

The HTK Book

Steve Young
Dan Kershaw
Julian Odell
Dave Ollason
Valtcho Valtchev
Phil Woodland

The HTK Book (for HTK Version 3.1)

©COPYRIGHT 1995-1999 Microsoft Corporation.

All Rights Reserved

First published December 1995

Reprinted March 1996

Revised for HTK Version 2.1 March 1997

Revised for HTK Version 2.2 January 1999

Revised for HTK Version 3.0 July 2000

Contents

I	Tutorial Overview	1
1	The Fundamentals of HTK	2
1.1	General Principles of HMMs	3
1.2	Isolated Word Recognition	3
1.3	Output Probability Specification	6
1.4	Baum-Welch Re-Estimation	6
1.5	Recognition and Viterbi Decoding	9
1.6	Continuous Speech Recognition	10
1.7	Speaker Adaptation	13
2	An Overview of the HTK Toolkit	14
2.1	HTK Software Architecture	14
2.2	Generic Properties of a HTK Tool	15
2.3	The Toolkit	16
2.3.1	Data Preparation Tools	16
2.3.2	Training Tools	16
2.3.3	Recognition Tools	19
2.3.4	Analysis Tool	19
2.4	Whats New In Version 2.2	19
2.4.1	Features Added To Version 2.1	20
3	A Tutorial Example of Using HTK	21
3.1	Data Preparation	22
3.1.1	Step 1 - the Task Grammar	22
3.1.2	Step 2 - the Dictionary	23
3.1.3	Step 3 - Recording the Data	25
3.1.4	Step 4 - Creating the Transcription Files	26
3.1.5	Step 5 - Coding the Data	28
3.2	Creating Monophone HMMs	29
3.2.1	Step 6 - Creating Flat Start Monophones	29
3.2.2	Step 7 - Fixing the Silence Models	32
3.2.3	Step 8 - Realigning the Training Data	33
3.3	Creating Tied-State Triphones	34
3.3.1	Step 9 - Making Triphones from Monophones	34
3.3.2	Step 10 - Making Tied-State Triphones	36
3.4	Recogniser Evaluation	39
3.4.1	Step 11 - Recognising the Test Data	39
3.5	Running the Recogniser Live	40
3.6	Adapting the HMMs	41
3.6.1	Step 12 - Preparation of the Adaptation Data	41
3.6.2	Step 13 - Generating the Transforms	42
3.6.3	Step 14 - Evaluation of the Adapted System	43
3.7	Summary	43

II	HTK in Depth	44
4	The Operating Environment	45
4.1	The Command Line	46
4.2	Script Files	46
4.3	Configuration Files	47
4.4	Standard Options	48
4.5	Error Reporting	49
4.6	Strings and Names	49
4.7	Memory Management	50
4.8	Input/Output via Pipes and Networks	51
4.9	Byte-swapping of HTK data files	51
4.10	Summary	52
5	Speech Input/Output	54
5.1	General Mechanism	54
5.2	Speech Signal Processing	56
5.3	Linear Prediction Analysis	58
5.4	Filterbank Analysis	59
5.5	Energy Measures	61
5.6	Delta and Acceleration Coefficients	61
5.7	Storage of Parameter Files	62
5.7.1	HTK Format Parameter Files	62
5.7.2	Esignal Format Parameter Files	64
5.8	Waveform File Formats	65
5.8.1	HTK File Format	65
5.8.2	Esignal File Format	65
5.8.3	TIMIT File Format	65
5.8.4	NIST File Format	65
5.8.5	SCRIBE File Format	66
5.8.6	SDES1 File Format	66
5.8.7	AIFF File Format	66
5.8.8	SUNAU8 File Format	66
5.8.9	OGI File Format	67
5.8.10	WAV File Format	67
5.8.11	ALIEN and NOHEAD File Formats	67
5.9	Direct Audio Input/Output	67
5.10	Multiple Input Streams	69
5.11	Vector Quantisation	70
5.12	Viewing Speech with HLIST	72
5.13	Copying and Coding using HCOPY	74
5.14	Version 1.5 Compatibility	75
5.15	Summary	76
6	Transcriptions and Label Files	79
6.1	Label File Structure	79
6.2	Label File Formats	80
6.2.1	HTK Label Files	80
6.2.2	ESPS Label Files	81
6.2.3	TIMIT Label Files	81
6.2.4	SCRIBE Label Files	81
6.3	Master Label Files	82
6.3.1	General Principles of MLFs	82
6.3.2	Syntax and Semantics	83
6.3.3	MLF Search	83
6.3.4	MLF Examples	84
6.4	Editing Label Files	86
6.5	Summary	89

7	HMM Definition Files	90
7.1	The HMM Parameters	91
7.2	Basic HMM Definitions	92
7.3	Macro Definitions	95
7.4	HMM Sets	99
7.5	Tied-Mixture Systems	103
7.6	Discrete Probability HMMs	103
7.7	Tee Models	105
7.8	Regression Class Trees for Adaptation	105
7.9	Binary Storage Format	106
7.10	The HMM Definition Language	107
8	HMM Parameter Estimation	111
8.1	Training Strategies	111
8.2	Initialisation using HINIT	114
8.3	Flat Starting with HCOMPV	117
8.4	Isolated Unit Re-Estimation using HREST	119
8.5	Embedded Training using HEREST	120
8.6	Single-Pass Retraining	123
8.7	Parameter Re-Estimation Formulae	123
8.7.1	Viterbi Training (HINIT)	124
8.7.2	Forward/Backward Probabilities	124
8.7.3	Single Model Reestimation(HREST)	126
8.7.4	Embedded Model Reestimation(HEREST)	127
9	HMM Adaptation	128
9.1	Model Adaptation using MLLR	129
9.1.1	Maximum Likelihood Linear Regression	129
9.1.2	MLLR and Regression Classes	129
9.1.3	Transform Model File Format	131
9.2	Model Adaptation using MAP	132
9.3	Using HEADAPT	134
9.4	MLLR Formulae	135
9.4.1	Estimation of the Mean Transformation Matrix	136
9.4.2	Estimation of the Variance Transformation Matrix	137
10	HMM System Refinement	138
10.1	Using HHED	139
10.2	Constructing Context-Dependent Models	139
10.3	Parameter Tying and Item Lists	140
10.4	Data-Driven Clustering	142
10.5	Tree-Based Clustering	144
10.6	Mixture Incrementing	146
10.7	Regression Class Tree Construction	147
10.8	Miscellaneous Operations	148
11	Discrete and Tied-Mixture Models	149
11.1	Modelling Discrete Sequences	149
11.2	Using Discrete Models with Speech	150
11.3	Tied Mixture Systems	152
11.4	Parameter Smoothing	154
12	Networks, Dictionaries and Language Models	155
12.1	How Networks are Used	156
12.2	Word Networks and Standard Lattice Format	157
12.3	Building a Word Network with HPARSE	159
12.4	Bigram Language Models	161
12.5	Building a Word Network with HBUILD	163
12.6	Testing a Word Network using HSGEN	164
12.7	Constructing a Dictionary	165

12.8	Word Network Expansion	167
12.9	Other Kinds of Recognition System	170
13	Decoding	172
13.1	Decoder Operation	172
13.2	Decoder Organisation	174
13.3	Recognition using Test Databases	176
13.4	Evaluating Recognition Results	177
13.5	Generating Forced Alignments	180
13.6	Decoding and Adaptation	181
13.6.1	Recognition with Adapted HMMs	181
13.6.2	Unsupervised Adaptation	181
13.7	Recognition using Direct Audio Input	182
13.8	N-Best Lists and Lattices	183
III	Reference Section	185
14	The HTK Tools	186
14.1	HBuild	187
14.1.1	Function	187
14.1.2	Use	187
14.1.3	Tracing	188
14.2	HCompV	189
14.2.1	Function	189
14.2.2	Use	189
14.2.3	Tracing	190
14.3	HCopy	191
14.3.1	Function	191
14.3.2	Use	191
14.3.3	Trace Output	193
14.4	HDMan	194
14.4.1	Function	194
14.4.2	Use	195
14.4.3	Tracing	196
14.5	HEAdapt	197
14.5.1	Function	197
14.5.2	Use	197
14.5.3	Tracing	198
14.6	HERest	200
14.6.1	Function	200
14.6.2	Use	200
14.6.3	Tracing	202
14.7	HHEd	203
14.7.1	Function	203
14.7.2	Use	210
14.7.3	Tracing	210
14.8	HInit	212
14.8.1	Function	212
14.8.2	Use	212
14.8.3	Tracing	213
14.9	HLEd	214
14.9.1	Function	214
14.9.2	Use	215
14.9.3	Tracing	216
14.10	HList	217
14.10.1	Function	217
14.10.2	Use	217
14.10.3	Tracing	217

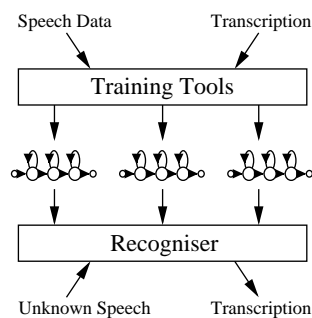
14.11HLStats	218
14.11.1 Function	218
14.11.2 Bigram Generation	218
14.11.3 Use	219
14.11.4 Tracing	219
14.12HParse	220
14.12.1 Function	220
14.12.2 Network Definition	220
14.12.3 Compatibility Mode	222
14.12.4 Use	222
14.12.5 Tracing	223
14.13HQuant	224
14.13.1 Function	224
14.13.2 VQ Codebook Format	224
14.13.3 Use	224
14.13.4 Tracing	225
14.14HRest	226
14.14.1 Function	226
14.14.2 Use	226
14.14.3 Tracing	227
14.15HResults	228
14.15.1 Function	228
14.15.2 Use	229
14.15.3 Tracing	230
14.16HSGen	231
14.16.1 Function	231
14.16.2 Use	231
14.16.3 Tracing	231
14.17HSLab	232
14.17.1 Function	232
14.17.2 Use	233
14.17.3 Tracing	235
14.18HSmooth	236
14.18.1 Function	236
14.18.2 Use	236
14.18.3 Tracing	237
14.19HVite	238
14.19.1 Function	238
14.19.2 Use	238
14.19.3 Tracing	240
15 Configuration Variables	241
15.1 Configuration Variables used in Library Modules	241
15.2 Configuration Variables used in Tools	244
16 Error and Warning Codes	245
16.1 Generic Errors	245
16.2 Summary of Errors by Tool and Module	246
17 HTK Standard Lattice Format (SLF)	261
17.1 SLF Files	261
17.2 Format	261
17.3 Syntax	262
17.4 Field Types	262
17.5 Example SLF file	263

Part I

Tutorial Overview

Chapter 1

The Fundamentals of HTK



HTK is a toolkit for building Hidden Markov Models (HMMs). HMMs can be used to model any time series and the core of HTK is similarly general-purpose. However, HTK is primarily designed for building HMM-based speech processing tools, in particular recognisers. Thus, much of the infrastructure support in HTK is dedicated to this task. As shown in the picture alongside, there are two major processing stages involved. Firstly, the HTK training tools are used to estimate the parameters of a set of HMMs using training utterances and their associated transcriptions. Secondly, unknown utterances are transcribed using the HTK recognition tools.

The main body of this book is mostly concerned with the mechanics of these two processes. However, before launching into detail it is necessary to understand some of the basic principles of HMMs. It is also helpful to have an overview of the toolkit and to have some appreciation of how training and recognition in HTK is organised.

This first part of the book attempts to provide this information. In this chapter, the basic ideas of HMMs and their use in speech recognition are introduced. The following chapter then presents a brief overview of HTK and, for users of older versions, it highlights the main differences in version 2.0 and later. Finally in this tutorial part of the book, chapter 3 describes how a HMM-based speech recogniser can be built using HTK. It does this by describing the construction of a simple small vocabulary continuous speech recogniser.

The second part of the book then revisits the topics skimmed over here and discusses each in detail. This can be read in conjunction with the third and final part of the book which provides a reference manual for HTK. This includes a description of each tool, summaries of the various parameters used to configure HTK and a list of the error messages that it generates when things go wrong.

Finally, note that this book is concerned only with HTK as a tool-kit. It does not provide information for using the HTK libraries as a programming environment.

1.1 General Principles of HMMs

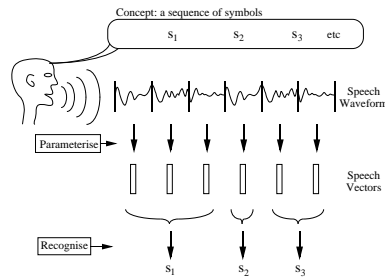


Fig. 1.1 Message Encoding/Decoding

Speech recognition systems generally assume that the speech signal is a realisation of some message encoded as a sequence of one or more symbols (see Fig. 1.1). To effect the reverse operation of recognising the underlying symbol sequence given a spoken utterance, the continuous speech waveform is first converted to a sequence of equally spaced discrete parameter vectors. This sequence of parameter vectors is assumed to form an exact representation of the speech waveform on the basis that for the duration covered by a single vector (typically 10ms or so), the speech waveform can be regarded as being stationary. Although this is not strictly true, it is a reasonable approximation. Typical parametric representations in common use are smoothed spectra or linear prediction coefficients plus various other representations derived from these.

The rôle of the recogniser is to effect a mapping between sequences of speech vectors and the wanted underlying symbol sequences. Two problems make this very difficult. Firstly, the mapping from symbols to speech is not one-to-one since different underlying symbols can give rise to similar speech sounds. Furthermore, there are large variations in the realised speech waveform due to speaker variability, mood, environment, etc. Secondly, the boundaries between symbols cannot be identified explicitly from the speech waveform. Hence, it is not possible to treat the speech waveform as a sequence of concatenated static patterns.

The second problem of not knowing the word boundary locations can be avoided by restricting the task to isolated word recognition. As shown in Fig. 1.2, this implies that the speech waveform corresponds to a single underlying symbol (e.g. word) chosen from a fixed vocabulary. Despite the fact that this simpler problem is somewhat artificial, it nevertheless has a wide range of practical applications. Furthermore, it serves as a good basis for introducing the basic ideas of HMM-based recognition before dealing with the more complex continuous speech case. Hence, isolated word recognition using HMMs will be dealt with first.

1.2 Isolated Word Recognition

Let each spoken word be represented by a sequence of speech vectors or *observations* \mathbf{O} , defined as

$$\mathbf{O} = \mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_T \quad (1.1)$$

where \mathbf{o}_t is the speech vector observed at time t . The isolated word recognition problem can then be regarded as that of computing

$$\arg \max_i \{P(w_i | \mathbf{O})\} \quad (1.2)$$

where w_i is the i 'th vocabulary word. This probability is not computable directly but using Bayes' Rule gives

$$P(w_i | \mathbf{O}) = \frac{P(\mathbf{O} | w_i) P(w_i)}{P(\mathbf{O})} \quad (1.3)$$

Thus, for a given set of prior probabilities $P(w_i)$, the most probable spoken word depends only on the likelihood $P(\mathbf{O} | w_i)$. Given the dimensionality of the observation sequence \mathbf{O} , the direct estimation of the joint conditional probability $P(\mathbf{o}_1, \mathbf{o}_2, \dots | w_i)$ from examples of spoken words is not practicable. However, if a parametric model of word production such as a Markov model is

assumed, then estimation from data is possible since the problem of estimating the class conditional observation densities $P(\mathbf{O}|w_i)$ is replaced by the much simpler problem of estimating the Markov model parameters.

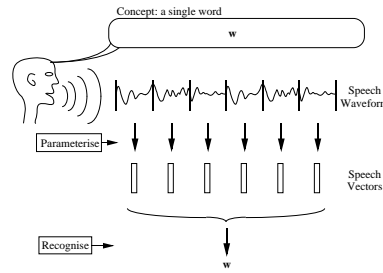


Fig. 1.2 Isolated Word Problem

In HMM based speech recognition, it is assumed that the sequence of observed speech vectors corresponding to each word is generated by a Markov model as shown in Fig. 1.3. A Markov model is a finite state machine which changes state once every time unit and each time t that a state j is entered, a speech vector \mathbf{o}_t is generated from the probability density $b_j(\mathbf{o}_t)$. Furthermore, the transition from state i to state j is also probabilistic and is governed by the discrete probability a_{ij} . Fig. 1.3 shows an example of this process where the six state model moves through the state sequence $X = 1, 2, 2, 3, 4, 4, 5, 6$ in order to generate the sequence \mathbf{o}_1 to \mathbf{o}_6 . Notice that in HTK, the entry and exit states of a HMM are non-emitting. This is to facilitate the construction of composite models as explained in more detail later.

The joint probability that \mathbf{O} is generated by the model M moving through the state sequence X is calculated simply as the product of the transition probabilities and the output probabilities. So for the state sequence X in Fig. 1.3

$$P(\mathbf{O}, X|M) = a_{12}b_2(\mathbf{o}_1)a_{22}b_2(\mathbf{o}_2)a_{23}b_3(\mathbf{o}_3) \dots \quad (1.4)$$

However, in practice, only the observation sequence \mathbf{O} is known and the underlying state sequence X is hidden. This is why it is called a *Hidden Markov Model*.

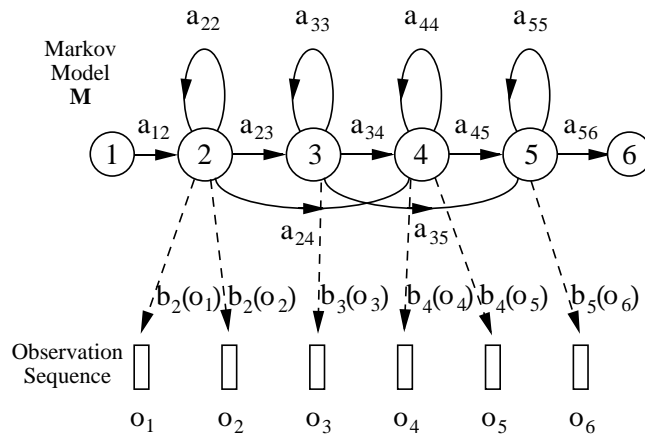


Fig. 1.3 The Markov Generation Model

Given that X is unknown, the required likelihood is computed by summing over all possible state sequences $X = x(1), x(2), x(3), \dots, x(T)$, that is

$$P(\mathbf{O}|M) = \sum_X a_{x(0)x(1)} \prod_{t=1}^T b_{x(t)}(\mathbf{o}_t) a_{x(t)x(t+1)} \quad (1.5)$$

where $x(0)$ is constrained to be the model entry state and $x(T + 1)$ is constrained to be the model exit state.

As an alternative to equation 1.5, the likelihood can be approximated by only considering the most likely state sequence, that is

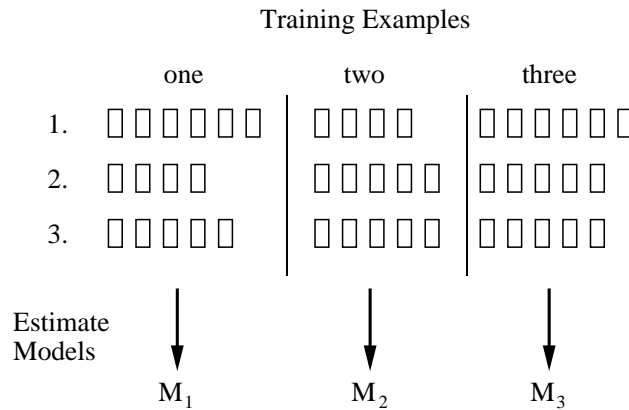
$$\hat{P}(\mathbf{O}|M) = \max_X \left\{ a_{x(0)x(1)} \prod_{t=1}^T b_{x(t)}(\mathbf{o}_t) a_{x(t)x(t+1)} \right\} \quad (1.6)$$

Although the direct computation of equations 1.5 and 1.6 is not tractable, simple recursive procedures exist which allow both quantities to be calculated very efficiently. Before going any further, however, notice that if equation 1.2 is computable then the recognition problem is solved. Given a set of models M_i corresponding to words w_i , equation 1.2 is solved by using 1.3 and assuming that

$$P(\mathbf{O}|w_i) = P(\mathbf{O}|M_i). \quad (1.7)$$

All this, of course, assumes that the parameters $\{a_{ij}\}$ and $\{b_j(\mathbf{o}_t)\}$ are known for each model M_i . Herein lies the elegance and power of the HMM framework. Given a set of training examples corresponding to a particular model, the parameters of that model can be determined automatically by a robust and efficient re-estimation procedure. Thus, provided that a sufficient number of representative examples of each word can be collected then a HMM can be constructed which implicitly models all of the many sources of variability inherent in real speech. Fig. 1.4 summarises the use of HMMs for isolated word recognition. Firstly, a HMM is trained for each vocabulary word using a number of examples of that word. In this case, the vocabulary consists of just three words: “one”, “two” and “three”. Secondly, to recognise some unknown word, the likelihood of each model generating that word is calculated and the most likely model identifies the word.

(a) Training



(b) Recognition

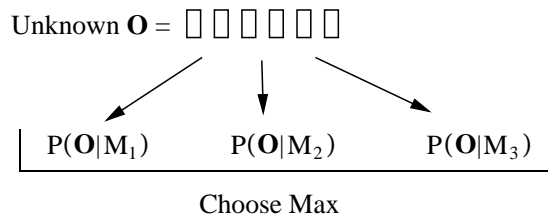


Fig. 1.4 Using HMMs for Isolated Word Recognition

1.3 Output Probability Specification

Before the problem of parameter estimation can be discussed in more detail, the form of the output distributions $\{b_j(\mathbf{o}_t)\}$ needs to be made explicit. HTK is designed primarily for modelling continuous parameters using continuous density multivariate output distributions. It can also handle observation sequences consisting of discrete symbols in which case, the output distributions are discrete probabilities. For simplicity, however, the presentation in this chapter will assume that continuous density distributions are being used. The minor differences that the use of discrete probabilities entail are noted in chapter 7 and discussed in more detail in chapter 11.

In common with most other continuous density HMM systems, HTK represents output distributions by Gaussian Mixture Densities. In HTK, however, a further generalisation is made. HTK allows each observation vector at time t to be split into a number of S independent data streams \mathbf{o}_{st} . The formula for computing $b_j(\mathbf{o}_t)$ is then

$$b_j(\mathbf{o}_t) = \prod_{s=1}^S \left[\sum_{m=1}^{M_s} c_{j_{sm}} \mathcal{N}(\mathbf{o}_{st}; \boldsymbol{\mu}_{j_{sm}}, \boldsymbol{\Sigma}_{j_{sm}}) \right]^{\gamma_s} \quad (1.8)$$

where M_s is the number of mixture components in stream s , $c_{j_{sm}}$ is the weight of the m 'th component and $\mathcal{N}(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is a multivariate Gaussian with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$, that is

$$\mathcal{N}(\mathbf{o}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^n |\boldsymbol{\Sigma}|}} e^{-\frac{1}{2}(\mathbf{o}-\boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1}(\mathbf{o}-\boldsymbol{\mu})} \quad (1.9)$$

where n is the dimensionality of \mathbf{o} .

The exponent γ_s is a stream weight¹. It can be used to give a particular stream more emphasis, however, it can only be set manually. No current HTK training tools can estimate values for it.

Multiple data streams are used to enable separate modelling of multiple information sources. In HTK, the processing of streams is completely general. However, the speech input modules assume that the source data is split into at most 4 streams. Chapter 5 discusses this in more detail but for now it is sufficient to remark that the default streams are the basic parameter vector, first (delta) and second (acceleration) difference coefficients and log energy.

1.4 Baum-Welch Re-Estimation

To determine the parameters of a HMM it is first necessary to make a rough guess at what they might be. Once this is done, more accurate (in the maximum likelihood sense) parameters can be found by applying the so-called Baum-Welch re-estimation formulae.

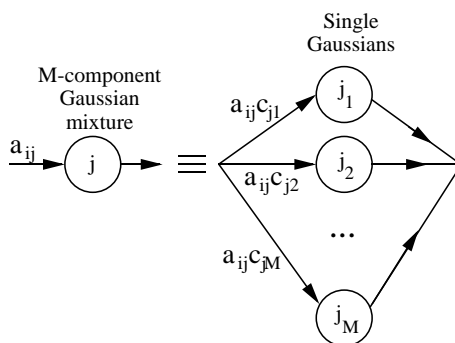


Fig. 1.5 Representing a Mixture

Chapter 8 gives the formulae used in HTK in full detail. Here the basis of the formulae will be presented in a very informal way. Firstly, it should be noted that the inclusion of multiple data streams does not alter matters significantly since each stream is considered to be statistically

¹often referred to as a codebook exponent.

independent. Furthermore, mixture components can be considered to be a special form of sub-state in which the transition probabilities are the mixture weights (see Fig. 1.5).

Thus, the essential problem is to estimate the means and variances of a HMM in which each state output distribution is a single component Gaussian, that is

$$b_j(\mathbf{o}_t) = \frac{1}{\sqrt{(2\pi)^n |\boldsymbol{\Sigma}_j|}} e^{-\frac{1}{2}(\mathbf{o}_t - \boldsymbol{\mu}_j)' \boldsymbol{\Sigma}_j^{-1} (\mathbf{o}_t - \boldsymbol{\mu}_j)} \quad (1.10)$$

If there was just one state j in the HMM, this parameter estimation would be easy. The maximum likelihood estimates of $\boldsymbol{\mu}_j$ and $\boldsymbol{\Sigma}_j$ would be just the simple averages, that is

$$\hat{\boldsymbol{\mu}}_j = \frac{1}{T} \sum_{t=1}^T \mathbf{o}_t \quad (1.11)$$

and

$$\hat{\boldsymbol{\Sigma}}_j = \frac{1}{T} \sum_{t=1}^T (\mathbf{o}_t - \boldsymbol{\mu}_j)(\mathbf{o}_t - \boldsymbol{\mu}_j)' \quad (1.12)$$

In practice, of course, there are multiple states and there is no direct assignment of observation vectors to individual states because the underlying state sequence is unknown. Note, however, that if some approximate assignment of vectors to states could be made then equations 1.11 and 1.12 could be used to give the required initial values for the parameters. Indeed, this is exactly what is done in the HTK tool called HINIT. HINIT first divides the training observation vectors equally amongst the model states and then uses equations 1.11 and 1.12 to give initial values for the mean and variance of each state. It then finds the maximum likelihood state sequence using the Viterbi algorithm described below, reassigns the observation vectors to states and then uses equations 1.11 and 1.12 again to get better initial values. This process is repeated until the estimates do not change.

Since the full likelihood of each observation sequence is based on the summation of all possible state sequences, each observation vector \mathbf{o}_t contributes to the computation of the maximum likelihood parameter values for each state j . In other words, instead of assigning each observation vector to a specific state as in the above approximation, each observation is assigned to every state in proportion to the probability of the model being in that state when the vector was observed. Thus, if $L_j(t)$ denotes the probability of being in state j at time t then the equations 1.11 and 1.12 given above become the following weighted averages

$$\hat{\boldsymbol{\mu}}_j = \frac{\sum_{t=1}^T L_j(t) \mathbf{o}_t}{\sum_{t=1}^T L_j(t)} \quad (1.13)$$

and

$$\hat{\boldsymbol{\Sigma}}_j = \frac{\sum_{t=1}^T L_j(t) (\mathbf{o}_t - \boldsymbol{\mu}_j)(\mathbf{o}_t - \boldsymbol{\mu}_j)'}{\sum_{t=1}^T L_j(t)} \quad (1.14)$$

where the summations in the denominators are included to give the required normalisation.

Equations 1.13 and 1.14 are the Baum-Welch re-estimation formulae for the means and covariances of a HMM. A similar but slightly more complex formula can be derived for the transition probabilities (see chapter 8).

Of course, to apply equations 1.13 and 1.14, the probability of state occupation $L_j(t)$ must be calculated. This is done efficiently using the so-called *Forward-Backward* algorithm. Let the forward probability² $\alpha_j(t)$ for some model M with N states be defined as

$$\alpha_j(t) = P(\mathbf{o}_1, \dots, \mathbf{o}_t, x(t) = j | M). \quad (1.15)$$

That is, $\alpha_j(t)$ is the joint probability of observing the first t speech vectors and being in state j at time t . This forward probability can be efficiently calculated by the following recursion

$$\alpha_j(t) = \left[\sum_{i=2}^{N-1} \alpha_i(t-1) a_{ij} \right] b_j(\mathbf{o}_t). \quad (1.16)$$

² Since the output distributions are densities, these are not really probabilities but it is a convenient fiction.

This recursion depends on the fact that the probability of being in state j at time t and seeing observation \mathbf{o}_t can be deduced by summing the forward probabilities for all possible predecessor states i weighted by the transition probability a_{ij} . The slightly odd limits are caused by the fact that states 1 and N are non-emitting³. The initial conditions for the above recursion are

$$\alpha_1(1) = 1 \quad (1.17)$$

$$\alpha_j(1) = a_{1j}b_j(\mathbf{o}_1) \quad (1.18)$$

for $1 < j < N$ and the final condition is given by

$$\alpha_N(T) = \sum_{i=2}^{N-1} \alpha_i(T)a_{iN}. \quad (1.19)$$

Notice here that from the definition of $\alpha_j(t)$,

$$P(\mathbf{O}|M) = \alpha_N(T). \quad (1.20)$$

Hence, the calculation of the forward probability also yields the total likelihood $P(\mathbf{O}|M)$.

The backward probability $\beta_j(t)$ is defined as

$$\beta_j(t) = P(\mathbf{o}_{t+1}, \dots, \mathbf{o}_T | x(t) = j, M). \quad (1.21)$$

As in the forward case, this backward probability can be computed efficiently using the following recursion

$$\beta_i(t) = \sum_{j=2}^{N-1} a_{ij}b_j(\mathbf{o}_{t+1})\beta_j(t+1) \quad (1.22)$$

with initial condition given by

$$\beta_i(T) = a_{iN} \quad (1.23)$$

for $1 < i < N$ and final condition given by

$$\beta_1(1) = \sum_{j=2}^{N-1} a_{1j}b_j(\mathbf{o}_1)\beta_j(1). \quad (1.24)$$

Notice that in the definitions above, the forward probability is a joint probability whereas the backward probability is a conditional probability. This somewhat asymmetric definition is deliberate since it allows the probability of state occupation to be determined by taking the product of the two probabilities. From the definitions,

$$\alpha_j(t)\beta_j(t) = P(\mathbf{O}, x(t) = j | M). \quad (1.25)$$

Hence,

$$\begin{aligned} L_j(t) &= P(x(t) = j | \mathbf{O}, M) \\ &= \frac{P(\mathbf{O}, x(t) = j | M)}{P(\mathbf{O} | M)} \\ &= \frac{1}{P} \alpha_j(t)\beta_j(t) \end{aligned} \quad (1.26)$$

where $P = P(\mathbf{O} | M)$.

All of the information needed to perform HMM parameter re-estimation using the Baum-Welch algorithm is now in place. The steps in this algorithm may be summarised as follows

1. For every parameter vector/matrix requiring re-estimation, allocate storage for the numerator and denominator summations of the form illustrated by equations 1.13 and 1.14. These storage locations are referred to as *accumulators*⁴.

³ To understand equations involving a non-emitting state at time t , the time should be thought of as being $t - \delta t$ if it is an entry state, and $t + \delta t$ if it is an exit state. This becomes important when HMMs are connected together in sequence so that transitions across non-emitting states take place *between frames*.

⁴ Note that normally the summations in the denominators of the re-estimation formulae are identical across the parameter sets of a given state and therefore only a single common storage location for the denominators is required and it need only be calculated once. However, HTK supports a generalised parameter tying mechanism which can result in the denominator summations being different. Hence, in HTK the denominator summations are always stored and calculated individually for each distinct parameter vector or matrix.

2. Calculate the forward and backward probabilities for all states j and times t .
3. For each state j and time t , use the probability $L_j(t)$ and the current observation vector \mathbf{o}_t to update the accumulators for that state.
4. Use the final accumulator values to calculate new parameter values.
5. If the value of $P = P(\mathbf{O}|M)$ for this iteration is not higher than the value at the previous iteration then stop, otherwise repeat the above steps using the new re-estimated parameter values.

All of the above assumes that the parameters for a HMM are re-estimated from a single observation sequence, that is a single example of the spoken word. In practice, many examples are needed to get good parameter estimates. However, the use of multiple observation sequences adds no additional complexity to the algorithm. Steps 2 and 3 above are simply repeated for each distinct training sequence.

One final point that should be mentioned is that the computation of the forward and backward probabilities involves taking the product of a large number of probabilities. In practice, this means that the actual numbers involved become very small. Hence, to avoid numerical problems, the forward-backward computation is computed in HTK using log arithmetic.

The HTK program which implements the above algorithm is called HREST. In combination with the tool HINIT for estimating initial values mentioned earlier, HREST allows isolated word HMMs to be constructed from a set of training examples using Baum-Welch re-estimation.

1.5 Recognition and Viterbi Decoding

The previous section has described the basic ideas underlying HMM parameter re-estimation using the Baum-Welch algorithm. In passing, it was noted that the efficient recursive algorithm for computing the forward probability also yielded as a by-product the total likelihood $P(\mathbf{O}|M)$. Thus, this algorithm could also be used to find the model which yields the maximum value of $P(\mathbf{O}|M_i)$, and hence, it could be used for recognition.

In practice, however, it is preferable to base recognition on the maximum likelihood state sequence since this generalises easily to the continuous speech case whereas the use of the total probability does not. This likelihood is computed using essentially the same algorithm as the forward probability calculation except that the summation is replaced by a maximum operation. For a given model M , let $\phi_j(t)$ represent the maximum likelihood of observing speech vectors \mathbf{o}_1 to \mathbf{o}_t and being in state j at time t . This partial likelihood can be computed efficiently using the following recursion (cf. equation 1.16)

$$\phi_j(t) = \max_i \{ \phi_i(t-1) a_{ij} \} b_j(\mathbf{o}_t). \quad (1.27)$$

where

$$\phi_1(1) = 1 \quad (1.28)$$

$$\phi_j(1) = a_{1j} b_j(\mathbf{o}_1) \quad (1.29)$$

for $1 < j < N$. The maximum likelihood $\hat{P}(\mathbf{O}|M)$ is then given by

$$\phi_N(T) = \max_i \{ \phi_i(T) a_{iN} \} \quad (1.30)$$

As for the re-estimation case, the direct computation of likelihoods leads to underflow, hence, log likelihoods are used instead. The recursion of equation 1.27 then becomes

$$\psi_j(t) = \max_i \{ \psi_i(t-1) + \log(a_{ij}) \} + \log(b_j(\mathbf{o}_t)). \quad (1.31)$$

This recursion forms the basis of the so-called Viterbi algorithm. As shown in Fig. 1.6, this algorithm can be visualised as finding the best path through a matrix where the vertical dimension represents the states of the HMM and the horizontal dimension represents the frames of speech (i.e. time). Each large dot in the picture represents the log probability of observing that frame at that time and each arc between dots corresponds to a log transition probability. The log probability of any path is computed simply by summing the log transition probabilities and the log output probabilities along that path. The paths are grown from left-to-right column-by-column. At time t , each partial path $\psi_i(t-1)$ is known for all states i , hence equation 1.31 can be used to compute $\psi_j(t)$ thereby extending the partial paths by one time frame.

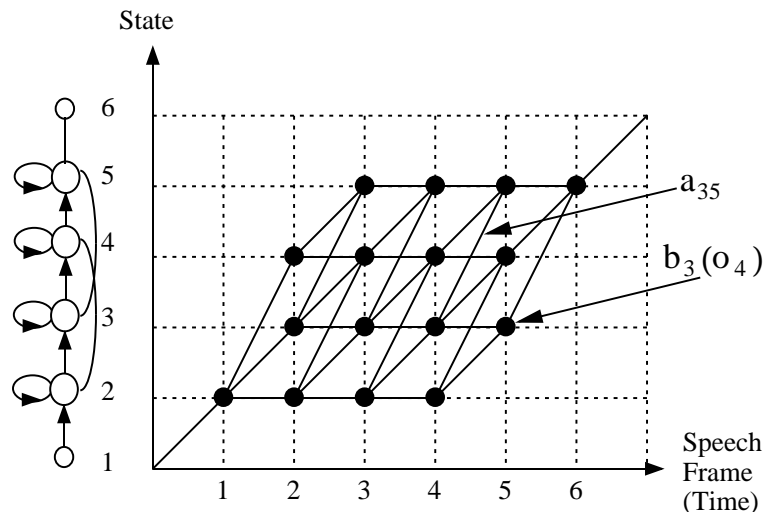


Fig. 1.6 The Viterbi Algorithm for Isolated Word Recognition

This concept of a path is extremely important and it is generalised below to deal with the continuous speech case.

This completes the discussion of isolated word recognition using HMMs. There is no HTK tool which implements the above Viterbi algorithm directly. Instead, a tool called HVITE is provided which along with its supporting libraries, HNET and HREC, is designed to handle continuous speech. Since this recogniser is syntax directed, it can also perform isolated word recognition as a special case. This is discussed in more detail below.

1.6 Continuous Speech Recognition

Returning now to the conceptual model of speech production and recognition exemplified by Fig. 1.1, it should be clear that the extension to continuous speech simply involves connecting HMMs together in sequence. Each model in the sequence corresponds directly to the assumed underlying symbol. These could be either whole words for so-called *connected speech recognition* or sub-words such as phonemes for *continuous speech recognition*. The reason for including the non-emitting entry and exit states should now be evident, these states provide the *glue* needed to join models together.

There are, however, some practical difficulties to overcome. The training data for continuous speech must consist of continuous utterances and, in general, the boundaries dividing the segments of speech corresponding to each underlying sub-word model in the sequence will not be known. In practice, it is usually feasible to mark the boundaries of a small amount of data by hand. All of the segments corresponding to a given model can then be extracted and the *isolated word* style of training described above can be used. However, the amount of data obtainable in this way is usually very limited and the resultant models will be poor estimates. Furthermore, even if there was a large amount of data, the boundaries imposed by hand-marking may not be optimal as far as the HMMs are concerned. Hence, in HTK the use of HINIT and HREST for initialising sub-word models is regarded as a *bootstrap* operation⁵. The main training phase involves the use of a tool called HEREST which does *embedded training*.

Embedded training uses the same Baum-Welch procedure as for the isolated case but rather than training each model individually all models are trained in parallel. It works in the following steps:

1. Allocate and zero accumulators for all parameters of all HMMs.
2. Get the next training utterance.

⁵ They can even be avoided altogether by using a *flat start* as described in section 8.3.

3. Construct a composite HMM by joining in sequence the HMMs corresponding to the symbol transcription of the training utterance.
4. Calculate the forward and backward probabilities for the composite HMM. The inclusion of intermediate non-emitting states in the composite model requires some changes to the computation of the forward and backward probabilities but these are only minor. The details are given in chapter 8.
5. Use the forward and backward probabilities to compute the probabilities of state occupation at each time frame and update the accumulators in the usual way.
6. Repeat from 2 until all training utterances have been processed.
7. Use the accumulators to calculate new parameter estimates for all of the HMMs.

These steps can then all be repeated as many times as is necessary to achieve the required convergence. Notice that although the location of symbol boundaries in the training data is not required (or wanted) for this procedure, the symbolic transcription of each training utterance is needed.

Whereas the extensions needed to the Baum-Welch procedure for training sub-word models are relatively minor⁶, the corresponding extensions to the Viterbi algorithm are more substantial.

In HTK, an alternative formulation of the Viterbi algorithm is used called the *Token Passing Model*⁷. In brief, the token passing model makes the concept of a state alignment path explicit. Imagine each state j of a HMM at time t holds a single moveable token which contains, amongst other information, the partial log probability $\psi_j(t)$. This token then represents a partial match between the observation sequence \mathbf{o}_1 to \mathbf{o}_t and the model subject to the constraint that the model is in state j at time t . The path extension algorithm represented by the recursion of equation 1.31 is then replaced by the equivalent *token passing algorithm* which is executed at each time frame t . The key steps in this algorithm are as follows

1. Pass a copy of every token in state i to all connecting states j , incrementing the log probability of the copy by $\log[a_{ij}] + \log[b_j(\mathbf{o}(t))]$.
2. Examine the tokens in every state and discard all but the token with the highest probability.

In practice, some modifications are needed to deal with the non-emitting states but these are straightforward if the tokens in entry states are assumed to represent paths extended to time $t - \delta t$ and tokens in exit states are assumed to represent paths extended to time $t + \delta t$.

The point of using the Token Passing Model is that it extends very simply to the continuous speech case. Suppose that the allowed sequence of HMMs is defined by a finite state network. For example, Fig. 1.7 shows a simple network in which each word is defined as a sequence of phoneme-based HMMs and all of the words are placed in a loop. In this network, the oval boxes denote HMM instances and the square boxes denote *word-end* nodes. This composite network is essentially just a single large HMM and the above Token Passing algorithm applies. The only difference now is that more information is needed beyond the log probability of the best token. When the best token reaches the end of the speech, the route it took through the network must be known in order to recover the recognised sequence of models.

⁶ In practice, a good deal of extra work is needed to achieve efficient operation on large training databases. For example, the HEREST tool includes facilities for pruning on both the forward and backward passes and parallel operation on a network of machines.

⁷ See "Token Passing: a Conceptual Model for Connected Speech Recognition Systems", SJ Young, NH Russell and JHS Thornton, CUED Technical Report F-INFENG/TR38, Cambridge University, 1989. Available by anonymous ftp from `svr-ftp.eng.cam.ac.uk`.

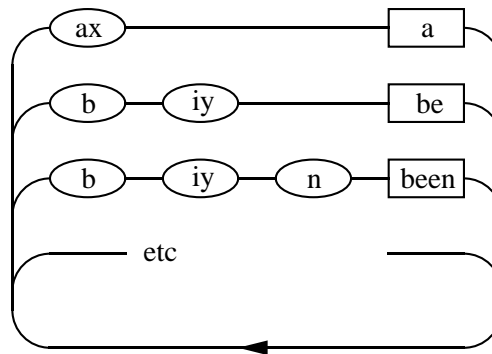


Fig. 1.7 Recognition Network for Continuously Spoken Word Recognition

The history of a token's route through the network may be recorded efficiently as follows. Every token carries a pointer called a *word end link*. When a token is propagated from the exit state of a word (indicated by passing through a word-end node) to the entry state of another, that transition represents a potential word boundary. Hence a record called a *Word Link Record* is generated in which is stored the identity of the word from which the token has just emerged and the current value of the token's link. The token's actual link is then replaced by a pointer to the newly created WLR. Fig. 1.8 illustrates this process.

Once all of the unknown speech has been processed, the WLRs attached to the link of the best matching token (i.e. the token with the highest log probability) can be traced back to give the best matching sequence of words. At the same time the positions of the word boundaries can also be extracted if required.

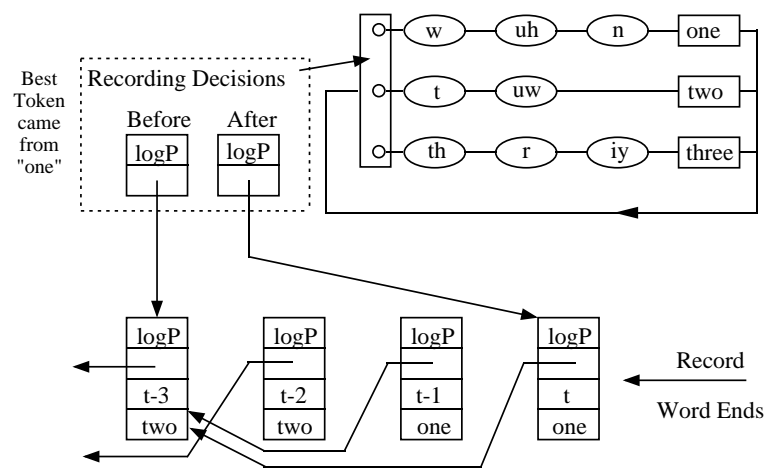


Fig. 1.8 Recording Word Boundary Decisions

The token passing algorithm for continuous speech has been described in terms of recording the word sequence only. If required, the same principle can be used to record decisions at the model and state level. Also, more than just the best token at each word boundary can be saved. This gives the potential for generating a lattice of hypotheses rather than just the single best hypothesis. Algorithms based on this idea are called *lattice N-best*. They are suboptimal because the use of a single token per state limits the number of different token histories that can be maintained. This limitation can be avoided by allowing each model state to hold multiple-tokens and regarding tokens as distinct if they come from different preceding words. This gives a class of algorithm called *word*

N-best which has been shown empirically to be comparable in performance to an optimal N-best algorithm.

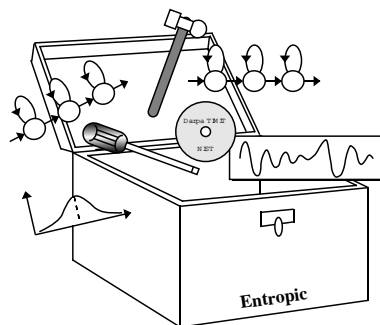
The above outlines the main idea of Token Passing as it is implemented within HTK. The algorithms are embedded in the library modules HNET and HREC and they may be invoked using the recogniser tool called HVITE. They provide single and multiple-token passing recognition, single-best output, lattice output, N-best lists, support for cross-word context-dependency, lattice rescoring and forced alignment.

1.7 Speaker Adaptation

Although the training and recognition techniques described previously can produce high performance recognition systems, these systems can be improved upon by customising the HMMs to the characteristics of a particular speaker. HTK provides the tools HEADAPT and HVITE to perform adaptation using a small amount of enrollment or adaptation data. The two tools differ in that HEADAPT performs offline supervised adaptation while HVITE recognises the adaptation data and uses the generated transcriptions to perform the adaptation. Generally, more robust adaptation is performed in a supervised mode, as provided by HEADAPT, but given an initial well trained model set, HVITE can still achieve noticeable improvements in performance. Full details of adaptation and how it is used in HTK can be found in Chapter 9.

Chapter 2

An Overview of the HTK Toolkit



The basic principles of HMM-based recognition were outlined in the previous chapter and a number of the key HTK tools have already been mentioned. This chapter describes the software architecture of a HTK tool. It then gives a brief outline of all the HTK tools and the way that they are used together to construct and test HMM-based recognisers. For the benefit of existing HTK users, the major changes in recent versions of HTK are listed. The following chapter will then illustrate the use of the HTK toolkit by working through a practical example of building a simple continuous speech recognition system.

2.1 HTK Software Architecture

Much of the functionality of HTK is built into the library modules. These modules ensure that every tool interfaces to the outside world in exactly the same way. They also provide a central resource of commonly used functions. Fig. 2.1 illustrates the software structure of a typical HTK tool and shows its input/output interfaces.

User input/output and interaction with the operating system is controlled by the library module `HSHELL` and all memory management is controlled by `HMEM`. Math support is provided by `HMATH` and the signal processing operations needed for speech analysis are in `HSIGP`. Each of the file types required by HTK has a dedicated interface module. `HLABEL` provides the interface for label files, `HLM` for language model files, `HNET` for networks and lattices, `HDICT` for dictionaries, `HVQ` for VQ codebooks and `HMODEL` for HMM definitions.

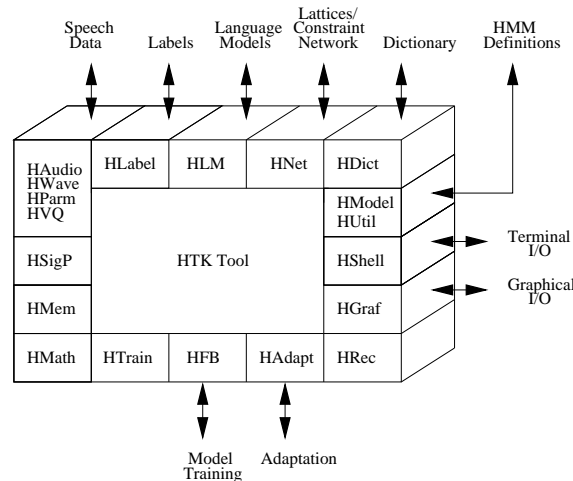


Fig. 2.1 Software Architecture

All speech input and output at the waveform level is via HWAVE and at the parameterised level via HPARM. As well as providing a consistent interface, HWAVE and HLABEL support multiple file formats allowing data to be imported from other systems. Direct audio input is supported by HAUDIO and simple interactive graphics is provided by HGRAF. HUTIL provides a number of utility routines for manipulating HMMs while HTRAIN and HFB contain support for the various HTK training tools. HADAPT provides support for the various HTK adaptation tools. Finally, HREC contains the main recognition processing functions.

As noted in the next section, fine control over the behaviour of these library modules is provided by setting configuration variables. Detailed descriptions of the functions provided by the library modules are given in the second part of this book and the relevant configuration variables are described as they arise. For reference purposes, a complete list is given in chapter 15.

2.2 Generic Properties of a HTK Tool

HTK tools are designed to run with a traditional command-line style interface. Each tool has a number of required arguments plus optional arguments. The latter are always prefixed by a minus sign. As an example, the following command would invoke the mythical HTK tool called HFOO

```
HFOO -T 1 -f 34.3 -a -s myfile file1 file2
```

This tool has two main arguments called `file1` and `file2` plus four optional arguments. Options are always introduced by a single letter option name followed where appropriate by the option value. The option value is always separated from the option name by a space. Thus, the value of the `-f` option is a real number, the value of the `-T` option is an integer number and the value of the `-s` option is a string. The `-a` option has no following value and it is used as a simple flag to enable or disable some feature of the tool. Options whose names are a capital letter have the same meaning across all tools. For example, the `-T` option is always used to control the trace output of a HTK tool.

In addition to command line arguments, the operation of a tool can be controlled by parameters stored in a configuration file. For example, if the command

```
HFOO -C config -f 34.3 -a -s myfile file1 file2
```

is executed, the tool HFOO will load the parameters stored in the configuration file `config` during its initialisation procedures. Configuration parameters can sometimes be used as an alternative to using command line arguments. For example, trace options can always be set within a configuration file. However, the main use of configuration files is to control the detailed behaviour of the library modules on which all HTK tools depend.

Although this style of command-line working may seem old-fashioned when compared to modern graphical user interfaces, it has many advantages. In particular, it makes it simple to write shell

scripts to control HTK tool execution. This is vital for performing large-scale system building and experimentation. Furthermore, defining all operations using text-based commands allows the details of system construction or experimental procedure to be recorded and documented.

Finally, note that a summary of the command line and options for any HTK tool can be obtained simply by executing the tool with no arguments.

2.3 The Toolkit

The HTK tools are best introduced by going through the processing steps involved in building a sub-word based continuous speech recogniser. As shown in Fig. 2.2, there are 4 main phases: data preparation, training, testing and analysis.

2.3.1 Data Preparation Tools

In order to build a set of HMMs, a set of speech data files and their associated transcriptions are required. Very often speech data will be obtained from database archives, typically on CD-ROMs. Before it can be used in training, it must be converted into the appropriate parametric form and any associated transcriptions must be converted to have the correct format and use the required phone or word labels. If the speech needs to be recorded, then the tool HSLAB can be used both to record the speech and to manually annotate it with any required transcriptions.

Although all HTK tools can parameterise waveforms *on-the-fly*, in practice it is usually better to parameterise the data just once. The tool HCOPY is used for this. As the name suggests, HCOPY is used to copy one or more source files to an output file. Normally, HCOPY copies the whole file, but a variety of mechanisms are provided for extracting segments of files and concatenating files. By setting the appropriate configuration variables, all input files can be converted to parametric form as they are read-in. Thus, simply copying each file in this manner performs the required encoding. The tool HLIST can be used to check the contents of any speech file and since it can also convert input on-the-fly, it can be used to check the results of any conversions before processing large quantities of data. Transcriptions will also need preparing. Typically the labels used in the original source transcriptions will not be exactly as required, for example, because of differences in the phone sets used. Also, HMM training might require the labels to be context-dependent. The tool HLED is a script-driven label editor which is designed to make the required transformations to label files. HLED can also output files to a single *Master Label File* MLF which is usually more convenient for subsequent processing. Finally on data preparation, HLSTATS can gather and display statistics on label files and where required, HQUANT can be used to build a VQ codebook in preparation for building discrete probability HMM system.

2.3.2 Training Tools

The second step of system building is to define the topology required for each HMM by writing a prototype definition. HTK allows HMMs to be built with any desired topology. HMM definitions can be stored externally as simple text files and hence it is possible to edit them with any convenient text editor. Alternatively, the standard HTK distribution includes a number of example HMM prototypes and a script to generate the most common topologies automatically. With the exception of the transition probabilities, all of the HMM parameters given in the prototype definition are ignored. The purpose of the prototype definition is only to specify the overall characteristics and topology of the HMM. The actual parameters will be computed later by the training tools. Sensible values for the transition probabilities must be given but the training process is very insensitive to these. An acceptable and simple strategy for choosing these probabilities is to make all of the transitions out of any state equally likely.

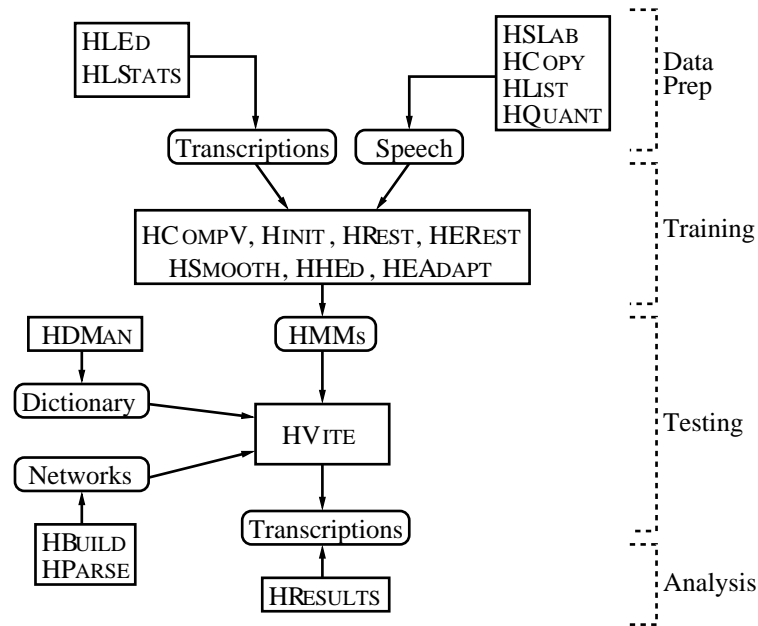


Fig. 2.2 HTK Processing Stages

The actual training process takes place in stages and it is illustrated in more detail in Fig. 2.3. Firstly, an initial set of models must be created. If there is some speech data available for which the location of the sub-word (i.e. phone) boundaries have been marked, then this can be used as *bootstrap data*. In this case, the tools HINIT and HREST provide *isolated word* style training using the fully labelled bootstrap data. Each of the required HMMs is generated individually. HINIT reads in all of the bootstrap training data and *cuts out* all of the examples of the required phone. It then iteratively computes an initial set of parameter values using a *segmental k-means* procedure. On the first cycle, the training data is uniformly segmented, each model state is matched with the corresponding data segments and then means and variances are estimated. If mixture Gaussian models are being trained, then a modified form of k-means clustering is used. On the second and successive cycles, the uniform segmentation is replaced by Viterbi alignment. The initial parameter values computed by HINIT are then further re-estimated by HREST. Again, the fully labelled bootstrap data is used but this time the segmental k-means procedure is replaced by the Baum-Welch re-estimation procedure described in the previous chapter. When no bootstrap data is available, a so-called *flat start* can be used. In this case all of the phone models are initialised to be identical and have state means and variances equal to the global speech mean and variance. The tool HCOMPV can be used for this.

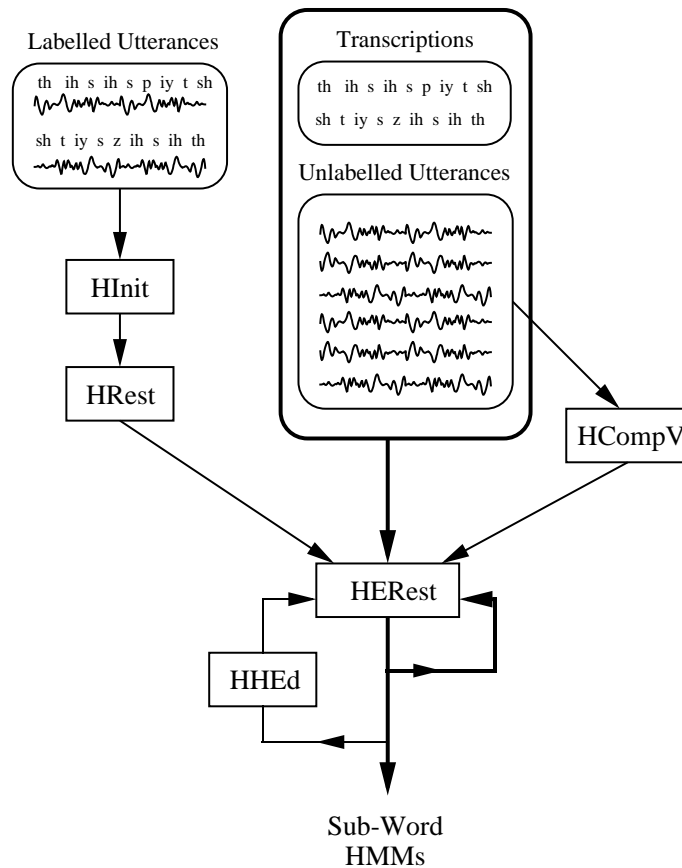


Fig. 2.3 Training Sub-word HMMs

Once an initial set of models has been created, the tool HEREST is used to perform *embedded training* using the entire training set. HEREST performs a single Baum-Welch re-estimation of the whole set of HMM phone models simultaneously. For each training utterance, the corresponding phone models are concatenated and then the forward-backward algorithm is used to accumulate the statistics of state occupation, means, variances, etc., for each HMM in the sequence. When all of the training data has been processed, the accumulated statistics are used to compute re-estimates of the HMM parameters. HEREST is the core HTK training tool. It is designed to process large databases, it has facilities for pruning to reduce computation and it can be run in parallel across a network of machines.

The philosophy of system construction in HTK is that HMMs should be refined incrementally. Thus, a typical progression is to start with a simple set of single Gaussian context-independent phone models and then iteratively refine them by expanding them to include context-dependency and use multiple mixture component Gaussian distributions. The tool HHED is a HMM definition editor which will clone models into context-dependent sets, apply a variety of parameter tyings and increment the number of mixture components in specified distributions. The usual process is to modify a set of HMMs in stages using HHED and then re-estimate the parameters of the modified set using HEREST after each stage. To improve performance for specific speakers the tools HEADAPT and HVITE can be used to adapt HMMs to better model the characteristics of particular speakers using a small amount of training or adaptation data. The end result of which is a speaker adapted system.

The single biggest problem in building context-dependent HMM systems is always data insufficiency. The more complex the model set, the more data is needed to make robust estimates of its parameters, and since data is usually limited, a balance must be struck between complexity and the available data. For continuous density systems, this balance is achieved by tying parameters together as mentioned above. Parameter tying allows data to be pooled so that the shared parameters can be robustly estimated. In addition to continuous density systems, HTK also supports

fully tied mixture systems and discrete probability systems. In these cases, the data insufficiency problem is usually addressed by smoothing the distributions and the tool HSMOOTH is used for this.

2.3.3 Recognition Tools

HTK provides a single recognition tool called HVITE which uses the token passing algorithm described in the previous chapter to perform Viterbi-based speech recognition. HVITE takes as input a network describing the allowable word sequences, a dictionary defining how each word is pronounced and a set of HMMs. It operates by converting the word network to a phone network and then attaching the appropriate HMM definition to each phone instance. Recognition can then be performed on either a list of stored speech files or on direct audio input. As noted at the end of the last chapter, HVITE can support cross-word triphones and it can run with multiple tokens to generate lattices containing multiple hypotheses. It can also be configured to rescore lattices and perform forced alignments.

The word networks needed to drive HVITE are usually either simple word loops in which any word can follow any other word or they are directed graphs representing a finite-state task grammar. In the former case, bigram probabilities are normally attached to the word transitions. Word networks are stored using the HTK standard lattice format. This is a text-based format and hence word networks can be created directly using a text-editor. However, this is rather tedious and hence HTK provides two tools to assist in creating word networks. Firstly, HBUILD allows sub-networks to be created and used within higher level networks. Hence, although the same low level notation is used, much duplication is avoided. Also, HBUILD can be used to generate word loops and it can also read in a backed-off bigram language model and modify the word loop transitions to incorporate the bigram probabilities. Note that the label statistics tool HLSTATS mentioned earlier can be used to generate a backed-off bigram language model.

As an alternative to specifying a word network directly, a higher level grammar notation can be used. This notation is based on the Extended Backus Naur Form (EBNF) used in compiler specification and it is compatible with the grammar specification language used in earlier versions of HTK. The tool HPARSE is supplied to convert this notation into the equivalent word network.

Whichever method is chosen to generate a word network, it is useful to be able to see examples of the *language* that it defines. The tool HSGEN is provided to do this. It takes as input a network and then randomly traverses the network outputting word strings. These strings can then be inspected to ensure that they correspond to what is required. HSGEN can also compute the empirical perplexity of the task.

Finally, the construction of large dictionaries can involve merging several sources and performing a variety of transformations on each sources. The dictionary management tool HDMAN is supplied to assist with this process.

2.3.4 Analysis Tool

Once the HMM-based recogniser has been built, it is necessary to evaluate its performance. This is usually done by using it to transcribe some pre-recorded test sentences and match the recogniser output with the correct reference transcriptions. This comparison is performed by a tool called HRESULTS which uses dynamic programming to align the two transcriptions and then count substitution, deletion and insertion errors. Options are provided to ensure that the algorithms and output formats used by HRESULTS are compatible with those used by the US National Institute of Standards and Technology (NIST). As well as global performance measures, HRESULTS can also provide speaker-by-speaker breakdowns, confusion matrices and time-aligned transcriptions. For word spotting applications, it can also compute *Figure of Merit* (FOM) scores and *Receiver Operating Curve* (ROC) information.

2.4 What's New In Version 2.2

This section lists the new features and refinements in HTK Version 2.2 compared to the preceding Version 2.1.

1. Speaker adaptation is now supported via the HEADAPT and HVITE tools, which adapt a current set of models to a new speaker and/or environment.

- HEADAPT performs offline supervised adaptation using maximum likelihood linear regression (MLLR) and/or maximum a-posteriori (MAP) adaptation.
- HVITE performs unsupervised adaptation using just MLLR.

Both tools can be used in a static mode, where all the data is presented prior to any adaptation, or in an incremental fashion.

2. Improved support for PC WAV files

In addition to 16-bit PCM linear, HTK can now read

- 8-bit CCITT mu-law
- 8-bit CCITT a-law
- 8-bit PCM linear

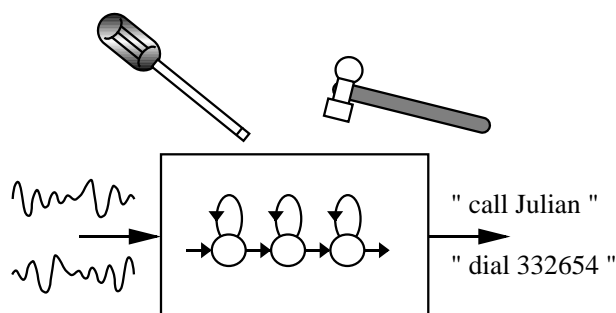
2.4.1 Features Added To Version 2.1

For the benefit of users of earlier versions of HTK this section lists the main changes in HTK Version 2.1 compared to the preceding Version 2.0.

1. The speech input handling has been partially re-designed and a new energy-based speech/silence detector has been incorporated into HPARM. The detector is robust yet flexible and can be configured through a number of configuration variables. Speech/silence detection can now be performed on waveform files. The calibration of speech/silence detector parameters is now accomplished by asking the user to speak an arbitrary sentence.
2. HPARM now allows random noise signal to be added to waveform data via the configuration parameter `ADDITHER`. This prevents numerical overflows which can occur with artificially created waveform data under some coding schemes.
3. HNET has been optimised for more efficient operation when performing forced alignments of utterances using HVITE. Further network optimisations tailored to biphone/triphone-based phone recognition have also been incorporated.
4. HVITE can now produce partial recognition hypothesis even when no tokens survive to the end of the network. This is accomplished by setting the HREC configuration parameter `FORCEOUT` to true.
5. Dictionary support has been extended to allow pronunciation probabilities to be associated with different pronunciations of the same word. At the same time, HVITE now allows the use of a pronunciation scale factor during recognition.
6. HTK now provides consistent support for reading and writing of HTK binary files (waveforms, binary MMFs, binary SLFs, HEREST accumulators) across different machine architectures incorporating automatic byte swapping. By default, all binary data files handled by the tools are now written/read in big-endian (`NONVAX`) byte order. The default behavior can be changed via the configuration parameters `NATURALREADORDER` and `NATURALWRITEORDER`.
7. HWAVE supports the reading of waveforms in Microsoft WAVE file format.
8. HAUDIO allows key-press control of live audio input.

Chapter 3

A Tutorial Example of Using HTK



This final chapter of the tutorial part of the book will describe the construction of a recogniser for simple voice dialling applications. This recogniser will be designed to recognise continuously spoken digit strings and a limited set of names. It is sub-word based so that adding a new name to the vocabulary involves only modification to the pronouncing dictionary and task grammar. The HMMs will be continuous density mixture Gaussian tied-state triphones with clustering performed using phonetic decision trees. Although the voice dialling task itself is quite simple, the system design is general-purpose and would be useful for a range of applications.

The system will be built from scratch even to the extent of recording training and test data using the HTK tool HSLAB. To make this tractable, the system will be speaker dependent¹, but the same design would be followed to build a speaker independent system. The only difference being that data would be required from a large number of speakers and there would be a consequential increase in model complexity.

Building a speech recogniser from scratch involves a number of inter-related subtasks and pedagogically it is not obvious what the best order is to present them. In the presentation here, the ordering is chronological so that in effect the text provides a recipe that could be followed to construct a similar system. The entire process is described in considerable detail in order to give a clear view of the range of functions that HTK addresses and thereby to motivate the rest of the book.

The HTK software distribution also contains an example of constructing a recognition system for the 1000 word ARPA Naval Resource Management Task. This is contained in the directory RMHTK of the HTK distribution. Further demonstration of HTK's capabilities can be found in the directory HTKDemo. Some example scripts that may be of assistance during the tutorial are available in the HTKTutorial directory.

At each step of the tutorial presented in this chapter, the user is advised to thoroughly read the entire section before executing the commands, and also to consult the reference section for each HTK tool being introduced (chapter 14), so that all command line options and arguments are clearly understood.

¹The final stage of the tutorial deals with adapting the speaker dependent models for new speakers

3.1 Data Preparation

The first stage of any recogniser development project is data preparation. Speech data is needed both for training and for testing. In the system to be built here, all of this speech will be recorded from scratch and to do this scripts are needed to prompt for each sentence. In the case of the test data, these prompt scripts will also provide the reference transcriptions against which the recogniser's performance can be measured and a convenient way to create them is to use the task grammar as a random generator. In the case of the training data, the prompt scripts will be used in conjunction with a pronunciation dictionary to provide the initial phone level transcriptions needed to start the HMM training process. Since the application requires that arbitrary names can be added to the recogniser, training data with good phonetic balance and coverage is needed. Here for convenience the prompt scripts needed for training are taken from the TIMIT acoustic-phonetic database.

It follows from the above that before the data can be recorded, a phone set must be defined, a dictionary must be constructed to cover both training and testing and a task grammar must be defined.

3.1.1 Step 1 - the Task Grammar

The goal of the system to be built here is to provide a voice-operated interface for phone dialling. Thus, the recogniser must handle digit strings and also personal name lists. Examples of typical inputs might be

Dial three three two six five four

Dial nine zero four one oh nine

Phone Woodland

Call Steve Young

HTK provides a grammar definition language for specifying simple task grammars such as this. It consists of a set of variable definitions followed by a regular expression describing the words to recognise. For the voice dialling application, a suitable grammar might be

```
$digit = ONE | TWO | THREE | FOUR | FIVE |
        SIX | SEVEN | EIGHT | NINE | OH | ZERO;
$name   = [ JOOP ] JANSEN |
          [ JULIAN ] ODELL |
          [ DAVE ] OLLASON |
          [ PHIL ] WOODLAND |
          [ STEVE ] YOUNG;
( SENT-START ( DIAL <$digit> | (PHONE|CALL) $name) SENT-END )
```

where the vertical bars denote alternatives, the square brackets denote optional items and the angle braces denote one or more repetitions. The complete grammar can be depicted as a network as shown in Fig. 3.1.

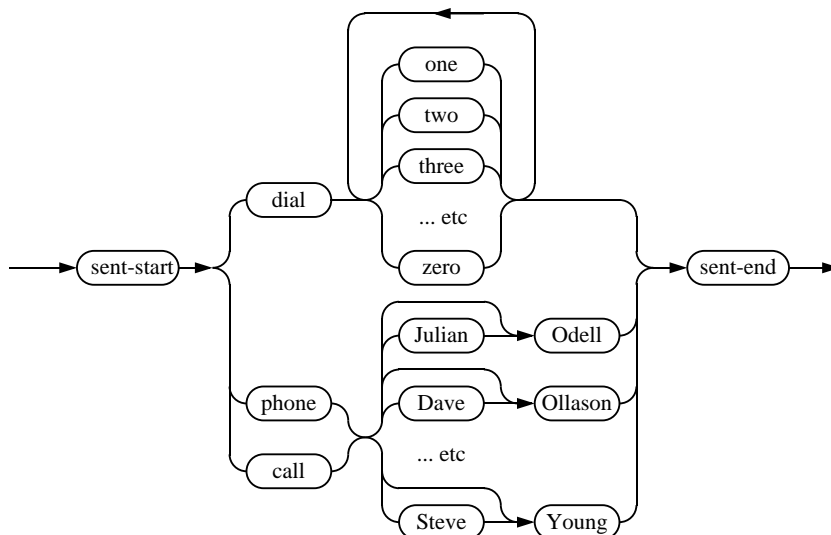


Fig. 3.1 Grammar for Voice Dialling

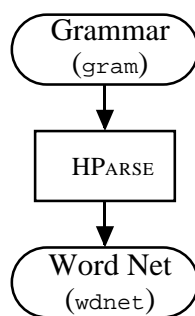


Fig. 3.2
Step 1

The above high level representation of a task grammar is provided for user convenience. The HTK recogniser actually requires a word network to be defined using a low level notation called HTK Standard Lattice Format (SLF) in which each word instance and each word-to-word transition is listed explicitly. This word network can be created automatically from the grammar above using the HPARSE tool, thus assuming that the file `gram` contains the above grammar, executing

```
HParse gram wtnet
```

will create an equivalent word network in the file `wtnet` (see Fig 3.2).

3.1.2 Step 2 - the Dictionary

The first step in building a dictionary is to create a sorted list of the required words. In the telephone dialling task pursued here, it is quite easy to create a list of required words by hand. However, if the task were more complex, it would be necessary to build a word list from the sample sentences present in the training data. Furthermore, to build robust acoustic models, it is necessary to train them on a large set of sentences containing many words and preferably phonetically balanced. For these reasons, the training data will consist of English sentences unrelated to the phone recognition task. Below, a short example of creating a word list from sentence prompts will be given. As noted above the training sentences given here are extracted from some prompts used with the TIMIT database and for convenience reasons they have been renumbered. For example, the first few items might be as follows

```

S0001 ONE VALIDATED ACTS OF SCHOOL DISTRICTS
S0002 TWO OTHER CASES ALSO WERE UNDER ADVISEMENT
S0003 BOTH FIGURES WOULD GO HIGHER IN LATER YEARS
S0004 THIS IS NOT A PROGRAM OF SOCIALIZED MEDICINE
etc

```

The desired training word list (`wlist`) could then be extracted automatically from these. Before using HTK, one would need to edit the text into a suitable format. For example, it would be necessary to change all white space to newlines and then to use the UNIX utilities `sort` and `uniq` to sort the words into a unique alphabetically ordered set, with one word per line. The script `prompts2wlist` from the `HTKTutorial` directory can be used for this purpose.

The dictionary itself can be built from a standard source using `HDMAN`. For this example, the British English BEEP pronouncing dictionary will be used². Its phone set will be adopted without modification except that the stress marks will be removed and a short-pause (`sp`) will be added to the end of every pronunciation. If the dictionary contains any silence markers then the `MP` command will merge the `sil` and `sp` phones into a single `sil`. These changes can be applied using `HDMAN` and an edit script (stored in `global.ded`) containing the three commands

```

AS sp
RS cmu
MP sil sil sp

```

where `cmu` refers to a style of stress marking in which the lexical stress level is marked by a single digit appended to the phone name (e.g. `eh2` means the phone `eh` with level 2 stress).

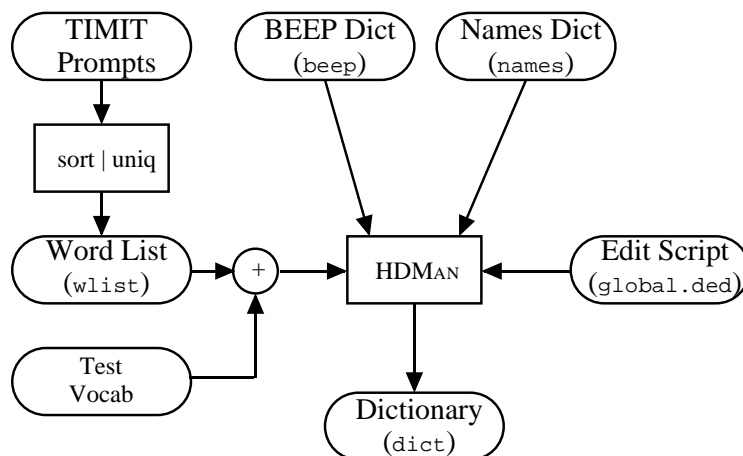


Fig. 3.3 Step 2

The command

```

HDMAN -m -w wlist -n monophones1 -l dlog dict beep names

```

will create a new dictionary called `dict` by searching the source dictionaries `beep` and `names` to find pronunciations for each word in `wlist` (see Fig 3.3). Here, the `wlist` in question needs only to be a sorted list of the words appearing in the task grammar given above.

Note that `names` is a manually constructed file containing pronunciations for the proper names used in the task grammar. The option `-l` instructs `HDMAN` to output a log file `dlog` which contains various statistics about the constructed dictionary. In particular, it indicates if there are words missing. `HDMAN` can also output a list of the phones used, here called `monophones1`. Once training and test data has been recorded, an HMM will be estimated for each of these phones.

The general format of each dictionary entry is

```

WORD [outsym] p1 p2 p3 ...

```

²Available by anonymous ftp from `svr-ftp.eng.cam.ac.uk/pub/comp.speech/dictionaries/beep.tar.gz`. Note that items beginning with unmatched quotes, found at the start of the dictionary, should be removed.

which means that the word `WORD` is pronounced as the sequence of phones `p1 p2 p3 ...`. The string in square brackets specifies the string to output when that word is recognised. If it is omitted then the word itself is output. If it is included but empty, then nothing is output.

To see what the dictionary is like, here are a few entries.

```
A          ah sp
A          ax sp
A          ey sp
CALL       k ao l sp
DIAL       d ay ax l sp
EIGHT      ey t sp
PHONE      f ow n sp
SENT-END   [] sil
SENT-START [] sil
SEVEN      s eh v n sp
TO         t ax sp
TO         t uw sp
ZERO       z ia r ow sp
```

Notice that function words such as `A` and `TO` have multiple pronunciations. The entries for `SENT-START` and `SENT-END` have a silence model `sil` as their pronunciations and null output symbols.

3.1.3 Step 3 - Recording the Data

The training and test data will be recorded using the HTK tool `HSLAB`. This is a combined waveform recording and labelling tool. In this example `HSLAB` will be used just for recording, as labels already exist. However, if you do not have pre-existing training sentences (such as those from the TIMIT database) you can create them either from pre-existing text (as described above) or by labelling your training utterances using `HSLAB`. `HSLAB` is invoked by typing

```
HSLab noname
```

This will cause a window to appear with a waveform display area in the upper half and a row of buttons, including a record button in the lower half. When the name of a normal file is given as argument, `HSLAB` displays its contents. Here, the special file name `noname` indicates that new data is to be recorded. `HSLAB` makes no special provision for prompting the user. However, each time the record button is pressed, it writes the subsequent recording alternately to a file called `noname_0.` and to a file called `noname_1.`. Thus, it is simple to write a shell script which for each successive line of a prompt file, outputs the prompt, waits for either `noname_0.` or `noname_1.` to appear, and then renames the file to the name prepending the prompt (see Fig. 3.4).

While the prompts for training sentences already were provided for above, the prompts for test sentences need to be generated before recording them. The tool `HSGEN` can be used to do this by randomly traversing a word network and outputting each word encountered. For example, typing

```
HSGen -l -n 200 wdnnet dict > testprompts
```

would generate 200 numbered test utterances, the first few of which would look something like:

```
1.  PHONE YOUNG
2.  DIAL OH SIX SEVEN SEVEN OH ZERO
3.  DIAL SEVEN NINE OH OH EIGHT SEVEN NINE NINE
4.  DIAL SIX NINE SIX TWO NINE FOUR ZERO NINE EIGHT
5.  CALL JULIAN ODELL
... etc
```

These can be piped to construct the prompt file `testprompts` for the required test data.

3.1.4 Step 4 - Creating the Transcription Files

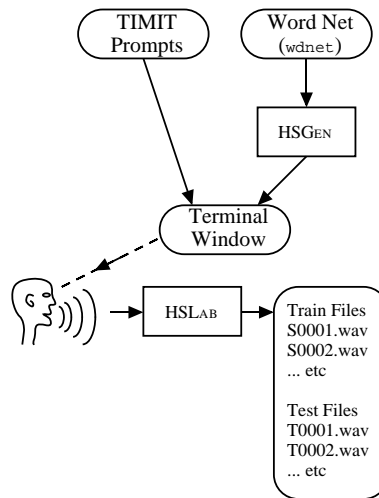


Fig. 3.4 Step 3

To train a set of HMMs, every file of training data must have an associated phone level transcription. Since there is no hand labelled data to bootstrap a set of models, a flat-start scheme will be used instead. To do this, two sets of phone transcriptions will be needed. The set used initially will have no short-pause (`sp`) models between words. Then once reasonable phone models have been generated, an `sp` model will be inserted between words to take care of any pauses introduced by the speaker.

The starting point for both sets of phone transcription is an orthographic transcription in HTK label format. This can be created fairly easily using a text editor or a scripting language. An example of this is found in the RM Demo at point 0.4. Alternatively, the script `prompts2mlf` has been provided in the `HTKTutorial` directory. The effect should be to convert the prompt utterances exemplified above into the following form:

```

#!MLF!#
"*/S0001.lab"
ONE
VALIDATED
ACTS
OF
SCHOOL
DISTRICTS
.
"*/S0002.lab"
TWO
OTHER
CASES
ALSO
WERE
UNDER
ADVISEMENT
.
"*/S0003.lab"
BOTH
FIGURES
(etc.)

```

As can be seen, the prompt labels need to be converted into path names, each word should be written on a single line and each utterance should be terminated by a single period on its own.

The first line of the file just identifies the file as a *Master Label File* (MLF). This is a single file containing a complete set of transcriptions. HTK allows each individual transcription to be stored in its own file but it is more efficient to use an MLF.

The form of the path name used in the MLF deserves some explanation since it is really a *pattern* and not a name. When HTK processes speech files, it expects to find a transcription (or *label file*) with the same name but a different extension. Thus, if the file `/root/sjy/data/S0001.wav` was being processed, HTK would look for a label file called `/root/sjy/data/S0001.lab`. When MLF files are used, HTK scans the file for a pattern which matches the required label file name. However, an asterisk will match any character string and hence the pattern used in the example is in effect path independent. It therefore allows the same transcriptions to be used with different versions of the speech data to be stored in different locations.

Once the word level MLF has been created, phone level MLFs can be generated using the label editor HLED. For example, assuming that the above word level MLF is stored in the file `words.mlf`, the command

```
HLEd -l '*' -d dict -i phones0.mlf mkphones0.led words.mlf
```

will generate a phone level transcription of the following form where the `-l` option is needed to generate the path `'*'` in the output patterns.

```
#!MLF!#
"*/S0001.lab"
sil
w
ah
n
v
ae
l
ih
d
.. etc
```

This process is illustrated in Fig. 3.5.

The HLED edit script `mkphones0.led` contains the following commands

```
EX
IS sil sil
DE sp
```

The expand `EX` command replaces each word in `words.mlf` by the corresponding pronunciation in the dictionary file `dict`. The `IS` command inserts a silence model `sil` at the start and end of every utterance. Finally, the delete `DE` command deletes all short-pause `sp` labels, which are not wanted in the transcription labels at this point.

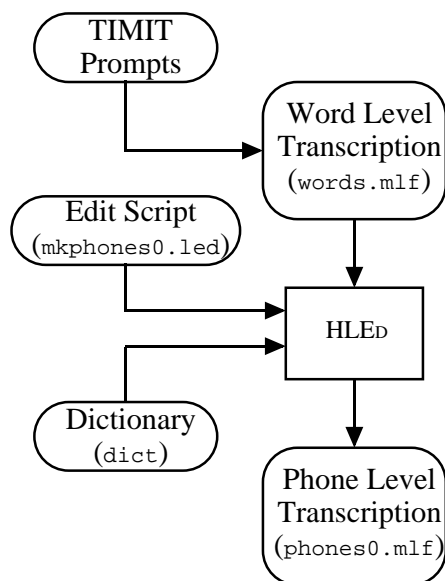


Fig. 3.5 Step 4

3.1.5 Step 5 - Coding the Data

The final stage of data preparation is to parameterise the raw speech waveforms into sequences of feature vectors. HTK support both FFT-based and LPC-based analysis. Here Mel Frequency Cepstral Coefficients (MFCCs), which are derived from FFT-based log spectra, will be used.

Coding can be performed using the tool `HCOPY` configured to automatically convert its input into MFCC vectors. To do this, a configuration file (`config`) is needed which specifies all of the conversion parameters. Reasonable settings for these are as follows

```

# Coding parameters
TARGETKIND = MFCC_0
TARGETRATE = 100000.0
SAVECOMPRESSED = T
SAVEWITHCRC = T
WINDOWSIZE = 250000.0
USEHAMMING = T
PREEMCOEF = 0.97
NUMCHANS = 26
CEPLIFTER = 22
NUMCEPS = 12
ENORMALISE = F
  
```

Some of these settings are in fact the default setting, but they are given explicitly here for completeness. In brief, they specify that the target parameters are to be MFCC using C_0 as the energy component, the frame period is 10msec (HTK uses units of 100ns), the output should be saved in compressed format, and a crc checksum should be added. The FFT should use a Hamming window and the signal should have first order preemphasis applied using a coefficient of 0.97. The filterbank should have 26 channels and 12 MFCC coefficients should be output. The variable `ENORMALISE` is by default true and performs energy normalisation on recorded audio files. It cannot be used with live audio and since the target system is for live audio, this variable should be set to false.

Note that explicitly creating coded data files is not necessary, as coding can be done "on-the-fly" from the original waveform files by specifying the appropriate configuration file (as above) with the relevant HTK tools. However, creating these files reduces the amount of preprocessing required during training, which itself can be a time-consuming process.

To run `HCOPY`, a list of each source file and its corresponding output file is needed. For example, the first few lines might look like

```

/root/sjy/waves/S0001.wav /root/sjy/train/S0001.mfc
/root/sjy/waves/S0002.wav /root/sjy/train/S0002.mfc
/root/sjy/waves/S0003.wav /root/sjy/train/S0003.mfc
/root/sjy/waves/S0004.wav /root/sjy/train/S0004.mfc
(etc.)

```

Files containing lists of files are referred to as script files³ and by convention are given the extension `scp` (although HTK does not demand this). Script files are specified using the standard `-S` option and their contents are read simply as extensions to the command line. Thus, they avoid the need for command lines with several thousand arguments⁴.

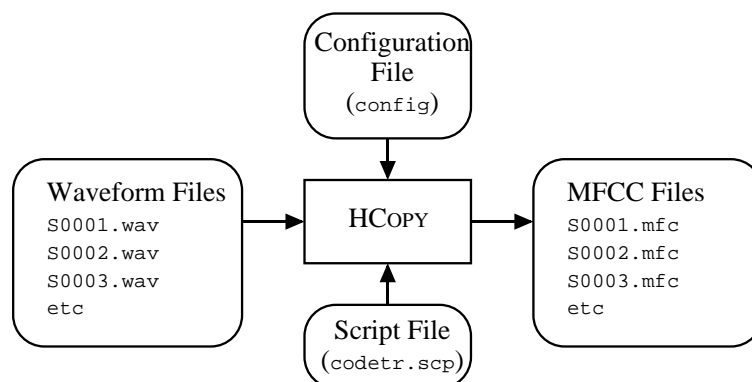


Fig. 3.6 Step 5

Assuming that the above script is stored in the file `codetr.scp`, the training data would be coded by executing

```
HCOPY -T 1 -C config -S codetr.scp
```

This is illustrated in Fig. 3.6. A similar procedure is used to code the test data (using `TARGETKIND = MFCC_0_D_A` in `config`) after which all of the pieces are in place to start training the HMMs.

3.2 Creating Monophone HMMs

In this section, the creation of a well-trained set of single-Gaussian monophone HMMs will be described. The starting point will be a set of identical monophone HMMs in which every mean and variance is identical. These are then retrained, short-pause models are added and the silence model is extended slightly. The monophones are then retrained.

Some of the dictionary entries have multiple pronunciations. However, when HLED was used to expand the word level MLF to create the phone level MLFs, it arbitrarily selected the first pronunciation it found. Once reasonable monophone HMMs have been created, the recogniser tool HVITE can be used to perform a *forced alignment* of the training data. By this means, a new phone level MLF is created in which the choice of pronunciations depends on the acoustic evidence. This new MLF can be used to perform a final re-estimation of the monophone HMMs.

3.2.1 Step 6 - Creating Flat Start Monophones

The first step in HMM training is to define a prototype model. The parameters of this model are not important, its purpose is to define the model topology. For phone-based systems, a good topology to use is 3-state left-right with no skips such as the following

```

~o <VecSize> 39 <MFCC_0_D_A>
~h "proto"

```

³ Not to be confused with files containing *edit* scripts

⁴ Most UNIX shells, especially the C shell, only allow a limited and quite small number of arguments.

```

<BeginHMM>
  <NumStates> 5
  <State> 2
    <Mean> 39
      0.0 0.0 0.0 ...
    <Variance> 39
      1.0 1.0 1.0 ...
  <State> 3
    <Mean> 39
      0.0 0.0 0.0 ...
    <Variance> 39
      1.0 1.0 1.0 ...
  <State> 4
    <Mean> 39
      0.0 0.0 0.0 ...
    <Variance> 39
      1.0 1.0 1.0 ...
  <TransP> 5
    0.0 1.0 0.0 0.0 0.0
    0.0 0.6 0.4 0.0 0.0
    0.0 0.0 0.6 0.4 0.0
    0.0 0.0 0.0 0.7 0.3
    0.0 0.0 0.0 0.0 0.0
<EndHMM>

```

where each ellipsed vector is of length 39. This number, 39, is computed from the length of the parameterised static vector (MFCC_0 = 13) plus the delta coefficients (+13) plus the acceleration coefficients (+13).

The HTK tool HCOMPV will scan a set of data files, compute the global mean and variance and set all of the Gaussians in a given HMM to have the same mean and variance. Hence, assuming that a list of all the training files is stored in `train.scp`, the command

```
HCompV -C config -f 0.01 -m -S train.scp -M hmm0 proto
```

will create a new version of `proto` in the directory `hmm0` in which the zero means and unit variances above have been replaced by the global speech means and variances. Note that the prototype HMM defines the parameter `kind` as `MFCC_0_D_A` (Note: 'zero' not 'oh'). This means that delta and acceleration coefficients are to be computed and appended to the static MFCC coefficients computed and stored during the coding process described above. To ensure that these are computed during loading, the configuration file `config` should be modified to change the target kind, i.e. the configuration file entry for `TARGETKIND` should be changed to

```
TARGETKIND = MFCC_0_D_A
```

HCOMPV has a number of options specified for it. The `-f` option causes a variance floor macro (called `vFloors`) to be generated which is equal to 0.01 times the global variance. This is a vector of values which will be used to set a floor on the variances estimated in the subsequent steps. The `-m` option asks for means to be computed as well as variances. Given this new prototype model stored in the directory `hmm0`, a *Master Macro File* (MMF) called `hmmdefs` containing a copy for each of the required monophone HMMs is constructed by manually copying the prototype and relabeling it for each required monophone (including "sil"). The format of an MMF is similar to that of an MLF and it serves a similar purpose in that it avoids having a large number of individual HMM definition files (see Fig. 3.7).

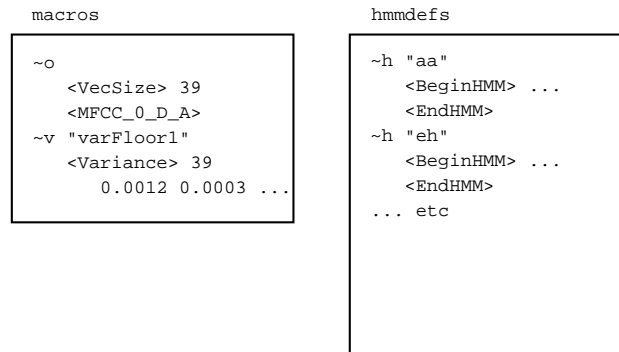


Fig. 3.7 Form of Master Macro Files

The flat start monophones stored in the directory `hmm0` are re-estimated using the embedded re-estimation tool `HEREST` invoked as follows

```

HERest -C config -I phones0.mlf -t 250.0 150.0 1000.0 \
-S train.scp -H hmm0/macros -H hmm0/hmmdefs -M hmm1 monophones0

```

The effect of this is to load all the models in `hmm0` which are listed in the model list `monophones0` (`monophones1` less the short pause (`sp`) model). These are then re-estimated them using the data listed in `train.scp` and the new model set is stored in the directory `hmm1`. Most of the files used in this invocation of `HEREST` have already been described. The exception is the file `macros`. This should contain a so-called *global options* macro and the variance floor macro `vFloors` generated earlier. The global options macro simply defines the HMM parameter kind and the vector size i.e.

```
~o <MFCC_0_D_A> <VecSize> 39
```

See Fig. 3.7. This can be combined with `vFloors` into a text file called `macros`.

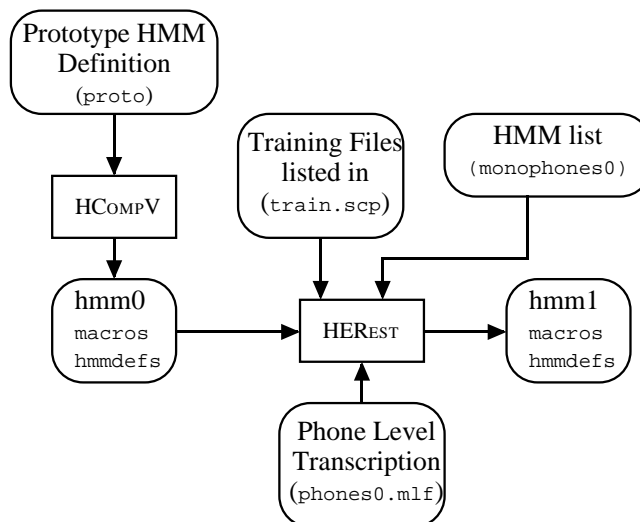


Fig. 3.8 Step 6

The `-t` option sets the pruning thresholds to be used during training. Pruning limits the range of state alignments that the forward-backward algorithm includes in its summation and it can reduce the amount of computation required by an order of magnitude. For most training files, a very tight pruning threshold can be set, however, some training files will provide poorer acoustic matching

and in consequence a wider pruning beam is needed. HEREST deals with this by having an auto-incrementing pruning threshold. In the above example, pruning is normally 250.0. If re-estimation fails on any particular file, the threshold is increased by 150.0 and the file is reprocessed. This is repeated until either the file is successfully processed or the pruning limit of 1000.0 is exceeded. At this point it is safe to assume that there is a serious problem with the training file and hence the fault should be fixed (typically it will be an incorrect transcription) or the training file should be discarded. The process leading to the initial set of monophones in the directory `hmm0` is illustrated in Fig. 3.8.

Each time HEREST is run it performs a single re-estimation. Each new HMM set is stored in a new directory. Execution of HEREST should be repeated twice more, changing the name of the input and output directories (set with the options `-H` and `-M`) each time, until the directory `hmm3` contains the final set of initialised monophone HMMs.

3.2.2 Step 7 - Fixing the Silence Models

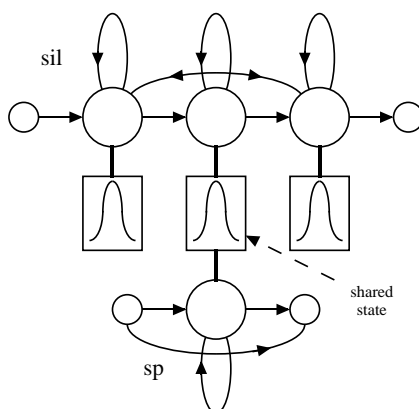


Fig. 3.9 Silence Models

The previous step has generated a 3 state left-to-right HMM for each phone and also a HMM for the silence model `sil`. The next step is to add extra transitions from states 2 to 4 and from states 4 to 2 in the silence model. The idea here is to make the model more robust by allowing individual states to absorb the various impulsive noises in the training data. The backward skip allows this to happen without committing the model to transit to the following word.

Also, at this point, a 1 state short pause `sp` model should be created. This should be a so-called *tee-model* which has a direct transition from entry to exit node. This `sp` has its emitting state tied to the centre state of the silence model. The required topology of the two silence models is shown in Fig. 3.9.

These silence models can be created in two stages

- Use a text editor on the file `hmm3/hmmdefs` to copy the centre state of the `sil` model to make a new `sp` model and store the resulting MMF `hmmdefs`, which includes the new `sp` model, in the new directory `hmm4`.
- Run the HMM editor HHED to add the extra transitions required and tie the `sp` state to the centre `sil` state

HHED works in a similar way to HLED. It applies a set of commands in a script to modify a set of HMMs. In this case, it is executed as follows

```
HHed -H hmm4/macros -H hmm4/hmmdefs -M hmm5 sil.hed monophones1
```

where `sil.hed` contains the following commands

```
AT 2 4 0.2 {sil.transP}
AT 4 2 0.2 {sil.transP}
AT 1 3 0.3 {sp.transP}
TI silst {sil.state[3],sp.state[2]}
```

The AT commands add transitions to the given transition matrices and the final TI command creates a tied-state called `silst`. The parameters of this tied-state are stored in the `hmmdefs` file and within each silence model, the original state parameters are replaced by the name of this macro. Macros are described in more detail below. For now it is sufficient to regard them simply as the mechanism by which HTK implements parameter sharing. Note that the phone list used here has been changed, because the original list `monophones0` has been extended by the new `sp` model. The new file is called `monophones1` and has been used in the above `HHED` command.

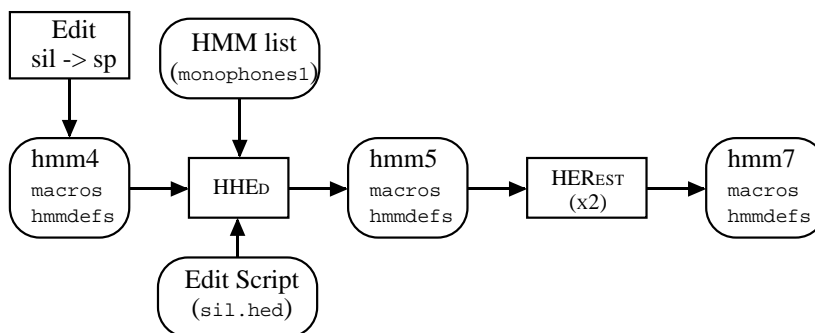


Fig. 3.10 Step 7

Finally, another two passes of `HEREST` are applied using the phone transcriptions with `sp` models between words. This leaves the set of monophone HMMs created so far in the directory `hmm7`. This step is illustrated in Fig. 3.10

3.2.3 Step 8 - Realigning the Training Data

As noted earlier, the dictionary contains multiple pronunciations for some words, particularly function words. The phone models created so far can be used to *realign* the training data and create new transcriptions. This can be done with a single invocation of the HTK recognition tool `HVITE`, viz

```

HVite -l '*' -o SWT -b silence -C config -a -H hmm7/macros \
      -H hmm7/hmmdefs -i aligned.mlf -m -t 250.0 -y lab \
      -I words.mlf -S train.scp dict monophones1
  
```

This command uses the HMMs stored in `hmm7` to transform the input word level transcription `words.mlf` to the new phone level transcription `aligned.mlf` using the pronunciations stored in the dictionary `dict` (see Fig 3.11). The key difference between this operation and the original word-to-phone mapping performed by `HLED` in step 4 is that the recogniser considers all pronunciations for each word and outputs the pronunciation that best matches the acoustic data.

In the above, the `-b` option is used to insert a silence model at the start and end of each utterance. The name `silence` is used on the assumption that the dictionary contains an entry

```
silence sil
```

Note that the dictionary should be sorted firstly by case (upper case first) and secondly alphabetically. The `-t` option sets a pruning level of 250.0 and the `-o` option is used to suppress the printing of scores, word names and time boundaries in the output MLF.

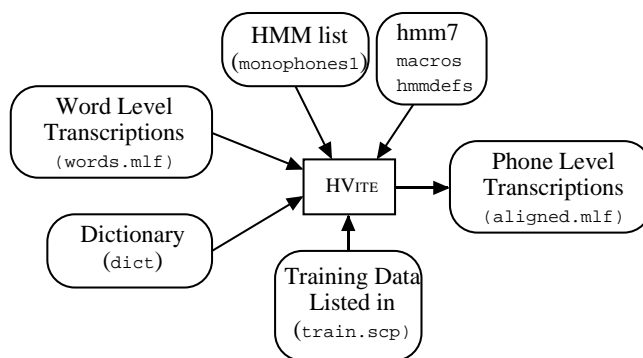


Fig. 3.11 Step 8

Once the new phone alignments have been created, another 2 passes of HEREST can be applied to reestimate the HMM set parameters again. Assuming that this is done, the final monophone HMM set will be stored in directory `hmm9`.

3.3 Creating Tied-State Triphones

Given a set of monophone HMMs, the final stage of model building is to create context-dependent triphone HMMs. This is done in two steps. Firstly, the monophone transcriptions are converted to triphone transcriptions and a set of triphone models are created by copying the monophones and re-estimating. Secondly, similar acoustic states of these triphones are tied to ensure that all state distributions can be robustly estimated.

3.3.1 Step 9 - Making Triphones from Monophones

Context-dependent triphones can be made by simply cloning monophones and then re-estimating using triphone transcriptions. The latter should be created first using HLED because a side-effect is to generate a list of all the triphones for which there is at least one example in the training data. That is, executing

```
HLEd -n triphones1 -l '*' -i wintri.mlf mktri.led aligned.mlf
```

will convert the monophone transcriptions in `aligned.mlf` to an equivalent set of triphone transcriptions in `wintri.mlf`. At the same time, a list of triphones is written to the file `triphones1`. The edit script `mktri.led` contains the commands

```
WB sp
WB sil
TC
```

The two `WB` commands define `sp` and `sil` as *word boundary symbols*. These then block the addition of context in the `TI` command, seen in the following script, which converts all phones (except word boundary symbols) to triphones `.`. For example,

```
sil th ih s sp m ae n sp ...
```

becomes

```
sil th+ih th-ih+s ih-s sp m+ae m-ae+n ae-n sp ...
```

This style of triphone transcription is referred to as *word internal*. Note that some biphones will also be generated as contexts at word boundaries will sometimes only include two phones.

The cloning of models can be done efficiently using the HMM editor HHED:

```
HHEd -B -H hmm9/macros -H hmm9/hmmdefs -M hmm10
mktri.hed monophones1
```


where the edit script `mktri.hed` contains a clone command `CL` followed by `TI` commands to tie all of the transition matrices in each triphone set, that is:

```
CL triphones1
TI T_ah {(*-ah+*,ah+*,*-ah).transP}
TI T_ax {(*-ax+*,ax+*,*-ax).transP}
TI T_ey {(*-ey+*,ey+*,*-ey).transP}
TI T_b {(*-b+*,b+*,*-b).transP}
TI T_ay {(*-ay+*,ay+*,*-ay).transP}
...
```

The file `mktri.hed` can be generated using the *Perl* script `maketrihed` included in the `HTKTutorial` directory. When running the `HHED` command you will get warnings about trying to tie transition matrices for the `sil` and `sp` models. Since neither model is context-dependent there aren't actually any matrices to tie.

The clone command `CL` takes as its argument the name of the file containing the list of triphones (and biphones) generated above. For each model of the form `a-b+c` in this list, it looks for the monophone `b` and makes a copy of it. Each `TI` command takes as its argument the name of a macro and a list of HMM components. The latter uses a notation which attempts to mimic the hierarchical structure of the HMM parameter set in which the transition matrix `transP` can be regarded as a sub-component of each HMM. The list of items within brackets are patterns designed to match the set of triphones, right biphones and left biphones for each phone.

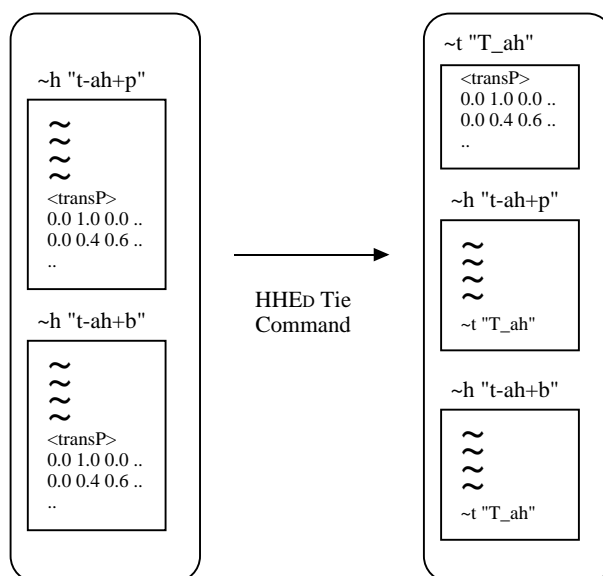


Fig. 3.12 Tying Transition Matrices

Up to now macros and tying have only been mentioned in passing. Although a full explanation must wait until chapter 7, a brief explanation is warranted here. Tying means that one or more HMMs share the same set of parameters. On the left side of Fig. 3.12, two HMM definitions are shown. Each HMM has its own individual transition matrix. On the right side, the effect of the first `TI` command in the edit script `mktri.hed` is shown. The individual transition matrices have been replaced by a reference to a *macro* called `T_ah` which contains a matrix shared by both models. When reestimating tied parameters, the data which would have been used for each of the original untied parameters is pooled so that a much more reliable estimate can be obtained.

Of course, tying could affect performance if performed indiscriminately. Hence, it is important to only tie parameters which have little effect on discrimination. This is the case here where the transition parameters do not vary significantly with acoustic context but nevertheless need to be estimated accurately. Some triphones will occur only once or twice and so very poor estimates would be obtained if tying was not done. These problems of data insufficiency will affect the output distributions too, but this will be dealt with in the next step.

Hitherto, all HMMs have been stored in text format and could be inspected like any text file. Now however, the model files will be getting larger and space and load/store times become an issue. For increased efficiency, HTK can store and load MMFs in binary format. Setting the standard `-B` option causes this to happen.

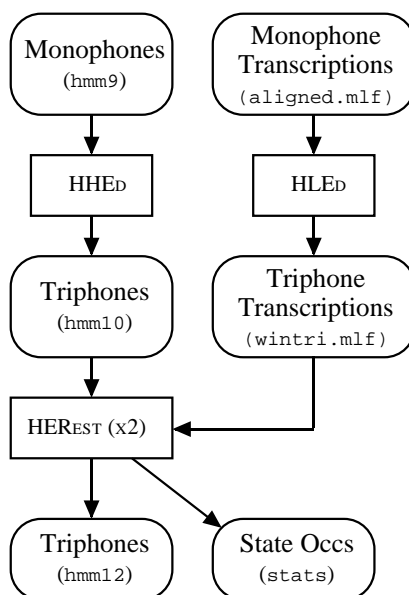


Fig. 3.13 Step 9

Once the context-dependent models have been cloned, the new triphone set can be re-estimated using HEREST. This is done as previously except that the monophone model list is replaced by a triphone list and the triphone transcriptions are used in place of the monophone transcriptions.

For the final pass of HEREST, the `-s` option should be used to generate a file of state occupation statistics called `stats`. In combination with the means and variances, these enable likelihoods to be calculated for clusters of states and are needed during the state-clustering process described below. Fig. 3.13 illustrates this step of the HMM construction procedure. Re-estimation should be again done twice, so that the resultant model sets will ultimately be saved in `hmm12`.

```

HERest -B -C config -I wintri.mlf -t 250.0 150.0 1000.0 -s stats \
-S train.scp -H hmm11/macros -H hmm11/hmmdefs -M hmm12 triphones1
  
```

3.3.2 Step 10 - Making Tied-State Triphones

The outcome of the previous stage is a set of triphone HMMs with all triphones in a phone set sharing the same transition matrix. When estimating these models, many of the variances in the output distributions will have been floored since there will be insufficient data associated with many of the states. The last step in the model building process is to tie states within triphone sets in order to share data and thus be able to make robust parameter estimates.

In the previous step, the `TI` command was used to explicitly tie all members of a set of transition matrices together. However, the choice of which states to tie requires a bit more subtlety since the performance of the recogniser depends crucially on how accurate the state output distributions capture the statistics of the speech data.

HHED provides two mechanisms which allow states to be clustered and then each cluster tied. The first is data-driven and uses a similarity measure between states. The second uses decision trees and is based on asking questions about the left and right contexts of each triphone. The decision tree attempts to find those contexts which make the largest difference to the acoustics and which should therefore distinguish clusters.

Decision tree state tying is performed by running HHED in the normal way, i.e.

```
HHEd -B -H hmm12/macros -H hmm12/hmmdefs -M hmm13 \
tree.hed triphones1 > log
```

Notice that the output is saved in a log file. This is important since some tuning of thresholds is usually needed.

The edit script `tree.hed`, which contains the instructions regarding which contexts to examine for possible clustering, can be rather long and complex. A script for automatically generating this file, `mkclscript`, is found in the RM Demo. A version of the `tree.hed` script, which can be used with this tutorial, is included in the `HTKTutorial` directory. Note that this script is only capable of creating the TB commands (decision tree clustering of states). The questions (QS) still need defining by the user. There is, however, an example list of questions which may be suitable to some tasks (or at least useful as an example) supplied with the RM demo (`lib/quests.hed`). The entire script appropriate for clustering English phone models is too long to show here in the text, however, its main components are given by the following fragments:

```
RO 100.0 stats
TR 0
QS "L_Class-Stop" {p-*,b-*,t-*,d-*,k-*,g-*}
QS "R_Class-Stop" {*+p,*+b,*+t,*+d,*+k,*+g}
QS "L_Nasal" {m-*,n-*,ng-*}
QS "R_Nasal" {*+m,*+n,*+ng}
QS "L_Glide" {y-*,w-*}
QS "R_Glide" {*+y,*+w}
...
QS "L_w" {w-*}
QS "R_w" {*+w}
QS "L_y" {y-*}
QS "R_y" {*+y}
QS "L_z" {z-*}
QS "R_z" {*+z}

TR 2

TB 350.0 "aa_s2" {(aa, *-aa, *-aa+*, aa+*).state[2]}
TB 350.0 "ae_s2" {(ae, *-ae, *-ae+*, ae+*).state[2]}
TB 350.0 "ah_s2" {(ah, *-ah, *-ah+*, ah+*).state[2]}
TB 350.0 "uh_s2" {(uh, *-uh, *-uh+*, uh+*).state[2]}
...
TB 350.0 "y_s4" {(y, *-y, *-y+*, y+*).state[4]}
TB 350.0 "z_s4" {(z, *-z, *-z+*, z+*).state[4]}
TB 350.0 "zh_s4" {(zh, *-zh, *-zh+*, zh+*).state[4]}

TR 1

AU "fulllist"
CO "tiedlist"

ST "trees"
```

Firstly, the `RO` command is used to set the outlier threshold to 100.0 and load the statistics file generated at the end of the previous step. The outlier threshold determines the minimum occupancy of any cluster and prevents a single outlier state forming a singleton cluster just because it is acoustically very different to all the other states. The `TR` command sets the trace level to zero in preparation for loading in the questions. Each `QS` command loads a single question and each question is defined by a set of contexts. For example, the first `QS` command defines a question called `L_Class-Stop` which is true if the left context is either of the stops `p`, `b`, `t`, `d`, `k` or `g`.

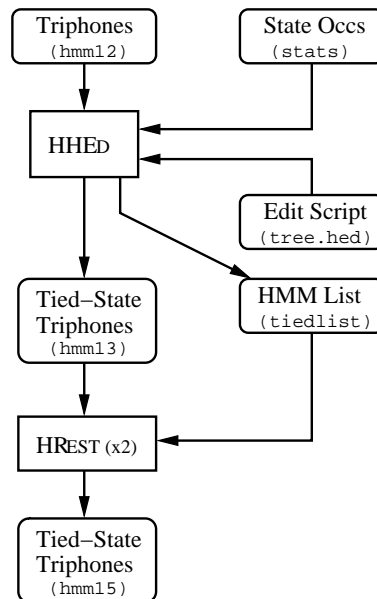


Fig. 3.14 Step 10

Notice that for a triphone system, it is necessary to include questions referring to both the right and left contexts of a phone. The questions should progress from wide, general classifications (such as consonant, vowel, nasal, diphthong, etc.) to specific instances of each phone. Ideally, the full set of questions loaded using the `QS` command would include every possible context which can influence the acoustic realisation of a phone, and can include any linguistic or phonetic classification which may be relevant. There is no harm in creating extra unnecessary questions, because those which are determined to be irrelevant to the data will be ignored.

The second `TR` command enables intermediate level progress reporting so that each of the following `TB` commands can be monitored. Each of these `TB` commands clusters one specific set of states. For example, the first `TB` command applies to the first emitting state of all context-dependent models for the phone `aa`.

Each `TB` command works as follows. Firstly, each set of states defined by the final argument is pooled to form a single cluster. Each question in the question set loaded by the `QS` commands is used to split the pool into two sets. The use of two sets rather than one, allows the log likelihood of the training data to be increased and the question which maximises this increase is selected for the first branch of the tree. The process is then repeated until the increase in log likelihood achievable by any question at any node is less than the threshold specified by the first argument (350.0 in this case).

Note that the values given in the `RO` and `TB` commands affect the degree of tying and therefore the number of states output in the clustered system. The values should be varied according to the amount of training data available. As a final step to the clustering, any pair of clusters which can be merged such that the decrease in log likelihood is below the threshold is merged. On completion, the states in each cluster i are tied to form a single shared state with macro name `xxx_i` where `xxx` is the name given by the second argument of the `TB` command.

The set of triphones used so far only includes those needed to cover the training data. The `AU` command takes as its argument a new list of triphones expanded to include all those needed for recognition. This list can be generated, for example, by using `HDMAN` on the entire dictionary (not just the training dictionary), converting it to triphones using the command `TC` and outputting a list of the distinct triphones to a file using the option `-n`

```
HDMAN -b sp -n fulllist -g global.ded -l flog beep-tri beep
```

The `-b sp` option specifies that the `sp` phone is used as a word boundary, and so is excluded from triphones. The effect of the `AU` command is to use the decision trees to synthesise all of the new previously unseen triphones in the new list.

Once all state-tying has been completed and new models synthesised, some models may share exactly the same 3 states and transition matrices and are thus identical. The `CO` command is used to compact the model set by finding all identical models and tying them together⁵, producing a new list of models called `tiedlist`.

One of the advantages of using decision tree clustering is that it allows previously unseen triphones to be synthesised. To do this, the trees must be saved and this is done by the `ST` command. Later if new previously unseen triphones are required, for example in the pronunciation of a new vocabulary item, the existing model set can be reloaded into `HHED`, the trees reloaded using the `LT` command and then a new extended list of triphones created using the `AU` command.

After `HHED` has completed, the effect of tying can be studied and the thresholds adjusted if necessary. The log file will include summary statistics which give the total number of physical states remaining and the number of models after compacting.

Finally, and for the last time, the models are re-estimated twice using `HEREST`. Fig. 3.14 illustrates this last step in the HMM build process. The trained models are then contained in the file `hmm15/hmmdefs`.

3.4 Recogniser Evaluation

The recogniser is now complete and its performance can be evaluated. The recognition network and dictionary have already been constructed, and test data has been recorded. Thus, all that is necessary is to run the recogniser and then evaluate the results using the HTK analysis tool `HRESULTS`

3.4.1 Step 11 - Recognising the Test Data

Assuming that `test.scp` holds a list of the coded test files, then each test file will be recognised and its transcription output to an MLF called `recout.mlf` by executing the following

```
HVite -H hmm15/macros -H hmm15/hmmdefs -S test.scp \
-l '*' -i recout.mlf -w wdnnet \
-p 0.0 -s 5.0 dict tiedlist
```

The options `-p` and `-s` set the *word insertion penalty* and the *grammar scale factor*, respectively. The word insertion penalty is a fixed value added to each token when it transits from the end of one word to the start of the next. The grammar scale factor is the amount by which the language model probability is scaled before being added to each token as it transits from the end of one word to the start of the next. These parameters can have a significant effect on recognition performance and hence, some tuning on development test data is well worthwhile.

The dictionary contains monophone transcriptions whereas the supplied HMM list contains word internal triphones. `HVITE` will make the necessary conversions when loading the word network `wdnet`. However, if the HMM list contained both monophones and context-dependent phones then `HVITE` would become confused. The required form of word-internal network expansion can be forced by setting the configuration variable `FORCECXTEXP` to true and `ALLOWXWRDEXP` to false (see chapter 12 for details).

Assuming that the MLF `testref.mlf` contains word level transcriptions for each test file⁶, the actual performance can be determined by running `HRESULTS` as follows

```
HResults -I testref.mlf tiedlist recout.mlf
```

the result would be a print-out of the form

```
===== HTK Results Analysis =====
Date: Sun Oct 22 16:14:45 1995
Ref : testrefs.mlf
Rec : recout.mlf
----- Overall Results -----
SENT: %Correct=98.50 [H=197, S=3, N=200]
```

⁵ Note that if the transition matrices had not been tied, the `CO` command would be ineffective since all models would be different by virtue of their unique transition matrices.

⁶The `HLED` tool may have to be used to insert silences at the start and end of each transcription or alternatively `HRESULTS` can be used to ignore silences (or any other symbols) using the `-e` option

```
WORD: %Corr=99.77, Acc=99.65 [H=853, D=1, S=1, I=1, N=855]
```

The line starting with **SENT:** indicates that of the 200 test utterances, 197 (98.50%) were correctly recognised. The following line starting with **WORD:** gives the word level statistics and indicates that of the 855 words in total, 853 (99.77%) were recognised correctly. There was 1 deletion error (D), 1 substitution error (S) and 1 insertion error (I). The accuracy figure (Acc) of 99.65% is lower than the percentage correct (Cor) because it takes account of the insertion errors which the latter ignores.

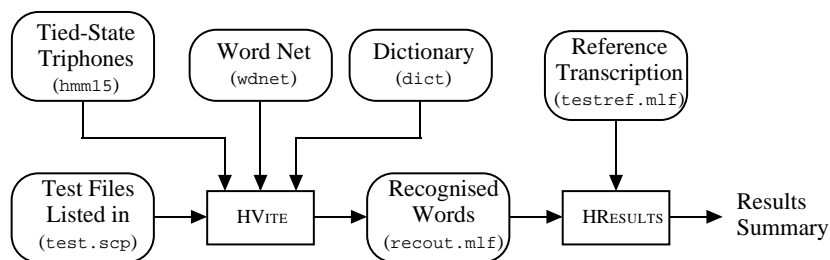


Fig. 3.15 Step 11

3.5 Running the Recogniser Live

The recogniser can also be run with live input. To do this it is only necessary to set the configuration variables needed to convert the input audio to the correct form of parameterisation. Specifically, the following needs to be appended to the configuration file `config` to create a new configuration file `config2`

```
# Waveform capture
SOURCERATE=625.0
SOURCEKIND=HAUDIO
SOURCEFORMAT=HTK
ENORMALISE=F
USESILDET=T
MEASURESIL=F
OUTSILWARN=T
```

These indicate that the source is direct audio with sample period 62.5 μ secs. The silence detector is enabled and a measurement of the background speech/silence levels should be made at start-up. The final line makes sure that a warning is printed when this silence measurement is being made.

Once the configuration file has been set-up for direct audio input, HVITE can be run as in the previous step except that no files need be given as arguments

```
HVite -H hmm15/macros -H hmm15/hmmdefs -C config2 \
-w wdnet -p 0.0 -s 5.0 dict tiedlist
```

On start-up, HVITE will prompt the user to speak an arbitrary sentence (approx. 4 secs) in order to measure the speech and background silence levels. It will then repeatedly recognise and, if trace level bit 1 is set, it will output each utterance to the terminal. A typical session is as follows

```
Read 1648 physical / 4131 logical HMMs
Read lattice with 26 nodes / 52 arcs
Created network with 123 nodes / 151 links

READY[1]>
Please speak sentence - measuring levels
Level measurement completed
DIAL FOUR SIX FOUR TWO FOUR OH
```

```

== [303 frames] -95.5773 [Ac=-28630.2 LM=-329.8] (Act=21.8)

READY[2]>
DIAL ZERO EIGHT SIX TWO
== [228 frames] -99.3758 [Ac=-22402.2 LM=-255.5] (Act=21.8)

READY[3]>
etc

```

During loading, information will be printed out regarding the different recogniser components. The physical models are the distinct HMMs used by the system, while the logical models include all model names. The number of logical models is higher than the number of physical models because many logically distinct models have been determined to be physically identical and have been merged during the previous model building steps. The lattice information refers to the number of links and nodes in the recognition syntax. The network information refers to actual recognition network built by expanding the lattice using the current HMM set, dictionary and any context expansion rules specified. After each utterance, the numerical information gives the total number of frames, the average log likelihood per frame, the total acoustic score, the total language model score and the average number of models active.

Note that if it was required to recognise a new name, then the following two changes would be needed

1. the grammar would be altered to include the new name
2. a pronunciation for the new name would be added to the dictionary

If the new name required triphones which did not exist, then they could be created by loading the existing triphone set into HHED, loading the decision trees using the LT command and then using the AU command to generate a new complete triphone set.

3.6 Adapting the HMMs

The previous sections have described the stages required to build a simple voice dialling system. To simplify this process, speaker dependent models were developed using training data from a single user. Consequently, recognition accuracy for any other users would be poor. To overcome this limitation, a set of speaker independent models could be constructed, but this would require large amounts of training data from a variety of speakers. An alternative is to adapt the current speaker dependent models to the characteristics of a new speaker using a small amount of training or adaptation data. In general, adaptation techniques are applied to well trained speaker independent model sets to enable them to better model the characteristics of particular speakers.

HTK supports both supervised adaptation, where the true transcription of the data is known and unsupervised adaptation where the transcription is hypothesised. In HTK supervised adaptation is performed offline by HEADAPT using maximum likelihood linear regression (MLLR) and/or maximum a-posteriori (MAP) techniques to estimate a series of transforms or a transformed model set, that reduces the mismatch between the current model set and the adaptation data. Unsupervised adaptation is provided by HVITE (see section 13.6.2), using just MLLR.

The following sections describe offline supervised adaptation (using MLLR) with the use of HEADAPT.

3.6.1 Step 12 - Preparation of the Adaptation Data

As in normal recogniser development, the first stage in adaptation involves data preparation. Speech data from the new user is required for both adapting the models and testing the adapted system. The data can be obtained in a similar fashion to that taken to prepare the original test data. Initially, prompt lists for the adaptation and test data will be generated using HSGEN. For example, typing

```

HSGen -l -n 20 wdnet dict > promptsAdapt
HSGen -l -n 20 wdnet dict > promptsTest

```

would produce two prompt files for the adaptation and test data. The amount of adaptation data required will normally be found empirically, but a performance improvement should be observable after just 30 seconds of speech. In this case, around 20 utterances should be sufficient. HSLAB can be used to record the associated speech.

Assuming that the script files `codeAdapt.scp` and `codeTest.scp` list the source and output files for the adaptation and test data respectively then both sets of speech can then be coded using the HCPY commands given below.

```
HCopy -C config -S codeAdapt.scp
HCopy -C config -S codeTest.scp
```

The final stage of preparation involves generating context dependent phone transcriptions of the adaptation data and word level transcriptions of the test data for use in adapting the models and evaluating their performance. The transcriptions of the test data can be obtained using `prompts2mlf`. To minimize the problem of multiple pronunciations the phone level transcriptions of the adaptation data can be obtained by using HVITE to perform a *forced alignment* of the adaptation data. Assuming that word level transcriptions are listed in `adaptWords.mlf`, then the following command will place the phone transcriptions in `adaptPhones.mlf`.

```
HVite -l '*' -o SWT -b silence -C config -a -H hmm15/macros \
-H hmm15/hmmdefs -i adaptPhones.mlf -m -t 250.0 \
-I adaptWords.mlf -y lab -S adapt.scp dict tiedlist
```

3.6.2 Step 13 - Generating the Transforms

HEADAPT provides two forms of MLLR adaptation depending on the amount of adaptation data available. If only small amounts are available a global transform can be generated for every output distribution of every model. As more adaptation data becomes available more specific transforms can be generated for specific groups of Gaussians. To identify the number of transforms that can be estimated using the current adaptation data, HEADAPT uses a regression class tree to cluster together groups of output distributions that are to undergo the same transformation. The HTK tool HHED can be used to build a regression class tree and store it as part of the HMM set. For example,

```
HHed -B -H hmm15/macros -H hmm15/hmmdefs -M hmm16 regtree.hed tiedlist
```

creates a regression class tree using the models stored in `hmm15`. The models are written out to the `hmm16` directory together with the regression class tree information. The HHED edit script `regtree.hed` contains the following commands

```
RN "models"
LS "stats"
RC 32 "rtree"
```

The RN command assigns an identifier to the HMM set. The LS command loads the state occupation statistics file `stats` generated by the last application of HEREST which created the models in `hmm15`. The RC command then attempts to build a regression class tree with 32 terminal or leaf nodes using these statistics.

HEADAPT can be used to perform either static adaptation, where all the adaptation data is processed in a single block or incremental adaptation, where adaptation is performed after a specified number of utterances and this is controlled by the `-i` option. In this tutorial the default setting of static adaptation will be used.

A typical use of HEADAPT involves two passes. On the first pass a global adaptation is performed. The second pass then uses the global transformation to transform the model set, producing better frame/state alignments which are then used to estimate a set of more specific transforms, using a regression class tree. After estimating the transforms, HEADAPT can output either the newly adapted model set or the transformations themselves in a transform model file (TMF). The latter can be advantageous if storage is an issue since the TMFs are significantly smaller than MMFs and the computational overhead incurred when transforming a model set using a TMF is negligible.

The two applications of HEADAPT below demonstrate a static two-pass adaptation approach where the global and regression class transformations are stored in the `global.tmf` and `rc.tmf` files respectively. The standard `-J` and `-K` options are used to load and save the TMFs respectively.


```
HEAdapt -C config -g -S adapt.scp -I adaptPhones.mlf -H hmm16/macros \
        -H hmm16/hmmdefs -K global.tmf tiedlist
```

```
HEAdapt -C config -S adapt.scp -I adaptPhones.mlf -H hmm16/macros \
        -H hmm16/hmmdefs -J global.tmf -K rc.tmf tiedlist
```

3.6.3 Step 14 - Evaluation of the Adapted System

To evaluate the performance of the adaptation, the test data previously recorded is recognised using HVITE. Assuming that `testAdapt.scp` contains a list of all of the coded test files, then HVITE can be invoked in much the same way as before but with the additional `-J` argument used to load the model transformation file `rc.tmf`.

```
HVite -H hmm16/macros -H hmm16/hmmdefs -S testAdapt.scp -l '*' \
      -J rc.tmf -i recoutAdapt.mlf -w wnet \
      -p 0.0 -s 5.0 dict tiedlist
```

The results of the adapted model set can then be observed using HRESULTS in the usual manner.

The RM Demo contains a section on speaker adaptation (point 5.6) and the recognition results obtained using an adapted model set are given below.

```
===== HTK Results Analysis =====
Date: Wed Jan 06 21:09:23 1999
Ref : usr/local/htk/RMHTK_V2.1/RMLib/wlabs/dms0_tst.mlf
Rec : adapt/dms0_tst.mlf
----- Overall Results -----
SENT: %Correct=66.33 [H=65, S=33, N=98]
WORD: %Corr=94.25, Acc=93.10 [H=738, D=11, S=34, I=9, N=783]
=====
```

The performance improvement gained by the adapted models can be evaluated by recognising the test data using the unadapted model set and comparing the two results. For the RM Demo task the following results were obtained with an unadapted model set.

```
===== HTK Results Analysis =====
Date: Mon Dec 14 10:59:28 1998
Ref : usr/local/htk/RMHTK_V2.1/RMLib/wlabs/dms0_tst.mlf
Rec : unadapt/dms0_tst.mlf
----- Overall Results -----
SENT: %Correct=46.00 [H=46, S=54, N=100]
WORD: %Corr=89.04, Acc=86.43 [H=715, D=26, S=62, I=21, N=803]
=====
```

3.7 Summary

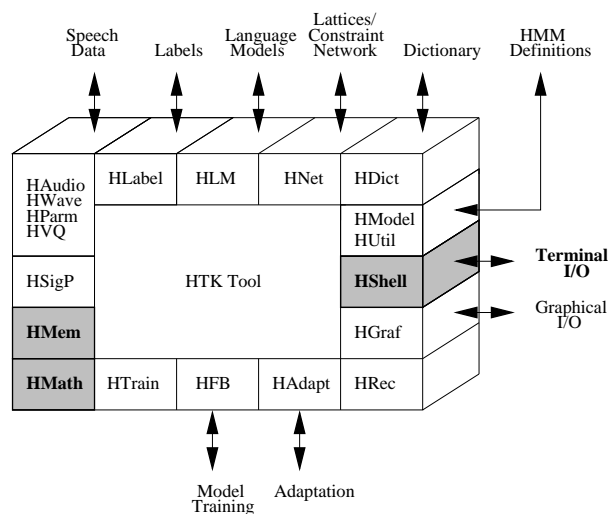
This chapter has described the construction of a tied-state phone-based continuous speech recogniser and in so doing, it has touched on most of the main areas addressed by HTK: recording, data preparation, HMM definitions, training tools, adaptation tools, networks, decoding and evaluating. The rest of this book discusses each of these topics in detail.

Part II

HTK in Depth

Chapter 4

The Operating Environment



This chapter discusses the various ways of controlling the operation of HTK tools along with related aspects of file system organisation, error reporting and memory management. All of the operating system and user interface functions are provided by the HTK module HShell. Memory management is a low level function which is largely invisible to the user, but it is useful to have a basic understanding of it in order to appreciate memory requirements and interpret diagnostic output from tools. Low level memory management in HTK is provided by HMEM and the management of higher level structures such as vectors and matrices is provided by HMATH.

The behaviour of a HTK tool depends on three sources of information. Firstly, all HTK tools are executed by issuing commands to the operating system shell. Each command typically contains the names of the various files that the tool needs to function and a number of optional arguments which control the detailed behaviour of the tool. Secondly, as noted in chapter 2 and shown in the adjacent figure, every HTK tool uses a set of standard library modules to interface to the various file types and to connect with the outside world. Many of these modules can be customised by setting parameters in a *configuration file*. Thirdly, a small number of parameters are specified using environment variables.

Terminal output mostly depends on the specific tool being used, however, there are some generic output functions which are provided by the library modules and which are therefore common across tools. These include version reporting, memory usage and error reporting.

Finally, HTK can read and write most data sources through pipes as an alternative to direct input and output from data files. This allows filters to be used, and in particular, it allows many of the external files used by HTK to be stored directly in compressed form and then decompressed *on-the-fly* when the data is read back in.

All of the above is discussed in more detail in the following sections.

4.1 The Command Line

The general form of command line for invoking a tool is¹

```
tool [options] files ...
```

Options always consist of a dash followed by a single letter. Some options are followed by an argument as follows

```
-i           - a switch option
-t 3        - an integer valued option
-a 0.01     - a float valued option
-s hello    - a string valued option
```

Option names consisting of a capital letter are common across all tools (see section 4.4). Integer arguments may be given in any of the standard C formats, for example, 13, 0xD and 015 all represent the same number. Typing the name of a tool on its own always causes a short summary of the command line options to be printed in place of its normal operation. For example, typing

```
HERest
```

would result in the following output

```
USAGE: HERest [options] hmmList dataFiles...
```

Option		Default
-c f	Mixture pruning threshold	10.0
-d s	dir to find hmm definitions	current
-m N	set min examples needed per model	3
-o s	extension for new hmm files	as src
-p N	set parallel mode to N	off
...		

The first line shows the names of the required files and the rest consists of a listing of each option, its meaning, and its default value.

The precise naming convention for specifying files depends on the operating system being used, but HTK always assumes the existence of a hierarchical file system and it maintains a distinction between directory paths and file names.

In general, a file will be located either in the current directory, some subdirectory of the current directory or some subdirectory of the root directory. For example, in the command

```
HList s1 dir/s2 /users/sjy/speech/s3
```

file *s1* must be in the current directory, *s2* must be in the directory *dir* within the current directory and *s3* must be in the directory */users/sjy/speech*.

Some tools allow directories to be specified via configuration parameters and command line options. In all cases, the final path character (eg / in UNIX, : in Mac OS) need not (but may be) included. For example, both of the following are acceptable and have equivalent effect

```
HInit -l i: -L mymodels/new/ hmmfile data*
HInit -l i: -L mymodels/new hmmfile data*
```

where the *-L* option specifies the directory in which to find the label files associated with the data files.

4.2 Script Files

Tools which require a potentially very long list of files (e.g. training tools) always allow the files to be specified in a script file via the *-S* option instead of via the command line. This is particularly useful when running under an OS with limited file name expansion capability. Thus, for example, HINIT may be invoked by either

¹All of the examples in this book assume the UNIX Operating System and the C Shell but the principles apply to any OS which supports hierarchical files and command line arguments

```
HInit hmmfile s1 s2 s3 s4 s5 ....
```

or

```
HInit -S filelist hmmfile
```

where `filelist` holds the list of files `s1`, `s2`, etc. Each file listed in a script should be separated by white space or a new line. Usually, files are listed on separate lines, however, when using `HCOPY` which read pairs of files as its arguments, it is normal to write each pair on a single line. Script files should only be used for storing ellipsed file list arguments. Note that shell meta-characters should not be used in script files and will not be interpreted by the HTK tools.

4.3 Configuration Files

Configuration files are used for customising the HTK working environment. They consist of a list of parameter-values pairs along with an optional prefix which limits the scope of the parameter to a specific module or tool.

The name of a configuration file can be specified explicitly on the command line using the `-C` command. For example, when executing

```
HERest ... -C myconfig s1 s2 s3 s4 ...
```

The operation of `HEREST` will depend on the parameter settings in the file `myconfig`.

When an explicit configuration file is specified, only those parameters mentioned in that file are actually changed and all other parameters retain their default values. These defaults are built-in. However, user-defined defaults can be set by assigning the name of a default configuration file to the environment variable `HCONFIG`. Thus, for example, using the UNIX C Shell, writing

```
setenv HCONFIG myconfig
HERest ... s1 s2 s3 s4 ...
```

would have an identical effect to the preceding example. However, in this case, a further refinement of the configuration values is possible since the opportunity to specify an explicit configuration file on the command line remains. For example, in

```
setenv HCONFIG myconfig
HERest ... -C xconfig s1 s2 s3 s4 ...
```

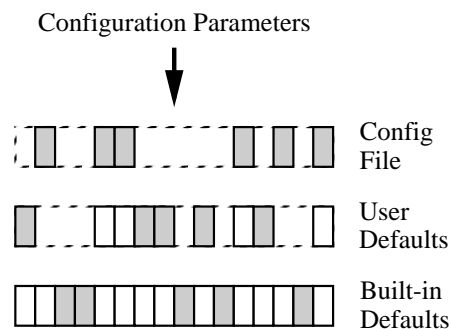


Fig. 4.1 Defining a Configuration

the parameter values in `xconfig` will over-ride those in `myconfig` which in turn will over-ride the built-in defaults. In practice, most HTK users will set general-purpose default configuration values using `HCONFIG` and will then over-ride these as required for specific tasks using the `-C` command line option. This is illustrated in Fig. 4.1 where the darkened rectangles indicate *active* parameter definitions. Viewed from above, all of the remaining parameter definitions can be seen to be masked by higher level over-rides.

The configuration file itself consists of a sequence of parameter definitions of the form

```
[MODULE:] PARAMETER = VALUE
```

One parameter definition is written per line and square brackets indicate that the module name is optional. Parameter definitions are not case sensitive but by convention they are written in upper case. A # character indicates that the rest of the line is a comment.

As an example, the following is a simple configuration file

```
# Example config file
  TARGETKIND = MFCC
  NUMCHANS   = 20
  WINDOWSIZE = 250000.0    # ie 25 msecs
  PREEMCOEF = 0.97
  ENORMALISE = T
HSHELL: TRACE = 02          # octal
HPARM:  TRACE = 0101
```

The first three lines contain no module name and hence they apply globally, that is, any library module or tool which is interested in the configuration parameter `NUMCHANS` will read the given parameter value. In practice, this is not a problem with library modules since nearly all configuration parameters have unique names. The final two lines show the same parameter name being given different values within different modules. This is an example of a parameter which every module responds to and hence does not have a unique name.

This example also shows each of the four possible types of value that can appear in a configuration file: string, integer, float and Boolean. The configuration parameter `TARGETKIND` requires a string value specifying the name of a speech parameter kind. Strings not starting with a letter should be enclosed in double quotes. `NUMCHANS` requires an integer value specifying the number of filter-bank channels to use in the analysis. `WINDOWSIZE` actually requires a floating-point value specifying the window size in units of 100ns. However, an integer can always be given wherever a float is required. `PREEMCOEF` also requires a floating-point value specifying the pre-emphasis coefficient to be used. Finally, `ENORMALISE` is a Boolean parameter which determines whether or not energy normalisation is to be performed, its value must be `T`, `TRUE` or `F`, `FALSE`. Notice also that, as in command line options, integer values can use the C conventions for writing in non-decimal bases. Thus, the trace value of `0101` is equal to decimal 65. This is particularly useful in this case because trace values are typically interpreted as bit-strings by HTK modules and tools.

If the name of a configuration variable is mis-typed, there will be no warning and the variable will simply be ignored. To help guard against this, the standard option `-D` can be used. This displays all of the configuration variables before and after the tool runs. In the latter case, all configuration variables which are still unread are marked by a hash character.

4.4 Standard Options

As noted in section 4.1, options consisting of a capital letter are common across all tools. Many are specific to particular file types and they will be introduced as they arise. However, five of the standard options are actually built-in to `HSHELL`. Two of these have been mentioned already. The option `-C` is used to specify a configuration file name and the option `-S` is used to specify a script file name.

The three remaining standard options provided directly by `HSHELL` are `-A`, `-D` and `-V`. The option `-A` causes the current command line arguments to be printed. When running experiments via scripts, it is a good idea to use this option to record in a log file the precise settings used for each tool. The `-D` option causes the current configuration parameters set by the user to be displayed both before and after the tool has executed. The initial display allows the configuration values to be checked before potentially wasting a large amount of cpu time through incorrectly set parameters. The final display shows which configuration variables were actually used during the execution of the tool. Its form is shown by the following example

```
HTK Configuration Parameters [3]
  Module/Tool   Parameter          Value
# HMODEL       SAVEBINARY        TRUE
HPARM          TARGETRATE        256000.000000
               TARGETKIND        MFCC_0
```

Here three configuration parameters have been set but the hash (#) indicates that `SAVEBINAR`y has not been used. Finally, the option `-V` causes version information for the tool and each module used by that tool to be listed. These should always be quoted when making bug reports.

Finally, all tools implement the trace option `-T`. Trace values are typically bit strings and the meaning of each bit is described in the reference section for each tool. Setting a trace option via the command line overrides any setting for that same trace option in a configuration file. This is a general rule, command line options always override defaults set in configuration files.

All of the standard options are listed in the final summary section of this chapter.

4.5 Error Reporting

The `HShell` module provides a standard mechanism for reporting errors and warnings. A typical error message is as follows

```
HList: ERROR [+1110]
IsWave: cannot open file speech.dat
```

This indicates that the tool `HList` is reporting an error number `+1110`. All errors have positive error numbers and always result in the tool terminating. Warnings have negative error numbers and the tool does not terminate. The first two digits of an error number indicate the module or tool in which the error is located (`HList` in this case) and the last two digits define the class of error. The second line of the error message names the actual routine in which the error occurred (here `IsWave`) and the actual error message. All errors and warnings are listed in the reference section at the end of this book indexed by error/warning number. This listing contains more details on each error or warning along with suggested causes.

Error messages are sent to the standard error stream but warnings are sent to the standard output stream. The reason for the latter is that most `HTK` tools are run with progress tracing enabled. Sending warnings to the standard output stream ensures that they are properly interleaved with the trace of progress so that it is easy to determine the point at which the warning was issued. Sending warnings to standard error would lose this information.

The default behaviour of a `HTK` tool on terminating due to an error is to exit normally returning the error number as exit status. If, however, the configuration variable `ABORTONERR` is set true, then the tool will core dump. This is a debugging facility which should not concern most users.

4.6 Strings and Names

Many `HTK` definition files include names of various types of objects: for example labels, model names, words in a grammar, etc. In order to achieve some uniformity, `HTK` applies standard rules for reading strings which are names.

A name string consists of a single white space delimited word or a quoted string. Either the single quote `'` or the double quote `"` can be used to quote strings but the start and end quotes must be matched. The backslash `\` character can also be used to introduce otherwise reserved characters. The character following a backslash is inserted into the string without special processing unless that character is a digit in the range 0 to 7. In that case, the three characters following the backslash are read and interpreted as an octal character code. When the three characters are not octal digits the result is not well defined.

In summary the special processing is

Notation	Meaning
<code>\\</code>	<code>\</code>
<code>_</code>	represents a space that will not terminate a string
<code>\'</code>	<code>'</code> (and will not end a quoted string)
<code>\"</code>	<code>"</code> (and will not end a quoted string)
<code>\nnn</code>	the character with octal code <code>\nnn</code>

Note that the above allows the same effect to be achieved in a number of different ways. For example,

```
"\"QUOTE"
\"QUOTE
```

```
'"QUOTE'
\042QUOTE
```

all produce the string "QUOTE.

The only exceptions to the above general rules are:

- Where models are specified in HHED scripts, commas (,), dots (.), and closing brackets ()) are all used as extra delimiters to allow HHED scripts created for earlier versions of HTK to be used unchanged. Hence for example, (a,b,c,d) would be split into 4 distinct name strings a, b, c and d.
- When the configuration variable RAWMITFORMAT is set true, each word in a language model definition file consists of a white space delimited string with no special processing being performed.
- Source dictionaries read by HDMAN are read using the standard HTK string conventions, however, the command IR can be used in a HDMAN source edit script to switch to using this raw format.
- To ensure that the general definition of a name string works properly in HTK master label files, all MLFs must have the reserved . and /// terminators alone on a line with no surrounding white space. If this causes problems reading old MLF files, the configuration variable V1COMPAT should be set true in the module HLABEL. In this case, HTK will attempt to simulate the behaviour of the older version 1.5.
- To force numbers to be interpreted as strings rather than times or scores in a label file, they must be quoted. If the configuration variable QUOTECHAR is set to ' or " then output labels will be quoted with the specified quote character. If QUOTECHAR is set to \, then output labels will be escaped. The default is to select the simplest quoting mechanism.

Note that under some versions of Unix HTK can support the 8-bit character sets used for the representation of various orthographies. In such cases the shell environment variable \$LANG usually governs which ISO character set is in use.

4.7 Memory Management

Memory management is a very low level function and is mostly invisible to HTK users. However, some applications require very large amounts of memory. For example, building the models for a large vocabulary continuous speech dictation system might require 150Mb or more. Clearly, when memory demands become this large, a proper understanding of the impact of system design decisions on memory usage is important. The first step in this is to have a basic understanding of memory allocation in HTK.

Many HTK tools dynamically construct large and complex data structures in memory. To keep strict control over this and to reduce memory allocation overheads to an absolute minimum, HTK performs its own memory management. Thus, every time that a module or tool wishes to allocate some memory, it does so by calling routines in HMEM. At a slightly higher level, math objects such as vectors and matrices are allocated by HMATH but using the primitives provided by HMEM.

To make memory allocation and de-allocation very fast, tools create specific memory allocators for specific objects or groups of objects. These memory allocators are divided into a sequence of blocks, and they are organised as either Stacks, M-heaps or C-heaps. A Stack constrains the pattern of allocation and de-allocation requests to be made in a last-allocated first-deallocated order but allows objects of any size to be allocated. An M-heap allows an arbitrary pattern of allocation and de-allocation requests to be made but all allocated objects must be the same size. Both of these memory allocation disciplines are more restricted than the general mechanism supplied by the operating system, and as a result, such memory operations are faster and incur no storage overhead due to the need to maintain hidden housekeeping information in each allocated object. Finally, a C-heap uses the underlying operating system and allows arbitrary allocation patterns, and as a result incurs the associated time and space overheads. The use of C-heaps is avoided wherever possible.

Most tools provide one or more trace options which show how much memory has been allocated. The following shows the form of the output


```

----- Heap Statistics -----
nblk=1, siz= 100000*1, used= 32056, alloc= 100000 : Global Stack[S]
nblk=1, siz= 200*28, used= 100, alloc= 5600 : cellHeap[M]
nblk=1, siz= 10000*1, used= 3450, alloc= 10000 : mlfHeap[S]
nblk=2, siz= 7504*1, used= 9216, alloc= 10346 : nameHeap[S]
-----

```

Each line describes the status of each memory allocator and gives the number of blocks allocated, the current block size (number of elements in block \times the number of bytes in each element)², the total number of bytes in use by the tool and the total number of bytes currently allocated to that allocator. The end of each line gives the name of the allocator and its type: Stack[S], M-heap[M] or C-heap[M]. The element size for Stacks will always be 1 but will be variable in M-heaps. The documentation for the memory intensive HTK tools indicates what each of the main memory allocators are used for and this information allows the effects of various system design choices to be monitored.

4.8 Input/Output via Pipes and Networks

Most types of file in HTK can be input or output via a pipe instead of directly from or to disk. The mechanism for doing this is to assign the required input or output filter command to a configuration parameter or to an environment variable, either can be used. Within this command, any occurrence of the dollar symbol \$ will be replaced by the name of the required file. The output of the command will then be input to or output from the HTK tool via a pipe.

For example, the following command will normally list the contents of the speech waveform file `spfile`

```
HList spfile
```

However, if the value of the environment variable `HWAVERFILTER` is set as follows

```
setenv HWAVERFILTER 'gunzip -c $'
```

then the effect is to invoke the decompression filter `gunzip` with its input connected to the file `spfile` and its output connected to `HLIST` via a pipe. Each different type of file has a unique associated variable so that multiple input and/or filters can be used. The full list of these is given in the summary section at the end of this chapter.

HTK is often used to process large amounts of data and typically this data is distributed across a network. In many systems, an attempt to open a file can fail because of temporary network *glitches*. In the majority of cases, a second or third attempt to open the file a few seconds later will succeed and all will be well. To allow this to be done automatically, HTK tools can be configured to retry opening a file several times before giving up. This is done simply by setting the configuration parameter `MAXTRYOPEN` to the required number of retries³.

4.9 Byte-swapping of HTK data files

Virtually all HTK tools can read and write data to and from binary files. The use of binary format as opposed to text can speed up the performance of the tools and at the same time reduce the file size when manipulating large quantities of data. Typical binary files used by the HTK tools are speech waveform/parameter files, binary master model files (MMF), binary accumulator files used in HMM parameter estimation and binary lattice files. However, the use of binary data format often introduces incompatibilities between different machine architectures due to the different byte ordering conventions used to represent numerical quantities. In such cases, byte swapping of the data is required. To avoid incompatibilities across different machine architectures, all HTK binary data files are written out using big-endian (`NONVAX`) representation of numerical values. Similarly, during loading HTK binary format files are assumed to be in `NONVAX` byte order. The default behavior can be altered using the configuration parameters `NATURALREADORDER` and `NATURALWRITEORDER`. Setting `NATURALREADORDER` to true will instruct the HTK tools to interpret the binary input data in the machine's natural byte order (byte swapping will never take place).

² Block sizes typically grow as more blocks are allocated

³ This does not work if input filters are used.

Similarly, setting `NATURALWRITEORDER` to true will instruct the tools to write out data using the machine's natural byte order. The default value of these two configuration variables is false which is the appropriate setting when using HTK in a multiple machine architecture environment. In an environment comprising entirely of machines with VAX byte order both configuration parameters can be set true which will disable the byte swapping procedure during reading and writing of data.

4.10 Summary

This section summarises the environment variables and configuration parameters recognised by `HSHELL`, `HMEM` and `HMATH`. It also provides a list of all the standard command line options used with HTK.

Table 4.1 lists all of the configuration parameters along with a brief description. A missing module name means that it is recognised by more than one module. Table 4.2 lists all of the environment parameters used by these modules. Finally, table 4.3 lists all of the standard options.

Module	Name	Description
HSHELL	ABORTONERR	core dump on error (for debugging)
HSHELL	HWAVEFILTER	Filter for waveform file input
HSHELL	HPARMFILTER	Filter for parameter file input
HSHELL	HLANGMODFILTER	Filter for language model file input
HSHELL	HMMLISTFILTER	Filter for HMM list file input
HSHELL	HMMDEFFILTER	Filter for HMM definition file input
HSHELL	HLABELFILTER	Filter for Label file input
HSHELL	HNETFILTER	Filter for Network file input
HSHELL	HDICTFILTER	Filter for Dictionary file input
HSHELL	HWAVEOFILTER	Filter for waveform file output
HSHELL	HPARMOFILTER	Filter for parameter file output
HSHELL	HLANGMODOFILTER	Filter for language model file output
HSHELL	HMMLISTOFILTER	Filter for HMM list file output
HSHELL	HMMDEFOFILTER	Filter for HMM definition file output
HSHELL	HLABELOFILTER	Filter for Label file output
HSHELL	HNETOFILTER	Filter for Network file output
HSHELL	HDICTOFILTER	Filter for Dictionary file output
HSHELL	MAXTRYOPEN	number of file open retries
HSHELL	NONUMESCAPES	prevent string output using \012 format
	NATURALREADORDER	Enable natural read order for HTK binary files
	NATURALWRITEORDER	Enable natural write order for HTK binary files
HMEM	PROTECTSTAKS	warn if stack is cut-back (debugging)
	TRACE	trace control (default=0)

Table 4.1 Configuration Parameters used in Operating Environment

Env Variable	Meaning
HCONFIG	Name of default configuration file
HxxxFILTER	Input/Output filters as above

Table 4.2 Environment Variables used in Operating Environment

Standard Option	Meaning
-A	Print command line arguments
-B	Store output HMM macro files in binary
-C <i>cf</i>	Configuration file is <i>cf</i>
-D	Display configuration variables
-F <i>fmt</i>	Set source data file format to <i>fmt</i>
-G <i>fmt</i>	Set source label file format to <i>fmt</i>
-H <i>mmf</i>	Load HMM macro file <i>mmf</i>
-I <i>mlf</i>	Load master label file <i>mlf</i>
-J <i>tmf</i>	Load transform model file <i>tmf</i>
-K <i>tmf</i>	Save transform model file <i>tmf</i>
-L <i>dir</i>	Look for label files in directory <i>dir</i>
-M <i>dir</i>	Store output HMM macro files in directory <i>dir</i>
-O <i>fmt</i>	Set output data file format to <i>fmt</i>
-P <i>fmt</i>	Set output label file format to <i>fmt</i>
-Q	Print command summary info
-S <i>scp</i>	Use command line script file <i>scp</i>
-T <i>N</i>	Set trace level to <i>N</i>
-V	Print version information
-X <i>ext</i>	Set label file extension to <i>ext</i>

Table. 4.3 Summary of Standard Options

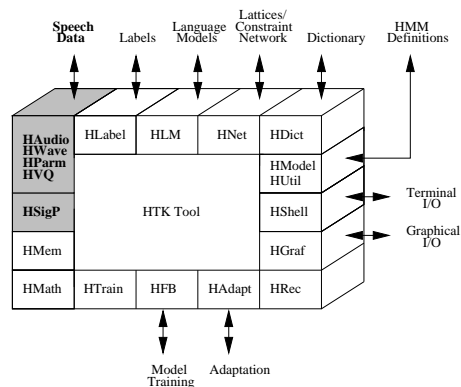
Chapter 5

Speech Input/Output

Many tools need to input parameterised speech data and HTK provides a number of different methods for doing this:

- input from a previously encoded speech parameter file
- input from a waveform file which is encoded as part of the input processing
- input from an audio device which is encoded as part of the input processing.

For input from a waveform file, a large number of different file formats are supported, including all of the commonly used CD-ROM formats. Input/output for parameter files is limited to the standard HTK file format and the new Entropic Esignal format.



All HTK speech input is controlled by configuration parameters which give details of what processing operations to apply to each input speech file or audio source. This chapter describes speech input/output in HTK. The general mechanisms are explained and the various configuration parameters are defined. The facilities for signal pre-processing, linear prediction-based processing, Fourier-based processing and vector quantisation are presented and the supported file formats are given. Also described are the facilities for augmenting the basic speech parameters with energy measures, delta coefficients and acceleration (delta-delta) coefficients and for splitting each parameter vector into multiple data streams to form *observations*. The chapter concludes with a brief description of the tools HLIST and HCOPY which are provided for viewing, manipulating and encoding speech files.

5.1 General Mechanism

The facilities for speech input and output in HTK are provided by five distinct modules: HAUDIO, HWAVE, HPARM, HVQ and HSIGP. The interconnections between these modules are shown in Fig. 5.1.

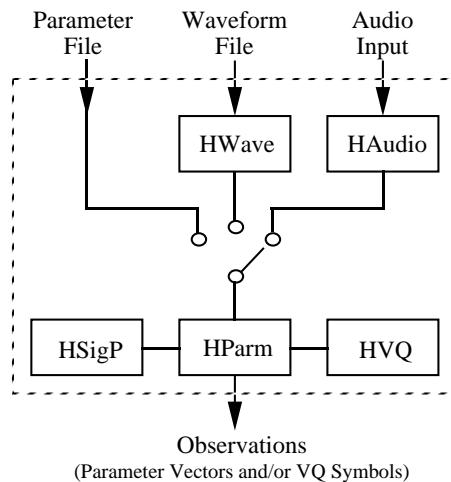


Fig. 5.1 Speech Input Subsystem

Waveforms are read from files using `HWAVE`, or are input direct from an audio device using `HAUDIO`. In a few rare cases, such as in the display tool `HSLAB`, only the speech waveform is needed. However, in most cases the waveform is wanted in parameterised form and the required encoding is performed by `HPARM` using the signal processing operations defined in `HSIGP`. The parameter vectors are output by `HPARM` in the form of observations which are the basic units of data processed by the HTK recognition and training tools. An observation contains all components of a raw parameter vector but it may be possibly split into a number of independent parts. Each such part is regarded by a HTK tool as a statistically independent data stream. Also, an observation may include VQ indices attached to each data stream. Alternatively, VQ indices can be read directly from a parameter file in which case the observation will contain only VQ indices.

Usually a HTK tool will require a number of speech data files to be specified on the command line. In the majority of cases, these files will be required in parameterised form. Thus, the following example invokes the HTK embedded training tool `HEREST` to re-estimate a set of models using the speech data files `s1`, `s2`, `s3`, ... These are input via the library module `HPARM` and they must be in exactly the form needed by the models.

```
HERest ... s1 s2 s3 s4 ...
```

However, if the external form of the speech data files is not in the required form, it will often be possible to convert them automatically during the input process. To do this, configuration parameter values are specified whose function is to define exactly how the conversion should be done. The key idea is that there is a *source parameter kind* and *target parameter kind*. The source refers to the natural form of the data in the external medium and the target refers to the form of the data that is required internally by the HTK tool. The principle function of the speech input subsystem is to convert the source parameter kind into the required target parameter kind.

Parameter kinds consist of a base form to which one or more qualifiers may be attached where each qualifier consists of a single letter preceded by an underscore character. Some examples of parameter kinds are

```
WAVEFORM    simple waveform
LPC         linear prediction coefficients
LPC_D_E    LPC with energy and delta coefficients
MFCC_C     compressed mel-cepstral coefficients
```

The required source and target parameter kinds are specified using the configuration parameters `SOURCEKIND` and `TARGETKIND`. Thus, if the following configuration parameters were defined

```
SOURCEKIND = WAVEFORM
TARGETKIND = MFCC_E
```

then the speech input subsystem would expect each input file to contain a speech waveform and it would convert it to mel-frequency cepstral coefficients with log energy appended.

The source need not be a waveform. For example, the configuration parameters

```
SOURCEKIND = LPC
TARGETKIND = LPREFC
```

would be used to read in files containing linear prediction coefficients and convert them to reflection coefficients.

For convenience, a special parameter kind called `ANON` is provided. When the source is specified as `ANON` then the actual kind of the source is determined from the input file. When `ANON` is used in the target kind, then it is assumed to be identical to the source. For example, the effect of the following configuration parameters

```
SOURCEKIND = ANON
TARGETKIND = ANON_D
```

would simply be to add delta coefficients to whatever the source form happened to be. The source and target parameter kinds default to `ANON` to indicate that by default no input conversions are performed. Note, however, that where two or more files are listed on the command line, the meaning of `ANON` will not be re-interpreted from one file to the next. Thus, it is a general rule, that any tool reading multiple source speech files requires that all the files have the same parameter kind.

The conversions applied by HTK's input subsystem can be complex and may not always behave exactly as expected. There are two facilities that can be used to help check and debug the set-up of the speech i/o configuration parameters. Firstly, the tool `HLIST` simply displays speech data by listing it on the terminal. However, since `HLIST` uses the speech input subsystem like all HTK tools, if a value for `TARGETKIND` is set, then it will display the target form rather than the source form. This is the simplest way to check the form of the speech data that will actually be delivered to a HTK tool. `HLIST` is described in more detail in section 5.12 below.

Secondly, trace output can be generated from the `HPARM` module by setting the `TRACE` configuration file parameter. This is a bit-string in which individual bits cover different parts of the conversion processing. The details are given in the reference section.

To summarise, speech input in HTK is controlled by configuration parameters. The key parameters are `SOURCEKIND` and `TARGETKIND` which specify the source and target parameter kinds. These determine the end-points of the required input conversion. However, to properly specify the detailed steps in between, more configuration parameters must be defined. These are described in subsequent sections.

5.2 Speech Signal Processing

In this section, the basic mechanisms involved in transforming a speech waveform into a sequence of parameter vectors will be described. Throughout this section, it is assumed that the `SOURCEKIND` is `WAVEFORM` and that data is being read from a HTK format file via `HWAVE`. Reading from different format files is described below in section 5.8. Much of the material in this section also applies to data read direct from an audio device, the additional features needed to deal with this latter case are described later in section 5.9.

The overall process is illustrated in Fig. 5.2 which shows the sampled waveform being converted into a sequence of parameter blocks. In general, HTK regards both waveform files and parameter files as being just sample sequences, the only difference being that in the former case the samples are 2-byte integers and in the latter they are multi-component vectors. The sample rate of the input waveform will normally be determined from the input file itself. However, it can be set explicitly using the configuration parameter `SOURCERATE`. The period between each parameter vector determines the output sample rate and it is set using the configuration parameter `TARGETRATE`. The segment of waveform used to determine each parameter vector is usually referred to as a window and its size is set by the configuration parameter `WINDOWSIZE`. Notice that the window size and frame rate are independent. Normally, the window size will be larger than the frame rate so that successive windows overlap as illustrated in Fig. 5.2.

For example, a waveform sampled at 16kHz would be converted into 100 parameter vectors per second using a 25 msec window by setting the following configuration parameters.

```

SOURCERATE = 625
TARGETRATE = 100000
WINDOWSIZE = 250000

```

Remember that all durations are specified in 100 nsec units¹.

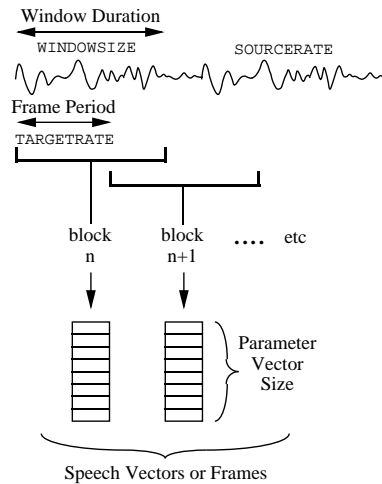


Fig. 5.2 Speech Encoding Process

Independent of what parameter kind is required, there are some simple pre-processing operations that can be applied prior to performing the actual signal analysis. Firstly, the DC mean can be removed from the source waveform by setting the Boolean configuration parameter `ZMEANSOURCE` to true (i.e. T). This is useful when the original analogue-digital conversion has added a DC offset to the signal. It is applied to each window individually so that it can be used both when reading from a file and when using direct audio input².

Secondly, it is common practice to pre-emphasise the signal by applying the first order difference equation

$$s'_n = s_n - k s_{n-1} \quad (5.1)$$

to the samples $\{s_n, n = 1, N\}$ in each window. Here k is the pre-emphasis coefficient which should be in the range $0 \leq k < 1$. It is specified using the configuration parameter `PREEMCOEF`. Finally, it is usually beneficial to taper the samples in each window so that discontinuities at the window edges are attenuated. This is done by setting the Boolean configuration parameter `USEHAMMING` to true. This applies the following transformation to the samples $\{s_n, n = 1, N\}$ in the window

$$s'_n = \left\{ 0.54 - 0.46 \cos \left(\frac{2\pi(n-1)}{N-1} \right) \right\} s_n \quad (5.2)$$

When both pre-emphasis and Hamming windowing are enabled, pre-emphasis is performed first.

In practice, all three of the above are usually applied. Hence, a configuration file will typically contain the following

```

ZMEANSOURCE = T
USEHAMMING = T
PREEMCOEF = 0.97

```

¹ The somewhat bizarre choice of 100nsec units originated in Version 1 of HTK when times were represented by integers and this unit was the best compromise between precision and range. Times are now represented by doubles and hence the constraints no longer apply. However, the need for backwards compatibility means that 100nsec units have been retained. The names `SOURCERATE` and `TARGETRATE` are also non-ideal, `SOURCEPERIOD` and `TARGETPERIOD` would be better.

² This method of applying a zero mean is different to HTK Version 1.5 where the mean was calculated and subtracted from the whole speech file in one operation. The configuration variable `V1COMPAT` can be set to revert to this older behaviour.

Certain types of artificially generated waveform data can cause numerical overflows with some coding schemes. In such cases adding a small amount of random noise to the waveform data solves the problem. The noise is added to the samples using

$$s'_n = s_n + qRND() \quad (5.3)$$

where $RND()$ is a uniformly distributed random value over the interval $[-1.0, +1.0]$ and q is the scaling factor. The amount of noise added to the data (q) is set with the configuration parameter `ADDDITHER` (default value 0.0). A positive value causes the noise signal added to be the same every time (ensuring that the same file always gives exactly the same results). With a negative value the noise is random and the same file may produce slightly different results in different trials.

One problem that can arise when processing speech waveform files obtained from external sources, such as databases on CD-ROM, is that the byte-order may be different to that used by the machine on which HTK is running. To deal with this problem, `HWAVE` can perform automatic byte-swapping in order to preserve proper byte order. HTK assumes by default that speech waveform data is encoded as a sequence of 2-byte integers as is the case for most current speech databases³. If the source format is known, then `HWAVE` will also make an assumption about the byte order used to create speech files in that format. It then checks the byte order of the machine that it is running on and automatically performs byte-swapping if the order is different. For unknown formats, proper byte order can be ensured by setting the configuration parameter `BYTEORDER` to `VAX` if the speech data was created on a little-endian machine such as a VAX or an IBM PC, and to anything else (e.g. `NONVAX`) if the speech data was created on a big-endian machine such as a SUN, HP or Macintosh machine.

The reading/writing of HTK format waveform files can be further controlled via the configuration parameters `NATURALREADORDER` and `NATURALWRITEORDER`. The effect and default settings of these parameters are described in section 4.9. Note that `BYTEORDER` should not be used when `NATURALREADORDER` is set to true. Finally, note that HTK can also byte-swap parameterised files in a similar way provided that only the byte-order of each 4 byte float requires inversion.

5.3 Linear Prediction Analysis

In linear prediction (LP) analysis, the vocal tract transfer function is modelled by an all-pole filter with transfer function⁴

$$H(z) = \frac{1}{\sum_{i=0}^p a_i z^{-i}} \quad (5.4)$$

where p is the number of poles and $a_0 \equiv 1$. The filter coefficients $\{a_i\}$ are chosen to minimise the mean square filter prediction error summed over the analysis window. The HTK module `HSIGP` uses the *autocorrelation method* to perform this optimisation as follows.

Given a window of speech samples $\{s_n, n = 1, N\}$, the first $p + 1$ terms of the autocorrelation sequence are calculated from

$$r_i = \sum_{j=1}^{N-i} s_j s_{j+i} \quad (5.5)$$

where $i = 0, p$. The filter coefficients are then computed recursively using a set of auxiliary coefficients $\{k_i\}$ which can be interpreted as the reflection coefficients of an equivalent acoustic tube and the prediction error E which is initially equal to r_0 . Let $\{k_j^{(i-1)}\}$ and $\{a_j^{(i-1)}\}$ be the reflection and filter coefficients for a filter of order $i - 1$, then a filter of order i can be calculated in three steps. Firstly, a new set of reflection coefficients are calculated.

$$k_j^{(i)} = k_j^{(i-1)} \quad (5.6)$$

for $j = 1, i - 1$ and

$$k_i^{(i)} = \left\{ r_i + \sum_{j=1}^{i-1} a_j^{(i-1)} r_{i-j} \right\} / E^{(i-1)} \quad (5.7)$$

³Many of the more recent speech databases use compression. In these cases, the data may be regarded as being logically encoded as a sequence of 2-byte integers even if the actual storage uses a variable length encoding scheme.

⁴Note that some textbooks define the denominator of equation 5.4 as $1 - \sum_{i=1}^p a_i z^{-i}$ so that the filter coefficients are the negatives of those computed by HTK.

Secondly, the prediction energy is updated.

$$E^{(i)} = (1 - k_i^{(i)} k_i^{(i)}) E^{(i-1)} \quad (5.8)$$

Finally, new filter coefficients are computed

$$a_j^{(i)} = a_j^{(i-1)} - k_i^{(i)} a_{i-j}^{(i-1)} \quad (5.9)$$

for $j = 1, i - 1$ and

$$a_i^{(i)} = -k_i^{(i)} \quad (5.10)$$

This process is repeated from $i = 1$ through to the required filter order $i = p$.

To effect the above transformation, the target parameter kind must be set to either `LPC` to obtain the LP filter parameters $\{a_i\}$ or `LPREFC` to obtain the reflection coefficients $\{k_i\}$. The required filter order must also be set using the configuration parameter `LPCORDER`. Thus, for example, the following configuration settings would produce a target parameterisation consisting of 12 reflection coefficients per vector.

```
TARGETKIND = LPREFC
LPCORDER = 12
```

An alternative LPC-based parameterisation is obtained by setting the target kind to `LPCEPSTRA` to generate linear prediction cepstra. The cepstrum of a signal is computed by taking a Fourier (or similar) transform of the log spectrum. In the case of linear prediction cepstra, the required spectrum is the linear prediction spectrum which can be obtained from the Fourier transform of the filter coefficients. However, it can be shown that the required cepstra can be more efficiently computed using a simple recursion

$$c_n = -a_n + \frac{1}{n} \sum_{i=1}^{n-1} (n-i) a_i c_{n-i} \quad (5.11)$$

The number of cepstra generated need not be the same as the number of filter coefficients, hence it is set by a separate configuration parameter called `NUMCEPS`.

The principal advantage of cepstral coefficients is that they are generally decorrelated and this allows diagonal covariances to be used in the HMMs. However, one minor problem with them is that the higher order cepstra are numerically quite small and this results in a very wide range of variances when going from the low to high cepstral coefficients. HTK does not have a problem with this but for pragmatic reasons such as displaying model parameters, flooring variances, etc., it is convenient to re-scale the cepstral coefficients to have similar magnitudes. This is done by setting the configuration parameter `CEPLIFTER` to some value L to *lifter* the cepstra according to the following formula

$$c'_n = \left(1 + \frac{L}{2} \sin \frac{\pi n}{L} \right) c_n \quad (5.12)$$

As an example, the following configuration parameters would use a 14'th order linear prediction analysis to generate 12 liftered LP cepstra per target vector

```
TARGETKIND = LPCEPSTRA
LPCORDER = 14
NUMCEPS = 12
CEPLIFTER = 22
```

These are typical of the values needed to generate a good front-end parameterisation for a speech recogniser based on linear prediction.

Finally, note that the conversions supported by HTK are not limited to the case where the source is a waveform. HTK can convert any LP-based parameter into any other LP-based parameter.

5.4 Filterbank Analysis

The human ear resolves frequencies non-linearly across the audio spectrum and empirical evidence suggests that designing a front-end to operate in a similar non-linear manner improves recognition performance. A popular alternative to linear prediction based analysis is therefore filterbank

analysis since this provides a much more straightforward route to obtaining the desired non-linear frequency resolution. However, filterbank amplitudes are highly correlated and hence, the use of a cepstral transformation in this case is virtually mandatory if the data is to be used in a HMM based recogniser with diagonal covariances.

HTK provides a simple Fourier transform based filterbank designed to give approximately equal resolution on a mel-scale. Fig. 5.3 illustrates the general form of this filterbank. As can be seen, the filters used are triangular and they are equally spaced along the mel-scale which is defined by

$$\text{Mel}(f) = 2595 \log_{10}\left(1 + \frac{f}{700}\right) \quad (5.13)$$

To implement this filterbank, the window of speech data is transformed using a Fourier transform and the magnitude is taken. The magnitude coefficients are then *binned* by correlating them with each triangular filter. Here binning means that each FFT magnitude coefficient is multiplied by the corresponding filter gain and the results accumulated. Thus, each bin holds a weighted sum representing the spectral magnitude in that filterbank channel. As an alternative, the Boolean configuration parameter USEPOWER can be set true to use the power rather than the magnitude of the Fourier transform in the binning process.

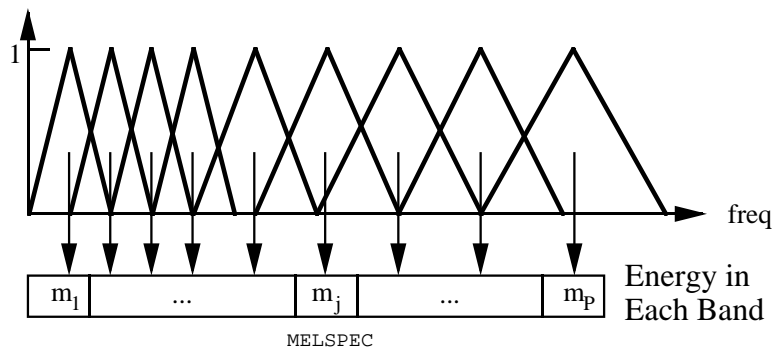


Fig. 5.3 Mel-Scale Filter Bank

Normally the triangular filters are spread over the whole frequency range from zero upto the Nyquist frequency. However, band-limiting is often useful to reject unwanted frequencies or avoid allocating filters to frequency regions in which there is no useful signal energy. For filterbank analysis only, lower and upper frequency cut-offs can be set using the configuration parameters LOFREQ and HIFREQ. For example,

```
LOFREQ = 300
HIFREQ = 3400
```

might be used for processing telephone speech. When low and high pass cut-offs are set in this way, the specified number of filterbank channels are distributed equally on the mel-scale across the resulting pass-band such that the lower cut-off of the first filter is at LOFREQ and the upper cut-off of the last filter is at HIFREQ.

If mel-scale filterbank parameters are required directly, then the target kind should be set to MELSPEC. Alternatively, log filterbank parameters can be generated by setting the target kind to FBANK. Most often, however, cepstral parameters are required and these are indicated by setting the target kind to MFCC standing for Mel-Frequency Cepstral Coefficients (MFCCs). These are calculated from the log filterbank amplitudes $\{m_j\}$ using the Discrete Cosine Transform

$$c_i = \sqrt{\frac{2}{N}} \sum_{j=1}^N m_j \cos\left(\frac{\pi i}{N}(j - 0.5)\right) \quad (5.14)$$

where N is the number of filterbank channels set by the configuration parameter NUMCHANS. The required number of cepstral coefficients is set by NUMCEPS as in the linear prediction case. Liftering can also be applied to MFCCs using the CEPLIFTER configuration parameter (see equation 5.12).

MFCCs are the parameterisation of choice for many speech recognition applications. They give good discrimination and lend themselves to a number of manipulations. In particular, the effect of inserting a transmission channel on the input speech is to multiply the speech spectrum by the channel transfer function. In the log cepstral domain, this multiplication becomes a simple addition which can be removed by subtracting the cepstral mean from all input vectors. In practice, of course, the mean has to be estimated over a limited amount of speech data so the subtraction will not be perfect. Nevertheless, this simple technique is very effective in practice where it compensates for long-term spectral effects such as those caused by different microphones and audio channels. To perform this so-called *Cepstral Mean Normalisation* in HTK it is only necessary to add the `_Z` qualifier to the target parameter kind. The mean is estimated by computing the average of each cepstral parameter across each input speech file. Since this cannot be done with live audio, cepstral mean compensation is not supported for this case.

5.5 Energy Measures

To augment the spectral parameters derived from linear prediction or mel-filterbank analysis, an energy term can be appended by including the qualifier `_E` in the target kind. The energy is computed as the log of the signal energy, that is, for speech samples $\{s_n, n = 1, N\}$

$$E = \log \sum_{n=1}^N s_n^2 \quad (5.15)$$

This log energy measure can be normalised to the range $-E_{min}..1.0$ by setting the Boolean configuration parameter `ENORMALISE` to true (default setting). This normalisation is implemented by subtracting the maximum value of E in the utterance and adding 1.0. Note that energy normalisation is incompatible with live audio input and in such circumstances the configuration variable `ENORMALISE` should be explicitly set false. The lowest energy in the utterance can be clamped using the configuration parameter `SILFLOOR` which gives the ratio between the maximum and minimum energies in the utterance in dB. Its default value is 50dB. Finally, the overall log energy can be arbitrarily scaled by the value of the configuration parameter `ESCALE` whose default is 0.1.

When calculating energy for LPC-derived parameterisations, the default is to use the zeroth delay autocorrelation coefficient (r_0). However, this means that the energy is calculated after windowing and pre-emphasis. If the configuration parameter `RAWENERGY` is set true, however, then energy is calculated separately before any windowing or pre-emphasis regardless of the requested parameterisation⁵.

In addition to, or in place of, the log energy, the qualifier `_0` can be added to a target kind to indicate that the 0th cepstral parameter C_0 is to be appended. This qualifier is only valid if the target kind is MFCC. Unlike earlier versions of HTK scaling factors set by the configuration variable `ESCALE` are not applied to C_0 ⁶.

5.6 Delta and Acceleration Coefficients

The performance of a speech recognition system can be greatly enhanced by adding time derivatives to the basic static parameters. In HTK, these are indicated by attaching qualifiers to the basic parameter kind. The qualifier `_D` indicates that first order regression coefficients (referred to as delta coefficients) are appended, and the qualifier `_A` indicates that second order regression coefficients (referred to as acceleration coefficients) are appended. The `_A` qualifier cannot be used without also using the `_D` qualifier.

The delta coefficients are computed using the following regression formula

$$d_t = \frac{\sum_{\theta=1}^{\Theta} \theta (c_{t+\theta} - c_{t-\theta})}{2 \sum_{\theta=1}^{\Theta} \theta^2} \quad (5.16)$$

where d_t is a delta coefficient at time t computed in terms of the corresponding static coefficients $c_{t-\theta}$ to $c_{t+\theta}$. The value of Θ is set using the configuration parameter `DELTAWINDOW`. The same

⁵ In any event, setting the compatibility variable `V1COMPAT` to true in HPARM will ensure that the calculation of energy is compatible with that computed by the Version 1 tool HCODE.

⁶ Unless `V1COMPAT` is set to true.

formula is applied to the delta coefficients to obtain acceleration coefficients except that in this case the window size is set by `ACCWINDOW`. Since equation 5.16 relies on past and future speech parameter values, some modification is needed at the beginning and end of the speech. The default behaviour is to replicate the first or last vector as needed to fill the regression window.

In older version 1.5 of HTK and earlier, this end-effect problem was solved by using simple first order differences at the start and end of the speech, that is

$$d_t = c_{t+1} - c_t, \quad t < \Theta \quad (5.17)$$

and

$$d_t = c_t - c_{t-1}, \quad t \geq T - \Theta \quad (5.18)$$

where T is the length of the data file. If required, this older behaviour can be restored by setting the configuration variable `V1COMPAT` to true in `HPARM`.

For some purposes, it is useful to use simple differences throughout. This can be achieved by setting the configuration variable `SIMPLEDIFFS` to true in `HPARM`. In this case, just the end-points of the delta window are used, i.e.

$$d_t = \frac{(c_{t+\Theta} - c_{t-\Theta})}{2\Theta} \quad (5.19)$$

When delta and acceleration coefficients are requested, they are computed for all static parameters including energy if present. In some applications, the absolute energy is not useful but time derivatives of the energy may be. By including the `_E` qualifier together with the `_N` qualifier, the absolute energy is suppressed leaving just the delta and acceleration coefficients of the energy.

5.7 Storage of Parameter Files

Whereas HTK can handle waveform data in a variety of file formats, all parameterised speech data is stored externally in either native HTK format data files or Entropic Esignal format files. Entropic ESPS format is no longer supported directly, but input and output filters can be used to convert ESPS to Esignal format on input and Esignal to ESPS on output.

5.7.1 HTK Format Parameter Files

HTK format files consist of a contiguous sequence of *samples* preceded by a header. Each sample is a vector of either 2-byte integers or 4-byte floats. 2-byte integers are used for compressed forms as described below and for vector quantised data as described later in section 5.11. HTK format data files can also be used to store speech waveforms as described in section 5.8.

The HTK file format header is 12 bytes long and contains the following data

<code>nSamples</code>	– number of samples in file (4-byte integer)
<code>sampPeriod</code>	– sample period in 100ns units (4-byte integer)
<code>sampSize</code>	– number of bytes per sample (2-byte integer)
<code>parmKind</code>	– a code indicating the sample kind (2-byte integer)

The parameter kind consists of a 6 bit code representing the basic parameter kind plus additional bits for each of the possible qualifiers. The basic parameter kind codes are

0	<code>WAVEFORM</code>	sampled waveform
1	<code>LPC</code>	linear prediction filter coefficients
2	<code>LPREFC</code>	linear prediction reflection coefficients
3	<code>LPCEPSTRA</code>	LPC cepstral coefficients
4	<code>LPDELCEP</code>	LPC cepstra plus delta coefficients
5	<code>IREFC</code>	LPC reflection coef in 16 bit integer format
6	<code>MFCC</code>	mel-frequency cepstral coefficients
7	<code>FBANK</code>	log mel-filter bank channel outputs
8	<code>MELSPEC</code>	linear mel-filter bank channel outputs
9	<code>USER</code>	user defined sample kind
10	<code>DISCRETE</code>	vector quantised data

and the bit-encoding for the qualifiers (in octal) is

_E	000100	has energy
_N	000200	absolute energy suppressed
_D	000400	has delta coefficients
_A	001000	has acceleration coefficients
_C	002000	is compressed
_Z	004000	has zero mean static coef.
_K	010000	has CRC checksum
_O	020000	has 0'th cepstral coef.

The `_A` qualifier can only be specified when `_D` is also specified. The `_N` qualifier is only valid when both energy and delta coefficients are present. The sample kind `LPDELCEP` is identical to `LPCEPSTRA_D` and is retained for compatibility with older versions of HTK. The `_C` and `_K` only exist in external files. Compressed files are always decompressed on loading and any attached CRC is checked and removed. An external file can contain both an energy term and a 0'th order cepstral coefficient. These may be retained on loading but normally one or the other is discarded⁷.

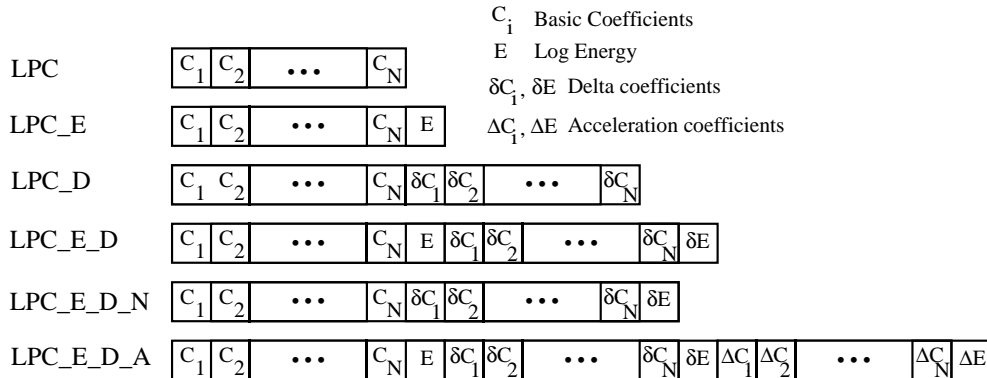


Fig. 5.4 Parameter Vector Layout in HTK Format Files

All parameterised forms of HTK data files consist of a sequence of vectors. Each vector is organised as shown by the examples in Fig 5.4 where various different qualified forms are listed. As can be seen, an energy value if present immediately follows the base coefficients. If delta coefficients are added, these follow the base coefficients and energy value. Note that the base form `LPC` is used in this figure only as an example, the same layout applies to all base sample kinds. If the 0'th order cepstral coefficient is included as well as energy then it is inserted immediately before the energy coefficient, otherwise it replaces it.

For external storage of speech parameter files, two compression methods are provided. For LP coding only, the `IREFC` parameter kind exploits the fact that the reflection coefficients are bounded by ± 1 and hence they can be stored as scaled integers such that $+1.0$ is stored as 32767 and -1.0 is stored as -32767 . For other types of parameterisation, a more general compression facility indicated by the `_C` qualifier is used. HTK compressed parameter files consist of a set of compressed parameter vectors stored as shorts such that for parameter x

$$x_{short} = A * x_{float} - B$$

The coefficients A and B are defined as

$$A = 2 * I / (x_{max} - x_{min})$$

$$B = (x_{max} + x_{min}) * I / (x_{max} - x_{min})$$

⁷ Some applications may require the 0'th order cepstral coefficient in order to recover the filterbank coefficients from the cepstral coefficients.

where x_{max} is the maximum value of parameter x in the whole file and x_{min} is the corresponding minimum. I is the maximum range of a 2-byte integer i.e. 32767. The values of A and B are stored as two floating point vectors prepended to the start of the file immediately after the header.

When a HTK tool writes out a speech file to external storage, no further signal conversions are performed. Thus, for most purposes, the target parameter kind specifies both the required internal representation and the form of the written output, if any. However, there is a distinction in the way that the external data is actually stored. Firstly, it can be compressed as described above by setting the configuration parameter `SAVECOMPRESSED` to true. If the target kind is `LPREFC` then this compression is implemented by converting to `IREFC` otherwise the general compression algorithm described above is used. Secondly, in order to avoid data corruption problems, externally stored HTK parameter files can have a cyclic redundancy checksum appended. This is indicated by the qualifier `.K` and it is generated by setting the configuration parameter `SAVEWITHCRC` to true. The principle tool which uses these output conversions is `HCOPY` (see section 5.13).

5.7.2 Esignal Format Parameter Files

The default for parameter files is native HTK format. However, HTK tools also support the Entropic Esignal format for both input and output. Esignal replaces the Entropic ESPS file format. To ensure compatibility Entropic provides conversion programs from ESPS to ESIG and vice versa.

To indicate that a source file is in Esignal format the configuration variable `SOURCEFORMAT` should be set to `ESIG`. Alternatively, `-F ESIG` can be specified as a command-line option. To generate Esignal format output files, the configuration variable `TARGETFORMAT` should be set to `ESIG` or the command line option `-O ESIG` should be set.

ESIG files consist of three parts: a preamble, a sequence of field specifications called the field list and a sequence of records. The preamble and the field list together constitute the header. The preamble is purely ASCII. Currently it consists of 6 information items that are all terminated by a new line. The information in the preamble is the following:

```

line 1      - identification of the file format
line 2      - version of the file format
line 3      - architecture (ASCII, EDR1, EDR2, machine name)
line 4      - preamble size (48 bytes)
line 5      - total header size
line 6      - record size

```

All ESIG files that are output by HTK programs contain the following global fields:

commandLine the command-line used to generate the file;

recordFreq a double value that indicates the sample frequency in Herz;

startTime a double value that indicates a time at which the first sample is presumed to be starting;

parmKind a character string that indicates the full type of parameters in the file, e.g: `MFCC.E.D`.

source.1 if the input file was an ESIG file this field includes the header items in the input file.

After that there are field specifiers for the records. The first specifier is for the basekind of the parameters, e.g: `MFCC`. Then for each available qualifier there are additional specifiers. Possible specifiers are:

```

zeroc
energy
delta
delta.zeroc
delta.energy
accs
accs.zeroc
accs.energy

```

The data segments of the ESIG files have exactly the same format as the the corresponding HTK files. This format was described in the previous section.

HTK can only input parameter files that have a valid parameter kind as value of the header field `parmKind`. If this field does not exist or if the value of this field does not contain a valid parameter kind, the file is rejected. After the header has been read the file is treated as an HTK file.

5.8 Waveform File Formats

For reading waveform data files, HTK can support a variety of different formats and these are all briefly described in this section. The default speech file format is HTK. If a different format is to be used, it can be specified by setting the configuration parameter `SOURCEFORMAT`. However, since file formats need to be changed often, they can also be set individually via the `-F` command-line option. This over-rides any setting of the `SOURCEFORMAT` configuration parameter.

Similarly for the output of waveforms, the format can be set using either the configuration parameter `TARGETFORMAT` or the `-O` command-line option. However, for output only native HTK format (HTK), Esignal format (ESIG) and headerless (NOHEAD) waveform files are supported.

The following sub-sections give a brief description of each of the waveform file formats supported by HTK.

5.8.1 HTK File Format

The HTK file format for waveforms is identical to that described in section 5.7 above. It consists of a 12 byte header followed by a sequence of 2 byte integer speech samples. For waveforms, the `sampSize` field will be 2 and the `parmKind` field will be 0. The `sampPeriod` field gives the sample period in 100ns units, hence for example, it will have the value 1000 for speech files sampled at 10kHz and 625 for speech files sampled at 16kHz.

5.8.2 Esignal File Format

The Esignal file format for waveforms is similar to that described in section 5.7 above with the following exceptions. When reading an ESIG waveform file the HTK programs only check whether the record length equals 2 and whether the datatype of the only field in the data records is `SHORT`. The data field that is created on output of a waveform is called `WAVEFORM`.

5.8.3 TIMIT File Format

The TIMIT format has the same structure as the HTK format except that the 12-byte header contains the following

```

hdrSize      - number of bytes in header ie 12 (2-byte integer)
version      - version number (2-byte integer)
numChannels  - number of channels (2-byte integer)
sampRate     - sample rate (2-byte integer)
nSamples     - number of samples in file (4-byte integer)

```

TIMIT format data is used only on the prototype TIMIT CD ROM.

5.8.4 NIST File Format

The NIST file format is also referred to as the Sphere file format. A NIST header consists of ASCII text. It begins with a label of the form `NISTxx` where `xx` is a version code followed by the number of bytes in the header. The remainder of the header consists of name value pairs of which HTK decodes the following

```

sample_rate      - sample rate in Hz
sample_n_bytes   - number of bytes in each sample
sample_count     - number of samples in file
sample_byte_format - byte order
sample_coding    - speech coding eg pcm,  $\mu$ law, shortpack
channels_interleaved - for 2 channel data only

```

The current NIST Sphere data format subsumes a variety of internal data organisations. HTK currently supports interleaved μ law used in Switchboard, Shortpack compression used in the original version of WSJ0 and standard 16bit linear PCM as used in Resource Management, TIMIT, etc. It does not currently support the Shorten compression format as used in WSJ1 due to licensing restrictions. Hence, to read WSJ1, the files must be converted using the NIST supplied decompression routines into standard 16 bit linear PCM. This is most conveniently done under UNIX by using the decompression program as an input filter set via the environment variable `HWAVEFILTER` (see section 4.8).

For interleaved μ law as used in Switchboard, the default is to add the two channels together. The left channel only can be obtained by setting the environment variable `STEREOMODE` to `LEFT` and the right channel only can be obtained by setting the environment variable `STEREOMODE` to `RIGHT`.

5.8.5 SCRIBE File Format

The SCRIBE format is a subset of the standard laid down by the European Esprit Programme SAM Project. SCRIBE data files are headerless and therefore consist of just a sequence of 16 bit sample values. HTK assumes by default that the sample rate is 20kHz. The configuration parameter `SOURCERATE` should be set to over-ride this. The byte ordering assumed for SCRIBE data files is `VAX` (little-endian).

5.8.6 SDES1 File Format

The SDES1 format refers to the “Sound Designer I” format defined by Digidesign Inc in 1985 for multimedia and general audio applications. It is used for storing short monoaural sound samples. The SDES1 header is complex (1336 bytes) since it allows for associated display window information to be stored in it as well as providing facilities for specifying repeat loops. The HTK input routine for this format just picks out the following information

```

headerSize  – size of header ie 1336 (2 byte integer)
              (182 byte filler)
fileSize    – number of bytes of sampled data (4 byte integer)
              (832 byte filler)
sampRate    – sample rate in Hz (4 byte integer)
sampPeriod  – sample period in microseconds (4 byte integer)
sampSize    – number of bits per sample ie 16 (2 byte integer)

```

5.8.7 AIFF File Format

The AIFF format was defined by Apple Computer for storing monoaural and multichannel sampled sounds. An AIFF file consists of a number of *chunks*. A *Common* chunk contains the fundamental parameters of the sound (sample rate, number of channels, etc) and a *Sound Data* chunk contains sampled audio data. HTK only partially supports AIFF since some of the information in it is stored as floating point numbers. In particular, the sample rate is stored in this form and to avoid portability problems, HTK ignores the given sample rate and assumes that it is 16kHz. If this default rate is incorrect, then the true sample period should be specified by setting the `SOURCERATE` configuration parameter. Full details of the AIFF format are available from Apple Developer Technical Support.

5.8.8 SUNAU8 File Format

The SUNAU8 format defines a subset of the “.au” and “.snd” audio file format used by Sun and NeXT. An SUNAU8 speech data file consists of a header followed by 8 bit μ law encoded speech samples. The header is 28 bytes and contains the following fields, each of which is 4 bytes

```

magicNumber  – magic number 0x2e736e64
dataLocation – offset to start of data
dataSize     – number of bytes of data
dataFormat   – data format code which is 1 for 8 bit  $\mu$ law
sampRate     – a sample rate code which is always 8012.821 Hz
numChan      – the number of channels
info         – arbitrary character string min length 4 bytes

```


No default byte ordering is assumed for this format. If the data source is known to be different to the machine being used, then the environment variable `BYTEORDER` must be set appropriately. Note that when used on Sun Sparc machines with 16 bit audio device the sampling rate of 8012.821Hz is not supported and playback will be performed at 8KHz.

5.8.9 OGI File Format

The OGI format is similar to TIMIT. The header contains the following

<code>hdrSize</code>	– number of bytes in header
<code>version</code>	– version number (2-byte integer)
<code>numChannels</code>	– number of channels (2-byte integer)
<code>sampRate</code>	– sample rate (2-byte integer)
<code>nSamples</code>	– number of samples in file (4-byte integer)
<code>lendian</code>	– used to test for byte swapping (4-byte integer)

5.8.10 WAV File Format

The WAV file format is a subset of Microsoft’s RIFF specification for the storage of multimedia files. A RIFF file starts out with a file header followed by a sequence of data “chunks”. A WAV file is often just a RIFF file with a single “WAVE” chunk which consists of two sub-chunks - a “fmt” chunk specifying the data format and a “data” chunk containing the actual sample data. The WAV file header contains the following

<code>'RIFF'</code>	– RIFF file identification (4 bytes)
<code><length></code>	– length field (4 bytes)
<code>'WAVE'</code>	– WAVE chunk identification (4 bytes)
<code>'fmt '</code>	– format sub-chunk identification (4 bytes)
<code>flength</code>	– length of format sub-chunk (4 byte integer)
<code>format</code>	– format specifier (2 byte integer)
<code>chans</code>	– number of channels (2 byte integer)
<code>sampsRate</code>	– sample rate in Hz (4 byte integer)
<code>bpssec</code>	– bytes per second (4 byte integer)
<code>bpsample</code>	– bytes per sample (2 byte integer)
<code>bpchan</code>	– bits per channel (2 byte integer)
<code>'data'</code>	– data sub-chunk identification (4 bytes)
<code>dlength</code>	– length of data sub-chunk (4 byte integer)

Support is provided for 8-bit CCITT mu-law, 8-bit CCITT a-law, 8-bit PCM linear and 16-bit PCM linear - all in stereo or mono (use of `STEREOMODE` parameter as per NIST). The default byte ordering assumed for WAV data files is VAX (little-endian).

5.8.11 ALIEN and NOHEAD File Formats

HTK tools can read speech waveform files with alien formats provided that their overall structure is that of a header followed by data. This is done by setting the format to `ALIEN` and setting the environment variable `HEADERSIZE` to the number of bytes in the header. HTK will then attempt to infer the rest of the information it needs. However, if input is from a pipe, then the number of samples expected must be set using the environment variable `NSAMPLES`. The sample rate of the source file is defined by the configuration parameter `SOURCERATE` as described in section 5.2. If the file has no header then the format `NOHEAD` may be specified instead of `ALIEN` in which case `HEADERSIZE` is assumed to be zero.

5.9 Direct Audio Input/Output

Many HTK tools, particularly recognition tools, can input speech waveform data directly from an audio device. The basic mechanism for doing this is to simply specify the `SOURCEKIND` as being `HAUDIO` following which speech samples will be read directly from the host computer’s audio input device.

Note that for live audio input, the configuration variable `ENORMALISE` should be set to false both during training and recognition. Energy normalisation cannot be used with live audio input, and the default setting for this variable is `TRUE`. When training models for live audio input, be sure to set `ENORMALISE` to false. If you have existing models trained with `ENORMALISE` set to true, you can retrain them using *single-pass retraining* (see section 8.6).

When using direct audio input, the input sampling rate may be set explicitly using the configuration parameter `SOURCERATE`, otherwise HTK will assume that it has been set by some external means such as an audio control panel. In the latter case, it must be possible for `HAUDIO` to obtain the sample rate from the audio driver otherwise an error message will be generated.

Although the detailed control of audio hardware is typically machine dependent, HTK provides a number of Boolean configuration variables to request specific input and output sources. These are indicated by the following table

Variable	Source/Sink
<code>LINEIN</code>	line input
<code>MICIN</code>	microphone input
<code>LINEOUT</code>	line output
<code>PHONESOUT</code>	headphones output
<code>SPEAKEROUT</code>	speaker output

The major complication in using direct audio is in starting and stopping the input device. The simplest approach to this is for HTK tools to take direct control and, for example, enable the audio input for a fixed period determined via a command line option. However, the `HAUDIO/HPARM` modules provides two more powerful built-in facilities for audio input control.

The first method of audio input control involves the use of an automatic energy-based speech/silence detector which is enabled by setting the configuration parameter `USESILDET` to true. Note that the speech/silence detector can also operate on waveform input files.

The automatic speech/silence detector uses a two level algorithm which first classifies each frame of data as either speech or silence and then applies a heuristic to determine the start and end of each utterance. The detector classifies each frame as speech or silence based solely on the log energy of the signal. When the energy value exceeds a threshold the frame is marked as speech otherwise as silence. The threshold is made up of two components both of which can be set by configuration variables. The first component represents the mean energy level of silence and can be set explicitly via the configuration parameter `SILEENERGY`. However, it is more usual to take a measurement from the environment directly. Setting the configuration parameter `MEASURESIL` to true will cause the detector to calibrate its parameters from the current acoustic environment just prior to sampling. The second threshold component is the level above which frames are classified as speech (`SPEECHTHRESH`). Once each frame has been classified as speech or silence they are grouped into windows consisting of `SPCSEQCOUNT` consecutive frames. When the number of frames marked as silence within each window falls below a glitch count the whole window is classed as speech. Two separate glitch counts are used, `SPCGLCHCOUNT` before speech onset is detected and `SILGLCHCOUNT` whilst searching for the end of the utterance. This allows the algorithm to take account of the tendency for the end of an utterance to be somewhat quieter than the beginning. Finally, a top level heuristic is used to determine the start and end of the utterance. The heuristic defines the start of speech as the beginning of the first window classified as speech. The actual start of the processed utterance is `SILMARGIN` frames before the detected start of speech to ensure that when the speech detector triggers slightly late the recognition accuracy is not affected. Once the start of the utterance has been found the detector searches for `SILSEQCOUNT` windows all classified as silence and sets the end of speech to be the end of the last window classified as speech. Once again the processed utterance is extended `SILMARGIN` frames to ensure that if the silence detector has triggered slightly early the whole of the speech is still available for further processing.

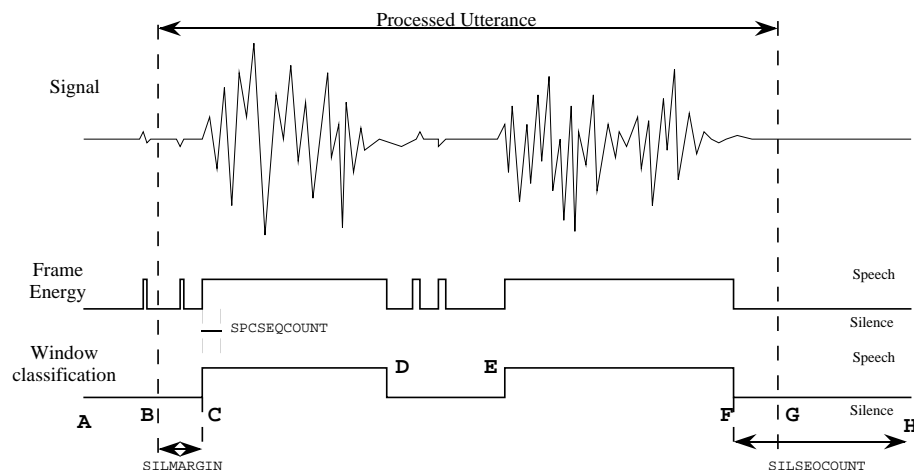


Fig. 5.5 Endpointer Parameters

Fig 5.5 shows an example of the speech/silence detection process. The waveform data is first classified as speech or silence at frame and then at window level before finally the start and end of the utterance are marked. In the example, audio input starts at point A and is stopped automatically at point H. The start of speech, C, occurs when a window of `SPCSEQCOUNT` frames are classified as speech and the start of the utterance occurs `SILMARGIN` frames earlier at B. The period of silence from D to E is not marked as the end of the utterance because it is shorter than `SILSEQCOUNT`. However after point F no more windows are classified as speech (although a few frames are) and so this is marked as the end of speech with the end of the utterance extended to G.

The second built-in mechanism for controlling audio input is by arranging for a signal to be sent from some other process. Sending the signal for the first time starts the audio device. If the speech detector is not enabled then sampling starts immediately and is stopped by sending the signal a second time. If automatic speech/silence detection is enabled, then the first signal starts the detector. Sampling stops immediately when a second signal is received or when silence is detected. The signal number is set using the configuration parameter `AUDIOSIG`. Keypress control operates in a similar fashion and is enabled by setting the configuration parameter `AUDIOSIG` to a negative number. In this mode an initial keypress will be required to start sampling/speech detection and a second keypress will stop sampling immediately.

Audio output is also supported by HTK. There are no generic facilities for output and the precise behaviour will depend on the tool used. It should be noted, however, that the audio input facilities provided by `HAUDIO` include provision for attaching a *replay buffer* to an audio input channel. This is typically used to store the last few seconds of each input to a recognition tool in a circular buffer so that the last utterance input can be replayed on demand.

5.10 Multiple Input Streams

As noted in section 5.1, HTK tools regard the input observation sequence as being divided into a number of independent *data streams*. For building continuous density HMM systems, this facility is of limited use and by far the most common case is that of a single data stream. However, when building tied-mixture systems or when using vector quantisation, a more uniform coverage of the acoustic space is obtained by separating energy, deltas, etc., into separate streams.

This separation of parameter vectors into streams takes place at the point where the vectors are extracted from the converted input file or audio device and transformed into an observation. The tools for HMM construction and for recognition thus view the input data as a sequence of observations but note that this is entirely internal to HTK. Externally data is always stored as a single sequence of parameter vectors.

When multiple streams are required, the division of the parameter vectors is performed automatically based on the parameter kind. This works according to the following rules.

1 stream single parameter vector. This is the default case.

2 streams if the parameter vector contains energy terms, then they are extracted and placed in stream 2. Stream 1 contains the remaining static coefficients and their deltas and accelerations, if any. Otherwise, the parameter vector must have appended delta coefficients and no appended acceleration coefficients. The vector is then split so that the static coefficients form stream 1 and the corresponding delta coefficients form stream 2.

3 streams if the parameter vector has acceleration coefficients, then vector is split with static coefficients plus any energy in stream 1, delta coefficients plus any delta energy in stream 2 and acceleration coefficients plus any acceleration energy in stream 3. Otherwise, the parameter vector must include log energy and must have appended delta coefficients. The vector is then split into three parts so that the static coefficients form stream 1, the delta coefficients form stream 2, and the log energy and delta log energy are combined to form stream 3.

4 streams the parameter vector must include log energy and must have appended delta and acceleration coefficients. The vector is split into 4 parts so that the static coefficients form stream 1, the delta coefficients form stream 2, the acceleration coefficients form stream 3 and the log energy, delta energy and acceleration energy are combined to form stream 4.

In all cases, the static log energy can be suppressed (via the `_N` qualifier). If none of the above rules apply for some required number of streams, then the parameter vector is simply incompatible with that form of observation. For example, the parameter kind `LPC_D_A` cannot be split into 2 streams, instead 3 streams should be used.

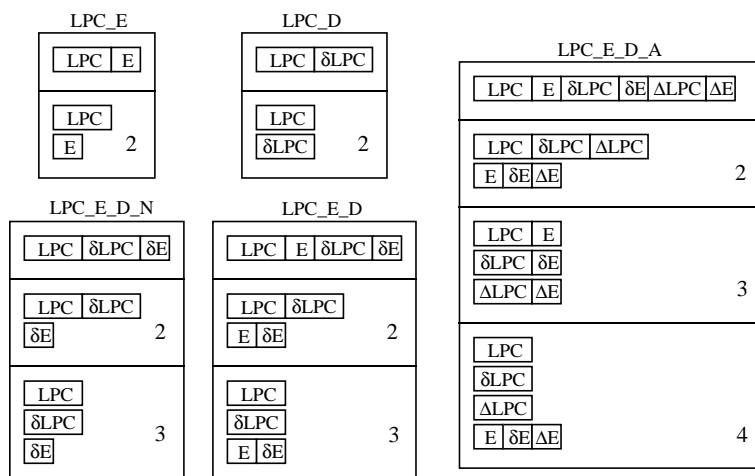


Fig. 5.6 Example Stream Construction

Fig. 5.6 illustrates the way that streams are constructed for a number of common cases. As earlier, the choice of LPC as the static coefficients is purely for illustration and the same mechanism applies to all base parameter kinds.

As discussed further in the next section, multiple data streams are often used with vector quantised data. In this case, each VQ symbol per input sample is placed in a separate data stream.

5.11 Vector Quantisation

Although HTK was designed primarily for building continuous density HMM systems, it also supports discrete density HMMs. Discrete HMMs are particularly useful for modelling data which is naturally symbolic. They can also be used with continuous signals such as speech by quantising each speech vector to give a unique VQ symbol for each input frame. The HTK module HVQ provides a basic facility for performing this vector quantisation. The VQ table (or codebook) can be constructed using the HTK tool HQANT.

When used with speech, the principle justification for using discrete HMMs is the much reduced computation. However, the use of vector quantisation introduces errors and it can lead to rather

fragile systems. For this reason, the use of continuous density systems is generally preferred. To facilitate the use of continuous density systems when there are computational constraints, HTK also allows VQ to be used as the basis for pre-selecting a subset of Gaussian components for evaluation at each time frame.

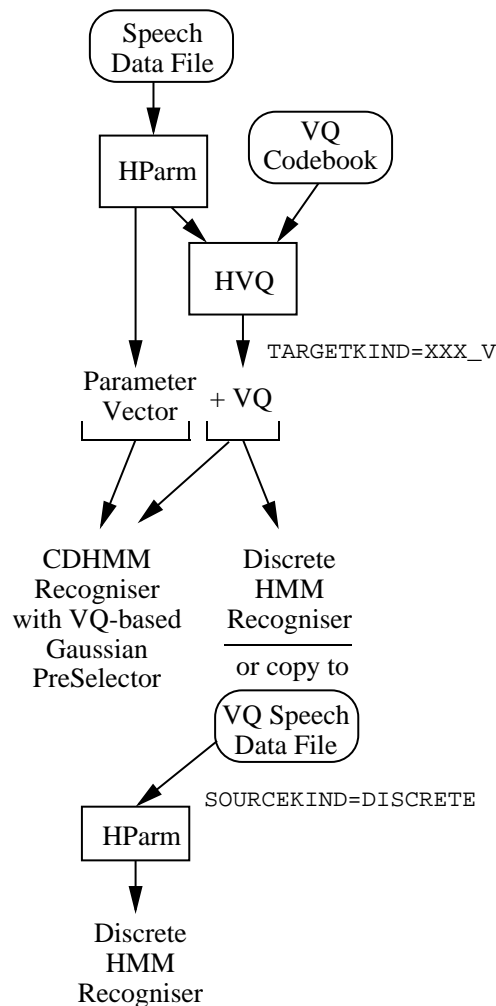


Fig. 5.7 Using Vector Quantisation

Fig. 5.7 illustrates the different ways that VQ can be used in HTK for a single data stream. For multiple streams, the same principles are applied to each stream individually. A converted speech waveform or file of parameter vectors can have VQ indices attached simply by specifying the name of a VQ table using the configuration parameter `VQTABLE` and by adding the `_V` qualifier to the target kind. The effect of this is that each *observation* passed to a recogniser can include both a conventional parameter vector and a VQ index. For continuous density HMM systems, a possible use of this might be to preselect Gaussians for evaluation (but note that HTK does not currently support this facility).

When used with a discrete HMM system, the continuous parameter vectors are ignored and only the VQ indices are used. For training and evaluating discrete HMMs, it is convenient to store speech data in vector quantised form. This is done using the tool `HCOPY` to read in and vector quantise each speech file. Normally, `HCOPY` copies the target form directly into the output file. However, if the configuration parameter `SAVEASVQ` is set, then it will store only the VQ indices and mark the kind of the newly created file as `DISCRETE`. Discrete files created in this way can be read directly by `HPARM` and the VQ symbols passed directly to a tool as indicated by the lower part of Fig. 5.7.

`HVQ` supports three types of distance metric and two organisations of VQ codebook. Each codebook consists of a collection of nodes where each node has a mean vector and optionally a

covariance matrix or diagonal variance vector. The corresponding distance metric used for each of these is simple Euclidean, full covariance Mahalanobis or diagonal covariance Mahalanobis. The codebook nodes are arranged in the form of a simple linear table or as a binary tree. In the linear case, the input vector is compared with every node in turn and the nearest determines the VQ index. In the binary tree case, each non-terminal node has a left and a right daughter. Starting with the top-most root node, the input is compared with the left and right daughter node and the nearest is selected. This process is repeated until a terminal node is reached.

VQ Tables are stored externally in text files consisting of a header followed by a sequence of node entries. The header consists of the following information

```

magic      - a magic number usually the original parameter kind
type       - 0 = linear tree, 1 = binary tree
mode       - 1 = diagonal covariance Mahalanobis
              2 = full covariance Mahalanobis
              5 = Euclidean
numNodes  - total number of nodes in the codebook
numS      - number of independent data streams
sw1,sw2,... - width of each data stream

```

Every node has a unique integer identifier and consists of the following

```

stream    - stream number for this node
vqidx     - VQ index for this node (0 if non-terminal)
nodeId    - integer id of this node
leftId    - integer id of left daughter node
rightId   - integer id of right daughter node
mean      - mean vector
cov       - diagonal variance or full covariance

```

The inclusion of the optional variance vector or covariance matrix depends on the mode in the header. If present they are stored in inverse form. In a binary tree, the root id is always 1. In linear codebooks, the left and right daughter node id's are ignored.

5.12 Viewing Speech with HList

As mentioned in section 5.1, the tool HLIST provides a dual rôle in HTK. Firstly, it can be used for examining the contents of speech data files. In general, HLIST displays three types of information

1. *source header*: requested using the `-h` option
2. *target header*: requested using the `-t` option
3. *target data*: printed by default. The begin and end samples of the displayed data can be specified using the `-s` and `-e` options.

When the default configuration parameters are used, no conversions are applied and the target data is identical to the contents of the file.

As an example, suppose that the file called `timit.wav` holds speech waveform data using the TIMIT format. The command

```
HList -h -e 49 -F TIMIT timit.wav
```

would display the source header information and the first 50 samples of the file. The output would look something like the following

```

----- Source: timit.wav -----
Sample Bytes:  2      Sample Kind:  WAVEFORM
Num Comps:    1      Sample Period: 62.5 us
Num Samples:  31437  File Format:   TIMIT
----- Samples: 0->49 -----
      0:      8      -4      -1      0      -2      -1      -3      -2      0      0

```

```

10:  -1    0   -1   -2   -1    1    0   -1   -2    1
20:  -2    0    0    0    2    1   -2    2    1    0
30:   1    0    0   -1    4    2    0   -1    4    0
40:   2    2    1   -1   -1    1    1    2    1    1
----- END -----

```

The source information confirms that the file contains WAVEFORM data with 2 byte samples and 31437 samples in total. The sample period is $62.5\mu s$ which corresponds to a 16kHz sample rate. The displayed data is numerically small because it corresponds to leading silence. Any part of the file could be viewed by suitable choice of the begin and end sample indices. For example,

```
HList -s 5000 -e 5049 -F TIMIT timit.wav
```

would display samples 5000 through to 5049. The output might look like the following

```

----- Samples: 5000->5049 -----
5000:   85  -116  -159  -252   23   99   69   92   79  -166
5010: -100  -123  -111   48  -19   15  111   41 -126  -304
5020: -189   91  162  255   80 -134 -174  -55   57  155
5030:   90   -1   33  154   68 -149  -70   91  165  240
5040:  297   50   13   72  187  189  193  244  198  128
----- END -----

```

The second use of HLIST is to check that input conversions are being performed properly. Suppose that the above TIMIT format file is part of a database to be used for training a recogniser and that mel-frequency cepstra are to be used along with energy and the first differential coefficients. Suitable configuration parameters needed to achieve this might be as follows

```

# Wave -> MFCC config file
SOURCEFORMAT = TIMIT      # same as -F TIMIT
TARGETKIND   = MFCC_E_D   # MFCC + Energy + Deltas
TARGETRATE   = 100000     # 10ms frame rate
WINDOWSIZE   = 200000    # 20ms window
NUMCHANS     = 24         # num filterbank chans
NUMCEPS     = 8           # compute c1 to c8

```

HLIST can be used to check this. For example, typing

```
HList -C config -o -h -t -s 100 -e 104 -i 9 timit.wav
```

will cause the waveform file to be converted, then the source header, the target header and parameter vectors 100 through to 104 to be listed. A typical output would be as follows

```

----- Source: timit.wav -----
Sample Bytes:  2      Sample Kind:  WAVEFORM
Num Comps:    1      Sample Period: 62.5 us
Num Samples:  31437  File Format:   TIMIT
----- Target -----
Sample Bytes:  72      Sample Kind:  MFCC_E_D
Num Comps:    18      Sample Period: 10000.0 us
Num Samples:  195     File Format:   HTK
----- Observation Structure -----
x:  MFCC-1 MFCC-2 MFCC-3 MFCC-4 MFCC-5 MFCC-6 MFCC-7 MFCC-8      E
     Del-1  Del-2  Del-3  Del-4  Del-5  Del-6  Del-7  Del-8  DelE
----- Samples: 100->104 -----
100:  3.573 -19.729 -1.256 -6.646 -8.293 -15.601 -23.404 10.988 0.834
      3.161 -1.913  0.573 -0.069 -4.935  2.309 -5.336  2.460 0.080
101:  3.372 -16.278 -4.683 -3.600 -11.030 -8.481 -21.210 10.472 0.777
      0.608 -1.850 -0.903 -0.665 -2.603 -0.194 -2.331  2.180 0.069
102:  2.823 -15.624 -5.367 -4.450 -12.045 -15.939 -22.082 14.794 0.830
      -0.051  0.633 -0.881 -0.067 -1.281 -0.410  1.312  1.021 0.005
103:  3.752 -17.135 -5.656 -6.114 -12.336 -15.115 -17.091 11.640 0.825
      -0.002 -0.204  0.015 -0.525 -1.237 -1.039  1.515  1.007 0.015

```

```

104:   3.127 -16.135  -5.176  -5.727 -14.044 -14.333 -18.905  15.506   0.833
      -0.034  -0.247   0.103  -0.223  -1.575   0.513   1.507   0.754   0.006
----- END -----

```

The target header information shows that the converted data consists of 195 parameter vectors, each vector having 18 components and being 72 bytes in size. The structure of each parameter vector is displayed as a simple sequence of floating-point numbers. The layout information described in section 5.7 can be used to interpret the data. However, including the `-o` option, as in the example, causes HLIST to output a schematic of the observation structure. Thus, it can be seen that the first row of each sample contains the static coefficients and the second contains the delta coefficients. The energy is in the final column. The command line option `-i 9` controls the number of values displayed per line and can be used to aid in the visual interpretation of the data. Notice finally that the command line option `-F TIMIT` was not required in this case because the source format was specified in the configuration file.

It should be stressed that when HLIST displays parameterised data, it does so in exactly the form that *observations* are passed to a HTK tool. So, for example, if the above data was input to a system built using 3 data streams, then this can be simulated by using the command line option `-n` to set the number of streams. For example, typing

```
HList -C config -n 3 -o -s 100 -e 101 -i 9 timit.wav
```

would result in the following output

```

----- Observation Structure -----
nTotal=18 nStatic=8 nDel=16 eSep=T
x.1:   MFCC-1 MFCC-2 MFCC-3 MFCC-4 MFCC-5 MFCC-6 MFCC-7 MFCC-8
x.2:   Del-1  Del-2  Del-3  Del-4  Del-5  Del-6  Del-7  Del-8
x.3:           E    DelE
----- Samples: 100->101 -----
100.1:  3.573 -19.729 -1.256 -6.646 -8.293 -15.601 -23.404  10.988
100.2:  3.161 -1.913  0.573 -0.069 -4.935  2.309 -5.336  2.460
100.3:  0.834  0.080
101.1:  3.372 -16.278 -4.683 -3.600 -11.030 -8.481 -21.210  10.472
101.2:  0.608 -1.850 -0.903 -0.665 -2.603 -0.194 -2.331  2.180
101.3:  0.777  0.069
----- END -----

```

Notice that the data is identical to the previous case, but it has been re-organised into separate streams.

5.13 Copying and Coding using HCopy

HCOPY is a general-purpose tool for copying and manipulating speech files. The general form of invocation is

```
HCopy src tgt
```

which will make a new copy called `tgt` of the file called `src`. HCOPY can also concatenate several sources together as in

```
HCopy src1 + src2 + src3 tgt
```

which concatenates the contents of `src1`, `src2` and `src3`, storing the results in the file `tgt`. As well as putting speech files together, HCOPY can also take them apart. For example,

```
HCopy -s 100 -e -100 src tgt
```

will extract samples 100 through to N-100 of the file `src` to the file `tgt` where N is the total number of samples in the source file. The range of samples to be copied can also be specified with reference to a label file, and modifications made to the speech file can be tracked in a copy of the label file. All of the various options provided by HCOPY are given in the reference section and in total they provide a powerful facility for manipulating speech data files.

However, the use of HCOPY extends beyond that of copying, chopping and concatenating files. HCOPY reads in all files using the speech input/output subsystem described in the preceding sections. Hence, by specifying an appropriate configuration file, HCOPY is also a speech coding tool. For example, if the configuration file `config` was set-up to convert waveform data to MFCC coefficients, the command

```
HCOPY -C config -s 100 -e -100 src.wav tgt.mfc
```

would parameterise the file waveform file `src.wav`, excluding the first and last 100 samples, and store the result in `tgt.mfc`.

HCOPY will process its arguments in pairs, and as with all HTK tools, argument lists can be written in a script file specified via the `-S` option. When coding a large database, the separate invocation of HCOPY for each file needing to be processed would incur a very large overhead. Hence, it is better to create a file, `flist` say, containing a list of all source and target files, as in for example,

```
src1.wav tgt1.mfc
src2.wav tgt2.mfc
src3.wav tgt3.mfc
src4.wav tgt4.mfc
etc
```

and then invoke HCOPY by

```
HCOPY -C config -s 100 -e -100 -S flist
```

which would encode each file listed in `flist` in a single invocation.

Normally HCOPY makes a direct copy of the target speech data in the output file. However, if the configuration parameter `SAVECOMPRESSED` is set true then the output is saved in compressed form and if the configuration parameter `SAVEWITHCRC` is set true then a checksum is appended to the output (see section 5.7). If the configuration parameter `SAVEASVQ` is set true then only VQ indices are saved and the kind of the target file is changed to `DISCRETE`. For this to work, the target kind must have the qualifier `_V` attached (see section 5.11).

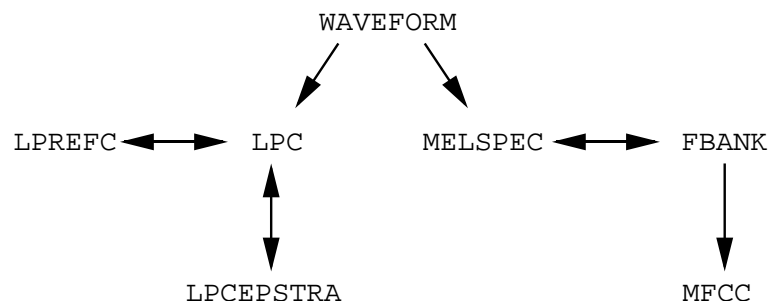


Fig. 5.8 Valid Parameter Kind Conversions

5.14 Version 1.5 Compatibility

The redesign of the HTK front-end in version 2 has introduced a number of differences in parameter encoding. The main changes are

1. Source waveform zero mean processing is now performed on a frame-by-frame basis.
2. Delta coefficients use a modified form of regression rather than simple differences at the start and end of the utterance.
3. Energy scaling is no longer applied to the zero'th MFCC coefficient.

If a parameter encoding is required which is as close as possible to the version 1.5 encoding, then the compatibility configuration variable `V1COMPAT` should be set to true.

Note also in this context that the default values for the various configuration values have been chosen to be consistent with the defaults or recommended practice for version 1.5.

5.15 Summary

This section summarises the various file formats, parameter kinds, qualifiers and configuration parameters used by HTK. Table 5.1 lists the audio speech file formats which can be read by the HWAVE module. Table 5.2 lists the basic parameter kinds supported by the HPARM module and Fig. 5.8 shows the various automatic conversions that can be performed by appropriate choice of source and target parameter kinds. Table 5.3 lists the available qualifiers for parameter kinds. The first 6 of these are used to describe the target kind. The source kind may already have some of these, HPARM adds the rest as needed. Note that HPARM can also delete qualifiers when converting from source to target. The final two qualifiers in Table 5.3 are only used in external files to indicate compression and an attached checksum. HPARM adds these qualifiers to the target form during output and only in response to setting the configuration parameters SAVECOMPRESSED and SAVEWITHCRC. Adding the `_C` or `_K` qualifiers to the target kind simply causes an error. Finally, Tables 5.4 and 5.5 lists all of the configuration parameters along with their meaning and default values.

Name	Description
HTK	The standard HTK file format
TIMIT	As used in the original prototype TIMIT CD-ROM
NIST	The standard SPHERE format used by the US NIST
SCRIBE	Subset of the European SAM standard used in the SCRIBE CD-ROM
SDES1	The Sound Designer 1 format defined by Digidesign Inc.
AIFF	Audio interchange file format
SUNAU8	Subset of 8bit ".au" and ".snd" formats used by Sun and NeXT
OGI	Format used by Oregon Graduate Institute similar to TIMIT
WAV	Microsoft WAVE files used on PCs
ESIG	Entropic Esignal file format
AUDIO	Pseudo format to indicate direct audio input
ALIEN	Pseudo format to indicate unsupported file, the alien header size must be set via the environment variable HDSIZE
NOHEAD	As for the ALIEN format but header size is zero

Table 5.1 Supported File Formats

Kind	Meaning
WAVEFORM	scalar samples (usually raw speech data)
LPC	linear prediction coefficients
LPREFC	linear prediction reflection coefficients
LPCEPSTRA	LP derived cepstral coefficients
LPDELCEP	LP cepstra + delta coef (obsolete)
IREFC	LPREFC stored as 16bit (short) integers
MFCC	mel-frequency cepstral coefficients
FBANK	log filter-bank parameters
MELSPEC	linear filter-bank parameters
USER	user defined parameters
DISCRETE	vector quantised codebook symbols
ANON	matches actual parameter kind

Table 5.2 Supported Parameter Kinds

Qualifier	Meaning
.A	Acceleration coefficients appended
.C	External form is compressed
.D	Delta coefficients appended
.E	Log energy appended
.K	External form has checksum appended
.N	Absolute log energy suppressed
.V	VQ index appended
.Z	Cepstral mean subtracted
.0	Cepstral C0 coefficient appended

Table. 5.3 Parameter Kind Qualifiers

Module	Name	Default	Description
HAUDIO	LINEIN	T	Select line input for audio
HAUDIO	MICIN	F	Select microphone input for audio
HAUDIO	LINEOUT	T	Select line output for audio
HAUDIO	SPEAKEROUT	F	Select speaker output for audio
HAUDIO	PHONESOUT	T	Select headphones output for audio
	SOURCEKIND	ANON	Parameter kind of source
	SOURCEFORMAT	HTK	File format of source
	SOURCERATE	0.0	Sample period of source in 100ns units
HWAVE	NSAMPLES		Num samples in alien file input via a pipe
HWAVE	HEADERSIZE		Size of header in an alien file
HWAVE	STEREOMODE		Select channel: RIGHT or LEFT
HWAVE	BYTEORDER		Define byte order VAX or other
	NATURALREADORDER	F	Enable natural read order for HTK files
	NATURALWRITEORDER	F	Enable natural write order for HTK files
	TARGETKIND	ANON	Parameter kind of target
	TARGETFORMAT	HTK	File format of target
	TARGETRATE	0.0	Sample period of target in 100ns units
HPARM	SAVECOMPRESSED	F	Save the output file in compressed form
HPARM	SAVEWITHCRC	T	Attach a checksum to output parameter file
HPARM	ADDITHER	0.0	Level of noise added to input signal
HPARM	ZMEANSOURCE	F	Zero mean source waveform before analysis
HPARM	WINDOWSIZE	256000.0	Analysis window size in 100ns units
HPARM	USEHAMMING	T	Use a Hamming window
HPARM	PREEMCOEF	0.97	Set pre-emphasis coefficient
HPARM	LPCORDER	12	Order of LPC analysis
HPARM	NUMCHANS	20	Number of filterbank channels
HPARM	LOFREQ	-1.0	Low frequency cut-off in fbank analysis
HPARM	HIFREQ	-1.0	High frequency cut-off in fbank analysis
HPARM	USEPOWER	F	Use power not magnitude in fbank analysis
HPARM	NUMCEPS	12	Number of cepstral parameters
HPARM	CEPLIFTER	22	Cepstral liftering coefficient
HPARM	ENORMALISE	T	Normalise log energy
HPARM	ESCALE	0.1	Scale log energy
HPARM	SILFLOOR	50.0	Energy silence floor (dB)
HPARM	DELTAWINDOW	2	Delta window size
HPARM	ACCWINDOW	2	Acceleration window size
HPARM	VQTABLE	NULL	Name of VQ table
HPARM	SAVEASVQ	F	Save only the VQ indices
HPARM	AUDIOSIG	0	Audio signal number for remote control

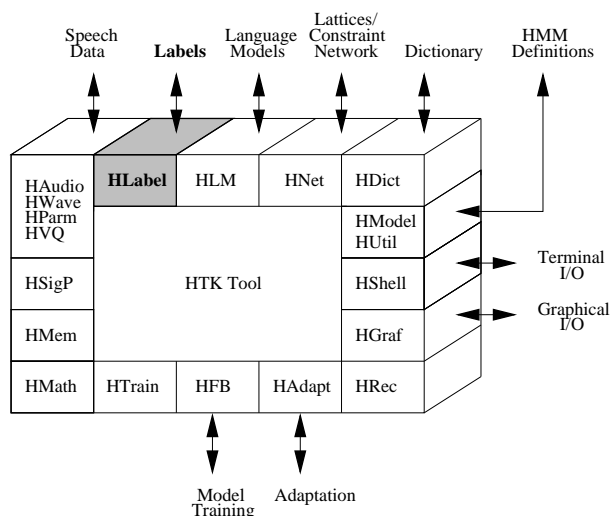
Table. 5.4 Configuration Parameters

Module	Name	Default	Description
HPARM	USESILDET	F	Enable speech/silence detector
HPARM	MEASURESIL	T	Measure background noise level prior to sampling
HPARM	OUTSILWARN	T	Print a warning message to <code>stdout</code> before measuring audio levels
HPARM	SPEECHTHRESH	9.0	Threshold for speech above silence level (dB)
HPARM	SILEENERGY	0.0	Average background noise level (dB)
HPARM	SPCSEQCOUNT	10	Window over which speech/silence decision reached
HPARM	SPCGLCHCOUNT	0	Maximum number of frames marked as silence in window which is classified as speech whilst expecting start of speech
HPARM	SILSEQCOUNT	100	Number of frames classified as silence needed to mark end of utterance
HPARM	SILGLCHCOUNT	2	Maximum number of frames marked as silence in window which is classified as speech whilst expecting silence
HPARM	SILMARGIN	40	Number of extra frames included before and after start and end of speech marks from the speech/silence detector
HPARM	V1COMPAT	F	Set Version 1.5 compatibility mode
	TRACE	0	Trace setting

Table. 5.5 Configuration Parameters (cont)

Chapter 6

Transcriptions and Label Files



Many of the operations performed by HTK which involve speech data files assume that the speech is divided into segments and each segment has a name or *label*. The set of labels associated with a speech file constitute a *transcription* and each transcription is stored in a separate *label file*. Typically, the name of the label file will be the same as the corresponding speech file but with a different extension. For convenience, label files are often stored in a separate directory and all HTK tools have an option to specify this. When very large numbers of files are being processed, label file access can be greatly facilitated by using *Master Label Files (MLFs)*. MLFs may be regarded as index files holding pointers to the actual label files which can either be embedded in the same index file or stored anywhere else in the file system. Thus, MLFs allow large sets of files to be stored in a single file, they allow a single transcription to be shared by many logical label files and they allow arbitrary file redirection.

The HTK interface to label files is provided by the module HLABEL which implements the MLF facility and support for a number of external label file formats. All of the facilities supplied by HLABEL, including the supported label file formats, are described in this chapter. In addition, HTK provides a tool called HLED for simple batch editing of label files and this is also described. Before proceeding to the details, however, the general structure of label files will be reviewed.

6.1 Label File Structure

Most transcriptions are single-alternative and single-level, that is to say, the associated speech file is described by a single sequence of labelled segments. Most standard label formats are of this kind. Sometimes, however, it is useful to have several levels of labels associated with the same basic segment sequence. For example, in training a HMM system it is useful to have both the word level transcriptions and the phone level transcriptions *side-by-side*.

Orthogonal to the requirement for multiple levels of description, a transcription may also need to include multiple alternative descriptions of the same speech file. For example, the output of a speech recogniser may be in the form of an *N-best* list where each word sequence in the list represents one possible interpretation of the input.

As an example, Fig. 6.1 shows a speech file and three different ways in which it might be labelled. In part (a), just a simple orthography is given and this single-level single-alternative type of transcription is the commonest case. Part (b) shows a 2-level transcription where the basic level consists of a sequence of phones but a higher level of word labels are also provided. Notice that there is a distinction between the basic level and the higher levels, since only the basic level has explicit boundary locations marked for every segment. The higher levels do not have explicit boundary information since this can always be inferred from the basic level boundaries. Finally, part (c) shows the case where knowledge of the contents of the speech file is uncertain and three possible word sequences are given.

HTK label files support multiple-alternative and multiple-level transcriptions. In addition to start and end times on the basic level, a label at any level may also have a score associated with it. When a transcription is loaded, all but one specific alternative can be discarded by setting the configuration variable `TRANSALT` to the required alternative `N`, where the first (i.e. normal) alternative is numbered 1. Similarly, all but a specified level can be discarded by setting the configuration variable `TRANSLEV` to the required level number where again the first (i.e. normal) level is numbered 1.

All non-HTK formats are limited to single-level single-alternative transcriptions.

6.2 Label File Formats

As with speech data files, HTK not only defines its own format for label files but also supports a number of external formats. Defining an external format is similar to the case for speech data files except that the relevant configuration variables for specifying a format other than HTK are called `SOURCELABEL` and `TARGETLABEL`. The source label format can also be specified using the `-G` command line option. As with using the `-F` command line option for speech data files, the `-G` option overrides any setting of `SOURCELABEL`.

6.2.1 HTK Label Files

The HTK label format is text based. As noted above, a single label file can contain multiple-alternatives and multiple-levels.

Each line of a HTK label file contains the actual label optionally preceded by start and end times, and optionally followed by a match score.

```
[start [end] ] name [score] { auxname [auxscore] } [comment]
```

where `start` denotes the start time of the labelled segment in 100ns units, `end` denotes the end time in 100ns units, `name` is the name of the segment and `score` is a floating point confidence score. All fields except the name are optional. If `end` is omitted then it is set equal to -1 and ignored. This case would occur with data which had been labelled frame synchronously. If `start` and `end` are both missing then both are set to -1 and the label file is treated as a simple symbolic transcription. The optional score would typically be a log probability generated by a recognition tool. When omitted the score is set to 0.0.

The following example corresponds to the transcription shown in part (a) of Fig. 6.1

```
0000000 3600000 ice
3600000 8200000 cream
```

Multiple levels are described by adding further names alongside the basic name. The lowest level (shortest segments) should be given first since only the lowest level has start and end times. The label file corresponding to the transcription illustrated in part (b) of Fig. 6.1 would be as follows.

```
0000000 2200000 ay      ice
2200000 3600000 s
3600000 4300000 k      cream
4300000 5000000 r
5000000 7400000 iy
7400000 8200000 m
```

Finally, multiple alternatives are written as a sequence of separate label lists separated by three slashes (///). The label file corresponding to the transcription illustrated in part (c) of Fig. 6.1 would therefore be as follows.

```
0000000 2200000 I
2200000 8200000 scream
///
0000000 3600000 ice
3600000 8200000 cream
///
0000000 3600000 eyes
3600000 8200000 cream
```

Actual label names can be any sequence of characters. However, the - and + characters are reserved for identifying the left and right context, respectively, in a context-dependent phone label. For example, the label N-aa+V might be used to denote the phone aa when preceded by a nasal and followed by a vowel. These context-dependency conventions are used in the label editor HLED, and are understood by all HTK tools.

6.2.2 ESPS Label Files

An *ESPS/waves+* label file is a text file with one label stored per line. Each label indicates a segment boundary. A complete description of the *ESPS/waves+* label format is given in the *ESPS/waves+* manual pages **xwaves (1-ESPS)** and **xlabel (1-ESPS)**. Only details required for use with HTK are given here.

The label data follows a header which ends with a line containing only a #. The header contents are generally ignored by HLABEL. The labels follow the header in the form

```
time ccode name
```

where **time** is a floating point number which denotes the boundary location in seconds, **ccode** is an integer color map entry used by *ESPS/waves+* in drawing segment boundaries and **name** is the name of the segment boundary. A typical value for **ccode** is 121.

While each HTK label can contain both a start and an end time which indicate the boundaries of a labeled segment, *ESPS/waves+* labels contain a single time in seconds which (by convention) refers to the end of the labeled segment. The starting time of the segment is taken to be the end of the previous segment and 0 initially.

ESPS/waves+ label files may have several boundary names per line. However, HLABEL only reads *ESPS/waves+* label files with a single name per boundary. Multiple-alternative and/or multiple-level HTK label data structures cannot be saved using *ESPS/waves+* format label files.

6.2.3 TIMIT Label Files

TIMIT label files are identical to single-alternative single-level HTK label files without scores except that the start and end times are given as sample numbers rather than absolute times. TIMIT label files are used on both the prototype and final versions of the TIMIT CD ROM.

6.2.4 SCRIBE Label Files

The SCRIBE label file format is a subset of the European SAM label file format. SAM label files are text files and each line begins with a label identifying the type of information stored on that line. The HTK SCRIBE format recognises just three label types

```
LBA      - acoustic label
LBB      - broad class label
UTS      - utterance
```

For each of these, the rest of the line is divided into comma separated fields. The LBA and LBB types have 4 fields: start sample, centre sample, end sample and label. HTK expects the centre sample to be blank. The UTS type has 3 fields: start sample, end sample and label. UTS labels may be multi-word since they can refer to a complete utterance. In order to make such labels usable within HTK tools, between word blanks are converted to underscore characters. The EX command in the HTK label editor HLED can then be used to split such a compound label into individual word labels if required.

6.3 Master Label Files

6.3.1 General Principles of MLFs

Logically, the organisation of data and label files is very simple. Every data file has a label file of the same name (but different extension) which is either stored in the same directory as the data file or in some other specified directory.

ice	cream
-----	-------



(a) 1-alternative, 1-level

ice		cream			
ay	s	k	r	iy	m



(b) 1-alternative, 2-level

I	scream
ice	cream
eyes	cream



(c) 3-alternative, 1-level

Fig. 6.1 Example Transcriptions

This scheme is sufficient for most needs and commendably simple. However, there are many cases where either it makes unnecessarily inefficient use of the operating system or it seriously inconveniences the user. For example, to use a training tool with isolated word data may require the generation of hundreds or thousands of label files each having just one label entry. Even where individual label files are appropriate (as in the phonetically transcribed TIMIT database), each label file must be stored in the same directory as the data file it transcribes, or all label files must be stored in the same directory. One cannot, for example, have a different directory of label files for each TIMIT dialect region and then run the HTK training tool HEREST on the whole database.

All of these problems can be solved by the use of Master Label Files (MLFs). Every HTK tool which uses label files has a `-I` option which can be used to specify the name of an MLF file. When an MLF has been loaded, the normal rules for locating a label file apply except that the MLF is searched first. If the required label file `f` is found via the MLF then that is loaded, otherwise the file `f` is opened as normal. If `f` does not exist, then an error is reported. The `-I` option may be repeated on the command line to open several MLF files simultaneously. In this case, each is searched in turn before trying to open the required file.

MLFs can do two things. Firstly, they can contain embedded label definitions so that many or all of the needed label definitions can be stored in the same file. Secondly, they can contain the names of sub-directories to search for label files. In effect, they allow multiple search paths to be defined. Both of these two types of definition can be mixed in a single MLF.

MLFs are quite complex to understand and use. However, they add considerable power and flexibility to HTK which combined with the `-S` and `-L` options mean that virtually any organisation of data and label files can be accommodated.

6.3.2 Syntax and Semantics

An MLF consists of one or more individual definitions. Blank lines in an MLF are ignored but otherwise the line structure is significant. The first line must contain just `#!MLF!#` to identify it as an MLF file. This is not necessary for use with the `-I` option but some HTK tools need to be able to distinguish an MLF from a normal label file. The following syntax of MLF files is described using an extended BNF notation in which alternatives are separated by a vertical bar `|`, parentheses `()` denote factoring, brackets `[]` denote options, and braces `{ }` denote zero or more repetitions.

```
MLF =      "#!MLF!#"
          MLFDef { MLFDef }
```

Each definition is either a transcription for immediate loading or a subdirectory to search.

```
MLFDef =   ImmediateTranscription | SubDirDef
```

An immediate transcription consists of a pattern on a line by itself immediately followed by a transcription which as far as the MLF is concerned is arbitrary text. It is read using whatever label file “driver” routines are installed in HLABEL. It is terminated by a period written on a line of its own.

```
ImmediateTranscription =
    Pattern
    Transcription
    “.”
```

A subdirectory definition simply gives the name of a subdirectory to search. If the required label file is found in that subdirectory then the label file is loaded, otherwise the next matching subdirectory definition is checked.

```
SubDirDef = Pattern SearchMode String
SearchMode = “->” | “=>”
```

The two types of search mode are described below. A pattern is just a string

```
Pattern =   String
```

except that the characters ‘?’ and ‘*’ embedded in the string act as wildcards such that ‘?’ matches any single character and ‘*’ matches 0 or more characters. A string is any sequence of characters enclosed in double quotes.

6.3.3 MLF Search

The names of label files in HTK are invariably reconstructed from an existing data file name and this means that the file names used to access label files can be partial or full path names in which the path has been constructed either from the path of the corresponding data file or by direct specification via the `-L` option. These path names are retained in the MLF search which proceeds as follows. The given label file specification `../d3/d2/d1/name` is matched against each pattern in the MLF. If a pattern matches, then either the named subdirectory is searched or an immediate definition is loaded. Pattern matching continues in this way until a definition is found. If no pattern matches then an attempt is made to open `../d3/d2/d1/name` directly. If this fails an error is reported.

The search of a sub-directory proceeds as follows. In simple search mode indicated by `->`, the file `name` must occur directly in the sub-directory. In full search mode indicated by `=>`, the files `name`, `d1/name`, `d2/d1/name`, etc. are searched for in that order. This full search allows a hierarchy of label files to be constructed which mirrors a hierarchy of data files (see Example 4 below).

Hashing is performed when the label file specification is either a full path name or in the form `*/file` so in these cases the search is very fast. Any other use of metacharacters invokes a linear

search with a full and relatively slow pattern match at each step. Note that all tools which generate label files have a `-l` option which is used to define the output directory in which to store individual label files. When outputting master label files, the `-l` option can be used to define the path in the output label file specifications. In particular, setting the option `-l '*'` causes the form `*/file` to be generated.

6.3.4 MLF Examples

1. Suppose a data set consisted of two training data files with corresponding label files:
`a.lab` contains

```
000000 590000 sil
600000 2090000 a
2100000 4500000 sil
```

`b.lab` contains

```
000000 990000 sil
1000000 3090000 b
3100000 4200000 sil
```

Then the above two individual label files could be replaced by a single MLF

```
#!MLF!#
"*/a.lab"
000000 590000 sil
600000 2090000 a
2100000 4500000 sil
.
"*/b.lab"
000000 990000 sil
1000000 3090000 b
3100000 4200000 sil
.
```

2. A digit data base contains training tokens `one.1.wav`, `one.2.wav`, `one.3.wav`, ..., `two.1.wav`, `two.2.wav`, `two.3.wav`, ..., etc. Label files are required containing just the name of the model so that HTK tools such as HEREST can be used. If MLFs are not used, individual label files are needed. For example, the individual label files `one.1.lab`, `one.2.lab`, `one.3.lab`, ... would be needed to identify instances of "one" even though each file contains the same entry, just

```
one
```

Using an MLF containing

```
#!MLF!#
"*/one.*.lab"
one
.
"*/two.*.lab"
two
.
"*/three.*.lab"
three
.
<etc.>
```

avoids the need for many duplicate label files.

3. A training database `/db` contains directories `dr1`, `dr2`, ..., `dr8`. Each directory contains a subdirectory called `labs` holding the label files for the data files in that directory. The following MLF would allow them to be found

```
#!MLF!#
"*" -> "/db/dr1/labs"
"*" -> "/db/dr2/labs"
...
"*" -> "/db/dr7/labs"
"*" -> "/db/dr8/labs"
```

Each attempt to open a label file will result in a linear search through `dr1` to `dr8` to find that file. If the sub-directory name is embedded into the label file name, then this searching can be avoided. For example, if the label files in directory `drx` had the form `drx_xxxx.lab`, then the MLF would be written as

```
#!MLF!#
"*/dr1_*" -> "/db/dr1/labs"
"*/dr2_*" -> "/db/dr2/labs"
...
"*/dr7_*" -> "/db/dr7/labs"
"*/dr8_*" -> "/db/dr8/labs"
```

4. A training database is organised as a hierarchy where `/disk1/db/dr1/sp2/u3.wav` is the data file for the third repetition from speaker 2 in dialect region `dr1` (see Figure 6.2).

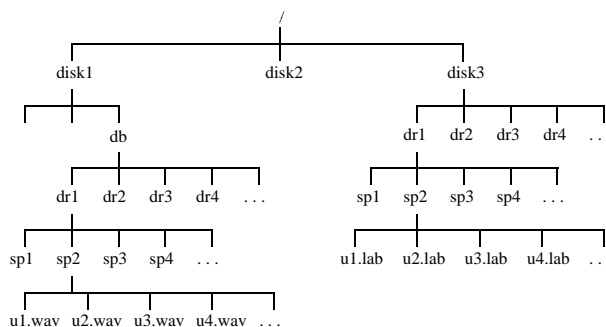


Fig. 6.2 Database Hierarchy: Data [Left]; Labels [Right].

Suppose that a similar hierarchy of label files was constructed on `disk3`. These label files could be found by any HTK tool by using an MLF containing just

```
#!MLF!#
"*" => "/disk3"
```

If for some reason all of the `drN` directories were renamed `ldrN` in the label hierarchy, then this could be handled by an MLF file containing

```
#!MLF!#
"*/dr1/*" => "/disk3/ldr1"
"*/dr2/*" => "/disk3/ldr2"
"*/dr3/*" => "/disk3/ldr3"
etc.
```

These few examples should illustrate the flexibility and power of MLF files. It should be noted, however, that when generating label names automatically from data file names, HTK sometimes discards path details. For example, during recognition, if the data files `/disk1/dr2/sx43.wav` and `/disk2/dr4/sx43.wav` are being recognised, and a single directory is specified for the output label files, then recognition results for both files will be written to a file called `sx43.lab`, and the latter occurrence will overwrite the former.

6.4 Editing Label Files

HTK training tools typically expect the labels used in transcription files to correspond directly to the names of the HMMs chosen to build an application. Hence, the label files supplied with a speech database will often need modifying. For example, the original transcriptions attached to a database might be at a fine level of acoustic detail. Groups of labels corresponding to a sequence of acoustic events (e.g. `pcl p'`) might need converting to some simpler form (e.g. `p`) which is more suitable for being represented by a HMM. As a second example, current high performance speech recognisers use a large number of context dependent models to allow more accurate acoustic modelling. For this case, the labels in the transcription must be converted to show the required contexts explicitly.

HTK supplies a tool called HLED for rapidly and efficiently converting label files. The HLED command invocation specifies the names of the files to be converted and the name of a script file holding the actual HLED commands. For example, the command

```
HLEd edfile.led 11 12 13
```

would apply the edit commands stored in the file `edfile.led` to each of the label files `11`, `12` and `13`. More commonly the new label files are stored in a new directory to avoid overwriting the originals. This is done by using the `-l` option. For example,

```
HLEd -l newlabs edfile.led 11 12 13
```

would have the same effect as previously except that the new label files would be stored in the directory `newlabs`.

Each edit command stored in an edit file is identified by a mnemonic consisting of two letters¹ and must be stored on a separate line. The supplied edit commands can be divided into two groups. The first group consists of commands which perform selective changes to specific labels and the second group contains commands which perform global transformations. The reference section defines all of these commands. Here a few examples will be given to illustrate the use of HLED.

As a first example, when using the TIMIT database, the original 61 phoneme symbol set is often mapped into a simpler 48 phoneme symbol set. The aim of this mapping is to delete all glottal stops, replace all closures preceding a voiced stop by a generic voiced closure (`vc1`), all closures preceding an unvoiced stop by a generic unvoiced closure (`c1`) and the different types of silence to a single generic silence (`sil`). A HLED script to do this might be

```
# Map 61 Phone Timit Set -> 48 Phones
S0
DE q
RE c1 pcl tcl kcl qcl
RE vc1 bcl dcl gcl
RE sil h# #h pau
```

The first line is a comment indicated by the initial hash character. The command on the second line is the *Sort* command `S0`. This is an example of a global command. Its effect is to sort all the labels into time order. Normally the labels in a transcription will already be in time order but some speech editors simply output labels in the order that the transcriber marked them. Since this would confuse the re-estimation tools, it is good practice to explicitly sort all label files in this way.

The command on the third line is the *Delete* command `DE`. This is a selective command. Its effect is to delete all of the labels listed on the rest of the command line, wherever they occur. In this case, there is just one label listed for deletion, the glottal stop `q`. Hence, the overall effect of this command will be to delete all occurrences of the `q` label in the edited label files.

The remaining commands in this example script are *Replace* commands `RE`. The effect of a *Replace* command is to substitute the first label following the `RE` for every occurrence of the remaining

¹ Some command names have single letter alternatives for compatibility with earlier versions of HTK.

labels on that line. Thus, for example, the command on the third line causes all occurrences of the labels `pc1`, `tc1`, `kc1` or `qc1` to be replaced by the label `c1`.

To illustrate the overall effect of the above HLED command script on a complete label file, the following TIMIT format label file

```
0000 2241 h#
2241 2715 w
2715 4360 ow
4360 5478 bc1
5478 5643 b
5643 6360 iy
6360 7269 tc1
7269 8313 t
8313 11400 ay
11400 12950 dc1
12950 14360 dh
14360 14640 h#
```

would be converted by the above script to the following

```
0 1400625 sil
1400625 1696875 w
1696875 2725000 ow
2725000 3423750 vc1
3423750 3526875 b
3526875 3975000 iy
3975000 4543125 c1
4543125 5195625 t
5195625 7125000 ay
7125000 8093750 vc1
8093750 8975000 dh
8975000 9150000 sil
```

Notice that label boundaries in TIMIT format are given in terms of sample numbers (16kHz sample rate), whereas the edited output file is in HTK format in which all times are in absolute 100ns units.

As well as the Replace command, there is also a *Merge* command `ME`. This command is used to replace a sequence of labels by a single label. For example, the following commands would merge the closure and release labels in the previous TIMIT transcription into single labels

```
ME b bc1 b
ME d dc1 dh
ME t tc1 t
```

As shown by this example, the label used for the merged sequence can be the same as occurs in the original but some care is needed since HLED commands are normally applied in sequence. Thus, a command on line n is applied to the label sequence that remains after the commands on lines 1 to $n - 1$ have been applied.

There is one exception to the above rule of sequential edit command application. The *Change* command `CH` provides for context sensitive replacement. However, when a sequence of Change commands occur in a script, the sequence is applied as a block so that the contexts which apply for each command are those that existed just prior to the block being executed. The Change command takes 4 arguments `X A Y B` such that every occurrence of label `Y` in the context of `A _ B` is changed to the label `X`. The contexts `A` and `B` refer to sets of labels and are defined by separate *Define Context* commands `DC`. The `CH` and `DC` commands are primarily used for creating context sensitive labels. For example, suppose that a set of context-dependent phoneme models are needed for TIMIT. Rather than treat all possible contexts separately and build separate triphones for each (see below), the possible contexts will be grouped into just 5 broad classes: `C` (consonant), `V` (vowel), `N` (nasal), `L` (liquid) and `S` (silence). The goal then is to translate a label sequence such as `sil b ah t iy n ...` into `sil+C S-b+V C-ah+C V-t+V C-iy+N V-n+ ...` where the `-` and `+` symbols within a label are recognised by HTK as defining the left and right context, respectively. To perform this transformation, it is necessary to firstly use `DC` commands to define the 5 contexts, that is

```
DC V iy ah ae eh ix ...
DC C t k d k g dh ...
DC L l r w j ...
DC N n m ng ...
DC S h# #h epi ...
```

Having defined the required contexts, a change command must be written for each context dependent triphone, that is

```
CH V-ah+V V ah V
CH V-ah+C V ah C
CH V-ah+N V ah N
CH V-ah+L V ah L
...
etc
```

This script will, of course, be rather long ($25 \times$ number of phonemes) but it can easily be generated automatically by a simple program or shell script.

The previous example shows how to transform a set of phonemes into a context dependent set in which the contexts are user-defined. For convenience, HLED provides a set of global transformation commands for converting phonemic transcriptions to conventional left or right biphones, or full triphones. For example, a script containing the single *Triphone Conversion* command TC will convert phoneme files to regular triphones. As an illustration, applying the TC command to a file containing the sequence `sil b ah t iy n ...` would give the transformed sequence `sil+b sil-b+ah b-ah+t ah-t+iy t-iy+n iy-n+ ...`. Notice that the first and last phonemes in the sequence cannot be transformed in the normal way. Hence, the left-most and right-most contexts of these start and end phonemes can be specified explicitly as arguments to the TC commands if required. For example, the command `TC # #` would give the sequence `#-sil+b sil-b+ah b-ah+t ah-t+iy t-iy+n iy-n+ ... +#`. Also, the contexts at pauses and word boundaries can be blocked using the WB command. For example, if `WB sp` was executed, the effect of a subsequent TC command on the sequence `sil b ah t sp iy n ...` would be to give the sequence `sil+b sil-b+ah b-ah+t ah-t sp iy+n iy-n+ ...`, where `sp` represents a short pause. Conversely, the NB command can be used to ignore a label as far as context is concerned. For example, if `NB sp` was executed, the effect of a subsequent TC command on the sequence `sil b ah t sp iy n ...` would be to give the sequence `sil+b sil-b+ah b-ah+t ah-t+iy sp t-iy+n iy-n+ ...`.

When processing HTK format label files with multiple levels, only the level 1 (i.e. left-most) labels are affected. To process a higher level, the *Move Level* command ML should be used. For example, in the script

```
ML 2
RE one 1
RE two 2
...
```

the Replace commands are applied to level 2 which is the first level above the basic level. The command `ML 1` returns to the base level. A complete level can be deleted by the *Delete Level* command DL. This command can be given a numeric argument to delete a specific level or with no argument, the current level is deleted. Multiple levels can also be split into single level alternatives by using the *Split Level* command SL.

When processing HTK format files with multiple alternatives, each alternative is processed as though it were a separate file.

Remember also that in addition to the explicit HLED commands, levels and alternatives can be filtered on input by setting the configuration variables `TRANSLEV` and `TRANSALT` (see section 6.1).

Finally, it should be noted that most HTK tools require all HMMs used in a system to be defined in a *HMM List*. HLED can be made to automatically generate such a list as a by-product of editing the label files by using the `-n` option. For example, the following command would apply the script `timit.led` to all files in the directory `t1labs`, write the converted files to the directory `h1labs` and also write out a list of all new labels in the edited files to `t1list`.

```
HLEd -n t1list -l h1labs -G TIMIT timit.led t1labs/*
```

Notice here that the `-G` option is used to inform HLED that the format of the source files is TIMIT. This could also be indicated by setting the configuration variable `SOURCELABEL`.

6.5 Summary

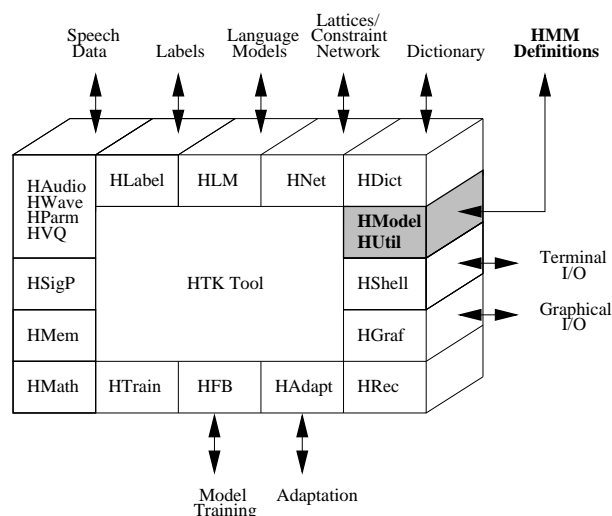
Table 6.1 lists all of the configuration parameters recognised by HLABEL along with a brief description. A missing module name means that it is recognised by more than one module.

Module	Name	Description
HLABEL	LABELSQQUOTE	Specify label quote character
HLABEL	SOURCELABEL	Source label format
HLABEL	SOURCERATE	Sample period for SCRIBE format
HLABEL	STRIPTRIPHONES	Remove triphone contexts on input
HLABEL	TARGETLABEL	Target label format
HLABEL	TRANSALT	Filter alternatives on input
HLABEL	TRANSLEV	Filter levels on input
HLABEL	V1COMPAT	Version 1.5 compatibility mode
	TRACE	trace control (default=0)

Table. 6.1 Configuration Parameters used with Labels

Chapter 7

HMM Definition Files



The principle function of HTK is to manipulate sets of hidden Markov models (HMMs). The definition of a HMM must specify the model topology, the transition parameters and the output distribution parameters. The HMM observation vectors can be divided into multiple independent data streams and each stream can have its own weight. In addition, a HMM can have ancillary information such as duration parameters. HTK supports both continuous mixture densities and discrete distributions. HTK also provides a generalised tying mechanism which allows parameters to be shared within and between models.

In order to encompass this rich variety of HMM types within a single framework, HTK uses a formal language to define HMMs. The interpretation of this language is handled by the library module HMODEL which is responsible for converting between the external and internal representations of HMMs. In addition, it provides all the basic probability function calculations. A second module HUTIL provides various additional facilities for manipulating HMMs once they have been loaded into memory.

The purpose of this chapter is to describe the HMM definition language in some detail. The chapter begins by describing how to write individual HMM definitions. HTK macros are then explained and the mechanisms for defining a complete model set are presented. The various flavours of HMM are then described and the use of binary files discussed. Finally, a formal description of the HTK HMM definition language is given.

As will be seen, the definition of a large HMM system can involve considerable complexity. However, in practice, HMM systems are built incrementally. The usual starting point is a single HMM definition which is then repeatedly cloned and refined using the various HTK tools (in particular, HEREST and HHEd). Hence, in practice, the HTK user rarely has to generate complex HMM definition files directly.

7.1 The HMM Parameters

A HMM consists of a number of states. Each state j has an associated observation probability distribution $b_j(\mathbf{o}_t)$ which determines the probability of generating observation \mathbf{o}_t at time t and each pair of states i and j has an associated transition probability a_{ij} . In HTK the entry state 1 and the exit state N of an N state HMM are non-emitting.

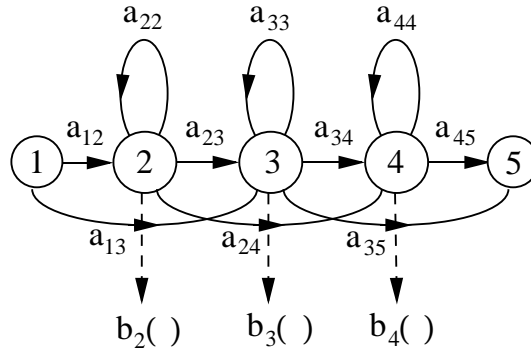


Fig. 7.1 Simple Left-Right HMM

Fig. 7.1 shows a simple left-right HMM with five states in total. Three of these are emitting states and have output probability distributions associated with them. The transition matrix for this model will have 5 rows and 5 columns. Each row will sum to one except for the final row which is always all zero since no transitions are allowed out of the final state.

HTK is principally concerned with continuous density models in which each observation probability distribution is represented by a mixture Gaussian density. In this case, for state j the probability $b_j(\mathbf{o}_t)$ of generating observation \mathbf{o}_t is given by

$$b_j(\mathbf{o}_t) = \prod_{s=1}^S \left[\sum_{m=1}^{M_{js}} c_{j_{sm}} \mathcal{N}(\mathbf{o}_{st}; \boldsymbol{\mu}_{j_{sm}}, \boldsymbol{\Sigma}_{j_{sm}}) \right]^{\gamma_s} \quad (7.1)$$

where M_{js} is the number of mixture components in state j for stream s , $c_{j_{sm}}$ is the weight of the m 'th component and $\mathcal{N}(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is a multivariate Gaussian with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$, that is

$$\mathcal{N}(\mathbf{o}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^n |\boldsymbol{\Sigma}|}} e^{-\frac{1}{2}(\mathbf{o}-\boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1}(\mathbf{o}-\boldsymbol{\mu})} \quad (7.2)$$

where n is the dimensionality of \mathbf{o} . The exponent γ_s is a stream weight and its default value is one. Other values can be used to emphasise particular streams, however, none of the standard HTK tools manipulate it.

HTK also supports discrete probability distributions in which case

$$b_j(\mathbf{o}_t) = \prod_{s=1}^S \{P_{js}[v_s(\mathbf{o}_{st})]\}^{\gamma_s} \quad (7.3)$$

where $v_s(\mathbf{o}_{st})$ is the output of the vector quantiser for stream s given input vector \mathbf{o}_{st} and $P_{js}[v]$ is the probability of state j generating symbol v in stream s .

In addition to the above, any model or state can have an associated vector of duration parameters $\{d_k\}^1$. Also, it is necessary to specify the kind of the observation vectors, and the width of the observation vector in each stream. Thus, the total information needed to define a single HMM is as follows

- type of observation vector
- number and width of each data stream

¹ No current HTK tool can estimate or use these.

- optional model duration parameter vector
- number of states
- for each emitting state and each stream
 - mixture component weights or discrete probabilities
 - if continuous density, then means and covariances
 - optional stream weight vector
 - optional duration parameter vector
- transition matrix

The following sections explain how these are defined.

7.2 Basic HMM Definitions

Some HTK tools require a single HMM to be defined. For example, the isolated-unit re-estimation tool HREST would be invoked as

```
HRest hmmdef s1 s2 s3 ....
```

This would cause the model defined in the file `hmmdef` to be input and its parameters re-estimated using the speech data files `s1`, `s2`, etc.

```
~h "hmm1"
<BeginHMM>
  <VecSize> 4 <MFCC>
  <NumStates> 5
  <State> 2
    <Mean> 4
      0.2 0.1 0.1 0.9
    <Variance> 4
      1.0 1.0 1.0 1.0
  <State> 3
    <Mean> 4
      0.4 0.9 0.2 0.1
    <Variance> 4
      1.0 2.0 2.0 0.5
  <State> 4
    <Mean> 4
      1.2 3.1 0.5 0.9
    <Variance> 4
      5.0 5.0 5.0 5.0
  <TransP> 5
    0.0 0.5 0.5 0.0 0.0
    0.0 0.4 0.4 0.2 0.0
    0.0 0.0 0.6 0.4 0.0
    0.0 0.0 0.0 0.7 0.3
    0.0 0.0 0.0 0.0 0.0
<EndHMM>
```

Fig. 7.2 Definition for Simple L-R HMM

HMM definition files consist of a sequence of symbols representing the elements of a simple language. These symbols are mainly keywords written within angle brackets and integer and floating point numbers. The full HTK definition language is presented more formally later in section 7.10. For now, the main features of the language will be described by some examples.

Fig 7.2 shows a HMM definition corresponding to the simple left-right HMM illustrated in Fig 7.1. It is a continuous density HMM with 5 states in total, 3 of which are emitting. The first symbol in the file `~h` indicates that the following string is the name of a macro of type `h` which means that it is a HMM definition (macros are explained in detail later). Thus, this definition describes a HMM called “hmm1”. Note that HMM names should be composed of alphanumeric characters only and must not consist solely of numbers. The HMM definition itself is bracketed by the symbols `<BeginHMM>` and `<EndHMM>`.

The first line of the definition proper specifies the *global* features of the HMM. In any system consisting of many HMMs, these features will be the same for all of them. In this case, the global definitions indicate that the observation vectors have 4 components (`<VecSize>` 4) and that they denote MFCC coefficients (`<MFCC>`).

The next line specifies the number of states in the HMM. There then follows a definition for each emitting state j , each of which has a single mean vector μ_j introduced by the keyword `<Mean>` and a diagonal variance vector Σ_j introduced by the keyword `<Variance>`. The definition ends with the transition matrix $\{a_{ij}\}$ introduced by the keyword `<TransP>`.

Notice that the dimension of each vector or matrix is specified explicitly before listing the component values. These dimensions must be consistent with the corresponding observation width (in the case of output distribution parameters) or number of states (in the case of transition matrices). Although in this example they could be inferred, HTK requires that they are included explicitly since, as will be described shortly, they can be detached from the HMM definition and stored elsewhere as a macro.

The definition for `hmm1` makes use of many defaults. In particular, there is no definition for the number of input data streams or for the number of mixture components per output distribution. Hence, in both cases, a default of 1 is assumed.

Fig 7.3 shows a HMM definition in which the emitting states are 2 component mixture Gaussians. The number of mixture components in each state j is indicated by the keyword `<NumMixes>` and each mixture component is prefixed by the keyword `<Mixture>` followed by the component index m and component weight c_{jm} . Note that there is no requirement for the number of mixture components to be the same in each distribution.

State definitions and the mixture components within them may be listed in any order. When a HMM definition is loaded, a check is made that all the required components have been defined. In addition, checks are made that the mixture component weights and each row of the transition matrix sum to one. If very rapid loading is required, this consistency checking can be inhibited by setting the Boolean configuration variable `CHKHMMDEFS` to false.

As an alternative to diagonal variance vectors, a Gaussian distribution can have a full rank covariance matrix. An example of this is shown in the definition for `hmm3` shown in Fig 7.4. Since covariance matrices are symmetric, they are stored in upper triangular form i.e. each row of the matrix starts at the diagonal element². Also, covariance matrices are stored in their inverse form i.e. HMM definitions contain Σ^{-1} rather than Σ . To reflect this, the keyword chosen to introduce a full covariance matrix is `<InvCovar>`.

² Covariance matrices are actually stored internally in lower triangular form

```

~h "hmm2"
<BeginHMM>
  <VecSize> 4 <MFCC>
  <NumStates> 4
  <State> 2 <NumMixes> 2
    <Mixture> 1 0.4
      <Mean> 4
        0.3 0.2 0.2 1.0
      <Variance> 4
        1.0 1.0 1.0 1.0
    <Mixture> 2 0.6
      <Mean> 4
        0.1 0.0 0.0 0.8
      <Variance> 4
        1.0 1.0 1.0 1.0
  <State> 3 <NumMixes> 2
    <Mixture> 1 0.7
      <Mean> 4
        0.1 0.2 0.6 1.4
      <Variance> 4
        1.0 1.0 1.0 1.0
    <Mixture> 2 0.3
      <Mean> 4
        2.1 0.0 1.0 1.8
      <Variance> 4
        1.0 1.0 1.0 1.0
  <TransP> 4
    0.0 1.0 0.0 0.0
    0.0 0.5 0.5 0.0
    0.0 0.0 0.6 0.4
    0.0 0.0 0.0 0.0
<EndHMM>

```

Fig. 7.3 Simple Mixture Gaussian HMM

Notice that only the second state has a full covariance Gaussian component. The first state has a mixture of two diagonal variance Gaussian components. Again, this illustrates the flexibility of HMM definition in HTK. If required the structure of every Gaussian can be individually configured.

Another possible way to store covariance information is in the form of the Choleski decomposition L of the inverse covariance matrix i.e. $\Sigma^{-1} = LL'$. Again this is stored externally in upper triangular form so L' is actually stored. It is distinguished from the normal inverse covariance matrix by using the keyword <LLTCovar> in place of <InvCovar>³.

The definition for `hmm3` also illustrates another macro type, that is, `~o`. This macro is used as an alternative way of specifying global options and, in fact, it is the format used by HTK tools when they write out a HMM definition. It is provided so that global options can be specified ahead of any other HMM parameters. As will be seen later, this is useful when using many types of macro.

As noted earlier, the observation vectors used to represent the speech signal can be divided into two or more statistically independent data streams. This corresponds to the splitting-up of the input speech vectors as described in section 5.10. In HMM definitions, the use of multiple data streams must be indicated by specifying the number of streams and the width (i.e dimension) of each stream as a global option. This is done using the keyword <StreamInfo> followed by the number of streams, and then a sequence of numbers indicating the width of each stream. The sum of these stream widths must equal the original vector size as indicated by the <VecSize> keyword.

³ The Choleski storage format is not used by default in HTK Version 2

```

~o <VecSize> 4 <MFCC>
~h "hmm3"
<BeginHMM>
  <NumStates> 4
  <State> 2 <NumMixes> 2
    <Mixture> 1 0.4
      <Mean> 4
        0.3 0.2 0.2 1.0
      <Variance> 4
        1.0 1.0 1.0 1.0
    <Mixture> 2 0.6
      <Mean> 4
        0.1 0.0 0.0 0.8
      <Variance> 4
        1.0 1.0 1.0 1.0
  <State> 3 <NumMixes> 1
    <Mean> 4
      0.10.20.61.4
    <InvCovar> 4
      1.00.10.00.0
      1.00.20.0
      1.00.1
      1.0
    <TransP> 4
      0.0 1.0 0.0 0.0
      0.0 0.5 0.5 0.0
      0.0 0.0 0.6 0.4
      0.0 0.0 0.0 0.0
<EndHMM>

```

Fig. 7.4 HMM with Full Covariance

An example of a HMM definition for multiple data streams is `hmm4` shown in Fig 7.5. This HMM is intended to model 2 distinct streams, the first has 3 components and the second has 1. This is indicated by the global option `<StreamInfo> 2 3 1`. The definition of each state output distribution now includes means and variances for each individual stream.

Thus, in Fig 7.5, each state is subdivided into 2 streams using the `<Stream>` keyword followed by the stream number. Note also, that each individual stream can be weighted. In state 2 of `hmm4`, the vector following the `<SWeights>` keyword indicates that stream 1 has a weight of 0.9 whereas stream 2 has a weight of 1.1. There is no stream weight vector in state 3 and hence the default weight of 1.0 will be assigned to each stream.

No HTK tools are supplied for estimating optimal stream weight values. Hence, they must either be set manually or derived from some outside source. However, once set, they are used in the calculation of output probabilities as specified in equations 7.1 and 7.3, and hence they will affect the operation of both the training and recognition tools.

7.3 Macro Definitions

So far, basic model definitions have been described in which all of the information required to define a HMM has been given directly between the `<BeginHMM>` and `<EndHMM>` keywords. As an alternative, HTK allows the internal parts of a definition to be written as separate units, possibly

```

~o <VecSize> 4 <MFCC>
  <StreamInfo> 2 3 1
~h "hmm4"
<BeginHMM>
  <NumStates> 4
  <State> 2
    <SWeights> 2 0.9 1.1
    <Stream> 1
      <Mean> 3
        0.2 0.1 0.1
      <Variance> 3
        1.0 1.0 1.0
    <Stream> 2
      <Mean> 1 0.0
      <Variance> 1 4.0
  <State> 3
    <Stream> 1
      <Mean> 3
        0.3 0.2 0.0
      <Variance> 3
        1.0 1.0 1.0
    <Stream> 2
      <Mean> 1 0.5
      <Variance> 1 3.0
  <TransP> 4
    0.0 1.0 0.0 0.0
    0.0 0.6 0.4 0.0
    0.0 0.0 0.4 0.6
    0.0 0.0 0.0 0.0
<EndHMM>

```

Fig. 7.5 HMM with 2 Data Streams

in several different files, and then referenced by name wherever they are needed. Such definitions are called *macros*.

```

~o <VecSize> 4 <MFCC>

~v "var"
  <Variance> 4
    1.0 1.0 1.0 1.0

```

Fig. 7.6 Simple Macro Definitions

HMM ($\sim h$) and global option macros ($\sim o$) have already been described. In fact, these are both rather special cases since neither is ever referenced explicitly by another definition. Indeed, the option macro is unusual in that since it is global and must be unique, it has no name. As an illustration of the use of macros, it may be observed that the variance vectors in the HMM definition `hmm2` given in Fig 7.3 are all identical. If this was intentional, then the variance vector could be defined as a macro as illustrated in Fig 7.6.

A macro definition consists of a macro type indicator followed by a user-defined macro name. In this case, the indicator is $\sim v$ and the name is `var`. Notice that a global options macro is included before the definition for `var`. HTK must know these before it can process any other definitions thus the first macro file specified on the command line of any HTK tool must have the global options macro. Global options macro need not be repeated at the head of every definition file, but it does no harm to do so.

```

~h "hmm5"
<BeginHMM>
  <NumStates> 4
  <State> 2 <NumMixes> 2
    <Mixture> 1 0.4
      <Mean> 4
        0.3 0.2 0.2 1.0
      ~v "var"
    <Mixture> 2 0.6
      <Mean> 4
        0.1 0.0 0.0 0.8
      ~v "var"
  <State> 3 <NumMixes> 2
    <Mixture> 1 0.7
      <Mean> 4
        0.1 0.2 0.6 1.4
      ~v "var"
    <Mixture> 2 0.3
      <Mean> 4
        2.1 0.0 1.0 1.8
      ~v "var"
  <TransP> 4
    0.0 1.0 0.0 0.0
    0.0 0.5 0.5 0.0
    0.0 0.0 0.6 0.4
    0.0 0.0 0.0 0.0
<EndHMM>

```

Fig. 7.7 A Definition Using Macros

Once defined, a macro is used simply by writing the type indicator and name exactly as written in the definition. Thus, for example, Fig 7.7 defines a HMM called `hmm5` which uses the variance macro `var` but is otherwise identical to the earlier HMM definition `hmm2`.

The definition for `hmm5` can be understood by substituting the textual body of the `var` macro everywhere that it is referenced. Textually this would make the definition for `hmm5` identical to that for `hmm2`, and indeed, if input to a recogniser, their effects would be similar. However, as will become clear in later chapters, the HMM definitions `hmm2` and `hmm5` differ in two ways. Firstly, if any attempt was made to re-estimate the parameters of `hmm2`, the values of the variance vectors would almost certainly diverge. However, the variance vectors of `hmm5` are tied together and are guaranteed to remain identical, even after re-estimation. Thus, in general, the use of a macro enforces a *tying* which results in the corresponding parameters being shared amongst all the HMM structures which reference that macro. Secondly, when used in a recognition tool, the computation required to decode using HMMs with tied parameters will often be reduced. This is particularly true when higher level parts of a HMM definition are tied such as whole states.

There are many different macro types. Some have special meanings but the following correspond to the various distinct points in the hierarchy of HMM parameters which can be tied.

- ~s shared state distribution
- ~m shared Gaussian mixture component
- ~u shared mean vector
- ~v shared diagonal variance vector
- ~i shared inverse full covariance matrix
- ~c shared choleski L' matrix
- ~x shared arbitrary transform matrix⁴
- ~t shared transition matrix
- ~d shared duration parameters
- ~w shared stream weight vector

Fig 7.8 illustrates these potential tie points graphically for the case of continuous density HMMs. In this figure, each solid black circle represents a potential tie point, and the associated macro type is indicated alongside it.

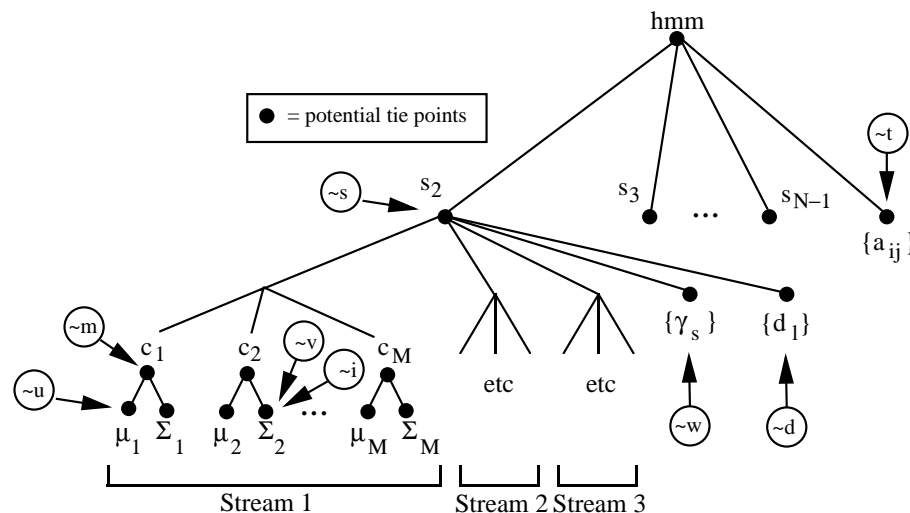


Fig. 7.8 HMM Hierarchy and Potential Tie Points

The tie points for discrete HMMs are identical except that the macro types `~m`, `~v`, `~c`, `~i` and `~u` are not relevant and are therefore excluded.

The macros with special meanings are as follows

- ~l logical HMM
- ~h physical HMM
- ~o global option values
- ~p tied mixture
- ~r regression class tree

The distinction between logical and physical HMMs will be explained in the next section and option macros have already been described. The `~p` macro is used by the HMM editor HHED for building tied mixture systems (see section 7.5). The `~l` or `~p` macros are special in the sense that they are created implicitly in order to represent specific kinds of parameter sharing and they never occur explicitly in HMM definitions.

7.4 HMM Sets

The previous sections have described how a single HMM definition can be specified. However, many HTK tools require complete model sets to be specified rather than just a single model. When this is the case, the individual HMMs which belong to the set are listed in a file rather than being enumerated explicitly on the command line. Thus, for example, a typical invocation of the tool HEREST might be as follows

```
HERest ... -H mf1 -H mf2 ... hlist
```

where each `-H` option names a macro file and `hlist` contains a list of HMM names, one per line. For example, it might contain

```
ha
hb
hc
```

In a case such as this, the macro files would normally contain definitions for the models `ha`, `hb` and `hc`, along with any lower level macro definitions that they might require.

As an illustration, Fig 7.9 and Fig 7.10 give examples of what the macro files `mf1` and `mf2` might contain. The first file contains definitions for three states and a transition matrix. The second file contains definitions for the three HMMs. In this example, each HMM shares the three states and the common transition matrix. A HMM set such as this is called a *tied-state* system.

The order in which macro files are listed on the command line and the order of definition within each file must ensure that all macro definitions are defined before they are referenced. Thus, macro files are typically organised such that all low level structures come first followed by states and transition matrices, with the actual HMM definitions coming last.

When the HMM list contains the name of a HMM for which no corresponding macro has been defined, then an attempt is made to open a file with the same name. This file is expected to contain a single definition corresponding to the required HMM. Thus, the general mechanism for loading a set of HMMs is as shown in Fig 7.11. In this example, the HMM list `hlist` contains the names of five HMMs of which only three have been predefined via the macro files. Hence, the remaining definitions are found in individual HMM definition files `hd` and `he`.

When a large number of HMMs must be loaded from individual files, it is common to store them in a specific directory. Most HTK tools allow this directory to be specified explicitly using a command line option. For example, in the command

```
HERest -d hdir ... hlist ....
```

the definitions for the HMM listed in `hlist` will be searched for in the subdirectory `hdir`.

After loading each HMM set, HMODEL marks it as belonging to one of the following categories (called the *HSK*ind

- PLAINHS
- SHAREDHS
- TIEDHS
- DISCRETEHS

Any HMM set containing discrete output distributions is assigned to the DISCRETEHS category (see section 7.6). If all mixture components are tied, then it is assigned to the TIEDHS category (see section 7.5). If it contains any shared states (`~s` macros) or Gaussians (`~m` macros) then it is SHAREDHS. Otherwise, it is PLAINHS. The category assigned to a HMM set determines which of several possible optimisations the various HTK tools can apply to it. As a check, the required kind

```

~o <VecSize> 4 <MFCC>
~s "stateA"
  <Mean> 4
    0.2 0.1 0.1 0.9
  <Variance> 4
    1.0 1.0 1.0 1.0
~s "stateB"
  <Mean> 4
    0.4 0.9 0.2 0.1
  <Variance> 4
    1.0 2.0 2.0 0.5
~s "stateC"
  <Mean> 4
    1.2 3.1 0.5 0.9
  <Variance> 4
    5.0 5.0 5.0 5.0
~t "tran"
  <TransP> 5
    0.0 0.5 0.5 0.0 0.0
    0.0 0.4 0.4 0.2 0.0
    0.0 0.0 0.6 0.4 0.0
    0.0 0.0 0.0 0.7 0.3
    0.0 0.0 0.0 0.0 0.0

```

Fig. 7.9 File mf1: shared state and transition matrix macros

of a HMM set can also be set via the configuration variable `HMMSETKIND`. For debugging purposes, this can also be used to re-categorise a `SHAREDHS` system as `PLAINHS`.

As shown in Figure 7.8, complete HMM definitions can be tied as well as their individual parameters. However, tying at the HMM level is defined in a different way. HMM lists have so far been described as simply a list of model names. In fact, every HMM has two names: a *logical* name and a *physical name*. The logical name reflects the rôle of the model and the physical name is used to identify the definition on disk. By default, the logical and physical names are identical. HMM tying is implemented by letting several logically distinct HMMs share the same physical definition. This is done by giving an explicit physical name immediately after the logical name in a HMM list.

```

~h "ha"
<BeginHMM>
  <NumStates> 5
  <State> 2
    ~s "stateA"
  <State> 3
    ~s "stateB"
  <State> 4
    ~s "stateB"
  ~t "tran"
<EndHMM>

~h "hb"
<BeginHMM>
  <NumStates> 5
  <State> 2
    ~s "stateB"
  <State> 3
    ~s "stateA"
  <State> 4
    ~s "stateC"
  ~t "tran"
<EndHMM>

~h "hc"
<BeginHMM>
  <NumStates> 5
  <State> 2
    ~s "stateC"
  <State> 3
    ~s "stateC"
  <State> 4
    ~s "stateB"
  ~t "tran"
<EndHMM>

```

Fig. 7.10 Simple Tied-State System

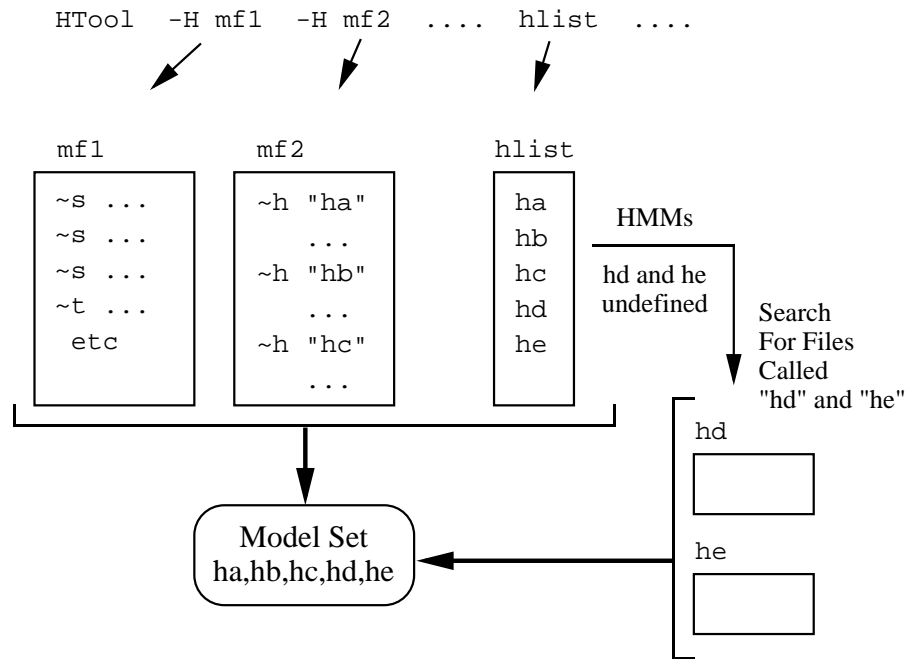


Fig. 7.11 Defining a Model Set

For example, in the HMM list shown in Fig 7.12, the logical HMMs `two`, `too` and `to` are tied and share the same physical HMM definition `tuw`. The HMMs `one` and `won` are also tied but in this case `won` shares `one`'s definition. There is, however, no subtle distinction here. The two different cases are given just to emphasise that the names used for the logical and physical HMMs can be the same or different, as is convenient. Finally, in this example, the models `three` and `four` are untied.

```
two    tuw
too    tuw
to     tuw
one
won    one
three
four
```

Fig. 7.12 HMM List with Tying

This mechanism is implemented internally by creating a `~l` macro definition for every HMM in the HMM list. If an explicit physical HMM is also given in the list, then the logical HMM is linked to that macro, otherwise a `~h` macro is created with the same name as the `~l` macro. Notice that this is one case where the “define before use” rule is relaxed. If an undefined `~h` is encountered then a dummy place-holder is created for it and, as explained above, HMODEL subsequently tries to find a HMM definition file of the same name.

Finally it should be noted that in earlier versions of HTK, there were no HMM macros. However, HMM definitions could be listed in a single *master macro file* or MMF. Each HMM definition began with its name written as a quoted string and ended with a period written on its own (just like master label files), and the first line of an MMF contained the string `#!MMF!#`. In HTK 3.1, the use of MMFs has been subsumed within the general macro definition facility using the `~h` type. However, for compatibility, the older MMF style of file can still be read by all HTK tools.

7.5 Tied-Mixture Systems

A Tied-Mixture System is one in which all Gaussian components are stored in a pool and all state output distributions share this pool. Fig 7.13 illustrates this for the case of single data stream.

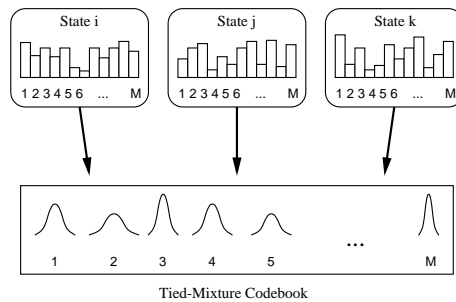


Fig. 7.13 Tied Mixture System

Each state output distribution is defined by M mixture component weights and since all states share the same components, all of the state-specific discrimination is encapsulated within these weights. The set of Gaussian components selected for the pool should be representative of the acoustic space covered by the feature vectors. To keep M manageable, multiple data streams are typically used with tied-mixture systems. For example, static parameters may be in one stream and delta parameters in another (see section 5.10). Each stream then has a separate pool of Gaussians which are often referred to as *codebooks*.

More formally, for S independent data streams, the output distribution for state j is defined as

$$b_j(\mathbf{o}_t) = \prod_{s=1}^S \left[\sum_{m=1}^{M_s} c_{j sm} \mathcal{N}(\mathbf{o}_{st}; \boldsymbol{\mu}_{sm}, \boldsymbol{\Sigma}_{sm}) \right]^{\gamma_s} \quad (7.4)$$

where the notation is identical to that used in equation 7.1. Note however that this equation differs from equation 7.1 in that the Gaussian component parameters and the number of mixture components per stream are state independent.

Tied-mixture systems lack the modelling accuracy of fully continuous density systems. However, they can often be implemented more efficiently since the total number of Gaussians which must be evaluated at each input frame is independent of the number of active HMM states and is typically much smaller.

A tied-mixture HMM system in HTK is defined by representing the pool of shared Gaussians as $\sim m$ macros with names “xxx1”, “xxx2”, ..., “xxxM” where “xxx” is an arbitrary name. Each HMM state definition is then specified by giving the name “xxx” followed by a list of the mixture weights. Multiple streams are identified using the `<Stream>` keyword as described previously.

As an example, Fig 7.14 shows a set of macro definitions which specify a 5 Gaussian component tied-mixture pool.

Fig 7.16 then shows a typical tied-mixture HMM definition which uses this pool. As can be seen, the mixture component weights are represented an array of real numbers as in the continuous density case.

The number of components in each tied-mixture codebook is typically of the order of 2 or 3 hundred. Hence, the list of mixture weights in each state is often long with many values being repeated, particularly floor values. To allow more efficient coding, successive identical values can be represented as a single value plus a repeat count in the form of an asterisk followed by an integer multiplier. For example, Fig 7.15 shows the same HMM definition as above but using repeat counts. When HTK writes out a tied-mixture definition, it uses repeat counts wherever possible.

7.6 Discrete Probability HMMs

Discrete probability HMMs model observation sequences which consist of symbols drawn from a discrete and finite set of size M . As in the case of tied-mixture systems described above, this set is

```

~o <VecSize> 2 <MFCC>
~m "mix1"
  <Mean> 2 0.0 0.1
  <Variance> 2 1.0 1.0
~m "mix2"
  <Mean> 2 0.2 0.3
  <Variance> 2 2.0 1.0
~m "mix3"
  <Mean> 2 0.0 0.1
  <Variance> 2 1.0 2.0
~m "mix4"
  <Mean> 2 0.4 0.1
  <Variance> 2 1.0 1.5
~m "mix5"
  <Mean> 2 0.9 0.7
  <Variance> 2 1.5 1.0

```

Fig. 7.14 Tied-Mixture Codebook

```

~h "htm"
<BeginHMM>
  <NumStates> 4
  <State> 2 <NumMixes> 5
    <TMix> mix 0.2 0.1 0.3*2 0.1
  <State> 3 <NumMixes> 5
    <TMix> mix 0.4 0.3 0.1*3
  <TransP> 4
  ...
<EndHMM>

```

Fig. 7.15 HMM using Repeat Counts

often referred to as a *codebook*.

The form of the output distributions in a discrete HMM was given in equation 7.3. It consists of a table giving the probability of each possible observation symbol. Each symbol is identified by an index in the range 1 to M and hence the probability of any symbol can be determined by a simple table look-up operation.

For speech applications, the observation symbols are generated by a vector quantiser which typically associates a prototype speech vector with each codebook symbol. Each incoming speech vector is then represented by the symbol whose associated prototype is closest. The prototypes themselves are chosen to cover the acoustic space and they are usually calculated by clustering a representative sample of speech vectors.

In HTK, discrete HMMs are specified using a very similar notation to that used for tied-mixture HMMs. A discrete HMM can have multiple data streams but the width of each stream must be 1. The output probabilities are stored as logs in a scaled integer format such that if $d_{js}[v]$ is the stored discrete probability for symbol v in stream s of state j , the true probability is given by

$$P_{js}[v] = \exp(-d_{js}[v]/2371.8) \quad (7.5)$$

Storage in the form of scaled logs allows discrete probability HMMs to be implemented very efficiently since HTK tools mostly use log arithmetic and direct storage in log form avoids the need for a run-time conversion. The range determined by the constant 2371.8 was selected to enable probabilities from 1.0 down to 0.000001 to be stored.

As an example, Fig 7.17 shows the definition of a discrete HMM called dhmm1. As can be seen, this has two streams. The codebook for stream 1 is size 10 and for stream 2, it is size 2. For

```

~h "htm"
<BeginHMM>
  <NumStates> 4
  <State> 2 <NumMixes> 5
    <TMix> mix 0.2 0.1 0.3 0.3 0.1
  <State> 3 <NumMixes> 5
    <TMix> mix 0.4 0.3 0.1 0.1 0.1
  <TransP> 4
    0.0 1.0 0.0 0.0
    0.0 0.5 0.5 0.0
    0.0 0.0 0.6 0.4
    0.0 0.0 0.0 0.0
<EndHMM>

```

Fig. 7.16 Tied-Mixture HMM

consistency with the representation used for continuous density HMMs, these sizes are encoded in the `<NumMixes>` specifier.

7.7 Tee Models

Normally, the transition probability from the non-emitting entry state to the non-emitting exit state of a HMM will be zero to ensure that the HMM aligns with at least one observation vector. Models which have a non-zero entry to exit transition probability are referred to as *tee-models*.

Tee-models are useful for modelling optional transient effects such as short pauses and noise bursts, particularly between words.

Although most HTK tools support tee-models, they are incompatible with those that work with isolated models such as HINIT and HREST. When a tee-model is loaded into one of these tools, its entry to exit transition probability is reset to zero and the first row of its transition matrix is renormalised.

7.8 Regression Class Trees for Adaptation

In order to perform adaptation HTK generally requires the use of a binary regression tree. Its use in the adaptation process is explained in further detail in chapter 9. After construction (see section 10.7) the terminal nodes of the binary regression tree contain mixture component groupings or clusters. These clusters are referred to as regression base classes. Each mixture component in an HMM set belongs to a unique regression base class. The binary regression tree is stored as part of the HMM set, since its structure is necessary for the dynamic adaptation procedure described in section 9.1.2. Also each mixture component has a regression base class identifier (the terminal node indices) stored with it.

An example is shown in figure 7.18 and corresponds with the tree shown in figure 9.1. The example shows the use of the keyword `<HMMSetId>` used to store an identifier for this HMM set. This is important because the regression tree is built based on this HMM set and is hence specific to it. Many sets of transforms may be built that can be applied to this HMM set, but only one HMM set can be transformed by a transform set that utilises the regression tree. The regression tree is described by non-terminal nodes `<Nodes>` and terminal nodes `<TNodes>`. Each node contains its index followed by either the indices of its children (if it is a non-terminal) or the number of mixture components clustered at a terminal. Each mixture component as defined by the keyword `<Mixture>` has an `<RClass>` keyword followed by the regression base class index. When an HMM definition is loaded, a check is made to see that all the regression classes have been defined and that the total number of mixture components loaded for each regression class matches the number of mixture components defined in the regression tree.

The regression tree together with the mixture regression base class numbers can be constructed automatically with the use of the tool HHED (see section 10.7).

```

~o <DISCRETE> <StreamInfo> 2 1 1
~h "dhmm1"
<BeginHMM>
  <NumStates> 4
  <State> 2
    <NumMixes> 10 2
    <SWeights> 2 0.9 1.1
    <Stream> 1
      <DProb> 3288*4 32767*6
    <Stream> 2
      <DProb> 1644*2
  <State> 3
    <NumMixes> 10 2
    <SWeights> 2 0.9 1.1
    <Stream> 1
      <DProb> 5461*10
    <Stream> 2
      <DProb> 1644*2
  <TransP> 4
    0.0 1.0 0.0 0.0
    0.0 0.5 0.5 0.0
    0.0 0.0 0.6 0.4
    0.0 0.0 0.0 0.0
<EndHMM>

```

Fig. 7.17 Discrete Probability HMM

7.9 Binary Storage Format

Throughout this chapter, a text-based representation has been used for the external storage of HMM definitions. For experimental work, text-based storage allows simple and direct access to HMM parameters and this can be invaluable. However, when using very large HMM sets, storage in text form is less practical since it is inefficient in its use of memory and the time taken to load can be excessive due to the large number of character to float conversions needed.

To solve these problems, HTK also provides a binary storage format. In binary mode, keywords are written as a single colon followed by an 8 bit code representing the actual keyword. Any subsequent numerical information following the keyword is then in binary. Integers are written as 16-bit shorts and all floating-point numbers are written as 32-bit single precision floats. The repeat factor used in the run-length encoding scheme for tied-mixture and discrete HMMs is written as a single byte. Its presence immediately after a 16-bit discrete log probability is indicated by setting the top bit to 1 (this is the reason why the range of discrete log probabilities is limited to 0 to 32767 i.e. only 15 bits are used for the actual value). For tied-mixtures, the repeat count is signalled by subtracting 2.0 from the weight.

Binary storage format and text storage format can be mixed within and between input files. Each time a keyword is encountered, its coding is used to determine whether the subsequent numerical information should be input in text or binary form. This means, for example, that binary files can be manually patched by replacing a binary-format definition by a text format definition⁵.

HTK tools provide a standard command line option (-B) to indicate that HMM definitions should be output in binary format. Alternatively, the Boolean configuration variable `SAVEBINARY` can be set to true to force binary format output.

⁵The fact that this is possible does not mean that it is recommended practice!


```

~o <HMMSetId> ecr1_us_mono
   <VecSize> 4 <MFCC>
~r "ecr1_us_mono_tree_4"
   <RegTree> 4
   <Node> 1 2 3
   <Node> 2 4 5
   <Node> 3 6 7
   <TNode> 4 30
   <TNode> 5 25
   <TNode> 6 40
   <TNode> 7 39
~s "stateA" <NumMixes> 3
   <Mixture> 1 0.34
   <RClass> 4
   ~u "mean51"
   ~v "var65"
   <Mixture> 2 0.52
   <RClass> 7
   ~u "mean32"
   ~v "var65"
   <Mixture> 3 0.14
   <RClass> 5
   ~u "mean12"
   ~v "var3"

```

Fig. 7.18 MMF with a regression tree and classes

7.10 The HMM Definition Language

To conclude this chapter, this section presents a formal description of the HMM definition language used by HTK. Syntax is described using an extended BNF notation in which alternatives are separated by a vertical bar |, parentheses () denote factoring, brackets [] denote options, and braces {} denote zero or more repetitions.

All keywords are enclosed in angle brackets⁶ and the case of the keyword name is not significant. White space is not significant except within double-quoted strings.

The top level structure of a HMM definition is shown by the following rule.

```

hmmdef = [ ~h macro ]
         <BeginHMM>
         [ globalOpts ]
         <NumStates> short
         state { state }
         [ regTree ]
         transP
         [ duration ]
         <EndHMM>

```

A HMM definition consists of an optional set of global options followed by the <NumStates> keyword whose following argument specifies the number of states in the model inclusive of the non-emitting entry and exit states⁷. The information for each state is then given in turn, followed by the parameters of the transition matrix and the model duration parameters, if any. The name of the HMM is given by the ~h macro. If the HMM is the only definition within a file, the ~h macro name can be omitted and the HMM name is assumed to be the same as the file name.

The global options are common to all HMMs. They can be given separately using a ~o option macro

⁶ This definition covers the textual version only. The syntax for the binary format is identical apart from the way that the lexical items are encoded.

⁷ Integer numbers are specified as either `char` or `short`. This has no effect on text-based definitions but for binary format it indicates the underlying C type used to represent the number.

```
optmacro = ~o globalOpts
```

or they can be included in one or more HMM definitions. Global options may be repeated but no definition can change a previous definition. All global options must be defined before any other macro definition is processed. In practice this means that any HMM system which uses parameter tying must have a `~o` option macro at the head of the first macro file processed.

The full set of global options is given below. Every HMM set must define the vector size (via `<VecSize>`), the stream widths (via `<StreamInfo>`) and the observation parameter kind. However, if only the stream widths are given, then the vector size will be inferred. If only the vector size is given, then a single stream of identical width will be assumed. All other options default to null.

```
globalOpts = option { option }
option = <HmSetId> string |
        <StreamInfo> short { short } |
        <VecSize> short |
        covkind |
        durkind |
        parmkind
```

The `<HmSetId>` option allows the user to give the MMF an identifier. This is used as a sanity check to make sure that a TMF can be safely applied to this MMF. The arguments to the `<StreamInfo>` option are the number of streams (default 1) and then for each stream, the width of that stream. The `<VecSize>` option gives the total number of elements in each input vector. If both `<VecSize>` and `<StreamInfo>` are included then the sum of all the stream widths must equal the input vector size.

The `covkind` defines the kind of the covariance matrix

```
covkind = <DiagC> | <InvDiagC> | <FullC> |
          <LLTC> | <XformC>
```

where `<InvDiagC>` is used internally. `<LLTC>` and `<XformC>` are not used in HTK Version 2.0. Setting the covariance kind as a global option forces all components to have this kind. In particular, it prevents mixing full and diagonal covariances within a HMM set.

The `durkind` denotes the type of duration model used according to the following rules

```
durkind = <nullD> | <poissonD> | <gammaD> | <genD>
```

For anything other than `<nullD>`, a duration vector must be supplied for the model or each state as described below. Note that no current HTK tool can estimate or use such duration vectors.

The parameter kind is any legal parameter kind including qualified forms (see section 5.1)

```
parmkind = <basekind{.D|.A|.E|.N|.Z|.O|.V|.C|.K}>
basekind = <discrete>|<lpc>|<lpcpstra>|<mfcc> | <fbank> |
          <melspec>|<lprefc>|<lpedelcep> | <user>
```

where the syntax rule for `parmkind` is non-standard in that no spaces are allowed between the base kind and any subsequent qualifiers. As noted in chapter 5, `<lpedelcep>` is provided only for compatibility with earlier versions of HTK and its further use should be avoided.

Each state of each HMM must have its own section defining the parameters associated with that state

```
state = <State: Exp > short stateinfo
```

where the short following `<State: Exp >` is the state number. State information can be defined in any order. The syntax is as follows

```
stateinfo = ~s macro |
           [ mixes ] [ weights ] stream { stream } [ duration ]
macro = string
```

A `stateinfo` definition consists of an optional specification of the number of mixtures, an optional set of stream weights, followed by a block of information for each stream, optionally terminated with a duration vector. Alternatively, `~s macro` can be written where `macro` is the name of a previously defined macro.

The optional `mixes` in a `stateinfo` definition specify the number of mixture components (or discrete codebook size) for each stream of that state

```
mixes =      <NumMixes> short {short}
```

where there should be one `short` for each stream. If this specification is omitted, it is assumed that all streams have just one mixture component.

The optional weights in a `stateinfo` definition define a set of exponent weights for each independent data stream. The syntax is

```
weights =    ~w macro | <SWeights> short vector
vector =     float { float }
```

where the `short` gives the number S of weights (which should match the value given in the `<StreamInfo>` option) and the `vector` contains the S stream weights γ_s (see section 7.1).

The definition of each `stream` depends on the kind of HMM set. In the normal case, it consists of a sequence of mixture component definitions optionally preceded by the stream number. If the stream number is omitted then it is assumed to be 1. For tied-mixture and discrete HMM sets, special forms are used.

```
stream =     [ <Stream> short ]
             (mixture { mixture } | tmixpdf | discpdf)
```

The definition of each mixture component consists of a Gaussian pdf optionally preceded by the mixture number and its weight

```
mixture =    [ <Mixture> short float ] mixpdf
```

If the `<Mixture>` part is missing then mixture 1 is assumed and the weight defaults to 1.0.

The `tmixpdf` option is used only for fully tied mixture sets. Since the `mixpdf` parts are all macros in a tied mixture system and since they are identical for every stream and state, it is only necessary to know the mixture weights. The `tmixpdf` syntax allows these to be specified in the following compact form

```
tmixpdf =    <TMix> macro weightList
weightList = repShort { repShort }
repShort =   short [ * char ]
```

where each `short` is a mixture component weight scaled so that a weight of 1.0 is represented by the integer 32767. The optional asterisk followed by a `char` is used to indicate a repeat count. For example, `0*5` is equivalent to 5 zeroes. The Gaussians which make-up the pool of tied-mixtures are defined using `~m` macros called `macro1`, `macro2`, `macro3`, etc.

Discrete probability HMMs are defined in a similar way

```
discpdf =    <DProb> weightList
```

The only difference is that the weights in the `weightList` are scaled log probabilities as defined in section 7.6.

The definition of a Gaussian pdf requires the mean vector to be given and one of the possible forms of covariance

```
mixpdf =     ~m macro | [ rclass ] mean cov [ <GConst> float ]
rclass =     <RClass> short
mean =       ~u macro | <Mean> short vector
cov =        var | inv | xform
var =        ~v macro | <Variance> short vector
inv =        ~i macro |
             (<InvCovar> | <LLTCovar>) short tmatrix
xform =      ~x macro | <Xform> short short matrix
matrix =     float {float}
tmatrix =    matrix
```

In `mean` and `var`, the `short` preceding the vector defines the length of the vector, in `inv` the `short` preceding the `tmatrix` gives the size of this square upper triangular matrix, and in `xform` the two `short`'s preceding the `matrix` give the number of rows and columns. The optional `<GConst>`⁸ gives

⁸specifically, in equation 7.2 the `GCONST` value seen in HMM sets is calculated by multiplying the determinant of the covariance matrix by $(2\pi)^n$

that part of the log probability of a Gaussian that can be precomputed. If it is omitted, then it will be computed during load-in, including it simply saves some time. HTK tools which output HMM definitions always include this field. The optional `<RClass>` stores the regression base class index that this mixture component belongs to, as specified by the regression class tree (which is also stored in the model set). HTK tools which output HMM definitions always include this field, and if there is no regression class tree then the regression identifier is set to zero.

In addition to defining the output distributions, a state can have a duration probability distribution defined for it. However, no current HTK tool can estimate or use these.

```
duration =    ~d macro | <Duration> short vector
```

Alternatively, as shown by the top level syntax for a `hmmdef`, duration parameters can be specified for a whole model.

A binary regression class tree (for the purposes of HMM adaptation as in chapter 9) may also exist for an HMM set. This is defined by

```
regTree =    ~r macro tree
tree =       <RegTree> short nodes
nodes =      (<Node> short short short | <TNode> short int) [ nodes ]
```

In `tree` the short preceding the `nodes` refers to the number of terminal nodes or leaves that the regression tree contains. Each node in `nodes` can either be a non-terminal `<Node>` or a terminal (leaf) `<TNode>`. For a `<Node>` the three following shorts refer to the node's index number and the index numbers of its children. For a `<TNode>`, the short refers to the leaf's index (which correspond to a regression base class index as stored at the component level in `RClass`, see above), while the `int` refers to the number of mixture components in this leaf cluster.

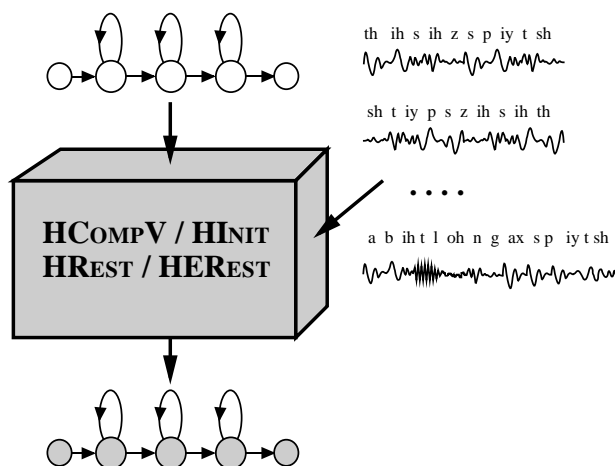
Finally, the transition matrix is defined by

```
transP =     ~t macro | <TransP> short matrix
```

where the short in this case should be equal to the number of states in the model.

Chapter 8

HMM Parameter Estimation



In chapter 7 the various types of HMM were described and the way in which they are represented within HTK was explained. Defining the structure and overall form of a set of HMMs is the first step towards building a recogniser. The second step is to estimate the parameters of the HMMs from examples of the data sequences that they are intended to model. This process of parameter estimation is usually called *training*. HTK supplies four basic tools for parameter estimation: HCOMPV, HINIT, HREST and HEREST. HCOMPV and HINIT are used for initialisation. HCOMPV will set the mean and variance of every Gaussian component in a HMM definition to be equal to the global mean and variance of the speech training data. This is typically used as an initialisation stage for *flat-start* training. Alternatively, a more detailed initialisation is possible using HINIT which will compute the parameters of a new HMM using a Viterbi style of estimation.

HREST and HEREST are used to refine the parameters of existing HMMs using Baum-Welch Re-estimation. Like HINIT, HREST performs *isolated-unit* training whereas HEREST operates on complete model sets and performs *embedded-unit* training. In general, whole word HMMs are built using HINIT and HREST, and continuous speech sub-word based systems are built using HEREST initialised by either HCOMPV or HINIT and HREST.

This chapter describes these training tools and their use for estimating the parameters of plain (i.e. untied) continuous density HMMs. The use of tying and special cases such as tied-mixture HMM sets and discrete probability HMMs are dealt with in later chapters. The first section of this chapter gives an overview of the various training strategies possible with HTK. This is then followed by sections covering initialisation, isolated-unit training, and embedded training. The chapter concludes with a section detailing the various formulae used by the training tools.

8.1 Training Strategies

As indicated in the introduction above, the basic operation of the HTK training tools involves reading in a set of one or more HMM definitions, and then using speech data to estimate the

parameters of these definitions. The speech data files are normally stored in parameterised form such as LPC or MFCC parameters. However, additional parameters such as delta coefficients are normally computed *on-the-fly* whilst loading each file.

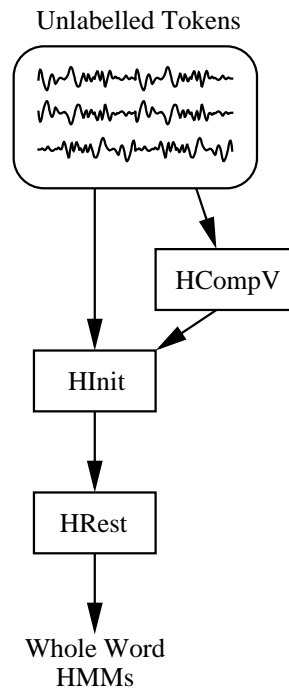


Fig. 8.1 Isolated Word Training

In fact, it is also possible to use waveform data directly by performing the full parameter conversion *on-the-fly*. Which approach is preferred depends on the available computing resources. The advantages of storing the data already encoded are that the data is more compact in parameterised form and pre-encoding avoids wasting compute time converting the data each time that it is read in. However, if the training data is derived from CD-ROMS and they can be accessed automatically on-line, then the extra compute may be worth the saving in magnetic disk storage.

The methods for configuring speech data input to HTK tools were described in detail in chapter 5. All of the various input mechanisms are supported by the HTK training tools except direct audio input.

The precise way in which the training tools are used depends on the type of HMM system to be built and the form of the available training data. Furthermore, HTK tools are designed to interface cleanly to each other, so a large number of configurations are possible. In practice, however, HMM-based speech recognisers are either whole-word or sub-word.

As the name suggests, whole word modelling refers to a technique whereby each individual word in the system vocabulary is modelled by a single HMM. As shown in Fig. 8.1, whole word HMMs are most commonly trained on examples of each word spoken in isolation. If these training examples, which are often called *tokens*, have had leading and trailing silence removed, then they can be input directly into the training tools without the need for any label information. The most common method of building whole word HMMs is to firstly use HINIT to calculate initial parameters for the model and then use HREST to refine the parameters using Baum-Welch re-estimation. Where there is limited training data and recognition in adverse noise environments is needed, so-called *fixed variance* models can offer improved robustness. These are models in which all the variances are set equal to the global speech variance and never subsequently re-estimated. The tool HCOMPV can be used to compute this global variance.

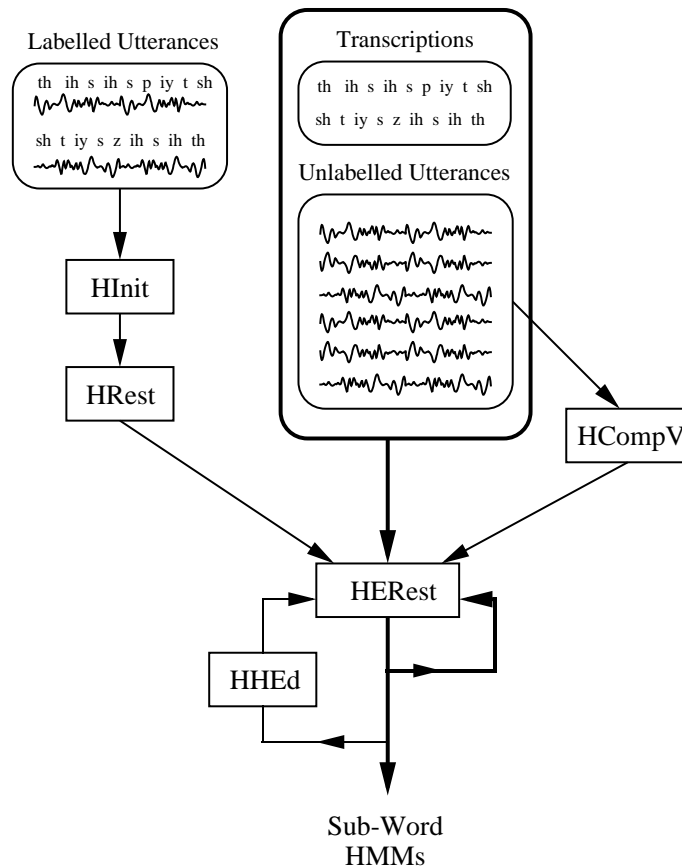


Fig. 8.2 Training Subword HMMs

Although HTK gives full support for building whole-word HMM systems, the bulk of its facilities are focussed on building sub-word systems in which the basic units are the individual sounds of the language called *phones*. One HMM is constructed for each such phone and continuous speech is recognised by joining the phones together to make any required vocabulary using a pronunciation dictionary.

The basic procedures involved in training a set of subword models are shown in Fig. 8.2. The core process involves the embedded training tool HEREST. HEREST uses continuously spoken utterances as its source of training data and simultaneously re-estimates the complete set of subword HMMs. For each input utterance, HEREST needs a transcription i.e. a list of the phones in that utterance. HEREST then joins together all of the subword HMMs corresponding to this phone list to make a single composite HMM. This composite HMM is used to collect the necessary statistics for the re-estimation. When all of the training utterances have been processed, the total set of accumulated statistics are used to re-estimate the parameters of all of the phone HMMs. It is important to emphasise that in the above process, the transcriptions are only needed to identify the sequence of phones in each utterance. No phone boundary information is needed.

The initialisation of a set of phone HMMs prior to embedded re-estimation using HEREST can be achieved in two different ways. As shown on the left of Fig. 8.2, a small set of hand-labelled *bootstrap* training data can be used along with the isolated training tools HINIT and HREST to initialise each phone HMM individually. When used in this way, both HINIT and HREST use the label information to extract all the segments of speech corresponding to the current phone HMM in order to perform isolated word training.

A simpler initialisation procedure uses HCOMPV to assign the global speech mean and variance to every Gaussian distribution in every phone HMM. This so-called *flat start* procedure implies that during the first cycle of embedded re-estimation, each training utterance will be uniformly segmented. The hope then is that enough of the phone models align with actual realisations of that phone so that on the second and subsequent iterations, the models align as intended.

One of the major problems to be faced in building any HMM-based system is that the amount of training data for each model will be variable and is rarely sufficient. To overcome this, HTK allows a variety of sharing mechanisms to be implemented whereby HMM parameters are tied together so that the training data is pooled and more robust estimates result. These tyings, along with a variety of other manipulations, are performed using the HTK HMM editor HHED. The use of HHED is described in a later chapter. Here it is sufficient to note that a phone-based HMM set typically goes through several refinement cycles of editing using HHED followed by parameter re-estimation using HEREST before the final model set is obtained.

Having described in outline the main training strategies, each of the above procedures will be described in more detail.

8.2 Initialisation using HInit

In order to create a HMM definition, it is first necessary to produce a prototype definition. As explained in Chapter 7, HMM definitions can be stored as a text file and hence the simplest way of creating a prototype is by using a text editor to manually produce a definition of the form shown in Fig 7.2, Fig 7.3 etc. The function of a prototype definition is to describe the form and topology of the HMM, the actual numbers used in the definition are not important. Hence, the vector size and parameter kind should be specified and the number of states chosen. The allowable transitions between states should be indicated by putting non-zero values in the corresponding elements of the transition matrix and zeros elsewhere. The rows of the transition matrix must sum to one except for the final row which should be all zero. Each state definition should show the required number of streams and mixture components in each stream. All mean values can be zero but diagonal variances should be positive and covariance matrices should have positive diagonal elements. All state definitions can be identical.

Having set up an appropriate prototype, a HMM can be initialised using the HTKtool HINIT. The basic principle of HINIT depends on the concept of a HMM as a generator of speech vectors. Every training example can be viewed as the output of the HMM whose parameters are to be estimated. Thus, if the state that generated each vector in the training data was known, then the unknown means and variances could be estimated by averaging all the vectors associated with each state. Similarly, the transition matrix could be estimated by simply counting the number of time slots that each state was occupied. This process is described more formally in section 8.7 below.

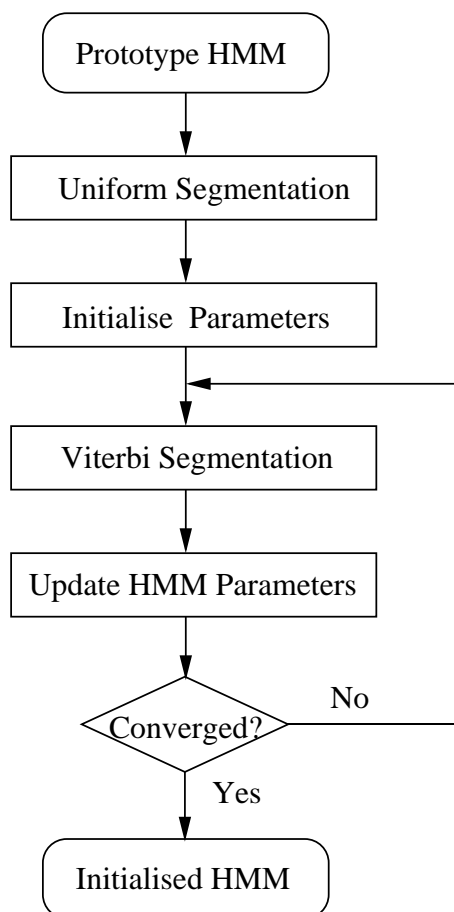


Fig. 8.3 HInit Operation

The above idea can be implemented by an iterative scheme as shown in Fig 8.3. Firstly, the Viterbi algorithm is used to find the most likely state sequence corresponding to each training example, then the HMM parameters are estimated. As a side-effect of finding the Viterbi state alignment, the log likelihood of the training data can be computed. Hence, the whole estimation process can be repeated until no further increase in likelihood is obtained.

This process requires some initial HMM parameters to get started. To circumvent this problem, HINIT starts by uniformly segmenting the data and associating each successive segment with successive states. Of course, this only makes sense if the HMM is left-right. If the HMM is ergodic, then the uniform segmentation can be disabled and some other approach taken. For example, HCOMPV can be used as described below.

If any HMM state has multiple mixture components, then the training vectors are associated with the mixture component with the highest likelihood. The number of vectors associated with each component within a state can then be used to estimate the mixture weights. In the uniform segmentation stage, a K-means clustering algorithm is used to cluster the vectors within each state.

Turning now to the practical use of HINIT, whole word models can be initialised by typing a command of the form

```
HInit hmm data1 data2 data3
```

where `hmm` is the name of the file holding the prototype HMM and `data1`, `data2`, etc. are the names of the speech files holding the training examples, each file holding a single example with no leading or trailing silence. The HMM definition can be distributed across a number of macro files loaded using the standard `-H` option. For example, in

```
HInit -H mac1 -H mac2 hmm data1 data2 data3 ...
```

then the macro files `mac1` and `mac2` would be loaded first. If these contained a definition for `hmm`, then no further HMM definition input would be attempted. If however, they did not contain a definition for `hmm`, then `HINIT` would attempt to open a file called `hmm` and would expect to find a definition for `hmm` within it. `HINIT` can in principle load a large set of HMM definitions, but it will only update the parameters of the single named HMM. On completion, `HINIT` will write out new versions of all HMM definitions loaded on start-up. The default behaviour is to write these to the current directory which has the usually undesirable effect of overwriting the prototype definition. This can be prevented by specifying a new directory for the output definitions using the `-M` option. Thus, typical usage of `HINIT` takes the form

```
HInit -H globals -M dir1 proto data1 data2 data3 ...
mv dir1/proto dir1/wordX
```

Here `globals` is assumed to hold a global options macro (and possibly others). The actual HMM definition is loaded from the file `proto` in the current directory and the newly initialised definition along with a copy of `globals` will be written to `dir1`. Since the newly created HMM will still be called `proto`, it is renamed as appropriate.

For most real tasks, the number of data files required will exceed the command line argument limit and a script file is used instead. Hence, if the names of the data files are stored in the file `trainlist` then typing

```
HInit -S trainlist -H globals -M dir1 proto
```

would have the same effect as previously.

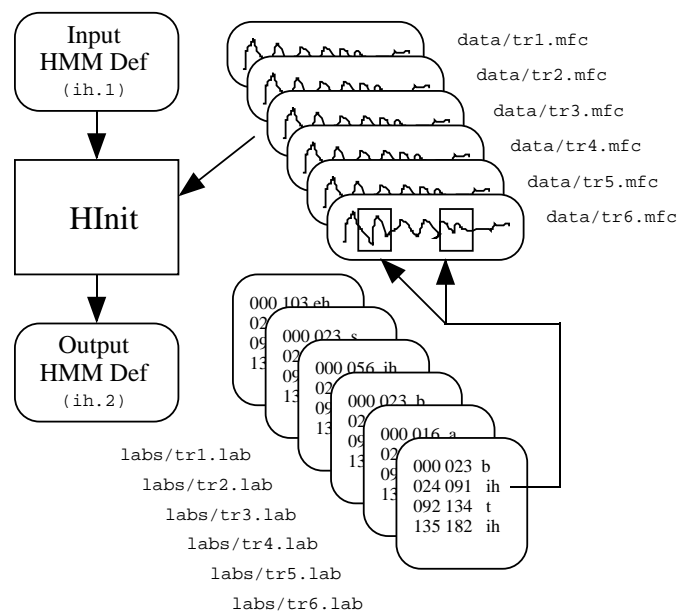


Fig. 8.4 File Processing in HInit

When building sub-word models, `HINIT` can be used in the same manner as above to initialise each individual sub-word HMM. However, in this case, the training data is typically continuous speech with associated label files identifying the speech segments corresponding to each sub-word. To illustrate this, the following command could be used to initialise a sub-word HMM for the phone `ih`

```
HInit -S trainlist -H globals -M dir1 -l ih -L labs proto
mv dir1/proto dir1/ih
```

where the option `-l` defines the name of the sub-word model, and the file `trainlist` is assumed to hold

```

data/tr1.mfc
data/tr2.mfc
data/tr3.mfc
data/tr4.mfc
data/tr5.mfc
data/tr6.mfc

```

In this case, HINIT will first try to find label files corresponding to each data file. In the example here, the standard `-L` option indicates that they are stored in a directory called `labs`. As an alternative, they could be stored in a Master Label File (MLF) and loaded via the standard option `-I`. Once the label files have been loaded, each data file is scanned and all segments corresponding to the label `ih` are loaded. Figure 8.4 illustrates this process.

All HTK tools support the `-T` trace option and although the details of tracing varies from tool to tool, setting the least significant bit (e.g. by `-T 1`), causes all tools to output top level progress information. In the case of HINIT, this information includes the log likelihood at each iteration and hence it is very useful for monitoring convergence. For example, enabling top level tracing in the previous example might result in the following being output

```

Initialising HMM proto . . .
States      : 2 3 4 (width)
Mixes s1:   1 1 1 ( 26 )
Num Using:  0 0 0
Parm Kind:  MFCC_E_D
Number of owners = 1
SegLab     :  ih
maxIter    :  20
epsilon    :  0.000100
minSeg     :  3
Updating   :  Means Variances MixWeights/DProbs TransProbs
16 Observation Sequences Loaded
Starting Estimation Process
Iteration 1: Average LogP = -898.24976
Iteration 2: Average LogP = -884.05402 Change = 14.19574
Iteration 3: Average LogP = -883.22119 Change = 0.83282
Iteration 4: Average LogP = -882.84381 Change = 0.37738
Iteration 5: Average LogP = -882.76526 Change = 0.07855
Iteration 6: Average LogP = -882.76526 Change = 0.00000
Estimation converged at iteration 7
Output written to directory :dir1:

```

The first part summarises the structure of the HMM, in this case, the data is single stream MFCC coefficients with energy and deltas appended. The HMM has 3 emitting states, each single Gaussian and the stream width is 26. The current option settings are then given followed by the convergence information. In this example, convergence was reached after 6 iterations, however if the `maxIter` limit was reached, then the process would terminate regardless.

HINIT provides a variety of command line options for controlling its detailed behaviour. The types of parameter estimated by HINIT can be controlled using the `-u` option, for example, `-u mtw` would update the means, transition matrices and mixture component weights but would leave the variances untouched. A variance floor can be applied using the `-v` to prevent any variance getting too small. This option applies the same variance floor to all speech vector elements. More precise control can be obtained by specifying a variance macro (i.e. a `v` macro) called `varFloor1` for stream 1, `varFloor2` for stream 2, etc. Each element of these variance vectors then defines a floor for the corresponding HMM variance components.

The full list of options supported by HINIT is described in the Reference Section.

8.3 Flat Starting with HCompV

One limitation of using HINIT for the initialisation of sub-word models is that it requires labelled training data. For cases where this is not readily available, an alternative initialisation strategy is to make all models equal initially and move straight to embedded training using HEREST. The

idea behind this so-called *flat start* training is similar to the uniform segmentation strategy adopted by HINIT since by making all states of all models equal, the first iteration of embedded training will effectively rely on a uniform segmentation of the data.

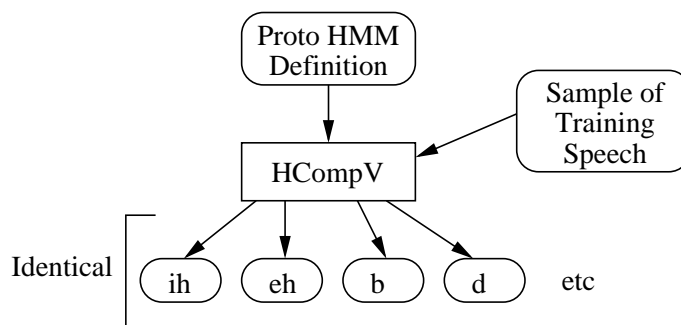


Fig. 8.5 Flat Start Initialisation

Flat start initialisation is provided by the HTK tool HCOMPV whose operation is illustrated by Fig 8.5. The input/output of HMM definition files and training files in HCOMPV works in exactly the same way as described above for HINIT. It reads in a prototype HMM definition and some training data and outputs a new definition in which every mean and covariance is equal to the global speech mean and covariance. Thus, for example, the following command would read a prototype definition called `proto`, read in all speech vectors from `data1`, `data2`, `data3`, etc, compute the global mean and covariance and write out a new version of `proto` in `dir1` with this mean and covariance.

```
HCompV -m -H globals -M dir1 proto data1 data2 data3 ...
```

The default operation of HCOMPV is only to update the covariances of the HMM and leave the means unchanged. The use of the `-m` option above causes the means to be updated too. This apparently curious default behaviour arises because HCOMPV is also used to initialise the variances in so-called *Fixed-Variance* HMMs. These are HMMs initialised in the normal way except that all covariances are set equal to the global speech covariance and never subsequently changed.

Finally, it should be noted that HCOMPV can also be used to generate variance floor macros by using the `-f` option.

8.4 Isolated Unit Re-Estimation using HRest

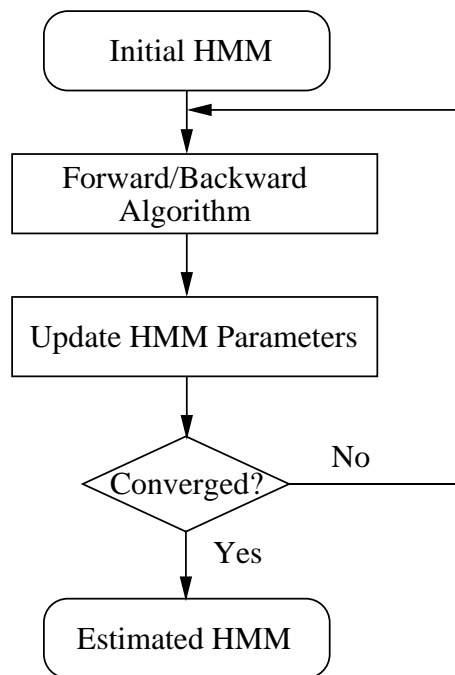


Fig. 8.6 HRest Operation

HREST is the final tool in the set designed to manipulate isolated unit HMMs. Its operation is very similar to HINIT except that, as shown in Fig 8.6, it expects the input HMM definition to have been initialised and it uses Baum-Welch re-estimation in place of Viterbi training. This involves finding the probability of being in each state at each time frame using the *Forward-Backward* algorithm. This probability is then used to form weighted averages for the HMM parameters. Thus, whereas Viterbi training makes a hard decision as to which state each training vector was “generated” by, Baum-Welch takes a soft decision. This can be helpful when estimating phone-based HMMs since there are no hard boundaries between phones in real speech and using a soft decision may give better results. The mathematical details of the Baum-Welch re-estimation process are given below in section 8.7.

HREST is usually applied directly to the models generated by HINIT. Hence for example, the generation of a sub-word model for the phone `ih` begun in section 8.2 would be continued by executing the following command

```
HRest -S trainlist -H dir1/globals -M dir2 -l ih -L labs dir1/ih
```

This will load the HMM definition for `ih` from `dir1`, re-estimate the parameters using the speech segments labelled with `ih` and write the new definition to directory `dir2`.

If HREST is used to build models with a large number of mixture components per state, a strategy must be chosen for dealing with *defunct mixture components*. These are mixture components which have very little associated training data and as a consequence either the variances or the corresponding mixture weight becomes very small. If either of these events happen, the mixture component is effectively deleted and provided that at least one component in that state is left, a warning is issued. If this behaviour is not desired then the variance can be floored as described previously using the `-v` option (or a variance floor macro) and/or the mixture weight can be floored using the `-w` option.

Finally, a problem which can arise when using HREST to initialise sub-word models is that of over-short training segments. By default, HREST ignores all training examples which have fewer frames than the model has emitting states. For example, suppose that a particular phone with 3 emitting states had only a few training examples with more than 2 frames of data. In this case, there would be two solutions. Firstly, the number of emitting states could be reduced. Since

HTK does not require all models to have the same number of states, this is perfectly feasible. Alternatively, some skip transitions could be added and the default reject mechanism disabled by setting the `-t` option. Note here that HINIT has the same reject mechanism and suffers from the same problems. HINIT, however, does not allow the reject mechanism to be suppressed since the uniform segmentation process would otherwise fail.

8.5 Embedded Training using HERest

Whereas isolated unit training is sufficient for building whole word models and initialisation of models using hand-labelled *bootstrap* data, the main HMM training procedures for building subword systems revolve around the concept of *embedded training*. Unlike the processes described so far, embedded training simultaneously updates all of the HMMs in a system using all of the training data. It is performed by HEREST which, unlike HREST, performs just a single iteration.

In outline, HEREST works as follows. On startup, HEREST loads in a complete set of HMM definitions. Every training file must have an associated label file which gives a transcription for that file. Only the sequence of labels is used by HEREST, however, and any boundary location information is ignored. Thus, these transcriptions can be generated automatically from the known orthography of what was said and a pronunciation dictionary.

HEREST processes each training file in turn. After loading it into memory, it uses the associated transcription to construct a composite HMM which spans the whole utterance. This composite HMM is made by concatenating instances of the phone HMMs corresponding to each label in the transcription. The Forward-Backward algorithm is then applied and the sums needed to form the weighted averages accumulated in the normal way. When all of the training files have been processed, the new parameter estimates are formed from the weighted sums and the updated HMM set is output.

The mathematical details of embedded Baum-Welch re-estimation are given below in section 8.7.

In order to use HEREST, it is first necessary to construct a file containing a list of all HMMs in the model set with each model name being written on a separate line. The names of the models in this list must correspond to the labels used in the transcriptions and there must be a corresponding model for every distinct transcription label. HEREST is typically invoked by a command line of the form

```
HERest -S trainlist -I labs -H dir1/hmacs -M dir2 hmmlist
```

where `hmmlist` contains the list of HMMs. On startup, HEREST will load the HMM master macro file (MMF) `hmacs` (there may be several of these). It then searches for a definition for each HMM listed in the `hmmlist`, if any HMM name is not found, it attempts to open a file of the same name in the current directory (or a directory designated by the `-d` option). Usually in large subword systems, however, all of the HMM definitions will be stored in MMFs. Similarly, all of the required transcriptions will be stored in one or more Master Label Files (MLFs), and in the example, they are stored in the single MLF called `labs`.

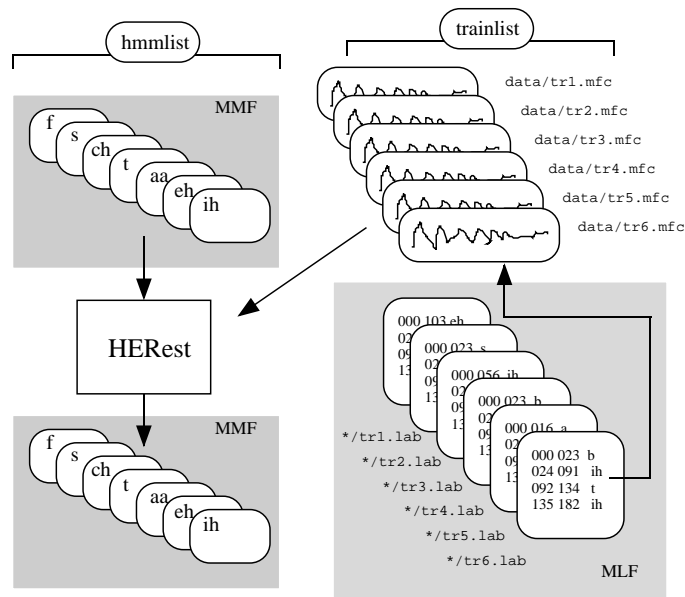


Fig. 8.7 File Processing in HERest

Once all MMFs and MLFs have been loaded, HEREST processes each file in the `trainlist`, and accumulates the required statistics as described above. On completion, an updated MMF is output to the directory `dir2`. If a second iteration is required, then HEREST is reinvoked reading in the MMF from `dir2` and outputting a new one to `dir3`, and so on. This process is illustrated by Fig 8.7.

When performing embedded training, it is good practice to monitor the performance of the models on unseen test data and stop training when no further improvement is obtained. Enabling top level tracing by setting `-T 1` will cause HEREST to output the overall log likelihood per frame of the training data. This measure could be used as a termination condition for repeated application of HEREST. However, repeated re-estimation to convergence may take an impossibly long time. Worse still, it can lead to over-training since the models can become too closely matched to the training data and fail to generalise well on unseen test data. Hence in practice around 2 to 5 cycles of embedded re-estimation are normally sufficient when training phone models.

In order to get accurate acoustic models, a large amount of training data is needed. Several hundred utterances are needed for speaker dependent recognition and several thousand are needed for speaker independent recognition. In the latter case, a single iteration of embedded training might take several hours to compute. There are two mechanisms for speeding up this computation. Firstly, HEREST has a pruning mechanism incorporated into its forward-backward computation. HEREST calculates the backward probabilities $\beta_j(t)$ first and then the forward probabilities $\alpha_j(t)$. The full computation of these probabilities for all values of state j and time t is unnecessary since many of these combinations will be highly improbable. On the forward pass, HEREST restricts the computation of the α values to just those for which the total log likelihood as determined by the product $\alpha_j(t)\beta_j(t)$ is within a fixed distance from the total likelihood $P(\mathbf{O}|M)$. This pruning is always enabled since it is completely safe and causes no loss of modelling accuracy.

Pruning on the backward pass is also possible. However, in this case, the likelihood product $\alpha_j(t)\beta_j(t)$ is unavailable since $\alpha_j(t)$ has yet to be computed, and hence a much broader *beam* must be set to avoid pruning errors. Pruning on the backward path is therefore under user control. It is set using the `-t` option which has two forms. In the simplest case, a fixed pruning beam is set. For example, using `-t 250.0` would set a fixed beam of 250.0. This method is adequate when there is sufficient compute time available to use a generously wide beam. When a narrower beam is used, HEREST will reject any utterance for which the beam proves to be too narrow. This can be avoided by using an incremental threshold. For example, executing

```
HERest -t 120.0 60.0 240.0 -S trainlist -I labs \
-H dir1/hmacs -M dir2 hmmlist
```

would cause HEREST to run normally at a beam width of 120.0. However, if a pruning error occurs, the beam is increased by 60.0 and HEREST reprocesses the offending training utterance. Repeated errors cause the beam width to be increased again and this continues until either the utterance is successfully processed or the upper beam limit is reached, in this case 240.0. Note that errors which occur at very high beam widths are often caused by transcription errors, hence, it is best not to set the upper limit too high.

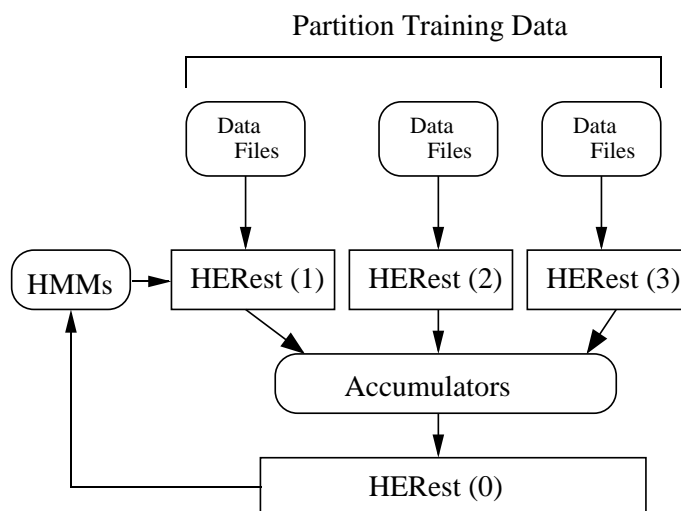


Fig. 8.8 HERest Parallel Operation

The second way of speeding-up the operation of HEREST is to use more than one computer in parallel. The way that this is done is to divide the training data amongst the available machines and then to run HEREST on each machine such that each invocation of HEREST uses the same initial set of models but has its own private set of data. By setting the option `-p N` where `N` is an integer, HEREST will dump the contents of all its accumulators into a file called `HERN.acc` rather than updating and outputting a new set of models. These dumped files are collected together and input to a new invocation of HEREST with the option `-p 0` set. HEREST then reloads the accumulators from all of the dump files and updates the models in the normal way. This process is illustrated in Figure 8.8.

To give a concrete example, suppose that four networked workstations were available to execute the HEREST command given earlier. The training files listed previously in `trainlist` would be split into four equal sets and a list of the files in each set stored in `trlist1`, `trlist2`, `trlist3`, and `trlist4`. On the first workstation, the command

```
HERest -S trlist1 -I labs -H dir1/hmacs -M dir2 -p 1 hmmlist
```

would be executed. This will load in the HMM definitions in `dir1/hmacs`, process the files listed in `trlist1` and finally dump its accumulators into a file called `HER1.acc` in the output directory `dir2`. At the same time, the command

```
HERest -S trlist2 -I labs -H dir1/hmacs -M dir2 -p 2 hmmlist
```

would be executed on the second workstation, and so on. When HEREST has finished on all four workstations, the following command will be executed on just one of them

```
HERest -H dir1/hmacs -M dir2 -p 0 hmmlist dir2/*.acc
```

where the list of training files has been replaced by the dumped accumulator files. This will cause the accumulated statistics to be reloaded and merged so that the model parameters can be reestimated and the new model set output to `dir2`. The time to perform this last phase of the operation is very small, hence the whole process will be around four times quicker than for the straightforward sequential case.

8.6 Single-Pass Retraining

In addition to re-estimating the parameters of a HMM set, HEREST also provides a mechanism for mapping a set of models trained using one parameterisation into another set based on a different parameterisation. This facility allows the front-end of a HMM-based recogniser to be modified without having to rebuild the models from scratch.

This facility is known as single-pass retraining. Given one set of well-trained models, a new set matching a different training data parameterisation can be generated in a single re-estimation pass. This is done by computing the forward and backward probabilities using the original models together with the original training data, but then switching to the new training data to compute the parameter estimates for the new set of models.

Single-pass retraining is enabled in HEREST by setting the `-r` switch. This causes the input training files to be read in pairs. The first of each pair is used to compute the forward/backward probabilities and the second is used to estimate the parameters for the new models. Very often, of course, data input to HTK is modified by the HPARM module in accordance with parameters set in a configuration file. In single-pass retraining mode, configuration parameters can be prefixed by the pseudo-module names HPARM1 and HPARM2. Then when reading in the first file of each pair, only the HPARM1 parameters are used and when reading the second file of each pair, only the HPARM2 parameters are used.

As an example, suppose that a set of models has been trained on data with MFCC_E_D parameterisation and a new set of models using Cepstral Mean Normalisation (_Z) is required. These two data parameterisations are specified in a configuration file (`config`) as two separate instances of the configuration variable TARGETKIND i.e.

```
# Single pass retraining
HPARM1: TARGETKIND = MFCC_E_D
HPARM2: TARGETKIND = MFCC_E_D_Z
```

HEREST would then be invoked with the `-r` option set to enable single-pass retraining. For example,

```
HERest -r -C config -S trainList -I labs -H dir1/hmacs -M dir2 hmmList
```

The script file `trainList` contains a list of data file pairs. For each pair, the first file should match the parameterisation of the original model set and the second file should match that of the required new set. This will cause the model parameter estimates to be performed using the new set of training data and a new set of models matching this data will be output to `dir2`. This process of single-pass retraining is a significantly faster route to a new set of models than training a fresh set from scratch.

8.7 Parameter Re-Estimation Formulae

For reference purposes, this section lists the various formulae employed within the HTK parameter estimation tools. All are standard, however, the use of non-emitting states and multiple data streams leads to various special cases which are usually not covered fully in the literature.

The following notation is used in this section

N	number of states
S	number of streams
M_s	number of mixture components in stream s
T	number of observations
Q	number of models in an embedded training sequence
N_q	number of states in the q 'th model in a training sequence
\mathbf{O}	a sequence of observations
\mathbf{o}_t	the observation at time t , $1 \leq t \leq T$
\mathbf{o}_{st}	the observation vector for stream s at time t
a_{ij}	the probability of a transition from state i to j
c_{jsm}	weight of mixture component m in state j stream s
$\boldsymbol{\mu}_{jsm}$	vector of means for the mixture component m of state j stream s
$\boldsymbol{\Sigma}_{jsm}$	covariance matrix for the mixture component m of state j stream s
λ	the set of all parameters defining a HMM

8.7.1 Viterbi Training (HInit)

In this style of model training, a set of training observations \mathbf{O}^r , $1 \leq r \leq R$ is used to estimate the parameters of a single HMM by iteratively computing Viterbi alignments. When used to initialise a new HMM, the Viterbi segmentation is replaced by a uniform segmentation (i.e. each training observation is divided into N equal segments) for the first iteration.

Apart from the first iteration on a new model, each training sequence \mathbf{O} is segmented using a state alignment procedure which results from maximising

$$\phi_N(T) = \max_i \phi_i(T) a_{iN}$$

for $1 < i < N$ where

$$\phi_j(t) = \left[\max_i \phi_i(t-1) a_{ij} \right] b_j(\mathbf{o}_t)$$

with initial conditions given by

$$\phi_1(1) = 1$$

$$\phi_j(1) = a_{1j} b_j(\mathbf{o}_1)$$

for $1 < j < N$. In this and all subsequent cases, the output probability $b_j(\cdot)$ is as defined in equations 7.1 and 7.2 in section 7.1.

If A_{ij} represents the total number of transitions from state i to state j in performing the above maximisations, then the transition probabilities can be estimated from the relative frequencies

$$\hat{a}_{ij} = \frac{A_{ij}}{\sum_{k=2}^N A_{ik}}$$

The sequence of states which maximises $\phi_N(T)$ implies an alignment of training data observations with states. Within each state, a further alignment of observations to mixture components is made. The tool HINIT provides two mechanisms for this: for each state and each stream

1. use clustering to allocate each observation \mathbf{o}_{st} to one of M_s clusters, or
2. associate each observation \mathbf{o}_{st} with the mixture component with the highest probability

In either case, the net result is that every observation is associated with a single unique mixture component. This association can be represented by the indicator function $\psi_{j_{sm}}^r(t)$ which is 1 if \mathbf{o}_{st}^r is associated with mixture component m of stream s of state j and is zero otherwise.

The means and variances are then estimated via simple averages

$$\hat{\boldsymbol{\mu}}_{j_{sm}} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \psi_{j_{sm}}^r(t) \mathbf{o}_{st}^r}{\sum_{r=1}^R \sum_{t=1}^{T_r} \psi_{j_{sm}}^r(t)}$$

$$\hat{\boldsymbol{\Sigma}}_{j_{sm}} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \psi_{j_{sm}}^r(t) (\mathbf{o}_{st}^r - \hat{\boldsymbol{\mu}}_{j_{sm}}) (\mathbf{o}_{st}^r - \hat{\boldsymbol{\mu}}_{j_{sm}})'}{\sum_{r=1}^R \sum_{t=1}^{T_r} \psi_{j_{sm}}^r(t)}$$

Finally, the mixture weights are based on the number of observations allocated to each component

$$\mathbf{c}_{j_{sm}} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \psi_{j_{sm}}^r(t)}{\sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{l=1}^{M_s} \psi_{j_{sl}}^r(t)}$$

8.7.2 Forward/Backward Probabilities

Baum-Welch training is similar to the Viterbi training described in the previous section except that the *hard* boundary implied by the ψ function is replaced by a *soft* boundary function L which represents the probability of an observation being associated any given Gaussian mixture component. This *occupation* probability is computed from the *forward* and *backward* probabilities.

For the isolated-unit style of training, the forward probability $\alpha_j(t)$ for $1 < j < N$ and $1 < t \leq T$ is calculated by the forward recursion

$$\alpha_j(t) = \left[\sum_{i=2}^{N-1} \alpha_i(t-1) a_{ij} \right] b_j(\mathbf{o}_t)$$

with initial conditions given by

$$\begin{aligned}\alpha_1(1) &= 1 \\ \alpha_j(1) &= a_{1j}b_j(\mathbf{o}_1)\end{aligned}$$

for $1 < j < N$ and final condition given by

$$\alpha_N(T) = \sum_{i=2}^{N-1} \alpha_i(T)a_{iN}$$

The backward probability $\beta_i(t)$ for $1 < i < N$ and $T > t \geq 1$ is calculated by the backward recursion

$$\beta_i(t) = \sum_{j=2}^{N-1} a_{ij}b_j(\mathbf{o}_{t+1})\beta_j(t+1)$$

with initial conditions given by

$$\beta_i(T) = a_{iN}$$

for $1 < i < N$ and final condition given by

$$\beta_1(1) = \sum_{j=2}^{N-1} a_{1j}b_j(\mathbf{o}_1)\beta_j(1)$$

In the case of embedded training where the HMM spanning the observations is a composite constructed by concatenating Q subword models, it is assumed that at time t , the α and β values corresponding to the entry state and exit states of a HMM represent the forward and backward probabilities at time $t - \Delta t$ and $t + \Delta t$, respectively, where Δt is small. The equations for calculating α and β are then as follows.

For the forward probability, the initial conditions are established at time $t = 1$ as follows

$$\begin{aligned}\alpha_1^{(q)}(1) &= \begin{cases} 1 & \text{if } q = 1 \\ \alpha_1^{(q-1)}(1)a_{1N_{q-1}}^{(q-1)} & \text{otherwise} \end{cases} \\ \alpha_j^{(q)}(1) &= a_{1j}^{(q)}b_j^{(q)}(\mathbf{o}_1) \\ \alpha_{N_q}^{(q)}(1) &= \sum_{i=2}^{N_q-1} \alpha_i^{(q)}(1)a_{iN_q}^{(q)}\end{aligned}$$

where the superscript in parentheses refers to the index of the model in the sequence of concatenated models. All unspecified values of α are zero. For time $t > 1$,

$$\begin{aligned}\alpha_1^{(q)}(t) &= \begin{cases} 0 & \text{if } q = 1 \\ \alpha_{N_{q-1}}^{(q-1)}(t-1) + \alpha_1^{(q-1)}(t)a_{1N_{q-1}}^{(q-1)} & \text{otherwise} \end{cases} \\ \alpha_j^{(q)}(t) &= \left[\alpha_1^{(q)}(t)a_{1j}^{(q)} + \sum_{i=2}^{N_q-1} \alpha_i^{(q)}(t-1)a_{ij}^{(q)} \right] b_j^{(q)}(\mathbf{o}_t) \\ \alpha_{N_q}^{(q)}(t) &= \sum_{i=2}^{N_q-1} \alpha_i^{(q)}(t)a_{iN_q}^{(q)}\end{aligned}$$

For the backward probability, the initial conditions are set at time $t = T$ as follows

$$\begin{aligned}\beta_{N_q}^{(q)}(T) &= \begin{cases} 1 & \text{if } q = Q \\ \beta_{N_{q+1}}^{(q+1)}(T)a_{1N_{q+1}}^{(q+1)} & \text{otherwise} \end{cases} \\ \beta_i^{(q)}(T) &= a_{iN_q}^{(q)}\beta_{N_q}^{(q)}(T) \\ \beta_1^{(q)}(T) &= \sum_{j=2}^{N_q-1} a_{1j}^{(q)}b_j^{(q)}(\mathbf{o}_T)\beta_j^{(q)}(T)\end{aligned}$$

where once again, all unspecified β values are zero. For time $t < T$,

$$\beta_{N_q}^{(q)}(t) = \begin{cases} 0 & \text{if } q = Q \\ \beta_1^{(q+1)}(t+1) + \beta_{N_{q+1}}^{(q+1)}(t)a_{1N_{q+1}}^{(q+1)} & \text{otherwise} \end{cases}$$

$$\beta_i^{(q)}(t) = a_{iN_q}^{(q)}\beta_{N_q}^{(q)}(t) + \sum_{j=2}^{N_q-1} a_{ij}^{(q)}b_j^{(q)}(\mathbf{o}_{t+1})\beta_j^{(q)}(t+1)$$

$$\beta_1^{(q)}(t) = \sum_{j=2}^{N_q-1} a_{1j}^{(q)}b_j^{(q)}(\mathbf{o}_t)\beta_j^{(q)}(t)$$

The total probability $P = \text{prob}(\mathbf{O}|\lambda)$ can be computed from either the forward or backward probabilities

$$P = \alpha_N(T) = \beta_1(1)$$

8.7.3 Single Model Reestimation(HRest)

In this style of model training, a set of training observations \mathbf{O}^r , $1 \leq r \leq R$ is used to estimate the parameters of a single HMM. The basic formula for the reestimation of the transition probabilities is

$$\hat{a}_{ij} = \frac{\sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r-1} \alpha_i^r(t)a_{ij}b_j(\mathbf{o}_{t+1}^r)\beta_j^r(t+1)}{\sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r} \alpha_i^r(t)\beta_i^r(t)}$$

where $1 < i < N$ and $1 < j < N$ and P_r is the total probability $P = \text{prob}(\mathbf{O}^r|\lambda)$ of the r 'th observation. The transitions from the non-emitting entry state are reestimated by

$$\hat{a}_{1j} = \frac{1}{R} \sum_{r=1}^R \frac{1}{P_r} \alpha_j^r(1)\beta_j^r(1)$$

where $1 < j < N$ and the transitions from the emitting states to the final non-emitting exit state are reestimated by

$$\hat{a}_{iN} = \frac{\sum_{r=1}^R \frac{1}{P_r} \alpha_i^r(T)\beta_i^r(T)}{\sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r} \alpha_i^r(t)\beta_i^r(t)}$$

where $1 < i < N$.

For a HMM with M_s mixture components in stream s , the means, covariances and mixture weights for that stream are reestimated as follows. Firstly, the probability of occupying the m 'th mixture component in stream s at time t for the r 'th observation is

$$L_{jsm}^r(t) = \frac{1}{P_r} U_j^r(t) c_{jsm} b_{jsm}(\mathbf{o}_{st}^r) \beta_j^r(t) b_{js}^*(\mathbf{o}_t^r)$$

where

$$U_j^r(t) = \begin{cases} a_{1j} & \text{if } t = 1 \\ \sum_{i=2}^{N-1} \alpha_i^r(t-1)a_{ij} & \text{otherwise} \end{cases} \quad (8.1)$$

and

$$b_{js}^*(\mathbf{o}_t^r) = \prod_{k \neq s} b_{jk}(\mathbf{o}_{kt}^r)$$

For single Gaussian streams, the probability of mixture component occupancy is equal to the probability of state occupancy and hence it is more efficient in this case to use

$$L_{jsm}^r(t) = L_j^r(t) = \frac{1}{P_r} \alpha_j(t)\beta_j(t)$$

Given the above definitions, the re-estimation formulae may now be expressed in terms of $L_{jsm}^r(t)$ as follows.

$$\hat{\boldsymbol{\mu}}_{jsm} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} L_{jsm}^r(t) \mathbf{o}_{st}^r}{\sum_{r=1}^R \sum_{t=1}^{T_r} L_{jsm}^r(t)}$$

$$\hat{\Sigma}_{j_{sm}} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} L_{j_{sm}}^r(t) (\mathbf{o}_{st}^r - \hat{\boldsymbol{\mu}}_{j_{sm}}) (\mathbf{o}_{st}^r - \hat{\boldsymbol{\mu}}_{j_{sm}})'}{\sum_{r=1}^R \sum_{t=1}^{T_r} L_{j_{sm}}^r(t)} \quad (8.2)$$

$$\mathbf{c}_{j_{sm}} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} L_{j_{sm}}^r(t)}{\sum_{r=1}^R \sum_{t=1}^{T_r} L_j^r(t)}$$

8.7.4 Embedded Model Reestimation(HERest)

The re-estimation formulae for the embedded model case have to be modified to take account of the fact that the entry states can be occupied at any time as a result of transitions out of the previous model. The basic formulae for the re-estimation of the transition probabilities is

$$\hat{a}_{ij}^{(q)} = \frac{\sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r-1} \alpha_i^{(q)r}(t) a_{ij}^{(q)} b_j^{(q)}(\mathbf{o}_{t+1}^r) \beta_j^{(q)r}(t+1)}{\sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r} \alpha_i^{(q)r}(t) \beta_i^{(q)r}(t)}$$

The transitions from the non-emitting entry states into the HMM are re-estimated by

$$\hat{a}_{1j}^{(q)} = \frac{\sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r-1} \alpha_1^{(q)r}(t) a_{1j}^{(q)} b_j^{(q)}(\mathbf{o}_t^r) \beta_j^{(q)r}(t)}{\sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r} \alpha_1^{(q)r}(t) \beta_1^{(q)r}(t) + \alpha_1^{(q)r}(t) a_{1N_q}^{(q)} \beta_1^{(q+1)r}(t)}$$

and the transitions out of the HMM into the non-emitting exit states are re-estimated by

$$\hat{a}_{iN_q}^{(q)} = \frac{\sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r-1} \alpha_i^{(q)r}(t) a_{iN_q}^{(q)} \beta_{N_q}^{(q)r}(t)}{\sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r} \alpha_i^{(q)r}(t) \beta_i^{(q)r}(t)}$$

Finally, the direct transitions from non-emitting entry to non-emitting exit states are re-estimated by

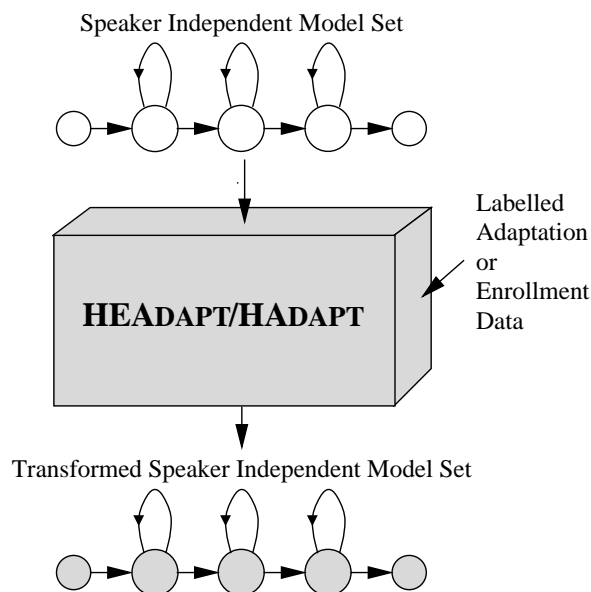
$$\hat{a}_{1N_q}^{(q)} = \frac{\sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r-1} \alpha_1^{(q)r}(t) a_{1N_q}^{(q)} \beta_1^{(q+1)r}(t)}{\sum_{r=1}^R \frac{1}{P_r} \sum_{t=1}^{T_r} \alpha_i^{(q)r}(t) \beta_i^{(q)r}(t) + \alpha_1^{(q)r}(t) a_{1N_q}^{(q)} \beta_1^{(q+1)r}(t)}$$

The re-estimation formulae for the output distributions are the same as for the single model case except for the obvious additional subscript for q . However, the probability calculations must now allow for transitions from the entry states by changing $U_j^r(t)$ in equation 8.1 to

$$U_j^{(q)r}(t) = \begin{cases} \alpha_1^{(q)r}(t) a_{1j}^{(q)} & \text{if } t = 1 \\ \alpha_1^{(q)r}(t) a_{1j}^{(q)} + \sum_{i=2}^{N_q-1} \alpha_i^{(q)r}(t-1) a_{ij}^{(q)} & \text{otherwise} \end{cases}$$

Chapter 9

HMM Adaptation



Chapter 8 described how the parameters are estimated for plain continuous density HMMs within HTK, primarily using the embedded training tool HEREST. Using the training strategy depicted in figure 8.2, together with other techniques can produce high performance speaker independent acoustic models for a large vocabulary recognition system. However it is possible to build improved acoustic models by tailoring a model set to a specific speaker. By collecting data from a speaker and training a model set on this speaker's data alone, the speaker's characteristics can be modelled more accurately. Such systems are commonly known as *speaker dependent* systems, and on a typical word recognition task, may have half the errors of a speaker independent system. The drawback of speaker dependent systems is that a large amount of data (typically hours) must be collected in order to obtain sufficient model accuracy.

Rather than training speaker dependent models, *adaptation* techniques can be applied. In this case, by using only a small amount of data from a new speaker, a good speaker independent system model set can be adapted to better fit the characteristics of this new speaker.

Speaker adaptation techniques can be used in various different modes. If the true transcription of the adaptation data is known then it is termed *supervised adaptation*, whereas if the adaptation data is unlabelled then it is termed *unsupervised adaptation*. In the case where all the adaptation data is available in one block, e.g. from a speaker enrollment session, then this termed *static adaptation*. Alternatively adaptation can proceed incrementally as adaptation data becomes available, and this is termed *incremental adaptation*.

HTK provides two tools to adapt continuous density HMMs. HEADAPT performs offline supervised adaptation using maximum likelihood linear regression (MLLR) and/or maximum a-posteriori (MAP) adaptation, while unsupervised adaptation is supported by HVITE (using only MLLR). In

this case HVITE not only performs recognition, but simultaneously adapts the model set as the data becomes available through recognition. Currently, MLLR adaptation can be applied in both incremental and static modes while MAP supports only static adaptation. If MLLR and MAP adaptation is to be performed simultaneously using HEADAPT in the same pass, then the restriction is that the entire adaptation must be performed statically¹.

This chapter describes the supervised adaptation tool HEADAPT. The first sections of the chapter give an overview of MLLR and MAP adaptation and this is followed by a section describing the general usages of HEADAPT to build simple and more complex adapted systems. The chapter concludes with a section detailing the various formulae used by the adaptation tool. The use of HVITE to perform unsupervised adaptation is discussed in section 13.6.2.

9.1 Model Adaptation using MLLR

9.1.1 Maximum Likelihood Linear Regression

Maximum likelihood linear regression or MLLR computes a set of transformations that will reduce the mismatch between an initial model set and the adaptation data². More specifically MLLR is a model adaptation technique that estimates a set of linear transformations for the mean and variance parameters of a Gaussian mixture HMM system. The effect of these transformations is to shift the component means and alter the variances in the initial system so that each state in the HMM system is more likely to generate the adaptation data. Note that due to computational reasons, MLLR is only implemented within HTK for diagonal covariance, single stream, continuous density HMMs.

The transformation matrix used to give a new estimate of the adapted mean is given by

$$\hat{\boldsymbol{\mu}} = \mathbf{W}\boldsymbol{\xi}, \quad (9.1)$$

where \mathbf{W} is the $n \times (n + 1)$ transformation matrix (where n is the dimensionality of the data) and $\boldsymbol{\xi}$ is the extended mean vector,

$$\boldsymbol{\xi} = [w \ \mu_1 \ \mu_2 \ \dots \ \mu_n]^T$$

where w represents a bias offset whose value is fixed (within HTK) at 1.

Hence \mathbf{W} can be decomposed into

$$\mathbf{W} = [\mathbf{b} \ \mathbf{A}] \quad (9.2)$$

where \mathbf{A} represents an $n \times n$ transformation matrix and \mathbf{b} represents a bias vector.

The transformation matrix \mathbf{W} is obtained by solving a maximisation problem using the *Expectation-Maximisation* (EM) technique. This technique is also used to compute the variance transformation matrix. Using EM results in the maximisation of a standard *auxiliary function*. (Full details are available in section 9.4.)

9.1.2 MLLR and Regression Classes

This adaptation method can be applied in a very flexible manner, depending on the amount of adaptation data that is available. If a small amount of data is available then a *global* adaptation transform can be generated. A global transform (as its name suggests) is applied to every Gaussian component in the model set. However as more adaptation data becomes available, improved adaptation is possible by increasing the number of transformations. Each transformation is now more specific and applied to certain groupings of Gaussian components. For instance the Gaussian components could be grouped into the broad phone classes: silence, vowels, stops, glides, nasals, fricatives, etc. The adaptation data could now be used to construct more specific broad class transforms to apply to these groupings.

Rather than specifying static component groupings or classes, a robust and dynamic method is used for the construction of further transformations as more adaptation data becomes available. MLLR makes use of a *regression class tree* to group the Gaussians in the model set, so that the set of transformations to be estimated can be chosen according to the amount and type of adaptation data that is available. The tying of each transformation across a number of mixture components

¹ By using two passes, one could perform incremental MLLR in the first pass (saving the new model or transform), followed by a second pass, this time using MAP adaptation.

² MLLR can also be used to perform environmental compensation by reducing the mismatch due to channel or additive noise effects.

makes it possible to adapt distributions for which there were no observations at all. With this process all models can be adapted and the adaptation process is dynamically refined when more adaptation data becomes available.

The regression class tree is constructed so as to cluster together components that are close in acoustic space, so that similar components can be transformed in a similar way. Note that the tree is built using the original speaker independent model set, and is thus independent of any new speaker. The tree is constructed with a centroid splitting algorithm, which uses a Euclidean distance measure. For more details see section 10.7. The terminal nodes or leaves of the tree specify the final component groupings, and are termed the *base (regression) classes*. Each Gaussian component of a model set belongs to one particular base class. The tool HHED can be used to build a binary regression class tree, and to label each component with a base class number. Both the tree and component base class numbers are saved automatically as part of the MMF. Please refer to section 7.8 and section 10.7 for further details.

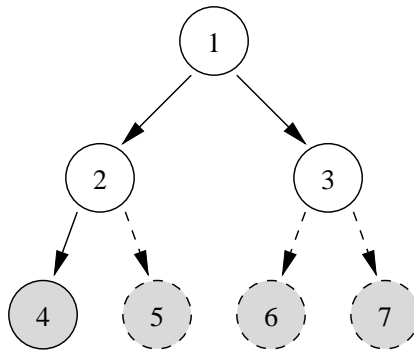


Fig. 9.1 A binary regression tree

Figure 9.1 shows a simple example of a binary regression tree with four base classes, denoted as $\{C_4, C_5, C_6, C_7\}$. During “dynamic” adaptation, the occupation counts are accumulated for each of the regression base classes. The diagram shows a solid arrow and circle (or node), indicating that there is sufficient data for a transformation matrix to be generated using the data associated with that class. A dotted line and circle indicates that there is insufficient data. For example neither node 6 or 7 has sufficient data; however when pooled at node 3, there is sufficient adaptation data. The amount of data that is “determined” as sufficient is set by the user as a command-line option to HEADAPT (see reference section 14.5).

HEADAPT uses a top-down approach to traverse the regression class tree. Here the search starts at the root node and progresses down the tree generating transforms only for those nodes which

1. have sufficient data **and**
2. are either terminal nodes (i.e. base classes) **or** have any children without sufficient data.

In the example shown in figure 9.1, transforms are constructed only for regression nodes 2, 3 and 4, which can be denoted as \mathbf{W}_2 , \mathbf{W}_3 and \mathbf{W}_4 . Hence when the transformed model set is required, the transformation matrices (mean and variance) are applied in the following fashion to the Gaussian components in each base class:-

$$\left\{ \begin{array}{l} \mathbf{W}_2 \rightarrow \{C_5\} \\ \mathbf{W}_3 \rightarrow \{C_6, C_7\} \\ \mathbf{W}_4 \rightarrow \{C_4\} \end{array} \right\}$$

At this point it is interesting to note that the global adaptation case is the same as a tree with just a root node, and is in fact treated as such.

9.1.3 Transform Model File Format

HEADAPT estimates the required transformation statistics and can either output a transformed MMF or a transform model file (TMF). The advantage in storing the transforms as opposed to an adapted MMF is that the TMFs are considerably smaller than MMFs (especially triphone MMFs). This section describes the format of the transform model file in detail.

The mean transformation matrix is stored as a block diagonal transformation matrix. The example block diagonal matrix \mathbf{A} shown below contains three blocks. The first block represents the transformation for only the static components of the feature vector, while the second represents the deltas and the third the accelerations. This block diagonal matrix example makes the assumption that for the transformation, there is no correlation between the statics, deltas and delta deltas. In practice this assumption works quite well.

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_s & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_\Delta & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{A}_{\Delta^2} \end{pmatrix}$$

This format reduces the number of transformation parameters required to be learnt, making the adaptation process faster. It also reduces the adaptation data required per transform when compared with the full case. When comparing the storage requirements, the 3 block diagonal matrix requires much less storage capacity than the full transform matrix. Note that for convenience a full transformation matrix is also stored as a block diagonal matrix, only in this case there is a single block.

The variance transformation is a diagonal matrix and as such is simply stored as a vector. Figure 9.2 shows a simple example of a TMF. In this case the feature vector has nine dimensions, and the mean transform has three diagonal blocks. The TMF can be saved in ASCII or binary format. The user header is always output in ascii. The first two fields are speaker descriptor fields. The next field <MMFID>, the MMF identifier, is obtained from the global options macro in the MMF, while the regression class tree identifier <RCID> is obtained from the regression tree macro name in the MMF. If global adaptation is being performed, then the <RCID> will contain the identifier `global`, since a tree is unnecessary in the global case. Note that the MMF and regression class tree identifiers are set within the MMF using the tool HHED. The final two fields are optional, but HEADAPT outputs these anyway for the user's convenience. These can be edited at any time (as can all the fields if desired, but editing <MMFID> and <RCID> fields should be avoided). The <CHAN> field should represent the adaptation data recording environment. Examples could be a particular microphone name, telephone channel or various background noise conditions. The <DESC> allow the user to enter any other information deemed useful. An example could be the speaker's dialect region.

```

<UID> djk
<NAME> Dan Kershaw
<MMFID> ECRL_UK_XWRD
<RCID> global
<CHAN> Standard
<DESC> None
<NBLOCKS> 3
<NODETHRESH> 700.0
<NODEOCC> 1 24881.8
<TRANSFORM> 1
  <MEAN_TR> 3
    <BLOCK> 1
      0.942 -0.032 -0.001
      -0.102 0.922 -0.015
      -0.016 0.045 0.910
    <BLOCK> 2
      1.021 -0.032 -0.011
      -0.017 1.074 -0.043
      -0.099 0.091 1.050
    <BLOCK> 3
      1.028 0.032 0.001
      -0.012 1.014 -0.011
      -0.091 -0.043 1.041
  <BIASOFFSET> 9
    -0.357 0.001 -0.002 0.132 0.072
    0.006 0.150 0.138 0.198
  <VARIANCE_TR> 9
    0.936 0.865 0.848 0.832 0.829
    0.786 0.947 0.869 0.912

```

Fig. 9.2 A Simple example of a TMF

Whenever a TMF is being used (in conjunction with an MMF), the MMF identifier in the MMF is checked against that in the TMF. These **must** match since the TMF is dependent on the model set it was constructed from. Also unless the `<RCID>` field is set to `global`, it is also checked for consistency against the regression tree identifier in the MMF.

The rest of the TMF contains a further information header, followed by all the transforms. The information header contains necessary transform set information such as the number of blocks used, node occupation threshold used, and the node occupation counts. Each transform has a regression class identifier number, the mean transformation matrix \mathbf{A} , an optional bias vector \mathbf{b} (as in equation 9.2) and an optional variance transformation diagonal matrix \mathbf{H} (stored as a vector). The example has both a bias offset and a variance transform.

9.2 Model Adaptation using MAP

Model adaptation can also be accomplished using a maximum a posteriori (MAP) approach. This adaptation process is sometimes referred to as Bayesian adaptation. MAP adaptation involves the use of prior knowledge about the model parameter distribution. Hence, if we know what the parameters of the model are likely to be (before observing any adaptation data) using the prior knowledge, we might well be able to make good use of the limited adaptation data, to obtain a decent MAP estimate. This type of prior is often termed an informative prior. Note that if the prior distribution indicates no preference as to what the model parameters are likely to be (a non-informative prior), then the MAP estimate obtained will be identical to that obtained using a maximum likelihood approach.

For MAP adaptation purposes, the informative priors that are generally used are the speaker independent model parameters. For mathematical tractability conjugate priors are used, which

results in a simple adaptation formula. The update formula for a single stream system for state j and mixture component m is

$$\hat{\boldsymbol{\mu}}_{jm} = \frac{N_{jm}}{N_{jm} + \tau} \bar{\boldsymbol{\mu}}_{jm} + \frac{\tau}{N_{jm} + \tau} \boldsymbol{\mu}_{jm} \quad (9.3)$$

where τ is a weighting of the a priori knowledge to the adaptation speech data and N is the occupation likelihood of the adaptation data, defined as,

$$N_{jm} = \sum_{r=1}^R \sum_{t=1}^{T_r} L_{jm}^r(t)$$

where $\boldsymbol{\mu}_{jm}$ is the speaker independent mean and $\bar{\boldsymbol{\mu}}_{jm}$ is the mean of the observed adaptation data and is defined as,

$$\bar{\boldsymbol{\mu}}_{jm} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} L_{jm}^r(t) \boldsymbol{o}_t^r}{\sum_{r=1}^R \sum_{t=1}^{T_r} L_{jm}^r(t)}$$

As can be seen, if the occupation likelihood of a Gaussian component (N_{jm}) is small, then the mean MAP estimate will remain close to the speaker independent component mean. With MAP adaptation, every single mean component in the system is updated with a MAP estimate, based on the prior mean, the weighting and the adaptation data. Hence, MAP adaptation requires a new “speaker-dependent” model set to be saved.

One obvious drawback to MAP adaptation is that it requires more adaptation data to be effective when compared to MLLR, because MAP adaptation is specifically defined at the component level. When larger amounts of adaptation training data become available, MAP begins to perform better than MLLR, due to this detailed update of each component (rather than the pooled Gaussian transformation approach of MLLR). In fact the two adaptation processes can be combined to improve performance still further, by using the MLLR transformed means as the priors for MAP adaptation (by replacing $\boldsymbol{\mu}_{jm}$ in equation 9.3 with the transformed mean of equation 9.1). In this case components that have a low occupation likelihood in the adaptation data, (and hence would not change much using MAP alone) have been adapted using a regression class transform in MLLR. An example usage is shown in the following section.

9.3 Using HEAdapt

At the outset HEADAPT operates in a very similar fashion to HEREST. Both use a frame/state alignment in order to accumulate various statistics about the data. In HEREST these statistics are used to estimate new model parameters whilst in HEADAPT they are used to estimate the transformations for each regression base class, or new model parameters. HEADAPT will currently only produce transforms with single stream data and PLAINHS or SHAREDHS HMM systems (see section 7.4 on HMM set kinds).

In outline, HEADAPT works as follows. On startup, HEADAPT loads in a complete set of HMM definitions, including, the regression class tree and the base class number of each Gaussian component. Note that HEADAPT requires the MMF to contain a regression class tree. Every training file must have an associated label file which gives a transcription for that file. Only the sequence of labels is used by HEADAPT, and any boundary location information is ignored. Thus, these transcriptions can be generated automatically from the known orthography of what was said and a pronunciation dictionary.

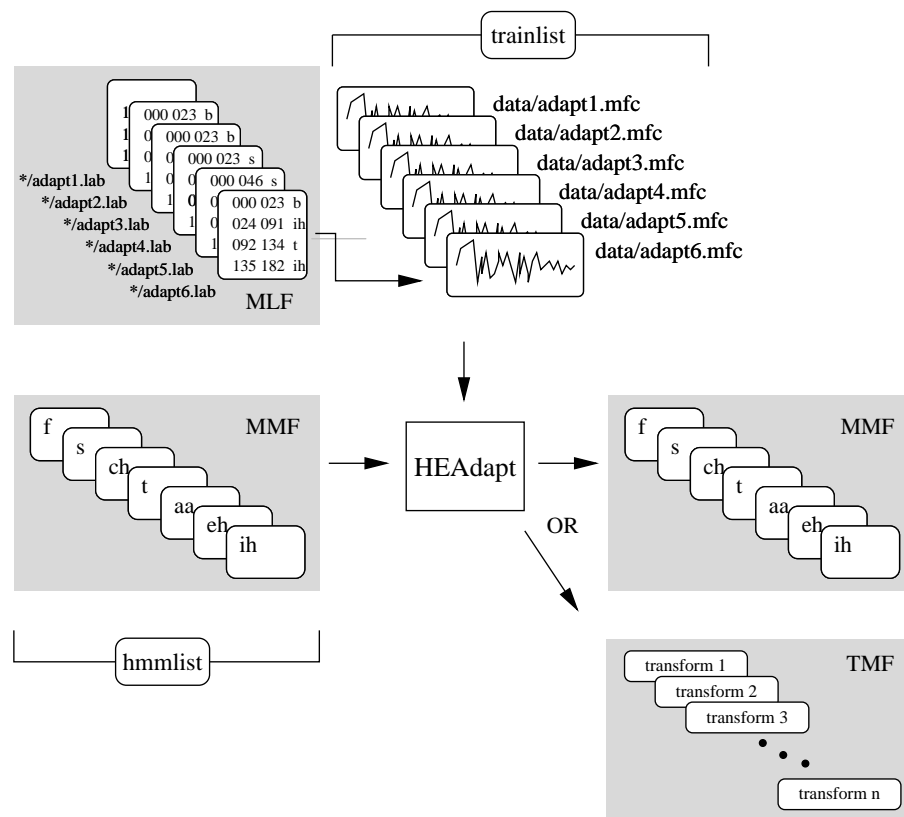


Fig. 9.3 File Processing in HEAdapt

HEADAPT processes each training file in turn. After loading it into memory, it uses the associated transcription to construct a composite HMM which spans the whole utterance. This composite HMM is made by concatenating instances of the phone HMMs corresponding to each label in the transcription. The Forward-Backward algorithm is then applied to obtain a frame/state alignment and the information necessary to form the standard auxiliary function is accumulated at the Gaussian component level. Note that this information is different from that required in HEREST (see section 9.4). When all of the training files have been processed (within the static or incremental block), the regression base class statistics are accumulated using the component level statistics. Next the regression class tree is traversed and the new regression class transformations are calculated for those regression classes containing a sufficient occupation count at the lowest level in the tree, as described in section 9.1.2. Finally either the updated (i.e. adapted) HMM set or the transformations are output. Note that HEADAPT produces a transforms model file (TMF) that contains

transforms that are estimated to *transform from the input MMF* to a new environment/speaker based on the adaptation data presented.

The mathematical details of the Forward-Backward algorithm are given in section 8.7, while the mathematical details for the MLLR mean and variance transformation calculations can be found in section 9.4.

HEADAPT is typically invoked by a command line of the form

```
HEAdapt -S adaptlist -I labs -H dir1/hmacs -M dir2 hmmlist
```

where `hmmlist` contains the list of HMMs.

Once all MMFs and MLFs have been loaded, HEADAPT processes each file in the `adaptlist`, and accumulates the required statistics as described above. On completion, an updated MMF is output to the directory `dir2`.

If the following form of the command is used

```
HEAdapt -S adaptlist -I labs -H dir1/hmacs -K dir2/tmf hmmlist
```

then on completion a transform model file (TMF) `tmf` is output to the directory `dir2`. This process is illustrated by Fig 9.3. Section 9.1.3 describes the TMF format in more detail. The output `tmf` contains transforms that transform the MMF `hmacs`. Once this is saved, HVITE can be used to perform recognition for the adapted speaker either using a transformed MMF or by using the speaker independent MMF together with a speaker specific TMF.

HEADAPT employs the same pruning mechanism as HEREST during the forward-backward computation. As such the pruning on the backward path is under the user's control, and the beam is set using the `-t` option.

HEADAPT can also be run several times in block or static fashion. For instance a first pass might entail a global adaptation (forced using the `-g` option), producing the TMF `global.tmf` by invoking

```
HEAdapt -g -S adaptlist -I labs -H mmf -K tmfs/global.tmf \
        hmmlist
```

The second pass could load in the global transformations (and transform the model set) using the `-J` option, performing a better frame/state alignment than the speaker independent model set, and output a set of regression class transformations,

```
HEAdapt -S adaptlist -I labs -H mmf -K tmfs/rc.tmf \
        -J tmfs/global.tmf hmmlist
```

Note again that the number of transformations is selected automatically and is dependent on the node occupation threshold setting and the amount of adaptation data available. Finally when producing a TMF, HEADAPT always generates a TMF to transform the input MMF in all cases. In the last example the input MMF is transformed by the global transform file `global.tmf` in order to obtain the frame/state alignment only. The final TMF that is output, `rc.tmf`, contains the set of transforms to transform the input MMF `mmf`, based on this frame/state alignment.

As an alternative, the second pass could entail MLLR together with MAP adaptation, outputting a new model set. Note that with MAP adaptation a transform can not be saved and a full HMM set must be output.

```
HEAdapt -S adaptlist -I labs -H mmf -M dir2 -k -j 12.0
        -J tmfs/global.tmf hmmlist
```

Note that MAP alone could be used by removing the `-k` option. The argument to the `-j` option represents the MAP adaptation scaling factor.

9.4 MLLR Formulae

For reference purposes, this section lists the various formulae employed within the HTK adaptation tool. It is assumed throughout that single stream data is used and that diagonal covariances are also used. All are standard and can be found in various literature.

The following notation is used in this section

\mathcal{M}	the model set
T	number of observations
m	a mixture component
\mathbf{O}	a sequence of observations
$\mathbf{o}(t)$	the observation at time t , $1 \leq t \leq T$
$\boldsymbol{\mu}_{m_r}$	mean vector for the mixture component m_r
$\boldsymbol{\xi}_{m_r}$	extended mean vector for the mixture component m_r
$\boldsymbol{\Sigma}_{m_r}$	covariance matrix for the mixture component m_r
$L_{m_r}(t)$	the occupancy probability for the mixture component m_r at time t

9.4.1 Estimation of the Mean Transformation Matrix

To enable robust transformations to be trained, the transform matrices are tied across a number of Gaussians. The set of Gaussians which share a transform is referred to as a regression class. For a particular transform case \mathbf{W}_m , the R Gaussian components $\{m_1, m_2, \dots, m_R\}$ will be tied together, as determined by the regression class tree (see section 9.1.2). By formulating the standard auxiliary function, and then maximising it with respect to the transformed mean, and considering only these tied Gaussian components, the following is obtained,

$$\sum_{t=1}^T \sum_{r=1}^R L_{m_r}(t) \boldsymbol{\Sigma}_{m_r}^{-1} \mathbf{o}(t) \boldsymbol{\xi}_{m_r}^T = \sum_{t=1}^T \sum_{r=1}^R L_{m_r}(t) \boldsymbol{\Sigma}_{m_r}^{-1} \mathbf{W}_m \boldsymbol{\xi}_{m_r} \boldsymbol{\xi}_{m_r}^T \quad (9.4)$$

and $L_{m_r}(t)$, the occupation likelihood, is defined as,

$$L_{m_r}(t) = p(q_{m_r}(t) \mid \mathcal{M}, \mathbf{O}_T)$$

where $q_{m_r}(t)$ indicates the Gaussian component m_r at time t , and $\mathbf{O}_T = \{\mathbf{o}(1), \dots, \mathbf{o}(T)\}$ is the adaptation data. The occupation likelihood is obtained from the forward-backward process described in section 8.7.

To solve for \mathbf{W}_m , two new terms are defined.

1. The left hand side of equation 9.4 is independent of the transformation matrix and is referred to as \mathbf{Z} , where

$$\mathbf{Z} = \sum_{t=1}^T \sum_{r=1}^R L_{m_r}(t) \boldsymbol{\Sigma}_{m_r}^{-1} \mathbf{o}(t) \boldsymbol{\xi}_{m_r}^T$$

2. A new variable $\mathbf{G}^{(i)}$ is defined with elements

$$g_{jq}^{(i)} = \sum_{r=1}^R v_{ii}^{(r)} d_{jq}^{(r)}$$

where

$$\mathbf{V}^{(r)} = \sum_{t=1}^T L_{m_r}(t) \boldsymbol{\Sigma}_{m_r}^{-1}$$

and

$$\mathbf{D}^{(r)} = \boldsymbol{\xi}_{m_r} \boldsymbol{\xi}_{m_r}^T$$

It can be seen that from these two new terms, \mathbf{W}_m can be calculated from

$$\mathbf{w}_i^T = \mathbf{G}_i^{-1} \mathbf{z}_i^T$$

where \mathbf{w}_i is the i^{th} vector of \mathbf{W}_m and \mathbf{z}_i is the i^{th} vector of \mathbf{Z} .

The regression class tree is used to generate the classes dynamically, so it is not known a-priori which regression classes will be used to estimate the transform. This does not present a problem, since $\mathbf{G}^{(i)}$ and \mathbf{Z} for the chosen regression class may be obtained from its child classes (as defined by the tree). If the parent node R has children $\{R_1, \dots, R_C\}$ then

$$\mathbf{Z} = \sum_{c=1}^C \mathbf{Z}^{(R_c)}$$

and

$$\mathbf{G}^{(i)} = \sum_{c=1}^C \mathbf{G}^{(iR_c)}$$

From this it is clear that it is only necessary to calculate $\mathbf{G}^{(i)}$ and \mathbf{Z} for only the most specific regression classes possible – i.e. the base classes.

9.4.2 Estimation of the Variance Transformation Matrix

Estimation of the variance transformation matrices is only available for diagonal covariance Gaussian systems. The Gaussian covariance is transformed using,

$$\hat{\Sigma}_m = \mathbf{B}_m^T \mathbf{H}_m \mathbf{B}_m$$

where \mathbf{H}_m is the linear transformation to be estimated and \mathbf{B}_m is the inverse of the Choleski factor of Σ_m^{-1} , so

$$\Sigma_m^{-1} = \mathbf{C}_m \mathbf{C}_m^T$$

and

$$\mathbf{B}_m = \mathbf{C}_m^{-1}$$

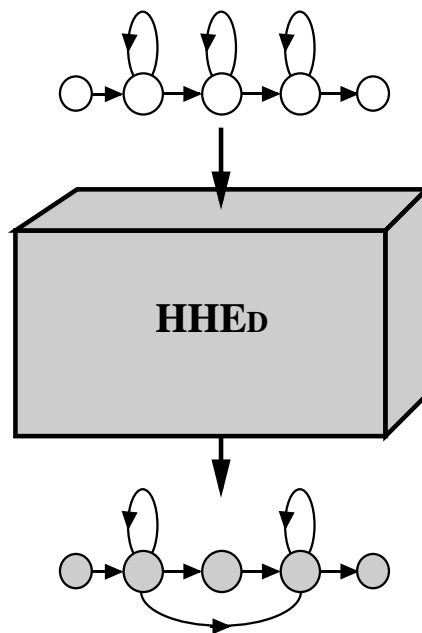
After rewriting the auxiliary function, the transform matrix \mathbf{H}_m is estimated from,

$$\mathbf{H}_m = \frac{\sum_{r=1}^{R_c} \mathbf{C}_{m_r}^T [L_{m_r}(t)(\mathbf{o}(t) - \boldsymbol{\mu}_{m_r})(\mathbf{o}(t) - \boldsymbol{\mu}_{m_r})^T] \mathbf{C}_{m_r}}{L_{m_r}(t)}$$

Here, \mathbf{H}_m is forced to be a diagonal transformation by setting the off-diagonal terms to zero, which ensures that $\hat{\Sigma}_m$ is also diagonal.

Chapter 10

HMM System Refinement



In chapter 8, the basic processes involved in training a continuous density HMM system were explained and examples were given of building a set of HMM phone models. In the practical application of these techniques to building real systems, there are often a number of problems to overcome. Most of these arise from the conflicting desire to have a large number of model parameters in order to achieve high accuracy, whilst at the same time having limited and uneven training data.

As mentioned previously, the HTK philosophy is to build systems incrementally. Starting with a set of context-independent monophone HMMs, a system can be refined in a sequence of stages. Each refinement step typically uses the HTK HMM definition editor HHED followed by re-estimation using HEREST. These incremental manipulations of the HMM set often involve parameter tying, thus many of HHED's operations involve generating new macro definitions.

The principle types of manipulation that can be performed by HHED are

- HMM cloning to form context-dependent model sets
- Generalised parameter tying
- Data driven and decision tree based clustering.
- Mixture component splitting
- Adding/removing state transitions
- Stream splitting, resizing and recasting

This chapter describes how the HTK tool HHED is used, its editing language and the main operations that can be performed.

10.1 Using HHed

The HMM editor HHED takes as input a set of HMM definitions and outputs a new modified set, usually to a new directory. It is invoked by a command line of the form

```
HHed -H MMF1 -H MMF2 ... -M newdir cmds.hed hmmlist
```

where `cmds.hed` is an edit script containing a list of edit commands. Each command is written on a separate line and begins with a 2 letter command name.

The effect of executing the above command line would be to read in the HMMs listed in `hmmlist` and defined by files `MMF1`, `MMF2`, etc., apply the editing operations defined in `cmds.hed` and then write the resulting system out to the directory `newdir`. As with all tools, HTK will attempt to replicate the file structure of the input in the output directory. By default, any new macros generated by HHED will be written to one or more of the existing MMFs. In doing this, HTK will attempt to ensure that the “definition before use” rule for macros is preserved, but it cannot always guarantee this. Hence, it is usually best to define explicit target file names for new macros. This can be done in two ways. Firstly, explicit target file names can be given in the edit script using the `UF` command. For example, if `cmds.hed` contained

```
....
UF smacs
# commands to generate state macros
....
UF vmacs
# commands to generate variance macros
....
```

then the output directory would contain an MMF called `smacs` containing a set of state macro definitions and an MMF called `vmacs` containing a set of variance macro definitions, these would be in addition to the existing MMF files `MMF1`, `MMF2`, etc.

Alternatively, the whole HMM system can be written to a single file using the `-w` option. For example,

```
HHed -H MMF1 -H MMF2 ... -w newMMF cmds.hed hmmlist
```

would write the whole of the edited HMM set to the file `newMMF`.

As mentioned previously, each execution of HHED is normally followed by re-estimation using HEREST. Normally, all the information needed by HHED is contained in the model set itself. However, some clustering operations require various statistics about the training data (see sections 10.4 and 10.5). These statistics are gathered by HEREST and output to a *stats file*, which is then read in by HHED. Note, however, that the statistics file generated by HEREST refers to the input model set not the re-estimated set. Thus for example, in the following sequence, the HHED edit script in `cmds.hed` contains a command (see the `R0` command in section 10.4) which references a statistics file (called `stats`) describing the HMM set defined by `hmm1/MMF`.

```
HERest -H hmm1/MMF -M hmmx -s stats hmmlist train1 train2 ....
HHed -H hmm1/MMF -M hmm2 cmds.hed hmmlist
```

The required statistics file is generated by HEREST but the re-estimated model set stored in `hmmx/MMF` is ignored and can be deleted.

10.2 Constructing Context-Dependent Models

The first stage of model refinement is usually to convert a set of initialised and trained context-independent monophone HMMs to a set of context dependent models. As explained in section 6.4, HTK uses the convention that a HMM name of the form `l-p+r` denotes the context-dependent version of the phone `p` which is to be used when the left neighbour is the phone `l` and the right neighbour is the phone `r`. To make a set of context dependent phone models, it is only necessary to construct a HMM list, called say `cdlist`, containing the required context-dependent models and then execute HHED with a single command in its edit script

```
CL cdlist
```

The effect of this command is that for each model `l-p+r` in `cdlist` it makes a copy of the monophone `p`.

The set of context-dependent models output by the above must be reestimated using HEREST. To do this, the training data transcriptions must be converted to use context-dependent labels and the original monophone hmm list must be replaced by `cdlist`. In fact, it is best to do this conversion before cloning the monophones because if the `HLED TC` command is used then the `-n` option can be used to generate the required list of context dependent HMMs automatically.

Before building a set of context-dependent models, it is necessary to decide whether or not cross-word triphones are to be used. If they are, then word boundaries in the training data can be ignored and all monophone labels can be converted to triphones. If, however, word internal triphones are to be used, then word boundaries in the training transcriptions must be marked in some way (either by an explicit marker which is subsequently deleted or by using a short pause *tee-model*). This word boundary marker is then identified to `HLED` using the `WB` command to make the `TC` command use biphones rather than triphones at word boundaries (see section 6.4).

All HTK tools can read and write HMM definitions in text or binary form. Text is good for seeing exactly what the tools are producing, but binary is much faster to load and store, and much more compact. Binary output is enabled either using the standard option `-B` or by setting the configuration variable `SAVEBINARY`. In the above example, the HMM set input to `HHED` will contain a small set of monophones whereas the output will be a large set of triphones. In order, to save storage and computation, this is usually a good point to switch to binary storage of MMFs.

10.3 Parameter Tying and Item Lists

As explained in Chapter 7, HTK uses macros to support a generalised parameter tying facility. Referring again to Fig. 7.7.8, each of the solid black circles denotes a potential *tie-point* in the hierarchy of HMM parameters. When two or more parameter sets are tied, the same set of parameter values are shared by all the *owners* of the tied set. Externally, tied parameters are represented by macros and internally they are represented by structure sharing. The accumulators needed for the numerators and denominators of the Baum-Welch re-estimation formulae given in section 8.7 are attached directly to the parameters themselves. Hence, when the values of a tied parameter set are re-estimated, all of the data which would have been used to estimate each individual untied parameter are effectively pooled leading to more robust parameter estimation.

Note also that although parameter tying is implemented in a way which makes it transparent to the HTK re-estimation and recognition tools, in practice, these tools do notice when a system has been tied and try to take advantage of it by avoiding redundant computations.

Although macro definitions could be written by hand, in practice, tying is performed by executing `HHED` commands and the resulting macros are thus generated automatically. The basic `HHED` command for tying a set of parameters is the `TI` command which has the form

```
TI macroname itemlist
```

This causes all items in the given `itemlist` to be tied together and output as a macro called `macroname`. Macro names are written as a string of characters optionally enclosed in double quotes. The latter are necessary if the name contains one or more characters which are not letters or digits.

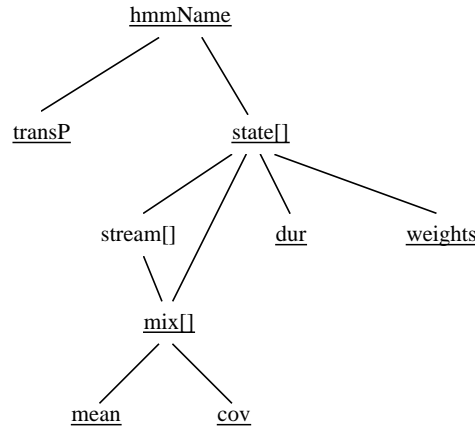


Fig. 10.1 Item List Construction

Item lists use a simple language to identify sets of points in the HMM parameter hierarchy illustrated in Fig. 7.7.8. This language is defined fully in the reference entry for HHED. The essential idea is that item lists represent paths down the hierarchical parameter tree where the direction *down* should be regarded as travelling from the *root* of the tree to towards the *leaves*. A path can be unique, or more usually, it can be a pattern representing a set of paths down the tree. The point at which each path stops identifies one member of the set represented by the item list. Fig. 10.1 shows the possible paths down the tree. In text form the branches are replaced by dots and the underlined node names are possible terminating points. At the topmost level, an item list is a comma separated list of paths enclosed in braces.

Some examples, should make all this clearer. Firstly, the following is a legal but somewhat long-winded way of specifying the set of items comprising states 2, 3 and 4 of the HMM called **aa**

```
{ aa.state[2],aa.state[3],aa.state[4] }
```

however in practice this would be written much more compactly as

```
{ aa.state[2-4] }
```

It must be emphasised that indices in item lists are really *patterns*. The set represented by an item list consists of all those elements which match the patterns. Thus, if **aa** only had two emitting states, the above item list would not generate an error. It would simply only match two items. The reason for this is that the same pattern can be applied to many different objects. For example, the HMM name can be replaced by a list of names enclosed in brackets, furthermore each HMM name can include ‘?’ characters which match any single character and ‘*’ characters which match zero or more characters. Thus

```
{ (aa+*,iy+*,eh+*).state[2-4] }
```

represents states 2, 3 and 4 of all biphone models corresponding to the phonemes **aa**, **iy** and **eh**. If **aa** had just 2 emitting states and the others had 4 emitting states, then this item list would include 2 states from each of the **aa** models and 3 states from each of the others. Moving further down the tree, the item list

```
{ *.state[2-4].stream[1].mix[1,3].cov }
```

denotes the set of all covariance vectors (or matrices) of the first and third mixture components of stream 1, of states 2 to 4 of all HMMs. Since many HMM systems are single stream, the **stream** part of the path can be omitted if its value is 1. Thus, the above could have been written

```
{ *.state[2-4].mix[1,3].cov }
```

These last two examples also show that indices can be written as comma separated lists as well as ranges, for example, [1,3,4-6,9] is a valid index list representing states 1, 3, 4, 5, 6, and 9.

When item lists are used as the argument to a TI command, the kind of items represented by the list determines the macro type in a fairly obvious way. The only non-obvious cases are firstly that lists ending in `cov` generate $\sim v$, $\sim i$, $\sim c$, or $\sim x$ macros as appropriate. If an explicit set of mixture components is defined as in

```
{ *.state[2].mix[1-5] }
```

then $\sim m$ macros are generated but omitting the indices altogether denotes a special case of mixture tying which is explained later in Chapter 11.

To illustrate the use of item lists, some example TI commands can now be given. Firstly, when a set of context-dependent models is created, it can be beneficial to share one transition matrix across all variants of a phone rather than having a distinct transition matrix for each. This could be achieved by adding TI commands immediately after the CL command described in the previous section, that is

```
CL cdlist
TI T_ah {*-ah+.transP}
TI T_eh {*-eh+.transP}
TI T_ae {*-ae+.transP}
TI T_ih {*-ih+.transP}
... etc
```

As a second example, a so-called Grand Variance HMM system can be generated very easily with the following HHed command

```
TI "gvar" { *.state[2-4].mix[1].cov }
```

where it is assumed that the HMMs are 3-state single mixture component models. The effect of this command is to tie all state distributions to a single global variance vector. For applications, where there is limited training data, this technique can improve performance, particularly in noise.

Speech recognition systems will often have distinct models for silence and short pauses. A silence model `sil` may have the normal 3 state topology whereas a short pause model may have just a single state. To avoid the two models *competing* with each other, the `sp` model state can be tied to the centre state of the `sil` model thus

```
TI "silst" { sp.state[2], sil.state[3] }
```

So far nothing has been said about how the parameters are actually determined when a set of items is replaced by a single shared representative. When states are tied, the state with the broadest variances and as few as possible zero mixture component weights is selected from the pool and used as the representative. When mean vectors are tied, the average of all the mean vectors in the pool is used and when variances are tied, the largest variance in the the pool is used. In all other cases, the last item in the tie-list is arbitrarily chosen as representative. All of these selection criteria are *ad hoc*, but since the tie operations are always followed by explicit re-estimation using HEREST, the precise choice of representative for a tied set is not critical.

Finally, tied parameters can be untied. For example, subsequent refinements of the context-dependent model set generated above with tied transition matrices might result in a much more compact set of models for which individual transition parameters could be robustly estimated. This can be done using the UT command whose effect is to untie all of the items in its argument list. For example, the command

```
UT {*-iy+.transP}
```

would untie the transition parameters in all variants of the `iy` phoneme. This untying works by simply making unique copies of the tied parameters. These untied parameters can then subsequently be re-estimated.

10.4 Data-Driven Clustering

In section 10.2, a method of triphone construction was described which involved cloning all monophones and then re-estimating them using data for which monophone labels have been replaced by triphone labels. This will lead to a very large set of models, and relatively little training data for

each model. Applying the argument that context will not greatly affect the centre states of triphone models, one way to reduce the total number of parameters without significantly altering the models' ability to represent the different contextual effects might be to tie all of the centre states across all models derived from the same monophone. This tying could be done by writing an edit script of the form

```
TI "iyS3" {*-iy+*.state[3]}
TI "ihS3" {*-ih+*.state[3]}
TI "ehS3" {*-eh+*.state[3]}
.... etc
```

Each TI command would tie all the centre states of all triphones in each phone group. Hence, if there were an average of 100 triphones per phone group then the total number of states per group would be reduced from 300 to 201.

Explicit tyings such as these can have some positive effect but overall they are not very satisfactory. Tying all centre states is too severe and worse still, the problem of undertraining for the left and right states remains. A much better approach is to use clustering to decide which states to tie. HHED provides two mechanisms for this. In this section a data-driven clustering approach will be described and in the next section, an alternative decision tree-based approach is presented.

Data-driven clustering is performed by the TC and NC commands. These both invoke the same top-down hierarchical procedure. Initially all states are placed in individual clusters. The pair of clusters which when combined would form the smallest resultant cluster are merged. This process repeats until either the size of the largest cluster reaches the threshold set by the TC command or the total number of clusters has fallen to that specified by the NC command. The size of cluster is defined as the greatest distance between any two states. The distance metric depends on the type of state distribution. For single Gaussians, a weighted Euclidean distance between the means is used and for tied-mixture systems a Euclidean distance between the mixture weights is used. For all other cases, the average probability of each component mean with respect to the other state is used. The details of the algorithm and these metrics are given in the reference section for HHED.

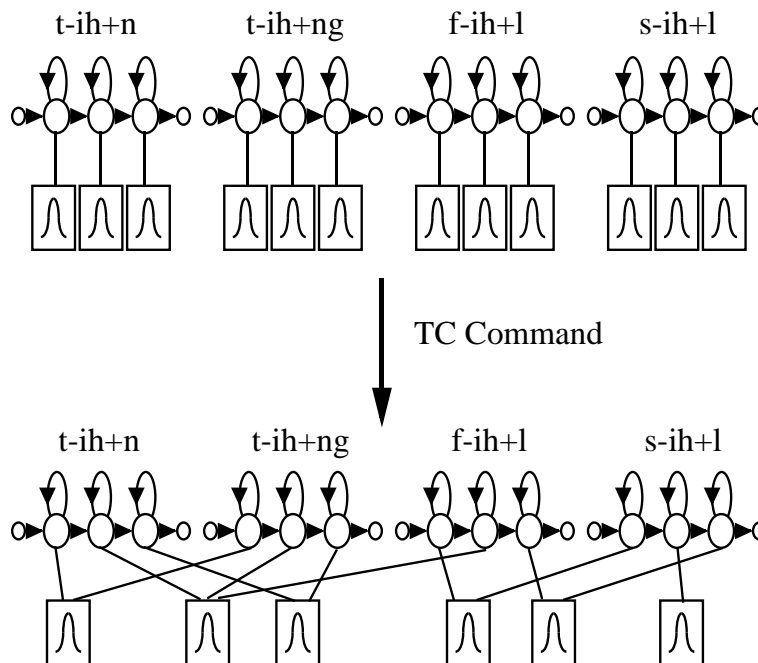


Fig. 10.2 Data-driven state tying

As an example, the following HHED script would cluster and tie the corresponding states of the triphone group for the phone ih

```
TC 100.0 "ihS2" {*-ih+*.state[2]}
```

```
TC 100.0 "ihS3" {*-ih+*.state[3]}
TC 100.0 "ihS4" {*-ih+*.state[4]}
```

In this example, each TC command performs clustering on the specified set of states, each cluster is tied and output as a macro. The macro name is generated by appending the cluster index to the macro name given in the command. The effect of this command is illustrated in Fig. 10.2. Note that if a word-internal triphone system is being built, it is sensible to include biphones as well as triphones in the item list, for example, the first command above would be written as

```
TC 100.0 "ihS2" {(*-ih,ih+*,*-ih+*).state[2]}
```

If the above TC commands are repeated for all phones, the resulting set of tied-state models will have far fewer parameters in total than the original untied set. The numeric argument immediately following the TC command name is the cluster threshold. Increasing this value will allow larger and hence, fewer clusters. The aim, of course, is to strike the right balance between compactness and the acoustic accuracy of the individual models. In practice, the use of this command requires some experimentation to find a good threshold value. HHED provides extensive trace output for monitoring clustering operations. Note in this respect that as well as setting tracing from the command line and the configuration file, tracing in HHED can be set by the TR command. Thus, tracing can be controlled at the command level. Further trace information can be obtained by including the SH command at strategic points in the edit script. The effect of executing this command is to list out all of the parameter tyings currently in force.

A potential problem with the use of the TC and NC commands is that *outlier* states will tend to form their own singleton clusters for which there is then insufficient data to properly train. One solution to this is to use the RO command to remove outliers. This command has the form

```
RO thresh "statsfile"
```

where *statsfile* is the name of a statistics file output using the `-s` option of HEREST. This statistics file holds the *occupation counts* for all states of the HMM set being trained. The term *occupation count* refers to the number of frames allocated to a particular state and can be used as a measure of how much training data is available for estimating the parameters of that state. The RO command must be executed *before* the TC or NC commands used to do the actual clustering. Its effect is to simply read in the statistics information from the given file and then to set a flag instructing the TC or NC commands to remove any outliers remaining at the conclusion of the normal clustering process. This is done by repeatedly finding the cluster with the smallest total occupation count and merging it with its nearest neighbour. This process is repeated until all clusters have a total occupation count which exceeds *thresh*, thereby ensuring that every cluster of states will be properly trained in the subsequent re-estimation performed by HEREST.

On completion of the above clustering and tying procedures, many of the models may be effectively identical, since acoustically similar triphones may share common clusters for all their emitting states. They are then, in effect, so-called *generalised triphones*. State tying can be further exploited if the HMMs which are effectively equivalent are identified and then tied via the physical-logical mapping¹ facility provided by HMM lists (see section 7.4). The effect of this would be to reduce the total number of HMM definitions required. HHED provides a compaction command to do all of this automatically. For example, the command

```
CO newList
```

will compact the currently loaded HMM set by identifying equivalent models and then tying them via the new HMM list output to the file *newList*. Note, however, that for two HMMs to be tied, they must be identical in all respects. This is one of the reasons why transition parameters are often tied across triphone groups otherwise HMMs with identical states would still be left distinct due to minor differences in their transition matrices.

10.5 Tree-Based Clustering

One limitation of the data-driven clustering procedure described above is that it does not deal with triphones for which there are no examples in the training data. When building word-internal triphone systems, this problem can often be avoided by careful design of the training database but when building large vocabulary cross-word triphone systems *unseen* triphones are unavoidable.

¹The physical HMM which corresponding to several logical HMMs will be arbitrarily named after one of them.

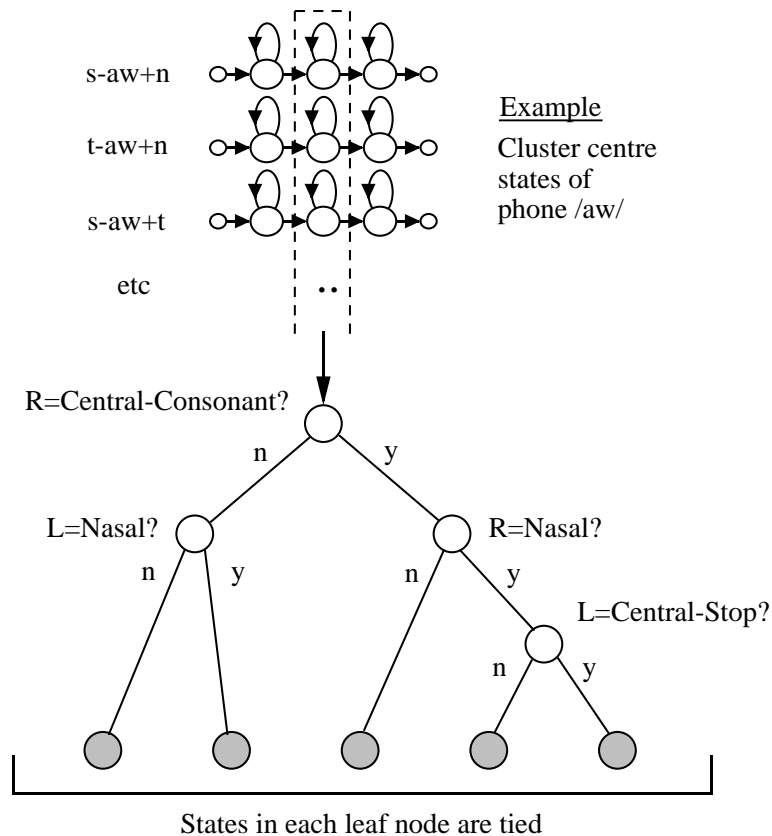


Fig. 10.3 Decision tree-based state tying

HHED provides an alternative decision tree based clustering mechanism which provides a similar quality of clustering but offers a solution to the unseen triphone problem. Decision tree-based clustering is invoked by the command TB which is analogous to the TC command described above and has an identical form, that is

```
TB thresh macroname itemlist
```

Apart from the clustering mechanism, there are some other differences between TC and TB. Firstly, TC uses a distance metric between states whereas TB uses a log likelihood criterion. Thus, the threshold values are not directly comparable. Furthermore, TC supports any type of output distribution whereas TB only supports single-Gaussian continuous density output distributions. Secondly, although the following describes only state clustering, the TB command can also be used to cluster whole models.

A phonetic decision tree is a binary tree in which a yes/no phonetic question is attached to each node. Initially all states in a given item list (typically a specific phone state position) are placed at the root node of a tree. Depending on each answer, the pool of states is successively split and this continues until the states have trickled down to leaf-nodes. All states in the same leaf node are then tied. For example, Fig 10.3 illustrates the case of tying the centre states of all triphones of the phone /aw/ (as in “out”). All of the states trickle down the tree and depending on the answer to the questions, they end up at one of the shaded terminal nodes. For example, in the illustrated case, the centre state of **s-aw+n** would join the second leaf node from the right since its right context is a central consonant, and its right context is a nasal but its left context is not a central stop.

The question at each node is chosen to (locally) maximise the likelihood of the training data given the final set of state tyings. Before any tree building can take place, all of the possible phonetic questions must be loaded into HHED using QS commands. Each question takes the form “Is the left or right context in the set P?” where the context is the model context as defined by its logical name. The set P is represented by an item list and for convenience every question is given a name. As an example, the following command

```
QS "L_Nasal" { ng-*,n-*,m-* }
```

defines the question “Is the left context a nasal?”.

It is possible to calculate the log likelihood of the training data given any pool of states (or models). Furthermore, this can be done without reference to the training data itself since for single Gaussian distributions the means, variances and state occupation counts (input via a stats file) form sufficient statistics. Splitting any pool into two will increase the log likelihood since it provides twice as many parameters to model the same amount of data. The increase obtained when each possible question is used can thus be calculated and the question selected which gives the biggest improvement.

Trees are therefore built using a top-down sequential optimisation process. Initially all states (or models) are placed in a single cluster at the root of the tree. The question is then found which gives the best split of the root node. This process is repeated until the increase in log likelihood falls below the threshold specified in the TB command. As a final stage, the decrease in log likelihood is calculated for merging terminal nodes with differing parents. Any pair of nodes for which this decrease is less than the threshold used to stop splitting are then merged.

As with the TC command, it is useful to prevent the creation of clusters with very little associated training data. The R0 command can therefore be used in tree clustering as well as in data-driven clustering. When used with trees, any split which would result in a total occupation count falling below the value specified is prohibited. Note that the R0 command can also be used to load the required stats file. Alternatively, the stats file can be loaded using the LS command.

As with data-driven clustering, using the trace facilities provided by HHED is recommended for monitoring and setting the appropriate thresholds. Basic tracing provides the following summary data for each tree

```
TB 350.00 aw_s3 {}
Tree based clustering
Start aw[3] : 28 have LogL=-86.899 occ=864.2
Via aw[3] : 5 gives LogL=-84.421 occ=864.2
End aw[3] : 5 gives LogL=-84.421 occ=864.2
TB: Stats 28->5 [17.9%] { 4537->285 [6.3%] total }
```

This example corresponds to the case illustrated in Fig 10.3. The TB command has been invoked with a threshold of 350.0 to cluster the centre states of the triphones of the phone *aw*. At the start of clustering with all 28 states in a single pool, the average log likelihood per unit of occupation is -86.9 and on completion with 5 clusters this has increased to -84.4. The middle line labelled “via” gives the position after the tree has been built but before terminal nodes have been merged (none were merged in this case). The last line summarises the overall position. After building this tree, a total of 4537 states were reduced to 285 clusters.

As noted at the start of this section, an important advantage of tree-based clustering is that it allows triphone models which have no training data to be synthesised. This is done in HHED using the AU command which has the form

```
AU hmmlist
```

Its effect is to scan the given `hmmlist` and any physical models listed which are not in the currently loaded set are synthesised. This is done by descending the previously constructed trees for that phone and answering the questions at each node based on the new unseen context. When each leaf node is reached, the state representing that cluster is used for the corresponding state in the unseen triphone.

The AU command can be used within the same edit script as the tree building commands. However, it will often be the case that a new set of triphones is needed at a later date, perhaps as a result of vocabulary changes. To make this possible, a complete set of trees can be saved using the ST command and then later reloaded using the LT command.

10.6 Mixture Incrementing

When building sub-word based continuous density systems, the final system will typically consist of multiple mixture component context-dependent HMMs. However, as indicated previously, the early stages of triphone construction, particularly state tying, are best done with single Gaussian models. Indeed, if tree-based clustering is to be used there is no option.

In HTK therefore, the conversion from single Gaussian HMMs to multiple mixture component HMMs is usually one of the final steps in building a system. The mechanism provided to do this is the HHED MU command which will increase the number of components in a mixture by a process called *mixture splitting*. This approach to building a multiple mixture component system is extremely flexible since it allows the number of mixture components to be repeatedly increased until the desired level of performance is achieved.

The MU command has the form

```
MU n itemList
```

where `n` gives the new number of mixture components required and `itemList` defines the actual mixture distributions to modify. This command works by repeatedly splitting the mixture with the largest mixture weight until the required number of components is obtained. The actual split is performed by copying the mixture, dividing the weights of both copies by 2, and finally perturbing the means by plus or minus 0.2 standard deviations. For example, the command

```
MU 3 {aa.state[2].mix}
```

would increase the number of mixture components in the output distribution for state 2 of model `aa` to 3. Normally, however, the number of components in all mixture distributions will be increased at the same time. Hence, a command of the form is more usual

```
MU 3 {*.state[2-4].mix}
```

It is usually a good idea to increment mixture components in stages, for example, by incrementing by 1 or 2 then re-estimating, then incrementing by 1 or 2 again and re-estimating, and so on until the required number of components are obtained. This also allows recognition performance to be monitored to find the optimum.

One final point with regard to multiple mixture component distributions is that all HTK tools ignore mixture components whose weights fall below a threshold value called MINMIX (defined in `HModel.h`). Such mixture components are called *defunct*. Defunct mixture components can be prevented by setting the `-w` option in `HEREST` so that all mixture weights are floored to some level above MINMIX. If mixture weights are allowed to fall below MINMIX then the corresponding Gaussian parameters will not be written out when the model containing that component is saved. It is possible to recover from this, however, since the MU command will replace defunct mixtures before performing any requested mixture component increment.

10.7 Regression Class Tree Construction

In order to perform most model adaptation tasks (see chapter 9), it will be necessary to produce a binary regression class tree. This tree is stored in the MMF, along with a regression base class identifier for each mixture component. An example regression tree and how it may be used is shown in subsection 9.1.2. HHED provides the means to construct a regression class tree for a given MMF, and is invoked using the RC command. It is also necessary to supply a statistics file, which is output using the `-s` option of `HERest`. The statistics file can be loaded by invoking the LS command.

A centroid-splitting algorithm using a Euclidean distance measure is used to grow the binary regression class tree to cluster the model set's mixture components. Each leaf node therefore specifies a particular mixture component cluster. This algorithm proceeds as follows until the requested number of terminals has been achieved.

- Select a terminal node that is to be split.
- Calculated the mean and variance from the mixture components clustered at this node.
- Create two children. Initialise their means to the parent mean perturbed in opposite directions (for each child) by a fraction of the variance.
- For each component at the parent node assign the component to one of the children by using a Euclidean distance measure to ascertain which child mean the component is closest to.
- Once all the components have been assigned, calculate the new means for the children, based on the component assignments.

- Keep re-assigning components to the children and re-estimating the child means until there is no change in assignments from one iteration to the next. Now finalise the split.

As an example, the following HHED script would produce a regression class tree with 32 terminal nodes, or regression base classes:-

```
LS "statsfile"  
RC 32 "rtree"
```

A further optional argument is possible with the RC command. This argument allows the user to specify the non-speech class mixture components using an `itemlist`, such as the silence mixture components.

```
LS "statsfile"  
RC 32 "rtree" sil.state[2-4].mix
```

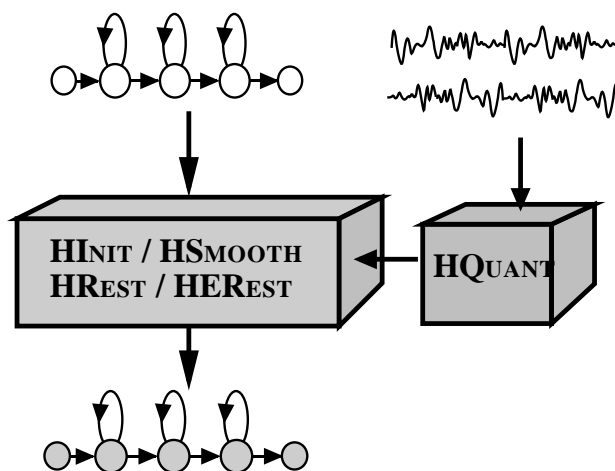
In this case the first split that will be made in the regression class tree will be to split the speech and non-speech sounds, after which the tree building continues as usual.

10.8 Miscellaneous Operations

The preceding sections have described the main HHED commands used for building continuous density systems with tied parameters. A further group of commands (`JO`, `TI` and `HK`) are used to build tied-mixture systems and these are described in Chapter 11. Those remaining cover a miscellany of functions. They are documented in the reference entry for HHED and include commands to add and remove state transitions (`AT`, `RT`); synthesise triphones from biphones (`MT`); change the parameter kind of a HMM (`SK`); modify stream dimensions (`SS`, `SU`, `SW`); change/add an identifier name to an MMF (`RN` command); and expand HMM sets by duplication, for example, as needed in making gender dependent models (`DP`).

Chapter 11

Discrete and Tied-Mixture Models



Most of the discussion so far has focussed on using HTK to model sequences of continuous-valued vectors. In contrast, this chapter is mainly concerned with using HTK to model sequences of discrete symbols. Discrete symbols arise naturally in modelling many types of data, for example, letters and words, bitmap images, and DNA sequences. Continuous signals can also be converted to discrete symbol sequences by using a quantiser and in particular, speech vectors can be *vector quantised* as described in section 5.11. In all cases, HTK expects a set of N discrete symbols to be represented by the contiguous sequence of integer numbers from 1 to N .

In HTK discrete probabilities are regarded as being closely analogous to the mixture weights of a continuous density system. As a consequence, the representation and processing of discrete HMMs shares a great deal with continuous density models. It follows from this that most of the principles and practice developed already are equally applicable to discrete systems. As a consequence, this chapter can be quite brief.

The first topic covered concerns building HMMs for discrete symbol sequences. The use of discrete HMMs with speech is then presented. The tool HQUANT is described and the method of converting continuous speech vectors to discrete symbols is reviewed. This is followed by a brief discussion of tied-mixture systems which can be regarded as a compromise between continuous and discrete density systems. Finally, the use of the HTK tool HSMOOTH for parameter smoothing by deleted interpolation is presented.

11.1 Modelling Discrete Sequences

Building HMMs for discrete symbol sequences is essentially the same as described previously for continuous density systems. Firstly, a prototype HMM definition must be specified in order to fix the model topology. For example, the following is a 3 state ergodic HMM in which the emitting states are fully connected.

```
~o <DISCRETE> <StreamInfo> 1 1
```

```

~h "dproto"
<BeginHMM>
  <NumStates> 5
  <State> 2 <NumMixes> 10
    <DProb> 5461*10
  <State> 3 <NumMixes> 10
    <DProb> 5461*10
  <State> 4 <NumMixes> 10
    <DProb> 5461*10
  <TransP> 5
    0.0 1.0 0.0 0.0 0.0
    0.0 0.3 0.3 0.3 0.1
    0.0 0.3 0.3 0.3 0.1
    0.0 0.3 0.3 0.3 0.1
    0.0 0.0 0.0 0.0 0.0
<EndHMM>

```

As described in chapter 7, the notation for discrete HMMs borrows heavily on that used for continuous density models by equating mixture components with symbol indices. Thus, this definition assumes that each training data sequence contains a single stream of symbols indexed from 1 to 10. In this example, all symbols in each state have been set to be equally likely¹. If prior information is available then this can of course be used to set these initial values.

The training data needed to build a discrete HMM can take one of two forms. It can either be discrete (SOURCEKIND=DISCRETE) in which case it consists of a sequence of 2-byte integer symbol indices. Alternatively, it can consist of continuous parameter vectors with an associated VQ codebook. This latter case is dealt with in the next section. Here it will be assumed that the data is symbolic and that it is therefore stored in discrete form. Given a set of training files listed in the script file `train.scp`, an initial HMM could be estimated using

```
HInit -T 1 -w 1.0 -o dhmm -S train.scp -M hmm0 dproto
```

This use of HINIT is identical to that which would be used for building whole word HMMs where no associated label file is assumed and the whole of each training sequence is used to estimate the HMM parameters. Its effect is to read in the prototype stored in the file `dproto` and then use the training examples to estimate initial values for the output distributions and transition probabilities. This is done by firstly uniformly segmenting the data and for each segment counting the number of occurrences of each symbol. These counts are then normalised to provide output distributions for each state. HINIT then uses the Viterbi algorithm to resegment the data and recompute the parameters. This is repeated until convergence is achieved or an upper limit on the iteration count is reached. The transition probabilities at each step are estimated simply by counting the number of times that each transition is made in the Viterbi alignments and normalising. The final model is renamed `dhmm` and stored in the directory `hmm0`.

When building discrete HMMs, it is important to floor the discrete probabilities so that no symbol has a zero probability. This is achieved using the `-w` option which specifies a floor value as a multiple of a global constant called MINMIX whose value is 10^{-5} .

The initialised HMM created by HINIT can then be further refined if desired by using HREST to perform Baum-Welch re-estimation. It would be invoked in a similar way to the above except that there is now no need to rename the model. For example,

```
HRest -T 1 -w 1.0 -S train.scp -M hmm1 hmm0/dhmm
```

would read in the model stored in `hmm0/dhmm` and write out a new model of the same name to the directory `hmm1`.

11.2 Using Discrete Models with Speech

As noted in section 5.11, discrete HMMs can be used to model speech by using a vector quantiser to map continuous density vectors into discrete symbols. A vector quantiser depends on a so-called

¹ Remember that discrete probabilities are scaled such that 32767 is equivalent to a probability of 0.000001 and 0 is equivalent to a probability of 1.0

codebook which defines a set of partitions of the vector space. Each partition is represented by the mean value of the speech vectors belonging to that partition and optionally a variance representing the spread. Each incoming speech vector is then matched with each partition and assigned the index corresponding to the partition which is closest using a Mahalanobis distance metric.

In HTK such a codebook can be built using the tool HQUANT. This tool takes as input a set of continuous speech vectors, clusters them and uses the centroid and optionally the variance of each cluster to define the partitions. HQUANT can build both linear and tree structured codebooks. To build a linear codebook, all training vectors are initially placed in one cluster and the mean calculated. The mean is then perturbed to give two means and the training vectors are partitioned according to which mean is nearest to them. The means are then recalculated and the data is repartitioned. At each cycle, the total distortion (i.e. total distance between the cluster members and the mean) is recorded and repartitioning continues until there is no significant reduction in distortion. The whole process then repeats by perturbing the mean of the cluster with the highest distortion. This continues until the required number of clusters have been found.

Since all training vectors are reallocated at every cycle, this is an expensive algorithm to compute. The maximum number of iterations within any single cluster increment can be limited using the configuration variable MAXCLUSTITER and although this can speed-up the computation significantly, the overall training process is still computationally expensive. Once built, vector quantisation is performed by scanning all codebook entries and finding the nearest entry. Thus, if a large codebook is used, the run-time VQ look-up operation can also be expensive.

As an alternative to building a linear codebook, a tree-structured codebook can be used. The algorithm for this is essentially the same as above except that every cluster is split at each stage so that the first cluster is split into two, they are split into four and so on. At each stage, the means are recorded so that when using the codebook for vector quantising a fast binary search can be used to find the appropriate leaf cluster. Tree-structured codebooks are much faster to build since there is no repeated reallocation of vectors and much faster in use since only $O(\log_2 N)$ distance need to be computed where N is the size of the codebook. Unfortunately, however, tree-structured codebooks will normally incur higher VQ distortion for a given codebook size.

When delta and acceleration coefficients are used, it is usually best to split the data into multiple streams (see section 5.10). In this case, a separate codebook is built for each stream.

As an example, the following invocation of HQUANT would generate a linear codebook in the file `linvq` using the data stored in the files listed in `vq.scp`.

```
HQuant -C config -s 4 -n 3 64 -n 4 16 -S vq.scp linvq
```

Here the configuration file `config` specifies the TARGETKIND as being MFCC_E_D_A i.e. static coefficients plus deltas plus accelerations plus energy. The `-s` options requests that this parameterisation be split into 4 separate streams. By default, each individual codebook has 256 entries, however, the `-n` option can be used to specify alternative sizes.

If a tree-structured codebook was wanted rather than a linear codebook, the `-t` option would be set. Also the default is to use Euclidean distances both for building the codebook and for subsequent coding. Setting the `-d` option causes a diagonal covariance Mahalanobis metric to be used and the `-f` option causes a full covariance Mahalanobis metric to be used.

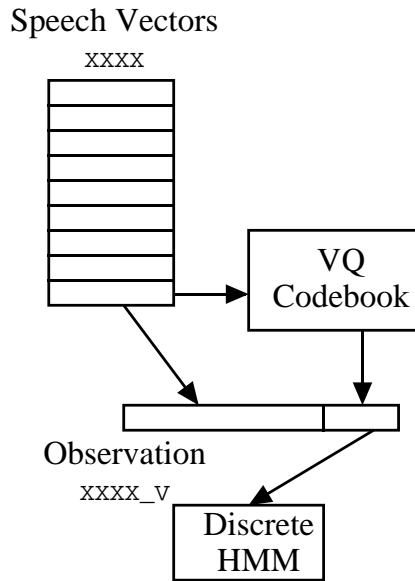


Fig. 11.1 VQ Processing

Once the codebook is built, normal speech vector files can be converted to discrete files using `HCOPY`. This was explained previously in section 5.11. The basic mechanism is to add the qualifier `_V` to the `TARGETKIND`. This causes `HPARM` to append a codebook index to each constructed observation vector. If the configuration variable `SAVEASVQ` is set true, then the output routines in `HPARM` will discard the original vectors and just save the VQ indices in a `DISCRETE` file. Alternatively, `HTK` will regard any speech vector with `_V` set as being compatible with discrete HMMs. Thus, it is not necessary to explicitly create a database of discrete training files if a set of continuous speech vector parameter files already exists. Fig. 11.1 illustrates this process.

Once the training data has been configured for discrete HMMs, the rest of the training process is similar to that previously described. The normal sequence is to build a set of monophone models and then clone them to make triphones. As in continuous density systems, state tying can be used to improve the robustness of the parameter estimates. However, in the case of discrete HMMs, alternative methods based on interpolation are possible. These are discussed in section 11.4.

11.3 Tied Mixture Systems

Discrete systems have the advantage of low run-time computation. However, vector quantisation reduces accuracy and this can lead to poor performance. As an intermediate between discrete and continuous, a fully tied-mixture system can be used. Tied-mixtures are conceptually just another example of the general parameter tying mechanism built-in to `HTK`. However, to use them effectively in speech recognition systems a number of storage and computational optimisations must be made. Hence, they are given special treatment in `HTK`.

When specific mixtures are tied as in

```
TI "mix" {*.state[2].mix[1]}
```

then a Gaussian mixture component is shared across all of the owners of the tie. In this example, all models will share the same Gaussian for the first mixture component of state 2. However, if the mixture component index is missing, then all of the mixture components participating in the tie are *joined* rather than tied. More specifically, the commands

```
JO 128 2.0
TI "mix" {*.state[2-4].mix}
```

has the following effect. All of the mixture components in states 2 to 4 of all models are collected into a pool. If the number of components in the pool exceeds 128, as set by the preceding join command

JO, then components with the smallest weights are removed until the pool size is exactly 128. Similarly, if the size of the initial pool is less than 128, then mixture components are split using the same algorithm as for the Mix-Up MU command. All states then share all of the mixture components in this pool. The new mixture weights are chosen to be proportional to the log probability of the corresponding new mixture component mean with respect to the original distribution for that state. The log is used here to give a wider spread of mixture weights. All mixture weights are floored to the value of the second argument of the JO command times MINMIX.

The net effect of the above two commands is to create a set of `tied-mixture` HMMs² where the same set of mixture components is shared across all states of all models. However, the type of the HMM set so created will still be `SHARED` and the internal representation will be the same as for any other set of parameter tyings. To obtain the optimised representation of the tied-mixture weights described in section 7.5, the following HHED HK command must be issued

```
HK TIEDHS
```

This will convert the internal representation to the special tied-mixture form in which all of the tied mixtures are stored in a global table and referenced implicitly instead of being referenced explicitly using pointers.

Tied-mixture HMMs work best if the information relating to different sources such as delta coefficients and energy are separated into distinct data streams. This can be done by setting up multiple data stream HMMs from the outset. However, it is simpler to use the SS command in HHED to split the data streams of the currently loaded HMM set. Thus, for example, the command

```
SS 4
```

would convert the currently loaded HMMs to use four separate data streams rather than one. When used in the construction of tied-mixture HMMs this is analogous to the use of multiple codebooks in discrete density HMMs.

The procedure for building a set of tied-mixture HMMs may be summarised as follows

1. Choose a *codebook* size for each data stream and then decide how many Gaussian components will be needed from an initial set of monophones to approximately fill this codebook. For example, suppose that there are 48 three state monophones. If codebook sizes of 128 are chosen for streams 1 and 2, and a codebook size of 64 is chosen for stream 3 then single Gaussian monophones would provide enough mixtures in total to fill the codebooks.
2. Train the initial set of monophones.
3. Use HHED to first split the HMMs into the required number of data streams, tie each individual stream and then convert the tied-mixture HMM set to have the kind `TIEDHS`. A typical script to do this for four streams would be

```
SS 4
JO 256 2.0
TI st1 {*.state[2-4].stream[1].mix}
JO 128 2.0
TI st2 {*.state[2-4].stream[2].mix}
JO 128 2.0
TI st3 {*.state[2-4].stream[3].mix}
JO 64 2.0
TI st4 {*.state[2-4].stream[4].mix}
HK TIEDHS
```

4. Re-estimate the models using HEREST in the normal way.

Once the set of retrained tied-mixture models has been produced, context dependent models can be constructed using similar methods to those outlined previously.

When evaluating probabilities in tied-mixture systems, it is often sufficient to sum just the most likely mixture components since for any particular input vector, its probability with respect to many of the Gaussian components will be very low. HTK tools recognise `TIEDHS` HMM sets as being special in the sense that additional optimisations are possible. When full tied-mixtures are

²Also called *semi-continuous* HMMs in the the literature.

used, then an additional layer of pruning is applied. At each time frame, the log probability of the current observation is computed for each mixture component. Then only those components which lie within a threshold of the most likely component are retained. This pruning is controlled by the `-c` option in HREST, HEREST and HVITE.

11.4 Parameter Smoothing

When large sets of context-dependent triphones are built using discrete models or tied-mixture models, under-training can be a severe problem since each state has a large number of mixture weight parameters to estimate. The HTK tool HSMOOTH allows these discrete probabilities or mixture component weights to be smoothed with the monophone weights using a technique called deleted interpolation.

HSMOOTH is used in combination with HEREST working in parallel mode. The training data is split into blocks and each block is used separately to re-estimate the HMMs. However, since HEREST is in parallel mode, it outputs a dump file of accumulators instead of updating the models. HSMOOTH is then used in place of the second pass of HEREST. It reads in the accumulator information from each of the blocks, performs deleted interpolation smoothing on the accumulator values and then outputs the re-estimated HMMs in the normal way.

HSMOOTH implements a conventional deleted interpolation scheme. However, optimisation of the smoothing weights uses a fast binary chop scheme rather than the more usual Baum-Welch approach. The algorithm for finding the optimal interpolation weights for a given state and stream is as follows where the description is given in terms of tied-mixture weights but the same applies to discrete probabilities.

Assume that HEREST has been set-up to output N separate blocks of accumulators. Let $w_i^{(n)}$ be the i 'th mixture weight based on the accumulator blocks 1 to N but excluding block n , and let $\bar{w}_i^{(n)}$ be the corresponding context independent weight. Let $x_i^{(n)}$ be the i 'th mixture weight count for the deleted block n . The derivative of the log likelihood of the deleted block, given the probability distribution with weights $c_i = \lambda w_i + (1 - \lambda)\bar{w}_i$ is given by

$$D(\lambda) = \sum_{n=1}^N \sum_{i=1}^M x_i^{(n)} \left[\frac{w_i^{(n)} - \bar{w}_i^{(n)}}{\lambda w_i^{(n)} + (1 - \lambda)\bar{w}_i^{(n)}} \right] \quad (11.1)$$

Since the log likelihood is a convex function of λ , this derivative allows the optimal value of λ to be found by a simple binary chop algorithm, viz.

```
function FindLambdaOpt:
  if (D(0) <= 0) return 0;
  if (D(1) >= 0) return 1;
  l=0; r=1;
  for (k=1; k<=maxStep; k++){
    m = (l+r)/2;
    if (D(m) == 0) return m;
    if (D(m) > 0) l=m; else r=m;
  }
  return m;
```

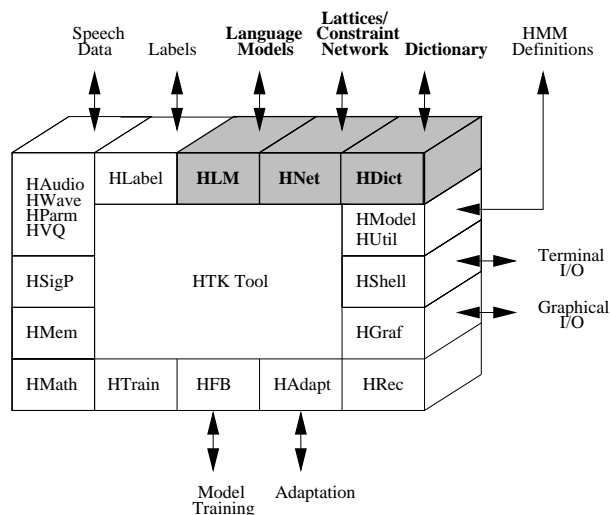
HSMOOTH is invoked in a similar way to HEREST. For example, suppose that the directory `hmm2` contains a set of accumulator files output by the first pass of HEREST running in parallel mode using as source the HMM definitions listed in `hlist` and stored in `hmm1/HMMDefs`. Then the command

```
HSmooth -c 4 -w 2.0 -H hmm1/HMMDefs -M hmm2 hlist hmm2/*.acc
```

would generate a new smoothed HMM set in `hmm2`. Here the `-w` option is used to set the minimum mixture component weight in any state to twice the value of `MINMIX`. The `-c` option sets the maximum number of iterations of the binary chop procedure to be 4.

Chapter 12

Networks, Dictionaries and Language Models



The preceding chapters have described how to process speech data and how to train various types of HMM. This and the following chapter are concerned with building a speech recogniser using HTK. This chapter focuses on the use of networks and dictionaries. A network describes the sequence of words that can be recognised and, for the case of sub-word systems, a dictionary describes the sequence of HMMs that constitute each word. A word level network will typically represent either a *Task Grammar* which defines all of the legal word sequences explicitly or a *Word Loop* which simply puts all words of the vocabulary in a loop and therefore allows any word to follow any other word. Word-loop networks are often augmented by a stochastic language model. Networks can also be used to define phone recognisers and various types of word-spotting systems.

Networks are specified using the HTK *Standard Lattice Format* (SLF) which is described in detail in Chapter 17. This is a general purpose text format which is used for representing multiple hypotheses in a recogniser output as well as word networks. Since SLF format is text-based, it can be written directly using any text editor. However, this can be rather tedious and HTK provides two tools which allow the application designer to use a higher-level representation. Firstly, the tool HPARSE allows networks to be generated from a source text containing extended BNF format grammar rules. This format was the only grammar definition language provided in earlier versions of HTK and hence HPARSE also provides backwards compatibility.

HPARSE task grammars are very easy to write, but they do not allow fine control over the actual network used by the recogniser. The tool HBUILD works directly at the SLF level to provide this detailed control. Its main function is to enable a large word network to be decomposed into a set of small self-contained sub-networks using as input an extended SLF format. This enhances the design process and avoids the need for unnecessary repetition.

HBUILD can also be used to perform a number of special-purpose functions. Firstly, it can construct word-loop and word-pair grammars automatically. Secondly, it can incorporate a statistical bigram language model into a network. These can be generated from label transcriptions using HLSTATS. However, HTK supports the standard ARPA MIT-LL text format for backed-off N-gram language models, and hence, import from other sources is possible.

Whichever tool is used to generate a word network, it is important to ensure that the generated network represents the intended grammar. It is also helpful to have some measure of the difficulty of the recognition task. To assist with this, the tool HSGEN is provided. This tool will generate example word sequences from an SLF network using random sampling. It will also estimate the perplexity of the network.

When a word network is loaded into a recogniser, a dictionary is consulted to convert each word in the network into a sequence of phone HMMs. The dictionary can have multiple pronunciations in which case several sequences may be joined in parallel to make a word. Options exist in this process to automatically convert the dictionary entries to context-dependent triphone models, either within a word or cross-word. Pronouncing dictionaries are a vital resource in building speech recognition systems and, in practice, word pronunciations can be derived from many different sources. The HTK tool HDMAN enables a dictionary to be constructed automatically from different sources. Each source can be individually edited and translated and merged to form a uniform HTK format dictionary.

The various facilities for describing a word network and expanding into a HMM level network suitable for building a recogniser are implemented by the HTK library module HNET. The facilities for loading and manipulating dictionaries are implemented by the HTK library module HDICT and for loading and manipulating language models are implemented by HLM. These facilities and those provided by HPARSE, HBUILD, HSGEN, HLSTATS and HDMAN are the subject of this chapter.

12.1 How Networks are Used

Before delving into the details of word networks and dictionaries, it will be helpful to understand their rôle in building a speech recogniser using HTK. Fig 12.1 illustrates the overall recognition process. A word network is defined using HTK Standard Lattice Format (SLF). An SLF word network is just a text file and it can be written directly with a text editor or a tool can be used to build it. HTK provides two such tools, HBUILD and HPARSE. These both take as input a textual description and output an SLF file. Whatever method is chosen, word network SLF generation is done *off-line* and is part of the system build process.

An SLF file contains a list of nodes representing words and a list of arcs representing the transitions between words. The transitions can have probabilities attached to them and these can be used to indicate *preferences* in a grammar network. They can also be used to represent bigram probabilities in a back-off bigram network and HBUILD can generate such a bigram network automatically. In addition to an SLF file, a HTK recogniser requires a dictionary to supply pronunciations for each word in the network and a set of acoustic HMM phone models. Dictionaries are input via the HTK interface module HDICT.

The dictionary, HMM set and word network are input to the HTK library module HNET whose function is to generate an equivalent network of HMMs. Each word in the dictionary may have several pronunciations and in this case there will be one branch in the network corresponding to each alternative pronunciation. Each pronunciation may consist either of a list of phones or a list of HMM names. In the former case, HNET can optionally expand the HMM network to use either word internal triphones or cross-word triphones. Once the HMM network has been constructed, it can be input to the decoder module HREC and used to recognise speech input. Note that HMM network construction is performed *on-line* at recognition time as part of the initialisation process.

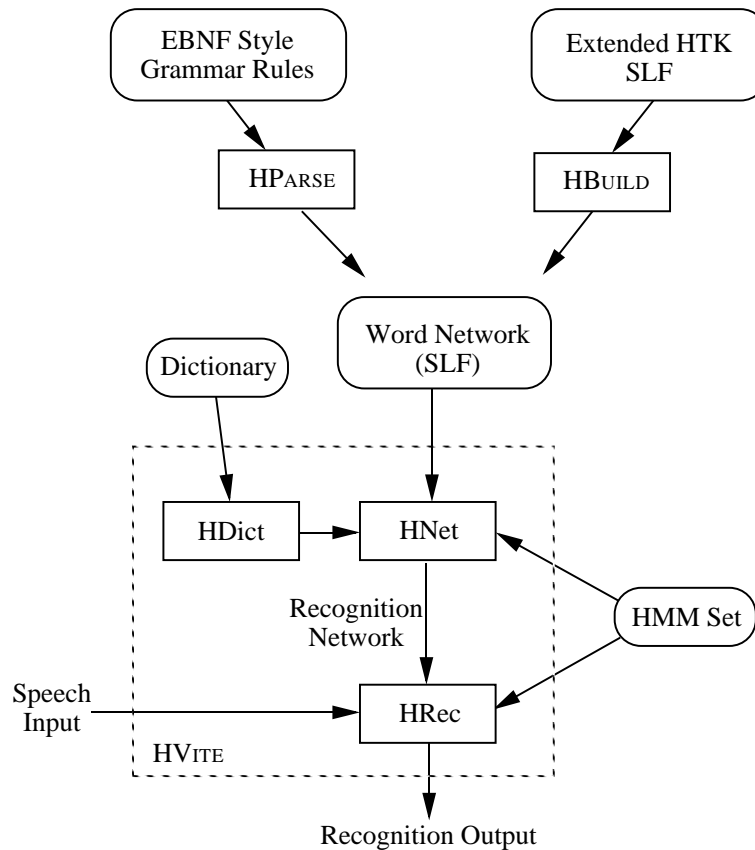


Fig. 12.1 Overview of the Recognition Process

For convenience, HTK provides a recognition tool called HVITE to allow the functions provided by HNET and HREC to be invoked from the command line. HVITE is particularly useful for running experimental evaluations on test speech stored in disk files and for basic testing using live audio input. However, application developers should note that HVITE is just a shell containing calls to load the word network, dictionary and models; generate the recognition network and then repeatedly recognise each input utterance. For embedded applications, it may well be appropriate to dispense with HVITE and call the functions in HNET and HREC directly from the application. The use of HVITE is explained in the next chapter.

12.2 Word Networks and Standard Lattice Format

This section provides a basic introduction to the HTK Standard Lattice Format (SLF). SLF files are used for a variety of functions some of which lie beyond the scope of the standard HTK package. The description here is limited to those features of SLF which are required to describe word networks suitable for input to HNET. The following Chapter describes the further features of SLF used for representing the output of a recogniser. For reference, a full description of SLF is given in Chapter 17.

A word network in SLF consists of a list of nodes and a list of arcs. The nodes represent words and the arcs represent the transition between words¹. Each node and arc definition is written on a single line and consists of a number of fields. Each field specification consists of a “name=value” pair. Field names can be any length but all commonly used field names consist of a single letter. By convention, field names starting with a capital letter are mandatory whereas field names starting with a lower-case letter are optional. Any line beginning with a # is a comment and is ignored.

¹More precisely, nodes represent the ends of words and arcs represent the transitions between word ends. This distinction becomes important when describing recognition output since acoustic scores are attached to arcs not nodes.

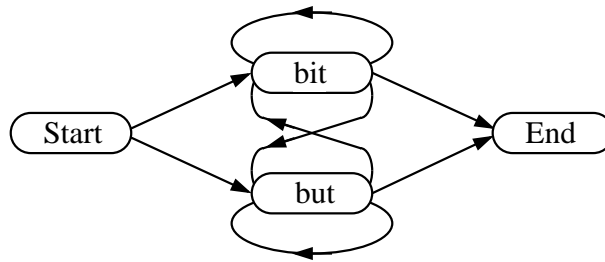


Fig. 12.2 A Simple Word Network

The following example should illustrate the basic format of an SLF word network file. It corresponds to the network illustrated in Fig 12.2 which represents all sequences consisting of the words “bit” and “but” starting with the word “start” and ending with the word “end”. As will be seen later, the start and end words will be mapped to a silence model so this grammar allows speakers to say “bit but but bit bit ...etc”.

```
# Define size of network: N=num nodes and L=num arcs
N=4 L=8
# List nodes: I=node-number, W=word
I=0 W=start
I=1 W=end
I=2 W=bit
I=3 W=but
# List arcs: J=arc-number, S=start-node, E=end-node
J=0 S=0 E=2
J=1 S=0 E=3
J=2 S=3 E=1
J=3 S=2 E=1
J=4 S=2 E=3
J=5 S=3 E=3
J=6 S=3 E=2
J=7 S=2 E=2
```

Notice that the first line which defines the size of the network must be given before any node or arc definitions. A node is a *network start node* if it has no predecessors, and a node is *network end node* if it has no successors. There must be one and only one network start node and one network end node. In the above, node 0 is a network start node and node 1 is a network end node. The choice of the names “start” and “end” for these nodes has no significance.

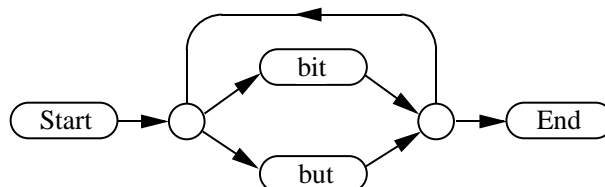


Fig. 12.3 A Word Network Using Null Nodes

A word network can have null nodes indicated by the special predefined word name !NULL. Null nodes are useful for reducing the number of arcs required. For example, the *Bit-But* network could be defined as follows

```
# Network using null nodes
```

```

N=6 L=7
I=0 W=start
I=1 W=end
I=2 W=bit
I=3 W=but
I=4 W=!NULL
I=5 W=!NULL
J=0 S=0 E=4
J=1 S=4 E=2
J=2 S=4 E=3
J=3 S=2 E=5
J=4 S=3 E=5
J=5 S=5 E=4
J=6 S=5 E=1

```

In this case, there is no significant saving, however, if there were many words in parallel, the total number of arcs would be much reduced by using null nodes to form common start and end points for the loop-back connections.

By default, all arcs are equally likely. However, the optional field `l=x` can be used to attach the log transition probability `x` to an arc. For example, if the word “but” was twice as likely as “bit”, the arcs numbered 1 and 2 in the last example could be changed to

```

J=1 S=4 E=2 l=-1.1
J=2 S=4 E=3 l=-0.4

```

Here the probabilities have been normalised to sum to 1, however, this is not necessary. The recogniser simply adds the scaled log probability to the path score and hence it can be regarded as an additive word transition penalty.

12.3 Building a Word Network with HParse

Whilst the construction of a word level SLF network file by hand is not difficult, it can be somewhat tedious. In earlier versions of HTK, a high level grammar notation based on extended Backus-Naur Form (EBNF) was used to specify recognition grammars. This *HParse* format was read-in directly by the recogniser and compiled into a finite state recognition network at run-time.

In HTK 3.1, *HParse* format is still supported but in the form of an *off-line* compilation into an SLF word network which can subsequently be used to drive a recogniser.

A HParse format grammar consists of an extended form of regular expression enclosed within parentheses. Expressions are constructed from sequences of words and the metacharacters

```

| denotes alternatives
[ ] encloses options
{ } denotes zero or more repetitions
< > denotes one or more repetitions
<< >> denotes context-sensitive loop

```

The following examples will illustrate the use of all of these except the last which is a special-purpose facility provided for constructing context-sensitive loops as found in for example, context-dependent phone loops and word-pair grammars. It is described in the reference entry for HPARSE.

As a first example, suppose that a simple isolated word single digit recogniser was required. A suitable syntax would be

```

(
  one | two | three | four | five |
  six | seven | eight | nine | zero
)

```

This would translate into the network shown in part (a) of Fig. 12.4. If this HParse format syntax definition was stored in a file called `digitsyn`, the equivalent SLF word network would be generated in the file `digitnet` by typing

```
HParse digitsyn digitnet
```

The above digit syntax assumes that each input digit is properly end-pointed. This requirement can be removed by adding a silence model before and after the digit

```
(
  sil (one | two | three | four | five |
      six | seven | eight | nine | zero) sil
)
```

As shown by graph (b) in Fig. 12.4, the allowable sequence of models now consists of silence followed by a digit followed by silence. If a sequence of digits needed to be recognised then angle brackets can be used to indicate one or more repetitions, the HParse grammar

```
(
  sil < one | two | three | four | five |
      six | seven | eight | nine | zero > sil
)
```

would accomplish this. Part (c) of Fig. 12.4 shows the network that would result in this case.

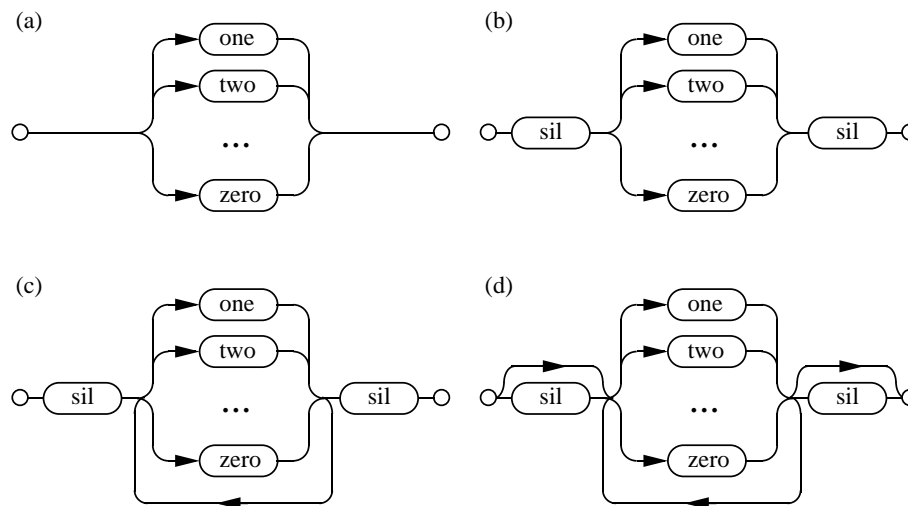


Fig. 12.4 Example Digit Recognition Networks

HParse grammars can define variables to represent sub-expressions. Variable names start with a dollar symbol and they are given values by definitions of the form

```
$var = expression ;
```

For example, the above connected digit grammar could be rewritten as

```
$digit = one | two | three | four | five |
        six | seven | eight | nine | zero;
(
  sil < $digit > sil
)
```

Here `$digit` is a variable whose value is the expression appearing on the right hand side of the assignment. Whenever the name of a variable appears within an expression, the corresponding expression is substituted. Note however that variables must be defined before use, hence, recursion is prohibited.

As a final refinement of the digit grammar, the start and end silence can be made optional by enclosing them within square brackets thus

```

$digit = one | two | three | four | five |
        six | seven | eight | nine | zero;
(
  [sil] < $digit > [sil]
)

```

Part (d) of Fig. 12.4 shows the network that would result in this last case.

HParse format grammars are a convenient way of specifying task grammars for interactive voice interfaces. As a final example, the following defines a simple grammar for the control of a telephone by voice.

```

$digit = one | two | three | four | five |
        six | seven | eight | nine | zero;
$number = $digit { [pause] $digit};
$score = shortcode $digit $digit;
$telnum = $score | $number;
$cmd = dial $telnum |
       enter $score for $number |
       redial | cancel;
$noise = lipsmack | breath | background;
( < $cmd | $noise > )

```

The dictionary entries for `pause`, `lipsmack`, `breath` and `background` would reference HMMs trained to model these types of noise and the corresponding output symbols in the dictionary would be null.

Finally, it should be noted that when the HParse format was used in earlier versions of HTK, word grammars contained word pronunciations embedded within them. This was done by using the reserved node names `WD_BEGIN` and `WD_END` to delimit word boundaries. To provide backwards compatibility, HPARSE can process these old format networks but when doing so it outputs a dictionary as well as a word network. This compatibility mode is defined fully in the reference section, to use it the configuration variable `V1COMPAT` must be set true or the `-c` option set.

Finally on the topic of word networks, it is important to note that any network containing an unbroken loop of one or more tee-models will generate an error. For example, if `sp` is a single state tee-model used to represent short pauses, then the following network would generate an error

```
( sil < sp | $digit > sil )
```

the intention here is to recognise a sequence of digits which may optionally be separated by short pauses. However, the syntax allows an endless sequence of `sp` models and hence, the recogniser could traverse this sequence without ever consuming any input. The solution to problems such as these is to rearrange the network. For example, the above could be written as

```
( sil < $digit sp > sil )
```

12.4 Bigram Language Models

Before continuing with the description of network generation and, in particular, the use of `HBUILD`, the use of bigram language models needs to be described. Support for statistical language models in HTK is provided by the library module `HLM`. Although the interface to `HLM` can support general N-grams, the facilities for constructing and using N-grams are limited to bigrams.

A bigram language model can be built using `HLSTATS` invoked as follows where it is assumed that all of the label files used for training are stored in an MLF called `labs`

```
HLStats -b bigfn -o wordlist labs
```

All words used in the label files must be listed in the `wordlist`. This command will read all of the transcriptions in `labs`, build a table of bigram counts in memory, and then output a back-off bigram to the file `bigfn`. The formulae used for this are given in the reference entry for `HLSTATS`. However, the basic idea is encapsulated in the following formula

$$p(i, j) = \begin{cases} (N(i, j) - D)/N(i) & \text{if } N(i, j) > t \\ b(i)p(j) & \text{otherwise} \end{cases}$$

where $N(i, j)$ is the number of times word j follows word i and $N(i)$ is the number of times that word i appears. Essentially, a small part of the available probability mass is deducted from the higher bigram counts and distributed amongst the infrequent bigrams. This process is called *discounting*. The default value for the discount constant D is 0.5 but this can be altered using the configuration variable DISCOUNT. When a bigram count falls below the threshold t , the bigram is backed-off to the unigram probability suitably scaled by a back-off weight in order to ensure that all bigram probabilities for a given history sum to one.

Backed-off bigrams are stored in a text file using the standard ARPA MIT-LL format which as used in HTK is as follows

```
\data\  
ngram 1=<num 1-grams>  
ngram 2=<num 2-ngrams>  
  
\1-grams:  
P(!ENTER)      !ENTER  B(!ENTER)  
P(W1)          W1      B(W1)  
P(W2)          W2      B(W2)  
...  
P(!EXIT)       !EXIT   B(!EXIT)  
  
\2-grams:  
P(W1 | !ENTER) !ENTER  W1  
P(W2 | !ENTER) !ENTER  W2  
P(W1 | W1)     W1      W1  
P(W2 | W1)     W1      W2  
P(W1 | W2)     W2      W1  
...  
P(!EXIT | W1)  W1      !EXIT  
P(!EXIT | W2)  W2      !EXIT  
  
\end\  

```

where all probabilities are stored as base-10 logs. The default start and end words, !ENTER and !EXIT can be changed using the HLSTATS -s option.

For some applications, a simple matrix style of bigram representation may be more appropriate. If the -o option is omitted in the above invocation of HLSTATS, then a simple full bigram matrix will be output using the format

```
!ENTER    0  P(W1 | !ENTER) P(W2 | !ENTER) .....  
W1        0  P(W1 | W1)      P(W2 | W1)      .....  
W2        0  P(W1 | W2)      P(W2 | W2)      .....  
...  
!EXIT    0  PN              PN              .....
```

where the probability $P(w_j|w_i)$ is given by row i, j of the matrix. If there are a total of N words in the vocabulary then PN in the above is set to $1/(N + 1)$, this ensures that the last row sums to one. As a very crude form of smoothing, a floor can be set using the -f minp option to prevent any entry falling below minp. Note, however, that this does not affect the bigram entries in the first column which are zero by definition. Finally, as with the storage of tied-mixture and discrete probabilities, a run-length encoding scheme is used whereby any value can be followed by an asterisk and a repeat count (see section 7.5).

12.5 Building a Word Network with HBuild

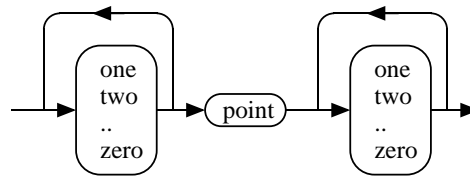


Fig. 12.5 Decimal Syntax

As mentioned in the introduction, the main function of HBUILD is allow a word-level network to be constructed from a main lattice and a set of sub-lattices. Any lattice can contain node definitions which refer to other lattices. This allows a word-level recognition network to be decomposed into a number of sub-networks which can be reused at different points in the network.

For example, suppose that decimal number input was required. A suitable network structure would be as shown in Fig. 12.5. However, to write this directly in an SLF file would require the digit loop to be written twice. This can be avoided by defining the digit loop as a sub-network and referencing it within the main *decimal* network as follows

```
# Digit network
SUBLAT=digits
N=14 L=21
# define digits
I=0 W=zero
I=1 W=one
I=2 W=two
...
I=9 W=nine
# enter/exit & loop-back null nodes
I=10 W=!NULL
I=11 W=!NULL
I=12 W=!NULL
I=13 W=!NULL
# null->>null->digits
J=0 S=10 E=11
J=1 S=11 E=0
J=2 S=11 E=1
...
J=10 S=11 E=9
# digits->>null->>null
J=11 S=0 E=12
...
J=19 S=9 E=12
J=20 S=12 E=13
# finally add loop back
J=21 S=12 E=11
.

# Decimal network
N=5 L=4
# digits -> point -> digits
I=0 W=start
I=1 L=digits
I=2 W=pause
I=3 L=digits
I=4 W=end
```

```

# digits -> point -> digits
J=0 S=0 E=1
J=1 S=1 E=2
J=2 S=2 E=3
J=3 S=3 E=4

```

The sub-network is identified by the field `SUBLAT` in the header and it is terminated by a single period on a line by itself. The main body of the sub-network is written as normal. Once defined, a sub-network can be substituted into a higher level network using an `L` field in a node definition, as in nodes 1 and 3 of the decimal network above.

Of course, this process can be continued and a higher level network could reference the decimal network wherever it needed decimal number entry.

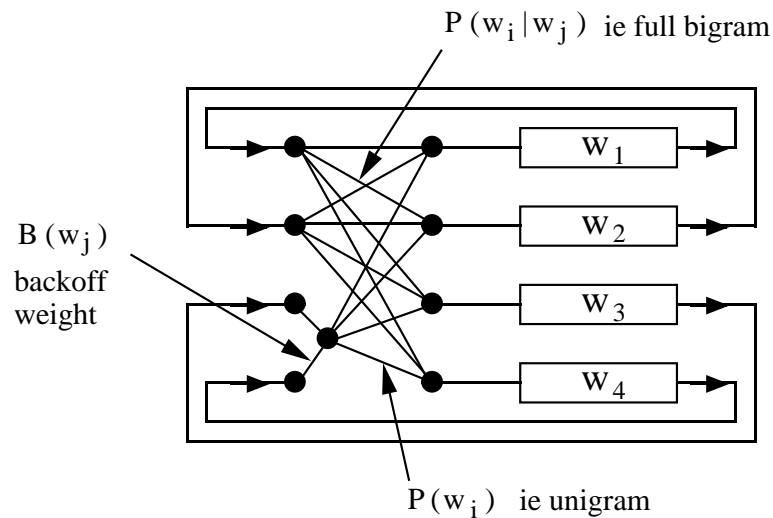


Fig. 12.6 Back-off Bigram Word-Loop Network

One of the commonest form of recognition network is the word-loop where all vocabulary items are placed in parallel with a loop-back to allow any word sequence to be recognised. This is the basic arrangement used in most dictation or transcription applications. `HBUILD` can build such a loop automatically from a list of words. It can also read in a bigram in either ARPA MIT-LL format or HTK matrix format and attach a bigram probability to each word transition. Note, however, that using a full bigram language model means that every distinct pair of words must have its own unique loop-back transition. This increases the size of the network considerably and slows down the recogniser. When a back-off bigram is used, however, backed-off transitions can share a common loop-back transition. Fig. 12.6 illustrates this. When backed-off bigrams are input via an ARPA MIT-LL format file, `HBUILD` will exploit this where possible.

Finally, `HBUILD` can automatically construct a word-pair grammar as used in the ARPA Naval Resource Management task.

12.6 Testing a Word Network using HSGen

When designing task grammars, it is useful to be able to check that the language defined by the final word network is as envisaged. One simple way to check this is to use the network as a generator by randomly traversing it and outputting the name of each word node encountered. HTK provides a very simple tool called `HSGEN` for doing this.

As an example if the file `bnet` contained the simple Bit-But network described above and the file `bdic` contained a corresponding dictionary then the command

```
HSGen bnet bdic
```

would generate a random list of examples of the language defined by `bnet`, for example,

```

start bit but bit bit bit end
start but bit but but end
start bit bit but but end
.... etc

```

This is perhaps not too informative in this case but for more complex grammars, this type of output can be quite illuminating.

HSGEN will also estimate the empirical entropy by recording the probability of each sentence generated. To use this facility, it is best to suppress the sentence output and generate a large number of examples. For example, executing

```
HSGen -s -n 1000 -q bnet bdic
```

where the `-s` option requests statistics, the `-q` option suppresses the output and `-n 1000` asks for 1000 sentences would generate the following output

```

Number of Nodes = 4 [0 null], Vocab Size = 4
Entropy = 1.156462, Perplexity = 2.229102
1000 Sentences: average len = 5.1, min=3, max=19

```

12.7 Constructing a Dictionary

As explained in section 12.1, the word level network is expanded by HNET to create the network of HMM instances needed by the recogniser. The way in which each word is expanded is determined from a dictionary.

A dictionary for use in HTK has a very simple format. Each line consists of a single word pronunciation with format

```
WORD [OUTSYM] PRONPROB P1 P2 P3 P4 ....
```

where `WORD` represents the word, `OUTSYM` is the symbol to output when that word is recognised, `PRONPROB` is the pronunciation probability (0.0 - 1.0) and `P1`, `P2`, ... are the sequence of phones or HMMs to be used in recognising that word. The output symbol and the pronunciation probability are optional. If an output symbol is not specified, the name of the word itself is output. Empty square brackets, `[]`, can be used to suppress any output when that word is recognised. For example, a dictionary for the But-Bit grammar might contain

```

bit          b i t
but          b u t
start []     sil
end []       sil

```

where for simplicity the phones are being represented by the letters of the alphabet.

If any word has more than one pronunciation, then the word has a repeated entry, for example,

```

the          th iy
the          th ax

```

corresponding to the stressed and unstressed forms of the word "the".

The pronunciations in a dictionary are normally at the phone level as in the above examples. However, if context-dependent models are wanted, these can be included directly in the dictionary. For example, the Bit-But entries might be written as

```

bit          b+i b-i+t i-t
but          b+u b-u+t u-t

```

In principle, this is never necessary since HNET can perform context expansion automatically, however, it saves computation to do this off-line as part of the dictionary construction process. Of course, this is only possible for word-internal context dependencies. Cross-word dependencies can only be generated by HNET.

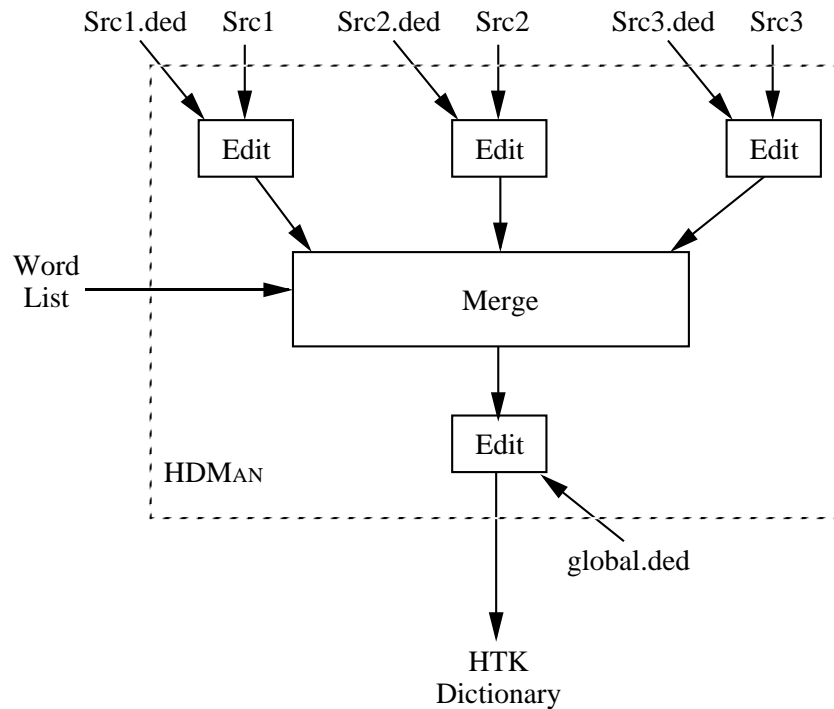


Fig. 12.7 Dictionary Construction using HDMAN

Pronouncing dictionaries are a valuable resource and if produced manually, they can require considerable investment. There are a number of commercial and public domain dictionaries available, however, these will typically have differing formats and will use different phone sets. To assist in the process of dictionary construction, HTK provides a tool called HDMAN which can be used to edit and merge differing source dictionaries to form a single uniform dictionary. The way that HDMAN works is illustrated in Fig. 12.7.

Each source dictionary file must have one pronunciation per line and the words must be sorted into alphabetical order. The word entries must be valid HTK strings as defined in section 4.6. If an arbitrary character sequence is to be allowed, then the input edit script should have the command `IM RAW` as its first command.

The basic operation of HDMAN is to scan the input streams and for each new word encountered, copy the entry to the output. In the figure, a word list is also shown. This is optional but if included HDMAN only copies words in the list. Normally, HDMAN copies just the first pronunciation that it finds for any word. Thus, the source dictionaries are usually arranged in order of *reliability*, possibly preceded by a small dictionary of special word pronunciations. For example, in Fig. 12.7, the main dictionary might be `Src2`. `Src1` might be a small dictionary containing correct pronunciations for words in `Src2` known to have errors in them. Finally, `Src3` might be a large poor quality dictionary (for example, it could be generated by a rule-based text-to-phone system) which is included as a last resort source of pronunciations for words not in the main dictionary.

As shown in the figure, HDMAN can apply a set of editing commands to each source dictionary and it can also edit the output stream. The commands available are described in full in the reference section. They operate in a similar way to those in HLED. Each set of commands is written in an edit script with one command per line. Each input edit script has the same name as the corresponding source dictionary but with the extension `.ded` added. The output edit script is stored in a file called `global.ded`. The commands provided include replace and delete at the word and phone level, context-sensitive replace and automatic conversions to left biphones, right biphones and word internal triphones.

When HDMAN loads a dictionary it adds word boundary symbols to the start and end of each pronunciation and then deletes them when writing out the new dictionary. The default for these word boundary symbols is `#` but it can be redefined using the `-b` option. The reason for this is to allow context-dependent edit commands to take account of word-initial and word-final phone positions. The examples below will illustrate this.

Rather than go through each HDMAN edit command in detail, some examples will illustrate the typical manipulations that can be performed by HDMAN. Firstly, suppose that a dictionary transcribed unstressed “-ed” endings as `ih0 d` but the required dictionary does not mark stress but uses a schwa in such cases, that is, the transformations

```
ih0 d #   ->  ax d
ih0       ->  ih (otherwise)
```

are required. These could be achieved by the following 3 commands

```
MP axd0 ih0 d #
SP axd0 ax d #
RP ih ih0
```

The context sensitive replace is achieved by merging all sequences of `ih0 d #` and then splitting the result into the sequence `ax d #`. The final `RP` command then unconditionally replaces all occurrences of `ih0` by `ih`. As a second similar example, suppose that all examples of `ax l` (as in “bottle”) are to be replaced by the single phone `e1` provided that the immediately following phone is a non-vowel. This requires the use of the `DC` command to define a context consisting of all non-vowels, then a merge using `MP` as above followed by a context-sensitive replace

```
DC nonv l r w y . . . . m n ng #
MP axl ax l
CR e1 * axl nonv
SP axl ax l
```

the final step converts all non-transformed cases of `ax l` back to their original form.

As a final example, a typical output transformation applied via the edit script `global.ded` will convert all phones to context-dependent form and append a short pause model `sp` at the end of each pronunciation. The following two commands will do this

```
TC
AS sp
```

For example, these commands would convert the dictionary entry

```
BAT b ah t
```

into

```
BAT b+ah b-ah+t ah-t sp
```

Finally, if the `-l` option is set, HDMAN will generate a log file containing a summary of the pronunciations used from each source and how many words, if any are missing. It is also possible to give HDMAN a phone list using the `-n` option. In this case, HDMAN will record how many times each phone was used and also, any phones that appeared in pronunciations but are not in the phone list. This is useful for detecting errors and unexpected phone symbols in the source dictionary.

12.8 Word Network Expansion

Now that word networks and dictionaries have been explained, the conversion of word level networks to model-based recognition networks will be described. Referring again to Fig 12.1, this expansion is performed automatically by the module HNET. By default, HNET attempts to infer the required expansion from the contents of the dictionary and the associated list of HMMs. However, 5 configurations parameters are supplied to apply more precise control where required: `ALLOWCXTEXP`, `ALLOWWRDEXP`, `FORCECXTEXP`, `FORCELEFTBI` and `FORCERIGHTBI`.

The expansion proceeds in four stages.

1. Context definition

The first step is to determine how model names are constructed from the dictionary entries and whether cross-word context expansion should be performed. The dictionary is scanned and each distinct phone is classified as either

(a) *Context Free*

In this case, the phone is skipped when determining context. An example is a model (**sp**) for short pauses. This will typically be inserted at the end of every word pronunciation but since it tends to cover a very short segment of speech it should not block context-dependent effects in a cross-word triphone system.

(b) *Context Independent*

The phone only exists in context-independent form. A typical example would be a silence model (**sil**). Note that the distinction that would be made by HNET between **sil** and **sp** is that whilst both would only appear in the HMM set in context-independent form, **sil** would appear in the contexts of other phones whereas **sp** would not.

(c) *Context Dependent*

This classification depends on whether a phone appears in the context part of the name and whether any context dependent versions of the phone exist in the HMMSet. Context Dependent phones will be subject to model name expansion.

2. *Determination of network type*

The default behaviour is to produce the simplest network possible. If the dictionary is closed (every phone name appears in the HMM list), then no expansion of phone names is performed. The resulting network is generated by straightforward substitution of each dictionary pronunciation for each word in the word network. If the dictionary is not closed, then if word internal context expansion would find each model in the HMM set then word internal context expansion is used. Otherwise, full cross-word context expansion is applied.

The determination of the network type can be modified by using the configuration parameters mentioned earlier. By default **ALLOWCXTEXP** is set true. If **ALLOWCXTEXP** is set false, then no expansion of phone names is performed and each phone corresponds to the model of the same name. The default value of **ALLOWWRDEXP** is false thus preventing context expansion across word boundaries. This also limits the expansion of the phone labels in the dictionary to word internal contexts only. If **FORCECXTEXP** is set true, then context expansion will be performed. For example, if the HMM set contained all monophones, all biphones and all triphones, then given a monophone dictionary, the default behaviour of HNET would be to generate a monophone recognition network since the dictionary would be closed. However, if **FORCECXTEXP** is set true and **ALLOWWRDEXP** is set false then word internal context expansion will be performed. If **FORCECXTEXP** is set true and **ALLOWWRDEXP** is set true then full cross-word context expansion will be performed.

3. *Network expansion*

Each word in the word network is transformed into a *word-end* node preceded by the sequence of model nodes corresponding to the word's pronunciation. For cross word context expansion, the initial and final context dependent phones (and any preceding/following context independent ones) are duplicated as many times as is necessary to cater for each different cross word context. Each duplicated word-final phone is followed by a similarly duplicated word-end node. Null words are simply transformed into word-end nodes with no preceding model nodes.

4. *Linking of models to network nodes*

Each model node is linked to the corresponding HMM definition. In each case, the required HMM model name is determined from the phone name and the surrounding context names. The algorithm used for this is

- (a) Construct the context-dependent name and see if the corresponding model exists.
- (b) Construct the context-independent name and see if the corresponding model exists.

If the configuration variable **ALLOWCXTEXP** is false (a) is skipped and if the configuration variable **FORCECXTEXP** is true (b) is skipped. If no matching model is found, an error is generated. When the right context is a boundary or **FORCELEFTBI** is true, then the context-dependent name takes the form of a left biphone, that is, the phone **p** with left context **l** becomes **l-p**. When the left context is a boundary or **FORCERIGHTBI** is true, then the context-dependent name takes the form of a right biphone, that is, the phone **p** with right context **r** becomes **p+r**. Otherwise, the context-dependent name is a full triphone, that is, **l-p+r**. Context-free phones are skipped in this process so

sil aa r sp y uw sp sil

would be expanded as

sil sil-aa+r aa-r+y sp r-y+uw y-uw+sil sp sil

assuming that `sil` is context-independent and `sp` is context-free. For word-internal systems, the context expansion can be further controlled via the configuration variable `CFWORDBOUNDARY`. When set true (default setting) context-free phones will be treated as word boundaries so

aa r sp y uw sp

would be expanded to

aa+r aa-r sp y+uw y-uw sp

Setting `CFWORDBOUNDARY` false would produce

aa+r aa-r+y sp r-y+uw y-uw sp

Note that in practice, stages (3) and (4) above actually proceed concurrently so that for the first and last phone of context-dependent models, logical models which have the same underlying physical model can be merged.

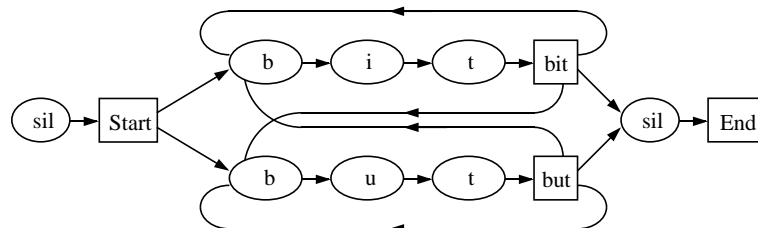


Fig. 12.8 Monophone Expansion of Bit-But Network

Having described the expansion process in some detail, some simple examples will help clarify the process. All of these are based on the Bit-But word network illustrated in Fig. 12.2. Firstly, assume that the dictionary contains simple monophone pronunciations, that is

bit	b i t
but	b u t
start	sil
end	sil

and the HMM set consists of just monophones

b i t u sil

In this case, HNET will find a closed dictionary. There will be no expansion and it will directly generate the network shown in Fig 12.8. In this figure, the rounded boxes represent model nodes and the square boxes represent word-end nodes.

Similarly, if the dictionary contained word-internal triphone pronunciations such as

bit	b+i b-i+t i-t
but	b+u b-u+t u-t
start	sil
end	sil

and the HMM set contains all the required models

b+i b-i+t i-t b+u b-u+t u-t sil

then again HNET will find a closed dictionary and the network shown in Fig. 12.9 would be generated.

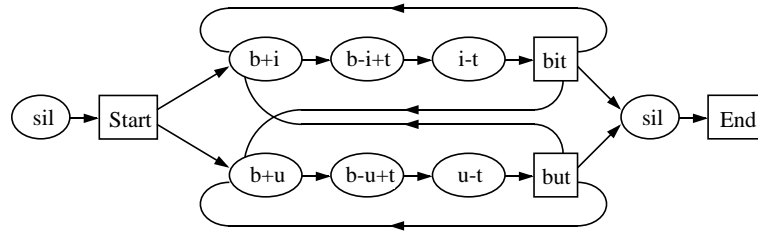


Fig. 12.9 Word Internal Triphone Expansion of Bit-But Network

If however the dictionary contained just the simple monophone pronunciations as in the first case above, but the HMM set contained just triphones, that is

```
sil-b+i  t-b+i  b-i+t  i-t+sil  i-t+b
sil-b+u  t-b+u  b-u+t  u-t+sil  u-t+b  sil
```

then HNET would perform full cross-word expansion and generate the network shown in Fig. 12.10.

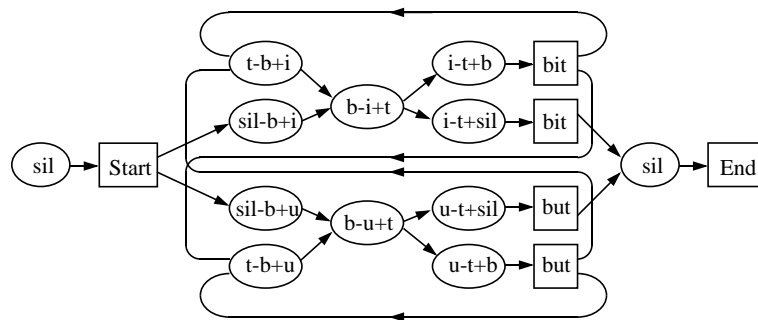


Fig. 12.10 Cross-Word Triphone Expansion of Bit-But Network

Now suppose that still using the simple monophone pronunciations, the HMM set contained all monophones, biphones and triphones. In this case, the default would be to generate the monophone network of Fig 12.8. If `FORCECXTEXP` is true but `ALLOWXWRDEXP` is set false then the word-internal network of Fig. 12.9 would be generated. Finally, if both `FORCECXTEXP` and `ALLOWXWRDEXP` are set true then the cross-word network of Fig. 12.10 would be generated.

12.9 Other Kinds of Recognition System

Although the recognition facilities of HTK are aimed primarily at sub-word based connected word recognition, it can nevertheless support a variety of other types of recognition system.

To build a phoneme recogniser, a word-level network is defined using an SLF file in the usual way except that each “word” in the network represents a single phone. The structure of the network will typically be a loop in which all phones loop back to each other.

The dictionary then contains an entry for each “word” such that the word and the pronunciation are the same, for example, the dictionary might contain

```
ih ih
eh eh
ah ah
... etc
```

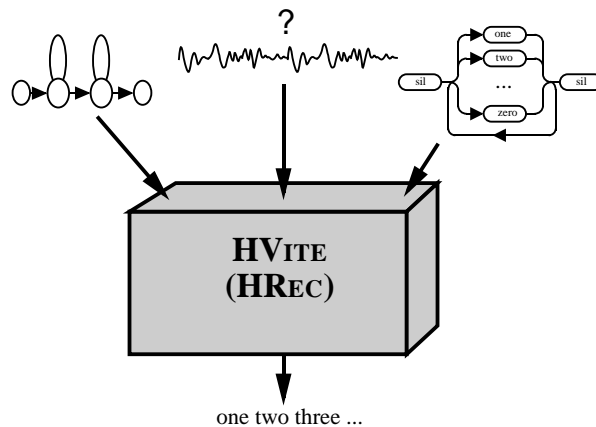

Phoneme recognisers often use biphones to provide some measure of context-dependency. Provided that the HMM set contains all the necessary biphones, then HNET will expand a simple phone loop into a context-sensitive biphone loop simply by setting the configuration variable FORCELEFTBI or FORCERIGHTBI to true, as appropriate.

Whole word recognisers can be set-up in a similar way. The word network is designed using the same considerations as for a sub-word based system but the dictionary gives the name of the whole-word HMM in place of each word pronunciation.

Finally, word spotting systems can be defined by placing each keyword in a word network in parallel with the appropriate filler models. The keywords can be whole-word models or subword based. Note in this case that word transition penalties placed on the transitions can be used to gain fine control over the false alarm rate.

Chapter 13

Decoding



The previous chapter has described how to construct a recognition network specifying what is allowed to be spoken and how each word is pronounced. Given such a network, its associated set of HMMs, and an unknown utterance, the probability of any path through the network can be computed. The task of a decoder is to find those paths which are the most likely.

As mentioned previously, decoding in HTK is performed by a library module called HREC. HREC uses the token passing paradigm to find the best path and, optionally, multiple alternative paths. In the latter case, it generates a lattice containing the multiple hypotheses which can if required be converted to an N-best list. To drive HREC from the command line, HTK provides a tool called HVITE. As well as providing basic recognition, HVITE can perform forced alignments, lattice rescoring and recognise direct audio input.

To assist in evaluating the performance of a recogniser using a test database and a set of reference transcriptions, HTK also provides a tool called HRESULTS to compute word accuracy and various related statistics. The principles and use of these recognition facilities are described in this chapter.

13.1 Decoder Operation

As described in Chapter 12 and illustrated by Fig. 12.1, decoding in HTK is controlled by a recognition network compiled from a word-level network, a dictionary and a set of HMMs. The recognition network consists of a set of nodes connected by arcs. Each node is either a HMM model instance or a word-end. Each model node is itself a network consisting of states connected by arcs. Thus, once fully compiled, a recognition network ultimately consists of HMM states connected by transitions. However, it can be viewed at three different levels: word, model and state. Fig. 13.1 illustrates this hierarchy.

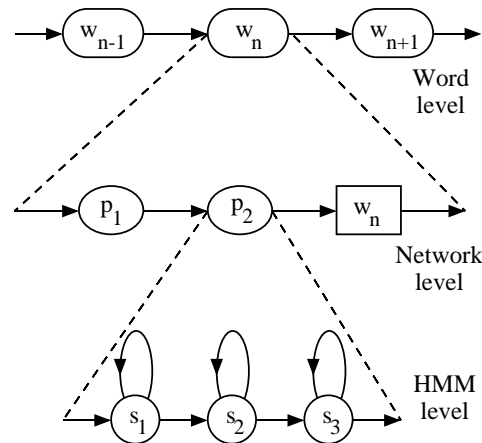


Fig. 13.1 Recognition Network Levels

For an unknown input utterance with T frames, every path from the start node to the exit node of the network which passes through exactly T emitting HMM states is a potential recognition hypothesis. Each of these paths has a log probability which is computed by summing the log probability of each individual transition in the path and the log probability of each emitting state generating the corresponding observation. Within-HMM transitions are determined from the HMM parameters, between-model transitions are constant and word-end transitions are determined by the language model likelihoods attached to the word level networks.

The job of the decoder is to find those paths through the network which have the highest log probability. These paths are found using a *Token Passing* algorithm. A token represents a partial path through the network extending from time 0 through to time t . At time 0, a token is placed in every possible start node.

Each time step, tokens are propagated along connecting transitions stopping whenever they reach an emitting HMM state. When there are multiple exits from a node, the token is copied so that all possible paths are explored in parallel. As the token passes across transitions and through nodes, its log probability is incremented by the corresponding transition and emission probabilities. A network node can hold at most N tokens. Hence, at the end of each time step, all but the N best tokens in any node are discarded.

As each token passes through the network it must maintain a history recording its route. The amount of detail in this history depends on the required recognition output. Normally, only word sequences are wanted and hence, only transitions out of word-end nodes need be recorded. However, for some purposes, it is useful to know the actual model sequence and the time of each model to model transition. Sometimes a description of each path down to the state level is required. All of this information, whatever level of detail is required, can conveniently be represented using a lattice structure.

Of course, the number of tokens allowed per node and the amount of history information requested will have a significant impact on the time and memory needed to compute the lattices. The most efficient configuration is $N = 1$ combined with just word level history information and this is sufficient for most purposes.

A large network will have many nodes and one way to make a significant reduction in the computation needed is to only propagate tokens which have some chance of being amongst the eventual winners. This process is called *pruning*. It is implemented at each time step by keeping a record of the best token overall and de-activating all tokens whose log probabilities fall more than a *beam-width* below the best. For efficiency reasons, it is best to implement primary pruning at the model rather than the state level. Thus, models are deactivated when they have no tokens in any state within the beam and they are reactivated whenever active tokens are propagated into them. State-level pruning is also implemented by replacing any token by a null (zero probability) token if it falls outside of the beam. If the pruning beam-width is set too small then the most likely path might be pruned before its token reaches the end of the utterance. This results in a *search error*. Setting the beam-width is thus a compromise between speed and avoiding search errors.

When using word loops with bigram probabilities, tokens emitted from word-end nodes will have a language model probability added to them before entering the following word. Since the range of language model probabilities is relatively small, a narrower beam can be applied to word-end nodes without incurring additional search errors. This beam is calculated relative to the best word-end token and it is called a *word-end beam*. In the case, of a recognition network with an arbitrary topology, word-end pruning may still be beneficial but this can only be justified empirically.

Finally, a third type of pruning control is provided. An upper-bound on the allowed use of compute resource can be applied by setting an upper-limit on the number of models in the network which can be active simultaneously. When this limit is reached, the pruning beam-width is reduced in order to prevent it being exceeded.

13.2 Decoder Organisation

The decoding process itself is performed by a set of core functions provided within the library module HREC. The process of recognising a sequence of utterances is illustrated in Fig. 13.2.

The first stage is to create a *recogniser-instance*. This is a data structure containing the compiled recognition network and storage for storing tokens. The point of encapsulating all of the information and storage needed for recognition into a single object is that HREC is re-entrant and can support multiple recognisers simultaneously. Thus, although this facility is not utilised in the supplied recogniser HVITE, it does provide applications developers with the capability to have multiple recognisers running with different networks.

Once a recogniser has been created, each unknown input is processed by first executing a *start recogniser* call, and then processing each observation one-by-one. When all input observations have been processed, recognition is completed by generating a lattice. This can be saved to disk as a standard lattice format (SLF) file or converted to a transcription.

The above decoder organisation is extremely flexible and this is demonstrated by the HTK tool HVITE which is a simple shell program designed to allow HREC to be driven from the command line.

Firstly, input control in the form of a recognition network allows three distinct modes of operation

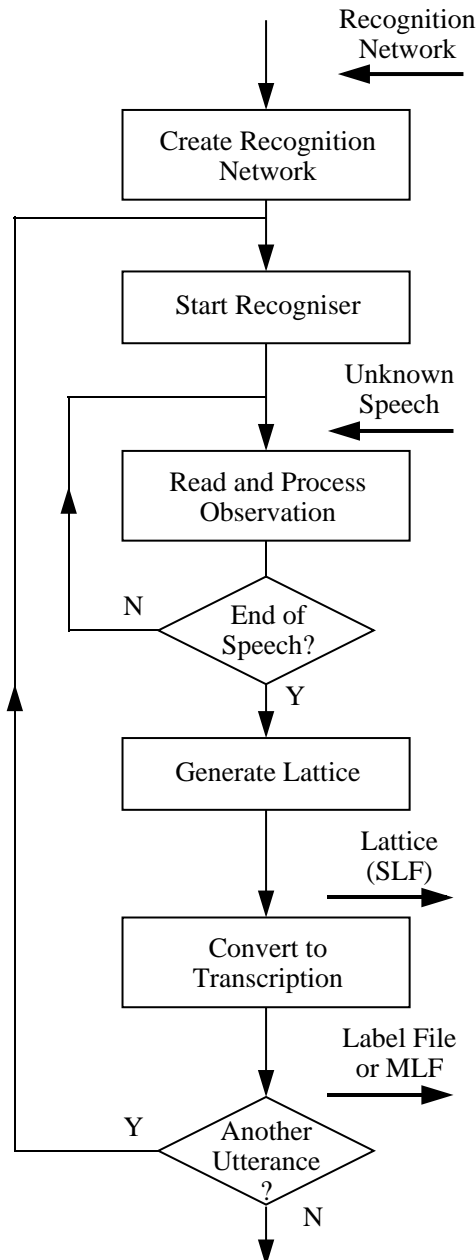


Fig. 13.2 Recognition Processing

1. *Recognition*

This is the conventional case in which the recognition network is compiled from a task level word network.

2. *Forced Alignment*

In this case, the recognition network is constructed from a word level transcription (i.e. orthography) and a dictionary. The compiled network may include optional silences between words and pronunciation variants. Forced alignment is often useful during training to automatically derive phone level transcriptions. It can also be used in automatic annotation systems.

3. *Lattice-based Rescoring*

In this case, the input network is compiled from a lattice generated during an earlier recognition run. This mode of operation can be extremely useful for recogniser development since rescoring can be an order of magnitude faster than normal recognition. The required lattices are usually generated by a basic recogniser running with multiple tokens, the idea being to

generate a lattice containing both the correct transcription plus a representative number of confusions. Rescoring can then be used to quickly evaluate the performance of more advanced recognisers and the effectiveness of new recognition techniques.

The second source of flexibility lies in the provision of multiple tokens and recognition output in the form of a lattice. In addition to providing a mechanism for rescoring, lattice output can be used as a source of multiple hypotheses either for further recognition processing or input to a natural language processor. Where convenient, lattice output can easily be converted into N-best lists.

Finally, since HREC is explicitly driven step-by-step at the observation level, it allows fine control over the recognition process and a variety of traceback and on-the-fly output possibilities.

For application developers, HREC and the HTK library modules on which it depends can be linked directly into applications. It will also be available in the form of an industry standard API. However, as mentioned earlier the HTK toolkit also supplies a tool called HVITE which is a shell program designed to allow HREC to be driven from the command line. The remainder of this chapter will therefore explain the various facilities provided for recognition from the perspective of HVITE.

13.3 Recognition using Test Databases

When building a speech recognition system or investigating speech recognition algorithms, performance must be monitored by testing on databases of test utterances for which reference transcriptions are available. To use HVITE for this purpose it is invoked with a command line of the form

```
HVite -w wdnnet dict hmmlist testf1 testf2 ....
```

where `wdnnet` is an SLF file containing the word level network, `dict` is the pronouncing dictionary and `hmmlist` contains a list of the HMMs to use. The effect of this command is that HVITE will use HNET to compile the word level network and then use HREC to recognise each test file. The parameter kind of these test files must match exactly with that used to train the HMMs. For evaluation purposes, test files are normally stored in parameterised form but only the basic static coefficients are saved on disk. For example, delta parameters are normally computed during loading. As explained in Chapter 5, HTK can perform a range of parameter conversions on loading and these are controlled by configuration variables. Thus, when using HVITE, it is normal to include a configuration file via the `-C` option in which the required target parameter kind is specified. Section 13.7 below on processing direct audio input explains the use of configuration files in more detail.

In the simple default form of invocation given above, HVITE would expect to find each HMM definition in a separate file in the current directory and each output transcription would be written to a separate file in the current directory. Also, of course, there will typically be a large number of test files.

In practice, it is much more convenient to store HMMs in master macro files (MMFs), store transcriptions in master label files (MLFs) and list data files in a script file. Thus, a more common form of the above invocation would be

```
HVite -T 1 -S test.scp -H hmmset -i results -w wdnnet dict hmmlist
```

where the file `test.scp` contains the list of test file names, `hmmset` is an MMF containing the HMM definitions¹, and `results` is the MLF for storing the recognition output.

As shown, it is usually a good idea to enable basic progress reporting by setting the trace option as shown. This will cause the recognised word string to be printed after processing each file. For example, in a digit recognition task the trace output might look like

```
File: testf1.mfc
SIL ONE NINE FOUR SIL
[178 frames] -96.1404 [Ac=-16931.8 LM=-181.2] (Act=75.0)
```

¹ Large HMM sets will often be distributed across a number of MMF files, in this case, the `-H` option will be repeated for each file.

where the information listed after the recognised string is the total number of frames in the utterance, the average log probability per frame, the total acoustic likelihood, the total language model likelihood and the average number of active models.

The corresponding transcription written to the output MLF form will contain an entry of the form

```
"testf1.rec"
  0 6200000 SIL -6067.333008
 6200000 9200000 ONE -3032.359131
 9200000 12300000 NINE -3020.820312
12300000 17600000 FOUR -4690.033203
17600000 17800000 SIL -302.439148
.
```

This shows the start and end time of each word and the total log probability. The fields output by HVITE can be controlled using the `-o`. For example, the option `-o ST` would suppress the scores and the times to give

```
"testf1.rec"
SIL
ONE
NINE
FOUR
SIL
.
```

In order to use HVITE effectively and efficiently, it is important to set appropriate values for its pruning thresholds and the language model scaling parameters. The main pruning beam is set by the `-t` option. Some experimentation will be necessary to determine appropriate levels but around 250.0 is usually a reasonable starting point. Word-end pruning (`-v`) and the maximum model limit (`-u`) can also be set if required, but these are not mandatory and their effectiveness will depend greatly on the task.

The relative levels of insertion and deletion errors can be controlled by scaling the language model likelihoods using the `-s` option and adding a fixed *penalty* using the `-p` option. For example, setting `-s 10.0 -p -20.0` would mean that every language model log probability x would be converted to $10x - 20$ before being added to the tokens emitted from the corresponding word-end node. As an extreme example, setting `-p 100.0` caused the digit recogniser above to output

```
SIL OH OH ONE OH OH OH NINE FOUR OH OH OH OH SIL
```

where adding 100 to each word-end transition has resulted in a large number of insertion errors. The word inserted is “oh” primarily because it is the shortest in the vocabulary. Another problem which may occur during recognition is the inability to arrive at the final node in the recognition network after processing the whole utterance. The user is made aware of the problem by the message “No tokens survived to final node of network”. The inability to match the data against the recognition network is usually caused by poorly trained acoustic models and/or very tight pruning beam-widths. In such cases, partial recognition results can still be obtained by setting the HREC configuration variable `FORCEOUT` true. The results will be based on the most likely partial hypothesis found in the network.

13.4 Evaluating Recognition Results

Once the test data has been processed by the recogniser, the next step is to analyse the results. The tool HRESULTS is provided for this purpose. HRESULTS compares the transcriptions output by HVITE with the original reference transcriptions and then outputs various statistics. HRESULTS matches each of the recognised and reference label sequences by performing an optimal string match using dynamic programming. Except when scoring word-spotter output as described later, it does not take any notice of any boundary timing information stored in the files being compared. The optimal string match works by calculating a score for the match with respect to the reference such that identical labels match with score 0, a label insertion carries a score of 7, a deletion carries a

score of 7 and a substitution carries a score of 10². The optimal string match is the label alignment which has the lowest possible score.

Once the optimal alignment has been found, the number of substitution errors (S), deletion errors (D) and insertion errors (I) can be calculated. The percentage correct is then

$$\text{Percent Correct} = \frac{N - D - S}{N} \times 100\% \quad (13.1)$$

where N is the total number of labels in the reference transcriptions. Notice that this measure ignores insertion errors. For many purposes, the percentage accuracy defined as

$$\text{Percent Accuracy} = \frac{N - D - S - I}{N} \times 100\% \quad (13.2)$$

is a more representative figure of recogniser performance.

HRESULTS outputs both of the above measures. As with all HTK tools it can process individual label files and files stored in MLFs. Here the examples will assume that both reference and test transcriptions are stored in MLFs.

As an example of use, suppose that the MLF `results` contains recogniser output transcriptions, `refs` contains the corresponding reference transcriptions and `wlist` contains a list of all labels appearing in these files. Then typing the command

```
HResults -I refs wlist results
```

would generate something like the following

```
===== HTK Results Analysis =====
Date: Sat Sep  2 14:14:22 1995
Ref : refs
Rec : results
----- Overall Results -----
SENT: %Correct=98.50 [H=197, S=3, N=200]
WORD: %Corr=99.77, Acc=99.65 [H=853, D=1, S=1, I=1, N=855]
=====
```

The first part shows the date and the names of the files being used. The line labelled `SENT` shows the total number of complete sentences which were recognised correctly. The second line labelled `WORD` gives the recognition statistics for the individual words³.

It is often useful to visually inspect the recognition errors. Setting the `-t` option causes aligned test and reference transcriptions to be output for all sentences containing errors. For example, a typical output might be

```
Aligned transcription: testf9.lab vs testf9.rec
LAB: FOUR      SEVEN NINE THREE
REC: FOUR OH SEVEN FIVE THREE
```

here an “oh” has been inserted by the recogniser and “nine” has been recognised as “five”

If preferred, results output can be formatted in an identical manner to NIST scoring software by setting the `-h` option. For example, the results given above would appear as follows in NIST format

```
,-----,
| HTK Results Analysis at Sat Sep  2 14:42:06 1995 |
| Ref: refs |
| Rec: results |
|=====|
|          # Snt | Corr  Sub  Del  Ins  Err  S. Err | |
|---|---|---|
| Sum/Avg | 200 | 99.77  0.12  0.12  0.12  0.35  1.50 |
|-----|
```

²The default behaviour of HRESULTS is slightly different to the widely used US NIST scoring software which uses weights of 3,3 and 4 and a slightly different alignment algorithm. Identical behaviour to NIST can be obtained by setting the `-n` option.

³All the examples here will assume that each label corresponds to a word but in general the labels could stand for any recognition unit such as phones, syllables, etc. HRESULTS does not care what the labels mean but for human consumption, the labels `SENT` and `WORD` can be changed using the `-a` and `-b` options.

When computing recognition results it is sometimes inappropriate to distinguish certain labels. For example, to assess a digit recogniser used for voice dialing it might be required to treat the alternative vocabulary items “oh” and “zero” as being equivalent. This can be done by making them equivalent using the `-e` option, that is

```
HResults -e ZERO OH .....
```

If a label is equated to the special label `???`, then it is ignored. Hence, for example, if the recognition output had silence marked by `SIL`, the setting the option `-e ??? SIL` would cause all the `SIL` labels to be ignored.

`HRESULTS` contains a number of other options. Recognition statistics can be generated for each file individually by setting the `-f` option and a confusion matrix can be generated by setting the `-p` option. When comparing phone recognition results, `HRESULTS` will strip any triphone contexts by setting the `-s` option. `HRESULTS` can also process N-best recognition output. Setting the option `-d N` causes `HRESULTS` to search the first `N` alternatives of each test output file to find the most accurate match with the reference labels.

When analysing the performance of a speaker independent recogniser it is often useful to obtain accuracy figures on a per speaker basis. This can be done using the option `-k mask` where `mask` is a pattern used to extract the speaker identifier from the test label file name. The pattern consists of a string of characters which can include the pattern matching metacharacters `*` and `?` to match zero or more characters and a single character, respectively. The pattern should also contain a string of one or more `%` characters which are used as a mask to identify the speaker identifier.

For example, suppose that the test filenames had the following structure

```
DIGITS_spkr_nnnn.rec
```

where `spkr` is a 4 character speaker id and `nnnn` is a 4 digit utterance id. Then executing `HRESULTS` by

```
HResults -h -k '*_%%%-????.*' .....
```

would give output of the form

```
,-----,
| HTK Results Analysis at Sat Sep  2 15:05:37 1995 |
| Ref: refs |
| Rec: results |
|-----|
|   SPKR | # Snt |  Corr   Sub   Del   Ins   Err   S. Err |
|-----|
|   dgo1 |   20 | 100.00  0.00  0.00  0.00  0.00  0.00 |
|-----|
|   pcw1 |   20 |  97.22  1.39  1.39  0.00  2.78 10.00 |
|-----|
|.....|
|=====|
| Sum/Avg |  200 |  99.77  0.12  0.12  0.12  0.35  1.50 |
|-----|
```

In addition to string matching, `HRESULTS` can also analyse the results of a recogniser configured for word-spotting. In this case, there is no DP alignment. Instead, each recogniser label w is compared with the reference transcriptions. If the start and end times of w lie either side of the mid-point of an identical label in the reference, then that recogniser label represents a *hit*, otherwise it is a *false-alarm* (FA).

The recogniser output must include the log likelihood scores as well as the word boundary information. These scores are used to compute the *Figure of Merit* (FOM) defined by NIST which is an upper-bound estimate on word spotting accuracy averaged over 1 to 10 false alarms per hour. The FOM is calculated as follows where it is assumed that the total duration of the test speech is T hours. For each word, all of the spots are ranked in score order. The percentage of true hits p_i found before the i 'th false alarm is then calculated for $i = 1 \dots N + 1$ where N is the first integer $\geq 10T - 0.5$. The figure of merit is then defined as

$$\text{FOM} = \frac{1}{10T}(p_1 + p_2 + \dots + p_N + ap_{N+1}) \quad (13.3)$$

where $a = 10T - N$ is a factor that interpolates to 10 false alarms per hour.

Word spotting analysis is enabled by setting the `-w` option and the resulting output has the form

```
----- Figures of Merit -----
Keyword:   #Hits   #FAs  #Actual   FOM
  BADGE:      92     83    102    73.56
  CAMERA:     20      2     22    89.86
  WINDOW:     84      8     92    86.98
  VIDEO:      72      6     72    99.81
  Overall:   268     99    188    87.55
-----
```

If required the standard time unit of 1 hour as used in the above definition of FOM can be changed using the `-u` option.

13.5 Generating Forced Alignments

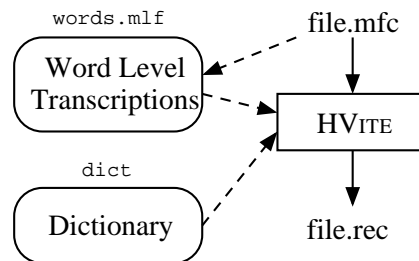


Fig. 13.3 Forced Alignment

HVITE can be made to compute forced alignments by not specifying a network with the `-w` option but by specifying the `-a` option instead. In this mode, HVITE computes a new network for each input utterance using the word level transcriptions and a dictionary. By default, the output transcription will just contain the words and their boundaries. One of the main uses of forced alignment, however, is to determine the actual pronunciations used in the utterances used to train the HMM system in this case, the `-m` option can be used to generate model level output transcriptions. This type of forced alignment is usually part of a *bootstrap* process, initially models are trained on the basis of one fixed pronunciation per word⁴. Then HVITE is used in forced alignment mode to select the best matching pronunciations. The new phone level transcriptions can then be used to retrain the HMMs. Since training data may have leading and trailing silence, it is usually necessary to insert a silence model at the start and end of the recognition network. The `-b` option can be used to do this.

As an illustration, executing

```
HVite -a -b sil -m -o SWT -I words.mlf \
      -H hmmset dict hmmlist file.mfc
```

would result in the following sequence of events (see Fig. 13.3). The input file name `file.mfc` would have its extension replaced by `lab` and then a label file of this name would be searched for. In this case, the MLF file `words.mlf` has been loaded. Assuming that this file contains a word level transcription called `file.lab`, this transcription along with the dictionary `dict` will be used to construct a network equivalent to `file.lab` but with alternative pronunciations included in parallel. Since `-b` option has been set, the specified `sil` model will be inserted at the start and end of the network. The decoder then finds the best matching path through the network and constructs a lattice which includes model alignment information. Finally, the lattice is converted to a transcription and output to the label file `file.rec`. As for testing on a database, alignments

⁴ The HLED EX command can be used to compute phone level transcriptions when there is only one possible phone transcription per word

will normally be computed on a large number of input files so in practice the input files would be listed in a `.scp` file and the output transcriptions would be written to an MLF using the `-i` option.

When the `-m` option is used, the transcriptions output by HVITE would by default contain both the model level and word level transcriptions. For example, a typical fragment of the output might be

```

7500000 8700000 f -1081.604736 FOUR 30.000000
8700000 9800000 ao -903.821350
9800000 10400000 r -665.931641
10400000 10400000 sp -0.103585
10400000 11700000 s -1266.470093 SEVEN 22.860001
11700000 12500000 eh -765.568237
12500000 13000000 v -476.323334
13000000 14400000 n -1285.369629
14400000 14400000 sp -0.103585

```

Here the score alongside each model name is the acoustic score for that segment. The score alongside the word is just the language model score.

Although the above information can be useful for some purposes, for example in bootstrap training, only the model names are required. The formatting option `-o SWT` in the above suppresses all output except the model names.

13.6 Decoding and Adaptation

Speaker adaptation techniques allow speaker independent model sets to be adapted to better fit the characteristics of individual speakers using a small amount of adaptation data. Chapter 9 described how the HEADAPT tool can be used to perform offline supervised adaptation (using the true transcription of the data).

This section describes how adapted model sets are used in the recognition process and also how HVITE can be used to perform unsupervised adaptation on a model set (when no transcription is available).

13.6.1 Recognition with Adapted HMMs

As described in section 9.1.3, HEADAPT can produce either a MMF containing the newly adapted model set or a TMF containing just the adaptation transform. If a transformed MMF has been constructed, then HVITE can be used in the usual way. If a TMF has been produced however, this needs to be passed to HVITE (using the `-J` option) along with the model set from which the transform was estimated. HVITE then transforms the model set using the TMF and recognises the input speech using the transformed model set. Thus, a common form of invocation would be

```

HVite -S test.scp -H hmmset -J trans.tmf -i results \
      -w wdnnet dict hmmlist

```

13.6.2 Unsupervised Adaptation

Unsupervised adaptation occurs when no transcription of the adaptation data exists and one must be generated. In this case HVITE can be used to create a transcription of the adaptation data and use this to estimate a transformation using MLLR. The transformation can then be saved to a TMF using the `-K` option.

Unsupervised adaptation is signalled by the use of the `-j` option and this also controls the mode of adaptation by specifying the number of utterances to be processed before a transform is estimated. Thus, the adaptation can be varied between static (adaptation only performed after recognition of all utterances) and incremental adaptation. As soon as a transform has been estimated during incremental adaptation, it is used to adapt the model set to improve performance for any subsequent utterances. Note however that only the final transformation is saved. To use HVITE for this purpose it is invoked with a command line of the form

```

HVite -S adapt.scp -H hmmset -K trans.tmf -j 10 -i results \
      -w wdnnet dict hmmlist

```

where `adapt.scp` contains a list of coded adaptation sentences, adaptation is being performed incrementally every 10 utterances and the final transform is stored in `trans.tmf`

13.7 Recognition using Direct Audio Input

In all of the preceding discussion, it has been assumed that input was from speech files stored on disk. These files would normally have been stored in parameterised form so that little or no conversion of the source speech data was required. When `HVITE` is invoked with no files listed on the command line, it assumes that input is to be taken directly from the audio input. In this case, configuration variables must be used to specify firstly how the speech waveform is to be captured and secondly, how the captured waveform is to be converted to parameterised form.

Dealing with waveform capture first, as described in section 5.9, `HTK` provides two main forms of control over speech capture: signals/keypress and an automatic speech/silence detector. To use the speech/silence detector alone, the configuration file would contain the following

```
# Waveform capture
SOURCERATE=625.0
SOURCEKIND=HAUDIO
SOURCEFORMAT=HTK
USESILDET=T
MEASURESIL=F
OUTSILWARN=T
ENORMALISE=F
```

where the source sampling rate is being set to 16kHz. Notice that the `SOURCEKIND` must be set to `HAUDIO` and the `SOURCEFORMAT` must be set to `HTK`. Setting the Boolean variable `USESILDET` causes the speech/silence detector to be used, and the `MEASURESIL` `OUTSILWARN` variables result in a measurement being taken of the background silence level prior to capturing the first utterance. To make sure that each input utterance is being captured properly, the `HVITE` option `-g` can be set to cause the captured wave to be output after each recognition attempt. Note that for a live audio input system, the configuration variable `ENORMALISE` should be explicitly set to `FALSE` both when training models and when performing recognition. Energy normalisation cannot be used with live audio input, and the default setting for this variable is `TRUE`.

As an alternative to using the speech/silence detector, a signal can be used to start and stop recording. For example,

```
# Waveform capture
SOURCERATE=625.0
SOURCEKIND=HAUDIO
SOURCEFORMAT=HTK
AUDIOSIG=2
```

would result in the Unix interrupt signal (usually the Control-C key) being used as a start and stop control⁵. Key-press control of the audio input can be obtained by setting `AUDIOSIG` to a negative number.

Both of the above can be used together, in this case, audio capture is disabled until the specified signal is received. From then on control is in the hands of the speech/silence detector.

The captured waveform must be converted to the required target parameter kind. Thus, the configuration file must define all of the parameters needed to control the conversion of the waveform to the required target kind. This process is described in detail in Chapter 5. As an example, the following parameters would allow conversion to Mel-frequency cepstral coefficients with delta and acceleration parameters.

```
# Waveform to MFCC parameters
TARGETKIND=MFCC_O_D_A
TARGETRATE=100000.0
WINDOWSIZE=250000.0
ZMEANSOURCE=T
```

⁵ The underlying signal number must be given, `HTK` cannot interpret the standard Unix signal names such as `SIGINT`

```

USEHAMMING = T
PREEMCOEF = 0.97
USEPOWER = T
NUMCHANS = 26
CEPLIFTER = 22
NUMCEPS = 12

```

Many of these variable settings are the default settings and could be omitted, they are included explicitly here as a reminder of the main configuration options available.

When HVITE is executed in direct audio input mode, it issues a prompt prior to each input and it is normal to enable basic tracing so that the recognition results can be seen. A typical terminal output might be

```

READY[1]>
Please speak sentence - measuring levels
Level measurement completed
DIAL ONE FOUR SEVEN
== [258 frames] -97.8668 [Ac=-25031.3 LM=-218.4] (Act=22.3)

READY[2]>
CALL NINE TWO EIGHT
== [233 frames] -97.0850 [Ac=-22402.5 LM=-218.4] (Act=21.8)

```

etc

If required, a transcription of each spoken input can be output to a label file or an MLF in the usual way by setting the `-e` option. However, to do this a file name must be synthesised. This is done by using a counter prefixed by the value of the HVITE configuration variable `RECOUTPREFIX` and suffixed by the value of `RECOUTSUFFIX`. For example, with the settings

```

RECOUTPREFIX = sjy
RECOUTSUFFIX = .rec

```

then the output transcriptions would be stored as `sjy0001.rec`, `sjy0002.rec` etc.

13.8 N-Best Lists and Lattices

As noted in section 13.1, HVITE can generate lattices and N-best outputs. To generate an N-best list, the `-n` option is used to specify the number of N-best tokens to store per state and the number of N-best hypotheses to generate. The result is that for each input utterance, a multiple alternative transcription is generated. For example, setting `-n 4 20` with a digit recogniser would generate an output of the form

```

"testf1.rec"
FOUR
SEVEN
NINE
OH
///
FOUR
SEVEN
NINE
OH
OH
///

etc

```

The lattices from which the N-best lists are generated can be output by setting the option `-z ext`. In this case, a lattice called `testf.ext` will be generated for each input test file `testf.xxx`. By default, these lattices will be stored in the same directory as the test files, but they can be redirected to another directory using the `-l` option.

The lattices generated by HVITE have the following general form

```

VERSION=1.0
UTTERANCE=testf1.mfc
lmname=wdnet
lmscale=20.00 wdpenalty=-30.00
vocab=dict
N=31 L=56
I=0 t=0.00
I=1 t=0.36
I=2 t=0.75
I=3 t=0.81
... etc
I=30 t=2.48
J=0 S=0 E=1 W=SILENCE v=0 a=-3239.01 l=0.00
J=1 S=1 E=2 W=FOUR v=0 a=-3820.77 l=0.00
... etc
J=55 S=29 E=30 W=SILENCE v=0 a=-246.99 l=-1.20

```

The first 5 lines comprise a header which records names of the files used to generate the lattice along with the settings of the language model scale and penalty factors. Each node in the lattice represents a point in time measured in seconds and each arc represents a word spanning the segment of the input starting at the time of its start node and ending at the time of its end node. For each such span, *v* gives the number of the pronunciation used, *a* gives the acoustic score and *l* gives the language model score.

The language model scores in output lattices do not include the scale factors and penalties. These are removed so that the lattice can be used as a constraint network for subsequent recogniser testing. When using HVITE normally, the word level network file is specified using the *-w* option. When the *-w* option is included but no file name is included, HVITE constructs the name of a lattice file from the name of the test file and inputs that. Hence, a new recognition network is created for each input file and recognition is very fast. For example, this is an efficient way of experimentally determining optimum values for the language model scale and penalty factors.

Part III

Reference Section

Chapter 14

The HTK Tools

14.1 HBuild

14.1.1 Function

This program is used to convert input files that represent language models in a number of different formats and output a standard HTK lattice. The main purpose of HBUILD is to allow the expansion of HTK multi-level lattices and the conversion of bigram language models (such as those generated by HLSTATS) into lattice format.

The specific input file types supported by HBUILD are:

1. HTK multi-level lattice files.
2. Back-off bigram files in ARPA/MIT-LL format.
3. Matrix bigram files produced by HLSTATS.
4. Word lists (to generate a word-loop grammar).
5. Word-pair grammars in ARPA Resource Management format.

The formats of both types of bigram supported by HBUILD are described in Chapter 12. The format for multi-level HTK lattice files is described in Chapter 17.

14.1.2 Use

HBUILD is invoked by the command line

```
HBuild [options] wordList outLatFile
```

The `wordList` should contain a list of all the words used in the input language model. The options specify the type of input language model as well as the source filename. If none of the flags specifying input language model type are given a simple word-loop is generated using the `wordList` given. After processing the input language model, the resulting lattice is saved to file `outLatFile`.

The operation of HBUILD is controlled by the following command line options

- b Output the lattice in binary format. This increases speed of subsequent loading (default ASCII text lattices).
- m `fn` The matrix format bigram in `fn` forms the input language model.
- n `fn` The ARPA/MIT-LL format back-off bigram in `fn` forms the input language model.
- s `st en` Set the bigram entry and exit words to `st` and `en`. (Default `!ENTER` and `!EXIT`). Note that no words will follow the exit word, or precede the entry word. Both the entry and exit word must be included in the `wordList`. This option is only effective in conjunction with the `-n` option.
- t `st en` This option is used with word-loops and word-pair grammars. An output lattice is produced with an initial word-symbol `st` (before the loop) and a final word-symbol `en` (after the loop). This allows initial and final silences to be specified. (Default is that the initial and final nodes are labelled with `!NULL`). Note that `st` and `en` shouldn't be included in the `wordList` unless they occur elsewhere in the network. This is only effective for word-loop and word-pair grammars.
- u `s` The unknown word is `s` (default `!NULL`). This option only has an effect when bigram input language models are specified. It can be used in conjunction with the `-z` flag to delete the symbol for unknown words from the output lattice.
- w `fn` The word-pair grammar in `fn` forms the input language model. The file must be in the format used for the ARPA Resource Management grammar.
- x `fn` The extended HTK lattice in `fn` forms the input language model. This option is used to expand a multi-level lattice into a single level lattice that can be processed by other HTK tools.
- z Delete (zap) any references to the unknown word (see `-u` option) in the output lattice.

HBUILD also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

14.1.3 Tracing

HBUILD supports the following trace options where each trace flag is given using an octal base
0001 basic progress reporting.

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

14.2 HCompV

14.2.1 Function

This program will calculate the global mean and covariance of a set of training data. It is primarily used to initialise the parameters of a HMM such that all component means and all covariances are set equal to the global data mean and covariance. This might form the first stage of a *flat start* training scheme where all models are initially given the same parameters. Alternatively, the covariances may be used as the basis for Fixed Variance and Grand Variance training schemes. These can sometimes be beneficial in adverse conditions where a fixed covariance matrix can give increased robustness.

When training large model sets from limited data, setting a floor is often necessary to prevent variances being badly underestimated through lack of data. One way of doing this is to define a variance macro called `varFloorN` where N is the stream index. HCOMPV can also be used to create these variance floor macros with values equal to a specified fraction of the global variance.

14.2.2 Use

HCOMPV is invoked via the command line

```
HCompV [options] hmm trainFiles ...
```

where `hmm` is the name of the physical HMM whose parameters are to be initialised. The effect of this command is to compute the covariance of the speech training data and then copy it into every Gaussian component of the given HMM definition. If there are multiple data streams, then a separate covariance is estimated for each stream. The HMM can have a mix of diagonal and full covariances and an option exists to update the means also. The HMM definition can be contained within one or more macro files loaded via the standard `-H` option. Otherwise, the definition will be read from a file called `hmm`. Any tyings in the input definition will be preserved in the output. By default, the new updated definition overwrites the existing one. However, a new definition file including any macro files can be created by specifying an appropriate target directory using the standard `-M` option.

In addition to the above, an option `-f` is provided to compute variance floor macros equal to a specified fraction of the global variance. In this case, the newly created macros are written to a file called `vFloors`. For each stream N defined for `hmm`, a variance macro called `varFloorN` is created. If a target directory is specified using the standard `-M` option then the new file will be written there, otherwise it is written in the current directory.

The list of train files can be stored in a script file if required. Furthermore, the data used for estimating the global covariance can be limited to that corresponding to a specified label.

The detailed operation of HCOMPV is controlled by the following command line options

- `-f f` Create variance floor macros with values equal to `f` times the global variance. One macro is created for each input stream and the output is stored in a file called `vFloors`.
- `-l s` The string `s` must be the name of a segment label. When this option is used, HCOMPV searches through all of the training files and uses only the speech frames from segments with the given label. When this option is not used, HCOMPV uses all of the data in each training file.
- `-m` The covariances of the output HMM are always updated however updating the means must be specifically requested. When this option is set, HCOMPV updates all the HMM component means with the sample mean computed from the training files.
- `-o s` The string `s` is used as the name of the output HMM in place of the source name.
- `-v f` This sets the minimum variance (i.e. diagonal elements of the covariance matrix) to the real value `f` (default value 0.0).
- `-B` Output HMM definition files in binary format.
- `-F fmt` Set the source data format to `fmt`.
- `-G fmt` Set the label file format to `fmt`.

- H *mmf* Load HMM macro model file *mmf*. This option may be repeated to load multiple MMFs.
- I *mlf* This loads the master label file *mlf*. This option may be repeated to load several MLFs.
- L *dir* Search directory *dir* for label files (default is to search current directory).
- M *dir* Store output HMM macro model files in the directory *dir*. If this option is not given, the new HMM definition will overwrite the existing one.
- X *ext* Set label file extension to *ext* (default is *lab*).

HCOMPV also supports the standard options *-A*, *-C*, *-D*, *-S*, *-T*, and *-V* as described in section 4.4.

14.2.3 Tracing

HCOMPV supports the following trace options where each trace flag is given using an octal base

00001 basic progress reporting.

00002 show covariance matrices.

00004 trace data loading.

00010 list label segments.

Trace flags are set using the *-T* option or the *TRACE* configuration variable.

14.3 HCopy

14.3.1 Function

This program will copy one or more data files to a designated output file, optionally converting the data into a parameterised form. While the source files can be in any supported format, the output format is always HTK. By default, the whole of the source file is copied to the target but options exist to only copy a specified segment. Hence, this program is used to convert data files in other formats to the HTK format, to concatenate or segment data files, and to parameterise the result. If any option is set which leads to the extraction of a segment of the source file rather than all of it, then segments will be extracted from all source files and concatenated to the target.

Labels will be copied/concatenated if any of the options indicating labels are specified (`-i -l -x -G -I -L -P -X`). In this case, each source data file must have an associated label file, and a target label file is created. The name of the target label file is the root name of the target data file with the extension `.lab`, unless the `-X` option is used. This new label file will contain the appropriately copied/truncated/concatenated labels to correspond with the target data file; all start and end boundaries are recalculated if necessary.

When used in conjunction with HSLAB, HCopy provides a facility for tasks such as cropping silence surrounding recorded utterances. Since input files may be coerced, HCopy can also be used to convert the parameter kind of a file, for example from WAVEFORM to MFCC, depending on the configuration options. Not all possible conversions can actually be performed; see Table 14.1 for a list of valid conversions. Conversions must be specified via a configuration file as described in chapter 5. Note also that the parameterisation qualifier `_N` cannot be used when saving files to disk, and is meant only for on-the-fly parameterisation.

14.3.2 Use

HCOPY is invoked by typing the command line

```
HCOPY [options] sa1 [ + sa2 + ... ] ta [ sb1 [ + sb2 + ... ] tb ... ]
```

This causes the contents of the one or more source files `sa1`, `sa2`, ... to be concatenated and the result copied to the given target file `ta`. To avoid the overhead of reinvoking the tool when processing large databases, multiple sources and targets may be specified, for example

```
HCOPY srcA.wav + srcB.wav tgtAB.wav srcC.wav tgtD.wav
```

will create two new files `tgtAB.wav` and `tgtD.wav`. HCopy takes file arguments from a script specified using the `-S` option exactly as from the command line, except that any newlines are ignored.

The allowable options to HCopy are as follows where all times and durations are given in 100 ns units and are written as floating-point numbers.

- `-a i` Use level `i` of associated label files with the `-n` and `-x` options. Note that this is not the same as using the `TRANSLEVEL` configuration variable since the `-a` option still allows all levels to be copied through to the output files.
- `-e f` End copying from the source file at time `f`. The default is the end of the file. If `f` is negative or zero, it is interpreted as a time relative to the end of the file, while a positive value indicates an absolute time from the start of the file.
- `-i mlf` Output label files to master file `mlf`.
- `-l s` Output label files to the directory `s`. The default is to output to the current directory.
- `-m t` Set a margin of duration `t` around the segments defined by the `-n` and `-x` options.
- `-n i [j]` Extract the speech segment corresponding to the `i`'th label in the source file. If `j` is specified, then the segment corresponding to the sequence of labels `i` to `j` is extracted. Labels are numbered from their position in the label file. A negative index can be used to count from the end of the label list. Thus, `-n 1 -1` would specify the segment starting at the first label and ending at the last.
- `-s f` Start copying from the source file at time `f`. The default is 0.0, ie the beginning of the file.

- t n Set the line width to n chars when formatting trace output.
- x s [n] Extract the speech segment corresponding to the first occurrence of label s in the source file. If n is specified, then the n'th occurrence is extracted. If multiple files are being concatenated, segments are extracted from each file in turn, and the label must exist for each concatenated file.
- F fmt Set the source data format to fmt.
- G fmt Set the label file format to fmt.
- I mlf This loads the master label file mlf. This option may be repeated to load several MLFs.
- L dir Search directory dir for label files (default is to search current directory).
- O fmt Set the target data format to fmt.
- P fmt Set the target label format to fmt.
- X ext Set label file extension to ext (default is lab).

HCOPY also supports the standard options -A, -C, -D, -S, -T, and -V as described in section 4.4.

Note that the parameter kind conversion mechanisms described in chapter 5 will be applied to all source files. In particular, if an automatic conversion is requested via the configuration file, then HCopy will copy or concatenate the converted source files, not the actual contents. Similarly, automatic byte swapping may occur depending on the source format and the configuration variable BYTEORDER. Because the sampling rate may change during conversions, the options that specify a position within a file i.e. -s and -e use absolute times rather than sample index numbers. All times in HTK are given in units of 100ns and are written as floating-point numbers. To save writing long strings of zeros, standard exponential notation may be used, for example -s 1E6 indicates a start time of 0.1 seconds from the beginning of the file.

	<i>Outputs</i>								
	W	A	V	E	F	O	R	M	D
			L						
			P						
			C					M	
			E	I				E	
			P	R			F	L	
			R	S	E	M	B	S	U
	L	E	T	E	F	A	P	S	E
	P	F	R	F	C	N	E	E	T
<i>Inputs</i>	M	C	C	A	C	C	K	C	R
WAVEFORM	✓	✓	✓	✓	✓	✓	✓	✓	✓
LPC		✓	✓	✓	✓				✓
LPREFC		✓	✓	✓	✓				✓
LPCEPSTRA		✓	✓	✓	✓				✓
IREFC		✓	✓	✓	✓				✓
MFCC						✓			✓
FBANK						✓	✓		✓
MELSPEC						✓	✓	✓	✓
USER									✓
DISCRETE									✓

Table 14.1 Valid Parameter Conversions

Note that truncations are performed *after* any desired coding, which may result in a loss of time resolution if the target file format has a lower sampling rate. Also, because of windowing effects, truncation, coding, and concatenation operations are not necessarily interchangeable. If in doubt, perform all truncation/concatenation in the waveform domain and then perform parameterisation as a last, separate invocation of HCopy.

14.3.3 Trace Output

HCOPY supports the following trace options where each trace flag is given using an octal base

00001 basic progress reporting.

00002 source and target file formats and parameter kinds.

00004 segment boundaries computed from label files.

00010 display memory usage after processing each file.

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

14.4 HDMAN

14.4.1 Function

The HTK tool HDMAN is used to prepare a pronunciation dictionary from one or more sources. It reads in a list of *editing* commands from a script file and then outputs an edited and merged copy of one or more dictionaries.

Each source pronunciation dictionary consists of comment lines and definition lines. Comment lines start with the # character (or optionally any one of a set of specified comment chars) and are ignored by HDMAN. Each definition line starts with a word and is followed by a sequence of symbols (phones) that define the pronunciation. The words and the phones are delimited by spaces or tabs, and the end of line delimits each definition.

Dictionaries used by HDMAN are read using the standard HTK string conventions (see section 4.6), however, the command IR can be used in a HDMAN source edit script to switch to using this raw format. Note that in the default mode, words and phones should not begin with unmatched quotes (they should be escaped with the backslash). All dictionary entries must already be alphabetically sorted before using HDMAN.

Each edit command in the script file must be on a separate line. Lines in the script file starting with a # are comment lines and are ignored. The commands supported are listed below. They can be displayed by HDMAN using the -Q option.

When no edit files are specified, HDMAN simply merges all of the input dictionaries and outputs them in sorted order. All input dictionaries must be sorted. Each input dictionary *xxx* may be processed by its own private set of edit commands stored in *xxx.ded*. Subsequent to the processing of the input dictionaries by their own unique edit scripts, the merged dictionary can be processed by commands in *global.ded* (or some other specified global edit file name).

Dictionaries are processed on a word by word basis in the order that they appear on the command line. Thus, all of the pronunciations for a given word are loaded into a buffer, then all edit commands are applied to these pronunciations. The result is then output and the next word loaded.

Where two or more dictionaries give pronunciations for the same word, the default behaviour is that only the first set of pronunciations encountered are retained and all others are ignored. An option exists to override this so that all pronunciations are concatenated.

Dictionary entries can be filtered by a word list such that all entries not in the list are ignored. Note that the word identifiers in the word list should match exactly (e.g. same case) their corresponding entries in the dictionary.

The edit commands provided by HDMAN are as follows

- AS A B ... Append silence models A, B, etc to each pronunciation.
- CR X A Y B Replace phone Y in the context of A.B by X. Contexts may include an asterisk * to denote any phone or a defined context set defined using the DC command.
- DC X A B ... Define the set A, B, ... as the context X.
- DD X A B ... Delete the definition for word X starting with phones A, B, ...
- DP A B C ... Delete any occurrences of phones A or B or C ...
- DS src Delete each pronunciation from source *src* unless it is the only one for the current word.
- DW X Y Z ... Delete words (& definitions) X, Y, Z, ...
- FW X Y Z ... Define X, Y, Z, ... as function words and change each phone in the definition to a function word specific phone. For example, in word W phone A would become W.A.
- IR Set the input mode to raw. In raw mode, words are regarded as arbitrary sequences of printing chars. In the default mode, words are strings as defined in section 4.6.
- LC [X] Convert all phones to be left-context dependent. If X is given then the 1st phone *a* in each word is changed to X-*a* otherwise it is unchanged.
- LP Convert all phones to lowercase.
- LW Convert all words to lowercase.

- MP X A B ... Merge any sequence of phones A B ... and rename as X.
- RC [X] Convert all phones to be right-context dependent. If X is given then the last phone z in each word is changed to z+X otherwise it is unchanged.
- RP X A B ... Replace all occurrences of phones A or B ... by X.
- RS system Remove stress marking. Currently the only stress marking system supported is that used in the dictionaries produced by Carnegie Mellon University (system = cmu).
- RW X A B ... Replace all occurrences of word A or B ... by X.
- SP X A B ... Split phone X into the sequence A B C ...
- TC [X [Y]] Convert phones to triphones. If X is given then the first phone a is converted to X-a+b otherwise it is unchanged. If Y is given then the last phone z is converted to y-z+Y otherwise if X is given then it is changed to y-z+X otherwise it is unchanged.
- UP Convert all phones to uppercase.
- UW Convert all words to uppercase.

14.4.2 Use

HDMAN is invoked by typing the command line

```
HDMan [options] newDict srcDict1 srcDict2 ...
```

This causes HDMAN read in the source dictionaries `srcDict1`, `srcDict2`, etc. and generate a new dictionary `newDict`. The available options are

- a s Each character in the string s denotes the start of a comment line. By default there is just one comment character defined which is #.
- b s Define s to be a word boundary symbol.
- e dir Look for edit scripts in the directory dir.
- g f File f holds the global edit script. By default, HDMAN expects the global edit script to be called `global.ded`.
- h i j Skip the first i lines of the j'th listed source dictionary.
- i Include word output symbols in the output dictionary.
- l s Write a log file to s. The log file will include dictionary statistics and a list of the number of occurrences of each phone.
- m Merge pronunciations from all source dictionaries. By default, HDMAN generates a single pronunciation for each word. If several input dictionaries have pronunciations for a word, then the first encountered is used. Setting this option causes all distinct pronunciations to be output for each word.
- n f Output a list of all distinct phones encountered to file f.
- o Disable dictionary output.
- p f Load the phone list stored in file f. This enables a check to be made that all output phones are in the supplied list. You need to create a log file (-l) to view the results of this check.
- t Tag output words with the name of the source dictionary which provided the pronunciation.
- w f Load the word list stored in file f. Only pronunciations for the words in this list will be extracted from the source dictionaries.
- Q Print a summary of all commands supported by this tool.

HDMAN also supports the standard options -A, -C, -D, -S, -T, and -V as described in section 4.4.

14.4.3 Tracing

HDMAN supports the following trace options where each trace flag is given using an octal base

00001 basic progress reporting

00002 word buffer operations

00004 show valid inputs

00010 word level editing

00020 word level editing in detail

00040 print edit scripts

00100 new phone recording

00200 pron deletions

00400 word deletions

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

14.5 HEAdapt

14.5.1 Function

This program is used to perform adaptation of a set of HMMs using either maximum likelihood linear regression (MLLR), maximum a-posteriori (MAP) or both. The default is MLLR. In order to perform the adaptation, the first stage requires a state/frame alignment. As such the initial operation of HEADAPT follows HEREST closely. The adaptation training data consists of one or more utterances each of which has a transcription in the form of a standard label file (segment boundaries are ignored). For each training utterance, a composite model is effectively synthesised by concatenating the phoneme models given by the transcription. Each mixture component's accumulators in HEADAPT are updated simultaneously by performing a standard Baum-Welch pass over each training utterance using the composite model. HEADAPT will also prune the α and β matrices, just as HEREST.

HEADAPT is intended to operate on HMMs which have been fully trained using HCOMPV, HINIT, HREST, HEREST. HEADAPT supports multiple mixture diagonal covariance Gaussian HMMs only (i.e. PLAINHS and SHAREDHS systems), with a single data stream only, and parameter tying within and between models. HEADAPT also supports tee-models (see section 7.7), for handling optional silence and non-speech sounds. These may be placed between the units (typically words or phones) listed in the transcriptions, but they cannot be used at the start or end of a transcription. Furthermore, chains of tee-models are not permitted.

After accumulating statistics, HEADAPT estimates the mean and (optionally) the variance transforms. HEADAPT will output either the adapted HMM set (as an MMF), or a transform model file (TMF). The TMF can then be applied to the original model set (for instance when using HVITE). Note that with MAP adaptation a transform is not available and a full HMM set must be output.

When HEADAPT is being run to calculate multiple regression transforms, the model set being adapted must contain a regression class tree. The regression class tree is constructed using the RC edit command in HHED.

14.5.2 Use

HEADAPT is invoked via the command line

```
HEAdapt [options] hmmList adaptFile ...
```

This causes the set of HMMs given in `hmmList` to be loaded. The given list of adaptation training files is then used to perform one adaptation cycle. As always, the list of training files can be stored in a script file if required. On completion, HEADAPT outputs new updated versions of each HMM definition or a new transform models file.

The detailed operation of HEADAPT is controlled by the following command line options

- b *N* Set the number of blocks to be used in the block diagonal matrix representation of the mean transformation. This option will override the config setting `HADAPT:BLOCKS`, or if this is not set the default number of blocks is 1.
- c *f* Set the minimum forward probability fixed distance for the alpha pruning to *f*. Restrict the computation of the α values to just those for which the total log likelihood $\alpha_j(t)\beta_j(t)$ is within distance *f* of the total likelihood (default 10.0).
- d *dir* Normally HEADAPT looks for HMM definitions (not already loaded via MMF files) in the current directory. This option tells HEADAPT to look in the directory *dir* to find them.
- f *field desc* Set the description field *field* in the transform model file to *desc*. Currently the choices for *field* are `uid`, `uname`, `chan` and `desc`.
- g Perform global adaptation only.
- i *N* Update the transforms (incrementally) after accumulating statistics every *N* utterances. The default operation is static adaptation, i.e. after seeing ALL the adaptation data.
- j *f* MAP adaptation with scaling factor *f*. The default operation is MLLR adaptation. If MAP adaptation is to be performed the default value of *f* is 15.0
- k Use MLLR to transform the HMM model set before performing MAP

- m *f* Set the minimum threshold occupation count for a regression class to *f*. A separate regression transformation will be generated at the lowest level in the tree for which there is sufficient occupancy (data). This option will override the config setting `HADAPT:OCCTHRESH`. The default setting is 700.0.
- o *ext* This causes the file name extensions of the original models (if any) to be replaced by *ext*.
- t *f* [*i* *l*] Set the pruning threshold to *f*. During the backward probability calculation, at each time *t* all (log) β values falling more than *f* below the maximum β value at that time are ignored. During the subsequent forward pass, (log) α values are only calculated if there are corresponding valid β values. Furthermore, if the ratio of the $\alpha\beta$ product divided by the total probability (as computed on the backward pass) falls below a fixed threshold then those values of α and β are ignored. Setting *f* to zero disables pruning (default value 0.0). Tight pruning thresholds can result in HEADAPT failing to process an utterance. If the *i* and *l* options are given, then a pruning error results in the threshold being increased by *i* and utterance processing restarts. If errors continue, this procedure will be repeated until the limit *l* is reached.
- u *flags* By default HEADAPT creates transforms for the means only. This option causes the parameters indicated by the *flags* to be created; this argument is a string containing one or more of the letters *m* (mean) and *v* (variance). The presence of a letter enables the creation of the corresponding part of the transform.
- w *f* This sets the minimum variance (i.e. diagonal element of the covariance matrix) to the real value *f* (default value 0.0).
- x *ext* By default, HEADAPT expects a HMM definition for the label *X* to be stored in a file called *X*. This option causes HEADAPT to look for the HMM definition in the file *X.ext*.
 - B Output the HMM definition files in binary format. If outputting a *tmf*, then this option specifies binary output for the *tmf*.
- F *fmt* Set the source data format to *fmt*.
- G *fmt* Set the label file format to *fmt*.
- H *mmf* Load HMM macro model file *mmf*. This option may be repeated to load multiple MMFs.
- I *m1f* This loads the master label file *m1f*. This option may be repeated to load several MLFs.
- J *tmf* Load a transform set from the transform model file *tmf*. The *tmf* is used to transform the *mmf* before performing the state/frame alignment, and a transform is calculated based on this state/frame alignment and the *mmf*.
- K *tmf* Save the transform set in the transform model file *tmf*.
- L *dir* Search directory *dir* for label files (default is to search current directory).
- M *dir* Store output HMM macro model files in the directory *dir*. If this option is not given, the new HMM definition will overwrite the existing one.
- X *ext* Set label file extension to *ext* (default is *lab*).

HEADAPT also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

14.5.3 Tracing

HEADAPT supports the following trace options where each trace flag is given using an octal base 00001 basic progress reporting.

00002 show the logical/physical HMM map.

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

The library that HEADAPT utilises, called HADAPT supports other useful trace options. For library modules, tracing has to be performed via the config file and the module name must prefix the trace (e.g `HADAPT:TRACE=0001`). The following are HADAPT trace options where each trace flag is given using an octal base

- 00001 basic progress reporting.
- 00002 trace on the accumulations.
- 00004 trace on the transformations.
- 00010 output the auxiliary function score.
- 00020 regression classes input/output tracing.
- 00040 regression class tree usage.
- 00200 detailed trace for accumulations at the class level.

14.6 HERest

14.6.1 Function

This program is used to perform a single re-estimation of the parameters of a set of HMMs using an *embedded training* version of the Baum-Welch algorithm. Training data consists of one or more utterances each of which has a transcription in the form of a standard label file (segment boundaries are ignored). For each training utterance, a composite model is effectively synthesised by concatenating the phoneme models given by the transcription. Each phone model has the same set of accumulators allocated to it as are used in HRest but in HEREST they are updated simultaneously by performing a standard Baum-Welch pass over each training utterance using the composite model.

HEREST is intended to operate on HMMs with initial parameter values estimated by HInit/HRest. HEREST supports multiple mixture Gaussians, discrete and tied-mixture HMMs, multiple data streams, parameter tying within and between models, and full or diagonal covariance matrices. HEREST also supports tee-models (see section 7.7), for handling optional silence and non-speech sounds. These may be placed between the units (typically words or phones) listed in the transcriptions but they cannot be used at the start or end of a transcription. Furthermore, chains of tee-models are not permitted.

HEREST includes features to allow parallel operation where a network of processors is available. When the training set is large, it can be split into separate chunks that are processed in parallel on multiple machines/processors, consequently speeding up the training process.

Like all re-estimation tools, HEREST allows a floor to be set on each individual variance by defining a variance floor macro for each data stream (see chapter 8).

HEREST also supports *single pass retraining*. Given a set of well-trained models, a set of new models using a different parameterisation of the training data can be generated in a single pass. This is done by computing the forward and backward probabilities using the original well-trained models and the original training data, but then switching to a new set of training data to compute the new parameter estimates.

HEREST operates in two distinct stages.

1. In the first stage, one of the following two options applies
 - (a) Each input data file contains training data which is processed and the accumulators for state occupation, state transition, means and variances are updated.
 - (b) Each data file contains a dump of the accumulators produced by previous runs of the program. These are read in and added together to form a single set of accumulators.
2. In the second stage, one of the following options applies
 - (a) The accumulators are used to calculate new estimates for the HMM parameters.
 - (b) The accumulators are dumped into a file.

Thus, on a single processor the default combination 1(a) and 2(a) would be used. However, if N processors are available then the training data would be split into N equal groups and HEREST would be set to process one data set on each processor using the combination 1(a) and 2(b). When all processors had finished, the program would then be run again using the combination 1(b) and 2(a) to load in the partial accumulators created by the N processors and do the final parameter re-estimation. The choice of which combination of operations HEREST will perform is governed by the `-p` option switch as described below.

As a further performance optimisation, HEREST will also prune the α and β matrices. By this means, a factor of 3 to 5 speed improvement and a similar reduction in memory requirements can be achieved with negligible effects on training performance (see the `-t` option below).

14.6.2 Use

HEREST is invoked via the command line

```
HERest [options] hmmList trainFile ...
```

This causes the set of HMMs given in `hmmList` to be loaded. The given list of training files is then used to perform one re-estimation cycle. As always, the list of training files can be stored in a script file if required. On completion, HEREST outputs new updated versions of each HMM definition. If the number of training examples falls below a specified threshold for some particular HMM, then the new parameters for that HMM are ignored and the original parameters are used instead.

The detailed operation of HEREST is controlled by the following command line options

- c *f* Set the threshold for tied-mixture observation pruning to *f*. For tied-mixture TIEDHS systems, only those mixture component probabilities which fall within *f* of the maximum mixture component probability are used in calculating the state output probabilities (default 10.0).
- d *dir* Normally HEREST looks for HMM definitions (not already loaded via MMF files) in the current directory. This option tells HEREST to look in the directory *dir* to find them.
- m *N* Set the minimum number of training examples required for any model to *N*. If the actual number falls below this value, the HMM is not updated and the original parameters are used for the new version (default value 3).
- o *ext* This causes the file name extensions of the original models (if any) to be replaced by *ext*.
- p *N* This switch is used to set parallel mode operation. If *p* is set to a positive integer *N*, then HEREST will process the training files and then dump all the accumulators into a file called `HERN.acc`. If *p* is set to 0, then it treats all file names input on the command line as the names of `.acc` dump files. It reads them all in, adds together all the partial accumulations and then re-estimates all the HMM parameters in the normal way.
- r This enables single-pass retraining. The list of training files is processed pair-by-pair. For each pair, the first file should match the parameterisation of the original model set. The second file should match the parameterisation of the required new set. All speech input processing is controlled by configuration variables in the normal way except that the variables describing the old parameterisation are qualified by the name `HPARM1` and the variables describing the new parameterisation are qualified by the name `HPARM2`. The stream widths for the old and the new must be identical.
- s *file* This causes statistics on occupation of each state to be output to the named file. This file is needed for the `RO` command of HHEd but it is also generally useful for assessing the amount of training material available for each HMM state.
- t *f* [*i* *l*] Set the pruning threshold to *f*. During the backward probability calculation, at each time *t* all (log) β values falling more than *f* below the maximum β value at that time are ignored. During the subsequent forward pass, (log) α values are only calculated if there are corresponding valid β values. Furthermore, if the ratio of the $\alpha\beta$ product divided by the total probability (as computed on the backward pass) falls below a fixed threshold then those values of α and β are ignored. Setting *f* to zero disables pruning (default value 0.0). Tight pruning thresholds can result in HEREST failing to process an utterance. if the *i* and *l* options are given, then a pruning error results in the threshold being increased by *i* and utterance processing restarts. If errors continue, this procedure will be repeated until the limit *l* is reached.
- u *flags* By default, HEREST updates all of the HMM parameters, that is, means, variances, mixture weights and transition probabilities. This option causes just the parameters indicated by the *flags* argument to be updated, this argument is a string containing one or more of the letters *m* (mean), *v* (variance), *t* (transition) and *w* (mixture weight). The presence of a letter enables the updating of the corresponding parameter set.
- v *f* This sets the minimum variance (i.e. diagonal element of the covariance matrix) to the real value *f* (default value 0.0).
- w *f* Any mixture weight which falls below the global constant `MINMIX` is treated as being zero. When this parameter is set, all mixture weights are floored to *f* * `MINMIX`.
- x *ext* By default, HEREST expects a HMM definition for the label *X* to be stored in a file called `X`. This option causes HEREST to look for the HMM definition in the file `X.ext`.

- B Output HMM definition files in binary format.
- F *fmt* Set the source data format to *fmt*.
- G *fmt* Set the label file format to *fmt*.
- H *mmf* Load HMM macro model file *mmf*. This option may be repeated to load multiple MMFs.
- I *mlf* This loads the master label file *mlf*. This option may be repeated to load several MLFs.
- L *dir* Search directory *dir* for label files (default is to search current directory).
- M *dir* Store output HMM macro model files in the directory *dir*. If this option is not given, the new HMM definition will overwrite the existing one.
- X *ext* Set label file extension to *ext* (default is *lab*).

HEREST also supports the standard options *-A*, *-C*, *-D*, *-S*, *-T*, and *-V* as described in section 4.4.

14.6.3 Tracing

HEREST supports the following trace options where each trace flag is given using an octal base

00001 basic progress reporting.

00002 show the logical/physical HMM map.

00004 report statistics on pruning.

00010 show the alpha/beta matrices.

00020 show the occupation counters.

00040 show the transition counters.

00100 show the mixture weight counters.

00200 show the calculation of the output probabilities.

00400 list the updated model parameters.

01000 show the average percentage utilisation of tied mixture components.

Trace flags are set using the *-T* option or the *TRACE* configuration variable.

14.7 HHEd

14.7.1 Function

HHEd is a script driven editor for manipulating sets of HMM definitions. Its basic operation is to load in a set of HMMs, apply a sequence of edit operations and then output the transformed set. HHEd is mainly used for applying tyings across selected HMM parameters. It also has facilities for cloning HMMs, clustering states and editing HMM structures.

Many HHEd commands operate on sets of similar items selected from the set of currently loaded HMMs. For example, it is possible to define a set of all final states of all vowel models, or all mean vectors of all mixture components within the model X, etc. Sets such as these are defined by item lists using the syntax rules given below. In all commands, all of the items in the set defined by an item list must be of the same type where the possible types are

s	- state	t	- transition matrix
p	- pdf	w	- stream weights
m	- mixture component	d	- duration parameters
u	- mean vector	x	- transform matrix
v	- variance vector	i	- inverse covariance matrix
h	- HMM definition		

Most of the above correspond directly to the tie points shown in Fig 7.8. There is just one exception. The type “p” corresponds to a pdf (ie a sum of Gaussian mixtures). Pdf’s cannot be tied, however, they can be named in a Tie (TI) command (see below) in which case, the effect is to join all of the contained mixture components into one pool of mixtures and then all of the mixtures in the pool are shared across all pdf’s. This allows conventional *tied-mixture* or *semi-continuous* HMM systems to be constructed.

The syntax rules for item lists are as follows. An item list consists of a comma separated list of item sets.

```
itemList      =  "{" itemSet { "," itemSet } "}"
```

Each `itemSet` consists of the name of one or more HMMs (or a pattern representing a set of HMMs) followed by a specification which represents a set of paths down the parameter hierarchy each terminating at one of the required parameter items.

```
itemSet       =  hmmName . ["transP" | "state" state ]
hmmName       =  ident | identList
identList     =  "(" ident { "," ident } ")"
ident         =  < char | metachar >
metachar      =  "?" | "*"
```

A `hmmName` consists of a single `ident` or a comma separated list of `ident`’s. The following examples are all valid `hmmName`’s:

```
aa three model001 (aa,iy,ah,uh) (one,two,three)
```

In addition, an `ident` can contain the metacharacter “?” which matches any single character and the metacharacter “*” which matches a string of zero or more characters. For example, the item list

```
{*-aa+*.transP}
```

would represent the set of transition matrices of all loaded triphone variations of `aa`.

Items within states require the state indices to be specified

```
state         =  index [ "." stateComp ]
index         =  "[" intRange { "," intRange } "]"
intRange      =  integer [ "-" integer ]
```

For example, the item list

```
{*.state[1,3-5,9]}
```

represents the set of all states 1, 3 to 5 inclusive and 9 of all currently loaded HMMs. Items within states include durational parameters, stream weights, pdf’s and all items within mixtures

```
stateComp = "dur" | "weights" | [ " stream" index ] "." "mix" [ mix ]
```

For example,

```
{(aa,ah,ax).state[2].dur}
```

denotes the set of durational parameter vectors from state 2 of the HMMs **aa**, **ah** and **ax**. Similarly,

```
{*.state[2-4].weights}
```

denotes the set of stream weights for states 2 to 4 of all currently loaded HMMs. The specification of pdf's may optionally include a list of the relevant streams, if omitted, stream 1 is assumed. For example,

```
{three.state[3].mix}
```

and

```
{three.state[3].stream[1].mix}
```

both denote a list of the single pdf belonging to stream 1 of state 3 of the HMM **three**.

Within a pdf, the possible item types are mixture components, mean vectors, and the various possible forms of covariance parameters

```
mix = index [ "." ( "mean" | "cov" ) ]
```

For example,

```
{*.state[2].mix[1-3]}
```

denotes the set of mixture components 1 to 3 from state 2 of all currently loaded HMMs and

```
{(one,two).state[4].stream[3].mix[1].mean}
```

denotes the set of mean vectors from mixture component 1, stream 3, state 4 of the HMMs **one** and **two**. When **cov** is specified, the type of the covariance item referred to is determined from the **CovKind** of the loaded models. Thus, for diagonal covariance models, the item list

```
{*.state[2-4].mix[1].cov}
```

would denote the set of variance vectors for mixture 1, states 2 to 4 of all loaded HMMs.

Note finally, that it is not an error to specify non-existent models, states, mixtures, etc. All item list specifications are regarded as patterns which are matched against the currently loaded set of models. All and only those items which match are included in the set. However, both a null result and a set of items of mixed type do result in errors.

All HHEd commands consist of a 2 character command name followed by zero or more arguments. In the following descriptions, item lists are shown as **itemList(c)** where the character **c** denotes the type of item expected by that command. If this type indicator is missing then the command works for all item types.

The HHEd commands are as follows

AT i j prob itemList(t)

Add a transition from state **i** to state **j** with probability **prob** for all transition matrices in **itemList**. The remaining transitions out of state **i** are rescaled so that $\sum_k a_{ik} = 1$. For example,

```
AT 1 3 0.1 {*.transP}
```

would add a skip transition to all loaded models from state 1 to state 3 with probability 0.1.

AU hmmList

Use a set of decision trees to create a new set of models specified by the **hmmList**. The decision trees may be made as a result of either the **TB** or **LT** command.

Each model in **hmmList** is constructed in the following manner. If a model with the same logical name already exists in the current HMM set this is used unchanged, otherwise the model is synthesised from the decision trees. If the trees cluster at the model level the synthesis results in a logical model sharing the physical model from the tree that matches the new context. If the clustering was performed at the state level a prototype model (an example of the same phone model occurring in a different context) is found and a new HMM is constructed that shares the transition matrix with the prototype model but consists of tied states selected using the decision tree.

CL `hmmList`

Clone a HMM list. The file `hmmList` should hold a list of HMMs all of whose logical names are either the same as, or are context-dependent versions of the currently loaded set of HMMs. For each name in `hmmList`, the corresponding HMM in the loaded set is cloned. On completion, the currently loaded set is discarded and replaced by the new set. For example, if the file `mylist` contained

```
A-A+A
A-A+B
B-A+A
B-B+B
B-B+A
```

and the currently loaded HMMs were just A and B, then A would be cloned 3 times to give the models A-A+A, A-A+B and B-A+A, and B would be cloned 2 times to give B-B+B and B-B+A. On completion, the original definitions for A and B would be deleted (they could be retained by including them in the new `hmmList`).

CO `newList`

Compact a set of HMMs. The effect of this command is to scan the currently loaded set of HMMs and identify all identical definitions. The physical name of the first model in each identical set is then assigned to all models in that set and all model definitions are replaced by a pointer to the first model definition. On completion, a new list of HMMs which includes the new model tyings is written out to file `newList`. For example, suppose that models A, B, C and D were currently loaded and A and B were identical. Then the command

```
CO tlist
```

would tie HMMs A and B, set the physical name of B to A and output the new HMM list

```
A
B A
C
D
```

to the file `tlist`. This command is used mainly after performing a sequence of parameter tying commands.

DP `s n id ...`

Duplicates a set of HMMs. This command is used to replicate a set of HMMs whilst allowing control over which structures will be shared between them. The first parameter controls duplication of tied structures. Any macros whose type appears in string `s` are duplicated with new names and only used in the duplicate model set. The remaining shared structures are common through all the model sets (original and duplicates). The second parameter defines the number of times the current HMM set should be duplicated with the remaining `n` parameters providing suffices to make the original macro identifiers unique in each duplicated HMM set.

For instance the following script could be used to duplicate a set of tied state models to produce gender dependent ones with tied variances.

```
MM "v_" { (*).state[2-4].mix[1-2].cov }
DP "v" 2 ":m" ":f"
```

The MM command converts all variances into macros (with each macro referring to only one variance). The DP command then duplicates the current HMM set twice. Each of the duplicate sets will share the tied variances with the original set but will have new mixture means, weights and state macros. The new macro names will be constructed by appending the id `":m"` or `":f"` to the original macro name whilst the model names have the id appended after the base phone name (so `ax-b+d` becomes `ax-b:m+d` or `ax-b:f+d`).

HK hsetkind

Converts model set from one kind to another. Although `hsetkind` can take the value `PLAINHS`, `SHAREDHS`, `TIEDHS` or `DISCRETEHS`, the `HK` command is most likely to be used when building tied-mixture systems (`hsetkind=TIEDHS`).

JO size minw

Set the size and minimum mixture weight for subsequent `Tie` (`TI`) commands applied to pdf's. The value of `size` sets the total number of mixtures in the tied mixture set (*codebook*) and `minw` sets a floor on the mixture weights as a multiple of `MINMIX`. This command only applies to tying item lists of type "p" (see the `Tie` `TI` command below).

LS statsfile

This command is used to read in the `HEREST` statistics file (see the `HEREST -s` option) stored in `statsfile`. These statistics are needed for certain clustering operations. The statistics file contains the occupation count for every HMM state.

LT treesfile

This command reads in the decision trees stored in `treesfile`. The trees file will consist of a set of questions defining contexts that may appear in the subsequent trees. The trees are used to identify either the state or the model that should be used in a particular context. The file would normally be produced by `ST` after tree based clustering has been performed.

MM macro itemList

This command makes each item ($I=1..N$) in `itemList` into a macro with name `nameI` and a usage of one. This command can prevent unnecessary duplication of structures when HMMs are cloned or duplicated.

MT triList newTriList

Make a set of triphones by merging the currently loaded set of biphones. This is a very specialised command. All currently loaded HMMs must have 3 emitting states and be either left or right context-dependent biphones. The list of HMMs stored in `triList` should contain one or more triphones. For each triphone in `triList` of the form `X-Y+Z`, there must be currently loaded biphones `X-Y` and `Y+Z`. A new triphone `X-Y+Z` is then synthesised by first cloning `Y+Z` and then replacing the state information for the initial emitting state by the state information for the initial emitting state of `X-Y`. Note that the underlying physical names of the biphones used to create the triphones are recorded so that where possible, triphones generated from tied biphones are also tied. On completion, the new list of triphones including aliases is written to the file `newTriList`.

MU m itemList(p)

Increase the number of non-defunct mixture components in each pdf in the `itemList` to `m` (when `m` is just a number) or by `m` (when `m` is a number preceded by a `+` sign). A defunct mixture is one for which the weight has fallen below `MINMIX`. This command works in two steps. Firstly, the weight of each mixture in each pdf is checked. If any defunct mixtures are discovered, then each is successively replaced by a non-defunct mixture component until either the required total number of non-defunct mixtures is reached or there are no defunct mixtures left. This replacement works by first deleting the defunct mixture and then finding the mixture with the largest weight and splitting it. The split operation is as follows. The weight of the mixture component is first halved and then the mixture is cloned. The two identical mean vectors are then perturbed by adding 0.2 standard deviations to one and subtracting the same amount from the other.

In the second step, the mixture component with the largest weight is split as above. This is repeated until the required number of mixture components are obtained. Whenever, a mixture is split, a count is incremented for that mixture so that splitting occurs evenly across the mixtures. Furthermore, a mixture whose *gconst* value falls more than four standard deviations below the mean is not split.

As an example, the command

```
MU 6 {*-aa*.state[3].mix}
```

would increase the number of mixture components in state 3 of all triphones of `aa` to 6.

```
NC N macro itemList(s)
```

N-cluster the states listed in the `itemList` and tie each cluster `i` as macro `macroi` where `i` is 1,2,3,...,N. The set of states in the `itemList` are divided into N clusters using the following furthest neighbour hierarchical cluster algorithm:

```
create 1 cluster for each state;
n = number of clusters;
while (n>N) {
    find i and j for which g(i,j) is minimum;
    merge clusters i and j;
}
```

Here $g(i, j)$ is the inter-group distance between clusters `i` and `j` defined as the maximum distance between any state in cluster `i` and any state in cluster `j`. The calculation of the inter-state distance depends on the type of HMMs involved. Single mixture Gaussians use

$$d(i, j) = \frac{1}{S} \sum_{s=1}^S \left[\frac{1}{V_s} \sum_{k=1}^{V_s} \frac{(\mu_{isk} - \mu_{j sk})^2}{\sigma_{isk} \sigma_{j sk}} \right]^{\frac{1}{2}} \quad (14.1)$$

where V_s is the dimensionality of stream s . Fully tied mixture systems (ie TIEDHS) use

$$d(i, j) = \frac{1}{S} \sum_{s=1}^S \left[\frac{1}{M_s} \sum_{m=1}^{M_s} (c_{ism} - c_{j sm})^2 \right]^{\frac{1}{2}} \quad (14.2)$$

and all others use

$$d(i, j) = -\frac{1}{S} \sum_{s=1}^S \frac{1}{M_s} \sum_{m=1}^{M_s} \log[b_{js}(\mu_{ism})] + \log[b_{is}(\mu_{j sm})] \quad (14.3)$$

where $b_{js}(\mathbf{x})$ is as defined in equation 7.1 for the continuous case and equation 7.3 for the discrete case. The actual tying of the states in each cluster is performed exactly as for the Tie (TI) command below. The macro for the `i`'th tied cluster is called `macroi`.

```
QS name itemList(h)
```

Define a question `name` which is true for all the models in `itemList`. These questions can subsequently be used as part of the decision tree based clustering procedure (see TB command below).

```
RC N identifier [itemlist]
```

This command is used to grow a regression class tree for adaptation purposes. A regression class tree is grown with N terminal or leaf nodes, using the centroid splitting algorithm with a Euclidean distance measure to cluster the model set's mixture components. Hence each leaf node specifies a particular mixture component cluster. The regression class tree is saved with the macro identifier `identifier_N`. Each Gaussian component is also labelled with a regression class number (corresponding to the leaf node number that the Gaussian component resides in). In order to grow the regression class tree it is necessary to load in a `statsfile` using the LS command. It is also possible to specify an `itemlist` containing the "non-speech" sound components such as the silence mixture components. If this is included then the first split made will result in one leaf containing the specified non-speech sound components, while the other leaf will contain the rest of the model set components. Tree construction then continues as usual.

```
RN hmmIdName
```

Rename or add the hmm set identifier in the global options macro to `hmmIdName`.

RM *hmmFile*

Load the *hmm* from *hmmFile* and subtract the mean from state 2, mixture 1 of the model from every loaded model. Every component of the mean is subtracted including deltas and accelerations.

RO *f [statsfile]*

This command is used to remove outlier states during clustering with subsequent *NC* or *TC* commands. If *statsfile* is present it first reads in the *HEREST* statistics file (see *LS*) otherwise it expects a separate *LS* command to have already been used to read in the statistics. Any subsequent *NC*, *TC* or *TB* commands are extended to ensure that the occupancy clusters produced exceeds the threshold *f*. For *TB* this is used to choose which questions are allowed to be used to split each node. Whereas for *NC* and *TC* a final merging pass is used and for as long the smallest cluster count falls below the threshold *f*, then that cluster is merged with its nearest neighbour.

RT *i j itemList(t)*

Remove the transition from state *i* to *j* in all transition matrices given in the *itemList*. After removal, the remaining non-zero transition probabilities for state *i* are rescaled so that $\sum_k a_{ik} = 1$.

SH

Show the current HMM set. This command can be inserted into edit scripts for debugging. It prints a summary of each loaded HMM identifying any tied parameters.

SK *skind*

Change the sample kind of all loaded HMMs to *skind*. This command is typically used in conjunction with the *SW* command. For example, to add delta coefficients to a set of models, the *SW* command would be used to double the stream widths and then this command would be used to add the *_D* qualifier.

SS *N*

Split into *N* independent data streams. This command causes the currently loaded set of HMMs to be converted from 1 data stream to *N* independent data streams. The widths of each stream are determined from the single stream vector size and the sample kind as described in section 5.10. Execution of this command will cause any tyings associated with the split stream to be undone.

ST *filename*

Save the currently defined questions and trees to file *filename*. This allows subsequent construction of models using for new contexts using the *LT* and *AU* commands.

SU *N w1 w2 w3 .. wN*

Split into *N* independent data streams with stream widths as specified. This command is similar to the *SS* command except that the width of each stream is defined explicitly by the user rather than using the built-in stream splitting rules. Execution of this command will cause any tyings associated with the split stream to be undone.

SW *s n*

Change the width of stream *s* of all currently loaded HMMs to *n*. Changing the width of stream involves changing the dimensions of all mean and variance vectors or covariance matrices. If *n* is greater than the current width of stream *s*, then mean vectors are extended with zeroes and variance vectors are extended with 1's. Covariance matrices are extended with zeroes everywhere except for the diagonal elements which are set to 1. This command preserves any tyings which may be in force.

TB **f** macro **itemList(s or h)**

Decision tree cluster all states in the given **itemList** and tie them as **macro_i** where **i** is 1,2,3,... This command performs a top down clustering of the states or models appearing in **itemList**. This clustering starts by placing all items in a single root node and then choosing a question from the current set to split the node in such a way as to maximise the likelihood of a single diagonal covariance Gaussian at each of the child nodes generating the training data. This splitting continues until the increase in likelihood falls below threshold **f** or no questions are available which do not pass the outlier threshold test. This type of clustering is only implemented for single mixture, diagonal covariance untied models.

TC **f** macro **itemList(s)**

Cluster all states in the given **itemList** and tie them as **macro_i** where **i** is 1,2,3,... This command is identical to the **NC** command described above except that the number of clusters is varied such that the maximum within cluster distance is less than the value given by **f**.

TI macro **itemList**

Tie the items in **itemList** and assign them to the specified **macro** name. This command applies to any item type but all of the items in **itemList** must be of the same type. The detailed method of tying depends on the item type as follows:

state(s) the state with the largest total value of **gConst** in stream 1 (indicating broad variances) and the minimum number of defunct mixture weights (see **MU** command) is selected from the item list and all states are tied to this typical state.

transitions(t) all transition matrices in the item list are tied to the last in the list.

mixture(m) all mixture components in the item list are tied to the last in the list.

mean(u) the average vector of all the mean vectors in the item list is calculated and all the means are tied to this average vector.

variance(v) a vector is constructed for which each element is the maximum of the corresponding elements from the set of variance vectors to be tied. All of the variances are then tied to this maximum vector.

covariance(i) all covariance matrices in the item list are tied to the last in the list.

xform(x) all transform matrices in the item list are tied to the last in the list.

duration(d) all duration vectors in the item list are tied to the last in the list.

stream weights(w) all stream weight vectors in the item list are tied to the last in the list.

pdf(p) as noted earlier, pdf's are tied to create tied mixture sets rather than to create a shared pdf. The procedure for tying pdf's is as follows

1. All mixtures from all pdf's in the item list are collected together in order of mixture weight.
2. If the number of mixtures exceeds the join size J [see the **Join (JO)** command above], then all but the first J mixtures are discarded.
3. If the number of mixtures is less than J , then the mixture with the largest weight is repeatedly split until there are exactly J mixture components. The split procedure used is the same as for the **MixUp (MU)** command described above.
4. All pdf's in the item list are made to share all J mixture components. The weight for each mixture is set proportional to the log likelihood of the mean vector of that mixture with respect to the original pdf.
5. Finally, all mixture weights below the floor set by the **Join** command are raised to the floor value and all of the mixture weights are renormalised.

TR n

Change the level of detail for tracing and consists of a number of separate flags which can be added together. Values 0001, 0002, 0004, 0008 have the same meaning as the command line trace level but apply only to a single block of commands (a block consisting of a set of commands of the name). A value of 0010 can be used to show current memory usage.

UT itemList

Untie all items in `itemList`. For each item in the item list, if the usage counter for that item is greater than 1 then it is cloned, the original shared item is replaced by the cloned copy and the usage count of the shared item is reduced by 1. If the usage count is already 1, the associated macro is simply deleted and the usage count set to 0 to indicate an unshared item. Note that it is not possible to untie a pdf since these are not actually shared [see the Tie (TI) command above].

14.7.2 Use

HHED is invoked by typing the command line

```
HHed [options] edCmdFile hmmList
```

where `edCmdFile` is a text file containing a sequence of edit commands as described above and `hmmList` defines the set of HMMs to be edited (see HMODEL for the format of HMM list). If the models are to be kept in separate files rather than being stored in an MMF, the configuration variable `KEEPDISTINCT` should be set to true. The available options for HHED are

- d `dir` This option tells HHED to look in the directory `dir` to find the model definitions.
- o `ext` This causes the file name extensions of the original models (if any) to be replaced by `ext`.
- w `mmf` Save all the macros and model definitions in a single master macro file `mmf`.
- x `s` Set the extension for the edited output files to be `s` (default is to use the original names unchanged).
- z Setting this option causes all aliases in the loaded HMM set to be deleted (zapped) immediately before loading the definitions. The result is that all logical names are ignored and the actual HMM list consists of just the physically distinct HMMs.
- B Output HMM definition files in binary format.
- H `mmf` Load HMM macro model file `mmf`. This option may be repeated to load multiple MMFs.
- M `dir` Store output HMM macro model files in the directory `dir`. If this option is not given, the new HMM definition will overwrite the existing one.
- Q Print a summary of all commands supported by this tool.

HHED also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

14.7.3 Tracing

HHED supports the following trace options where each trace flag is given using an octal base

- 00001 basic progress reporting.
- 00002 intermediate progress reporting.
- 00004 detailed progress reporting.
- 00010 show item lists used for each command.
- 00020 show memory usage.
- 00100 show changes to macro definitions.
- 00200 show changes to stream widths.

00400 show clusters.

00800 show questions.

01000 show tree filtering.

02000 show tree splitting.

04000 show tree merging.

10000 show good question scores.

20000 show all question scores.

40000 show all merge scores.

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

14.8 HInit

14.8.1 Function

HINIT is used to provide initial estimates for the parameters of a single HMM using a set of observation sequences. It works by repeatedly using Viterbi alignment to segment the training observations and then recomputing the parameters by pooling the vectors in each segment. For mixture Gaussians, each vector in each segment is aligned with the component with the highest likelihood. Each cluster of vectors then determines the parameters of the associated mixture component. In the absence of an initial model, the process is started by performing a uniform segmentation of each training observation and for mixture Gaussians, the vectors in each uniform segment are clustered using a modified K-Means algorithm¹.

HINIT can be used to provide initial estimates of whole word models in which case the observation sequences are realisations of the corresponding vocabulary word. Alternatively, HINIT can be used to generate initial estimates of *seed* HMMs for phoneme-based speech recognition. In this latter case, the observation sequences will consist of segments of continuously spoken training material. HINIT will *cut* these out of the training data automatically by simply giving it a segment label.

In both of the above applications, HINIT normally takes as input a prototype HMM definition which defines the required HMM topology i.e. it has the form of the required HMM except that means, variances and mixture weights are ignored². The transition matrix of the prototype specifies both the allowed transitions and their initial probabilities. Transitions which are assigned zero probability will remain zero and hence denote non-allowed transitions. HINIT estimates transition probabilities by counting the number of times each state is visited during the alignment process.

HINIT supports multiple mixtures, multiple streams, parameter tying within a single model, full or diagonal covariance matrices, tied-mixture models and discrete models. The output of HInit is typically input to HRest.

Like all re-estimation tools, HINIT allows a floor to be set on each individual variance by defining a variance floor macro for each data stream (see chapter 8).

14.8.2 Use

HINIT is invoked via the command line

```
HInit [options] hmm trainFiles ...
```

This causes the means and variances of the given `hmm` to be estimated repeatedly using the data in `trainFiles` until either a maximum iteration limit is reached or the estimation converges. The HMM definition can be contained within one or more macro files loaded via the standard `-H` option. Otherwise, the definition will be read from a file called `hmm`. The list of train files can be stored in a script file if required.

The detailed operation of HINIT is controlled by the following command line options

- `-e f` This sets the convergence factor to the real value `f`. The convergence factor is the relative change between successive values of $P_{max}(O|\lambda)$ computed as a by-product of the Viterbi alignment stage (default value 0.0001).
- `-i N` This sets the maximum number of estimation cycles to `N` (default value 20).
- `-l s` The string `s` must be the name of a segment label. When this option is used, HINIT searches through all of the training files and cuts out all segments with the given label. When this option is not used, HINIT assumes that each training file is a single token.
- `-m N` This sets the minimum number of training examples so that if fewer than `N` examples are supplied an error is reported (default value 3).
- `-n` This flag suppresses the initial uniform segmentation performed by HINIT allowing it to be used to update the parameters of an existing model.

¹This algorithm is significantly different from earlier versions of HTK where K-means clustering was used at every iteration and the Viterbi alignment was limited to states

²Prototypes should either have `GConst` set (the value does not matter) to avoid HTK trying to compute it or variances should be set to a positive value such as 1.0 to ensure that `GConst` is computable

- o *s* The string *s* is used as the name of the output HMM in place of the source name. This is provided in `HINIT` since it is often used to initialise a model from a prototype input definition. The default is to use the source name.
- u *flags* By default, `HINIT` updates all of the HMM parameters, that is, means, variances, mixture weights and transition probabilities. This option causes just the parameters indicated by the *flags* argument to be updated, this argument is a string containing one or more of the letters *m* (mean), *v* (variance), *t* (transition) and *w* (mixture weight). The presence of a letter enables the updating of the corresponding parameter set.
- v *f* This sets the minimum variance (i.e. diagonal element of the covariance matrix) to the real value *f*. The default value is 0.0.
- w *f* Any mixture weight or discrete observation probability which falls below the global constant `MINMIX` is treated as being zero. When this parameter is set, all mixture weights are floored to $f * \text{MINMIX}$.
- B Output HMM definition files in binary format.
- F *fmt* Set the source data format to *fmt*.
- G *fmt* Set the label file format to *fmt*.
- H *mmf* Load HMM macro model file *mmf*. This option may be repeated to load multiple MMFs.
- I *mlf* This loads the master label file *mlf*. This option may be repeated to load several MLFs.
- L *dir* Search directory *dir* for label files (default is to search current directory).
- M *dir* Store output HMM macro model files in the directory *dir*. If this option is not given, the new HMM definition will overwrite the existing one.
- X *ext* Set label file extension to *ext* (default is `lab`).

`HINIT` also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

14.8.3 Tracing

`HINIT` supports the following trace options where each trace flag is given using an octal base

- 000001 basic progress reporting.
- 000002 file loading information.
- 000004 segments within each file.
- 000010 uniform segmentation.
- 000020 Viterbi alignment.
- 000040 state alignment.
- 000100 mixture component alignment.
- 000200 count updating.
- 000400 output probabilities.

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

14.9 HLEd

14.9.1 Function

This program is a simple editor for manipulating label files. Typical examples of its use might be to merge a sequence of labels into a single composite label or to expand a set of labels into a context sensitive set. HLEd works by reading in a list of *editing* commands from an edit script file and then makes an edited copy of one or more label files. For multiple level files, edit commands are applied to the *current level* which is initially the first (i.e. 1). Other levels may be edited by moving to the required level using the ML Move Level command.

Each edit command in the script file must be on a separate line. The first two-letter mnemonic on each line is the command name and the remaining letters denote labels³. The commands supported may be divided into two sets. Those in the first set are used to edit individual labels and they are as follows

CH X A Y B Change Y in the context of A_B to X. A and/or B may be a * to match any context, otherwise they must be defined by a DC command (see below). A block of consecutive CH commands are effectively executed in parallel so that the contexts are those that exist before any of the commands in the block are applied.

DC A B C .. define the context A as the set of labels B, C, etc.

DE A B .. Delete any occurrences of labels A or B etc.

FI A Y B Find Y in the context of A_B and count the number of occurrences.

ME X A B .. Merge any sequence of labels A B C etc. and call the new segment X.

ML N Move to label level N.

RE X A B .. Replace all occurrences of labels A or B etc. by the label X.

The commands in the second set perform global operations on whole transcriptions. They are as follows.

DL [N] Delete all labels in the current level. If the optional integer arg is given, then level N is deleted.

EX Expand all labels either from words to phones using the first pronunciation from a dictionary when it is specified on the command line otherwise expand labels of the form A_B.C.D... into a sequence of separate labels A B C D .. This is useful for label formats which include a complete orthography as a single label or for creating a set of sub-word labels from a word orthography for a sub-word based recogniser. When a label is expanded in this way, the label duration is divided into equal length segments. This can only be performed on the root level of a multiple level file.

FG X Mark all unlabelled segments of the input file of duration greater than T_g msec with the label X. The default value for T_g is 50000.0 (=5msecs) but this can be changed using the -g command line option. This command is mainly used for explicitly labelling inter-word silences in data files for which only the actual speech has been transcribed.

IS A B Insert label A at the start of every transcription and B at the end. This command is usually used to insert silence labels.

IT Ignore triphone contexts in CH and FI commands.

LC [X] Convert all phoneme labels to left context dependent. If X is given then the first phoneme label a becomes X-a otherwise it is left unchanged.

NB X The label X (typically a short pause) should be ignored at word boundaries when using the context commands LC, RC and TC.

³In earlier versions of HTK, HLEd command names consisted of a single letter. These are still supported for backwards compatibility and they are included in the command summary produced using the -Q option. However, commands introduced since version 2.0 have two letter names.

- RC [X] Convert all phoneme labels to right context dependent. If X is given then the last phoneme label **z** becomes **z+X** otherwise it is left unchanged.
- SB X Define the label X to be a sentence boundary marker. This label can then be used in context-sensitive change commands.
- SO Sort all labels into time order.
- SP Split multiple levels into multiple alternative label lists.
- TC [X[Y]] Convert all phoneme labels to Triphones, that is left and right context dependent. If X is given then the first phoneme label **a** becomes **X-a+b** otherwise it is left unchanged. If Y is given then the last phoneme label **z** becomes **y-z+Y** otherwise if X is given then it becomes **y-z+X** otherwise it is left unchanged.
- WB X Define X to be an inter-word label. This command affects the operation of the LC, RC and TC commands. The expansion of context labels is blocked wherever an inter-word label occurs.

The source and target label file formats can be defined using the `-G` and `-P` command line arguments. They can also be set using the configuration variables `SOURCELABEL` and `TARGETLABEL`. The default for both cases is the HTK format.

14.9.2 Use

HLEd is invoked by typing the command line

```
HLEd [options] edCmdFile labFiles ..
```

This causes HLEd to be applied to each `labFile` in turn using the edit commands listed in `edCmdFile`. The `labFiles` may be master label files. The available options are

- b Suppress label boundary times in output files.
- d *s* Read a dictionary from file *s* and use this for expanding labels when the EX command is used.
- i *mlf* This specifies that the output transcriptions are written to the master label file *mlf*.
- g *t* Set the minimum gap detected by the FG to be *t* (default 50000.0 = 5msecs). All gaps of shorter duration than *t* are ignored and not labelled.
- l *s* Directory to store output label files (default is current directory). When output is directed to an MLF, this option can be used to add a path to each output file name. In particular, setting the option `-l '*'` will cause a label file named *xxx* to be prefixed by the pattern `"/xxx"` in the output MLF file. This is useful for generating MLFs which are independent of the location of the corresponding data files.
- m Strip all labels to monophones on loading.
- n *fn* This option causes a list of all new label names created to be output to the file *fn*.
- G *fmt* Set the label file format to *fmt*.
- I *mlf* This loads the master label file *mlf*. This option may be repeated to load several MLFs.
- P *fmt* Set the target label format to *fmt*.
- X *ext* Set label file extension to *ext* (default is `lab`).

HLEd also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

14.9.3 Tracing

HLEd supports the following trace options where each trace flag is given using an octal base

000001 basic progress reporting.

000002 edit script details.

000004 general command operation.

000010 change operations.

000020 level split/merge operations.

000040 delete level operation.

000100 edit file input.

000200 memory usage.

000400 dictionary expansion in EX command

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

14.10 HList

14.10.1 Function

This program will list the contents of one or more data sources in any HTK supported format. It uses the full HTK speech input facilities described in chapter 5 and it can thus read data from a waveform file, from a parameter file and direct from an audio source. HLIST provides a dual rôle in HTK. Firstly, it is used for examining the contents of speech data files. For this function, the TARGETKIND configuration variable should not be set since no conversions of the data are required. Secondly, it is used for checking that input conversions are being performed properly. In the latter case, a configuration designed for a recognition system can be used with HLIST to make sure that the translation from the source data into the required observation structure is exactly as intended. To assist this, options are provided to split the input data into separate data streams (-n) and to explicitly list the identity of each parameter in an observation (-o).

14.10.2 Use

HList is invoked by typing the command line

```
HList [options] file ...
```

This causes the contents of each `file` to be listed to the standard output. If no files are given and the source format is `HAUDIO`, then the audio source is listed. The source form of the data can be converted and listed in a variety of target forms by appropriate settings of the configuration variables, in particular `TARGETKIND`⁴.

The allowable options to HLIST are

- d Force each observation to be listed as discrete VQ symbols. For this to be possible the source must be either `DISCRETE` or have an associated VQ table specified via the `VQTABLE` configuration variable.
- e N End listing samples at sample index N.
- h Print the source header information.
- i N Print N items on each line.
- n N Display the data split into N independent data streams.
- o Show the observation structure. This identifies the rôle of each item in each sample vector.
- p Playback the audio. When sourcing from an audio device, this option enables the playback buffer so that after displaying the sampled data, the captured audio is replayed.
- r Print the raw data only. This is useful for exporting a file into a program which can only accept simple character format data.
- s N Start listing samples from sample index N. The first sample index is 0.
- t Print the target header information.
- F `fmt` Set the source data format to `fmt`.

HLIST also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

14.10.3 Tracing

HLIST supports the following trace options where each trace flag is given using an octal base 00001 basic progress reporting.

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

⁴The `TARGETKIND` is equivalent to the `HCOERCE` environment variable used in earlier versions of HTK

14.11 HLStats

14.11.1 Function

This program will read in a HMM list and a set of HTK format transcriptions (label files). It will then compute various statistics which are intended to assist in analysing acoustic training data and generating simple language models for recognition. The specific functions provided by HLSTATS are:

1. number of occurrences of each distinct logical HMM and/or each distinct physical HMM. The list printed can be limited to the N most infrequent models.
2. minimum, maximum and average durations of each logical HMM in the transcriptions.
3. a matrix of bigram probabilities
4. an ARPA/MIT-LL format text file of back-off bigram probabilities
5. a list of labels which cover the given set of transcriptions.

14.11.2 Bigram Generation

When using the bigram generating options, each transcription is assumed to have a unique entry and exit label which by default are !ENTER and !EXIT. If these labels are not present they are inserted. In addition, any label occurring in a transcription which is not listed in the HMM list is mapped to a unique label called !NULL.

HLSTATS processes all input transcriptions and maps all labels to a set of unique integers in the range 1 to L , where L is the number of distinct labels. For each adjacent pair of labels i and j , it counts the total number of occurrences $N(i, j)$. Let the total number of occurrences of label i be $N(i) = \sum_{j=1}^L N(i, j)$.

For matrix bigrams, the bigram probability $p(i, j)$ is given by

$$p(i, j) = \begin{cases} \alpha N(i, j)/N(i) & \text{if } N(i) > 0 \\ 1/L & \text{if } N(i) = 0 \\ f & \text{otherwise} \end{cases}$$

where f is a floor probability set by the `-f` option and α is chosen to ensure that $\sum_{j=1}^L p(i, j) = 1$.

For back-off bigrams, the unigram probabilities $p(i)$ are given by

$$p(i) = \begin{cases} N(i)/N & \text{if } N(i) > u \\ u/N & \text{otherwise} \end{cases}$$

where u is unigram floor count set by the `-u` option and $N = \sum_{i=1}^L \max[N(i), u]$.

The backed-off bigram probabilities are given by

$$p(i, j) = \begin{cases} (N(i, j) - D)/N(i) & \text{if } N(i, j) > t \\ b(i)p(j) & \text{otherwise} \end{cases}$$

where D is a discount and t is a bigram count threshold set by the `-t` option. The discount D is fixed at 0.5 but can be changed via the configuration variable `DISCOUNT`. The back-off weight $b(i)$ is calculated to ensure that $\sum_{j=1}^L p(i, j) = 1$, i.e.

$$b(i) = \frac{1 - \sum_{j \in B} p(i, j)}{1 - \sum_{j \in B} p(j)}$$

where B is the set of all words for which $p(i, j)$ has a bigram.

The formats of matrix and ARPA/MIT-LL format bigram files are described in Chapter 12.

14.11.3 Use

HLSTATS is invoked by the command line

```
HLStats [options] hmmList labFiles ....
```

The `hmmList` should contain a list of all the labels (ie model names) used in the following label files for which statistics are required. Any labels not appearing in the list are ignored and assumed to be out-of-vocabulary. The list of labels is specified in the same way as for a HMM list (see HMODEL) and the logical \Rightarrow physical mapping is preserved to allow statistics to be collected about physical names as well as logical ones. The `labFiles` may be master label files. The available options are

- b `fn` Compute bigram statistics and store result in the file `fn`.
- c `N` Count the number of occurrences of each logical model listed in the `hmmList` and on completion list all models for which there are `N` or less occurrences.
 - d Compute minimum, maximum and average duration statistics for each label.
- f `f` Set the matrix bigram floor probability to `f` (default value 0.0). This option should be used in conjunction with the `-b` option.
- h `N` Set the bigram hashtable size to medium(`N=1`) or large (`N=2`). This option should be used in conjunction with the `-b` option. The default is small(`N=0`).
- l `fn` Output a list of covering labels to file `fn`. Only labels occurring in the `labList` are counted (others are assumed to be out-of-vocabulary). However, this list may contain labels that do not occur in any of the label files. The list of labels written to `fn` will however contain only those labels which occur at least once.
 - o Produce backed-off bigrams rather than matrix ones. These are output in the standard ARPA/MIT-LL textual format.
- p `N` Count the number of occurrences of each physical model listed in the `hmmList` and on completion list all models for which there are `N` or less occurrences.
- s `st en` Set the sentence start and end labels to `st` and `en`. (Default `!ENTER` and `!EXIT`).
- t `n` Set the threshold count for including a bigram in a backed-off bigram language model. This option should be used in conjunction with the `-b` and `-o` options.
- u `f` Set the unigram floor probability to `f` when constructing a back-off bigram language model. This option should be used in conjunction with the `-b` and `-o` options.
- G `fmt` Set the label file format to `fmt`.
- I `mlf` This loads the master label file `mlf`. This option may be repeated to load several MLFs.

HLSTATS also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

14.11.4 Tracing

HLSTATS supports the following trace options where each trace flag is given using an octal base

00001 basic progress reporting.

00002 trace memory usage.

00004 show basic statistics whilst calculating bigrams. This includes the global training data entropy and the entropy for each distinct label.

00010 show file names as they are processed.

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

14.12 HParse

14.12.1 Function

The HPARSE program generates word level lattice files (for use with e.g. HVITE) from a text file syntax description containing a set of rewrite rules based on extended Backus-Naur Form (EBNF). The EBNF rules are used to generate an internal representation of the corresponding finite-state network where HPARSE network nodes represent the words in the network, and are connected via sets of links. This HPARSE network is then converted to HTK V2 word level lattice. The program provides one convenient way of defining such word level lattices.

HPARSE also provides a *compatibility mode* for use with HPARSE syntax descriptions used in HTK V1.5 where the same format was used to define both the word level syntax and the dictionary. In compatibility mode HPARSE will output the word level portion of such a syntax as an HTK V2 lattice file (via HNET) and the pronunciation information as an HTK V2 dictionary file (via HDICT).

The lattice produced by HPARSE will often contain a number of !NULL nodes in order to reduce the number of arcs in the lattice. The use of such !NULL nodes can both reduce size and increase efficiency when used by recognition programs such as HVITE.

14.12.2 Network Definition

The syntax rules for the textual definition of the network are as follows. Each node in the network has a **nodename**. This node name will normally correspond to a word in the final syntax network. Additionally, for use in compatibility mode, each node can also have an external name.

```
name      = char{char}
nodename  = name [ "%" ( "%" | name ) ]
```

Here **char** represents any character except one of the meta chars { } [] < > | = \$ () ; \ / *. The latter may, however, be escaped using a backslash. The first name in a **nodename** represents the name of the node (“internal name”), and the second optional name is the “external” name. This is used only in compatibility mode, and is, by default the same as the internal name.

Network definitions may also contain variables

```
variable  = $name
```

Variables are identified by a leading \$ character. They stand for sub-networks and must be defined before they appear in the RHS of a rule using the form

```
subnet    = variable “=” expr “;”
```

An **expr** consists of a set of alternative sequences representing parallel branches of the network.

```
expr      = sequence { “|” sequence }
sequence  = factor{factor}
```

Each sequence is composed of a sequence of factors where a factor is either a node name, a variable representing some sub-network or an expression contained within various sorts of brackets.

```
factor    = “(” expr “)”
           “{” expr “}”
           “<” expr “>”
           “[” expr “]”
           “<<” expr “>>”
           nodename
           variable
```

Ordinary parentheses are used to bracket sub-expressions, curly braces { } denote zero or more repetitions and angle brackets <> denote one or more repetitions. Square brackets [] are used to enclose optional items. The double angle brackets are a special feature included for building context dependent loops and are explained further below. Finally, the complete network is defined by a list of sub-network definitions followed by a single expression within parentheses.

```
network   = {subnet} “(” expr “)”
```

Note that C style comments may be placed anywhere in the text of the network definition.

As an example, the following network defines a syntax for some simple edit commands

```
$dir  = up | down | left | right;
$mvcmd = move $dir | top | bottom;
$item  = char | word | line | page;
$dlcmd = delete [$item]; /* default is char */
$incmd = insert;
$encmd = end [insert];
$cmd  = $mvcmd|$dlcmd|$incmd|$encmd;
({sil} < $cmd {sil} > quit)
```

Double angle brackets are used to construct contextually consistent context-dependent loops such as a word-pair grammar.⁵ This function can also be used to generate consistent triphone loops for phone recognition⁶. The entry and exit conditions to a context-dependent loop can be controlled by the invisible pseudo-words TLOOP_BEGIN and TLOOP_END. The right context of TLOOP_BEGIN defines the legal loop start nodes, and the left context of TLOOP_END defines the legal loop finishers. If TLOOP_BEGIN/TLOOP_END are not present then all models are connected to the entry/exit of the loop.

A word-pair grammar simply defines the legal set of words that can follow each word in the vocabulary. To generate a network to represent such a grammar a right context-dependent loop could be used. The legal sentence set of sentence start and end words are defined as above using TLOOP_BEGIN/TLOOP_END.

For example, the following lists the legal followers for each word in a 7 word vocabulary

```
ENTRY      - show, tell, give
show       - me, all
tell       - me, all
me         - all
all        - names, addresses
names     - and, names, addresses, show, tell, EXIT
addresses - and, names, addresses, show, tell, EXIT
and        - names, addresses, show, tell
```

HPARSE can generate a suitable lattice to represent this word-pair grammar by using the following specification:

```
$TLOOP_BEGIN_FLLWRS = show|tell|give;
$TLOOP_END_PREDS    = names|addresses;
$show_FLLWRS        = me|all;
$tell_FLLWRS        = me|all;
$me_FLLWRS          = all;
$all_FLLWRS         = names|addresses;
$names_FLLWRS       = and|names|addresses|show|tell|TLOOP_END;
$addresses_FLLWRS   = and|names|addresses|show|tell|TLOOP_END;
$and_FLLWRS         = names|addresses|show|tell;
```

```
( sil <<
  TLOOP_BEGIN+TLOOP_BEGIN_FLLWRS |
  TLOOP_END_PREDS-TLOOP_END |
  show+show_FLLWRS |
  tell+tell_FLLWRS |
  me+me_FLLWRS |
  all+all_FLLWRS |
  names+names_FLLWRS |
  addresses+addresses_FLLWRS |
  and+and_FLLWRS
  >> sil )
```

⁵The expression between double angle brackets must be a simple list of alternative node names or a variable which has such a list as its value

⁶In HTK V2 it is preferable for these context-loop expansions to be done automatically via HNET, to avoid requiring a dictionary entry for every context-dependent model

where it is assumed that each utterance begins and ends with `sil` model.

In this example, each set of contexts is defined by creating a variable whose alternatives are the individual contexts. The actual context-dependent loop is indicated by the `<< >>` brackets. Each element in this loop is a single variable name of the form `A-B+C` where `A` represents the left context, `C` represents the right context and `B` is the actual word. Each of `A`, `B` and `C` can be nodenames or variable names but note that this is the only case where variable names are expanded automatically and the usual `$` symbol is not used⁷. Both `A` and `C` are optional, and left and right contexts can be mixed in the same triphone loop.

14.12.3 Compatibility Mode

In HPARSE compatibility mode, the interpretation of the ENBF network is that used by the HTK V1.5 HVITE program. In which HPARSE ENBF notation was used to define both the word level syntax and the dictionary. Compatibility mode is aimed at converting files written for HTK V1.5 into their equivalent HTK V2 representation. Therefore HPARSE will output the word level portion of such a ENBF syntax as an HTK V2 lattice file and the pronunciation information is optionally stored in an HTK V2 dictionary file. When operating in compatibility mode and not generating dictionary output, the pronunciation information is discarded.

In compatibility mode, the reserved node names `WD_BEGIN` and `WD_END` are used to delimit word boundaries—nodes between a `WD_BEGIN/WD_END` pair are called “word-internal” while all other nodes are “word-external”. All `WD_BEGIN/WD_END` nodes must have an “external name” attached that denotes the word. It is a requirement that the number of `WD_BEGIN` and the number of `WD_END` nodes are equal and furthermore that there isn’t a direct connection from a `WD_BEGIN` node to a `WD_END`. For example a portion of such an HTK V1.5 network could be

```

$A      =  WD_BEGIN%A ax WD_END%A;
$ABDOMEN =  WD_BEGIN%ABDOMEN ae b d ax m ax n WD_END%ABDOMEN;
$ABIDES  =  WD_BEGIN%ABIDES ax b ay d z WD_END%ABIDES;
$ABOLISH =  WD_BEGIN%ABOLISH ax b aa l ih sh WD_END%ABOLISH;
... etc

( <
  $A | $ABDOMEN | $ABIDES | $ABOLISH | ... etc
> )

```

HPARSE will output the connectivity of the words in an HTK V2 word lattice format file and the pronunciation information in an HTK V2 dictionary. Word-external nodes are treated as words and stored in the lattice with corresponding entries in the dictionary.

It should be noted that in HTK V1.5 any ENBF network could appear between a `WD_BEGIN/WD_END` pair, which includes loops. Care should therefore be taken with syntaxes that define very complex sets of alternative pronunciations. It should also be noted that each dictionary entry is limited in length to 100 phones. If multiple instances of the same word are found in the expanded HParse network, a dictionary entry will be created for only the first instance and subsequent instances are ignored (a warning is printed). If words with a `NULL` external name are present then the dictionary will contain a `NULL` output symbol.

Finally, since the implementation of the generation of the HPARSE network has been revised⁸ the semantics of variable definition and use has been slightly changed. Previously variables could be redefined during network definition and each use would follow the most recent definition. In HTK V2 only the final definition of any variable is used in network expansion.

14.12.4 Use

HPARSE is invoked via the command line

```
HParse [options] syntaxFile latFile
```

⁷If the base-names or left/right context of the context-dependent names in a context-dependent loop are variables, no `$` symbols are used when writing the context-dependent nodename.

⁸With the added benefit of rectifying some residual bugs in the HTK V1.5 implementation

HPARSE will then read the set of ENBF rules in `syntaxFile` and produce the output lattice in `latFile`.

The detailed operation of HPARSE is controlled by the following command line options

- b Output the lattice in binary format. This increases speed of subsequent loading (default ASCII text lattices).
- c Set V1.5 compatibility mode. Compatibility mode can also be set by using the configuration variable V1COMPAT (default compatibility mode disabled).
- d *s* Output dictionary to file *s*. This is only a valid option when operating in compatibility mode. If not set no dictionary will be produced.
- l Include language model log probabilities in the output These log probabilities are calculated as $-\log(\text{number of followers})$ for each network node.

HPARSE also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

14.12.5 Tracing

HPARSE supports the following trace options where each trace flag is given using an octal base

0001 basic progress reporting.

0002 show final HParse network (before conversion to a lattice)

0004 print memory statistics after HParse lattice generation

0010 show progress through glue node removal.

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

14.13 HQuant

14.13.1 Function

This program will construct a HTK format VQ table consisting of one or more codebooks, each corresponding to an independent data stream. A codebook is a collection of indexed reference vectors that are intended to represent the structure of the training data. Ideally, every compact cluster of points in the training data is represented by one reference vector. VQ tables are used directly when building systems based on Discrete Probability HMMs. In this case, the continuous-valued speech vectors in each stream are replaced by the index of the closest reference vector in each corresponding codebook. Codebooks can also be used to attach VQ indices to continuous observations. A typical use of this is to preselect Gaussians in probability computations. More information on the use of VQ tables is given in section 5.11.

Codebook construction consists of finding clusters in the training data, assigning a unique reference vector (the cluster centroid) to each, and storing the resultant reference vectors and associated codebook data in a VQ table. HQUANT uses a top-down clustering process whereby clusters are iteratively split until the desired number of clusters are found. The optimal codebook size (number of reference vectors) depends on the structure and amount of the training data, but a value of 256 is commonly used.

HQUANT can construct both linear (i.e. flat) and tree-structured (i.e. binary) codebooks. Linear codebooks can have lower quantisation errors but tree-structured codebooks have $\log_2 N$ access times compared to N for the linear case. The distance metric can either be Euclidean, diagonal covariance Mahalanobis or full covariance Mahalanobis.

14.13.2 VQ Codebook Format

Externally, a VQ table is stored in a text file consisting of a header followed by a sequence of entries representing each tree node. One tree is built for each stream and linear codebooks are represented by a tree in which there are only right branches.

The header consists of a magic number followed by the covariance kind, the number of following nodes, the number of streams and the width of each stream.

```
header =      magic type covkind numNodes numS swidth1 swidth2 ...
```

where *magic* is a magic number which is usually the code for the parameter kind of the data. The *type* defines the type of codebook

```
type =      linear (0) , binary tree-structured (1)
```

The covariance kind determines the type of distance metric to be used

```
covkind =   diagonal covariance (1), full covariance (2), euclidean (5)
```

Within the file, these covariances are stored in inverse form.

Each node entry has the following form

```
node-entry =  stream vqidx nodeld leftld rightld
              mean-vector
              [inverse-covariance-matrix | inverse-variance-vector]
```

Stream is the stream index for this entry. *Vqidx* is the VQ index corresponding to this entry. This is the number that appears in vector quantised speech files. In tree-structured code-books, it is zero for non-terminal nodes. Every node has a unique integer identifier (distinct from the VQ index) given by *nodeld*. The left and right daughter of the current node are given by *leftld* and *rightld*. In a linear codebook, the left identifier is always zero.

Some examples of VQ tables are given in Chapter 11.

14.13.3 Use

HQUANT is invoked via the command line

```
HQuant [options] vqFile trainFiles ...
```

where `vqFile` is the name of the output VQ table file. The effect of this command is to read in training data from each `trainFile`, cluster the data and write the final cluster centres into the VQ table file.

The list of training files can be stored in a script file if required. Furthermore, the data used for training the codebook can be limited to that corresponding to a specified label. This can be used, for example, to train phone specific codebooks. When constructing a linear codebook, the maximum number of iterations per cluster can be limited by setting the configuration variable `MAXCLUSTITER`. The minimum number of samples in any one cluster can be set using the configuration variable `MINCLUSTSIZE`.

The detailed operation of `HQUANT` is controlled by the following command line options

- d Use a diagonal-covariance Mahalanobis distance metric for clustering (default is to use a Euclidean distance metric).
- f Use a full-covariance Mahalanobis distance metric for clustering (default is to use a Euclidean distance metric).
- g Output the global covariance to a codebook. Normally, covariances are computed individually for each cluster using the data in that cluster. This option computes a global covariance across all the clusters.
- l *s* The string *s* must be the name of a segment label. When this option is used, `HQUANT` searches through all of the training files and uses only the speech frames from segments with the given label. When this option is not used, `HQUANT` uses all of the data in each training file.
- n *S N* Set size of codebook for stream *S* to *N* (default 256). If tree-structured codebooks are required then *N* must be a power of 2.
- s *N* Set number of streams to *N* (default 1). Unless the `-w` option is used, the width of each stream is set automatically depending on the size and parameter kind of the training data.
- t Create tree-structured codebooks (default linear).
- w *S N* Set width of stream *S* to *N*. This option overrides the default decomposition that `HTK` normally uses to divide a parameter file into streams. If this option is used, it must be repeated for each individual stream.
- F *fmt* Set the source data format to *fmt*.
- G *fmt* Set the label file format to *fmt*.
- I *mlf* This loads the master label file *mlf*. This option may be repeated to load several MLFs.
- L *dir* Search directory *dir* for label files (default is to search current directory).
- X *ext* Set label file extension to *ext* (default is `lab`).

`HQUANT` also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

14.13.4 Tracing

`HQUANT` supports the following trace options where each trace flag is given using an octal base

- 00001 basic progress reporting.
- 00002 dump global mean and covariance
- 00004 trace data loading.
- 00010 list label segments.
- 00020 dump clusters.
- 00040 dump VQ table.

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

14.14 HRest

14.14.1 Function

HREST performs basic Baum-Welch re-estimation of the parameters of a single HMM using a set of observation sequences. HREST can be used for normal isolated word training in which the observation sequences are realisations of the corresponding vocabulary word.

Alternatively, HREST can be used to generate *seed* HMMs for phoneme-based recognition. In this latter case, the observation sequences will consist of segments of continuously spoken training material. HREST will *cut* these out of the training data automatically by simply giving it a segment label.

In both of the above applications, HREST is intended to operate on HMMs with initial parameter values estimated by HINIT.

HREST supports multiple mixture components, multiple streams, parameter tying within a single model, full or diagonal covariance matrices, tied-mixture models and discrete models. The outputs of HREST are often further processed by HEREST.

Like all re-estimation tools, HREST allows a floor to be set on each individual variance by defining a variance floor macro for each data stream (see chapter 8). If any diagonal covariance component falls below 0.00001, then the corresponding mixture weight is set to zero. A warning is issued if the number of mixtures is greater than one, otherwise an error occurs. Applying a variance floor via the `-v` option or a variance floor macro can be used to prevent this.

14.14.2 Use

HREST is invoked via the command line

```
HRest [options] hmm trainFiles ...
```

This causes the parameters of the given `hmm` to be re-estimated repeatedly using the data in `trainFiles` until either a maximum iteration limit is reached or the re-estimation converges. The HMM definition can be contained within one or more macro files loaded via the standard `-H` option. Otherwise, the definition will be read from a file called `hmm`. The list of train files can be stored in a script file if required.

The detailed operation of HREST is controlled by the following command line options

- c *f* Set the threshold for tied-mixture observation pruning to *f*. When all mixtures of all models are tied to create a full tied-mixture system, the calculation of output probabilities is treated as a special case. Only those mixture component probabilities which fall within *f* of the maximum mixture component probability are used in calculating the state output probabilities (default 10.0).
- e *f* This sets the convergence factor to the real value *f*. The convergence factor is the relative change between successive values of $P(O|\lambda)$ (default value 0.0001).
- i *N* This sets the maximum number of re-estimation cycles to *N* (default value 20).
- l *s* The string *s* must be the name of a segment label. When this option is used, HREST searches through all of the training files and cuts out all segments with the given label. When this option is not used, HREST assumes that each training file is a single token.
- m *N* Sets the minimum number of training examples to be *N*. If fewer than *N* examples are supplied then an error is reported (default value 3).
- t Normally, training sequences are rejected if they have fewer frames than the number of emitting states in the HMM. Setting this switch disables this reject mechanism⁹.
- u *flags* By default, HREST updates all of the HMM parameters, that is, means, variances, mixture weights and transition probabilities. This option causes just the parameters indicated by the *flags* argument to be updated, this argument is a string containing one or more of the letters *m* (mean), *v* (variance), *t* (transition) and *w* (mixture weight). The presence of a letter enables the updating of the corresponding parameter set.

⁹Using this option only makes sense if the HMM has skip transitions

- v *f* This sets the minimum variance (i.e. diagonal element of the covariance matrix) to the real value *f*. This is ignored if an explicit variance floor macro is defined. The default value is 0.0.
- w *f* Any mixture weight or discrete observation probability which falls below the global constant MINMIX is treated as being zero. When this parameter is set, all mixture weights are floored to $f * \text{MINMIX}$.
- B Output HMM definition files in binary format.
- F *fmt* Set the source data format to *fmt*.
- G *fmt* Set the label file format to *fmt*.
- H *mmf* Load HMM macro model file *mmf*. This option may be repeated to load multiple MMFs.
- I *mlf* This loads the master label file *mlf*. This option may be repeated to load several MLFs.
- L *dir* Search directory *dir* for label files (default is to search current directory).
- M *dir* Store output HMM macro model files in the directory *dir*. If this option is not given, the new HMM definition will overwrite the existing one.
- X *ext* Set label file extension to *ext* (default is *lab*).

HREST also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

14.14.3 Tracing

HREST supports the following trace options where each trace flag is given using an octal base

000001 basic progress reporting.

000002 output information on the training data loaded.

000004 the observation probabilities.

000010 the alpha matrices.

000020 the beta matrices.

000040 the occupation counters.

000100 the transition counters.

000200 the mean counters.

000400 the variance counters.

001000 the mixture weight counters.

002000 the re-estimated transition matrix.

004000 the re-estimated mixture weights.

010000 the re-estimated means.

020000 the re-estimated variances.

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

14.15 HResults

14.15.1 Function

HRESULTS is the HTK performance analysis tool. It reads in a set of label files (typically output from a recognition tool such as HVITE) and compares them with the corresponding reference transcription files. For the analysis of speech recognition output, the comparison is based on a Dynamic Programming-based string alignment procedure. For the analysis of word-spotting output, the comparison uses the standard US NIST FOM metric.

When used to calculate the sentence accuracy using DP the basic output is recognition statistics for the whole file set in the format

```
----- Overall Results -----
SENT: %Correct=13.00 [H=13, S=87, N=100]
WORD: %Corr=53.36, Acc=44.90 [H=460,D=49,S=353,I=73,N=862]
=====
```

The first line gives the sentence-level accuracy based on the total number of label files which are identical to the transcription files. The second line is the word accuracy based on the DP matches between the label files and the transcriptions¹⁰. In this second line, H is the number of correct labels, D is the number of deletions, S is the number of substitutions, I is the number of insertions and N is the total number of labels in the defining transcription files. The percentage number of labels correctly recognised is given by

$$\%Correct = \frac{H}{N} \times 100\% \quad (14.4)$$

and the accuracy is computed by

$$Accuracy = \frac{H - I}{N} \times 100\% \quad (14.5)$$

In addition to the standard HTK output format, HRESULTS provides an alternative similar to that used in the US NIST scoring package, i.e.

```
|=====|
|          # Snt | Corr   Sub   Del   Ins   Err   S. Err | |
|---|---|---|
| Sum/Avg |   87 | 53.36 40.95  5.68  8.47 55.10  87.00 |
|-----|
```

Optional extra outputs available from HRESULTS are

- recognition statistics on a per file basis
- recognition statistics on a per speaker basis
- recognition statistics from best of N alternatives
- time-aligned transcriptions
- confusion matrices

For comparison purposes, it is also possible to assign two labels to the same equivalence class (see `-e` option). Also, the *null* label `???` is defined so that making any label equivalent to the null label means that it will be ignored in the matching process. Note that the order of equivalence labels is important, to ensure that label `X` is ignored, the command line option `-e ??? X` would be used. Label files containing triphone labels of the form `A-B+C` can be optionally stripped down to just the class name `B` via the `-s` switch.

The word spotting mode of scoring can be used to calculate hits, false alarms and the associated figure of merit for each of a set of keywords. Optionally it can also calculate ROC information over a range of false alarm rates. A typical output is as follows

¹⁰ The choice of “Sentence” and “Word” here is the usual case but is otherwise arbitrary. HRESULTS just compares label sequences. The sequences could be paragraphs, sentences, phrases or words, and the labels could be phrases, words, syllables or phones, etc. Options exist to change the output designations ‘SENT’ and ‘WORD’ to whatever is appropriate.

----- Figures of Merit -----				
Keyword:	#Hits	#FAs	#Actual	FOM
A:	8	1	14	30.54
B:	4	2	14	15.27
Overall:	12	3	28	22.91

which shows the number of hits and false alarms (FA) for two keywords A and B. A label in the test file with start time t_s and end time t_e constitutes a hit if there is a corresponding label in the reference file such that $t_s < t_m < t_e$ where t_m is the mid-point of the reference label.

Note that for keyword scoring, the test transcriptions must include a score with each labelled word spot and all transcriptions must include boundary time information.

The FOM gives the % of hits averaged over the range 1 to 10 FA's per hour. This is calculated by first ordering all spots for a particular keyword according to the match score. Then for each FA rate f , the number of hits are counted starting from the top of the ordered list and stopping when f have been encountered. This corresponds to a *posteriori* setting of the keyword detection threshold and effectively gives an upper bound on keyword spotting performance.

14.15.2 Use

HRESULTS is invoked by typing the command line

```
HResults [options] hmmList recFiles ...
```

This causes HRESULTS to be applied to each `recFile` in turn. The `hmmList` should contain a list of all model names for which result information is required. Note, however, that since the context dependent parts of a label can be stripped, this list is not necessarily the same as the one used to perform the actual recognition. For each `recFile`, a transcription file with the same name but the extension `.lab` (or some user specified extension - see the `-X` option) is read in and matched with it. The `recfiles` may be master label files (MLFs), but note that even if such an MLF is loaded using the `-I` option, the list of files to be checked still needs to be passed, either as individual command line arguments or via a script with the `-S` option. For this reason, it is simpler to pass the `recFile` MLF as one of the command line filename arguments. For loading reference label file MLF's, the `-I` option must be used. The reference labels and the recognition labels must have different file extensions. The available options are

- a `s` change the label SENT in the output to `s`.
- b `s` change the label WORD in the output to `s`.
- c when comparing labels convert to upper case. Note that case is still significant for equivalences (see `-e` below).
- d `N` search the first `N` alternatives for each test label file to find the most accurate match with the reference labels. Output results will be based on the most accurate match to allow `NBest` error rates to be found.
- e `s t` the label `t` is made equivalent to the label `s`. More precisely, `t` is assigned to an equivalence class of which `s` is the identifying member.
- f Normally, HRESULTS accumulates statistics for all input files and just outputs a summary on completion. This option forces match statistics to be output for each input test file.
- g `fmt` This sets the test label format to `fmt`. If this is not set, the `recFiles` should be in the same format as the reference files.
- h Output the results in the same format as US NIST scoring software.
- k `s` Collect and output results on a speaker by speaker basis (as well as globally). `s` defines a pattern which is used to extract the speaker identifier from the test label file name. In addition to the pattern matching metacharacters `*` and `?` (which match zero or more characters and a single character respectively), the character `%` matches any character whilst including it as part of the speaker identifier.

- m *N* Terminate after collecting statistics from the first *N* files.
- n Set US NIST scoring software compatibility.
- p This option causes a phoneme confusion matrix to be output.
- s This option causes all phoneme labels with the form *A-B+C* to be converted to *B*. It is useful for analysing the results of phone recognisers using context dependent models.
- t This option causes a time-aligned transcription of each test file to be output provided that it differs from the reference transcription file.
- u *f* Changes the time unit for calculating false alarm rates (for word spotting scoring) to *f* hours (default is 1.0).
- w Perform word spotting analysis rather than string accuracy calculation.
- z *s* This redefines the null class name to *s*. The default null class name is *???*, which may be difficult to manage in shell script programming.
- G *fmt* Set the label file format to *fmt*.
- I *mlf* This loads the master label file *mlf*. This option may be repeated to load several MLFs.
- L *dir* Search directory *dir* for label files (default is to search current directory).
- X *ext* Set label file extension to *ext* (default is *lab*).

HRESULTS also supports the standard options *-A*, *-C*, *-D*, *-S*, *-T*, and *-V* as described in section 4.4.

14.15.3 Tracing

HRESULTS supports the following trace options where each trace flag is given using an octal base

00001 basic progress reporting.

00002 show error rate for each test alternative.

00004 show speaker identifier matches.

00010 warn about non-keywords found during word spotting.

00020 show detailed word spotting scores.

00040 show memory usage.

Trace flags are set using the *-T* option or the *TRACE* configuration variable.

14.16 HSGen

14.16.1 Function

This program will read in a word network definition in standard HTK lattice format representing a Regular Grammar G and randomly generate sentences from the language $L(G)$ of G . The sentences are written to standard output, one per line and an option is provided to number them if required.

The empirical entropy H_e can also be calculated using the formula

$$H_e = \frac{\sum_k P(S_k)}{\sum_k |S_k|} \quad (14.6)$$

where S_k is the k 'th sentence generated and $|S_k|$ is its length. The perplexity P_e is computed from H_e by

$$P_e = 2^{H_e} \quad (14.7)$$

The probability of each sentence $P(S_k)$ is computed from the product of the individual branch probabilities.

14.16.2 Use

HSGEN is invoked by the command line

```
HSGen [options] wnet dictfile
```

where `dictfile` is a dictionary containing all of the words used in the word network stored in `wnet`. This dictionary is only used as a word list, the pronunciations are ignored.

The available options are

- l When this option is set, each generated sentence is preceded by a line number.
- n N This sets the total number of sentences generated to be N (default value 100).
- q Set quiet mode. This suppresses the printing of sentences. It is useful when estimating the entropy of $L(G)$ since the accuracy of the latter depends on the number of sentences generated.
- s Compute word network statistics. When set, the number of network nodes, the vocabulary size, the empirical entropy, the perplexity, the average sentence length, the minimum sentence length and the maximum sentence length are computed and printed on the standard output.

HSLAB also supports the standard options -A, -C, -D, -S, -T, and -V as described in section 4.4.

14.16.3 Tracing

HSLAB supports the following trace options where each trace flag is given using an octal base

00001 basic progress reporting

00002 detailed trace of lattice traversal

Trace flags are set using the -T option or the TRACE configuration variable.

14.17 HSLab

14.17.1 Function

HSLAB is an interactive label editor for manipulating speech label files. An example of using HSLAB would be to load a sampled waveform file, determine the boundaries of the speech units of interest and assign labels to them. Alternatively, an existing label file can be loaded and edited by changing current label boundaries, deleting and creating new labels. HSLAB is the only tool in the HTK package which makes use of the graphics library HGRAF.

When started HSLAB displays a window which is split into two parts: a display section and a control section (see Fig 14.1). The display section contains the plotted speech waveform with the associated labels. The control section consists of a palette of buttons which are used to invoke the various facilities available in the tool. The buttons are laid out into three different groups depending on the function they perform. Group one (top row) contains buttons related to basic input/output commands. Group two (middle row) implements the viewing and record/playback functions. The buttons in group three (bottom row) are used for labelling. To invoke a particular function, place the mouse pointer onto the corresponding button and click once. All commands which require further interaction with the user after invocation will display guiding text in the message area telling the user what he or she is expected to do next. For example, to delete a label, the user will click on **Delete**, the message “Please select label to delete” will appear in the message area and the user will be expected to click in that part of the display section corresponding to the label to be deleted (not on the label itself).

A *marked region* is a slice of the waveform currently visible in the window. A region is marked by clicking on **Mark** and specifying two boundaries by clicking in the display section. When marked, a region will be displayed in inverse colours. In the presence of a marked region the commands **Play**, **Label** and **Label as** will be applied to the specified region rather than to the whole of the waveform visible on the screen. Part of the waveform can also be made into a marked region with the commands **Zoom Out** and **Select**. **Zoom Out** will take the user back to the previous level of magnification and the waveform being displayed before the execution of the command will become a marked region. **Select** will make the part of the waveform corresponding to a particular label into a marked region. This can be useful for playing back existing labels.

Labelling is carried out with **Label** and **Label as**. **Label** will assign *The Current Label* to a specified slice of the waveform, whilst **Label as** will prompt the user to type-in the labelling string. *The Current Label* is shown in the button in the bottom right corner of the control section. It defaults to “Speech” and it can be changed by clicking on the button it resides in. Multiple alternative transcriptions are manipulated using the **Set [?]** and **New** buttons. The former is used to select the desired transcription, the latter is used to create a new alternative transcription.

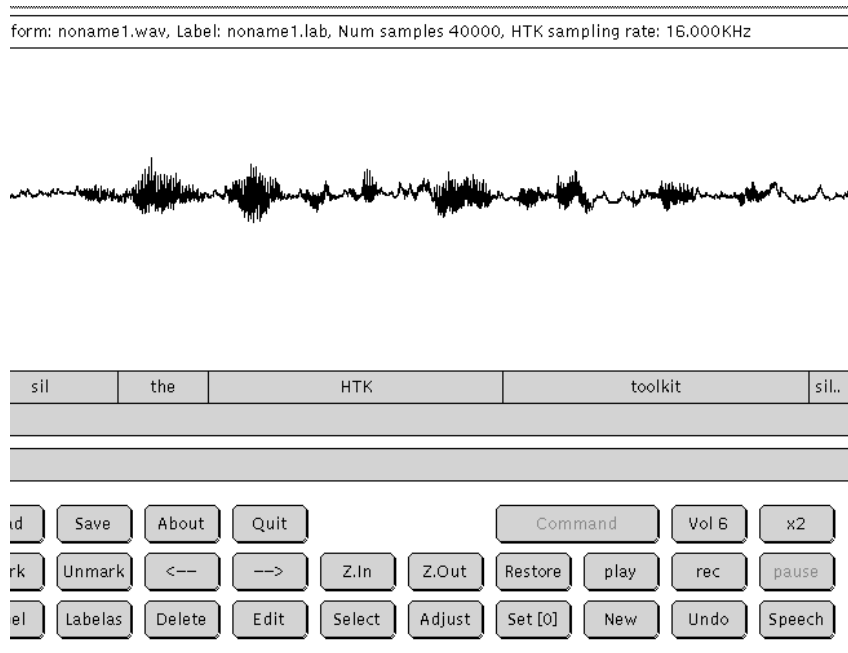


Fig. 14.1 HSLab display window

14.17.2 Use

HSLAB is invoked by typing the command line

```
HSLab [options] dataFile
```

where `dataFile` is a data file in any of the supported formats with a `WAVEFORM` sample kind. If the given data file does not exist, then HSLAB will assume that a new file is to be recorded with this name.

The available options for HSLAB are

- a With this switch present, the numeric part of the global labelling string is automatically incremented after every `Label` operation.
- i file This option allows transcription files to be output to the named master label file (MLF).
- n Normally HSLAB expects to load an existing label file whose name is derived from the speech data file. This option tells HSLAB that a new empty transcription is to be created for the loaded data-file.
- s string This option allows the user to set the string displayed in the “command” button used to trigger external commands.
- F fmt Set the source data format to `fmt`.
- G fmt Set the label file format to `fmt`.
- I mlf This loads the master label file `mlf`. This option may be repeated to load several MLFs.
- L dir Search directory `dir` for label files (default is to search current directory).
- X ext Set label file extension to `ext` (default is `lab`).

HSLAB also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

The following is a summary of the function of each HSLAB button.

- Load** Load a speech data file and associated transcription. If changes have been made to the currently loaded transcription since it was last saved the user will be prompted to save these changes before loading the new file.
- Save** Save changes made to the transcription into a specified label file.
- About** Print information about HSLab.
- Quit** Exit from HSLAB. If alterations have been made to the currently loaded transcription since it was last saved, the user will be prompted to save these changes before exiting.
- Command** This button is used to trigger an external command which can process the waveform file currently loaded in HSLAB. This is accomplished by setting the environment variable `HSLABCMD` to the shell command required to perform the processing. When the **Command** button is pressed, any occurrence of `$` in the shell command is replaced by the name of the currently loaded waveform file. Note that only the filename without its extension is substituted. The string displayed in the “command” button can be changed using the `-s` option.
- Mark** Mark a region of the displayed waveform. The user is prompted to specify the start and the end point of a region with the mouse pointer. The marked region will be displayed in inverse colours. Only one region can be marked at a time.
- Unmark** Unmark a previously marked region.
- <--** Scroll the display to the left.
- >** Scroll the display to the right.
- Z.In** Zoom into a part of the displayed waveform. If there is a currently marked region then that region will be zoomed into, otherwise, the user will be prompted to select a slice of the waveform by specifying two points using the mouse pointer.
- Z.Out** Restore the previous viewing level.
- Restore** Display the complete waveform into the window. Any marked regions will be unmarked.
- Play** If there is a marked region of the waveform then that portion of the signal will be played through the internal speaker. Otherwise, the command will apply to the waveform visible on the screen.
- Rec** This initiates recording from the audio input device. The maximum duration of a recording is limited to 2 mins at 16KHz sampling rate. Two bar-graphs are displayed: the first (red) shows the number of samples recorded, the second bar (green) displays the energy of the incoming signal. Once pressed, the **Rec** button changes into **Stop** which, in turn, is used to terminate the operation. When finished, the audio data stored in the buffer is written out to disk. Each recording is stored in alternating files `dataFile_0` and `dataFile_1`.
- Pause** Clicking on this button pauses/un-pauses the recording operation.
- Volume** This button is used to select the playback volume of the audio device.
- x1** This button selects the current level of waveform magnification. The available factors are $\times 1$, $\times 2$, $\times 4$, $\times 8$, $\times 16$, and $\times 32$.
- Label** If a marked region exists, then the waveform contained in the region will be labelled with *The Current Label*. Otherwise, the command will be applied to the waveform visible on the screen.
- Labelas** Same as above, however, the user is prompted to type in the labelling string.
- Delete** Delete a label.
- Edit** Edit the string of a label.
- Select** Select a label as a marked region.

Adjust Adjust the boundaries of a label. To select the label boundary to adjust, click in the display near to the label boundary.

Set [?] This button is used to select the current alternative transcription displayed and used in HSLAB.

New Creates a new alternative transcription. If an empty alternative transcription already exists, then a new transcription is not created.

Undo Single level undo operation for labelling commands.

Speech Change the current labelling string (the button in the bottom right of the control area).

The following “mouse” shortcuts are provided. To mark a region position the pointer at one of the desired boundaries, then press the left mouse button and while holding it down position the pointer at the other region boundary. Upon releasing the mouse button the marked region will be highlighted. To play a label position the mouse cursor anywhere within the corresponding label “slice” in the label area of the display and click the left mouse button.

14.17.3 Tracing

HSLAB does not provide any trace options.

14.18 HSmooth

14.18.1 Function

This program is used to smooth a set of context-dependent tied mixture or discrete HMM's using deleted interpolation. The program operates as a replacement for the second pass of HEREST when working in *parallel mode*¹¹. It reads in the N sets of accumulator files containing the statistics accumulated during the first pass and then interpolates the mixture weights between each context dependent model and its corresponding context independent version. The interpolation weights are chosen to maximise the likelihood of each deleted block with respect to the probability distributions estimated from all other blocks.

14.18.2 Use

HSMOOTH is invoked via the command line

```
HSmooth [options] hmmList accFile ...
```

where `hmmList` contains the list of context dependent models to be smoothed and each `accFile` is a file of the form `HERN.acc` dumped by a previous execution of HEREST with the `-p` option set to N . The HMM definitions are loaded and then for every state and stream of every context dependent model X , the optimal interpolation weight is found for smoothing between the mixture weights determined from the X accumulators alone and those determined from the context independent version of X . The latter is computed simply by summing the accumulators across all context dependent allophones of X .

The detailed operation of HSMOOTH is controlled by the following command line options

- b *f* Set the value of epsilon for convergence in the binary chop optimisation procedure to *f*. The binary chop optimisation procedure for each interpolation weight terminates when the gradient is within epsilon of zero (default 0.001).
- c *N* Set maximum number of interpolation iterations for the binary chop optimisation procedure to be N (default 16).
- d *dir* Normally HSMOOTH expects to find the HMM definitions in the current directory. This option tells HSMOOTH to look in the directory *dir* to find them.
- m *N* Set the minimum number of training examples required for any model to N . If the actual number falls below this value, the HMM is not updated and the original parameters are used for the new version (default value 1).
- o *ext* This causes the file name extensions of the original models (if any) to be replaced by *ext*.
- s *file* This causes statistics on occupation of each state to be output to the named file.
- u *flags* By default, HSMOOTH updates all of the HMM parameters, that is, means, variances and transition probabilities. This option causes just the parameters indicated by the *flags* argument to be updated, this argument is a string containing one or more of the letters *m* (mean), *v* (variance), *t* (transition) and *w* (mixture weight). The presence of a letter enables the updating of the corresponding parameter set.
- v *f* This sets the minimum variance (i.e. diagonal element of the covariance matrix) to the real value *f* (default value 0.0).
- w *f* Any mixture weight which falls below the global constant MINMIX is treated as being zero. When this parameter is set, all mixture weights are floored to $f * \text{MINMIX}$.
- x *ext* By default, HSMOOTH expects a HMM definition for the model X to be stored in a file called X . This option causes HSMOOTH to look for the HMM definition in the file $X.\text{ext}$.
- B Output HMM definition files in binary format.
- H *mmf* Load HMM macro model file *mmf*. This option may be repeated to load multiple MMFs.

¹¹It is not, of course, necessary to have multiple processors to use this program since each 'parallel' activation can be executed sequentially on a single processor

`-M dir` Store output HMM macro model files in the directory `dir`. If this option is not given, the new HMM definition will overwrite the existing one.

HSMOOTH also supports the standard options `-A`, `-C`, `-D`, `-S`, `-T`, and `-V` as described in section 4.4.

14.18.3 Tracing

HSMOOTH supports the following trace options where each trace flag is given using an octal base

00001 basic progress reporting.

00002 show interpolation weights.

00004 give details of optimisation algorithm.

Trace flags are set using the `-T` option or the `TRACE` configuration variable.

14.19 HVite

14.19.1 Function

HVITE is a general-purpose Viterbi word recogniser. It will match a speech file against a network of HMMs and output a transcription for each. When performing N-best recognition a word level lattice containing multiple hypotheses can also be produced.

Either a word level lattice or a label file is read in and then expanded using the supplied dictionary to create a model based network. This allows arbitrary finite state word networks and simple forced alignment to be specified.

This expansion can be used to create context independent, word internal context dependent and cross word context dependent networks. The way in which the expansion is performed is determined automatically from the dictionary and HMMList. When all labels appearing in the dictionary are defined in the HMMList no expansion of model names is performed. Otherwise if all the labels in the dictionary can be satisfied by models dependent only upon word internal context these will be used else cross word context expansion will be performed. These defaults can be overridden by HNET configuration parameters.

HVITE supports shared parameters and appropriately pre-computes output probabilities. For increased processing speed, HVITE can optionally perform a beam search controlled by a user specified threshold (see `-t` option). When fully tied mixture models are used, observation pruning is also provided (see the `-c` option). Speaker adaptation is also supported by HVITE both in terms of recognition using an adapted model set or a TMF (see the `-J` option), and in the estimation of a transform by unsupervised adaptation using maximum likelihood linear regression (MLLR) (see the `-j` option).

14.19.2 Use

HVITE is invoked via the command line

```
HVite [options] dictFile hmmList testFiles ...
```

HVite will then either load a single network file and match this against each of the test files `-w netFile`, or create a new network for each test file either from the corresponding label file `-a` or from a word lattice `-w`. When a new network is created for each test file the path name of the label (or lattice) file to load is determined from the test file name and the `-L` and `-X` options described below.

If no `testFiles` are specified the `-w s` option must be specified and recognition will be performed from direct audio.

The `hmmList` should contain a list of the models required to construct the network from the word level representation.

The recogniser output is written in the form of a label file whose path name is determined from the test file name and the `-l` and `-x` options described below. The list of test files can be stored in a script file if required.

When performing N-best recognition (see `-n N` option described below) the output label file can contain multiple alternatives `-n N M` and a lattice file containing multiple hypotheses can be produced.

The detailed operation of HVITE is controlled by the following command line options

- `-a` Perform alignment. HVITE will load a label file and create an alignment network for each test file.
- `-b s` Use `s` as the sentence boundary during alignment.
- `-c f` Set the tied-mixture observation pruning threshold to `f`. When all mixtures of all models are tied to create a full tied-mixture system, the calculation of output probabilities is treated as a special case. Only those mixture component probabilities which fall within `f` of the maximum mixture component probability are used in calculating the state output probabilities (default 10.0).
- `-d dir` This specifies the directory to search for the HMM definition files corresponding to the labels used in the recognition network.

- e When using direct audio input, output transcriptions are not normally saved. When this option is set, each output transcription is written to a file called PnS where n is an integer which increments with each output file, P and S are strings which are by default empty but can be set using the configuration variables RECOUTPREFIX and RECOUTSUFFIX.
- f During recognition keep track of full state alignment. The output label file will contain multiple levels. The first level will be the state number and the second will be the word name (not the output symbol).
- g When using direct audio input, this option enables audio replay of each input utterance after it has been recognised.
- i s Output transcriptions to MLF s.
- j i Perform incremental MLLR adaptation every i utterances
- k s1 s2 Set the description field s1 in the transform model file to s2. Currently the choices for field are uid, unname, chan and desc.
- l dir This specifies the directory to store the output label files. If this option is not used then HVITE will store the label files in the same directory as the data. When output is directed to an MLF, this option can be used to add a path to each output file name. In particular, setting the option -l '*' will cause a label file named xxx to be prefixed by the pattern "*/xxx" in the output MLF file. This is useful for generating MLFs which are independent of the location of the corresponding data files.
- m During recognition keep track of model boundaries. The output label file will contain multiple levels. The first level will be the model number and the second will be the word name (not the output symbol).
- n i [N] Use i tokens in each state to perform N-best recognition. The number of alternative output hypotheses N defaults to 1.
- o s Choose how the output labels should be formatted. s is a string with certain letters (from NSCTWM) indicating binary flags that control formatting options. N normalise acoustic scores by dividing by the duration (in frames) of the segment. S remove scores from output label. By default scores will be set to the total likelihood of the segment. C Set the transcription labels to start and end on frame centres. By default start times are set to the start time of the frame and end times are set to the end time of the frame. T Do not include times in output label files. W Do not include words in output label files when performing state or model alignment. M Do not include model names in output label files when performing state and model alignment.
- p f Set the word insertion log probability to f (default 0.0).
- q s Choose how the output lattice should be formatted. s is a string with certain letters (from ABtvaldmn) indicating binary flags that control formatting options. A attach word labels to arcs rather than nodes. B output lattices in binary for speed. t output node times. v output pronunciation information. a output acoustic likelihoods. l output language model likelihoods. d output word alignments (if available). m output within word alignment durations. n output within word alignment likelihoods.
- r f Set the dictionary pronunciation probability scale factor to f. (default value 1.0).
- s f Set the grammar scale factor to f. This factor post-multiplies the language model likelihoods from the word lattices. (default value 1.0).
- t f [i l] Enable beam searching such that any model whose maximum log probability token falls more than f below the maximum for all models is deactivated. Setting f to 0.0 disables the beam search mechanism (default value 0.0). In alignment mode two extra parameters i and l can be specified. If the alignment fails at the initial pruning threshold f, then the threshold will be increased by i and the alignment will be retried. This procedure is repeated until the alignment succeeds or the threshold limit l is reached.
- u i Set the maximum number of active models to i. Setting i to 0 disables this limit (default 0).

- v *f* Enable word end pruning. Do not propagate tokens from word end nodes that fall more than *f* below the maximum word end likelihood. (default 0.0).
- w [*s*] Perform recognition from word level networks. If *s* is included then use it to define the network used for every file.
- x *ext* This sets the extension to use for HMM definition files to *ext*.
- y *ext* This sets the extension for output label files to *ext* (default *rec*).
- z *ext* Enable output of lattices (if performing NBest recognition) with extension *ext* (default off).
- L *dir* This specifies the directory to find input label (when *-a* is specified) or network files (when *-w* is specified).
- X *s* Set the extension for the input label or network files to be *s* (default value *lab*).
- F *fmt* Set the source data format to *fmt*.
- G *fmt* Set the label file format to *fmt*.
- H *mmf* Load HMM macro model file *mmf*. This option may be repeated to load multiple MMFs.
- I *m1f* This loads the master label file *m1f*. This option may be repeated to load several MLFs.
- J *tmf* Load a transform set from the transform model file *tmf*.
- K *tmf* Save the transform set in the transform model file *tmf*.
- P *fmt* Set the target label format to *fmt*.

HVITE also supports the standard options *-A*, *-C*, *-D*, *-S*, *-T*, and *-V* as described in section 4.4.

14.19.3 Tracing

HVITE supports the following trace options where each trace flag is given using an octal base

- 0001 enable basic progress reporting.
- 0002 list observations.
- 0004 frame-by-frame best token.
- 0010 show memory usage at start and finish.
- 0020 show memory usage after each utterance.

Trace flags are set using the *-T* option or the *TRACE* configuration variable.

Chapter 15

Configuration Variables

This chapter tabulates all configuration variables used in HTK.

15.1 Configuration Variables used in Library Modules

Table 15.1: Library Module Configuration Variables

Module	Name	Default	Description
HPARM	SOURCEFORMAT	HTK	File format of source
HWAVE	TARGETFORMAT	HTK	File format of target
HLABEL HAUDIO HWAVE HPARM	SOURCERATE	0.0	Sample rate of source in 100ns units
HPARM HWAVE	TARGETRATE	0.0	Sample rate of target in 100ns units
HAUDIO	LINEOUT	T	Enable audio output to machine line output
	PHONESOUT	T	Enable audio output to machine phones output
	SPEAKEROUT	F	Enable audio output to machine internal speaker
	LINEIN	T	Enable audio input from machine line input
	MICIN	F	Enable audio input from machine mic input
HWAVE	NSAMPLES		Num samples in alien file input via a pipe
	HEADERSIZE		Size of header in an alien file
	BYTEORDER		Define byte order VAX or other
	STEREOMODE		Select channel: RIGHT or LEFT
HPARM	SOURCEKIND	ANON	Parameter kind of source
	TARGETKIND	ANON	Parameter kind of target
	SAVECOMPRESSED	F	Save the output file in compressed form
	SAVEWITHCRC	T	Attach a checksum to output parameter file
	ADDITHER	0.0	Level of noise added to input signal
	ZMEANSOURCE	F	Zero mean source waveform before analysis
	WINDOWSIZE	256000.0	Analysis window size in 100ns units
	USEHAMMING	T	Use a Hamming window
	PREEMCOEF	0.97	Set pre-emphasis coefficient
	LPCORDER	12	Order of lpc analysis
	NUMCHANS	20	Number of filterbank channels
	LOFREQ	-1.0	Low frequency cut-off in fbank analysis
	HIFREQ	-1.0	High frequency cut-off in fbank analysis
	WARPFREQ	1.0	Frequency warping factor
WARPLCUTOFF		Lower frequency threshold for non-linear warping	

Module	Name	Default	Description
	CMEANDIR		Directory to find cepstral mean vecotrs
	CMEANMASK		Filename mask for cepstral mean vecotrs
	VARSCALEDIR		Directory to find cepstral variance vecotrs
	VARSCALEMASK		Filename mask for cepstral variance vecotrs
	VARSCALEFN		Filename of global variance scaling vector
	COMPRESSFACT	0.33	Amplitude compression factor for PLP
HLabel HPARM	V1COMPAT	F	HTK V1 compatibility setting
HWAVE	NATURALREADORDER	F	Enable natural read order for binary files
HShell	NATURALWRITEORDER	F	Enable natural write order for binary files
HPARM	USEPOWER	F	Use power not magnitude in fbank analysis
	NUMCEPS	12	Number of cepstral parameters
	CEPLIFTER	22	Cepstral liftering coefficient
	ENORMALISE	T	Normalise log energy
	ESCALE	0.1	Scale log energy
	SILFLOOR	50.0	Energy silence floor in dBs
	DELTAWINDOW	2	Delta window size
	ACCWINDOW	2	Acceleration window size
	VQTABLE	NULL	Name of VQ table
	SIMPLEDIFFS	F	Use simple differences for delta calculations
	RAWENERGY	T	Use raw energy
	AUDIOSIG	0	Audio signal number for remote control
	USESILDET	F	Enable speech/silence detector
	MEASURESIL	T	Measure background silence level
	OUTSILWARN	T	Print a warning message to <code>stdout</code> before measuring audio levels
	SPEECHTHRESH	9.0	Threshold for speech above silence level (in dB)
	SILENERGY	0.0	Average background noise level (in dB) - will normally be measured rather than supplied in configuration
	SPCSEQCOUNT	10	Window over which speech/silence decision reached
SPCGLCHCOUNT	0	Maximum number of frames marked as silence in window which is classified as speech whilst expecting start of speech	
SILSEQCOUNT	100	Number of frames classified as silence needed to mark end of utterance	
SILGLCHCOUNT	2	Maximum number of frames marked as silence in window which is classified as speech whilst expecting silence	
SILMARGIN	40	Number of extra frames included before and after start and end of speech marks from the speech/silence detector	
HLabel	STRIPTRIPHONES	F	Enable triphone stripping
	TRANSALT	0	Filter all but specified label alternative
	TRANSLEV	0	Filter all but specified label level
	LABELSQUOTE	NULL	Select method for quoting in label files
	SOURCELABEL	HTK	Source label format
	TARGETLABEL	HTK	Target label format
HMem	PROTECTSTAKS	F	Enable stack protection
	CHKHMMDEFS	T	Check consistency of HMM defs
	SAVEBINARY	F	Save HMM defs in binary format

Module	Name	Default	Description
HMODEL	KEEPDISTINCT	F	Keep orphan HMMs in distinct files
	SAVEGLOBOPTS	T	Save $\sim o$ with HMM defs
	ORPHANMACFILE	NULL	Last resort file for new macros
	HMMSETKIND	NULL	Kind of HMM Set
	ALLOWOTHERHMMS	T	Allow MMFs to contain HMM definitions which are not listed in the HMM List
	DISCRETELZERO	F	Map DLOGZERO to LZERO in output probability calculations
HNET	FORCECXTEXP	F	Force triphone context expansion to get model names (is overridden by ALLOWCXTEXP)
	FORCELEFTBI	F	Force left biphone context expansion to get model names ie. don't try triphone names
	FORCERIGHTBI	F	Force right biphone context expansion to get model names ie. don't try triphone names
	ALLOWCXTEXP	T	Allow context expansion to get model names
	ALLOWXWRDEXP	F	Allow context expansion across words
	FACTORLM	F	Factor language model likelihoods throughout words rather than applying all at transition into word. This can increase accuracy when pruning is tight and language model likelihoods are relatively high.
	CFWORDBOUNDARY	T	In word-internal triphone systems, context-free phones will be treated as word boundaries
HREC	FORCEOUT	F	Forces the most likely partial hypothesis to be used as the recognition result even when no token reaches the end of the network by the last frame of the utterance
HSHELL	ABORTONERR	F	Causes HError to abort rather than exit
	NONUMESCAPES	F	Prevent writing in 012 format
	MAXTRYOPEN	1	Maximum number of attempts which will be made to open the same file
	EXTENDFILENAME	T	Support for extended filenames
HTRAIN	MAXCLUSTITER	10	Maximum number of cluster iterations
	MINCLUSTSIZE	3	Minimum number of elements in any one cluster
	BINARYACCFORMAT	T	Save accumulator files in binary format
HFB	HSKIPSTART	-1	Start of skip over region (debugging only)
	HSKIPEND	-1	End of skip over region (debugging only)
HADAPT	USEVAR	F	Compute variance transform
	ADPTSIL	T	Transform the silence
	BLOCKS	1	Number of blocks used in the block diagonal matrix implementation
	SAVEBINARY	F	Save HMMs/transforms in binary format
	OCCTHRESH	700	Minimum occupation before computing a regression class transform for a node
	TRACE	0	Trace setting

15.2 Configuration Variables used in Tools

Module	Name	Default	Description
HCOMPV	UPDATEMEANS	F	Update means
	SAVEBINARY	F	Load/Save in binary format
	MINVARFLOOR	0.0	Minimum variance floor
HCOPY	NSTREAMS	1	Number of streams
	SAVEASVQ	F	Save only the VQ indices
	SOURCEFORMAT	HTK	File format of source
	TARGETFORMAT	HTK	File format of target
	SOURCEKIND	ANON	Parameter kind of source
	TARGETKIND	ANON	Parameter kind of target
HEREST	SAVEBINARY	F	Load/Save in binary format
	BINARYACFORMAT	T	Load/Save accumulators in binary format
HEADAPT	SAVEBINARY	F	Load/Save in binary format
HHED	TREEMERGE	T	After tree splitting, merge leaves
	TIEDMIXNAME	TM	Tied mixture base name
HPARSE	V1COMPAT	F	Enable compatibility with HTK V1.X
HRESULTS	REFLEVEL	0	Label level to be used as reference
	TESTLEVEL	0	Label level to be scored
	STRIPCONTEXT	F	Strip triphone contexts
	IGNORECASE	F	If enabled, converts labels to uppercase
	NISTSCORE	F	Use NIST formatting
	PHRASELABEL	SENT	Label for phrase level statistics
	PHONELABEL	WORD	Label for word level statistics
	SPEAKERMASK	NULL	If set then report on a per speaker basis
HVITE	RECOUTPREFIX	NULL	Prefix for direct audio output name
	RECOUTSUFFIX	NULL	Suffix for direct audio output name
	SAVEBINARY	F	Save transforms as binary
HLSTATS	DISCOUNT	0.5	Discount constant for backoff bigrams
HLIST	AUDIOSIG	0	Audio signal number for remote control
	SOURCERATE	0.0	Sample rate of source in 100ns units
	TRACE	0	Trace setting

Table 15.2: Tool Specific Configuration Variables

Chapter 16

Error and Warning Codes

When a problem occurs in any HTK tool, either error or warning messages are printed.

If a warning occurs then a message is sent to standard output and execution continues. The format of this warning message is as follows:

```
WARNING [-nnnn] Function: 'Brief Explanation' in HTool
```

The message consists of four parts. On the first line is the tool name and the error number. On the second line is the function in which the problem occurred and a brief textual explanation. The reason for sending warnings to standard output is so that they are synchronised with any trace output.

If an error occurs a number of error messages may be produced on standard error. Many of the functions in the HTK Library do not exit immediately when an error condition occurs, but instead print a message and return a failure value back to their calling function. This process may be repeated several times. When the HTK Tool that called the function receives the failure value, it exits the program with a fatal error message. Thus the displayed output has a typical format as follows:

```
ERROR [+nnnn] FunctionA: 'Brief explanation'  
ERROR [+nnnn] FunctionB: 'Brief explanation'  
ERROR [+nnnn] FunctionC: 'Brief explanation'  
FATAL ERROR - Terminating program HTool
```

Error numbers in HTK are allocated on a module by module and tool by tool basis in blocks of 100 as shown by the table shown overleaf. Within each block of 100 numbers the first 20 (0 - 19) and the final 10 (90-99) are reserved for standard types of error which are common to all tools and library modules.

All other codes are module or tool specific.

16.1 Generic Errors

- +??00 Initialisation failed
The initialisation procedure for the tool produced an error. This could be due to errors in the command line arguments or configuration file.
- +??01 Facility not implemented
HTK does not support the operation requested.
- +??05 Available memory exhausted
The operation requires more memory than is available.
- +??06 Audio not available
The audio device is not available, either there is no driver for the current machine, the library was compiled with NO_AUDIO set or another process has exclusive access to the audio device.

HCOPY	1000-1099	HShell	5000-5099
HList	1100-1199	HMem	5100-5199
HLEd	1200-1299	HMath	5200-5299
HLStats	1300-1399	HSigP	5300-5399
HDMan	1400-1499		
HSLab	1500-1599	HAudio	6000-6099
		HVQ	6100-6199
		HWave	6200-6299
HCompV	2000-2099	HParm	6300-6399
HInit	2100-2199	HLabel	6500-6599
HRest	2200-2299		
HERest	2300-2399	HGraf	6800-6899
HSmooth	2400-2499		
HQuant	2500-2599	HModel	7000-7099
HHEd	2600-2699	HTrain	7100-7199
HEAdapt	2700-2799	HUtil	7200-7299
		HFB	7300-7399
		HAdapt	7400-7499
HBuild	3000-3099		
HParse	3100-3199	HDict	8000-8099
HVite	3200-3299	HLM	8100-8199
HResults	3300-3399	HNet	8200-8299
HSGen	3400-3499	HRec	8500-8599

- +??10 Cannot open file for reading
Specified file could not be opened for reading. The file may not exist or the filter through which it is read may not be set correctly.
- +??11 Cannot open file for writing
Specified file could not be opened for writing. The directory may not exist or be writable by the user or the filter through which the file is written may not be set correctly.
- +??13 Cannot read from file
Cannot read data from file. The file may have been truncated, incorrectly formatted or the filter process may have died.
- +??14 Cannot write to file
Cannot write data to file. The file system is full or the filter process has died.
- +??15 Required function parameter not set
You have called a library routine without setting one of the arguments.
- +??16 Memory heap of incorrect type
Some library routines require you to pass them a heap of a particular type.
- +??19 Command line syntax error
The command line is badly formed, refer to the manual or the command summary printed when the command is executed without arguments.
- +??9? Sanity check failed
Several functions perform checks that structures are self consistent and that everything is functioning correctly. When these sanity checks fail they indicate the code is not functioning as intended. These errors should not occur and are not correctable by the user.

16.2 Summary of Errors by Tool and Module

HCOPY

- +1030 Non-existent part of file specified
HCOPY needed to access a non-existent part of the input file. Check that the times are specified correctly, that the label file contains enough labels and that it corresponds to the data file.

- ±1031 Label file formatted incorrectly
HCOPY is only able to properly copy label files with the same number of levels/alternatives. When using labels with multiple alternatives only the first one is used to determine segment boundaries.
- +1032 Appending files of different type/size/rate
Files that are joined together must have the same parameter kind and sample rate.
- 1089 ALIEN format set
Input/output format has been set to ALIEN, ensure that this was intended.

HLIST

HLED

- +1230 Edit script syntax error
The HLED command script contains a syntax error, check the input script against the descriptions of each command in section 14.9 or obtained by running HLEd -Q.
- ±1231 Operation invalid
You have either exceeded HLED limits on the number of boundaries that can be specified, tried to perform an operation on a non-existent level or tried to sort an auxiliary level into time order. None of these operations are supported.
- +1232 Cannot find pronunciation
The dictionary does not contain a valid pronunciation (only occurs when attempting expansion from a dictionary).
- 1289 ALIEN format set
Input/output format has been set to ALIEN, ensure that this was intended.

HLSTATS

- +1328 Load/Make HMMSet failed
The model set could not be loaded due to either an error opening the file or the data within being inconsistent.
- ±1330 No operation specified
You have invoked HLSTATS but have not specified an operation to be performed.
- 1389 ALIEN format set
Input format has been set to ALIEN, ensure that this was intended.

HDMAN

- ±1430 Limit exceeded
HDMAN has several built in limits on the number of different pronunciation, phones, contexts and command arguments. This error occurs when you try to exceed one of them.
- ±1431 Item not found
Could not find item for deletion. Check that it actually occurs in the dictionary.
- ±1450 Edit script file syntax error
The HDMAN command script contains a syntax error, check the input script against the descriptions of each command in section 14.4 or obtained by running HDMan -Q.
- ±1451 Dictionary file syntax error
One of the input dictionaries contained a syntax error. Ensure that it is in a HTK readable form (see section 12.7).
- ±1452 Word out of order in dictionary error
Entries in the dictionary must be sorted into alphabetical (ASCII) order.

HSLAB

- 1589 ALIEN format set
Input/output format has been set to ALIEN, ensure that this was intended.

HCOMPV

- +2020 HMM does not appear in HMMSet
Supplied HMM filename does not appear in HMMSet. Check correspondence between HMM filename and HMMSet.
- +2021 Not enough data to calculate variance
There are not enough frames of data to evaluate a reliable estimate of variance. Use more data.
- +2028 Load/Make HMMSet failed
The model set could not be loaded due to either an error opening the file or the data within being inconsistent.
- +2030 Needs continuous models
HCompV can only operate on models with an HMM set kind of `PLAINHS` or `SHAREDHS`.
- +2050 Data does not match HMM
An aspect of the data does not match the equivalent aspect in the HMMSet. Check the parameter kind of the data.
- 2089 ALIEN format set
Input format has been set to `ALIEN`, ensure that this was intended.

HINIT

- +2120 Unknown update flag
Unknown flag set by `-u` option, use combinations of `tmvw`.
- +2121 Too little data
Not enough data to reliably estimate parameters. Use more training data.
- +2122 Segment with fewer frames than model states
Segment may be too short to be matched to model, do not use this segment for training.
- +2123 Cannot mix covariance kind in a single mix
Covariance kind of all mixture components in any one state must be the same.
- +2124 Bad covariance kind
Covariance kind of mixture component must be either `FULLC` or `DIAGC`.
- +2125 No best mix found
The Viterbi mixture component allocation failed to find a most likely component with this data. Check that data is not corrupt and that parameter values produced by the initial uniform segmentation are reasonable.
- +2126 No path through segment
The Viterbi segmentation failed to find a path through model with this data. Check that data is not corrupt and that a valid path exists through the model.
- +2127 Zero occurrence count
Parameter has had no data assigned to it and cannot be updated. Ensure that each parameter can be estimated by using more training data or fewer parameters.
- +2128 Load/Make HMMSet failed
The model set could not be loaded due to either an error opening the file or the data within being inconsistent.
- +2129 HMM not found
HMM missing from HMMSet. Check that the HMMSet is complete and has not been corrupted.
- +2150 Data does not match HMM
An aspect of the data does not match the equivalent aspect in the HMMSet. Check the parameter kind of the data.
- +2170 Index out of range
Trying to access a mixture component or VQ index beyond the range in the current HMM.
- 2189 ALIEN format set
Input format has been set to `ALIEN`, ensure that this was intended.

HREST

- +2220 Unknown update flag
Unknown flag set by `-u` option, use combinations of `tmvw`.
- +2221 Too few training examples
There are fewer training examples than the minimum set by the `-m` option (default 3). Either reduce the value specified by `-m` or use more training examples.
- +2222 Zero occurrence count
Parameter has had no data assigned to it and cannot be updated. Ensure that each parameter can be estimated by using more training data or fewer parameters.
- +2223 Floor too high
Mix weight floor has been set so high that the sum over all mixture components exceeds unity. Reduce the floor value.
- 2225 Defunct Mix X.Y.Z
Not enough training data to re-estimate the covariance vector of mixture component Z in stream Y of state X. The weight of the mixture component is set to 0.0 and it will never recover even with further training.
- +2226 No training data
None of the supplied training data could be used to re-estimate the model. Data may be corrupt or has been floored.
- +2228 Load/Make HMMSet failed
The model set could not be loaded due to either an error opening the file or the data within being inconsistent.
- +2250 Data does not match HMM
An aspect of the data does not match the equivalent aspect in the HMMSet. Check the parameter kind of the data.
- 2289 ALIEN format set
Input format has been set to `ALIEN`, ensure that this was intended.

HEREST

- +2320 Unknown update flag
Unknown flag set by `-u` option, use combinations of `tmvw`.
- +2321 Load/Make HMMSet failed
The model set could not be loaded due to either an error opening the file or the data within being inconsistent.
- 2326 No transitions
No transition out of an emitting state, ensure that there is a transition path from beginning to end of model.
- +2327 Floor too high
Mix weight floor has been set so high that the sum over all mixture components exceeds unity. Reduce the floor value.
- +2328 No mixtures above floor
None of the mixture component weights are greater than the floor value, reduce the floor value.
- 2330 Zero occurrence count
Parameter has had no data assigned to it and cannot be updated. Ensure that each parameter can be estimated by using more training data or fewer parameters.
- 2331 Not enough training examples
Model was not updated as there were not enough training examples. Either reduce the minimum specified by `-m` or use more data.
- 2389 ALIEN format set
Input format has been set to `ALIEN`, ensure that this was intended.

HSMOOTH

- +2420 Unknown update flag
Unknown flag set by `-u` option, use combinations of `tmvw`.
- +2421 Invalid HMM set kind
HSMOOTH can only be used if HMM set kind is either DISCRETE or TIED.
- +2422 Too many monophones in list
HSMOOTH is limited to HMMSets containing fewer than 500 monophones.
- +2423 Different number of states for smoothing
Monophones and context-dependent models have differing numbers of states.
- 2424 No transitions
No transition out of an emitting state, ensure that there is a transition path from beginning to end of model.
- +2425 Floor too high
Mix weight floor has been set so high that the sum over all mixture components exceeds unity. Reduce the floor value.
- 2427 Zero occurrence count
Parameter has had no data assigned to it and cannot be updated. Ensure that each parameter can be estimated by using more training data or fewer parameters.
- 2428 Not enough training examples
Model was not updated as there were not enough training examples. Either reduce the minimum specified by `-m` or use more data.
- +2429 Load/Make HMMSet failed
The model set could not be loaded due to either an error opening the file or the data within being inconsistent.

HQUANT

- +2530 Stream widths invalid
The chosen stream widths are invalid. Check that these match the parameter kind and are specified correctly.
- +2531 Data does not match codebook
Ensure that the parameter kind of the data matches that of the codebook being generated.

HHEd

- +2628 Load/Make HMMSet failed
The model set could not be loaded due to either an error opening the file or the data within being inconsistent.
- ±2630 Tying null or different sized items
You have executed a tie command on items which do not have the appropriate structure or the structures are not matched. Ensure that the item list refers only to the items that you wish to tie together.
- 2631 Performing operation on no items
The item list was empty, no operation is performed.
- +2632 Command parameter invalid
The parameters to the command are invalid either because they refer to parts of the model that do not exist (for instance a state that does not appear in the model) or because they do not represent an acceptable value (for instance HMMSet kind is not PLAINHS, SHAREDHS, TIEDHS or DISCRETEHS).
- +2634 Join parameters invalid or not set
Make sure than the join parameters (set by the `JO` command) are reasonable. In particular take care that the floor is low enough to ensure that when summed over all the mixture components the sum is below 1.0.
- +2635 Cannot find matching item
Search for specified item was unsuccessful. When this occurs with the `CL` or `MT` commands ensure that the appropriate monophone/biphone models are in the current HMMSet.

- 2637 Small gConst
A small gConst indicates a very low variance in that particular Gaussian. This could be indicative of over-training of the models.
- 2638 No typical state
When tying states together a search is performed for the distribution with largest variance and all tied states share this distribution. If this cannot be found the first in the list will be used instead.
- 2639 Long macro name
In general macro names should not exceed 20 characters in length.
- +2640 Not implemented
You have asked HHED to perform a function that is not implemented.
- +2641 Invalid stream split
The specified number/width of the streams does not agree with the parameter kind/vector size of the models.
- +2650 Edit script syntax error
The HHED command script contains a syntax error, check the input script against the descriptions of each command in section 14.7 or obtained by running HHED -Q.
- +2651 Command range error
The value specified in the command script is out of range. Ensure that the specified state exists and the the value given is valid.
- ±2655 Stats file load error
Either loading occupation statistics for the second time or executing an operation that needs the statistics loaded without loading them.
- +2660 Trees file syntax error
The trees file format did not correspond to that expected. Ensure that the file is complete and has not been corrupted.
- +2661 Trees file macro/question not recognised
The question or macro referred to does not exist. Ensure that the file is complete and has not been corrupted.
- +2662 Trying to sythesize for unknown model
There is no tree or prototype model for the new context. Ensure that a tree has been constructed for the base phone.
- ±2663 Invalid types to tree cluster
Tree clustering will only work for single Gaussian diagonal covariance untied models of similar topology.

HEADAPT

- +2720 Invalid update flag
Means (m) and variances (v) can be updated, check that the update flags are correct.
- 2721 No update requested
Mean transformation will be updated even though none specified.
- +2728 Load/Make HMMSet failed
The model set could not be loaded due to either an error opening the file or the data within being inconsistent.
- +2730 Invalid HMM kind
MLLR adaptation only available for PLAIN or SHARED systems.
- +2731 Invalid HMM kind
MLLR adaptation only currently available for single stream data.
- 2740 Invalid Adaptation Mode
Incremental MAP adaptation is not supported, use in static mode.
- 2741 Invalid Adaptation Mode
MLLR and MAP adaptation does not support incremental adaptation. If incremental MLLR and static MAP adaptation is required, perform the steps separately.

- 2789 Alien format set
Input format has been set to ALIEN, ensure that this was intended.
- HBUILD
- ±3030 Mismatch between command line and language model
Ensure that the !ENTER and !EXIT words are correctly defined and that the supplied files are of the appropriate type.
- ±3031 Unknown word
Ensure that the word list corresponds to the language model/lattice supplied.
- HPARSE
- ±3130 Variable not defined
You have referenced a network that has not yet been defined. Check that all networks are defined before they are referenced.
- ±3131 Loop or word expansion error
There is either a mismatch between the WD_BEGIN WD_END pairs or a triphone loop is badly formed.
- ±3132 Dictionary error
When generating a dictionary a word exceeded the maximum number of phones, a word occurred twice or no dictionary was produced.
- ±3150 Syntax error in HParse file
The HPARSE network definition contains a syntax error, check the input file against the network description in section 14.12.
- HVITE
- +3228 Load/Make HMMSet failed
The model set could not be loaded due to either an error opening the file or the data within being inconsistent.
- ±3230 Unsupported operation
HVITE is not able to perform the operation requested
- ±3231 Data does not match HMMs
There is a mismatch between the data file and the HMMSet. Ensure that the data is parameterised in the correct format and the configuration parameters match those used during training.
- +3232 MMF Load Error
The HMMSet does not contain a well-formed regression class tree.
- 3289 ALIEN format set
Input/output format has been set to ALIEN, ensure that this was intended.
- HRESULTS
- 3330 Empty file
The file was empty and will be skipped.
- +3331 Unknown label
The label did not appear in the list supplied to HResults. This error will only occur if calculating confusion matrices so normally the contents of the word list file will have no effect on results.
- +3332 Too many labels
HRESULTS will only generate confusion statistics for a small number of labels.
- ±3333 Cannot calculate word spot results
When calculating word spotting results the label files need to have both times and scores present.
- 3389 ALIEN format set
Input format has been set to ALIEN, ensure that this was intended.

HSGEN

- 3420 Network malformed
The word network is malformed. The information in a node (word and following arcs) is set incorrectly.

HSHELL

- +5020 Command line processing error
- +5021 Command line argument type error
- +5022 Command line argument range error
The command line is badly formed. Ensure that it matches the syntax and values expected by the command (check the manual page or the syntax obtained by running HTOOL without any arguments).
- +5050 Configuration file format error
HSHELL was unable to parse the configuration. Check that it is of the format described in section 4.3.
- +5051 Script file format error
Check that the script file is just a list of file names and that if any file names are quoted that the quotes occur in pairs.
- +5070 Module version syntax error
A module registered with HShell with an incorrectly formatted version string (which should be of the form "`!HVER!HModule: Vers.str [WHO DD/MM/YY]`").
- +5071 Too many configuration parameters
The size of the buffer used by one of the tools or modules to read its configuration parameters was exceeded. Either reduce the total number of configuration parameters in the file or make more of them specific to their particular module rather than global.
- +5072 Configuration parameter of wrong type
The configuration parameter is of the wrong type. Check that its type agrees with that shown in chapter 15.
- +5073 Configuration parameter out of range
The configuration parameter is out of range.

HMEM

- +5170 Heap parameters invalid
You have tried to create a heap with unreasonable parameters. Adjust these so that the growth factor is positive and the initial block size is no larger than the maximum. For MSTAK the element size should be 1.
- +5171 Heap not found
The specified heap could not be found, ensure that it has not been deleted or memory overwritten.
- +5172 Heap does not support operation
The heap is of the wrong type to support the requested operation. In particular it is not possible to `Reset` or `Delete` a `CHEAP`.
- +5173 Wrong element size for MHEAP
You have tried to allocate an item of the wrong size from a `MHEAP`. All items on a `MHEAP` must be of the same size.
- +5174 Heap not initialised
You have tried to allocate an item on a heap that has not yet been created. Ensure that `CreateHeap` is called to initialise the heap before any items are allocated from it.
- +5175 Freeing unseen item
You have tried to free an item from the wrong heap. This can occur if the wrong heap is specified, the item pointer has been corrupted or the item has already been freed implicitly by a `Reset/DeleteHeap` call.

HMATH

- +5220 Singular covariance matrix
The covariance matrix was not invertible. This may indicate a lack of training data or linearly dependent parameters.
- +5270 Size mismatch
The input parameters were of incompatible sizes.
- +5271 Log of negative
Result would be logarithm of a negative number.

HSigP

- +5320 No results for WaveToLPC
Call did not include Vectors for the results.
- +5321 Vector size mismatch
Input vectors were of mismatched sizes.
- 5322 Clamped samples during zero mean
During a zero mean operation samples were clipped as they were outside the allowable range.

HAUDIO

- +6020 Replay buffer not active
Attempt to access a replay buffer when one was not attached.
- +6021 Cannot StartAudio without measuring silence
An attempt was made to start audio input through the silence detector without first measuring or supplying the background silence values.
- +6070 Audio frame size/rate invalid
The choice of frame period and window duration are invalid. Check both these and the sample rate.
- 6071 Setting speech threshold below silence
The thresholds used in the speech detector have been set so that the threshold for detecting speech is set below that of detecting silence.

HVQ

- +6150 VQ file format error
The VQ file was incorrectly formatted. Ensure that the file is complete and has not been corrupted.
- +6151 VQ file range error
A value from the VQ file was out of range. Ensure that the file is complete and has not been corrupted.
- +6170 Magic number mismatch
The VQ magic number (normally based on parameter kind) does not match that expected. Check that the parameter kind used to quantise the data and create the VQ table matches the current parameter kind.
- +6171 VQ table already exists
All VQ tables must have distinct names. This error will occur if you try to create or load a VQ table with the same name as one already loaded.
- +6172 Invalid covariance kind
Entries in VQ tables must have either NULLC, FULLC or INVDIAGC covariance kind.
- +6173 Node not in table
A node was missing from the VQ table. Ensure that the VQ table was properly created or that the file was complete.
- +6174 Stream codebook mismatch
The number or size of streams in the VQ table does not match that requested.

HWAVE

- +6220 Cannot fseek/ftell
Unless the wave file is read through a pipe fseek and ftell are expected to work correctly so that HWAVE can calculate the file size. If this error occurs when using an input pipe, supply the number of samples in the file using the configuration variable NSAMPLES.
- +6221 File appears to be a infinite
HWAVE cannot determine the size of the file.
- +6230 Config parameter not set
A necessary configuration parameter has not been set. Determine the correct value and place this in the configuration file before re-invoking the tool.
- +6250 Premature end of header
HWAVE could not read the complete file header.
- +6251 Header contains invalid data
HWAVE was unable to successfully parse the header. The header is invalid, of the wrong type or be a variation that HWAVE does not handle.
- +6252 Header missing essential data
The header was missing a piece of information necessary for HWAVE to load the file. Check the processing of the input file and re-process if necessary.
- +6253 Premature end of data
The file ended before all the data was read correctly. Check that the file is complete, has not been corrupted and where necessary NSAMPLES is set correctly.
- +6254 Data formatted incorrectly
The data could not be decoded properly. Check that the file was complete and processed correctly.
- +6270 File format invalid
The file format is not valid for the operation requested.
- +6271 Attempt to read outside file
You have tried to read a sample outside of the waveform file.

HPARM

- +6320 Configuration mismatch
The data file does not match the configuration. Check the configuration file is correct.
- +6321 Invalid parameter kind
Parameter kind is not valid. Check the configuration file.
- +6322 Conversion not possible
The specified conversion is not possible. Check the configuration is correct and re-code the data from waveform files if necessary.
- +6323 Audio error
An audio error has been detected. Check the HAUDIO configuration and the audio device.
- +6324 Buffer not initialised
Ensure that the buffer is used in the correct manner.
- +6325 Silence detection failed
The silence detector was not initialised correctly before use.
- +6328 Load/Make HMMSet failed
The model set could not be loaded due to either an error opening the file or the data within being inconsistent.
- +6350 CRC error
The CRC does not match that of the data. Check the data file is complete and has not been corrupted.
- 6351 Byte swapping not possible
HPARM will attempt to byte swap parameter files but this may not work if the floating point representation of the machine that generated the file is different from that which is reading it.

- +6352 File too short to parameterise
The file does not contain enough data to produce a single observation. Check the file is complete and not corrupt. If it is, it should be discarded.
- +6370 Unknown parameter kind
The specified parameter kind is not recognised. Refer to section 5.15 for a list of allowable parameter kinds and qualifiers.
- +6371 Invalid parameters for coding
The chosen parameters are not valid for coding. Choose different ones.
- +6372 Stream widths not valid
Cannot split the data into the specified number of streams. Check that the parameter kind is correct and matches any models used.
- +6373 Buffer/observation mismatch
The observation parameter kind should match that of the input buffer. Check that the configuration parameter kind is correct and matches that of any models used.
- +6374 Buffer size too small for window
Calculation of delta parameters requires a window larger than the buffer size chosen. Increase the size of the buffer.
- +6375 Frame not in buffer
An attempt was made to access a frame that does not appear in the buffer. Make sure that the file actually contains the specified frame.
- +6376 Mean/Variance normalisation failed
The mean or variance normalisation vector from the file specified by the normalisation dir and mask cannot be applied. Make sure the file format is correct and the vectors are of the right dimension.

HLABEL

- +6520 MLF index out of range
An attempt was made to access an MLF that has not been loaded or to load too many MLFs.
- +6521 fseek/ftell not possible
HLABEL needs random access to MLFs. This error is generated when this is not possible (for instance if access is via a pipe).
- +6550 HTK format error
- +6551 MLF format error
- +6552 TIMIT format error
- ±6553 ESPS format error
- +6554 SCRIBE format error
A label file was formatted incorrectly. Label file formats are described in chapter 6.
- +6570 Level out of range
Attempted to access a non-existent label level. Check that the correct label file has been loaded.
- +6571 Label out of range
Attempted to access a non-existent label. Check that the correct label file has been loaded and that the correct level is being accessed.
- +6572 Invalid format
The specified file format is not valid for the particular operation.

HMODEL

- +7020 Cannot find physical HMM
No physical HMM exists for a particular logical model. Check that the HMMSet was loaded or created correctly.
- +7021 INVDIAG internal format
Attempts to load or save models with INVDIAG covariance kind will fail as this is a purely internal model format.

- ±7023 varFloor should be variance floor
HMODEL reserves the macro name varFloorN as the variance floor for stream N. These should be variance macros (type v) of the correct size for the particular stream.
- +7024 Variance tending to 0.0
A variance has become too low. Start using a variance floor or increase the amount of training data.
- +7030 HMM set incomplete or inconsistent
The HMMSet contained missing or inconsistent data. Check that the file is complete and has not been corrupted.
- +7031 HMM parameters inconsistent
Some model parameters were inconsistent. Check that the file is complete and has not been corrupted.
- ±7032 Option mismatch
All HMMs in a particular set must have consistent options.
- +7035 Unknown macro
Macro does not exist. Check that the name is correct and appears in the HMMSet.
- +7036 Duplicate macro
Attempted to create a macro with the same name as one already present. Choose a different name.
- +7037 Invalid macro
Macro had invalid type. See section 7.3 describes the allowable macro types.
- +7050 Model file format error
- +7060 HMM List format error
The file was formatted incorrectly. Check the file is complete and has not been corrupted.
- +7070 Invalid HMM kind
Invalid HMMSet kind. Check that this is specified correctly.
- +7071 Observation not compatible with HMMSet
Attempted to calculate an observation likelihood for an observation not compatible with the HMMSet. Check that the parameter kind is set correctly.

HTRAIN

- +7120 Clustering failed
Almost certainly due to a lack of data, reduce the number of clusters requested or increase amount of data.
- +7150 Accumulator file format error
Cannot read an item from an accumulator file. Check that file is complete and not corrupted.
- +7170 Unsupported covariance kind
Covariance kind must be either FULLC, DIAGC or INVDIAGC.
- +7171 Item out of range
Attempt made to access data beyond expected range. Check that the item number is correct.
- +7172 Tree size must be power of 2
Requested codebook size must be a power of 2 when using tree based clustering.
- 7173 Segment empty
Empty data segment in file. Check that file has not become corrupted and that the start and end segment times are correct.

HUTIL

- +7220 HMMSet empty
A scan was initiated for a HMMSet with no members.
- +7230 Item list parse error
The item list syntax was incorrect. Check the item list specification in section 14.7.

- +7231 Item list type error
Each item in a particular list should be of the same type and size.
- +7250 Stats file format error
Stats file is of wrong format. Note the format of the stats file has changed in HTK_V2.0 and old files will need converting to the new format.
- +7251 Stats file model error
A model name encountered in the stats file is invalid check that the model set corresponds to that used to generate the stats file and that the stats file is complete and has not been corrupted.
- +7270 Accessing non-existent macro
Attempt to perform operation on non-existent macro.
- +7271 Member id out of range
Attempt to perform set operation on out of range member.

HFB

- +7321 Unknown model
Model in HMM List not found in HMMSet, check that the correct HMM List is being used.
- +7322 Invalid output probability
Mixture component probability has not been set. This should not occur in normal use.
- +7323 Beta prune failed on taper
Utterance is possibly too short for minimum duration of model sequence. Check transcription.
- 7324 No path through utterance
No path was found on the beta training pass, relax the pruning threshold.
- 7325 Empty label file
No labels found in label file, check label file.
- +7326 Single-pass retraining data mismatch
Paired training files must contain the same number of observations. Use original data to re-parameterise.
- ±7332 HMM with unreachable states
HMM has an unreachable state, check transition matrix.
- 7333 Transition matrix with discontinuity
Check transition matrix.
- +7350 Data does not match HMM
An aspect of the data does not match the equivalent aspect in the HMMSet. Check the parameter kind of the data.

HADAPT

- ±7410 TMF load error
Can't open the input TMF.
- 7420 Can't calculate the auxilliary function
Only possible with mean and variance transforms, check update flags.
- +7421 MMF load error
MMF does not contain a HMM identifier, use HHed to generate one.
- +7422 MMF load error
MMF does not contain a regression class, use HHed to generate one
- +7425 Mismatch number of Gaussian components
The number of Gaussian components found for a regression class does not much the number expected. Check the MMF is not corrupted.
- +7430 MMF load error
Can't find the regression class tree.

- +7431 Can't add regression class accumulate
Problem adding a frame of accumulation at the component level.
- ±7440 Symbol not found
Could not read in symbol.
- +7450 Invalid HMM kind
Can only adapt PLAIN and SHARED systems.
- +7460 Can't create block transformation matrix
Check number of blocks compatible with vector size.
- +7470 MAP weight hook is NULL

HDICT

- +8050 Dictionary file format error
The dictionary file is not correctly formatted. Section 12.7 describes the HTK dictionary file format.

HLM

- +8150 LM syntax error
The language model file was formatted incorrectly. Check the file is complete and has not been corrupted.
- ±8151 LM range error
The specified value(s) for the language model probability are not valid. Check the input files are correct.

HNET

- +8220 No such word
The specified word does not exist or does not have a valid pronunciation.
- 8221 Duplicate pronunciations removed
During network generations duplicate identical pronunciations of the same word are removed.
- +8230 Contexts not consistent
HNET can only deal with the standard HTK method for specifying context `left-phone+right` and will only allow context free phones if they are context independent and only form part of the word. This may be indicative of an inconsistency between the symbols in the dictionary and the hmms as defined. There may be a model/phone in the dictionary that has not been defined in the HMM list or may not have a corresponding model. See also section 12.8 on context expansion.
- +8231 No such model
A particular model could not be found. Make sure that the network is being expanded in the correct fashion and then ensure that your HMM list will cover all required contexts.
- +8232 Lattice badly formed
Could not convert lattice to network. The lattice should have a single well defined start and a single well defined end. When cross word expansion is being performed the number of !NULL words that can be concatenated in a string is limited.
- +8250 Lattice format error
The lattice file is formatted incorrectly. Ensure that the lattice is of the format described in chapter 17.
- +8251 Lattice file data error
The value specified in the lattice file is invalid.
- +8252 Lattice file with multiple start/end nodes
A lattice should have only one well defined start node and one well defined end node.
- +8253 Lattice with invalid sub lattices
The sub lattices referred to by the main lattices are malformed.

HREC

- ±8520 Invalid HMM
One of the HMMs in the network is invalid. Check that the HMMSet has been correctly initialised.
- +8521 Network structure invalid
The network is incorrectly structured. Take care to avoid loops that can be traversed without consuming observations (this may occur if you introduce any 'tee' words in which all the models making up that word contain tee-transitions). Also ensure that the recogniser and the network have been created and initialised correctly.
- +8522 Lattice structure invalid
The lattice was incorrectly formed. Ensure that the lattice was created properly.
- ±8570 Recogniser not initialised correctly
Ensure the recogniser is initialised and used correctly.
- +8571 Data does not match HMMs
The observation does not match the HMM structure. Check the parameter kind of the data and ensure that the data is matched to the HMMs.

HGRAF

- +6870 X11 error
Ensure that the DISPLAY variable is set and that the X11 window system is configured correctly.

Chapter 17

HTK Standard Lattice Format (SLF)

17.1 SLF Files

Lattices in HTK are used for storing multiple hypotheses from the output of a speech recogniser and for specifying finite state syntax networks for recognition. The HTK standard lattice format (SLF) is designed to be extensible and to be able to serve a variety of purposes. However, in order to facilitate the transfer of lattices, it incorporates a core set of common features.

An SLF file can contain zero or more sub-lattices followed by a main lattice. Sub-lattices are used for defining sub-networks prior to their use in subsequent sub-lattices or the main lattice. They are identified by the presence of a `SUBLAT` field and they are terminated by a single period on a line by itself.

A lattice consists of optional header information followed by a sequence of node definitions and a sequence of links (arc) definitions. Nodes and links are numbered and the first definition line must give the total number of each.

Each link represents a word instance occurring between two nodes, however for more compact storage the nodes often hold the word labels since these are frequently common to all words entering a node (the node effectively represents the end of several word instances). This is also used in lattices representing word-level networks where each node is a word end, and each arc is a word transition.

Each node may optionally be labelled with a word hypothesis and with a time. Each link has a start and end node number and may optionally be labelled with a word hypothesis (including the pronunciation variant, acoustic score and segmentation of the word hypothesis) and a language model score.

17.2 Format

The format is designed to allow optional information that at its most detailed gives full identity, alignment and score (log likelihood) information at the word and phone level to allow calculation of the alignment and likelihood of an individual hypothesis. However, without scores or times the lattice is just a word graph. The format is designed to be extendible. Further field names can be defined to allow arbitrary information to be added to the lattice without making the resulting file unreadable by others.

The lattices are stored in a text file as a series of fields that form two blocks:

- A header, specifying general information about the lattice.
- The node and link definitions.

Either block may contain comment lines, for which the first character is a `#` and the rest of the line is ignored.

All non-comment lines consist of fields, separated by white space. Fields consist of an alphanumeric field name, followed by a delimiter (the character `=` or `'`) and a (possibly “quoted”) field value. All letters in the name are significant. Single character field names are reserved for fields

defined in the specification and single character abbreviations may be used for many of the fields defined below.

The convention used to define the current field names is that lower case is used for optional fields and upper case is used for required fields. The meaning of field names can be dependent on the context in which they appear.

The header must include a field specifying which utterance was used to generate the lattice and a field specifying the version of the lattice specification used. The header is terminated by a line which defines the number of nodes and links in the lattice.

The node definitions are optional but if included each node definition consists of a single line which specifies the node number followed by optional fields that may (for instance) define the time of the node or the word hypothesis ending at that node.

The link definitions are required and each link definition consists of a single line which specifies the link number as well as the start and end node numbers that it connects to and optionally other information about the link such as the word identity and language model score. If word identity information is not present in node definitions then it must appear in link definitions.

17.3 Syntax

The following rules define the syntax of a lattice using the first character of the field names defined below to represent the entire field. Any unrecognised fields may be ignored and no user defined fields may share first characters with pre-defined field names. parentheses () indicate that terms within may appear in any order and braces {} indicate that terms within are optional.

```
lattice_file = header size_def terms
header       = ( V {\n} U {\n} S {\n} ) { base {\n} start {\n} end {\n} dir {\n}
              tscale {\n} lmscale {\n} lmpen {\n} }
size_def     = ( L N ) \n
terms        = node_def | link_def { terms }
node_def     = I { t W v d a } \n
link_def     = J ( S E { W v d a n l} ) \n
```

17.4 Field Types

The currently defined fields are as follows:-

Field	abbr	o c	Description
Header fields			
VERSION=%s	V	o	Lattice specification adhered to
UTTERANCE=%s	U	o	Utterance identifier
SUBLAT=%s	S	o	Sub-lattice name
base=%f		o	LogBase for Likelihoods (0.0 not logs,default base e)
lmname=%s		o	Name of Language model
lmscale=%f		o	Scaling factor for language model
wdpenalty=%f		o	Word insertion penalty
Lattice Size fields			
NODES=%d	N	c	Number of nodes in lattice
LINKS=%d	L	c	Number of links in lattice
Node Fields			
I=%d			Node identifier. Starts node information
time=%f	t	o	Time from start of utterance (in seconds)
WORD=%s	W	wc	Word (If lattice labels nodes rather than links)
L=%s		wc	Substitute named sub-lattice for this node
var=%d	v	wo	Pronunciation variant number
Link Fields			

```

J=%d          Link identifier.  Starts link information
START=%d      S c Start node number (of the link)
END=%d        E c End node number (of the link)
WORD=%s       W wc Word (If lattice labels links rather than nodes)
var=%d        v wo Pronunciation variant number
div=%s        d wo Segmentation [ d=(label[,duration,like]:)+ ]
acoustic=%f   a wo Acoustic likelihood of link
ngram=%f      n o NGram likelihood of link
language=%f   l o General language model likelihood of link

```

Note. The word identity (and associated 'w' fields var,div and acoustic) must appear on either the link or the end node (forwards lattices) or the start node (backwards lattices) of the link.
abbr is a possible single character abbreviation for the field name
o|c indicates whether field is optional or compulsory.

17.5 Example SLF file

The following is a real lattice (generated by the HTK Large Vocabulary Recogniser with a 20k WSJ bigram) with word labels occurring on the links.

Note that the !ENTER and !EXIT "words" model initial and final silence.

```

# Lattice generated by CU-HTK 23 Feb 94
#
# File : "/data/wsj/wsj1/si_dt_20/4k0/4k0c030t.wv2"
#
# Best hypothesis "!ENTER IT DIDN'T ELABORATE !EXIT" Score=-20218.25
#
# Language model scores from "/lib/baseline-lm/bg-boc-lm20o.nvp".
# Dictionary used "/lib/dictionaries/dragon/wsj.sls".
# Acoustic scores from "/models/htk2/hmm11".
#
# Header
#
VERSION=1.1
UTTERANCE=4k0c030t
lmscale=16.0
wdpenalty=0.0
#
# Size line
#
N=24 L=39
#
# Node definitions
#
I=0 t=0.00
I=1 t=0.25
I=2 t=0.26
I=3 t=0.61
I=4 t=0.62
I=5 t=0.62
I=6 t=0.71
I=7 t=0.72
I=8 t=0.72
I=9 t=0.72
I=10 t=0.72
I=11 t=0.72

```

```

I=12  t=0.72
I=13  t=0.73
I=14  t=0.78
I=15  t=0.78
I=16  t=0.80
I=17  t=0.80
I=18  t=0.81
I=19  t=0.81
I=20  t=1.33
I=21  t=2.09
I=22  t=2.09
I=23  t=2.85
#
# Link definitions.
#
J=0    S=0    E=1    W=!ENTER    v=0  a=-1432.27  l=0.00
J=1    S=0    E=2    W=!ENTER    v=0  a=-1500.93  l=0.00
J=2    S=0    E=3    W=!ENTER    v=0  a=-3759.32  l=0.00
J=3    S=0    E=4    W=!ENTER    v=0  a=-3829.60  l=0.00
J=4    S=1    E=5    W=TO        v=3  a=-2434.05  l=-87.29
J=5    S=2    E=5    W=TO        v=1  a=-2431.55  l=-87.29
J=6    S=4    E=6    W=AND       v=3  a=-798.30   l=-69.71
J=7    S=4    E=7    W=IT        v=0  a=-791.79   l=-62.05
J=8    S=4    E=8    W=AND       v=2  a=-836.88   l=-69.71
J=9    S=3    E=9    W=BUT       v=0  a=-965.47   l=-51.14
J=10   S=4    E=10   W=A.        v=0  a=-783.36   l=-105.95
J=11   S=4    E=11   W=IN        v=0  a=-835.98   l=-49.01
J=12   S=4    E=12   W=A         v=0  a=-783.36   l=-59.66
J=13   S=4    E=13   W=AT        v=0  a=-923.59   l=-77.95
J=14   S=4    E=14   W=THE       v=0  a=-1326.40  l=-27.96
J=15   S=4    E=15   W=E.        v=0  a=-1321.67  l=-121.96
J=16   S=4    E=16   W=A         v=2  a=-1451.38  l=-59.66
J=17   S=4    E=17   W=THE       v=2  a=-1490.78  l=-27.96
J=18   S=4    E=18   W=IT        v=0  a=-1450.07  l=-62.05
J=19   S=5    E=18   W=IT        v=0  a=-1450.07  l=-110.42
J=20   S=6    E=18   W=IT        v=0  a=-775.76   l=-85.12
J=21   S=7    E=18   W=IT        v=0  a=-687.68   l=-125.32
J=22   S=8    E=18   W=IT        v=0  a=-687.68   l=-85.12
J=23   S=9    E=18   W=IT        v=0  a=-687.68   l=-50.28
J=24   S=10   E=18   W=IT        v=0  a=-689.67   l=-108.91
J=25   S=11   E=18   W=IT        v=0  a=-706.89   l=-113.78
J=26   S=12   E=18   W=IT        v=0  a=-689.67   l=-194.91
J=27   S=13   E=18   W=IT        v=0  a=-619.20   l=-100.24
J=28   S=4    E=19   W=IT        v=1  a=-1567.49  l=-62.05
J=29   S=14   E=20   W=DIDN'T    v=0  a=-4452.87  l=-195.48
J=30   S=15   E=20   W=DIDN'T    v=0  a=-4452.87  l=-118.62
J=31   S=16   E=20   W=DIDN'T    v=0  a=-4303.97  l=-189.88
J=32   S=17   E=20   W=DIDN'T    v=0  a=-4303.97  l=-195.48
J=33   S=18   E=20   W=DIDN'T    v=0  a=-4222.70  l=-78.74
J=34   S=19   E=20   W=DIDN'T    v=0  a=-4235.65  l=-78.74
J=35   S=20   E=21   W=ELABORATE v=2  a=-5847.54  l=-62.72
J=36   S=20   E=22   W=ELABORATE v=0  a=-5859.59  l=-62.72
J=37   S=21   E=23   W=!EXIT     v=0  a=-4651.00  l=-13.83
J=38   S=22   E=23   W=!EXIT     v=0  a=-4651.00  l=-13.83

```

Index

- ABORTONERR, 49
- accumulators, 8, 122
- accuracy figure, 39
- ACCWINDOW, 62
- adaptation, 13, 41
 - adaptation modes, 128
 - generating transforms, 42
 - global transforms, 42, 129
 - MAP, 41, 132
 - MLLR, 41, 129
 - MLLR formulae, 135
 - regression tree, 42, 105, 129, 147
 - supervised adaptation, 41, 128
 - transform model file, 42, 131
 - unsupervised adaptation, 41, 128, 181
- ADDITHER, 58
- ALIEN, 67
- all-pole filter, 58
- ALLOWCXTEXP, 167
- ALLOWWRDEXP, 39, 167
- analysis
 - FFT-based, 28
 - LPC-based, 28
- ANON, 56
- AT command, 33, 148
- AU command, 38, 39, 41, 146
- audio output, 68, 69
- audio source, 68
- AUDIOSIG, 69
- average log probability, 177

- back-off bigrams, 161
 - ARPA MIT-LL format, 162
- backward probability, 8
- Baum-Welch algorithm, 8
- Baum-Welch re-estimation, 6, 7
 - embedded unit, 120
 - isolated unit, 119
- Bayes' Rule, 3
- beam width, 122, 173
- <BeginHMM>, 93
- binary chop, 154
- binary storage, 106, 140
- binning, 60
- Boolean values, 48
- bootstrapping, 10, 17, 113
- byte swapping, 51, 58
- byte-order, 58
- BYTEORDER, 58

- C-heaps, 50

- CEPLIFTER, 59, 60
- cepstral analysis
 - filter bank, 60
 - liftering coefficient, 59
 - LPC based, 59
 - power vs magnitude, 60
- cepstral coefficients
 - liftering, 59
- cepstral mean normalisation, 61
- CFWORDBOUNDARY, 169
- CH command, 87
- check sums, 75
- CHKHMMDEFS, 93
- Choleski decomposition, 94
- CL command, 35, 140
- cloning, 34, 35, 140
- cluster merging, 38
- clustering
 - data-driven, 143
 - tracing in, 146
 - tree-based, 144
- CO command, 39, 144
- codebook, 104
- codebook exponent, 6
- codebooks, 6
- coding, 28
- command line
 - arguments, 29, 46
 - ellipsed arguments, 47
 - integer argument formats, 46
 - options, 15, 46
 - script files, 29
- compression, 75
- configuration files, 15, 47–48
 - default, 47
 - display, 49
 - format, 47
 - types, 48
- configuration parameters
 - operating environment, 52
 - switching, 123
- configuration variables, 15
 - display of, 48
 - summary, 241
- confusion matrix, 179
- context dependent models, 139
- continuous speech recognition, 113
- covariance matrix, 91
- cross-word network expansion, 170
- cross-word triphones, 140

- data insufficiency, 36
- data preparation, 16, 22
- DC command, 87, 167
- DE command, 86
- decision tree-based clustering, 145
- decision trees, 36
 - loading and storing, 146
- decoder, 172
 - alignment mode, 175
 - evaluation, 176
 - forced alignment, 180
 - live input, 182
 - N-best, 183
 - operation, 172
 - organisation, 174
 - output formatting, 181
 - output MLF, 177
 - progress reporting, 176
 - recognition mode, 175
 - rescoring mode, 176
 - results analysis, 177
 - trace output, 177
 - using adapted HMMs, 181
- decompression filter, 51
- defunct mixture components, 119
- defunct mixtures, 147
- deleted interpolation, 154
- deletion errors, 177
- delta coefficients, 61
- DELTAWINDOW, 61
- dictionaries, 155
- dictionary
 - construction, 24, 165
 - edit commands, 166
 - entry, 24
 - format, 24
 - formats, 165
 - output symbols, 165
- digit recogniser, 159
- direct audio input, 68
 - signal control
 - keypress, 69
 - silence detector
 - speech detector, 68
- DISCOUNT, 162
- DISCRETE, 71
- discrete data, 150
- discrete HMM
 - output probability scaling, 104
- discrete HMMs, 104, 149
- discrete probability, 91, 103
- DISCRETEHS, 99
- DP command, 148
- duration parameters, 91
- duration vector, 108
- EBNF, 19, 159
- edit commands
 - single letter, 86
- edit file, 86
- embedded re-estimation, 31
- embedded training, 10, 18, 113, 120
- <EndHMM>, 93
- energy suppression, 70
- ENORMALISE, 48, 61
- environment variables, 52
- error message
 - format, 245
- error number
 - structure of, 245
- error numbers
 - structure of, 49
- errors, 49
 - full listing, 245
- ESCALE, 61
- EX command, 81, 180
- extended Backus-Naur Form, 159
- extensions
 - mfc, 28
 - scp, 29
 - wav, 25
- Figure of Merit, 19, 179
- file formats
 - ALIEN, 67
 - Audio Interchange (AIFF), 66
 - Esignal, 64, 65
 - HTK, 62, 65
 - NIST, 65
 - NOHEAD, 67
 - OGI, 67
 - SCRIBE, 66
 - Sound Designer(SDES1), 66
 - Sun audio (SUNAU8), 66
 - TIMIT, 65
 - WAV, 67
- files
 - adding checksums, 75
 - compressing, 75
 - configuration, 47
 - copying, 74
 - listing contents, 72
 - network problems, 51
 - opening, 51
 - script, 46
 - VQ codebook, 72
- filters, 51
- fixed-variance, 118
- flat start, 17, 26, 30, 113, 118
- float values, 48
- FOM, 19, 179
- FORCECXP, 39, 167
- forced alignment, 13, 29, 180
- FORCELEFTBI, 167
- FORCEOUT, 177
- FORCERIGHTBI, 167
- forward probability, 7
- forward-backward

- embedded, 120
 - isolated unit, 119
- Forward-Backward algorithm, 7
- full rank covariance, 93
- Gaussian mixture, 6
- Gaussian pre-selection, 71
- GCONST value, 109
- generalised triphones, 144
- `global.ded`, 166
- global options, 107
- global options macro, 116
- global speech variance, 112
- grammar, 159
- grammar scale factor, 39
- grand variance, 142
- Hamming Window, 57
- HAUDIO, 15
- HAUDIO, 67
- HBUILD, 19, 161, 163, 187–188
- HCOMPV, 17, 30, 112, 118, 189–190
- HCONFIG, 47
- HCOPY, 16, 28, 47, 74, 151, 191–193
- HDICT, 14
- HDMAN, 19, 24, 166, 194–196
- HEADAPT, 13, 18, 42, 128, 181, 197–199
- HEADERSIZE, 67
- HEREST, 10, 18, 31, 113, 120, 139, 200–202
- HGRAF, 15
- HHED, 18, 32, 41, 114, 138, 203–211
- HIFREQ, 60
- HINIT, 7, 9, 17, 46, 112, 212–213
- HK command, 153
- HLABEL, 14
- HHED, 35
- HLED, 16, 27, 34, 86, 180, 214–216
- HLIST, 16, 72, 217
- HLM, 14, 161
- HLSTATS, 16, 161, 218–219
- HMATH, 14, 45
- HMEM, 14, 45, 50
- HMM
 - binary storage, 36
 - build philosophy, 18
 - cloning, 34
 - definition files, 30
 - definitions, 4, 90
 - editor, 18
 - instance of, 11
 - parameters, 91
 - triphones, 34
- HMM definition
 - stream weight, 95
 - basic form, 92
 - binary storage, 106
 - covariance matrix, 93
 - formal syntax, 107
 - global features, 93
 - global options, 107
 - global options macro, 94
 - macro types, 98
 - macros, 97
 - mean vector, 93
 - mixture components, 93
 - multiple data streams, 95
 - stream weight, 95
 - symbols in, 92
 - tied-mixture, 103
 - transition matrix, 93
- HMM lists, 88, 99, 100, 120
- HMM name, 93
- HMM refinement, 138
- HMM sets, 99
 - types, 99
- HMM tying, 100
- HMODEL, 14
- HNET, 13, 14
- HParm
 - SILENERGY, 68
 - SILGLCHCOUNT, 68
 - SILMARGIN, 68
 - SILSEQCOUNT, 68
 - SPCGLCHCOUNT, 68
 - SPCSEQCOUNT, 68
 - SPEECHTHRESH, 68
- HPARM, 15
- HPARSE, 19, 23, 159, 220–223
- HParse format, 159
 - compatibility mode, 161
 - in V1.5, 161
 - variables, 160
- HQUANT, 16, 224–225
- HREC, 13, 15, 174
- HREST, 9, 17, 112, 226–227
- HRESULTS, 19, 177, 228–230
- HSGEN, 19, 25, 164, 231
- HSHELL, 14, 45
- HSIGP, 14
- HSK*ind, 99
- HSLAB, 16, 25, 232–235
- HSMOOTH, 19, 154, 236–237
- HTRAIN, 15
- HUTIL, 15
- HEADAPT, 134
- HVITE, 10, 13, 18, 19, 33, 39, 174, 181, 238–240
- HVQ, 14
- HWAVE, 15
- HWAVEFILTER, 66
- insertion errors, 177
- integer values, 48
- <InvCovar>, 93
- isolated word training, 112
- item lists, 35, 141
 - indexing, 141
 - pattern matching, 141
- JO command, 153

- K-means clustering, 115
- label files, 79
 - ESPS format, 81
 - HTK format, 80
 - SCRIBE format, 81
 - TIMIT format, 81
- labels
 - changing, 87
 - context dependent, 88
 - context markers, 81
 - deleting, 86
 - editing, 86
 - external formats, 80
 - merging, 87
 - moving level, 88
 - multiple level, 80
 - replacing, 87
 - side-by-side, 79
 - sorting, 86
- language model scaling, 177
- language models
 - bigram, 161
- lattice
 - comment lines, 261
 - field names, 262
 - format, 19, 261
 - header, 261
 - language model scale factor, 184
 - link, 261
 - N-best, 12
 - node, 261
 - rescoring, 13
 - syntax, 262
- lattice generation, 183
- lattices, 261
- library modules, 14
- likelihood computation, 4
- linear prediction, 58
 - cepstra, 59
- LINEIN, 68
- LINEOUT, 68
- live input, 40
- <LLTCovar>, 94
- LOFREQ, 60
- log arithmetic, 9
- LPC, 59
- LPCEPSTRA, 59
- LPCORDER, 59
- LPREFC, 59
- LS command, 42, 146, 147
- LT command, 39, 41, 146
- M-heaps, 50
- macro definition, 97
- macro substitution, 98
- macros, 33, 97
 - special meanings, 98
 - types, 98
- marking word boundaries, 140
- master label files, 27, 79, 82, 117, 120
 - embedded label definitions, 82
 - examples, 84
 - multiple search paths, 82
 - pattern matching, 83
 - patterns, 27, 84
 - search, 83
 - sub-directory search, 83
 - syntax, 83
 - wildcards, 83
- master macro file, 102
- master macro files, 30
 - input/output, 139
- matrix dimensions, 93
- MAXCLUSTITER, 151
- maximum model limit, 177
- MAXTRYOPEN, 51
- ME command, 87
- <Mean>, 93
- mean vector, 91
- MEASURESIL, 182
- mel scale, 60
- MELSPEC, 60
- memory
 - allocators, 50
 - element sizes, 51
 - statistics, 50
- memory management, 50
- MFCC coefficients, 28, 93
- MICIN, 68
- minimum occupancy, 37
- MINMIX, 147, 153, 154
- <Mixture>, 93, 109
- mixture component, 91
- mixture incrementing, 146
- mixture splitting, 147
- mixture tying, 153
- mixture weight floor, 147
- ML command, 88
- MLF, 27, 79
- MMF, 30, 102
- model compaction, 39
- model training
 - clustering, 143
 - compacting, 144
 - context dependency, 139
 - embedded, 120
 - embedded subword formulae, 127
 - forward/backward formulae, 124
 - HMM editing, 139
 - in parallel, 122
 - initialisation, 114
 - isolated unit formulae, 126
 - mixture components, 115
 - pruning, 121
 - re-estimation formulae, 123
 - sub-word initialisation, 117
 - tying, 140
 - update control, 117

- Viterbi formulae, 124
- whole word, 116
- monitoring convergence, 117, 121
- monophone HMM
 - construction of, 29
- MP command, 167
- MT command, 148
- MU command, 147
- mu law encoded files, 66
- multiple alternative transcriptions, 183
- multiple hypotheses, 261
- multiple recognisers, 174
- multiple streams, 69
 - rules for, 69
- multiple-tokens, 13

- N-best, 13, 183
- N-grams, 161
- NATURALREADORDER, 51
- NATURALWRITEORDER, 51
- NC command, 143
- network type, 168
- networks, 155
 - in recognition, 156
 - word-internal, 39
- new features
 - in Version 2.1, 20
 - in Version 2.2, 19
- NIST, 19
- NIST format, 178
- NIST scoring software, 178
- NIST Sphere data format, 66
- non-emitting states, 10
- non-printing chars, 49
- NSAMPLES, 67
- NUMCEPS, 59, 60
- NUMCHANS, 48, 60
- <NumMixes>, 93, 105
- <NumStates>, 107

- observations
 - displaying structure of, 74
- operating system, 45
- outlier threshold, 37
- output filter, 51
- output lattice format, 183
- output probability
 - continuous case, 6, 91
 - discrete case, 91
- OUTSILWARN, 182
- over-short training segments, 119

- parameter estimation, 111
- parameter kind, 62
- parameter tie points, 98
- parameter tying, 35
- parameterisation, 28
- partial results, 177
- path, 10
 - as a token, 11
- partial, 9
- phone alignment, 33
- phone mapping, 33
- phone model initialisation, 113
- phone recognition, 170
- phones, 113
- PHONESOUT, 68
- phonetic questions, 145
- pipes, 45, 51
- PLAINHS, 100
- pre-emphasis, 57
- PREEMCOEF, 57
- prompt script
 - generation of, 25
- prototype definition, 16
- pruning, 18, 31, 173, 177
 - in tied mixtures, 153
- pruning errors, 122

- QS command, 37, 145
- qualifiers, 55, 62
 - _A, 61
 - _C, 63, 76
 - _D, 61
 - _E, 61
 - _K, 63, 64, 76
 - _N, 62, 70
 - _O, 61
 - _V, 71, 152
 - _V, 75
 - _Z, 61
 - codes, 63
 - ESIG field specifiers, 65
 - summary, 77

- RAWENERGY, 61
- RC command, 42, 147
- RE command, 86
- realignment, 33
- recogniser evaluation, 39
- recogniser performance, 178
- recognition
 - direct audio input, 40
 - errors, 178
 - hypothesis, 173
 - network, 172
 - output, 40
 - overall process, 157
 - results analysis, 40
 - statistics, 178
 - tools, 19
- recording speech, 25
- RECOUTPREFIX, 183
- RECOUTSUFFIX, 183
- reflection coefficients, 58
- regression formula, 61
- removing outliers, 144
- results analysis, 19
- RN, 148
- RN command, 42

- RO command, 37, 139, 144
- RP command, 167
- RT command, 148

- SAVEASVQ, 71, 152
- SAVEBINARy, 106, 140
- SAVECOMPRESSED, 75
- SAVETHCRC, 75
- script files, 46, 116
- search errors, 174
- segmental k-means, 17
- sentence generation, 165
- SH command, 144
- SHAREDHS, 99
- short pause, 32
- signals
 - for recording control, 182
- silence floor, 61
- silence model, 32, 33, 142
- SILFLOOR, 61
- simple differences, 62
- SIMPLEDIFFS, 62
- single-pass retraining, 123
- singleton clusters, 144
- SK command, 148
- SLF, 19, 23, 155, 157
 - arc probabilities, 159
 - format, 157
 - null nodes, 158
 - word network, 158
- SO command, 86
- software architecture, 14
- SOURCEFORMAT, 64, 65
- SOURCEKIND, 55, 182
- SOURCELABEL, 80, 88
- SOURCERATE, 56, 68
- SP command, 167
- speaker identifier, 179
- SPEAKEROUT, 68
- speech input, 54
 - automatic conversion, 55
 - bandpass filtering, 60
 - blocking, 56
 - byte order, 58
 - DC offset, 57
 - direct audio, 67
 - dynamic coefficients, 61
 - energy measures, 61
 - filter bank, 60
 - general mechanism, 54
 - Hamming window function, 57
 - monitoring, 72
 - pre-emphasis, 57
 - pre-processing, 57
 - summary of variables, 76
 - target kind, 55
- speech/silence detector, 182
- SS command, 148, 153
- ST command, 39, 146

- stacks, 50
- standard lattice format, 19, 23, 155, 157
 - definition, 261
- standard options, 48
 - A, 48
 - C, 47, 48
 - D, 48
 - F, 64, 65, 80
 - G, 80
 - I, 82
 - L, 83, 117
 - S, 46, 48, 83
 - T, 49, 117
 - V, 49
 - summary, 53
- state clustering, 36
- state transitions
 - adding/removing, 148
- state tying, 36, 144
- statistics
 - state occupation, 36
- statistics file, 37, 139, 144
- <Stream>, 95, 103
- stream weight, 6, 91
- <StreamInfo>, 94, 108
- streams, 6
- stress marking, 24
- string matching, 177
- string values, 48
- strings
 - metacharacters in, 49
 - output of, 50
 - rules for, 49
- SU command, 148
- sub-lattices, 163, 261
- SUBLAT, 164, 261
- SW command, 148
- <SWeights>, 95

- TARGETKIND, 48, 55, 152
- TARGETRATE, 56
- task grammar, 22, 161
- TB command, 38, 145
- TC command, 88, 140, 143
- tee-models, 32, 105
 - in networks, 161
- termination, 49
- TI command, 33, 142
- tied parameters, 140
- tied-mixture system, 103
- tied-mixtures, 142, 152
 - build procedure, 153
 - output distribution, 103
- tied-state, 99
- TIMIT database, 23, 86
- token history, 173
- token passing, 173
- Token Passing Model, 11
- total likelihood, 8, 9

- TR command, 37
- tracing, 49
- training
 - sub-word, 113
 - whole-word, 112
- training tools, 16
- TRANSALT, 80, 88
- transcription
 - orthographic, 26
- transcriptions
 - model level, 181
 - phone level, 181
 - word level, 181
- transitions
 - adding them, 32
- TRANSLEV, 80, 88
- <TransP>, 93
- tree building, 38
- tree optimisation, 146
- triphones
 - by cloning, 34
 - from monophones, 34
 - notation, 35
 - synthesising unseen, 41
 - word internal, 34
- tying
 - examples of, 142
 - exemplar selection, 142
 - states, 36
 - transition matrices, 35

- UF command, 139
- under-training, 154
- uniform segmentation, 115
- unseen triphones, 39, 144
 - synthesising, 146
- up-mixing, 146
- upper triangular form, 93
- USEHAMMING, 57
- USEPOWER, 60
- USESILDET, 68, 182
- UT command, 142

- V1COMPAT, 62, 161
- varFloorN, 117
- variance
 - flooring problems, 36
- <Variance>, 93
- variance floor macros, 30
 - generating, 118
- variance floors, 117
- <VecSize>, 93, 108
- vector dimensions, 93
- vector quantisation, 70
 - code book external format, 72
 - distance metrics, 71
 - type of, 72
 - uses of, 71
- Viterbi training, 7, 115
- VQTABLE, 71

- warning codes
 - full listing, 245
- warning message
 - format, 245
- warnings, 49
- waveform capture, 182
- WB command, 34, 140
- whole word modelling, 112
- whole word recognition, 171
- WINDOWSIZE, 48, 56
- WLR, 12
- word equivalence, 179
- word insertion penalty, 39
- word internal, 34
- Word Link Record, 12
- word list, 24
- word N-best, 13
- expansion rules, 167
- word networks
 - tee-models in, 161
- word spotting, 171
- word-end nodes, 11, 12, 173
- word-internal network expansion, 170
- word-loop network, 164
- word-pair grammar, 164

- ZMEANSOURCE, 57