

# Sentence Fusion for Multidocument News Summarization

Regina Barzilay\*  
Massachusetts Institute of Technology

Kathleen R. McKeown†  
Columbia University

*A system that can produce informative summaries, highlighting common information found in many online documents, will help web users to pinpoint information that they need without extensive reading. In this paper, we introduce sentence fusion, a novel text-to-text generation technique for synthesizing common information across documents. Sentence fusion involves bottom-up local multisequence alignment to identify phrases conveying similar information; and statistical generation to combine common phrases into a sentence. Sentence fusion moves the summarization field from the use of purely extractive methods to the generation of abstracts, which contain sentences not found in any of the input documents and which can synthesize information across sources.*

## 1. Introduction

Redundancy in large text collections, such as the web, creates both problems and opportunities for natural language systems. On the one hand, the presence of numerous sources conveying the same information causes difficulties for end users of search engines and news providers; they must read the same information over and over again. On the other hand, redundancy can be exploited to identify important and accurate information for applications such as summarization and question answering (Mani and Bloedorn, 1997; Radev and McKeown, 1998; Radev, Prager, and Samn, 2000; Clarke, Cormack, and Lynam, 2001; Dumais et al., 2002; Chu-Carroll et al., 2003). Clearly, it would be highly desirable to have a mechanism that could identify common information among multiple related documents and fuse it into a coherent text. In this paper, we present a method for sentence fusion that exploits redundancy to achieve this task in the context of multidocument summarization.

A straightforward approach for approximating sentence fusion can be found in the use of sentence extraction for multidocument summarization (Carbonell and Goldstein, 1998; Radev, Jing, and Budzikowska, 2000; Marcu and Gerber, 2001; Lin and Hovy, 2002). Once a system finds a set of sentences that convey similar information (e.g., by clustering), one of these sentences is selected to represent the set. This is a robust approach that is always guaranteed to output a grammatical sentence. However, extraction is only a coarse approximation of fusion. An extracted sentence may include not only common information, but additional information specific to the article from which it came, leading to source bias and aggravating fluency problems in the extracted summary. Attempting to solve this problem by including more sentences might lead to a verbose and repetitive summary.

---

\* Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA 02458

† Department of Computer Science, Columbia University, New York, NY 10027

Instead, we want a fine-grained approach that can identify only those pieces of sentences that are common. Language generation offers an appealing approach to the problem, but the use of generation in this context raises significant research challenges. In particular, generation for sentence fusion must be able to operate in a domain-independent fashion, scalable to handle a large variety of input documents with various degrees of overlap. In the past, generation systems were developed for limited domains and required a rich semantic representation as input. In contrast, for this task we require *text-to-text generation*, the ability to produce a new text given a set of related texts as input. If language generation can be scaled to take fully formed text as input without semantic interpretation, selecting content and producing well-formed English sentences as output, then generation has a large potential payoff.

In this paper, we present the concept of sentence fusion, a novel text-to-text generation technique which, given a set of similar sentences, produces a new sentence containing the information common to most sentences in the set. The research challenges in developing such an algorithm lie in two areas: identification of the fragments conveying common information and combination of the fragments into a sentence. To identify common information, we have developed a method for aligning syntactic trees of input sentences, incorporating paraphrasing information. Our alignment problem poses unique challenges: we only want to match a subset of the subtrees in each sentence and are given few constraints on permissible alignments (e.g., arising from constituent ordering, start or end points). Our algorithm meets these challenges through bottom-up local multisequence alignment, using words and paraphrases as anchors. Combination of fragments is addressed through construction of a fusion tree encompassing the resulting alignment and linearization of the tree into a sentence using a language model. Our approach to sentence fusion thus features the integration of robust statistical techniques, such as local, multi-sequence alignment and language modeling, with linguistic representations automatically derived from input documents.

Sentence fusion is a significant first step towards the generation of abstracts, as opposed to extracts (Borko and Bernier, 1975), for multi-document summarization. While there has been research on sentence reduction for single document summarization (Grefenstette, 1998; Mani, Gates, and Bloedorn, 1999; Knight and Marcu, 2001; Jing and McKeown, 2000; Reizler et al., 2003), analysis of human-written multi-document summaries shows that most sentences contain information drawn from multiple documents (Banko and Vanderwende, 2004). Unlike extraction methods (used by the vast majority of summarization researchers), sentence fusion allows for the true synthesis of information from a set of input documents. It generates summary sentences by reusing and altering phrases from input sentences, combining common information from several sources. Consequently, one summary sentence may include information conveyed in several input sentences. Our evaluation shows that our approach is promising, with sentence fusion outperforming sentence extraction for the task of content selection.

This paper describes the implementation of the sentence fusion method within the multidocument summarization system MultiGen, which daily summarizes multiple news articles on the same event as part<sup>1</sup> of Columbia's news browsing system Newsblaster.<sup>2</sup> Analysis of the system's output reveals the capabilities and the weaknesses of our text-to-text generation method and identifies interesting challenges that will require new insights.

In the next section, we provide an overview of the multidocument summarization

---

<sup>1</sup> In addition to MultiGen, Newsblaster utilizes another summarizer DEMS (Schiffman, Nenkova, and McKeown, 2002) to summarize heterogeneous sets of articles.

<sup>2</sup> <http://newsblaster.cs.columbia.edu/>.

system MultiGen focusing on components that produce input or operate over output of sentence fusion. In Section 3, we provide an overview of our fusion algorithm and detail on its main steps: identification of common information (Section 3.1), fusion lattice computation (Section 3.2), and lattice linearization (Section 3.3). Evaluation results and their analysis are presented in Section 4. An overview of related work and discussion of future directions conclude the paper.

## 2. Framework for Sentence Fusion: MultiGen

Sentence fusion is the central technique used within the MultiGen summarization system. MultiGen takes as input a cluster of news stories on the same event and produces a summary which synthesizes common information across input stories. An example of a MultiGen summary is shown in Figure 1. The input clusters are automatically produced from a large quantity of news articles that are retrieved by Newsblaster from 30 news sites each day.

### Agency Suspends Smallpox Vaccines for People With Heart Disease

**Summary from the U.S.**


A second health care worker has died of a heart attack (3) after receiving a smallpox vaccination (9) and officials are investigating whether vaccinations are to blame (3) for cardiac problems. (6) The vaccine never has been associated with heart trouble but as a precaution (3) the U.s. centers for Disease Control and Prevention (14) is advising people with a history of heart disease to be vaccinated (3) until further notice. (14) Strom suggested that the Bush administration reassess whether it necessary and safe to continue with its aggressive plan to inoculate millions of health care workers and emergency responders. (1)

**Story keywords**

vaccine, Heart, Smallpox, vaccinated, Disease

**Source articles**

1. [Vaccination program in peril after second death](#) (seattletimes.nwsource.com, 03/28/2003, 319 words)
2. [Wired News: Smallpox Shots: Proceed With Care](#) (Wired, 03/27/2003, 559 words)
3. [2nd worker dies after smallpox vaccination](#) (suntimes.com, 03/28/2003, 358 words)
4. [2nd worker dies after smallpox vaccine](#) (dallasnews.com, 03/28/2003, 499 words)
5. [Smallpox vaccine is reviewed after second fatal heart attack](#) (boston.com, 03/28/2003, 732 words)
6. [Second Smallpox Vaccine Death Evad](#) (CRS News 03/28/2003, 865 words)



**Figure 1**

An example of MultiGen summary as shown in the Columbia Newsblaster Interface. Summary phrases are followed by parenthesized numbers indicating their source articles. The last sentence is extracted since it was repeated verbatim in several input articles.

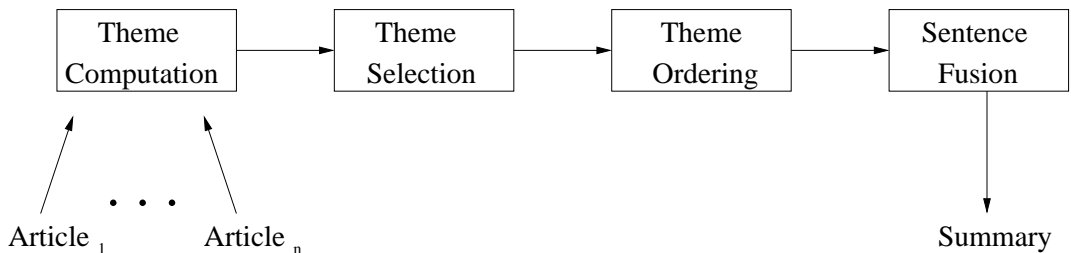
In order to understand the role of sentence fusion within summarization, we overview the MultiGen architecture, providing details on the processes that precede sentence fusion and thus, the input that the fusion component requires. Fusion itself is discussed in the following sections of the paper.

MultiGen follows a pipeline architecture, shown in Figure 2. The analysis component of the system, Simfinder (Hatzivassiloglou, Klavans, and Eskin, 1999) clusters sentences of input documents into *themes*, groups of sentences that convey similar information

(Section 2.1). Once themes are constructed, the system selects a subset of the groups to be included in the summary, depending on the desired compression length (Section 2.2). The selected groups are passed to the ordering component, which selects a complete order among themes (Section 2.3).

**2.1 Theme Construction**

The analysis component of MultiGen, Simfinder, identifies *themes*, groups of sentences from different documents that each say roughly the same thing. Each theme will ultimately correspond to at most one sentence in the output summary, generated by the fusion component, and there may be many themes for a set of articles. An example of a theme is shown in Table 1. As this set illustrates, sentences within a theme are not exact repetitions of each other; they usually include phrases expressing information that is *not* common to all sentences in the theme. Information that is common across sentences is shown in bold; other portions of the sentence are specific to individual articles. If one of these sentences were used as is to represent the theme, the summary would contain extraneous information. Also, errors in clustering may result in the inclusion of some unrelated sentences. Evaluation involving human judges revealed that Simfinder identifies similar sentences with 49.3% precision at 52.9% recall (Hatzivassiloglou, Klavans, and Eskin, 1999). We will discuss later how this error rate influences sentence fusion.



**Figure 2**  
MultiGen Architecture

1. IDF Spokeswoman did not confirm this, but said <b>the Palestinians fired an anti-tank missile at a bulldozer</b> .
2. The clash erupted when <b>Palestinian militants fired machine-guns and anti-tank missiles at a bulldozer</b> that was building an embankment in the area to better protect Israeli forces.
3. The army expressed "regret at the loss of innocent lives" but a senior commander said troops had shot in self-defense <b>after being fired at while using bulldozers</b> to build a new embankment at an army base in the area.
<b>fusion sentence:</b> Palestinians fired an anti-tank missile at a bulldozer.

**Table 1**  
A theme with the corresponding fusion sentence.

To identify themes, Simfinder extracts linguistically motivated features for each sentence, including WordNet synsets (Miller et al., 1990) and syntactic dependencies, such as subject-verb and verb-object relations. A log-linear regression model is used to combine the evidence from the various features to a single similarity value. The model was trained on a large set of sentences which were manually marked for similarity. The output of the

model is a listing of real-valued similarity values on sentence pairs. These similarity values are fed into a clustering algorithm that partitions the sentences into closely related groups.

## 2.2 Theme Selection

To generate a summary of predetermined length, we induce a ranking on the themes and select the  $n$  highest.<sup>3</sup> This ranking is based on three features of the theme: size measured as the number of sentences, similarity of sentences in a theme, and salience score. The first two of these scores are produced by Simfinder, and the salience score of the theme is computed using *lexical chains* (Morris and Hirst, 1991; Barzilay and Elhadad, 1997) as described below. Since each of these scores has a different range of values, we perform ranking based on each score separately, and then, induce total ranking by summing ranks from individual categories:

$$\begin{aligned} \text{Rank (theme)} = & \text{Rank (Number of sentences in theme)} + \\ & \text{Rank (Similarity of sentences in theme)} + \\ & \text{Rank (Sum of lexical chain scores in theme)} \end{aligned}$$

Lexical chains — sequences of semantically related words — are tightly connected to the lexical cohesive structure of the text and have been shown to be useful for determining which sentences are important for single document summarization (Barzilay and Elhadad, 1997; Silber and McCoy, 2002). In the multidocument scenario, lexical chains can be adapted for theme ranking based on the salience of theme sentences within their original documents. Specifically, a theme that has many sentences ranked high by lexical chains as important for a single document summary, is, in turn, given a higher salience score for the multidocument summary. In our implementation, a salience score for a theme is computed as the sum of lexical chain scores of each sentence in a theme.

## 2.3 Theme Ordering

Once we filter out the themes that have a low rank, the next task is to order the selected themes into coherent text. Our ordering strategy aims to capture chronological order of the main events and ensure coherence. To implement this strategy in MultiGen, we select for each theme the sentence which has the earliest publication time (*theme time stamp*). To increase the coherence of the output text, we identify blocks of topically-related themes and then apply chronological ordering on blocks of themes using theme time stamps (Barzilay, Elhadad, and McKeown, 2002). These stages thus produce a sorted set of themes which are passed as input to the sentence fusion component, described in the next section.

## 3. Sentence Fusion

Given a group of similar sentences—a theme—the problem is to create a concise and fluent fusion of information with this theme, reflecting facts common to all sentences. An example of a fusion sentence is shown in Table 1. To achieve this goal we need to identify phrases common to most theme sentences, and then combine them into a new sentence.

At one extreme, we might consider a shallow approach to the fusion problem, adapting the “bag of words” approach. However, sentence intersection in a set-theoretic sense

---

<sup>3</sup> Typically, Simfinder produces at least 20 themes given an average Newsblaster cluster of nine articles. The length of a generated summary typically does not exceed seven sentences.

produces poor results. For example, the intersection of the first two sentences from the theme shown in Table 1 is (*the, fired, anti-tank, at, a, bulldozer*). Besides being ungrammatical, it is impossible to understand what event this intersection describes. The inadequacy of the “bag of words” method to the fusion task demonstrates the need for a more linguistically motivated approach. At the other extreme, previous approaches (Radev and McKeown, 1998) have demonstrated that this task is feasible when a detailed semantic representation of the input sentences is available. However, these approaches operate in a limited domain (e.g., terrorist events), where information extraction systems can be used to interpret the source text. The task of mapping input text into a semantic representation in a domain-independent setting extends well beyond the ability of current analysis methods. These considerations suggest that we need a new method for the sentence fusion task. Ideally, such a method would not require a full semantic representation. Rather, it would rely on input texts and shallow linguistic knowledge (such as parse trees) that can be automatically derived from a corpus to generate a fusion sentence.

In our approach, sentence fusion is modeled after the typical generation pipeline: content selection (what to say) and surface realization (how to say it). In contrast to traditional generation systems where a content selection component chooses content from semantic units, our task is complicated by the lack of semantics in the textual input. At the same time, we can benefit from the textual information given in the input sentences for the tasks of syntactic realization, phrasing, and ordering; in many cases, constraints on text realization are already present in the input.

The algorithm operates in three phases:

- *Identification of common information* (Section 3.1)
- *Fusion lattice computation* (Section 3.2)
- *Lattice linearization* (Section 3.3)

Content selection occurs primarily in the first phase, in which our algorithm uses local alignment across pairs of parsed sentences, from which we select fragments to be included in the fusion sentence. Instead of examining all possible ways to combine these fragments, we select a sentence in the input which contains most of the fragments and transform its parsed tree into the fusion lattice by eliminating non-essential information and augmenting it with information from other input sentences. This construction of the fusion lattice targets content selection but, in the process, alternative verbalizations are selected and thus, some aspects of realization are also carried out in this phase. Finally, we generate a sentence from this representation based on a language model derived from a large body of texts. This approach generates a fusion sentence by reusing and altering phrases from the input sentences, performing text-to-text generation.

### 3.1 Identification of Common Information

Our task is to identify information shared between sentences. We do this by aligning constituents in the syntactic parse trees for the input sentences. Our alignment process differs considerably from alignment for other NL tasks, such as machine translation, because we cannot expect a complete alignment. Rather, a subset of the subtrees in one sentence will match different subsets of the subtrees in the others. Furthermore, order across trees is not preserved; there is no natural starting point for alignment; and there are no constraints on crosses. For these reasons we have developed a bottom-up local multisequence alignment algorithm that uses words and phrases as anchors for matching. This algorithm operates on the dependency trees for pairs of input sentences. We use a dependency-based representation because it abstracts over features irrelevant

for comparison such as constituent ordering. In the paragraphs that follow, we first describe how this representation is computed, then describe how dependency subtrees are aligned, and finally describe how we choose between constituents conveying overlapping information.

In this section we first describe an algorithm which, given a pair of sentences, determines which sentence constituents convey information appearing in both sentences. This algorithm will be applied to pairwise combinations of sentences in the input set of related sentences.

The intuition behind the algorithm is to compare all constituents of one sentence to those of another, and to select the most similar ones. Of course, how this comparison is done depends on the particular sentence representation used. A good sentence representation would emphasize sentence features that are relevant for comparison, such as dependencies between sentence constituents, while ignoring irrelevant features, such as constituent ordering. A representation which fits these requirements is a dependency-based representation (Melcuk, 1988). We first detail how this representation is computed, then describe a method for aligning dependency subtrees, and finally present a method for selecting components conveying overlapping information.

**3.1.1 Sentence Representation.** In many NLP applications, the structure of a sentence is represented using phrase structure trees. An alternative representation is a *dependency tree*, which describes the sentence structure in terms of dependencies between words. The similarity of the dependency tree to a predicate-argument structure makes it a natural representation for our comparison.<sup>4</sup> This representation can be constructed from the output of a traditional parser. In fact, we have developed a rule-based component that transforms the phrase-structure output of Collins’ parser (Collins, 1997) into a representation where a node has a direct link to its dependents. We also mark verb-subject and verb-node dependencies in the tree.

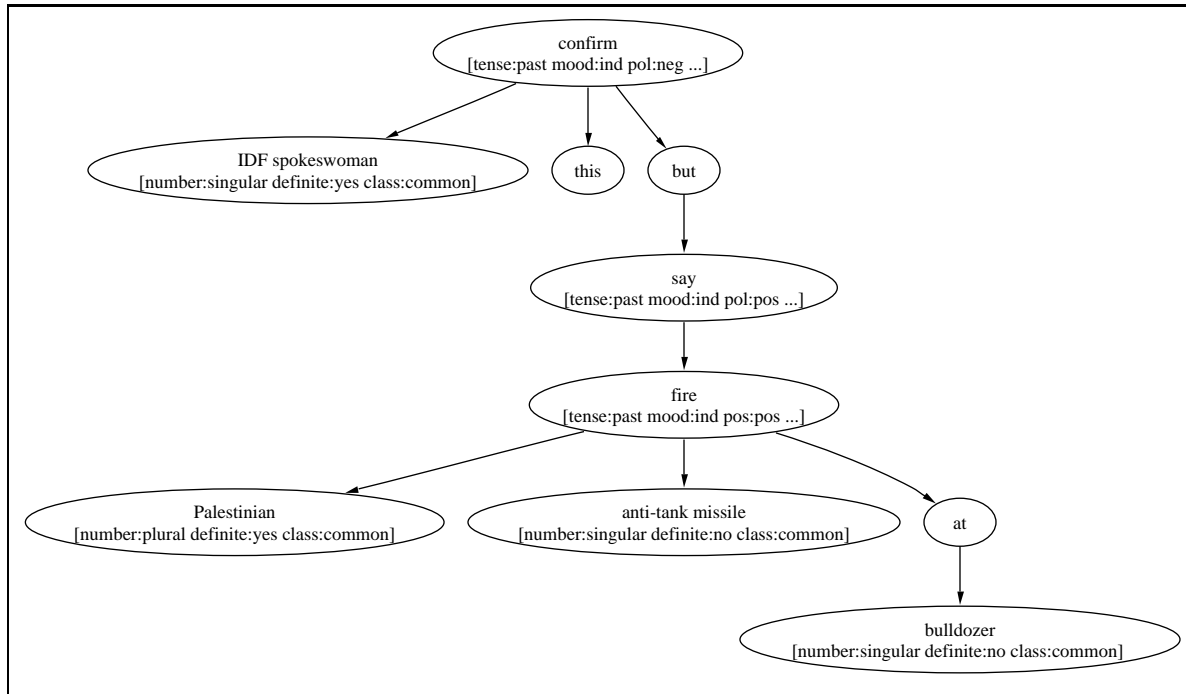
The process of comparing trees can be further facilitated if the dependency tree is abstracted to a canonical form which eliminates features irrelevant to the comparison. We hypothesize that the difference in grammatical features such as auxiliaries, number, and tense have a secondary effect when comparing the meaning of sentences. Therefore, we represent in the dependency tree only non-auxiliary words with their associated grammatical features. For nouns, we record their number, articles, and class (common or proper). For verbs, we record tense, mood (indicative, conditional or infinitive), voice, polarity, aspect (simple or continuous), and taxis (perfect or none). The eliminated auxiliary words can be recreated using these recorded features. We also transform all the passive voice sentences to the active voice, changing the order of affected children.

While the alignment algorithm described in Section 3.1.2 produces one-to-one mappings, in practice some paraphrases are not decomposable to words, forming one-to-many or many-to-many paraphrases. Our manual analysis of paraphrased sentences (Barzilay, 2003) revealed that such alignments most frequently occur in pairs of noun phrases (e.g., “*faculty member*” and “*professor*”) and pairs including verbs with particles (e.g., “*stand up*”, “*rise*”). To correctly align such phrases, we flatten subtrees containing noun phrases and verbs with particles into one node. We subsequently determine matches between flattened sentences using statistical metrics.

An example of a sentence and its dependency tree with associated features is shown in

---

<sup>4</sup> Two paraphrasing sentences which differ in word order may have significantly different trees in phrase-based format. For instance, this phenomenon occurs when an adverbial is moved from a position in the middle of a sentence to the beginning of a sentence. In contrast, dependency representations of these sentences are very similar.

**Figure 3**

Dependency tree of the sentence “*The IDF spokeswoman did not confirm this, but said the Palestinians fired an anti-tank missile at a bulldozer on the site.*” The features of the node “confirm” are explicitly marked in the graph.

Figure 3. (In figures of dependency trees hereafter, node features are omitted for clarity.)

**3.1.2 Alignment.** Our alignment of dependency trees is driven by two sources of information: the similarity between the structure of the dependency trees, and the similarity between two given words. In determining the structural similarity between two trees, we take into account the types of edges (which indicate the relationships between nodes). An edge is labeled by the syntactic function of the two nodes it connects (e.g., subject-verb). It is unlikely that an edge connecting a subject and verb in one sentence, for example, corresponds to an edge connecting a verb and an adjective in another sentence.

The word similarity measures take into account more than word identity: they also identify pairs of paraphrases, using WordNet and a paraphrasing dictionary. We automatically constructed the paraphrasing dictionary from a large comparable news corpus using the co-training method described in (Barzilay and McKeown, 2001). The dictionary contains pairs of word-level paraphrases as well as phrase-level paraphrases.<sup>5</sup> Several examples of automatically extracted paraphrases are given in Table 2. During alignment, each pair of non-identical words that do not comprise a synset in WordNet is looked up

<sup>5</sup> The comparable corpus and the derived dictionary are available at <http://www.cs.cornell.edu/~regina/thesis-data/comp/input/processed.tbz2> and <http://www.cs.cornell.edu/~regina/thesis-data/comp/output/comp2-ALL.txt>. For details on the corpus collection and evaluation of the paraphrase quality see (Barzilay, 2003).



in the paraphrasing dictionary; in the case of a match, the pair is considered to be a paraphrase.

(auto, automobile), (closing, settling), (rejected, does not accept), (military, army), (IWC, International Whaling Commission), (Japan, country), (researching, examining), (harvesting, killing), (mission-control office, control centers), (father, pastor), (past 50 years, four decades), (Wangler, Wanger), (teacher, pastor), (fondling, groping), (Kalkilya, Qalqilya), (accused, suspected), (language, terms), (head, president), (U.N., United Nations), (Islamabad, Kabul), (goes, travels), (said, testified), (article, report), (chaos, upheaval), (Gore, Lieberman), (revolt, uprising), (more restrictive local measures, stronger local regulations) (countries, nations), (barred, suspended), (alert, warning), (declined, refused), (anthrax, infection), (expelled, removed), (White House, White House spokesman Ari Fleischer), (gunmen, militants)

**Table 2**

Lexical paraphrases extracted by the algorithm from the comparable news corpus.

We now give an intuitive explanation of how our tree similarity function, denoted by  $Sim$ , is computed. If the optimal alignment of two trees is known, then the value of the similarity function is the sum of the similarity scores of aligned nodes and aligned edges. Since the best alignment of given trees is not known *a priori*, we select the maximal score among plausible alignments of the trees. Instead of exhaustively traversing the space of all possible alignments, we recursively construct the best alignment for trees of given depths, assuming that we know how to find an optimal alignment for trees of shorter depths. More specifically, at each point of the traversal we consider two cases, shown in Figure 4. In the first case, two top nodes are aligned to each other and their children are aligned in an optimal way by applying the algorithm to shorter trees. In the second case, one tree is aligned with one of the children of the top node of the other tree; again we can apply our algorithm for this computation, since we decrease the height of one of the trees.

Before giving the precise definition of  $Sim$ , we introduce some notation. When  $T$  is a tree with root node  $v$ , we let  $c(T)$  denote the set containing all children of  $v$ . For a tree  $T$  containing a node  $s$ , the subtree of  $T$  which has  $s$  as its root node is denoted by  $T_s$ .

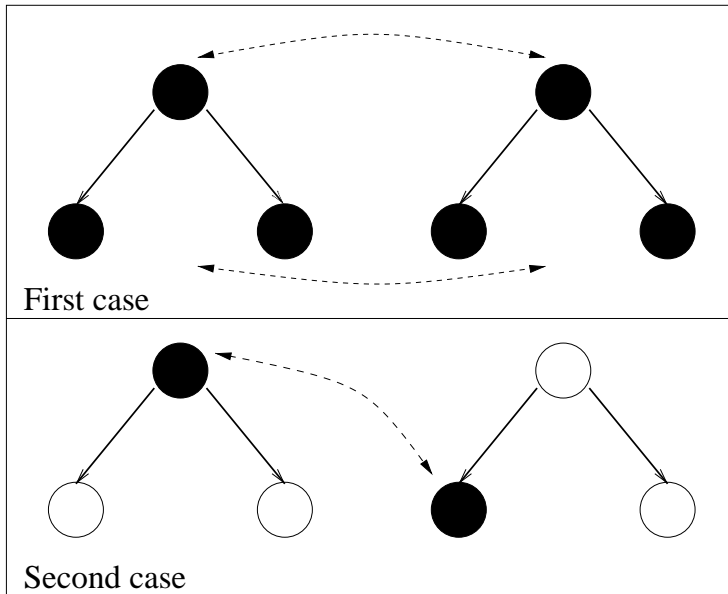
Given two trees  $T$  and  $T'$  with root nodes  $v$  and  $v'$ , respectively, the similarity  $Sim(T, T')$  between the trees is defined to be the maximum of the three expressions  $NodeCompare(T, T')$ ,  $\max\{Sim(T_s, T') : s \in c(T)\}$ , and  $\max\{Sim(T, T_{s'}) : s' \in c(T')\}$ . The upper part of Figure 4 depicts the computation of  $NodeCompare(T, T')$ , where two top nodes are aligned to each other. The remaining expressions,  $\max\{Sim(T_s, T') : s \in c(T)\}$ , and  $\max\{Sim(T, T_{s'}) : s' \in c(T')\}$ , capture mappings in which the top of one tree is aligned with one of the children of the top node of the other tree (the bottom of the Figure 4).

The maximization in the  $NodeCompare$  formula searches for the best possible alignment for the child nodes of the given pair of nodes and is defined by

$$NodeCompare(T, T') =$$

$$NodeSim(v, v') + \max_{m \in M(c(T), c(T'))} \left[ \sum_{(s, s') \in m} (EdgeSim((v, s), (v', s')) + Sim(T_s, T_{s'})) \right]$$

where  $M(A, A')$  is the set of all possible matchings between  $A$  and  $A'$ , and a matching (between  $A$  and  $A'$ ) is a subset  $m$  of  $A \times A'$  such that for any two distinct elements

**Figure 4**

Tree alignment computation. In the first case two tops are aligned, while in the second case the top of one tree is aligned to a child of another tree.

$(a, a'), (b, b') \in m$ , both  $a \neq b$  and  $a' \neq b'$ . In the base case, when one of the trees has depth one,  $NodeCompare(T, T')$  is defined to be  $NodeSim(v, v')$ .

The similarity score  $NodeSim(v, v')$  of atomic nodes depends on whether the corresponding words are identical, paraphrases or unrelated. The similarity scores for pairs of identical words, pairs of synonyms, pairs of paraphrases or edges (given in Table 3) are manually derived using a small development corpus. While learning of the similarity scores automatically is an appealing alternative, its application in the fusion context is challenging due to the absence of a large training corpus and the lack of an automatic evaluation function.<sup>6</sup> The similarity of nodes containing flattened subtrees,<sup>7</sup> such as noun phrases, is computed as the score of their intersection normalized by the length of the longest phrase. For instance, the similarity score of the noun phrases “*anti-tank missile*” and “*machine gun and anti-tank missile*” is computed as a ratio between the score of their intersection “*anti-tank missile*” (2), divided by the length of the latter phrase (4).

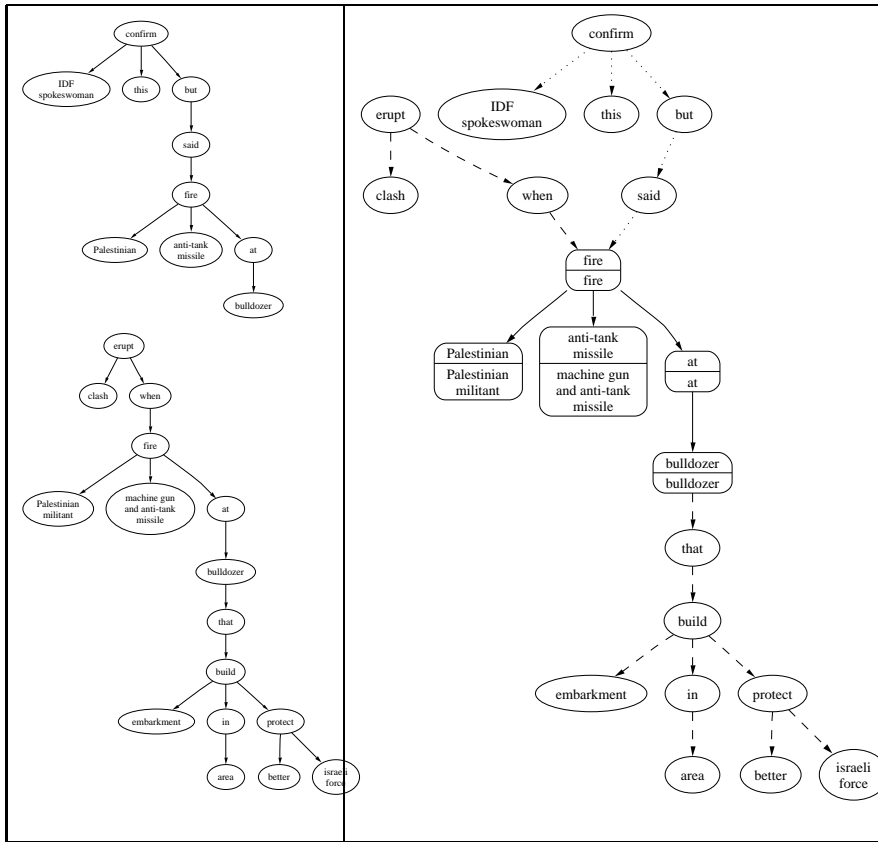
The computation of the similarity function  $Sim$  is performed using bottom-up dynamic programming, where the shortest subtrees are processed first. The alignment algorithm returns the similarity score of the trees as well as the optimal mapping between the subtrees of input trees. In the resulting tree mapping, the pairs of nodes whose  $NodeSim$

<sup>6</sup> Our preliminary experiments with n-gram-based overlap measures, such as BLEU (Papineni et al., 2002) and ROUGE (Lin and Hovy, 2003), show that these metrics do not correlate with human judgments on the fusion task, when tested against two reference outputs. This is to be expected: as lexical variability across input sentences grows, the number of possible ways to fuse them by machine as well by human also grows. The accuracy of match between the system output and the reference sentences is largely depends on the features of the input sentences, rather than on the underlying fusion method.

<sup>7</sup> Pairs of phrases that form an entry in the paraphrasing dictionary are compared as pairs of atomic entries.

Category	Node Sim	Category	NodeSim
Identical words	1	Edges are subject-verb	0.03
Synonyms	1	Edges are verb-object	0.03
Paraphrases	0.5	Edges are same type	0.02
Other	-0.1	Other	0

**Table 3**  
Node and edge similarity scores used by the alignment algorithm.



**Figure 5** Two dependency trees and their alignment tree. Solid lines represent aligned edges. Dotted and dashed edges represent unaligned edges of the theme sentences.

positively contributed to the alignment are considered as parallel.

Figure 5 shows two dependency trees and their alignment.

As it is evident from the *Sim* definition, we are only considering one-to-one node “matchings”: every node in one tree is mapped to at most one node in another tree. This restriction is necessary, since the problem of optimizing many-to-many alignment is NP-hard.<sup>8</sup> The subtree flattening, performed during the preprocessing stage, aims to minimize the negative effect of the restriction on alignment granularity.

Another important property of our algorithm is that it produces a local alignment. Local alignment maps local regions with high similarity to each other rather than creating an overall optimal global alignment of the entire tree. This strategy is more meaningful when only partial meaning overlap is expected between input sentences, as in typical sen-

<sup>8</sup> The complexity of our algorithm is polynomial in the number of nodes.

Let  $n_1$  denote the number of nodes in the first tree, and  $n_2$  denote the number of nodes in the second tree. We assume that the branching factor of a parse tree is bounded above by a constant. The function *NodeCompare* is evaluated only once on each node pair. Therefore, it is evaluated  $n_1 * n_2$  times totally. Each evaluation is computed in constant time, assuming that values of the function for node children are known. Since we use memoization, the total time of the procedure is  $O(n_1 * n_2)$ .

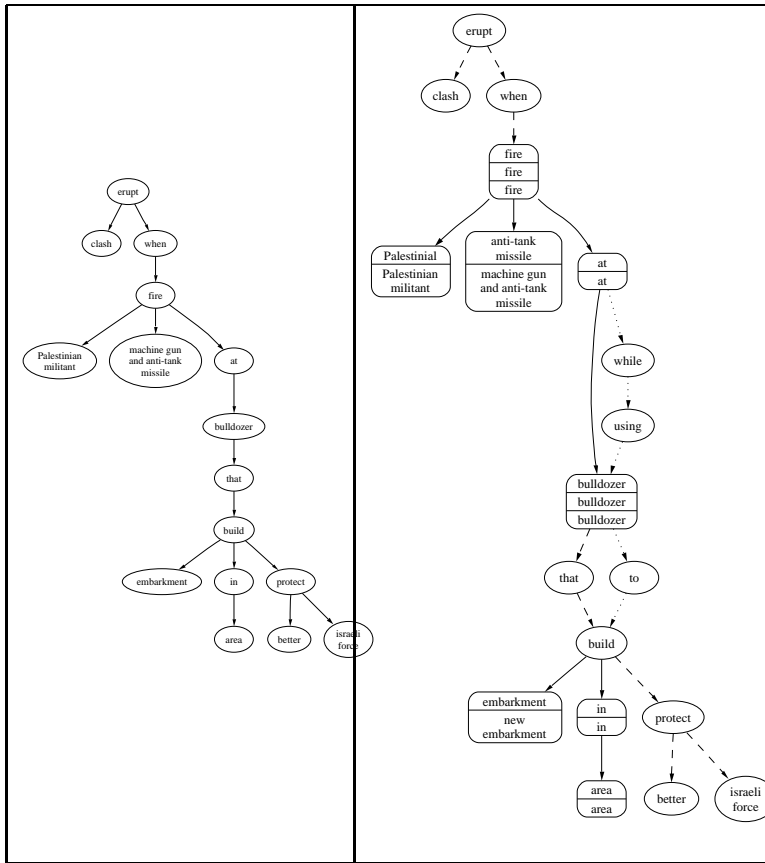
tence fusion input. Only these high similarity regions, which we call *intersection subtrees*, are included in the fusion sentence.

### 3.2 Fusion Lattice Computation

The next question we address is how to put together intersection subtrees. During this process, the system will remove phrases from a selected sentence, add phrases from other sentences, and replace words with the paraphrases that annotate each node. Obviously, among the many possible combinations, we are interested only in those combinations which yield semantically sound sentences and do not distort the information presented in the input sentences. We cannot explore every possible combination, since the lack of semantic information in the trees prohibits us from assessing the quality of the resulting sentences. Instead, we select a combination already present in the input sentences as a basis, and transform it into a fusion sentence by removing extraneous information and augmenting the fusion sentence with information from other sentences. The advantage of this strategy is that, when the initial sentence is semantically correct and the applied transformations aim to preserve semantic correctness, the resulting sentence is a semantically correct one. In fact, early experimentation with generation from constituent phrases (e.g., NPs, VPs, etc.) demonstrated that it was difficult to ensure that semantically anomalous or ungrammatical sentences would not be generated. Our generation strategy is reminiscent of earlier work on revision for summarization (Robin and McKeown, 1996), although Robin and McKeown used a three-tiered representation of each sentence, including its semantics, deep, and surface syntax, all of which were used as triggers for revision.

The three steps of the fusion lattice computation are ordered as follows: selection of the *basis tree*, augmentation of the tree with alternative verbalizations, and pruning of the extraneous subtrees. Alignment is essential for all the steps. The selection of the basis tree is guided by the number of intersection subtrees it includes; in the best case, it contains all such subtrees. The basis tree is the centroid of the input sentences — a sentence which is the most similar to the other sentences in the input. Using the alignment-based similarity score described in Section 3.1.2, we identify a centroid by computing for each sentence the average similarity score between the sentence and the rest of the input sentences, and then selecting a sentence with a maximal score.

Next, we augment the basis tree with information present in the other input sentences. More specifically, we add alternative verbalizations for the nodes in the basis tree and the intersection subtrees which are not part of the basis tree. The alternative verbalizations are readily available from the pairwise alignments of the basis tree with other trees in the input computed in the previous section. For each node of the basis tree we record all verbalizations from the nodes of the other input trees aligned with a given node. A verbalization can be a single word, or it can be a phrase, if a node represents a noun compound or a verb with a particle. An example of a fusion lattice, augmented with alternative verbalizations, is given in Figure 6. Even after this augmentation, the fusion lattice may not include all of the intersection subtrees. The main difficulty in subtree insertion is finding its acceptable placement, which is often determined by various sources of knowledge: syntactic, semantic and idiosyncratic. Therefore, we only insert subtrees whose top node aligns with one of the nodes in a basis tree. We further constrain the insertion procedure by only inserting trees that appear in at least half of the sentences of a theme. These two restrictions prevent the algorithm from generating overly long,

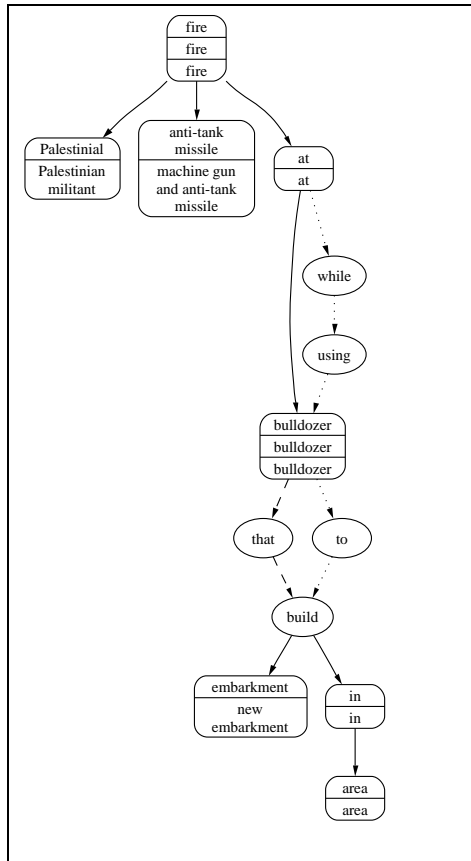


**Figure 6**  
 A basis lattice before and after the augmentation. Solid lines represent aligned edges of the basis tree. Dashed edges represent unaligned edges of the basis tree, and dotted edges represent insertions from other theme sentences.

unreadable sentences.<sup>9</sup>

Finally, subtrees which are not part of the intersection are pruned off the basis tree. However, removing all such subtrees may result in an ungrammatical or semantically flawed sentence; for example, we might create a sentence without a subject. This over-pruning may happen if either the input to the fusion algorithm is noisy, or the alignment has failed to identify the similarity between some subtrees. Therefore, we perform more conservative pruning, deleting self-contained components which can be removed without leaving ungrammatical sentences. As previously observed in the literature (Mani, Gates, and Bloedorn, 1999; Jing and McKeown, 2000), such components include a clause in the clause conjunction, relative clauses, and some elements within a clause (such as adverbs and prepositions). For example, this procedure transforms the lattice in Figure 6 into the pruned basis lattice shown in Figure 7 by deleting the clause “*the clash erupted*” and the verb phrase “*to better protect Israeli forces.*” These phrases are eliminated because they do not appear in other sentences of a theme, and at the same time their removal

<sup>9</sup> The preference for shorter fusion sentences is further enforced during the linearization stage because our scoring function monotonically decreases with the length of a sentence.



**Figure 7**  
A pruned basis lattice.

does not interfere with the well-formedness of the fusion sentence. Once these subtrees are removed, the fusion lattice construction is completed.

### 3.3 Generation

The final stage in sentence fusion is linearization of the fusion lattice. Sentence generation includes selection of a tree traversal order, lexical choice among available alternatives and placement of auxiliaries, such as determiners. Our generation method utilizes information given in the input sentences to restrict the search space and then chooses among remaining alternatives based using a language model derived from a large text collection. We first motivate the need in reordering and rephrasing, and then discuss our implementation.

For the word ordering task, we do not have to consider all the possible traversals, since the number of valid traversals is limited by ordering constraints encoded in the fusion lattice. However, the basis lattice does not uniquely determine the ordering: the placement of trees inserted in the basis lattice from other theme sentences are not restricted by the original basis tree. While the ordering of many sentence constituents is determined by their syntactic roles, some constituents, such as time, location and manner circumstantials, are free to move (Elhadad et al., 2001). Therefore, the algorithm still

has to select an appropriate order from among different orders of the inserted trees.

The process so far produces a sentence that can be quite different from the extracted sentence; while the basis sentences provides guidance for the generation process, constituents may be removed, added in, or reordered. Wording can also be modified during this process. Although the selection of words and phrases which appear in the basis tree is a safe choice, enriching the fusion sentence with alternative verbalizations has several benefits. In applications such as summarization, where the length of the produced sentence is a factor, a shorter alternative is desirable. This goal can be achieved by selecting the shortest paraphrase among available alternatives. Alternate verbalizations can also be used to replace anaphoric expressions, for instance, when the basis tree contains a noun phrase with anaphoric expressions (e.g., “*his visit*”) and one of the other verbalizations is anaphora-free. Substitution of the latter for the anaphoric expression may increase the clarity of the produced sentence, since frequently the antecedent of the anaphoric expression is not present in a summary. In addition to cases where substitution is preferable but not mandatory, there are cases where it is a required step for generation of a fluent sentence. As a result of subtree insertions and deletions, the words used in the basis tree may not be a good choice after the transformations and the best verbalization would be achieved by using their paraphrase from another theme sentence. As an example, consider the case of two paraphrasing verbs with different subcategorization frames, such as “tell” and “say”. If the phrase “*our correspondent*” is removed from the sentence “*Sharon told our correspondent that the elections were delayed . . .*”, a replacement of the verb “*told*” with “*said*” yields a more readable sentence.

The task of auxiliary placement is alleviated by the presence of features stored in the input nodes. In most cases, aligned words stored in the same node have the same feature values, which uniquely determine an auxiliary selection and conjugation. However, in some cases, aligned words have different grammatical features, in which case the linearization algorithm needs to select among available alternatives.

Linearization of the fusion sentence involves the selection of the best phrasing and placement of auxiliaries as well as the determination of optimal ordering. Since we do not have sufficient semantic information to perform such selection, our algorithm is driven by corpus-derived knowledge. We generate all possible sentences<sup>10</sup> from the valid traversals of the fusion lattice, and score their likelihood according to statistics derived from a corpus. This approach, originally proposed by (Knight and Hatzivassiloglou, 1995; Langkilde and Knight, 1998), is a standard method used in statistical generation. We trained a trigram model with Good-Turing smoothing over 60 Megabyte of news articles collected by Newsblaster using the second version CMU-Cambridge Statistical Language Modeling toolkit (Clarkson and Rosenfeld, 1997). The sentence with the lowest length-normalized entropy (the best score) is selected as the verbalization of the fusion lattice. Table 4 shows several verbalizations produced by our algorithm from the central tree in Table 7. Here, we can see that the lowest scoring sentence is both grammatical and concise.

This table also illustrates that entropy-based scoring does not always correlate with the quality of the generated sentence. For example, the fifth sentence in Table 4 — “*Palestinians fired anti-tank missile at a bulldozer to build a new embankment in the area.*” — is not a well-formed sentence; however, our language model gave it a better score than its well-formed alternatives, the second and the third sentences (See Section 4 for further discussion.) Despite these shortcomings, we preferred entropy-based scoring to symbolic linearization. In the next section, we motivate our choice.

---

<sup>10</sup> Due to the efficiency constraints imposed by Newsblaster, we sample only a subset of 20,000 paths. The sample is selected randomly.



Sentence	Entropy
Palestinians fired an anti-tank missile at a bulldozer.	4.25
Palestinian militants fired machine-guns and anti-tank missiles at a bulldozer.	5.86
Palestinian militants fired machine-guns and anti-tank missiles at a bulldozer that was building an embankment in the area.	6.22
Palestinians fired anti-tank missiles at while using a bulldozer.	7.04
Palestinians fired anti-tank missile at a bulldozer to build a new embankment in the area.	5.46

**Table 4**

Alternative linearizations of the fusion lattice with the corresponding entropy values

**3.3.1 Statistical versus Symbolic Linearization.** In the previous version of the system (Barzilay, McKeown, and Elhadad, 1999), we performed linearization of a fusion dependency structure using the language generator FUF/SURGE (Elhadad and Robin, 1996). As a large-scale linearizer used in many traditional semantic-to-text generation systems, FUF/SURGE could be an appealing solution to the task of surface realization. Because the input structure and the requirements on the linearizer are quite different in text-to-text generation, we had to design rules for mapping between dependency structures produced by the fusion component and FUF/SURGE input. For instance, FUF/SURGE requires that the input contain a semantic role for prepositional phrases, such as *manner*, *purpose* or *location*, which is not present in our dependency representation; thus we had to augment the dependency representation with this information. In the case of inaccurate prediction or the lack of relevant semantic information, the linearizer scrambles the order of sentence constituents, selects wrong prepositions or even fails to generate an output. Another feature of the FUF/SURGE system that negatively influences system performance is its limited ability to reuse phrases readily available in the input, instead of generating every phrase from scratch. This makes the generation process more complex and, thus, prone to error.

While the initial experiments conducted on a set of manually constructed themes seemed promising, the system performance deteriorated significantly when it was applied to automatically constructed themes. Our experience led us to believe that transformation of an arbitrary sentence into FUF/SURGE input representation is similar in its complexity to semantic parsing, a challenging problem on its own right. Rather than refining the mapping mechanism, we modified MultiGen to use a statistical linearization component, which handles uncertainty and noise in the input in a more robust way.

## 4. Sentence Fusion Evaluation

In our previous work, we evaluated the overall summarization strategy of MultiGen in multiple experiments, including comparisons with human-written summaries in the DUC<sup>11</sup> evaluation (McKeown et al., 2001; McKeown et al., 2002) and quality assessment in the context of a particular information access task in the Newsblaster framework (McKeown et al., 2002).

In this paper, we aim to evaluate the sentence fusion algorithm in isolation from other system components; we analyze the algorithm performance in terms of content selection and the grammaticality of the produced sentences. We will first present our evaluation methodology (Section 4.1), then we describe our data (Section 4.2), the results (Section 4.3) and their analysis (Section 4.4).

### 4.1 Methods

**Construction of a reference sentence** We evaluated content selection by comparing an automatically generated sentence with a reference sentence. The reference sentence was produced by a human (hereafter RFA) who was instructed to generate a sentence conveying information common to many sentences in a theme. The RFA was not familiar with the fusion algorithm. The RFA was provided with the list of the theme sentences; the original documents were not included. The instructions given to the RFA included several examples of themes with fusion sentences generated by the authors. Even though the RFA was not instructed to use phrases from input sentences, the sentences presented

---

<sup>11</sup> DUC (Document Understanding Conference) is a community-based evaluation of summarization systems organized by DARPA.

as examples reused many phrases from the input sentences. We believe that phrase reuse elucidates the connection between input sentences and a resulting fusion sentence. An example of a theme, a reference sentence and a system output are shown in Table 5.

#1	The forest is about 70 miles west of Portland.
#2	Their bodies were found Saturday in a remote part of Tillamook State Forest, about 40 miles west of Portland.
#3	Elk hunters found their bodies Saturday in the Tillamook State Forest, about 60 miles west of the family’s hometown of Portland.
#4	The area where the bodies were found is in a mountainous forest about 70 miles west of Portland.
Reference	The bodies were found Saturday in the forest area west of Portland.
System	The bodies <sub>4</sub> were found <sub>2</sub> Saturday <sub>2</sub> in <sub>3</sub> the Tillamook <sub>3</sub> State <sub>3</sub> Forest <sub>3</sub> west <sub>2</sub> of <sub>2</sub> Portland <sub>2</sub> .
#1	Four people including an Islamic cleric have been detained in Pakistan after a fatal attack on a church on Christmas Day.
#2	Police detained six people on Thursday following a grenade attack on a church that killed three girls and wounded 13 people on Christmas Day.
#3	A grenade attack on a Protestant church in Islamabad killed five people, including a U.S. Embassy employee and her 17 - year - old daughter.
Reference	A grenade attack on a church killed several people.
System	A <sub>3</sub> grenade <sub>3</sub> attack <sub>3</sub> on <sub>3</sub> a protestant <sub>3</sub> church <sub>3</sub> in <sub>3</sub> Islamabad <sub>3</sub> killed <sub>3</sub> six <sub>2</sub> people <sub>2</sub> .

**Table 5**

Examples from the test set. Each example contains a theme, a reference sentence generated by the RFA and a sentence generated by the system. Subscripts in the system-generated sentence represent a theme sentence from which a word was extracted.

**Data Selection** We wanted to test the performance of the fusion component on automatically computed inputs which reflect the accuracy of the existing preprocessing tools. For this reason, the test data was selected randomly from material collected by Newsblaster. To remove themes irrelevant for fusion evaluation, we introduced two additional filters. First, we excluded themes that contain identical or nearly identical sentences (with cosine similarity higher than 0.8). When processing such sentences, our algorithm reduces to sentence extraction which does not allow us to evaluate generation abilities of our algorithm. Second, themes for which the RFA was unable to create a reference sentence were also removed from the test set. As we mentioned above, Simfinder does not always produce accurate themes,<sup>12</sup> and therefore, the RFA could choose not to generate a reference sentence if the theme sentences had little in common. An example of a theme for which no sentence was generated is shown in Table 6. As a result of this filtering, 34% of the sentences were removed.

**Baselines** In addition to the system-generated sentence, we also included in the evaluation a fusion sentence generated by another human (hereafter, RFA2) and three baselines. (Following the DUC terminology, we refer to the baselines, our system and the RFA2 peers.) The first baseline is the shortest sentence among the theme sentences, which is obviously grammatical, and also it has a good chance of being representative of

<sup>12</sup> To mitigate the effects of Simfinder noise in MultiGen, we induced a similarity threshold on input trees — trees which are not similar to the basis tree are not used in the fusion process.

The shares have fallen 60 percent this year.
They said Qwest was forcing them to exchange their bonds at a fraction of face value — between 52.5 percent and 82.5 percent, depending on the bond — or else fall lower in the pecking order for repayment in case Qwest went broke.
Qwest had offered to exchange up to \$12.9 billion of the old bonds, which carried interest rates between 5.875 percent and 7.9 percent.
The new debt carries rates between 13 percent and 14 percent.
Their yield fell to about 15.22 percent from 15.98 percent.

**Table 6**

An example of noisy Simfinder output.

common topics conveyed in the input. The second baseline is produced by a simplification of our algorithm, where paraphrase information is omitted during the alignment process. This baseline is included to capture the contribution of paraphrase information to the performance of the fusion algorithm. The third baseline consists of the basis sentence. The comparison with this baseline reveals the contribution of the insertion and deletion stages in the fusion algorithm. The comparison against an RFA2 sentence provides an upper boundary on the system and baselines performance. In addition, this comparison will shed light on the human agreement on this task.

**Comparison against a reference sentence** The judge is given a peer sentence along with the corresponding reference sentence. The judge also has access to the original theme from which these sentences were generated. The order of the presentation is randomized across themes and peer systems. Reference and peer sentences are divided into clauses by the authors. The judges assess overlap on the clause-level between reference and peer sentences. The wording of the instructions was inspired by the DUC instructions for clause comparison. For each clause in the reference sentence, the judge decides whether the meaning of a corresponding clause is conveyed in a peer sentence. In addition to 0 score for “No Overlap” and 1 for “Full Overlap,” this framework allows for “Partial Overlap” with a score of 0.5. From the overlap data, we compute weighted recall and precision based on fractional counts (Hatzivassiloglou and McKeown, 1993). Recall is a ratio of weighted clause overlap between a peer and a reference sentence, and the number of clauses in a reference sentence. Precision is a ratio of weighted clause overlap between a peer and a reference sentence, and the number of clauses in a peer sentence.

**Grammaticality assessment** Grammaticality is rated in three categories: “Grammatical” (3), “Partially Grammatical” (2), and “Not Grammatical” (1). The judges were instructed to rate a sentence in the “Grammatical” category if it didn’t contain any grammatical mistakes. The “Partially Grammatical” included sentences that contain at most one mistake in agreement, articles and tense realization. The “Non Grammatical” category includes sentences that are corrupted by multiple mistakes of the former type, order sentence components in erroneous fashion or miss important components (e.g., subject).

Punctuation is one issue in assessing grammaticality. Proper placement of punctuation is a limitation specific<sup>13</sup> to our implementation of the sentence fusion algorithm that we are well aware of. Therefore, in our grammaticality evaluation (following the DUC procedure), the judge was asked to ignore punctuation.

## 4.2 Data

To evaluate our sentence fusion algorithm, we selected 100 themes following the procedure described in the previous section. Each set varied from two to seven sentences, with 3.82 sentences on average. The generated fusion sentences consisted of 1.91 clauses on average. None of the sentences in the test set were fully extracted; on average, each sentence fused fragments from 2.14 theme sentences. Out of 100 sentence, 57 sentences produced by the algorithm combined phrases from several sentences, while the rest of the sentences comprised subsequences of one of the theme sentences. Note that compression is different from sentence extraction. We included these sentences in the evaluation, because they reflect both content selection and realization capacities of the algorithm.

---

<sup>13</sup> We were unable to develop a set of rules which works in most cases. Punctuation placement is determined by a variety of features; considering all possible interactions of these features is hard. We believe that corpus-based algorithms for automatic restoration of punctuation developed for speech recognition applications (Beeferman, Berger, and Lafferty, 1998; Shieber and Tao, 2003) could help in our task, and we plan to experiment with them in the future.

Table 5 shows two sentences from the test corpus, along with input sentences. The examples are chosen so as to reflect good and bad performance cases. Note that the first example results in inclusion of the essential information (the fact that bodies were found, along with time and place) and leaves out details (that it was a remote location or how many miles west it was, a fact that is in dispute in any case). The problematic example incorrectly selects the number of people killed as six, even though this number is not repeated and different numbers are referred to in the text. This mistake is caused by a noisy entry in our paraphrasing dictionary which erroneously identifies “five” and “six” as paraphrases of each other.

### 4.3 Results

Table 7 shows the compression rate, precision, recall, F-measure and grammaticality score for each algorithm. The compression rate of a sentence was computed as the ratio of its output length to the average length of the theme input sentences.

We use  $\chi^2$  tests to determine whether the performance of our method in terms of F-measure gain is significantly different from RFA2 and other baselines (see Table 8). The presence of the diacritics \* and \*\* in Table 7 indicates significant differences (at  $p < 0.05$  and  $p < 0.01$ , respectively).

Peer	Compression	Precision	Recall	F-measure	Grammaticality
RFA2	54%	98%	94%	96%**	2.9
System	78%	65%	72%	68%	2.3
Baseline 1	69%	52%	38%	44%**	3
Baseline 2	111%	41%	67%	51%*	3
Baseline 3	73%	63%	64%	63%	2.4

**Table 7**

Evaluation results for a human-crafted fusion sentence(RFA2), our system output, the shortest sentence in the theme (baseline 1), the basis sentence(baseline 2) and a simplified version of our algorithm without paraphrasing information (baseline 3).

Peer	$\chi^2$	Confidence Level
Fusion/RFA2	30.49	< 0.01
Fusion/Baseline 1	10.71	< 0.01
Fusion/Baseline 2	5.29	< 0.05
Fusion/Baseline 3	0.35	not significantly different

**Table 8**

$\chi^2$  test on F-measure.

### 4.4 Discussion

The results in Table 7 demonstrate that sentences manually generated by the second human participant (RFA2) are not only the shortest, but are also closest to the reference sentence in terms of selected information. The tight connection<sup>14</sup> between sentences generated by RFAs establishes a high upper boundary for the fusion task. While neither our system nor the baselines were able to reach this performance, the fusion algorithm clearly outperforms all the baselines in terms of content selection, at a reasonable level

<sup>14</sup> We cannot apply Kappa statistics (Siegel and Castellan, 1988) for measuring agreement in the content selection task since an event space is not well-defined. This prevents us from computing the probability of the random agreement.

of compression. The performance of baseline 1 and baseline 2 demonstrates that neither the shortest sentence nor the basis sentence are adequate substitutions for fusion in terms of content selection; both precision and recall are significantly lower, as measure by  $\chi^2$  test (see Table 8). The gap in recall between our system and baseline 3 confirms our hypothesis about the importance of paraphrasing information for the fusion process. Omission of paraphrases (baseline 2) causes a 8% drop in recall due to the inability to match equivalent phrases with different wording.

Table 7 also reveals a downside of the fusion algorithm: automatically generated sentences contain grammatical errors, unlike fully extracted, human-written sentences. Given the high sensitivity of humans to processing ungrammatical sentences, one has to consider the benefits of flexible information selection against the decrease in readability of the generated sentences. Sentence fusion may not be a worthy direction to pursue, if low grammaticality is intrinsic to the algorithm and its correction requires knowledge which cannot be automatically acquired. In the remainder of the section, we show that this is not the case. Our manual analysis of generated sentences revealed that most of the grammatical mistakes are caused by the linearization component, or, more specifically, by suboptimal scoring of the language model. Language modeling is an active area of research, and we believe that advancement in this direction will be able to drastically boost the linearization capacity of our algorithm.

**4.4.1 Error Analysis.** In this section, we discuss the results of our manual analysis of mistakes in content selection and surface realization. Note that in some cases multiple errors are entwined in one sentence, which makes it hard to distinguish between a sequence of independent mistakes and a cause-and-effect chains. Therefore, the presented counts should be viewed as approximations, rather than precise numbers.

We start with the analysis of the test set, and continue with the description of some interesting mistakes that we encountered during system development.

**Mistakes in Content Selection** Most of the mistakes in content selection can be attributed to problems with alignment. In most cases (17 cases), erroneous alignments missed relevant word mappings due to the lack of a corresponding entry in our paraphrasing resources. At the same time, mapping of unrelated words (as shown in Table 5) is quite rare (two cases). This performance level is quite predictable given the accuracy of an automatically constructed dictionary and limited coverage of WordNet. Noise in lexical resources was exacerbated by the simplicity of our weighting scheme supports limited forms of mapping typology, and also uses manually assigned weights. Even in the presence of accurate lexical information, the algorithm occasionally produced suboptimal alignments (four cases).

Another source of errors (two cases) is the algorithm’s inability to handle many-to-many alignments. Namely, two trees conveying the same meaning may not be decomposable into the node level mappings which our algorithm aims to compute. For example, the mapping between the sentences in Table 9 expressed by the rule “*X denied claims by Y*”  $\leftrightarrow$  “*X said that Y’s claim was untrue*” cannot be decomposed into smaller matching units. At least two mistakes resulted from noisy preprocessing (tokenization and parsing).

Syria denied claims by Israeli Prime Minister Ariel Sharon. . .
The Syrian spokesman said that Sharon’s claim was untrue. . .

**Table 9**

A pair of sentences which cannot be fully decomposed.

In addition to alignment, overcutting during the lattice pruning caused omission of

three clauses that were present in the corresponding reference sentences. The sentence “*Conservatives were cheering language.*” is an example of an incomplete sentence derived from the following input sentence: “*Conservatives were cheering language in the final version that insures that one-third of all funds for prevention programs be used to promote abstinence.*” The omission of a relative clause was possible because some sentences in the input theme contain a noun “*language*” without any relative clauses.

**Mistakes in Surface Realization** Grammatical mistakes included incorrect selection of determiners, erroneous word ordering, omission of essential sentence constituents, incorrect realization of negation constructions and tense. These mistakes (42) originated during linearization of the lattice, and were caused either by incompleteness of the linearizer or by suboptimal scoring of language model. Mistakes of the first type are caused by missing rules for generating auxiliaries given node features. An example of this phenomenon is the sentence “*The coalition to have play a central role.*”, which verbalizes the verb construction “*will have to play*” incorrectly. Our linearizer lacks the completeness of existing application-independent linearizers, such as the unification based FUF/SURGE (Elhadad and Robin, 1996) and the probabilistic Fergus (Bangalore and Rambow, 2000). Unfortunately, we were unable to reuse any of the existing large-scale linearizers due to significant structural differences between input expected by these linearizers and the format of a fusion lattice. We are currently working on adapting Fergus for the sentence fusion task.

Mistakes related to suboptimal scoring are more common — 33 out of 42; in these cases, a language model selected ill-formed sentences, assigning a worse score to a better sentence. The sentence “*The diplomats were given to leave the country in 10 days.*” illustrates a suboptimal linearization of the fusion lattice. The correct linearizations — “*The diplomats were given 10 days to leave the country.*” and “*The diplomats were ordered to leave the country in 10 days.*” — were present in the fusion lattice, but the language model picked the incorrect verbalization. We found that in 27 cases the optimal verbalizations (in the authors’ view) were ranked below the top ten sentences ranked by the language model. We believe that more powerful language models, which incorporate more linguistic knowledge (such as syntax-based models), can improve the quality of generated sentences.

**4.4.2 Further Analysis.** In addition to analyzing errors found in this particular study, we also regularly track the quality of generated summaries on Newsblaster’s web page. We have noted a number of interesting errors that crop up from time to time, which seem to require information about the full syntactic parse, semantics or even discourse. Consider, for example, the last sentence from a summary entitled “*Estrogen-Progestin Supplements Now Linked to Dementia*” is shown in Table 10. This sentence was created by sentence fusion and clearly, there is a problem. Certainly, there was a study “*finding the risk of dementia in women who took one type of combined hormone pill*” but it was not the government study which was abruptly halted last summer. In looking at the two sentences from which this summary sentence was drawn, we can see that there is a good amount of overlap between the two, but the component does not have enough information about the referents of the different terms to know that two different studies are involved and that fusion should not take place. One topic of our future work (Section 6) is on the problem of reference and summarization.

Another example is shown in Table 11. Here again, the problem is reference. The first error is in the references to “the segments”. The two uses of “segments” in the first source document sentence do not refer to the same entity and thus, when the modifier is dropped, we get an anomaly. The second, more unusual problem is in the equation of “Clinton/Dole”, “Dole/Clinton” and “Clinton and Dole”.



#1	Last summer, a government study was abruptly halted after finding an increased risk of breast cancer, heart attacks and strokes in women who took one type of combined hormone pill.
#2	The most common form of hormone replacement therapy, already linked to breast cancer, stroke and heart disease, does not improve mental functioning as some earlier studies suggested and may increase the risk of dementia, researchers said on Tuesday.
System	Last <sub>1</sub> summer <sub>1</sub> a <sub>1</sub> government <sub>1</sub> study <sub>1</sub> abruptly <sub>1</sub> was <sub>1</sub> halted <sub>1</sub> after <sub>1</sub> finding <sub>1</sub> the <sub>2</sub> risk <sub>2</sub> of <sub>2</sub> dementia <sub>2</sub> in <sub>1</sub> women <sub>1</sub> who <sub>1</sub> took <sub>1</sub> one <sub>1</sub> type <sub>1</sub> of <sub>1</sub> combined <sub>1</sub> hormone <sub>1</sub> pill <sub>1</sub> .

**Table 10**

An example of wrong reference selection. Subscripts in the generated sentence represent a theme sentence from which a word was extracted.

#1	The segments will revive the “Point-Counterpoint” segments popular until they stopped airing in 1979, but will instead be called “Clinton/Dole” one week and “Dole/Clinton” the next week.
#2	Clinton and Dole have signed up to do the segment for the next 10 weeks, Hewitt said.
#3	The segments will be called “Clinton Dole” one week and “Dole Clinton” the next.
System	The <sub>1</sub> segments <sub>1</sub> will <sub>1</sub> revive <sub>1</sub> the <sub>3</sub> segments <sub>3</sub> until <sub>1</sub> they <sub>1</sub> stopped <sub>1</sub> airing <sub>1</sub> in <sub>1</sub> 1979 <sub>1</sub> but <sub>1</sub> instead <sub>1</sub> will <sub>1</sub> be <sub>1</sub> called <sub>1</sub> Clinton <sub>2</sub> and <sub>2</sub> Dole <sub>2</sub> .

**Table 11**

An example of incorrect reference selection. Subscripts in the generated sentence represent a theme sentence from which a word was extracted.

## 5. Related Work

Text-to-text generation is an emerging area of NLP. Unlike traditional concept-to-text generation approaches, text-to-text generation methods take text as input, and transform it into a new text satisfying some constraints (e.g., length or level of sophistication). In addition to sentence fusion, compression algorithms (Grefenstette, 1998; Mani, Gates, and Bloedorn, 1999; Knight and Marcu, 2001; Jing and McKeown, 2000; Reizler et al., 2003) and methods for expansion of a multiparallel corpus (Pang, Knight, and Marcu, 2003) are other examples of such methods.

Compression methods were developed for single-document summarization, and they aim to reduce a sentence by eliminating constituents which are not crucial for its understanding nor salient enough to include in the summary. These approaches are based on the observation that the “importance” of a sentence constituent can often be determined based on shallow features, such as its syntactic role and the words it contains. For example, in many cases a relative clause that is peripheral to the central point of the document can be removed from a sentence without significantly distorting its meaning. While earlier approaches for text compression were based on symbolic reduction rules (Grefenstette, 1998; Mani, Gates, and Bloedorn, 1999), more recent approaches use an aligned corpus of documents and their human written summaries to determine which constituents can be reduced (Knight and Marcu, 2001; Jing and McKeown, 2000; Reizler et al., 2003). Alignment is made between the summary sentences, which have been manually compressed, and the original sentences from which they were drawn.

Knight and Marcu (2000) treat reduction as a translation process using a noisy-

channel model (Brown et al., 1993). In this model, a short (compressed) string is treated as a source and additions to this string are considered to be noise. The probability of a source string  $s$  is computed by the combination of a standard probabilistic context-free grammar score, which is derived from the grammar rules that yielded tree  $s$ , and a word-bigram score, computed over the leaves of the tree. The stochastic channel model creates a large tree  $t$  from a smaller tree  $s$  by choosing an extension template for each node based on the labels of the node and its children. In the decoding stage, the system searches for the short string  $s$  that maximizes  $P(s|t)$ , which (for fixed  $t$ ) is equivalent to maximizing  $P(s) * P(t|s)$ .

While this approach exploits only syntactic and lexical information, (Jing and McKeown, 2000) also rely on cohesion information, derived from word distribution in a text: phrases that are linked to a local context are kept, while phrases that have no such links are dropped. Another difference between these two methods is the extensive use of domain-independent knowledge sources in the latter. For example, a lexicon is used to identify which components of the sentence are obligatory to keep it grammatically correct. The corpus in this approach is used to estimate the degree to which the fragment is extraneous and can be omitted from a summary. A phrase is removed only if it is not grammatically obligatory, is not linked to a local context, and has a reasonable probability of being removed by humans. In addition to reducing the original sentences, (Jing and McKeown, 2000) use a number of manually compiled rules to aggregate reduced sentences; for example, reduced clauses might be conjoined with “*and*”.

Sentence fusion exhibits similarities with compression algorithms in the ways in which it copes with the lack of semantic data in the generation process, relying on shallow analysis of the input and statistics derived from a corpus. Clearly, the difference in the nature of both tasks and in the type of input they expect (single sentence versus multiple sentences) dictates the use of different methods. Having multiple sentences in the input poses new challenges — such as a need for sentence comparison — but at the same time it opens up new possibilities for generation. While the output of existing compression algorithms is always a substring of the original sentence, sentence fusion may generate a new sentence which is not a substring of any of the input sentences. This is achieved by arranging fragments of several input sentences into one sentence.

The only other text-to-text generation approach with a capability of producing novel utterances is that of Pang, Knight and Marcu (2003). Their method operates over multiple English translations of the same foreign sentence, and is intended to generate novel paraphrases of the input sentences. Like sentence fusion, their method aligns parse trees of the input sentences and then uses a language model to linearize the derived lattice. The main difference between the two methods is in the type of the alignment: our algorithm performs local alignment, while the algorithm of (Pang, Knight, and Marcu, 2003) performs global alignment. The differences in alignment are caused by differences in input: their method expects semantically equivalent sentences, while our algorithm operates over sentences with only partial meaning overlap. The presence of deletions and insertions in input sentences makes their alignment a particularly new significant challenge.

## 6. Conclusions and Future Work

In this paper, we presented sentence fusion, a novel method for text-to-text generation which, given a set of similar sentences, produces a new sentence containing the information common to most sentences. Unlike traditional generation methods, sentence fusion does not require an elaborate semantic representation of the input, but instead relies on the shallow linguistic representation automatically derived from the input documents

and knowledge acquired from a large text corpus. Generation is performed by reusing and altering phrases from input sentences.

As the evaluation described in Section 4 shows, our method accurately identifies common information and in most of the cases generates a well-formed fusion sentence. Our algorithm outperforms the shortest sentence baseline in terms of content selection, without a significant drop in grammaticality. We also show that augmenting the fusion process with paraphrasing knowledge improves the output by both measures. However, there is still a gap between our system and human performance.

An important goal for future work on sentence fusion is to increase the flexibility of this component. In our current implementation, we took a conservative approach which eliminates some valid combinations of input phrases in order to ensure a well-formed output. Therefore, we eliminate “high-risk” transformations, which reduces the generative power of the algorithm. This approach permits the possibility of a noisy alignment; furthermore, the language model does not effectively discriminate between grammatical and ungrammatical sentences. We believe that the process of aligning theme sentences can be improved by learning the similarity function, instead of using manually assigned weights. An interesting question is how such a similarity function can be induced in an unsupervised fashion. We can also improve the flexibility of the fusion algorithm by using a more powerful language model. Recent research (Daume et al., 2002) showed that syntax-based language models are more suitable for language-generation tasks; the study of such models is a promising direction to explore.

An important feature of the sentence fusion algorithm is its ability to generate multiple verbalizations of a given fusion lattice. In our implementation, this property is only utilized to produce grammatical texts in the changed syntactic context, but it can also be used to increase coherence of the text at the discourse level by taking context into account. In our current system, each sentence is generated in isolation, independently from what is said before and what will be said after. Clear evidence of the limitation of this approach is found in the selection of referring expressions. For example, all summary sentences may contain the full description of a named entity (e.g., “*President of Columbia University Lee Bollinger*”), while the use of shorter descriptions such as “*Bollinger*” or anaphoric expressions in some summary sentences would increase its readability (Schiffman, Nenkova, and McKeown, 2002; Nenkova and McKeown, 2003). These constraints can be incorporated into the sentence fusion algorithm, since our alignment-based representation of themes often contains several alternative descriptions of the same object.

Beyond the problem of referring expression generation, we found that by selecting appropriate paraphrases of each summary sentence, we can significantly improve the coherence of an output summary. An important research direction for future work is to develop a probabilistic text model that can capture properties of well-formed texts, just as a language model captures properties of sentence grammaticality. Ideally, such a model would be able to discriminate between cohesive fluent texts and ill-formed texts, guiding the selection of sentence paraphrases to achieve an optimal sentence sequence.

## References

- Bangalore, Srinivas and Owen Rambow. 2000. Exploiting a probabilistic hierarchical model for generation. In *Proceedings of COLING*, pages 42–48.
- Banko, Michele and Lucy Vanderwende. 2004. Using n-grams to understand the nature of summaries. In *Proceedings of HLT-NAACL*, pages 1–4.
- Barzilay, Regina. 2003. *Information Fusion for Multi-Document Summarization: Paraphrasing and Generation*. Ph.D. thesis, Columbia University.
- Barzilay, Regina and Michael Elhadad. 1997. Using lexical chains for text summarization. In

- Proceedings of the ACL Workshop on Intelligent Scalable Text Summarization*, pages 10–17, Madrid, Spain, August.
- Barzilay, Regina, Noemie Elhadad, and Kathleen McKeown. 2002. Inferring strategies for sentence ordering in multi-document news summarization. *Journal of Artificial Intelligence Research*, 17:35–55.
- Barzilay, Regina and Kathleen McKeown. 2001. Extracting paraphrases from a parallel corpus. In *Proceedings of the ACL/EACL*, pages 50–57.
- Barzilay, Regina, Kathleen McKeown, and Michael Elhadad. 1999. Information fusion in the context of multi-document summarization. In *Proceedings of the ACL*, pages 550–557.
- Beeferman, Doug, Adam Berger, and John Lafferty. 1998. Cyberpunc: A lightweight punctuation annotation system for speech. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 689–692.
- Borko, Harold and Charles Bernier. 1975. *Abstracting Concepts and Methods*. Academic Press, New York.
- Brown, Peter F., Stephen Della Pietra, Vincent Della Pietra, and Robert Mercer. 1993. The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*, 19(2):263–311.
- Carbonell, Jaime and Jade Goldstein. 1998. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of SIGIR*, pages 335–336.
- Chu-Carroll, Jennifer, Krzysztof Czuba, John Prager, and Abraham Ittycheriah. 2003. In question answering: Two heads are better than one. In *Proceedings of HLT-NAACL*, pages 24–31.
- Clarke, Charles, Gordon Cormack, and Thomas Lynam. 2001. Exploiting redundancy in question answering. In *Proceedings of SIGIR*, pages 358–365.
- Clarkson, Philip and R. Rosenfeld. 1997. Statistical language modeling using the cmu-cambridge toolkit. In *Proceedings ESCA Eurospeech*, volume 5, pages 2707–2710.
- Collins, Michael. 1997. Three generative, lexicalised models for statistical parsing. In *Proceedings of the ACL/EACL*, pages 16–23, Madrid, Spain.
- Daume, Hal, Kevin Knight, Irene Langkilde-Geary, Daniel Marcu, and Kenji Yamada. 2002. The importance of lexicalized syntax models for natural language generation tasks. In *Proceedings of INLG*, Arden House, NJ.
- Dumais, Susan, Michele Banko, Eric Brill, Jimmy Lin, and Andrew Ng. 2002. Web question answering: Is more always better? In *Proceedings of SIGIR*.
- Elhadad, Michael, Yael Netzer, Regina Barzilay, and Kathleen McKeown. 2001. Ordering circumstantials for multi-document summarization. In *Proceedings of BISFAI*.
- Elhadad, Michael and Jacques Robin. 1996. An overview of surge: A reusable comprehensive syntactic realization component. Technical Report 96-03, Dept of Mathematics and Computer Science, Ben Gurion University, Beer Sheva, Israel.
- Grefenstette, Gregory. 1998. Producing intelligent telegraphic text reduction to provide an audio scanning service for the blind. In *Proceedings of the AAAI Spring Workshop on Intelligent Text Summarization*, pages 111–115.
- Hatzivassiloglou, V., J. Klavans, and E. Eskin. 1999. Detecting text similarity over short passages: Exploring linguistic feature combinations via machine learning. In *Proceedings of the Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*.
- Hatzivassiloglou, Vasileios and Kathleen McKeown. 1993. Towards the automatic identification of adjectival scales: Clustering adjectives according to meaning. In *Proceedings of the ACL*, pages 172–182.
- Jing, Hongyang and Kathleen McKeown. 2000. Cut and paste based summarization. In *Proceedings of the First Conference of the North American Chapter of the Association of Computational Linguistics*, pages 178–185, Seattle, Washington.
- Knight, Kevin and Vasileios Hatzivassiloglou. 1995. Two-level, many-path generation. In *Proceedings of the ACL*, pages 252–260.
- Knight, Kevin and Daniel Marcu. 2001. Statistics-based summarization - step one: Sentence compression. In *Proceeding of the 17th National Conference of the American Association for Artificial Intelligence AAAI*, pages 703–710, Austin, Texas.

- Langkilde, Irene and Kevin Knight. 1998. Generation that exploits corpus-based statistical knowledge. In *Proceedings of the ACL/COLING*, pages 704–710.
- Lin, Chin-Yew and Eduard Hovy. 2002. From single to multi-document summarization: A prototype system and its evaluation. In *Proceedings of the ACL*, pages 457–464.
- Lin, Chin-Yew and Eduard Hovy. 2003. Automatic evaluation of summaries using n-gram co-occurrence statistics. In *Proceedings of HLT-NAACL*, pages 150–157.
- Mani, Inderjeet and Eric Bloedorn. 1997. Multi-document summarization by graph search and matching. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence (AAAI-97)*, pages 622–628, Providence, Rhode Island. AAAI.
- Mani, Inderjeet, Barbara Gates, and Eric Bloedorn. 1999. Improving summaries by revising them. In *Proceedings of the ACL*, pages 558–565.
- Marcu, Daniel and Laurie Gerber. 2001. An inquiry into the nature of multidocument abstracts, extracts, and their evaluation. In *Proceedings of the NAACL Workshop on Automatic Summarization*, pages 2–11.
- McKeown, Kathleen, Regina Barzilay, David Evans, Vasileios Hatzivassiloglou, Judith Klavans, Ani Nenkova, Carl Sable, Barry Schiffman, and Sergey Sigelman. 2002. Tracking and Summarizing News on a Daily Basis with Columbia’s Newsblaster. In *Proceedings of the Human Language Technology Conference (HLT-02)*.
- McKeown, Kathleen R., Regina Barzilay, David Evans, Vasileios Hatzivassiloglou, Min Yen Kan, Barry Schiffman, and Simone Teufel. 2001. Columbia multi-document summarization: Approach and evaluation. In *Proceedings of the Document Understanding Conference (DUC01)*.
- Melcuk, Igor. 1988. *Dependency Syntax: Theory and Practice*. Albany: State University of New York Press.
- Miller, G.A., R. Beckwith, C. Fellbaum, D. Gross, and K.J. Miller. 1990. Introduction to WordNet: An on-line lexical database. *International Journal of Lexicography (special issue)*, 3(4):235–245.
- Morris, Jane and Graeme Hirst. 1991. Lexical cohesion, the thesaurus, and the structure of text. *Computational Linguistics*, 17(1):21–48, March.
- Nenkova, Ani and Kathleen R. McKeown. 2003. References to named entities: A corpus study. In *Proceedings of the Human Language Technology Conference, Companion Volume*, pages 70–73.
- Pang, Bo, Kevin Knight, and Daniel Marcu. 2003. Syntax-based alignment of multiple translations: Extracting paraphrases and generating new sentences. In *Proceedings of HLT-NAACL*, pages 180–187.
- Papineni, Kishore A., Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: A method for automatic evaluation of machine translation. In *Proceedings of the ACL*, pages 311–318.
- Radev, Dragomir, Hongyan Jing, and Malgorzata Budzikowska. 2000. Centroid-based summarization of multiple documents: sentence extraction, utility-based evaluation, and user studies. In *Proceedings of the ANLP/NAACL 2000 Workshop on Automatic Summarization*.
- Radev, Dragomir and Kathleen R. McKeown. 1998. Generating natural language summaries from multiple on-line sources. *Computational Linguistics*, 24(3):469–500, September.
- Radev, Dragomir, John Prager, and Valerie Samn. 2000. Ranking suspected answers to natural language questions using predictive annotation. In *Proceedings of 6th Conference on Applied Natural Language Processing (ANLP)*, pages 150–157.
- Reizler, Stefan, Tracy H. King, Richard Crouch, and Annie Zaenen. 2003. Statistical sentence condensation using ambiguity packing and stochastic disambiguation methods for lexical-functional grammar. In *Proceedings of HLT-NAACL*, pages 197–204.
- Robin, Jacques and Kathleen McKeown. 1996. Empirically designing and evaluating a new revision-based model for summary generation. *Artificial Intelligence*, 85(1–2):135–179.
- Schiffman, Barry, Ani Nenkova, and Kathleen R. McKeown. 2002. Experiments in multidocument summarization. In *Proceedings of HLT*.
- Shieber, Stuart and Xiapong Tao. 2003. Comma restoration using constituency information. In *Proceedings of HLT-NAACL*, pages 142–148.
- Siegel, Sidney and N.John Castellan. 1988. *Non Parametric Statistics for Behavioral Sciences*. McGraw-Hill.
- Silber, Gregory and Kathleen McCoy. 2002. Computed lexical chains as an intermediate representation for automatic text summarization. *Computational Linguistics*, 28(4):487–496.