

COMS 4705 (Fall 2011)
Machine Translation Part I

Overview

- Challenges in machine translation
- Classical machine translation
- A brief introduction to statistical MT
- Evaluation of MT systems

Lexical Ambiguity

Example 1:

book the flight \Rightarrow reservar

read the book \Rightarrow libro

Example 2:

the box was in the pen

the pen was on the table

Example 3:

kill a man \Rightarrow matar

kill a process \Rightarrow acabar

Differing Word Orders

- English word order is *subject – verb – object*
- Japanese word order is *subject – object – verb*

English: IBM bought Lotus

Japanese: *IBM Lotus bought*

English: Sources said that IBM bought Lotus yesterday

Japanese: *Sources yesterday IBM Lotus bought that said*

Syntactic Structure is not Preserved Across Translations

The bottle floated into the cave



La botella entro a la cuerva flotando
(the bottle entered the cave floating)

Syntactic Ambiguity Causes Problems

John hit the dog with the stick



John golpeo el perro con el palo/que tenia el palo

Pronoun Resolution

The computer outputs the data; it is fast.



La computadora imprime los datos; **es** rapida

The computer outputs the data; it is stored in ascii.



La computadora imprime los datos; **están** almacenados en ascii

Differing Treatments of Tense

From Dorr et. al 1998:

Mary **went** to Mexico. During her stay she learned Spanish.

Went \Rightarrow iba (simple past/preterit)

Mary **went** to Mexico. When she returned she started to speak Spanish.

Went \Rightarrow fue (ongoing past/imperfect)

The Best Translation May not be 1-1

(From Manning and Schuetze):

According to our survey, 1988 sales of mineral water and soft drinks were much higher than in 1987, reflecting the growing popularity of these products. Cola drink manufacturers in particular achieved above average growth rates.

⇒

Quant aux eaux minerales et aux limonades, elles recontrent toujours plus d'adeptes. En effet notre sondage fait ressortir des ventes nettement superieures a celles de 1987, pour les boissons a base de cola notamment.

With regard to the mineral waters and the lemonades (soft drinks) they encounter still more users. Indeed our survey makes stand out the sales clearly superior to those in 1987 for cola-based drinks especially.

From Babel Fish:

Aznar ha premiado a Rodrigo Rato (vicepresidente primero), Javier Arenas (vicepresidente segundo y ministro de la Presidencia) y Eduardo Zaplana (ministro portavoz y titular de Trabajo) en la septima remodelacion de Gobierno en sus dos legislaturas. Las caras nuevas del Ejecutivo son las de Juan Costa, al frente del Ministerio de Ciencia y Tecnologia, y la de Julia Garcia Valdecasas, que ocupara la cartera de Administraciones Publicas.



Aznar has awarded to Rodrigo Short while (vice-president first), Javier Sands (vice-president second and minister of the Presidency) and Eduardo Zaplana (minister spokesman and holder of Work) in the seventh remodeling of Government in its two legislatures. The new faces of the Executive are those of Juan Coast, to the front of the Ministry of Science and Technology, and the one of Julia Garci'a Valdecasas, who will occupy the portfolio of Public Administrations.

An Example: Google Translation from Arabic

Stock prices retreated in the stock markets again with increasing concern about the circumstances surrounding the credit markets in the world, due mostly to the problems it faces American mortgage lending market, which raised concern among investors.

The index retreated Vuciji / 100 on the London Stock Exchange at the beginning of a percentage point in the dealings of up to 6082 points, while the Nikkei index retreated / 225 Japanese rate of 2.2% to close at the lowest level in eight months. The American Jones index has lost about 1.6 points Tuesday to reach 13029 points, the Nasdaq index had lost 1.7 of its value.

These declines came despite statements by the American Federal Reserve Bank (Central Bank), in which he said that the process of pumping more funds into capital markets when necessary.

The American Federal Reserve Board, for the purposes of relaxation of tension in global financial markets, resulting in the Gaza backtrackings American real estate lending, have pumped billions of dollars of emergency funds allocation to the banking sector during the past few days, on Friday and Monday. As the European Central Bank did the same.

Overview

- Challenges in machine translation
- **Classical machine translation**
- A brief introduction to statistical MT
- Evaluation of MT systems

Direct Machine Translation

- Translation is word-by-word
- Very little analysis of the source text (e.g., no syntactic or semantic analysis)
- Relies on a large bilingual dictionary. For each word in the source language, the dictionary specifies a set of rules for translating that word
- After the words are translated, simple reordering rules are applied (e.g., move adjectives after nouns when translating from English to French)

An Example of a set of Direct Translation Rules

(From Jurafsky and Martin, edition 2, chapter 25. Originally from a system from Panov 1960)

Rules for translating *much* or *many* into Russian:

if preceding word is *how* **return** *skol'ko*

else if preceding word is *as* **return** *stol'ko zhe*

else if word is *much*

if preceding word is *very* **return** *nil*

else if following word is a noun **return** *mnogo*

else (word is *many*)

if preceding word is a preposition and following word is noun **return** *mnogii*

else return *mnogo*

Some Problems with Direct Machine Translation

- Lack of any analysis of the source language causes several problems, for example:

- Difficult or impossible to capture long-range reorderings

English: Sources said that IBM bought Lotus yesterday
Japanese: *Sources yesterday IBM Lotus bought that said*

- Words are translated without disambiguation of their syntactic role

e.g., *that* can be a complementizer or determiner, and will often be translated differently for these two cases

They said *that* ...

They like *that* ice-cream

Transfer-Based Approaches

- Three phases in translation:

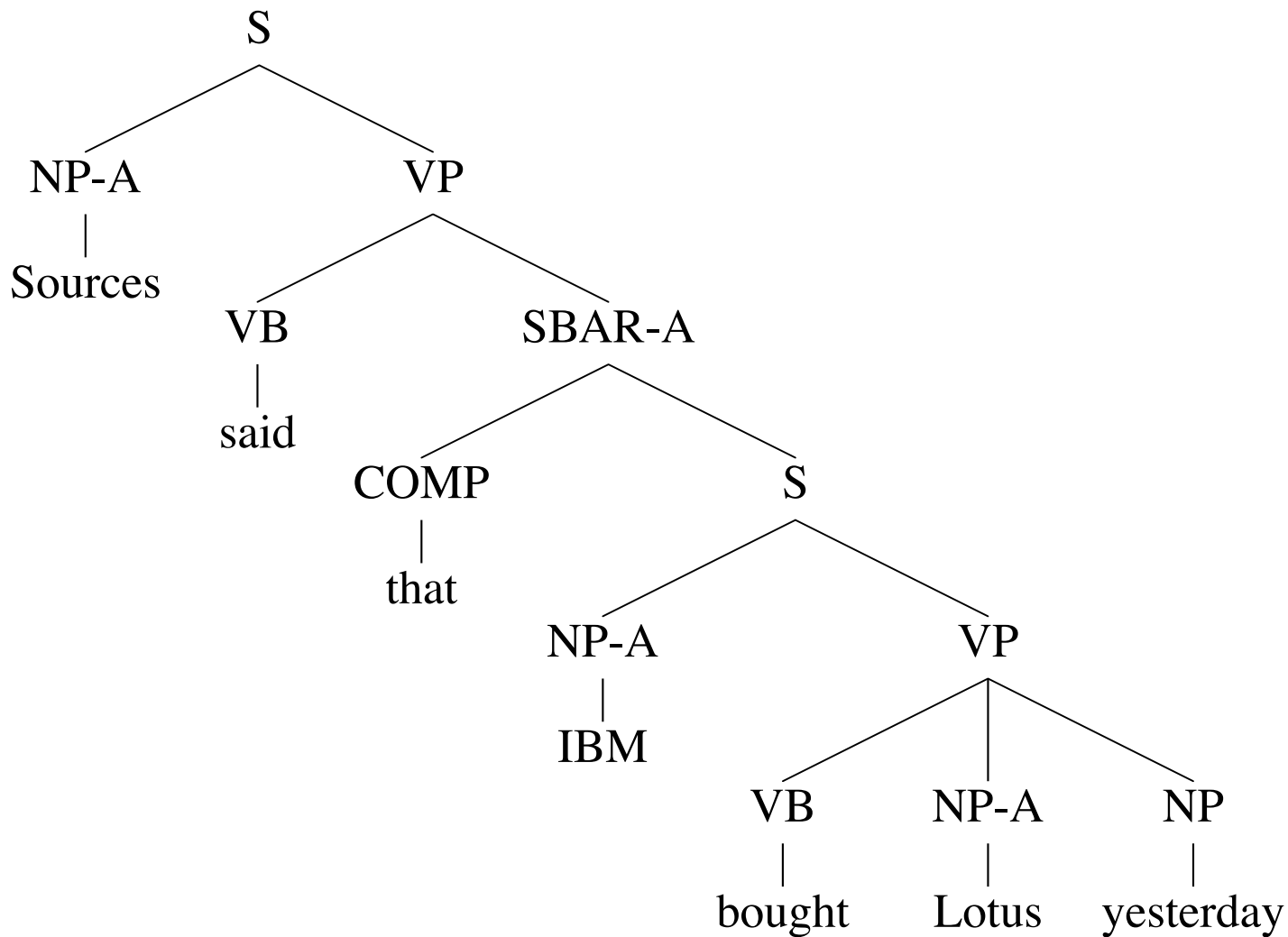
Analysis: Analyze the source language sentence; for example, build a syntactic analysis of the source language sentence.

Transfer: Convert the source-language parse tree to a target-language parse tree.

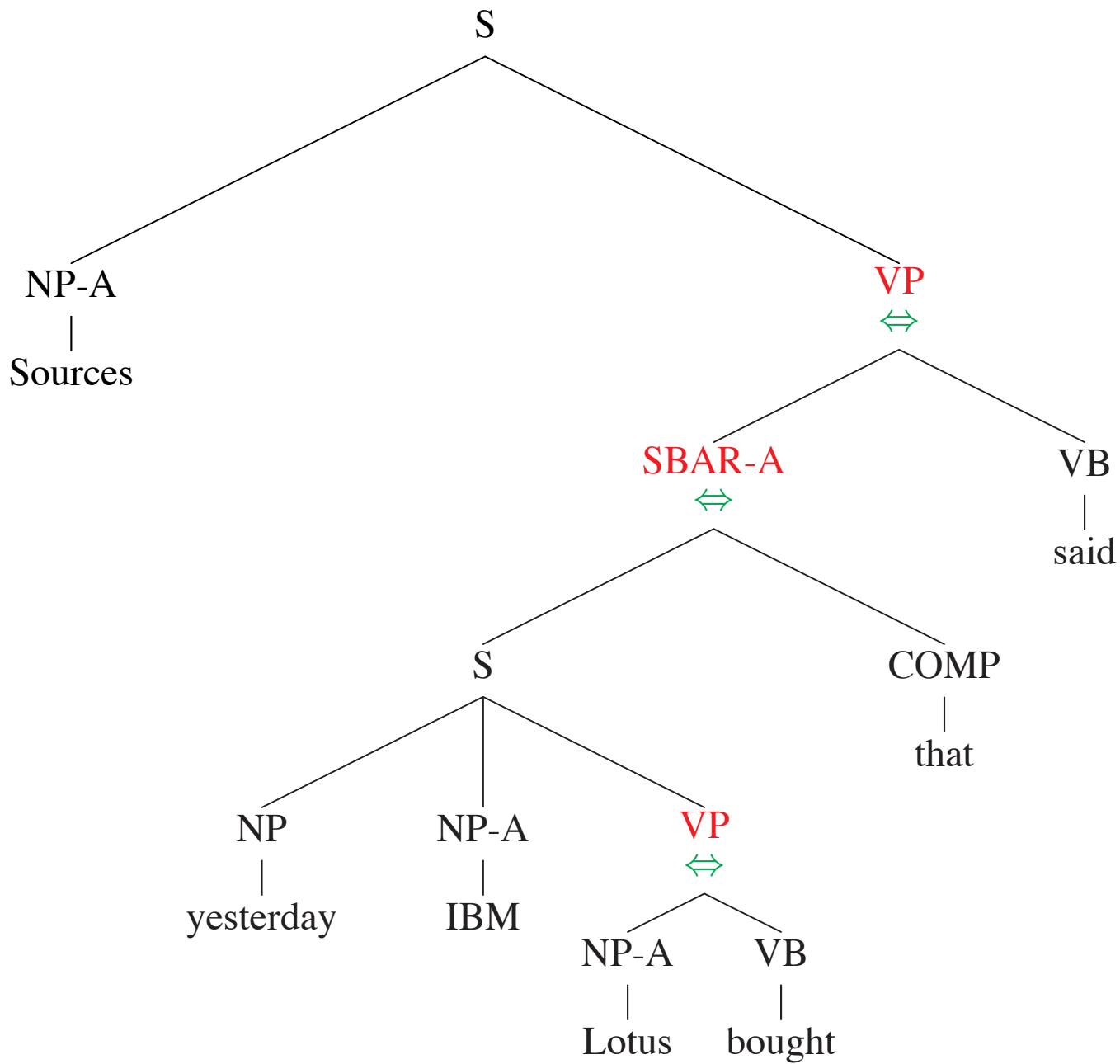
Generation: Convert the target-language parse tree to an output sentence.

Transfer-Based Approaches

- The “parse trees” involved can vary from shallow analyses to much deeper analyses (even semantic representations).
- The transfer rules might look quite similar to the rules for direct translation systems. But they can now operate on syntactic structures.
- It's easier with these approaches to handle long-distance reorderings
- The *Systran* systems are a classic example of this approach



⇒ Japanese: *Sources yesterday IBM Lotus bought that said*



Interlingua-Based Translation

- Two phases in translation:

Analysis: Analyze the source language sentence into a (language-independent) representation of its meaning.

Generation: Convert the meaning representation into an output sentence.

Interlingua-Based Translation

One Advantage: If we want to build a translation system that translates between n languages, we need to develop n analysis and generation systems. With a transfer based system, we'd need to develop $O(n^2)$ sets of translation rules.

Disadvantage: What would a language-independent representation look like?

Interlingua-Based Translation

- How to represent different concepts in an interlingua?
- Different languages break down concepts in quite different ways:

German has two words for *wall*: one for an internal wall, one for a wall that is outside

Japanese has two words for *brother*: one for an elder brother, one for a younger brother

Spanish has two words for *leg*: *pierna* for a human's leg, *pata* for an animal's leg, or the leg of a table

- An interlingua might end up simple being an intersection of these different ways of breaking down concepts, but that doesn't seem very satisfactory...

Overview

- Challenges in machine translation
- Classical machine translation
- A brief introduction to statistical MT
- Evaluation of MT systems

A Brief Introduction to Statistical MT

- Parallel corpora are available in several language pairs
- Basic idea: use a parallel corpus as a training set of translation examples
- Classic example: IBM work on French-English translation, using the Canadian Hansards. (1.7 million sentences of 30 words or less in length).
- Idea goes back to Warren Weaver (1949): suggested applying statistical and cryptanalytic techniques to translation.

The Noisy Channel Model

- Goal: translation system from French to English
- Have a model $p(e | f)$ which estimates conditional probability of any English sentence e given the French sentence f . Use the training corpus to set the parameters.
- A Noisy Channel Model has two components:

$p(e)$ **the language model**

$p(f | e)$ **the translation model**

- Giving:

$$p(e | f) = \frac{p(e, f)}{p(f)} = \frac{p(e)p(f | e)}{\sum_e p(e)p(f | e)}$$

and

$$\operatorname{argmax}_e p(e | f) = \operatorname{argmax}_e p(e)p(f | e)$$

More About the Noisy Channel Model

- The **language model** $p(e)$ could be a trigram model, estimated from any data (parallel corpus not needed to estimate the parameters)
- The **translation model** $p(f | e)$ is trained from a parallel corpus of French/English pairs.
- Note:
 - The translation model is backwards!
 - The language model can make up for deficiencies of the translation model.
 - Later we'll talk about how to build $p(f | e)$
 - Decoding, i.e., finding

$$\operatorname{argmax}_e p(e)p(f | e)$$

is also a challenging problem.

Example from Koehn and Knight tutorial

Translation from Spanish to English, candidate translations based on $p(\text{Spanish} | \text{English})$ alone:

Que hambre tengo yo

→

What hunger have $p(S|E) = 0.000014$

Hungry I am so $p(S|E) = 0.000001$

I am so hungry $p(S|E) = 0.0000015$

Have i that hunger $p(S|E) = 0.000020$

...

With $p(\text{Spanish} | \text{English}) \times p(\text{English})$:

Que hambre tengo yo

→

What hunger have $p(S|E)p(E) = 0.000014 \times 0.000001$

Hungry I am so $p(S|E)p(E) = 0.000001 \times 0.0000014$

I am so hungry $p(S|E)p(E) = 0.0000015 \times 0.0001$

Have i that hunger $p(S|E)p(E) = 0.000020 \times 0.00000098$

...