

# Introduction to Machine Translation

Joan Andreu Sánchez

Departamento Sistemas Informáticos y Computación  
Instituto Tecnológico de Informática  
Universidad Politécnica Valencia

PASCAL 2 Ghana Bootcamp 2011

**URL:** <http://www.dsic.upv.es/~jandreu>

**e-mail:** [jandreu@dsic.upv.es](mailto:jandreu@dsic.upv.es)

## Index

### 1. Introduction

- 1.1 Objectives of MT
- 1.2 Approaches to MT
- 1.3 Linguistic resources
- 1.4 Assessment

### 2. Statistical alignment models

- 2.1 Statistical framework to MT
- 2.2 Alignments
- 2.3 Statistical alignment models
- 2.4 Categorization in MT

### 3. Advanced statistical alignment models

- 3.1 Fertility-based models
- 3.2 The search problem
- 3.3 Using linguistic knowledge

### 4. Phrase-based models

- 4.1 Beyond word models
- 4.2 Phrase-based models

### 5. Syntax-based translation models

- 5.1 Introduction
- 5.2 ITG for MT
- 5.3 Tree-to-string models
- 5.4 Hierarchical MT

Slides in Sections 1, 2 and 3 have been prepared from slides supplied by F. Casacuberta.

## Index

### 1. Introduction

- 1.1 Objectives of MT
- 1.2 Approaches to MT
- 1.3 Linguistic resources
- 1.4 Assessment

### 2. Statistical alignment models

- 2.1 Statistical framework to MT
- 2.2 Alignments
- 2.3 Statistical alignment models
- 2.4 Categorization in MT

### 3. Advanced statistical alignment models

- 3.1 Fertility-based models
- 3.2 The search problem
- 3.3 Using linguistic knowledge

### 4. Phrase-based models

- 4.1 Beyond word models
- 4.2 Phrase-based models

### 5. Syntax-based translation models

- 5.1 Introduction
- 5.2 ITG for MT
- 5.3 Tree-to-string models
- 5.4 Hierarchical MT

## MT objectives: Erroneous conceptions

- MT is a waste of time because a machine never will translate Shakespeare
- In general, the quality of translation you can get from an MT system is very low
- MT threatens the jobs of translators
- There is an MT system that translates what you say into Japanese and translates the other speaker's replies in English

## MT objectives: Facts

- There are many situations that a MT systems produce reliable, if less than perfect, translations at high speed
- In some circumstances, MT systems can produce good quality outputs
- MT does not threaten translators' jobs: High demand of translations and too repetitive translation jobs
- Speech-to-speech MT is still a research topic
- There are many open research problems in MT
- Building a traditional MT system is a time consuming job
- A user will typically have to invest a considerable amount of effort in customizing an MT system

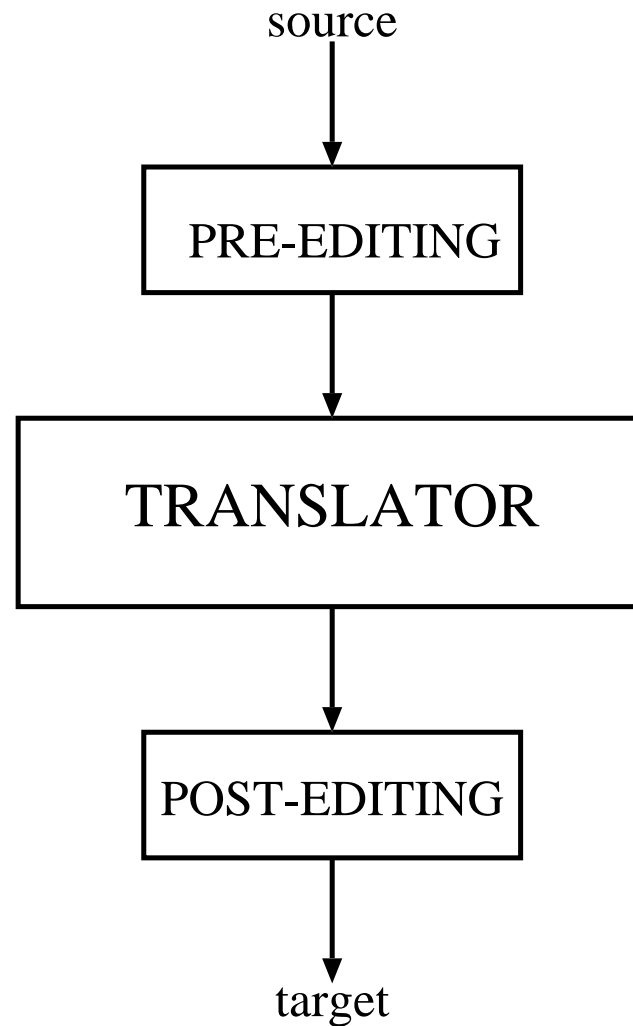
### MT objectives: need of pre/post-editing

- While the number of errors and bad constructions is high, “post-editing” can make the result useful.
- Many problems could have been avoided by making the source text “simpler”.
- Simplification of the translation problem by using adequate rules to produce “controlled” (i.e., simple and regular) source text.

# 1.1 OBJECTIVES OF MT

---

## General scheme for MT

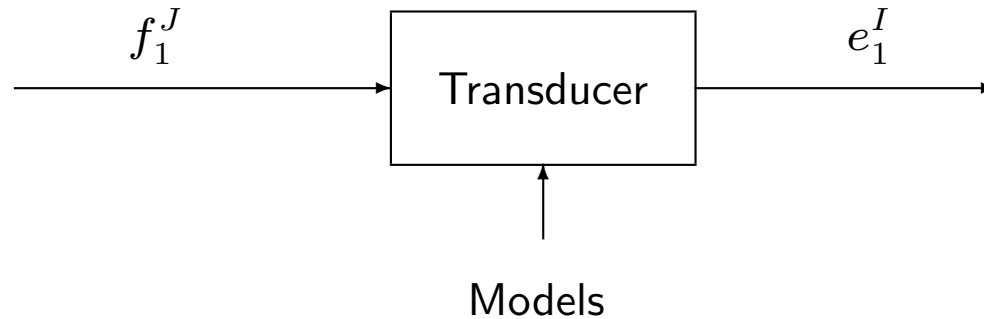


## Technologies

- (Linguistic) knowledge-based methods
- (Memorized) example-based methods
  - Translation memories
- Statistical models
  - Alignment models
  - Syntax-based models
  - Finite-State models
- Hybrid models



## Statistical MT



- Inverse approach (noisy channel)

$$\hat{e}_1^I = \arg \max_{e_1^I} \Pr(e_1^I | f_1^J) = \arg \max_{e_1^I} \Pr(f_1^J | e_1^I) \Pr(e_1^I)$$

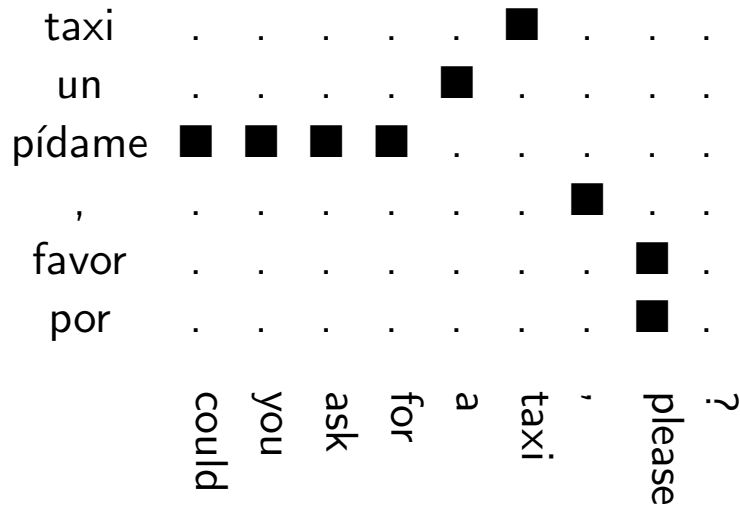
- Direct approach

$$\hat{e}_1^I = \arg \max_{e_1^I} \Pr(e_1^I | f_1^J)$$

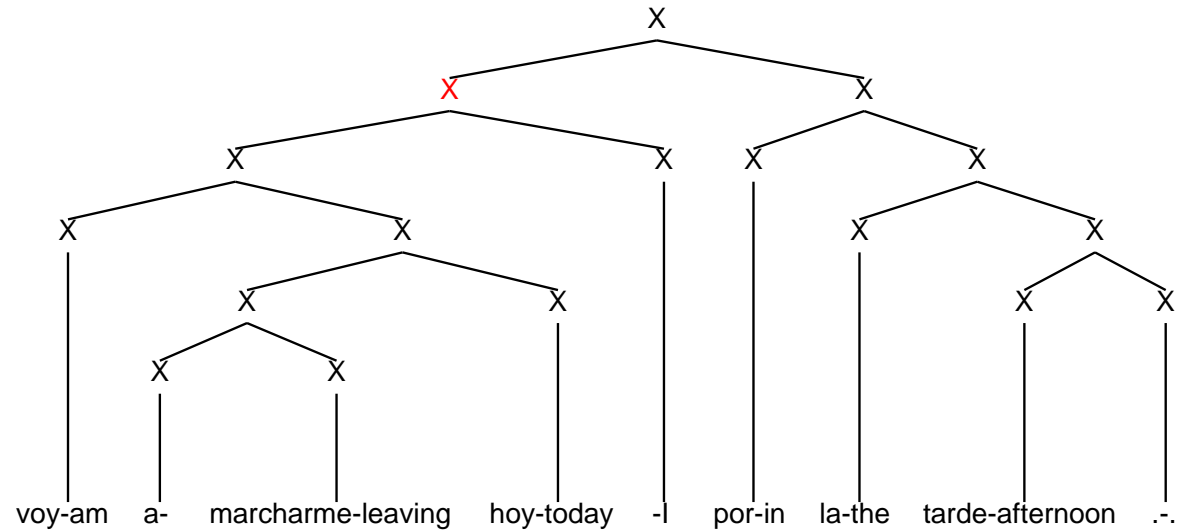
# 1.2 APPROACHES TO MT

## Statistical approaches to MT

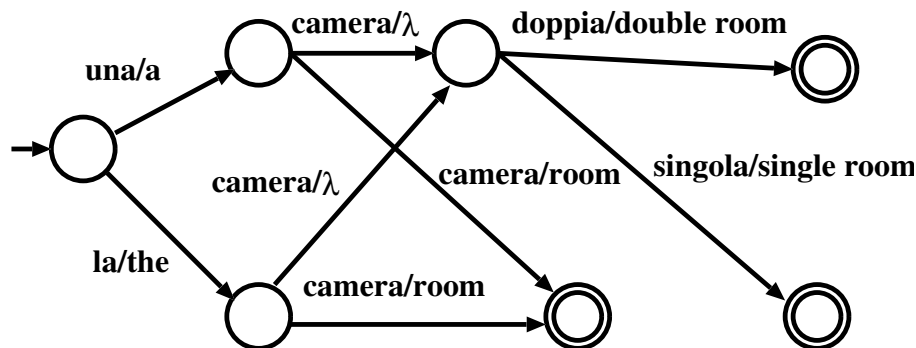
→ Word-alignment approaches



→ Syntactic approaches



→ Finite-state approaches



# 1.3 LINGUISTIC RESOURCES

---

## Resources

- Dictionaries
- Grammars
- Corpora
- Paragraph-aligned and Labeled Corpora

## 1.4 ASSESMENT

---

- Test sentences with reference translation
- Automatic assessment
  - Editing Distances:
    - Translation Word Error Rate (TWER)
    - Translation Error Rate (TER)
  - Multireference TWER
  - N-Gram based: *BLUE* and *NIST* score
  - . . . .

## Index

### 1. Introduction

- 1.1 Objectives of MT
- 1.2 Approaches to MT
- 1.3 Linguistic resources
- 1.4 Assessment

### 2. **Statistical alignment models**

- 2.1 Statistical framework to MT
- 2.2 Alignments
- 2.3 Statistical alignment models
- 2.4 Categorization in MT

### 3. Advanced statistical alignment models

- 3.1 Fertility-based models
- 3.2 The search problem
- 3.3 Using linguistic knowledge

### 4. Phrase-based models

- 4.1 Beyond word models
- 4.2 Phrase-based models

### 5. Syntax-based translation models

- 5.1 Introduction
- 5.2 ITG for MT
- 5.3 Tree-to-string models
- 5.4 Hierarchical MT

## 2.1 STATISTICAL FRAMEWORK FOR MT

---

### General framework

- Every sentence  $y$  in one language is a translation of any sentence  $x$  in another language
- For each possible pair of sentences,  $y$  and  $x$ , there is a probability  $\Pr(y | x)$
- The probability of pairs of sentences as  
*quiero una habitación doble con vistas al mar # are all expenses included in the bill ?*  
should be low
- The probability of pairs of sentences as  
*¿ hay alguna habitación tranquila libre ? # is there a quiet room available ?*  
should be high

### General framework

Given a source sentence  $x$ , search for the sentence  $\hat{y}$

$$\hat{y} = \arg \max_y \Pr(y | x)$$

### Approaches

- A direct approach: *maximum entropy models*
- An inverse approach: *channel models*

## 2.1 STATISTICAL FRAMEWORK FOR MT

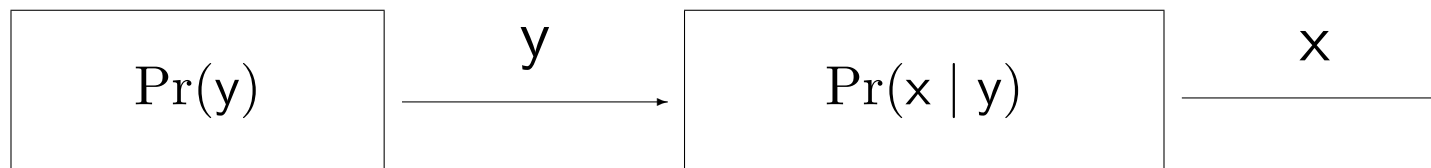
---

### An inverse approach

*Given a source sentence  $x$ , search for the sentence  $\hat{y}$*

$$\hat{y} = \arg \max_y \Pr(y | x) = \arg \max_y \Pr(x | y) \cdot \Pr(y)$$

A channel model



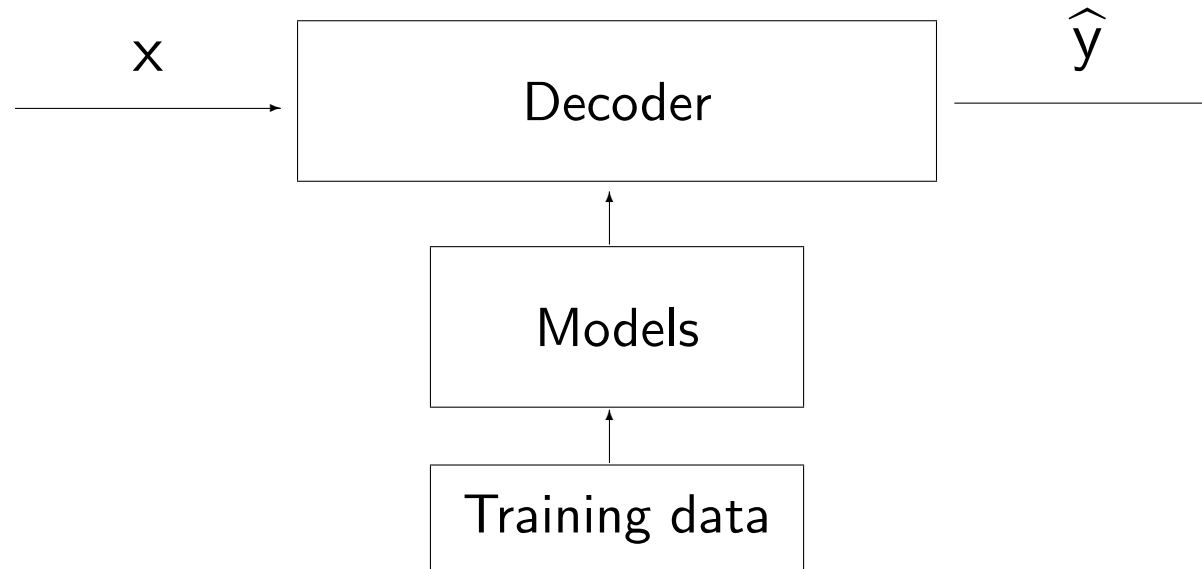
A target-language model + alignment and lexicon models



## 2.1 STATISTICAL FRAMEWORK FOR MT

---

### Translation search

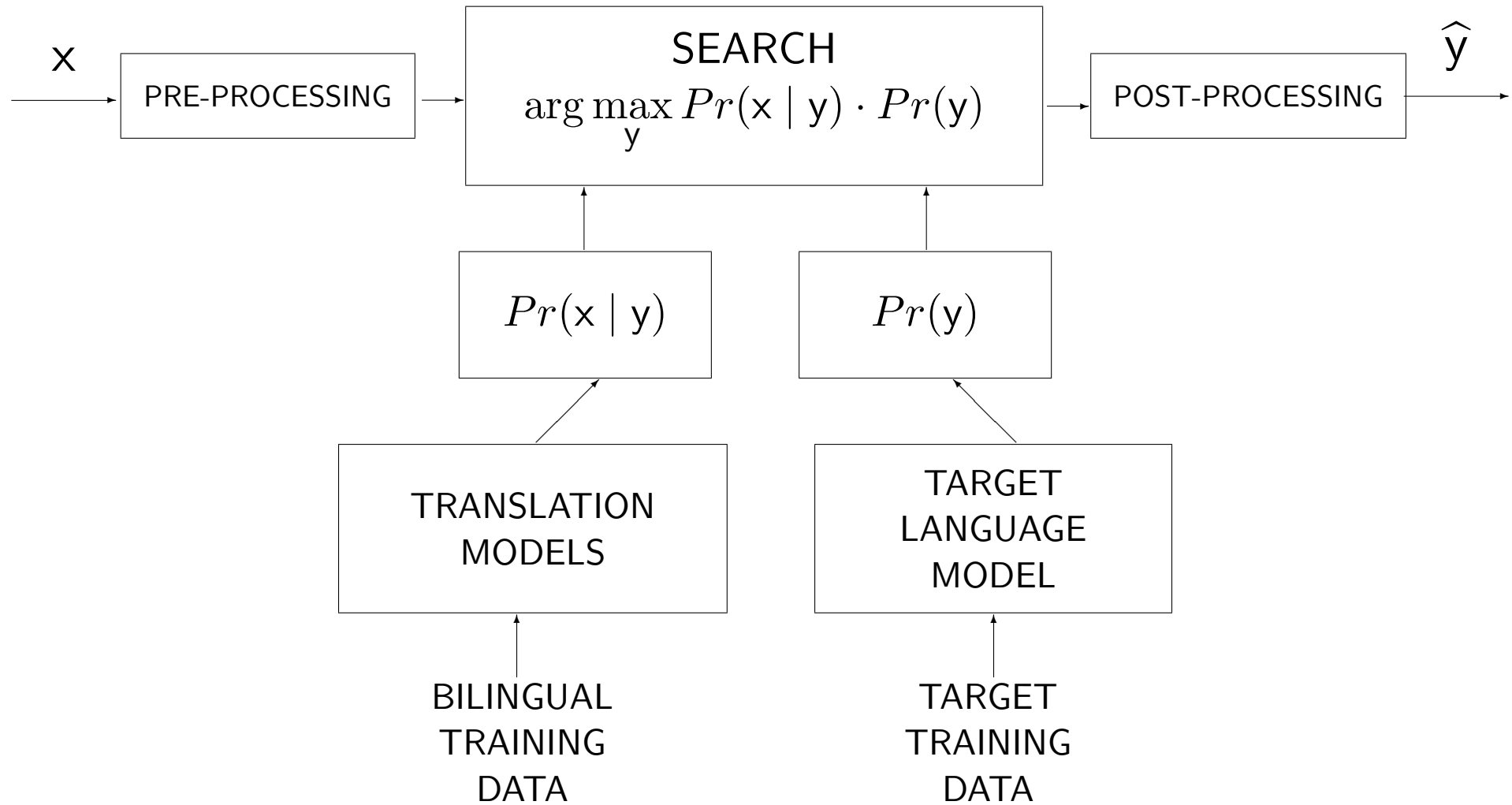


Inverse approach:

- A target-language model:  $\Pr(y) \approx Pr(y)$
- Translation models (alignment and lexicon models):  $\Pr(x | y) \approx Pr(x | y)$
- Search procedure:  $\hat{y} = \arg \max_y \Pr(x | y) \cdot Pr(y)$

## 2.1 STATISTICAL FRAMEWORK FOR MT

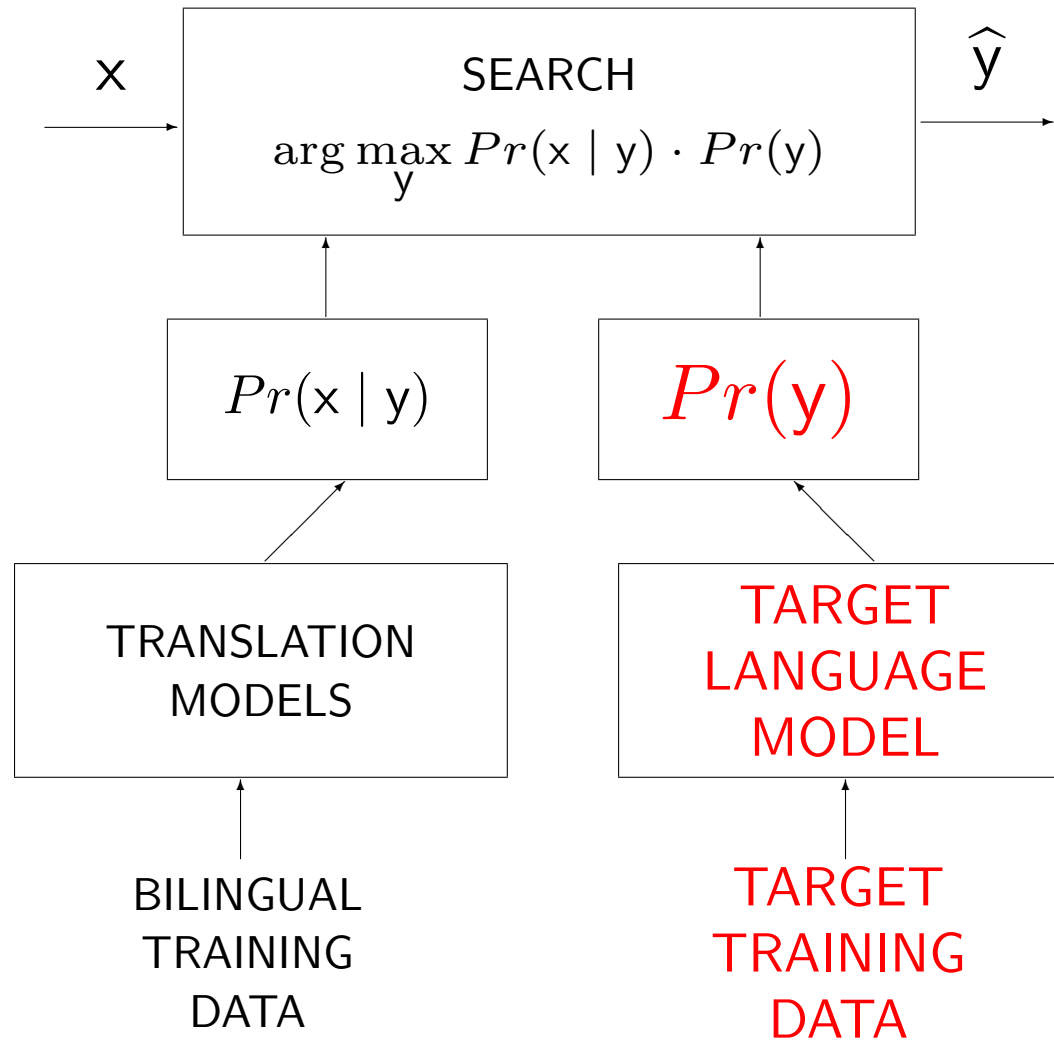
### An inverse approach



## 2.1 STATISTICAL FRAMEWORK FOR MT

---

An inverse approach: The target language model



### Language models

#### Word n-grams

$$\Pr(y) = \prod_{i=1}^{|y|} \Pr(y_i | y_1 \dots y_{i-1}) \approx Pr(y) = \prod_{i=1}^{|y|} p_n(y_i | y_{i-n+1} \dots y_{i-1})$$

#### n-grams of categories

$$\Pr(y) \approx Pr(y) = \prod_{i=1}^{|y|} p_n(C_i | C_{i-n+1} \dots C_{i-1}) \cdot p(y_i | C_i)$$

#### Regular or context-free grammars

$$\Pr(y) \approx Pr(y) = \sum_{d(y)} p_G(d(y)) \approx \max_{d(y)} p_G(d(y))$$

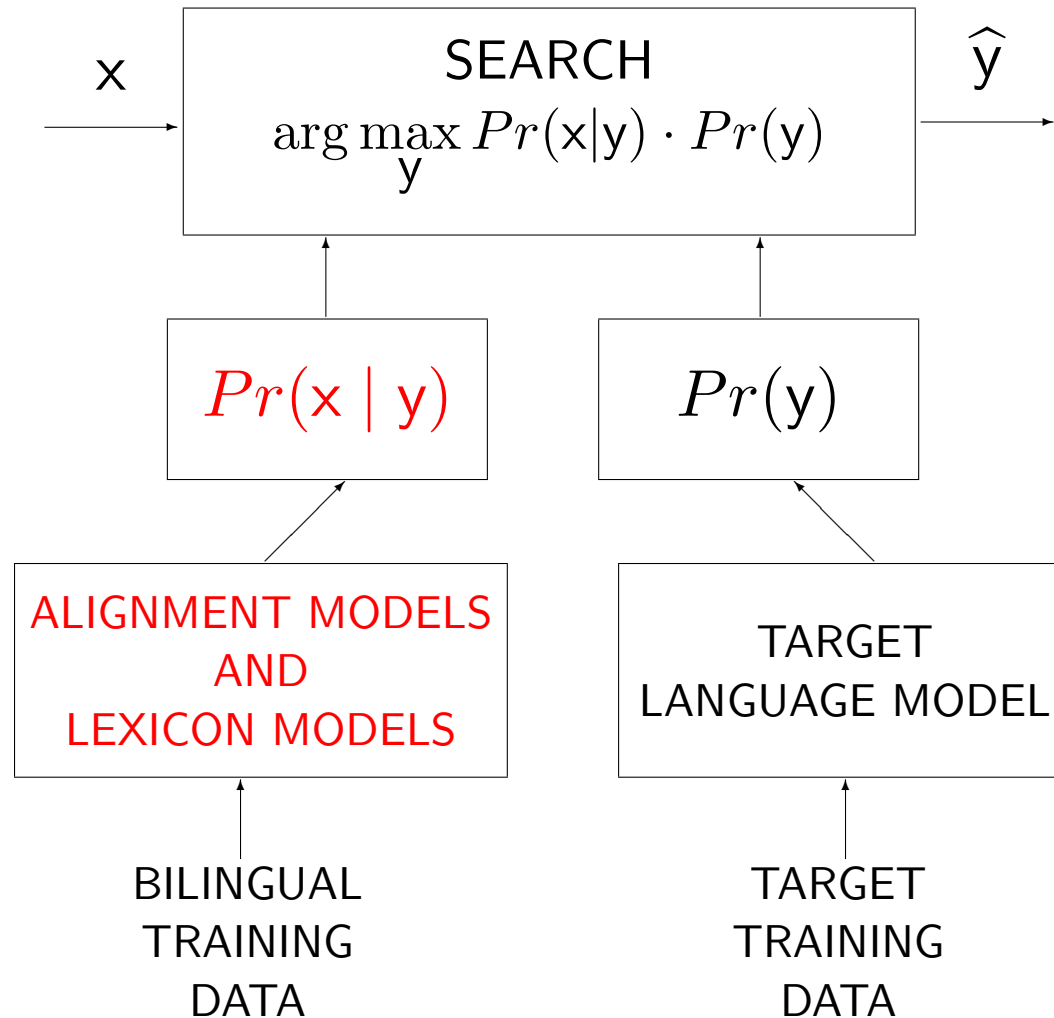
### Learning language models

- Probabilistic estimation techniques.
  - Maximum likelihood
  - Maximum entropy.
- SMOOTHING.
- Extensions: cache, triggers, categories, etc.
- Widely used toolkits for  $n$ -grams:
  - SRILM - The SRI Language Modeling Toolkit  
<http://www.speech.sri.com/projects/srilm/>
  - The CMU Statistical Language Modeling (SLM) Toolkit  
[http://www.speech.cs.cmu.edu/SLM\\_info.html](http://www.speech.cs.cmu.edu/SLM_info.html)

## 2.2 ALIGNMENTS

---

### An inverse approach



## 2.2 ALIGNMENTS

### Example of word alignments

.	.	.	.	.	.	.	.	.	.	.	.	■
Cabedo	.	.	.	.	.	.	.	.	.	.	■	.
Rosario	.	.	.	.	.	.	.	.	.	■	.	.
de	.	.	.	.	.	.	.	.	■	.	.	.
nombre	.	.	.	.	.	.	.	.	■	.	.	.
a	.	.	.	.	.	.	.	.	■	.	.	.
tranquila	.	.	.	.	.	.	■	.	.	.	.	.
habitación	.	.	.	.	.	.	.	■	.	.	.	.
una	.	.	.	.	.	■	.	.	.	.	.	.
de	.	.	.	.	■	.	.	.	.	.	.	.
reserva	.	.	.	■	.	.	.	.	.	.	.	.
la	.	.	■	.	.	.	.	.	.	.	.	.
hecho	.	■	.	.	.	.	.	.	.	.	.	.
he	■	.	.	.	.	.	.	.	.	.	.	.
		have	made	a	reservation	for	a	quiet	room	for	Rosario	Cabedo

## 2.2 ALIGNMENTS

### Example of word alignments [Ney 03a]

?	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	■		
proposal	.	.	.	.	■	.	.	.	.	.	.	.	.	.	.	.	.	.	.		
new	.	.	.	■	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.		
the	.	.	■	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.		
under	■	■	■	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.		
fees	.	.	.	.	.	.	.	.	.	.	.	.	.	.	■	■	■	.	.		
collecting	.	.	.	.	.	.	.	.	.	.	.	.	■	■	.	.	.	.	.		
and	.	.	.	.	.	.	.	.	.	.	.	■	.	.	.	.	.	.	.		
administering	.	.	.	.	.	.	.	.	.	.	■	.	.	.	.	.	.	.	.		
of	.	.	.	.	.	.	.	.	.	■	.	.	.	.	.	.	.	.	.		
cost	.	.	.	.	.	.	.	.	■	.	.	.	.	.	.	.	.	.	.		
anticipated	.	.	.	.	.	.	.	.	■	.	.	.	.	.	.	.	.	.	.		
the	.	.	.	.	.	.	.	■	.	.	.	.	.	.	.	.	.	.	.		
is	.	.	.	.	.	.	.	■	.	.	.	.	.	.	.	.	.	.	.		
What	.	.	.	.	.	.	■	■	.	.	.	.	.	.	.	.	.	.	.		
	En	vertu	de	les	nouvelles	propositions	,	quel	est	le	cout	prevu	de	administration	et	de	perception	de	les	droits	?



## 2.2 ALIGNMENTS

---

➤ **Alignments** [Brown 90]:  $J = |x|$  y  $I = |y|$

$$a \subseteq \{1, \dots, J\} \times \{1, \dots, I\}$$

- Number of connections:  $I \cdot J$
- Number of alignments:  $2^{I \cdot J}$

➤ Constrain:  $a : \{1, \dots, J\} \rightarrow \{0, \dots, I\}$ , ( $a_j = 0 \Rightarrow j$  in  $x$  is not aligned with any position in  $y$ ).

- Number of alignments:  $(I + 1)^J$

➤ Set of possible alignments:  $\mathcal{A}(x, y)$

➤ The probability of translation  $y$  to  $x$  through an alignment  $a$  is  $\Pr(x, a \mid y)$

$$\Pr(x \mid y) = \sum_{a \in \mathcal{A}(y, x)} \Pr(x, a \mid y)$$

## 2.2 ALIGNMENTS

---

$$\begin{aligned}\Pr(\mathbf{x}, \mathbf{a} \mid \mathbf{y}) &= \Pr(J \mid \mathbf{y}) \cdot \Pr(\mathbf{x}, \mathbf{a} \mid J, \mathbf{y}) \\ &= \Pr(J \mid \mathbf{y}) \cdot \Pr(\mathbf{a} \mid J, \mathbf{y}) \cdot \Pr(\mathbf{x} \mid \mathbf{a}, J, \mathbf{y})\end{aligned}$$

- Length probability:  $\Pr(J \mid \mathbf{y})$
- **Alignment probability:**  $\Pr(\mathbf{a} \mid J, \mathbf{y})$
- **Lexicon probability:**  $\Pr(\mathbf{x} \mid \mathbf{a}, J, \mathbf{y})$

$$\Pr(\mathbf{a} \mid J, \mathbf{y}) = \prod_{j=1}^J \Pr(a_j \mid \mathbf{a}_1^{j-1}, J, \mathbf{y}) \qquad \Pr(\mathbf{x} \mid \mathbf{a}, J, \mathbf{y}) = \prod_{j=1}^J \Pr(x_j \mid \mathbf{x}_1^{j-1}, \mathbf{a}, J, \mathbf{y})$$

$$\Pr(\mathbf{x}, \mathbf{a} \mid \mathbf{y}) = \Pr(J \mid \mathbf{y}) \cdot \prod_{j=1}^J \Pr(a_j \mid \mathbf{a}_1^{j-1}, \mathbf{x}_1^{j-1}, J, \mathbf{y}) \cdot \Pr(x_j \mid \mathbf{a}_1^j, \mathbf{x}_1^{j-1}, J, \mathbf{y})$$

### Zero-order models

- Model 1
- Model 2
- The Viterbi approximation
- The search problem

## 2.3 STATISTICAL ALIGNMENTS MODELS

---

### Model 1

$$\Pr(\mathbf{x}, \mathbf{a} \mid \mathbf{y}) = \Pr(J \mid \mathbf{y}) \cdot \prod_{j=1}^J \Pr(\mathbf{a}_j \mid \mathbf{a}_1^{j-1}, \mathbf{x}_1^{j-1}, J, \mathbf{y}) \cdot \Pr(\mathbf{x}_j \mid \mathbf{a}_1^j, \mathbf{x}_1^{j-1}, J, \mathbf{y})$$

- $\Pr(J \mid \mathbf{y}) \approx n(J|I)$
- $\Pr(\mathbf{a}_j \mid \mathbf{a}_1^{j-1}, \mathbf{x}_1^{j-1}, J, \mathbf{y}) \approx \frac{1}{(I+1)^J}$
- $\Pr(\mathbf{x}_j \mid \mathbf{a}_1^j, \mathbf{x}_1^{j-1}, J, \mathbf{y}) \approx l(\mathbf{x}_j \mid y_{\mathbf{a}_j})$

$l(\mathbf{x}_j \mid y_i)$  defines a **statistical lexicon**

$$\Pr(\mathbf{x} \mid \mathbf{y}) \approx P_{M1}(\mathbf{x} \mid \mathbf{y}) = \frac{n(J|I)}{(I+1)^J} \prod_{j=1}^J \sum_{i=0}^I l(\mathbf{x}_j \mid y_i)$$

### Model 1

- $\Pr(J | y) \approx n(J|I)$
- $\Pr(a_j | a_1^{j-1}, x_1^{j-1}, J, y) \approx \frac{1}{(I+1)^J}$
- $\Pr(x_j | a_1^j, x_1^{j-1}, J, y) \approx l(x_j | y_{a_j})$

Generative process: Given a target sentence  $y$  of length  $I$ ,

1. Choose the length of the source sentence  $J$  according to  $n(J|I)$
2. For each  $1 \leq j \leq J$ , choose a position  $a_j$  in the target sentence according to a uniform distribution.
3. For each  $1 \leq j \leq J$  choose a source word  $x_j$  according to  $l(x_j | y_{a_j})$

## 2.3 STATISTICAL ALIGNMENTS MODELS

---

### Model 1: An example

Given $y$ :	$a$	$double$	$room$	$(I = 3)$	
Choose $J$ ( $n(J   3)$ ): ( $J = 5$ )	1	2	3	4	5
Choose $a_j$ (uniform)	1	3	2	2	2
	$a$	room	double	double	double
Choose $x_j$ ( $l(x_j   y_i)$ )	Una	habitación	con	dos	camas

## 2.3 STATISTICAL ALIGNMENTS MODELS

---

### Model 2

$$\Pr(\mathbf{x}, \mathbf{a} \mid \mathbf{y}) = \Pr(J \mid \mathbf{y}) \cdot \prod_{j=1}^J \Pr(\mathbf{a}_j \mid \mathbf{a}_1^{j-1}, \mathbf{x}_1^{j-1}, J, \mathbf{y}) \cdot \Pr(\mathbf{x}_j \mid \mathbf{a}_1^j, \mathbf{x}_1^{j-1}, J, \mathbf{y})$$

- $\Pr(J \mid \mathbf{y}) \approx n(J|I)$
- $\Pr(\mathbf{a}_j \mid \mathbf{a}_1^{j-1}, \mathbf{x}_1^{j-1}, J, \mathbf{y}) \approx a(\mathbf{a}_j \mid j, J, I)$
- $\Pr(\mathbf{x}_j \mid \mathbf{a}_1^j, \mathbf{x}_1^{j-1}, J, \mathbf{y}) \approx l(\mathbf{x}_j \mid y_{\mathbf{a}_j})$

$l(\mathbf{x}_j \mid y_i)$  defines a **statistical lexicon**

$a(i \mid j, J, I)$  defines **statistical alignments**

$$\Pr(\mathbf{x} \mid \mathbf{y}) \approx P_{M2}(\mathbf{x} \mid \mathbf{y}) = n(J|I) \cdot \prod_{j=1}^J \sum_{i=0}^I a(i \mid j, J, I) \cdot l(\mathbf{x}_j \mid y_i)$$

## 2.3 STATISTICAL ALIGNMENTS MODELS

---

### Model 2

- $\Pr(J | y) \approx n(J|I)$
- $\Pr(a_j | a_1^{j-1}, x_1^{j-1}, J, y) \approx a(a_j | j, J, I)$
- $\Pr(x_j | a_1^j, x_1^{j-1}, J, y) \approx l(x_j | y_{a_j})$

Generative process: Given a target sentence  $y$  of length  $I$ ,

1. Choose the length of the source sentence  $J$  according to  $n(J|I)$ .
2. For each  $1 \leq j \leq J$ , choose a position  $a_j$  in the target sentence according to  $a(a_j | j, J, I)$ .
3. For each  $1 \leq j \leq J$  choose a source word  $x_j$  according to  $l(x_j | y_{a_j})$ .



## 2.3 STATISTICAL ALIGNMENTS MODELS

---

### Model 2: An example

Given $y$ :	$a$	$double$	$room$	$(I = 3)$	
Choose $J$ ( $n(J   3)$ ): ( $J = 5$ )	1	2	3	4	5
Choose $a_j$ ( $a(a_j   j, I, J)$ )	1	3	2	2	2
	$a$	room	double	double	double
Choose $x_j$ ( $l(x_j   y_i)$ )	Una	habitación	con	dos	camas

The translation process: searching

$$\arg \max_y Pr(x | y) \cdot Pr(y)$$

**A computational difficult problem [Knight 99]**

ALGORITHMIC SOLUTIONS:

- Dynamic Programming like [Ney 00a]
- Stack-Decoding:  $A^*$  or Branch & Bound [Brown 90]

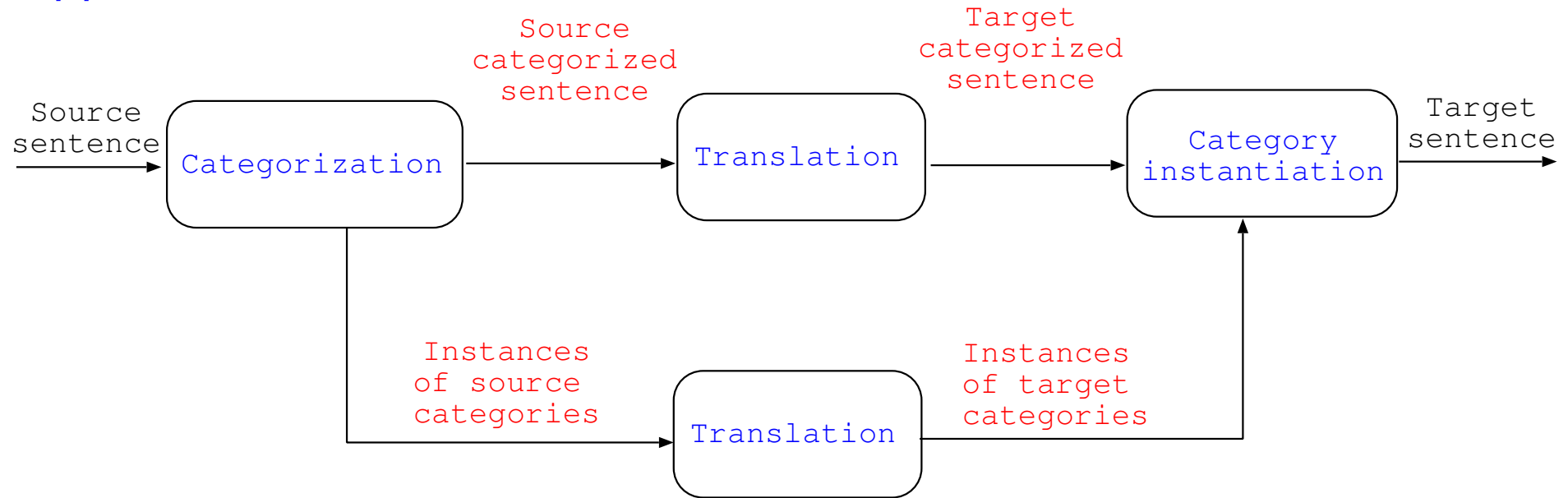
## 2.4 CATEGORIZATION IN MT

---

- Too many parameters to be estimated
- Many words play the same role: names, dates, etc.
- Substitution of words by categories:
  - The vocabulary size decreases.
  - Easy word addition to the vocabulary.
- Examples:
  - `mi nombre es $NAME.masc $SURNAME . # my name is $NAME.masc $SURNAME .`
  - `nos vamos a ir el $DATE a $HOUR . # we are leaving on $DATE at $HOUR .`
- Given a bilingual corpus:
  - Automatic extraction of bilingual categories.
  - Manual extraction of bilingual categories.

## 2.4 CATEGORIZATION IN MT

### An approach



1. **CATEGORIZATION:** Translate the source sentence into an source categorized sentence and obtain the source instances of each category.
2. **CATEGORIZED TRANSLATION:** Translate the source categorized sentence into a target categorized sentence.
3. **TRANSLATION OF EACH CATEGORY:** Translate the source instances of each category detected.
4. **CATEGORY RESOLUTION:** Substitution of each target category by the corresponding instance translation.

## 2.4 CATEGORIZATION IN MT

---

### An example

me voy a ir el dia veintiseis de abril a las doce en punto de la mañana

*Statistical  
Categorization*

*Viterbi  
Alignment*

me voy a ir el dia \$DATE a \$HOUR de la mañana

\$DATE = veintiseis de abril  
\$HOUR = las doce en punto

*Statistical Translation*

I am leaving on \$DATE at \$HOUR in the morning

\$DATE = April the twenty-sixth  
\$HOUR = twelve o'clock

*Category Resolution*

I am leaving on April the twenty-sixth at twelve o'clock in the morning

### Automatic categorization

➤ *Extended word categories* [Barrachina 99]

1. Align a bilingual corpus
2. Build extended words using the alignments
3. Apply a clustering algorithm to the corpus of extended word sentences

➤ *Statistical bilingual categories* [Och 99]

1. Align a bilingual corpus
2. Apply a clustering algorithm to the target corpus.
3. Apply a clustering algorithm to the source corpus taking into account the categories of target words aligned to the source words.

## Index

### 1. Introduction

- 1.1 Objectives of MT
- 1.2 Approaches to MT
- 1.3 Linguistic resources
- 1.4 Assessment

### 2. Statistical alignment models

- 2.1 Statistical framework to MT
- 2.2 Alignments
- 2.3 Statistical alignment models
- 2.4 Categorization in MT

### 3. **Advanced statistical alignment models**

- 3.1 Fertility-based models
- 3.2 The search problem
- 3.3 Using linguistic knowledge

### 4. Phrase-based models

- 4.1 Beyond word models
- 4.2 Phrase-based models

### 5. Syntax-based translation models

- 5.1 Introduction
- 5.2 ITG for MT
- 5.3 Tree-to-string models
- 5.4 Hierarchical MT

### Alignments

$$\Pr(\mathbf{x} \mid \mathbf{y}) = \sum_{\mathbf{a} \in \mathcal{A}(\mathbf{y}, \mathbf{x})} \Pr(\mathbf{x}, \mathbf{a} \mid \mathbf{y}) = \Pr(J \mid \mathbf{y}) \cdot \sum_{\mathbf{a} \in \mathcal{A}(\mathbf{y}, \mathbf{x})} \Pr(\mathbf{x}, \mathbf{a} \mid J, \mathbf{y})$$

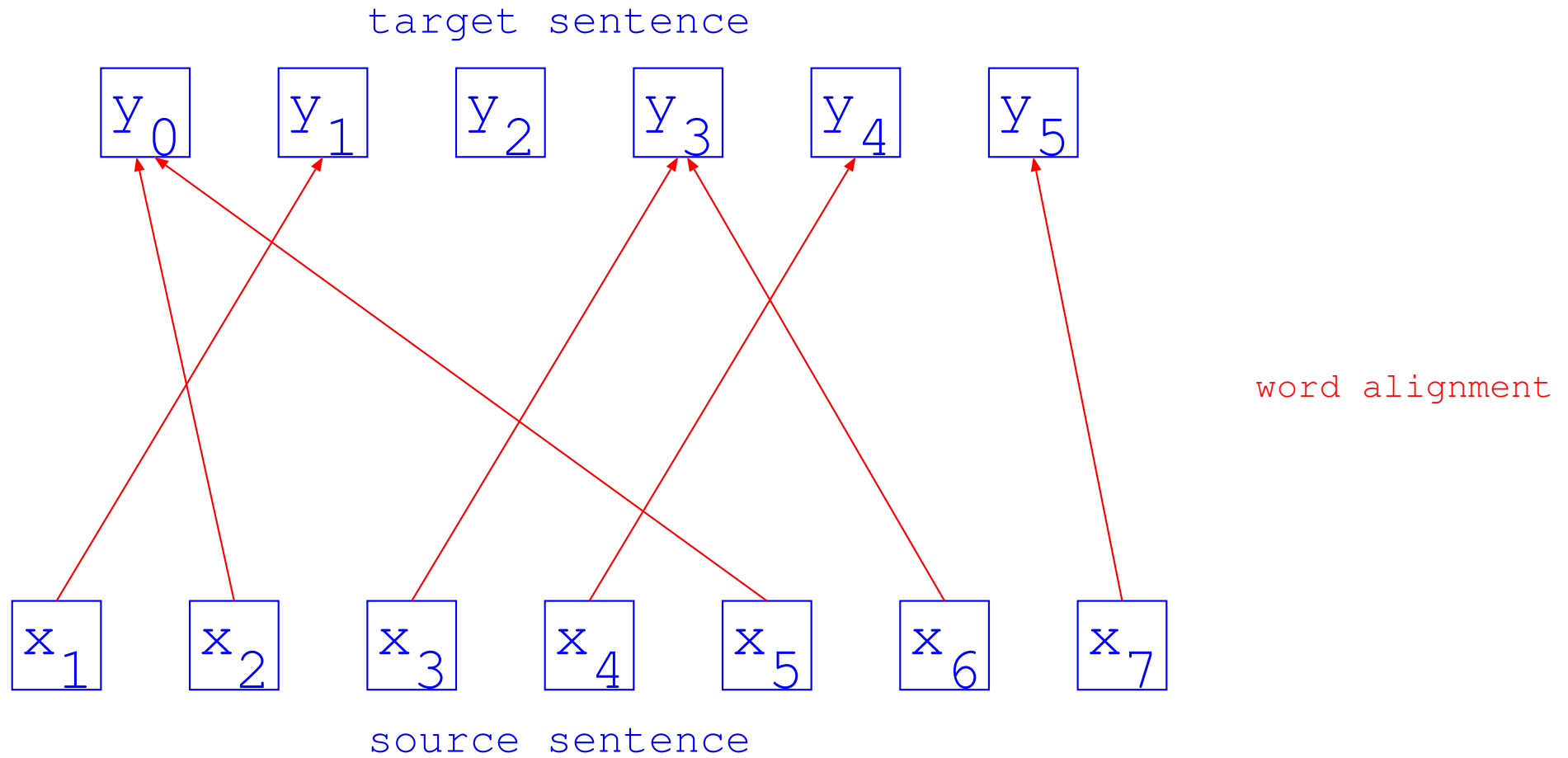
Alignment probabilities and lexicon probabilities

- Model 1
- Model 2



# 3.1 FERTILITY-BASED MODELS

## Models 1, 2 or HMM



## 3.1 FERTILITY-BASED MODELS

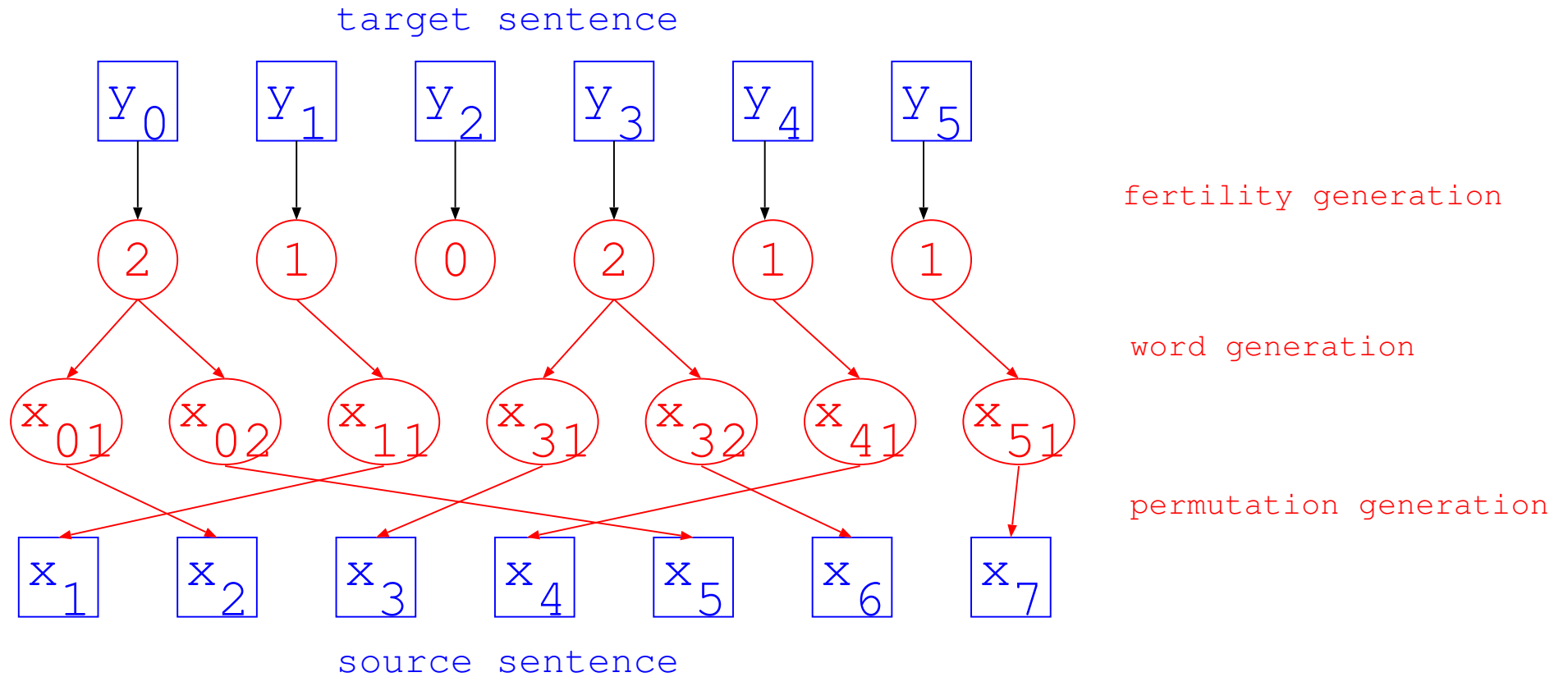
---

### Models 3, 4 and 5

- Model 3: Lexicon, fertility and distortion models
- Model 4 is a refined version of distortion distribution in Model 3
- Model 5 is a consistent version of distortion distribution in Model 4

# 3.1 FERTILITY-BASED MODELS

## Fertility



## 3.1 FERTILITY-BASED MODELS

---

**Fertility**  $\phi$  of  $y_i \in \Delta$ : number of the source words connected to an target word  $y_i$

1. Choose how many source words are connected to a target word  $y_i$ : *fertility* of  $y_i$
2. Choose a set of the source words, a *tablet*  $\tau_i$ , that is connected to  $i$ -th target word
3. Choose the *position*  $\pi_{i,k}$  in the source sentence of the  $k$ -th word  $\tau_{i,k}$  that is connected to the  $i$ -th target word

## 3.1 FERTILITY-BASED MODELS

---

### Model 3

Given a target sentence  $y$  of length  $I$ :

1. For each  $1 \leq i \leq I$  choose a length  $\phi_i$
2. Choose a length  $\phi_0$
3.  $J = \sum_{i=0}^I \phi_i$ .
4. For each  $1 \leq i \leq I$  and  $1 \leq k \leq \phi_i$ , choose a source word
5. For each  $1 \leq i \leq I$  and  $1 \leq k \leq \phi_i$ , choose a position
6. If any position has been chosen then **error** (*inconsistent model*).
7. For each  $1 \leq k \leq \phi_0$  choose a position from the vacant positions according to a uniform distribution.

## 3.1 FERTILITY-BASED MODELS

---

### Example

Given $y$ :	$a$	<i>double</i>	<i>room</i>	$(I = 3)$	
$i$	1	2	3		
Choose $\phi(y_i) = \phi$ using $f(\phi y_i)$	1	3	1		
Choose $\tau_{i,k} = x$ using $l(x y_i)$	{una}	{con,	camas, dos}	{habitación}	
Choose $\pi_{i,k} = j$ using $d(j i, I, J)$	1	3	5	4	2
$j$	1	2	3	4	5
x	una	habitación	con	dos	camas

## 3.1 FERTILITY-BASED MODELS

---

### Examples of alignments

#### Corpus EUTRANS-I: Spanish-English

1	2	3	4	5	6	7	8	9	10
por	favor	,	¿	podría	ver	alguna	habitación	tranquila	?

- MODEL 1, ITERATION 5  
could (5) | (6) see (6) a (7) quiet (9) room (8) , (3) please (2) ? (4)
- MODEL 2, ITERATION 2  
could (5) | (6) see (6) a (7) quiet (9) room (8) , (3) please (3) ? (10)
- MODEL 3, ITERATION 2  
could (5) | (5) see (6) a (7) quiet (9) room (8) , (3) please (2) ? (10)

### Conventional IBM Models Training

- Every model has a specific set of free parameters.
- To train the model parameters  $\theta$ : A maximum likelihood criterium, using a parallel training corpus consisting of  $S$  sentence pairs  $\{(x^{(n)}, y^{(n)}) : n = 1, \dots, N\}$ :

$$\hat{\theta} = \arg \max_{\theta} \prod_{n=1}^N \sum_{\mathbf{a}} p_{\theta}(x^{(n)}, \mathbf{a} | y^{(n)}) \quad .$$

- The training is carried out using the Expectation-Maximization (EM) algorithm.



## 3.2 THE SEARCH PROBLEM

---

$$\hat{y} = \arg \max_y Pr(x | y) \cdot Pr(y)$$

- Search is a **NP-Hard** problem. [Knight 99]
- **Algorithmic solutions:** (+ heuristics for efficient suboptimal solutions)
  - *Dynamic Programming* [Tillmann 03]
  - *Stack-decoding, A\* or Branch & Bound* (Ortiz , 2003)

### Some stack-decoding proposals

- Candide systems from IBM [Berger et al. 96]: Multiple stacks, model 3.
- Multiple stack-decoding [Wang and Waibel 98]: Model 2.
- Algorithm  $A^*$  [Ueffing et al. 01]: model 4.
- Basic stack-decoding strategy:
  - Origin of the *stack decoding* or  $A^*$ : ASR
  - Optimal solution to the search problem (Jelinek, 1976)
  - Incremental development of practical hypothesis
  - The hypothesis are stored in a priority queue (a type of 'stack')
  - Selection and expansion of the top of the stack(s).

### A taxonomy of the stack-decoding algorithms

- Basic stack-decoding algorithm:
  - All the hypothesis are stored in a one stack
  - A hypothesis is selected in each iteration: the hypothesis with higher score in the stack
- Problem: hypothesis with a high number of aligned words are discarded.
- Possible solutions:
  - Use of heuristics: an estimation of the contribution to the set of the optimal score.
  - Multiple stacks.
- Taxonomy:
  - Single stack algorithms  $A^*$
  - Multiple stack algorithms

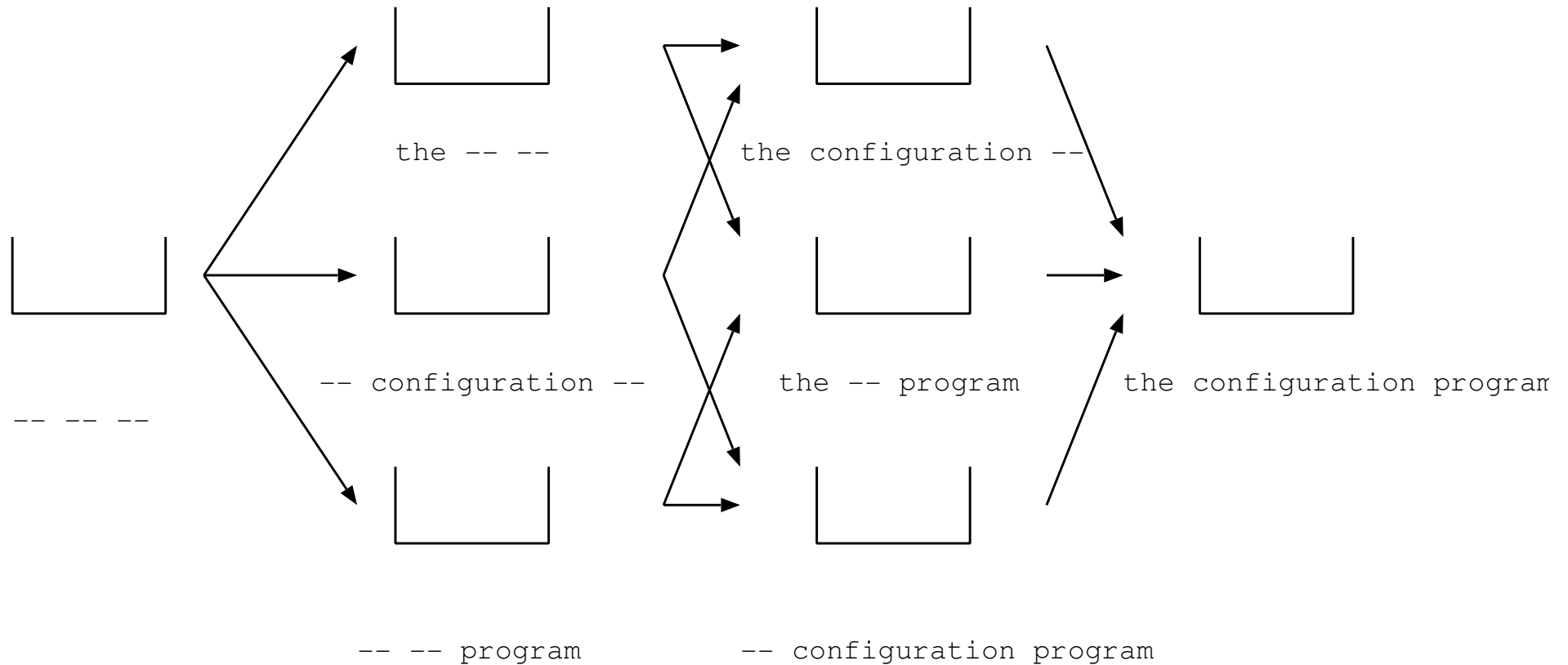
### Basic multiple stack decoding *StackDecoding*

- A hypothesis in a stack:
  - A prefix of the target sentence ( $y_1^i$ )
  - A coverage subset of source positions ( $\mathcal{C}$ )
  - A score ( $S$ ).
- There is one stack for each possible subset of source positions which words has already been translated.
- The possible number of stacks can be very high.
- In each iteration, the best hypothesis from each available stack is selected to generate new extended hypothesis.
- The new hypothesis is stored in the corresponding stack.

## 3.2 THE SEARCH PROBLEM

---

Source sentence: "the configuration program"



## 3.3 USING LINGUISTIC KNOWLEDGE

---

Is the linguistic knowledge needed for statistical machine translation?

➤ YES?

- There are many linguistic knowledge available.
- The bilingual training data can be better exploited.

➤ NOT?

- Many linguistic knowledge is hard to formalize.
- The generation of new linguistic knowledge requires great human effort.

## 3.3 USING LINGUISTIC KNOWLEDGE

---

### Linguistic knowledge that has been used in statistical machine translation

- Morpho-syntactic knowledge: lexicon, Part-of-Speech, etc...  
(Nießen and Ney, 2004)

Hybrid linguistic-statistical approaches have been used with success (i.e. *hidden markov models*)

- Others: Cognates (Kondrak, Marcu and Knight, 2003), named entities (Huang, Vogel and Waibel, 2003), ...
- Syntactic information: next topic!

### Morpho-syntactic knowledge in statistical machine translation

Nießen and Ney, 2004. *Statistical machine translation with scarce resources using morpho-syntactic information*. Computational Linguistics.

- Present statistical machine translation systems often treat different inflected forms of the same lemma as if they were independent of one another.
- The bilingual data can be better exploited by explicitly taking into account the interdependencies of related inflected forms.



## 3.3 USING LINGUISTIC KNOWLEDGE

---

### Morpho-syntactic knowledge in statistical machine translation

yo **como** pan

- Morphological and syntactic tags (POS, tense, person, ...)
- The *base form*

*comer* verb indicative present singular 1

## Index

### 1. Introduction

- 1.1 Objectives of MT
- 1.2 Approaches to MT
- 1.3 Linguistic resources
- 1.4 Assessment

### 2. Statistical alignment models

- 2.1 Statistical framework to MT
- 2.2 Alignments
- 2.3 Statistical alignment models
- 2.4 Categorization in MT

### 3. Advanced statistical alignment models

- 3.1 Fertility-based models
- 3.2 The search problem
- 3.3 Using linguistic knowledge

### 4. **Phrase-based models**

- 4.1 Beyond word models
- 4.2 Phrase-based models

### 5. Syntax-based translation models

- 5.1 Introduction
- 5.2 ITG for MT
- 5.3 Tree-to-string models
- 5.4 Hierarchical MT

## 4.1 BEYOND WORD MODELS

---

- The basic assumption in the current word-based models: Each source word is generated by only one target word.
- This assumption does not correspond to the nature of natural language. In some cases, it is necessary to know the context.
- Solutions:
  - *Context-dependent dictionaries*. The basic unit is the word.
  - *Word sequences*:
    - *Alignment templates*: A sequence of source (classes of) words is aligned with a sequence of target (classes of) words. Inside the templates there are word-to-word correspondences. The basic unit is the word.
    - *Phrase-based models*: A sequence of source words is aligned with a sequence of target words. The basic unit is the phrase.



### Word sequences

The statistical dictionaries of single word pairs are substituted by statistical dictionaries of *bilingual phrases*.

Bilingual phrases are related with a bilingual segmentation.

- Problem: The generalisation capability, since only sequences of segments that have been seen in the training corpus are accepted.
- Problem: The selection of adequate bilingual phrases.

## 4.2 PHRASE-BASED MODELS

### An example

		y: could you ask for a taxi , please ?									
	y	could	you	ask	for	a	taxi	,	please	?	
	i	1	2	3	4	5	6	7	8	9=I	
Segmentation	i				$i_1$		$i_2$			$i_3$	
Translation	x		[ pídame ]				[ un taxi . ]			[ por favor , ]	
Permutation	$\alpha$		$\alpha_1 = 2$			$\alpha_2 = 3$			$\alpha_3 = 1$		
		por	favor	,		pídame		un	taxi	.	
	j	1	2	3	4	5	6	7			
Segmentation	$\gamma$			$\gamma_1$		$\gamma_2$				$\gamma_3$	

x: por favor , pídame un taxi .

### Log-linear models

Search for a target sentence with maximum *posterior* probability:

$$\hat{y} = \arg \max_y \Pr(y | x)$$

$$\hat{y} = \arg \max_y \frac{\exp \left( \sum_{k=1}^K \lambda_k h_k(x, y) \right)}{\sum_{y'} \exp \left( \sum_{k=1}^K \lambda_k h_k(x, y') \right)} = \arg \max_y \sum_{k=1}^K \lambda_k h_k(x, y)$$

- $h_1(x, y) = \log Pr(y)$ , a language model
- $h_2(x, y) = \log Pr_{PB}(y | x)$ , phrase-based models
- $h_3(x, y) = \log Pr_{PB}(x | y)$ , phrase-based inverse model
- $h_4(x, y) = \log Pr_{M1}(x | y)$ , statistical dictionaries
- $h_5(x, y) = \log Pr_{M1}(y | x)$ , statistical inverse dictionaries
- ...

### Learning phrase-based models

Given a sentence-aligned corpus  $\mathcal{T}$ :

- A word-aligned corpus is generated using the GIZA++ toolkit with  $\mathcal{T}$   
<http://www-i6.informatik.rwth-aachen.de/Colleagues/och/software/GIZA++.html>
- A set of bilingual word sequences from the word aligned corpus is extracted.
- The parameters of the phrase-model are estimated.



### Estimating the parameters

#### Estimating the parameters

By relative frequencies, for each pair of segments  $(x, y)$ :

$$p(\tilde{x} | \tilde{y}) = \frac{N(\tilde{x}, \tilde{y})}{N(\tilde{y})}$$

where  $N(\tilde{y})$  denotes the number of times that phrase  $\tilde{y}$  has appeared, and  $N(\tilde{x}, \tilde{y})$  is the number of times that the bilingual phrase  $(\tilde{x}, \tilde{y})$  has appeared.

#### Distortion model

$$p(\alpha_k | \alpha_{k-1}) = p_0^{|\gamma\alpha_k - \gamma\alpha_{k-1}|},$$

where  $p_0$  is a parameter to be adjusted using a validation set.

## Index

### 1. Introduction

- 1.1 Objectives of MT
- 1.2 Approaches to MT
- 1.3 Linguistic resources
- 1.4 Assessment

### 2. Statistical alignment models

- 2.1 Statistical framework to MT
- 2.2 Alignments
- 2.3 Statistical alignment models
- 2.4 Categorization in MT

### 3. Advanced statistical alignment models

- 3.1 Fertility-based models
- 3.2 The search problem
- 3.3 Using linguistic knowledge

### 4. Phrase-based models

- 4.1 Beyond word models
- 4.2 Phrase-based models

### 5. **Syntax-based translation models**

- 5.1 Introduction
- 5.2 ITG for MT
- 5.3 Tree-to-string models
- 5.4 Hierarchical MT



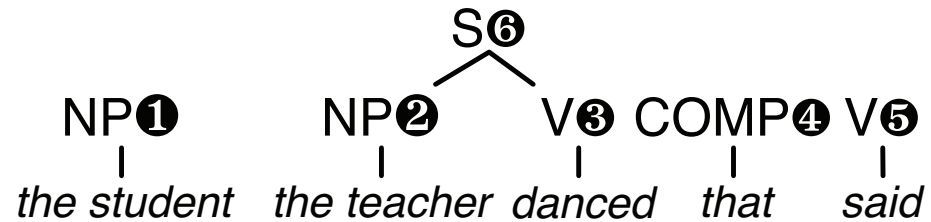
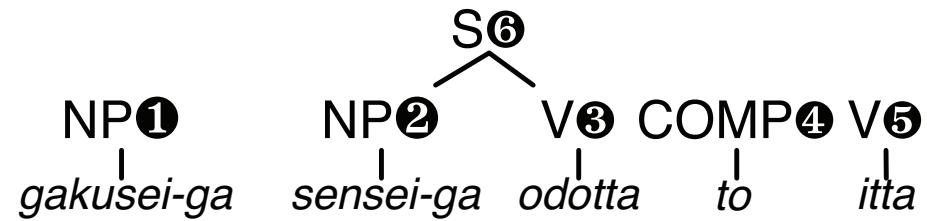
### Example SCFG\*

	Japanese	English
S →	NP① VP②	NP① VP②
S' →	S① COMP②	COMP② S①
VP →	NP① V②	V② NP①
NP →	<i>gakusei-ga</i>	<i>student</i>
NP →	<i>sensei-ga</i>	<i>teacher</i>
V →	<i>odotta</i>	<i>danced</i>
V →	<i>itta</i>	<i>said</i>
COMP →	<i>to</i>	<i>that</i>

\* Slide source: <http://www.mt-archive.info/MTMarathon-2009-Li-ppt.pdf>

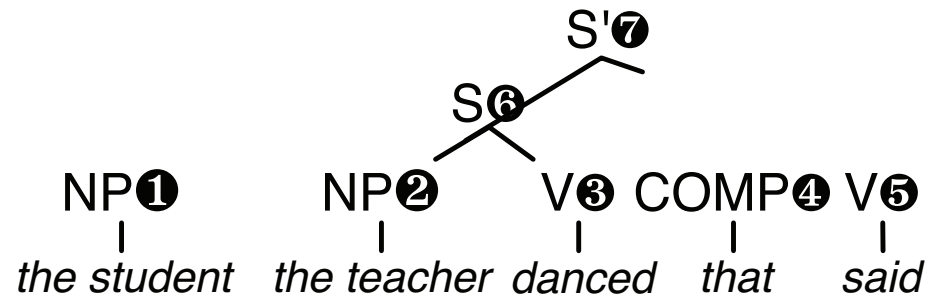
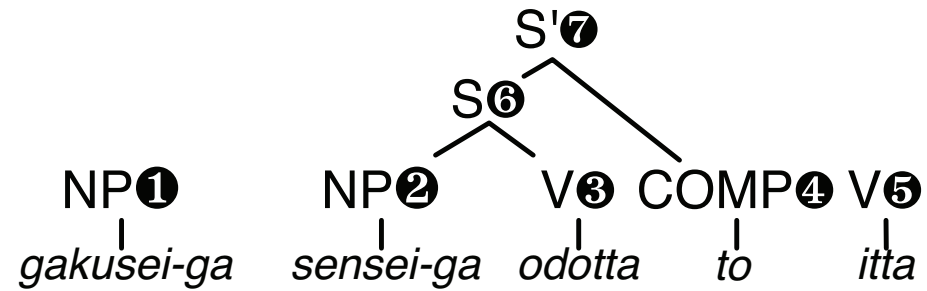
## 5.1 INTRODUCTION

---



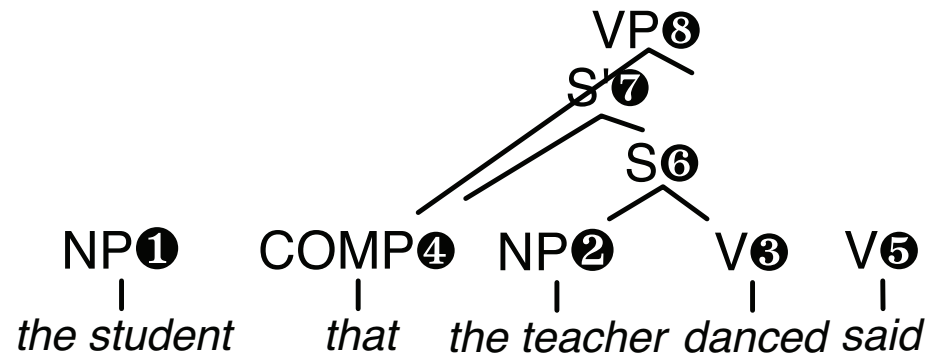
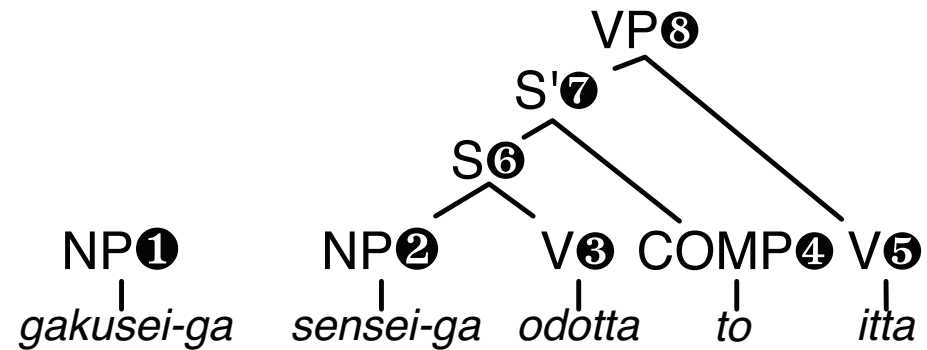
## 5.1 INTRODUCTION

---



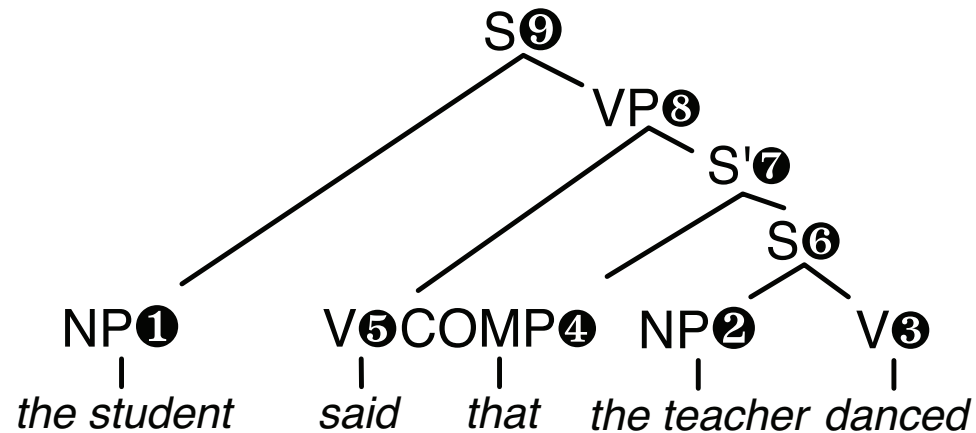
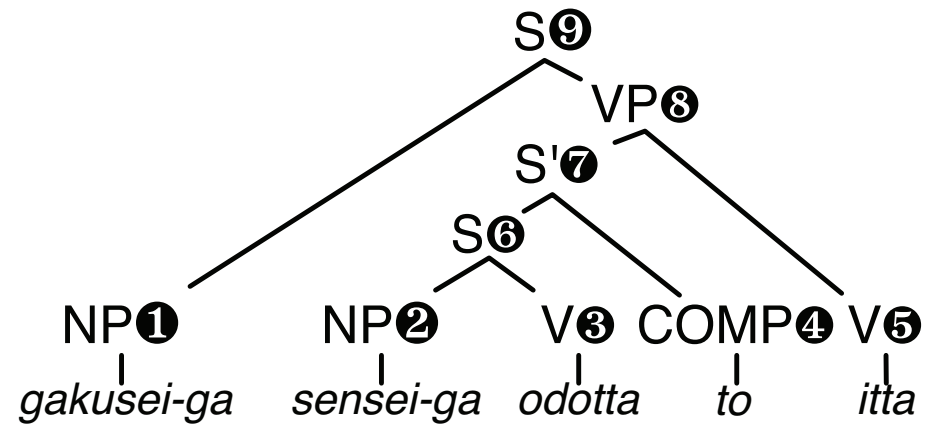
## 5.1 INTRODUCTION

---



## 5.1 INTRODUCTION

---



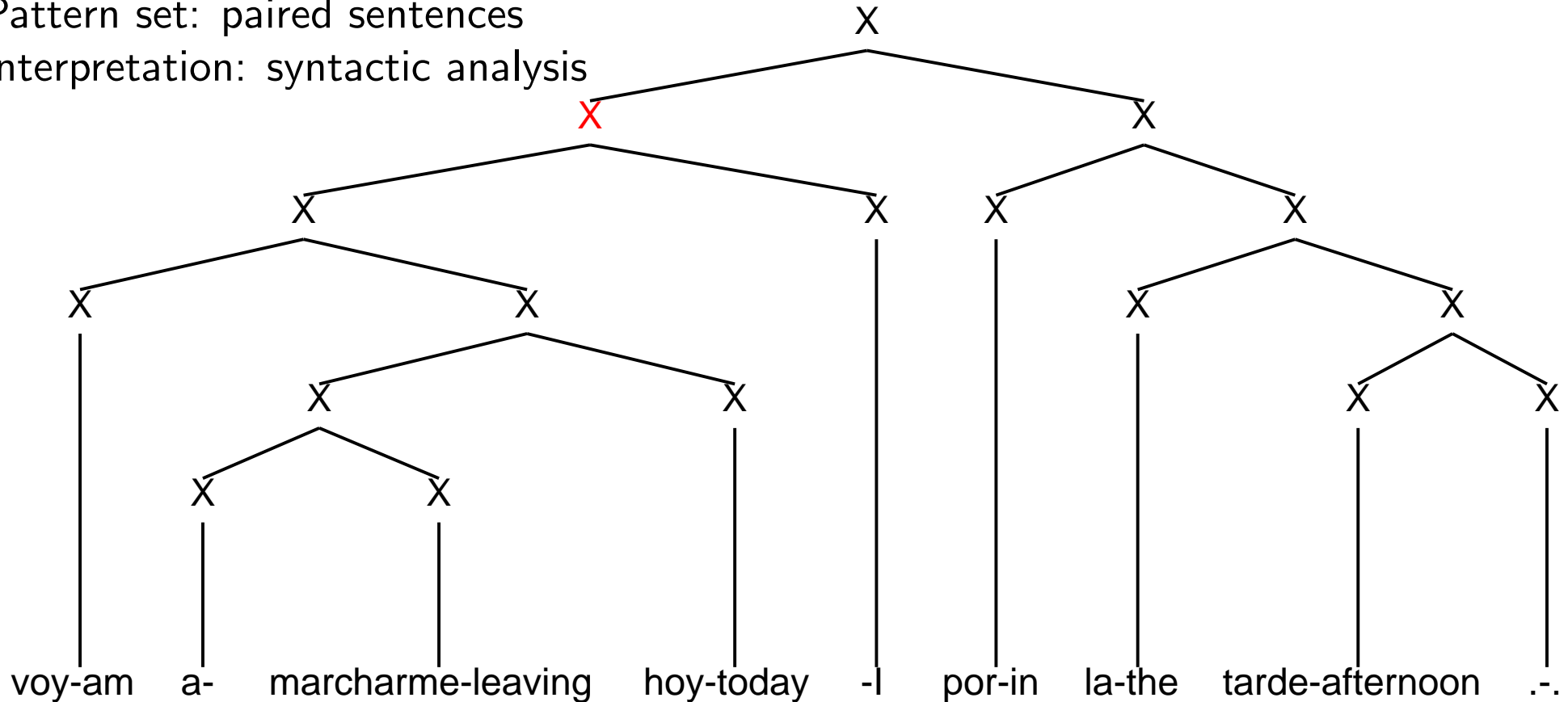
## 5.2 ITG FOR MT

### Stochastic inversion transduction grammars [Wu 97, Maryanski 79]

- Primitives: two alphabets (words, punctuation symbols, . . . )
- Object representation: two paired written sentences

“voy a marcharme hoy por la tarde”  $\iff$  “I am leaving today in the afternoon”

- Pattern set: paired sentences
- Interpretation: syntactic analysis





## 5.2 ITG FOR MT

---

➤ **ITG:**  $G = (N, W_1, W_2, R, S)$

$R$  is a finite set of straight orientation rules  $A \rightarrow [a_1 a_2 \dots a_r]$  and inverted orientation rules  $A \rightarrow \langle a_1 a_2 \dots a_r \rangle$ ,  $a_i \in N \cup X$  and  $X = (W_1 \cup \{\epsilon\}) \times (W_2 \cup \{\epsilon\})$

**Theorem.** For any ITG  $G$ , there exists an equivalent ITG  $G'$  in which every production takes one of the following forms:

$$\begin{array}{lll} S \rightarrow \epsilon/\epsilon & A \rightarrow x/\epsilon & A \rightarrow [BC] \\ S \rightarrow x/y & A \rightarrow \epsilon/y & A \rightarrow \langle BC \rangle \end{array}$$

➤ **SITG:**  $G_s = (G, p)$  where:

➤  $G$  is an ITG

➤  $p$  is a function that attaches a probability to each rule:

$$p : R \rightarrow ]0, 1] \quad \sum_{1 \leq j \leq n_i} p(A_i \rightarrow \alpha_j) = 1, \quad \forall A_i \in N$$

### Stochastic derivation for SITG

Given a sequence of stochastic events:

$$(S, S) = (\alpha_0, \beta_0) \xrightarrow{r_1} (\alpha_1, \beta_1) \xrightarrow{r_2} (\alpha_2, \beta_2) \cdots (\alpha_{m-1}, \beta_{m-1}) \xrightarrow{r_m} (\alpha_m, \beta_m) = (x, y)$$

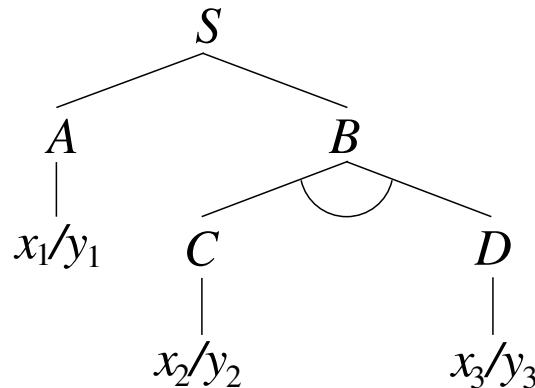
Probability of  $(x, y)$  being generated by  $G_s = (G, p)$  from the rule sequence

$d_x = (r_1, \dots, r_m)$ , is:

$$P_{G_s}((x, y), d_x) = \prod_{j=1 \dots m} p(r_j)$$

### Example

$S \rightarrow [AB]$   
 $A \rightarrow x_1/y_1$   
 $B \rightarrow \langle CD \rangle$   
 $C \rightarrow x_2/y_2$   
 $D \rightarrow x_3/y_3$



$$(S, S) \Rightarrow (AB, AB) \Rightarrow (x_1B, y_1B) \Rightarrow (x_1CD, y_1DC) \Rightarrow (x_1x_2D, y_1Dy_2) \Rightarrow (x_1x_2x_3, y_1y_3y_2)$$

Probability of a string pair

$$\Pr_{G_s}(x, y) = \sum_{d_x \in D_x} \Pr_{G_s}((x, y), d_x)$$

Probability of the best derivation

$$\widehat{\Pr}_{G_s}(x, y) = \max_{d_x \in D_x} \Pr_{G_s}((x, y), d_x)$$

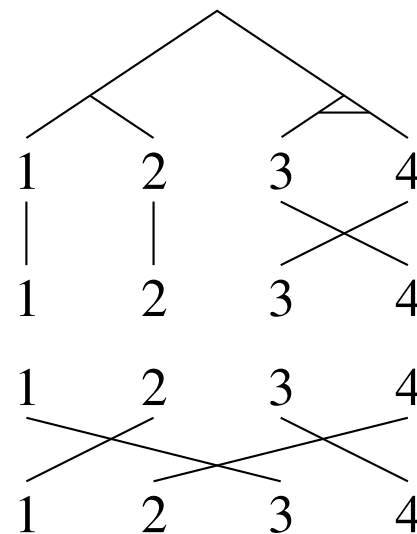
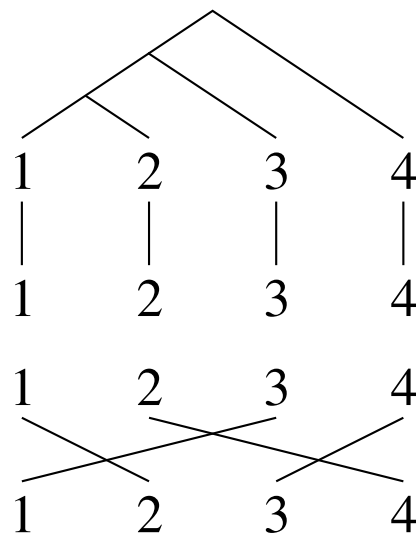
Language generated by a SITG

$$L(G_s) = \{(x, y) \mid \Pr_{G_s}(x, y) > 0\}$$

## 5.2 ITG FOR MT

### Expressiveness of ITGs

**YES**



**NO**

$r$	ITG	all matchings	ratio
1	1	1	1.000
2	2	2	1.000
3	6	6	1.000
4	22	24	0.917
5	90	120	0.750

$r$	ITG	all matchings	ratio
6	394	720	0.547
7	1,806	5,040	0.358
8	8,558	40,320	0.212
9	41,586	362,880	0.115
10	206,098	3,628,800	0.057

- Parsing:

- Inside algorithm
- Viterbi algorithm

- Learning:

- Structure learning
- Probabilistic estimation: Inside-outside estimation  
Viterbi-based estimation

- Translation:

- Adapted Cooke-Kasami-Younger parser algorithm

### Viterbi algorithm [Wu 97, Gascó 10b]

➤ Given  $(x, y) \in (W_1^*, W_2^*)$  and  $A \in N$

$$\delta_{i,j,k,l}(A) = \widehat{\text{Pr}}(A \xrightarrow{*} x_{i+1} \cdots x_j / y_{k+1} \cdots y_l)$$

➤ Initialization

$$\begin{array}{ll} \delta_{i-1,i,k-1,k}(A) = p(A \rightarrow x_i / y_k) & 1 \leq i \leq |x|, 1 \leq k \leq |y| \\ \delta_{i-1,i,k,k}(A) = p(A \rightarrow x_i / \epsilon) & 1 \leq i \leq |x|, 0 \leq k \leq |y| \\ \delta_{i,i,k-1,k}(A) = p(A \rightarrow \epsilon / y_k) & 0 \leq i \leq |x|, 1 \leq k \leq |y| \end{array}$$

## 5.2 ITG FOR MT

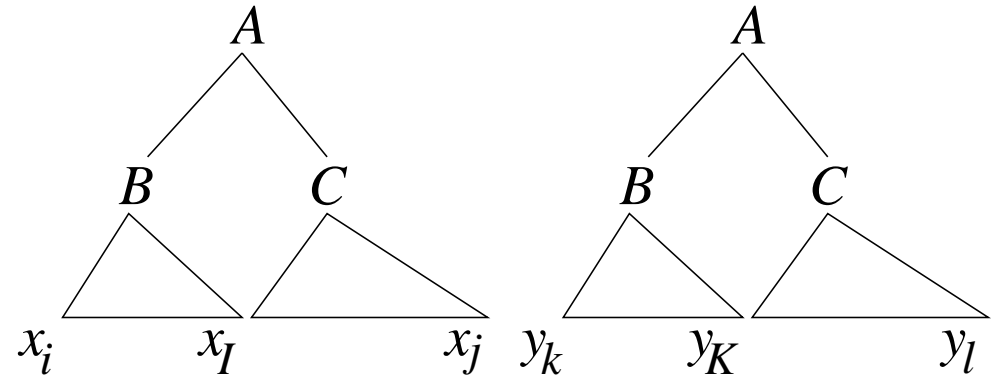
- Recursion. For all  $A \in N$ , and  $i, j, k, l$  such that  $0 \leq i < j \leq |x|$ ,  $0 \leq k < l \leq |y|$  and  $j - i + l - k \geq 2$ :

$$\delta_{ijkl}(A) = \max(\delta_{ijkl}^{\square}(A), \delta_{ijkl}^{\langle \rangle}(A))$$

$$\delta_{ijkl}^{\square}(A) = \max_{B, C \in N} p(A \rightarrow [BC]) \delta_{iIkK}(B) \delta_{IjKl}(C)$$

$$i \leq I \leq j, k \leq K \leq l$$

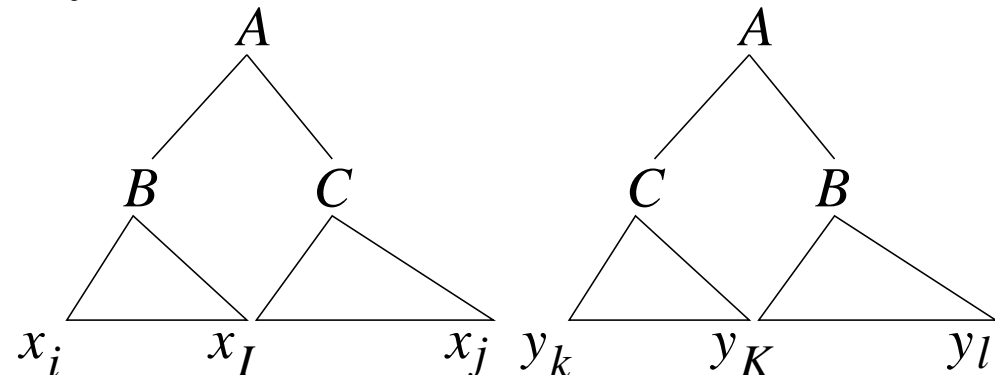
$$((j-I) + (l-K)) \times ((l-i) + (K-k)) \neq 0$$



$$\delta_{ijkl}^{\langle \rangle}(A) = \max_{B, C \in N} p(A \rightarrow \langle BC \rangle) \delta_{iIKl}(B) \delta_{IjKl}(C)$$

$$i \leq I \leq j, k \leq K \leq l$$

$$((j-I) + (K-k)) \times ((I-i) + (I-K)) \neq 0$$

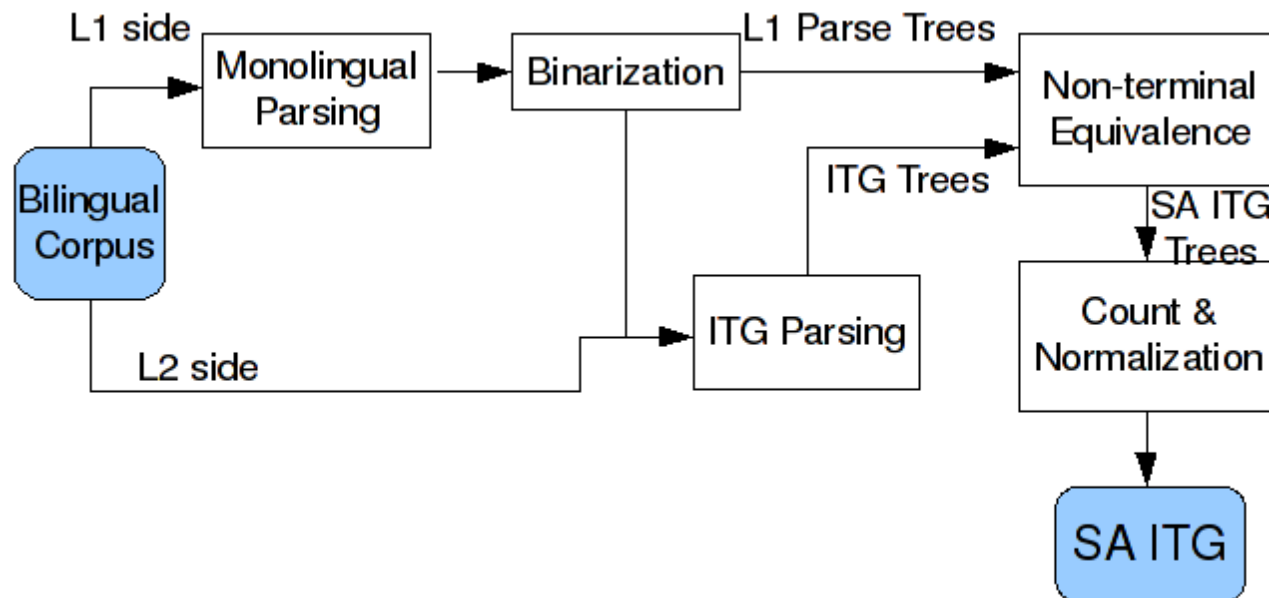


## 5.2 ITG FOR MT

---

[Gascó 10a]

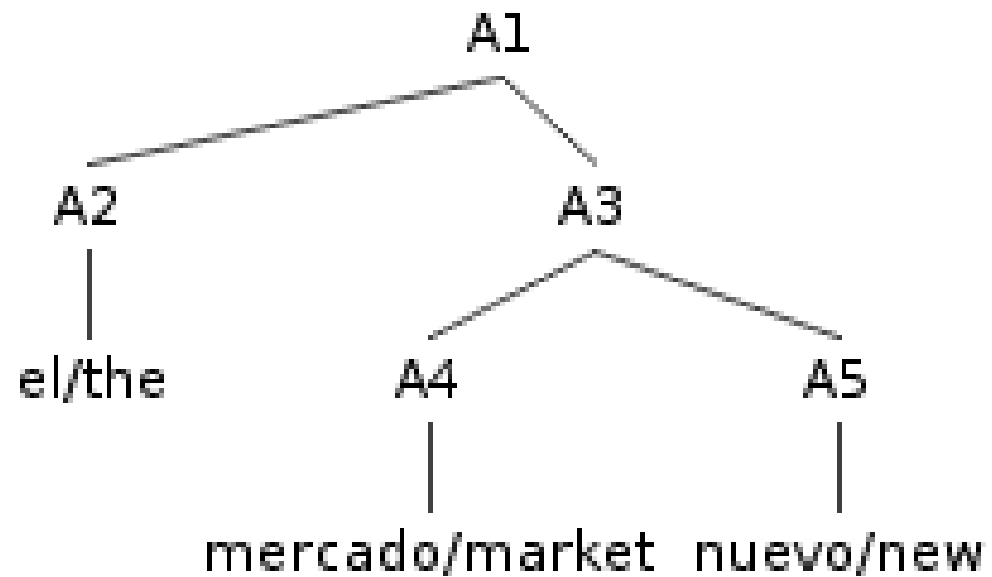
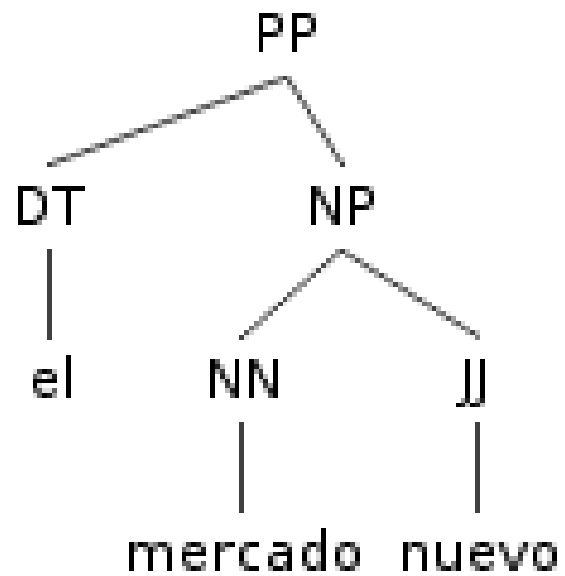
1. Create an initial SITG
2. Estimate the probabilities
3. Attach linguistic information to the non-terminal symbols





## 5.2 ITG FOR MT

---



## 5.2 ITG FOR MT

---

- IWSLT 2008 (Chinese-English BTEC)
- Standard tools: GIZA++, ZMERT
- Stanford parser for Chinese
- Baseline: Moses, 5-gram

Corpus Set	Statistic	Chinese	English
Training	Sentences	42,655	
	Words	330,163	380,431
	Voc. Size	8,773	8,387
DevSet	Sentences	489	
	Words	3,169	3,861
	OOV Words	111	115
Test	Sentences	507	
	Words	3,357	-
	OOV Words	97	-

System	%BLEU
Baseline PBT	41.1
Initial ITG	41.2
Re-estimated ITG	41.8
Source SAITG	42.9
Target SAITG	43.0

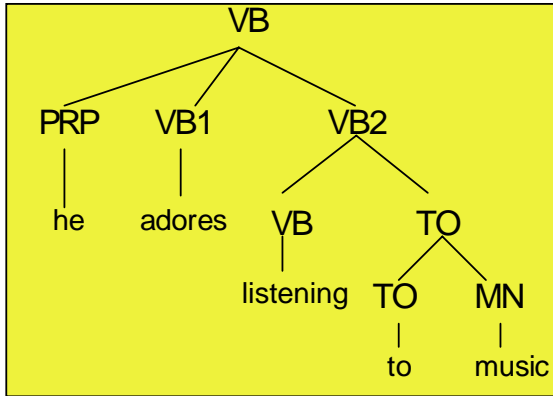
### Main ideas [Yamada 01]

- The input sentence is preprocessed by a syntactic parser
- The channel performs operations on each node of the parse tree:
  - reordering child nodes
  - inserting extra words at each node
  - translating leaf words
- The output of the the model is a string.

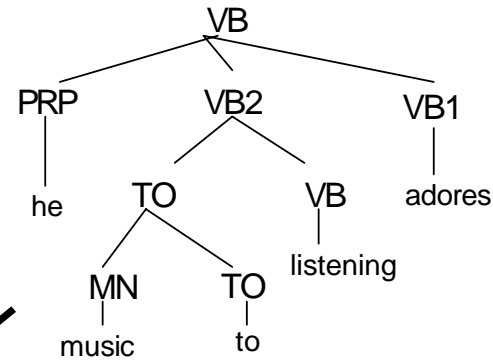
# 5.3 TREE-TO-STRING MODELS

## An example\*

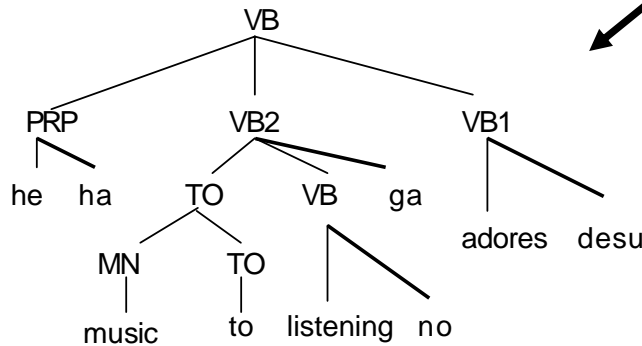
Parse Tree(E)



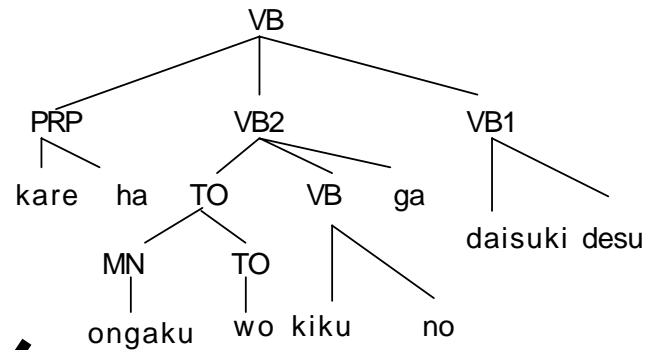
Reorder



Insert



Translate



Take Leaves

Sentence(J)

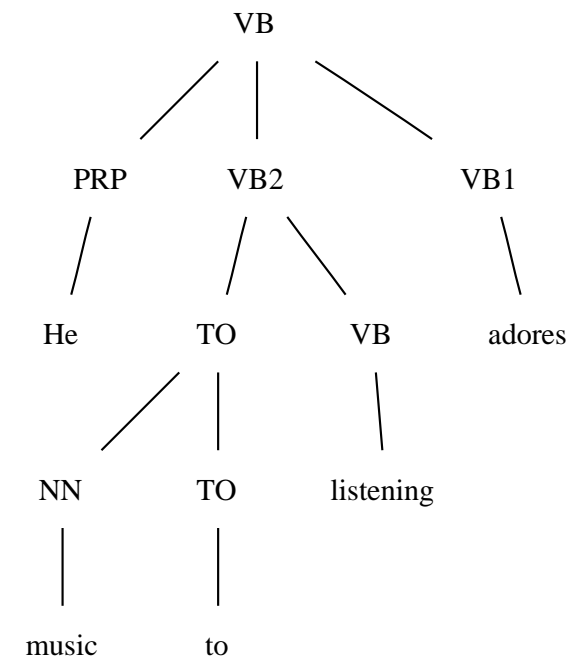
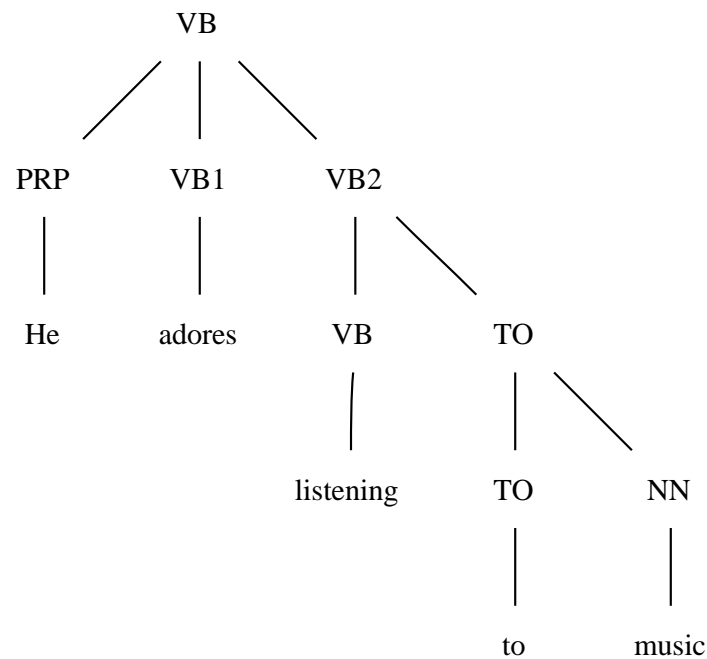
*Kare ha ongaku wo kiku no ga daisuki desu*

\*Source: <http://www.isi.edu/natural-language/people/cs562-8-22-06.pdf>

## 5.3 TREE-TO-STRING MODELS

⇒ The reordering is decided according to the *r-table*

original order	reordering	P(reorder)
PRP VB1 VB2	PRP VB1 VB2	0.074
	PRP VB2 VB1	0.723
	VB1 PRP VB2	0.061
	...	...
VB TO	VB TO	0.252
	TO VB	0.749
TO NN	TO NN	0.107
	NN TO	0.893
	...	...



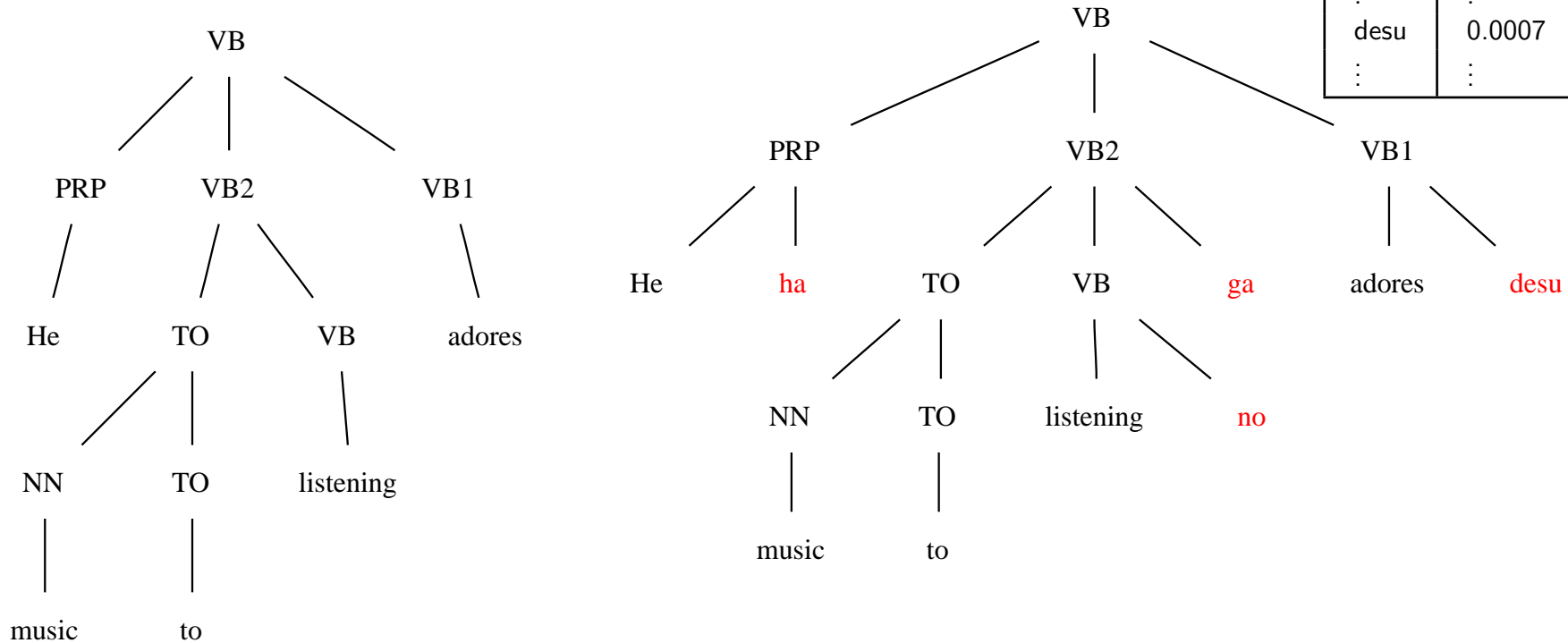
Reordering probability:  $0.723 \cdot 0.749 \cdot 0.893 = 0.484$

## 5.3 TREE-TO-STRING MODELS

⇒ The insertion of a new node is decided according to the *n-table*

parent	TOP	VB	VB	VB	TO	TO	...
node	VB	VB	PRP	TO	TO	NN	...
P(None)	0.735	0.687	0.344	0.709	0.900	0.800	...
P(Left)	0.004	0.061	0.004	0.030	0.003	0.096	...
P(right)	0.260	0.252	0.652	0.261	0.007	0.104	...

w	P(ins-w)
ha	0.219
ta	0.131
wo	0.099
no	0.094
ni	0.080
te	0.078
ga	0.062
⋮	⋮
desu	0.0007
⋮	⋮

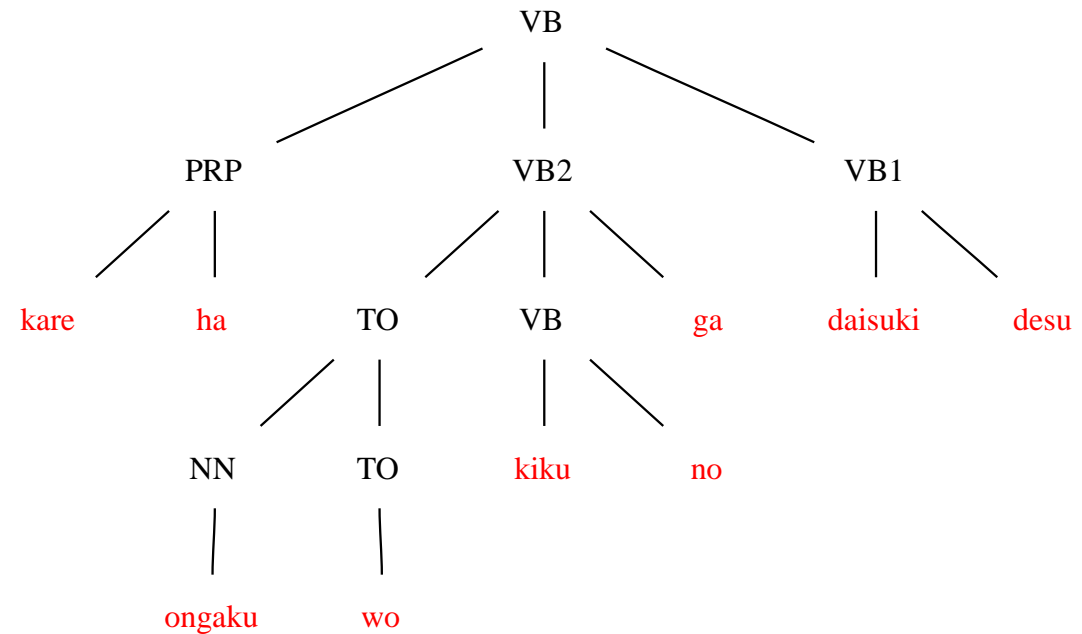
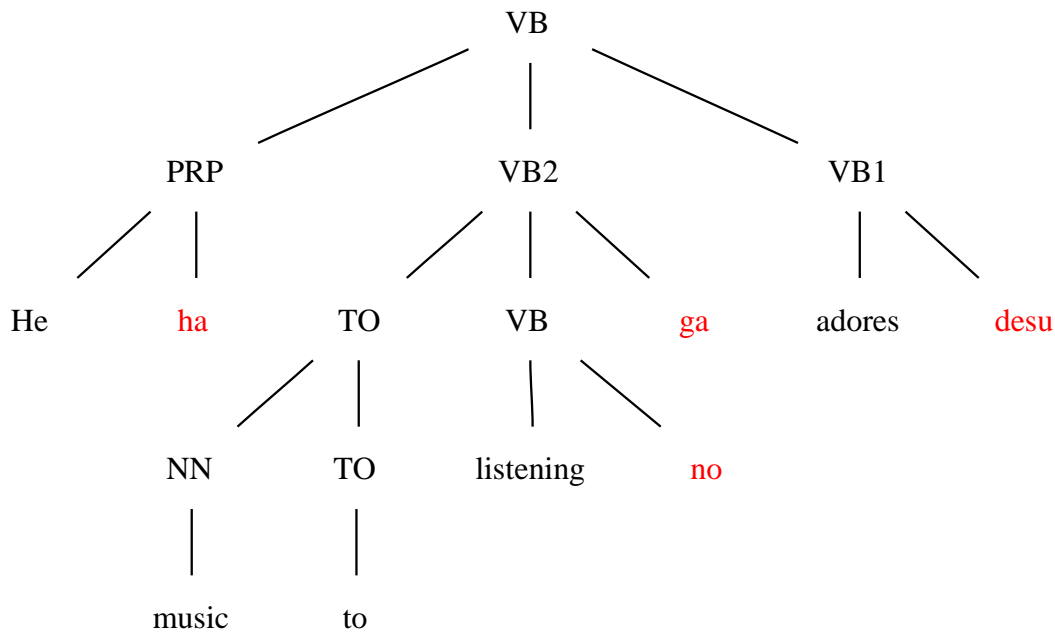


Insertion probability:  $(0.652 \cdot 0.219) \cdot (0.252 \cdot 0.094) \cdot (0.252 \cdot 0.062) \cdot (0.252 \cdot 0.0007) \cdot 0.735 \cdot 0.709 \cdot 0.900 \cdot 0.800 = 3.498e - 9$

## 5.3 TREE-TO-STRING MODELS

⇒ The translation is decided according to the *t-table*

adores		he		listening		music		to		...
daisuki	1.000	kare	0.952	kiku	0.333	ongaku	0.900	ni	0.216	...
		NULL	0.016	kii	0.333	naru	0.100	NULL	0.204	
		nani	0.005	mi	0.333			to	0.133	
		⋮	⋮	⋮	⋮			⋮	⋮	



Translation probability:  $0.952 \cdot 0.900 \cdot 0.038 \cdot 1.000 = 0.0108$

### Decoder description

- Given a French sentence, the decoder will find the most plausible English parse tree
- Idea: a mechanism similar to normal parsing is used
- Steps:
  1. Start from an English context-free grammar and incorporate to it the channel operations
  2. For each non-lexical rule (such as “VP → VB NP PP”), supplement the grammar with reordered rules and probabilities are taken from the r-table
  3. Rules such as “VP → VP X” and “X → *word*” are added and probabilities are taken from the n-table
  4. For each lexical rule in the English grammar, we add rules such as “englishWord → foreingWord”
  5. Parse a string of foreign words
  6. Undo reordering operations and remove leaf nodes with foreign words
  7. Among all possible tree, choose pick the best in which the product of the LM and the TM probability is the highest



### Main ideas [Chiang 07]

- It allows to capture difficult reordering
- Hierarchical phrases: phrases that can contain other phrases
- Related to Synchronous CFG: useful for specifying relations between languages.
- Rules are as follows:

$$X \rightarrow \langle \gamma, \alpha, \sim \rangle$$

where

- $X$  is a non-terminal symbol
- $\gamma, \alpha$  are strings of terminal and non-terminal symbols
- $\sim$  is one-to-one correspondence between non-terminal occurrences in  $\gamma$  and  $\alpha$

### Rule extraction

- Rules are extracted from word-alignments sentences
  - Extract a rule for each phrase pair
  - Replace phrase pairs in each rule by a non-terminal symbol if another rule produces that phrase pair.
- The set of rules of two word-aligned sentences  $\langle f, e, \sim \rangle$  is the smallest set satisfying the following:

- If  $\langle f_i^j, e_{i'}^{j'} \rangle$  is an initial phrase pair, then add the following rule:

$$X \rightarrow \langle f_i^j, e_{i'}^{j'} \rangle$$

- If  $(X \rightarrow \langle \gamma, \alpha \rangle)$  is a rule and  $\langle f_i^j, e_{i'}^{j'} \rangle$  is an initial phrase pair such that  $\gamma = \gamma_1 f_i^j \gamma_2$  and  $\alpha = \alpha_1 e_{i'}^{j'} \alpha_2$ , then add the following rule:

$$X \rightarrow \langle \gamma_1 X_k \gamma_2, \alpha_1 X_k \alpha_2 \rangle$$

- Glue rules:

$$S \rightarrow \langle S_1 X_2, S_1 X_2 \rangle$$

$$S \rightarrow \langle X_1, X_1 \rangle$$

### Translation model

- Log-linear model over derivations:

$$P(D) \propto \prod_i \Phi_i(D)^{\lambda_i}$$

where  $\Phi_i$  are features defined on derivations and  $\lambda_i$  are feature weights

- Features: functions on the rules and an additional LM function:

$$P(D) \propto P_{LM}(e)^{\lambda_{LM}} \prod_{i \neq LM} \prod_{(X \rightarrow \langle \gamma, \alpha \rangle) \in D} \Phi_i(X \rightarrow \langle \gamma, \alpha \rangle)^{\lambda_i}$$

- Features on rules:

- $P(\gamma \mid \alpha)$  and  $P(\alpha \mid \gamma)$
- Lexical weights:  $P_w(\gamma \mid \alpha)$  and  $P_w(\alpha \mid \gamma)$
- A penalty  $\exp(-1)$  to learn a preference for longer or shorter derivations
- Word penalty:  $\exp(-\#T(\alpha))$

### Training

- Rules probabilities obtained from frequencies
- $\lambda_i$ : minimum-error-rate training [Och 02]
- CKY-based algorithm

---

# REFERENCES

# REFERENCES

---

- [Barrachina 99] S. Barrachina and J.M Vilar. *Bilingual clustering using monolingual algorithms*. TMI. 1999.
- [Brown 90] P. F. Brown et al. *A statistical approach to machine translation*. Computational Linguistics, 16, 79–85, 1990.
- [Brown 93] P. F. Brown et al. *The mathematics of statistical machine translation: parameter estimation*. Computational Linguistics, 19(2), 263–310, 1993.
- [Chiang 07] D. Chiang *Hierarchical phrase-based translation*. Computational Linguistics, 33(2), 201–228, 2007.
- [Gascó 10a] G. Gascó and J.A. Sánchez. *Syntax augmented inversion transduction grammars for machine translation*. 11th International Conference on Intelligent Text Processing and Computational Linguistics (CICLING), March, 2010.
- [Gascó 10b] G. Gascó, J.A. Sánchez and J.M. Benedí. *Enlarged Search Space for SITG Parsing*, Proc. 11th Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL HLT), June, 2010, 653-656.
- [Knight 99] K. Knight. *Decoding Complexity in Word-Replacement Translation Models*, Computational Linguistics, Squibs & Discussion, 25(4), 1999.
- [Kumar 04] S. Kumar and W. Byrne. *A Weighted Finite State Transducer Implementation of the Alignment Template Model for Statistical Machine Translation*. Proceedings of HLT-NAACL 2003, May 2003.

# REFERENCES

---

- [Ney 00a] H. Ney, S. Nießen, F. Och, H. Sawaf, C. Tillmann and S. Vogel. *Algorithms for Statistical Translation of Spoken Language*. IEEE Transactions on Speech and Audio Processing, vol. 8(1), 24–36, 2000.
- [Ney 03a] H. Ney *Statistical Natural Language Processing*, 2003, Canadian Hansard.
- [Och 99] F.J. Och. *An Efficient Method for Determining Bilingual Word Classes*. EACL. 1999.
- [Och 02] F.J. Och. *Discriminative training and maximum entropy models for statistical machine translation*. Proc. of ACL, 295–302, 2002.
- [Tillmann 01] C. Tillmann. *Word re-ordering and DP based search for SMT*. PhD Thesis, 2001.
- [Wu 97] D. Wu. *Stochastic Inversion Transduction Grammars and Bilingual Parsing of Parallel Corpora*. Computational Linguistics, 23(3):377-403, 1997.
- [Yamada 01] K. Yamada and K. Knight. *A Syntax-Based Statistical Translation Model*. ACL, 2001.