



ĐẠI HỌC BÁCH KHOA HÀ NỘI
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

Machine translation

Lê Thanh Hương

School of Information and Communication Technology

Email: huonglt@soict.hust.edu.vn

An example

- Au sortir de la saison 97/98 et surtout au debut de cette saison 98/99...
- With leaving season 97/98 and especially at the beginning of this season 98/99...

Challenges

1. Capture variation and similarities amongst languages
 - Morphologically: # morphemes/word:
 - *Isolating* languages (Vietnamese, Cantonese) – 1 word/ 1 morpheme
 - *Polysynthetic languages* (Siberian Yupik), 1 word = a whole sentence
 - Degree to which morphemes are segmentable

Challenges

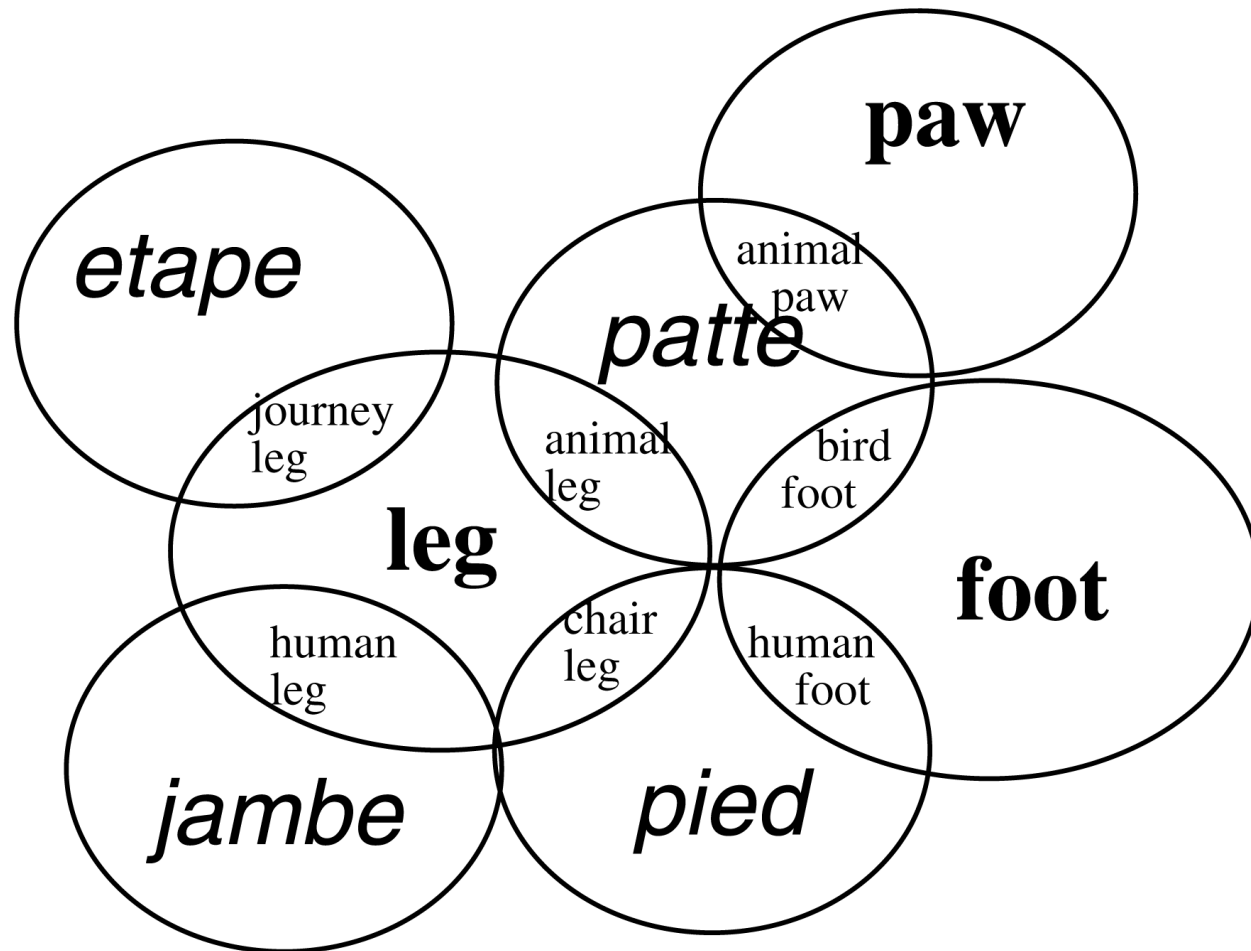
2. Syntax: order of words in a sentence

- *To Yukio; Yukio ne*
- English vs. Vietnamese:
 - *The* (affix1) *red* (affix2) *flag* (head)
 - *Lá cờ* (head) *đỏ* (affix2) *ấy* (affix1)

3. Differences in specificity

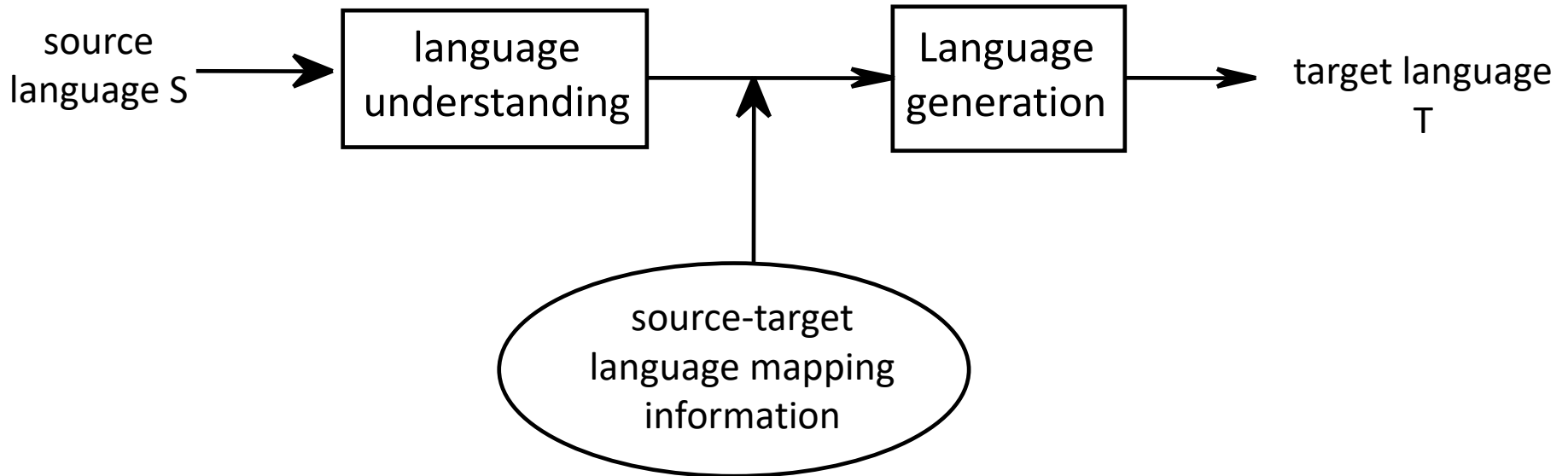
English	brother	Vietnamese	anh em
English	wall	German	wand (inside) mauer(outside)
German	berg	English	hill mountain

Conceptual space



Lexical gap: Jp, no word for *privacy*; Eng: no word for *oyakoko* (filial piety)

Three main blocks in machine translation



Language understanding

1. Lexical ambiguity:

English: *book* - Spanish *libro, reservar*

⇒ Use syntactic context

2. Syntactic ambiguity:

I saw the guy on the hill with the telescope

3. Semantic ambiguity

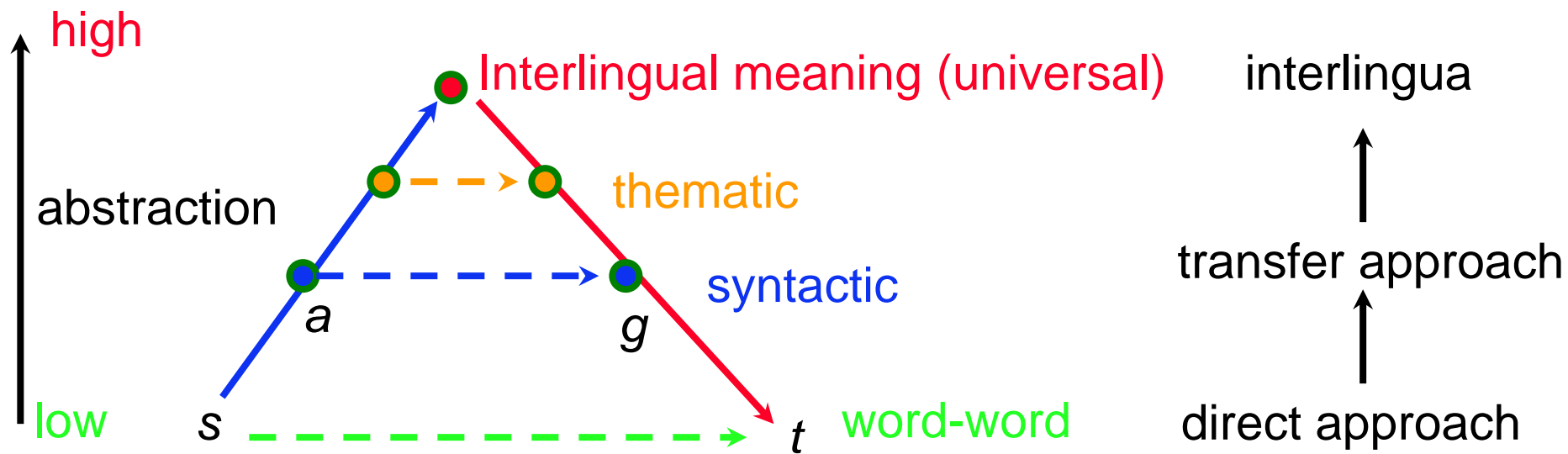
- *E: While driving, John swerved & hit a tree*



John's car

- *S: Mientras que John estaba manejando, se desvio y golpeo con un arbo*

Methods of machine translation

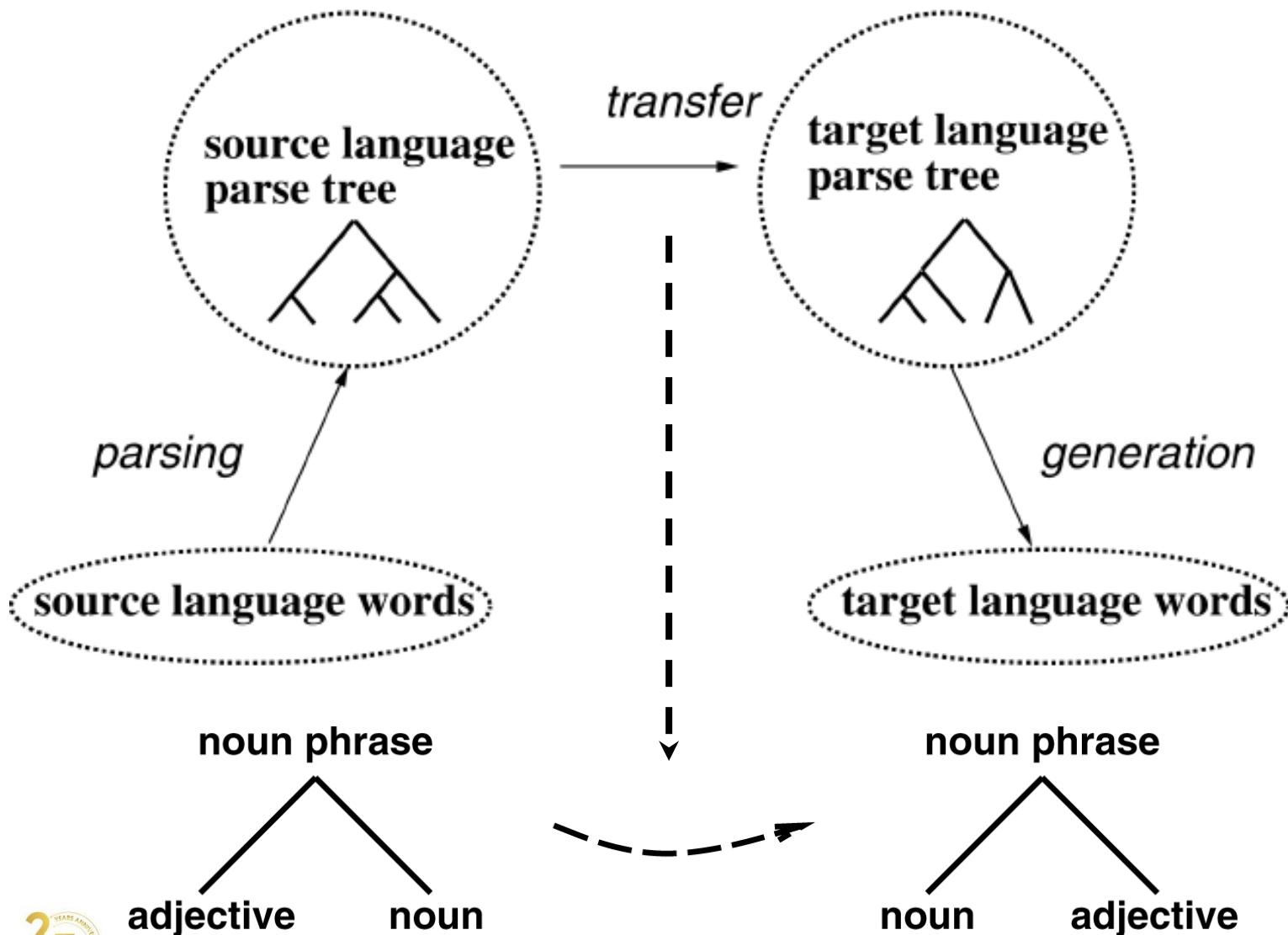


$$a = a(s)$$

$$g = f(a(s)); f - \text{transfer function}$$

$$t = g(f(a(s)))$$

The triangle/transfer diagram



List of transforms

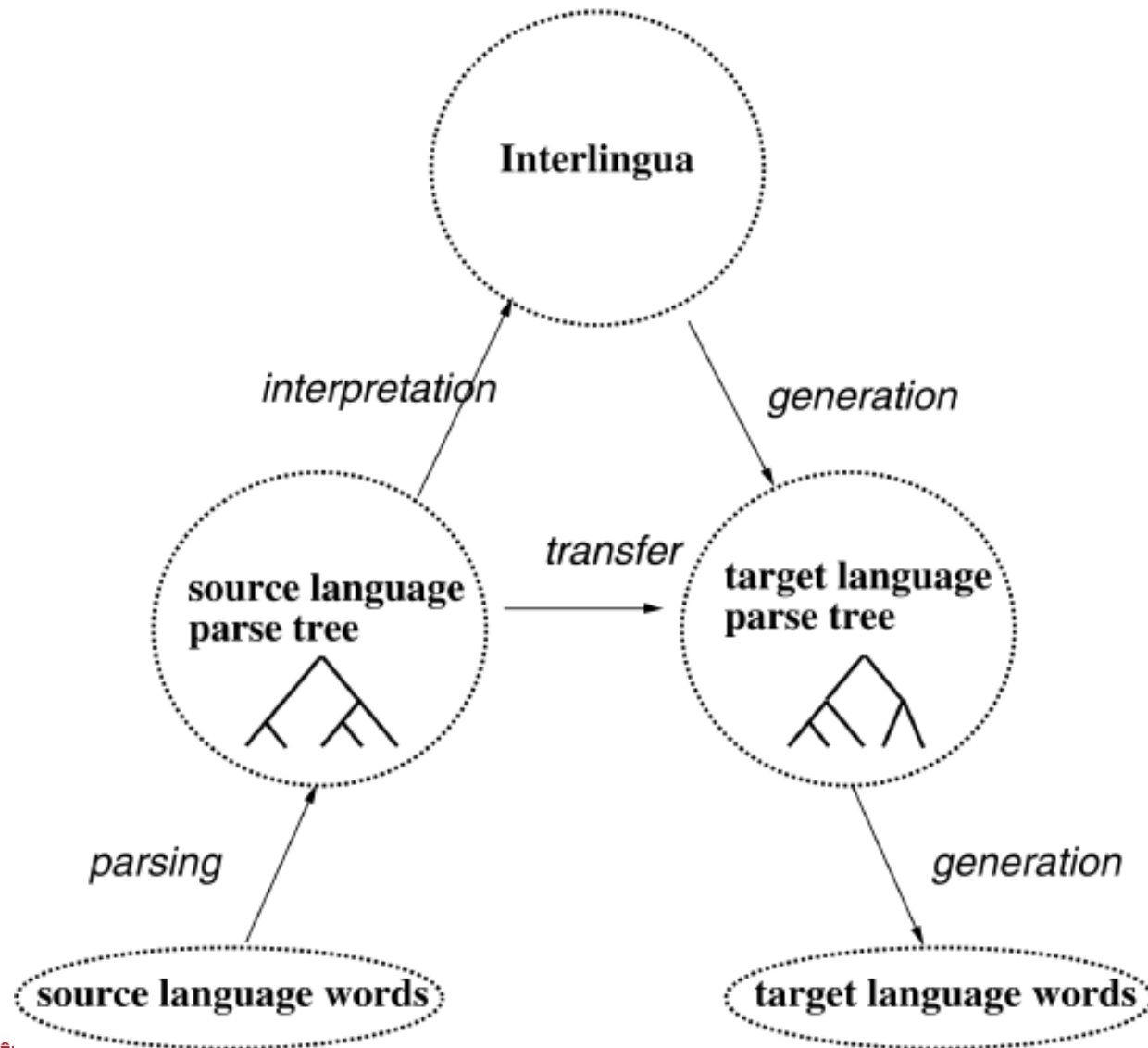
English to French:

1. $NP \rightarrow \text{Adjective}_1 \text{ Noun}_2$
 \Rightarrow
 $NP \rightarrow \text{Noun}_2 \text{ Adjective}_1$

Japanese to English:

2. $\text{Existential-There-Sentence} \rightarrow \text{There}_1 \text{ Verb}_2 \text{ NP}_3 \text{ Postnominal}_4$
 \Rightarrow
 $\text{Sentence} \rightarrow (\text{NP} \rightarrow \text{NP}_3 \text{ Relative-Clause}_4) \text{ Verb}_2$
3. $NP \rightarrow \text{NP}_1 \text{ Relative Clause}_2$
 \Rightarrow
 $NP \rightarrow \text{Relative-Clause}_2 \text{ NP}_1$

The triangle/transfer diagram



Interlingua approach: using meaning

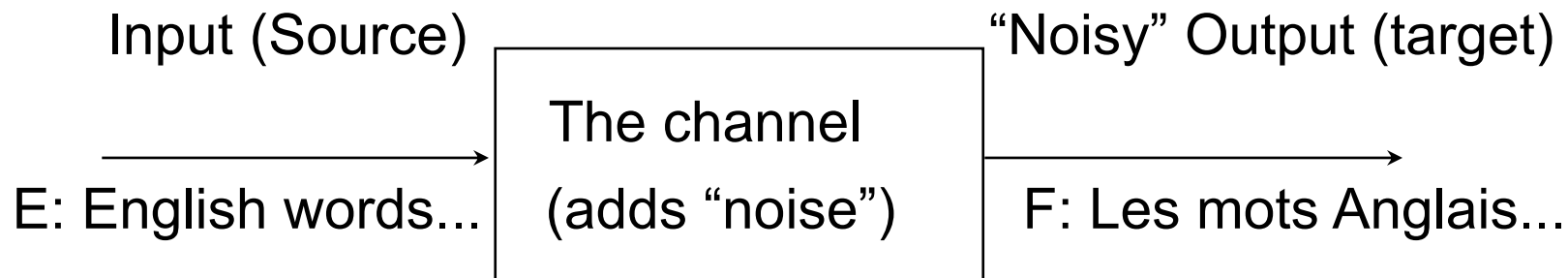
- Transfer: one pair of rules per language pair
- Objects/events (ontology)

event	gardening						
agent	<table><tr><td>man</td><td></td></tr><tr><td>number</td><td>sg</td></tr><tr><td>definiteness</td><td>indef</td></tr></table>	man		number	sg	definiteness	indef
man							
number	sg						
definiteness	indef						
aspect	progressive						
tense	past						

Statistical machine translation

The Main Idea

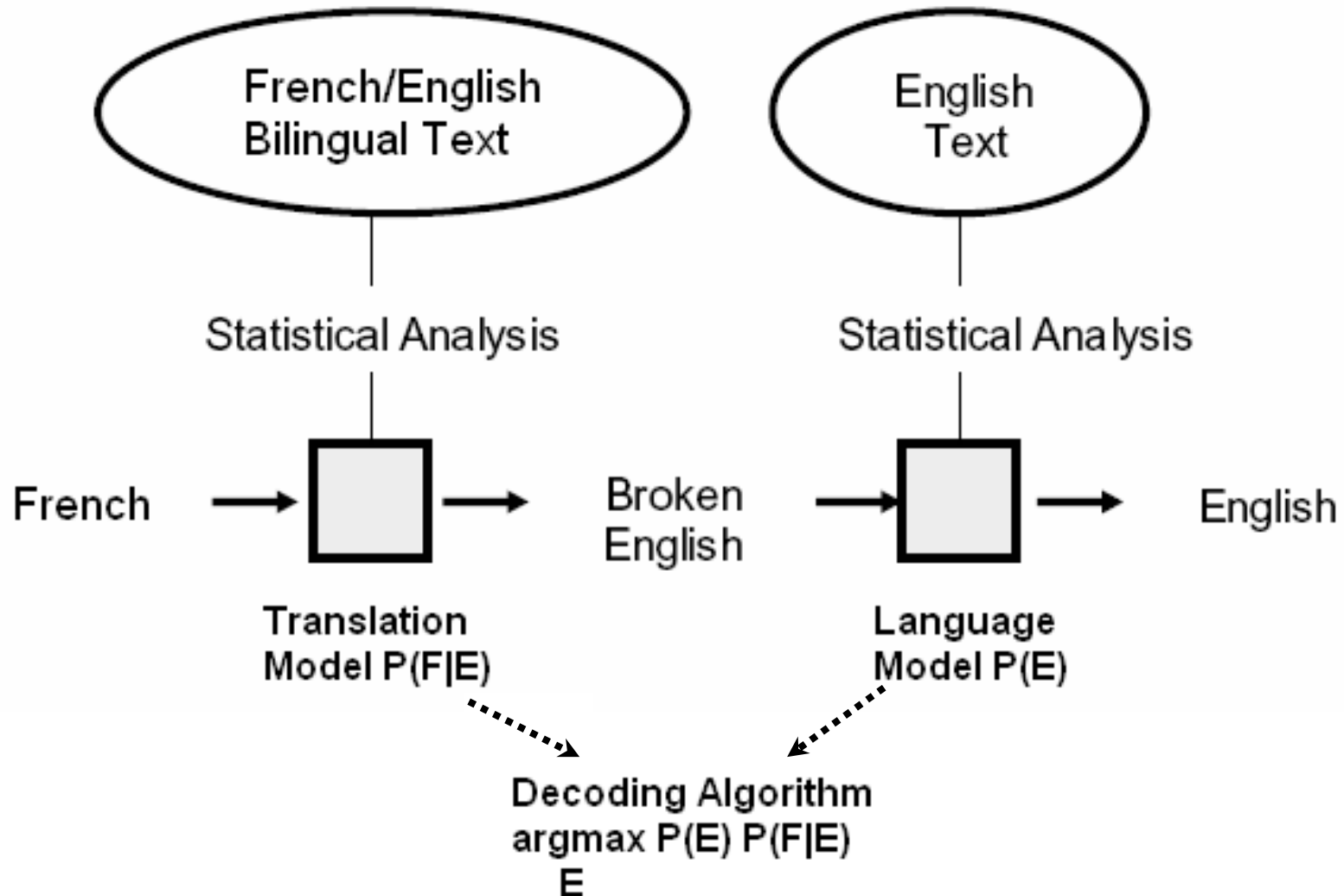
- Treat translation as a noisy channel problem



- The translation model: $P(E|F) = P(F|E) P(E) / P(F)$
- Interested in rediscovering \underline{E} given \underline{F} :
After the usual simplification ($P(F)$ fixed):

$$\operatorname{argmax}_E P(E|F) = \operatorname{argmax}_E P(F|E) P(E)$$

A Statistical MT System



The necessities

- **Language Model** (LM): our *expectation* of seeing a particular *English* (E) sentence, *a priori*:
 $P(E)$
- **Translation Model** (TM): Target sentence in *French* (F) *given* source sentence in English:
 $P(F|E)$
- Search procedure
 - Given F, find best E using the LM and TM distributions.
- Usual problem: sparse data!
 - We cannot create a “sentence dictionary” $E \leftrightarrow F$
 - Typically, we do not see a sentence even twice!

Alignment Idea

- TM doesn't care about correct strings of English words
- Use the “tagging” approach:
 - 1 English word (“tag”) ~ 1 French word (“word”)
 - not realistic: even #words in two sentences isn't equal
 - use “Alignment”.
- **Sentence alignment**: find some group of sentences in one language corresponds to some other group of sentences in another language.

Sentence Alignment

The old man is
happy. He has
fished many times.
His wife talks to
him. The fish are
jumping. The
sharks await.

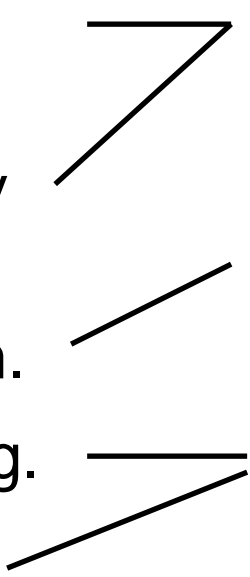
El viejo está feliz
porque ha pescado
muchos veces. Su
mujer habla con él.
Los tiburones
esperan.

Sentence Alignment

1. The old man is happy.
2. He has fished many times.
3. His wife talks to him.
4. The fish are jumping.
5. The sharks await.

1. El viejo está feliz porque ha pescado muchos veces.
2. Su mujer habla con él.
3. Los tiburones esperan.

Sentence Alignment

- | | | |
|------------------------------|--|--|
| 1. The old man is happy. | | 1. El viejo está feliz porque ha pescado muchos veces. |
| 2. He has fished many times. | | 2. Su mujer habla con él. |
| 3. His wife talks to him. | | 3. Los tiburones esperan. |
| 4. The fish are jumping. | | |
| 5. The sharks await. | | |
- 

Difficulties:

- ✿ **Crossing dependencies:** the order of sentences are changed in the translation.

Word alignment - easy

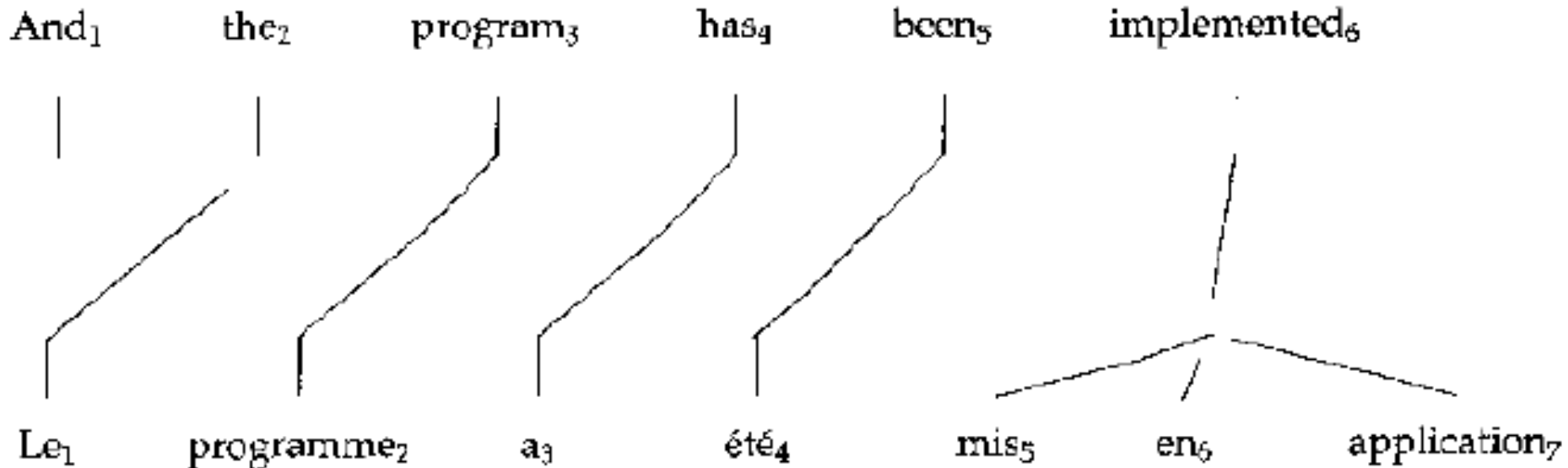
Japan shaken by two new quakes|

Le Japon secoué par deux nouveaux séismes

Extra word appears in French: “spurious” word

Word alignment - harder

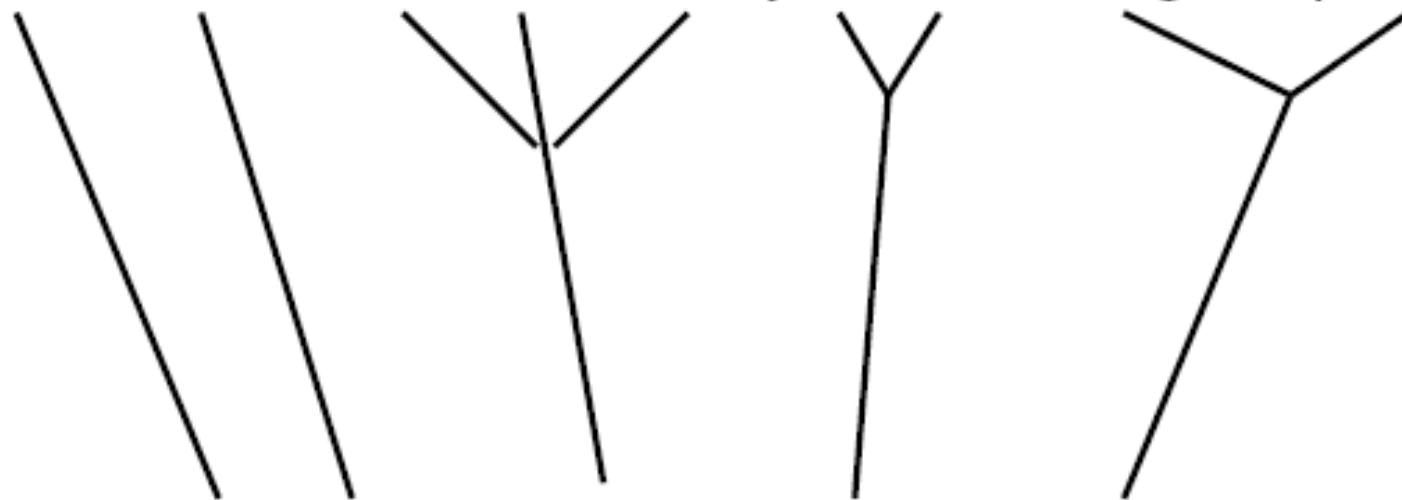
“Zero fertility” word: not translated



One word translated as several words

Word alignment - harder

The balance was the territory of the aboriginal people

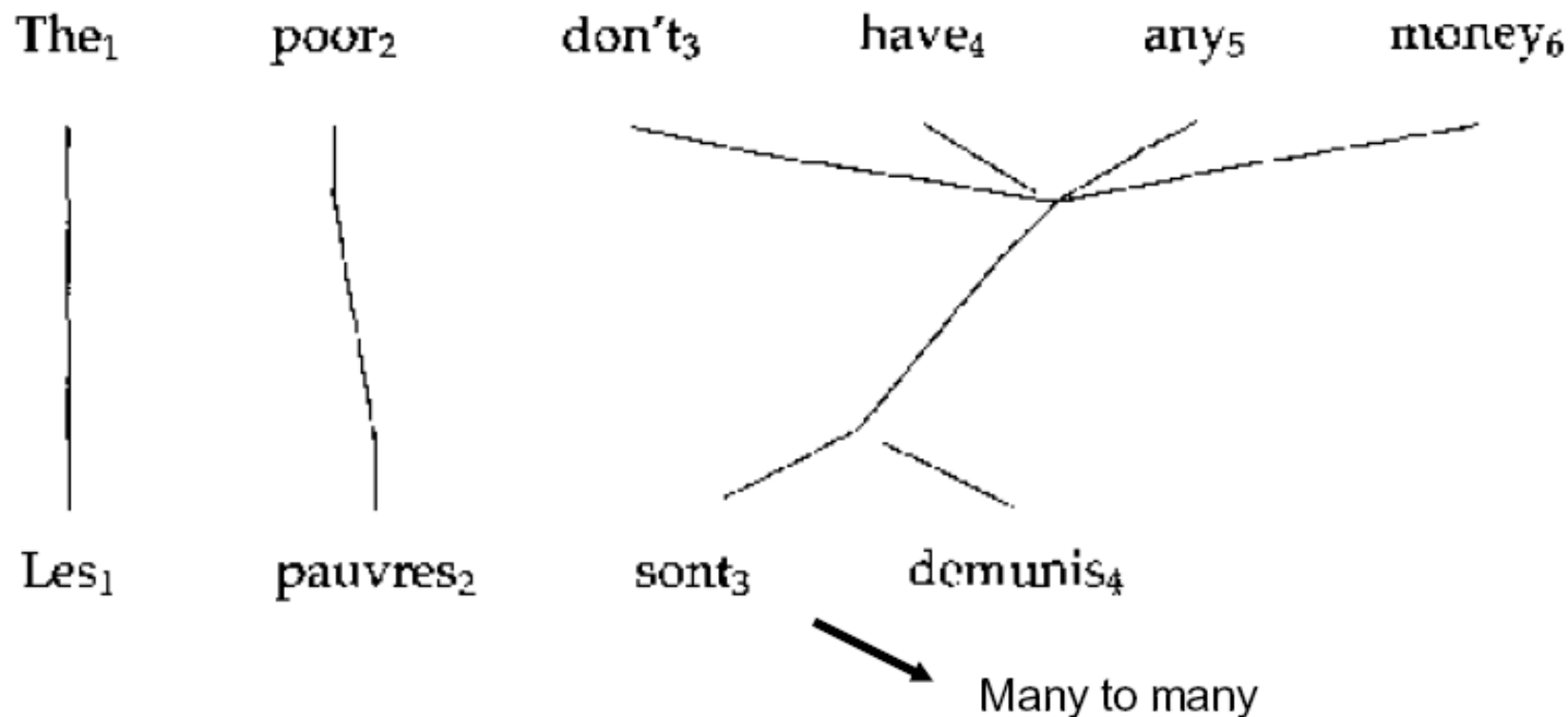


Le reste appartenait aux autochtones



Several words translated as one

Word alignment - hard



- A line group linking a minimal subset of words is called a 'ceptr' in the IBM work

Word alignment - encoding

0 1 2 3 4 5 6

- e_0 And the program has been implemented



- f_0 Le programme a été mis en application

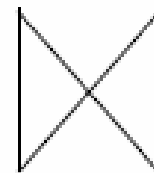
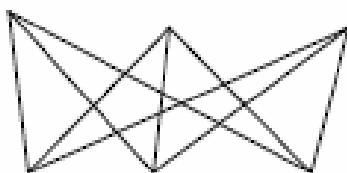
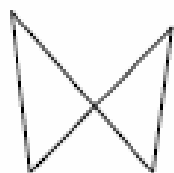
0 1 2 3 4 5 6 7

- Linear notation:

- $f_0(1)$ Le(2) programme(3) a(4) été(5) mis(6) en(6) application(6)
- e_0 And(0) the(1) program(2) has(3) been(4) implemented(5,6,7)

Word alignment learning with EM

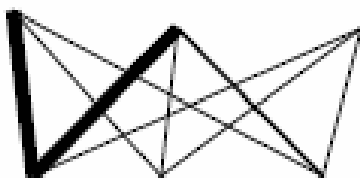
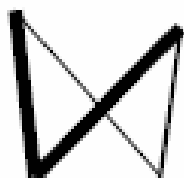
... la maison ... la maison bleue ... la fleur ...



... the house ... the blue house ... the flower ...



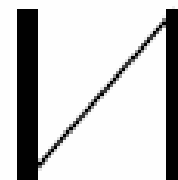
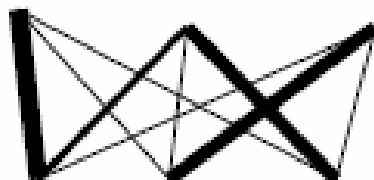
... la maison ... la maison bleue ... la fleur ...



... the house ... the blue house ... the flower ...

Word alignment learning with EM

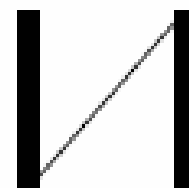
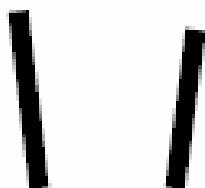
... la maison ... la maison bleue ... la fleur ...



... the house ... the blue house ... the flower ...



... la maison ... la maison bleue ... la fleur ...



... the house ... the blue house ... the flower ...

Word alignment learning with EM

... la maison ... la maison bleue ... la fleur ...

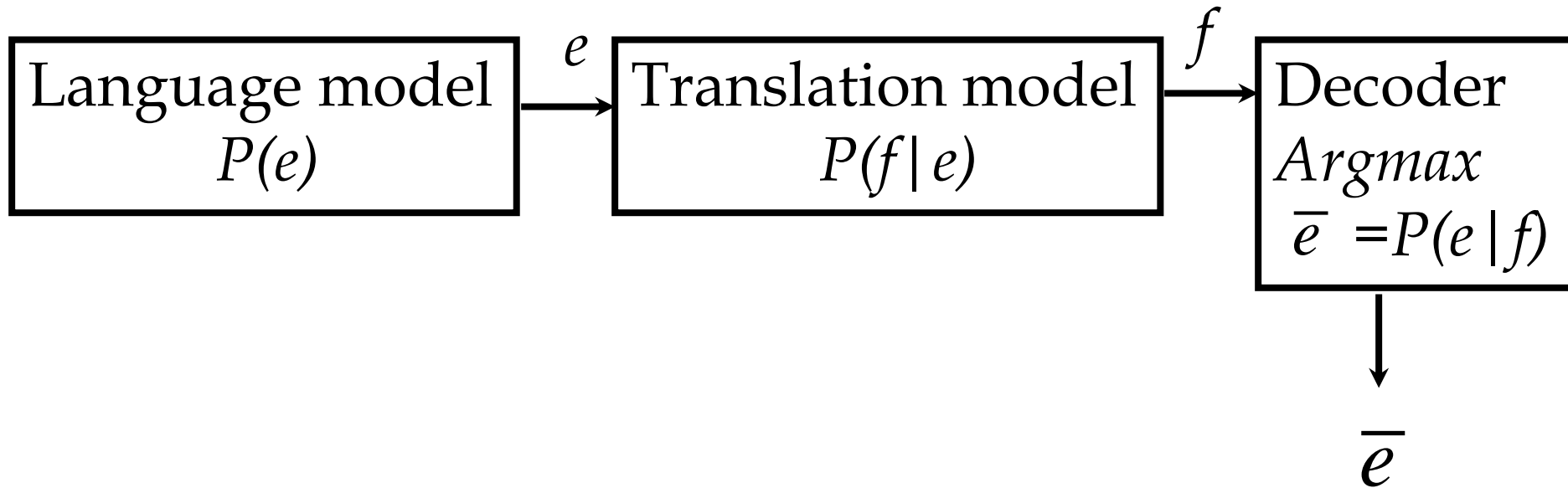
\ /

\ X

| |

... the house ... the blue house ... the flower ...

Noisy channel again



Another Alignment System

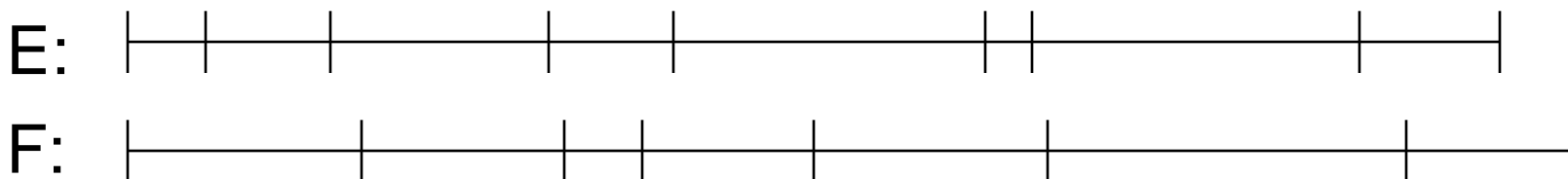
- Available corpus assumed:
 - parallel text (translation $E \leftrightarrow F$)
- Sentence alignment
 - sentence detection
 - sentence alignment
- Word alignment
 - tokenization
 - word alignment (with restrictions)

Sentence Boundary Detection

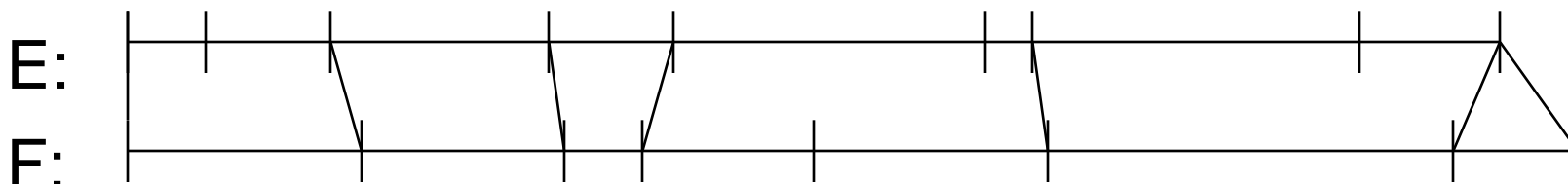
- Rules, lists:
 - Sentence breaks:
 - paragraphs (if marked)
 - certain characters: ?, !, ; (...almost sure)
 - Problem: period .
 - end of sentence (... left yesterday. He was heading to...)
 - decimal point: 3.6 (three-point-six)
 - thousand segment separator: 3.200
 - abbreviation: cf., e.g., Calif., Mt., Mr.
 - ellipsis: ...
 - other languages: ordinal number indication (2nd ~ 2.)
 - initials: A. B. Smith
- Statistical methods: e.g., Maximum Entropy

Sentence Alignment

- Problem: sentences detected only:



- **Desired output**: Segmentation with equal number of segments, spanning continuously the whole text.
- Original sentence boundaries kept:

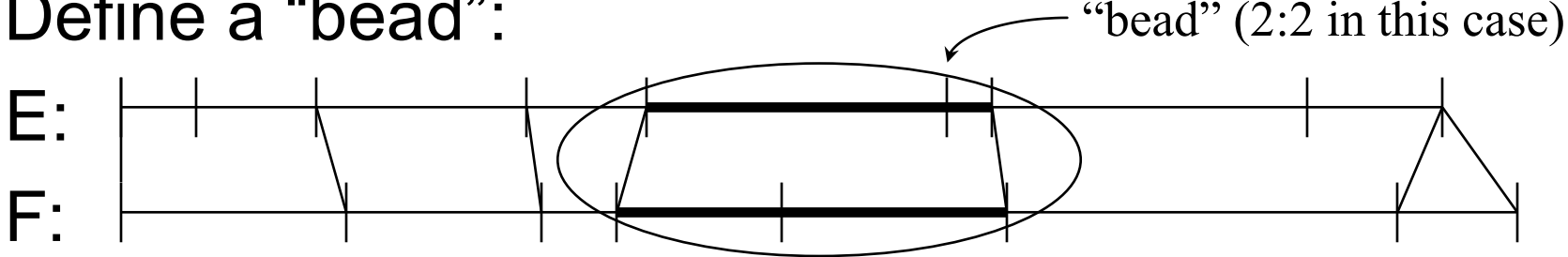


- Alignments obtained: 2-1, 1-1, 1-1, 2-2, 2-1, 0-1

Alignment Methods

- Several methods (probabilistic and not prob.)
 - character-length based
 - word-length based
 - “cognates” (word identity used)
 - using an existing dictionary (F: prendre ~ E: make, take)
 - using word “distance” (similarity): names, numbers, borrowed words, Latin origin words, ...
- Best performing:
 - statistical, word- or character- length based (with some words perhaps)

Length-based Alignment

- First, define the problem probabilistically:
$$\operatorname{argmax}_A P(A|E,F) = \operatorname{argmax}_A P(A,E,F) \quad (E,F \text{ fixed})$$
- Define a “bead”:


The diagram shows two horizontal lines, E and F, representing sequences. Vertical tick marks are placed along both lines. A region where the sequences overlap is highlighted with a thick black bar and an oval. An arrow points to this region with the label “bead” (2:2 in this case).
- Approximate:
$$P(A,E,F) \cong \prod_{i=1..n} P(B_i),$$

where B_i is a bead; $P(B_i)$ does not depend on the rest of E,F .

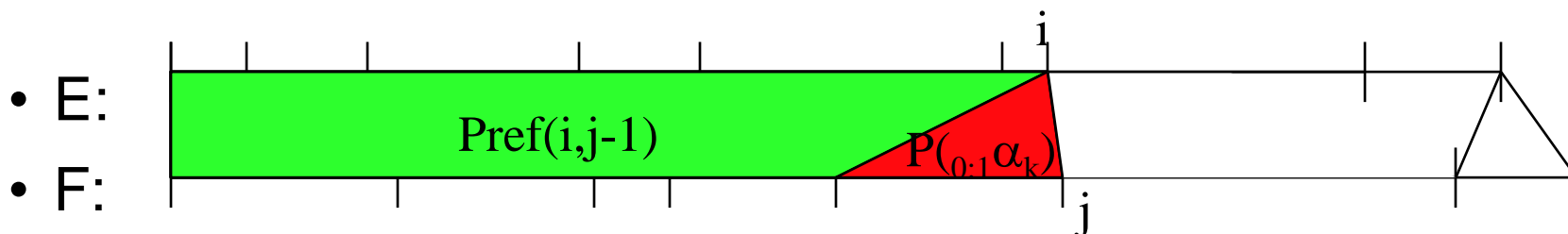
The Alignment Task

Some definitions:

- Given $P(A,E,F) \cong \prod_{i=1..n} P(B_i)$,
find the partitioning of (E,F) into n beads $B_{i=1..n}$,
that maximizes $P(A,E,F)$ over training data.
- $B_i =_{p:q} \alpha_i$, where $p:q \in \{0:1, 1:0, 1:1, 1:2, 2:1, 2:2\}$
describes the type of alignment
- $\text{Pref}(i,j)$ - probability of the best alignment from the
start of (E,F) data $(1,1)$ up to (i,j)

Recursive Definition

- Initialize: $\text{Pref}(0,0) = 0$.
- $\text{Pref}(i,j) = \max ($
 $\text{Pref}(i,j-1) P_{(0:1)}\alpha_k, \text{Pref}(i-1,j) P_{(1:0)}\alpha_k, \text{Pref}(i-1,j-1) P_{(1:1)}\alpha_k,$
 $\text{Pref}(i-1,j-2) P_{(1:2)}\alpha_k, \text{Pref}(i-2,j-1) P_{(2:1)}\alpha_k, \text{Pref}(i-2,j-2) P_{(2:2)}\alpha_k)$
- This is enough for a Viterbi-like search.



Probability of a Bead

- Remains to define $P_{(p:q)}\alpha_k$:
 - \underline{k} refers to the “next” bead, with segments of \underline{p} and \underline{q} sentences, lengths $l_{k,e}$ and $l_{k,f}$.
- Use normal distribution for length variation:
$$P_{(p:q)}\alpha_k = P(\delta(l_{k,e}, l_{k,f}, \mu, \sigma^2), p:q) \cong P(\delta(l_{k,e}, l_{k,f}, \mu, \sigma^2))P(p:q)$$
$$\delta(l_{k,e}, l_{k,f}, \mu, \sigma^2) = (l_{k,f} - \mu l_{k,e}) / \sqrt{l_{k,e} \sigma^2}$$
- Estimate $P(p:q)$ from small amount of data, or even guess and re-estimate after aligning some data.
- Words etc. might be used as better clues in $P_{(p:q)}\alpha_k$ def.

Word Alignment

- Length alone does not help:
 - words can be swapped, and mutual translations have often vastly different length.
- Idea:
 - Assume some (simple) translation model.
 - Find its parameters by considering virtually all alignments.
 - After we have the parameters, find the best alignment given those parameters.

Given the following bilingual dataset:
Mike yêu Jane. Mike loves Jane.
Jane yêu hoa. Jane loves flowers.
Jane thích đọc sách. Jane likes reading.
Compute the following translation probabilities:
 $P(\text{"Mike likes reading"} | \text{Mike thích đọc sách})$
 $P(\text{"Mike likes flowers"} | \text{Mike yêu hoa})$

- Start with sentence-aligned corpus.
- Let (E,F) be a pair of sentences (actually, a bead).
 1. Initialize $p(f|e)$ randomly, $f \in F$, $e \in E$.
 2. Compute expected counts over the corpus:

$$c(f,e) = \sum_{(E,F); e \in E, f \in F} p(f|e)$$

\forall aligned pair (E,F) , find if e in E and f in F ; if yes, add $p(f|e)$.

3. Reestimate:

$$p(f|e) = c(f,e) / c(e) \quad [c(e) = \sum_f c(f,e)]$$

4. Iterate until change of $p(f|e)$ is small.

Best Alignment

Select, for each (E,F),

$$A = \operatorname{argmax}_A P(A|F,E) = \operatorname{argmax}_A P(F,A|E)/P(F) = \\ \operatorname{argmax}_A P(F,A|E) = \operatorname{argmax}_A (\varepsilon / (l+1)^m \prod_{j=1..m} p(f_j|e_{a_j})) = \operatorname{argmax}_A \prod_{j=1..m} p(f_j|e_{a_j})$$

- Use dynamic programming, Viterbi-like algorithm.
- Recompute $p(f|e)$

Exercise

- Given the following bilingual dataset:
 - Mike yêu Jane. Mike loves Jane.
 - Jane yêu hoa. Jane loves flowers.
 - Jane thích đọc sách. Jane likes reading.
- Compute the following translation probabilities:
 - $P(\text{"Mike thích đọc sách"} | \text{Mike likes reading})$
 - $P(\text{"Mike yêu hoa"} | \text{Mike likes flowers})$

Evaluate

Evaluation based on Hansard corpus:

- 48% of French sentences are translated correctly
- 2 types of errors:
 - Mistranslation:
 - Permettez que je donne un exemple à chambre
 - Let me give an example in the House (incorrect decoding)
 - (Let me give the House an example)
 - Grammatical translation:
 - Vous avez besoin de toute l'aide disponible
 - You need all of the benefits available (ungrammatical decoding)
 - (You need all the help you can get)

Reason

- **Distortion**: English words at the beginning of a sentence are aligned with French words at the end of a sentence – this reduces the probability of alignment
- **Fertility** : the correspondence between English and French words (1-to-1, 1-to-2, 1-to-0, ...),
 - For example, fertility(**farmers**) in the corpus = 2, because this word when translated into English usually consists of 2 words: **les agriculteurs**
 - To go

Reason

- **Independent Assumptions:** Short sentences are preferred because there are fewer probabilities (when multiplying)

⇒ multiply the result by a constant proportional to the sentence length

- **Training data dependence:** a small change in the training data causes a large change in the parameter estimates

For example, $P(le/the)$ changed from 0.610 to 0.497

- **Efficiency.** Remove sentences > 30 words, because the search space increases exponentially
- **Lack of language knowledge**

Lack of language knowledge

- Can't save information about terms: example can't be aligned “ to go ” and “ aller ”
- No local binding :

Eg, is she a mathematician

- **Phonemic.** Words made up of different phonemes are considered separate symbols
- **Sparse data.** Ratings for uncommon words are incorrect

Open sources

- GIZA++: statistical machine translation tool to train IBM 1-5 model for word alignment
- MOSES: statistical machine translation tool
- Moses has two types of translation: **phrase-based** and **tree-based**

Machine translation using the syntax

Why Syntax?

- Need much more grammatical output
- Need accurate control over re-ordering
- Need accurate insertion of function words
- Word translations need to depend on grammatically-related word

Yamada and Knight (2001): The need for phrasal syntax

- He adores listening to music.

The diagram illustrates the need for phrasal syntax by showing incorrect connections between the Japanese sentence and its English translation. Lines connect the English words to the Japanese words in a way that does not reflect the actual syntactic structure. For example, a line connects 'He' to 'daisuki' (loves), and another connects 'listening' to 'kiku' (listen). These lines are crossed out with a large 'X', indicating that such a word-for-word mapping is incorrect. The correct syntactic structure is represented by the Japanese sentence itself, which is a phrasal sentence where the verb 'kiku' is part of a larger phrase 'kiku no ga daisuki' (listening is loved).

彼は音楽を聞くのが大好きです
Kare ha ongaku wo kiku no ga daisuki desu

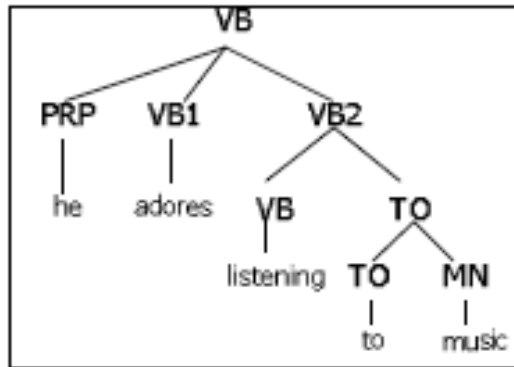
Syntax-based model



- Preprocess English by a parser
- Probabilistic operations on a parse-tree
 1. Reorder child nodes
 2. Insert extra nodes
 3. Translate leaf words

Parse Tree(E) → Sentence (J)

Parse Tree(E)



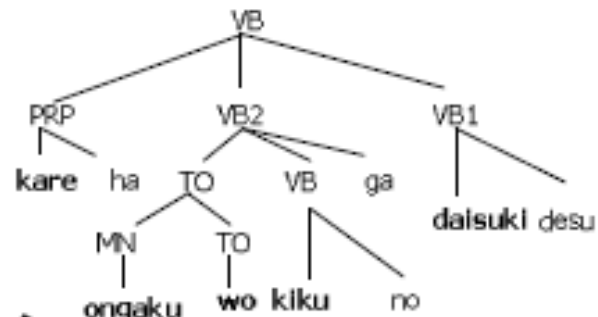
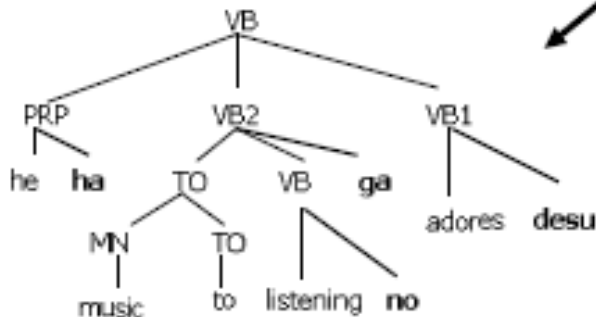
Reorder



Insert



Translate



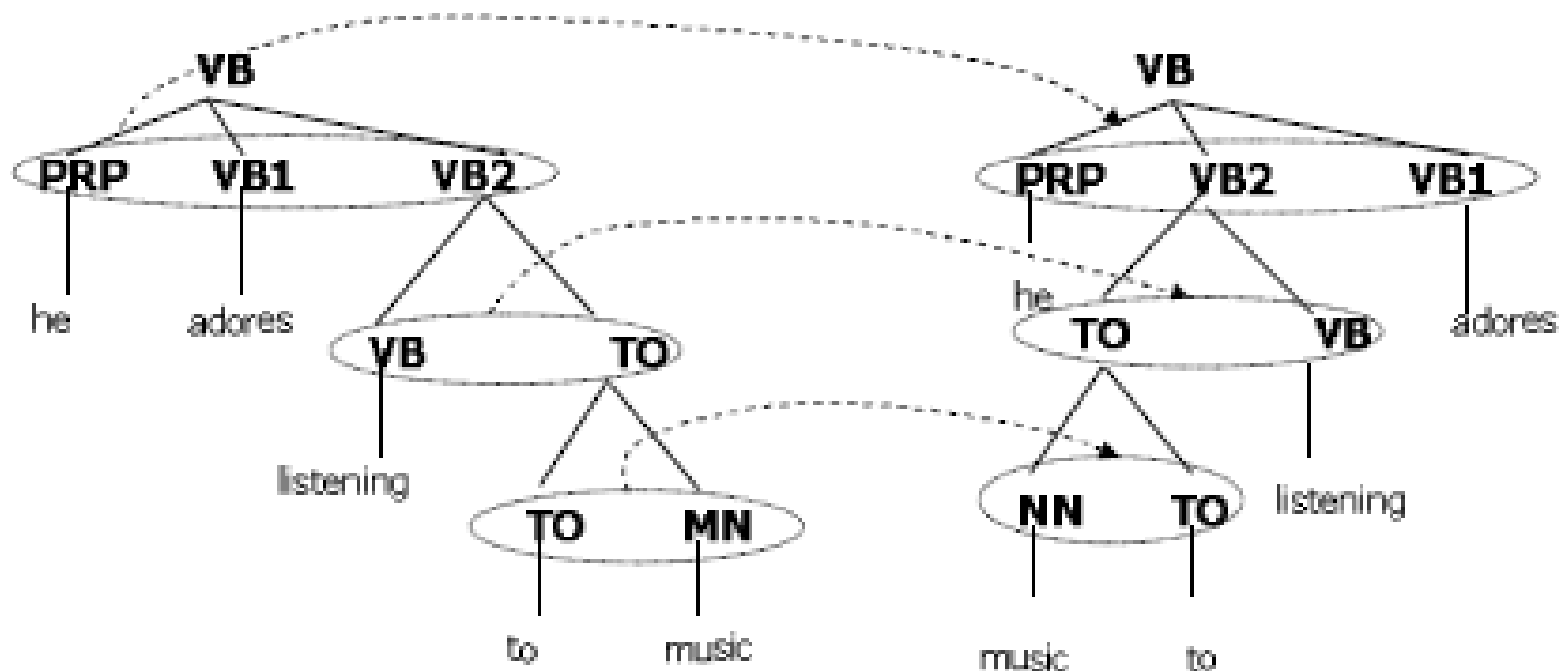
Take Leaves



Sentence(J)

Kare ha ongaku wo kiku no ga daisuki desu

1. Reorder



$$P(\text{PRP VB1 VB2} \mid \text{PRP VB2 VB1}) = 0.723$$

$$P(\text{VB TO} \mid \text{TO VB}) = 0.749$$

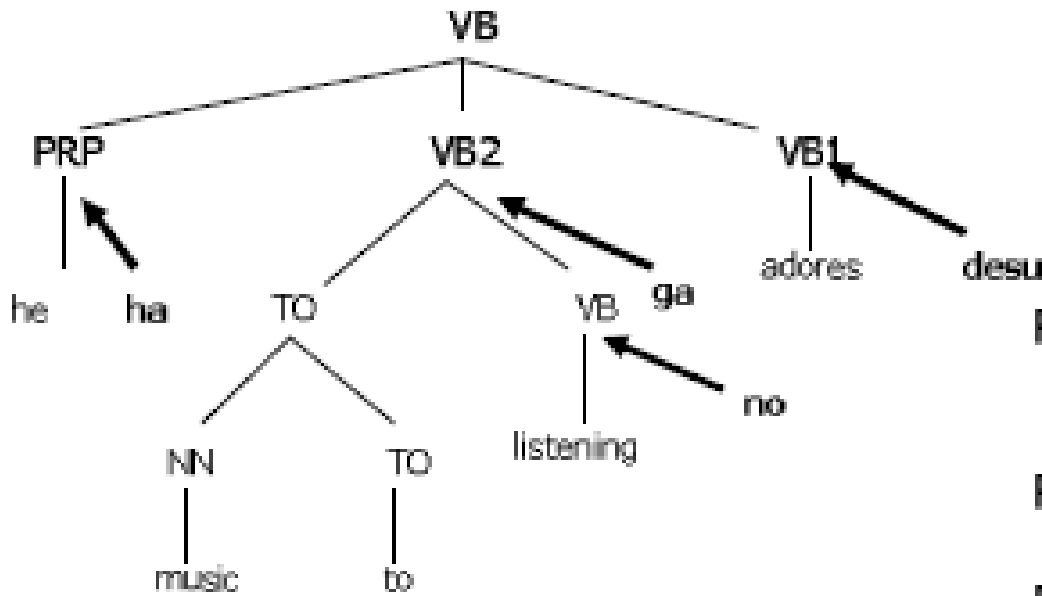
$$P(\text{TO NN} \mid \text{NN TO}) = 0.893$$

Conditional feature = child label sequence

Parameter Table: Reorder

Original order	Coming back	P(Reorder Original order)
PRP VB1 VB2	PRP VB1 VB2	0.074
	PRP VB2 VB1	0.723
	VB1 PRP VB2	0.061
	VB1 VB2 PRP	0.037
	VB2 PRP VB1	0.083
	VB2 VB1 PRP	0.021
VB TO	VB TO	0.107
	TO VB	0.893
TO NN	TO NN	0.251
	NN TO	0.749

2. Insert



$$P(\text{none}|\text{TOP-VB}) = 0.735$$

⋮

$$P(\text{right}|\text{VB-PRP}) * P(\text{ha}) = 0.652 * 0.219$$

$$P(\text{right}|\text{VB-VB}) * P(\text{ga}) = 0.252 * 0.062$$

⋮

$$P(\text{none}|\text{TO-TO}) = 0.900$$

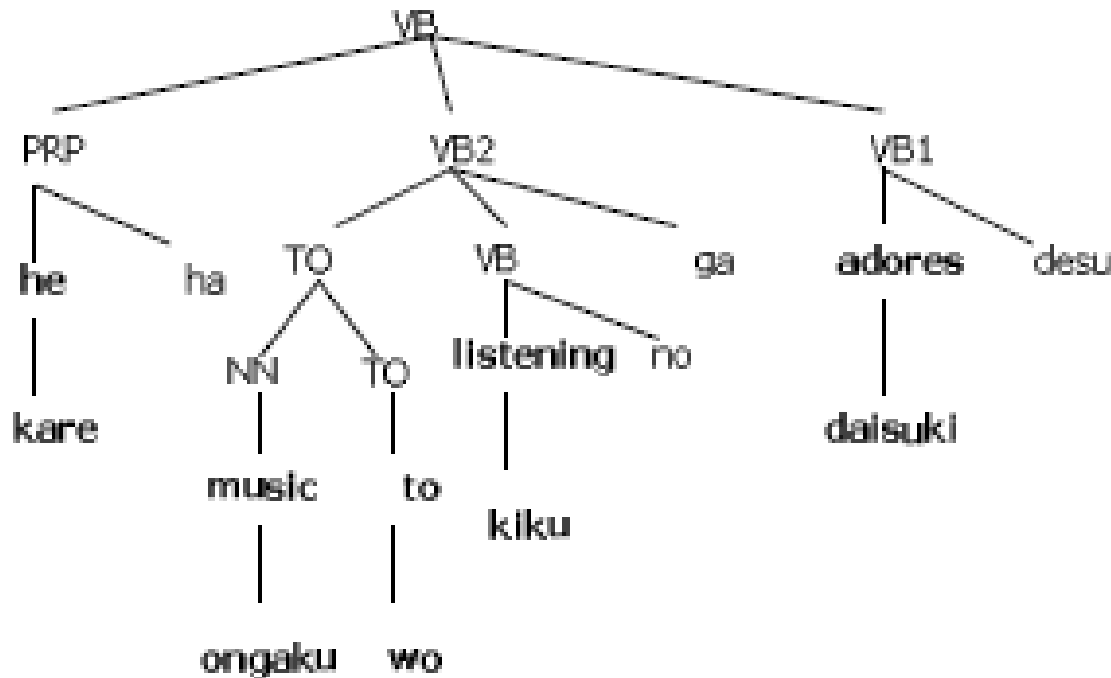
Conditional Feature = parent label & node label (position) & none (word selection)

Parameter table: insert

Parent label node level	TOP VB	VB VB	VB TO	TO TO	TO NN	TO NN
P (none)	0.735	0.687	0.344	0.700	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	0.097	0.104

W	P (insert-w)
ha	0.219
ta	0.131
wo	0.099
no	0.094
ni	0.090
te	0.078
ga	0.062
deu	0.0007

3. Translate



$P(\text{he} \text{---} \text{kare}) = 0.952$
 $P(\text{music} \text{---} \text{ongaku}) = 0.900$
 $P(\text{to} \text{---} \text{wo}) = 0.038$
 $P(\text{listening} \text{---} \text{kiku}) = 0.333$
 $P(\text{adore} \text{---} \text{daisuki}) = 1.000$

Conditional feature = word identity (English)

Parameter table: Translate

E	adores	he	listening	music	to
J	daisuki 1.000	kare 0.952 NULL 0.016 nani 0.005 da 0.003 shi 0.003 	kiku 0.333 kii 0.333 mi 0.333	ongaku 0.900 naru 0.100	ni 0.216 NULL 0.204 to 0.133 no 0.046 wo 0.038

Note: Translation to NULL = deletion

Experiment

- Training Corpus: J-E 2K sentence pairs
- J: Tokenized by Chasen [Matsumoto, et al., 1999]
- E: Parsed by Collins Parser [Collins, 1999]
 - Trained: 40K Treebank, Accuracy: ~90%
- E: Flatten parse tree
 - To Capture word-order difference (SVO->SOV)
- EM Training: 20 Iterations
 - 50 min/iter (Sparc 200Mhz 1-CPU) or
 - 30 sec/iter (Pentium3 700Mhz 30-CPU)

Result

	Average	#perf sent
Y/K models	0.582	10
IBM model 5	0.431	0

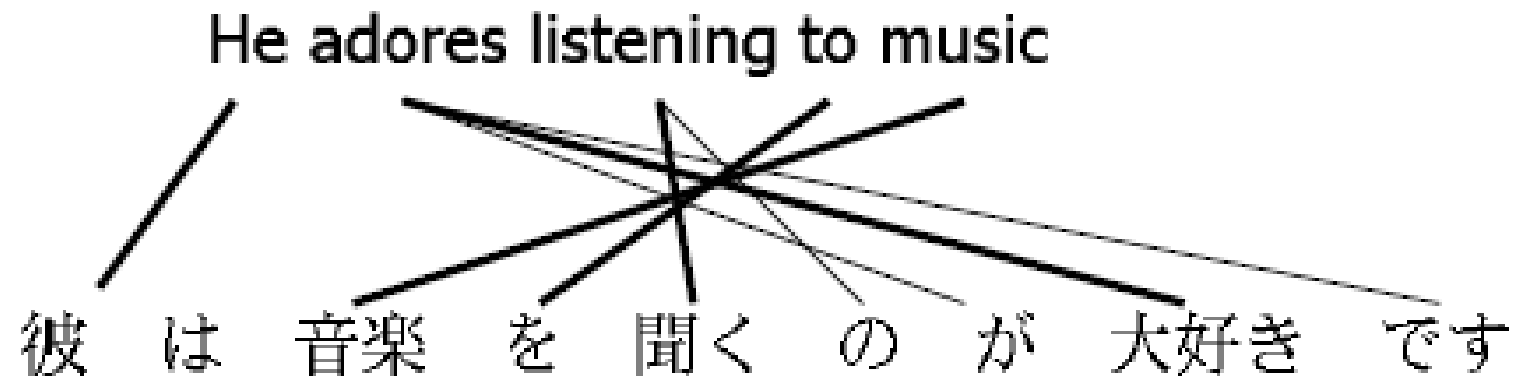
- Average by 3 humans for 50 sentences
- ok(1.0), not sure (0.5), wrong (0.0)
- Precision only

Result: Alignment 1

Syntax-based Model



IBM Model 3

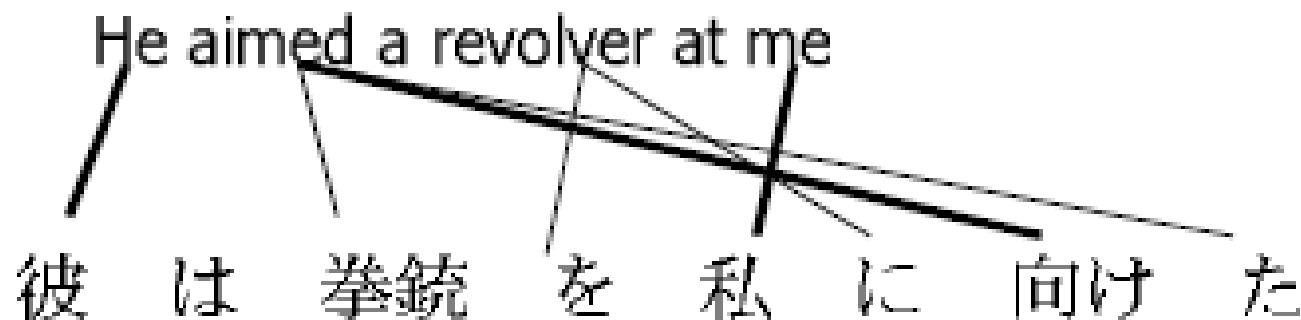


Result: Alignment 2

Syntax-based model



IBM Model 3



Some open sources

- See <http://fosmt.org/>
 - Moses
 - Giza++

Some MT systems on the web

- http://www.google.com/language_tools?hl=en
- <http://www.systransoft.com/index.html>
- <http://babelfish.altavista.digital.com/>