

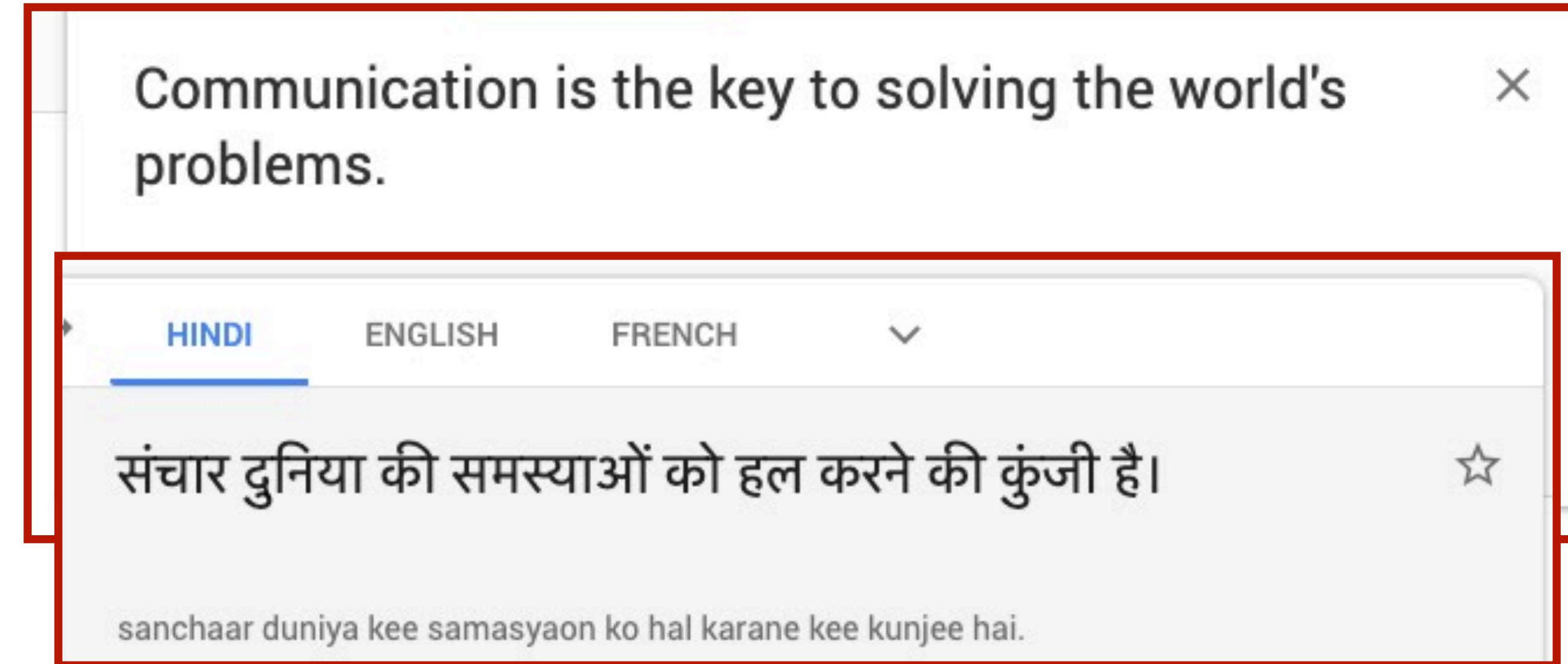
CMPT 825: Natural Language Processing

# Machine Translation

Spring 2020

Adapted from slides from Chris Manning, Abigail See, Matthew Lamm,  
Danqi Chen and Karthik Narasimhan

# Translation



- One of the “holy grail” problems in artificial intelligence
- Practical use case: Facilitate communication between people in the world
- Extremely challenging (especially for low-resource languages)

# Easy and not so easy translations

- Easy:
  - I like apples ↔ ich mag Äpfel (German)
- Not so easy:
  - I like apples ↔ J'aime les pommes (French)
  - I like red apples ↔ J'aime les pommes rouges (French)
  - /es ↔ the but /les pommes ↔ apples

# MT basics

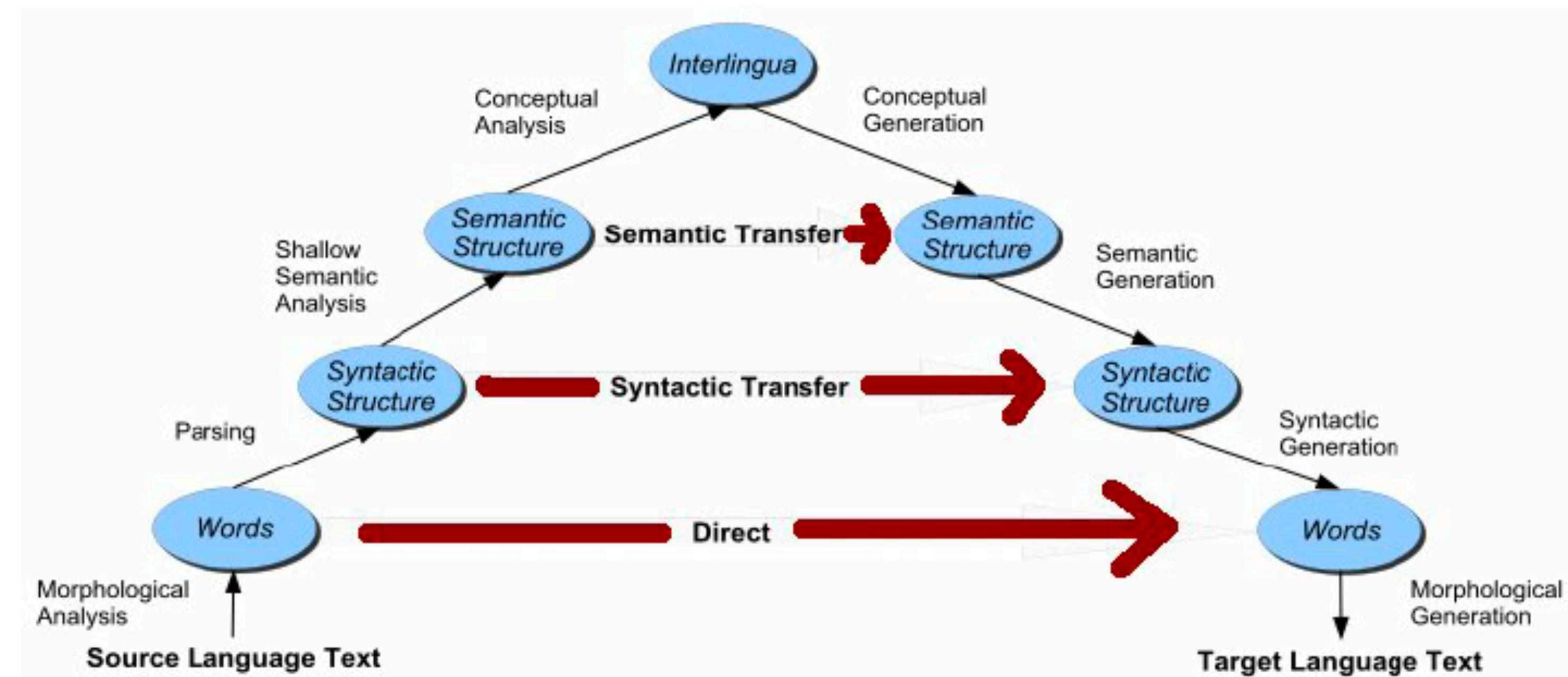
- **Goal:** Translate a sentence  $w^{(s)}$  in a **source language (input)** to a sentence in the **target language (output)**
- Can be formulated as an optimization problem:
  - $\hat{w}^{(t)} = \arg \max_{w^{(t)}} \psi(w^{(s)}, w^{(t)})$
  - where  $\psi$  is a scoring function over source and target sentences
- Requires **two** components:
  - Learning algorithm to compute parameters of  $\psi$
  - Decoding algorithm for computing the best translation  $\hat{w}^{(t)}$

# Why is MT challenging?

- Single words may be replaced with multi-word phrases
  - I like **apples**  $\leftrightarrow$  J'aime **les pommes**
- Reordering of phrases
  - I like **red apples**  $\leftrightarrow$  J'aime **les pommes rouges**
- Contextual dependence
  - *les*  $\leftrightarrow$  *the* but *les pommes*  $\leftrightarrow$  *apples*

Extremely large output space  $\implies$  Decoding is NP-hard

# Vauquois Pyramid



- Hierarchy of concepts and distances between them in different languages
- Lowest level: individual words/characters
- Higher levels: syntax, semantics
- Interlingua: Generic language-agnostic representation of meaning

# Evaluating translation quality

- Two main criteria:
  - **Adequacy:** Translation  $w^{(t)}$  should adequately reflect the linguistic content of  $w^{(s)}$
  - **Fluency:** Translation  $w^{(t)}$  should be fluent text in the target language

	Adequate?	Fluent?
<i>To Vinay it like Python</i>	yes	no
<i>Vinay debugs memory leaks</i>	no	yes
<i>Vinay likes Python</i>	yes	yes

Different translations of *A Vinay le gusta Python*

# Evaluation metrics

- Manual evaluation is most accurate, but expensive
- Automated evaluation metrics:
  - Compare system hypothesis with reference translations
  - BiLingual Evaluation Understudy (BLEU) (Papineni et al., 2002):
    - Modified n-gram precision

$$p_n = \frac{\text{number of } n\text{-grams appearing in both reference and hypothesis translations}}{\text{number of } n\text{-grams appearing in the hypothesis translation}}$$

# BLEU

$$\text{BLEU} = \exp \frac{1}{N} \sum_{n=1}^N \log p_n$$

Two modifications:

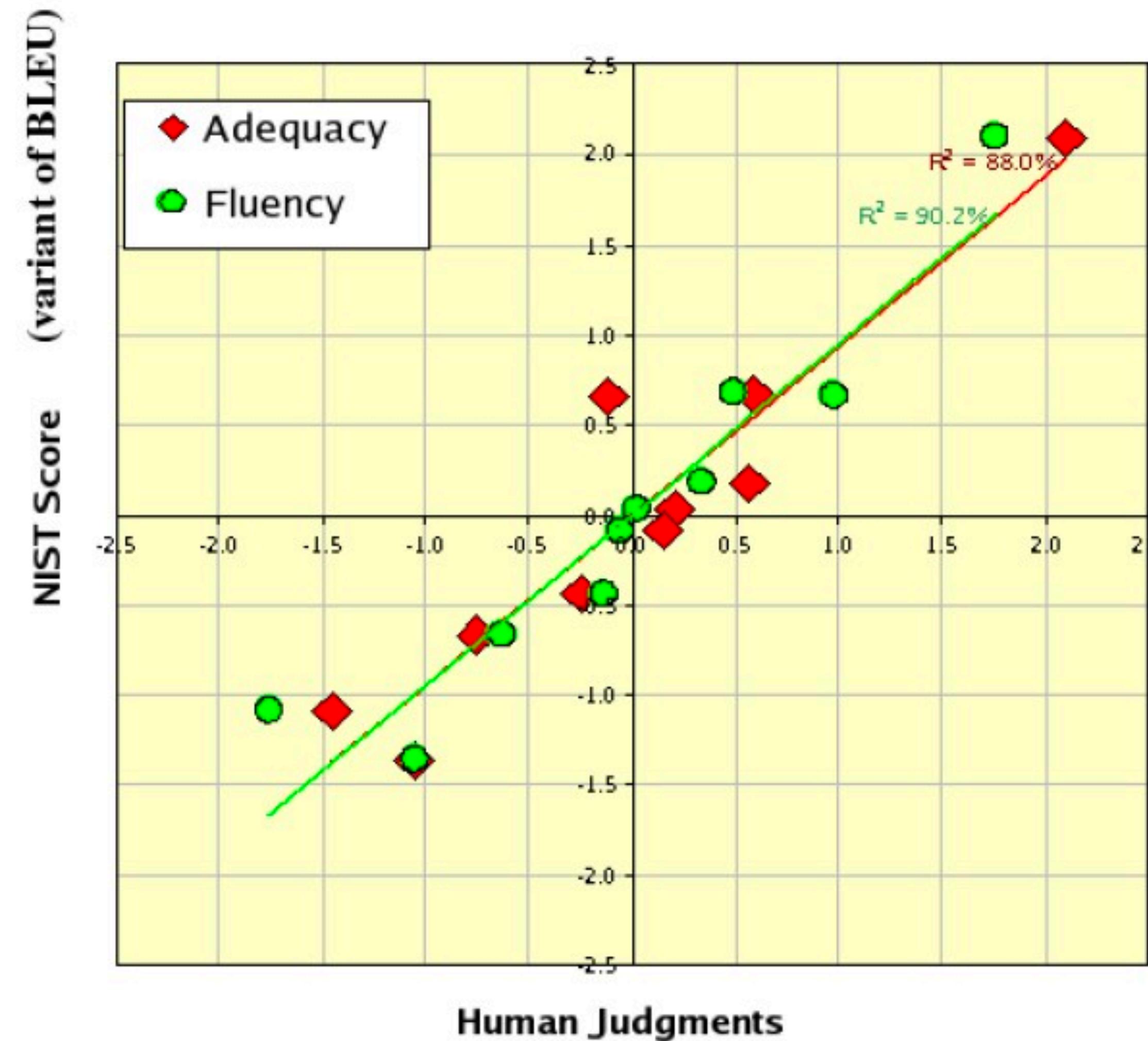
- To avoid  $\log 0$ , all precisions are smoothed
- Each n-gram in reference can be used at most once
  - Ex. **Hypothesis**: *to to to to to* vs **Reference**: *to be or not to be* should not get a unigram precision of 1

Precision-based metrics favor short translations

- Solution: Multiply score with a brevity penalty for translations shorter than reference,  $e^{1-r/h}$

# BLEU

- Correlates somewhat well with human judgements



(G. Doddington, NIST)

# BLEU scores

	<b>Translation</b>	$p_1$	$p_2$	$p_3$	$p_4$	BP	BLEU
<i>Reference</i>	<i>Vinay likes programming in Python</i>						
<i>Sys1</i>	<i>To Vinay it like to program Python</i>	$\frac{2}{7}$	0	0	0	1	.21
<i>Sys2</i>	<i>Vinay likes Python</i>	$\frac{3}{3}$	$\frac{1}{2}$	0	0	.51	.33
<i>Sys3</i>	<i>Vinay likes programming in his pajamas</i>	$\frac{4}{6}$	$\frac{3}{5}$	$\frac{2}{4}$	$\frac{1}{3}$	1	.76

Sample BLEU scores for various system outputs

- Alternatives have been proposed:  
    - METEOR: weighted F-measure
    - Translation Error Rate (TER): Edit distance between hypothesis and reference
- Issues?**

# Machine Translation (MT)

task of translating a sentence  
from one language (the **source language**)  
to a sentence in another language (the **target language**)

# History

- Started in the **1950s**: rule-based, tightly linked to formal linguistics theories
  - Russian → English (motivated by the Cold War!)
  - Systems were mostly **rule-based**, using a bilingual dictionary to map Russian words to their English counterparts



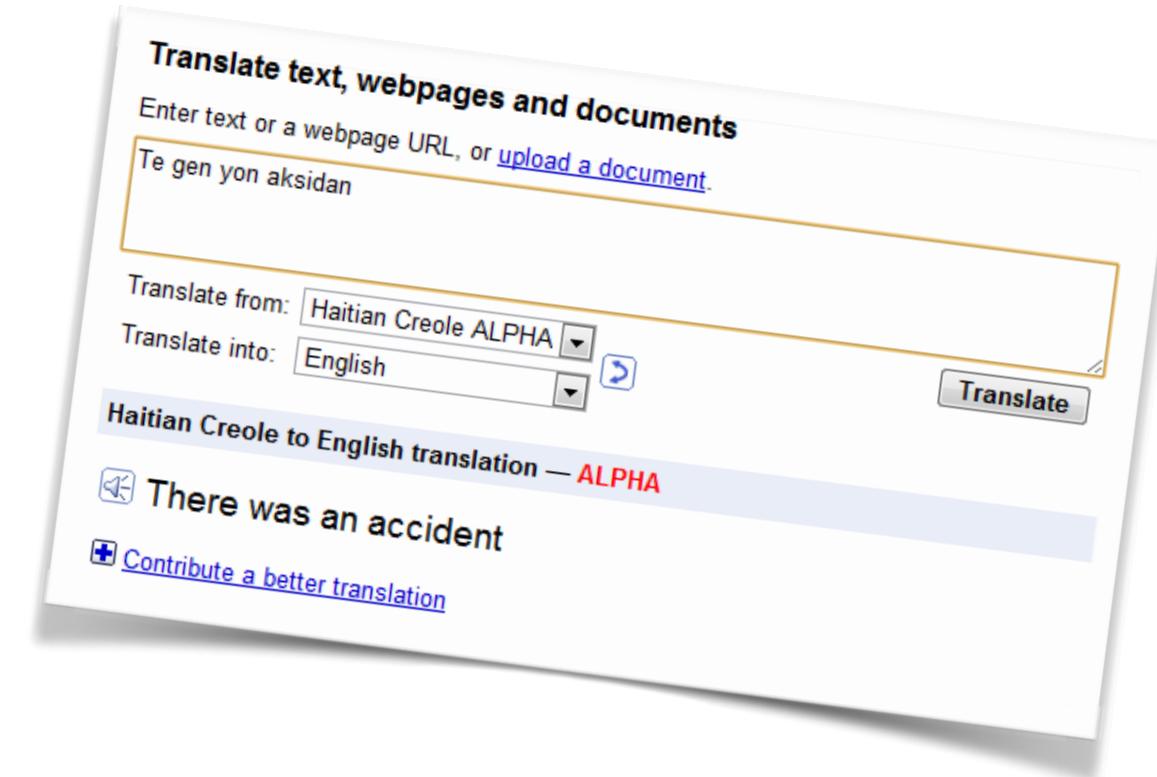
1 minute video showing 1954 MT:  
<https://youtu.be/K-HfpsHPmvw>

# History

- Started in the 1950s: rule-based, tightly linked to formal linguistics theories
- (late) 1980s to 2000s: Statistical MT
- 2000s-2014: Statistical Phrase-Based MT
- 2014-Present: Neural Machine Translation

# History

- Started in the 1950s: rule-based, tightly linked to formal linguistics theories
- 1980s: Statistical MT
- 2000s-2015: **Statistical Phrase-Based MT**
- 2015-Present: Neural Machine Translation



# Statistical MT

- Key Idea: Learn **probabilistic model** from **data**
- To find the best English sentence **y**, given French sentence **x**:

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

- Decompose using Bayes Rule:

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

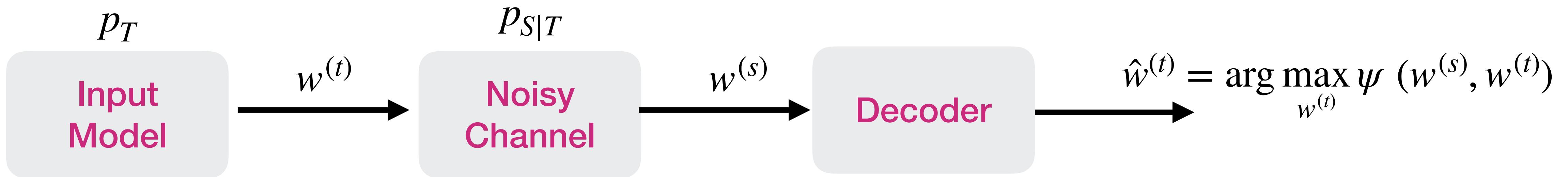
**Translation/Alignment Model**  
Models how words and phrases  
should be translated  
(*adequacy/fidelity*).

Learn from *parallel* data

**Language Model**  
Models how to write good  
English (*fluency*)

Learn from *monolingual* data

# Noisy channel model



$$\psi(w^{(s)}, w^{(t)}) = \psi_A(w^{(s)}, w^{(t)}) + \psi_F(w^{(t)})$$

$$\log P_{S,T}(w^{(s)}, w^{(t)}) = \log P_{S|T}(w^{(s)} | w^{(t)}) + \log P_T(w^{(t)})$$

- Generative process for source sentence
- Use Bayes rule to recover  $w^{(t)}$  that is maximally likely under the conditional distribution  $p_{T|S}$  (which is what we want)

# Data

- Statistical MT relies requires lots of **parallel corpora**

1. Chapter 4, Koch (DE)	de	es
<p><b>context</b> We would like to ensure that there is a reference to this <b>as early as the recitals</b> and that the period within which the Council has to make a decision - which is not clearly worded - is set at a maximum of three months .</p>	<p>Wir möchten sicherstellen , daß hierauf bereits in den Erwägungsgründen hingewiesen wird und die uneindeutig formulierte Frist , innerhalb der der Rat eine Entscheidung treffen muß , auf maximal drei Monate fixiert wird .</p>	<p>Quisiéramos asegurar que se aluda ya a esto en los considerandos y que el plazo , imprecisamente formulado , dentro del cual el Consejo ha de adoptar una decisión , se fije en tres meses como máximo .</p>
2. Chapter 3, FÄrm (SV)	de	es
<p><b>context</b> Our experience of modern administration tells us that openness , decentralisation of responsibility and qualified evaluation are often <b>as effective as detailed bureaucratic supervision</b> .</p>	<p>Unsere Erfahrungen mit moderner Verwaltung besagen , daß Transparenz , Dezentralisation der Verantwortlichkeiten und eine qualifizierte Auswertung oft ebenso effektiv sind wie bürokratische Detailkontrolle .</p>	<p>Nuestras experiencias en materia de administración moderna nos señalan que la apertura , la descentralización de las responsabilidades y las evaluaciones bien hechas son a menudo tan eficaces como los controles burocráticos detallados .</p>

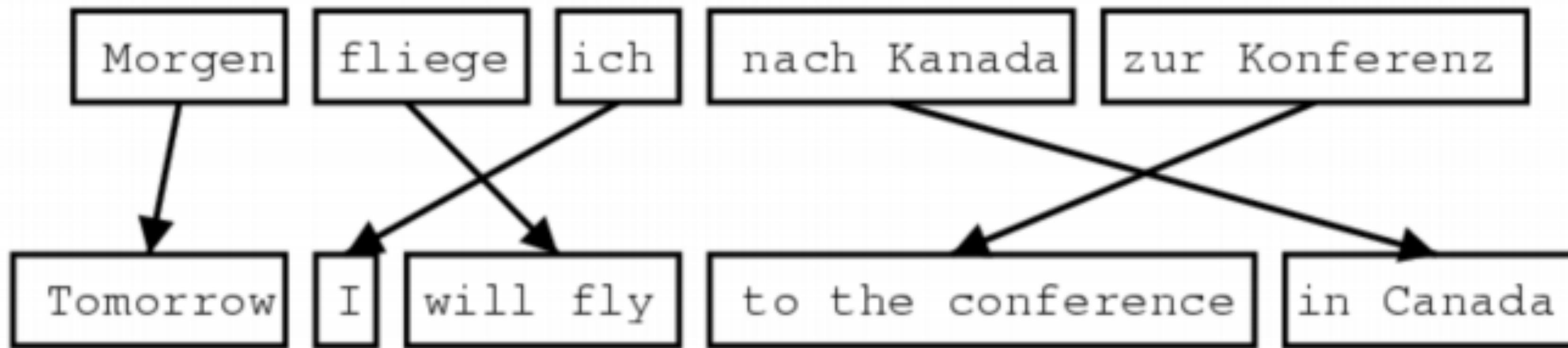
(Europarl, Koehn, 2005)

- Not available for many low-resource languages in the world

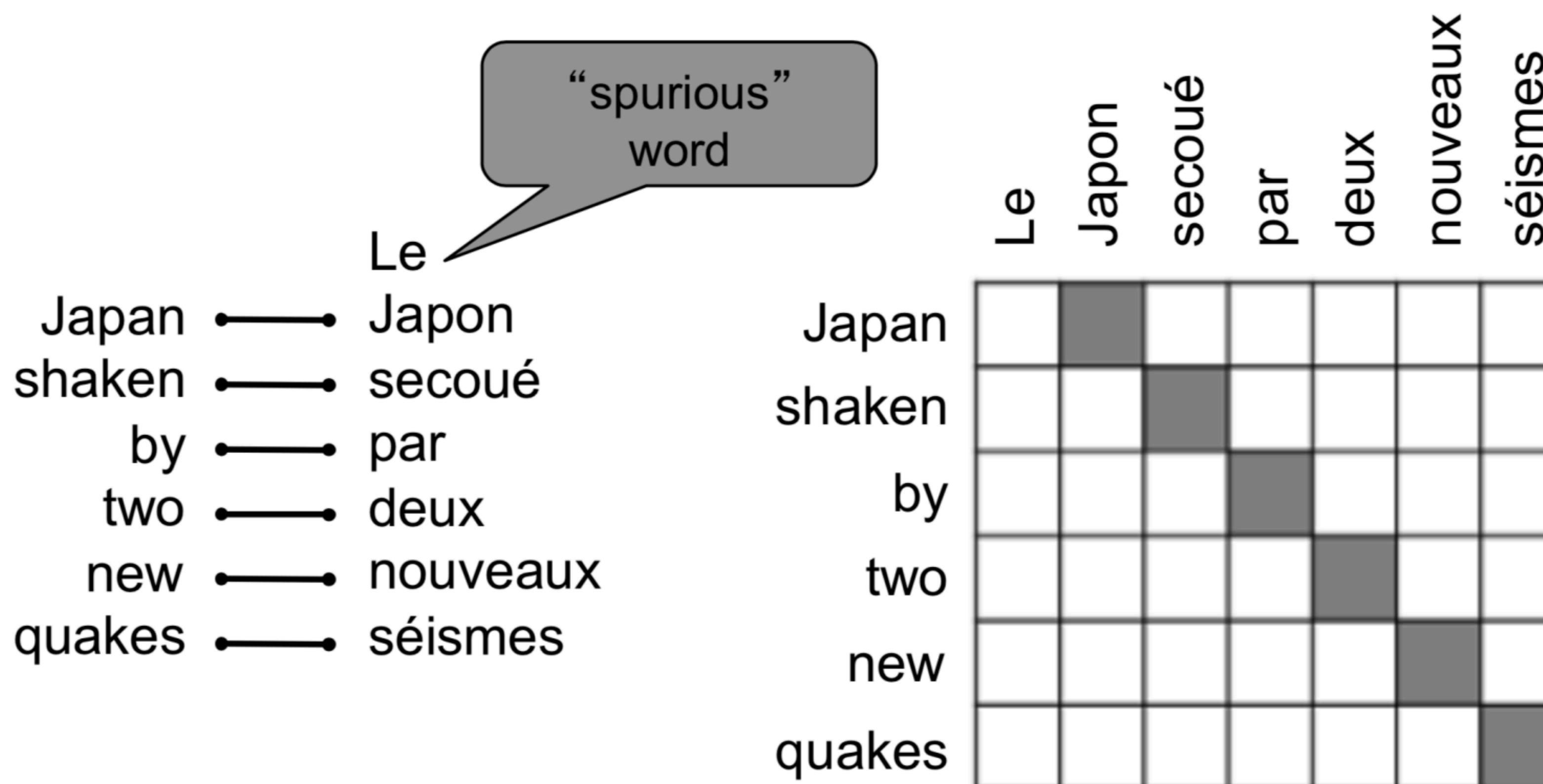
# How to define the translation model?

Introduce latent variable modeling the **alignment** (word-level correspondence) between the source sentence **x** and the target sentence **y**

$$P(x|y) = P(x, A|y)$$



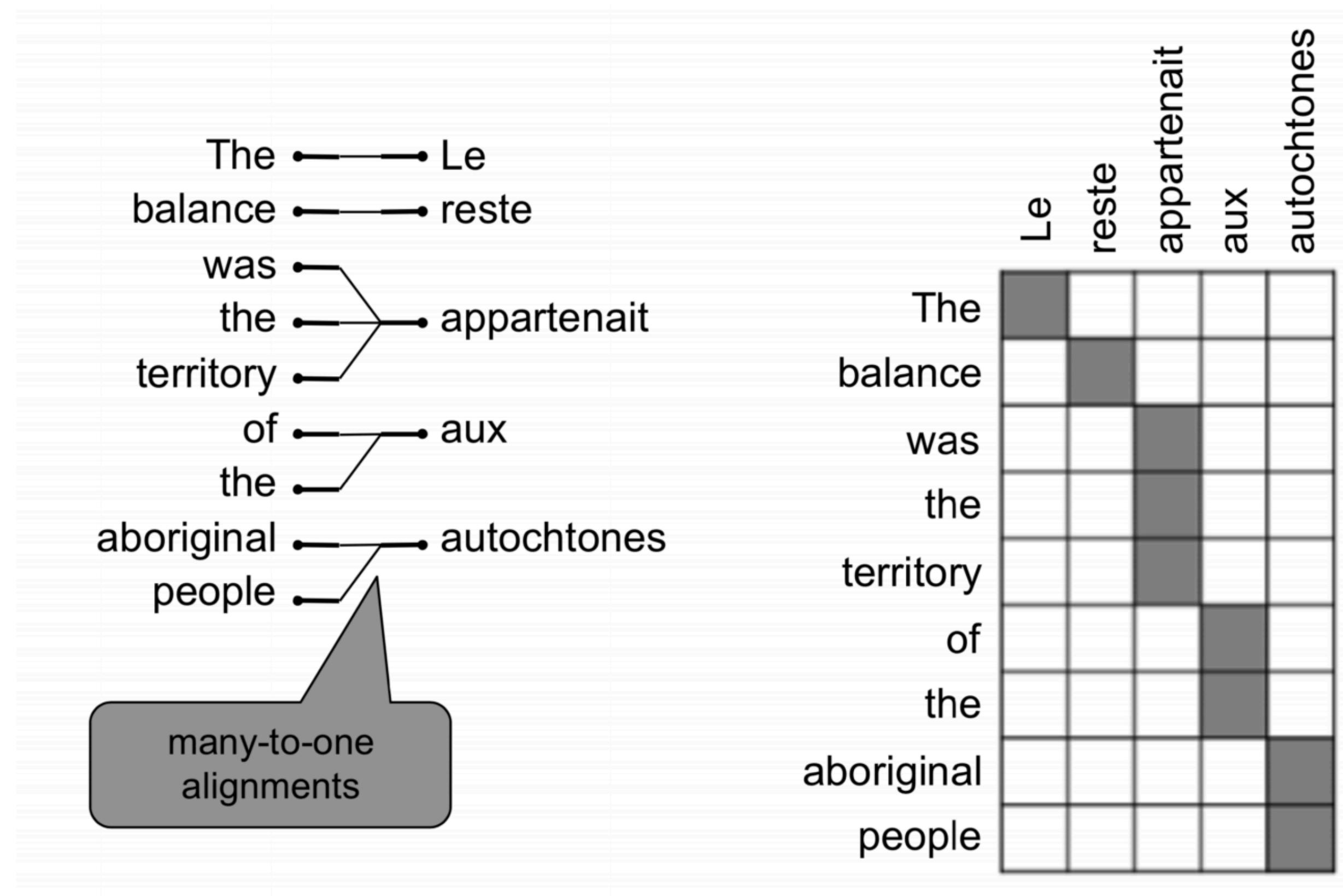
# What is alignment?



Examples from: “The Mathematics of Statistical Machine Translation: Parameter Estimation”, Brown et al, 1993. <http://www.aclweb.org/anthology/J93-2003>

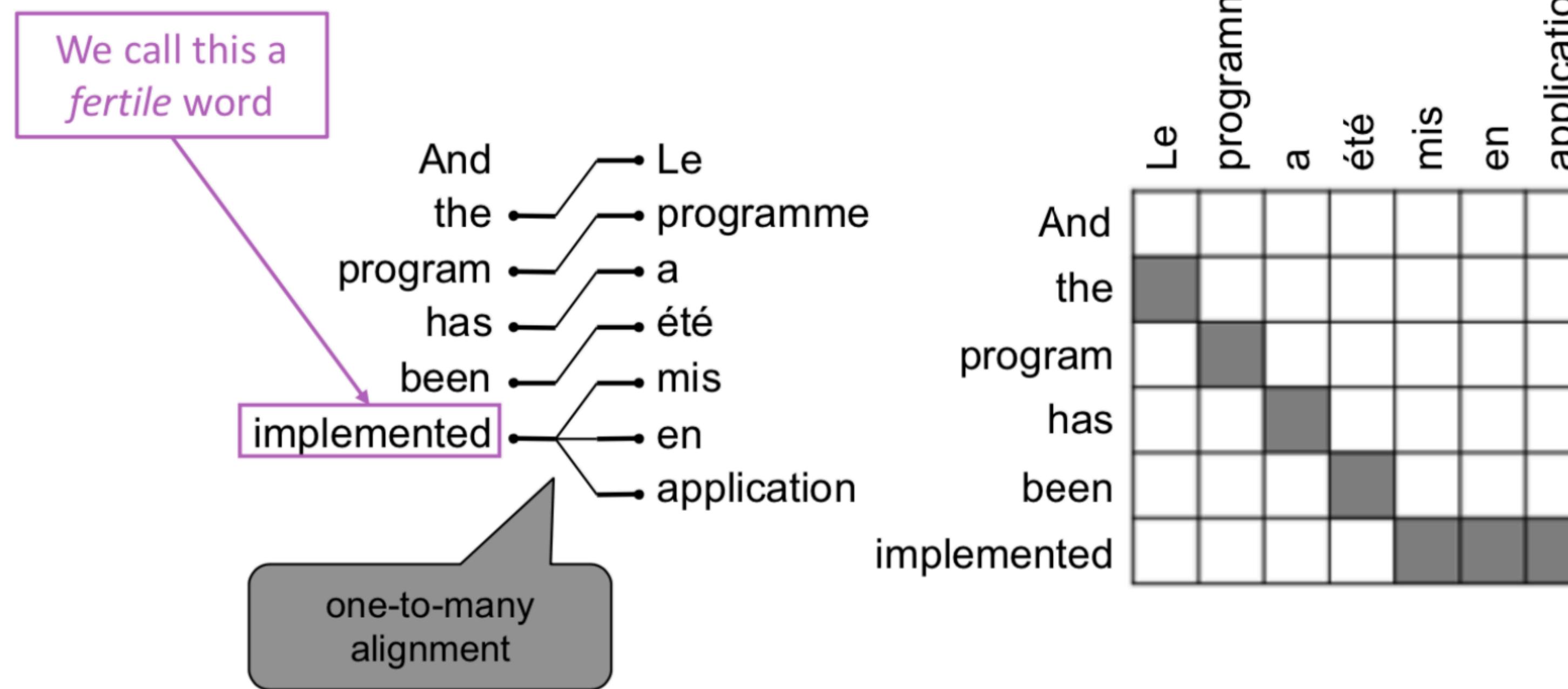
# Alignment is complex

Alignment can be many-to-one



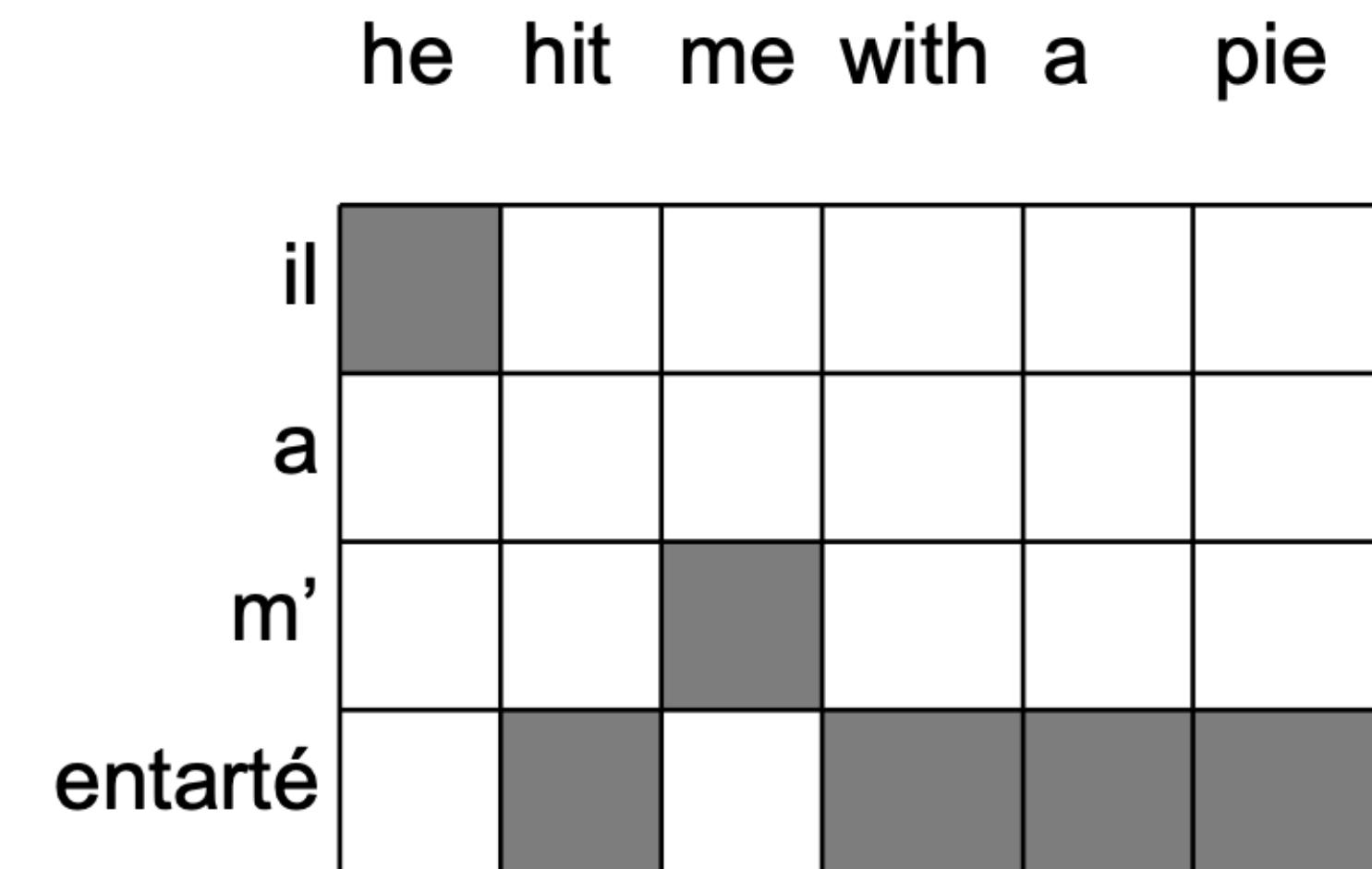
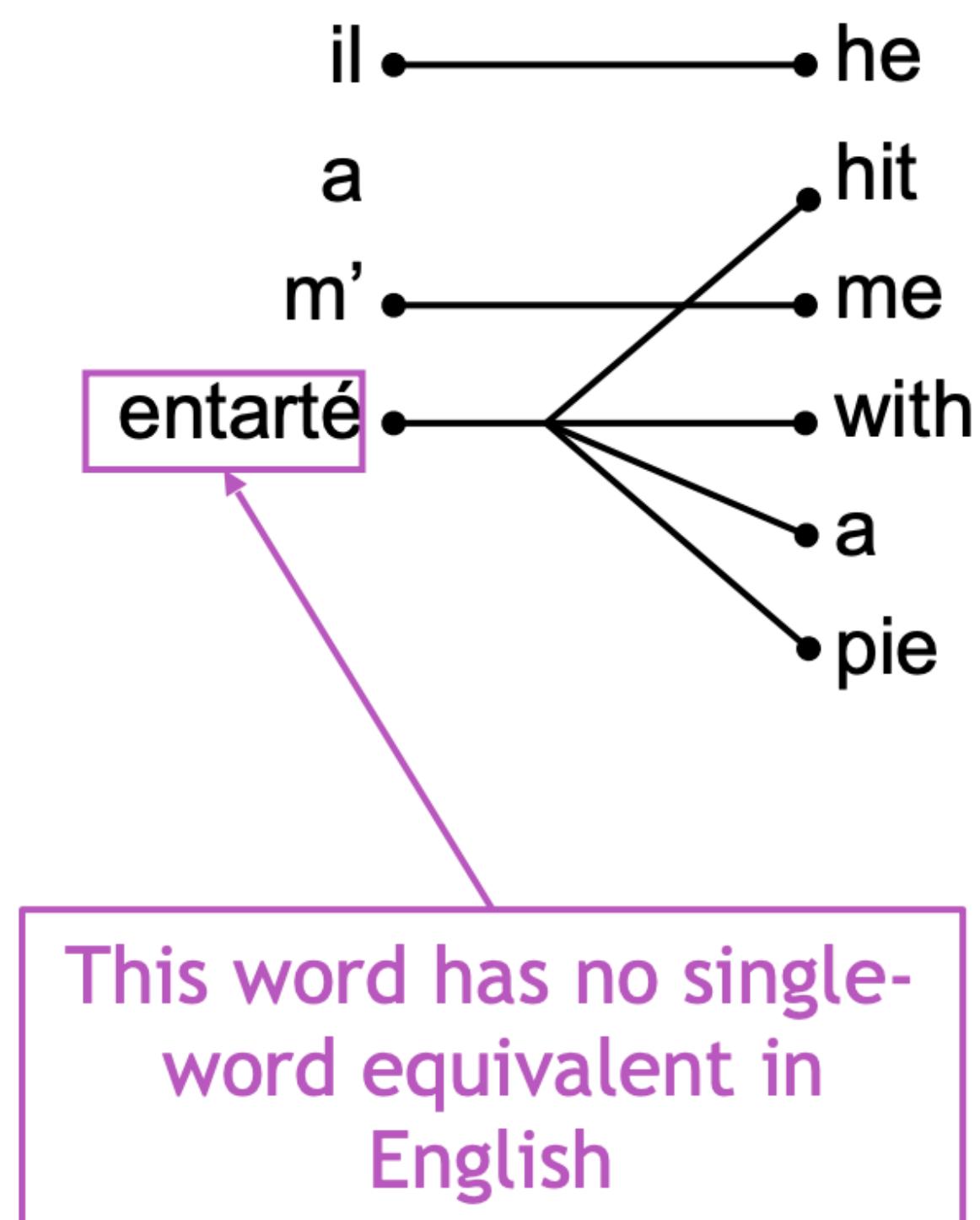
# Alignment is complex

Alignment can be **one-to-many**



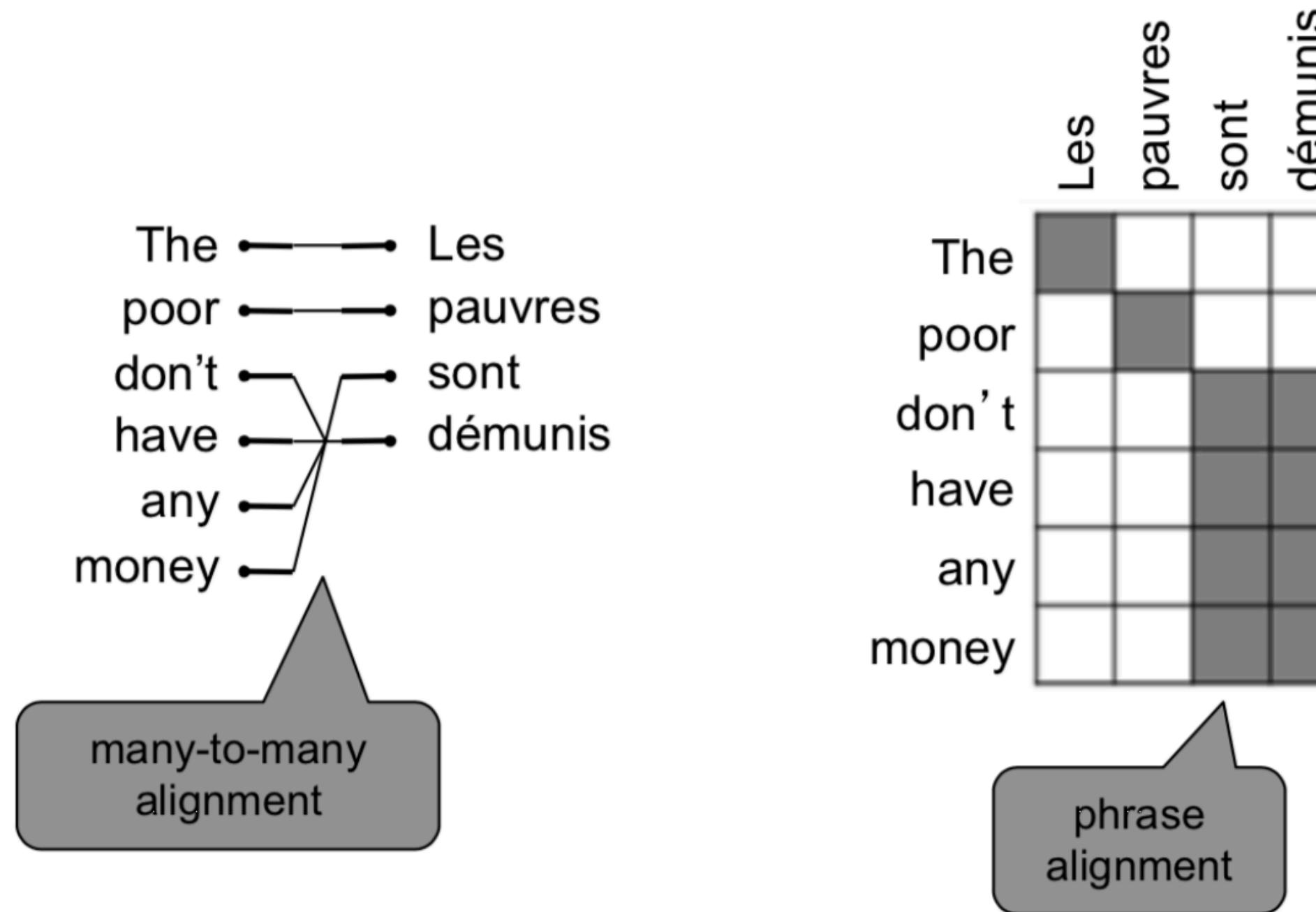
# Alignment is complex

Some words are very fertile!



# Alignment is complex

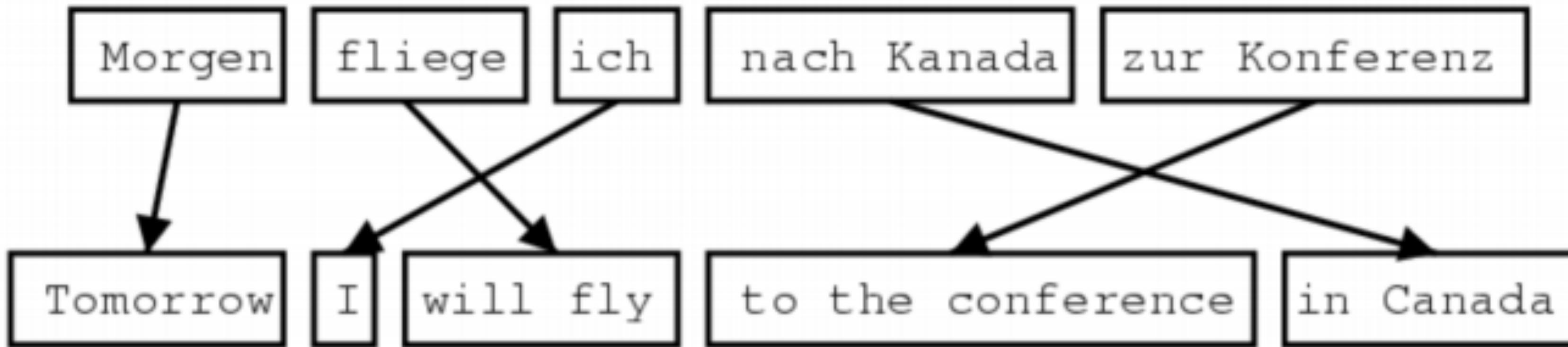
Alignment can be many-to-many (phrase-level)



# How to define the translation model?

Given the alignment, how do we incorporate in our model?

$$P(x|y) = P(x, A|y)$$



# Incorporating alignments

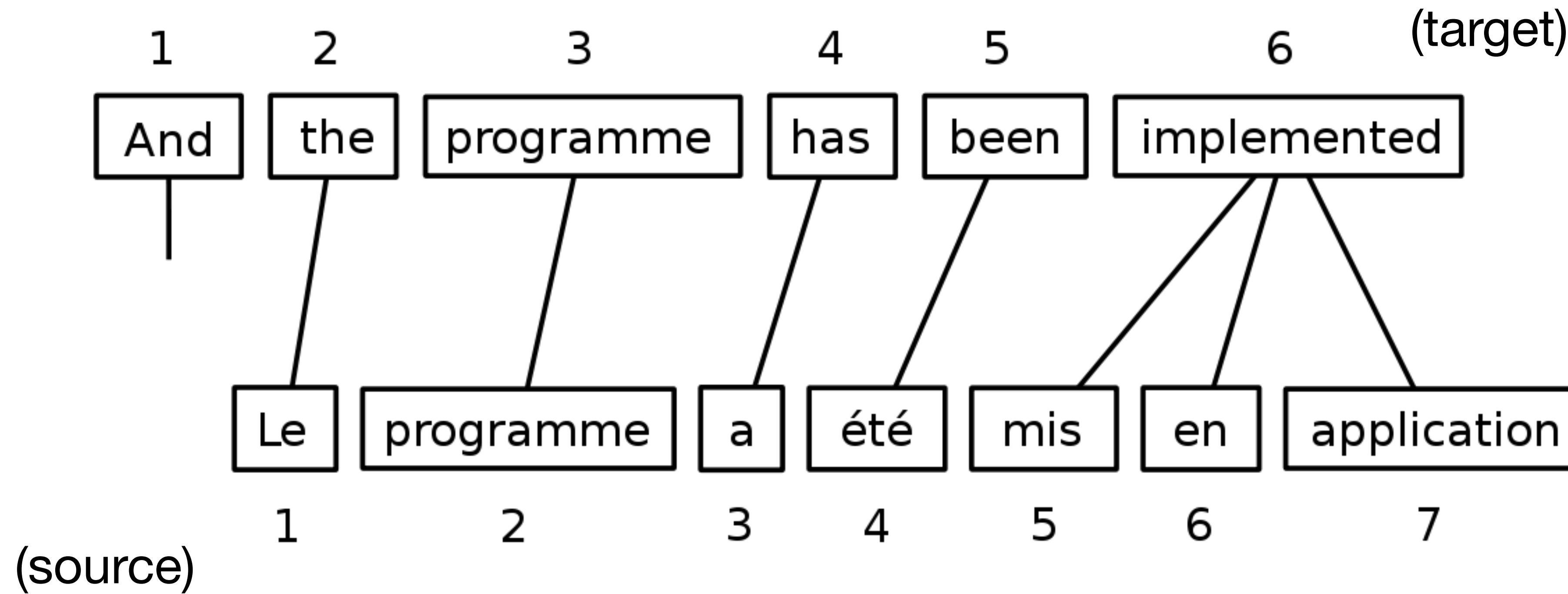
- Joint probability of alignment and translation can be defined as:

$$\begin{aligned} p(\mathbf{w}^{(s)}, \mathcal{A} \mid \mathbf{w}^{(t)}) &= \prod_{m=1}^{M^{(s)}} p(w_m^{(s)}, a_m \mid w_{a_m}^{(t)}, m, M^{(s)}, M^{(t)}) \\ &= \prod_{m=1}^{M^{(s)}} p(a_m \mid m, M^{(s)}, M^{(t)}) \times p(w_m^{(s)} \mid w_{a_m}^{(t)}). \end{aligned}$$

- $M^{(s)}, M^{(t)}$  are the number of words in source and target sentences
- $a_m$  is the alignment of the  $m^{th}$  word in the source sentence, i.e. it specifies that the  $m^{th}$  word is aligned to the  $a_m^{th}$  word in target

Is this sufficient?

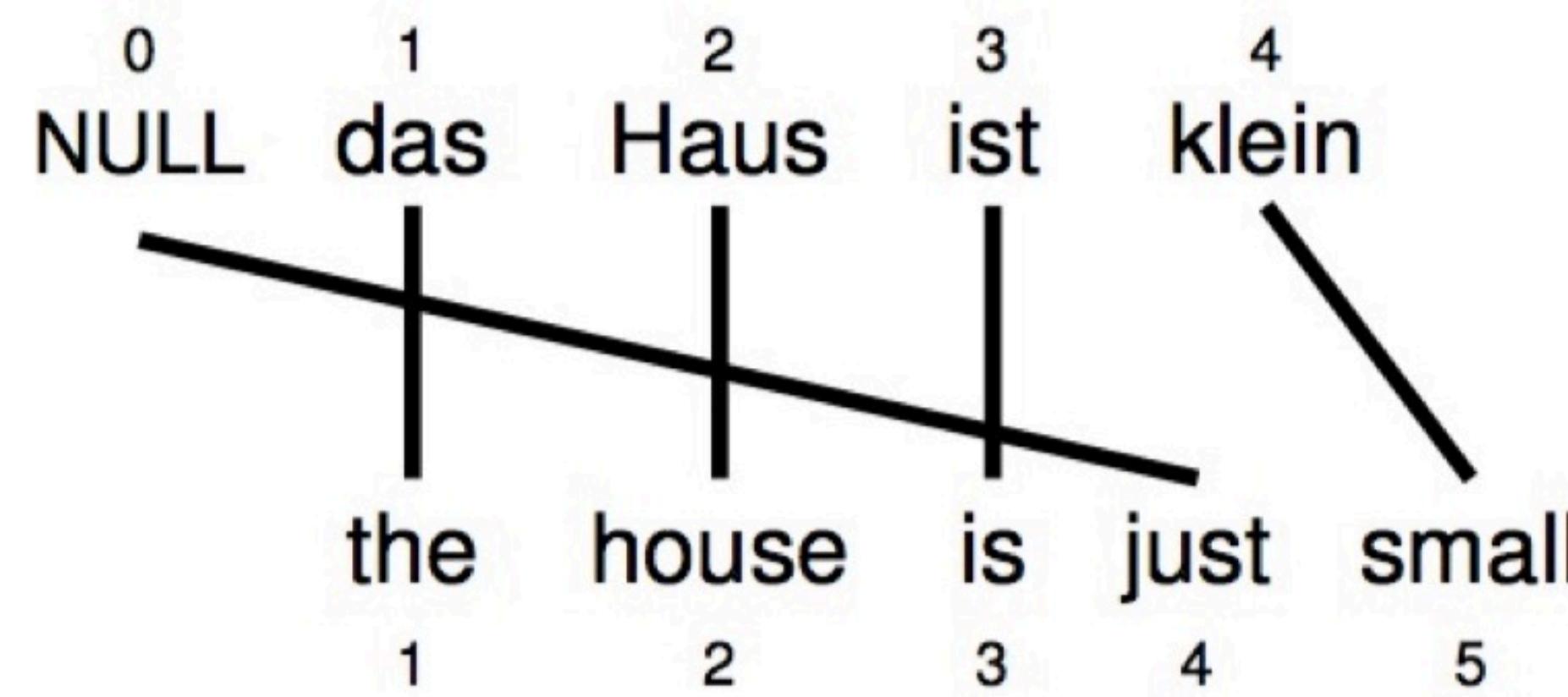
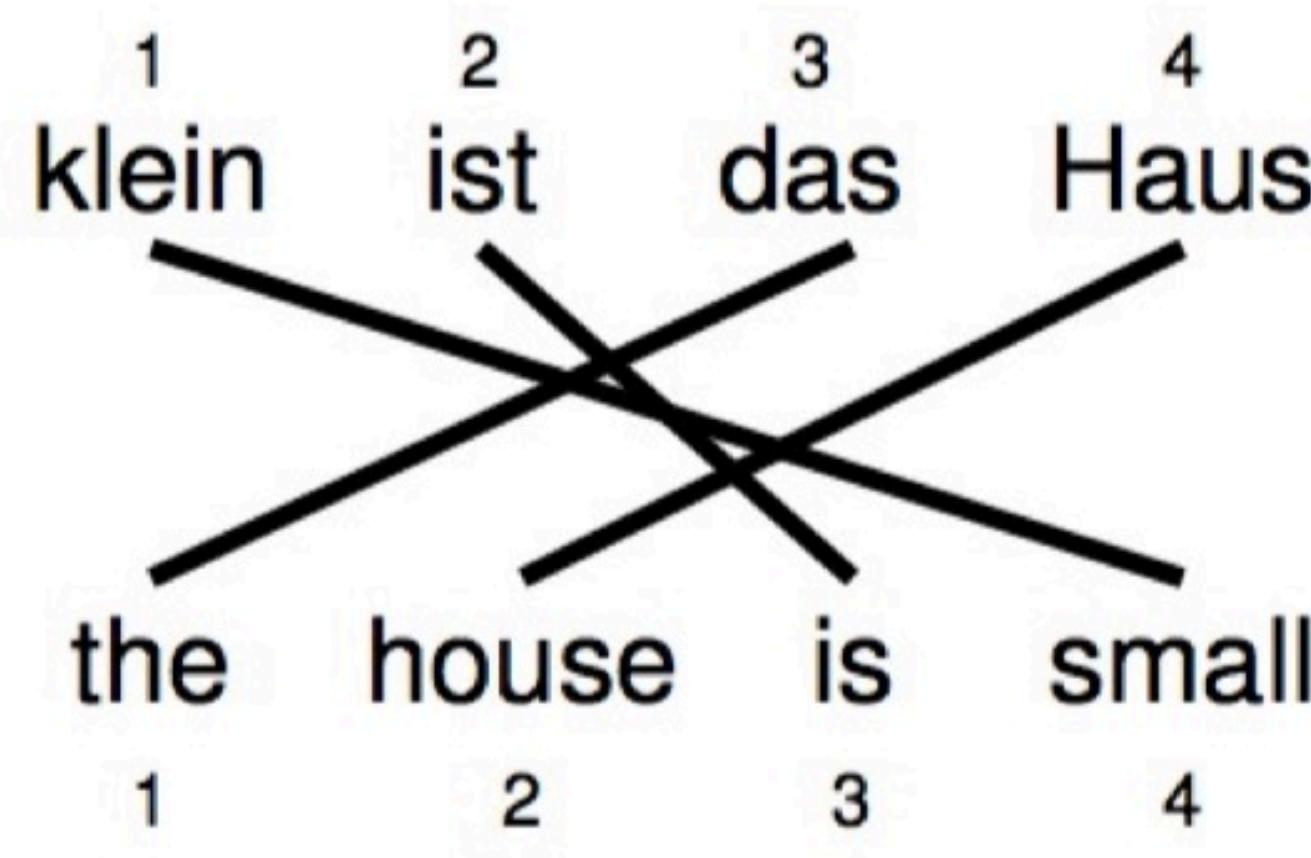
# Incorporating alignments



$$a_1 = 2, a_2 = 3, a_3 = 4, \dots$$

*Multiple source words may align to the same target word!*

# Reordering and word insertion



$$\mathbf{a} = (3, 4, 2, 1)^\top$$

$$\mathbf{a} = (1, 2, 3, 0, 4)^\top$$

Assume extra NULL token

# Independence assumptions

$$\begin{aligned} p(\mathbf{w}^{(s)}, \mathcal{A} \mid \mathbf{w}^{(t)}) &= \prod_{m=1}^{M^{(s)}} p(w_m^{(s)}, a_m \mid w_{a_m}^{(t)}, m, M^{(s)}, M^{(t)}) \\ &= \prod_{m=1}^{M^{(s)}} p(a_m \mid m, M^{(s)}, M^{(t)}) \times p(w_m^{(s)} \mid w_{a_m}^{(t)}). \end{aligned}$$

- Two independence assumptions:
  - Alignment probability factors across tokens:

$$p(\mathcal{A} \mid \mathbf{w}^{(s)}, \mathbf{w}^{(t)}) = \prod_{m=1}^{M^{(s)}} p(a_m \mid m, M^{(s)}, M^{(t)}).$$

- Translation probability factors across tokens:

$$p(\mathbf{w}^{(s)} \mid \mathbf{w}^{(t)}, \mathcal{A}) = \prod_{m=1}^{M^{(s)}} p(w_m^{(s)} \mid w_{a_m}^{(t)}),$$

# How do we translate?

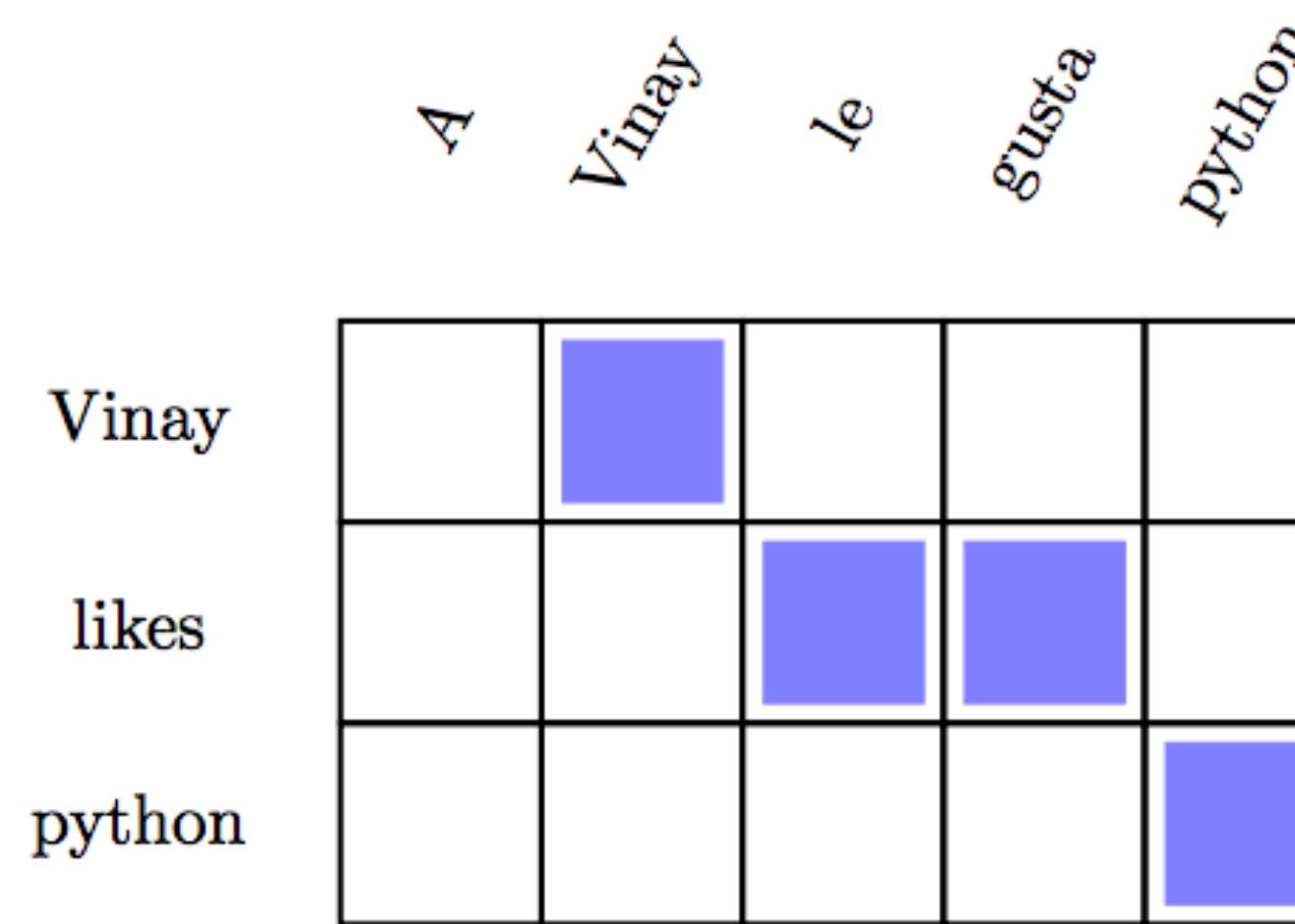
- We want:  $\arg \max_{w^{(t)}} p(w^{(t)} | w^{(s)}) = \arg \max_{w^{(t)}} \frac{p(w^{(s)}, w^{(t)})}{p(w^{(s)})}$
- Sum over all possible alignments:

$$\begin{aligned} p(w^{(s)}, w^{(t)}) &= \sum_{\mathcal{A}} p(w^{(s)}, w^{(t)}, \mathcal{A}) \\ &= p(w^{(t)}) \sum_{\mathcal{A}} p(\mathcal{A}) \times p(w^{(s)} | w^{(t)}, \mathcal{A}) \end{aligned}$$

- Alternatively, take the max over alignments

# Alignments

- **Key question:** How should we align words in source to words in target?

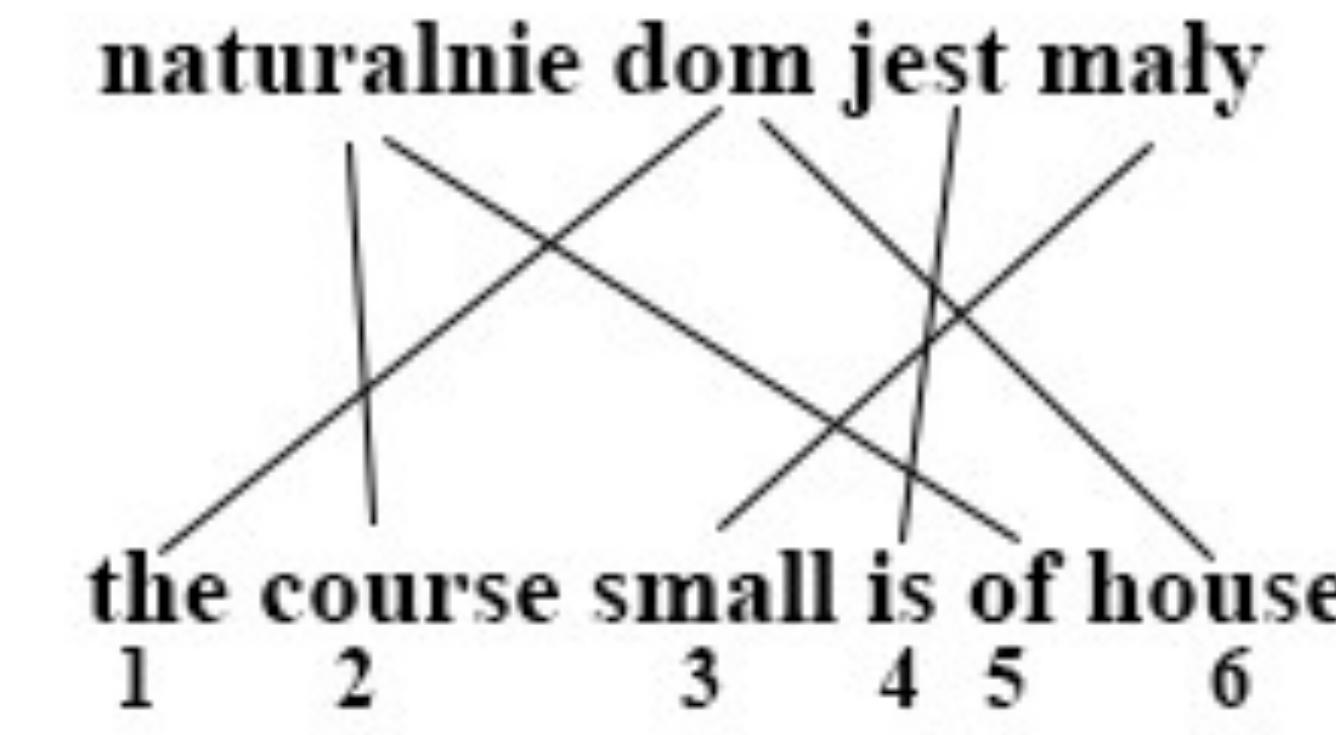


**good**  $\mathcal{A}(\mathbf{w}^{(s)}, \mathbf{w}^{(t)}) = \{(A, \emptyset), (Vinay, Vinay), (le, likes), (gusta, likes), (Python, Python)\}.$

**bad**  $\mathcal{A}(\mathbf{w}^{(s)}, \mathbf{w}^{(t)}) = \{(A, Vinay), (Vinay, likes), (le, Python), (gusta, \emptyset), (Python, \emptyset)\}.$

# IBM Model I

- Assume  $p(a_m | m, M^{(s)}, M^{(t)}) = \frac{1}{M^{(t)}}$
- Is this a good assumption?



Every alignment is equally likely!

# IBM Model I

- Each source word is aligned to at most one target word
- Further, assume  $p(a_m | m, M^{(s)}, M^{(t)}) = \frac{1}{M^{(t)}}$
- We then have:
$$p(w^{(s)}, w^{(t)}) = p(w^{(t)}) \sum_A \left(\frac{1}{M^{(t)}}\right)^{M^{(s)}} p(w^{(s)} | w^{(t)})$$
- How do we estimate  $p(w^{(s)} = v | w^{(t)} = u)$  ?

# IBM Model I

- If we had word-to-word alignments, we could compute the probabilities using the MLE:
- $p(v | u) = \frac{count(u, v)}{count(u)}$
- where  $count(u, v)$  = #instances where word  $u$  was aligned to word  $v$  in the training set
- However, word-to-word alignments are often hard to come by

What can we do?

# EM for Model I

- (E-Step) If we had an accurate translation model, we can estimate likelihood of each alignment as:

$$q_m(a_m \mid \mathbf{w}^{(s)}, \mathbf{w}^{(t)}) \propto p(a_m \mid m, M^{(s)}, M^{(t)}) \times p(w_m^{(s)} \mid w_{a_m}^{(t)}),$$

- (M Step) Use expected count to re-estimate translation parameters:

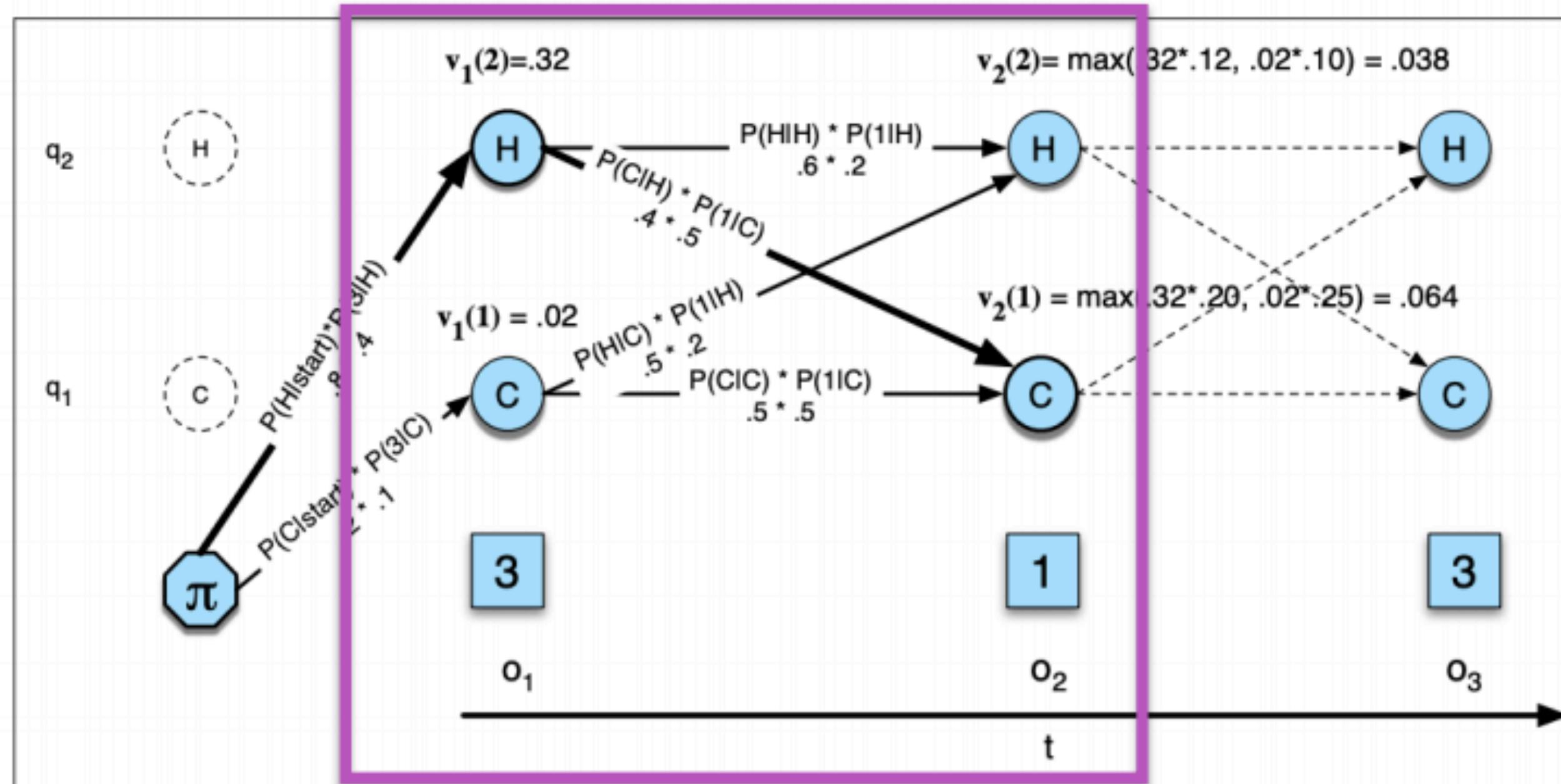
$$p(v \mid u) = \frac{E_q[\text{count}(u, v)]}{\text{count}(u)}$$

$$E_q [\text{count}(u, v)] = \sum_m q_m(a_m \mid \mathbf{w}^{(s)}, \mathbf{w}^{(t)}) \times \delta(w_m^{(s)} = v) \times \delta(w_{a_m}^{(t)} = u).$$

# Independence assumptions allow for Viterbi decoding

Impose strong **independence assumptions** in model:

$$p(x, a|y) = \prod_{j=1}^{l_x} p(x_j | f_{a(j)})$$



Source: "Speech and Language Processing", Chapter A, Jurafsky and Martin, 2019.

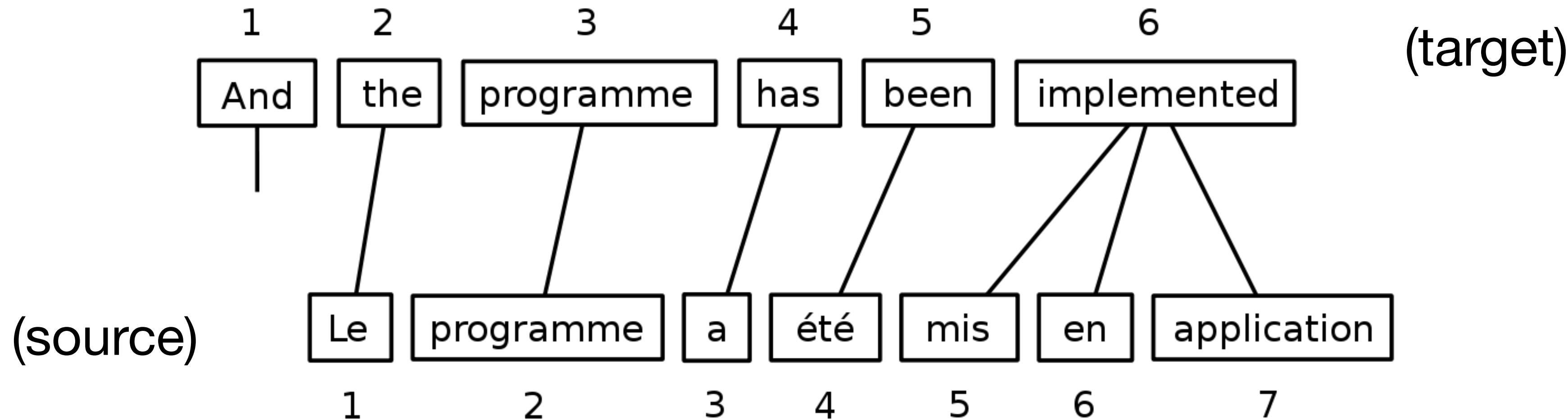
In general, use greedy or beam decoding

# Model I: Decoding

- Pick target sentence length  $M^{(t)}$
- Decode:  $\arg \max_{w^{(t)}} p(w^{(t)} | w^{(s)}) = \arg \max_{w^{(t)}} p(w^{(s)}, w^{(t)})$
- $p(w^{(s)}, w^{(t)}) = p(w^{(t)}) \sum_A \frac{1}{M^{(t)}} p(w^{(s)} | w^{(t)})$

$$\begin{aligned} p(\mathbf{w}^{(s)}, \mathcal{A} | \mathbf{w}^{(t)}) &= \prod_{m=1}^{M^{(s)}} p(w_m^{(s)}, a_m | w_{a_m}^{(t)}, m, M^{(s)}, M^{(t)}) \\ &= \prod_{m=1}^{M^{(s)}} p(a_m | m, M^{(s)}, M^{(t)}) \times p(w_m^{(s)} | w_{a_m}^{(t)}). \end{aligned}$$

# Model I: Decoding



At every step  $m$ , pick target word to maximize product of:

1. Language model:  
 $p_{LM}(w_m^{(t)} | w_{<m}^{(t)})$
  2. Translation model:  
 $p(w_{b_m}^{(s)} | w_m^{(t)})$

where  $b_m$  is the inverse alignment from target to source

# IBM Model 2

- Slightly relaxed assumption:
  - $p(a_m | m, M^{(s)}, M^{(t)})$  is also estimated, not set to constant
  - Original independence assumptions still required:
    - Alignment probability factors across tokens:

$$p(\mathcal{A} | \mathbf{w}^{(s)}, \mathbf{w}^{(t)}) = \prod_{m=1}^{M^{(s)}} p(a_m | m, M^{(s)}, M^{(t)}).$$

- Translation probability factors across tokens:

$$p(\mathbf{w}^{(s)} | \mathbf{w}^{(t)}, \mathcal{A}) = \prod_{m=1}^{M^{(s)}} p(w_m^{(s)} | w_{a_m}^{(t)}),$$

# Other IBM models

Model 1: lexical translation

Model 2: additional absolute alignment model

Model 3: extra fertility model

Model 4: added relative alignment model

Model 5: fixed deficiency problem.

Model 6: Model 4 combined with a [HMM](#) alignment model in a log linear way

- Models 3 - 6 make successively weaker assumptions
  - But get progressively harder to optimize
  - Simpler models are often used to ‘initialize’ complex ones
    - e.g train Model 1 and use it to initialize Model 2 parameters

# Phrase-based MT

- Word-by-word translation is not sufficient in many cases

(literal)

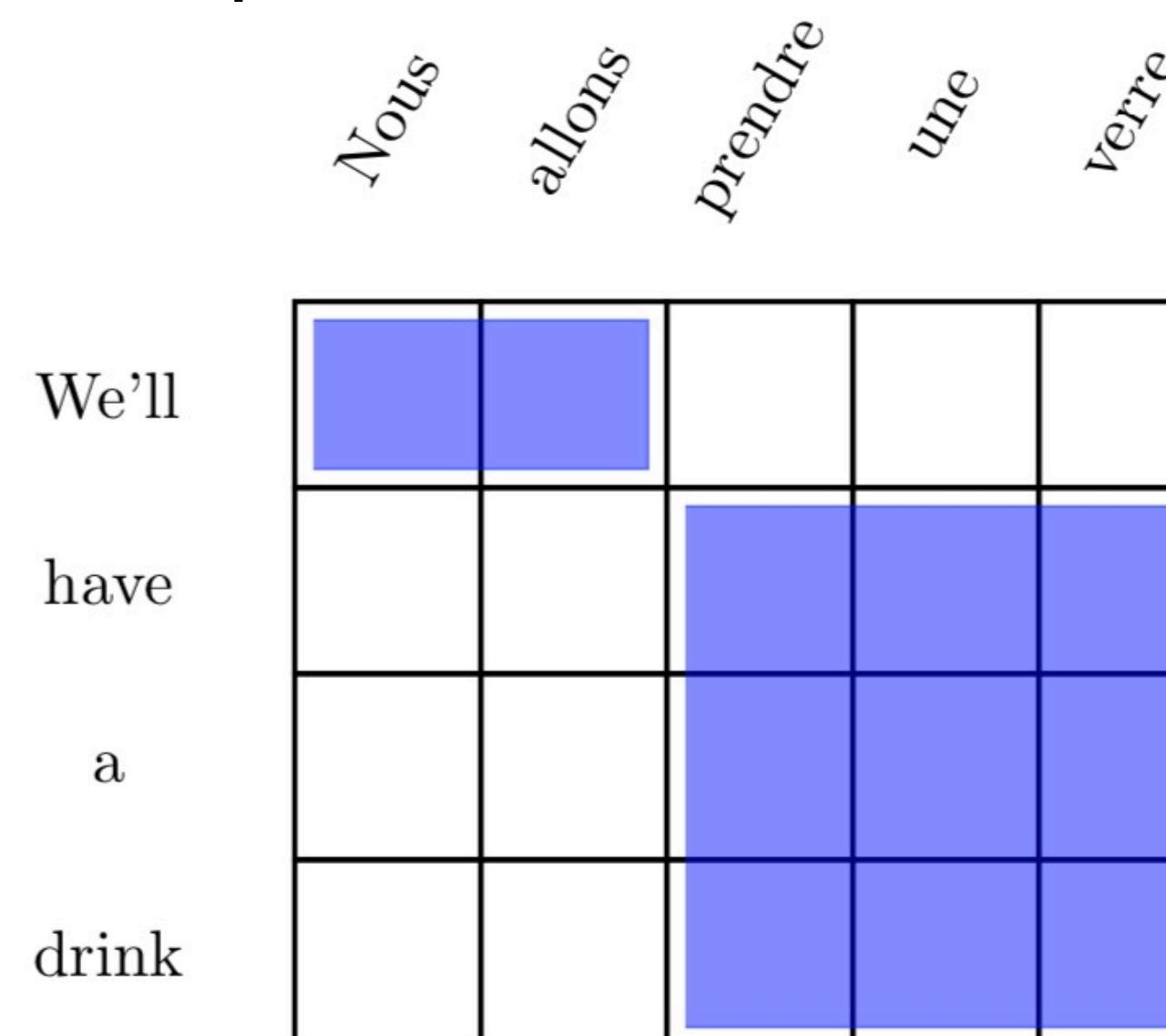
*Nous allons prendre un verre*

(actual)

We will take a glass

We'll have a drink

- Solution: build alignments and translation tables between multiword spans or “phrases”



# Phrase-based MT

- Solution: build alignments and translation tables between multiword spans or “phrases”
- Translations condition on multi-word units and assign probabilities to multi-word units
- Alignments map from spans to spans

$$p(\mathbf{w}^{(s)} \mid \mathbf{w}^{(t)}, \mathcal{A}) = \prod_{((i,j),(k,\ell)) \in \mathcal{A}} p_{w^{(s)}|w^{(t)}}(\{w_{i+1}^{(s)}, w_{i+2}^{(s)}, \dots, w_j^{(s)}\} \mid \{w_{k+1}^{(t)}, w_{k+2}^{(t)}, \dots, w_\ell^{(t)}\})$$

# Syntactic MT

- ▶ Rather than use phrases, use a *synchronous context-free grammar*: constructs “parallel” trees in two languages simultaneously

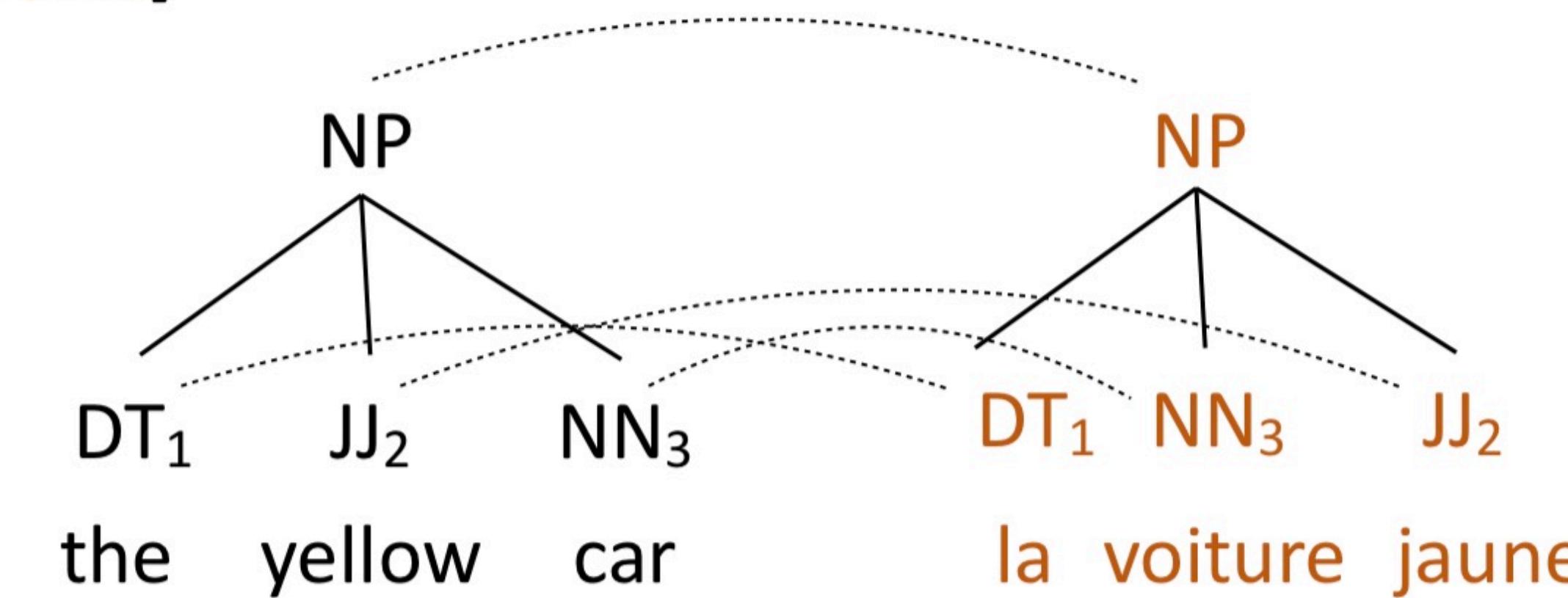
$NP \rightarrow [DT_1\ JJ_2\ NN_3; DT_1\ NN_3\ JJ_2]$

$DT \rightarrow [\text{the}, \text{la}]$

$DT \rightarrow [\text{the}, \text{le}]$

$NN \rightarrow [\text{car}, \text{voiture}]$

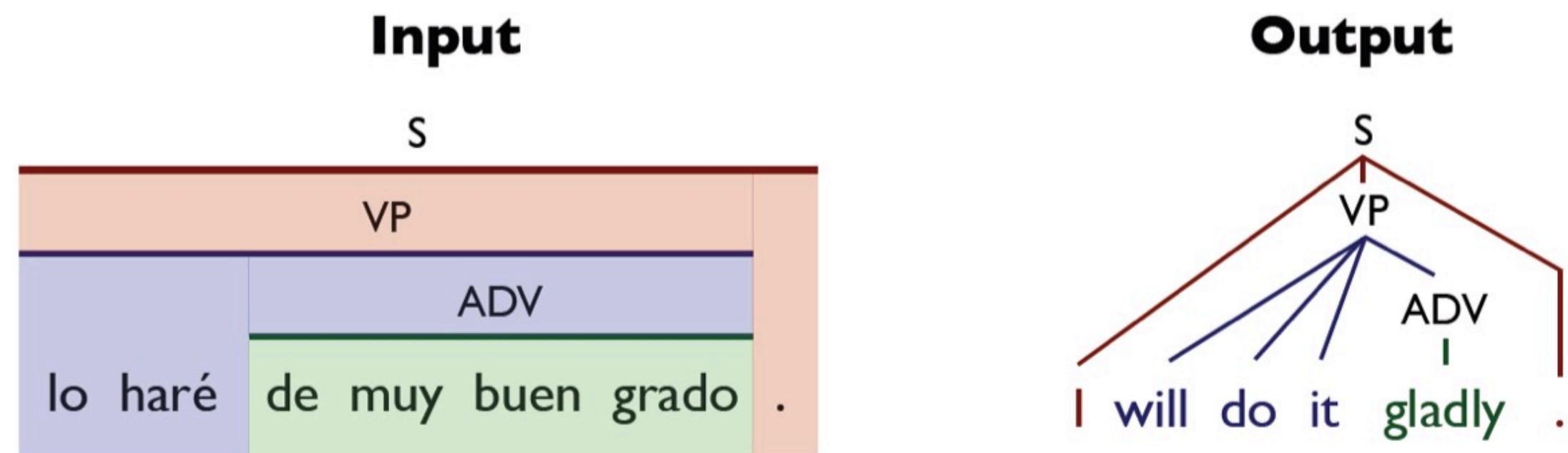
$JJ \rightarrow [\text{yellow}, \text{jaune}]$



- ▶ Assumes parallel syntax up to reordering
- ▶ Translation = parse the input with “half” the grammar, read off other half

(Slide credit: Greg Durrett)

# Syntactic MT

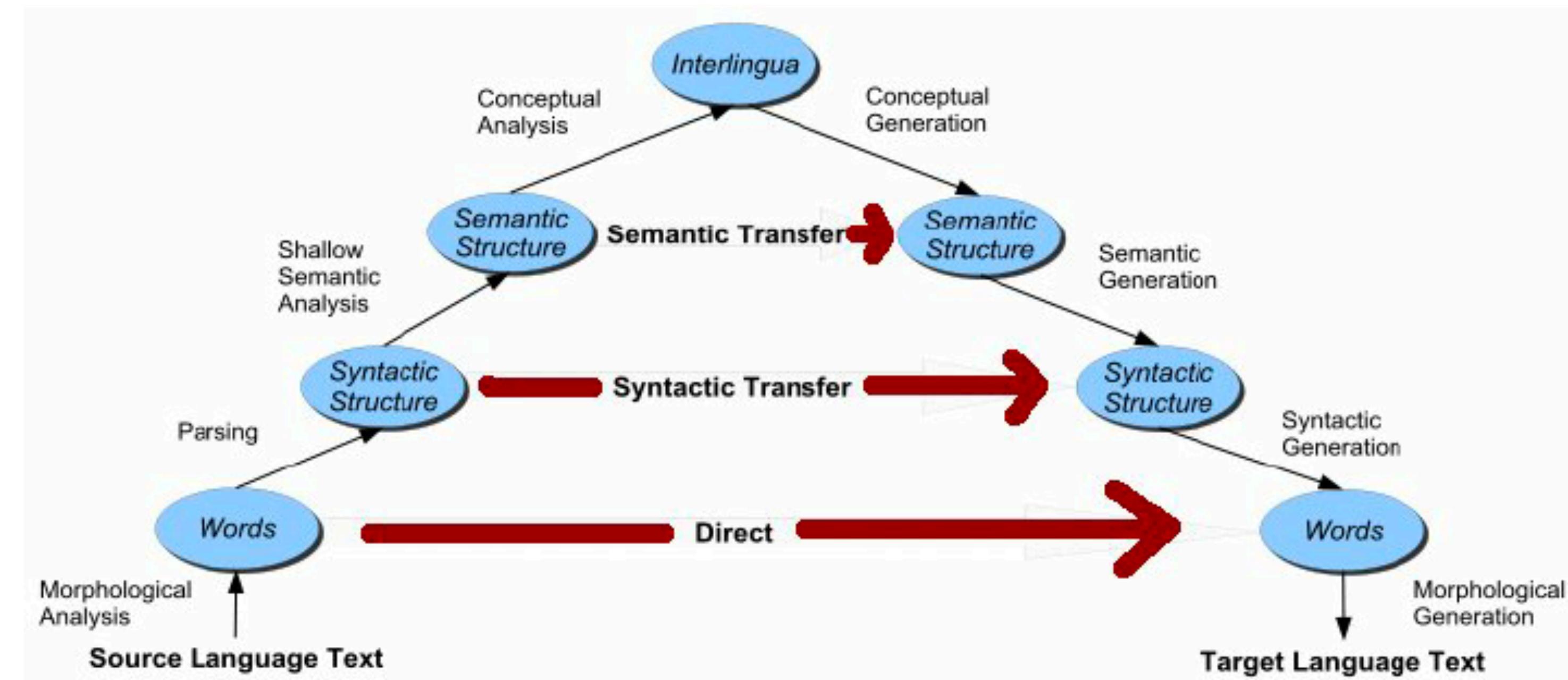


- Relax this by using lexicalized rules, like “syntactic phrases”
- Leads to HUGE grammars, parsing is slow

s → ⟨ VP . ; I VP . ⟩ OR s → ⟨ VP . ; you VP . ⟩  
VP → ⟨ lo haré ADV ; will do it ADV ⟩  
s → ⟨ lo haré ADV . ; I will do it ADV . ⟩  
ADV → ⟨ de muy buen grado ; gladly ⟩

Slide credit: Dan Klein

# Vauquois Pyramid



- Hierarchy of concepts and distances between them in different languages
- Lowest level: individual words/characters
- Higher levels: syntax, semantics
- Interlingua: Generic language-agnostic representation of meaning

# Statistical MT

- 1990s to 2010s: huge research area
- Extremely complex systems
  - Many separately-designed subcomponents
  - Lots of feature engineering
  - Needed to compile/maintain extra resources (phrase tables)
- Lots of human effort to maintain

# History

- Started in the 1950s: rule-based, tightly linked to formal linguistics theories
- 1980s: Statistical MT
- 2000s-2015: Statistical Phrase-Based MT
- 2015-Present: Neural Machine Translation
- ~2018-Present: Neural Machine Translation + PBMT Hybrid