

# CS388: Natural Language Processing

## Lecture 16: Machine Translation 1

Greg Durrett



Some slides adapted from Dan Klein, UC Berkeley



Star Wars The Third Gathers: The Backstroke of the West  
(subtitles machine translated from Chinese)



# Administrivia

---

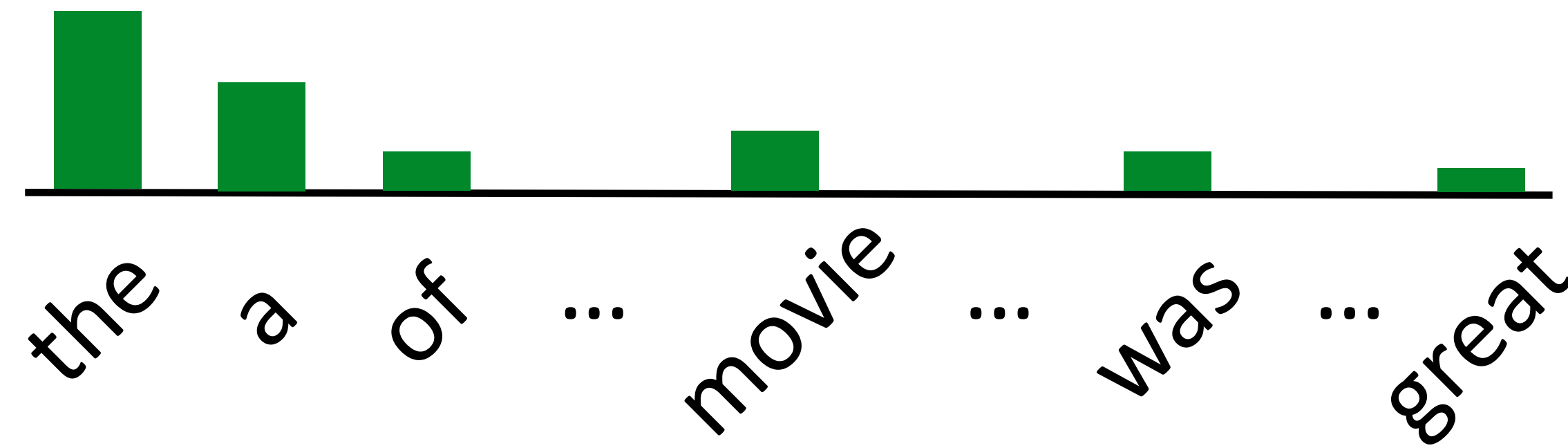
- ▶ Project 2 due in a week



# Recall: Pointer Networks

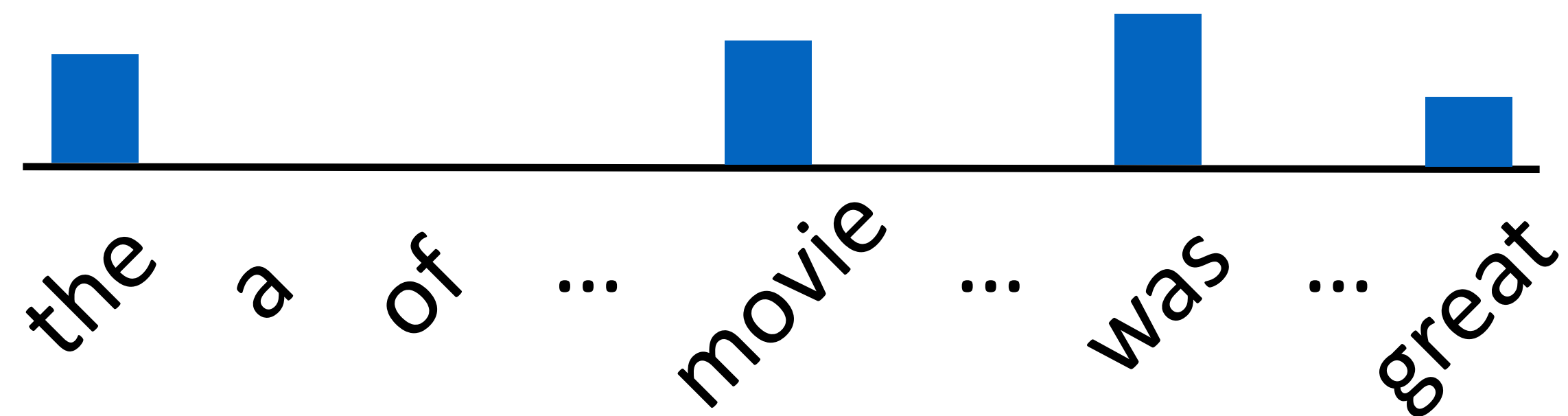
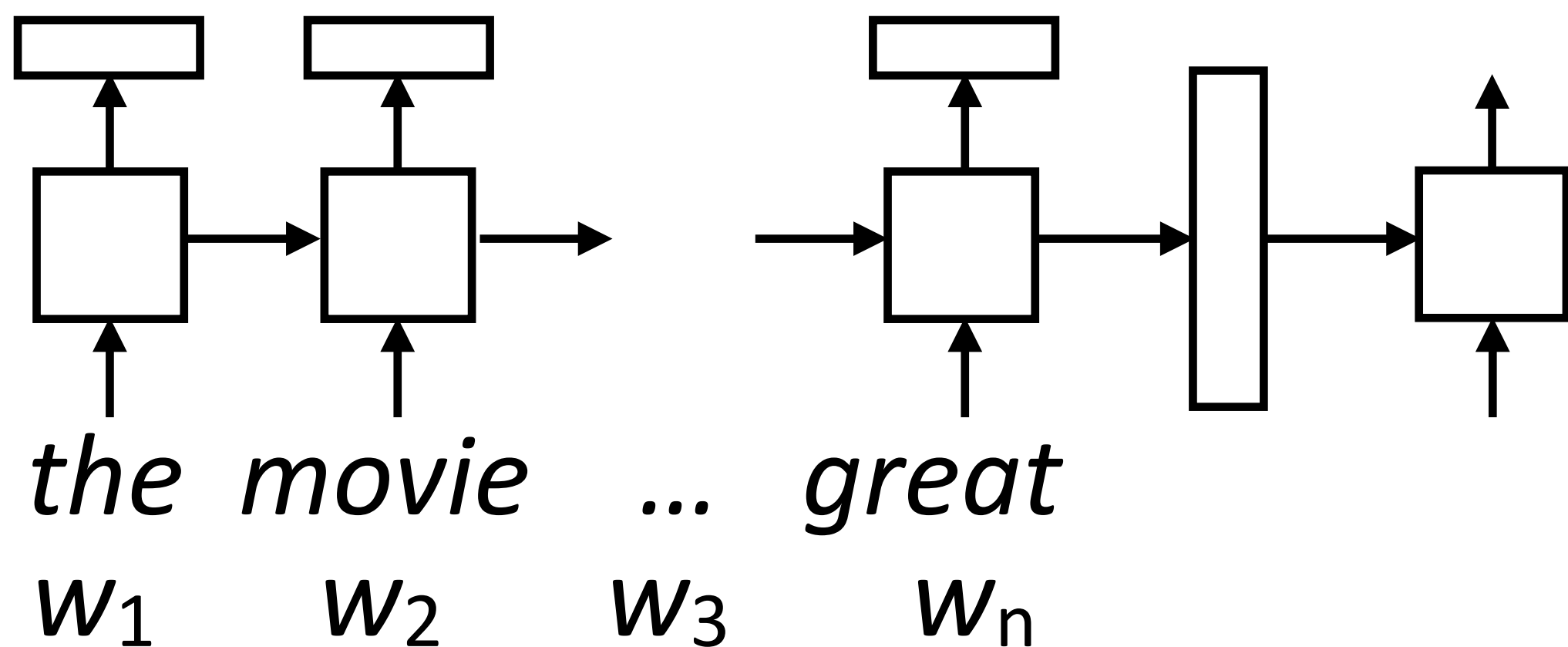
$$P(y_i | \mathbf{x}, y_1, \dots, y_{i-1}) = \text{softmax}(W[c_i; \bar{h}_i])$$

- ▶ Standard decoder ( $P_{\text{vocab}}$ ): softmax over vocabulary, all words get  $>0$  prob



- ▶ Pointer network: predict from *source words* instead of *target vocab*

$$P_{\text{pointer}}(y_i | \mathbf{x}, y_1, \dots, y_{i-1}) \propto \begin{cases} \exp(h_j^\top V \bar{h}_i) & \text{if } y_i = w_j \\ 0 & \text{otherwise} \end{cases}$$



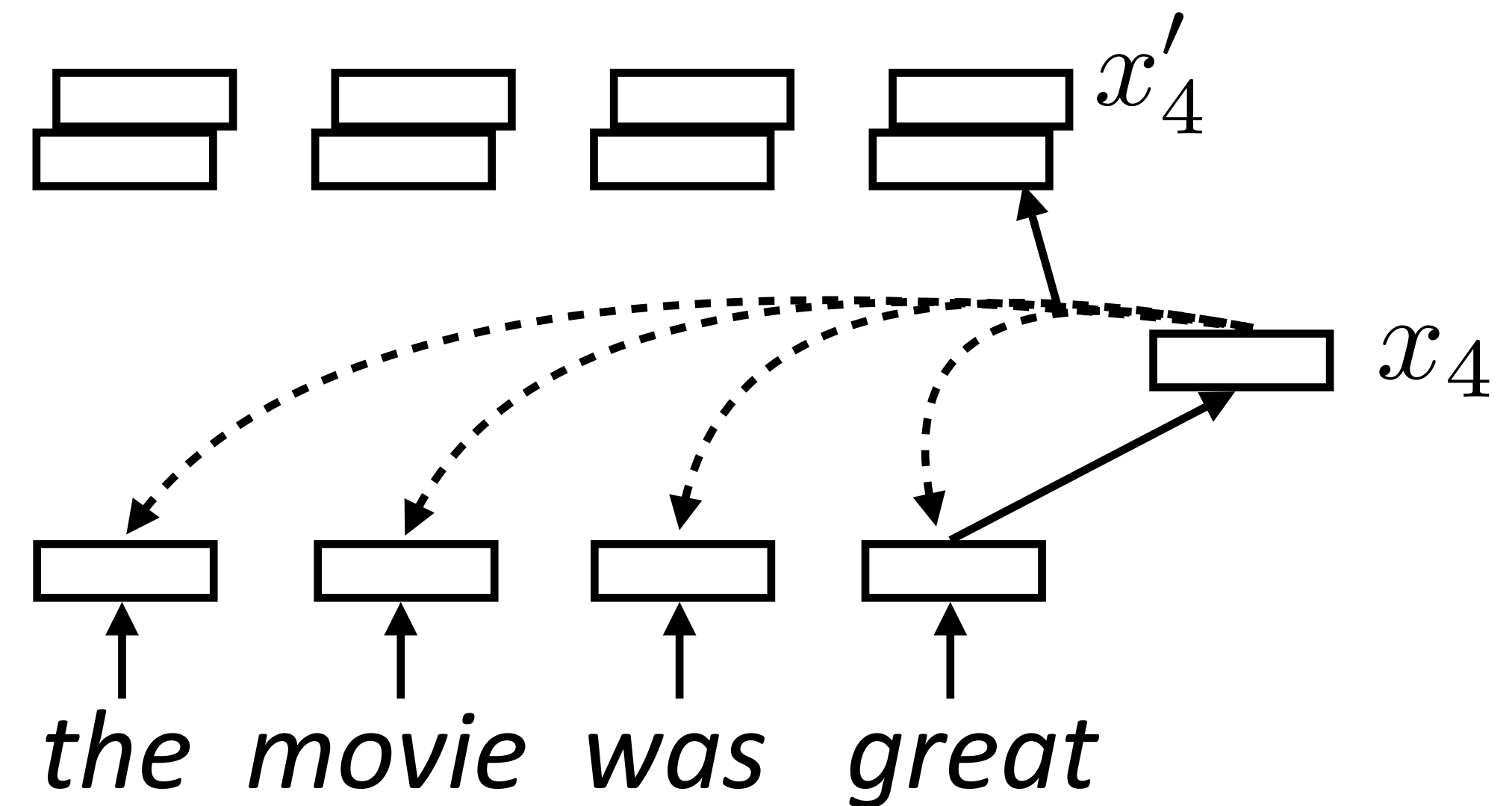


# Recall: Self-Attention/Transformers

- ▶ Each word forms a “query” which then computes attention over each word

$$\alpha_{i,j} = \text{softmax}(x_i^\top x_j) \quad \text{scalar}$$

$$x'_i = \sum_{j=1}^n \alpha_{i,j} x_j \quad \text{vector} = \text{sum of scalar} * \text{vector}$$



- ▶ Multiple “heads” analogous to different convolutional filters. Use parameters  $W_k$  and  $V_k$  to get different attention values + transform vectors

$$\alpha_{k,i,j} = \text{softmax}(x_i^\top W_k x_j) \quad x'_{k,i} = \sum_{j=1}^n \alpha_{k,i,j} V_k x_j$$



# This Lecture

---

- ▶ MT basics, evaluation
- ▶ Word alignment
- ▶ Phrase-based decoders
- ▶ Syntax-based decoders

# MT Basics





# MT Ideally

- ▶ *I have a friend*  $\Rightarrow \exists x \text{ friend}(x, \text{self}) \Rightarrow J'ai un ami$   
*J'ai une amie* (friend is female)
- ▶ May need information you didn't think about in your representation
- ▶ Hard for semantic representations to cover everything
- ▶ Everyone has a friend  $\Rightarrow \begin{matrix} \exists x \forall y \text{ friend}(x, y) \\ \forall x \exists y \text{ friend}(x, y) \end{matrix} \Rightarrow \text{Tous a un ami}$
- ▶ Can often get away without doing all disambiguation — same ambiguities may exist in both languages



# MT in Practice

---

- ▶ Bitext: this is what we learn translation systems from. What can you learn?

Je fais un bureau

I'm making a desk

Je fais une soupe

I'm making soup

Je fais un bureau

I make a desk

Qu'est-ce que tu fais?

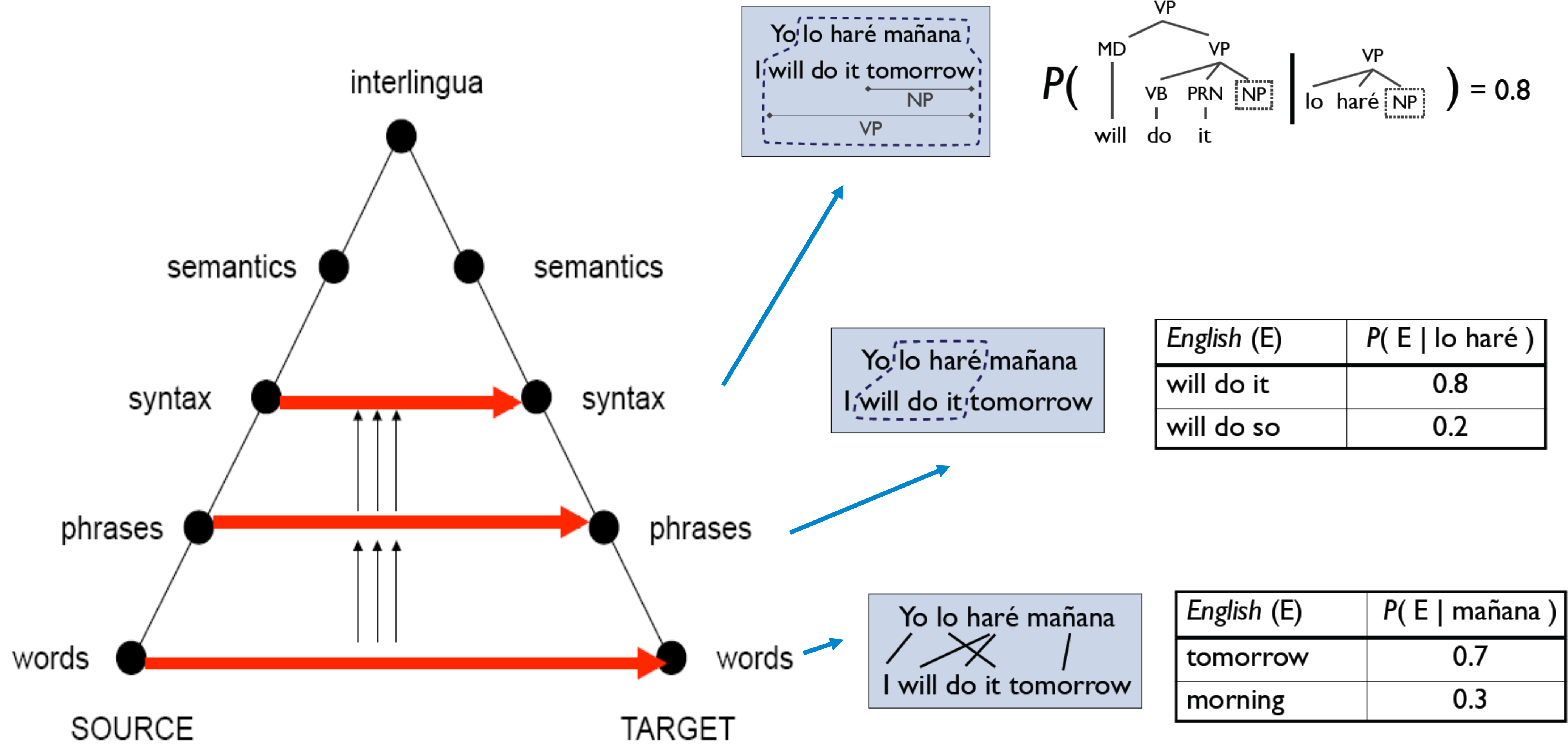
What are you making?

- ▶ What makes this hard?
  - Not word-to-word translation
  - Multiple translations of a single source (ambiguous)





# Levels of Transfer: Vauquois Triangle



► Today: mostly phrase-based, some syntax

Slide credit: Dan Klein



# Phrase-Based MT

---

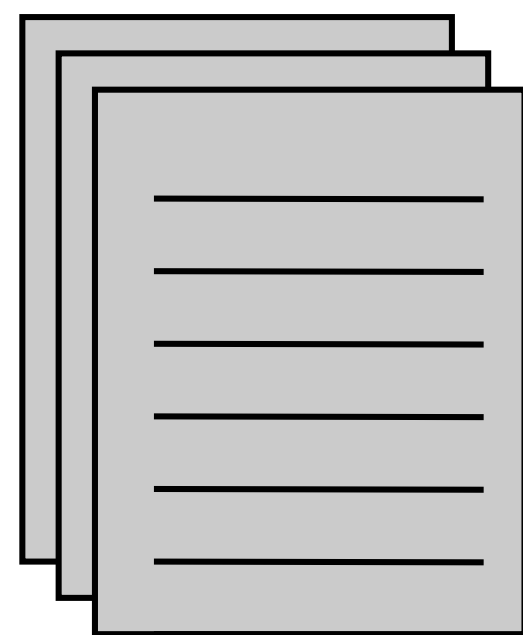
- ▶ Key idea: translation works better the bigger chunks you use
- ▶ Remember phrases from training data, translate piece-by-piece and stitch those pieces together to translate
  - ▶ How to identify phrases? Word alignment over source-target bitext
  - ▶ How to stitch together? Language model over target language
  - ▶ Decoder takes phrases and a language model and searches over possible translations
- ▶ NOT like standard discriminative models (take a bunch of translation pairs, learn a ton of parameters in an end-to-end way)



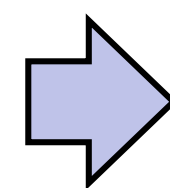
# Phrase-Based MT

cat ||| chat ||| 0.9  
the cat ||| le chat ||| 0.8  
dog ||| chien ||| 0.8  
house ||| maison ||| 0.6  
my house ||| ma maison ||| 0.9  
language ||| langue ||| 0.9  
...

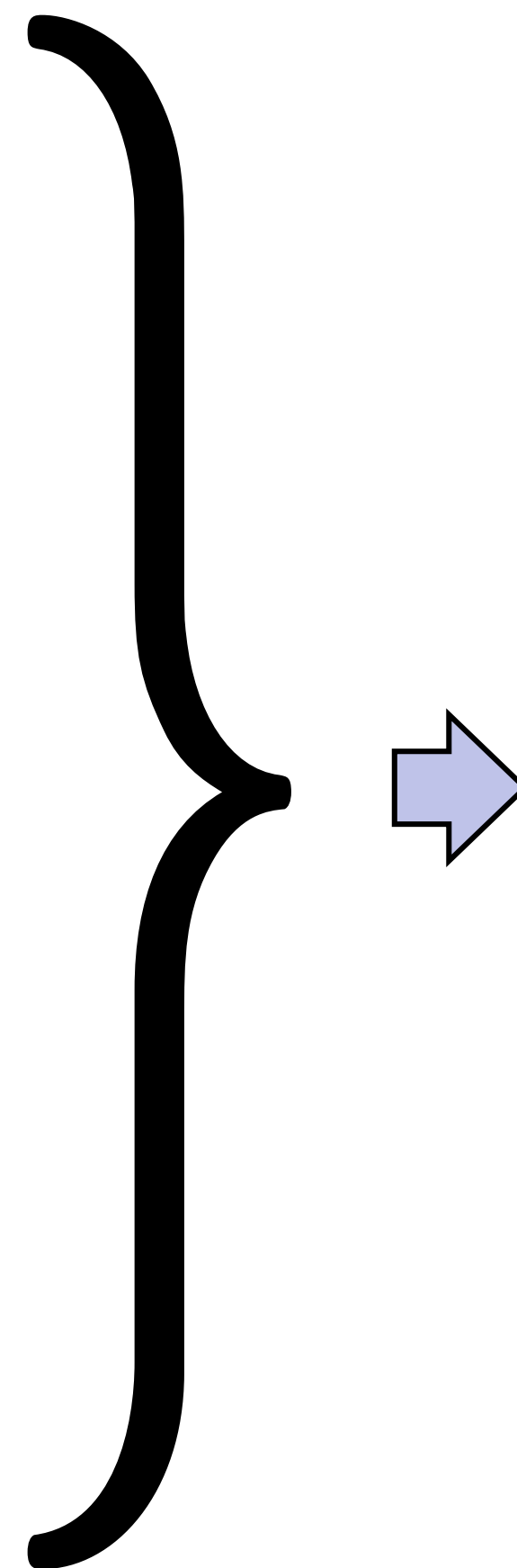
Phrase table  $P(f|e)$



Unlabeled English data



Language model  $P(e)$



$$P(e|f) \propto P(f|e)P(e)$$

Noisy channel model:  
combine scores from  
translation model +  
language model to  
translate foreign to  
English

“Translate faithfully but make fluent English”



# Evaluating MT

- ▶ Fluency: does it sound good in the target language?
- ▶ Fidelity/adequacy: does it capture the meaning of the original?
- ▶ BLEU score: geometric mean of 1-, 2-, 3-, and 4-gram *precision* vs. a reference, multiplied by brevity penalty (penalizes short translations)

$$\text{BLEU} = \text{BP} \cdot \exp \left( \sum_{n=1}^N w_n \log p_n \right) \quad \text{Typically } n = 4, w_i = 1/4$$

$$\text{BP} = \begin{cases} 1 & \text{if } c > r \\ e^{(1-r/c)} & \text{if } c \leq r \end{cases} \quad \begin{array}{l} r = \text{length of reference} \\ c = \text{length of prediction} \end{array}$$

# Word Alignment



# Word Alignment

- ▶ Input: a bitext, pairs of translated sentences

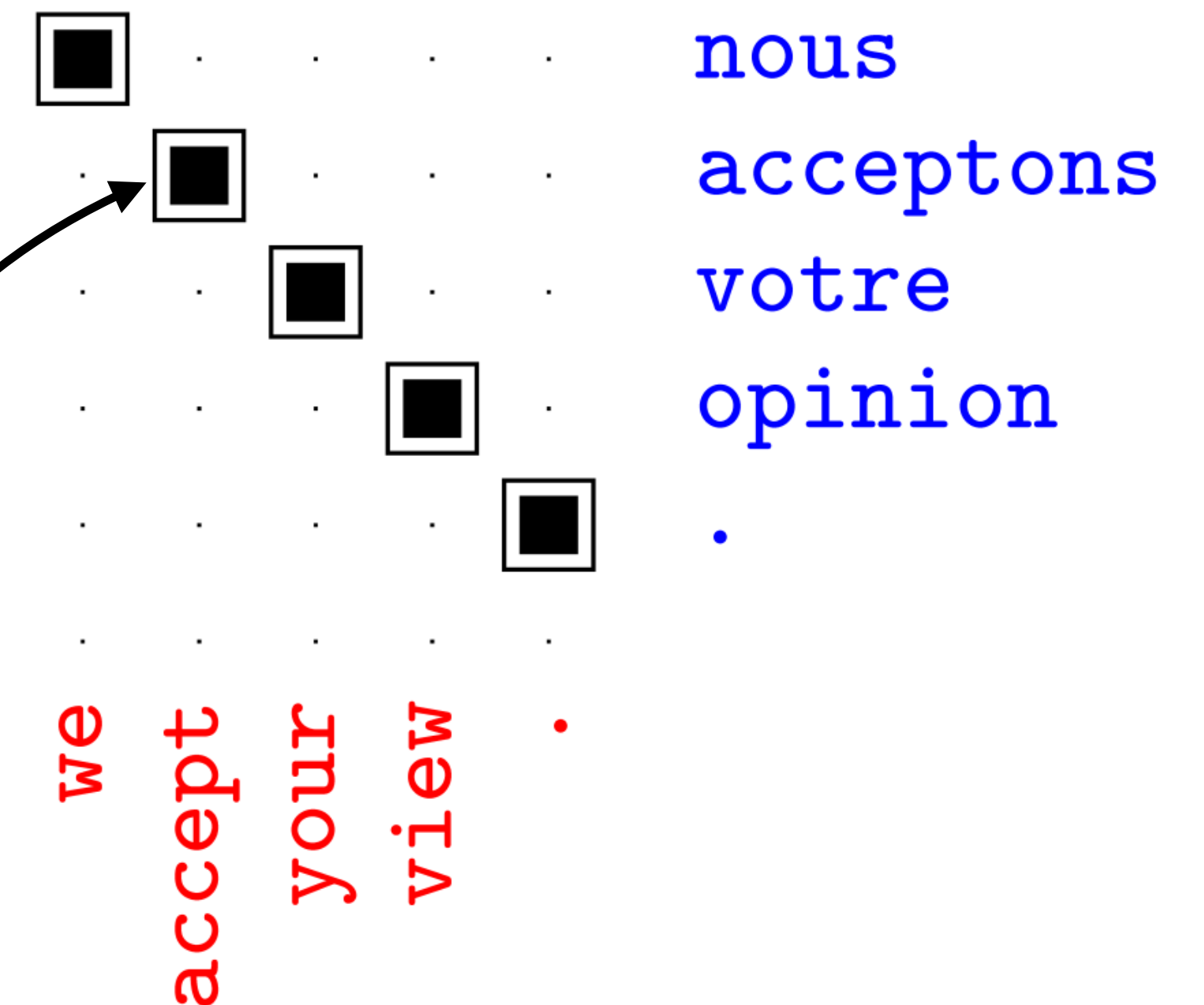
nous acceptons votre opinion . ||| we accept your view

nous allons changer d'avis ||| we are going to change our minds

- ▶ Output: alignments between words in each sentence

- ▶ We will see how to turn these into phrases

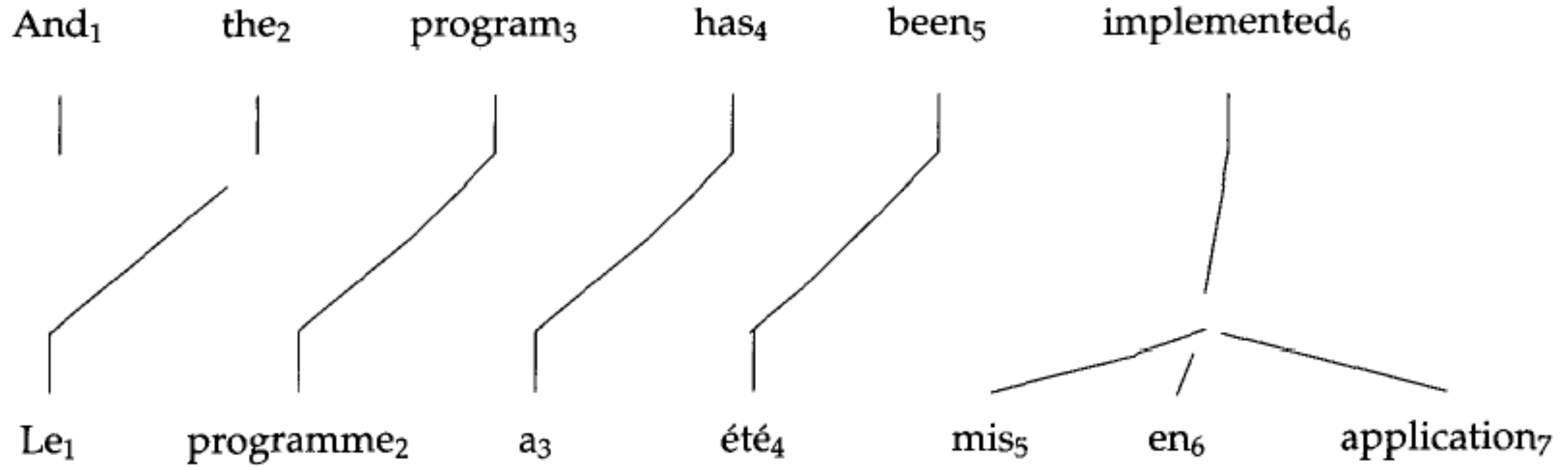
“accept and acceptons are aligned”







# 1-to-Many Alignments





# Word Alignment

---

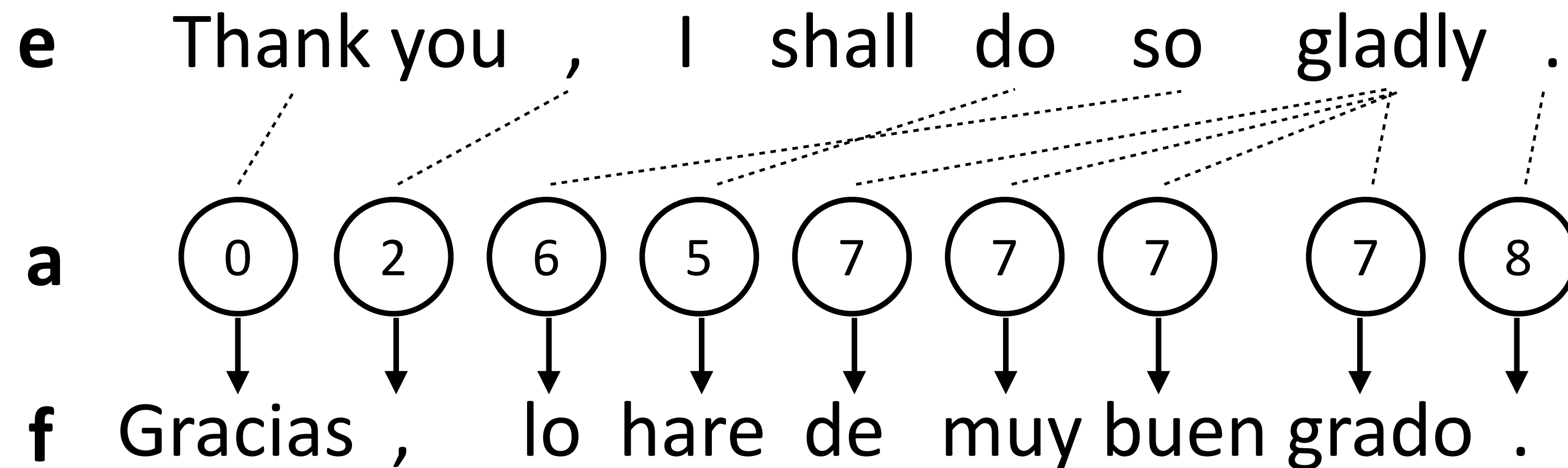
- ▶ Models  $P(\mathbf{f}|\mathbf{e})$ : probability of “French” sentence being generated from “English” sentence according to a model
- ▶ Latent variable model: 
$$P(\mathbf{f}|\mathbf{e}) = \sum_{\mathbf{a}} P(\mathbf{f}, \mathbf{a}|\mathbf{e}) = \sum_{\mathbf{a}} P(\mathbf{f}|\mathbf{a}, \mathbf{e})P(\mathbf{a})$$
- ▶ Correct alignments should lead to higher-likelihood generations, so by optimizing this objective we will learn correct alignments



# IBM Model 1

- ▶ Each French word is aligned to *at most* one English word

$$P(\mathbf{f}, \mathbf{a}|\mathbf{e}) = \prod_{i=1}^n P(f_i|e_{a_i})P(a_i)$$



- ▶ Set  $P(\mathbf{a})$  uniformly (no prior over good alignments)

- ▶  $P(f_i|e_{a_i})$ : word translation probability table. Learn with EM

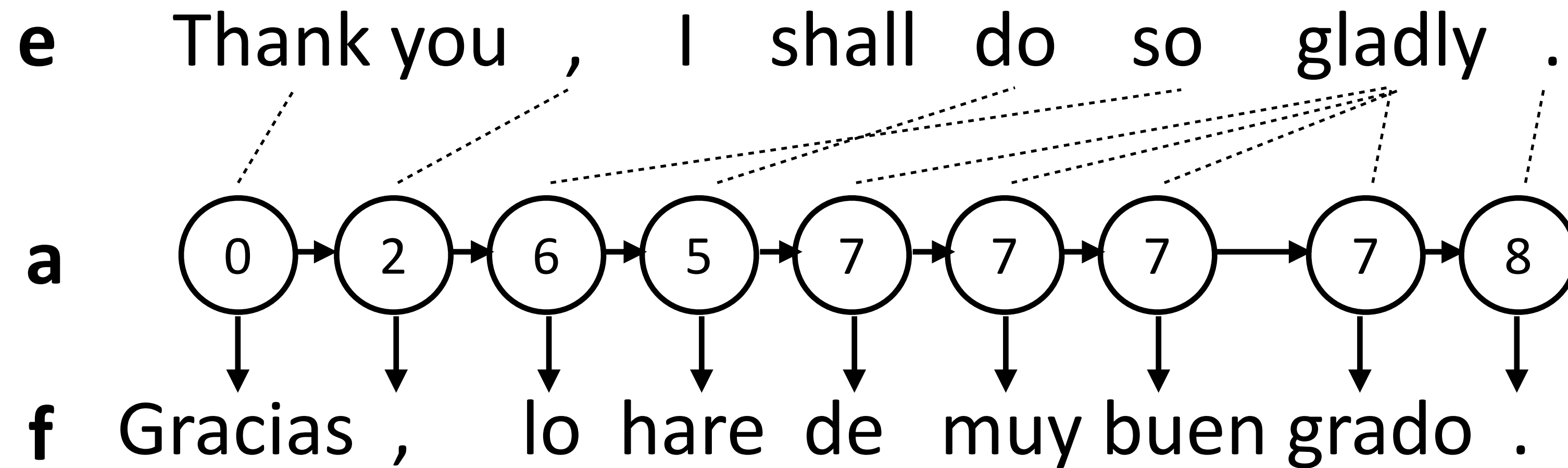
Brown et al. (1993)



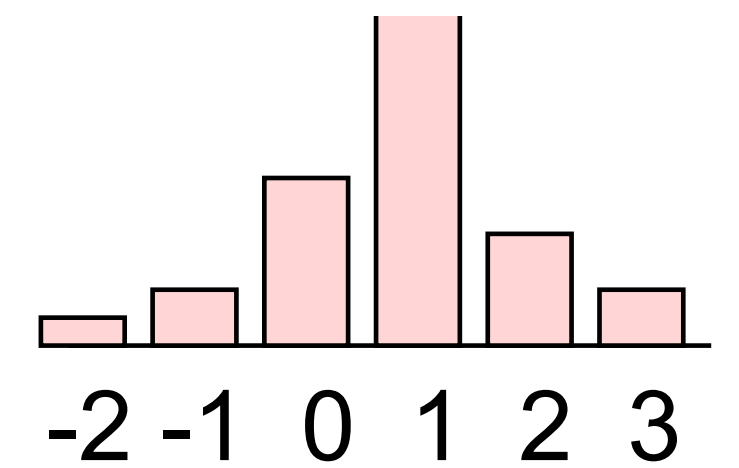
# HMM for Alignment

- ▶ Sequential dependence between a's to capture monotonicity

$$P(\mathbf{f}, \mathbf{a}|\mathbf{e}) = \prod_{i=1}^n P(f_i|e_{a_i})P(a_i|a_{i-1})$$



- ▶ Alignment dist parameterized by jump size:  $P(a_j - a_{j-1})$  →



- ▶  $P(f_i|e_{a_i})$ : same as before

Vogel et al. (1996)





# Phrase Extraction

- ▶ Find contiguous sets of aligned words in the two languages that don't have alignments to other words

d'assister à la reunion et ||| to attend the meeting and

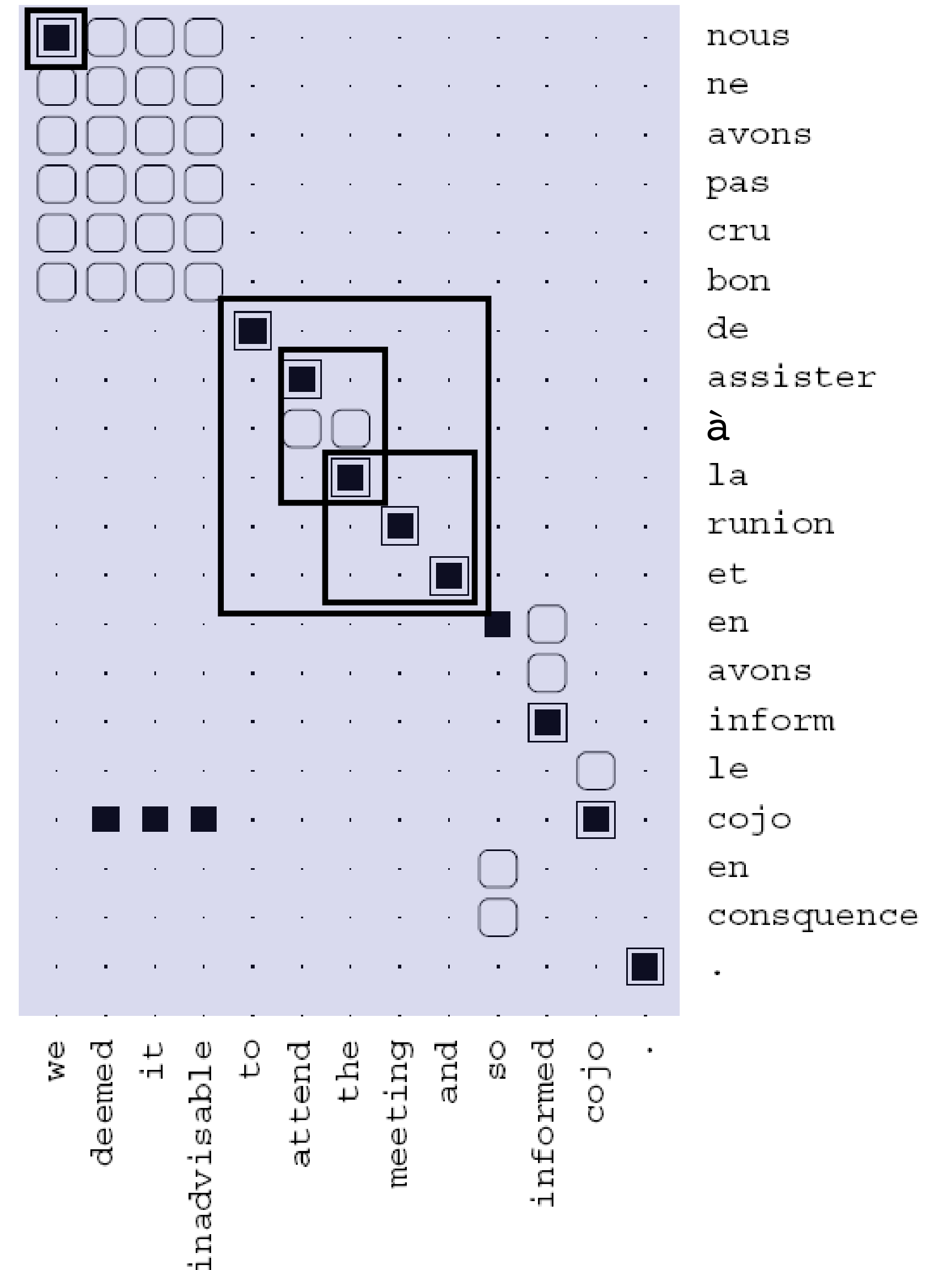
assister à la reunion ||| attend the meeting

la reunion and ||| the meeting and

nous ||| we

...

- ▶ Lots of phrases possible, count across all sentences and score by frequency





Decoding



# Recall: $n$ -gram Language Models

$$P(\mathbf{w}) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2) \dots$$

- ▶  $n$ -gram models: distribution of next word is a multinomial conditioned on previous  $n-1$  words  $P(w_i|w_1, \dots, w_{i-1}) = P(w_i|w_{i-n+1}, \dots, w_{i-1})$

I visited San \_\_\_\_\_ put a distribution over the next word

$$P(w|\text{visited San}) = \frac{\text{count}(\text{visited San}, w)}{\text{count}(\text{visited San})}$$

Maximum likelihood estimate of this 3-gram probability from a corpus

- ▶ Typically use  $\sim 5$ -gram language models for translation

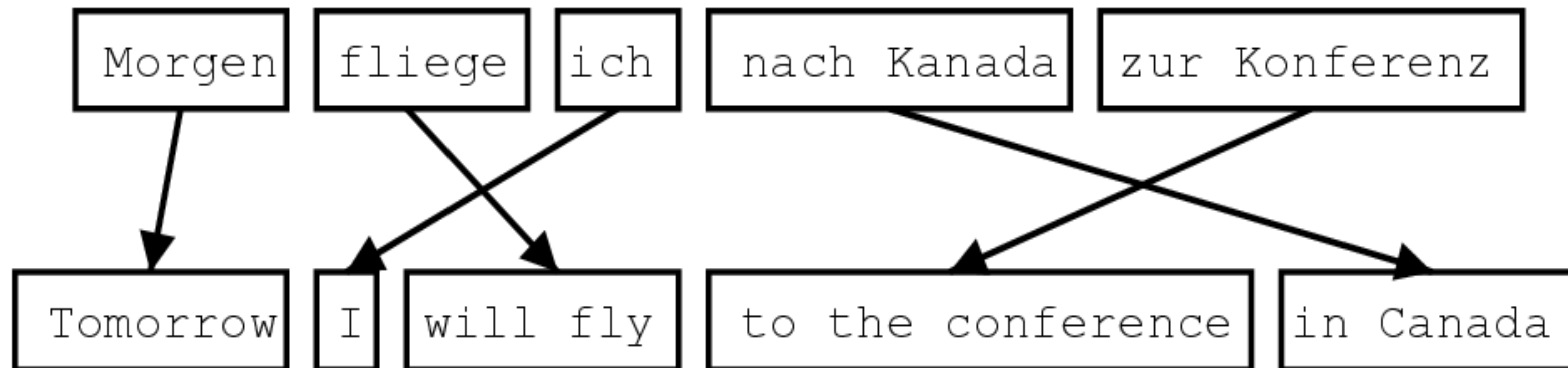


# Phrase-Based Decoding

- ▶ Inputs:

- ▶ n-gram language model:  $P(e_i|e_1, \dots, e_{i-1}) \approx P(e_i|e_{i-n-1}, \dots, e_{i-1})$
- ▶ Phrase table: set of phrase pairs  $(\mathbf{e}, \mathbf{f})$  with probabilities  $P(\mathbf{f}|\mathbf{e})$

- ▶ What we want to find:  $\mathbf{e}$  produced by a series of phrase-by-phrase translations from an input  $\mathbf{f}$ , possibly with reordering:







# Phrase lattices are big!

这 7人 中包括 来自 法国 和 俄罗斯 的 宇航 员 .

the	7 people	including	by some	and	the russian	the	the astronauts	,
it	7 people included		by france	and the	the russian	international astronautical	of rapporteur .	
this	7 out	including the	from	the french	and the russian	the fifth	.	
these	7 among	including from		the french and	of the russian	of	space	members .
that	7 persons	including from the		of france	and to	russian	of the	aerospace
	7 include		from the	of france and	russian	astronauts	.	the
	7 numbers include		from france		and russian	of astronauts who	.	"
	7 populations include		those from france		and russian	astronauts .		
	7 deportees included		come from	france	and russia	in	astronautical	personnel ;
	7 philtrum	including those from		france and	russia	a space	member	
		including representatives from		france and the	russia	astronaut		
		include	came from	france and russia		by cosmonauts		
		include representatives from		french	and russia	cosmonauts		
		include	came from france		and russia 's	cosmonauts .		
		includes	coming from	french and	russia 's	cosmonaut		
				french and russian	's	astronavigation	member .	
				french	and russia	astronauts		
					and russia 's		special rapporteur	
					, and russia		rapporteur	
					, and russia		rapporteur .	
					, and russia			
				or	russia 's			



# Phrase-Based Decoding

*The decoder...*

*tries different segmentations,*

*translates phrase by phrase,*

*and considers reorderings.*

▶ Input

lo haré | rápidamente |.

▶ Translations

I'll do it | quickly |.

quickly | I'll do it |.

$$\arg \max_{\mathbf{e}} [P(\mathbf{f}|\mathbf{e}) \cdot P(\mathbf{e})]$$

▶ Decoding objective (for 3-gram LM)

$$\arg \max_{\mathbf{e}} \left[ \prod_{\langle \bar{e}, \bar{f} \rangle} P(\bar{f}|\bar{e}) \cdot \prod_{i=1}^{|\mathbf{e}|} P(e_i|e_{i-1}, e_{i-2}) \right]$$

Slide credit: Dan Klein





# Monotonic Translation

María	no	dio	una	bofetada	a	la	bruja	verde
<u>Mary</u>	<u>not</u>	<u>give</u>	<u>a</u>	<u>slap</u>	<u>to</u>	<u>the</u>	<u>witch</u>	<u>green</u>
	<u>did not</u>		<u>a</u>	<u>slap</u>	<u>by</u>		<u>green</u>	<u>witch</u>
	<u>no</u>		<u>slap</u>		<u>to the</u>			
	<u>did not give</u>				<u>to</u>			
					<u>the</u>			
			<u>slap</u>			<u>the</u>	<u>witch</u>	

► If we translate with beam search, what state do we need to keep in the beam?

- What have we translated so far?  $\arg \max_e \left[ \prod_{\langle \bar{e}, \bar{f} \rangle} P(\bar{f} | \bar{e}) \cdot \prod_{i=1}^{|\bar{e}|} P(e_i | e_{i-1}, e_{i-2}) \right]$
- What words have we produced so far?  
(need to remember the last 2 words for a 3-gram LM)

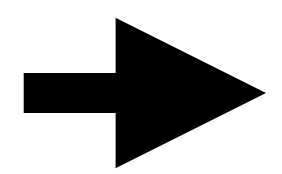




# Monotonic Translation

María	no	dio	una	bofetada	a	la	bruja	verde
<u>Mary</u>	<u>not</u>	<u>give</u>	<u>a</u>	<u>slap</u>	<u>to</u>	<u>the</u>	<u>witch</u>	<u>green</u>
	<u>did not</u>		<u>a</u>	<u>slap</u>	<u>by</u>		<u>green</u>	<u>witch</u>
	<u>no</u>		<u>slap</u>		<u>to the</u>			
	<u>did not give</u>				<u>to</u>			
				<u>slap</u>		<u>the</u>		
						<u>the</u>	<u>witch</u>	

Mary  
idx = 1      -1.1



...did not  
idx = 2      -0.3

Mary not  
idx = 2      -1.2

Mary no  
idx = 2      -2.9

- ▶ Beam state: where we're at, what the current translation so far is, and score of that translation
- ▶ Advancing state consists of trying each possible translation that could get us to this timestep



# Monotonic Translation

María	no	dio	una	bofetada	a	la	bruja	verde
<u>Mary</u>	<u>not</u>	<u>give</u>	<u>a</u>	<u>slap</u>	<u>to</u>	<u>the</u>	<u>witch</u>	<u>green</u>
	<u>did not</u>		<u>a</u>	<u>slap</u>	<u>by</u>		<u>green</u>	<u>witch</u>
	<u>no</u>		<u>slap</u>		<u>to the</u>			
	<u>did not give</u>				<u>to</u>			
				<u>slap</u>		<u>the</u>		
							<u>the</u>	<u>witch</u>

...did not idx = 2	-0.3
Mary not idx = 2	-1.2
Mary no idx = 2	-2.9

$$\text{score} = \log [ \underbrace{P(\text{Mary}) P(\text{not} | \text{Mary})}_{\text{LM}} \underbrace{P(\text{María} | \text{Mary}) P(\text{no} | \text{not})}_{\text{TM}} ]$$

In reality:  $\text{score} = \alpha \log P(\text{LM}) + \beta \log P(\text{TM})$

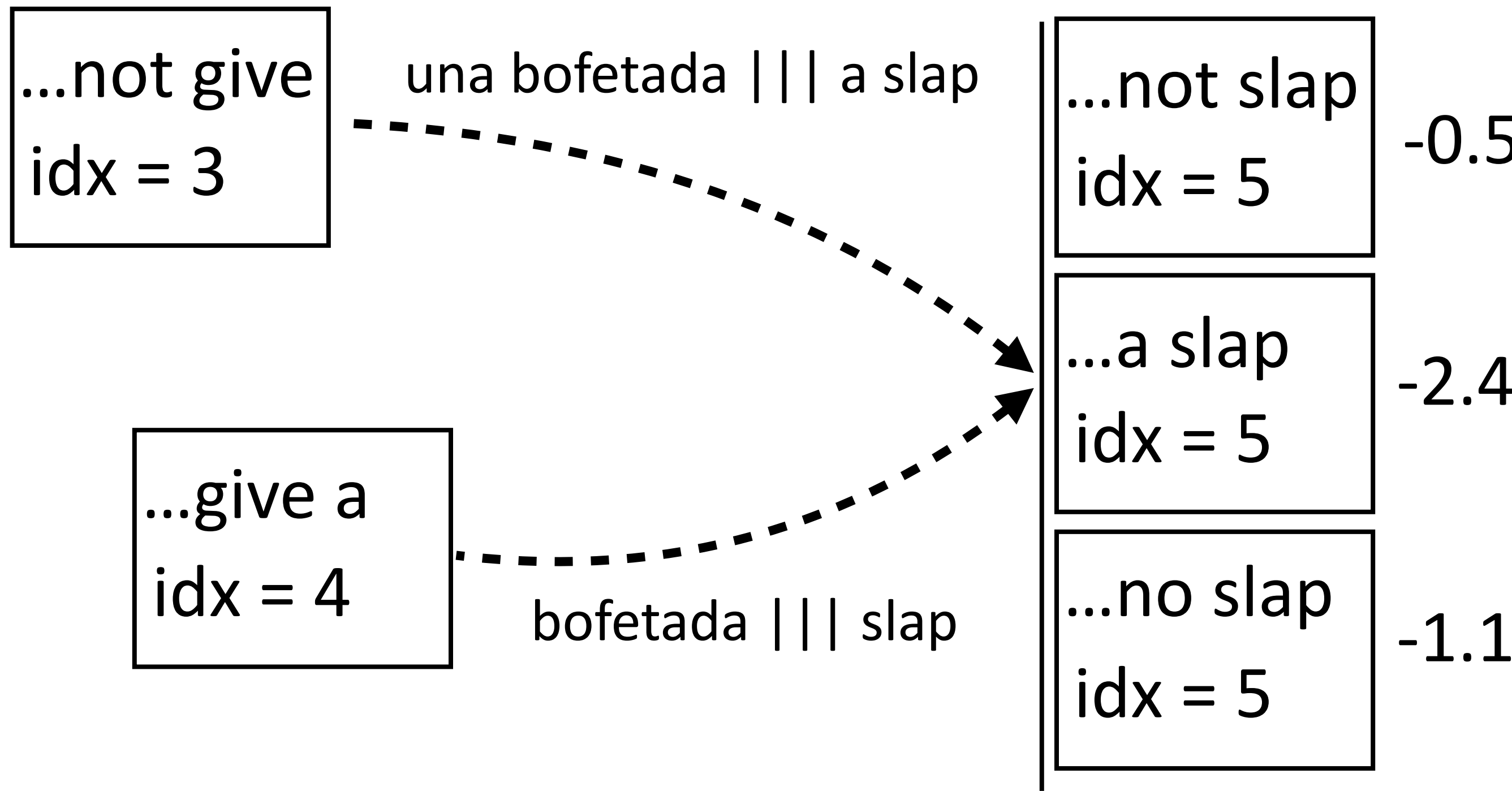
...and TM is broken down into several features

Koehn (2004)



# Monotonic Translation

Maria	no	dio	una	bofetada	a	la	bruja	verde
<u>Mary</u>	<u>not</u>	<u>give</u>	<u>a</u>	<u>slap</u>	<u>to</u>	<u>the</u>	<u>witch</u>	<u>green</u>
	<u>did not</u>		<u>a</u>	<u>slap</u>	<u>by</u>		<u>green</u>	<u>witch</u>
	<u>no</u>		<u>slap</u>		<u>to the</u>			
	<u>did not give</u>				<u>to</u>			
			<u>slap</u>		<u>the</u>			
				<u>slap</u>		<u>the</u>	<u>witch</u>	



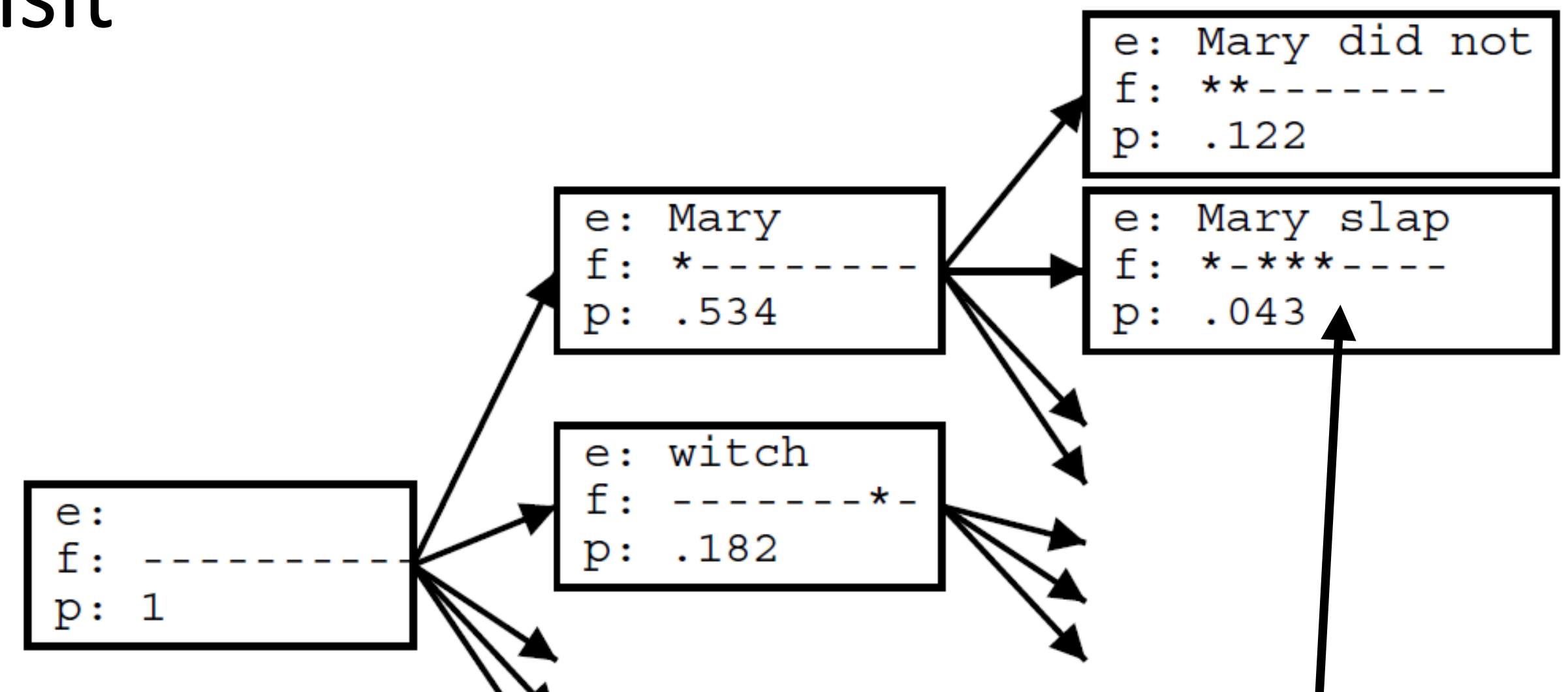
- ▶ Several paths can get us to this state, max over them (like Viterbi)
- ▶ Variable-length translation pieces = semi-HMM



# Non-Monotonic Translation

Maria	no	dio	una	bofetada	a	la	bruja	verde
<u>Mary</u>	<u>not</u>	<u>give</u>	<u>a</u>	<u>slap</u>	<u>to</u>	<u>the</u>	<u>witch</u>	<u>green</u>
	<u>did not</u>		<u>a</u>	<u>slap</u>	<u>by</u>		<u>green</u>	<u>witch</u>
	<u>no</u>		<u>slap</u>		<u>to the</u>			
	<u>did not give</u>				<u>to</u>			
				<u>slap</u>		<u>the</u>		
							<u>the</u>	<u>witch</u>

- ▶ Non-monotonic translation: can visit source sentence “out of order”
- ▶ State needs to describe which words have been translated and which haven’t
- ▶ Big enough phrases already capture lots of reorderings, so this isn’t as important as you think



translated: Maria, dio,  
una, bofetada



# Moses

---

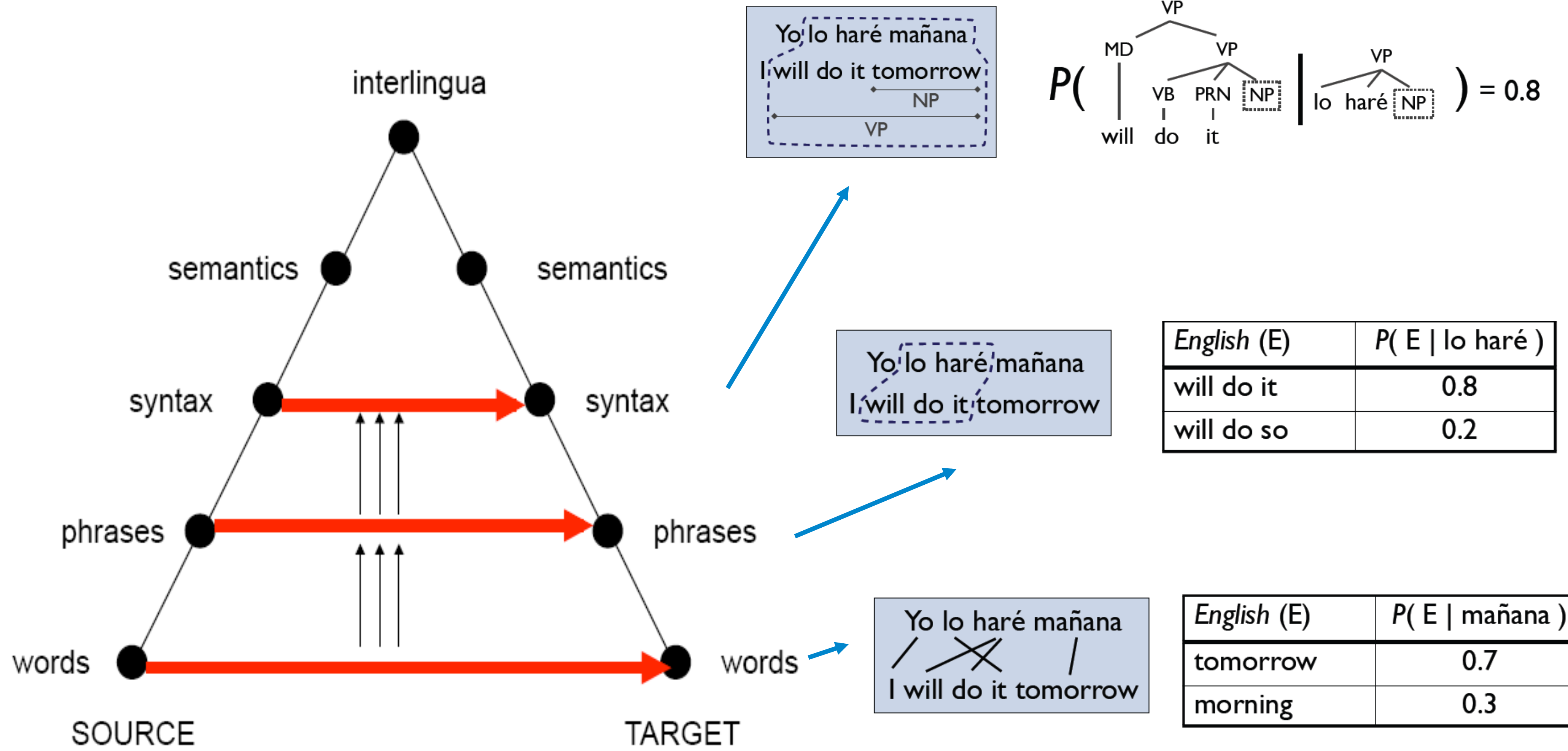
- ▶ Toolkit for machine translation due to Philipp Koehn + Hieu Hoang
  - ▶ Pharaoh (Koehn, 2004) is the decoder from Koehn's thesis
- ▶ Moses implements word alignment, language models, and this decoder, plus training regimes and more
  - ▶ Highly optimized and heavily engineered, could more or less build SOTA translation systems with this from 2007-2015
- ▶ Next time: results on these and comparisons to neural methods



Syntax



# Levels of Transfer: Vauquois Triangle



► Is syntax a “better” abstraction than phrases?



# Syntactic MT

- ▶ Rather than use phrases, use a *synchronous context-free grammar*

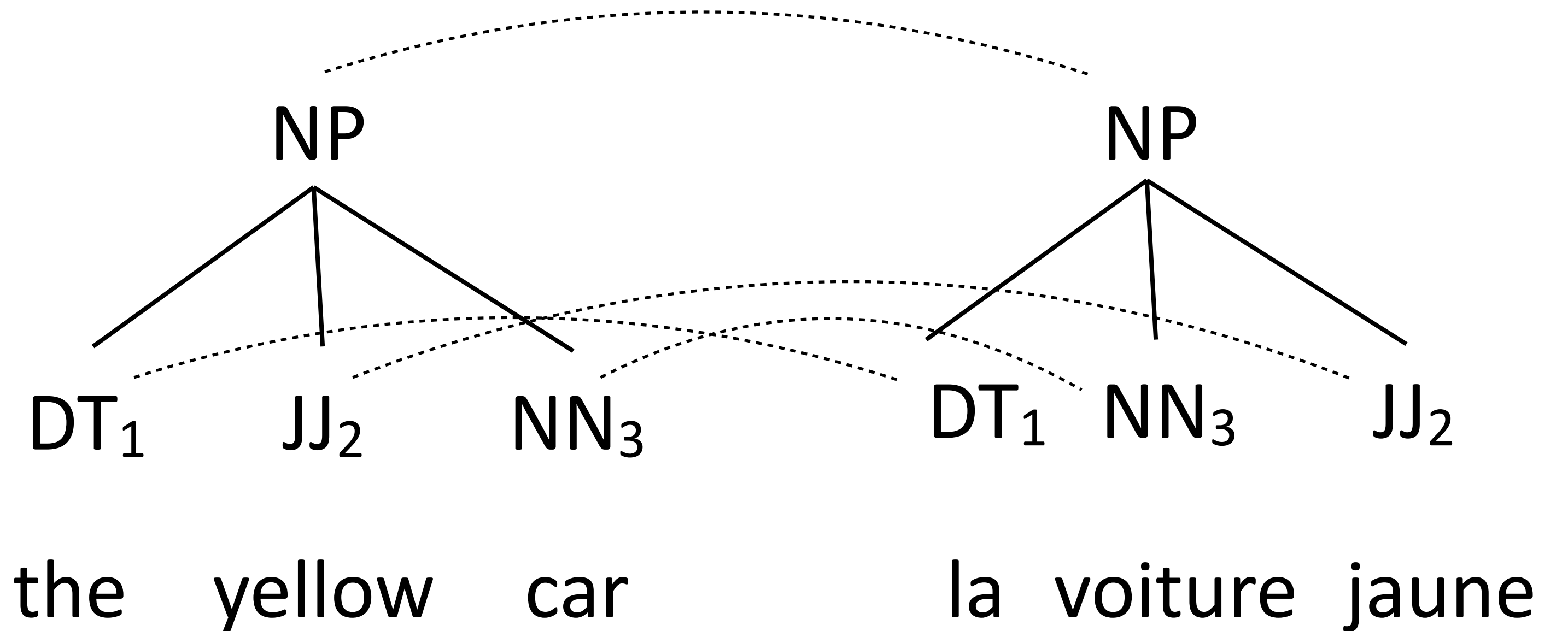
NP  $\rightarrow$  [DT<sub>1</sub> JJ<sub>2</sub> NN<sub>3</sub>; DT<sub>1</sub> NN<sub>3</sub> JJ<sub>2</sub>]

DT  $\rightarrow$  [the, la]

DT  $\rightarrow$  [the, le]

NN  $\rightarrow$  [car, voiture]

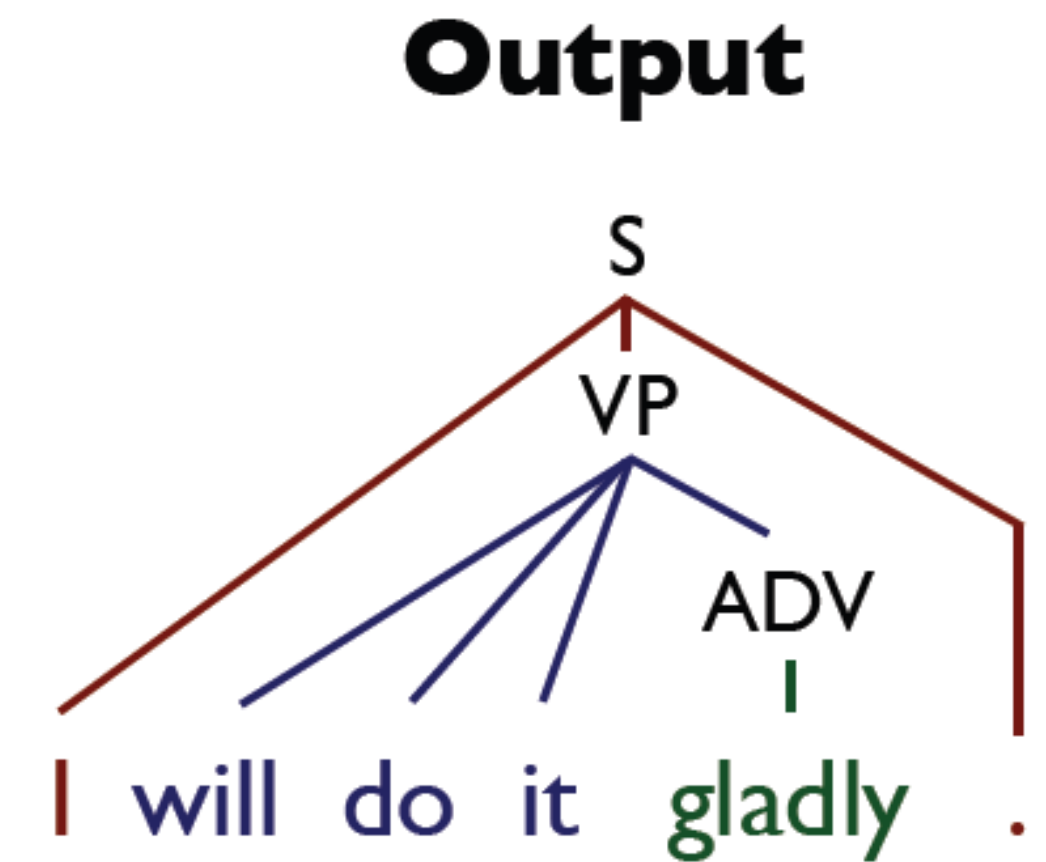
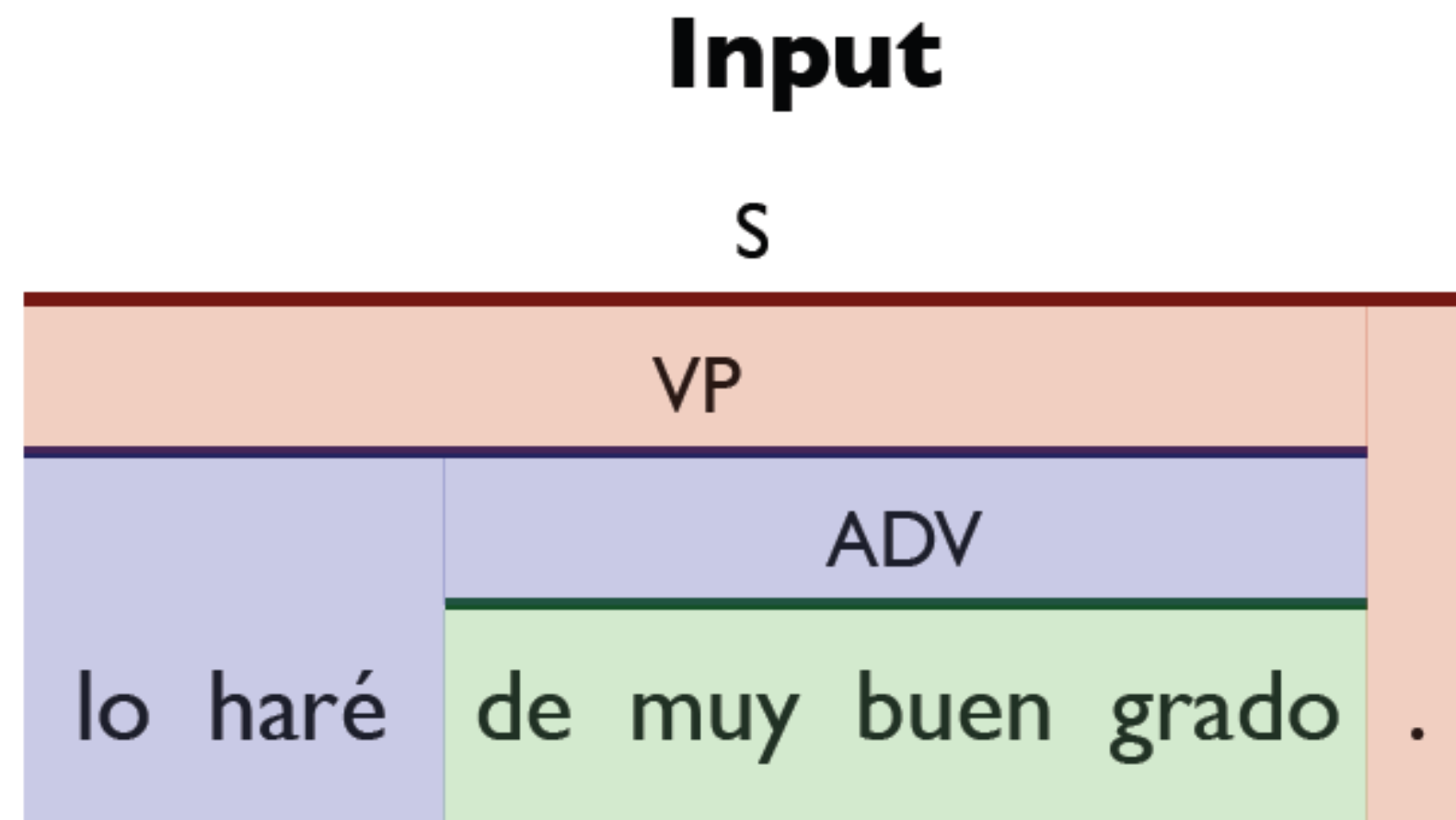
JJ  $\rightarrow$  [yellow, jaune]



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



# Syntactic MT



- ▶ Use lexicalized rules, look like “syntactic phrases”
- ▶ Leads to HUGE grammars, parsing is slow

## Grammar

$S \rightarrow \langle VP . ; I VP . \rangle$  **OR**  $S \rightarrow \langle VP . ; you VP . \rangle$

$VP \rightarrow \langle lo haré ADV ; will do it ADV \rangle$

$S \rightarrow \langle lo haré ADV . ; I will do it ADV . \rangle$

$ADV \rightarrow \langle de muy buen grado ; gladly \rangle$



# Takeaways

---

- ▶ Phrase-based systems consist of 3 pieces: aligner, language model, decoder
  - ▶ HMMs work well for alignment
  - ▶ N-gram language models are scalable and historically worked well
  - ▶ Decoder requires searching through a complex state space
- ▶ Lots of system variants incorporating syntax
- ▶ Next time: neural MT