

Decoding in Phrase-Based Machine Translation

(Building the translation)

Not required for the homework



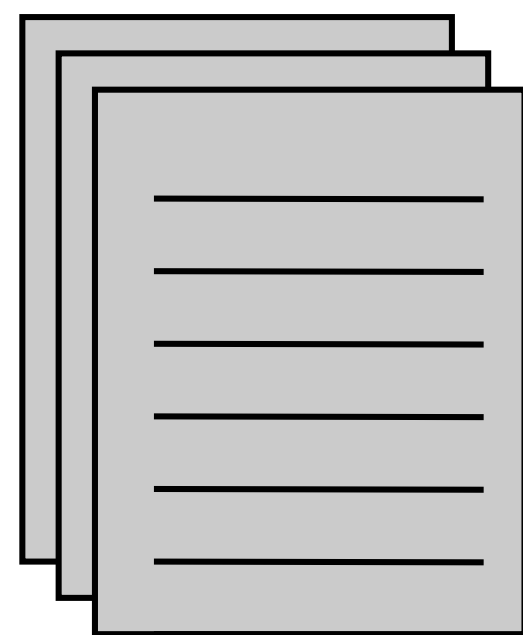
Phrase Extraction



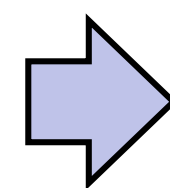
Phrase-Based MT

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
...

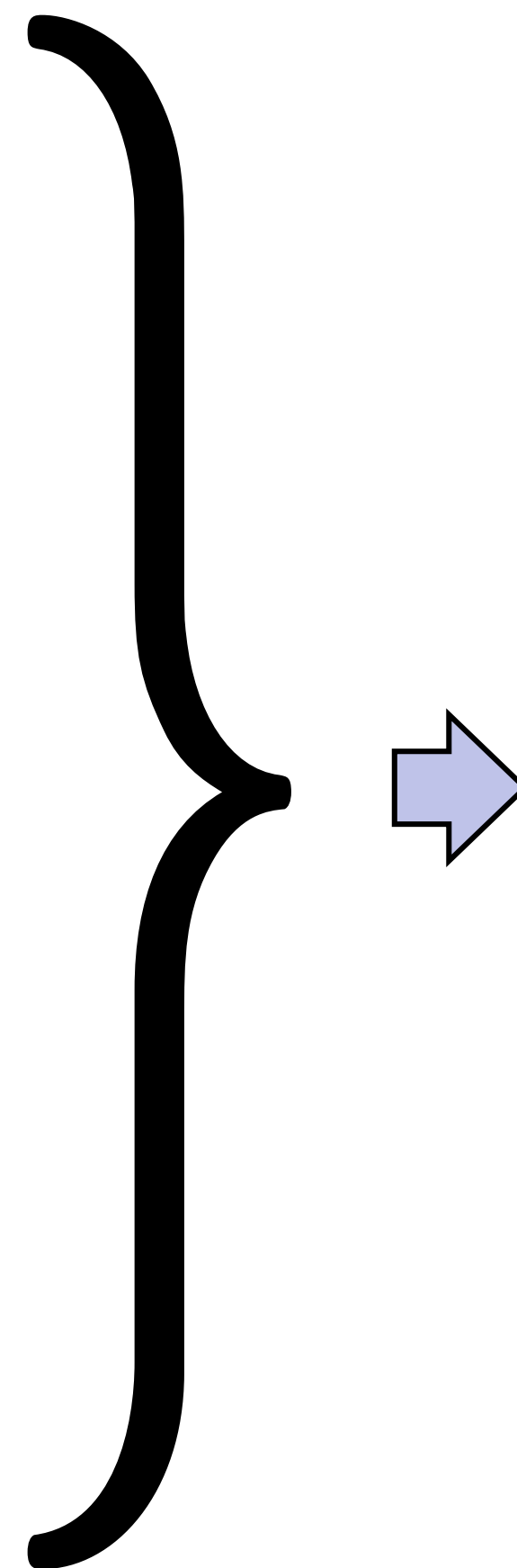
Phrase table $P(f|e)$



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”



Phrase-Based Decoding

- ▶ Noisy channel model: $P(\mathbf{e} | \mathbf{f}) \propto P(\mathbf{f} | \mathbf{e}) P(\mathbf{e})$ (ignore $P(\mathbf{f})$ term)
Translation model (TM) Language model (LM)
- ▶ Inputs needed
 - ▶ Language model that scores $P(e_i | e_1, \dots, e_{i-1}) \approx P(e_i | e_{i-n-1}, \dots, e_{i-1})$
 - ▶ 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}



Phrase Lattice

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>	

- ▶ Given an input sentence, look at our phrase table to find all possible translations of all possible spans
- ▶ Monotonic translation: need to translate each word in order, explore paths in the lattice that don't skip any words
- ▶ Looks like Viterbi, but the scoring is more complicated



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?

► Score

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

► Where are we in the sentence

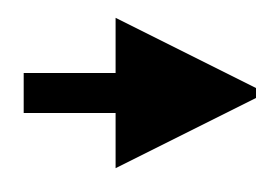
► What words have we produced so far (actually only need to remember the last 2 words when using a 3-gram LM)



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>the</u>			
			<u>slap</u>			<u>the</u>	<u>witch</u>	

Mary
idx = 1 -1.1



...did not
idx = 2 -0.1

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	-0.1
Mary not idx = 2	-1.2
Mary no idx = 2	-2.9

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

In reality: $\text{score} = \alpha \log P(\text{LM}) + \beta \log P(\text{TM})$
 ...and TM is broken down into several features



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 **a ton** more stuff
 - ▶ Highly optimized and heavily engineered, could more or less build SOTA translation systems with this from 2007-2013



Moses

SOURCE	Cela constituerait une solution transitoire qui permettrait de conduire à terme à une charte à valeur contraignante.
HUMAN	That would be an interim solution which would make it possible to work towards a binding charter in the long term .
1x DATA	[this] [constituerait] [assistance] [transitoire] [who] [permettrait] [licences] [to] [terme] [to] [a] [charter] [to] [value] [contraignante] [.]
10x DATA	[it] [would] [a solution] [transitional] [which] [would] [of] [lead] [to] [term] [to a] [charter] [to] [value] [binding] [.]
100x DATA	[this] [would be] [a transitional solution] [which would] [lead to] [a charter] [legally binding] [.]
1000x DATA	[that would be] [a transitional solution] [which would] [eventually lead to] [a binding charter] [.]

slide credit:
Dan Klein



Evaluating MT

- ▶ Fluency: does it sound good in the target language?
- ▶ Fidelity/adequacy: does it capture the meaning of the original?
- ▶ Automatic evaluation tries to approximate this...
- ▶ BLEU score: geometric mean of 1-, 2-, 3-, and 4-gram *precision* vs. a reference, multiplied by brevity penalty (penalizes short translations)
 - ▶ 1-gram precision: do you predict words that are in the reference?
 - ▶ 4-gram precision: to get this right, you need those words to be in the right order!
- ▶ Better metrics: human-in-the-loop variants

Syntactic MT



Syntactic MT

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

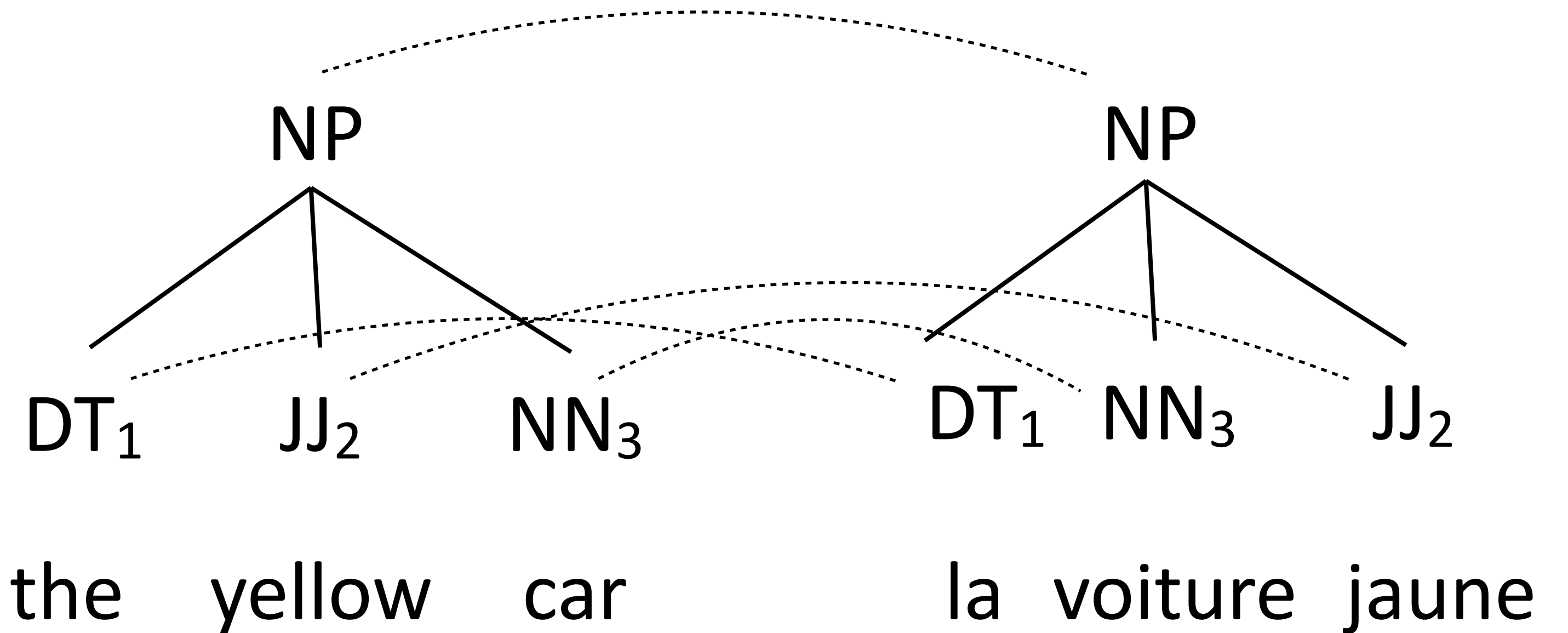
NP \rightarrow [DT₁ JJ₂ NN₃; DT₁ NN₃ JJ₂]

DT \rightarrow [the, la]

DT \rightarrow [the, le]

NN \rightarrow [car, voiture]

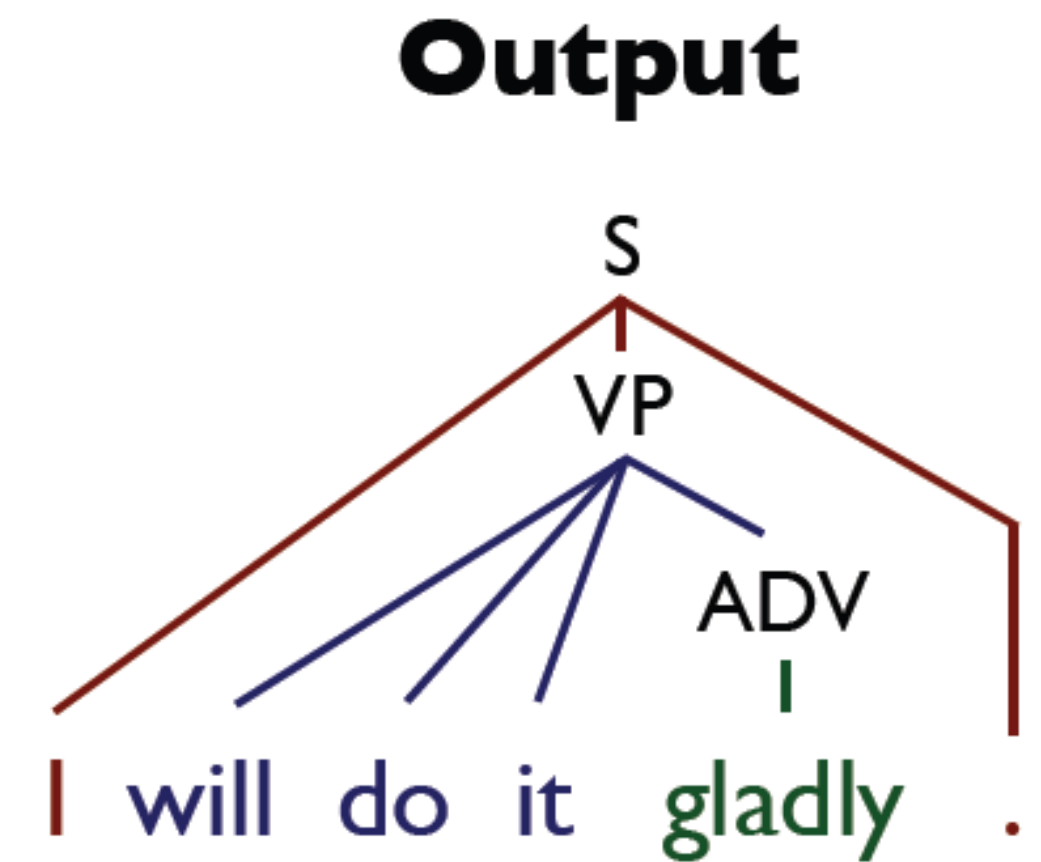
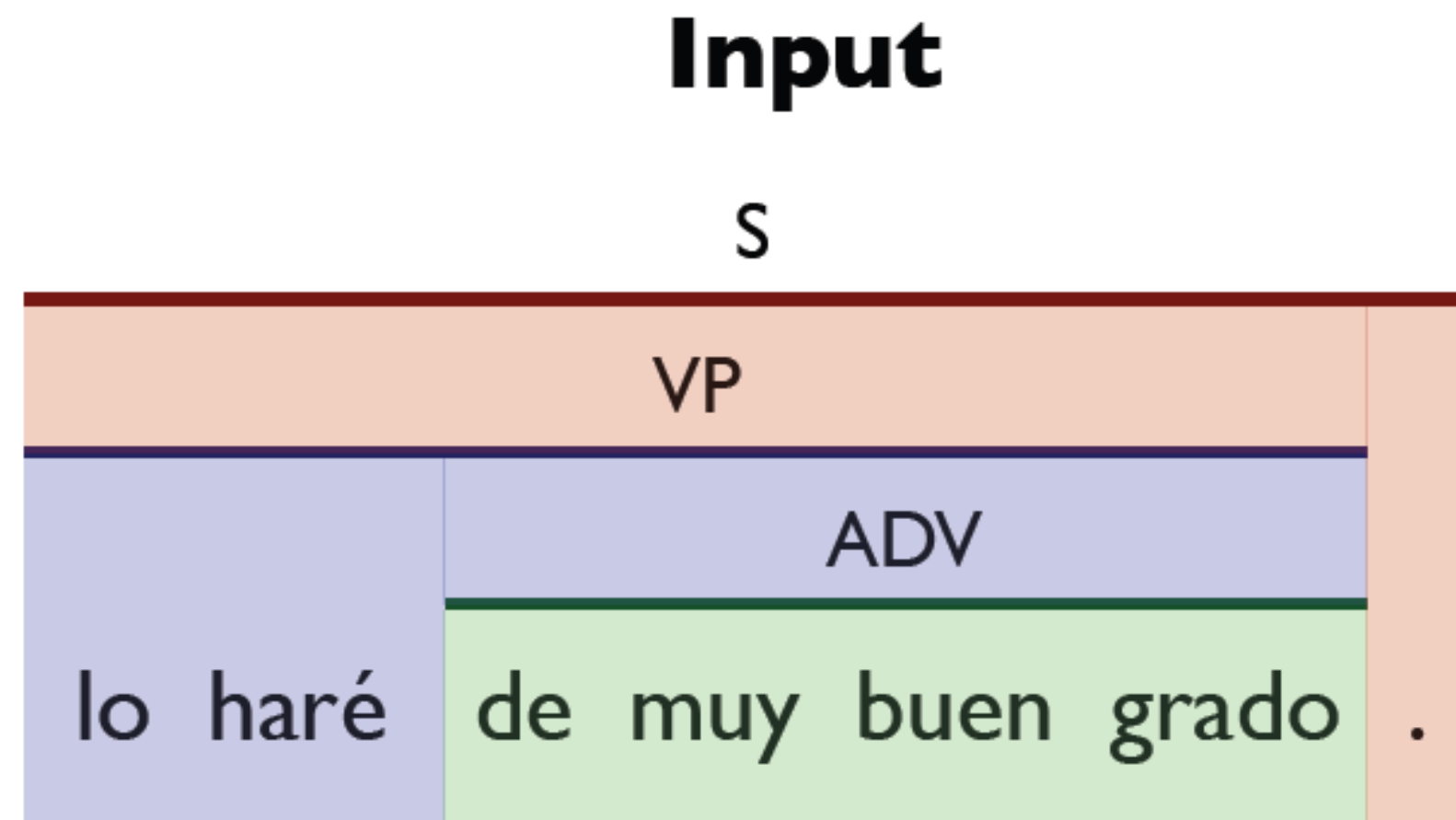
JJ \rightarrow [yellow, jaune]



- ▶ Translation = parse the input with “half” of the grammar, read off the other half
- ▶ Assumes parallel tree structures, but there can be 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** $S \rightarrow \langle VP . ; you VP . \rangle$

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

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

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