

# Natural Language Processing



## Machine Translation III

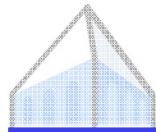
Dan Klein – UC Berkeley

# Phrase-Based MT



# Phrase-Based Translation Overview

<b>Input:</b>	lo haré   <b>rápidamente</b>  .	<i>The decoder... tries different segmentations, translates phrase by phrase, and considers reorderings.</i>
<b>Translations:</b>	I'll do it   <b>quickly</b>  .  quickly   I'll do it  .	
<b>Objective:</b>	$\arg \max_{\mathbf{e}} [P(\mathbf{f} \mathbf{e}) \cdot P(\mathbf{e})]$	$\arg \max_{\mathbf{e}} \left[ \prod_{\langle \bar{e}, \bar{f} \rangle} P(\bar{f} \bar{e}) \cdot \prod_{i=1}^{ \mathbf{e} } P(e_i e_{i-1}, e_{i-2}) \right]$

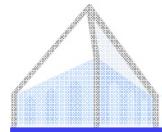


# Phrase-Based Decoding

这 7人 中包括 来自 法国 和 俄罗斯 的 宇航 员 .

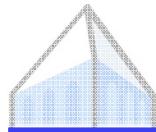
the	7 people	including	by some	and	the russian	the	the astronauts	,
it	7 people included	by france		and the	the russian	international astronomical	of rapporteur .	
this	7 out	including the	from	the french	and the russian	the fifth		.
these	7 among	including from		the french and	of the russian	of	space	members .
that	7 persons	including from the		of france	and to	russian	of the aerospace	members .
	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			

Decoder design is important: [Koehn et al. 03]



# Phrase-Based Decoding

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 slap</u>		<u>by</u>		<u>green 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>witch</u>		

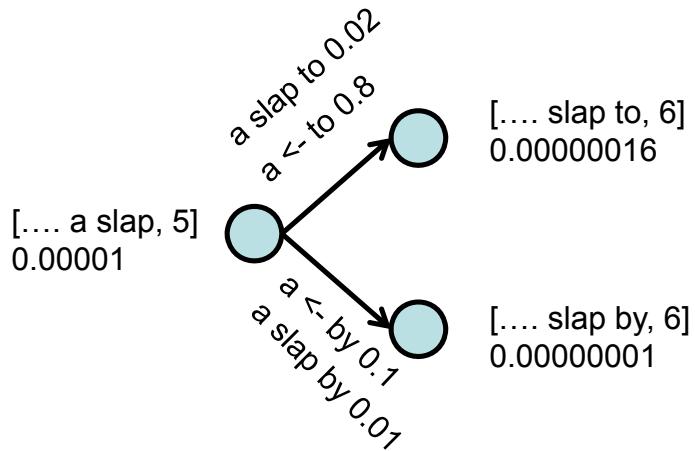


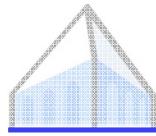
# Monotonic Word Translation

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

- Cost is  $LM * TM$
- It's an HMM?
  - $P(e|e_{-1}, e_{-2})$
  - $P(f|e)$
- State includes
  - Exposed English
  - Position in foreign
- Dynamic program loop?

```
for (fPosition in 1...|f|)  
    for (eContext in allEContexts)  
        for (eOption in translations[fPosition])  
            score = scores[fPosition-1][eContext] * LM(eContext+eOption) * TM(eOption, fWord[fPosition])  
            scores[fPosition][eContext[2]+eOption] =max score
```





# Beam Decoding

- For real MT models, this kind of dynamic program is a disaster (why?)
- Standard solution is beam search: for each position, keep track of only the best k hypotheses

```
for (fPosition in 1...|f|)
    for (eContext in bestEContexts[fPosition])
        for (eOption in translations[fPosition])
            score = scores[fPosition-1][eContext] * LM(eContext+eOption) * TM(eOption, fWord[fPosition])
            bestEContexts.maybeAdd(eContext[2]+eOption, score)
```

- Still pretty slow... why?
- Useful trick: cube pruning (Chiang 2005)

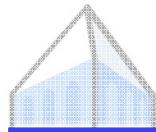
	1	4	7
1	2	5	8
2	3	6	9
6	7	10	13
10	11	14	17

	1	4	7
1	2	5	
2	3		
6			
10			

	1	4	7
2	5		
3			
7			

	1	4	7
2	5	8	
3			
7			

Example from David Chiang



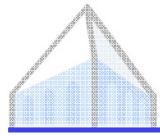
# Phrase Translation

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 slap</u>		<u>by</u>		<u>green 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 witch</u>	

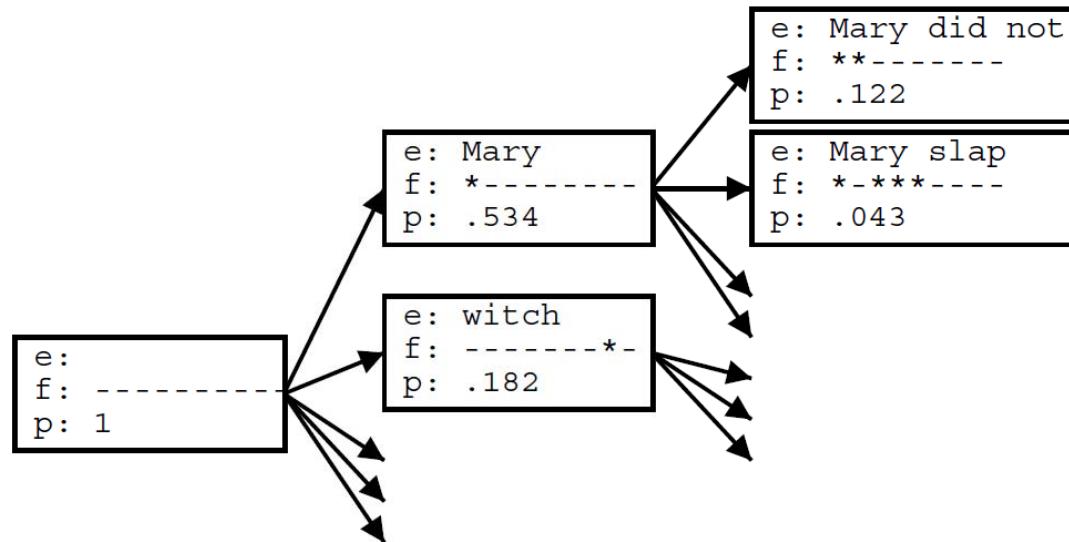
- If monotonic, almost an HMM; technically a semi-HMM

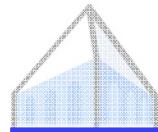
```
for (fPosition in 1...|f|)
    for (lastPosition < fPosition)
        for (eContext in eContexts)
            for (eOption in translations[fPosition])
                ... combine hypothesis for (lastPosition ending in eContext) with eOption
```

- If distortion... now what?



# Non-Monotonic Phrasal MT





## Pruning: Beams + Forward Costs

Maria no dio una bofetada a la bruja verde

—

e: Mary did not  
f: \*-----  
p: 0.154

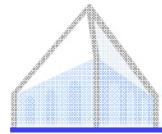
better  
partial  
translation

—

e: the  
f: -----\*--  
p: 0.354

covers  
easier part  
--> lower cost

- Problem: easy partial analyses are cheaper
  - Solution 1: use beams per foreign subset
  - Solution 2: estimate forward costs (A\*-like)



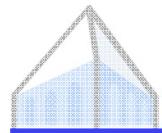
# The Pharaoh Decoder

Maria	no	dio	una	bofetada	a	la	bruja	verde
-------	----	-----	-----	----------	---	----	-------	-------

Mary    not    give    a    slap    to    the    witch    green  
did not                 a slap                 by                 green witch  
no                 slap                 to the  
did not give                 to  
                       the  
                       slap                 the witch

Maria	no	dio una bofetada	a la	bruja	verde
-------	----	------------------	------	-------	-------

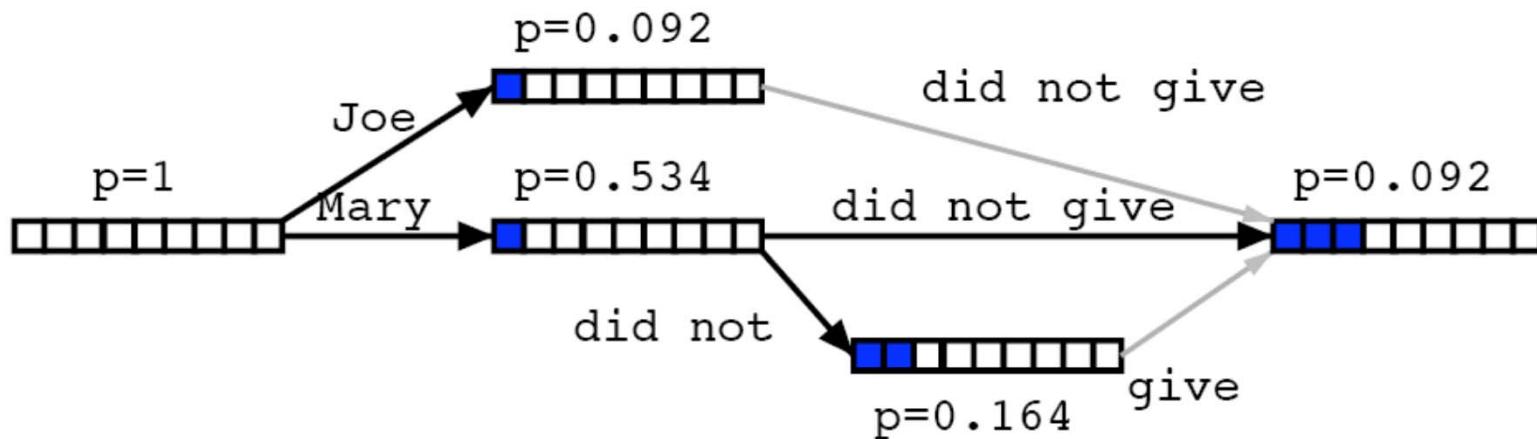
Mary	did not	slap	the	green	witch
------	---------	------	-----	-------	-------



# Hypothesis Lattices

Maria	no	dio	una	bofetada	a	la	bruja	verde
-------	----	-----	-----	----------	---	----	-------	-------

Mary    not    give    a    slap    to    the    witch    green  
did not    a slap    by    the    green  
no    slap    to the  
did not give    to  
                the  
                slap                      the witch



# Parameter Tuning

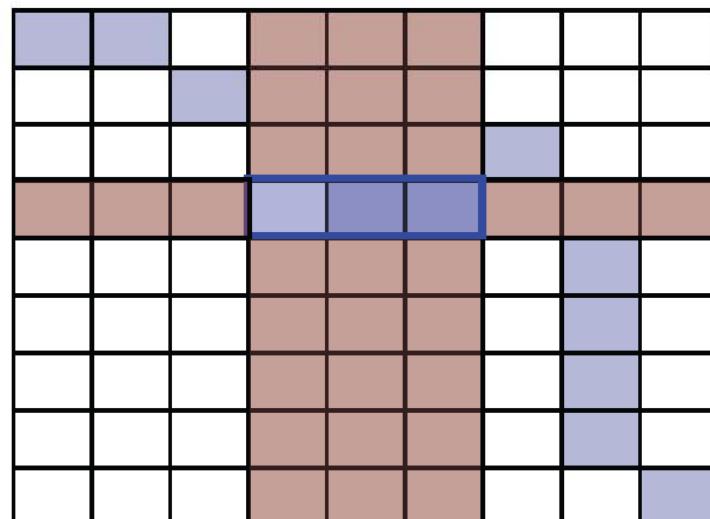
# Counting Phrase Pairs

**Input:**

Gracias , lo haré de muy buen grado .  
Thank you , I shall do so gladly .

*First, we learn word alignments,*

*then we infer aligned phrases.*



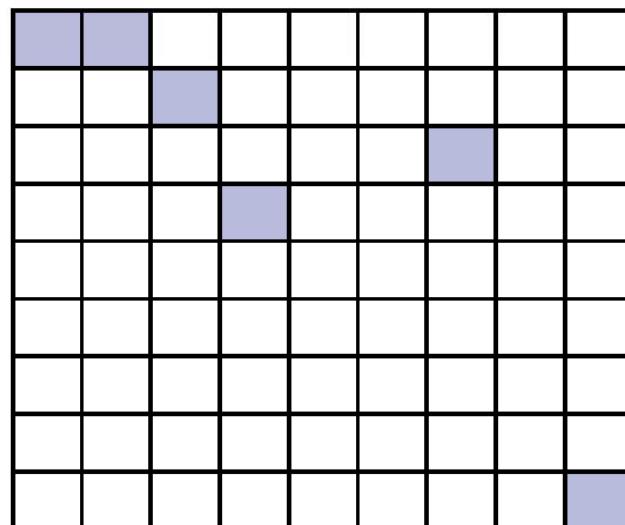
Thank you , I shall do so gladly .

## Gloss

Gracias	Thanks
,	,
lo	that
haré	do [first; future]
de	of
muy	very
buen	good
grado	degree
.	.

# What Happens in Practice

A real word alignment  
(GIZA++ Model 4 with  
grow-diag-final combination)

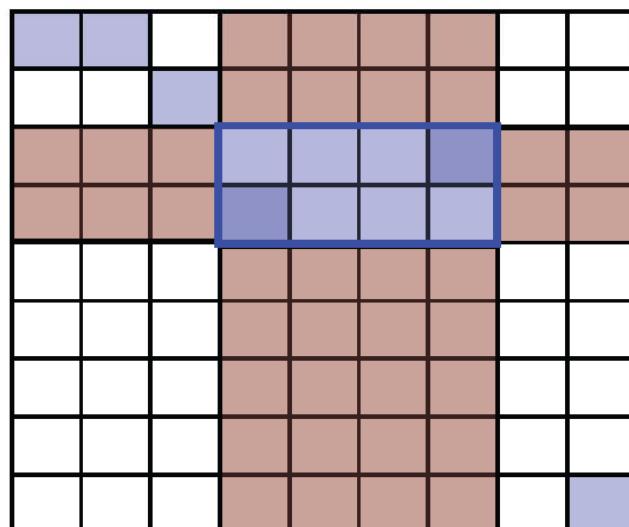


Thank you , I shall do so gladly .

<b>Gloss</b>	
Gracias	Thanks
,	,
lo	that
haré	do [first; future]
de	of
muy	very
buen	good
grado	degree
.	.

# What Happens in Practice

A real word alignment  
(GIZA++ Model 4 with  
grow-diag-final combination)

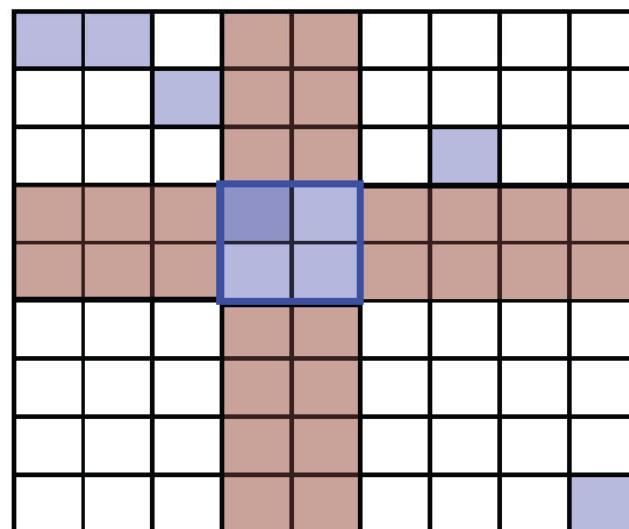


Thank you , I shall do so gladly .

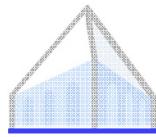
	<b>Gloss</b>
Gracias	Thanks
,	,
lo	that
haré	do [first; future]
de	of
muy	very
buen	good
grado	degree
.	.

# What Happens in Practice

A real word alignment  
(GIZA++ Model 4 with  
grow-diag-final combination)

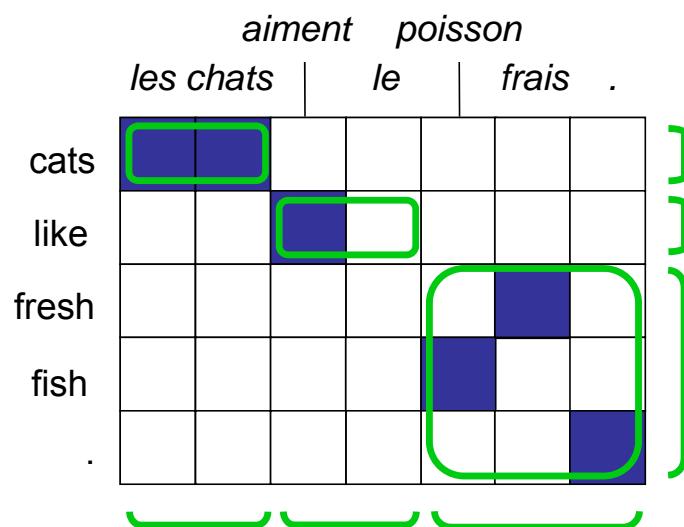


	<b>Gloss</b>
Gracias	<i>Thanks</i>
,	,
lo	<i>that</i>
haré	<i>do [first; future]</i>
de	<i>of</i>
muy	<i>very</i>
buen	<i>good</i>
grado	<i>degree</i>
.	.

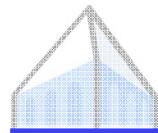


# Phrase Scoring

$$\phi_{new}(\bar{e}_j | \bar{f}_i) = \frac{c(\bar{f}_i, \bar{e}_j)}{c(\bar{f}_i)}$$

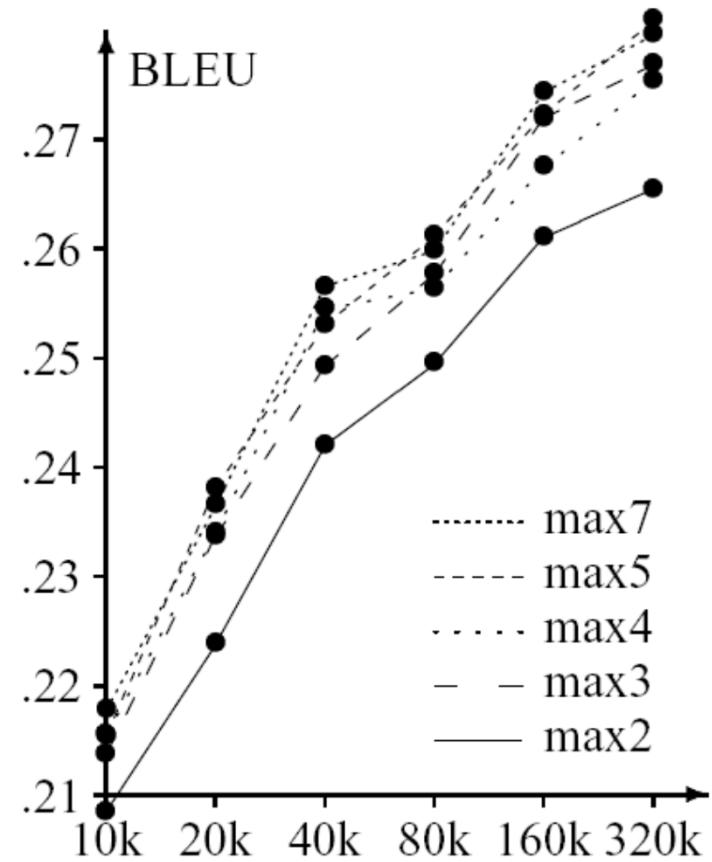
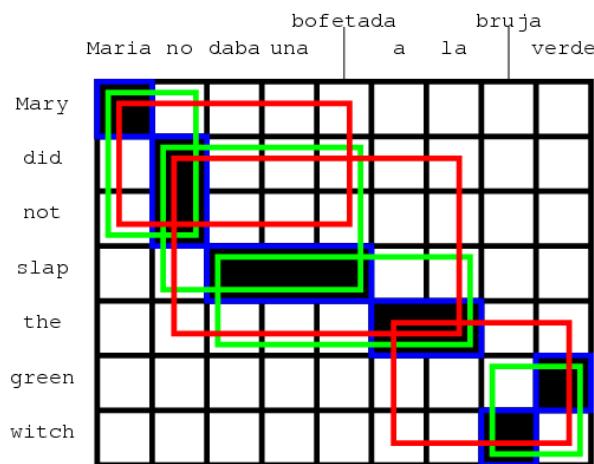


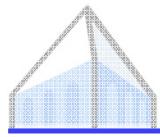
- Learning weights has been tried, several times:
  - [Marcu and Wong, 02]
  - [DeNero et al, 06]
  - ... and others
- Seems not to work well, for a variety of partially understood reasons
- Main issue: big chunks get all the weight, obvious priors don't help
  - Though, [DeNero et al 08]



# Phrase Size

- Phrases do help
  - But they don't need to be long
  - Why should this be?



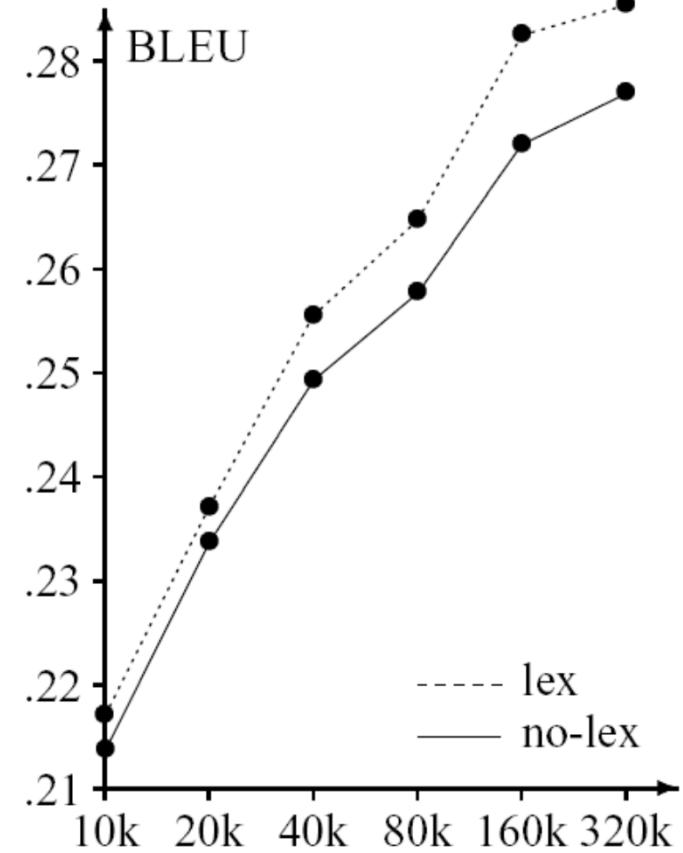


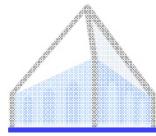
# Lexical Weighting

$$\phi(\bar{f}_i|\bar{e}_i) = \frac{\text{count}(\bar{f}_i, \bar{e}_i)}{\text{count}(\bar{e}_i)} p_w(\bar{f}_i|\bar{e}_i)$$

	f1	f2	f3	
NULL	--	--	##	
e1	##	--	--	
e2	--	##	--	
e3	--	##	--	

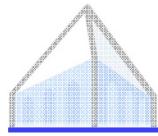
$$\begin{aligned} p_w(\bar{f}|\bar{e}, a) &= p_w(f_1 f_2 f_3 | e_1 e_2 e_3, a) \\ &= w(f_1|e_1) \\ &\quad \times \frac{1}{2}(w(f_2|e_2) + w(f_2|e_3)) \\ &\quad \times w(f_3|\text{NULL}) \end{aligned}$$





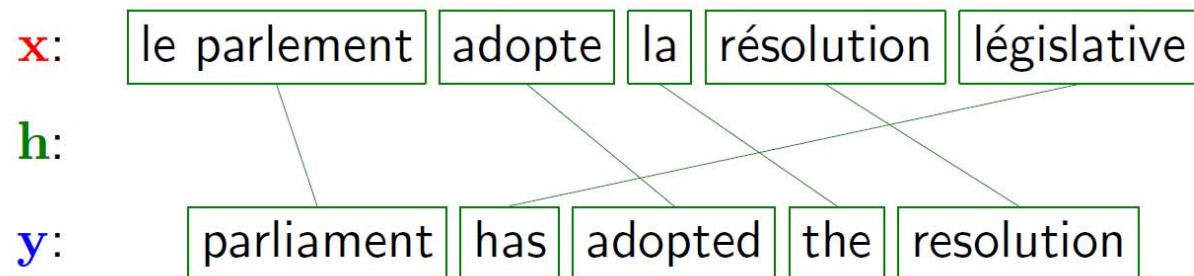
## Tuning for MT

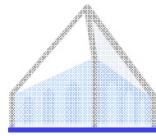
- Features encapsulate lots of information
  - Basic MT systems have around 6 features
  - $P(e|f)$ ,  $P(f|e)$ , lexical weighting, language model
- How to tune feature weights?
- Idea 1: Use your favorite classifier



# Why Tuning is Hard

- Problem 1: There are latent variables
  - Alignments and segmentations
  - Possibility: forced decoding (but it can go badly)

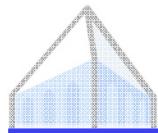




# Why Tuning is Hard

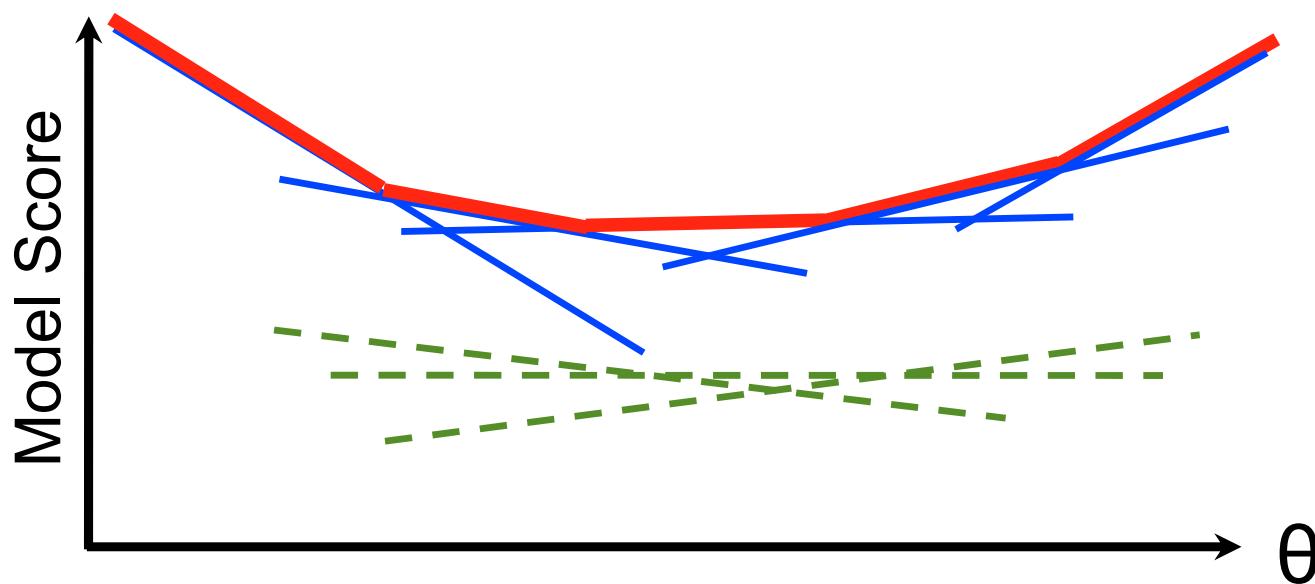
---

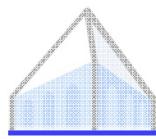
- **Problem 3: Computational constraints**
  - Discriminative training involves repeated decoding
  - Very slow! So people tune on sets much smaller than those used to build phrase tables



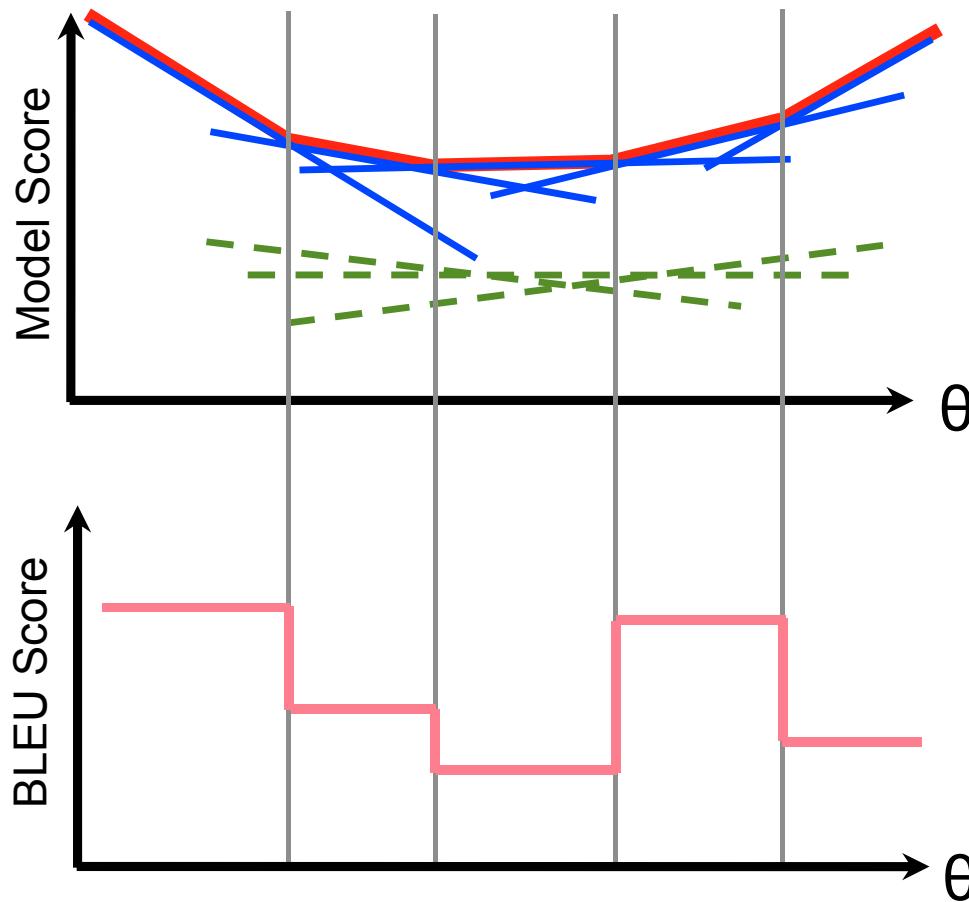
# Minimum Error Rate Training

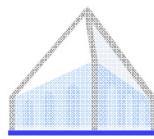
- Standard method: minimize BLEU directly (Och 03)
  - MERT is a discontinuous objective
  - Only works for max ~10 features, but works very well then
  - Here: k-best lists, but forest methods exist (Machery et al 08)
  - Recently, lots of alternatives being explored for more features



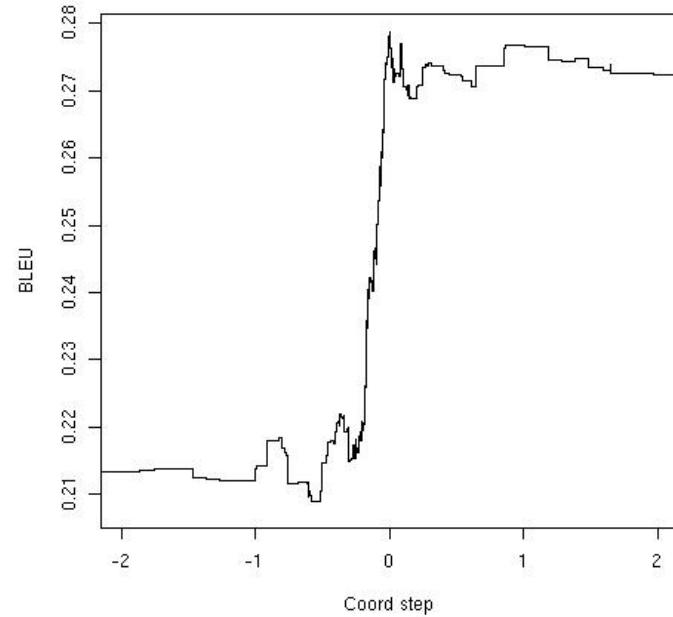
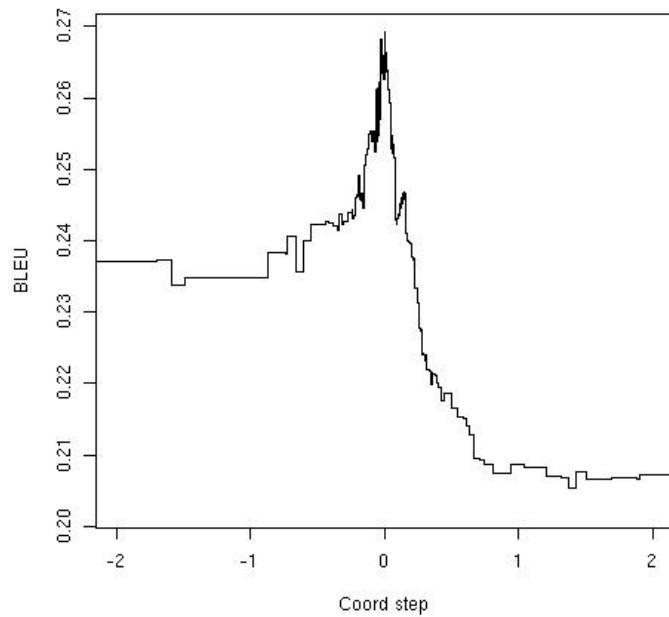


# MERT





# MERT





# Translating with Tree Transducers

**Input**

lo haré de muy buen grado .

**Output**

**Grammar**



# Translating with Tree Transducers

**Input**

lo haré de muy buen grado .

**Output**

**Grammar**

ADV → ⟨ de muy buen grado ; gladly ⟩

# Syntactic Models



# Translating with Tree Transducers

## Input

ADV  
lo haré de muy buen grado .

## Output

ADV  
I  
gladly

## Grammar

ADV → < de muy buen grado ; gladly >



# Translating with Tree Transducers

## Input

ADV  
lo haré de muy buen grado .

## Output

ADV  
I  
gladly

## Grammar

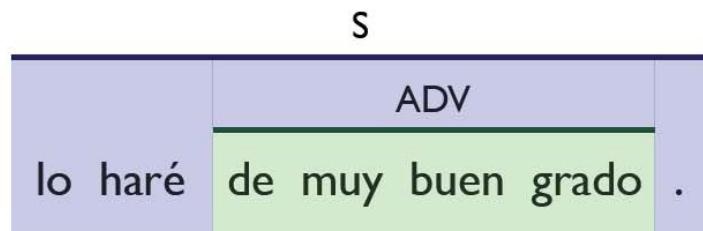
$s \rightarrow \langle \text{lo haré ADV .} ; \text{ I will do it ADV .} \rangle$

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

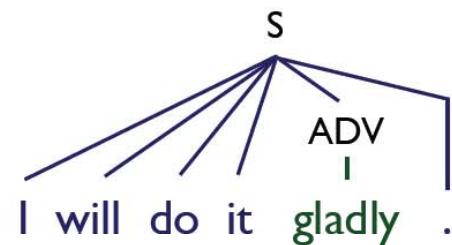


# Translating with Tree Transducers

## Input



## Output



## Grammar

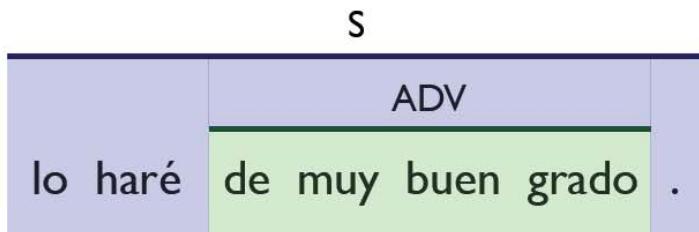
$s \rightarrow \langle \text{lo haré } \text{ADV} . ; \text{I will do it } \text{ADV} . \rangle$

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

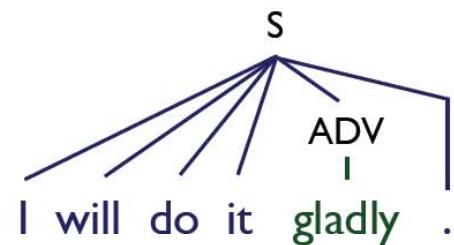


# Translating with Tree Transducers

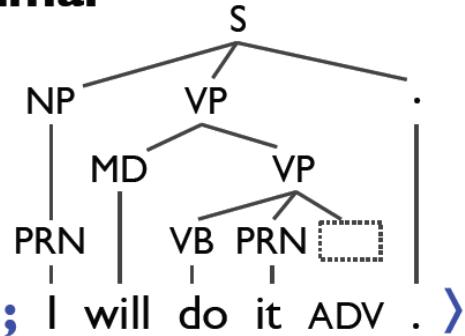
## Input



## Output



## Grammar



$s \rightarrow \langle \text{lo haré} \text{ ADV} . ; \text{ I will do it ADV} . \rangle$

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



# Translating with Tree Transducers

## Input

ADV  
lo haré de muy buen grado .

## Output

ADV  
I  
gladly

## Grammar

$s \rightarrow \langle \text{lo haré ADV .} ; \text{ I will do it ADV .} \rangle$

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



# Translating with Tree Transducers

## Input

ADV  
lo haré de muy buen grado .

## Output

ADV  
I  
gladly

## Grammar

VP → ⟨ lo haré ADV ; will do it ADV ⟩

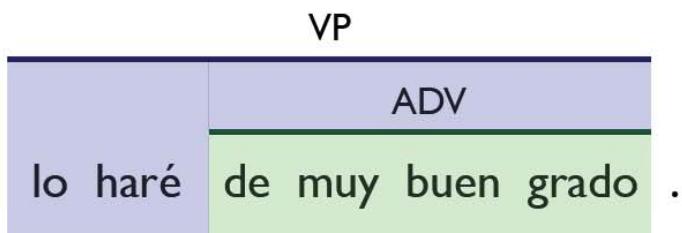
S → ⟨ lo haré ADV . ; I will do it ADV . ⟩

ADV → ⟨ de muy buen grado ; gladly ⟩

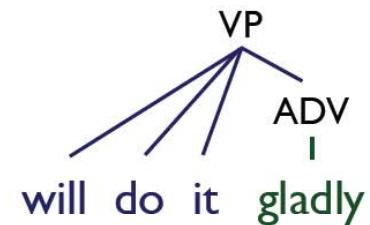


# Translating with Tree Transducers

## Input



## Output



## Grammar

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

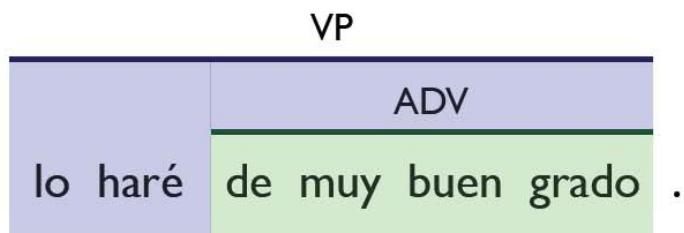
$\text{s} \rightarrow \langle \text{lo haré ADV .} ; \text{I will do it ADV .} \rangle$

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

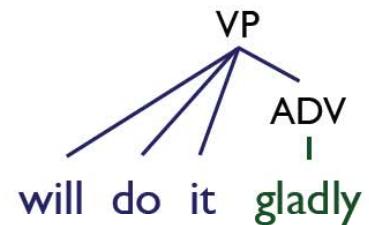


# Translating with Tree Transducers

## Input



## Output



## Grammar

$s \rightarrow \langle VP . ; | VP . \rangle$

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

$s \rightarrow \langle lo\haré\ ADV . ; | will\ do\ it\ ADV . \rangle$

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

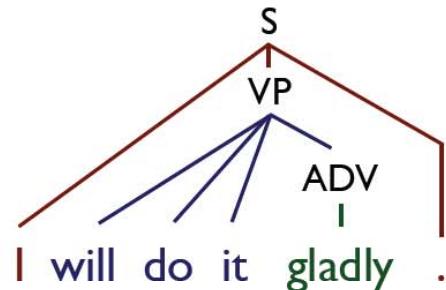


# Translating with Tree Transducers

## Input



## Output



## Grammar

$$S \rightarrow \langle VP . ; I VP . \rangle$$

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

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

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

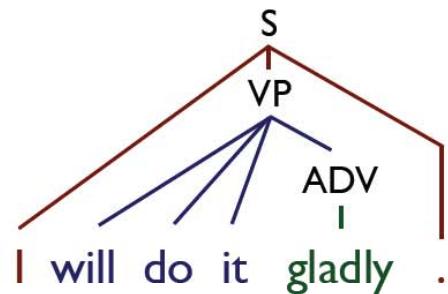


# Translating with Tree Transducers

## Input



## Output



## Grammar

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

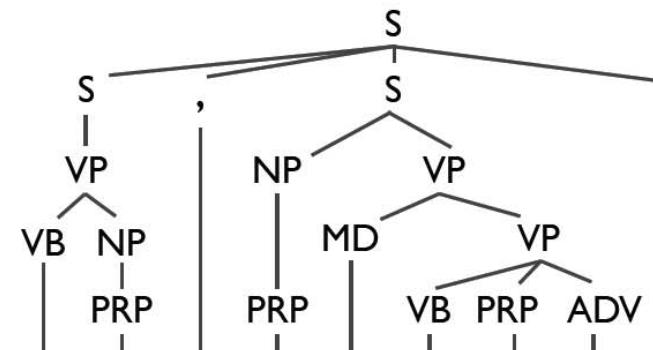
$VP \rightarrow \langle lo\ haré\ ADV ; will\ do\ it\ ADV \rangle$

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

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



# Learning Grammars for Translation



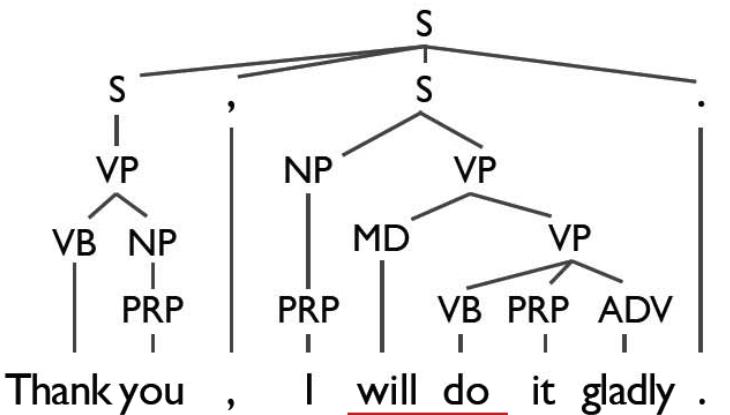
Thank you , I will do it gladly .


Gracias  
,  
lo  
haré  
de  
muy  
buen  
grado  
.

## Grammar Rules



# Learning Grammars for Translation



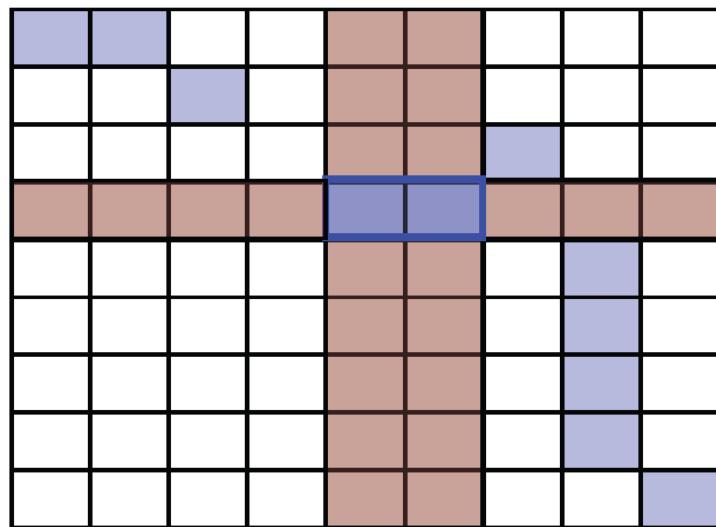
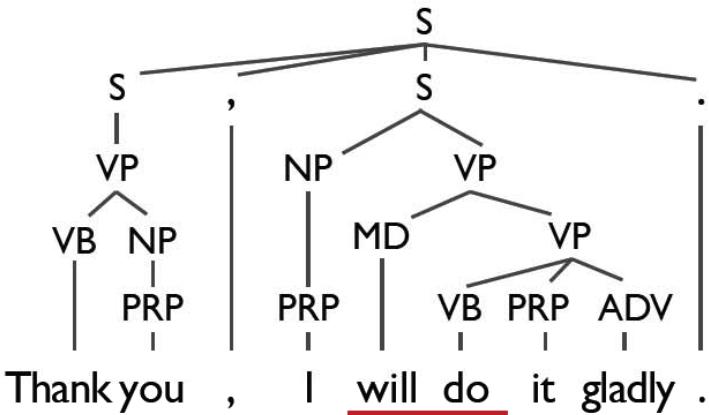
A grid-based diagram representing a neural network or matrix. It consists of several colored cells: light blue, dark blue, light brown, and dark brown. A horizontal blue bar is positioned across the middle row of the dark brown cells. The grid is composed of approximately 10 columns and 10 rows of cells.

Gracias  
,  
lo  
haré  
de  
muy  
buen  
grado  
.

## Grammar Rules



# Learning Grammars for Translation



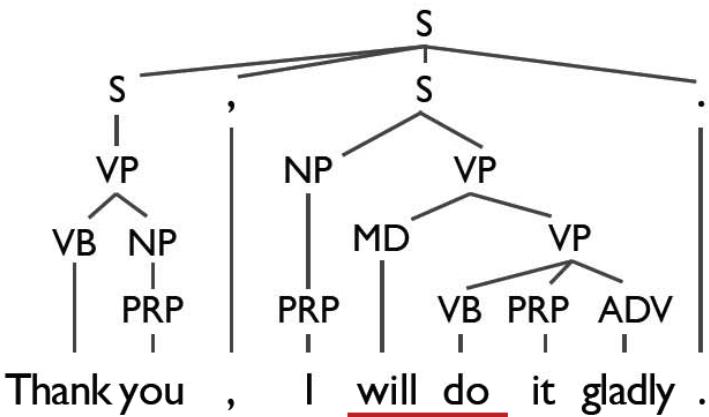
Gracias  
,  
lo  
haré  
de  
muy  
buen  
grado  
.

## Grammar Rules

⟨haré ; will do⟩



# Learning Grammars for Translation



A grid-based diagram representing a neural network or matrix. It consists of several rows and columns of squares. Some squares are colored light blue, some are light brown, and some are white. A prominent vertical column of light brown squares runs down the middle of the grid. A horizontal row of light blue squares spans across the middle of the grid. The grid is used to map words from one language to another, with arrows indicating the flow of information between specific cells.

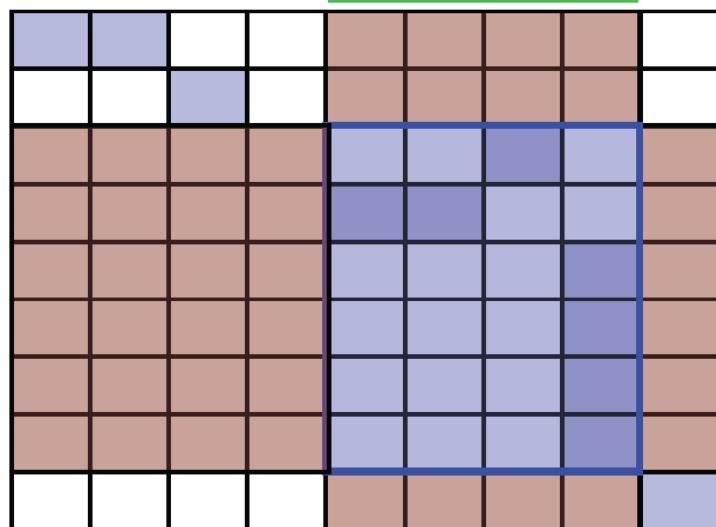
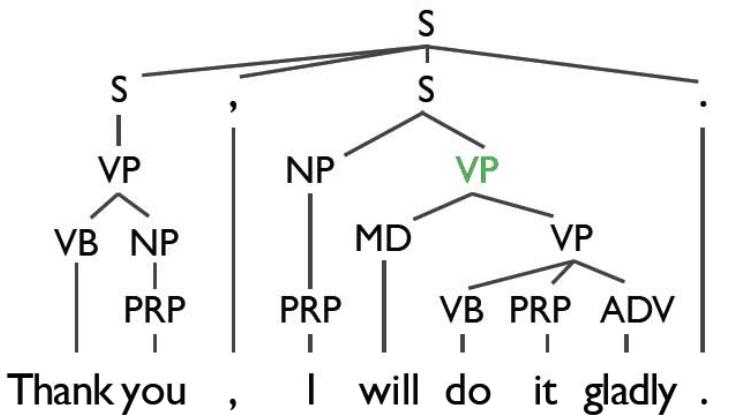
Gracias  
,  
lo  
haré  
de  
muy  
buen  
grado  
.

## Grammar Rules

~~haré ; will do~~



# Learning Grammars for Translation



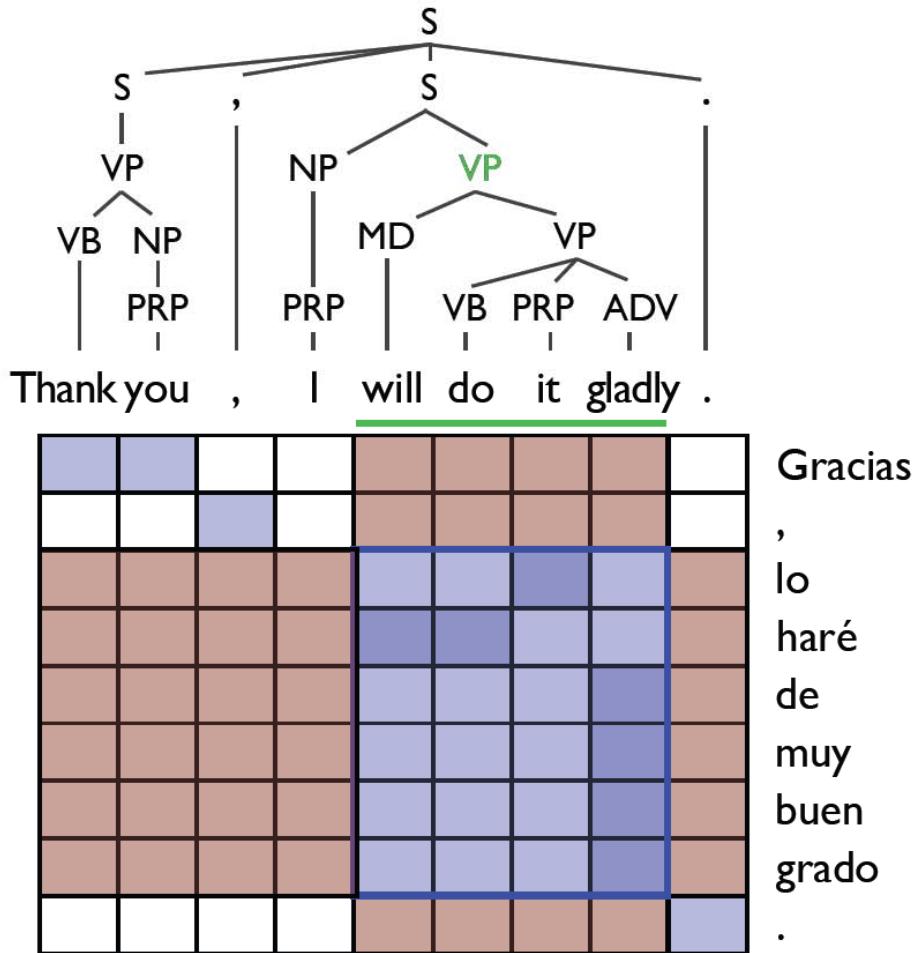
Gracias  
,  
lo  
haré  
de  
muy  
buen  
grado  
.

## Grammar Rules

~~haré ; will do~~



# Learning Grammars for Translation



## Grammar Rules

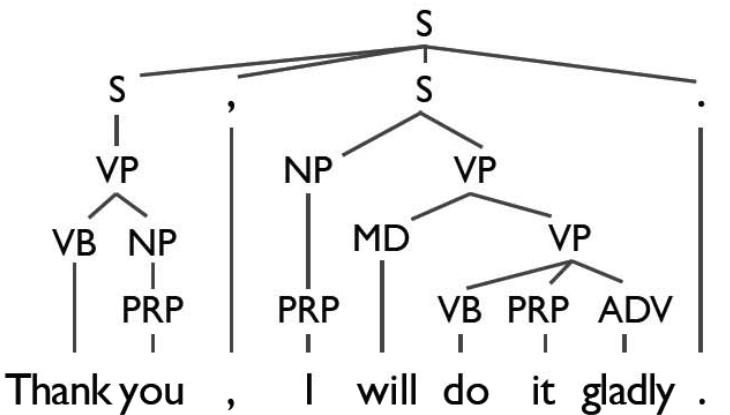
~~⟨haré ; will do⟩~~

VP →

⟨lo haré de ... grado ;  
will do it gladly⟩



# Learning Grammars for Translation




Gracias  
,  
lo  
haré  
de  
muy  
buen  
grado  
.

## Grammar Rules

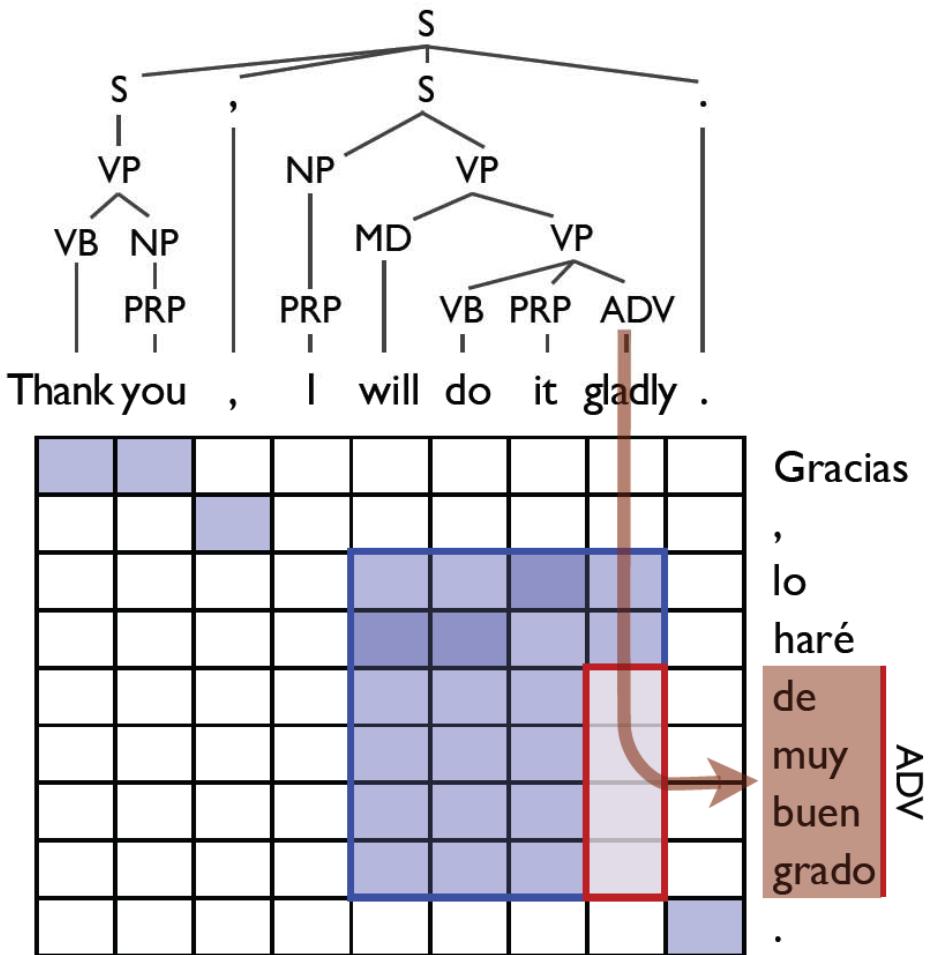
~~⟨haré ; will do⟩~~

VP →

⟨lo haré de ... grado ;  
will do it gladly⟩



# Learning Grammars for Translation

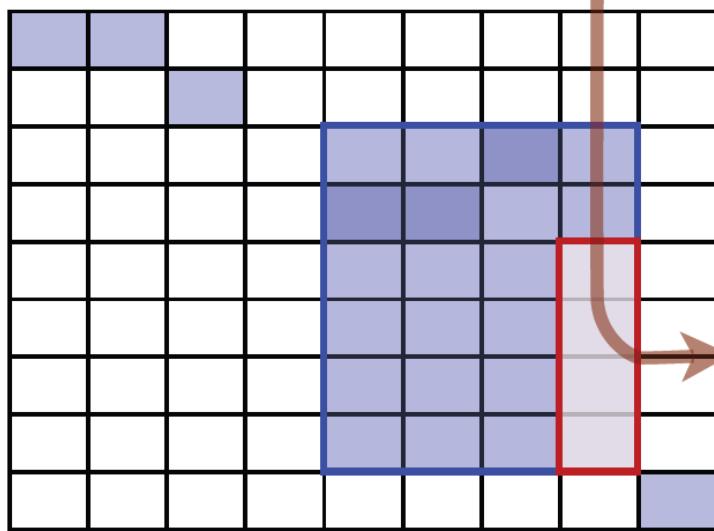
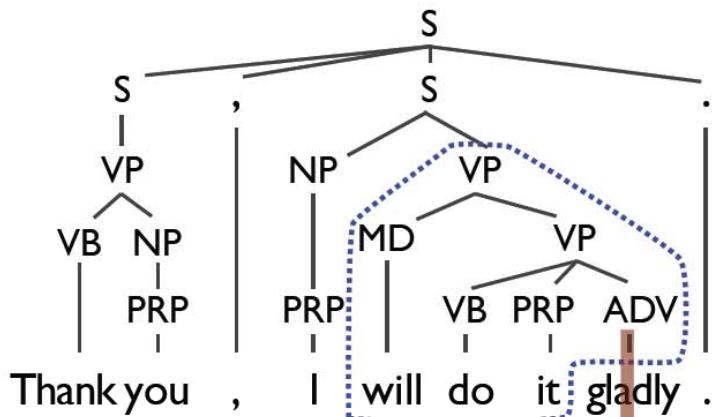


## Grammar Rules

~~haré ; will do~~  
VP →  
(lo haré de ... grado ;  
will do it gladly)



# Learning Grammars for Translation



Gracias

,  
lo  
haré  
de  
muy  
buen  
grado

ADV

## Grammar Rules

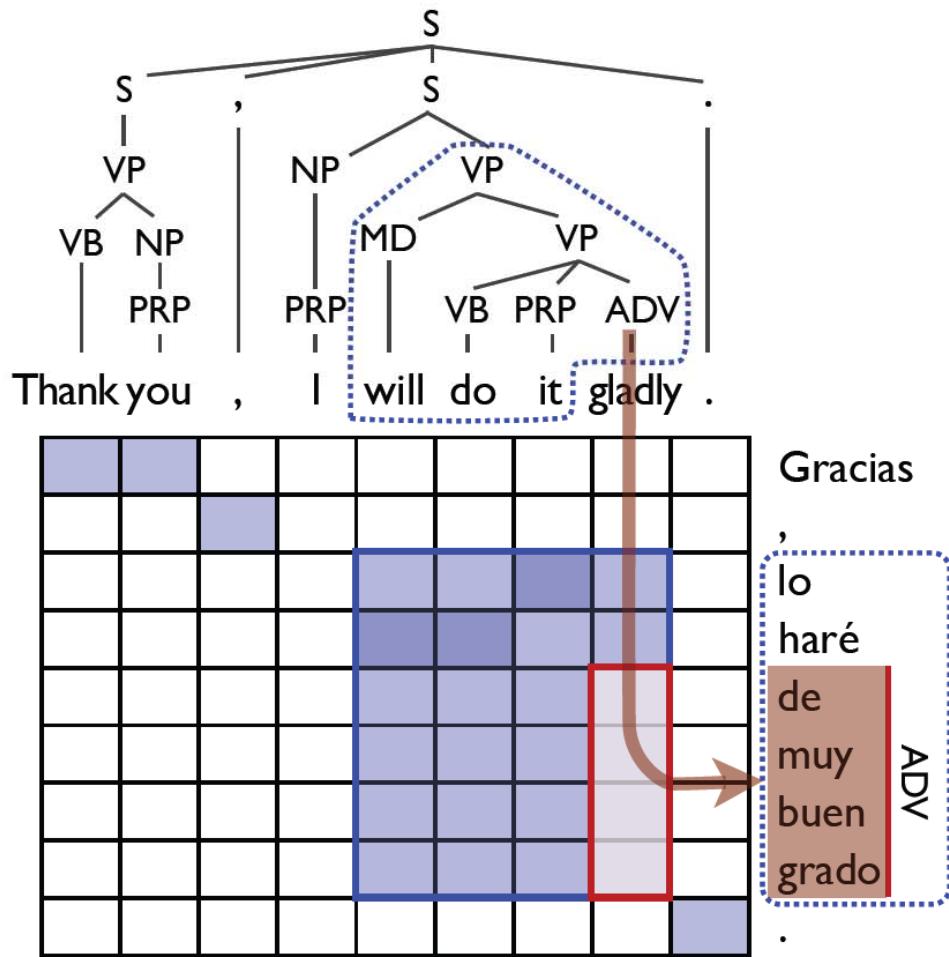
~~⟨haré ; will do⟩~~

VP →

⟨lo haré de ... grado ;  
will do it gladly⟩



# Learning Grammars for Translation



## Grammar Rules

~~⟨haré ; will do⟩~~

VP →  
 ⟨lo haré de ... grado ;  
 will do it gladly⟩

VP →  
 ⟨lo haré ADV ;  
 will do it ADV⟩



# The Size of Tree Transducer Grammars

Extracted a transducer grammar from a 220 million word bitext

Relativized the grammar to each test sentence

Kept all rules with at most 6 non-terminals



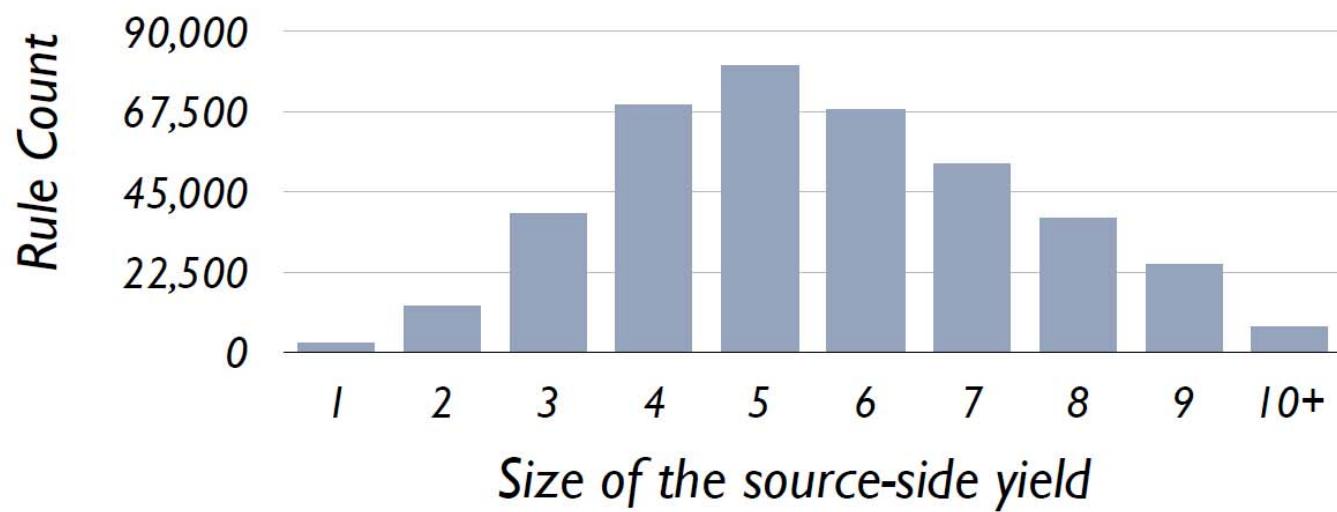
# The Size of Tree Transducer Grammars

Extracted a transducer grammar from a 220 million word bitext

Relativized the grammar to each test sentence

Kept all rules with at most 6 non-terminals

Rules matching an example 40-word sentence





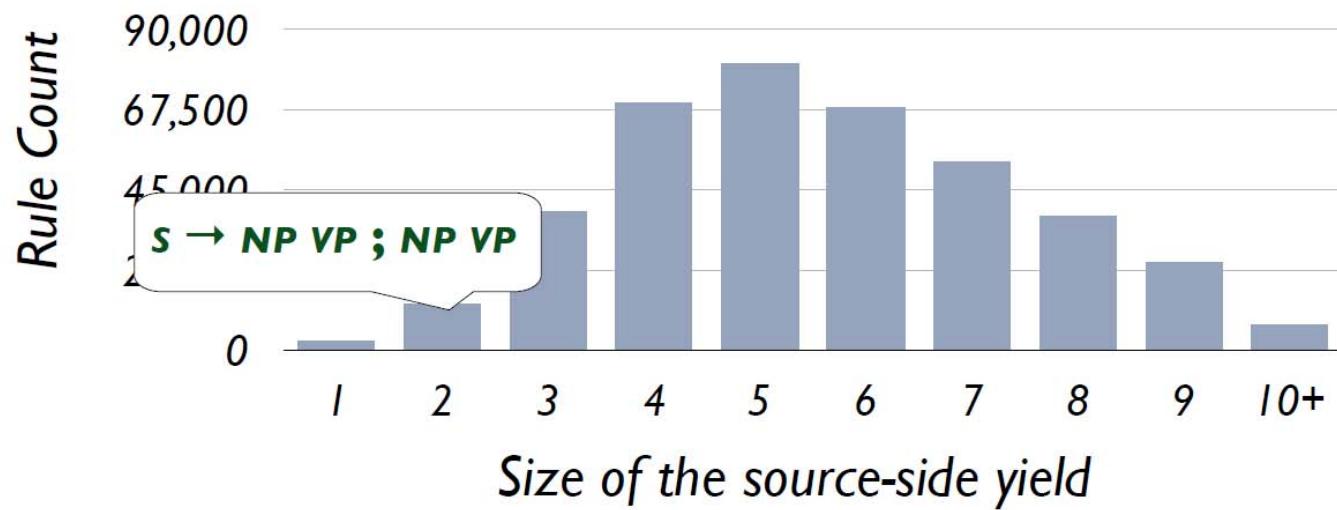
# The Size of Tree Transducer Grammars

Extracted a transducer grammar from a 220 million word bitext

Relativized the grammar to each test sentence

Kept all rules with at most 6 non-terminals

Rules matching an example 40-word sentence



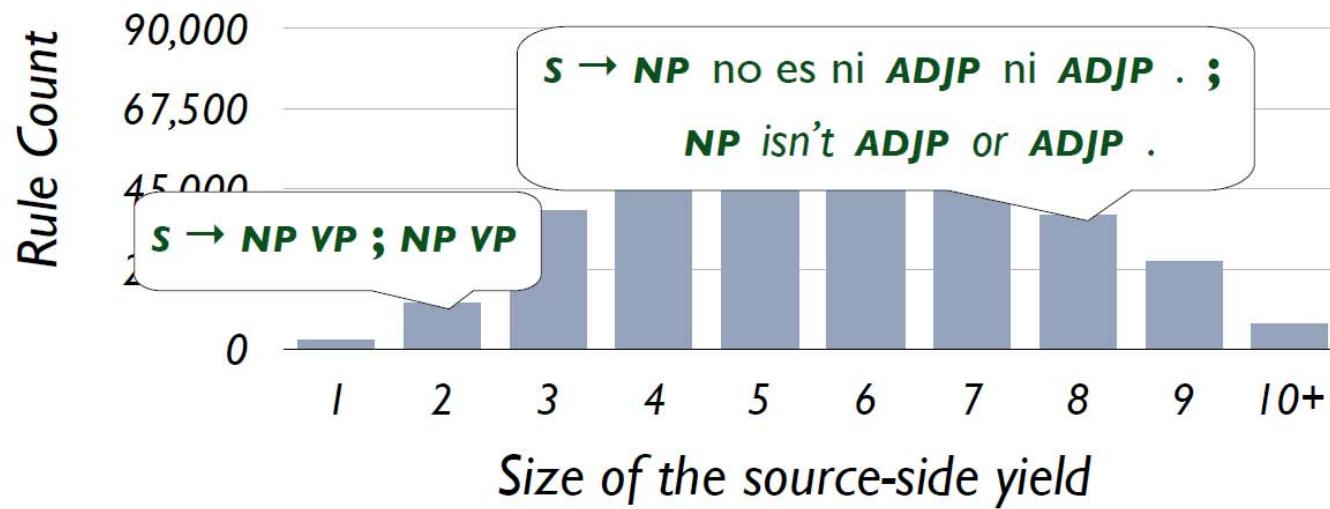
# The Size of Tree Transducer Grammars

Extracted a transducer grammar from a 220 million word bitext

Relativized the grammar to each test sentence

Kept all rules with at most 6 non-terminals

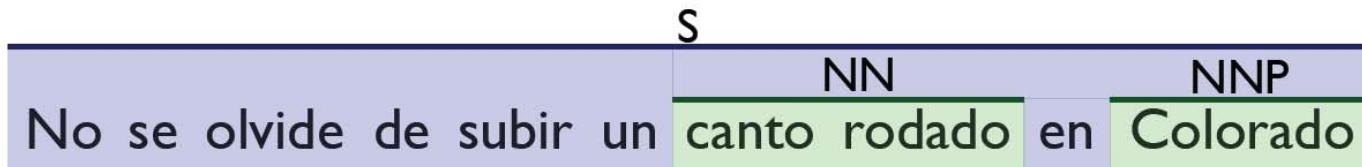
## Rules matching an example 40-word sentence



# Syntactic Decoding



# Tree Transducer Grammars



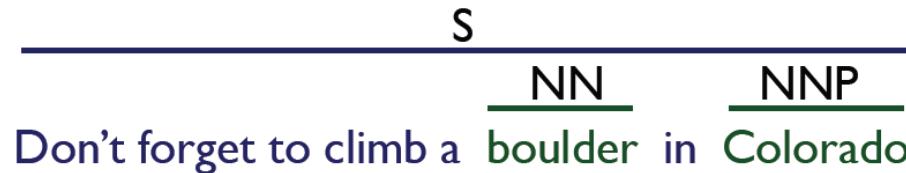
## Synchronous Grammar

**NNP** → Colorado ; *Colorado*

**NN** → canto rodado ; *boulder*

**S** → No se olvide de subir un **NN** en **NNP** ; *Don't forget to climb a NN in NNP*

## Output





# CKY-style Bottom-up Parsing

For each  
span length:



# CKY-style Bottom-up Parsing

For each span length:

For each span  $[i,j]$ :



# CKY-style Bottom-up Parsing

For each span length:

For each span  $[i,j]$ :

Apply all grammar rules to  $[i,j]$



# CKY-style Bottom-up Parsing

For each span length:

For each span  $[i,j]$ :

Apply all grammar rules to  $[i,j]$

Binary rule:  $X \rightarrow Y Z$



# CKY-style Bottom-up Parsing

For each span length:

For each span  $[i,j]$ :

Apply all grammar rules to  $[i,j]$

Binary rule:  $X \rightarrow Y Z$

Split points:  $i < k < j$

Operations:  $O(j - i)$

Time scales with: Grammar constant



## CKY-style Bottom-up Parsing

For each span length:

For each span  $[i,j]$ :

Apply all grammar rules to  $[i,j]$

$_i \text{ No se olvide de subir un canto rodado en Colorado }_j$



## CKY-style Bottom-up Parsing

For each span length:

For each span [i,j]:

Apply all grammar rules to [i,j]

**S** → No se **VB** de subir un **NN** en **NNP**

$_i$  No se olvide de subir un canto rodado en Colorado  $_j$



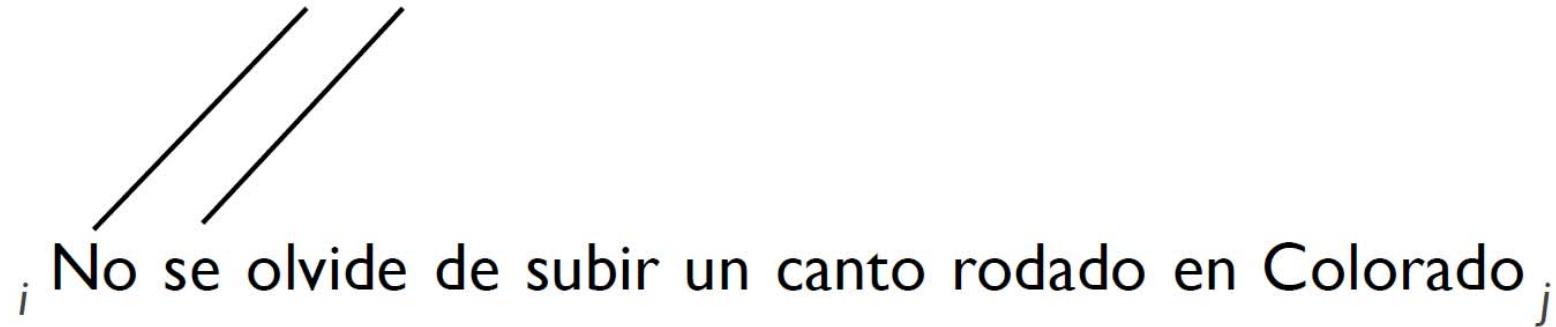
## CKY-style Bottom-up Parsing

For each span length:

For each span [i,j]:

Apply all grammar rules to [i,j]

**S** → No se **VB** de subir un **NN** en **NNP**



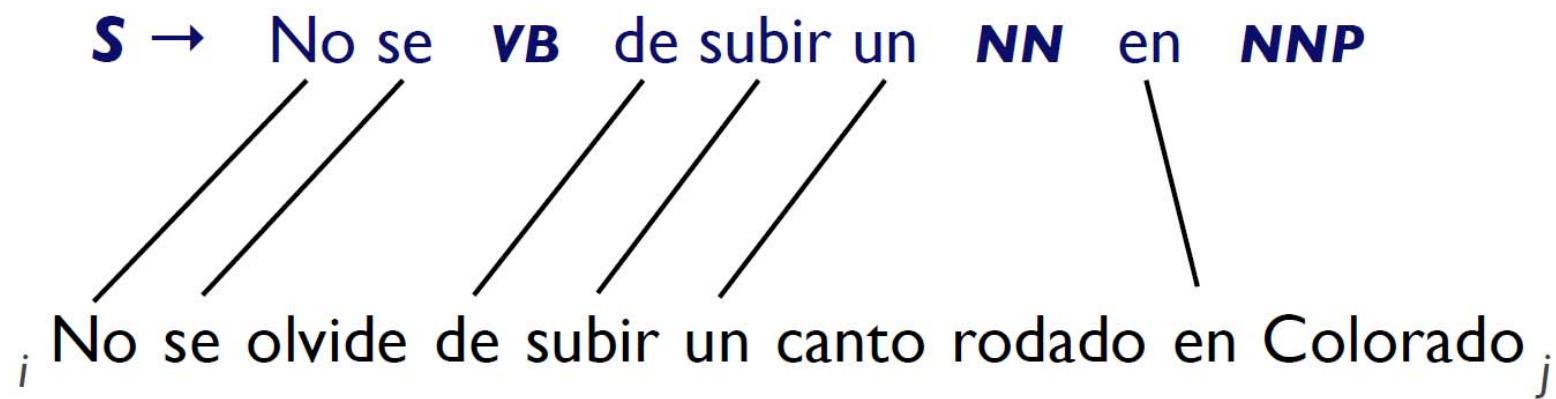


## CKY-style Bottom-up Parsing

For each span length:

For each span [i,j]:

Apply all grammar rules to [i,j]



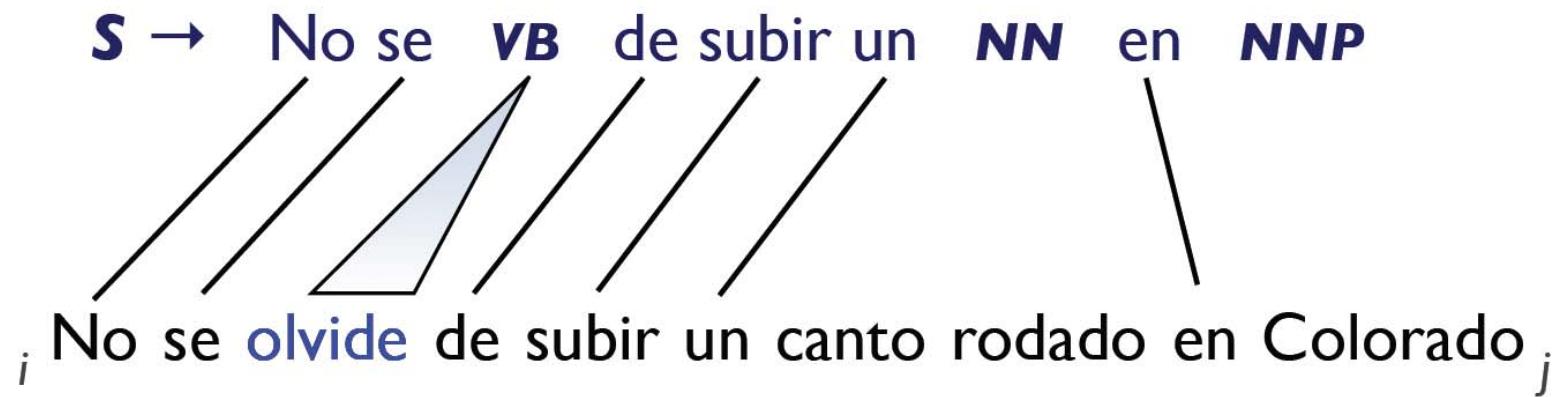


## CKY-style Bottom-up Parsing

For each span length:

For each span [i,j]:

Apply all grammar rules to [i,j]



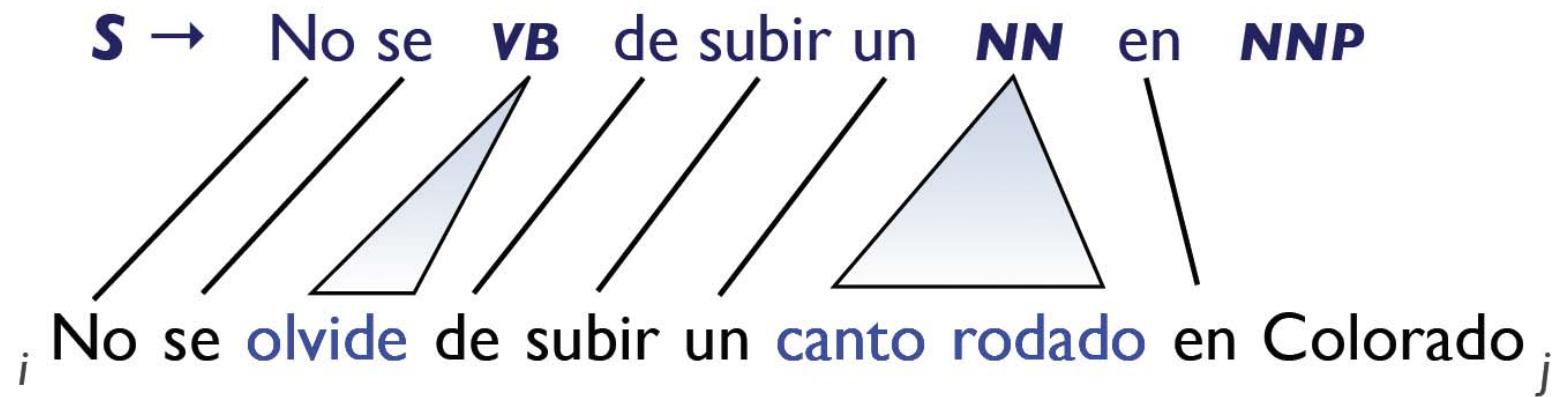


## CKY-style Bottom-up Parsing

For each span length:

For each span [i,j]:

Apply all grammar rules to [i,j]



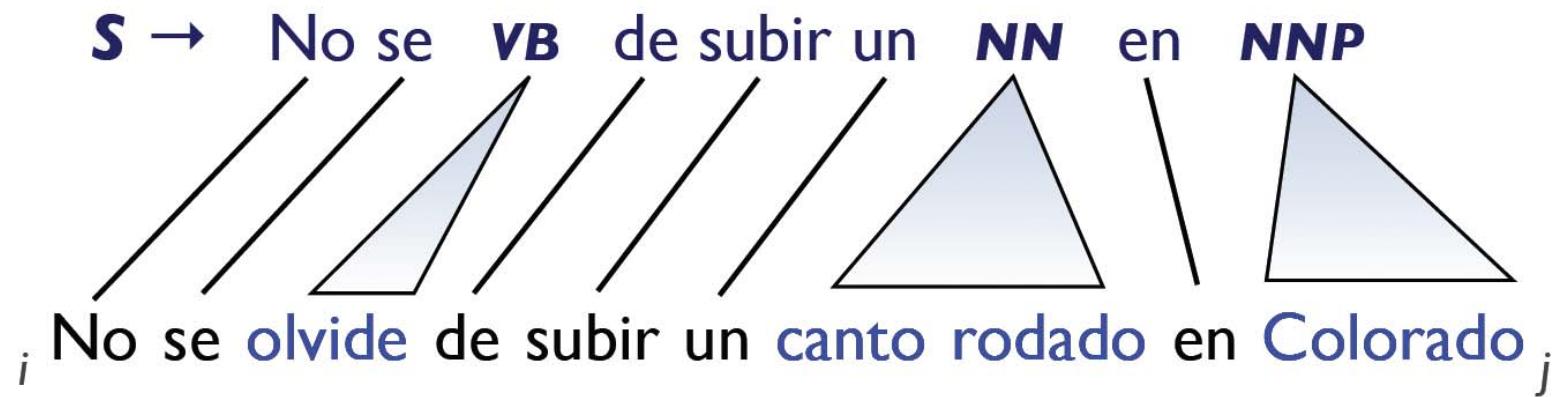


## CKY-style Bottom-up Parsing

For each span length:

For each span [i,j]:

Apply all grammar rules to [i,j]





## CKY-style Bottom-up Parsing

For each span length:

For each span [i,j]:

Apply all grammar rules to [i,j]



*Many untransformed lexical rules can be applied in linear time*



## CKY-style Bottom-up Parsing

For each span length:

For each span [i,j]:

Apply all grammar rules to [i,j]

**S** → No se **VP** **NP** **PP**

$i$  No se olvide de subir un canto rodado en Colorado  $j$



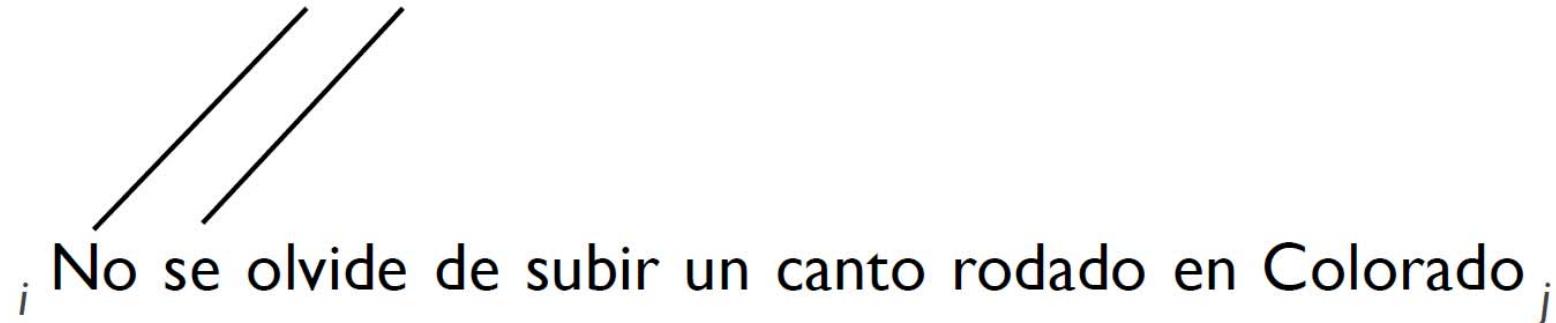
## CKY-style Bottom-up Parsing

For each span length:

For each span [i,j]:

Apply all grammar rules to [i,j]

**S** → No se **VP** **NP** **PP**

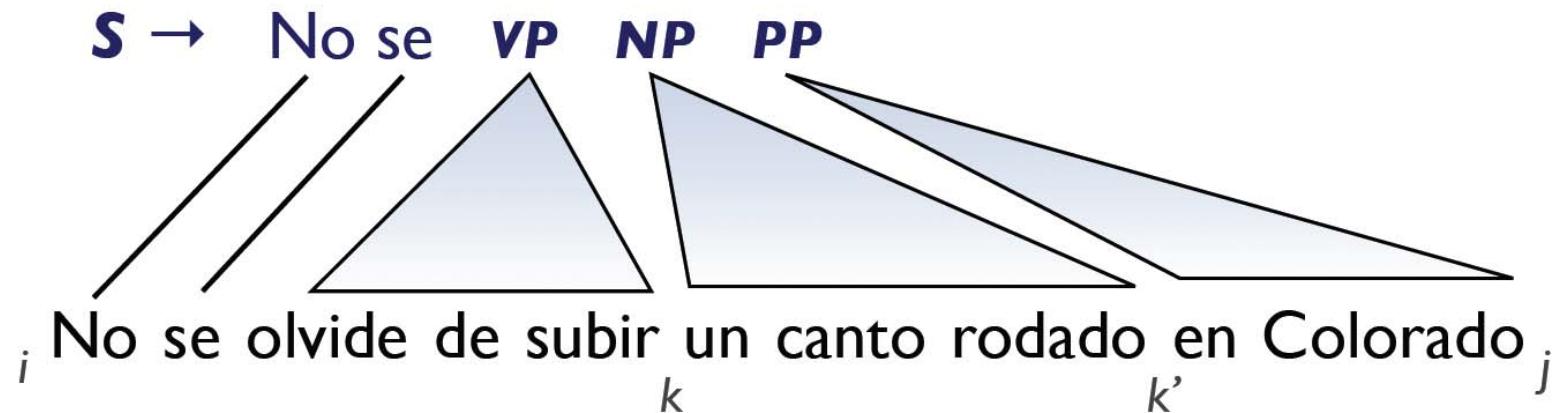


## CKY-style Bottom-up Parsing

For each span length:

For each span [i,j]:

Apply all grammar rules to [i,j]



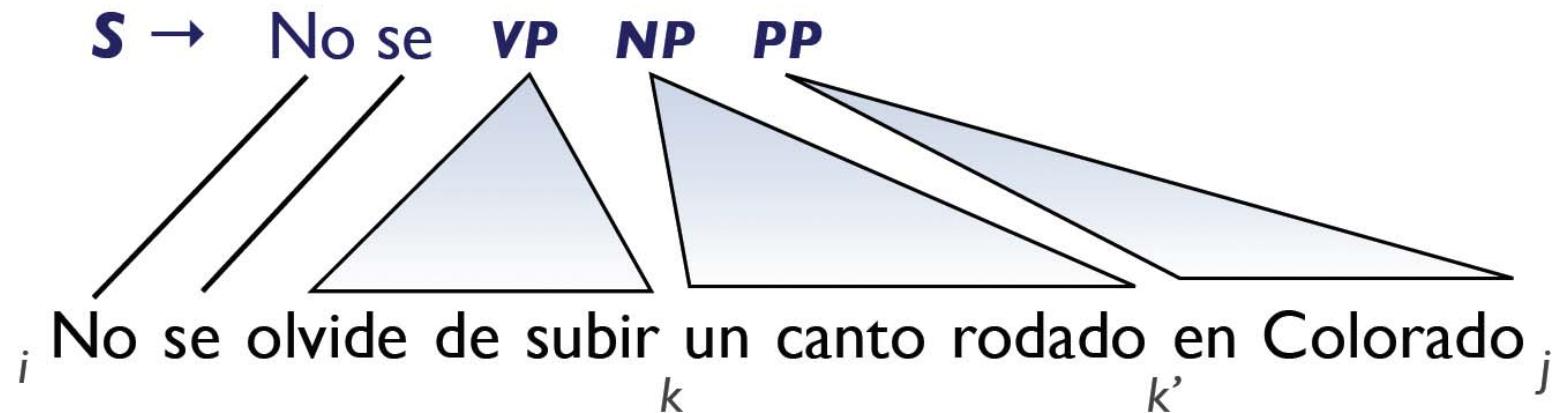


## CKY-style Bottom-up Parsing

For each span length:

For each span [i,j]:

Apply all grammar rules to [i,j]



**Problem:** Applying adjacent non-terminals is slow



# Eliminating Non-terminal Sequences

## Lexical Normal Form (LNF)

- (a) lexical rules have at most one adjacent non-terminal
- (b) all unlexicalized rules are binary.

Original rule:  $S \rightarrow \text{No se } VB\ VB \text{ un } NN\ PP$

Transformed rules:  $S \rightarrow \text{No se } VB\sim VB \text{ un } NN\sim PP$

$VB\sim VB \rightarrow VB\ VB$

$NN\sim PP \rightarrow NN\ PP$

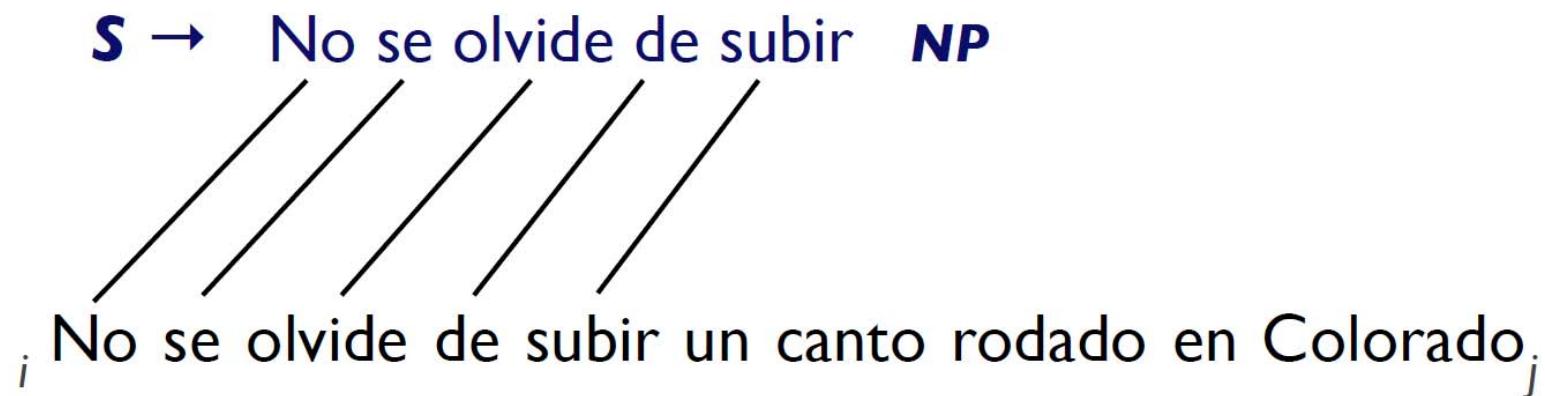
Parsing stages:

- Lexical rules are applied by matching
- Unlexicalized rules are applied by iterating over split points



# Speeding up Lexical Rule Application

**Problem:** Lexical rules can apply to many spans





# Speeding up Lexical Rule Application

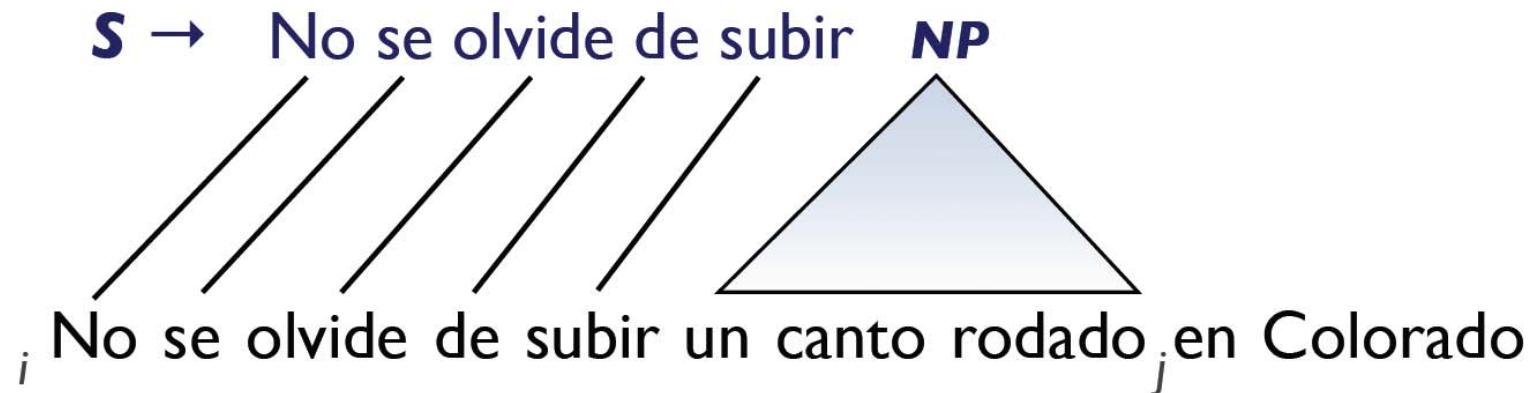
**Problem:** Lexical rules can apply to many spans





# Speeding up Lexical Rule Application

**Problem:** Lexical rules can apply to many spans





# Speeding up Lexical Rule Application

**Problem:** Lexical rules can apply to many spans



# Flexible Syntax

## Soft Syntactic MT: From Chiang 2010



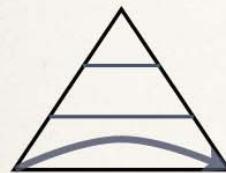
*reference:* An official from Japan 's science and technology ministry said , " We are highly encouraged by Abraham 's comment .

*Hiero:* Officials of the Japanese ministry of education and science , " said Abraham speeches , we are deeply encouraged by .

*string-to-tree:* Japan 's ministry of education , culture , sports , science and technology , " Abraham 's statement , which is most encouraging , " the official said .

# Previous work

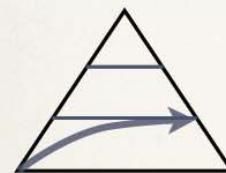
---



string-to-string

ITG (Wu 1997)

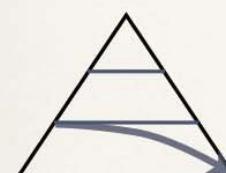
Hiero  
(Chiang 2005)



string-to-tree

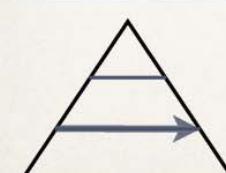
Yamada & Knight  
2001

Galley et al  
2004/2006



tree-to-string

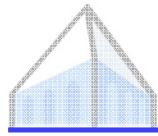
Huang et al 2006  
Y Liu et al 2006



tree-to-tree

DOT (Poutsma 2000)  
Eisner 2003

Stat-XFER (Lavie et al 2008)  
M Zhang et al. 2008  
Y Liu et al., 2009



## Hiero Rules

---

$S \rightarrow \langle S_{\textcolor{brown}{1}} X_{\textcolor{teal}{2}}, S_{\textcolor{brown}{1}} X_{\textcolor{teal}{2}} \rangle$

$S \rightarrow \langle X_{\textcolor{teal}{1}}, X_{\textcolor{teal}{1}} \rangle$

$X \rightarrow \langle \text{yu } X_{\textcolor{teal}{1}} \text{ you } X_{\textcolor{teal}{2}}, \text{have } X_{\textcolor{teal}{2}} \text{ with } X_{\textcolor{teal}{1}} \rangle$

$X \rightarrow \langle X_{\textcolor{teal}{1}} \text{ de } X_{\textcolor{teal}{2}}, \text{the } X_{\textcolor{teal}{2}} \text{ that } X_{\textcolor{teal}{1}} \rangle$

$X \rightarrow \langle X_{\textcolor{teal}{1}} \text{ zhiyi, one of } X_{\textcolor{teal}{1}} \rangle$

$X \rightarrow \langle \text{Aozhou, Australia} \rangle$

$X \rightarrow \langle \text{shi, is} \rangle$

$X \rightarrow \langle \text{shaoshu guojia, few countries} \rangle$

$X \rightarrow \langle \text{bangjiao, diplomatic relations} \rangle$

$X \rightarrow \langle \text{Bei Han, North Korea} \rangle$

From [Chiang et al, 2005]

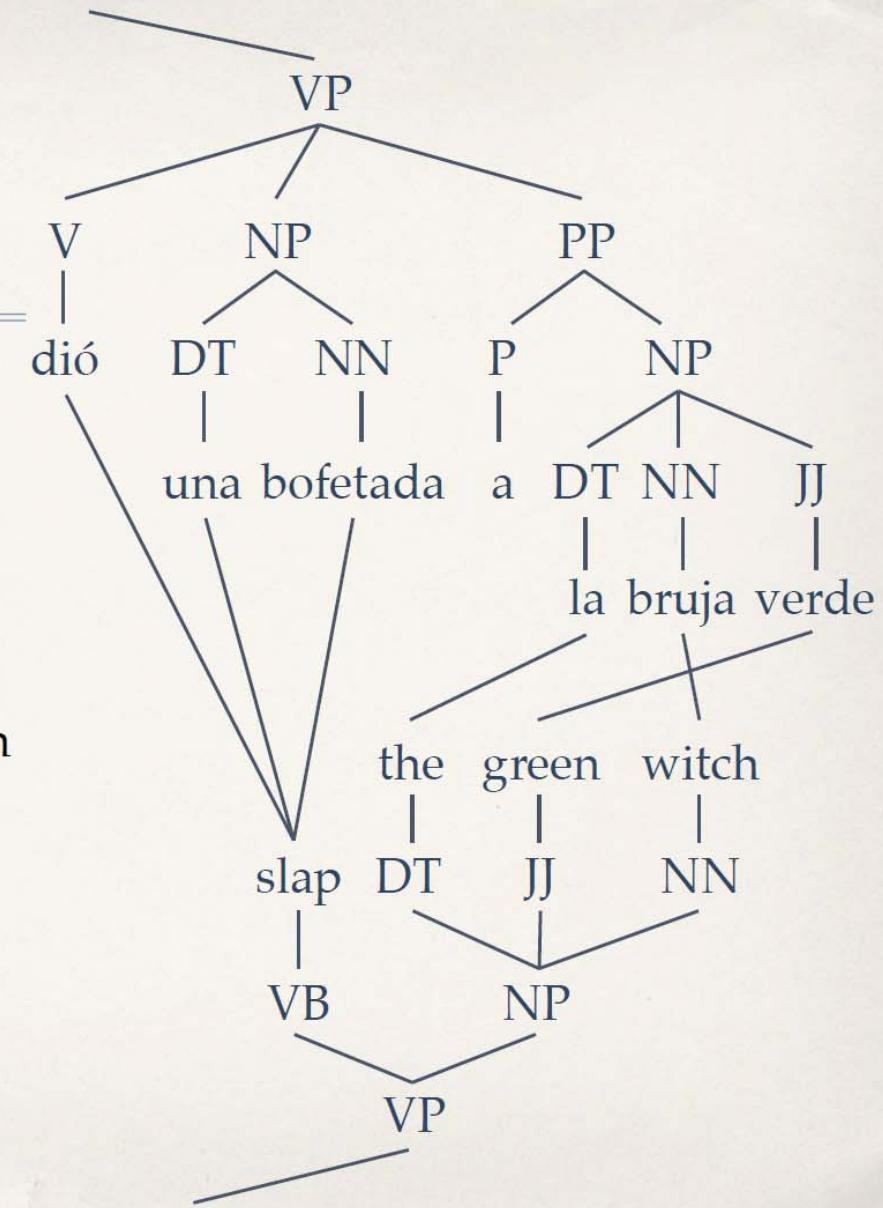
# STSG extraction

## 1. Phrases

- respect word alignments
- are syntactic constituents on *both* sides

## 2. Phrase pairs form rules

## 3. Subtract phrases to form rules



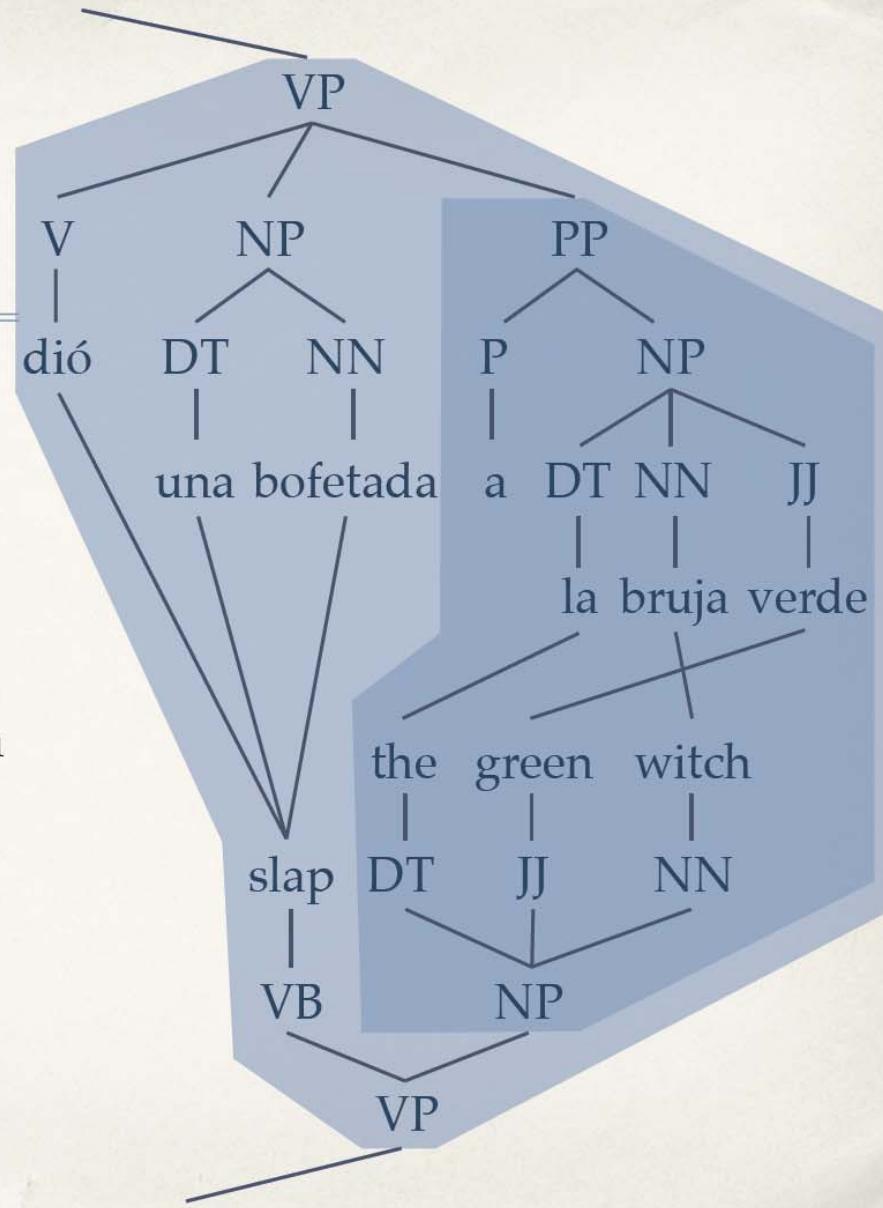
# STSG extraction

## 1. Phrases

- respect word alignments
- are syntactic constituents on *both* sides

## 2. Phrase pairs form rules

## 3. Subtract phrases to form rules



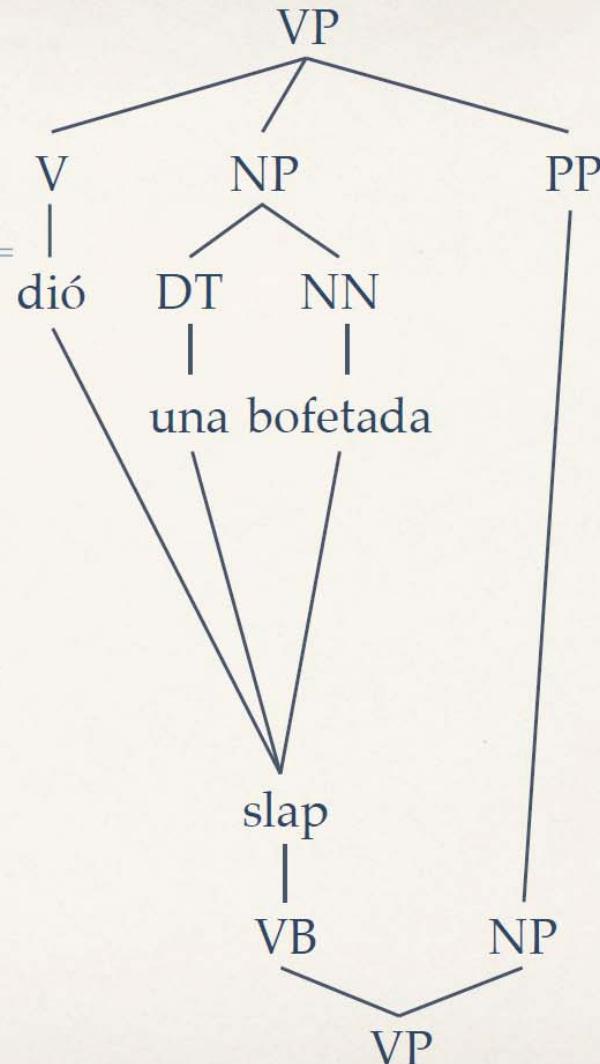
# STSG extraction

## 1. Phrases

- respect word alignments
- are syntactic constituents on *both* sides

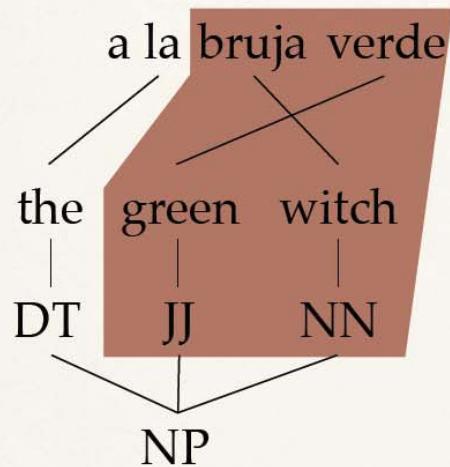
## 2. Phrase pairs form rules

## 3. Subtract phrases to form rules

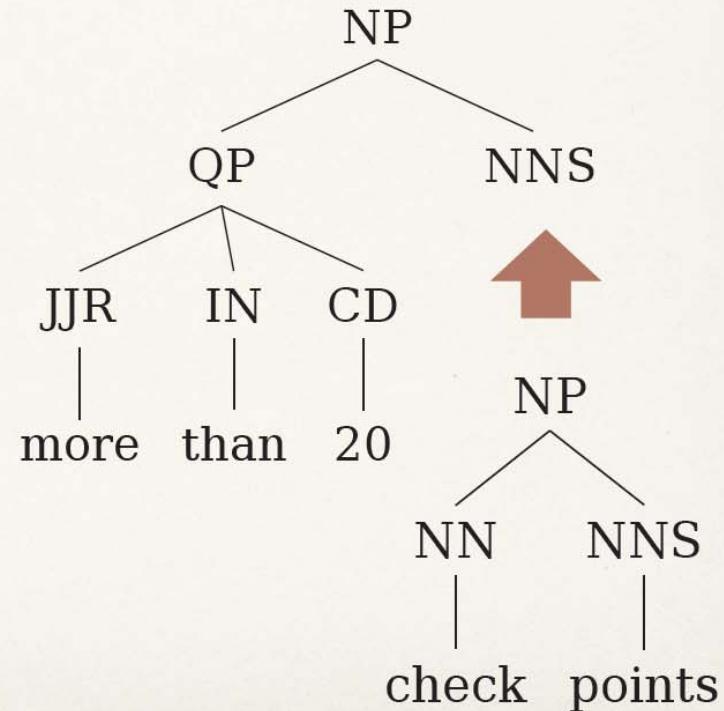


# Why is tree-to-tree hard?

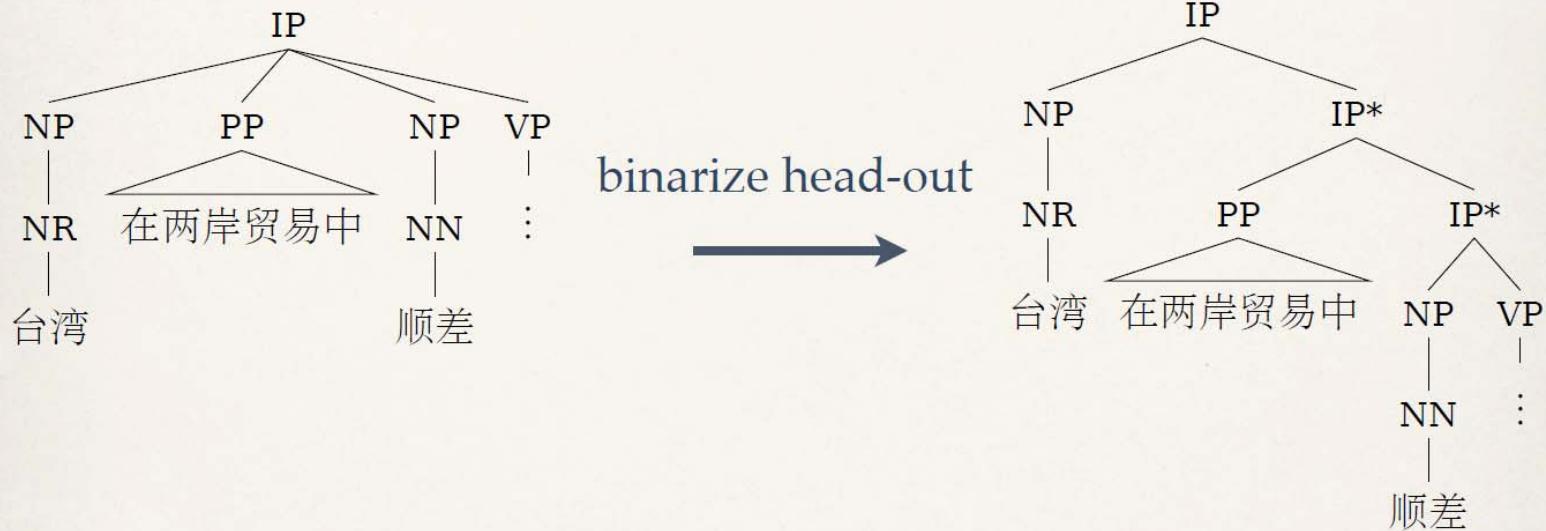
too few rules



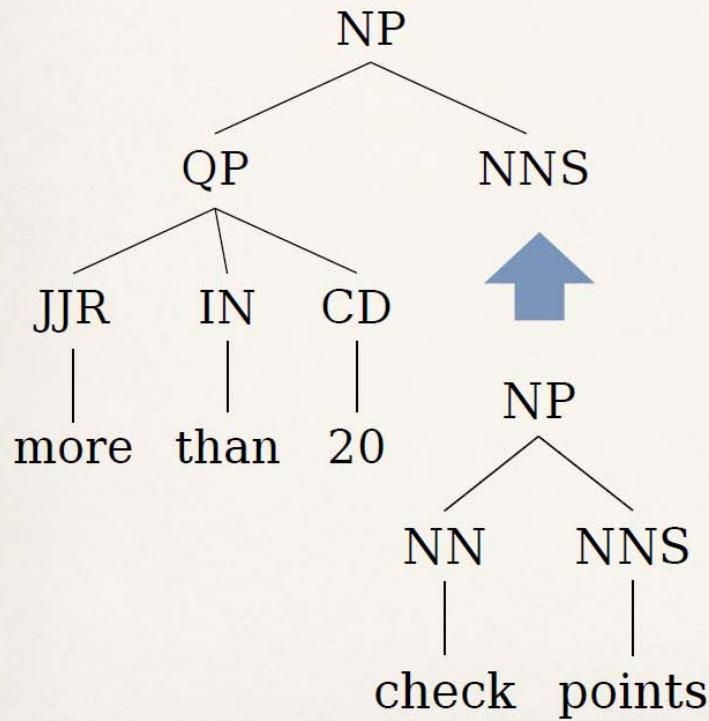
too few derivations



# Extracting more rules

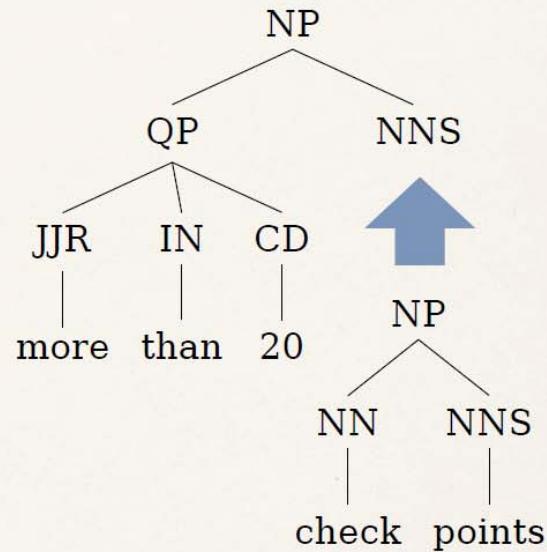
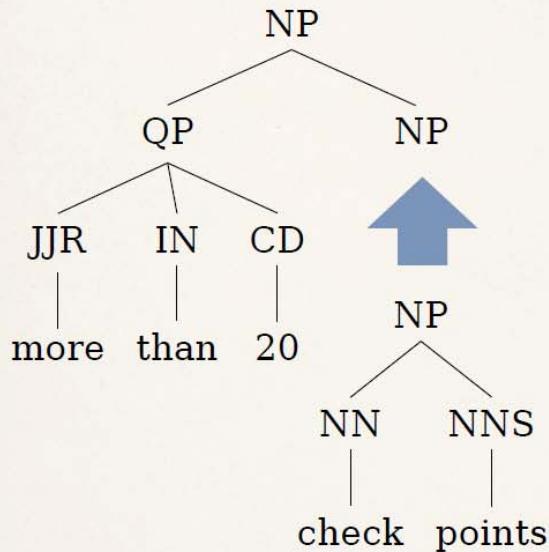


# Allow more derivations



- ❖ STSG: allow only matching substitutions
- ❖ Hiero-like: allow any substitutions
- ❖ Let the model learn to choose:
  - ❖ matching substitutions
  - ❖ mismatching substitutions
  - ❖ monotone phrase-based

# Allow more derivations



*fire subst:NP→NP*  
*fire subst:match*

*fire subst:NNS→NP*  
*fire subst:unmatch*

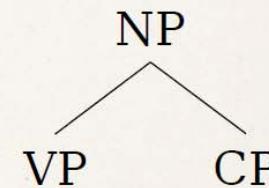
# Allow more derivations

Hiero-like decoding

$$[X, i, j] \quad [X, j+1, k]$$
$$[X, i, k]$$

$X \rightarrow X$  的  $X$

STSG decoding

$$[VP, i, j] \quad [NP, j+1, k]$$
$$[NP, i, k]$$


fuzzy STSG  
decoding

$$[A, i, j] \quad [B, j+1, k]$$
$$[NP, i, k]$$

```
graph TD; NP --- VP; NP --- CP; CP --- DEC; CP --- NP; NP --- 的;
```

# Results

extraction	Chinese-English			Arabic-English		
	rules	feats	BLEU	rules	feats	BLEU
Hiero	440M	1k	23.7	790M	1k	48.9
fuzzy STSG	50M	5k	23.9	38M	5k	47.5
fuzzy STSG +binarize	64M	5k	24.3	40M	6k	48.1
fuzzy STSG +SAMT	440M	160k	24.3	790M	130k	49.7

# Example tree-to-tree translation

---

日本 文部科学省官员 表示 , " 亚伯拉罕 的 发言 , 令 我们 深感 鼓舞  
Japan MEXT official said , " Abraham 's comment make us deeply-feel courage

*reference: An official from Japan 's science and technology ministry said , " We are highly encouraged by Abraham 's comment .*

*Hiero: Officials of the Japanese ministry of education and science , " said Abraham speeches , we are deeply encouraged by .*

*string-to-tree: Japan 's ministry of education , culture , sports , science and technology , " Abraham 's statement , which is most encouraging , " the official said .*

*Fuzzy STSG, binarize: Officials of the Japanese ministry of education , culture , sports , science and technology , said , " we are very encouraged by the speeches of Abraham .*

