

# Syntax-based decoding

JHU Machine Translation class  
April 1, 2014

# Administrative

- Homework 4 out, due April 14
- Final project proposals due today

# Where do grammars come from?

- We left off on Thursday with
  - a formalism for describing the relationship between two languages,
  - an loosely-sketched algorithm for producing translations
- Questions for today:
  - Where do synchronous grammars come from?
  - How do we decode with an ngram language model?

# Data-driven grammar extraction

- Grammar rules are not written by hand, they are extracted from bilingual parallel corpora

## Arabic

فالتعذيب لا يزال يمارس على نطاق واسع

وتتم عمليات الاعتقال والاحتجاز دون سبب بصورة روتينية

وحان وقت التحلى بالبصيرة والشجاعة السياسية .

...

## English

Torture is still being practised on a wide scale.

Arrest and detention without cause take place routinely.

This is a time for vision and political courage

...

## Chinese

我国 能源 原材料 工业 生产 大幅度 增长 .

非国大 要求 阻止 更多 被 拘留 人员 死亡 .

...

## English

China's energy and raw materials production up.

ANC calls for steps to prevent deaths in police custody .

...

# Hiero

- Consider the redundancy in this phrase table

Spanish	English
la bruja verde	the green witch
la bruja roja	the red witch
la bruja azul	the blue witch

- What generalization is missing?

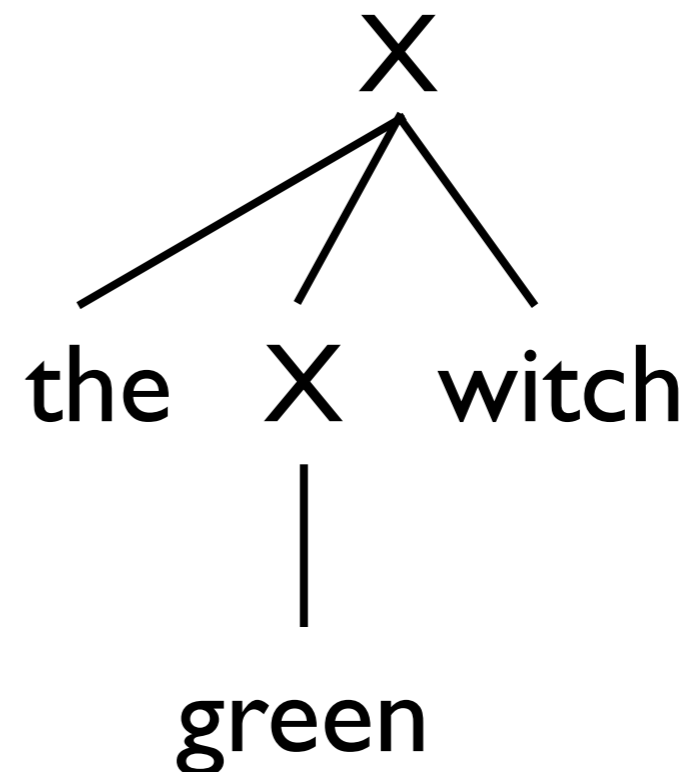
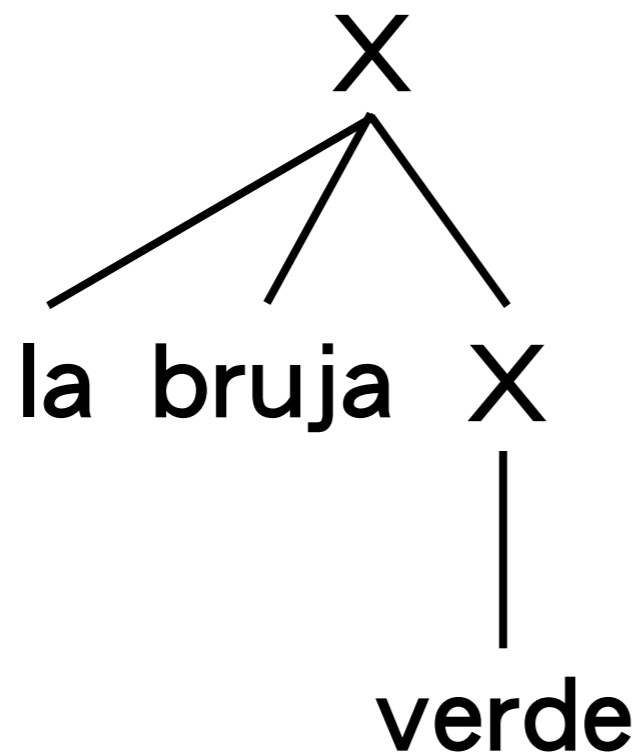
# Hiero

- Synchronous grammar rules

$X \rightarrow \text{la bruja } X_{(1)} \mid \mid \mid \text{the } X_{(1)} \text{ witch}$

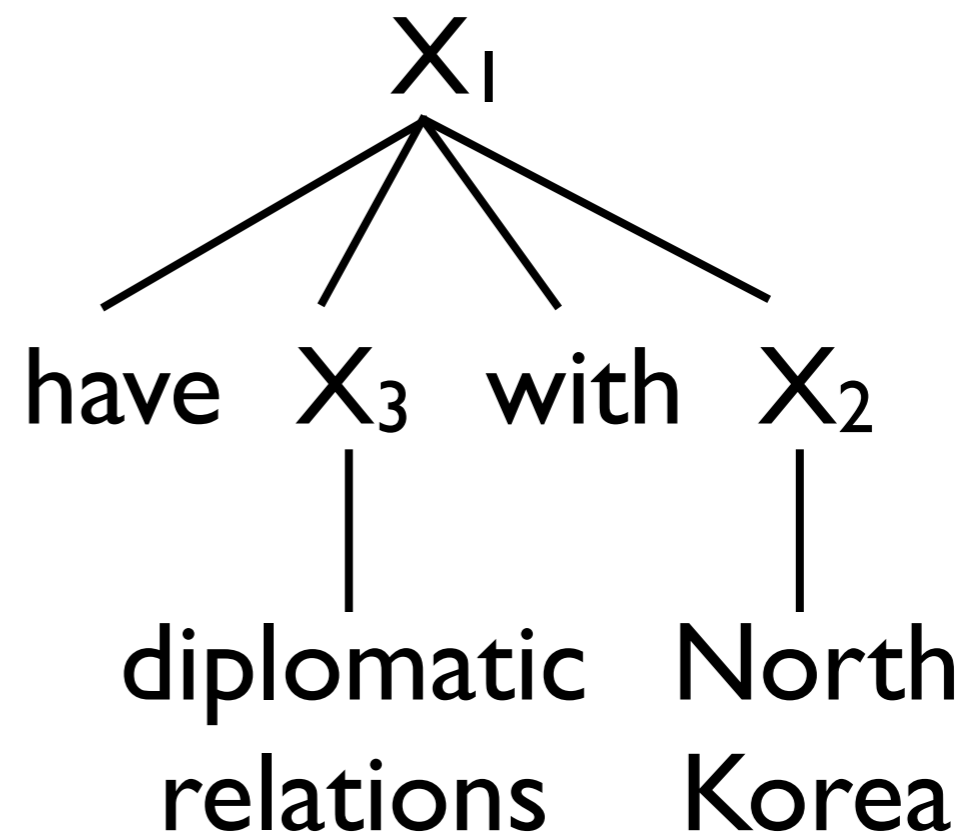
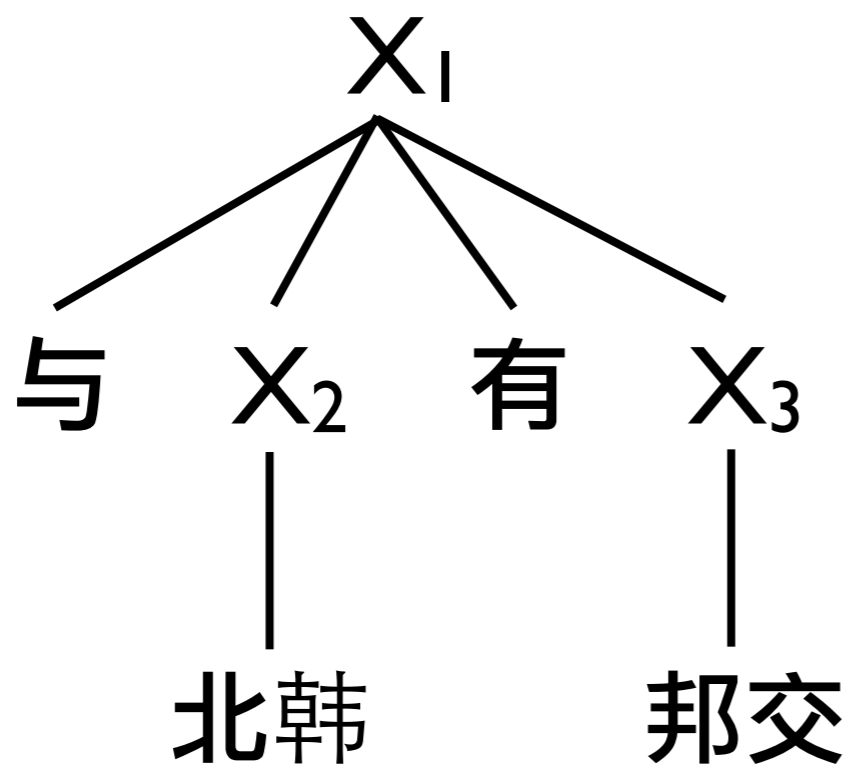
$X \rightarrow \text{verde} \mid \mid \mid \text{green}$

- As a tree



# Hiero-style SCFG rules

- Most common type of SCFG in SMT is Hiero which has rules w/one non-terminal symbol
- Not as nice as linguistically motivated rules, does not capture the reordering in Urdu



# Hiero

- Consider the redundancy in this phrase table

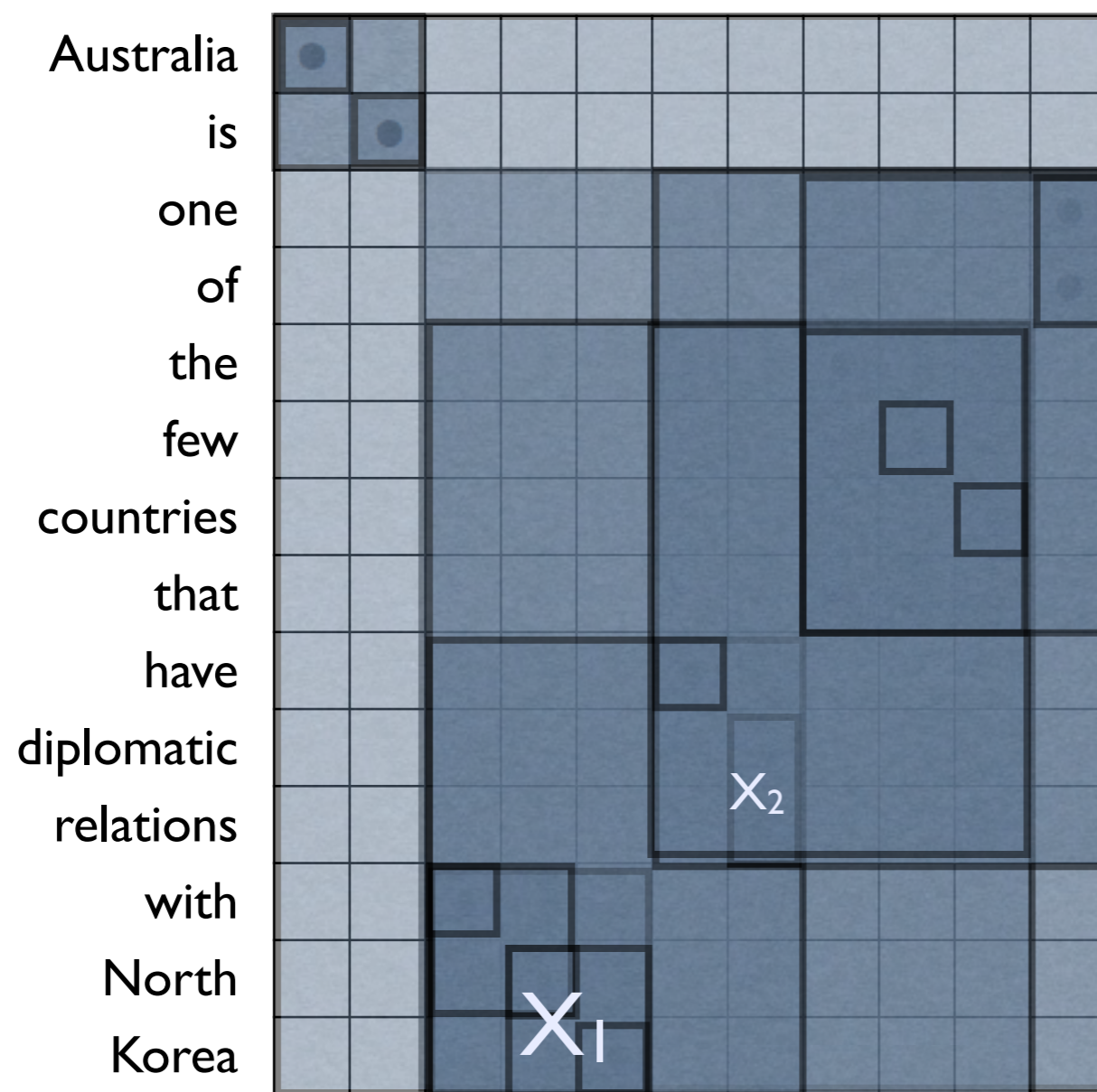
Spanish	English
la bruja verde	the green witch
la bruja roja	the red witch
la bruja azul	the blue witch

- What generalization is missing?
- Hiero abandons conventional English syntax
- Relies instead on evidence-based phrasal “subtractions”



# Extracting Hiero rules

澳 洲 是 与 北 韩 有 邦 交 的 少 数 国 家 之 一



$X \rightarrow$  与 北 韩 有 邦 交,  
have diplomatic relations  
with North Korea

$X \rightarrow$  邦 交,  
diplomatic relations

$X \rightarrow$  北 韩,  
North Korea

$X \rightarrow$  与  $X_1$  有  $X_2$ ,  
have  $X_2$  with  $X_1$

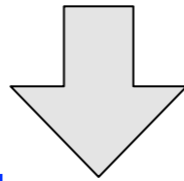
# Decoding

- We now have a way to obtain a synchronous grammar
- Last week, we sketched the decoding algorithm, which was based on parsing
- Today, we'll cover it in more detail, and correct a crucial omission (ngram language models)

# Review (I)

- We've discussed how syntactic differences between languages motivated reordering as a preprocessing step

Ich **werde** Ihnen den Report  
**aushaendigen**, damit Sie den  
eventuell **uebernehmen koennen**.



Ich **werde aushaendigen** Ihnen  
den Report, damit Sie **koennen**  
**uebernehmen** den eventuell.

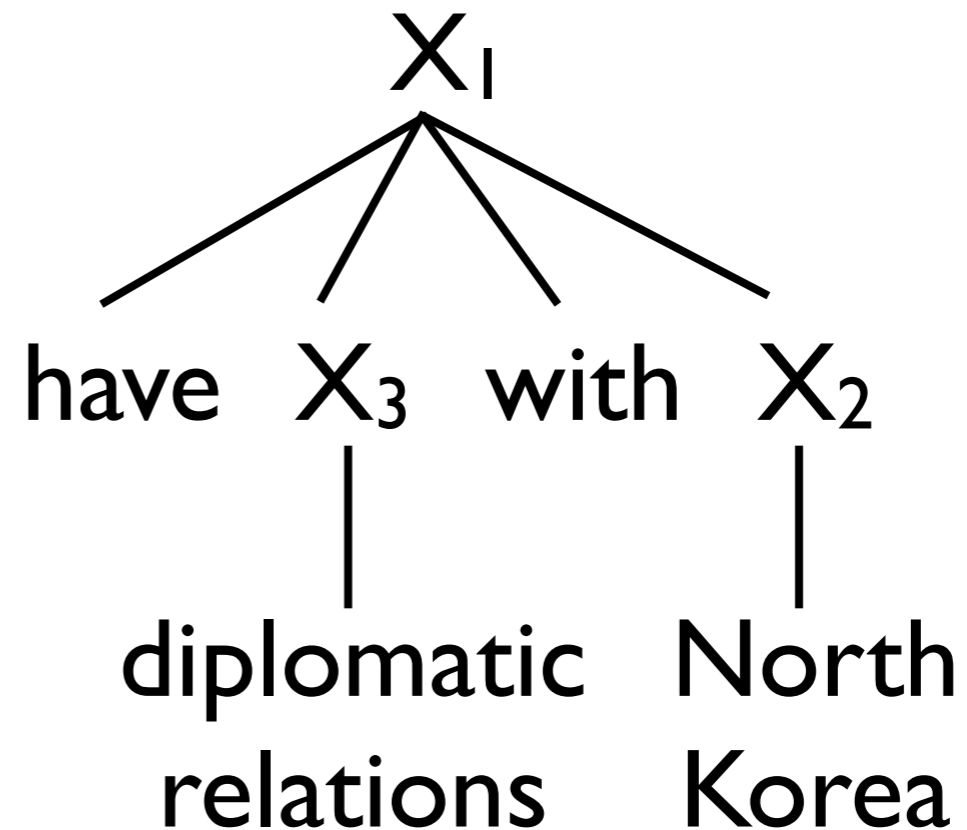
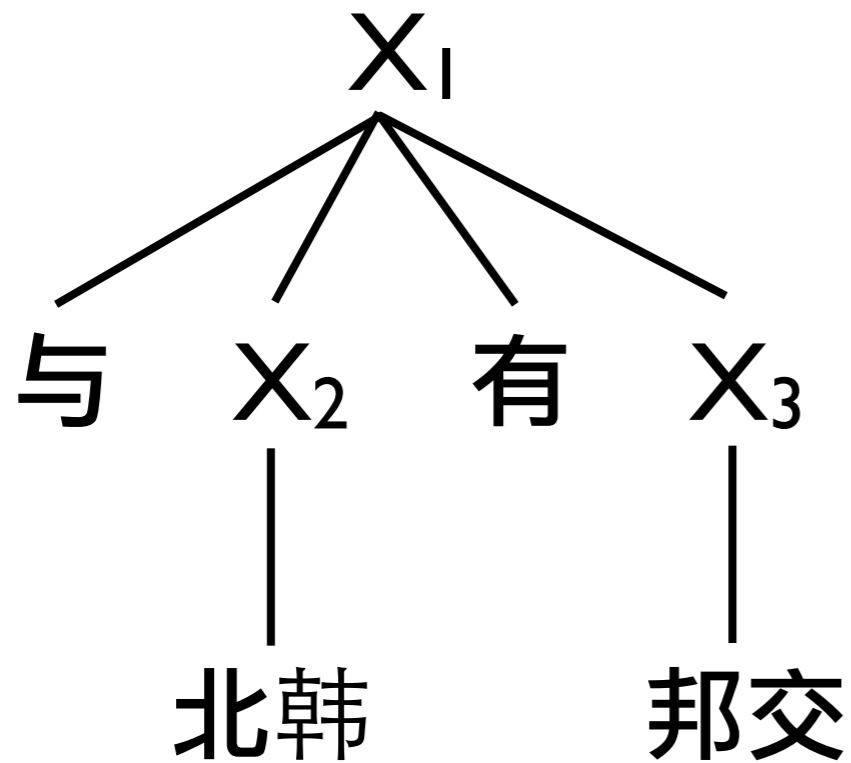
# Review (2)

- We've also discussed **synchronous grammar** rules, which describe the generation of sentences in pairs

	Urdu	English
S →	NP① VP②	NP① VP②
VP →	PP① VP②	VP② PP①
VP →	V① AUX②	AUX② V①
PP →	NP① P②	P② NP①
NP →	<i>hamd ansary</i>	<i>Hamid Ansari</i>
NP →	<i>na}b sdr</i>	<i>Vice President</i>
V →	<i>namzd</i>	<i>nominated</i>
P →	<i>kylye</i>	<i>for</i>
AUX →	<i>taa</i>	<i>was</i>

# Review (3)

- ...and how we could extract those rules automatically from text



# Today

- How do we actually *decode* with these grammars?

- The solution is the CKY / CYK algorithm

{CKY algorithm}	{CYK algorithm}
6,090	~ 13,700

- Outline

- Parsing in one language

- Parsing in two languages with *inversion transduction grammar (ITG)*

- Decoding as parsing with *synchronous context-free grammars (SCFG) and integrated language models*

- Time-permitting: advanced topics



# Review: monolingual parsing

Using the CKY algorithm to find (the best) structure for a sentence given a grammar

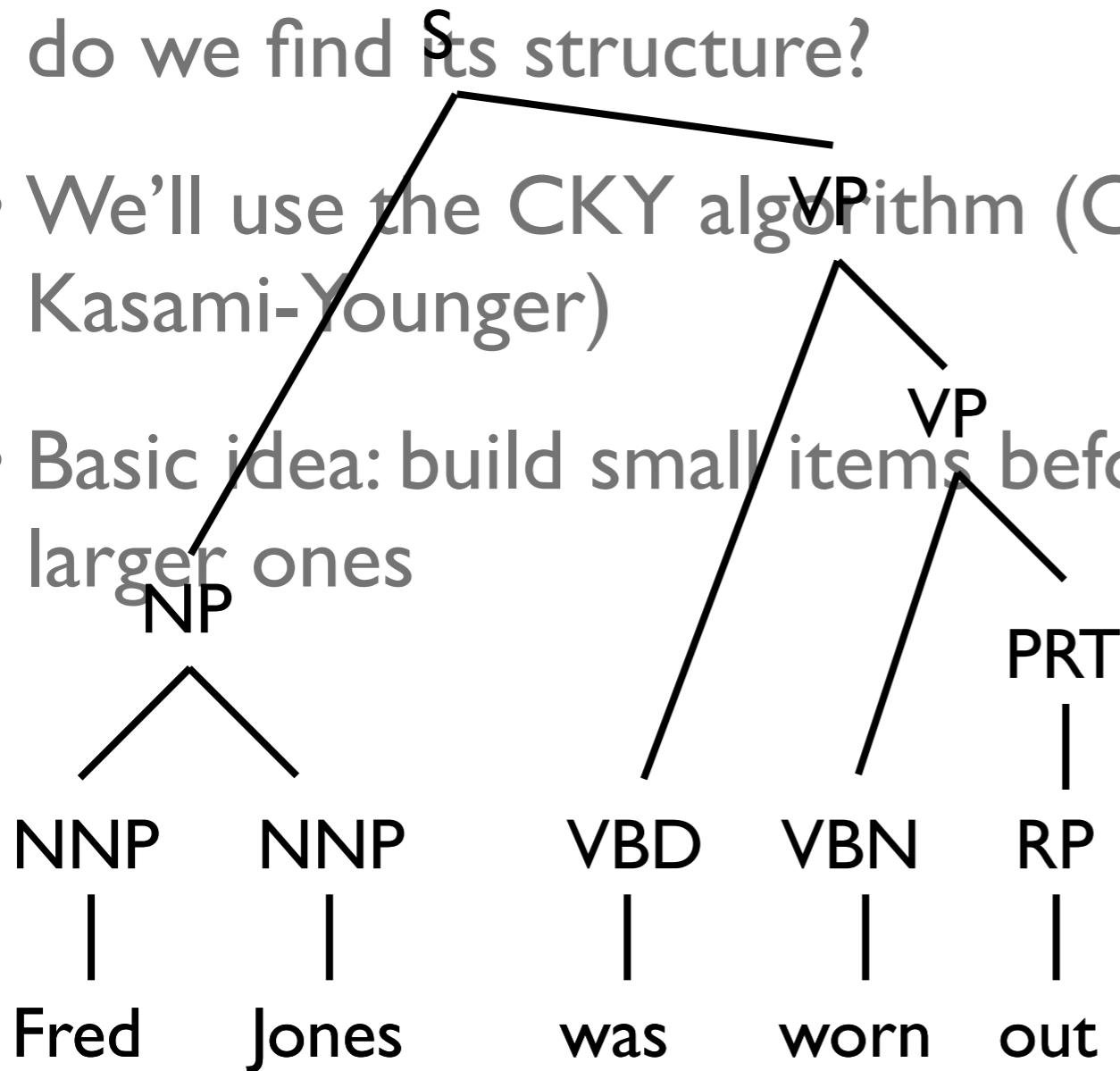
# Formal definitions

- **Formal languages** are (possibly infinite) sets of strings that are generated by a grammar
  - e.g.,  $\{a^+\}$  is a language of all strings with one or more *as*
  - Its grammar could be written as
$$A \rightarrow Aa$$
$$A \rightarrow a$$
- We can view **natural languages** in this manner, too
  - e.g., the **English language** is the set of word sequences that constitute valid English sentences
  - We believe there to be a grammar that generates those sentences
  - We don't know what it is, but we have some guesses and approximations



# Parsing

- Given a sentence and a grammar, how do we find its structure?
- We'll use the CKY algorithm (Cocke-Kasami-Younger)
- Basic idea: build small items before larger ones

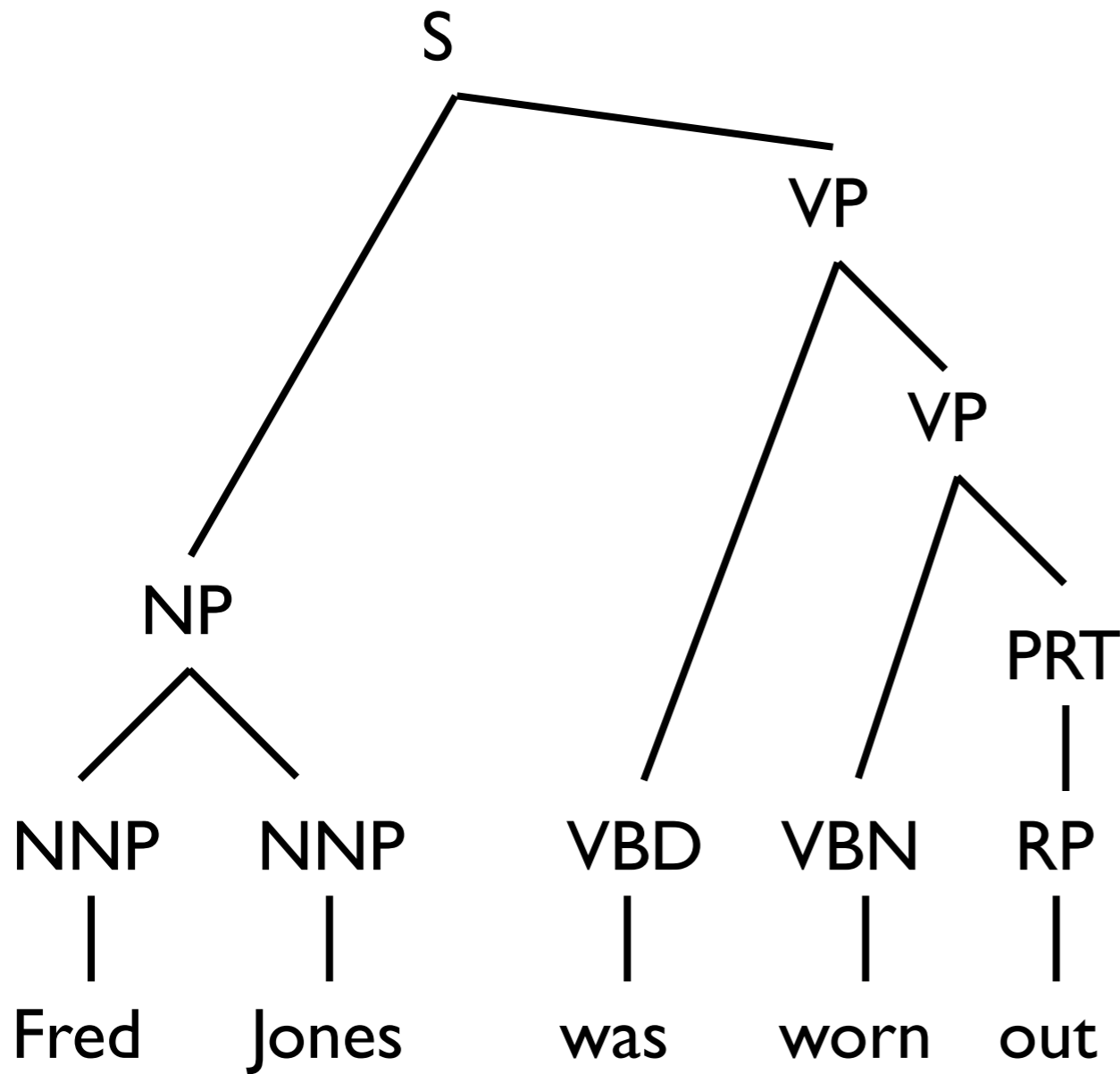


sentence

S	→	NP VP
VP	→	VBN PRT
PRT	→	RP
VP	→	VBD VP
NP	→	NNP NNP
NNP	→	Fred   Jones
VBD	→	was
VBN	→	worn
RP	→	out

grammar

# Parsing with CKY



sentence

S	→	NP VP
VP	→	VBN PRT
PRT	→	RP
VP	→	VBD VP
NP	→	NNP NNP
NNP	→	Fred   Jones
VBD	→	was
VBN	→	worn
RP	→	out

grammar

# Implementation details

- Dynamic programming maintains a **chart** of items

- Each cell item represents the **dynamic programming state**

- (NNP,1,1), (S,1,5)

- The **chart** is the collection of all items

- The **score** resolves alternate ways of constructing an item

- We also store **backpointers**: the items and rule used to construct each item

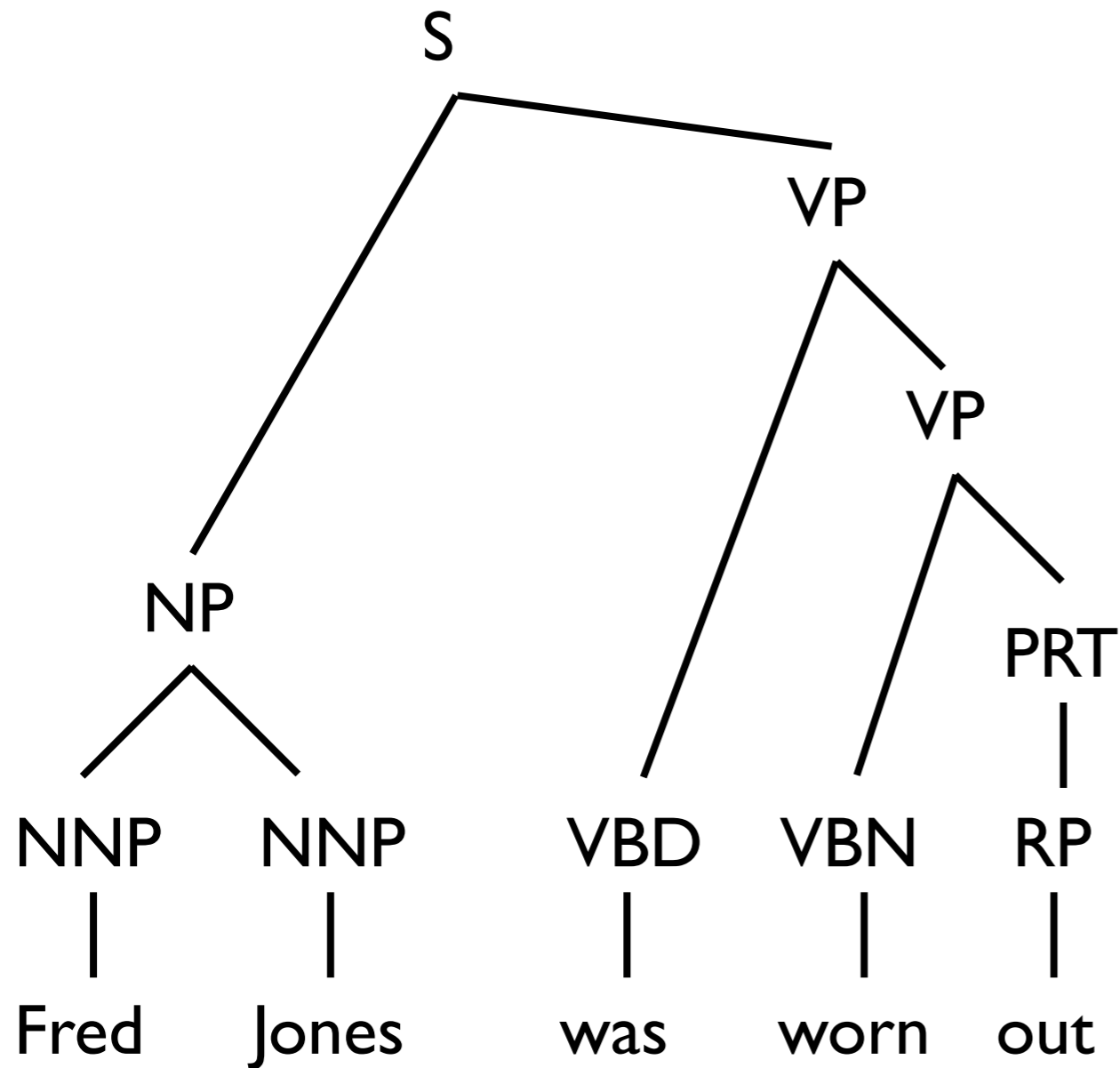
a.k.a. “predecessor”

```
struct item {
    // d.p. state
    string nt;
    int i, j;
    // backpointer
    float score;
    Rule* rule;
    item* rhs1,
        rhs2;
}
```

# CKY algorithm

```
input: words[1..N]
for i in 1..N
  for each unary rule  $X \rightarrow \text{words}[i]$ 
    add  $(X, i, i)$  to the chart
for span in 1..N
  for i in 1..(N-span)
    j = i + span
    for k in i..j
      for rule  $X \rightarrow Y Z$ 
        if  $(Y, i, k)$  and  $(Z, k, j)$ 
          add  $(X, i, j)$  to the chart
output:  $(S, 1, N)$ 
```

# Parsing with CKY



	1	2	3	4	5	
Fred	NNP					1
Jones	NP	NNP				2
was			VBD			3
worn				VBN		4
out	S		VP	VP	RP PRT	5
	Fred	Jones	was	worn	out	

```

item
  nt = "S";
  i = 1, j = 5;
  score = -42.5;
  Rule = &rule("S → NP VP")
  rhs1 = &item(NP, 1, 2);
  rhs2 = &item(VP, 3, 5);
  
```

# Reconstructing the best parse

- We can reconstruct the best parse by following backpointers

```

nodes.append(item(S, 1, N))
while nodes.size() > 0:
    item = nodes.pop()
    print item
    nodes.append(item.rhsr)
    nodes.append(item.rhsl)

```

nodes

~~((S, 1, 5))~~ ~~((NP, 1, 2))~~ ((NNP, 1, 1))

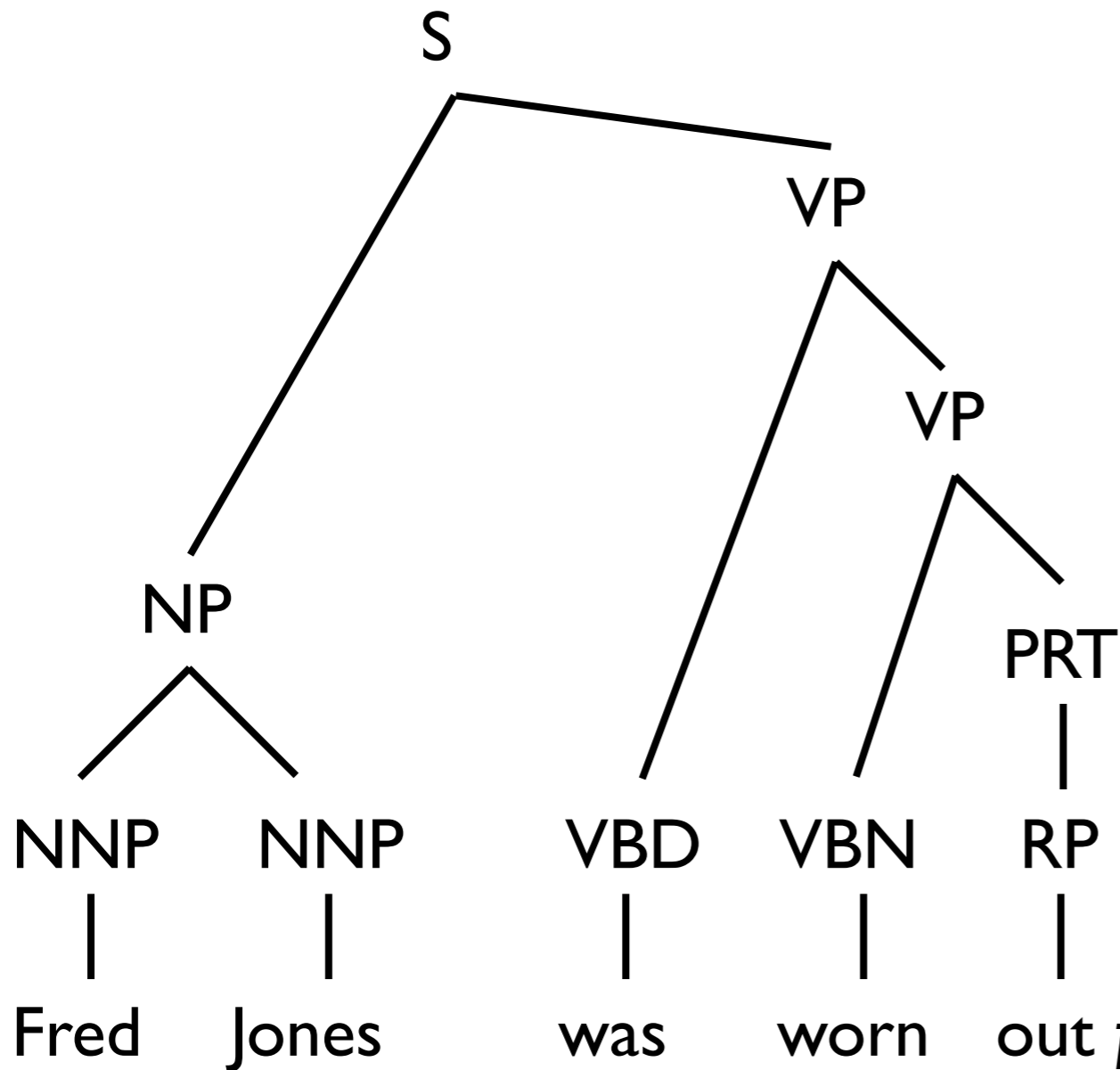
Fred	NNP				
Jones	NP	NNP			
was			VBD		
worn				VBN	
out	S		VP	VP	RP PRT
	Fred	Jones	was	worn	out

```

S → NP VP (1,5)
  NP → NNP NNP (1,2)
    NNP → Fred (1,1)
    NNP → Jones (2,2)
  VP → VBD VP (3,5)
    VBD → was (3,3)
    VP → VBN PRT (4,5)
      VBN → worn (4,4)
      PRT → RP (5,5)
        RP → out (5,5)

```

# Parsing with CKY

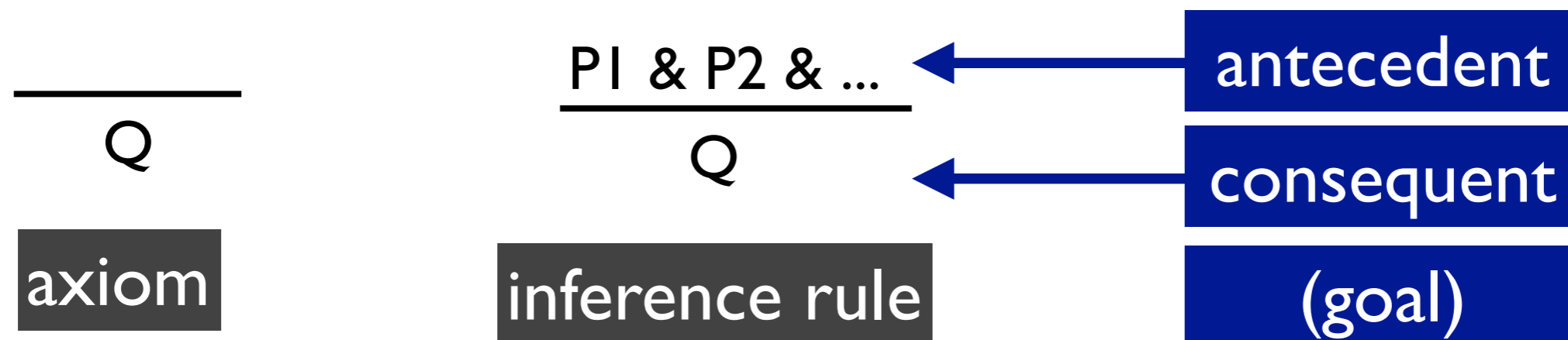


	1	2	3	4	5		
1	Fred	NNP					1
2	Jones	NP	NNP				2
3	was			VBD			3
4	worn				VBN		4
5	out	S		VP	VP	RP PRT	5
		Fred	Jones	was	worn	out	

*Fred Jones was worn out from caring for his often screaming and crying wife during the day but he couldn't sleep at night for she in a stupor from the drugs that didn't ease the pain would set the house ablaze with a cigarette*

# Parsing as (weighted) deduction

- Deductive reasoning:
  - **axioms**: statements that are true or false (“it is raining”)
  - **inference rules**: statements that are conditionally true (“If it is raining and I am outside, I’ll get wet”)
  - **goals**: statements that are licensed by combinations of axioms, inference rules, and other conclusions (“I am wet”)





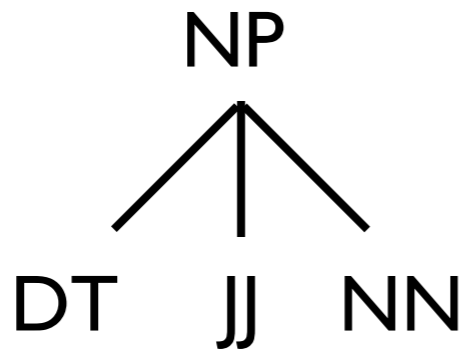
# Parsing as (weighted) deduction

- input: words  $w[1..N]$

Axioms	$\overline{\overline{X \rightarrow w[i]}}$	for all $(X \rightarrow w[i])$
Inference rules	$\frac{X \rightarrow w[i]}{(X, i, i)}$ $\frac{(B, i, j) \quad (C, j, k) \quad A \rightarrow BC}{(A, i, k)}$	in bottom-up order (smaller spans first)
Goal	$(S, 1, n)$	

# Complexity

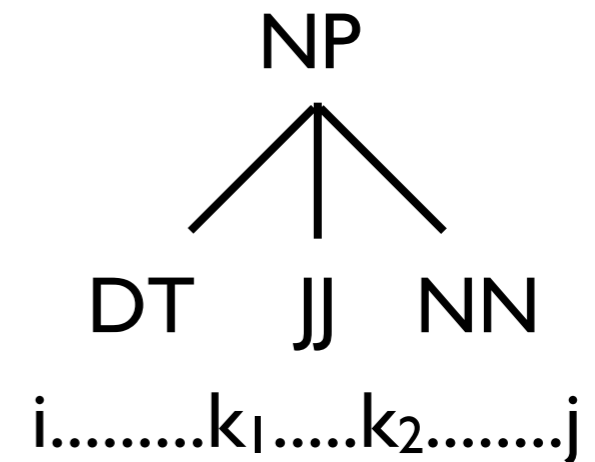
- Complexity of parsing is  $O(Gn^3)$ 
  - $G$  - number of (binarized) rules in the grammar
  - $n$  - length of the sentence
- All those rules were binary; what about longer rules?
  - e.g.,



- We have to enumerate every split point!

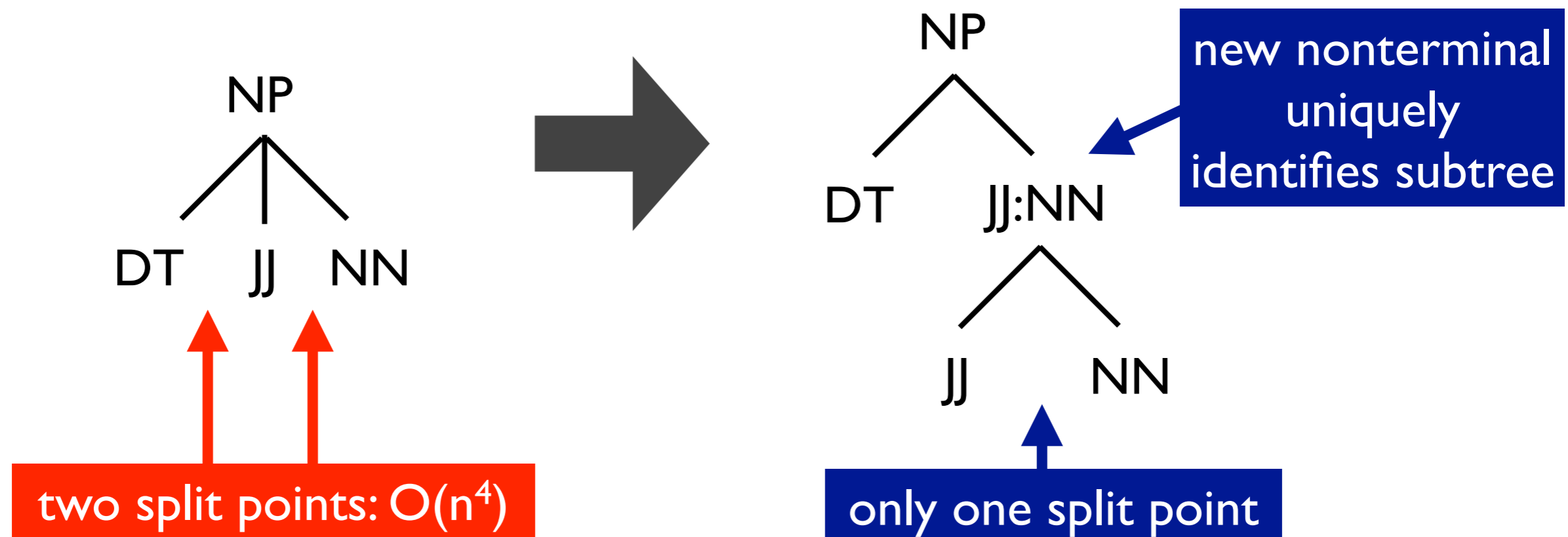
# CKY algorithm

```
input: words[1..N]
for i in 1..N
  for each unary rule X → words[i]
    add (X,i,i) to the chart
for span in 1..N
  for i in 1..(N-span)
    j = i + span
    for k1 in i..j-1
      for k2 in k1..j
        for rule X → W Y Z
          if (W,i,k1) and (Y,k1,k2) and (Z,k2,j)
            add (X,i,j) to the chart
output: (S,1,N)
```



# Binarization into Chomsky Normal Form

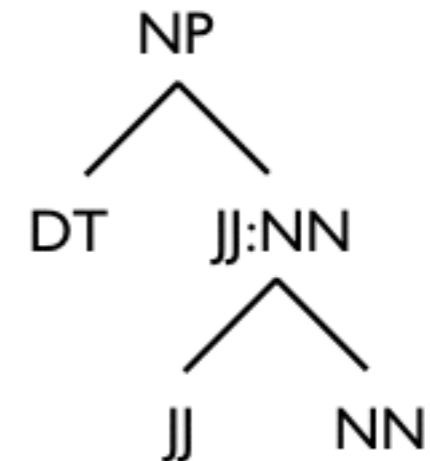
- In general, for a rule with  $k$  RHS items, complexity is  $O(n^{k+1})$  (and cumbersome, since you have to explicitly add inner loops to enumerate them)
- Fortunately, we can **binarize** rules to make them all have a rank of 2



# CKY algorithm

- In summary, monolingual parsing:
  - finds the best structure
  - works bottom-up, enumerating all spans, from small to large, building searching for applicable rules and building new chart items
  - works with the binarized form of a grammars (easily unbinarized afterward) for a complexity of  $O(Gn^3)$
  - all grammars are binarizable

	1	2	3	4	5	
Fred	NNP					1
Jones	NP	NNP				2
was			VBD			3
worn				VBN		4
out	S		VP	VP	RP	5
	Fred	Jones	was	worn	out	



# Synchronous parsing

# Synchronous parsing

- We can extend CKY to parse two languages at once!
- Consider the following grammar:

A  $\rightarrow$  fat, gordos (lexical)  
A  $\rightarrow$  thin, delgados  
N  $\rightarrow$  cats, gatos  
VP  $\rightarrow$  eat, comen  
NP  $\rightarrow$  A<sup>(1)</sup> N<sup>(2)</sup>, N<sup>(2)</sup> A<sup>(1)</sup> (inverted)  
S  $\rightarrow$  NP<sup>(1)</sup> VP<sup>(2)</sup>, NP<sup>(1)</sup> VP<sup>(2)</sup> (straight)

- and the following sentence pair:

fat cats eat / gatos gordos comen

# Synchronous parsing

- We now have to enumerate *pairs* of spans
  - instead of (i,j)...
  - ...we have (i,j) and (s,t)
- For each of the bilingual blocks, we attempt to match both **straight** and **inverted** rules



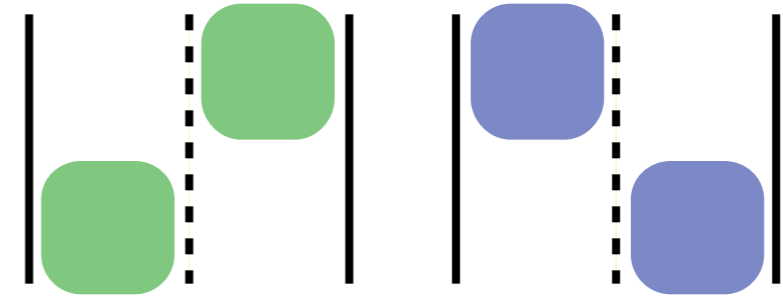
$A \rightarrow \text{fat, gordos}$   
 $N \rightarrow \text{cats, gatos}$   
 $VP \rightarrow \text{eat, comen}$   
 $VP \rightarrow \text{eat, como}$   
 $NP \rightarrow A^{(1)} N^{(2)}, N^{(2)} A^{(1)}$   
 $S \rightarrow NP^{(1)} VP^{(2)}, NP^{(1)} VP^{(2)}$

comen	(3,3,3,3)		(3,3,3,3)
gordos	(1,1,2,2)	(1,2,1,2)	
gatos		(2,2,1,1)	
	fat	cats	eat



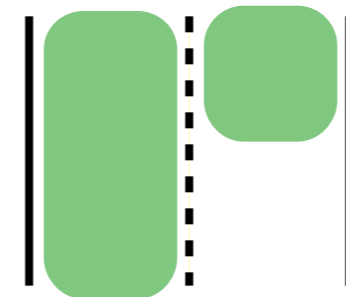
# Relation to monolingual parsing

- Why do we combine like this?
  - Think about monolingual CKY: combine adjacent spans

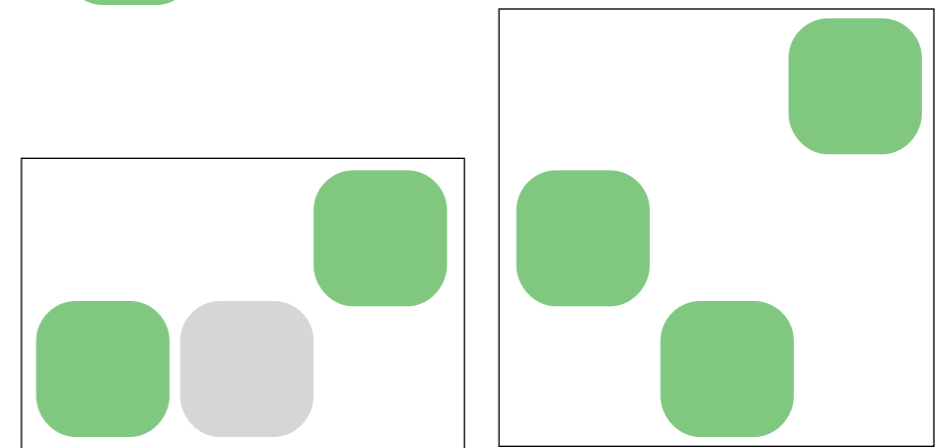


- These pieces are adjacent in both languages; it's only when we consider them *together* that reordering comes into play

- Why can't we do this?
  - It doesn't make sense!



- What about these?
  - Possible, but complex



gap

rank > 2

# CKY for synchronous parsing

```
input: source[1..N], target[1..M]
for span1 in 1..N
  for i in 1..(N-span1)
    j = i + span1
    for k in i..j
      for span2 in 1..M
        for s in 1..(M-span2)
          t = s + span2
          for u in s..t
            for rule X → [Y Z]
              if (Y, i, k, s, u) and
                 (Z, k, j, u, v) then
                add (X, i, j, s, t) to chart
output: (S, 1, N, 1, M)
```

comen				M
gordos				3
gatos				2
	fat	cats	eat	1
				N
				1
				2
				3

# Synchronous parsing

- Complexity:  
 $O(GN^3M^3) \approx O(GN^6)$
- Why?
  - We have to enumerate all valid combinations of six variables
  - This can be seen in the six nested loops of the algorithm

A → fat, gordos

N → cats, gatos

VP → eat, comen

VP → eat, como

NP → A<sup>(1)</sup> N<sup>(2)</sup>, N<sup>(2)</sup> A<sup>(1)</sup>

S → NP<sup>(1)</sup> VP<sup>(2)</sup>, NP<sup>(1)</sup> VP<sup>(2)</sup>

comen	(3,3,3,3)		(3,3,3,3)
gordos	(1,1,2,2)	(1,2,1,2)	
gatos		(2,2,1,1)	
	fat	cats	eat

# Visualization of $O(GN^6)$ complexity

```
input: source[1..N], target[1..M]
```

```
1 for span1 in 1..N
```

```
2   for i in 1..(N-span1)
```

```
     j = i + span1
```

```
3     for k in i..j
```

```
4       for span2 in 1..M
```

```
5         for s in 1..(M-span2)
```

```
           t = s + span2
```

```
6         for u in s..t
```

```
           for rule  $X \rightarrow [Y Z]$ 
```

```
             if (Y, i, k, s, u) and
```

```
                (Z, k, j, u, v) then
```

```
                 add (X, i, j, s, t) to chart
```

```
output: (S, 1, N, 1, M)
```

# Synchronous binarization

- In the above, we considered two nonterminals (per side)
- What if we want more (Zhang et al., 2006)?

$S \rightarrow NP^{(1)} VP^{(2)} PP^{(3)}, NP^{(1)} PP^{(3)} VP^{(2)}$   
 $NP \rightarrow \text{Powell, Baoweier}$   
 $VP \rightarrow \text{held a meeting, } juxing\ le\ huitan$   
 $PP \rightarrow \text{with Sharon, } yu\ Shalong$

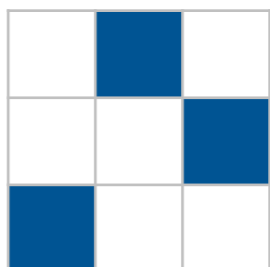
- Three nonterminals? No problem:

$S \rightarrow V_{NP-PP} VP$       or       $S \rightarrow NP V_{PP-VP}$   
 $V_{NP-PP} \rightarrow NP PP$        $V_{PP-VP} \rightarrow PP VP$

- More?

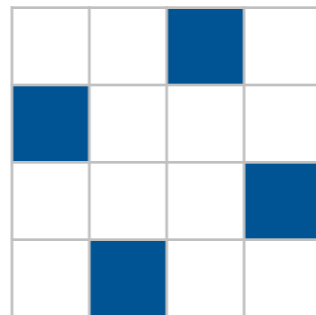
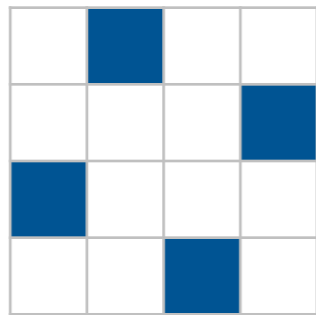
# Permutations

- The nonterminals in the right-hand side of a rule define a permutation between the languages
  - we assume the source language nonterminals are in order (wlog)
  - intermingled terminal symbols do not affect binarization ability
- Example:  $S \rightarrow NP^{(1)} VP^{(2)} PP^{(3)}, NP^{(1)} PP^{(3)} VP^{(2)}$ 
  - permutation: 1 3 2



# Synchronous binarization

- Bad news: synchronous grammars can't be binarized in the general case (Shapiro & Stephens, 1991; Wu, 1997) \*
- Famous examples: the (2,4,1,3) and (3,1,4,2) permutations

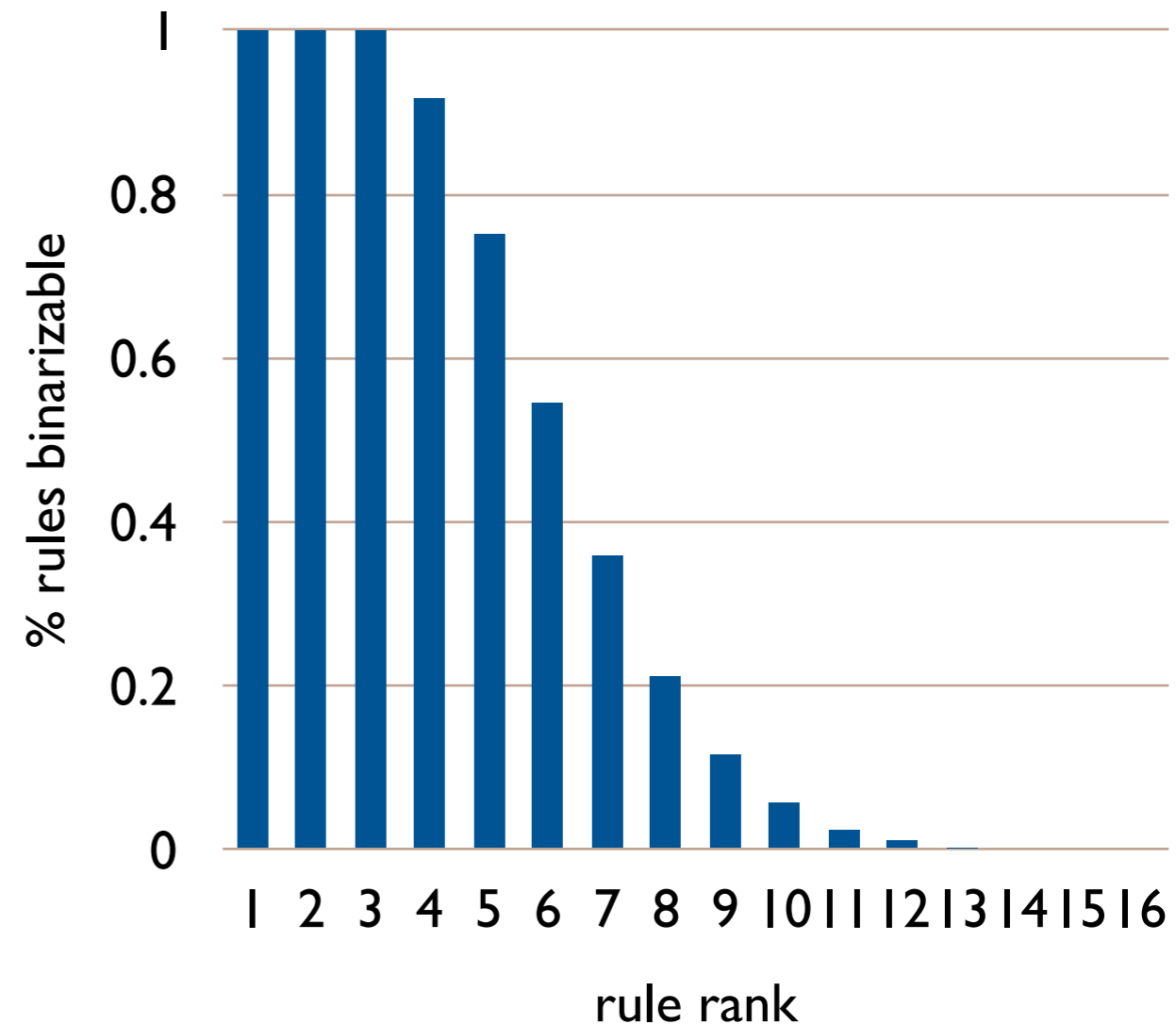


- What makes these unbinarizable?
  - Crucial: parsing works by combining *adjacent* elements
  - No pair of alignments here is adjacent in *both* languages simultaneously

(\* ) *Technically, you can binarize any synchronous grammar, but you may increase the fan-out, which mitigates the potential gains.*

# Synchronous binarization

- As the rank of a rule grows, the percentage of binarizable rules approaches 0



- In summary:
  - We can't binarize all rules
  - The first unbinarizable rule has rank 4



# Silver lining

- Empirically, we don't observe that many non-binarizable rules (Zhang et al., 2006):

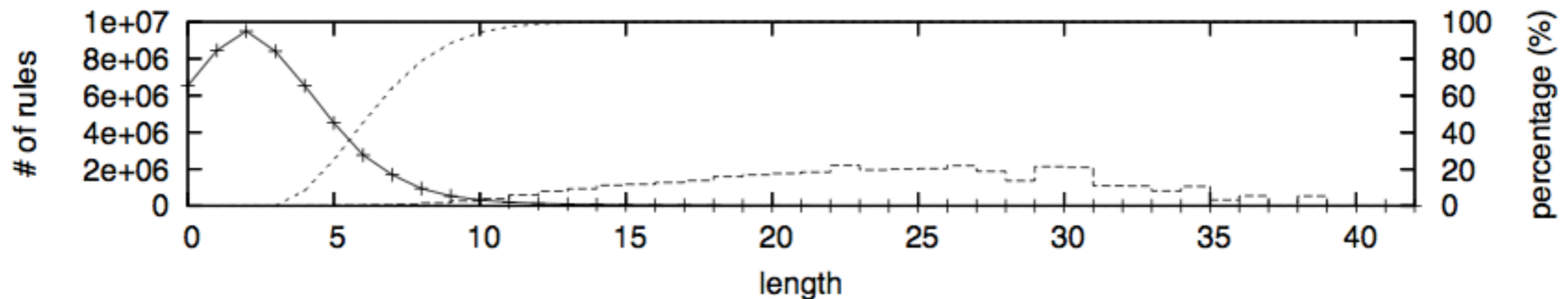


Figure 6: The solid-line curve represents the distribution of all rules against permutation lengths. The dashed-line stairs indicate the percentage of non-binarizable rules in our initial rule set while the dotted-line denotes that percentage among all permutations.

- ...and we can safely throw out the ones we do find
  - 99.7% of rules extracted were binarizable
  - many not were due to alignment errors

# Decoding as parsing

# Synchronous decoding

- Enough parsing; what we care about is decoding
- Parsing is relevant, though, because we can view decoding as a task where we are doing synchronous parsing but we don't happen to know the target side text
- This works by parsing with a **source-side projection** of the synchronous grammar rules
  - At the end, we can follow backpointers to discover the most probable target side

# Updated data structure

- Just like regular parsing, we combine items in pairs to produce new items over larger spans:

$$\frac{(A,1,1) \quad (N,2,2)}{(NP,1,2)}$$

- However, we also have to maintain our guess of the target side

A → fat, gordos

N → cats, gatos

VP → eat, comen

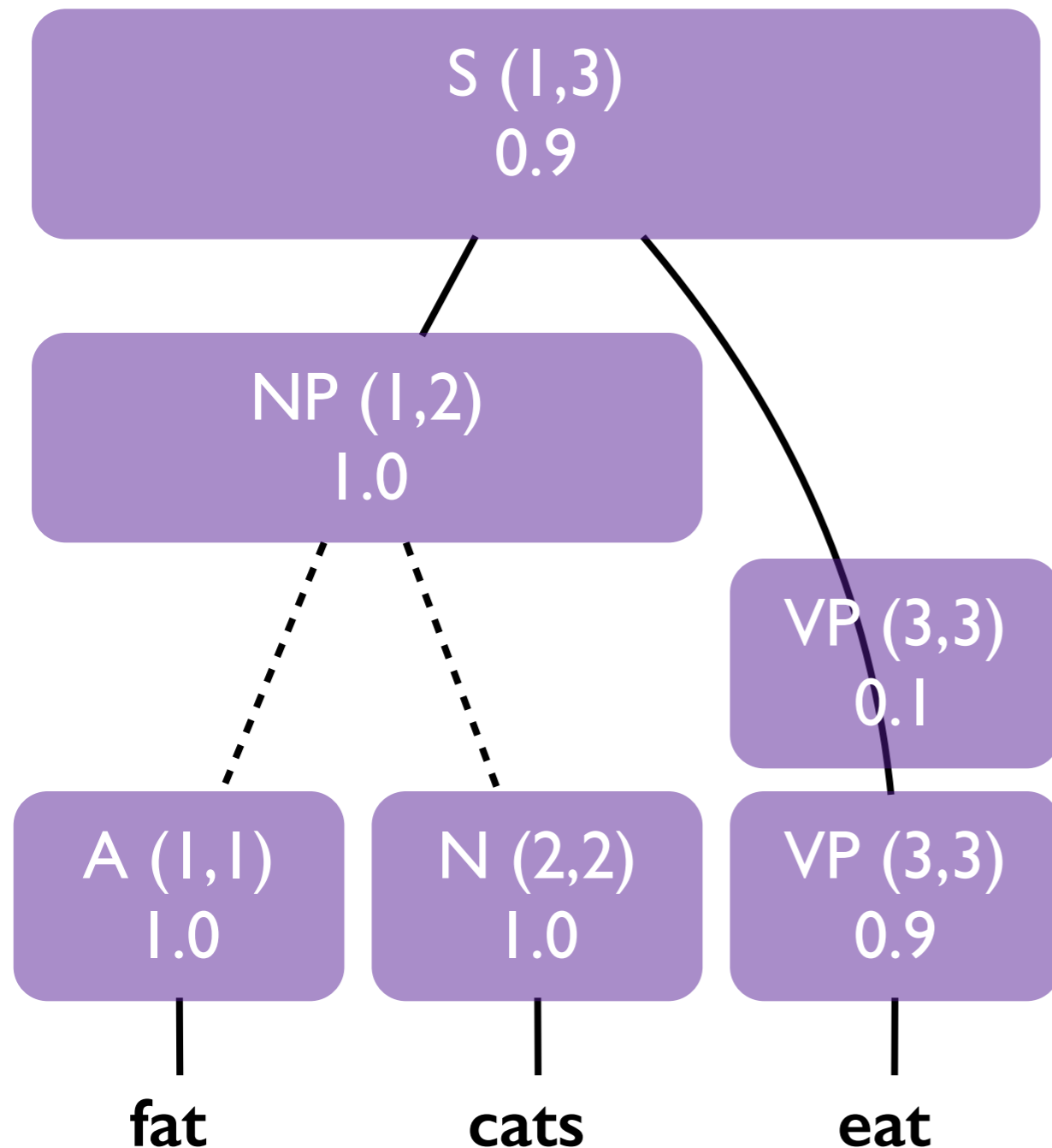
VP → eat, como

NP → A<sup>(1)</sup> N<sup>(2)</sup>, N<sup>(2)</sup> A<sup>(1)</sup>

S → NP<sup>(1)</sup> VP<sup>(2)</sup>, NP<sup>(1)</sup> VP<sup>(2)</sup>

# Decoding

- Again, a bottom-up process



A → fat, gordos

1.0

N → cats, gatos

1.0

VP → eat, comen

0.1

VP → eat, como

0.9

NP → A<sup>(1)</sup> N<sup>(2)</sup>, N<sup>(2)</sup> A<sup>(1)</sup>

1.0

S → NP<sup>(1)</sup> VP<sup>(2)</sup>, NP<sup>(1)</sup> VP<sup>(2)</sup>

1.0

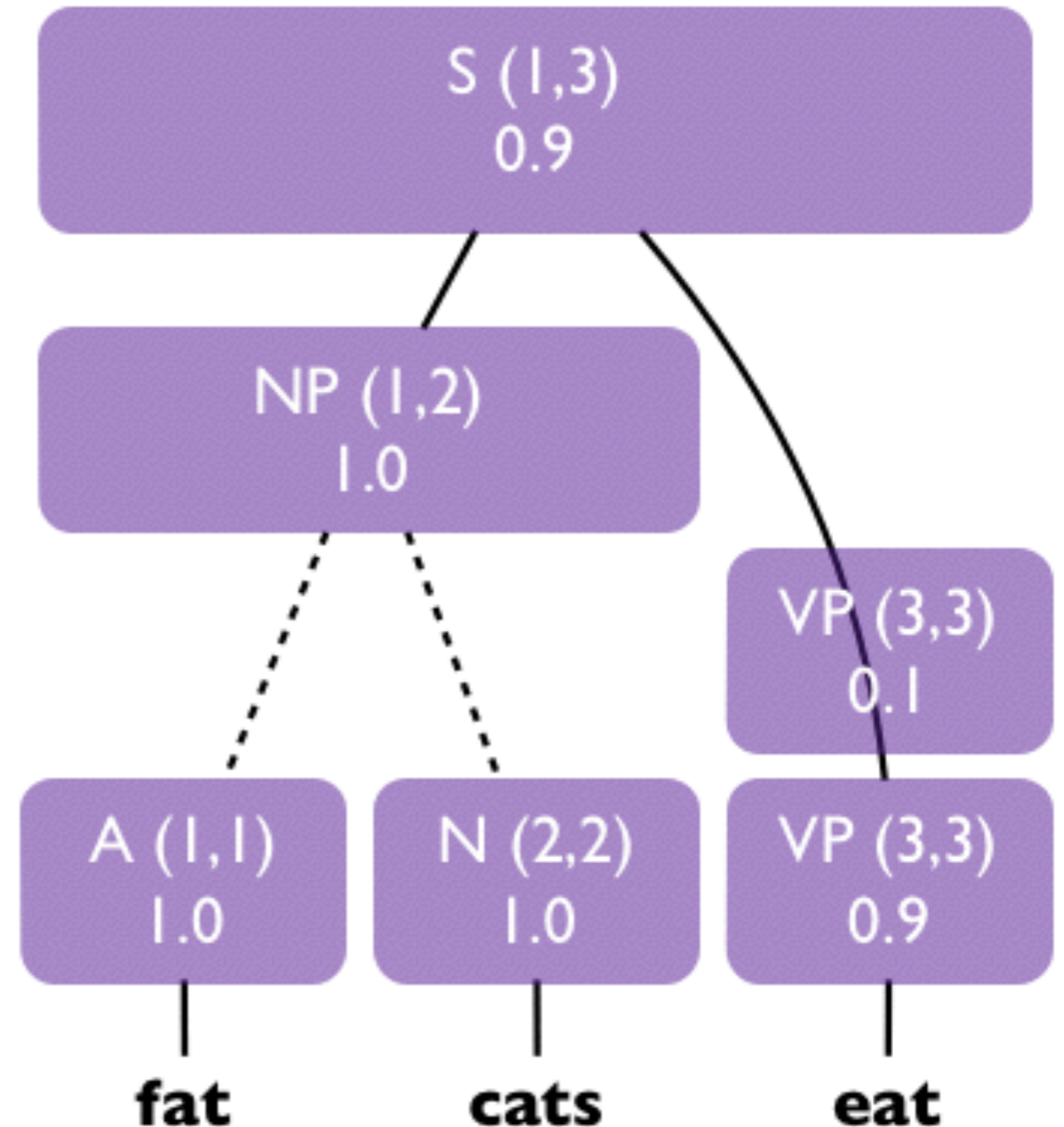
## Legend

— straight rule application

- - - inverted rule application

# Getting the translation

- Follow the backpointers
  - (S,1,3)
    - (NP,1,2)
      - (N,2,2) → gatos
      - (A,1,1) → gordos
    - (VP,3,3) → como
- translation:  
gatos gordos como  
\* *cats fat lps-eat*



# What happened?

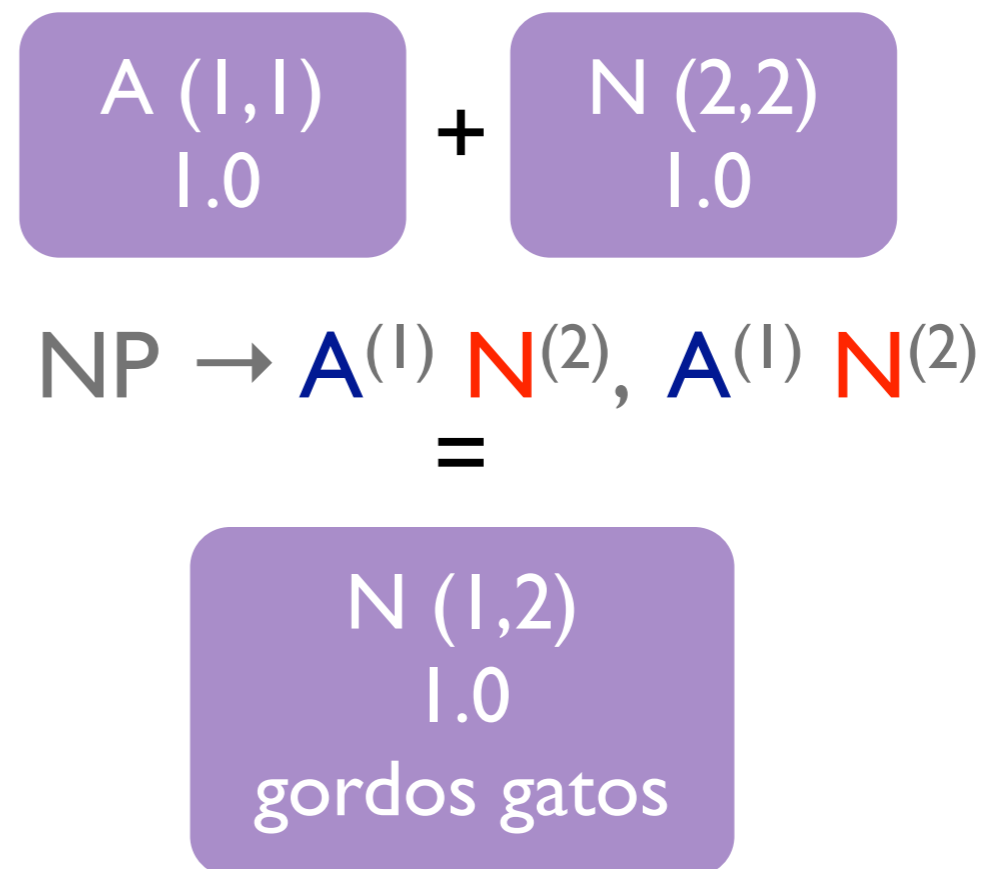
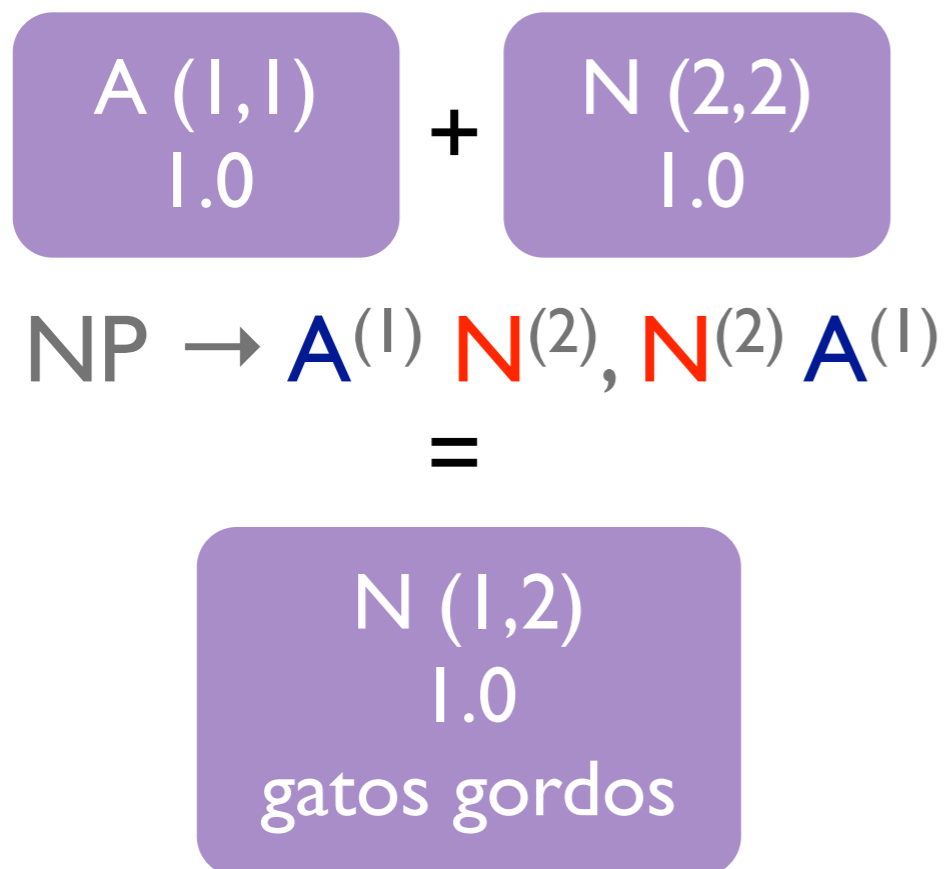
- We forgot the language model
- We're inventing the target side (which is what decoding does), so we need to incorporate it
- How?
  - Stack-based decoding: we maintained the last word
  - Integration was easy because hypotheses always extended to the right
  - Here, hypotheses are *merged* either straight or inverted

# Language model integration

## phrase-based



## synchronous grammars





# Language model integration

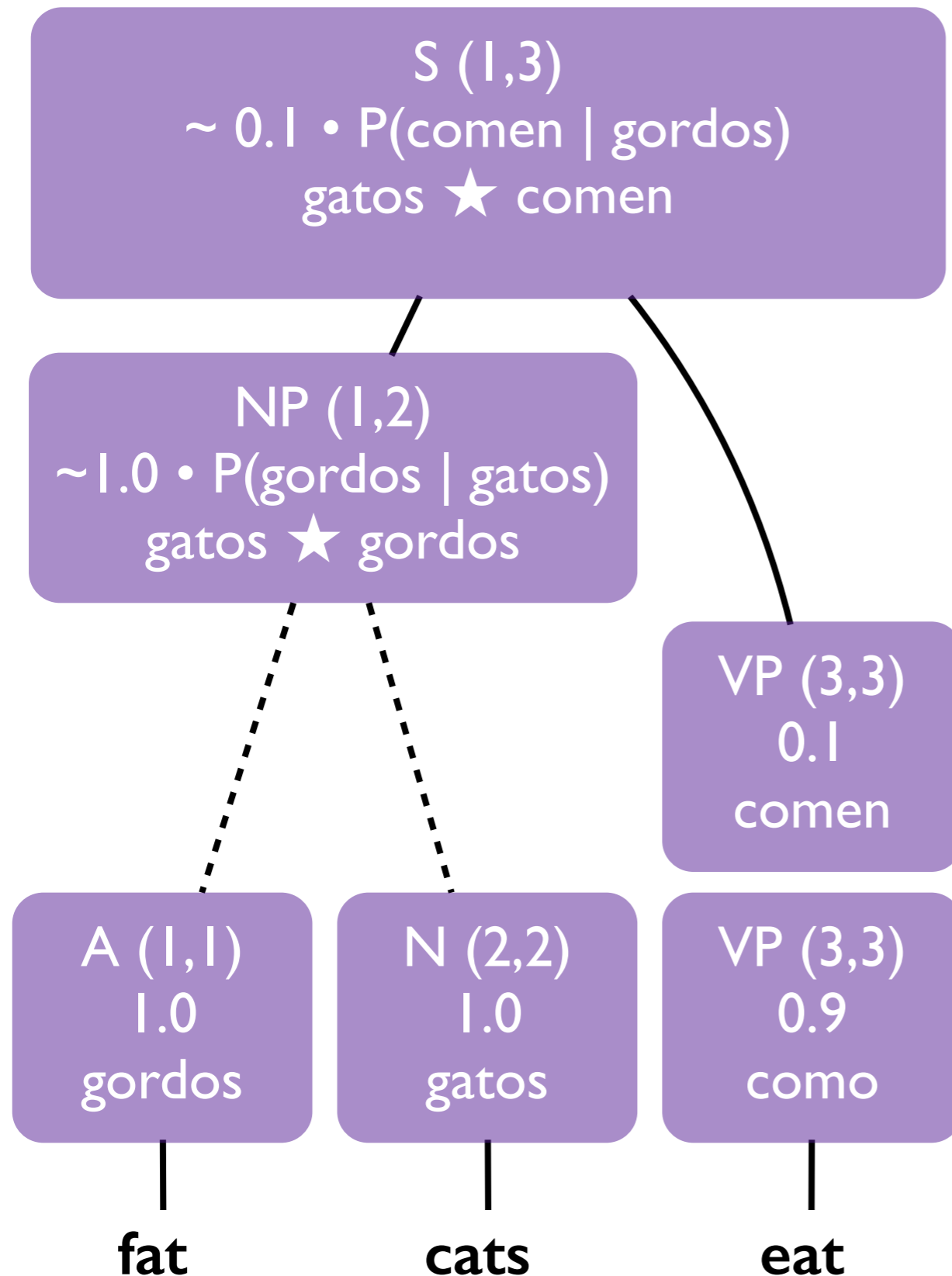
- We still maintain a chart of items, but now the items have to contain the target side words
- Just like regular parsing, we combine items in pairs to produce new items over larger spans
- When items are merged, we can use these words to compute a language model probability
- Formally, we are intersecting a *context-free grammar* (the translation model) with a *regular grammar* (Bar-Hillel et al., 1964; Wu, 1996)

# Updated data structure

- With dynamic programming, we only need a word on either side
  - (for bigram LMs; for the general case, see Chiang (2007, §5.3.2))
  - Following Chiang, we represent the elided middle portion with a ★
  - The complete string can be reconstructed by following the backpointers

```
struct item {  
    // d.p. state  
    string nt;  
    int i, j;  
    string left_words;  
    string right_words;  
    // backpointer  
    float score;  
    Rule* rule;  
    item* rhs1,  
        rhs2;  
}
```

# Decoding with an integrated LM



A → fat, gordos

N → cats, gatos

VP → eat, comen

VP → eat, como

NP → A<sup>(1)</sup> N<sup>(2)</sup>, N<sup>(2)</sup> A<sup>(1)</sup>

S → NP<sup>(1)</sup> VP<sup>(2)</sup>, NP<sup>(1)</sup> VP<sup>(2)</sup>

1.0

1.0

0.1

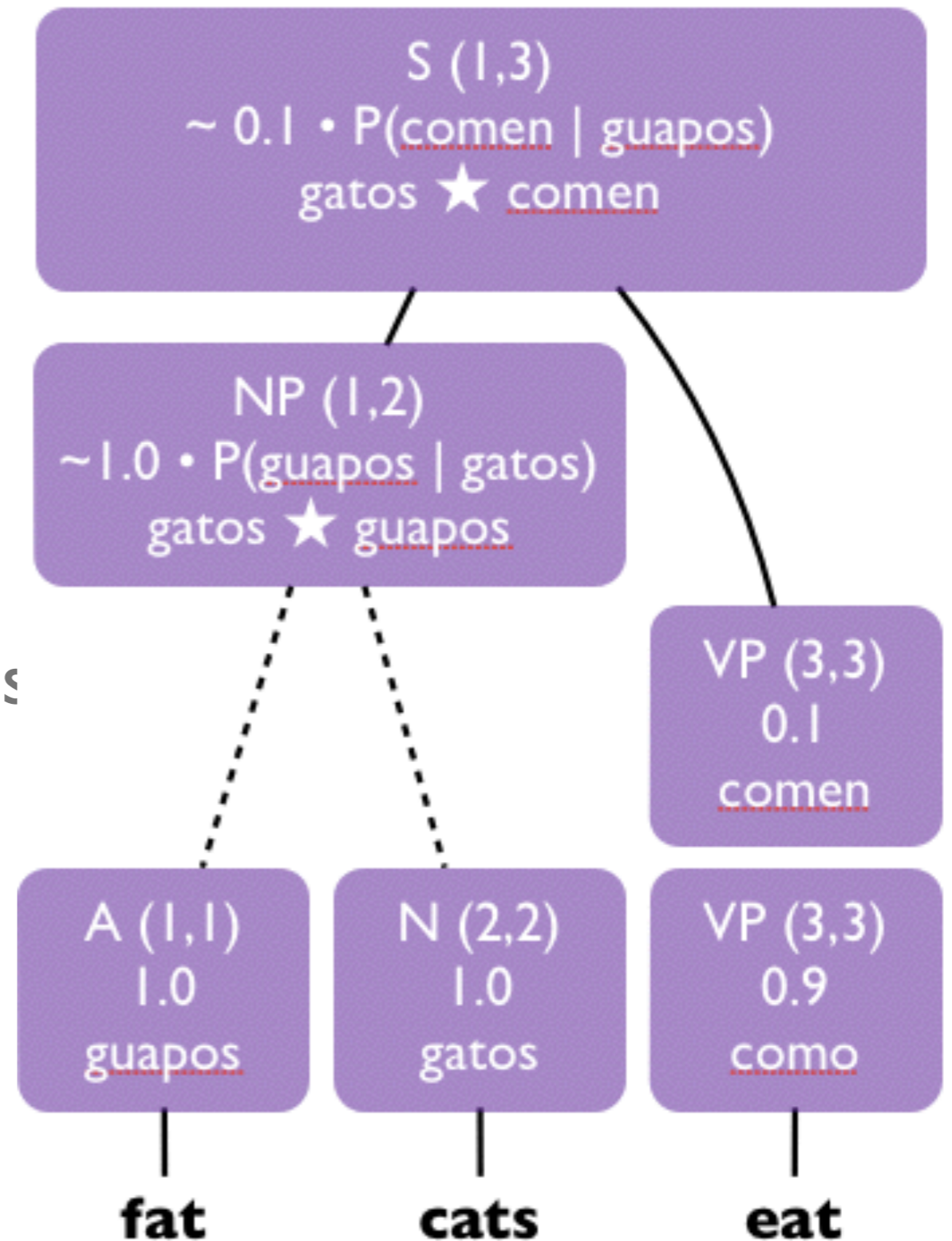
0.9

1.0

1.0

# Getting the translation

- Follow the backpointers
  - (S, 1, 3, gatos★comen)
    - (NP, 1, 2, gatos★gordos)
      - (N, 2, 2, gatos) → gatos
      - (A, 1, 1, gordos) → gordos
    - (VP, 3, 3, comen) → comen
- translation:  
gatos gordos comen  
*cats fat 3pp-eat*



# Pruning

- We have also not dealt much with ambiguity and competition amongst hypotheses
- In general, there are too many hypotheses to consider, so we keep only the top  $k$  of them (per input span  $(i,j)$ )
- When considering a span  $(i,j)$  and a split point  $k$ , we have a large number of ways to combine items
  - there can be any number of applicable rules
  - there can be up to  $k$  items located at span  $(i,k)$
  - there can be up to  $k$  items located at span  $(k,j)$

# Applying a unary rule

- The naive way is to consider the full cross product

		[X, 6, 8; the scheme]	[X, 6, 8; the plan]	[X, 6, 8; the project]	
		1	4	7	
$X \rightarrow \langle \text{cong } X_{\square}, \text{from } X_{\square} \rangle$	1	2.1	5.1	8.2	[X, 5, 8; from the $\star$ the scheme] : 2.1
$X \rightarrow \langle \text{cong } X_{\square}, \text{from the } X_{\square} \rangle$	2	5.5	8.5	11.5	[X, 5, 8; from the $\star$ the plan] : 5.1
$X \rightarrow \langle \text{cong } X_{\square}, \text{since } X_{\square} \rangle$	6	7.7	10.6	13.1	[X, 5, 8; from the $\star$ the scheme] : 5.5
$X \rightarrow \langle \text{cong } X_{\square}, \text{through } X_{\square} \rangle$	10	11.1	14.3	17.3	[X, 5, 8; since the $\star$ the scheme] : 7.7

# Cube pruning

- When considering a span  $(i,j)$  of a length- $N$  sentence:
  - *unary rules*: there are  $rk$  items to compute ( $r$  the number of rules,  $k$  the number of child items)
  - *binary rules*: there are  $Nrk^2$  items to compute (since there are  $O(N)$  split points)
- However, we're only going to be keeping the top  $k$  of them!
  - this problem gets worse as  $k$  gets larger
- We'd like to avoid computing all of these new items, which we accomplish with **cube pruning**

# Cube pruning

- We start with sorted lists of rules and the items they applied to
- Observation:
  - the best item comes from the **best rule** and the **best cell**
  - the next-best item uses either the **2nd best rule** or the **2nd-best cell**

	rule	rhsl	rhsr
1	1	1	1
2	2	7	3
3	4	9	4
...			

best item

	rule	rhsl	rhsr
1	1	1	1
2	2	7	3
3	4	9	4
...			

2nd-best

	rule	rhsl	rhsr
1	1	1	1
2	2	7	3
3	4	9	4
...			

3rd-best



# Applying a unary rule

- The Huang & Chiang (2005) way:

		[X, 6, 8; the scheme]	[X, 6, 8; the plan]	[X, 6, 8; the project]
		1	4	7
$X \rightarrow \langle \text{cong } X_{\square}, \text{from } X_{\square} \rangle$	1	2.1	5.1	
$X \rightarrow \langle \text{cong } X_{\square}, \text{from the } X_{\square} \rangle$	2	5.5		
$X \rightarrow \langle \text{cong } X_{\square}, \text{since } X_{\square} \rangle$	6			
$X \rightarrow \langle \text{cong } X_{\square}, \text{through } X_{\square} \rangle$	10			

		[X, 6, 8; the scheme]	[X, 6, 8; the plan]	[X, 6, 8; the project]
		1	4	7
		2.1	5.1	8.2
		5.5	8.5	

		[X, 6, 8; the scheme]	[X, 6, 8; the plan]	[X, 6, 8; the project]
		1	4	7
		2.1	5.1	8.2
		5.5	8.5	
		7.7		

# Cube pruning

- We haven't discussed the language model, which complicates this procedure by making it *nonmonotonic*
- But that's the basic idea

# Summary

- Today, we have reviewed
  - Monolingual parsing
  - Synchronous (bilingual) parsing
  - Decoding as parsing with an intersected bigram language model
- We have also briefly touched on efficiency considerations with cube pruning

# Advanced topics

# Advanced topics: implicit binarization

- We'd decoded in an ITG settings, where the rules all look like this:

$X \rightarrow \text{boy, chico}$  (lexical)

$X \rightarrow X^{(1)} X^{(2)}, X^{(2)} X^{(1)}$  (inverted)

$X \rightarrow X^{(1)} X^{(2)}, X^{(1)} X^{(2)}$  (straight)

- This is the closest thing to Chomsky Normal Form for synchronous grammars
- How do we decode with intermingled terminals and nonterminals?

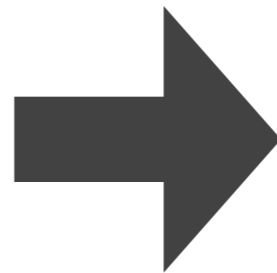
$X \rightarrow \text{the } X^{(1)} \text{ was } X^{(2)}, \text{el } X^{(1)} \text{ era } X^{(2)}$

# Advanced topics: implicit binarization

- One answer: binarize (terminals can always be binarized):

$X \rightarrow$  the  $X$  was  $X$ , el  $X$  era  $X$

$X \rightarrow$  the  $X_{174}$ , el  $X_{174}$



$X_{174} \rightarrow X X_{295}, X X_{295}$

$X_{295} \rightarrow$  was  $X$ , era  $X$

- However, this is inefficient:
  - it leads to a huge blowup in the number of nonterminals
  - it introduces a split point that has to be searched over (avoidable in this case, but not always)

# Advanced topics: implicit binarization

- Instead, we'd like to do implicit, **Earley-style binarization**

# Advanced topics: spurious ambiguity

- *Spurious ambiguity* - multiple structures leading to the same interpretation
- Especially problematic in ITG with its weak grammar
- This can be addressed in various ways
  - Grammar canonical forms