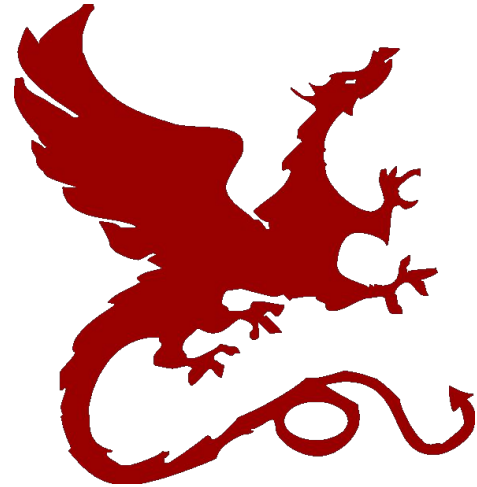


# Algorithms for NLP



## Parsing II

Yulia Tsvetkov – CMU

Slides: Ivan Titov – University of Edinburgh,  
Chris Dyer – Deepmind



# Announcements

---

- HW2 out
- Today: Sachin will give an overview of HW2
- Recitation on EM next week 10/12
- Recitation on HW2 the week after 10/19
- Yulia office hours
  - today: 3:30-4:00
  - next week Yulia is away, no office hours

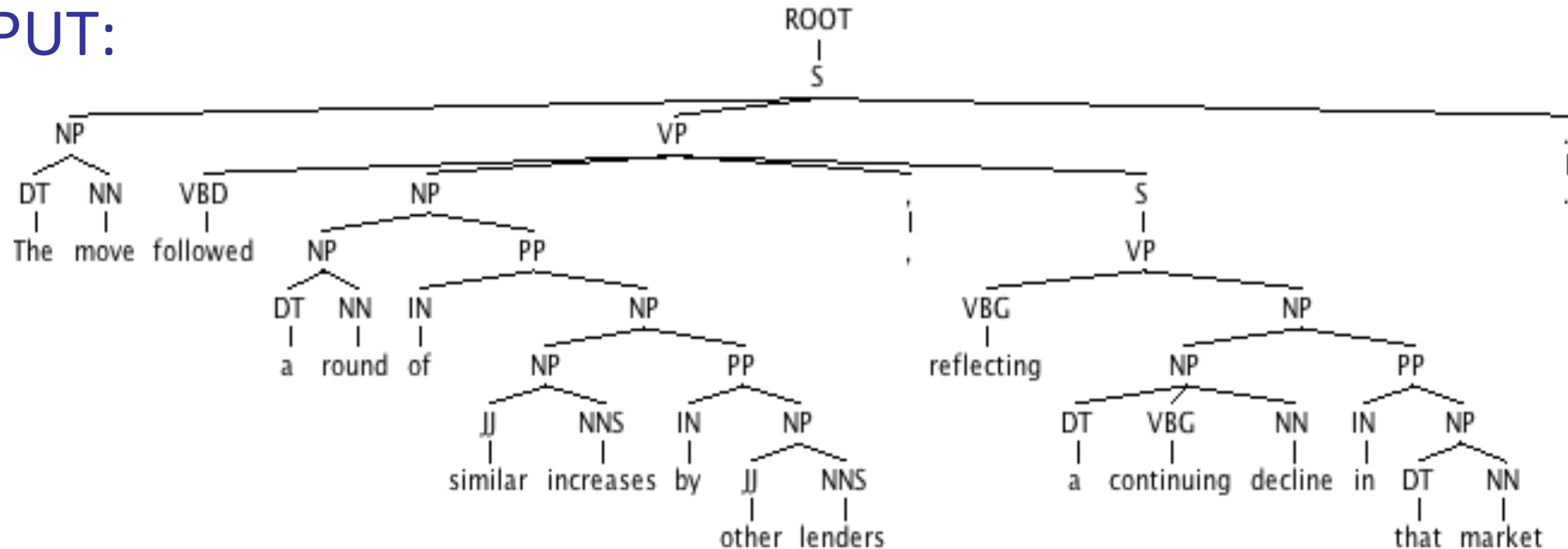


# Syntactic Parsing

- INPUT:

- The move followed a round of similar increases by other lenders, reflecting a continuing decline in that market

- OUTPUT:





# Context Free Grammar (CFG)

---

## Grammar (CFG)

ROOT  $\rightarrow$  S

S  $\rightarrow$  NP VP

NP  $\rightarrow$  DT NN

NP  $\rightarrow$  NN NNS

NP  $\rightarrow$  NP PP

VP  $\rightarrow$  VBP NP

VP  $\rightarrow$  VBP NP PP

PP  $\rightarrow$  IN NP

## Lexicon

NN  $\rightarrow$  interest

NNS  $\rightarrow$  raises

VBP  $\rightarrow$  interest

VBZ  $\rightarrow$  raises

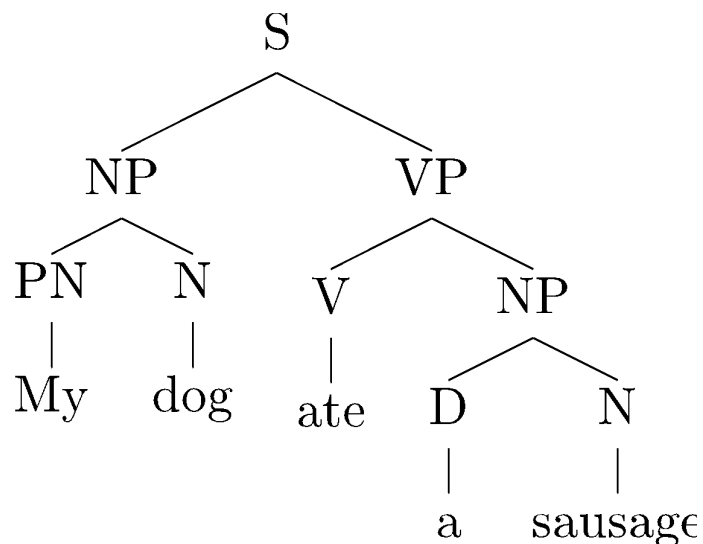
...

- Other grammar formalisms: LFG, HPSG, TAG, CCG...





# Constituent trees



- Internal nodes correspond to phrases
  - S – a sentence
  - NP (Noun Phrase): My dog, a sandwich, lakes, ..
  - VP (Verb Phrase): ate a sausage, barked, ...
  - PP (Prepositional phrases): with a friend, in a car, ...
- Nodes immediately above words are PoS tags (aka preterminals)
  - PN – pronoun
  - D – determiner
  - V – verb
  - N – noun
  - P – preposition

# Parsing with CKY

	lead		can		poison	
0		1		2		3

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

---

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules

lead	can	poison
0	1	2

max = 1      max = 2      max = 3

min = 0			S?
min = 1			
min = 2			

Chart (aka parsing triangle)

$S \rightarrow NP VP$

$VP \rightarrow M V$   
 $VP \rightarrow V$

$NP \rightarrow N$   
 $NP \rightarrow N NP$

---

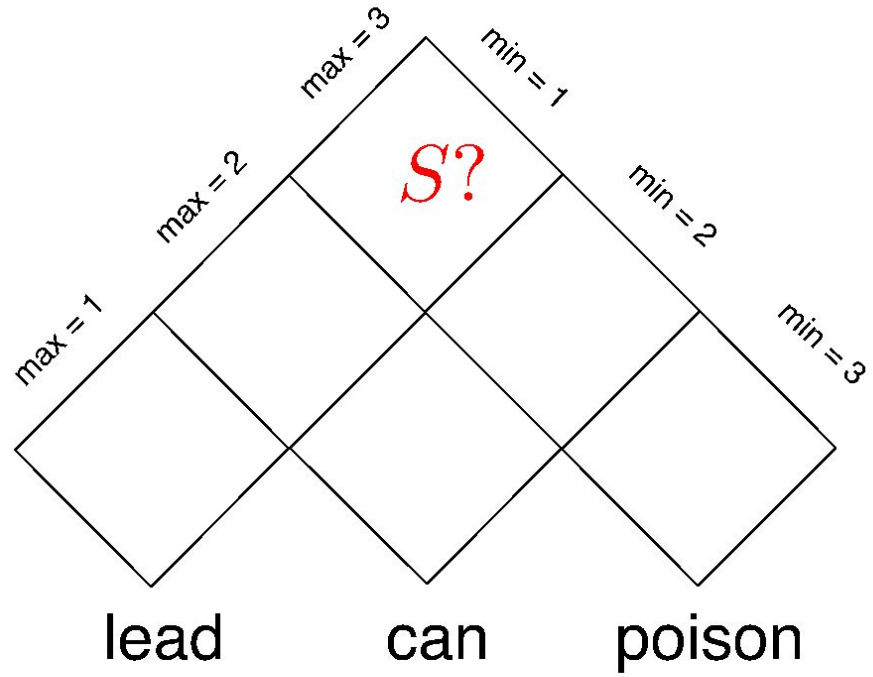
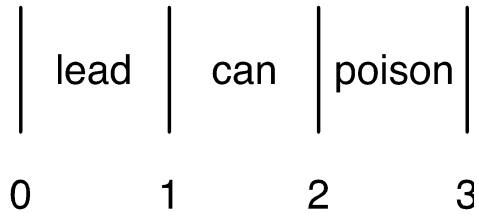
$N \rightarrow can$   
 $N \rightarrow lead$   
 $N \rightarrow poison$

$M \rightarrow can$   
 $M \rightarrow must$

$V \rightarrow poison$   
 $V \rightarrow lead$

Inner rules

Preterminal rules



$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

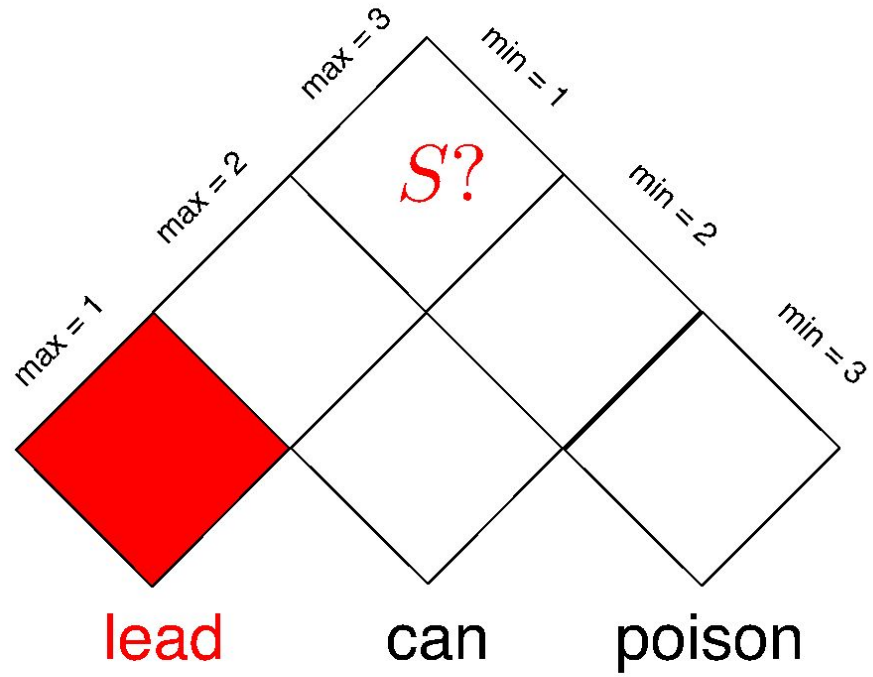
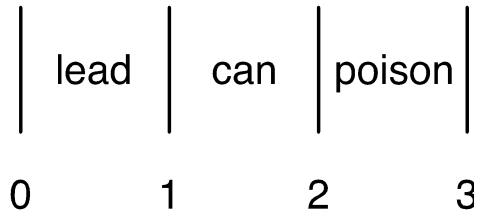
$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules



$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

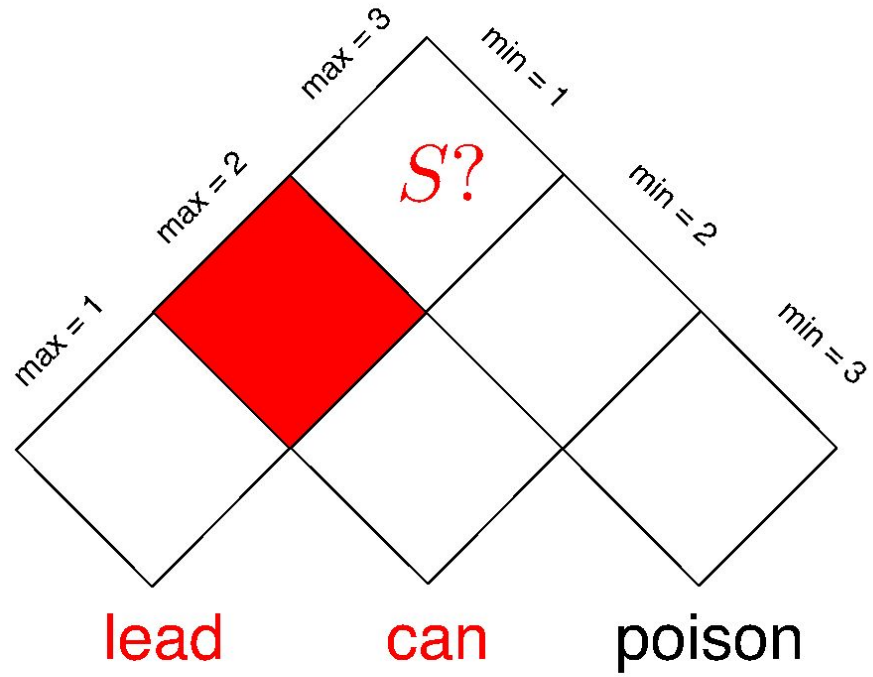
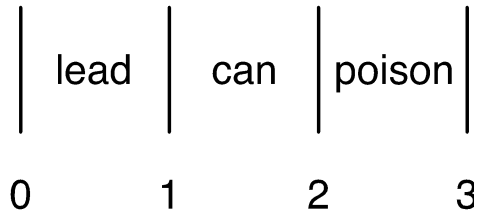
$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules



$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules

lead	can	poison
0	1	2

max = 1      max = 2      max = 3

min = 0			S?
min = 1			
min = 2			

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules

| lead | can | poison |  
 0      1      2      3

max = 1      max = 2      max = 3

min = 0	1	4	6 <i>S?</i>
min = 1		2	5
min = 2			3

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

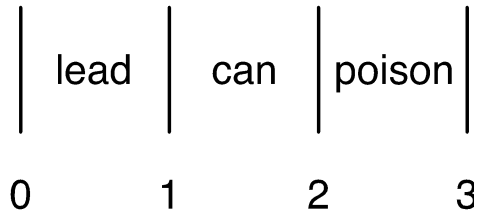
$V \rightarrow poison$

$V \rightarrow lead$

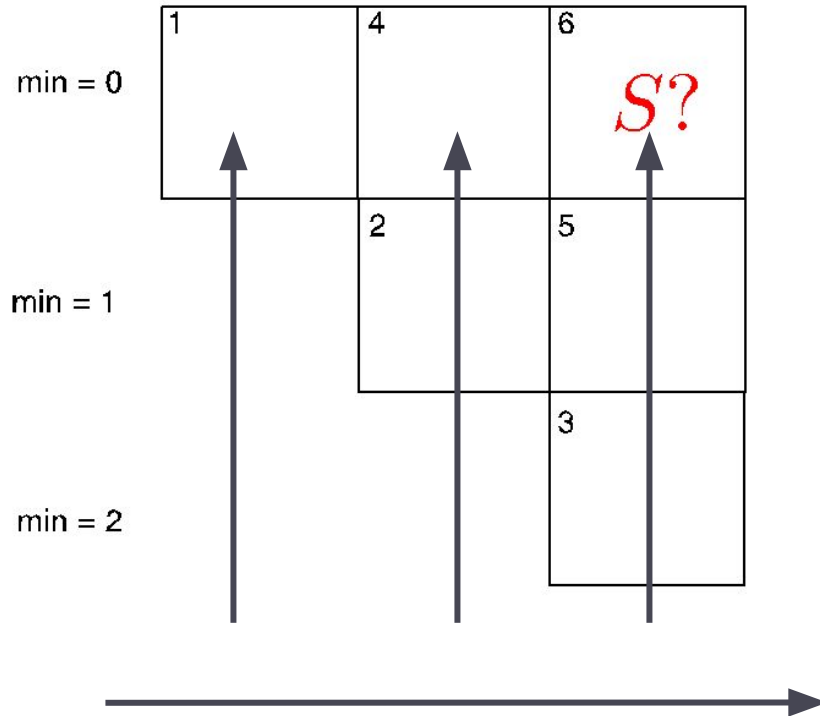
Inner rules

Preterminal rules





max = 1      max = 2      max = 3



$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

Inner rules

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

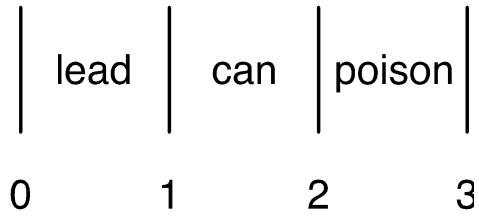
$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Preterminal rules



max = 1      max = 2      max = 3

min = 0	1 ?		
min = 1		2 ?	
min = 2			3 ?

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules

lead	can	poison
0	1	2

max = 1      max = 2      max = 3

min = 0	1 ?		
min = 1		2 ?	
min = 2			3 ?

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

- $N \rightarrow can$
  - $N \rightarrow lead$
  - $N \rightarrow poison$
  
  - $M \rightarrow can$
  - $M \rightarrow must$
  
  - $V \rightarrow poison$
  - $V \rightarrow lead$

Inner rules

Preterminal rules

| lead | can | poison |  
 0      1      2      3

		max = 1	max = 2	max = 3
min = 0	1 <i>N, V</i>			
min = 1		2 <i>N, M</i>		
min = 2			3 <i>N, V</i>	

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules

| lead | can | poison |  
 0      1      2      3

max = 1      max = 2      max = 3

min = 0	1 <i>N, V</i> <i>NP, VP</i>	4 ?	
min = 1		2 <i>N, M</i> <i>NP</i>	
min = 2			3 <i>N, V</i> <i>NP, VP</i>

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules

| lead | can | poison |  
 0      1      2      3

max = 1      max = 2      max = 3

min = 0	1 <i>N, V</i> <i>NP, VP</i>	4 ?	
min = 1		2 <i>N, M</i> <i>NP</i>	
min = 2			3 <i>N, V</i> <i>NP, VP</i>

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

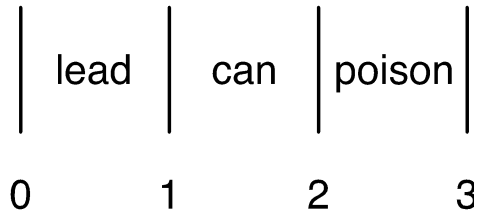
$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules



max = 1      max = 2      max = 3

min = 0	<div style="display: flex; flex-direction: column; justify-content: space-between;"> <span>1</span> <span><i>N, V</i></span> <span><i>NP, VP</i></span> </div>	<div style="display: flex; flex-direction: column; justify-content: space-between;"> <span>4</span> <span><i>NP</i></span> </div>	
min = 1	<div style="display: flex; flex-direction: column; justify-content: space-between;"> <span>2</span> <span><i>N, M</i></span> <span><i>NP</i></span> </div>		
min = 2		<div style="display: flex; flex-direction: column; justify-content: space-between;"> <span>3</span> <span><i>N, V</i></span> <span><i>NP, VP</i></span> </div>	

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules

lead	can	poison
0	1	2

max = 1      max = 2      max = 3

min = 0	<div style="display: flex; justify-content: space-between;"> <span>1</span> <span><i>N, V</i></span> </div> <div style="display: flex; justify-content: space-between;"> <span><i>NP, VP</i></span> <span>4</span> </div>	<div style="display: flex; justify-content: space-between;"> <span>4</span> <span><i>NP</i></span> </div>	
min = 1	<div style="display: flex; justify-content: space-between;"> <span>2</span> <span><i>N, M</i></span> </div> <div style="display: flex; justify-content: space-between;"> <span><i>NP</i></span> <span>5</span> </div>	?	
min = 2		<div style="display: flex; justify-content: space-between;"> <span>3</span> <span><i>N, V</i></span> </div> <div style="display: flex; justify-content: space-between;"> <span><i>NP, VP</i></span> </div>	

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules



lead	can	poison
0	1	2
1	2	3

max = 1      max = 2      max = 3

min = 0	<sup>1</sup> <i>N, V</i> <i>NP, VP</i>	<sup>4</sup> <i>NP</i>	
min = 1		<sup>2</sup> <i>N, M</i> <i>NP</i>	<sup>5</sup> <i>S, VP,</i> <i>NP</i>
min = 2			<sup>3</sup> <i>N, V</i> <i>NP, VP</i>

$S \rightarrow NP VP$

$VP \rightarrow M V$   
 $VP \rightarrow V$

$NP \rightarrow N$   
 $NP \rightarrow N NP$

---

*N* → *can*  
*N* → *lead*  
*N* → *poison*

*M* → *can*  
*M* → *must*

*V* → *poison*  
*V* → *lead*

Inner rules

Preterminal rules

| lead | can | poison |  
 0      1      2      3

max = 1      max = 2      max = 3

min = 0	1 <i>N, V</i> <i>NP, VP</i>	4 <i>NP</i>	6 ?
min = 1		2 <i>N, M</i> <i>NP</i>	5 <i>S, VP,</i> <i>NP</i>
min = 2			3 <i>N, V</i> <i>NP, VP</i>

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules

| lead | can | poison |  
 0      1      2      3

max = 1      max = 2      max = 3

min = 0	1 <i>N, V</i> <i>NP, VP</i>	4 <i>NP</i>	6 ?
min = 1		2 <i>N, M</i> <i>NP</i>	5 <i>S, VP,</i> <i>NP</i>
min = 2			3 <i>N V</i> <i>NP VP</i>

$S \rightarrow NP VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules

| lead | can | poison |  
 0      1      2      3

max = 1      max = 2      max = 3

min = 0	1 N, V NP, VP	4 NP	6 S, NP
min = 1		2 N, M NP	5 S, VP, NP
min = 2			3 N, V NP, VP

mid=1

$S \rightarrow NP VP$

$VP \rightarrow M V$   
 $VP \rightarrow V$

$NP \rightarrow N$   
 $NP \rightarrow N NP$

$N \rightarrow can$   
 $N \rightarrow lead$   
 $N \rightarrow poison$

$M \rightarrow can$   
 $M \rightarrow must$

$V \rightarrow poison$   
 $V \rightarrow lead$

Inner rules

Preterminal rules

$S \rightarrow NP VP$

0 | lead | 1 | can | 2 | poison | 3

max = 1    max = 2    max = 3

min = 0	1 $N, V$ $NP, VP$	4 $NP$	6 $S, NP$ $S(?!)$
min = 1		2 $N, M$ $NP$	5 $S, VP,$ $NP$
min = 2			3 $N, V$ $NP, VP$

*mid=2*

$VP \rightarrow M V$   
 $VP \rightarrow V$

$NP \rightarrow N$   
 $NP \rightarrow N NP$

$N \rightarrow can$   
 $N \rightarrow lead$   
 $N \rightarrow poison$

$M \rightarrow can$   
 $M \rightarrow must$

$V \rightarrow poison$   
 $V \rightarrow lead$

Inner rules

Preterminal rules

$S \rightarrow NP VP$

0 | lead | 1 | can | 2 | poison | 3

max = 1      max = 2      max = 3

min = 0	1 $N, V$ $NP, VP$	4 $NP$	6 $S, NP$ $S(?!)$
min = 1		2 $N, M$ $NP$	5 $S, VP,$ $NP$
min = 2			3 $N, V$ $NP, VP$

$VP \rightarrow M V$

$VP \rightarrow V$

$NP \rightarrow N$

$NP \rightarrow N NP$

$N \rightarrow can$

$N \rightarrow lead$

$N \rightarrow poison$

$M \rightarrow can$

$M \rightarrow must$

$V \rightarrow poison$

$V \rightarrow lead$

Inner rules

Preterminal rules

Apparently the sentence is ambiguous for the grammar: (as the grammar overgenerates)

# PCFGs

$S \rightarrow NP VP$  1.0

$N \rightarrow girl$  0.2

$N \rightarrow telescope$  0.7

$VP \rightarrow V$  0.2

$VP \rightarrow V NP$  0.4

$N \rightarrow sandwich$  0.1

$VP \rightarrow VP PP$  0.4

$PN \rightarrow I$  1.0

$V \rightarrow saw$  0.5

$V \rightarrow ate$  0.5

$NP \rightarrow NP PP$  0.3

$NP \rightarrow D N$  0.5

$P \rightarrow with$  0.6

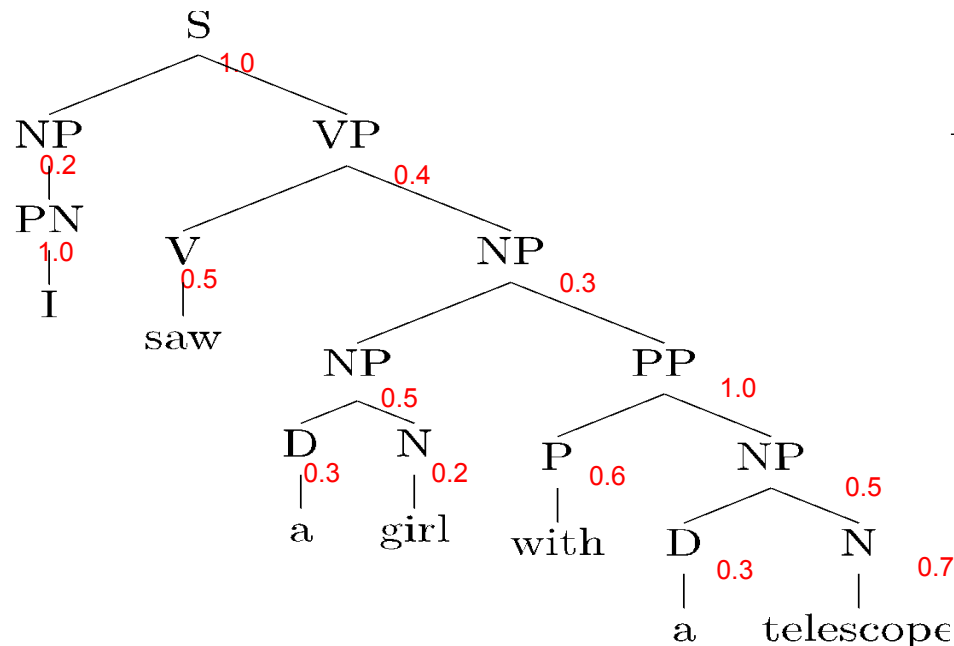
$NP \rightarrow PN$  0.2

$P \rightarrow in$  0.4

$PP \rightarrow P NP$  1.0

$D \rightarrow a$  0.3

$D \rightarrow the$  0.7



$$\begin{aligned} p(T) &= 1.0 \times 0.2 \times 1.0 \times 0.4 \times 0.5 \times 0.3 \times \\ &\quad 0.5 \times 0.3 \times 0.2 \times 1.0 \times 0.6 \times 0.5 \times 0.3 \times 0.7 \\ &= 2.26 \times 10^{-5} \end{aligned}$$



# CKY with PCFGs

---

- Chart is represented by a 3d array of **floats**

`chart [min] [max] [label]`

- It stores probabilities for the most probable subtree with a given signature
- `chart [0] [n] [S]` will store the probability of **the most probable full parse tree**





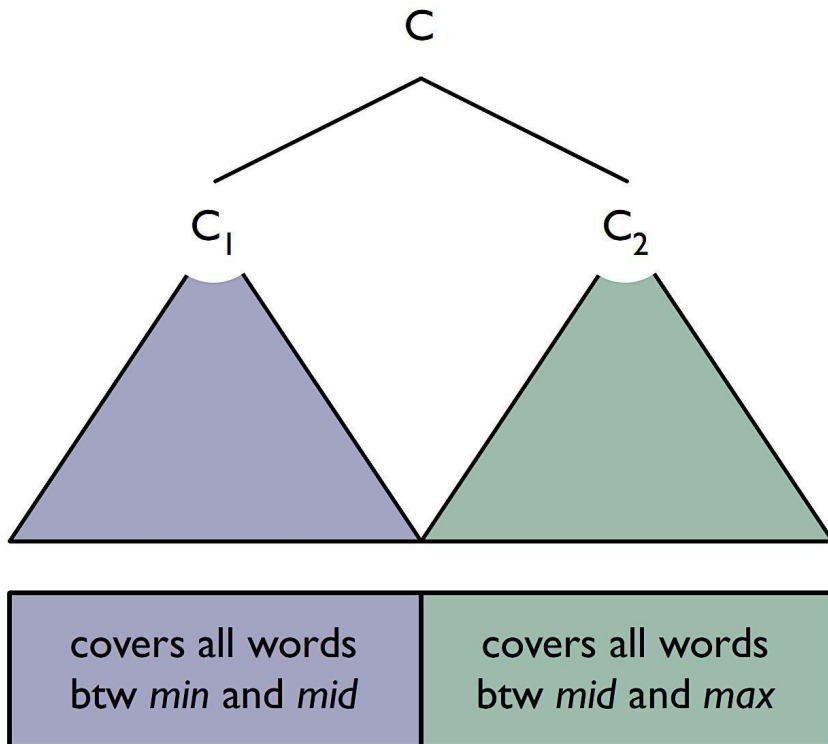
# Intuition

$$C \rightarrow C_1 C_2$$

For every  $C$  choose  $C_1, C_2$  and  $mid$  such that

$$P(T_1) \times P(T_2) \times P(C \rightarrow C_1 C_2)$$

is maximal, where  $T_1$  and  $T_2$  are left and right subtrees.





# Implementation: preterminal rules

---

for each  $w_i$  from left to right

for each preterminal rule  $C \rightarrow w_i$

$\text{chart}[i - 1][i][C] = p(C \rightarrow w_i)$



# Implementation: binary rules

```
for each max from 2 to n
```

```
for each min from max - 2 down to 0
```

```
for each syntactic category C
```

```
double best = undefined
```

```
for each binary rule C -> C1 C2
```

```
for each mid from min + 1 to max - 1
```

```
double t1 = chart[min][mid][C1]
```

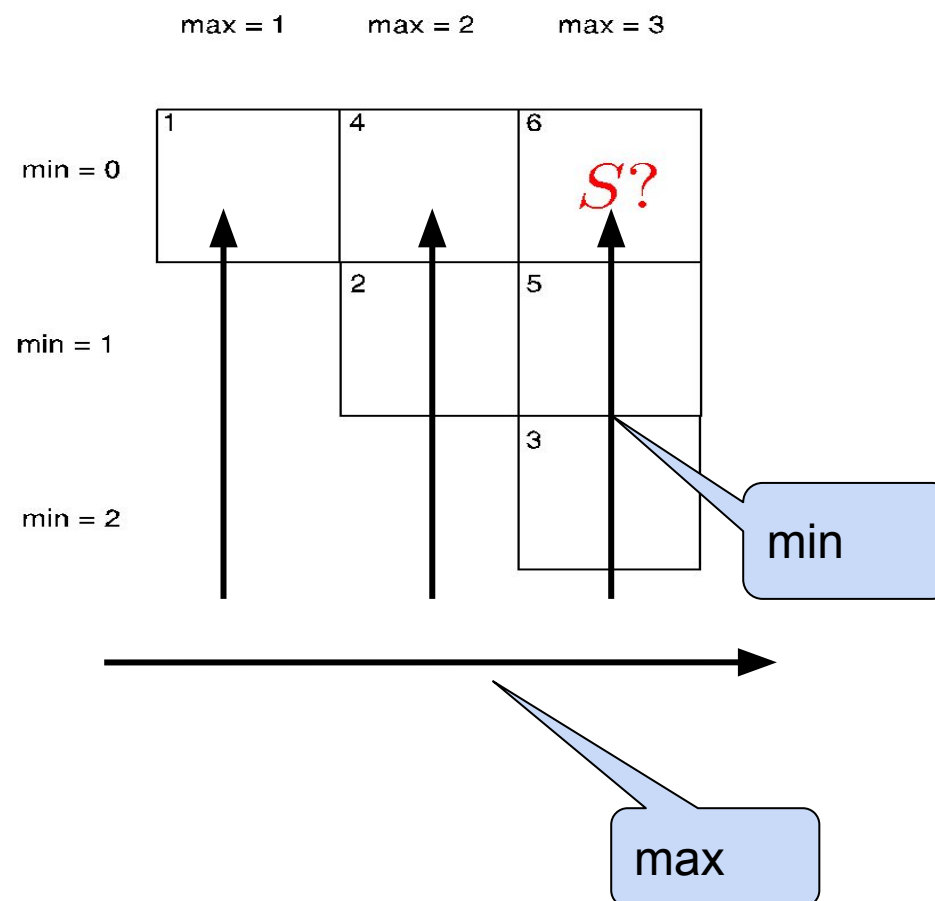
```
double t2 = chart[mid][max][C2]
```

```
double candidate = t1 * t2 * p(C -> C1 C2)
```

```
if candidate > best then
```

```
best = candidate
```

```
chart[min][max][C] = best
```





# Recovery of the tree

---

- For each signature we store backpointers to the elements from which it was built
  - start recovering from  $[0, n, S]$
- What backpointers do we store?



# Recovery of the tree

---

- For each signature we store backpointers to the elements from which it was built
  - start recovering from  $[0, n, S]$
- What backpointers do we store?
  - rule
  - for binary rules, midpoint



# Constraints on the grammar

---

- The basic CKY algorithm supports only rules in the **Chomsky Normal Form (CNF)**:

$$C \rightarrow x$$

$$C \rightarrow C_1 C_2$$

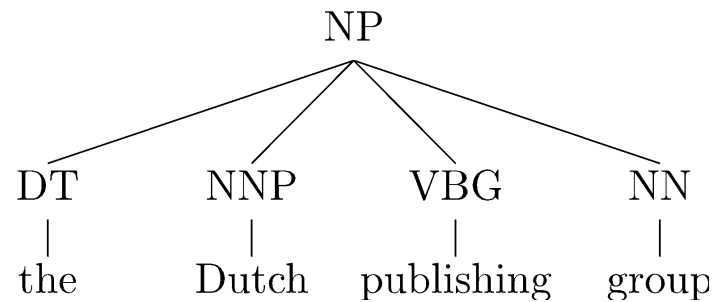
- Any CFG can be converted to an equivalent CNF
  - Equivalent means that they define **the same language**
  - However (syntactic) **trees will look differently**
  - It is possible to address it by defining such transformations that allows for easy **reverse transformation**



# Transformation to CNF form: binarization

---

- Consider  $NP \rightarrow DT\ NNP\ VBG\ NN$

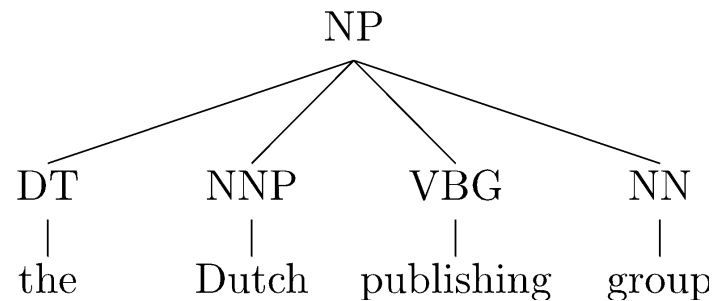


- How do we get a set of binary rules which are equivalent?



# Transformation to CNF form: binarization

- Consider  $NP \rightarrow DT \ NNP \ VBG \ NN$



- How do we get a set of binary rules which are equivalent?

$NP \rightarrow DT \ X$

$X \rightarrow NNP \ Y$

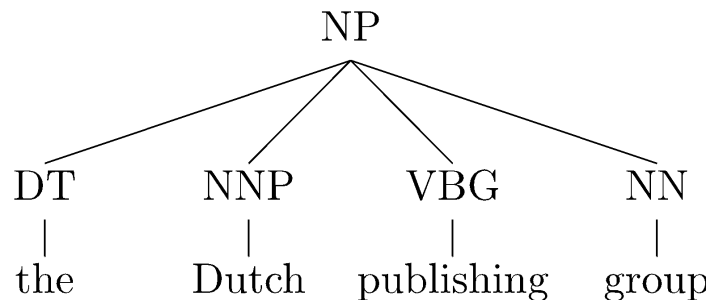
$Y \rightarrow VBG \ NN$





# Transformation to CNF form: binarization

- Consider  $NP \rightarrow DT\ NNP\ VBG\ NN$



- How do we get a set of binary rules which are equivalent?

$NP \rightarrow DT\ X$

$X \rightarrow NNP\ Y$

$Y \rightarrow VBG\ NN$

- A more systematic way to refer to new non-terminals

$NP \rightarrow DT\ @NP|DT$

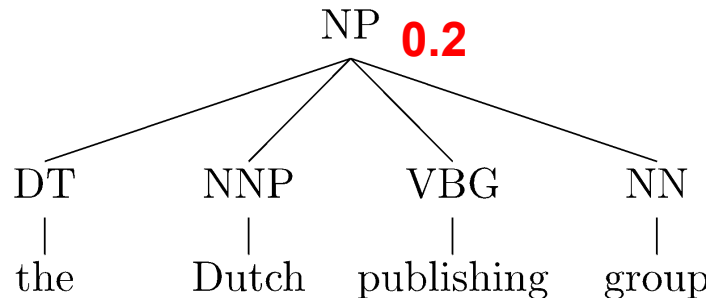
$@NP|DT \rightarrow NNP\ @NP|DT\_NNP$

$@NP|DT\_NNP \rightarrow VBG\ NN$



# Transformation to CNF form: binarization

- Consider  $NP \rightarrow DT\ NNP\ VBG\ NN$  **0.2**

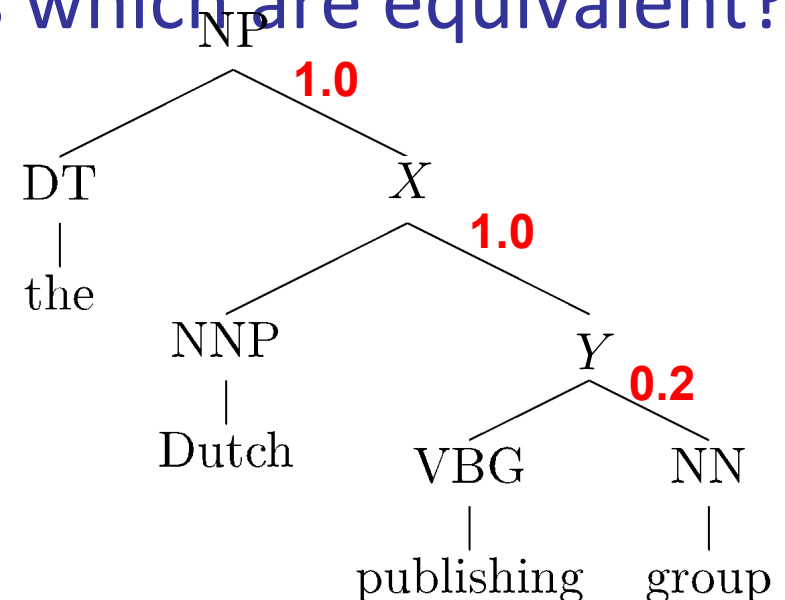


- How do we get a set of binary rules which are equivalent?

$NP \rightarrow DT\ X$  **1.0**

$X \rightarrow NNP\ Y$  **1.0**

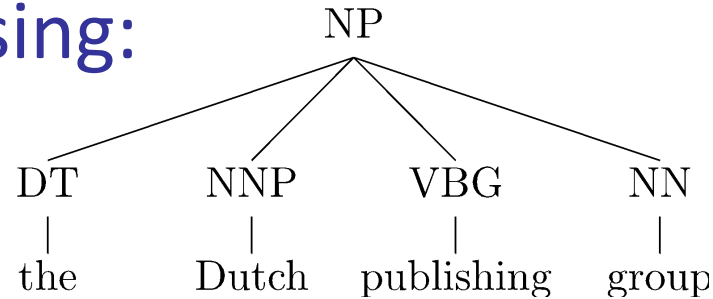
$Y \rightarrow VBG\ NN$  **0.2**



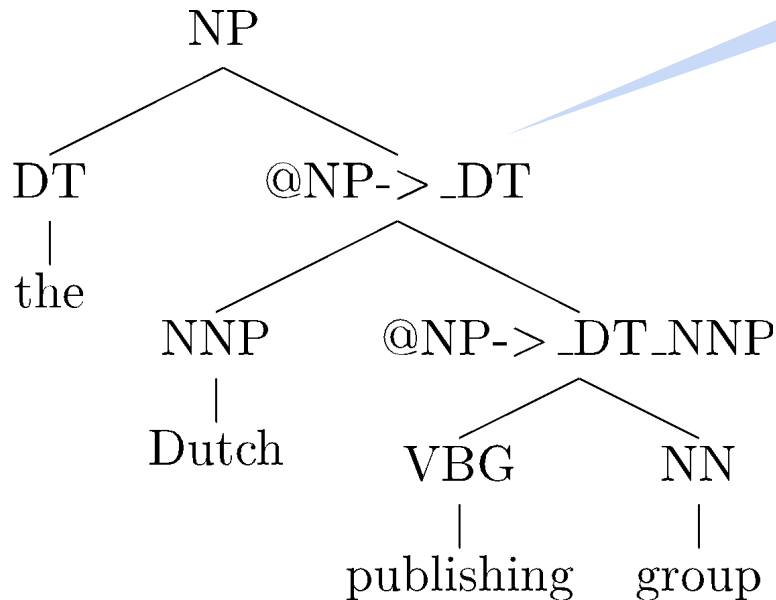


# Transformation to CNF form: binarization

- Instead of binarizing tuples we can binarize trees on preprocessing:



Also known as **lossless Markovization** in the context of PCFGs



Can be easily reversed on postprocessing



# Unary Rules

---

- CNF includes only two types of rules:

$$C \rightarrow x$$

$$C \rightarrow C_1 C_2$$

- What about unary rules:

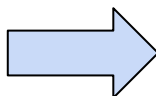
$$C \rightarrow C_1$$



# Unary Rules

## CFG

A → X  
B → X  
C → X  
...  
X → C<sub>1</sub>C<sub>2</sub>  
...  
X → run  
X → play  
X → sleep  
X → love



## CNF

A → run,    B → run,    C → run,    X → run,  
A → play,    B → play,    C → play,    X → play,  
A → sleep,    B → sleep,    C → sleep,    X → sleep,  
A → love    B → love    C → love    X → love  
...    ...    ...    ...  
A → C<sub>1</sub>C<sub>2</sub>    B → C<sub>1</sub>C<sub>2</sub>    C → C<sub>1</sub>C<sub>2</sub>    X → C<sub>1</sub>C<sub>2</sub>

- explode the grammar
- make it hard to reverse

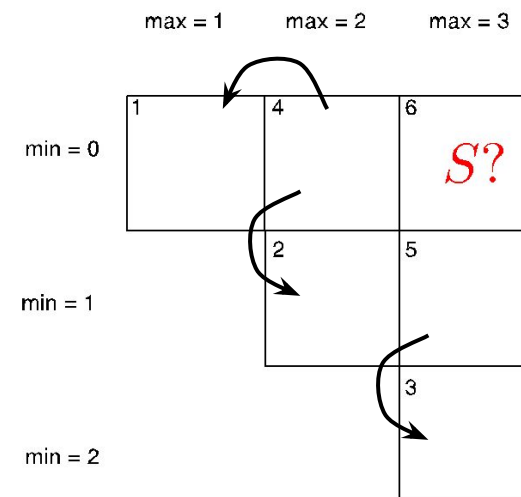


# Unary rules

- How to integrate unary rules  $C \rightarrow C_1$  ?

```
for each max from 1 to n
    for each min from max - 1 down to 0
        // First, try all binary rules as before.
        ...
        // Then, try all unary rules.
        for each syntactic category C
            for each unary rule  $C \rightarrow C_1$ 
                chart[min][max][C] = maximum (chart[min][max][C],
                                                chart[min][max][C1])
```

**new bounds!**





# Unary closure

- What if the grammar contained 2 rules:

$$A \rightarrow B$$

$$B \rightarrow C$$

- But C can be derived from A by a chain of rules:

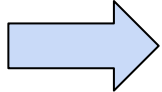
$$A \rightarrow B \rightarrow C$$

- One could support chains in the algorithm but it is easier to extend the grammar, to get the **transitive closure**

$$\begin{array}{l} A \rightarrow B \\ B \rightarrow C \end{array} \quad \Rightarrow \quad \begin{array}{l} A \rightarrow B \\ B \rightarrow C \\ A \rightarrow C \end{array}$$



# Why unary closure

$A \rightarrow B$   
 $B \rightarrow C$        $A \rightarrow C$

```
// Then, try all unary rules.
```

```
for each syntactic category C
```

```
  for each unary rule  $C \rightarrow C_1$ 
```

```
    if chart[min][max][C1] then
```

```
      chart[min][max][C] = true
```

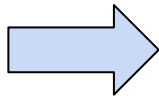




# Why unary closure

$A \rightarrow B$

$B \rightarrow C$



$A \rightarrow C$

**scenario 1**

1 C	4	6 <i>S?</i>
	2	5
		3

$B \rightarrow C$

1 C, B	4	6 <i>S?</i>
	2	5
		3

$A \rightarrow B$

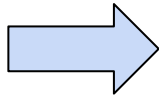
1 C, B, A	4	6 <i>S?</i>
	2	5
		3



# Why unary closure

$A \rightarrow B$

$B \rightarrow C$



$A \rightarrow C$

**scenario 1**

1	C	4		6	$S?$
		2		5	
			3		

$B \rightarrow C$

1	C, B	4		6	$S?$
		2		5	
			3		

$A \rightarrow B$

1	C, B, A	4		6	$S?$
		2		5	
			3		

**scenario 2**

1	C	4		6	$S?$
		2		5	
			3		

$A \rightarrow B$

1	C	4		6	$S?$
		2		5	
			3		

$B \rightarrow C$

1	C, B	4		6	$S?$
		2		5	
			3		



# Unary closure

- What if the grammar contained 2 rules:

$$A \rightarrow B$$

$$B \rightarrow C$$

- But C can be derived from A by a chain of rules:

$$A \rightarrow B \rightarrow C$$

- One could support chains in the algorithm but it is easier to extend the grammar, to get the **transitive closure**

$$\begin{array}{l} A \rightarrow B \\ B \rightarrow C \end{array} \quad \Rightarrow \quad \begin{array}{l} A \rightarrow B \\ B \rightarrow C \\ A \rightarrow C \end{array} \quad \begin{array}{l} A \rightarrow A \\ B \rightarrow B \\ C \rightarrow C \end{array}$$

Convenient for programming reasons in the PCFG case



# Unary (reflexive transitive) closure

$A \rightarrow B$	0.1	$\Rightarrow$	$A \rightarrow B$	0.1	$A \rightarrow A$	1
$B \rightarrow C$	0.2		$B \rightarrow C$	0.2	$B \rightarrow B$	1
...			$A \rightarrow C$	$0.2 \times 0.1$	$C \rightarrow C$	1
				...		

Note that this is not a PCFG anymore as the rules do not sum to 1 for each parent



# Unary (reflexive transitive) closure

The fact that the rule is composite needs to be stored to recover the true tree

$A \rightarrow B$	0.1	$\Rightarrow$	$A \rightarrow B$	0.1	$A \rightarrow A$	1
$B \rightarrow C$	0.2		$B \rightarrow C$	0.2	$B \rightarrow B$	1
...			$A \rightarrow C$	$0.2 \times 0.1$	$C \rightarrow C$	1
				...		

Note that this is not a PCFG anymore as the rules do not sum to 1 for each parent



# Unary (reflexive transitive) closure

The fact that the rule is composite needs to be stored to recover the true tree

$A \rightarrow B$	0.1	$\Rightarrow$	$A \rightarrow B$	0.1	$A \rightarrow A$	1
$B \rightarrow C$	0.2		$B \rightarrow C$	0.2	$B \rightarrow B$	1
...			$A \rightarrow C$	$0.2 \times 0.1$	$C \rightarrow C$	1
			...		...	

Note that this is not a PCFG anymore as the rules do not sum to 1 for each parent

$A \rightarrow B$	0.1	$\Rightarrow$	$A \rightarrow B$	0.1	$A \rightarrow A$	1
$B \rightarrow C$	0.2		$B \rightarrow C$	0.1	$B \rightarrow B$	1
$A \rightarrow C$	$1.e - 5$		$A \rightarrow C$	0.02	$C \rightarrow C$	1

What about loops, like:  $A \rightarrow B \rightarrow A \rightarrow C$  ?



# Recovery of the tree

---

- For each signature we store backpointers to the elements from which it was built
  - start recovering from  $[0, n, S]$
- What do we store in backpointers?
  - rule
  - for binary rules, midpoint
- Be careful with unary rules
  - Basically you can assume that you always used an unary rule from the closure (but it could be the trivial one  $C \rightarrow C$ )



# Speeding up the algorithm

---

- **Basic pruning (roughly):**
  - For every span  $(i,j)$  store only labels which have the probability at most  $N$  times smaller than the probability of the most probable label for this span
  - Check not all rules but only rules yielding subtree labels having non-zero probability
- **Coarse-to-fine pruning**
  - Parse with a smaller (simpler) grammar, and precompute (posterior) probabilities for each spans, and use only the ones with non-negligible probability from the previous grammar





# Parsing evaluation

---

- **Intrinsic evaluation:**
  - **Automatic:** evaluate against annotation provided by human experts (gold standard) according to some predefined measure
  - **Manual:** ... according to human judgment
- **Extrinsic evaluation:** score syntactic representation by comparing how well a system using this representation performs on some task
  - E.g., use syntactic representation as input for a semantic analyzer and compare results of the analyzer using syntax predicted by different parsers.



# Standard evaluation setting in parsing

---

- Automatic intrinsic evaluation is used: parsers are evaluated against gold standard by provided by linguists
  - There is a standard split into the parts:
    - training set: used for estimation of model parameters
    - development set: used for tuning the model (initial experiments)
    - test set: final experiments to compare against previous work



# Automatic evaluation of constituent parsers

---

- **Exact match:** percentage of trees predicted correctly
- **Bracket score:** scores how well individual phrases (and their boundaries) are identified

The most standard measure;  
we will focus on it



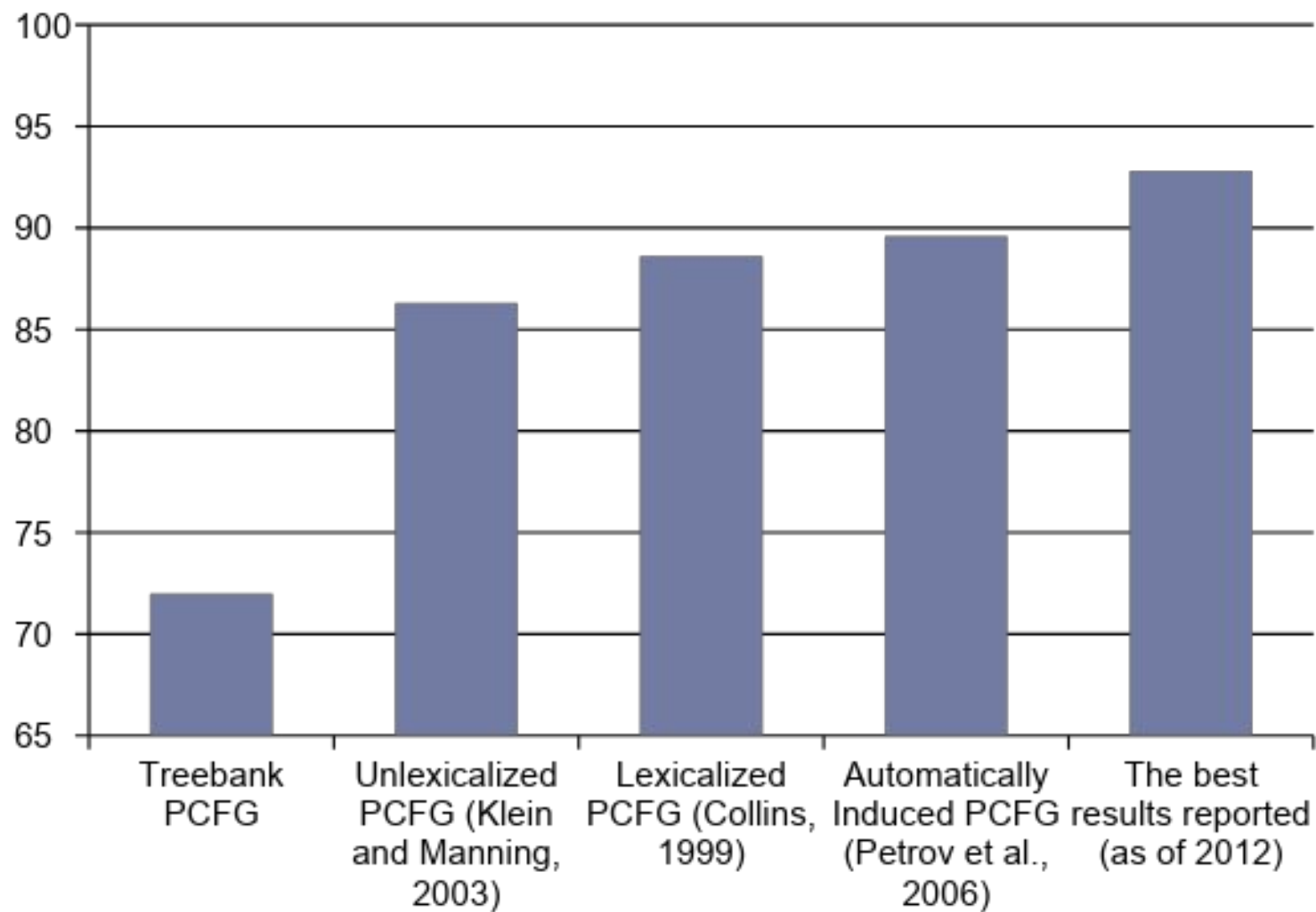
# Brackets scores

Subtree signatures for  
CKY

- The most standard score is **bracket score**
- It regards a tree as a collection of brackets:  $[min, max, C]$
- The set of brackets predicted by a parser is compared against the set of brackets in the tree annotated by a linguist
- **Precision, recall** and **F1** are used as scores



# Preview: F1 bracket score





---

# Estimating PCFGs



# Estimating PCFGs

Associate probabilities with the rules :  $p(X \rightarrow \alpha)$

$$\forall X \rightarrow \alpha \in R: 0 \leq p(X \rightarrow \alpha) \leq 1$$

$$\forall X \in N: \sum_{\alpha: X \rightarrow \alpha \in R} p(X \rightarrow \alpha) = 1$$

$S \rightarrow NP VP$	1.0	(NP A girl) (VP ate a sandwich)	$N \rightarrow girl$	0.2
$VP \rightarrow V$	0.2		$N \rightarrow telescope$	0.7
$VP \rightarrow V NP$	0.4	(VP ate) (NP a sandwich)	$N \rightarrow sandwich$	0.1
$VP \rightarrow VP PP$	0.4	(VP saw a girl) (PP with ...)	$PN \rightarrow I$	1.0
$NP \rightarrow NP PP$	0.3	(NP a girl) (PP with ....)	$V \rightarrow saw$	0.5
$NP \rightarrow D N$	0.5	(D a) (N sandwich)	$V \rightarrow ate$	0.5
$NP \rightarrow PN$	0.2		$P \rightarrow with$	0.6
$PP \rightarrow P NP$	1.0	(P with) (NP with a sandwich)	$P \rightarrow in$	0.4
			$D \rightarrow a$	0.3
			$D \rightarrow the$	0.7



# Estimating PCFGs: Intuition

---

- Probabilistic Regular Grammar

$$N^i \rightarrow w^j N^k$$

$$N^i \rightarrow w^j$$

Start state,  $N^1$





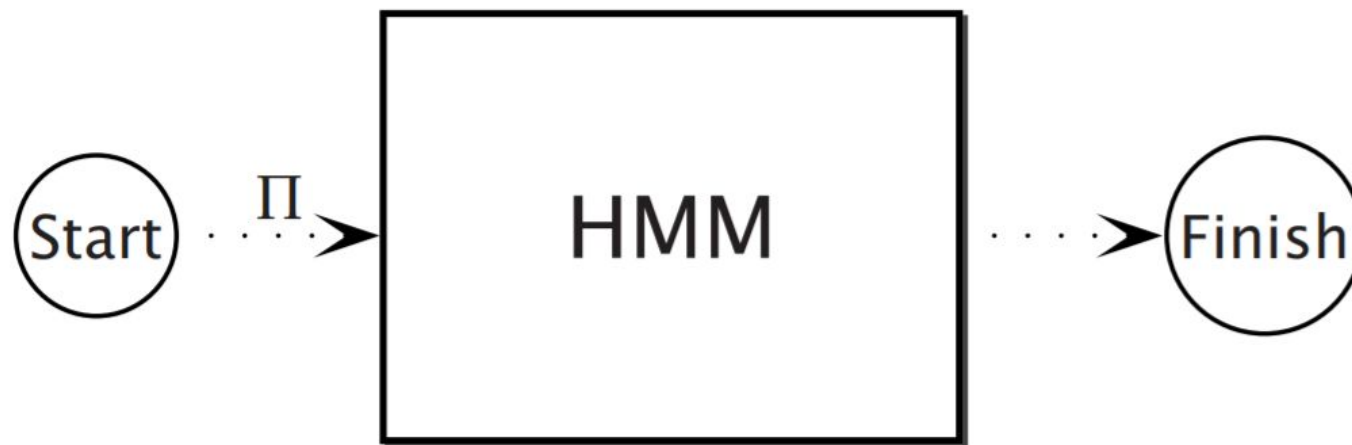
# Estimating PCFGs: Intuition

- Probabilistic Regular Grammar

$$N^i \rightarrow w^j N^k$$

$$N^i \rightarrow w^j$$

Start state,  $N^1$





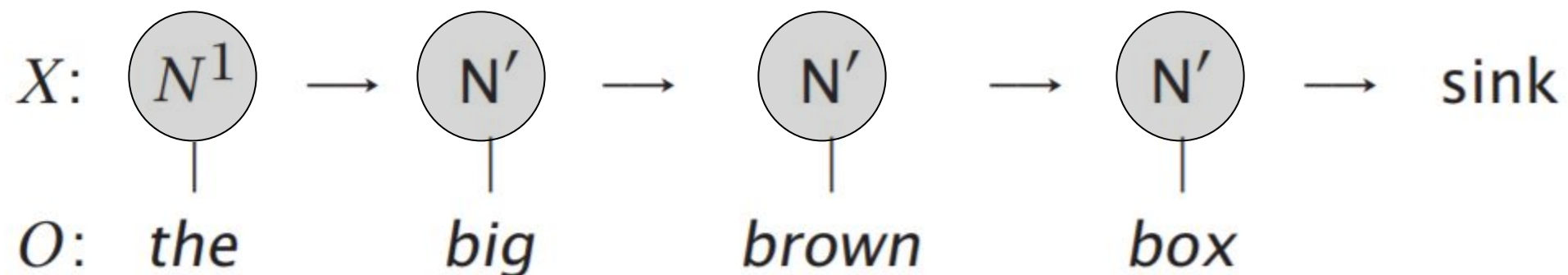
# Estimating PCFGs: Intuition

- Probabilistic Regular Grammar

$$N^i \rightarrow w^j N^k$$

$$N^i \rightarrow w^j$$

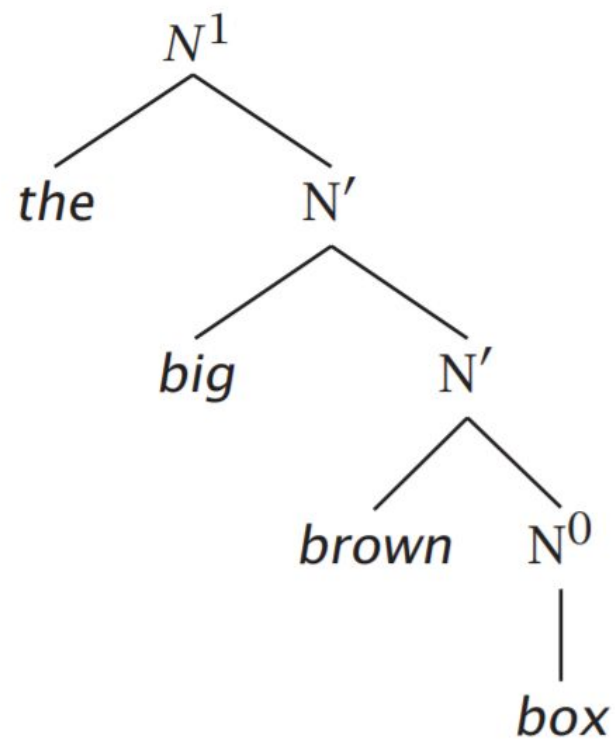
Start state,  $N^1$





# Estimating PCFGs: Intuition

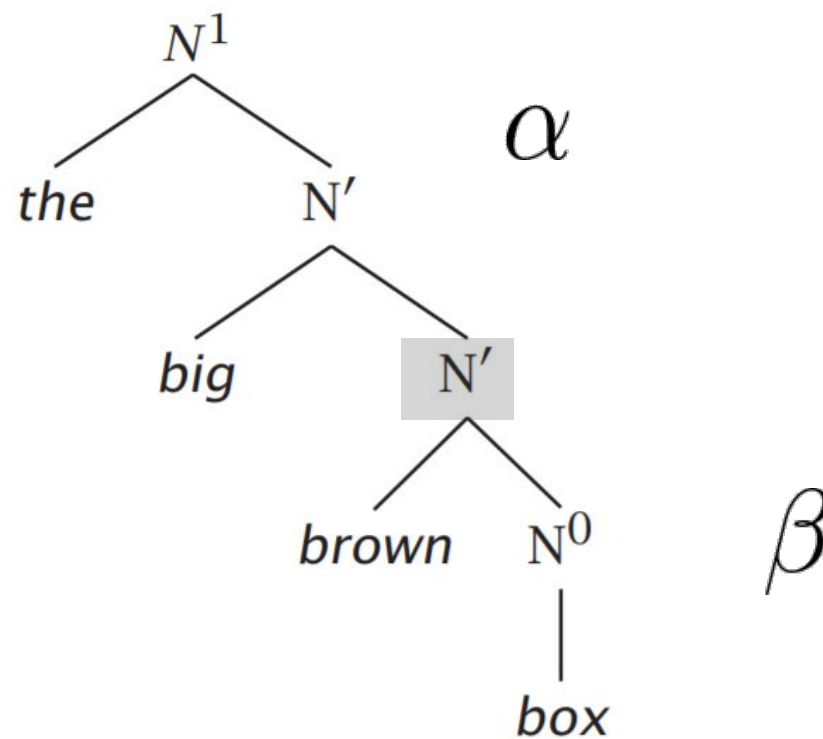
$X: \text{NP} \rightarrow \text{N}' \rightarrow \text{N}' \rightarrow \text{N}' \rightarrow \text{sink}$   
 $O: \textit{the} \quad \textit{big} \quad \textit{brown} \quad \textit{box}$





# Estimating PCFGs: Intuition

X: NP → N' → N' → N' → sink  
O: the big brown box





# Unsupervised estimation of PCFGs

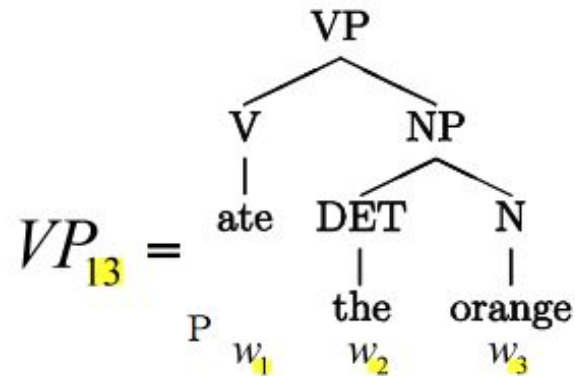
---

- Notation
- Calculating inside probabilities
- Calculating outside probabilities
- The inside-outside algorithm (EM) - preview



# Notation

- Non-terminal symbols (latent variables):  $\{N^1, \dots, N^n\}$
- Sentence (observed data):  $\{w_1, \dots, w_m\} = w_{1:m}$
- $N_{pq}^j$  denotes that  $N^j$  spans  $w_{pq}$  in the sentence





# Inside probability

- ▶ Definition (compare with backward prob for HMMs):

$$\beta_j(p, q) = P(w_p, \dots, w_q | N_{pq}^j, G) = P(N_{pq}^j \rightarrow w_{pq} | G)$$

- ▶ Computed recursively

- ▶ Base case:  $\beta_j(k, k) = P(w_k | N_{kk}^j, G) = P(N_j \rightarrow w_k | G)$
- ▶ Induction:

$$\beta_j(p, q) = \sum_{rs} \sum_{d=p}^{q-1} P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)$$

The grammar is binarized

let's draw...



# Implementation: PCFG parsing

---

```
for each max from 2 to n
  for each min from max - 2 down to 0
    for each syntactic category C
      double best = undefined
      for each binary rule C -> C1 C2
        for each mid from min + 1 to max - 1
          double t1 = chart[min][mid][C1]
          double t2 = chart[mid][max][C2]
          double candidate = t1 * t2 * p(C -> C1 C2)
          if candidate > best then
            best = candidate
      chart[min][max][C] = best
```





# Implementation: inside

---

```
for each max from 2 to n
```

```
  for each min from max - 2 down to 0
```

```
    for each syntactic category C
```

```
      double total = 0.0
```

```
        for each binary rule  $C \rightarrow C_1 C_2$ 
```

```
          for each mid from min + 1 to max - 1
```

```
            double  $t_1$  = chart[min][mid][ $C_1$ ]
```

```
            double  $t_2$  = chart[mid][max][ $C_2$ ]
```

```
            double candidate =  $t_1 * t_2 * p(C \rightarrow C_1 C_2)$ 
```

```
              total = total + candidate
```

```
    chart[min][max][C] = total
```



# Implementation: inside

```
for each max from 2 to n
```

```
  for each min from max - 2 down to 0
```

```
    for each syntactic category C
```

```
      double total = 0.0
```

```
        for each binary rule C -> C1 C2
```

```
          for each mid from min + 1 to max - 1
```

```
            double t1 = chart[min][mid][C1]
```

```
            double t2 = chart[mid][max][C2]
```

```
            double candidate = t1 * t2 * p(C -> C1 C2)
```

```
            total = total + candidate
```

```
          chart[min][max][C] = total
```

$$\beta_j(p, q) = \sum_{rs} \sum_{d=p}^{q-1} P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)$$



# Implementation: inside

```
for each max from 2 to n
```

```
  for each min from max - 2 down to 0
```

```
    for each syntactic category C
```

```
      double total = 0.0
```

```
        for each binary rule C -> C1 C2
```

```
          for each mid from min + 1 to max - 1
```

```
            double t1 = chart[min][mid][C1]
```

```
            double t2 = chart[mid][max][C2]
```

```
            double candidate = t1 * t2 * p(C -> C1 C2)
```

```
          total = total + candidate
```

```
        chart[min][max][C] = total
```

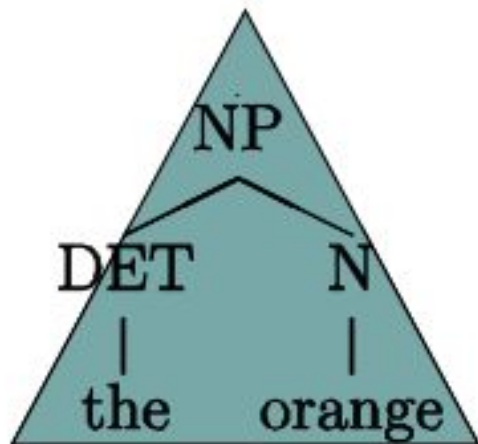
$$\beta_j(p, q) = \sum_{rs} \sum_{d=p}^{q-1} P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)$$



# Inside probability: example

NP→DET N 0.8  
DET→a 0.6  
N→apple 0.8

NP→N 0.2  
DET→the 0.4  
N→orange 0.2



$\beta_{DET}(1,1)$

$\beta_N(2,2)$

$\beta_{NP}(1,2)$

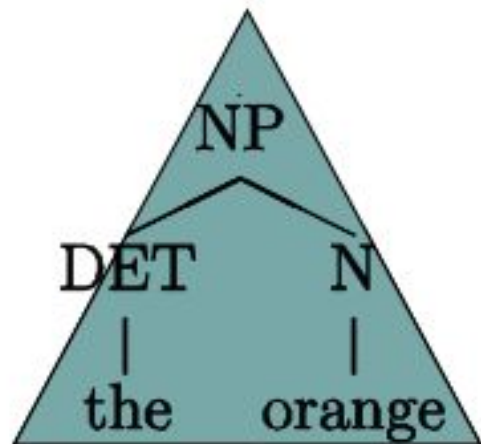
$\beta_{NP}(1,2)$



# Inside probability: example

NP → DET N    0.8  
DET → a        0.6  
N → apple      0.8

NP → N        0.2  
DET → the     0.4  
N → orange    0.2



$$\beta_{DET}(1,1) = P(the | DET_{11}, G) = P(DET \rightarrow the | G) = 0.4$$

$$\beta_N(2,2)$$

$$\beta_{NP}(1,2)$$

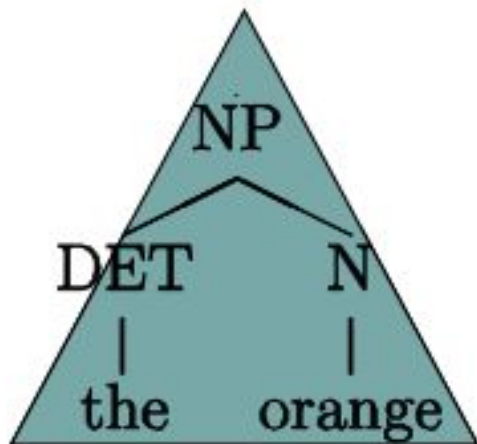
$$\beta_{NP}(1,2)$$



# Inside probability: example

NP → DET N    0.8  
DET → a        0.6  
N → apple      0.8

NP → N        0.2  
DET → the     0.4  
N → orange    0.2



$$\beta_{DET}(1,1) = P(the | DET_{11}, G) = P(DET \rightarrow the | G) = 0.4$$

$$\beta_N(2,2) = P(N \rightarrow orange | G) = 0.2$$

$$\beta_{NP}(1,2)$$

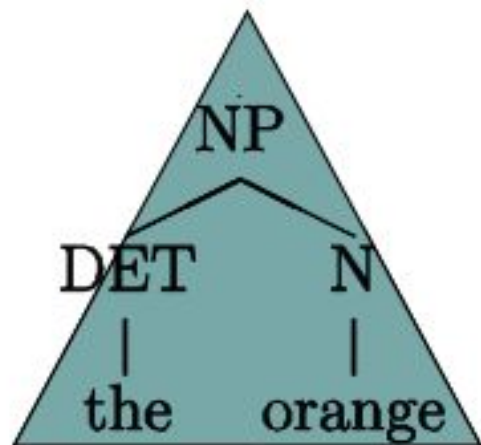
$$\beta_{NP}(1,2)$$



# Inside probability: example

NP → DET N    0.8  
DET → a        0.6  
N → apple      0.8

NP → N        0.2  
DET → the     0.4  
N → orange    0.2



$$\beta_{DET}(1,1) = P(the | DET_{11}, G) = P(DET \rightarrow the | G) = 0.4$$

$$\beta_N(2,2) = P(N \rightarrow orange | G) = 0.2$$

$$\begin{aligned} \beta_{NP}(1,2) &= P(NP \rightarrow DET \cdot N) \beta_{DET}(1,1) \beta_N(2,2) \\ &= 0.8 \quad \quad \quad \times 0.4 \quad \quad \times 0.2 \end{aligned}$$

$$\beta_{NP}(1,2) = 0.064$$

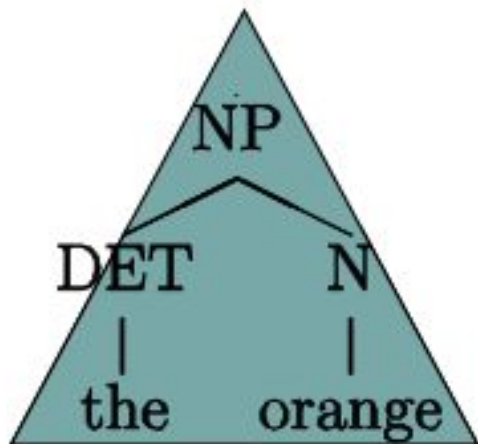




# Inside probability: example

NP → DET N    0.8  
DET → a        0.6  
N → apple      0.8

NP → N        0.2  
DET → the     0.4  
N → orange    0.2



$$\beta_{DET}(1,1) = P(the | DET_{11}, G) = P(DET \rightarrow the | G) = 0.4$$

$$\beta_N(2,2) = P(N \rightarrow orange | G) = 0.2$$

$$\begin{aligned} \beta_{NP}(1,2) &= P(NP \rightarrow DET \cdot N) \beta_{DET}(1,1) \beta_N(2,2) \\ &= 0.8 \quad \quad \quad \times 0.4 \quad \quad \times 0.2 \end{aligned}$$

$$\beta_{NP}(1,2) = 0.064$$

$$\beta_S(1, m) = P(S \rightarrow w_1, \dots, w_m | G)$$



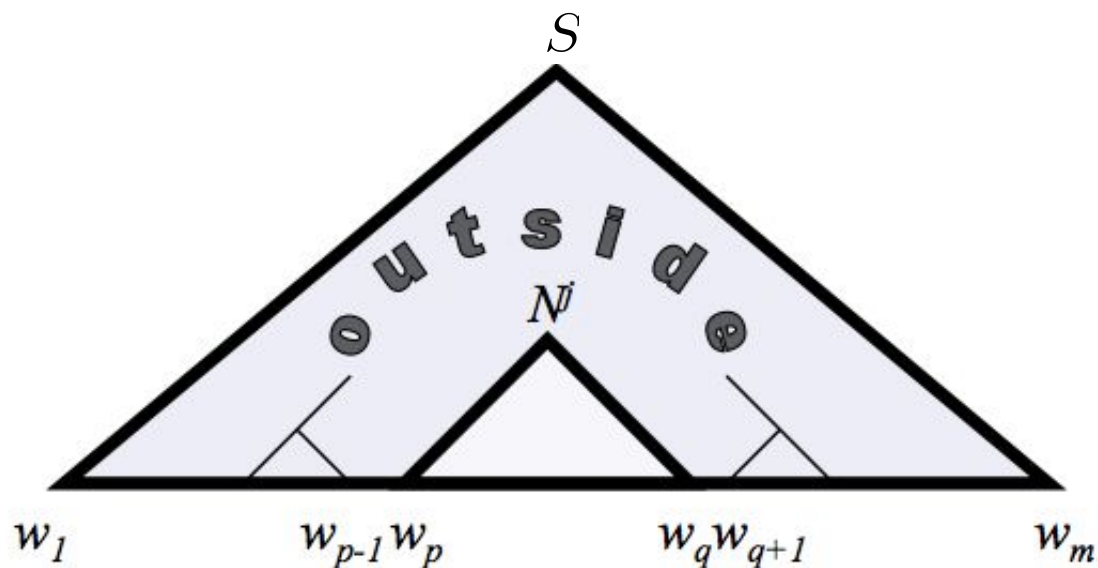


# Outside probability

- ▶ Definition (compare with forward prob for HMMs):

$$\alpha_j(p, q) = P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m} | G)$$

- ▶ The joint probability of starting with  $S$ , generating words  $w_1, \dots, w_{p-1}$ , the non terminal  $N^j$  and words  $w_{q+1}, \dots, w_m$ .





# Calculating outside probability

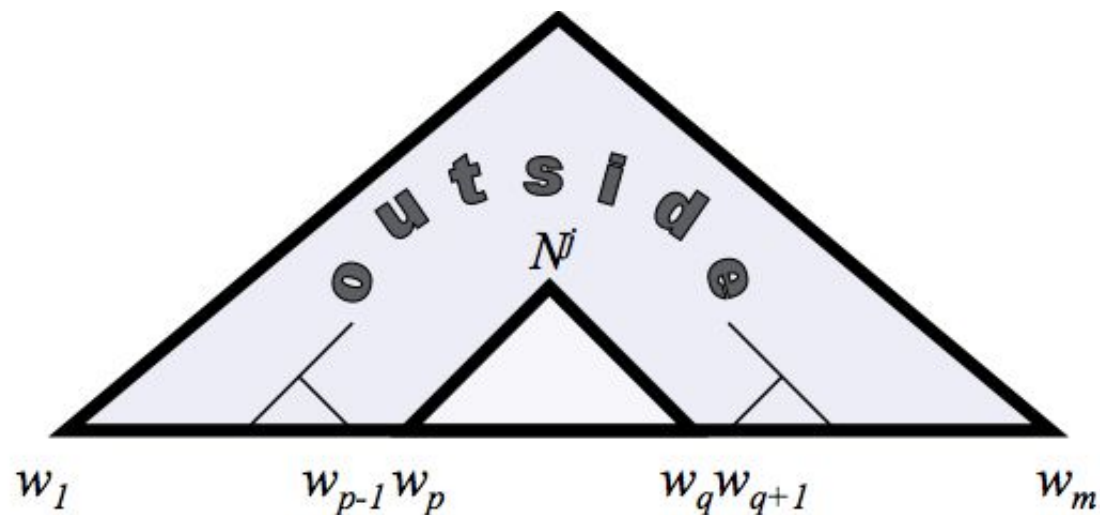
- ▶ Computed recursively, base case

$$\alpha_1(1, m) = \alpha_S(1, m) = 1$$

$$\alpha_{j \neq 1}(1, m) = 0$$

- ▶ Induction?

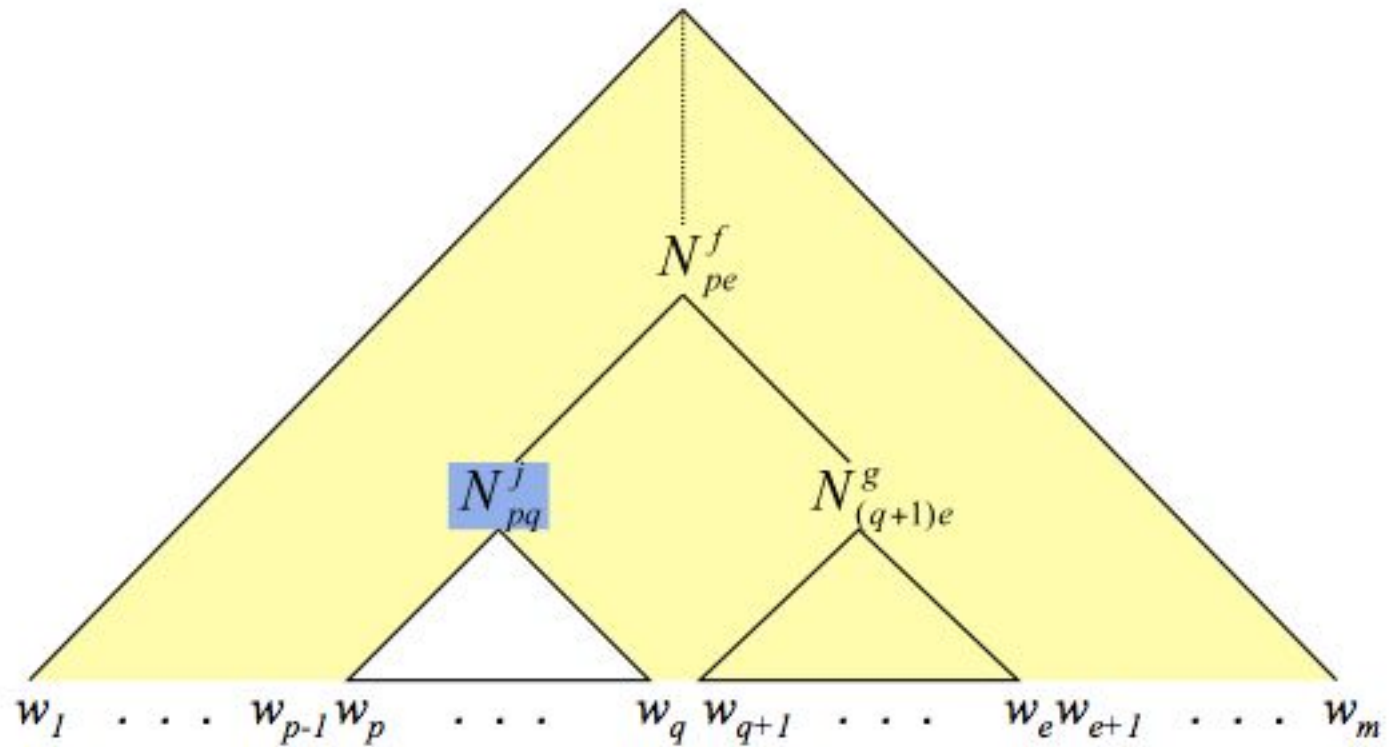
- ▶ Intuition:  $N_{pq}^j$  must be either the L or R child of a parent node. We first consider the case when it is the L child.





# Calculating outside probability

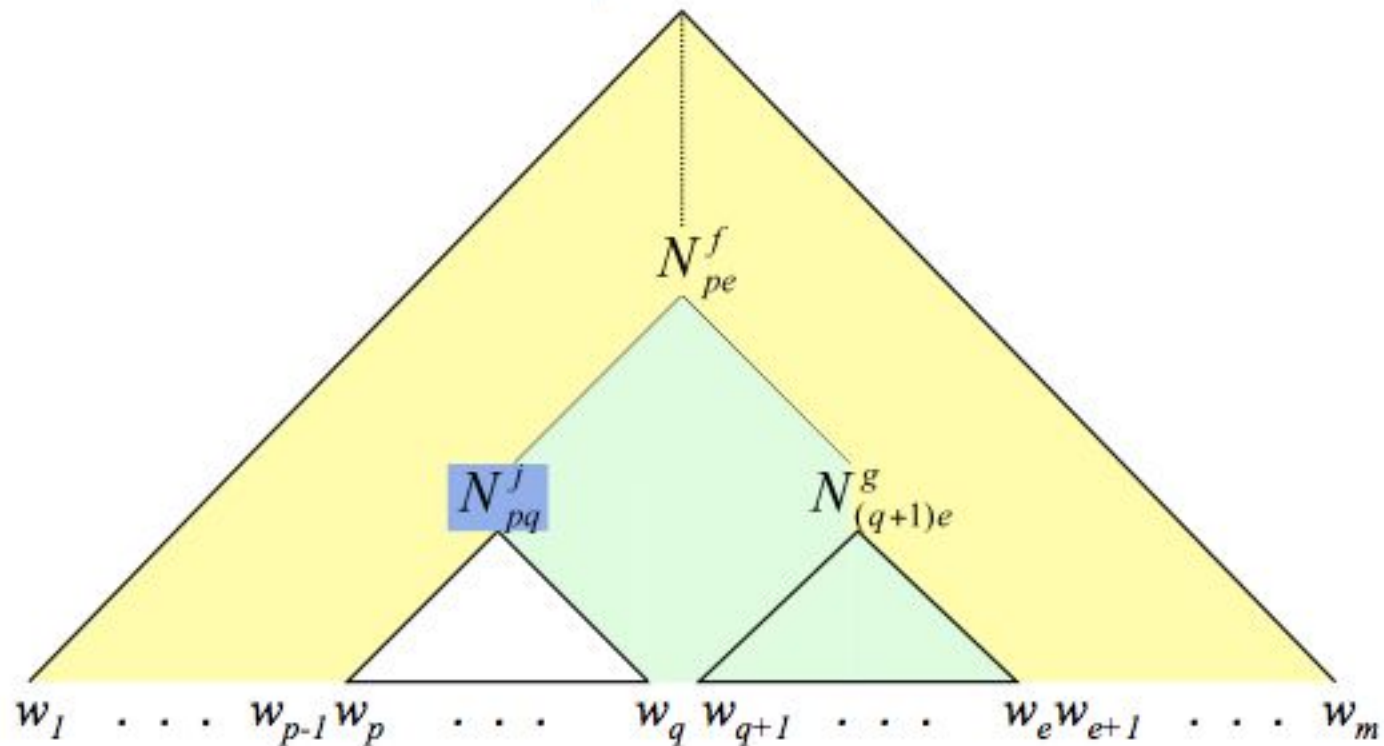
- ▶ The yellow area is the probability we would like to calculate
  - ▶ How do we decompose it?





# Calculating outside probability

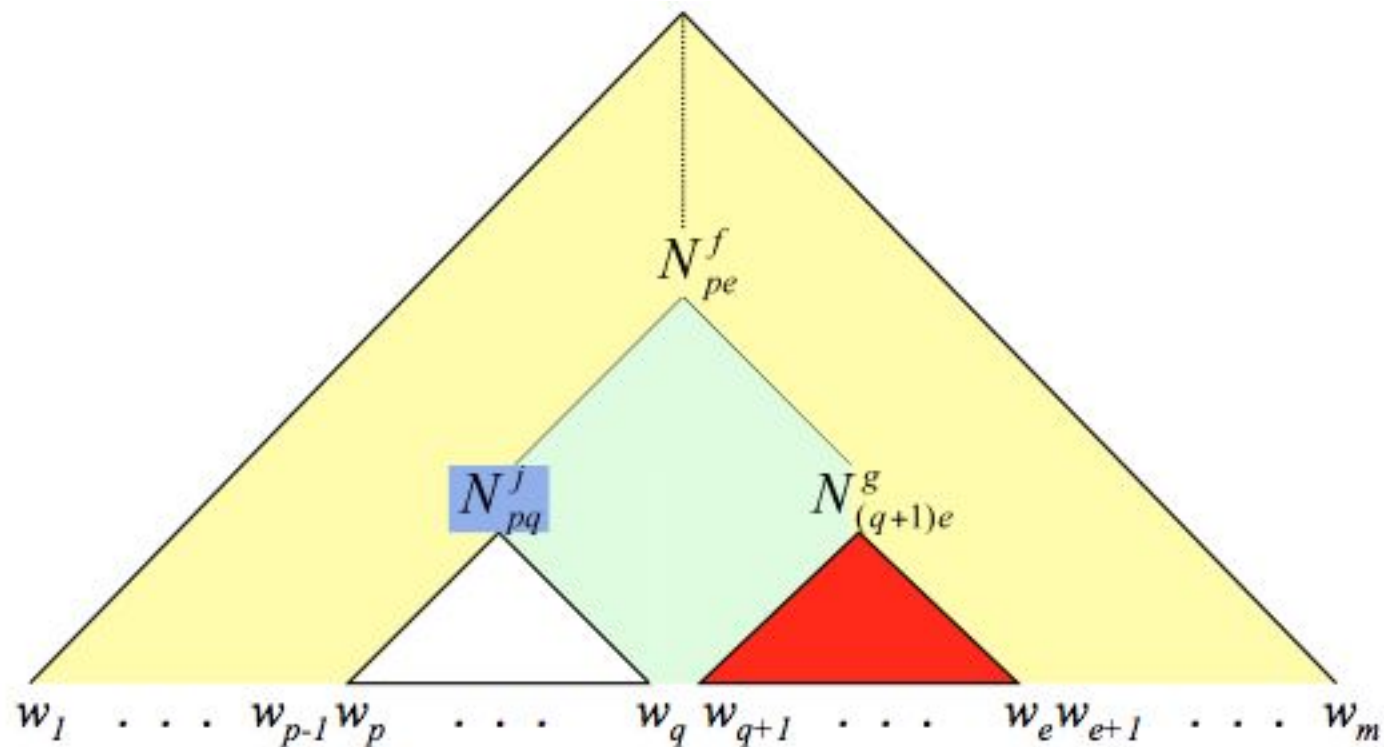
- ▶ Step 1: We assume that  $N_{pe}^f$  is the parent of  $N_{pq}^j$ . Its outside probability,  $\alpha_f(p, e)$  (represented by the yellow shading) is available recursively. But how do we compute the green part?





# Calculating outside probability

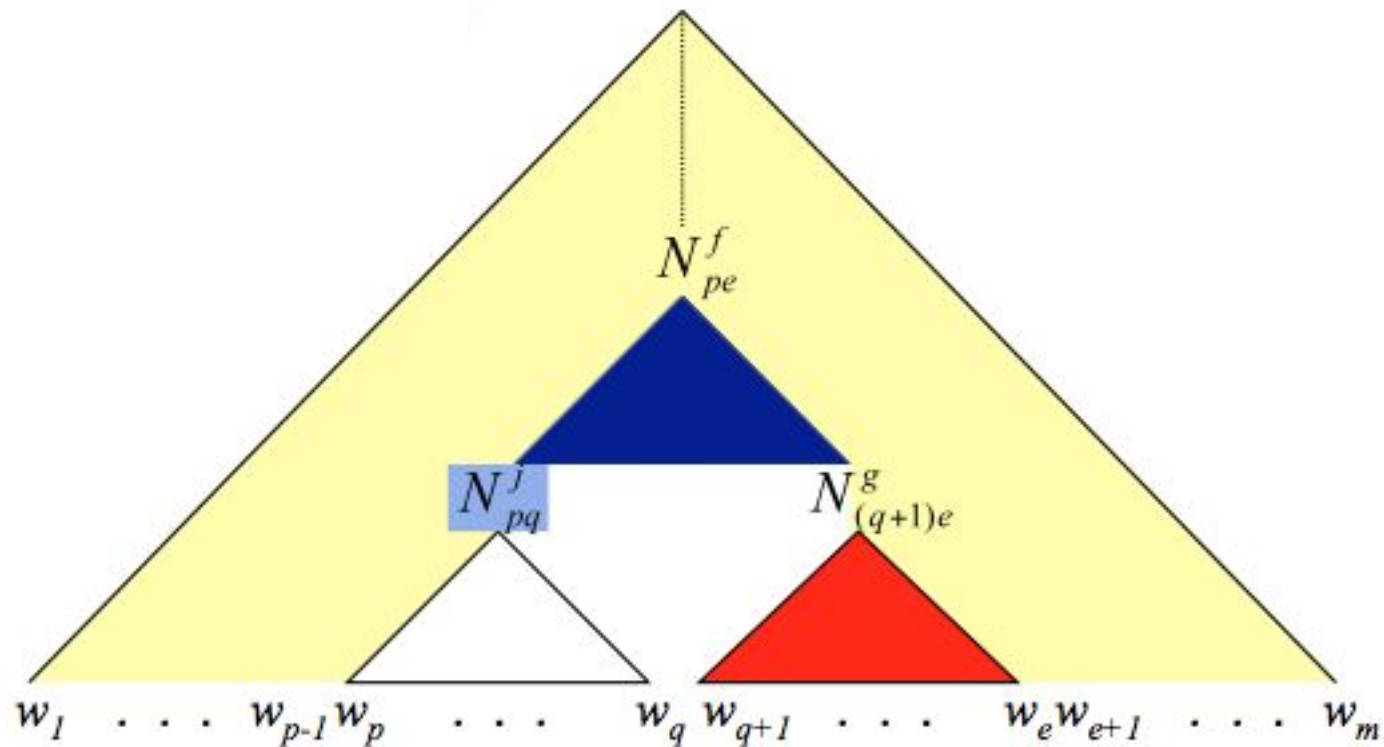
- ▶ Step 1: The red shaded area is the inside probability for  $N_{(q+1)e}^g$ , i.e.  $\beta_q(q+1, e)$





# Calculating outside probability

- ▶ Step 3: The blue shaded area is just the production  $N^f \rightarrow N^j N^g$ , the corresponding probability  $P(N^f \rightarrow N^j N^g | N^f, G)$



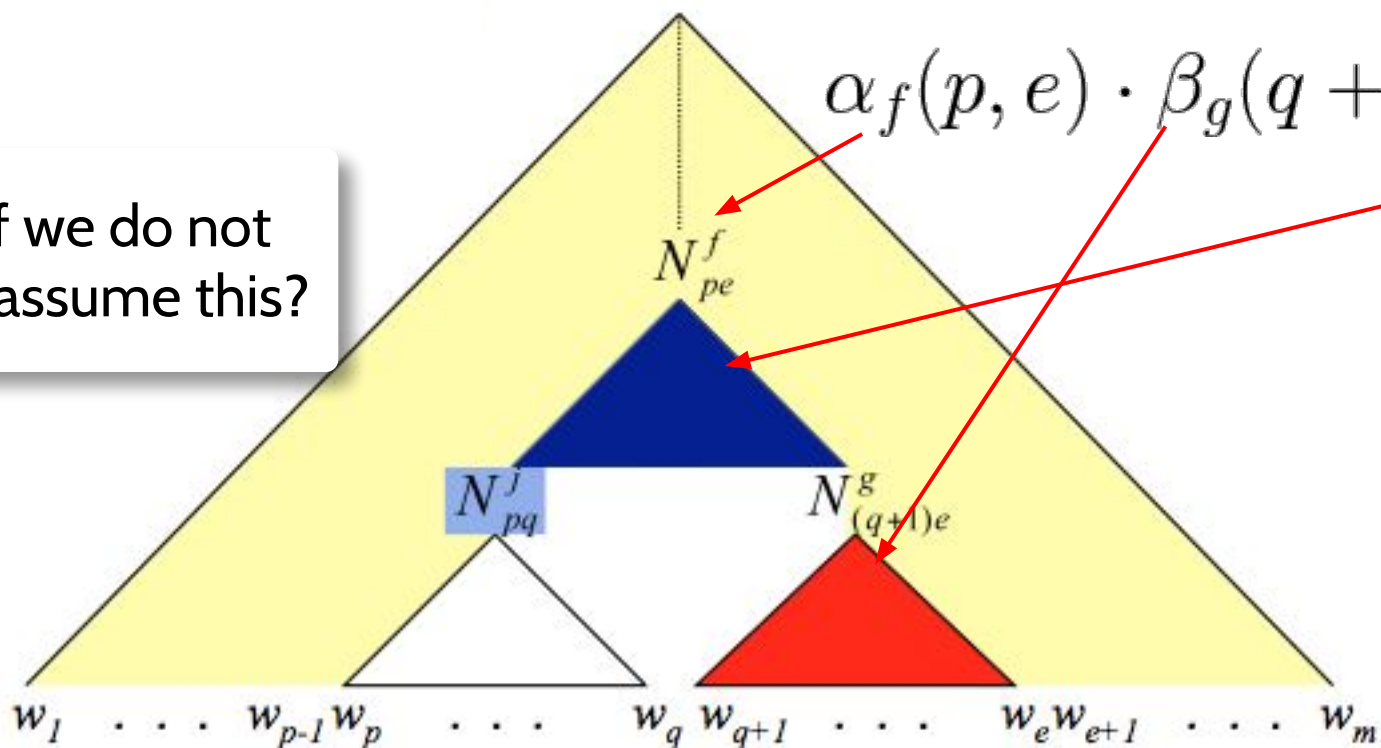


# Calculating outside probability

- ▶ If we multiply the terms together, we have the joint probability corresponding to the yellow, red and blue areas, **assuming**  $N^j$  was the L child of  $N^f$ , and give fixed non-terminals  $f$  and  $g$ , as well as a fixed partition  $e$

$$\alpha_f(p, e) \cdot \beta_g(q + 1, e) \cdot P(N^f \rightarrow N^j N^g)$$

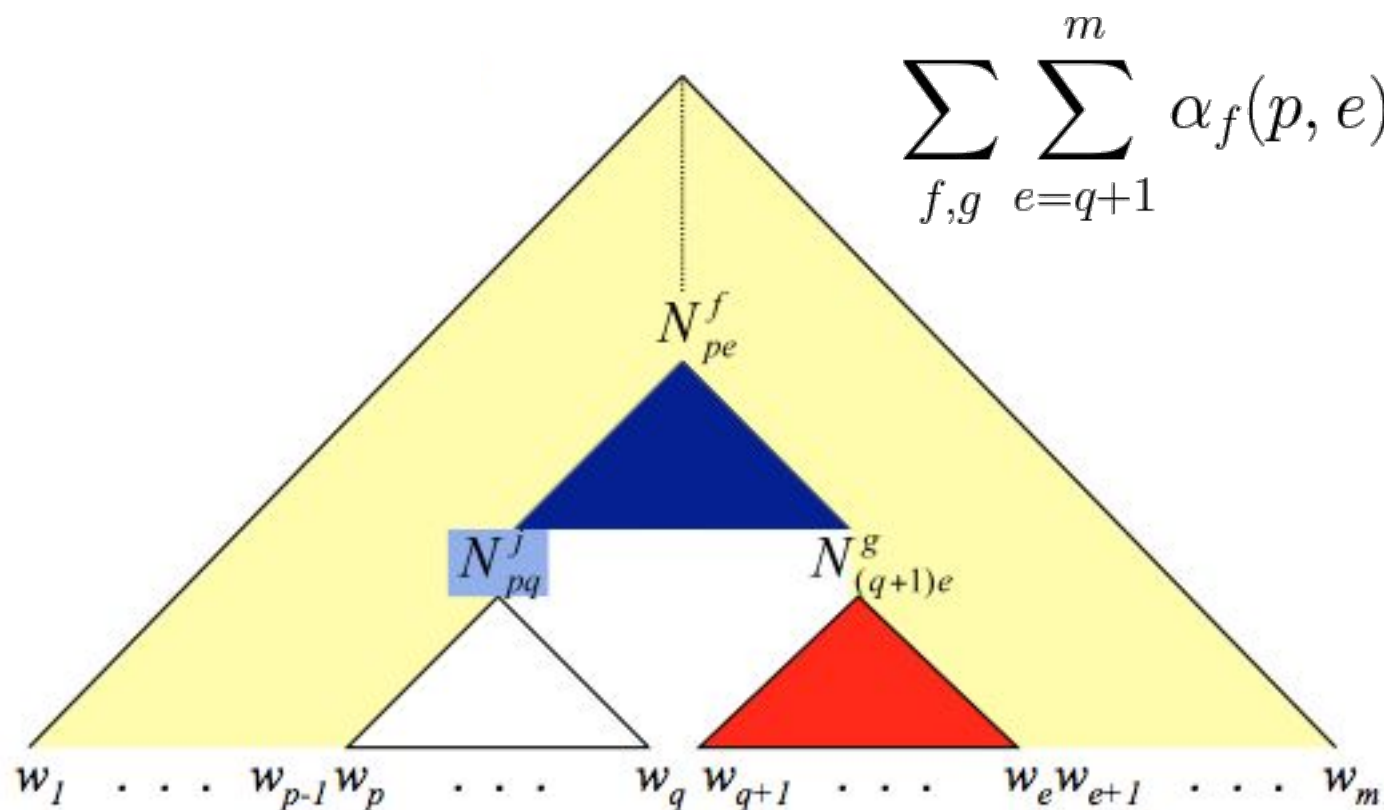
What if we do not want to assume this?





# Calculating outside probability

- ▶ The joint probability corresponding to the yellow, red and blue areas, **assuming**  $N^j$  was the **L** child of some non-terminal:



$$\sum_{f,g} \sum_{e=q+1}^m \alpha_f(p, e) \cdot \beta_g(q+1, e) \cdot P(N^f \rightarrow N^j N^g)$$

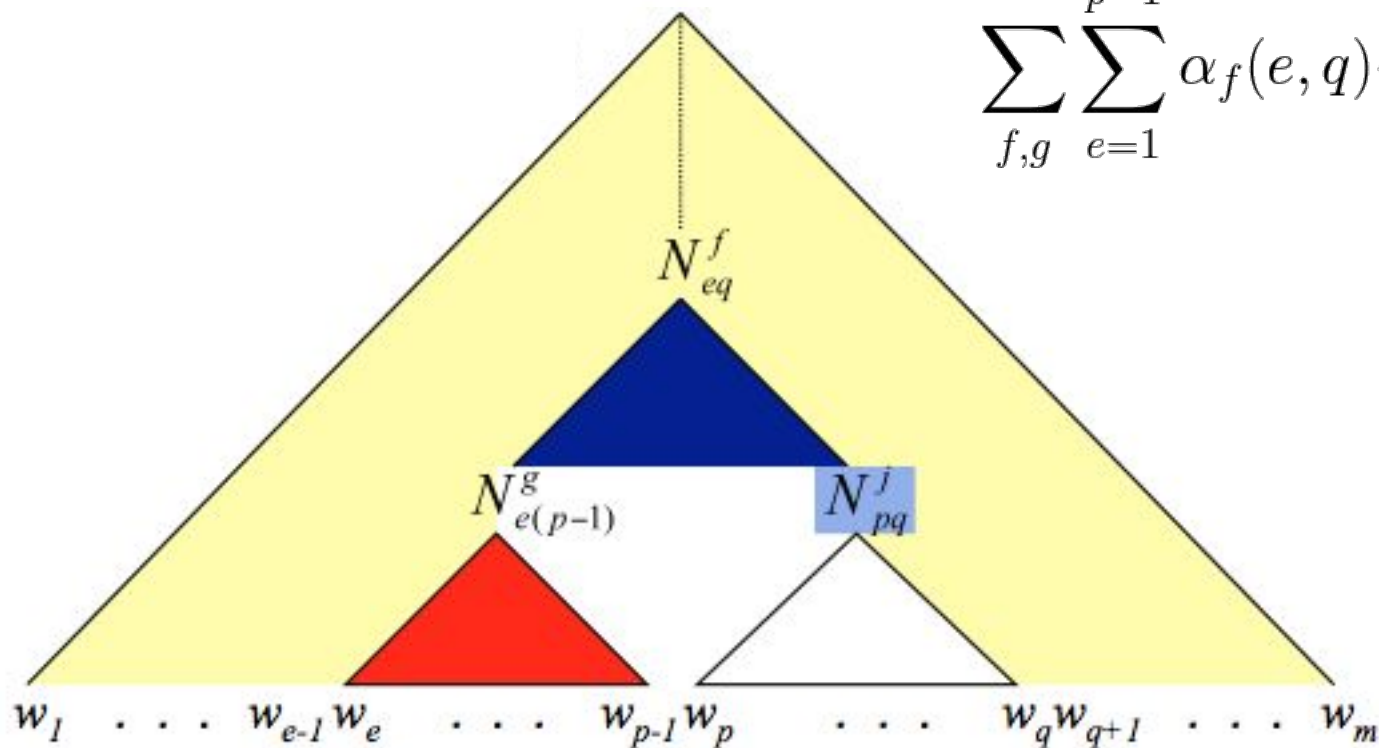




# Calculating outside probability

- ▶ The joint probability corresponding to the yellow, red and blue areas, **assuming**  $N^j$  was the **R** child of some non-terminal:

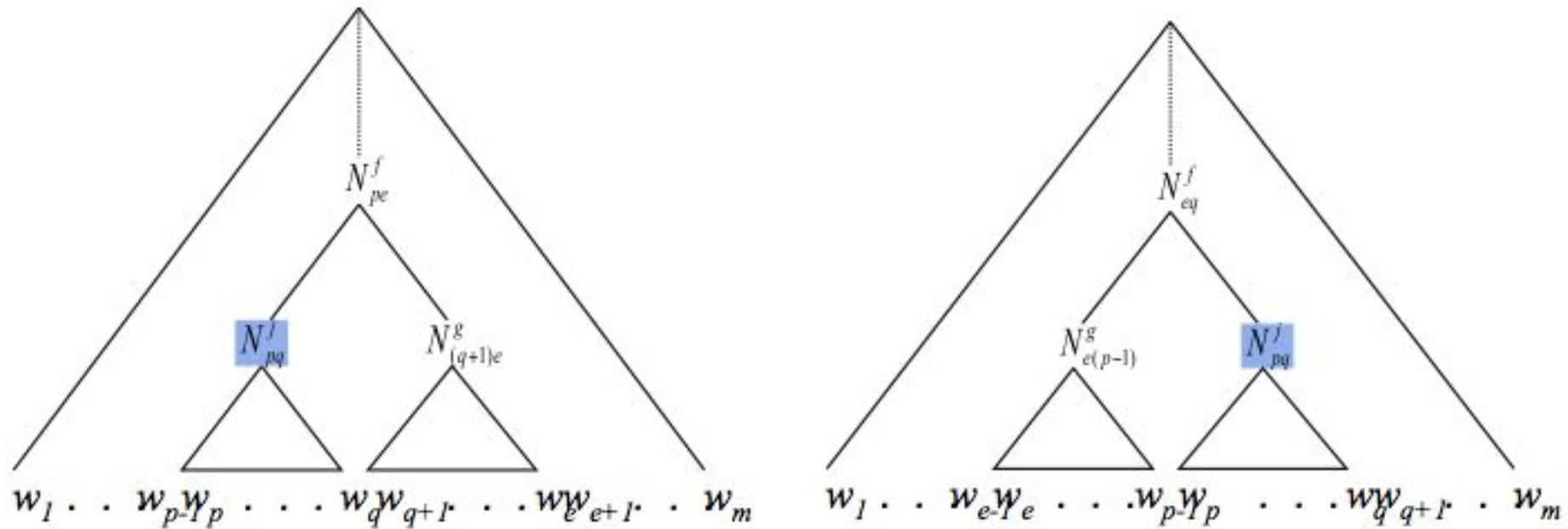
$$\sum_{f,g} \sum_{e=1}^{p-1} \alpha_f(e, q) \cdot \beta_g(e, p-1) \cdot P(N^f \rightarrow N^g N^j)$$





# Calculating outside probability

- ▶ The joint final joint probability (the sum over the L and R cases):

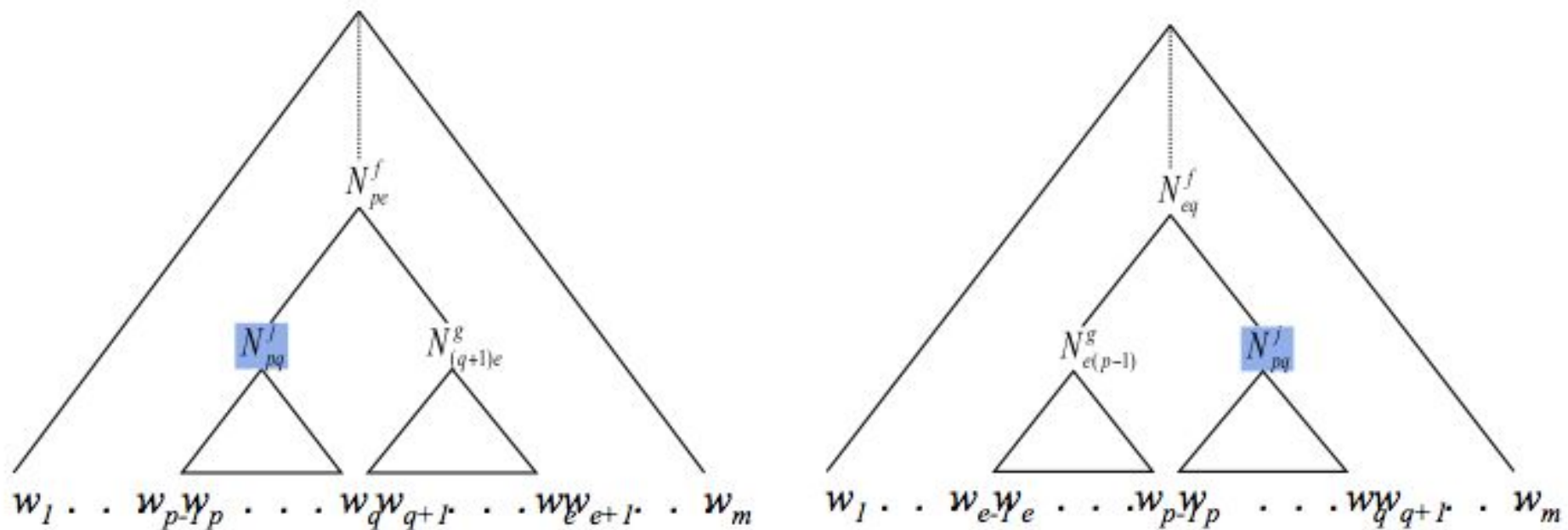


$$\alpha_j(p, q) = \sum_{f, g} \sum_{e=q+1}^m \alpha_f(p, e) \cdot \beta_g(q+1, e) \cdot P(N^f \rightarrow N^j N^g) + \sum_{f, g} \sum_{e=1}^{p-1} \alpha_f(e, q) \cdot \beta_g(e, p-1) \cdot P(N^f \rightarrow N^g N^j)$$



# Calculating outside probability

- ▶ The joint final joint probability (the sum over the L and R cases):



$$\alpha_j(p, q) = \sum_{\substack{f, g \neq j \\ \underline{f, g} \neq j}} \sum_{e=q+1}^m \alpha_f(p, e) \cdot \beta_g(q+1, e) \cdot P(N^f \rightarrow N^j N^g) + \sum_{f, g} \sum_{e=1}^{p-1} \alpha_f(e, q) \cdot \beta_g(e, p-1) \cdot P(N^f \rightarrow N^g N^j)$$



# Inside-outside algorithm

---

- ▶ For PCFGs we need to compute:

$$\theta^t = P(N^j \rightarrow N^r N^s | N^j)$$



# EM

---

- Given two events,  $x$  and  $y$ , the maximum likelihood estimation (MLE) for their conditional probability is:

$$P(x | y) = \frac{\textit{count}(x, y)}{\textit{count}(x)}$$

- If they are observable, it's easy to see what to do: just count the events in a representative corpus and use the MLE



# EM

- What these are hidden variables that cannot be observed directly?
- Use a model  $\mu$  and iteratively improve the model based on a corpus of observable data ( $O$ ) generated by the hidden variables:

$$P_{\hat{\mu}}(x | y) = \frac{E_{\mu}[\text{count}(x, y) | O]}{E_{\mu}[\text{count}(x) | O]}$$

- It is worth noting that if you know how to calculate the numerator, the denominator is trivially derivable.



# EM

---

- By updating  $\mu$  and iterating, the model converges to at least a local maximum
- This can be proven, but I will not do it here.



# The inside-outside algorithm

---

- Goal: estimate a model  $\mu$  that is a PCFG (in Chomsky normal form) that characterizes a corpus of text.
  
- Required input:
  - Size of non-terminal vocabulary,  $n$
  - At least one sentence to be modeled,  $O$





# The inside-outside algorithm

---

- Stated with the general schema described earlier, we seek to the MLE probabilities for productions in the grammar

$$P(N^j \rightarrow N^r N^s \mid N^j) = \frac{\text{count}(N^j \rightarrow N^r N^s, N^j)}{\text{count}(N^j)}$$

- (Observe that this would be trivially easy to calculate this with a treebank, since the non-terminals are observable in a treebank)



# The inside-outside algorithm

---

- Since the non-terminals are not visible, we can use EM to estimate the probabilities iteratively:

$$P_{\hat{\mu}}(N^j \rightarrow N^r N^s \mid N^j) = \frac{E_{\mu}[\text{count}(N^j \rightarrow N^r N^s, N^j) \mid O]}{E_{\mu}[\text{count}(N^j) \mid O]}$$



# To be continued...

---

- Next: recitation on EM