

CS447: Natural Language Processing

<http://courses.engr.illinois.edu/cs447>

# Lecture 16: PCFG Parsing (updated)

Julia Hockenmaier

*juliahmr@illinois.edu*

3324 Siebel Center

# Overview

# Where we're at

## Previous lecture:

### Standard **CKY** (for non-probabilistic CFGs)

The CKY algorithm finds all possible parse trees  $\tau$  for a sentence  $S = w^{(1)} \dots w^{(n)}$  under a CFG  $G$  in Chomsky Normal Form.

## Today's lecture:

### **Probabilistic Context-Free Grammars (PCFGs)**

- CFGs in which each rule is associated with a probability

### **CKY for PCFGs (Viterbi):**

- CKY for PCFGs finds the most likely parse tree  
 $\tau^* = \operatorname{argmax} P(\tau \mid S)$  for the sentence  $S$  under a PCFG.

### **Shortcomings of PCFGs (and ways to overcome them)**

### **Penn Treebank Parsing**

### **Evaluating PCFG parsers**

# CKY: filling the chart

w	...	...	...	$w_i$	...	w	
							w
							...
							..
							$w_i$
							...
							w

w	...	...	...	$w_i$	...	w	
							w
							...
							..
							$w_i$
							...
							w

w	...	...	...	$w_i$	...	w	
							w
							...
							..
							$w_i$
							...
							w

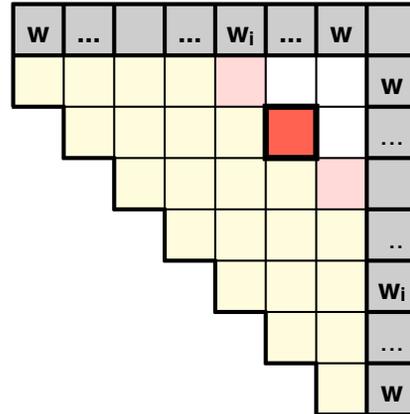
w	...	...	...	$w_i$	...	w	
							w
							...
							..
							$w_i$
							...
							w

w	...	...	...	$w_i$	...	w	
							w
							...
							..
							$w_i$
							...
							w

w	...	...	...	$w_i$	...	w	
							w
							...
							..
							$w_i$
							...
							w

w	...	...	...	$w_i$	...	w	
							w
							...
							..
							$w_i$
							...
							w

# CKY: filling one cell

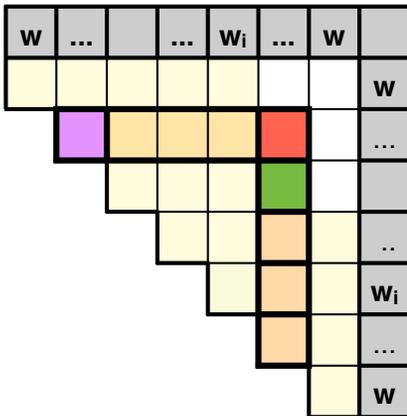


chart[2][6]:

$w_1$   **$w_2$**   **$w_3$**   **$w_4$**   **$w_5$**   **$w_6$**   $w_7$

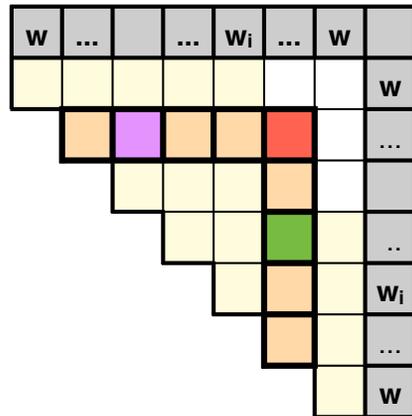
chart[2][6]:

$w_1$   **$w_2$**  **$w_3$**  **$w_4$**  **$w_5$**  **$w_6$**   $w_7$



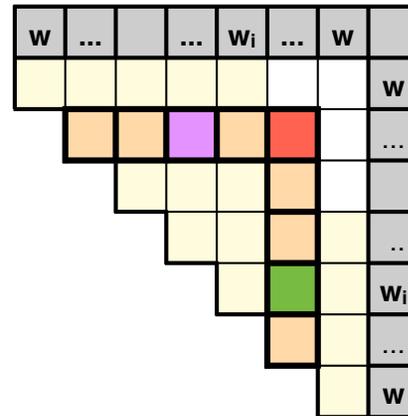
chart[2][6]:

$w_1$   **$w_2$**  **$w_3$**  **$w_4$**  **$w_5$**  **$w_6$**   $w_7$



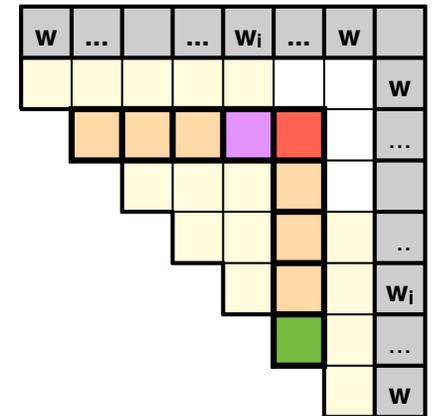
chart[2][6]:

$w_1$   **$w_2$**  **$w_3$**  **$w_4$**  **$w_5$**  **$w_6$**   $w_7$



chart[2][6]:

$w_1$   **$w_2$**  **$w_3$**  **$w_4$**  **$w_5$**  **$w_6$**   $w_7$



# CKY for standard CFGs

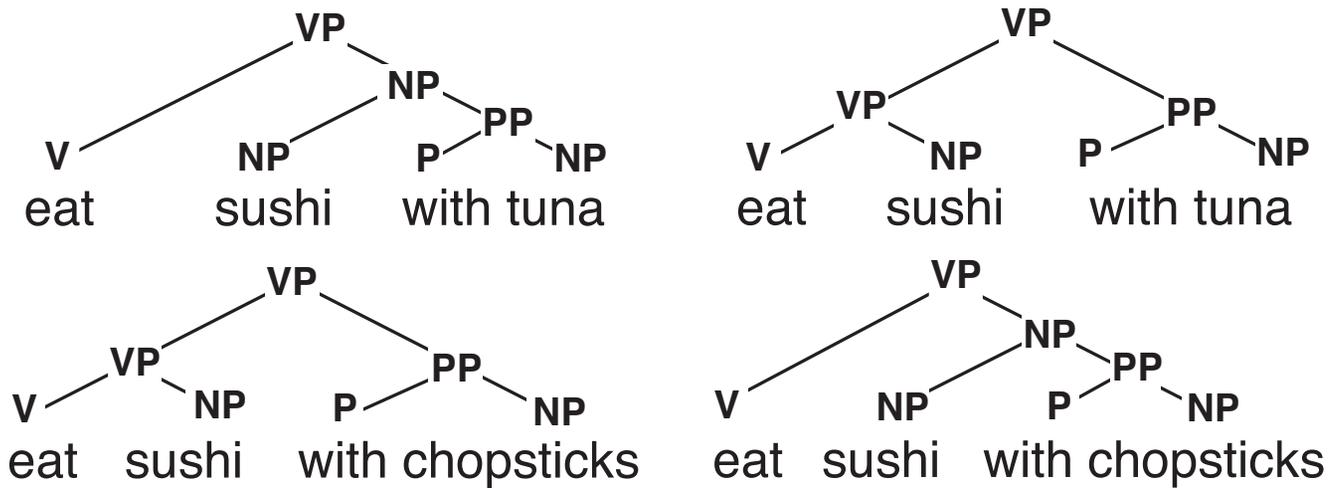
CKY is a bottom-up chart parsing algorithm that finds all possible parse trees  $\tau$  for a sentence  $S = w^{(1)} \dots w^{(n)}$  under a CFG  $G$  in Chomsky Normal Form (CNF).

- **CNF**:  $G$  has two types of rules:  $X \longrightarrow Y Z$  and  $X \longrightarrow w$  ( $X, Y, Z$  are nonterminals,  $w$  is a terminal)
- CKY is a **dynamic programming** algorithm
- The **parse chart** is an  $n \times n$  upper triangular matrix:  
Each cell  $\text{chart}[i][j]$  ( $i \leq j$ ) stores **all subtrees** for  $w^{(i)} \dots w^{(j)}$
- Each cell  $\text{chart}[i][j]$  has at most **one entry for each nonterminal  $X$**  (and **pairs of backpointers** to each pair of  $(Y, Z)$  entry in cells  $\text{chart}[i][k]$   $\text{chart}[k+1][j]$  from which an  $X$  can be formed
- Time Complexity:  $O(n^3 |G|)$

Probabilistic  
Context-Free  
Grammars (PCFGs)

# Grammars are ambiguous

A grammar might generate multiple trees for a sentence:



What's the most likely parse  $\tau$  for sentence  $S$  ?

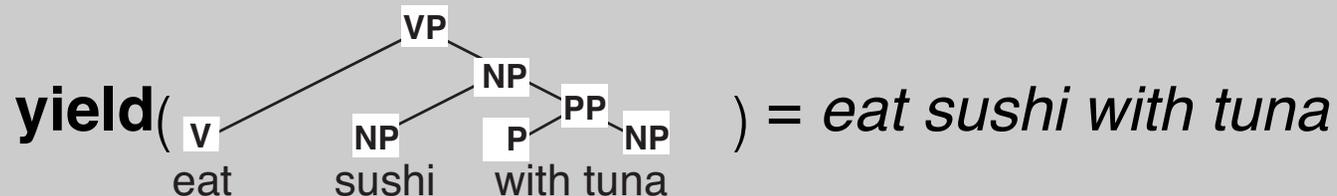
**We need a model of  $P(\tau | S)$**

# Computing $P(\tau | S)$

Using Bayes' Rule:

$$\begin{aligned} \arg \max_{\tau} P(\tau | S) &= \arg \max_{\tau} \frac{P(\tau, S)}{P(S)} \\ &= \arg \max_{\tau} P(\tau, S) \\ &= \arg \max_{\tau} P(\tau) \quad \text{if } S = \text{yield}(\tau) \end{aligned}$$

The **yield of a tree** is the string of terminal symbols that can be read off the leaf nodes



# Computing $P(\tau)$

$T$  is the (infinite) set of all trees in the language:

$$L = \{s \in \Sigma^* \mid \exists \tau \in T : \text{yield}(\tau) = s\}$$

We need to define  $P(\tau)$  such that:

$$\forall \tau \in T : 0 \leq P(\tau) \leq 1$$

$$\sum_{\tau \in T} P(\tau) = 1$$

The set  $T$  is generated by a context-free grammar

$S \rightarrow NP VP$	$VP \rightarrow Verb NP$	$NP \rightarrow Det Noun$
$S \rightarrow S conj S$	$VP \rightarrow VP PP$	$NP \rightarrow NP PP$
$S \rightarrow \dots\dots$	$VP \rightarrow \dots\dots$	$NP \rightarrow \dots\dots$

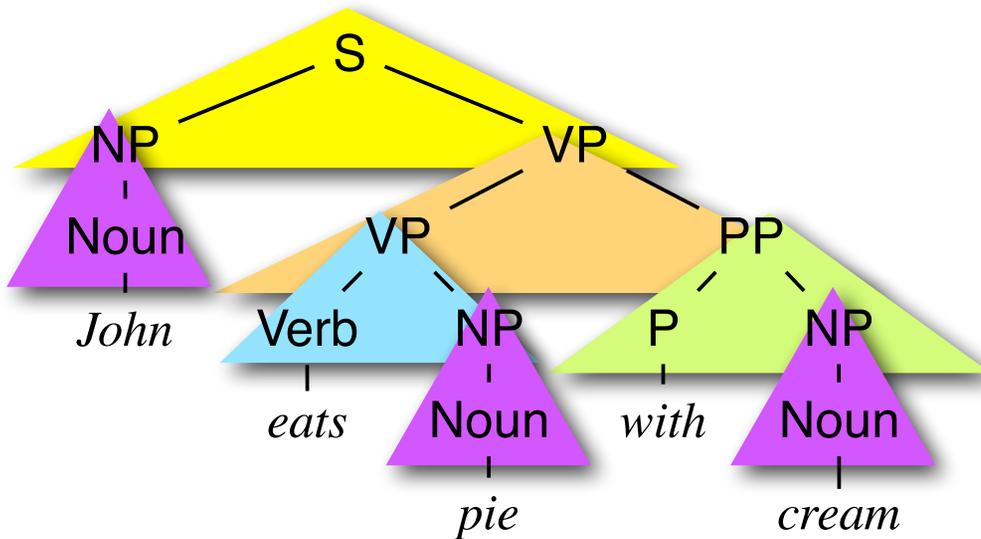
# Probabilistic Context-Free Grammars

For every nonterminal  $X$ , define a probability distribution  $P(X \rightarrow \alpha \mid X)$  over all rules with the same LHS symbol  $X$ :

S	→ NP VP	0.8
S	→ S conj S	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP conj NP	0.2
VP	→ Verb	0.4
VP	→ Verb NP	0.3
VP	→ Verb NP NP	0.1
VP	→ VP PP	0.2
PP	→ P NP	1.0

# Computing $P(\tau)$ with a PCFG

The probability of a tree  $\tau$  is the product of the probabilities of all its rules:



S	→ NP VP	0.8
S	→ S conj S	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP conj NP	0.2
VP	→ Verb	0.4
VP	→ Verb NP	0.3
VP	→ Verb NP NP	0.1
VP	→ VP PP	0.2
PP	→ P NP	1.0

$$P(\tau) = 0.8 \times 0.3 \times 0.2 \times 1.0 \times 0.2^3$$

$$= 0.00384$$

# Learning the parameters of a PCFG

If we have a treebank (a corpus in which each sentence is associated with a parse tree), we can just **count** the number of times each rule appears, e.g.:

$S \rightarrow NP VP .$  (count = 1000)

$S \rightarrow S \text{ conj } S .$  (count = 220)

$PP \rightarrow IN NP$  (count = 700)

and then we divide the **count** (observed frequency) of **each rule**  $X \rightarrow Y Z$  by the **sum of the frequencies of all rules with the same LHS**  $X$  to turn these counts into probabilities:

$S \rightarrow NP VP .$  ( $p = 1000/1220$ )

$S \rightarrow S \text{ conj } S .$  ( $p = 220/1220$ )

$PP \rightarrow IN NP$  ( $p = 700/700$ )

# More on probabilities:

## Computing $P(s)$ :

If  $P(\tau)$  is the probability of a tree  $\tau$ , the probability of a sentence  $s$  is the sum of the probabilities of all its parse trees:

$$P(s) = \sum_{\tau: \text{yield}(\tau) = s} P(\tau)$$

## How do we know that $P(L) = \sum_{\tau} P(\tau) = 1$ ?

If we have learned the PCFG from a corpus via MLE, this is guaranteed to be the case.

But if we set the probabilities by hand, we could run into trouble: In this PCFG, the probability mass of all finite trees is less than 1:

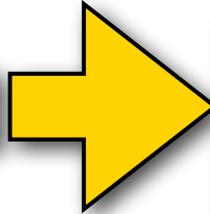
$$S \rightarrow S S \quad (0.9) \qquad S \rightarrow w \quad (0.1)$$

$$\begin{aligned} P(L) &= P("w") + P("ww") + P("w[ww]") + P("[ww]w") + \dots \\ &= .1 + .009 + 0.00081 + 0.00081 + \dots \ll 1 \end{aligned}$$

PCFG Decoding:  
CKY with Viterbi

# How do we handle flat rules?

S	→	NP VP	0.8
<b>S</b>	<b>→</b>	<b>S conj S</b>	<b>0.2</b>
NP	→	Noun	0.2
NP	→	Det Noun	0.4
NP	→	NP PP	0.2
<b>NP</b>	<b>→</b>	<b>NP conj NP</b>	<b>0.2</b>
VP	→	Verb	0.3
VP	→	Verb NP	0.3
<b>VP</b>	<b>→</b>	<b>Verb NP NP</b>	<b>0.1</b>
VP	→	VP PP	0.3
PP	→	PP NP	1.0
Prep	→	P	1.0
Noun	→	N	1.0
Verb	→	V	1.0



<b>S</b>	<b>→</b>	<b>S ConjS</b>	<b>0.2</b>
<b>ConjS</b>	<b>→</b>	<b>conj S</b>	<b>1.0</b>

**Binarize** each flat rule by adding a **unique dummy nonterminal** (ConjS), and **setting the probability** of the new rule with the dummy nonterminal on the LHS to **1**

# How do we handle flat rules?

S	→	NP VP	0.8
<b>S</b>	<b>→</b>	<b>S conj S</b>	<b>0.2</b>
NP	→	Noun	0.2
NP	→	Det Noun	0.4
NP	→	NP PP	0.2
<b>NP</b>	<b>→</b>	<b>NP conj NP</b>	<b>0.2</b>
VP	→	Verb	0.3
VP	→	Verb NP	0.3
<b>VP</b>	<b>→</b>	<b>Verb NP NP</b>	<b>0.1</b>
VP	→	VP PP	0.3
PP	→	PP NP	1.0
Prep	→	P	1.0
Noun	→	N	1.0
Verb	→	V	1.0

S	→	NP VP	0.8
<b>S</b>	<b>→</b>	<b>S ConjS</b>	<b>0.2</b>
NP	→	Noun	0.2
NP	→	Det Noun	0.4
NP	→	NP PP	0.2
<b>NP</b>	<b>→</b>	<b>NP ConjNP</b>	<b>0.2</b>
VP	→	Verb	0.3
VP	→	Verb NP	0.3
<b>VP</b>	<b>→</b>	<b>Verb NPNP</b>	<b>0.1</b>
VP	→	VP PP	0.3
PP	→	PP NP	1.0
Prep	→	P	1.0
Noun	→	N	1.0
Verb	→	V	1.0
<b>ConjS</b>	<b>→</b>	<b>conj S</b>	<b>1.0</b>
<b>ConjNP</b>	<b>→</b>	<b>conj NP</b>	<b>1.0</b>
<b>NPNP</b>	<b>→</b>	<b>NP NP</b>	<b>1.0</b>

# Probabilistic CKY: Viterbi

Like standard CKY, but with probabilities.

Finding the most likely tree is similar to Viterbi for HMMs:

## Initialization:

- [optional] Every chart entry that corresponds to a **terminal** (entry  $w$  in  $cell[i][i]$ ) has a Viterbi probability  $P_{VIT}(w_{[i][i]}) = 1$  (\*)
- Every entry for a **non-terminal**  $X$  in  $cell[i][i]$  has Viterbi probability  $P_{VIT}(X_{[i][i]}) = P(X \rightarrow w \mid X)$  [and a single backpointer to  $w_{[i][i]}$  (\*)]

**Recurrence:** For every entry that corresponds to a **non-terminal**  $x$  in  $cell[i][j]$ , keep only the highest-scoring pair of backpointers to any pair of children ( $Y$  in  $cell[i][k]$  and  $Z$  in  $cell[k+1][j]$ ):  
$$P_{VIT}(X_{[i][j]}) = \operatorname{argmax}_{Y,Z,k} P_{VIT}(Y_{[i][k]}) \times P_{VIT}(Z_{[k+1][j]}) \times P(X \rightarrow Y Z \mid X)$$

**Final step:** Return the Viterbi parse for the start symbol  $S$  in the top  $cell[1][n]$ .

\*this is unnecessary for simple PCFGs, but can be helpful for more complex probability models

# Probabilistic CKY

## Input: POS-tagged sentence

John\_N eats\_V pie\_N with\_P cream\_N

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
					<b>John</b>
					<b>eats</b>
					<b>pie</b>
					<b>with</b>
					<b>cream</b>

S	→	NP VP	0.8
S	→	S Conjs	0.2
NP	→	Noun	0.2
NP	→	Det Noun	0.4
NP	→	NP PP	0.2
NP	→	NP ConjNP	0.2
VP	→	Verb	0.3
VP	→	Verb NP	0.3
VP	→	Verb NPNP	0.1
VP	→	VP PP	0.3
PP	→	PP NP	1.0
Prep	→	P	1.0
Noun	→	N	1.0
Verb	→	V	1.0
Conjs	→	conj S	1.0
ConjNP	→	conj NP	1.0
NPNP	→	NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John\_N eats\_V pie\_N with\_P cream\_N

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
					<b>John</b>
					<b>eats</b>
					<b>pie</b>
					<b>with</b>
					<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0



# Probabilistic CKY

## Input: POS-tagged sentence

John\_N eats\_V pie\_N with\_P cream\_N

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2					<b>John</b>
					<b>eats</b>
					<b>pie</b>
					<b>with</b>
					<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0



# Probabilistic CKY

## Input: POS-tagged sentence

John\_N eats\_V pie\_N with\_P cream\_N

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2					<b>John</b>
	Verb VP 1.0 0.3				<b>eats</b>
					<b>pie</b>
					<b>with</b>
					<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
<b>VP</b>	<b>→ Verb</b>	<b>0.3</b>
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
<b>Verb</b>	<b>→ V</b>	<b>1.0</b>
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John\_N eats\_V pie\_N with\_P cream\_N

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2					<b>John</b>
	Verb VP 1.0 0.3				<b>eats</b>
		Noun NP 1.0 0.2			<b>pie</b>
					<b>with</b>
					<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0



# Probabilistic CKY

## Input: POS-tagged sentence

John\_N eats\_V pie\_N with\_P cream\_N

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2					<b>John</b>
	Verb VP 1.0 0.3				<b>eats</b>
		Noun NP 1.0 0.2			<b>pie</b>
			Prep 1.0		<b>with</b>
					<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John\_N eats\_V pie\_N with\_P cream\_N

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2					<b>John</b>
	Verb VP 1.0 0.3				<b>eats</b>
		Noun NP 1.0 0.2			<b>pie</b>
			Prep 1.0		<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0



# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun <b>NP</b> 1.0 0.2	<b>S</b> 0.8 · 0.2 · 0.3				<b>John</b>
	Verb <b>VP</b> 1.0 0.3				<b>eats</b>
		Noun <b>NP</b> 1.0 0.2			<b>pie</b>
			Prep 1.0		<b>with</b>
				Noun <b>NP</b> 1.0 0.2	<b>cream</b>

<b>S</b>	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3				<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06			<b>eats</b>
		Noun NP 1.0 0.2			<b>pie</b>
			Prep 1.0		<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John\_N eats\_V pie\_N with\_P cream\_N

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3				<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06			<b>eats</b>
		Noun NP 1.0 0.2			<b>pie</b>
			Prep 1.0		<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3				<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06			<b>eats</b>
		Noun NP 1.0 0.2			<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

**Input: POS-tagged sentence**

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun <b>NP</b> 1.0 0.2	S 0.8 · 0.2 · 0.3	<b>S</b> 0.8 · 0.2 · 0.06			<b>John</b>
	Verb <b>VP</b> 1.0 0.3	<b>VP</b> 1 · 0.3 · 0.2 = 0.06			<b>eats</b>
		Noun <b>NP</b> 1.0 0.2			<b>pie</b>
			Prep 1.0	<b>PP</b> 1 · 1 · 0.2	<b>with</b>
				Noun <b>NP</b> 1.0 0.2	<b>cream</b>

<b>S</b>	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06			<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06			<b>eats</b>
		Noun NP 1.0 0.2			<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06			<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06			<b>eats</b>
		Noun NP 1.0 0.2			<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0



# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06			<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06			<b>eats</b>
		Noun NP 1.0 0.2			<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06			<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06			<b>eats</b>
		Noun NP 1.0 0.2		NP 0.2 · 0.2 · 0.2 = 0.008	<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06			<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06			<b>eats</b>
		Noun NP 1.0 0.2		NP 0.2 · 0.2 · 0.2 = 0.008	<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06			<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06			<b>eats</b>
		Noun NP 1.0 0.2		NP 0.2 · 0.2 · 0.2 = 0.008	<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John\_N eats\_V pie\_N with\_P cream\_N

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06			<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06			<b>eats</b>
		Noun NP 1.0 0.2		NP 0.2 · 0.2 · 0.2 = 0.008	<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

John	eats	pie	with	cream	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06			John
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06			eats
		Noun NP 1.0 0.2		NP 0.2 · 0.2 · 0.2 = 0.008	pie
			Prep 1.0	PP 1 · 1 · 0.2	with
				Noun NP 1.0 0.2	cream

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06			<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06		VP max( 1.0 · 0.008 · 0.3, 0.06 · 0.2 · 0.3 )	<b>eats</b>
		Noun NP 1.0 0.2		NP 0.2 · 0.2 · 0.2 = 0.008	<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun <b>NP</b> 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06		<b>S</b> 0.2 · 0.0036 · 0.8	<b>John</b>
	Verb <b>VP</b> 1.0 0.3	<b>VP</b> 1 · 0.3 · 0.2 = 0.06		<b>VP</b> 0.0036	<b>eats</b>
		Noun <b>NP</b> 1.0 0.2		<b>NP</b> 0.2 · 0.2 · 0.2 = 0.008	<b>pie</b>
			Prep 1.0	<b>PP</b> 1 · 1 · 0.2	<b>with</b>
				Noun <b>NP</b> 1.0 0.2	<b>cream</b>

<b>S</b>	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0



# Probabilistic CKY

## Input: POS-tagged sentence

John\_N eats\_V pie\_N with\_P cream\_N

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06		S 0.2 · 0.0036 · 0.8	<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06		VP 0.0036	<b>eats</b>
		Noun NP 1.0 0.2		NP 0.2 · 0.2 · 0.2 = 0.008	<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

**Input: POS-tagged sentence**

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06		S 0.2 · 0.0036 · 0.8	<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06		VP 0.0036	<b>eats</b>
		Noun NP 1.0 0.2		NP 0.2 · 0.2 · 0.2 = 0.008	<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Probabilistic CKY

## Input: POS-tagged sentence

John<sub>N</sub> eats<sub>V</sub> pie<sub>N</sub> with<sub>P</sub> cream<sub>N</sub>

<b>John</b>	<b>eats</b>	<b>pie</b>	<b>with</b>	<b>cream</b>	
Noun NP 1.0 0.2	S 0.8 · 0.2 · 0.3	S 0.8 · 0.2 · 0.06		S 0.2 · 0.0036 · 0.8	<b>John</b>
	Verb VP 1.0 0.3	VP 1 · 0.3 · 0.2 = 0.06		VP 0.0036	<b>eats</b>
		Noun NP 1.0 0.2		NP 0.2 · 0.2 · 0.2 = 0.008	<b>pie</b>
			Prep 1.0	PP 1 · 1 · 0.2	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

S	→ NP VP	0.8
S	→ S Conjs	0.2
NP	→ Noun	0.2
NP	→ Det Noun	0.4
NP	→ NP PP	0.2
NP	→ NP ConjNP	0.2
VP	→ Verb	0.3
VP	→ Verb NP	0.3
VP	→ Verb NPNP	0.1
VP	→ VP PP	0.3
PP	→ PP NP	1.0
Prep	→ P	1.0
Noun	→ N	1.0
Verb	→ V	1.0
Conjs	→ conj S	1.0
ConjNP	→ conj NP	1.0
NPNP	→ NP NP	1.0

# Extracting the final tree

## Input: POS-tagged sentence

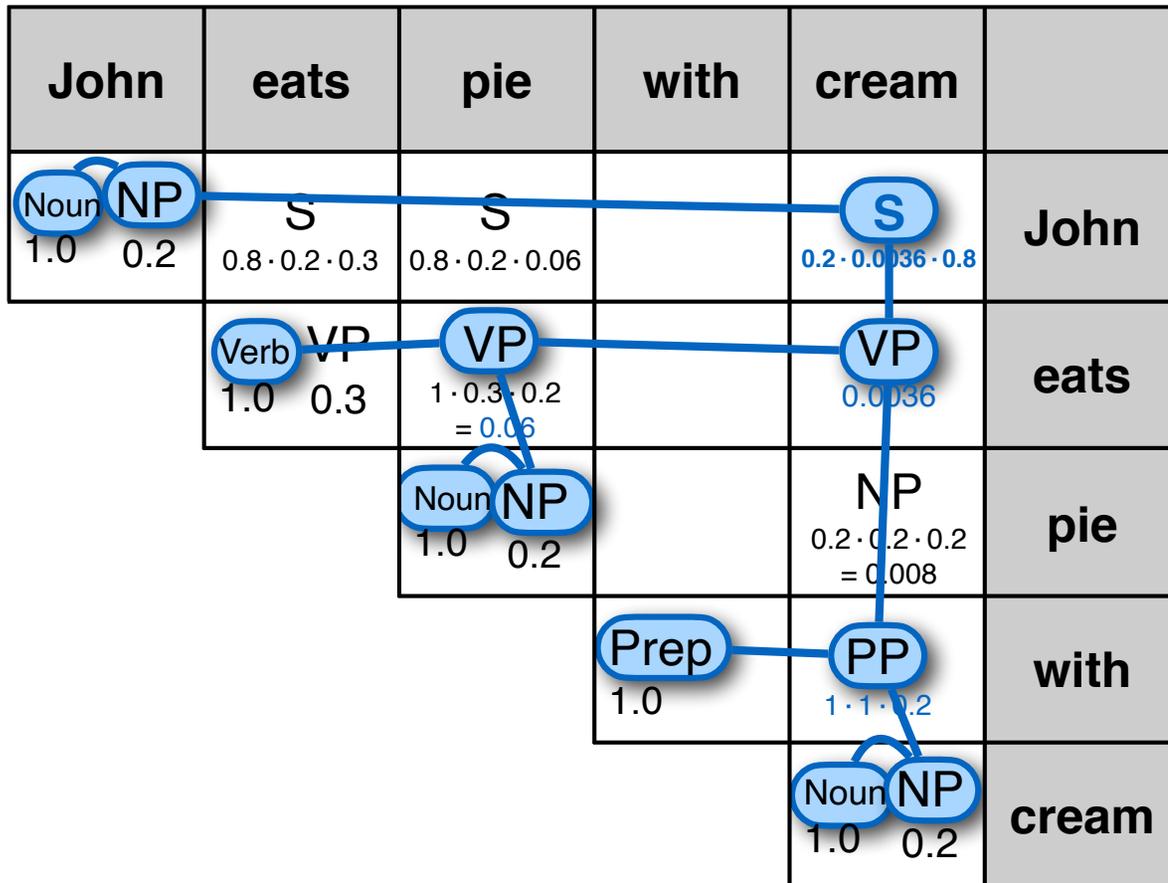
John\_N eats\_V pie\_N with\_P cream\_N

John	eats	pie	with	cream	
Noun NP 1.0 0.2	S $0.8 \cdot 0.2 \cdot 0.3$	S $0.8 \cdot 0.2 \cdot 0.06$		<b>S</b> $0.2 \cdot 0.0036 \cdot 0.8$	<b>John</b>
	Verb VP 1.0 0.3	VP $1 \cdot 0.3 \cdot 0.2$ = 0.06		VP 0.0036	<b>eats</b>
		Noun NP 1.0 0.2		NP $0.2 \cdot 0.2 \cdot 0.2$ = 0.008	<b>pie</b>
			Prep 1.0	PP $1 \cdot 1 \cdot 0.2$	<b>with</b>
				Noun NP 1.0 0.2	<b>cream</b>

# Extracting the final tree

## Input: POS-tagged sentence

John\_N eats\_V pie\_N with\_P cream\_N



```
(S (NP (N John))
  (VP (VP (Verb eats)
          (NP (Noun pie))))
  (PP (Prep with)
      (NP (Noun cream)))))
```

Shortcomings  
of PCFGs

# How well can a PCFG model the distribution of trees?

PCFGs make **independence assumptions**:

Only the label of a node determines what children it has.

Factors that influence these assumptions:

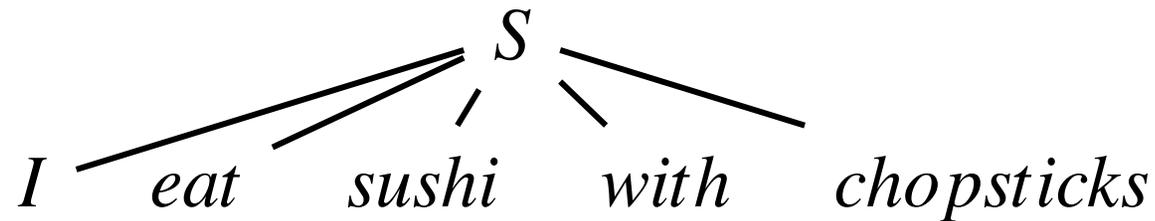
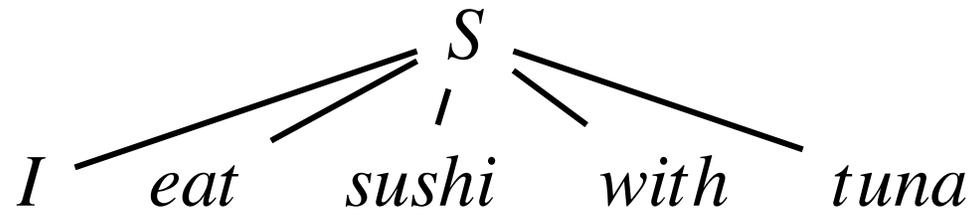
**Shape** of the trees:

A corpus with **flat trees** (i.e. few nodes/sentence) results in a model with few independence assumptions.

**Labeling** of the trees:

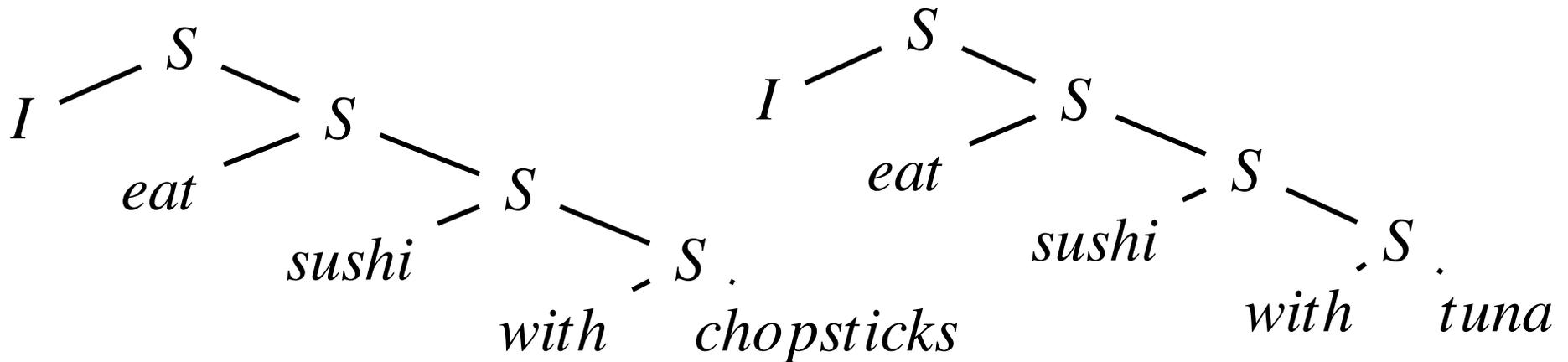
A corpus with **many node labels** (nonterminals) results in a model with few independence assumptions.

# Example 1: flat trees



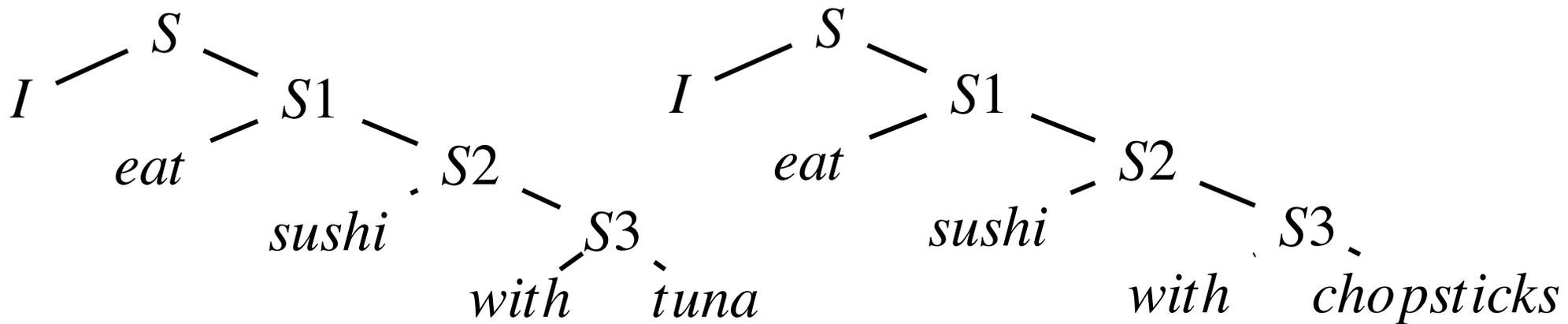
**What sentences would a PCFG estimated from this corpus generate?**

# Example 2: deep trees, few labels



**What sentences would a PCFG estimated from this corpus generate?**

# Example 3: deep trees, many labels



**What sentences would a PCFG estimated from this corpus generate?**

# Aside: Bias/Variance tradeoff

A probability model has **low bias** if it makes **few independence assumptions**.

⇒ It can capture the structures in the training data.

But: this typically leads to a **more fine-grained partitioning** of the training data.

Hence, **fewer data points are available** to estimate the model parameters.

This yields poor estimates of the distribution.

That is, such models have **high variance**



# Two ways to improve performance

## ... change the (internal) grammar:

### - Parent annotation/state splits:

Not all NPs/VPs/DTs/... are the same.

It matters where they are in the tree

## ... change the probability model:

### – Lexicalization:

Words matter!

### – Markovization:

Generalizing the rules

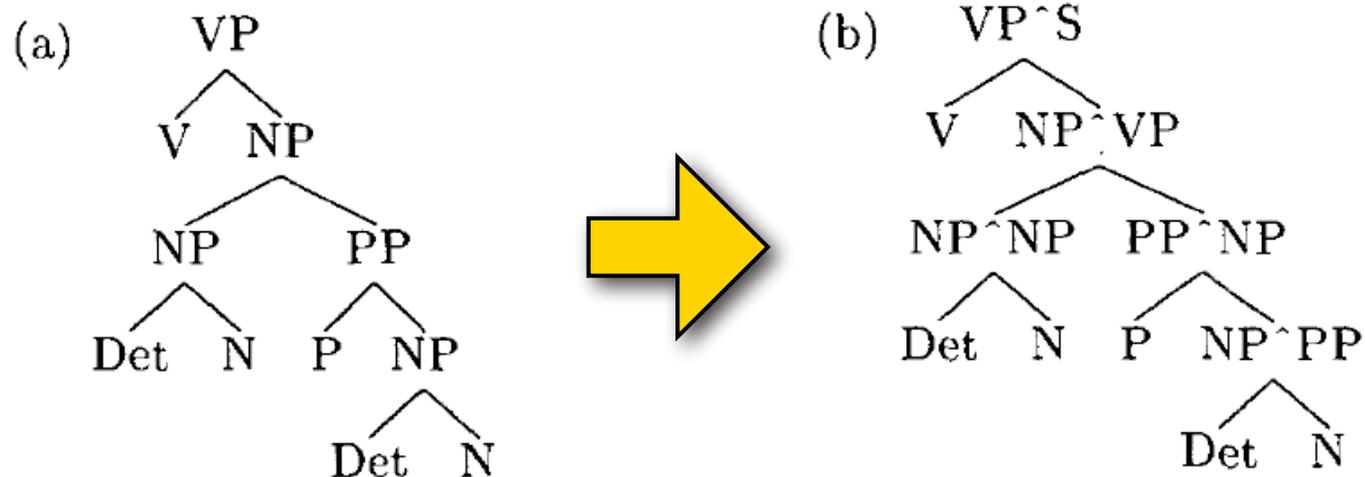


# The parent transformation

PCFGs assume the expansion of any nonterminal is independent of its parent.

But this is not true: NP subjects more likely to be modified than objects.

We can **change the grammar** by adding the name of the parent node to each nonterminal



# Lexicalization

PCFGs can't distinguish between  
*“eat sushi with chopsticks”* and *“eat sushi with tuna”*.

We need to take words into account!

$P(\text{VP}_{\text{eat}} \rightarrow \text{VP PP}_{\text{with chopsticks}} \mid \text{VP}_{\text{eat}})$

vs.  $P(\text{VP}_{\text{eat}} \rightarrow \text{VP PP}_{\text{with tuna}} \mid \text{VP}_{\text{eat}})$

Problem: sparse data ( $\text{PP}_{\text{with fatty|white|... tuna....}}$ )

Solution: only take **head words** into account!

Assumption: each constituent has one head word.

# Lexicalized PCFGs

At the root (start symbol  $S$ ), generate the head word of the sentence,  $w_s$ , with  $P(w_s)$

## Lexicalized rule probabilities:

Every nonterminal is lexicalized:  $X_{w_x}$

Condition rules  $X_{w_x} \rightarrow \alpha Y \beta$  on the lexicalized LHS  $X_{w_x}$

$P(X_{w_x} \rightarrow \alpha Y \beta \mid X_{w_x})$

## Word-word dependencies:

For each nonterminal  $Y$  in RHS of a rule  $X_{w_x} \rightarrow \alpha Y \beta$ , condition  $w_Y$  (the head word of  $Y$ ) on  $X$  and  $w_x$ :

$P(w_Y \mid Y, X, w_x)$

# Dealing with unknown words

A lexicalized PCFG assigns zero probability to any word that does not appear in the training data.

## Solution:

Training: Replace rare words in training data with a token 'UNKNOWN'.

Testing: Replace unseen words with 'UNKNOWN'

# Markov PCFGs (Collins parser)

The RHS of each CFG rule consists of:  
one head  $H_X$ ,  $n$  left sisters  $L_i$  and  $m$  right sisters  $R_i$ :

$$X \rightarrow \underbrace{L_n \dots L_1}_{\text{left sisters}} H_X \underbrace{R_1 \dots R_m}_{\text{right sisters}}$$

Replace rule probabilities with a generative process:

For each nonterminal  $X$

- generate its head  $H_X$  (nonterminal or terminal)
- then generate its left sisters  $L_{1..n}$  and a STOP symbol conditioned on  $H_X$
- then generate its right sisters  $R_{1..n}$  and a STOP symbol conditioned on  $H_X$

# Penn Treebank Parsing

# The Penn Treebank

The first publicly available syntactically annotated corpus

Wall Street Journal (50,000 sentences, 1 million words)

also Switchboard, Brown corpus, ATIS

The annotation:

- POS-tagged (Ratnaparkhi's MXPOST)
- Manually annotated with phrase-structure trees
- Richer than standard CFG: *Traces* and other *null elements* used to represent non-local dependencies (designed to allow extraction of predicate-argument structure), although these are typically removed when we do parsing  
[more on non-local dependencies and traces later in the semester]

The standard data set for English phrase-structure parsers



# The Treebank label set

## 48 preterminals (tags):

- 36 POS tags, 12 other symbols (punctuation etc.)
- Simplified version of Brown tagset (87 tags)  
(cf. Lancaster-Oslo/Bergen (LOB) tag set: 126 tags)

## 14 nonterminals:

Standard inventory (S, NP, VP, PP, ADJP, ADVP, SBAR,...)

Many nonterminals have function tags indicating their syntactic roles (NP-SBJ: subject NP) or what role they play

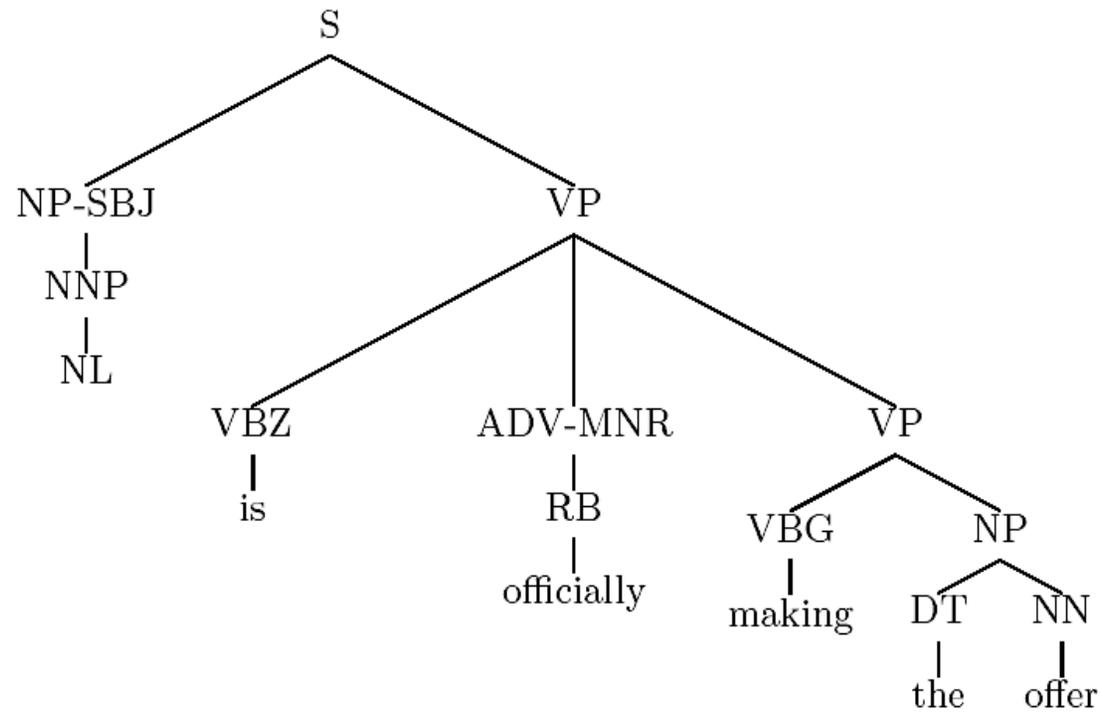
(e.g. PP-LOC: locative PP, i.e. indicating a location [“in NYC”]

PP-DIR: directional PP, indicating a direction [“to NYC”],

ADVP-MNR: manner adverb [“slowly”]).

For historical reasons, these function tags are typically removed before parsing.

# A simple example



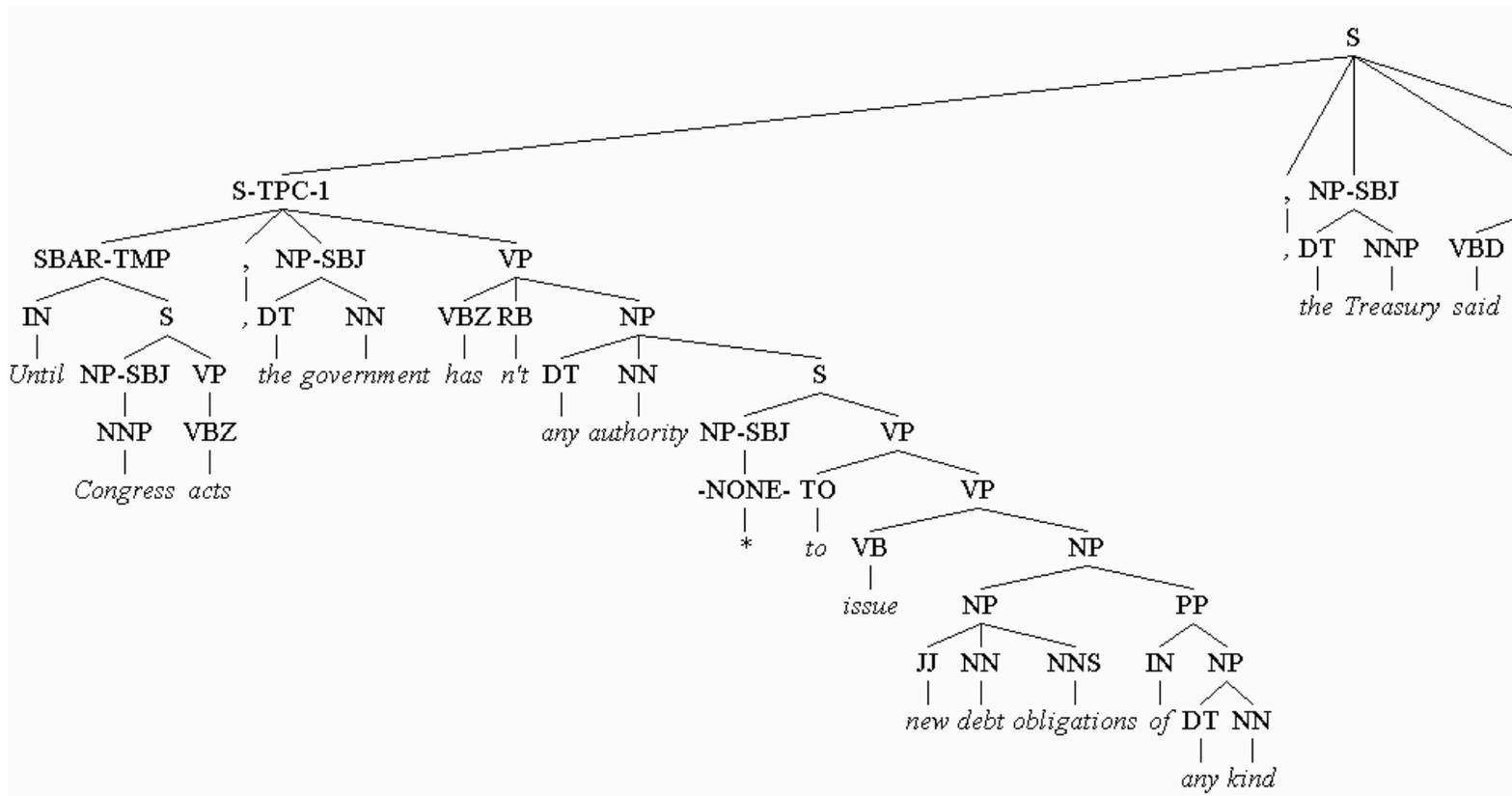
Relatively flat structures:

- There is no noun level
- VP arguments and adjuncts appear at the same level

Function tags, e.g. -SBJ (subject), -MNR (manner)

# A more realistic (partial) example

*Until Congress acts, the government hasn't any authority to issue new debt obligations of any kind, the Treasury said ....*



# The Penn Treebank CFG

The Penn Treebank uses very flat rules, e.g.:

```
NP → DT JJ NN
NP → DT JJ NNS
NP → DT JJ NN NN
NP → DT JJ JJ NN
NP → DT JJ CD NNS
NP → RB DT JJ NN NN
NP → RB DT JJ JJ NNS
NP → DT JJ JJ NNP NNS
NP → DT NNP NNP NNP NNP JJ NN
NP → DT JJ NNP CC JJ JJ NN NNS
NP → RB DT JJS NN NN SBAR
NP → DT VBG JJ NNP NNP CC NNP
NP → DT JJ NNS , NNS CC NN NNS NN
NP → DT JJ JJ VBG NN NNP NNP FW NNP
NP → NP JJ , JJ `` SBAR '' NNS
```

Basic PCFGs don't work well on the Penn Treebank

- Many of these rules appear only once.
- But many of these rules are very similar.

Can we generalize by not treating each rule as atomic?

# Summary

The Penn Treebank has a large number of very flat rules.

Accurate parsing requires modifications to basic PCFG models:

- Generalizing across similar rules (“Markov PCFGs”)
- Modeling word-word dependencies  
(although this does not help as much as people used to think)
- Refining the nonterminals to capture more context

How much of this transfers to other treebanks or languages?

# CFG Parser Evaluation

# Precision and recall

Precision and recall were originally developed as evaluation metrics for information retrieval:

- **Precision**: What percentage of retrieved documents are relevant to the query?
- **Recall**: What percentage of relevant documents were retrieved?

In NLP, they are often used in addition to accuracy:

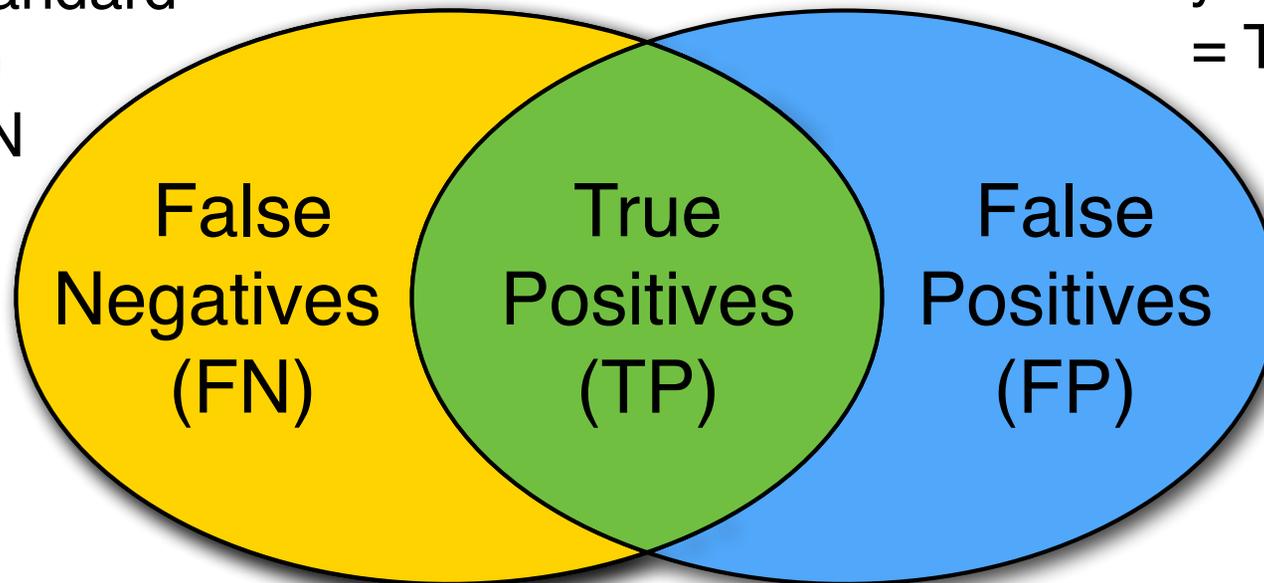
- **Precision**: What percentage of items that were assigned label X do actually have label X in the test data?
- **Recall**: What percentage of items that have label X in the test data were assigned label X by the system?

Particularly useful when there are more than two labels.

# True vs. false positives, false negatives

Items labeled X  
in the gold standard  
(‘truth’)  
= TP + FN

Items labeled X  
by the system  
= TP + FP

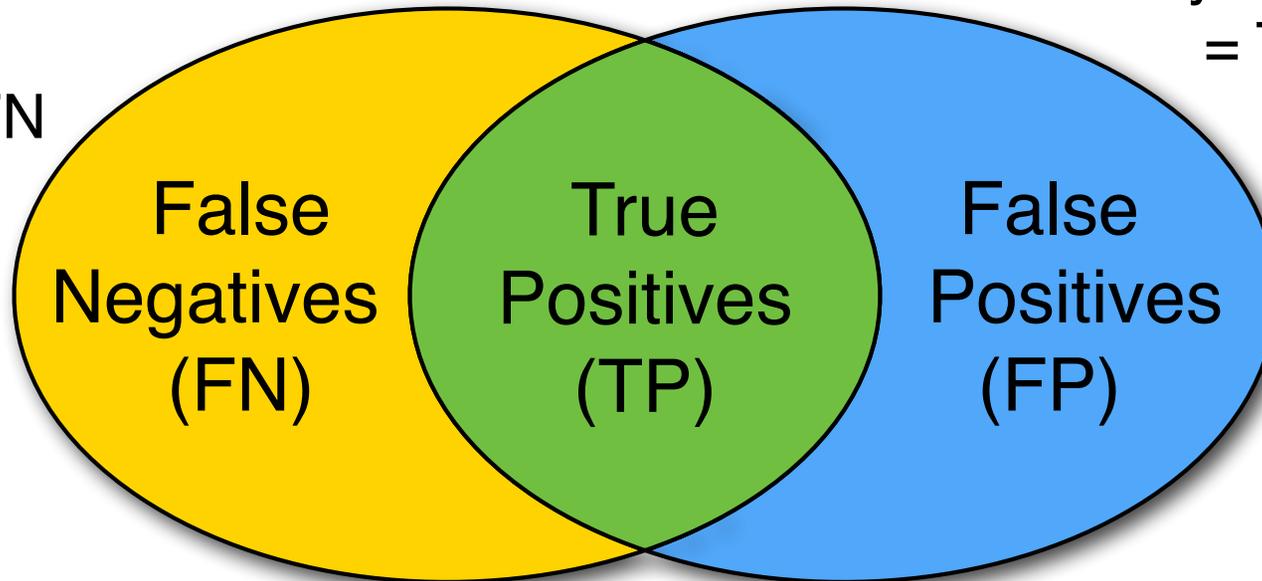


- True positives: Items that were labeled X by the system, and should be labeled X.
- False positives: Items that were labeled X by the system, but should not be labeled X.
- False negatives: Items that were not labeled X by the system, but should be labeled X

# Precision, recall, f-measure

Items labeled X  
in the gold standard  
(‘truth’)  
= TP + FN

Items labeled X  
by the system  
= TP + FP



**Precision:**  $P = \frac{TP}{TP + FP}$

**Recall:**  $R = \frac{TP}{TP + FN}$

**F-measure:** harmonic mean of precision and recall

$$F = \frac{2 \cdot P \cdot R}{P + R}$$

# Evalb (“Parseval”)

Measures recovery of phrase-structure trees.

**Labeled:** span and label (NP, PP,...) has to be right.

[Earlier variant— unlabeled: span of nodes has to be right]

Two aspects of evaluation

**Precision:** *How many of the predicted nodes are correct?*

**Recall:** *How many of the correct nodes were predicted?*

*Usually combined into one metric (**F-measure**):*

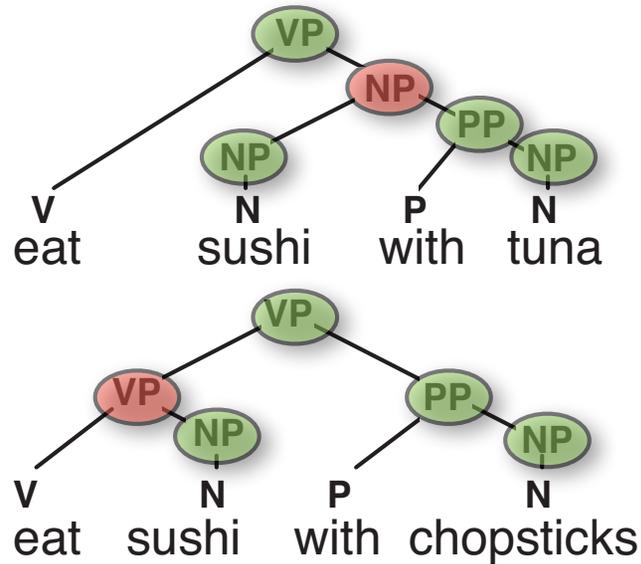
$$P = \frac{\text{\#correctly predicted nodes}}{\text{\#predicted nodes}}$$

$$R = \frac{\text{\#correctly predicted nodes}}{\text{\#correct nodes}}$$

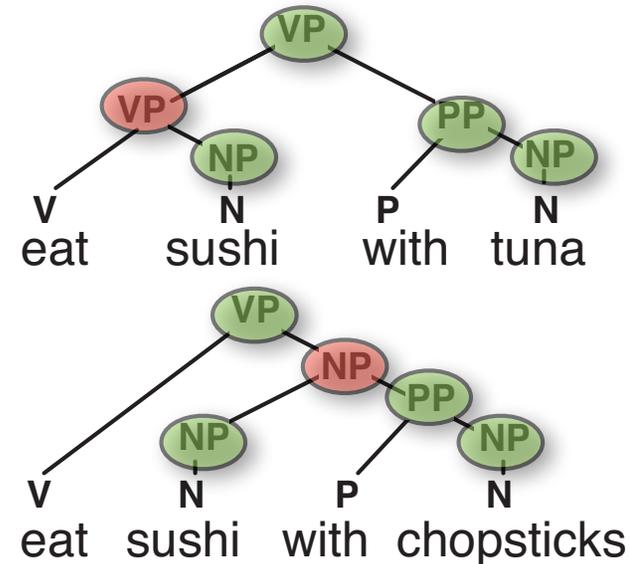
$$F = \frac{2PR}{P + R}$$

# Parseval in practice

## Gold standard



## Parser output



*eat sushi with tuna*: Precision: 4/5 Recall: 4/5

*eat sushi with chopsticks*: Precision: 4/5 Recall: 4/5