

# CMP462: Natural Language Processing



## Lecture 09: CKY Parser

Mohamed Alaa El-Dien Aly  
Computer Engineering Department  
Cairo University  
Spring 2013

# Agenda

- CKY Parsing Algorithm
- Worked example
- Parsing Evaluation
- PCFG Training

## **Acknowledgment:**

Most slides adapted from Chris Manning and Dan Jurafsky's NLP class on [Coursera](#).



# Review:

## Context Free Grammars (CFGs)

- $G = (T, N, S, R)$ 
  - T is a set of terminal symbols e.g. fish, people
  - N is a set of nonterminal symbols e.g. NP, VP
  - S is the start symbol ( $S \in N$ )
  - R is a set of rules/productions of the form  $X \rightarrow \gamma$  where  $X \in N$  and  $\gamma \in (N \cup T)^*$   
e.g.  $S \rightarrow NP VP$
  
- A grammar G generates a language L.



# Review:

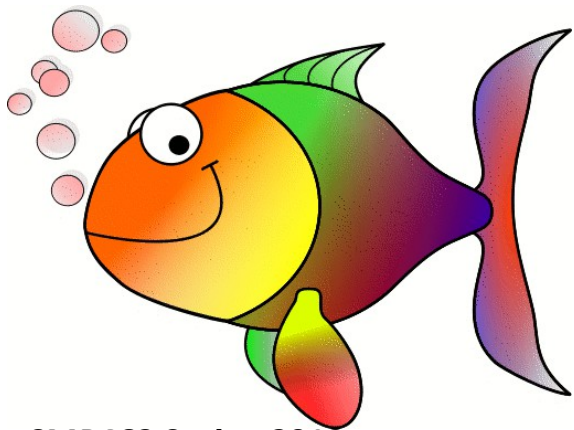
## Probabilistic Context Free Grammars (PCFGs)

- $G = (T, N, S, R, P)$ 
  - T is a set of terminal symbols
  - N is a set of nonterminal symbols
  - S is the start symbol ( $S \in N$ )
  - R is a set of rules/productions of the form  $X \rightarrow \gamma$
  - P is a probability function
    - $P: R \rightarrow [0,1]$
    - $\forall X \in N, \sum_{X \rightarrow \gamma \in R} P(X \rightarrow \gamma) = 1$



## Review: A PCFG

$S \rightarrow NP VP$	1.0		$N \rightarrow people$	0.5
$VP \rightarrow V NP$	0.6	} Sum to 1	$N \rightarrow fish$	0.2
$VP \rightarrow V NP PP$	0.4		$N \rightarrow tanks$	0.2
$NP \rightarrow NP NP$	0.1	} Sum to 1	$N \rightarrow rods$	0.1
$NP \rightarrow NP PP$	0.2		$V \rightarrow people$	0.1
$NP \rightarrow N$	0.7		$V \rightarrow fish$	0.6
$PP \rightarrow P NP$	1.0		$V \rightarrow tanks$	0.3
			$P \rightarrow with$	1.0



[With empty NP removed so less ambiguous]



## Review: The probability of trees and strings

- $P(t)$  – The probability of a tree  $t$  is the product of the probabilities of the rules used to generate it.
- $P(s)$  – The probability of the string  $s$  is the sum of the probabilities of the trees which have that string as their yield

$$P(s) = \sum_j P(t) \text{ where } t \text{ is a parse of } s$$



# Review: Chomsky Normal Form

- All rules are of the form  $X \rightarrow YZ$  or  $X \rightarrow w$ 
  - $X, Y, Z \in N$  and  $w \in T$
- A transformation to this form doesn't change the weak generative capacity of a CFG
  - That is, it recognizes the same language
    - But maybe with different trees
- Transformations:
  - Empties and unaries are removed recursively
  - n-ary rules are divided by introducing new nonterminals ( $n > 2$ )



## Original CFG

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$VP \rightarrow V NP PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$NP \rightarrow N$

$NP \rightarrow e$

$PP \rightarrow P NP$

$N \rightarrow \text{people}$

$N \rightarrow \text{fish}$

$N \rightarrow \text{tanks}$

$N \rightarrow \text{rods}$

$V \rightarrow \text{people}$

$V \rightarrow \text{fish}$

$V \rightarrow \text{tanks}$

$P \rightarrow \text{with}$

*people fish tanks*

*people fish with rods*





# CFG in Chomsky Normal Form

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V @VP\_V$

$@VP\_V \rightarrow NP PP$

$S \rightarrow V @S\_V$

$@S\_V \rightarrow NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$NP \rightarrow P NP$

$PP \rightarrow P NP$

$NP \rightarrow \textit{people}$

$NP \rightarrow \textit{fish}$

$NP \rightarrow \textit{tanks}$

$NP \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$S \rightarrow \textit{people}$

$VP \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$S \rightarrow \textit{fish}$

$VP \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$S \rightarrow \textit{tanks}$

$VP \rightarrow \textit{tanks}$

$P \rightarrow \textit{with}$

$PP \rightarrow \textit{with}$



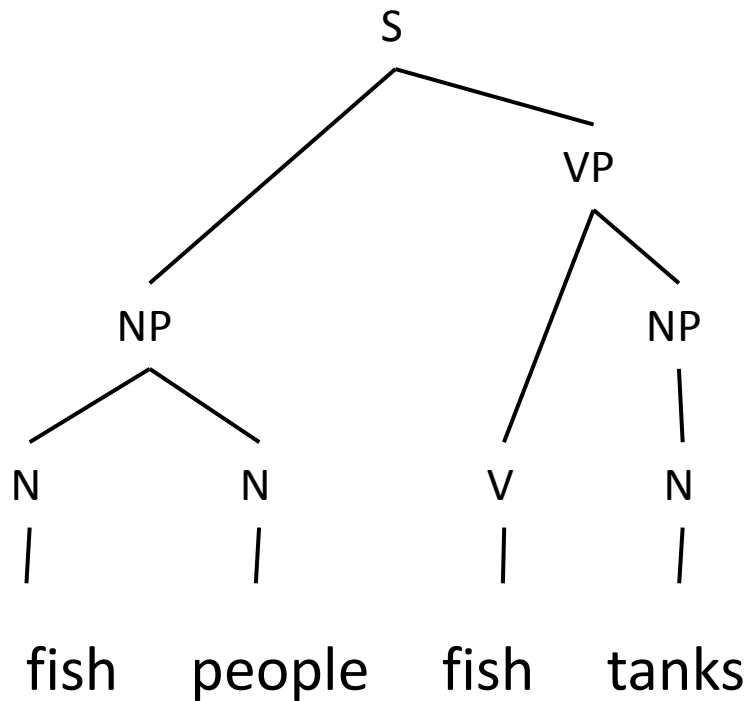
# CKY Parsing

Exact polynomial  
time parsing of  
(P)CFGs



# Constituency Parsing

Arrive at a parse tree consistent with the grammar



Start with a sentence

## PCFG

Rule Prob  $\theta_i$

$S \rightarrow NP VP \quad \theta_0$

$NP \rightarrow NP NP \quad \theta_1$

...

$N \rightarrow \text{fish} \quad \theta_{42}$

$N \rightarrow \text{people} \quad \theta_{43}$

$V \rightarrow \text{fish} \quad \theta_{44}$

...

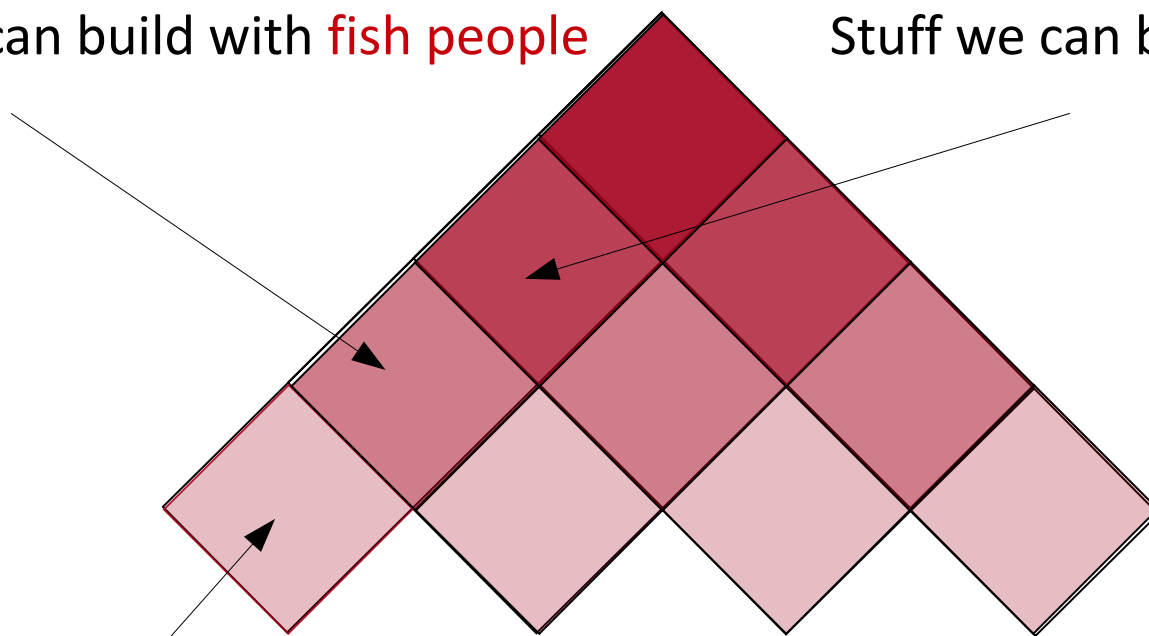


# Cocke-Kasami-Younger (CKY) Constituency Parsing

Parse Triangle or Chart

Stuff we can build with **fish people**

Stuff we can build with **fish people fish**



Fill chart bottom-up

fish people fish tanks

Stuff we can build with **fish**



# Viterbi (Max) Scores

Combine a constituent from the **left** with one from the **right** according to a grammar **rule** and compute its **probability**

$$\begin{aligned} \text{NP} &\rightarrow \text{NP NP} \\ 0.35 &\times 0.14 \times 0.1 \\ &= 0.0049 \end{aligned}$$

NP 0.35

V 0.1

N 0.5

people

VP 0.06

NP 0.14

V 0.6

N 0.2

fish

S  $\rightarrow$  NP VP 0.9

S  $\rightarrow$  VP 0.1

VP  $\rightarrow$  V NP 0.5

VP  $\rightarrow$  V 0.1

VP  $\rightarrow$  V @VP\_V 0.3

VP  $\rightarrow$  V PP 0.1

@VP\_V  $\rightarrow$  NP PP 1.0

NP  $\rightarrow$  NP NP 0.1

NP  $\rightarrow$  NP PP 0.2

NP  $\rightarrow$  N 0.7

PP  $\rightarrow$  P NP 1.0



# Viterbi (Max) Scores

Combine a constituent from the **left** with one from the **right** according to a grammar **rule** and compute its **probability**

$$\begin{aligned} \text{VP} &\rightarrow \text{V NP} \\ 0.1 \times 0.14 \times 0.5 \\ &= 0.007 \end{aligned}$$

NP 0.0049

NP 0.35

V 0.1

N 0.5

people

VP 0.06

NP 0.14

V 0.6

N 0.2

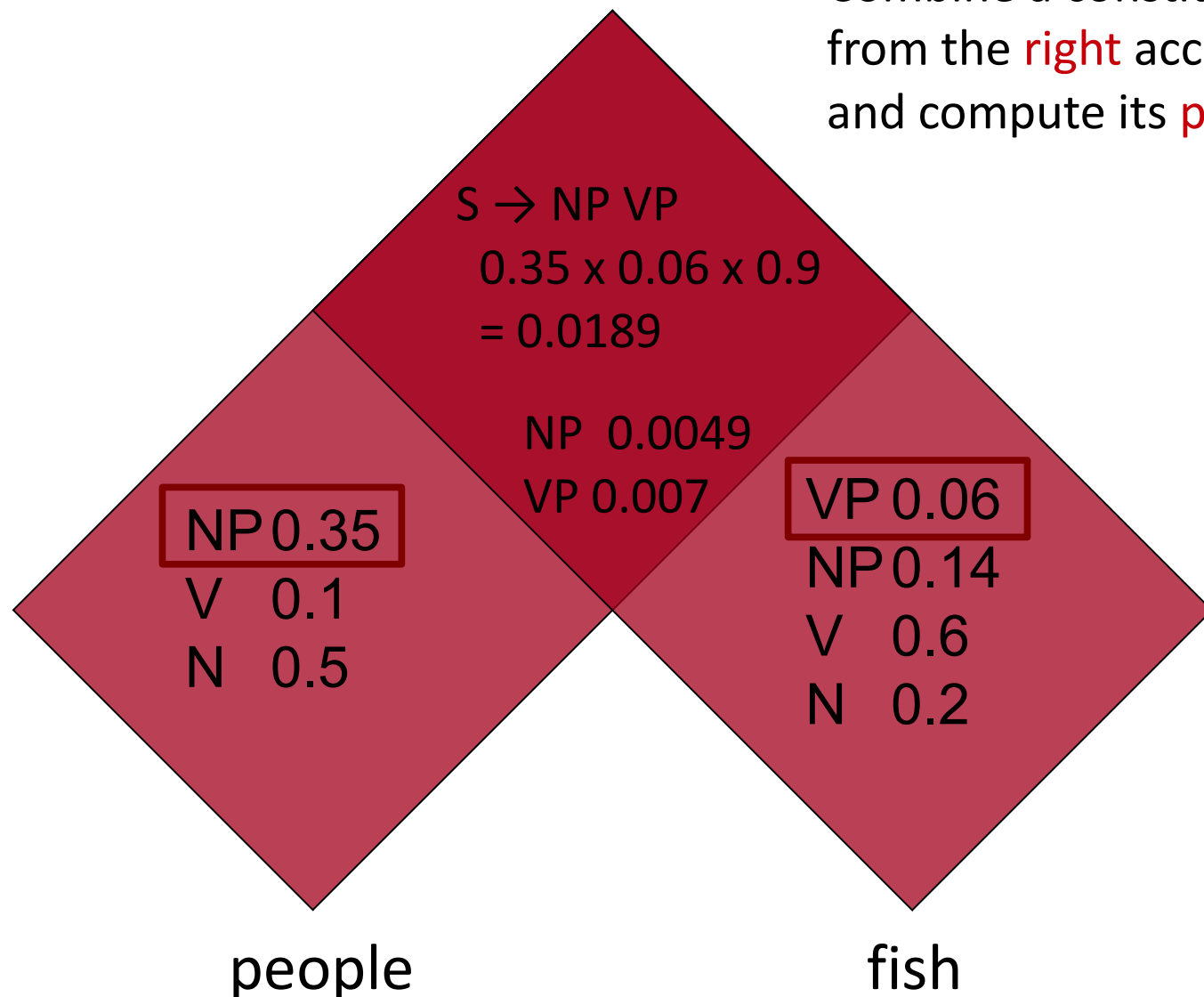
fish

$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP\_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP\_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0



# Viterbi (Max) Scores

Combine a constituent from the **left** with one from the **right** according to a grammar **rule** and compute its **probability**



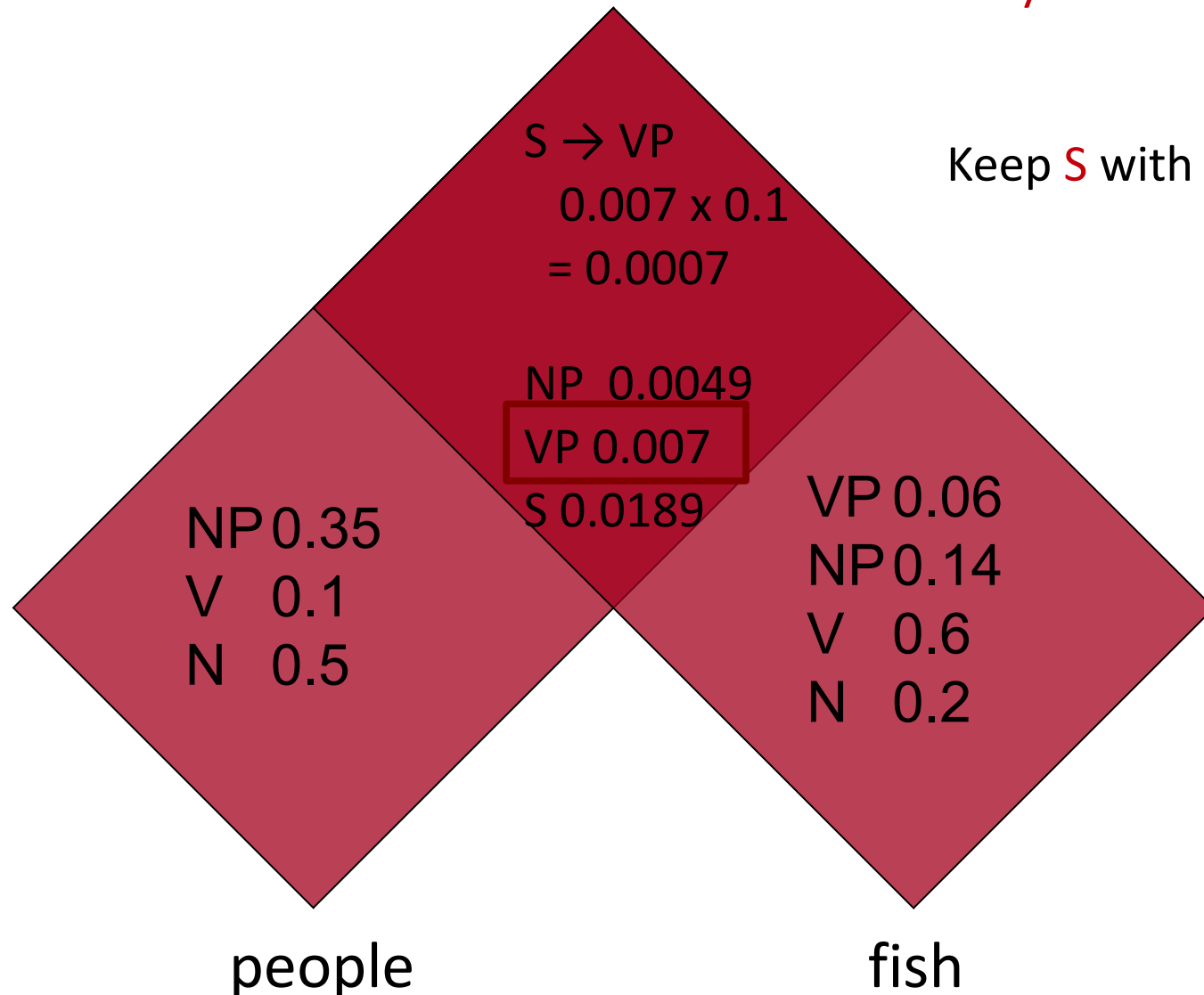
$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP\_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP\_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0



# Viterbi (Max) Scores

Use **unary** rules inside the cell

Keep **S** with the **highest** probability



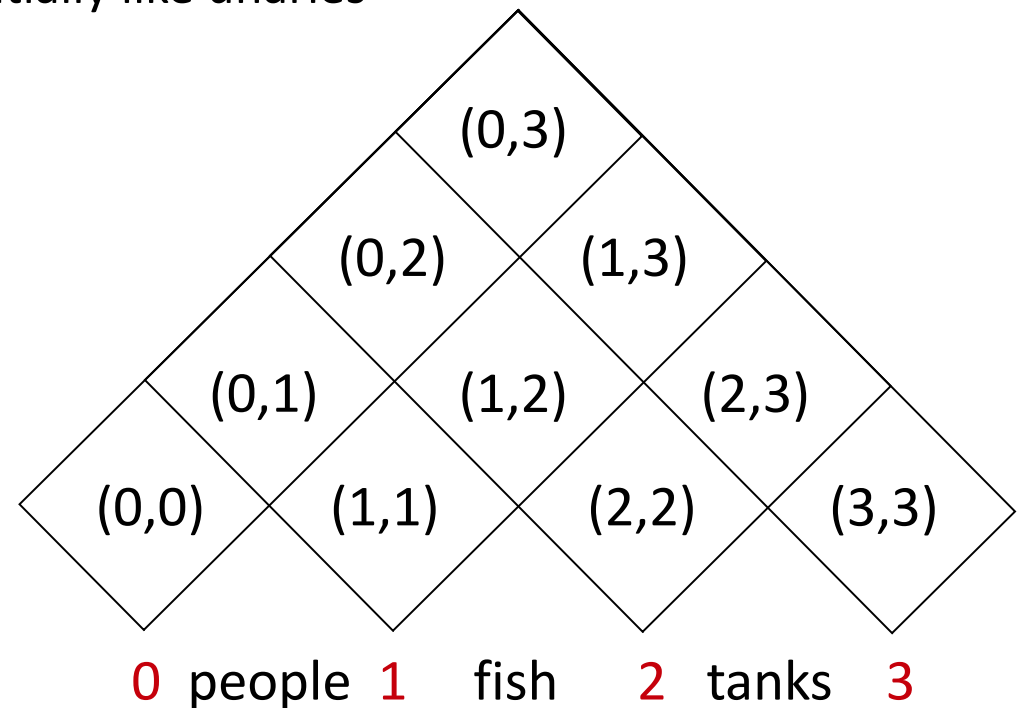
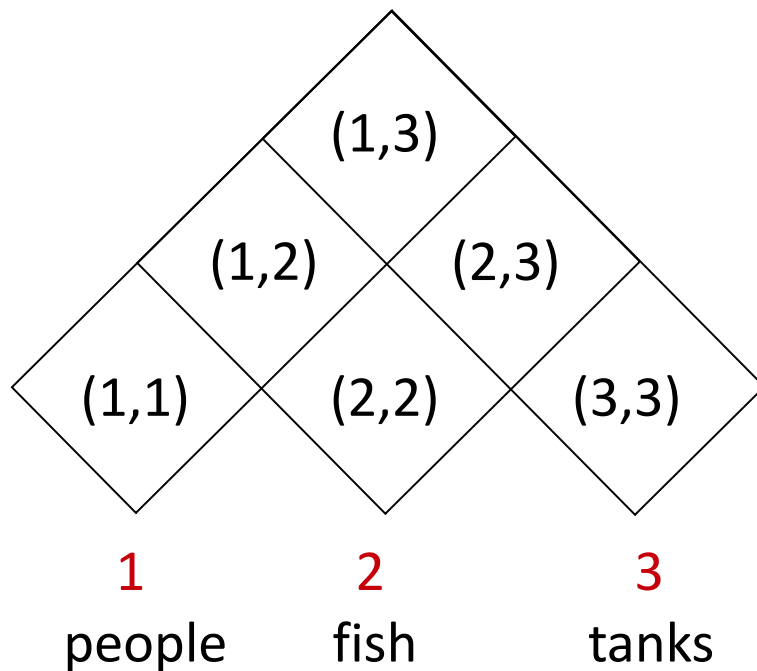
$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP\_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP\_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0





# Extended CKY parsing

- Unaries can be incorporated into the algorithm
  - Messy, but doesn't increase algorithmic complexity
- Empties can be incorporated
  - Use fenceposts
  - Doesn't increase complexity; essentially like unaries





# Extended CKY parsing

- Unaries can be incorporated into the algorithm
  - Messy, but doesn't increase algorithmic complexity
- Empties can be incorporated
  - Use fenceposts
  - Doesn't increase complexity; essentially like unaries
- Binarization is *vital*
  - Without binarization, you don't get parsing cubic in the length of the sentence and in the number of nonterminals in the grammar
    - Binarization may be an explicit transformation or implicit in how the parser works (Earley-style dotted rules), but it's always there.



# The CKY algorithm (1960/1965) ... extended to unaries

```

function CKY(words, grammar) returns [most_probable_parse, prob]
  score = new double[#(words)+1][#(words)+1][#(nonterms)]
  back = new Pair[#(words)+1][#(words)+1][#nonterms]
  for i=0; i<#(words); i++
    for A in nonterms
      if A -> words[i] in grammar
        score[i][i+1][A] = P(A -> words[i])
  //handle unaries
  boolean added = true
  while added
    added = false
    for A, B in nonterms
      if score[i][i+1][B] > 0 && A->B in grammar
        prob = P(A->B)*score[i][i+1][B]
        if prob > score[i][i+1][A]
          score[i][i+1][A] = prob
          back[i][i+1][A] = B
          added = true

```

Array of **scores**  
one element per  
**nonterminal** per **cell**

**back pointers** for  
reconstructing the  
best parse tree



# The CKY algorithm (1960/1965) ... extended to unaries

```

function CKY(words, grammar) returns [most_probable_parse, prob]
  score = new double[#(words)+1][#(words)+1][#(nonterms)]
  back = new Pair[#(words)+1][#(words)+1][#nonterms]]
  for i=0; i<#(words); i++
    for A in nonterms
      if A -> words[i] in grammar
        score[i][i+1][A] = P(A -> words[i])
      //handle unaries
      boolean added = true
      while added
        added = false
        for A, B in nonterms
          if score[i][i+1][B] > 0 && A->B in grammar
            prob = P(A->B)*score[i][i+1][B]
            if prob > score[i][i+1][A]
              score[i][i+1][A] = prob
              back[i][i+1][A] = B
            added = true

```

Lexicon

Rules  $A \rightarrow \text{word}$

Unary Rules  $B \rightarrow A$



# The CKY algorithm (1960/1965)

## ... extended to unaries

```

for span = 2 to #(words)
  for begin = 0 to #(words) - span
    end = begin + span
    for split = begin+1 to end-1
      for A, B, C in nonterms
        prob = score[begin][split][B] * score[split][end][C] * P(A -> BC)
        if prob > score[begin][end][A]
          score[begin][end][A] = prob
          back[begin][end][A] = new Triple(split, B, C)
      //handle unaries
      boolean added = true
      while added
        added = false
        for A, B in nonterms
          prob = P(A -> B) * score[begin][end][B];
          if prob > score[begin][end][A]
            score[begin][end][A] = prob
            back[begin][end][A] = B
            added = true
    return buildTree(score, back)
  
```

← Loop on **spans**

← Rules **A → B C**

← Rules **C → A**

Grammar





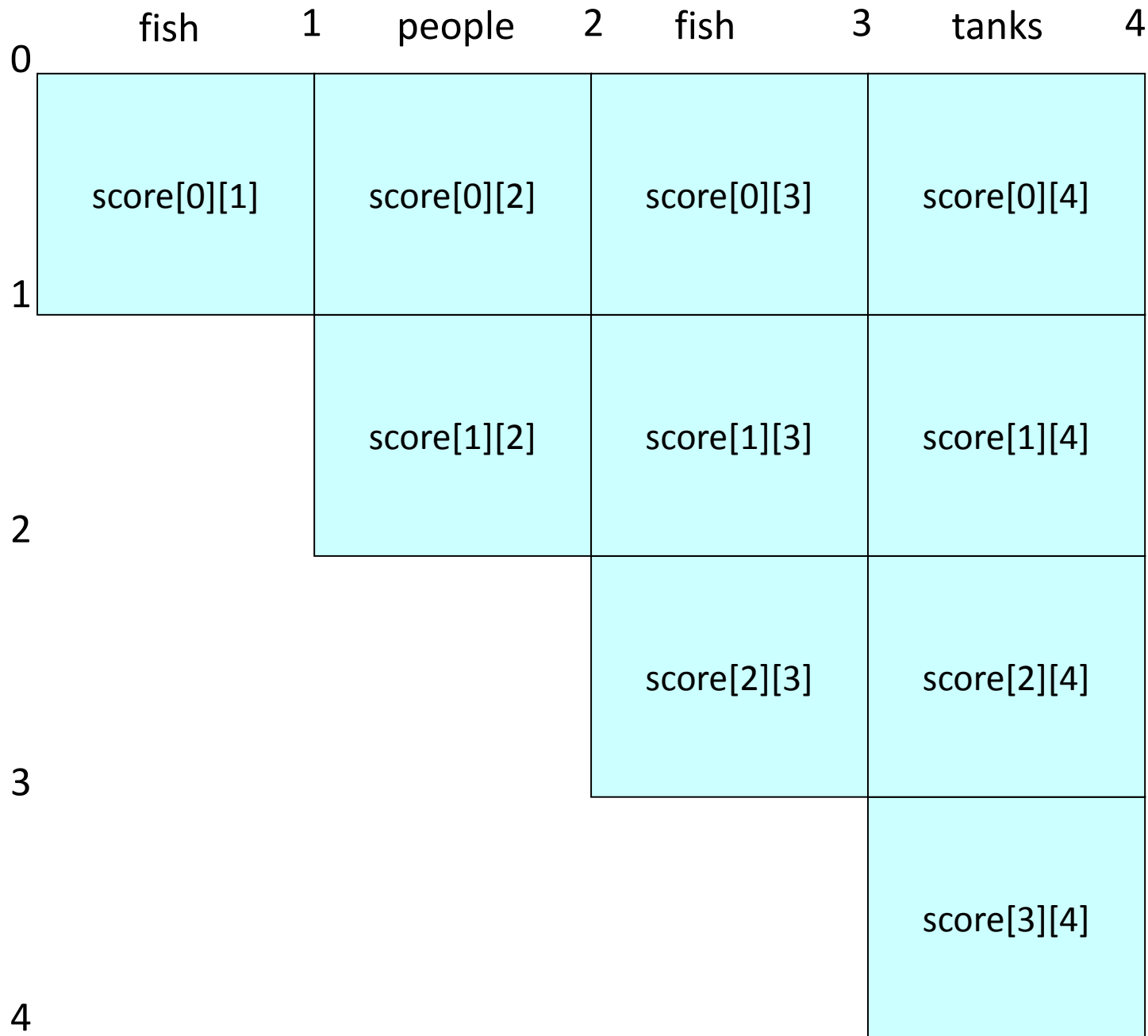
# The grammar: Binary, no epsilons,

$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP\_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP\_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0

$N \rightarrow people$	0.5
$N \rightarrow fish$	0.2
$N \rightarrow tanks$	0.2
$N \rightarrow rods$	0.1
$V \rightarrow people$	0.1
$V \rightarrow fish$	0.6
$V \rightarrow tanks$	0.3
$P \rightarrow with$	1.0



Rotate the triangle 45 degrees to the right!







- S → NP VP 0.9
- S → VP 0.1
- VP → V NP 0.5
- VP → V 0.1
- VP → V @VP\_V 0.3
- VP → V PP 0.1
- @VP\_V → NP PP 1.0
- NP → NP NP 0.1
- NP → NP PP 0.2
- NP → N 0.7
- PP → P NP 1.0
  
- N → *people* 0.5
- N → *fish* 0.2
- N → *tanks* 0.2
- N → *rods* 0.1
- V → *people* 0.1
- V → *fish* 0.6
- V → *tanks* 0.3
- P → *with* 1.0

	fish 1	people 2	fish 3	tanks 4
0	N → fish 0.2 V → fish 0.6			
1		N → people 0.5 V → people 0.1		
2			N → fish 0.2 V → fish 0.6	
3				N → tanks 0.2 V → tanks 0.1
4				

```

for i=0; i<#(words); i++
  for A in nonterms
    if A -> words[i] in grammar
      score[i][i+1][A] = P(A -> words[i]);
    
```



- S → NP VP 0.9
- S → VP 0.1**
- VP → V NP 0.5
- VP → V 0.1**
- VP → V @VP\_V 0.3
- VP → V PP 0.1
- @VP\_V → NP PP 1.0
- NP → NP NP 0.1
- NP → NP PP 0.2
- NP → N 0.7**
- PP → P NP 1.0

- N → *people* 0.5
- N → *fish* 0.2
- N → *tanks* 0.2
- N → *rods* 0.1
- V → *people* 0.1
- V → *fish* 0.6
- V → *tanks* 0.3
- P → *with* 1.0

	fish 1	people 2	fish 3	tanks 4
0	N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006			
1		N → people 0.5 V → people 0.1		
2			N → fish 0.2 V → fish 0.6	
				N → tanks 0.2 V → tanks 0.1

```
// handle unaries
boolean added = true
while added
    added = false
    for A, B in nonterms
        if score[i][i+1][B] > 0 && A->B in grammar
            prob = P(A->B)*score[i][i+1][B]
            if(prob > score[i][i+1][A])
                score[i][i+1][A] = prob
                back[i][i+1][A] = B
                added = true
```



- S → NP VP 0.9
- S → VP 0.1
- VP → V NP 0.5
- VP → V 0.1
- VP → V @VP\_V 0.3
- VP → V PP 0.1
- @VP\_V → NP PP 1.0
- NP → NP NP 0.1
- NP → NP PP 0.2
- NP → N 0.7
- PP → P NP 1.0
  
- N → *people* 0.5
- N → *fish* 0.2
- N → *tanks* 0.2
- N → *rods* 0.1
- V → *people* 0.1
- V → *fish* 0.6
- V → *tanks* 0.3
- P → *with* 1.0

	fish 1	people 2	fish 3	tanks 4
0	N → fish 0.2 <span style="border: 1px solid red; padding: 2px;">V → fish 0.6</span> NP → N 0.14 VP → V 0.06 S → VP 0.006	VP → V NP 0.105 S → NP VP 0.00126 NP → NP NP 0.0049		
1		N → people 0.5 V → people 0.1 <span style="border: 1px solid red; padding: 2px;">NP → N 0.35</span> VP → V 0.01 S → VP 0.001		
2			N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006	
3				N → tanks 0.2 V → tanks 0.1 NP → N 0.14 VP → V 0.03 S → VP 0.003
4				

```

prob=score[begin][split][B]*score[split][end][C]*P(A->BC)
if (prob > score[begin][end][A])
    score[begin][end][A] = prob
    back[begin][end][A] = new Triple(split,B,C)
    
```



- S → NP VP 0.9
- S → VP 0.1**
- VP → V NP 0.5
- VP → V 0.1
- VP → V @VP\_V 0.3
- VP → V PP 0.1
- @VP\_V → NP PP 1.0
- NP → NP NP 0.1
- NP → NP PP 0.2
- NP → N 0.7
- PP → P NP 1.0
  
- N → *people* 0.5
- N → *fish* 0.2
- N → *tanks* 0.2
- N → *rods* 0.1
- V → *people* 0.1
- V → *fish* 0.6
- V → *tanks* 0.3
- P → *with* 1.0

	fish	1	people	2	fish	3	tanks	4
0	N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006		<b>VP → V NP 0.105</b> S → NP VP 0.00126 NP → NP NP 0.0049					
1			<b>S → VP 0.0105</b>					
2			N → people 0.5 V → people 0.1 NP → N 0.35 VP → V 0.01 S → VP 0.001	NP → NP NP 0.0049 VP → V NP 0.007 S → NP VP 0.0189				
3				N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006	NP → NP NP 0.00196 VP → V NP 0.042 S → NP VP 0.00378			
4						N → tanks 0.2 V → tanks 0.1 NP → N 0.14 VP → V 0.03 S → VP 0.003		

Handle *unaries*

```
//handle unaries
boolean added = true
while added
  added = false
  for A, B in nonterms
    prob = P(A->B)*score[begin][end][B];
    if prob > score[begin][end][A]
      score[begin][end][A] = prob
      back[begin][end][A] = B
  added = true
```



- S → NP VP 0.9
- S → VP 0.1
- VP → V NP 0.5
- VP → V 0.1
- VP → V @VP\_V 0.3
- VP → V PP 0.1
- @VP\_V → NP PP 1.0
- NP → NP NP 0.1
- NP → NP PP 0.2
- NP → N 0.7
- PP → P NP 1.0
  
- N → *people* 0.5
- N → *fish* 0.2
- N → *tanks* 0.2
- N → *rods* 0.1
- V → *people* 0.1
- V → *fish* 0.6
- V → *tanks* 0.3
- P → *with* 1.0

	fish	1	people	2	fish	3	tanks	4
0	N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006		NP → NP NP 0.0049 VP → V NP 0.105 <i>S → VP</i> 0.0105					
1			N → people 0.5 V → people 0.1 NP → N 0.35 VP → V 0.01 S → VP 0.001	NP → NP NP 0.0049 VP → V NP 0.007 S → NP VP 0.0189				
2				N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006		NP → NP NP 0.00196 VP → V NP 0.042 S → NP VP 0.00378		
3						N → tanks 0.2 V → tanks 0.1 NP → N 0.14 VP → V 0.03 S → VP 0.003		
4								



- S → NP VP 0.9
- S → VP 0.1
- VP → V NP 0.5
- VP → V 0.1
- VP → V @VP\_V 0.3
- VP → V PP 0.1
- @VP\_V → NP PP 1.0
- NP → NP NP 0.1
- NP → NP PP 0.2
- NP → N 0.7
- PP → P NP 1.0
- N → *people* 0.5
- N → *fish* 0.2
- N → *tanks* 0.2
- N → *rods* 0.1
- V → *people* 0.1
- V → *fish* 0.6
- V → *tanks* 0.3
- P → *with* 1.0

	0	fish	1	people	2	fish	3	tanks	4
	0	<p>N → fish 0.2</p> <p>V → fish 0.6</p> <p>NP → N 0.14</p> <p>VP → V 0.06</p> <p>S → VP 0.006</p>	<p>NP → NP NP 0.0049</p> <p>VP → V NP 0.105</p> <p>S → VP 0.0105</p>	To fill this cell span (0, 3)					
span (0,1)	1		<p>N → people 0.5</p> <p>V → people 0.1</p> <p>NP → N 0.35</p> <p>VP → V 0.01</p> <p>S → VP 0.001</p>	<p>NP → NP NP 0.0049</p> <p>VP → V NP 0.007</p> <p>S → NP VP 0.0189</p>					
	2			span (1,3)	<p>N → fish 0.2</p> <p>V → fish 0.6</p> <p>NP → N 0.14</p> <p>VP → V 0.06</p> <p>S → VP 0.006</p>	<p>NP → NP NP 0.00196</p> <p>VP → V NP 0.042</p> <p>S → VP 0.0042</p>			
Find something from these cells	3					<p>N → tanks 0.2</p> <p>V → tanks 0.1</p> <p>NP → N 0.14</p> <p>VP → V 0.03</p> <p>S → VP 0.003</p>			
	4								

```

for split = begin+1 to end-1
  for A,B,C in nonterms
    prob=score[begin][split][B]*score[split][end][C]*P(A->BC)
    if prob > score[begin][end][A]
      score[begin][end][A] = prob
      back[begin][end][A] = new Triple(split,B,C)
    
```



- S → NP VP 0.9
- S → VP 0.1
- VP → V NP 0.5
- VP → V 0.1
- VP → V @VP\_V 0.3
- VP → V PP 0.1
- @VP\_V → NP PP 1.0
- NP → NP NP 0.1
- NP → NP PP 0.2
- NP → N 0.7
- PP → P NP 1.0
- N → *people* 0.5
- N → *fish* 0.2
- N → *tanks* 0.2
- N → *rods* 0.1
- V → *people* 0.1
- V → *fish* 0.6
- V → *tanks* 0.3
- P → *with* 1.0

	fish	1	people	2	fish	3	tanks	4
0	N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006	<div style="border: 2px solid red; border-radius: 50%; padding: 10px;">                     NP → NP NP 0.0049                      VP → V NP 0.105                      S → VP 0.0105                 </div>		To fill this cell span (0, 3)				
1								
2	span (0,2)		N → people 0.5 V → people 0.1 NP → N 0.35 VP → V 0.01 S → VP 0.001	NP → NP NP 0.0049 VP → V NP 0.007 S → NP VP 0.0189				
3	Or these cells		<div style="border: 2px solid red; border-radius: 50%; padding: 10px;">                     N → fish 0.2                      V → fish 0.6                      NP → N 0.14                      VP → V 0.06                      S → VP 0.006                 </div>		NP → NP NP 0.00196 VP → V NP 0.042 S → VP 0.0042			
4					N → tanks 0.2 V → tanks 0.1 NP → N 0.14 VP → V 0.03 S → VP 0.003			

```

for split = begin+1 to end-1
  for A,B,C in nonterms
    prob=score[begin][split][B]*score[split][end][C]*P(A->BC)
    if prob > score[begin][end][A]
      score[begin][end][A] = prob
      back[begin][end][A] = new Triple(split,B,C)
    
```



$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0
$N \rightarrow people$	0.5
$N \rightarrow fish$	0.2
$N \rightarrow tanks$	0.2
$N \rightarrow rods$	0.1
$V \rightarrow people$	0.1
$V \rightarrow fish$	0.6
$V \rightarrow tanks$	0.3
$P \rightarrow with$	1.0

	fish	1	people	2	fish	3	tanks	4
0	$N \rightarrow fish$ 0.2 $V \rightarrow fish$ 0.6 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.06 $S \rightarrow VP$ 0.006		$NP \rightarrow NP NP$ 0.0049 $VP \rightarrow V NP$ 0.105 $S \rightarrow VP$ 0.0105		$NP \rightarrow NP NP$ $S \rightarrow NP VP$			
1			$N \rightarrow people$ 0.5 $V \rightarrow people$ 0.1 $NP \rightarrow N$ 0.35 $VP \rightarrow V$ 0.01 $S \rightarrow VP$ 0.001		$NP \rightarrow NP NP$ 0.0049 $VP \rightarrow V NP$ 0.007 $S \rightarrow NP VP$ 0.0189			
2					$N \rightarrow fish$ 0.2 $V \rightarrow fish$ 0.6 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.06 $S \rightarrow VP$ 0.006		$NP \rightarrow NP NP$ 0.00196 $VP \rightarrow V NP$ 0.042 $S \rightarrow VP$ 0.0042	
3							$N \rightarrow tanks$ 0.2 $V \rightarrow tanks$ 0.1 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.03 $S \rightarrow VP$ 0.003	
4								

```

for split = begin+1 to end-1
  for A,B,C in nonterms
    prob=score[begin][split][B]*score[split][end][C]*P(A->BC)
    if prob > score[begin][end][A]
      score[begin][end][A] = prob
      back[begin][end][A] = new Triple(split,B,C)
    
```





$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$VP \rightarrow V NP$	0.5
$VP \rightarrow V$	0.1
$VP \rightarrow V @VP\_V$	0.3
$VP \rightarrow V PP$	0.1
$@VP\_V \rightarrow NP PP$	1.0
$NP \rightarrow NP NP$	0.1
$NP \rightarrow NP PP$	0.2
$NP \rightarrow N$	0.7
$PP \rightarrow P NP$	1.0
$N \rightarrow people$	0.5
$N \rightarrow fish$	0.2
$N \rightarrow tanks$	0.2
$N \rightarrow rods$	0.1
$V \rightarrow people$	0.1
$V \rightarrow fish$	0.6
$V \rightarrow tanks$	0.3
$P \rightarrow with$	1.0

	0	fish	1	people	2	fish	3	tanks	4
0	$N \rightarrow fish$ 0.2 $V \rightarrow fish$ 0.6 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.06 $S \rightarrow VP$ 0.006		$NP \rightarrow NP NP$ 0.0049 $VP \rightarrow V NP$ 0.105 $S \rightarrow VP$ 0.0105		$NP \rightarrow NP NP$ $S \rightarrow NP VP$ $NP \rightarrow NP NP$ $VP \rightarrow V NP$ $S \rightarrow NP VP$				
1				$N \rightarrow people$ 0.5 $V \rightarrow people$ 0.1 $NP \rightarrow N$ 0.35 $VP \rightarrow V$ 0.01 $S \rightarrow VP$ 0.001		$NP \rightarrow NP NP$ 0.0049 $VP \rightarrow V NP$ 0.007 $S \rightarrow NP VP$ 0.0189			
2						$N \rightarrow fish$ 0.2 $V \rightarrow fish$ 0.6 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.06 $S \rightarrow VP$ 0.006		$NP \rightarrow NP NP$ 0.00196 $VP \rightarrow V NP$ 0.042 $S \rightarrow VP$ 0.0042	
3								$N \rightarrow tanks$ 0.2 $V \rightarrow tanks$ 0.1 $NP \rightarrow N$ 0.14 $VP \rightarrow V$ 0.03 $S \rightarrow VP$ 0.003	
4									

Keep **NP** and **S** with highest probability

```

for split = begin+1 to end-1
  for A,B,C in nonterms
    prob=score[begin][split][B]*score[split][end][C]*P(A->BC)
    if prob > score[begin][end][A]
      score[begin][end][A] = prob
      back[begin][end][A] = new Triple(split,B,C)
    
```



- S → NP VP 0.9
- S → VP 0.1
- VP → V NP 0.5
- VP → V 0.1
- VP → V @VP\_V 0.3
- VP → V PP 0.1
- @VP\_V → NP PP 1.0
- NP → NP NP 0.1
- NP → NP PP 0.2
- NP → N 0.7
- PP → P NP 1.0
- N → *people* 0.5
- N → *fish* 0.2
- N → *tanks* 0.2
- N → *rods* 0.1
- V → *people* 0.1
- V → *fish* 0.6
- V → *tanks* 0.3
- P → *with* 1.0

	fish 1	people 2	fish 3	tanks 4
0	N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006	NP → NP NP 0.0049 VP → V NP 0.105 S → VP 0.0105	NP → NP NP 0.0000686 VP → V NP 0.00147 S → NP VP 0.000882	
1		N → people 0.5 V → people 0.1 NP → N 0.35 VP → V 0.01 S → VP 0.001	NP → NP NP 0.0049 VP → V NP 0.007 S → NP VP 0.0189	
2			N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006	NP → NP NP 0.00196 VP → V NP 0.042 S → VP 0.0042
3				N → tanks 0.2 V → tanks 0.1 NP → N 0.14 VP → V 0.03 S → VP 0.003
4				

Handle **unaries!**

```

for split = begin+1 to end-1
  for A,B,C in nonterms
    prob=score[begin][split][B]*score[split][end][C]*P(A->BC)
    if prob > score[begin][end][A]
      score[begin][end][A] = prob
      back[begin][end][A] = new Triple(split,B,C)
    
```



- S → NP VP 0.9
- S → VP 0.1
- VP → V NP 0.5
- VP → V 0.1
- VP → V @VP\_V 0.3
- VP → V PP 0.1
- @VP\_V → NP PP 1.0
- NP → NP NP 0.1
- NP → NP PP 0.2
- NP → N 0.7
- PP → P NP 1.0
- N → *people* 0.5
- N → *fish* 0.2
- N → *tanks* 0.2
- N → *rods* 0.1
- V → *people* 0.1
- V → *fish* 0.6
- V → *tanks* 0.3
- P → *with* 1.0

	fish	1	people	2	fish	3	tanks	4
0	N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006	NP → NP NP 0.0049 VP → V NP 0.105 S → VP 0.0105		NP → NP NP 0.0000686 VP → V NP 0.00147 S → NP VP 0.000882				
1			N → people 0.5 V → people 0.1 NP → N 0.35 VP → V 0.01 S → VP 0.001	NP → NP NP 0.0049 VP → V NP 0.007 S → NP VP 0.0189		Fill cell span (1,4)		
2				N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006				
3							N → tanks 0.2 V → tanks 0.1 NP → N 0.14 VP → V 0.03 S → VP 0.003	
4								

```

for split = begin+1 to end-1
  for A,B,C in nonterms
    prob=score[begin][split][B]*score[split][end][C]*P(A->BC)
    if prob > score[begin][end][A]
      score[begin][end][A] = prob
      back[begin][end][A] = new Triple(split,B,C)
    
```



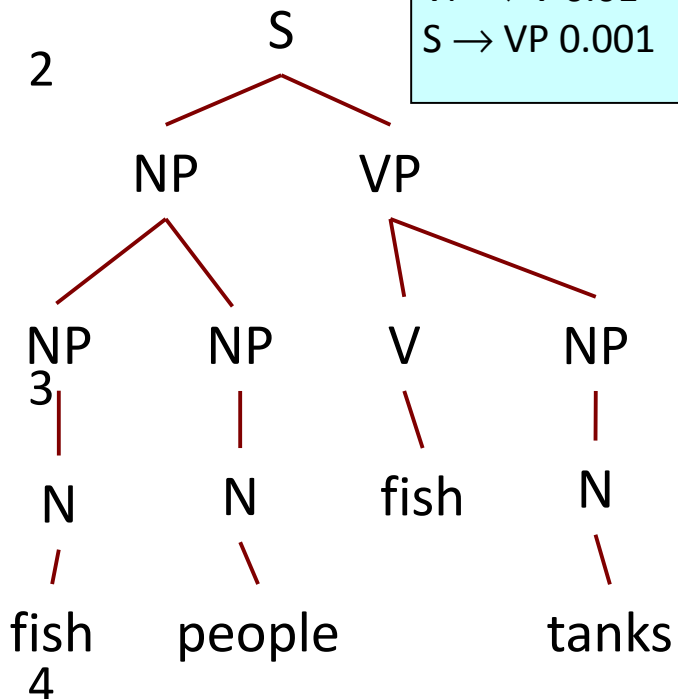
- S → NP VP 0.9
- S → VP 0.1
- VP → V NP 0.5
- VP → V 0.1
- VP → V @VP\_V 0.3
- VP → V PP 0.1
- @VP\_V → NP PP 1.0
- NP → NP NP 0.1
- NP → NP PP 0.2
- NP → N 0.7
- PP → P NP 1.0
- N → *people* 0.5
- N → *fish* 0.2
- N → *tanks* 0.2
- N → *rods* 0.1
- V → *people* 0.1
- V → *fish* 0.6
- V → *tanks* 0.3
- P → *with* 1.0

	fish	1	people	2	fish	3	tanks	4
0	N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006		NP → NP NP 0.0049 VP → V NP 0.105 S → VP 0.0105		NP → NP NP 0.0000686 VP → V NP 0.00147 S → NP VP 0.000882	span(0,4)		
1			N → people 0.5 V → people 0.1 NP → N 0.35 VP → V 0.01 S → VP 0.001		NP → NP NP 0.0049 VP → V NP 0.007 S → NP VP 0.0189			
2					N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006		NP → NP NP 0.00196 VP → V NP 0.042 S → VP 0.0042	
3							N → tanks 0.2 V → tanks 0.1 NP → N 0.14 VP → V 0.03 S → VP 0.003	
4	for split = begin+1 to end-1 for A,B,C in nonterms prob=score[begin][split][B]*score[split][end][C]*P(A->BC) if prob > score[begin][end][A] score[begin][end][A] = prob back[begin][end][A] = new Triple(split,B,C)							



- S → NP VP 0.9
- S → VP 0.1
- VP → V NP 0.5
- VP → V 0.1
- VP → V @VP\_V 0.3
- VP → V PP 0.1
- @VP\_V → NP PP 1.0
- NP → NP NP 0.1
- NP → NP PP 0.2
- NP → N 0.7
- PP → P NP 1.0
- N → *people* 0.5
- N → *fish* 0.2
- N → *tanks* 0.2
- N → *rods* 0.1
- V → *people* 0.1
- V → *fish* 0.6
- V → *tanks* 0.3
- P → *with* 1.0

	fish 1	people 2	fish 3	tanks 4
0	N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006	NP → NP NP 0.0049 VP → V NP 0.105 S → VP 0.0105	NP → NP NP 0.0000686 VP → V NP 0.00147 S → NP VP 0.000882	NP → NP NP 0.0000009604 VP → V NP 0.00002058 S → NP VP 0.00018522
1		N → people 0.5 V → people 0.1 NP → N 0.35 VP → V 0.01 S → VP 0.001	NP → NP NP 0.0049 VP → V NP 0.007 S → NP VP 0.0189	NP → NP NP 0.0000686 VP → V NP 0.000098 S → NP VP 0.01323
2			N → fish 0.2 V → fish 0.6 NP → N 0.14 VP → V 0.06 S → VP 0.006	NP → NP NP 0.00196 VP → V NP 0.042 S → VP 0.0042
3				N → tanks 0.2 V → tanks 0.1 NP → N 0.14 VP → V 0.03 S → VP 0.003



Call buildTree(score, back) to get the best parse

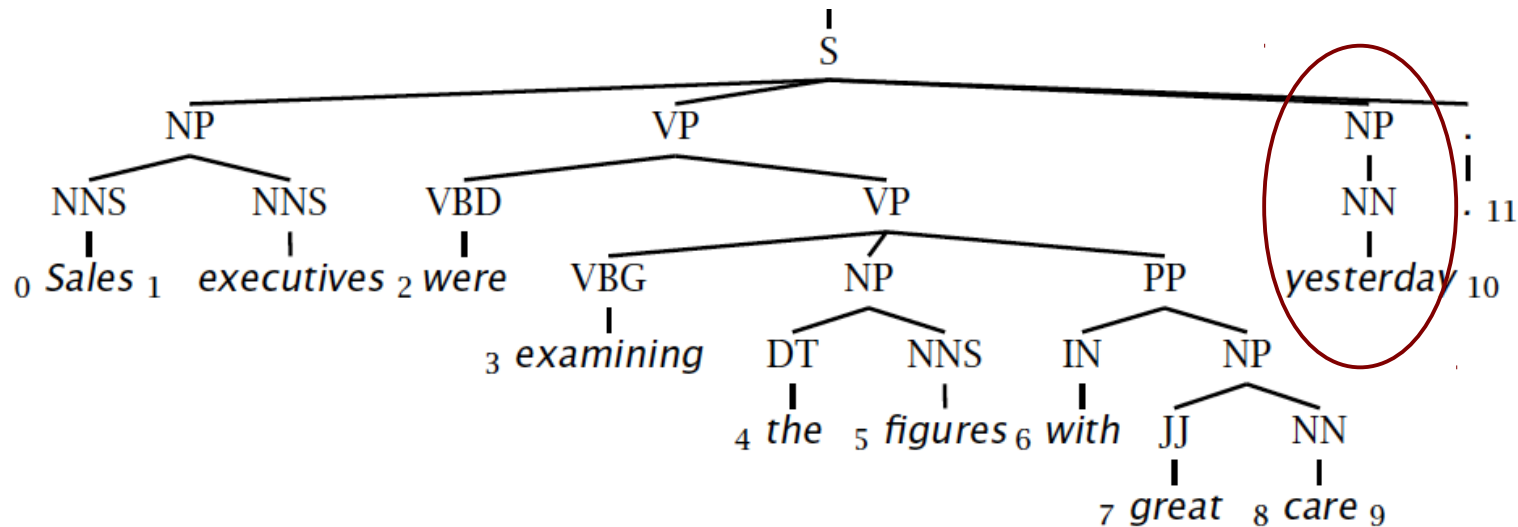






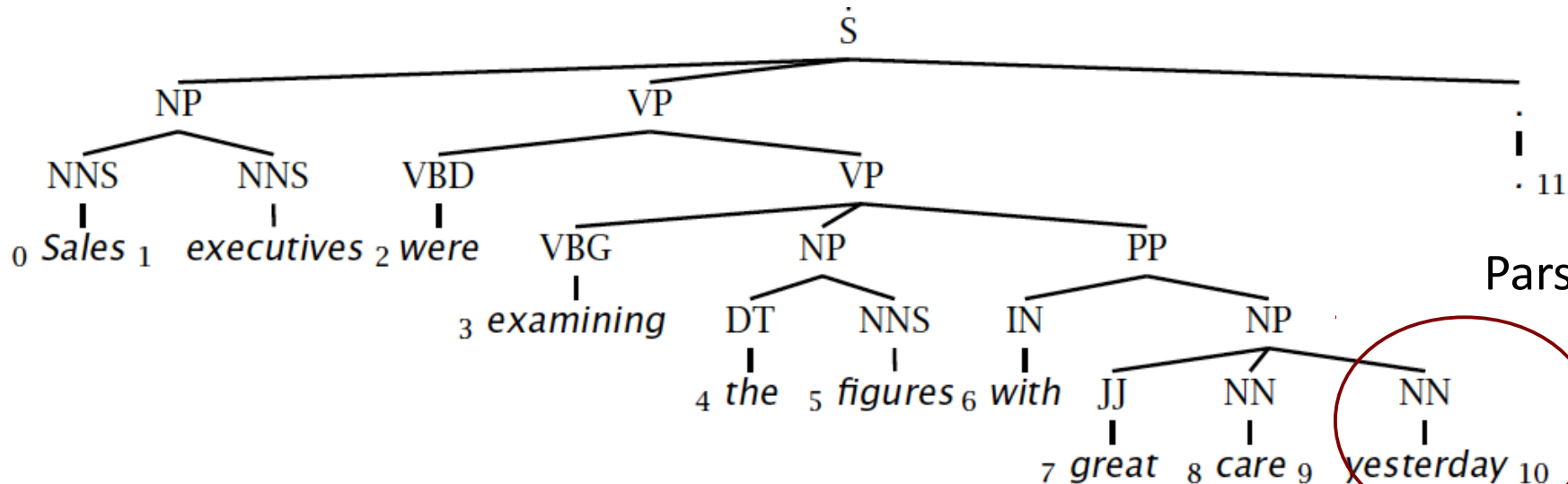
# Evaluating constituency parsing

Gold standard brackets: S-(0:11), NP-(0:2), VP-(2:9), VP-(3:9), NP-(4:6), PP-(6-9), NP-(7,9), NP-(9:10)



Gold standard

Candidate brackets: S-(0:11), NP-(0:2), VP-(2:10), VP-(3:10), NP-(4:6), PP-(6-10), NP-(7,10)



Parser Output



# Evaluating constituency parsing

## Gold standard brackets:

**S-(0:11)**, **NP-(0:2)**, VP-(2:9), VP-(3:9), **NP-(4:6)**, PP-(6-9), NP-(7,9), NP-(9:10)

## Candidate brackets:

**S-(0:11)**, **NP-(0:2)**, VP-(2:10), VP-(3:10), **NP-(4:6)**, PP-(6-10), NP-(7,10)

Labeled Precision	$3/7 = 42.9\%$
Labeled Recall	$3/8 = 37.5\%$
LP/LR F1	40.0%
POS Tagging Accuracy	$11/11 = 100.0\%$

One small mistake causes a big drop in the prec/recall!





# PCFG Training

- We saw how to transform an arbitrary PCFG into Chomsky Normal Form
- And we saw how to use a PCFG to get most probable parse of a sentence
- But, where do these PCFGs come from i.e. where to get the grammar **rules** and their **probabilities**?

# PCFG Training

- The answer is: **Treebanks**
- To train a PCFG:
  - Use all the hand-labeled rules in the treebank as grammar rules
  - The probability of a given rule is its relative count

$$P(N \rightarrow \xi) = \frac{\text{Count}(N \rightarrow \xi)}{\sum_{\gamma} \text{Count}(N \rightarrow \gamma)}$$

# Recap

- CKY Parsing Algorithm
- Worked example
- Parsing Evaluation
- PCFG Training