

# (Probabilistic) Context-Free Grammars

# A phrase structure grammar

$S \rightarrow NP VP$   
 $VP \rightarrow V NP$   
 $NP \rightarrow N$   
 $VP \rightarrow V NP PP$   
 $NP \rightarrow NP NP$   
 $NP \rightarrow NP PP$   
 $NP \rightarrow e$   
 $PP \rightarrow P NP$

***people fish tanks***

*people fish with rods*

$N \rightarrow \text{people}$

$V \rightarrow \text{fish}$

$N \rightarrow \text{fish}$

$N \rightarrow \text{tanks}$

$N \rightarrow \text{rods}$

$V \rightarrow \text{people}$

$V \rightarrow \text{tanks}$

$P \rightarrow \text{with}$

Ambiguous: People people people, fish fish fish

# Phrase structure grammars = context-free grammars (CFGs)

- $G = (T, N, S, R)$ 
  - $T$  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  $X \rightarrow \gamma$ 
    - $X \in N$  and  $\gamma \in (N \cup T)^*$
- A grammar  $G$  generates a language  $L$ .

# Phrase structure grammars in NLP

- $G = (T, C, N, S, L, R)$ 
  - $T$  is a set of terminal symbols
  - $C$  is a set of preterminal symbols
  - $N$  is a set of nonterminal symbols
  - $S$  is the start symbol ( $S \in N$ )
  - $L$  is the lexicon, a set of items of the form  $X \rightarrow x$ 
    - $X \in C$  and  $x \in T$
  - $R$  is the grammar, a set of items of the form  $X \rightarrow \gamma$ 
    - $X \in N$  and  $\gamma \in (N \cup C)^*$
- By usual convention,  $S$  is the start symbol, but in statistical NLP, we usually have an extra node at the top (ROOT, TOP)
- We usually write  $e$  for an empty sequence, rather than nothing

# A phrase structure grammar (empty, unary, binary)

## Grammar Rules

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$VP \rightarrow V NP PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$NP \rightarrow N$

$NP \rightarrow e$

$PP \rightarrow P NP$

*EMPTY fish tanks*

*people fish EMPTY*

## Lexicon

$N \rightarrow \textit{people}$

$N \rightarrow \textit{fish}$

$N \rightarrow \textit{tanks}$

$N \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$P \rightarrow \textit{with}$

# Probabilistic/stochastic – 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
    - $P: R \rightarrow [0,1]$
    - $$\forall X \in N, \sum_{X \rightarrow \gamma \in R} P(X \rightarrow \gamma) = 1$$

# A PCFG

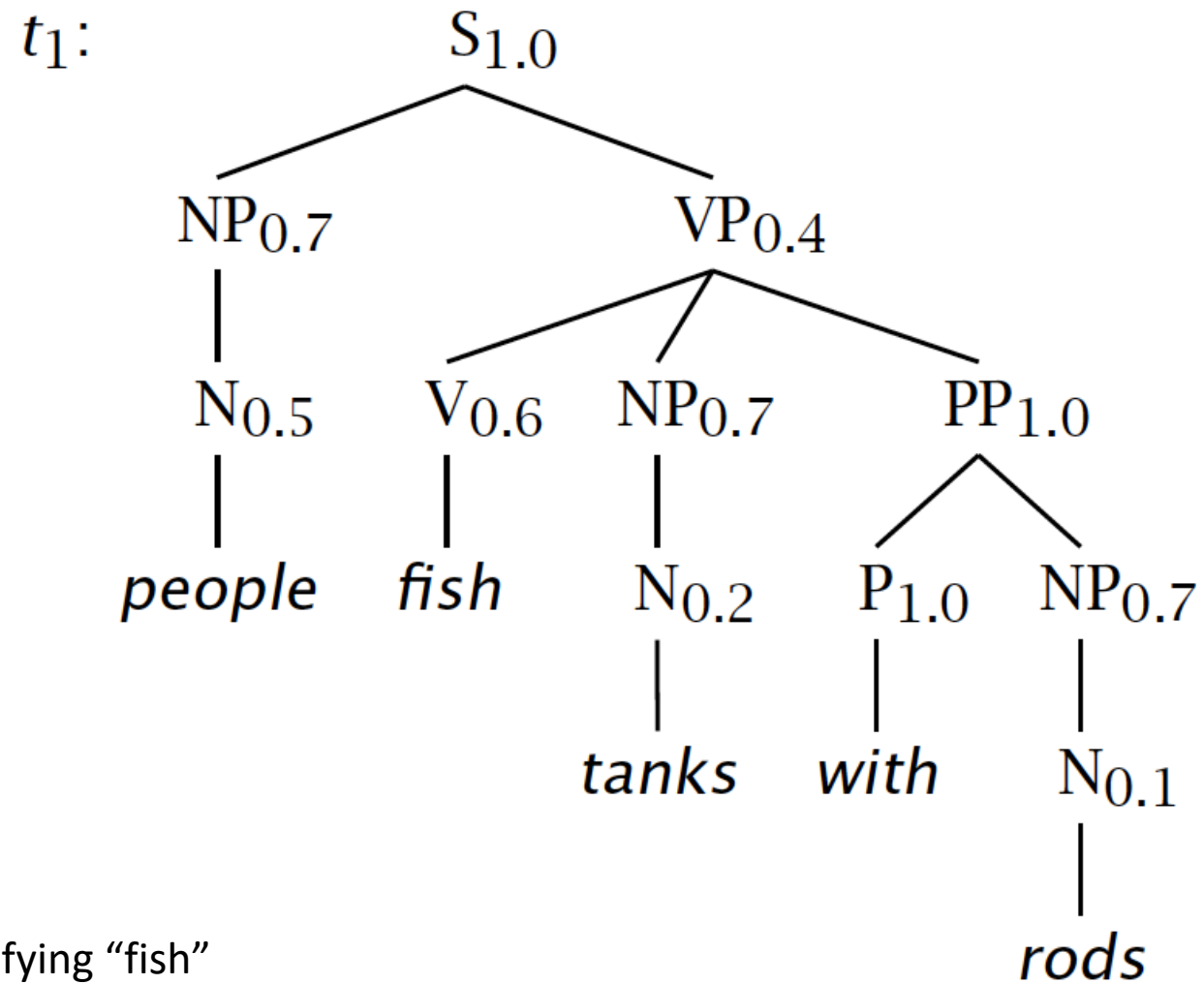
$S \rightarrow NP VP$	1.0	$N \rightarrow people$	0.5
$VP \rightarrow V NP$	0.6	$N \rightarrow fish$	0.2
$VP \rightarrow V NP PP$	0.4	$N \rightarrow tanks$	0.2
$NP \rightarrow NP NP$	0.1	$N \rightarrow rods$	0.1
$NP \rightarrow NP PP$	0.2	$V \rightarrow people$	0.1
$NP \rightarrow N$	0.7	$V \rightarrow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

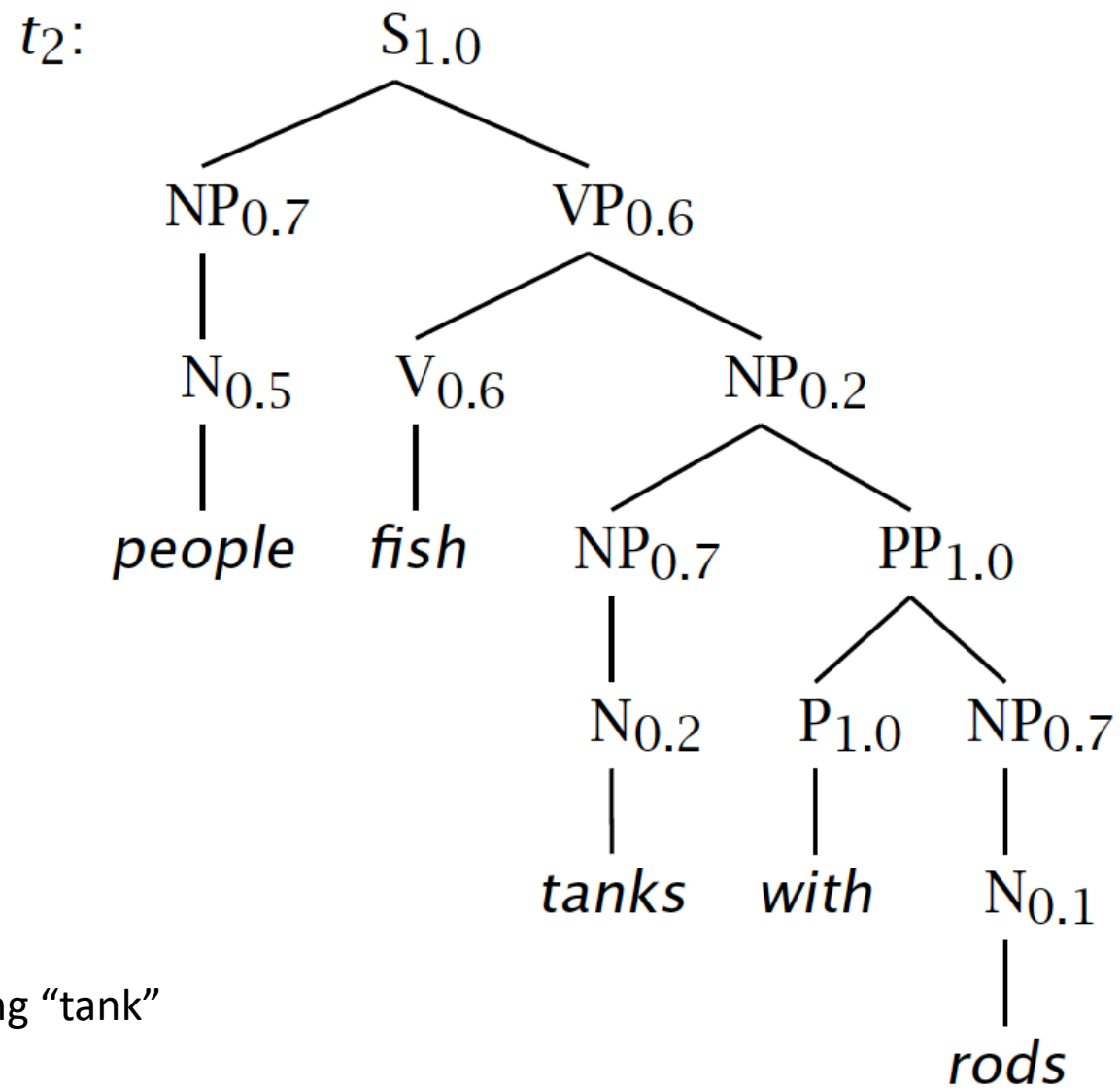
# The probability of trees and strings

- $P(t)$  – 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

$$\begin{aligned} P(s) &= \sum_t P(s, t) \text{ where } t \text{ is a parse of } s \\ &= \sum_t P(t) \end{aligned}$$







Preposition “with” modifying “tank”

# Tree and String Probabilities

- $s = \textit{people fish tanks with rods}$

- $P(t_1) = 1.0 \times 0.7 \times 0.4 \times 0.5 \times 0.6 \times 0.7$   
 $\times 1.0 \times 0.2 \times 1.0 \times 0.7 \times 0.1$   
 $= 0.0008232$

Verb attach

- $P(t_2) = 1.0 \times 0.7 \times 0.6 \times 0.5 \times 0.6 \times 0.2$   
 $\times 0.7 \times 1.0 \times 0.2 \times 1.0 \times 0.7 \times 0.1$   
 $= 0.00024696$

Noun attach

- $P(s) = P(t_1) + P(t_2)$   
 $= 0.0008232 + 0.00024696$   
 $= 0.00107016$

# Grammar Transforms

Restricting the grammar form for efficient parsing

# 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
- Empties and unaries are removed recursively
- $n$ -ary rules are divided by introducing new nonterminals ( $n > 2$ )

# A phrase structure grammar

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$VP \rightarrow V NP PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$NP \rightarrow N$

$NP \rightarrow e$

$PP \rightarrow P NP$

$N \rightarrow \textit{people}$

$N \rightarrow \textit{fish}$

$N \rightarrow \textit{tanks}$

$N \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$P \rightarrow \textit{with}$

Start discussing epsilon removal

# Chomsky Normal Form steps

$S \rightarrow NP VP$

$S \rightarrow VP$

$VP \rightarrow V NP$

$VP \rightarrow V$

$VP \rightarrow V NP PP$

$VP \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP$

$NP \rightarrow NP PP$

$NP \rightarrow PP$

$NP \rightarrow N$

$PP \rightarrow P NP$

$PP \rightarrow P$

$N \rightarrow \textit{people}$

$N \rightarrow \textit{fish}$

$N \rightarrow \textit{tanks}$

$N \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$P \rightarrow \textit{with}$

Start discussing unary removal downwards: remove  $S \rightarrow VP$

# Chomsky Normal Form steps

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V$

$S \rightarrow V$

$VP \rightarrow V NP PP$

$S \rightarrow V NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP$

$NP \rightarrow NP PP$

$NP \rightarrow PP$

$NP \rightarrow N$

$PP \rightarrow P NP$

$PP \rightarrow P$

$N \rightarrow \textit{people}$

$N \rightarrow \textit{fish}$

$N \rightarrow \textit{tanks}$

$N \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$P \rightarrow \textit{with}$

Remove more unaries, next  $S \rightarrow V$



# Chomsky Normal Form steps

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V$

$VP \rightarrow V NP PP$

$S \rightarrow V NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP$

$NP \rightarrow NP PP$

$NP \rightarrow PP$

$NP \rightarrow N$

$PP \rightarrow P NP$

$PP \rightarrow P$

$N \rightarrow \textit{people}$

$N \rightarrow \textit{fish}$

$N \rightarrow \textit{tanks}$

$N \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$S \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$S \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$S \rightarrow \textit{tanks}$

$P \rightarrow \textit{with}$

After remove  $S \rightarrow V$  get this, and then do  $VP \rightarrow V$

# Chomsky Normal Form steps

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V NP PP$

$S \rightarrow V NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP$

$NP \rightarrow NP PP$

$NP \rightarrow PP$

**$NP \rightarrow N$**

$PP \rightarrow P NP$

$PP \rightarrow P$

$N \rightarrow \textit{people}$

$N \rightarrow \textit{fish}$

$N \rightarrow \textit{tanks}$

$N \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$S \rightarrow \textit{people}$

$VP \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$S \rightarrow \textit{fish}$

$VP \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$S \rightarrow \textit{tanks}$

$VP \rightarrow \textit{tanks}$

$P \rightarrow \textit{with}$

# Chomsky Normal Form steps

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V NP PP$

$S \rightarrow V NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$NP \rightarrow P NP$

$PP \rightarrow P NP$

$NP \rightarrow \textit{people}$

$NP \rightarrow \textit{fish}$

$NP \rightarrow \textit{tanks}$

$NP \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$S \rightarrow \textit{people}$

$VP \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$S \rightarrow \textit{fish}$

$VP \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$S \rightarrow \textit{tanks}$

$VP \rightarrow \textit{tanks}$

$P \rightarrow \textit{with}$

$PP \rightarrow \textit{with}$

And then binarize

# Chomsky Normal Form steps

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V @VP\_V$

$@VP\_V \rightarrow NP PP$

$S \rightarrow V @S\_V$

$@S\_V \rightarrow NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$NP \rightarrow P NP$

$PP \rightarrow P NP$

$NP \rightarrow \textit{people}$

$NP \rightarrow \textit{fish}$

$NP \rightarrow \textit{tanks}$

$NP \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$S \rightarrow \textit{people}$

$VP \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$S \rightarrow \textit{fish}$

$VP \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$S \rightarrow \textit{tanks}$

$VP \rightarrow \textit{tanks}$

$P \rightarrow \textit{with}$

$PP \rightarrow \textit{with}$

# A phrase structure grammar

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$VP \rightarrow V NP PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$NP \rightarrow N$

$NP \rightarrow e$

$PP \rightarrow P NP$

$N \rightarrow \textit{people}$

$N \rightarrow \textit{fish}$

$N \rightarrow \textit{tanks}$

$N \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$P \rightarrow \textit{with}$

# Chomsky Normal Form steps

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$S \rightarrow V NP$

$VP \rightarrow V @VP\_V$

$@VP\_V \rightarrow NP PP$

$S \rightarrow V @S\_V$

$@S\_V \rightarrow NP PP$

$VP \rightarrow V PP$

$S \rightarrow V PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$NP \rightarrow P NP$

$PP \rightarrow P NP$

$NP \rightarrow \textit{people}$

$NP \rightarrow \textit{fish}$

$NP \rightarrow \textit{tanks}$

$NP \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$S \rightarrow \textit{people}$

$VP \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$S \rightarrow \textit{fish}$

$VP \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$S \rightarrow \textit{tanks}$

$VP \rightarrow \textit{tanks}$

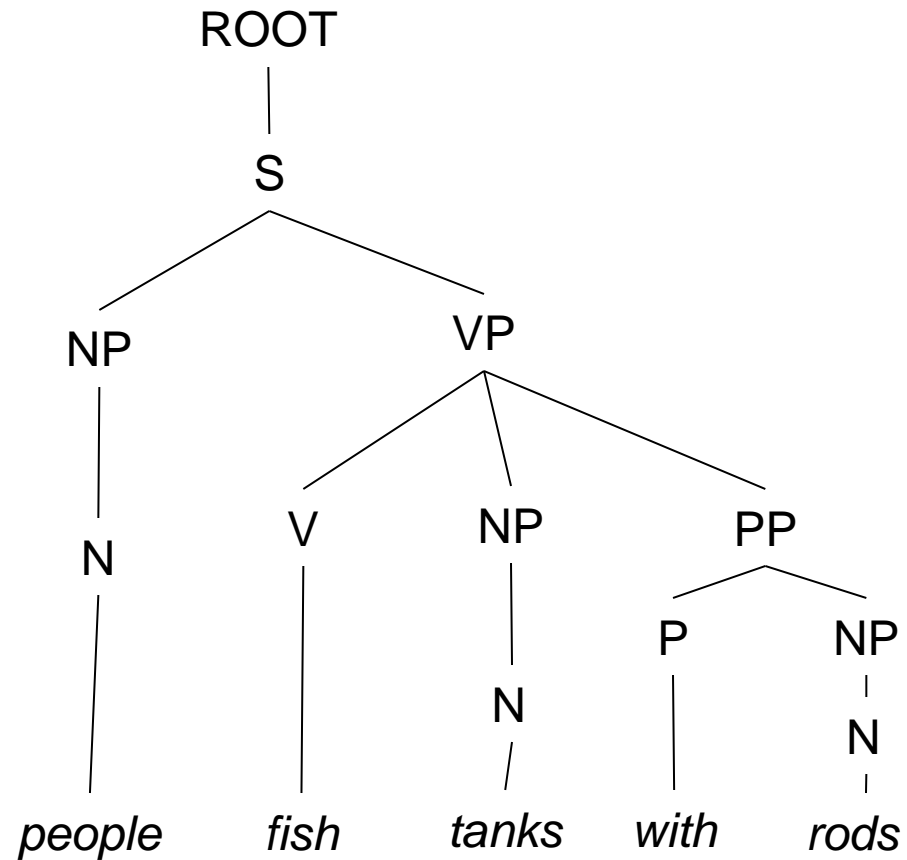
$P \rightarrow \textit{with}$

$PP \rightarrow \textit{with}$

# Chomsky Normal Form

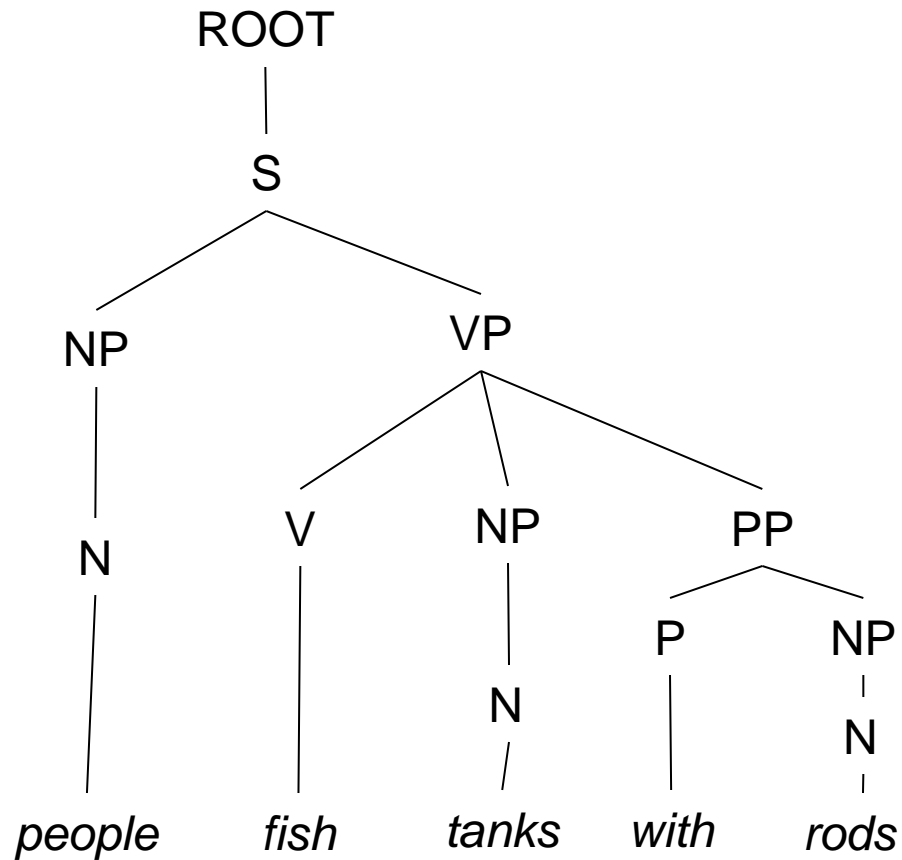
- You should think of this as a transformation for efficient parsing
- With some extra book-keeping in symbol names, you can even reconstruct the same trees with a detransform
- In practice full Chomsky Normal Form is a pain
  - Reconstructing n-aries is easy
  - Reconstructing unaries/empties is trickier
- **Binarization is crucial for cubic time CFG parsing**
- The rest isn't necessary; it just makes the algorithms cleaner and a bit quicker

## An example: before binarization...

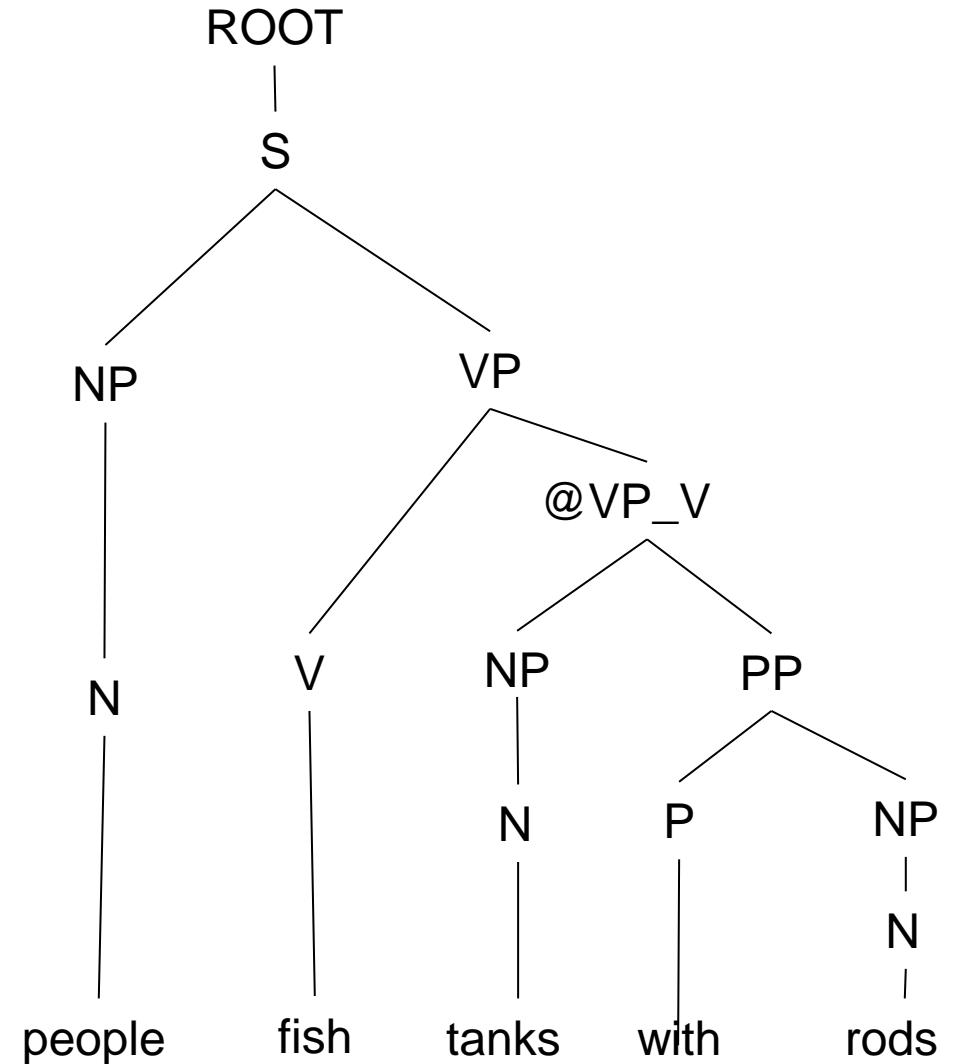




## An example: before binarization...

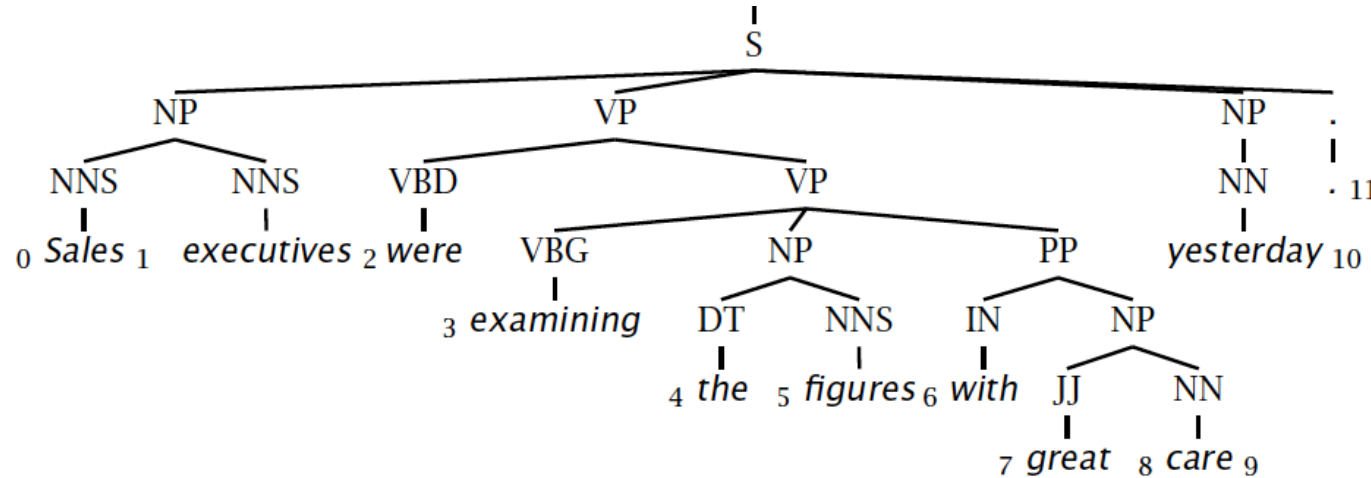


## After binarization...

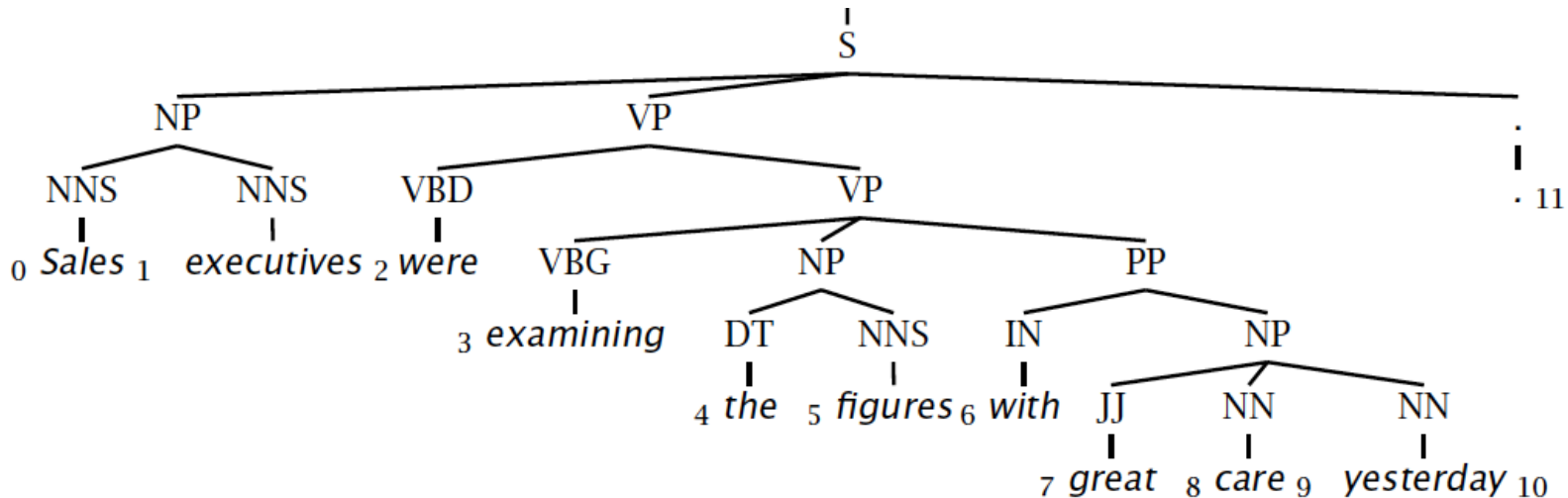


# Evaluating constituency parsing

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)



# 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%
Tagging Accuracy	$11/11 = 100.0\%$