



# CFGs and PCFGs

# (Probabilistic) Context-Free Grammars



# A phrase structure grammar

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$VP \rightarrow V NP PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$NP \rightarrow N$

$NP \rightarrow e$

$PP \rightarrow P NP$

$N \rightarrow \text{people}$

$N \rightarrow \text{fish}$

$N \rightarrow \text{tanks}$

$N \rightarrow \text{rods}$

$V \rightarrow \text{people}$

$V \rightarrow \text{fish}$

$V \rightarrow \text{tanks}$

$P \rightarrow \text{with}$

*people fish tanks*

*people fish with rods*



# 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  $X \rightarrow x$ 
    - $X \in P$  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

$S \rightarrow NP VP$

$VP \rightarrow V NP$

$VP \rightarrow V NP PP$

$NP \rightarrow NP NP$

$NP \rightarrow NP PP$

$NP \rightarrow N$

$NP \rightarrow e$

$PP \rightarrow P NP$

$N \rightarrow \textit{people}$

$N \rightarrow \textit{fish}$

$N \rightarrow \textit{tanks}$

$N \rightarrow \textit{rods}$

$V \rightarrow \textit{people}$

$V \rightarrow \textit{fish}$

$V \rightarrow \textit{tanks}$

$P \rightarrow \textit{with}$

*people fish tanks*

*people fish with rods*



# Probabilistic – or 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 grammar  $G$  generates a language model  $L$ .

$$\sum_{\gamma \in T^*} P(\gamma) = 1$$



# A PCFG

$S \rightarrow NP VP$  1.0

$VP \rightarrow V NP$  0.6

$VP \rightarrow V NP PP$  0.4

$NP \rightarrow NP NP$  0.1

$NP \rightarrow NP PP$  0.2

$NP \rightarrow N$  0.7

$PP \rightarrow P NP$  1.0

$N \rightarrow \textit{people}$  0.5

$N \rightarrow \textit{fish}$  0.2

$N \rightarrow \textit{tanks}$  0.2

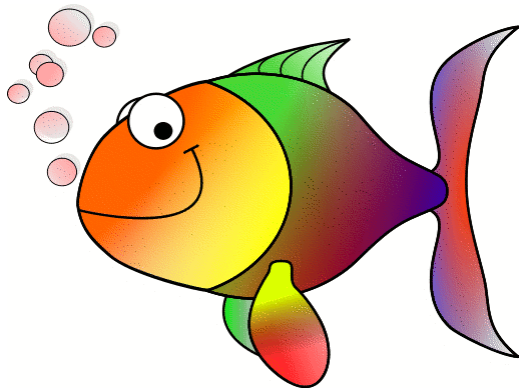
$N \rightarrow \textit{rods}$  0.1

$V \rightarrow \textit{people}$  0.1

$V \rightarrow \textit{fish}$  0.6

$V \rightarrow \textit{tanks}$  0.3

$P \rightarrow \textit{with}$  1.0



[With empty NP removed  
so less ambiguous]

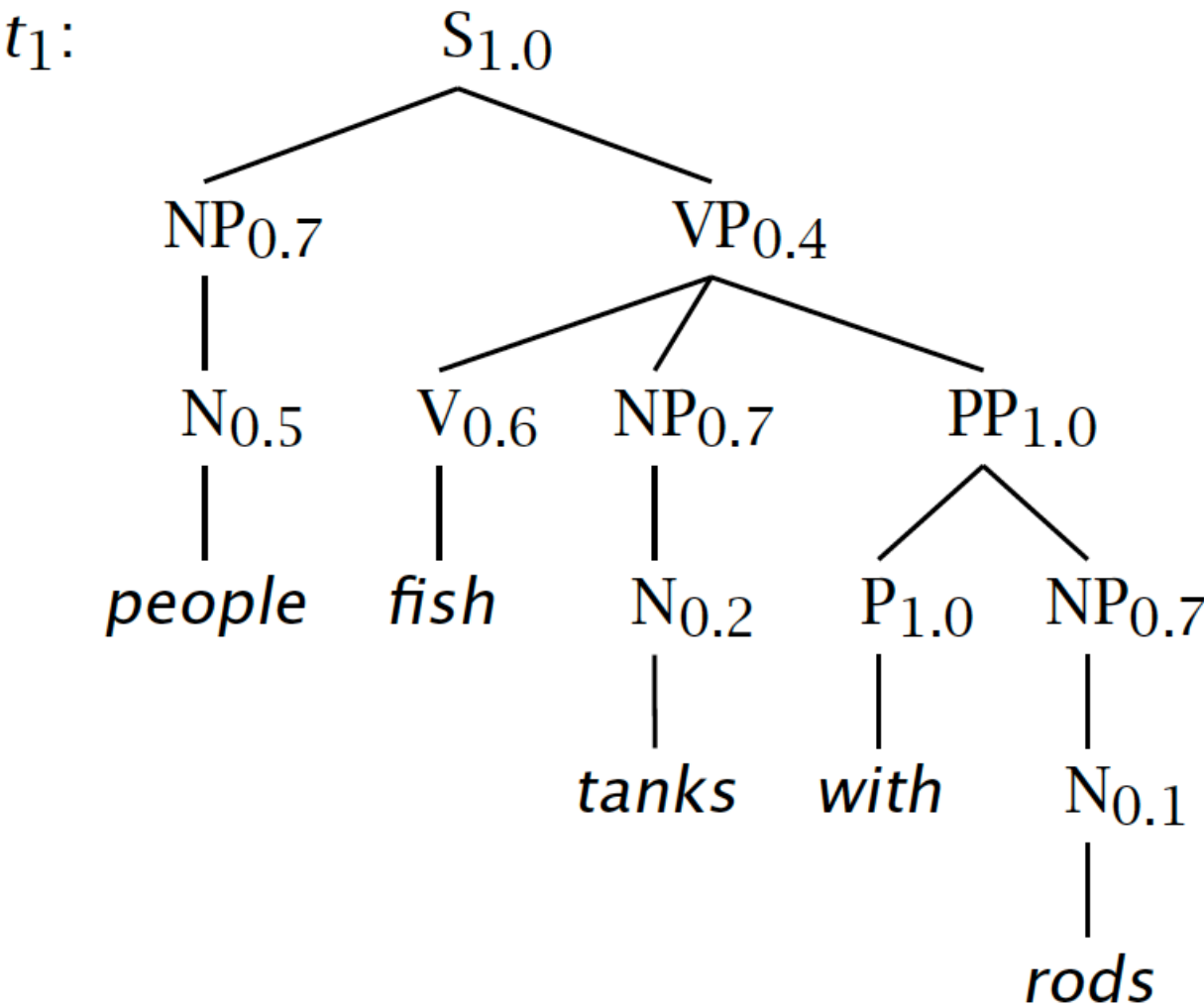


# The probability of trees and strings

- $P(t)$  – 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

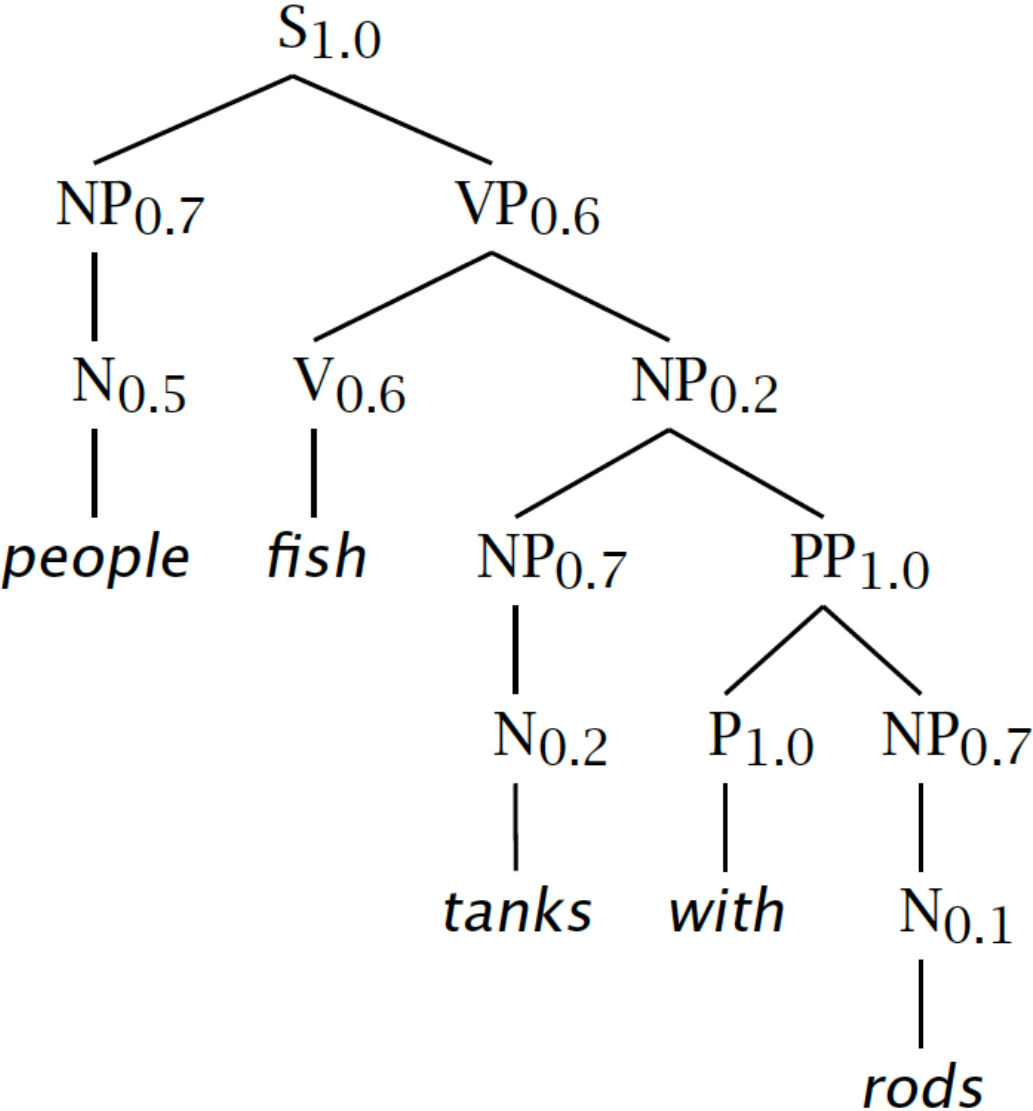
$$\begin{aligned} P(s) &= \sum_j P(s, t) \text{ where } t \text{ is a parse of } s \\ &= \sum_j P(t) \end{aligned}$$







$t_2$ :



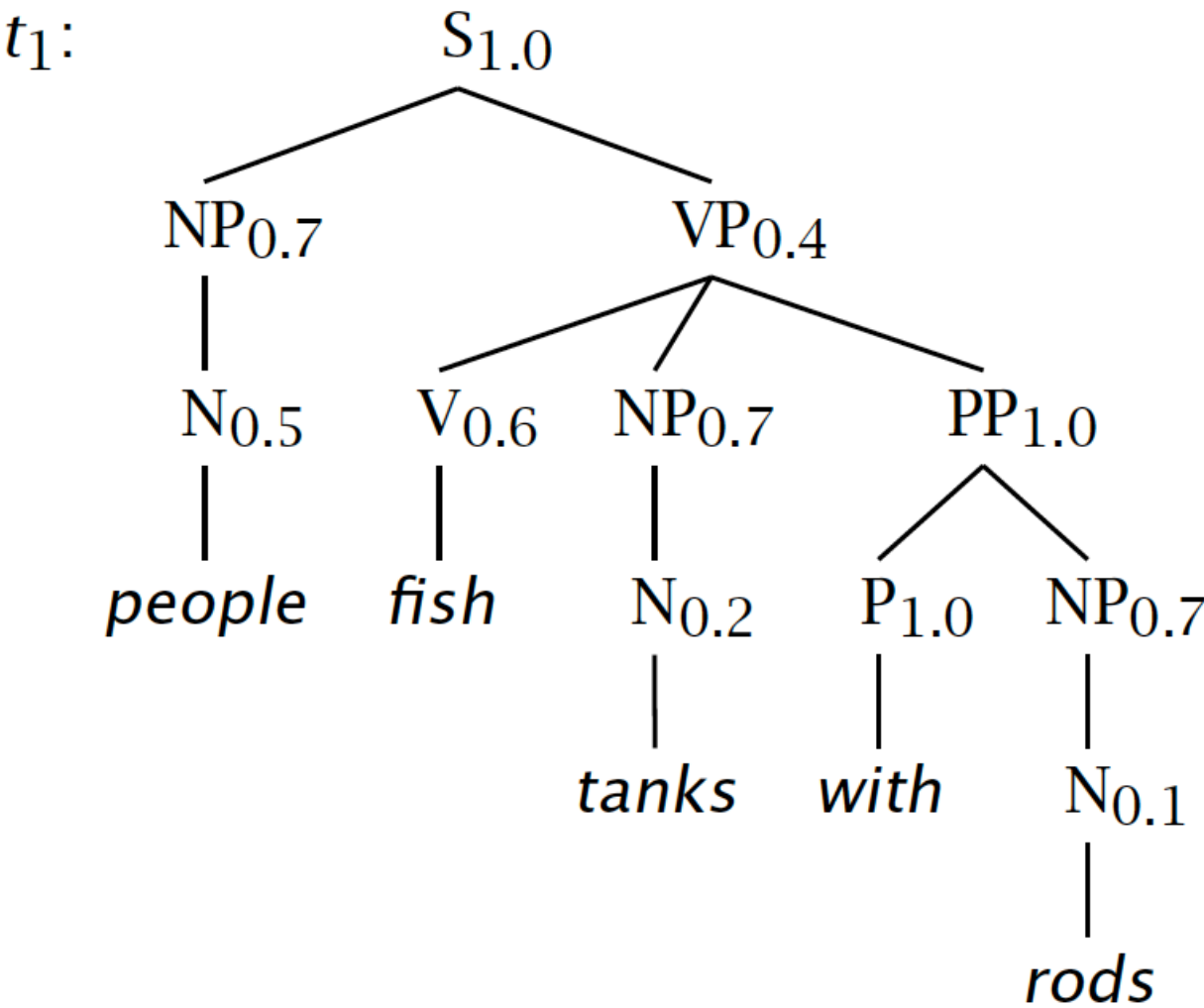


# Tree and String Probabilities

- $s = \textit{people fish tanks with rods}$
- $P(t_1) = 1.0 \times 0.7 \times 0.4 \times 0.5 \times 0.6 \times 0.7$   
 $\times 1.0 \times 0.2 \times 1.0 \times 0.7 \times 0.1$   
 $= 0.0008232$
- $P(t_2) = 1.0 \times 0.7 \times 0.6 \times 0.5 \times 0.6 \times 0.2$   
 $\times 0.7 \times 1.0 \times 0.2 \times 1.0 \times 0.7 \times 0.1$   
 $= 0.00024696$
- $P(s) = P(t_1) + P(t_2)$   
 $= 0.0008232 + 0.00024696$   
 $= 0.00107016$

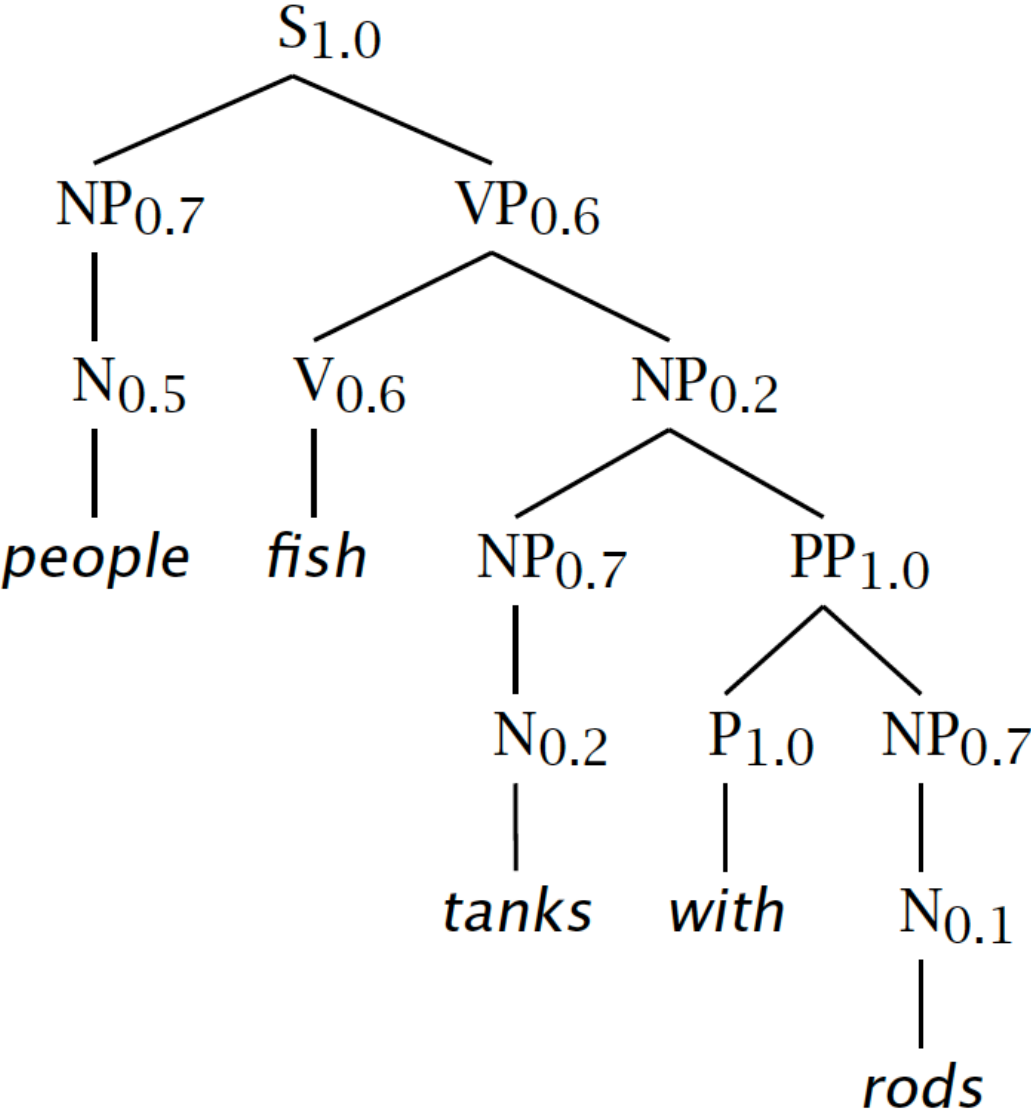
Verb attach

Noun attach





$t_2$ :





# CFGs and PCFGs

# (Probabilistic) Context-Free Grammars