# Context-Free Parsing

He He

New York University

October 27, 2020

# Logistics

- Homework 3
- Project proposal and group
- Project presentation

# Table of Contents

# Langauge is a set of strings

**Formal language**:

▶ A set of **strings** consisting of **words** from an **alphabet**

▶ Well-formed according to a set of rules

▶ Studies the syntactical aspects of a language

Examples:

▶ Formulas (logic): $(p_1 \wedge p_2) \vee (\neg p_3)$   $(( P_1$

▶ Programming languages: `int a, b = 0;`

▶ Sequences from the alphabet $\{a, b\}$ that ends with two $a$'s

Questions:

▶ Formal language theory: expressiveness power, recognizability etc.

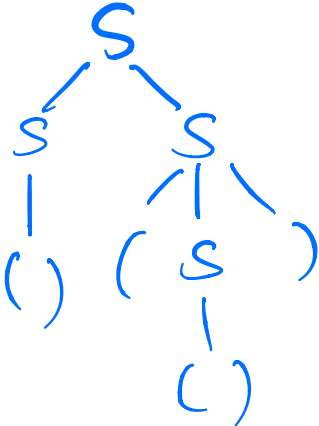▶ Can we design formal languages that capture as many properties of natural language as possible?

# Context-free language

**Context-free languages (CFL)** are generated by a **context-free grammar** $G = (\Sigma, N, R, S)$:

- a finite alphabet $\Sigma$ of **terminals** (words)
- a finite set of **non-terminals** $N$ disjoint from $\Sigma$ (word groups)
- a set of **production rules** $R$ of the form $A \to \beta$, where $A \in N, \beta \in (\Sigma \cup N)^*$ (how to group words)
- a start symbol $S \in N$ (root of derivation)

Example:

$S \to SS$
$S \to (S)$
$S \to ()$



string:

$() (())$

# Natural language syntax

Construct a formal language to represent the syntax of natural language

- ▶ Expressivity: how many syntactic phenomena can it cover?
- ▶ Computation: how fast can we parse a sentence?

Context-free grammars for natural language

- ▶ Captures nested structures which are common in natural language

    [I told Mary that [John told Jane that [Ted told Tom a secret]]].

- ▶ Captures long-range dependencies

    the burnt and badly-ground Italian coffee
    these burnt and badly-ground Italian coffees

- ▶ Strikes a good balance between expressivity and computation

# Phrase-structure grammar for English

Sentences are broken down into **constituents**.

A constituent works as a single unit in a sentence.

▶ Can be moved around or replaced without breaking grammaticality.
(Abigail) and (her younger brother) (bought a fish).

  *VP*

Construct CFG for English

▶ Each word is a terminal, derived from its POS tag.

▶ Each sentence is derived from the start symbol $S$.

▶ Each phrase type is a non-terminal.

▶ Each constituent is derived from a non-terminal.

Grammar design: choose the right set of non-terminals that produces different constituents.

# A toy example CFG

$N = \{S, NP, VP, PP, DT, Vi, Vt, NN, IN\}$
$S = S$
$\Sigma = \{sleeps, saw, man, woman, dog, telescope, the, with, in\}$

$R =$

| S | → | NP | VP |
|---|---|-----|-----|
| VP | → | Vi | |
| VP | → | Vt | NP |
| VP | → | VP | PP |
| NP | → | DT | NN |
| NP | → | NP | PP |
| PP | → | IN | NP |

| Vi | → | sleeps |
|----|---|--------|
| Vt | → | saw |
| NN | → | man |
| NN | → | woman |
| NN | → | telescope |
| NN | → | dog |
| DT | → | the |
| IN | → | with |
| IN | → | in |

``grammar''

Lexicon

**Lexicon**: rules that produce the terminals

# Parsing

$[_S [_{VP} [_{DT} \text{the}]\ [_{NN} \text{man}]][_{VP} [_{Vi} \text{sleeps}]]]$

$R =$

| S | $\rightarrow$ | NP | VP |
|---|---|---|---|
| VP | $\rightarrow$ | Vi | |
| VP | $\rightarrow$ | Vt | NP |
| VP | $\rightarrow$ | VP | PP |
| NP | $\rightarrow$ | DT | NN |
| NP | $\rightarrow$ | NP | PP |
| PP | $\rightarrow$ | IN | NP |

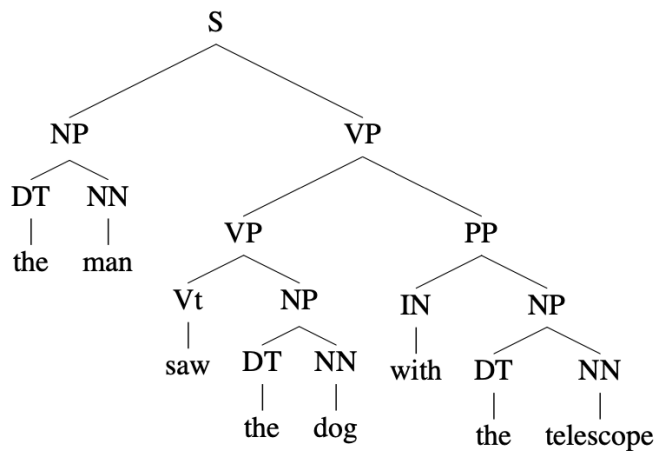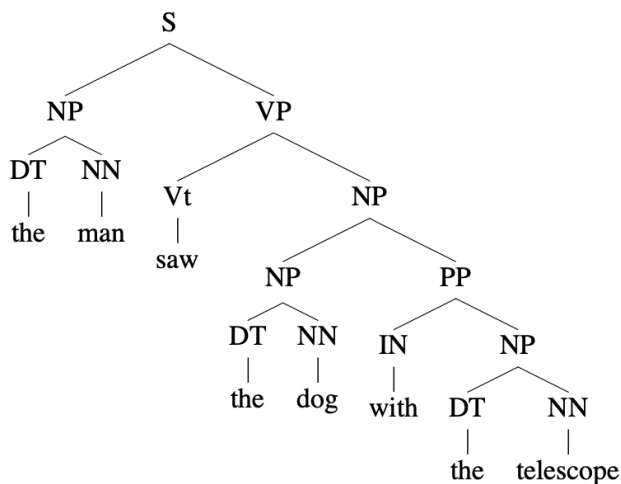| Vi | $\rightarrow$ | sleeps |
|---|---|---|
| Vt | $\rightarrow$ | saw |
| NN | $\rightarrow$ | man |
| NN | $\rightarrow$ | woman |
| NN | $\rightarrow$ | telescope |
| NN | $\rightarrow$ | dog |
| DT | $\rightarrow$ | the |
| IN | $\rightarrow$ | with |
| IN | $\rightarrow$ | in |

Can we derive the sentence "the man sleeps"?

$S \rightarrow NP\ VP$
$\rightarrow DT\ NN\ VP$
$\rightarrow \text{the } NN\ VP$
$\rightarrow \text{the man } VP$
$\rightarrow \text{the man } Vi$
$\rightarrow \text{the man sleeps}$

# Ambiguity

Can a sentence have multiple parse trees?



Exercise: find parse trees for
"She announced a program to promote safety in trucks and vans".

# Table of Contents

# PCFG

Notation: let $\mathcal{T}_G$ be the set of all possible left-most parse trees under the grammar $G$.

Goal: define a probability distribution $p(t)$ over parse trees $t \in \mathcal{T}_G$

Parsing: pick the most likely parse tree for a sentence $s$

$$\arg\max_{t \in \mathcal{T}_G(s)} p(t)$$

Three questions:

▶ Modeling: how to define $p(t)$ for trees?

▶ Learning: how to estimate parameters of the distribution $p(t)$?
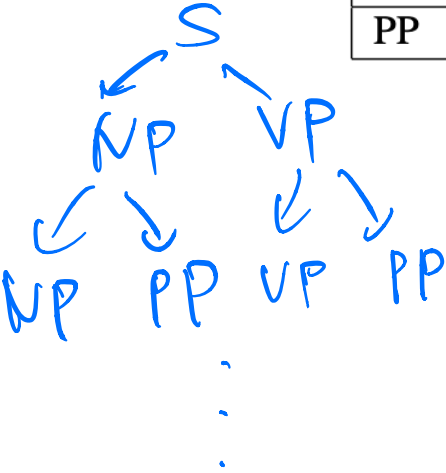
▶ Inference: how to find the most likely tree efficiently?

# Modeling

Generate parse trees: iteratively sample a production rule to expand a non-terminal

$R =$

| S | $\rightarrow$ | NP | VP |
|---|---|---|---|
| VP | $\rightarrow$ | Vi | |
| VP | $\rightarrow$ | Vt | NP |
| VP | $\rightarrow$ | VP | PP |
| NP | $\rightarrow$ | DT | NN |
| NP | $\rightarrow$ | NP | PP |
| PP | $\rightarrow$ | IN | NP |

| Vi | $\rightarrow$ | sleeps |
|---|---|---|
| Vt | $\rightarrow$ | saw |
| NN | $\rightarrow$ | man |
| NN | $\rightarrow$ | woman |
| NN | $\rightarrow$ | telescope |
| NN | $\rightarrow$ | dog |
| DT | $\rightarrow$ | the |
| IN | $\rightarrow$ | with |
| IN | $\rightarrow$ | in |

# PCFG

A **PCFG** consists of

▶ A CFG $G = (\Sigma, N, R, S)$

▶ Probabilities of production rules $q(\alpha \to \beta)$ for each $\alpha \to \beta \in R$ such that

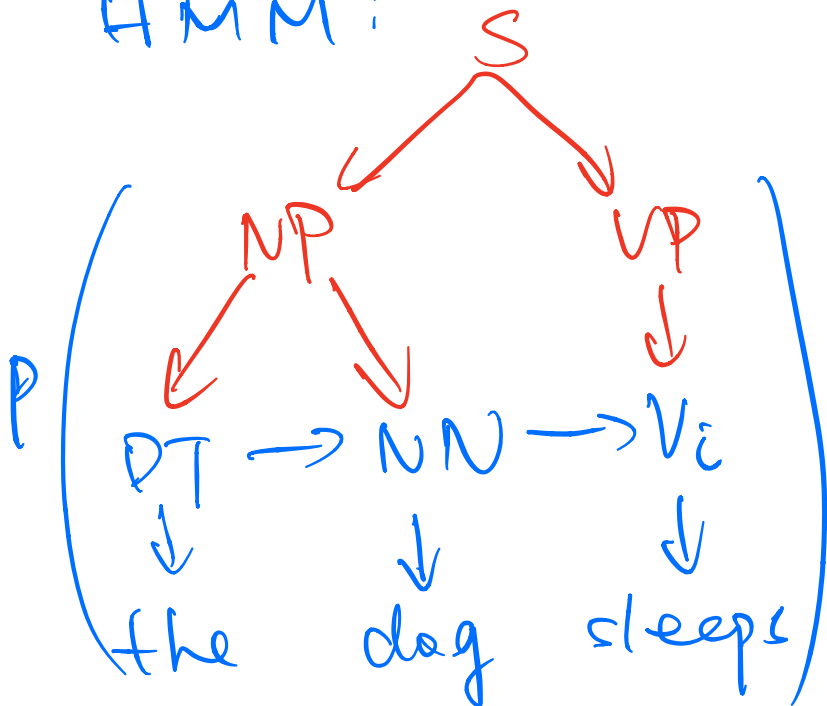$$\sum_{\beta:\, X \to \beta \in R} q(X \to \beta) = 1 \quad \forall X \in N$$

$R, q =$

| S | $\to$ | NP | VP | 1.0 |
|---|---|---|---|---|
| VP | $\to$ | Vi | | 0.3 |
| VP | $\to$ | Vt | NP | 0.5 |
| VP | $\to$ | VP | PP | 0.2 |
| NP | $\to$ | DT | NN | 0.8 |
| NP | $\to$ | NP | PP | 0.2 |
| PP | $\to$ | IN | NP | 1.0 |

| Vi | $\to$ | sleeps | 1.0 |
|---|---|---|---|
| Vt | $\to$ | saw | 1.0 |
| NN | $\to$ | man | 0.1 |
| NN | $\to$ | woman | 0.1 |
| NN | $\to$ | telescope | 0.3 |
| NN | $\to$ | dog | 0.5 |
| DT | $\to$ | the | 1.0 |
| IN | $\to$ | with | 0.6 |
| IN | $\to$ | in | 0.4 |

# From HMM to PCFG

HMM:

$$P\left( \begin{array}{c} \text{S} \\ \text{NP} \quad \text{VP} \\ \text{DT} \rightarrow \text{NN} \rightarrow V_i \\ \text{the} \quad \text{dog} \quad \text{sleeps} \end{array} \right)$$

$P(\text{the} \mid DT) \cdots$

$= P(DT, NN \mid NP)$

$P(V_i \mid VP)$

$P(VP, NP \mid S) \, P(S)$

# Probabilities of parse trees

Given a parse tree $t$ consisting of rules $\alpha_1 \to \beta_1, \ldots, \alpha_n \to \beta_n$, its probabilities under the PCFG is

$$p(t) = \prod_{i=1}^{n} q(\alpha_i \to \beta_i)$$

Example:

# Learning $P(S \to NP\ VP)$

Given a set of trees, we can estimate rule probabilities by MLE.

# Learning

Given a set of trees, we can estimate rule probabilities by MLE.

$$q(\alpha \to \beta) = \frac{\text{count}(\alpha \to \beta)}{\sum_{\beta' \,:\, \alpha \to \beta' \in R} \text{count}(\alpha \to \beta')}$$

Training data: treebanks

S → NP VP .          S → NP VP

```
((S
   (NP-SBJ (DT That)
    (JJ cold) (, ,)
    (JJ empty) (NN sky) )      ((S
   (VP (VBD was)                   (NP-SBJ The/DT flight/NN )
    (ADJP-PRD (JJ full)            (VP should/MD
     (PP (IN of)                    (VP arrive/VB
      (NP (NN fire)                  (PP-TMP at/IN
       (CC and)                       (NP eleven/CD a.m/RB ))
       (NN light) ))))                (NP-TMP tomorrow/NN )))))
   (. .) ))
         (a)                                 (b)
```

**Figure 12.7** Parsed sentences from the LDC Treebank3 version of the Brown (a) and ATIS (b) corpora.

# Parsing

Input: sentences, (P)CFG
Output: derivations / parse trees (with scores/probabilities)

Total number of parse trees for a sentence?

Consider a minimal CFG:

$X \rightarrow XX$

$X \rightarrow$ aardvark|abacus| . . . |zyther

# of parse trees = # of strings with balanced brackets

$((w_1 w_2)(w_3 w_4))$, $(((w_1 w_2)w_3)w_4)$, ...

# of strings with $n$ pairs of brackets:

$$\text{Catalan number } C_n = \frac{1}{n+1}\binom{2n}{n}$$

# Chomsky normal form (CNF)

A CFG is in **Chomsky normal form** if every production rule takes one of the following forms:

- ▶ Binary non-terminal production: $A \to BC$ where $A, B, C \in N$.
- ▶ Unary terminal production: $A \to a$ where $A \in N, a \in \Sigma$.

Grammars in CNF produces binary parse trees.

Convert a ~~production rule~~ CFG to CNF: VP → VBD NP PP

    VP → VBD @VP-VBD

    @VP-VBD → NP PP

                                 $VP \to V$

We assume the grammar are in CNF.       $V \to x \, Y \qquad VP \to x \, Y$

# Dynamic programming on the tree

$$p(t) = \underbrace{q(A \to BC)}_{\text{top rule}} \times \underbrace{q(t_B)}_{\text{left child}} \times \underbrace{q(t_C)}_{\text{right child}}$$



What are the variables when constructing a tree rooted at $A$ spanning $x_i, \ldots, x_j$?

- ▶ The production rule $A \to BC$
- ▶ The splitting point $s$: $B$ spans $x_i, \ldots, x_s$ and $C$ spans $x_{s+1}, \ldots, x_j$
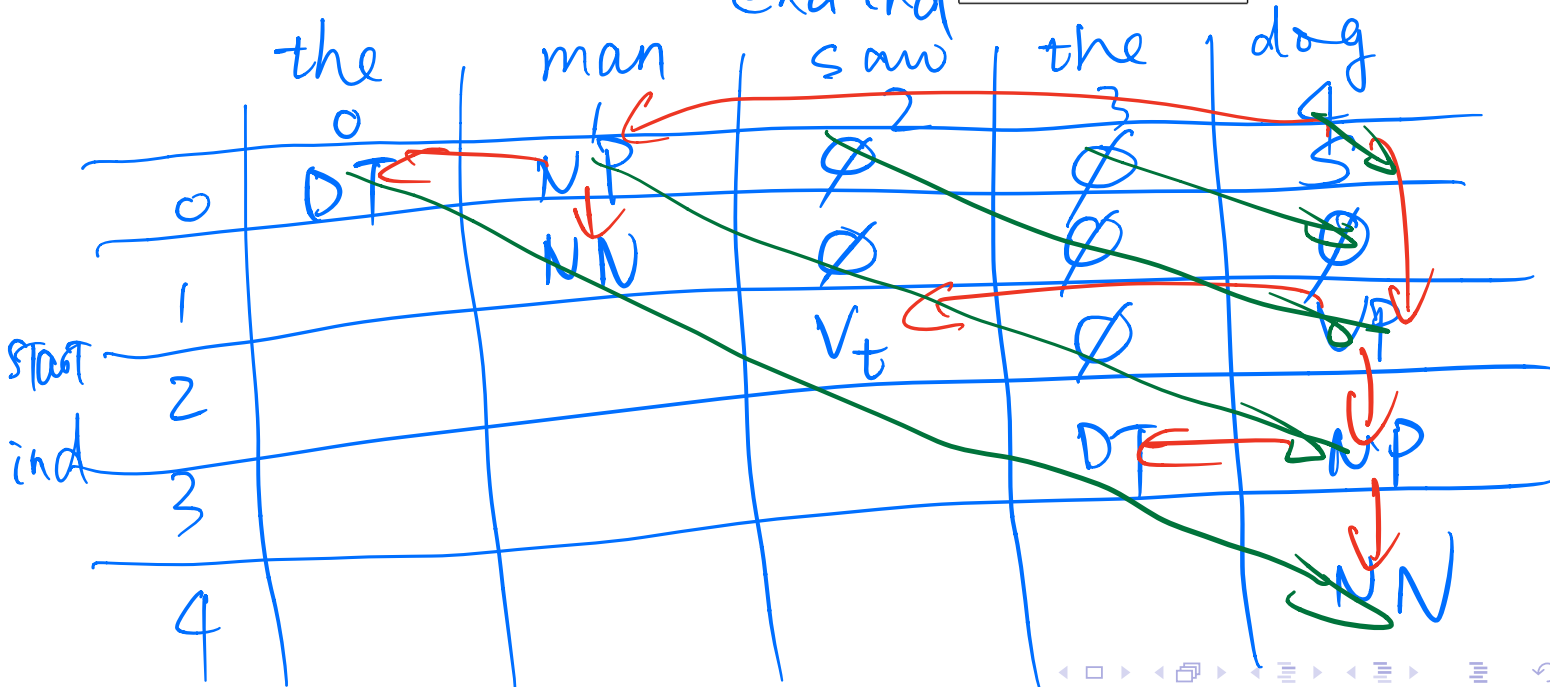
# Bottom-up parsing

$n \, |R|$ 　　　$O\left(n^2 \cdot n \cdot |R|\right)$

$R =$

| S | → | NP | VP |
|---|---|----|----|
| VP | → | Vi | |
| VP | → | Vt | NP |
| VP | → | VP | PP |
| NP | → | DT | NN |
| NP | → | NP | PP |
| PP | → | IN | NP |

| Vi | → | sleeps |
|----|---|--------|
| Vt | → | saw |
| NN | → | man |
| NN | → | woman |
| NN | → | telescope |
| NN | → | dog |
| DT | → | the |
| IN | → | with |
| IN | → | in |

the ∧ man ∧ saw

end ind

| | the | man | saw | the | dog |
|--|-----|-----|-----|-----|-----|
| | 0 | 1 | 2 | 3 | 4 |
| 0 | DT | NP / NN | ∅ | ∅ | |
| 1 | | | ∅ | ∅ | |
| 2 | | | Vt | ∅ | VP |
| 3 | | | | DT | NP |
| 4 | | | | | NN |

START

ind

# The CYK algorithm

Notation: $\mathcal{T}(i, j, X)$ is the set of trees with root node $X$ spanning $x_i, \ldots, x_j$

Subproblem:

$$\pi(i, j, X) = \max_{t \in \mathcal{T}(i,j,X)} p(t)$$

Base case:

$$\pi(i, i, X) = \begin{cases} q(X \to x_i) & \text{if } X \to x_i \in R \\ 0 & \text{otherwise} \end{cases}$$

Recursion:

$$\pi(i, j, X) = \max_{\substack{Y, Z \in N \\ s \in \{i, \ldots, j-1\}}} q(X \to YZ) \times \pi(i, s, Y) \times \pi(s+1, j)$$

Use backtracking to find the argmax tree.

# Variants of CYK

**Argmax**: find the most likely tree (analogous to Viterbi).

$$\pi(i,j,X) = \max_{\substack{Y,Z \in N \\ s \in \{i,\dots,j-1\}}} q(X \to YZ) \times \pi(i,s,Y) \times \pi(s+1,j)$$

**Recognition**: does the string belong to the language?

$$\pi(i,j,X) = \bigvee_{\substack{Y,Z \in N \\ s \in \{i,\dots,j-1\}}} \mathbb{I}[X \to YZ \in R] \wedge \pi(i,s,Y) \wedge \pi(s+1,j)$$

**Marginalization**: what's the probability of the string being generated from the grammar? (the **inside algorithm**)

$$\pi(i,j,X) = \sum_{\substack{Y,Z \in N \\ s \in \{i,\dots,j-1\}}} q(X \to YZ) \times \pi(i,s,Y) \times \pi(s+1,j)$$

Complexity?

# Summary

| | NB | HMM | PCFG |
|---|---|---|---|
| output structure | category | sequence | tree |
| learning | | MLE | |
| decoding | brute force | viterbi | CKY |
| marginalization | | $P(y_i \mid x)$ $P(y_{i+1}, y_i \mid x)$ | $P(i, j, N \mid x)$ |
| unsupervised learning | | EM | |

# Table of Contents

# CRF for trees

Input: sequence of words $x = (x_1, \ldots, x_n)$
Output: parse tree $y \in \mathcal{T}(x)$
Model: decompose by production rules

$$p(y \mid x; \theta) \propto \prod_{(r,s)} \psi(r, s \mid x; \theta)$$

- ▶ $r$: production rule
- ▶ $s$: start, split, end indices of the rule $r$



$$\psi(S \rightarrow NP\ VP, (0, 1, 2))$$

$$\psi(NP \rightarrow DT\ NN, (0, 0, 1))$$

# CRF parsing

Potential functions:

$$\psi(r, s \mid x; \theta) = \exp\left(\theta \cdot \phi(r, s, x)\right)$$

$$\prod_{(r,s)\in\mathcal{T}(x)} \psi(r, s \mid x; \theta) = \exp\left(\sum_{(r,s)\in\mathcal{T}(x)} \theta \cdot \phi(r, s, x)\right) \Bigg/ Z$$

Learning: MLE $\log P(y\mid x)$

1. Compute the partition function by the inside algorithm
2. Call autograd to compute the gradient (backpropagation)
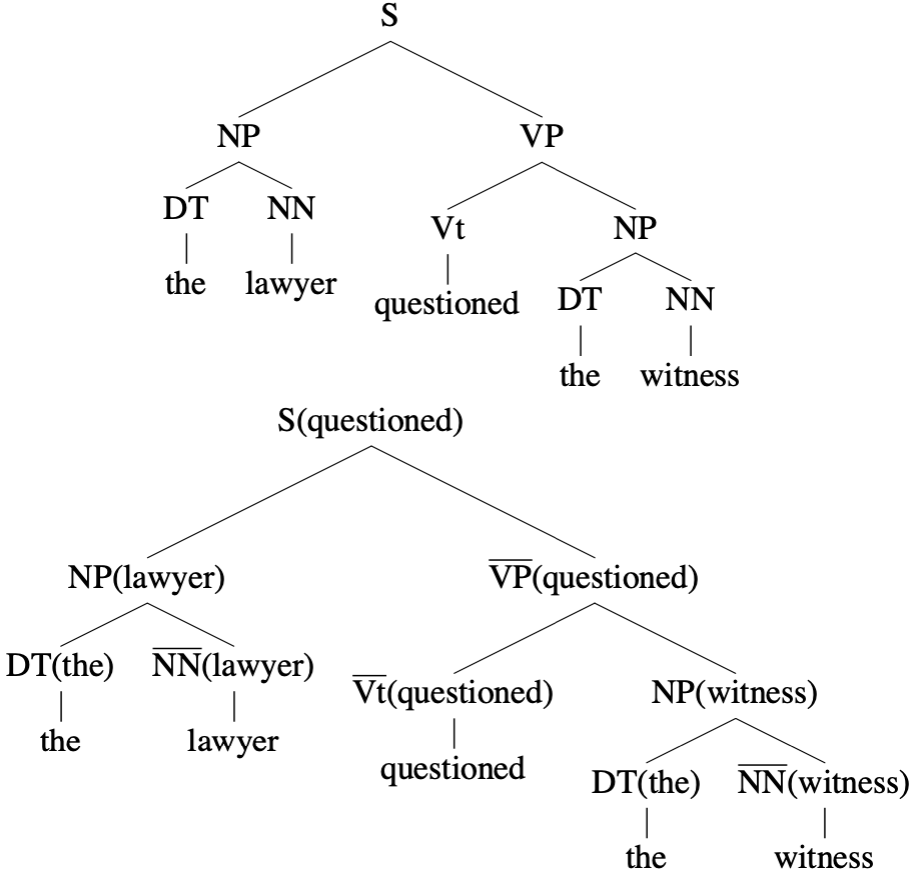
Inference: CYK

# Limitations of PCFG

# Limitations of PCFG



Limit

No lexical information

# Lexicalized PCFG

Attach the "head" of the span to each non-terminal

# Features

local score $= \theta \cdot \phi($VP $\rightarrow$ VBD NP$, (5, 6, 8), $...averted financial disaster...$)$



Figure: Less grammar, more features. [Hall+ 14]

# Neural CRF parser

a) $\phi = f_s^\top W f_o$

b) $\phi = g(Hv(f_w))^\top W f_o$
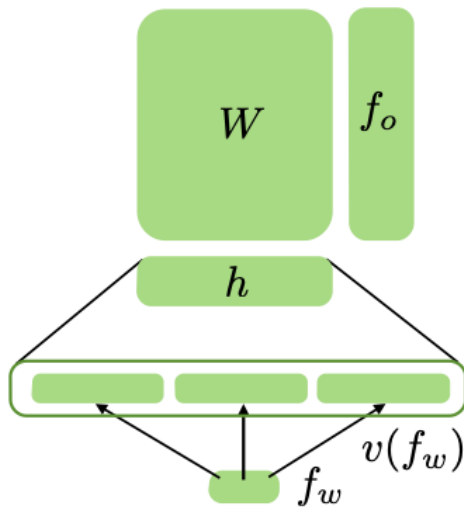


$$W_{ij} = \text{e ght}([[f_{s,i} \quad f_{o,j}]])$$

Figure: Neural CRF Parsing. [Durrett+ 15]

# Evaluation

$$( i, j, x )$$

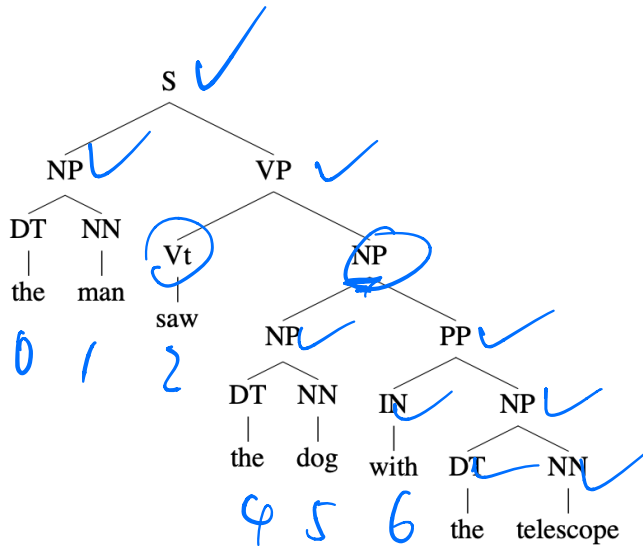$$\text{recall} = \frac{\#\text{correct constituents}}{\#\text{total constituents in gold trees}}$$

$$\text{precision} = \frac{\#\text{correct constituents}}{\#\text{total constituents in predicted trees}}$$

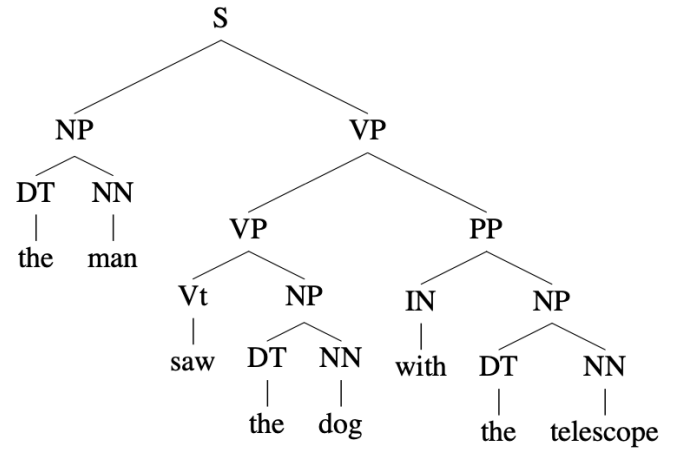$$\text{F1} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

▶ Labeled F1: the non-terminal node label must be correct
▶ Unlabeled F1: just consider the tree structure

# Example

$(0, 8, S) \quad (0, 1, NP) \cdots$



(a) Gold.

(b) Predicted.