# Tutorial on Probabilistic Context-Free Grammars

Raphael Hoffmann

590AI, Winter 2009

# Outline

- PCFGs: Inference and Learning
- Parsing English
- Discriminative CFGs
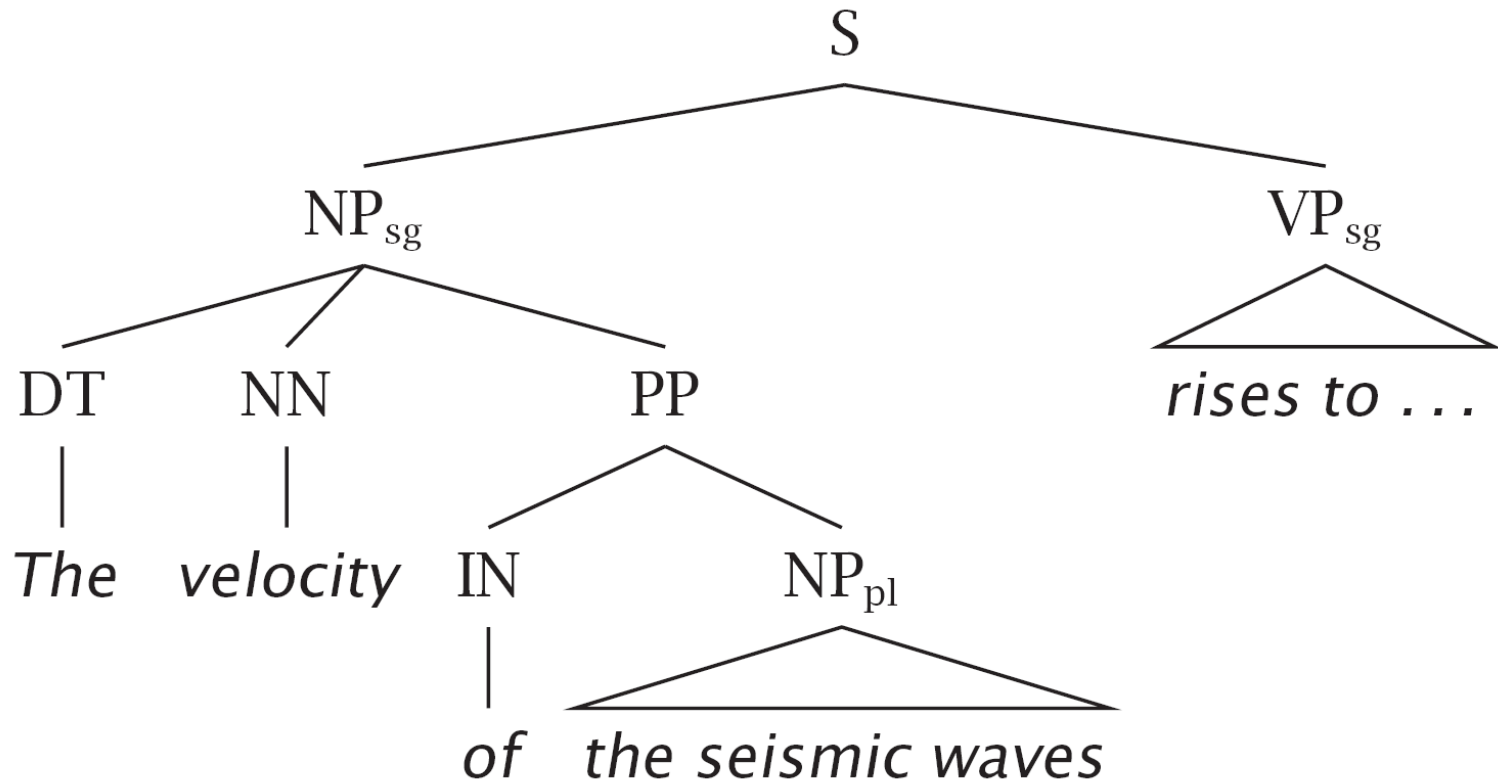- Grammar Induction

# Image Search for "pcfg"

# Outline

- PCFGs: Inference and Learning
- Parsing English
- Discriminative CFGs
- Grammar Induction

# The velocity of the seismic waves rises to …

# A CFG consists of

- Terminals $\qquad w^1, w^2, \ldots, w^V$

- Nonterminals $\qquad N^1, N^2, \ldots, N^n$

- Start symbol $\qquad N^1$

- Rules $\qquad N^i \longrightarrow \zeta^j$

  where $\zeta^j$ is a sequence of terminals and nonterminals
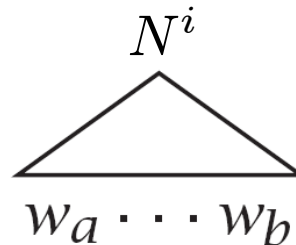
# A (generative) PCFG consists of

- Terminals $w^1, w^2, \ldots, w^V$

- Nonterminals $N^1, N^2, \ldots, N^n$

- Start symbol $N^1$

- Rules $N^i \longrightarrow \zeta^j$

  where $\zeta^j$ is a sequence of terminals and nonterminals

- Rule probabilities  such that
  $$\forall_i \sum_j P(N^i \longrightarrow \zeta^j) = 1$$

# Notation

sentence:  sequence of words   $w_1 w_2 \ldots w_m$

$w_{ab}$ :   the subsequence  $w_a \ldots w_b$

$N^i_{ab}$ :   nonterminal $N^i$ dominates  $w_a \ldots w_b$

$$
\begin{array}{c}
N^i \\
\triangle \\
w_a \cdots w_b
\end{array}
$$

$N^i \Longrightarrow^* \zeta$ :    repeated derivation from $N^i$ gives $\zeta$

# Probability of a Sentence

$$P(w_{1n}) = \sum_t P(w_{1n}, t)$$

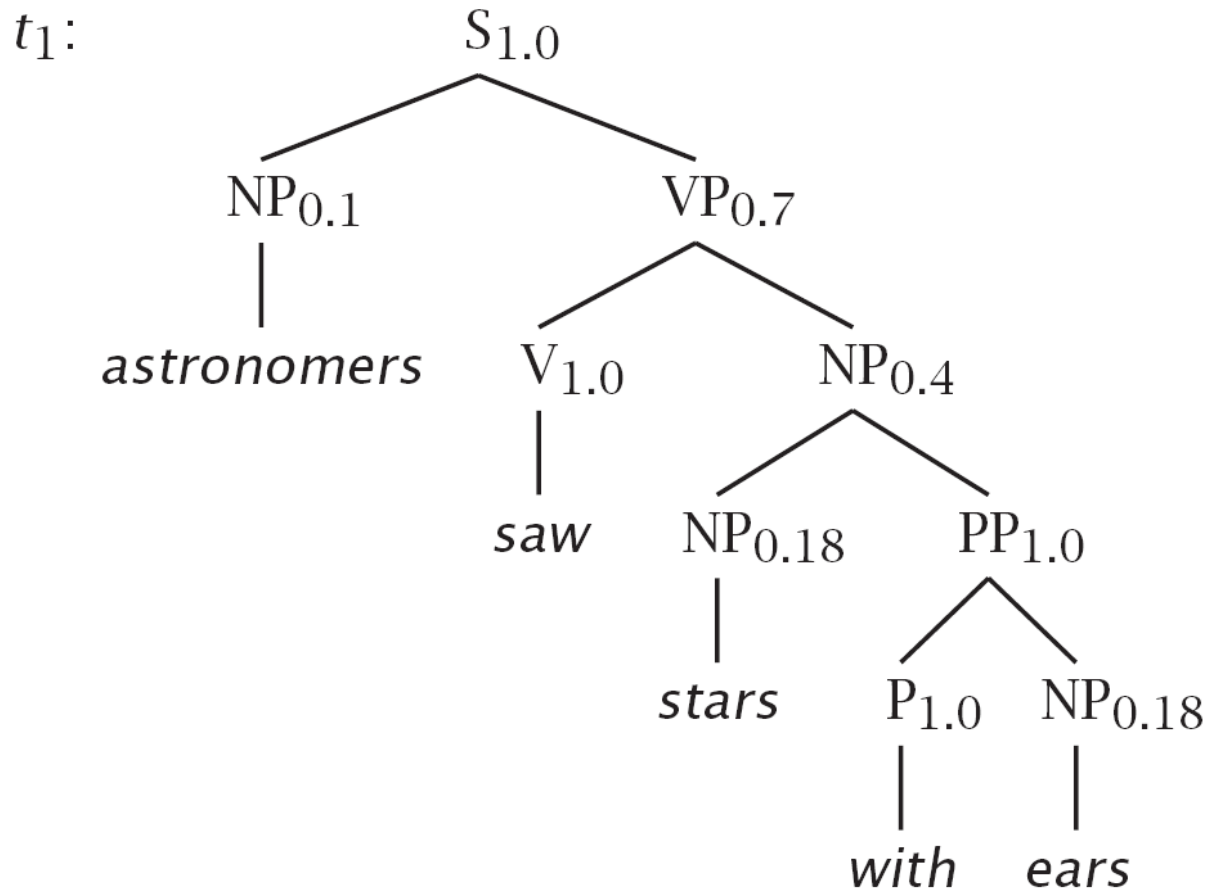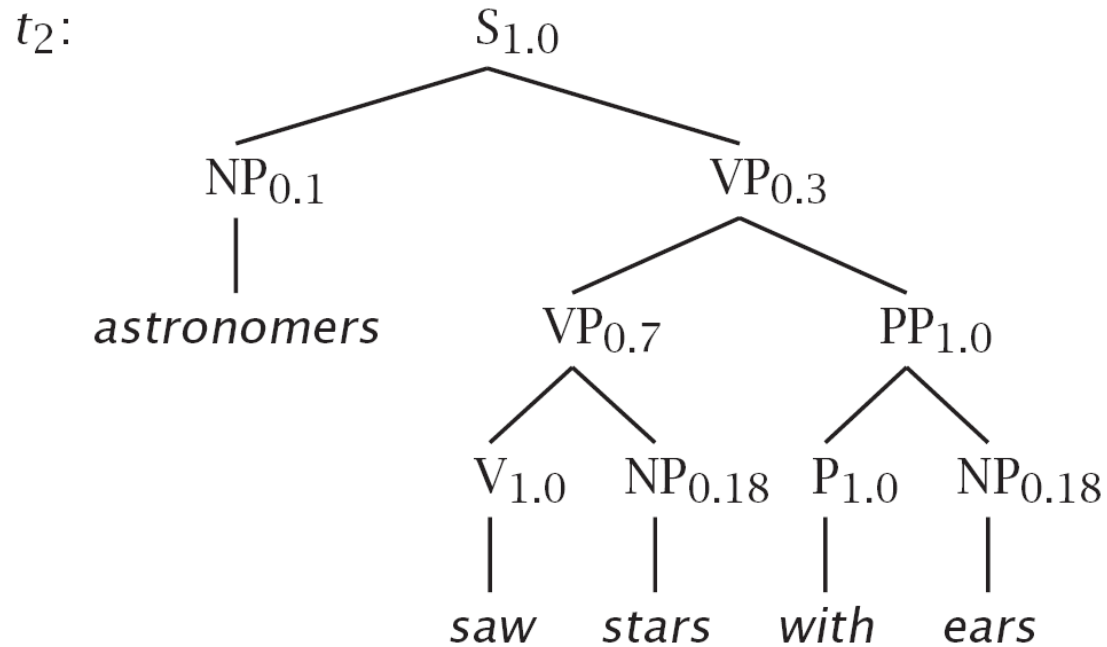where $t$ a parse tree of $w_{1n}$

# Example

| | | | | |
|---|---|---|---|---|
| S → NP VP | 1.0 | NP → NP PP | 0.4 |
| PP → P NP | 1.0 | NP → *astronomers* | 0.1 |
| VP → V NP | 0.7 | NP → *ears* | 0.18 |
| VP → VP PP | 0.3 | NP → *saw* | 0.04 |
| P → *with* | 1.0 | NP → *stars* | 0.18 |
| V → *saw* | 1.0 | NP → *telescopes* | 0.1 |

- Terminals        with, saw, astronomers, ears, stars, telescopes
- Nonterminals    S, PP, P, NP, VP, V
- Start symbol    S

# astronomers saw stars with ears

# astronomers saw stars with ears
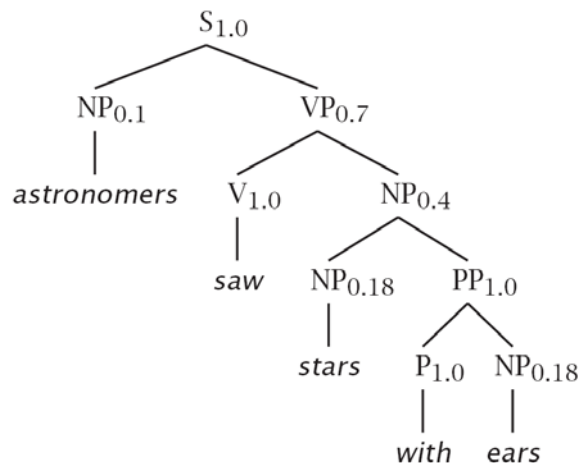
# Probabilities

$t_1$:



$$
\begin{aligned}
P(t_1) & = 1.0 \times 0.1 \times 0.7 \times 1.0 \times 0.4 \\
& \quad \times 0.18 \times 1.0 \times 1.0 \times 0.18 \\
& = 0.0009072 \\
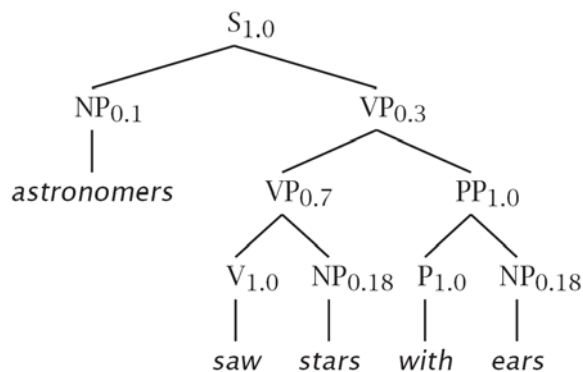P(t_2) & = 1.0 \times 0.1 \times 0.3 \times 0.7 \times 1.0 \\
& \quad \times 0.18 \times 1.0 \times 1.0 \times 0.18 \\
& = 0.0006804 \\
P(w_{15}) & = P(t_1) + P(t_2) = 0.0015876
\end{aligned}
$$

$t_2$:

# Assumptions of PCFGs

1.  Place invariance (like time invariance in HMMs)

$$\forall k \quad P(N^j_{k(k+c)} \longrightarrow \zeta) \quad \text{is the same}$$

# Assumptions of PCFGs

1. Place invariance (like time invariance in HMMs)

$$\forall k \quad P(N^j_{k(k+c)} \longrightarrow \zeta) \quad \text{is the same}$$

2. Context free

$$P(N^j_{kl} \longrightarrow \zeta \mid \text{words outside } w_k \ldots w_l) = P(N^j_{kl} \longrightarrow \zeta)$$

# Assumptions of PCFGs

1. Place invariance (like time invariance in HMMs)

$$\forall k \quad P(N^j_{k(k+c)} \longrightarrow \zeta) \quad \text{is the same}$$

2. Context free

$$P(N^j_{kl} \longrightarrow \zeta \mid \text{words outside } w_k \ldots w_l) = P(N^j_{kl} \longrightarrow \zeta)$$

3. Ancestor free

$$P(N^j_{kl} \longrightarrow \zeta \mid \text{ancestor nodes of } N^j_{kl}) = P(N^j_{kl} \longrightarrow \zeta)$$

# Some Features of PCFGs

- Partial solution for grammar ambiguity
- Can be learned from positive data alone
(but grammar induction difficult)
- Robustness
(admit everything with low probability)
- Gives a probabilistic language model
- Predictive power better than that for a HMM

# Some Features of PCFGs

- Not lexicalized (probabilities do not factor in lexical co-occurrence)

- PCFG is a worse language model for English than n-gram models

- Certain biases: smaller trees more probable (average WSJ sentence 23 words)

# Inconsistent Distributions

- S → *rhubarb*   $P = \frac{1}{3}$

  S → S S        $P = \frac{2}{3}$

- *rhubarb*                                    $\frac{1}{3}$

  *rhubarb rhubarb*                   $\frac{2}{3} \times \frac{1}{3} \times \frac{1}{3} = \frac{2}{27}$

  *rhubarb rhubarb rhubarb*   $\left(\frac{2}{3}\right)^2 \times \left(\frac{1}{3}\right)^3 \times 2 = \frac{8}{243}$

  . . .

- $P(\mathcal{L}) = \frac{1}{3} + \frac{2}{27} + \frac{8}{243} + \ldots = \frac{1}{2}$

- Improper/inconsistent distribution

- Not a problem if you estimate from parsed treebank: Chi and Geman 1998).

# Questions

Let $w_{1m}$ be a sentence, $G$ a grammar, $t$ a parse tree.

- What is the probability of a sentence?

$$P(w_{1m}|G)$$

# Questions

Let $w_{1m}$ be a sentence, $G$ a grammar, $t$ a parse tree.

- What is the probability of a sentence?

$$P(w_{1m}|G)$$

- What is the most likely parse of sentence?

$$\arg\max_t P(t|w_{1m}, G)$$

# Questions

Let $w_{1m}$ be a sentence, $G$ a grammar, $t$ a parse tree.

- What is the probability of a sentence?
$$P(w_{1m}|G)$$
- What is the most likely parse of sentence?
$$\arg\max_t P(t|w_{1m}, G)$$
- What rule probs. maximize probs. of sentences?
Find $G$ that maximizes $P(w_{1m}|G)$

# Chomsky Normal Form

- Any CFG grammar can be represented in CNF where all rules take the form

$$N^i \longrightarrow N^j N^k$$

$$N^i \longrightarrow w^j$$

# HMMs and PCFGs

- HMMs: distribution over strings of certain length
$$\forall n \sum_{w_{1n}} P(w_{1n}) = 1$$
- PCFGs: distribution over strings of language $L$
$$\sum_{w \in L} P(w) = 1$$

# HMMs and PCFGs

- HMMs: distribution over strings of certain length

$$\forall n \sum_{w_{1n}} P(w_{1n}) = 1$$

- PCFGs: distribution over strings of language *L*

$$\sum_{w \in L} P(w) = 1$$

- Consider

$$P(\text{John decided to bake a})$$

high probability in HMM, low probability in PCFG

# HMMs and PCFGs

- HMMs: distribution over strings of certain length

$$\forall n \sum_{w_{1n}} P(w_{1n}) = 1$$

- PCFGs: distribution over strings of language *L*

$$\sum_{w \in L} P(w) = 1$$

- Consider

$$P(\text{John decided to bake a})$$

high probability in HMM, low probability in PCFG

- Probabilistic Regular Grammar

$$N^i \longrightarrow w^j N^k$$
$$N^i \longrightarrow w^j$$

# HMMs and PCFGs

$X$:  NP  $\longrightarrow$  N′  $\longrightarrow$  N′  $\longrightarrow$  N′  $\longrightarrow$  sink

$O$:  *the*  *big*  *brown*  *box*

```
              NP
           /      \
        the        N′
                /      \
             big        N′
                     /      \
                  brown      N⁰
                             |
                            box
```

# Inside and Outside Probabilities

- For HMMs we have

Forwards $\quad \alpha_i(t) \quad = \quad P(w_{1(t-1)}, X_t = i)$

Backwards $\quad \beta_i(t) \quad = \quad P(w_{tT} | X_t = i)$

# Inside and Outside Probabilities

- ## For HMMs we have

Forwards $\quad \alpha_i(t) \quad = \quad P(w_{1(t-1)}, X_t = i)$

Backwards $\quad \beta_i(t) \quad = \quad P(w_{tT} | X_t = i)$

- ## For PCFGs we have

Outside $\quad \alpha_j(p, q) \quad = \quad P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m} | G)$

Inside $\quad \beta_j(p, q) \quad = \quad P(w_{pq} | N_{pq}^j, G)$

# Inside and Outside Probabilities



$N^1$

$\alpha$

$N^J$

$\beta$

$w_1 \quad \cdots \quad w_{p-1} w_p \quad \cdots \quad w_q w_{q+1} \quad \cdots \quad w_m$

# Probability of a sentence

Outside $\quad \alpha_j(p, q) \quad = \quad P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m}|G)$

Inside $\quad \beta_j(p, q) \quad = \quad P(w_{pq}|N_{pq}^j, G)$

$$P(w_{1m}|G) \quad = \quad \beta_1(1, m)$$

$$P(w_{1m}|G) \quad = \quad \sum_j \alpha_j(k, k) P(N^j \longrightarrow w_k)$$

# Inside Probabilities

$$\beta_j(p, q) \quad = \quad P(w_{pq}|N_{pq}^j, G)$$

- Base case

$$\beta_j(k, k) \quad = \quad P(w_{kk}|N_{kk}^j, G)$$
$$= \quad P(N^j \longrightarrow w_k|G)$$

# Inside Probabilities

$$\beta_j(p, q) \quad = \quad P(w_{pq}|N^j_{pq}, G)$$

- ## Base case

$$\beta_j(k, k) \quad = \quad P(w_{kk}|N^j_{kk}, G)$$
$$= \quad P(N^j \longrightarrow w_k|G)$$

- ## Induction

Want to find $\beta_j(p, q)$ for $p < q$

Since we assume Chomsky Normal Form,
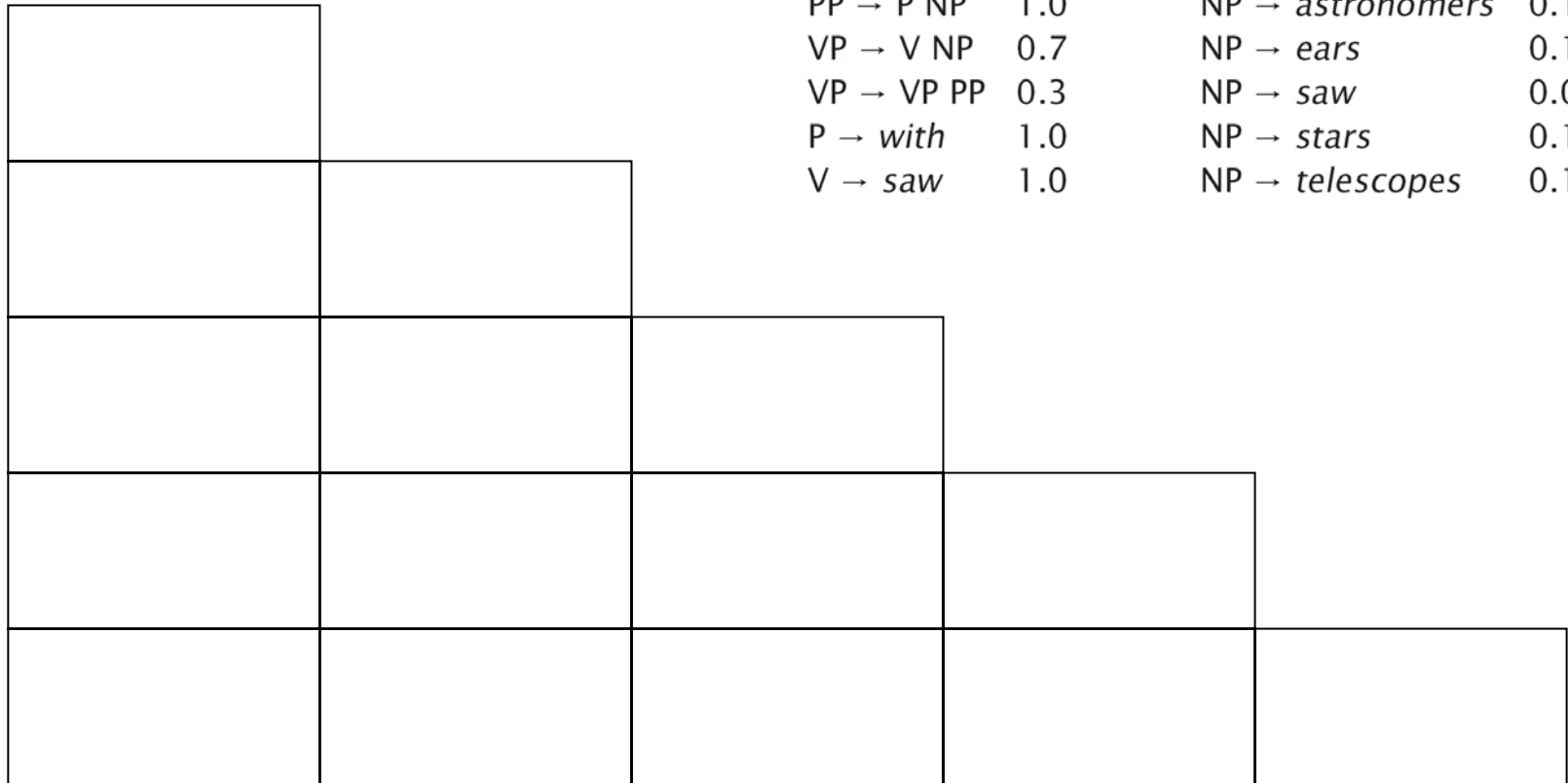the first rule must be of the form $N^j \longrightarrow N^r N^s$

So we can divide the sentence in two in
various places and sum the result

$$\beta_j(p, q) = \sum_{r,s} \sum_{d=p}^{q-1} P(N^j \longrightarrow N^r N^s)\beta_r(p, d)\beta_s(d + 1, q)$$

# CYK Algorithm

| | | | | |
|---|---|---|---|---|
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| astronomers | saw | stars | with | ears |

| | | | | |
|---|---|---|---|---|
| S → NP VP | 1.0 | | NP → NP PP | 0.4 |
| PP → P NP | 1.0 | | NP → *astronomers* | 0.1 |
| VP → V NP | 0.7 | | NP → *ears* | 0.18 |
| VP → VP PP | 0.3 | | NP → *saw* | 0.04 |
| P → *with* | 1.0 | | NP → *stars* | 0.18 |
| V → *saw* | 1.0 | | NP → *telescopes* | 0.1 |

# CYK Algorithm

| | | | | |
|---|---|---|---|---|
| S → NP VP | 1.0 | | NP → NP PP | 0.4 |
| PP → P NP | 1.0 | | NP → *astronomers* | 0.1 |
| VP → V NP | 0.7 | | NP → *ears* | 0.18 |
| VP → VP PP | 0.3 | | NP → *saw* | 0.04 |
| P → *with* | 1.0 | | NP → *stars* | 0.18 |
| V → *saw* | 1.0 | | NP → *telescopes* | 0.1 |

| | | | | |
|---|---|---|---|---|
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| $\beta_{NP} = 0.1$ ? | $\beta_V = 1.0$ $\beta_{NP} = 0.04$ ? | $\beta_{NP} = 0.18$ ? | $\beta_P = 1.0$ | $\beta_{NP} = 0.18$ |
| astronomers | saw | stars | with | ears |

# CYK Algorithm

$$
\begin{array}{llll}
S \rightarrow NP\ VP & 1.0 & NP \rightarrow NP\ PP & 0.4 \\
PP \rightarrow P\ NP & 1.0 & NP \rightarrow astronomers & 0.1 \\
VP \rightarrow V\ NP & 0.7 & NP \rightarrow ears & 0.18 \\
VP \rightarrow VP\ PP & 0.3 & NP \rightarrow saw & 0.04 \\
P \rightarrow with & 1.0 & NP \rightarrow stars & 0.18 \\
V \rightarrow saw & 1.0 & NP \rightarrow telescopes & 0.1 \\
\end{array}
$$

| | | | | |
|---|---|---|---|---|
| | | | | |
| | | | | |
| | | | | |
| ? | $\beta_{VP} = 0.126$ ? | ? | $\beta_{PP} = 0.18$ ? | |
| $\beta_{NP} = 0.1$ | $\beta_V = 1.0$ $\beta_{NP} = 0.04$ | $\beta_{NP} = 0.18$ | $\beta_P = 1.0$ | $\beta_{NP} = 0.18$ |
| astronomers | saw | stars | with | ears |

# CYK Algorithm

| | | | | |
|---|---|---|---|---|
| $\beta_S = 0.015876$ | | | | |
| | $\beta_{VP} = 0.015876$ | | | |
| $\beta_S = 0.0126$ **?** | | $\beta_{NP} = 0.01296$ | | |
| | $\beta_{VP} = 0.126$ | | $\beta_{PP} = 0.18$ | |
| $\beta_{NP} = 0.1$ | $\beta_V = 1.0$ $\beta_{NP} = 0.04$ | | $\beta_P = 1.0$ | $\beta_{NP} = 0.18$ |
| astronomers | saw | stars | with | ears |

| | | | | |
|---|---|---|---|---|
| S → NP VP | 1.0 | | NP → NP PP | 0.4 |
| PP → P NP | 1.0 | | NP → *astronomers* | 0.1 |
| VP → V NP | 0.7 | | NP → *ears* | 0.18 |
| VP → VP PP | 0.3 | | NP → *saw* | 0.04 |
| P → *with* | 1.0 | | NP → *stars* | 0.18 |
| V → *saw* | 1.0 | | NP → *telescopes* | 0.1 |

# CYK Algorithm

Worst case: O(m³r)
         m = length of sentence
         r = number of rules in grammar
         n = number of non-terminals
If we consider all possible CNF rules: O(m³n³)

| astronomers | saw | stars | with | ears |
|---|---|---|---|---|
| $\beta_S = 0.015876$ | | | | |
| | $\beta_{VP} = 0.015876$ | | | |
| $\beta_S = 0.0126$ | | $\beta_{NP} = 0.01296$ | | |
| | $\beta_{VP} = 0.126$ | | $\beta_{PP} = 0.18$ | |
| $\beta_{NP} = 0.1$ | $\beta_V = 1.0$ $\beta_{NP} = 0.04$ | $\beta_{NP} = 0.18$ | $\beta_P = 1.0$ | $\beta_{NP} = 0.18$ |

# Outside Probabilities

- Compute top-down (after inside probabilities)
- Base case

$$\alpha_1(1, m) = 1$$

$$\alpha_j(1, m) = 0, \text{ for } j \neq 1$$

- Induction

$$\alpha_j(p, q) = \left( \sum_{f,g} \sum_{e=q+1}^{m} \alpha_f(p, e) P(N^f \longrightarrow N^j N^g) \beta_g(q+1, e) \right)$$

$$+ \left( \sum_{f,g} \sum_{e=1}^{p-1} \alpha_f(e, q) P(N^f \longrightarrow N^g N^j) \beta_g(e, p-1) \right)$$

# Probability of a node existing

- As with a HMM, we can form a product of the inside and outside probabilities.

$$\alpha_j(p, q)\beta_j(p, q) = P(w_{1m}, N_{pq}^j | G)$$

# Probability of a node existing

- As with a HMM, we can form a product of the inside and outside probabilities.

$$\alpha_j(p,q)\beta_j(p,q) = P(w_{1m}, N_{pq}^j | G)$$

- Therefore,

$$P(w_{1m}, N_{pq} | G) = \sum_j \alpha_j(p,q)\beta_j(p,q)$$

# Probability of a node existing

- As with a HMM, we can form a product of the inside and outside probabilities.

$$\alpha_j(p,q)\beta_j(p,q) = P(w_{1m}, N_{pq}^j | G)$$

- Therefore,

$$P(w_{1m}, N_{pq} | G) = \sum_j \alpha_j(p,q)\beta_j(p,q)$$

- Just in the cases of the root node and the preterminals, we know there will be some such constituent.

# Training

- If have data → count

$$\hat{P}(N^j \longrightarrow \zeta) = \frac{C(N^j \longrightarrow \zeta)}{\sum_\gamma C(N^j \longrightarrow \gamma)}$$

# Training

- ## If have data → count

$$\hat{P}(N^j \longrightarrow \zeta) = \frac{C(N^j \longrightarrow \zeta)}{\sum_\gamma C(N^j \longrightarrow \gamma)}$$

- ## else use EM  (Inside-Outside-Algorithm)

**repeat**

    compute $\alpha_j$ 's and $\beta_j$ 's

    compute $\hat{P}$ 's

$$\hat{P}(N^j \longrightarrow N^r N^s) = \ldots$$

$$\hat{P}(N^j \longrightarrow w^k) = \ldots$$

**end**

two really long formulas with $\alpha_j$ 's and $\beta_j$ 's

# EM Problems

- Slow: $O(m^3 n^3)$ for each sentence and each iteration

- Local maxima (Charniak: 300 trials led to 300 different max.)

- In practice, need >3 times more non-terminals than are theoretically needed

- No guarantee that learned non-terminals correspond to NP, VP, …

# Bracketing helps

Pereira/Schabes '92:

- Train on sentences:
  37% of predicted brackets correct

- Train on sentences + brackets:
  90% of predicted brackets correct

# Grammar Induction

- Rules typically selected by linguist
- Automatic induction is very difficult for context-free languages
- It is easy to find *some* form of structure, but little resemblance to that of linguistics/NLP
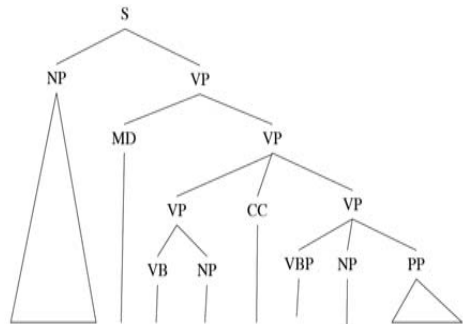
# Outline

- PCFGs: Inference and Learning
- Parsing English
- Discriminative CFGs
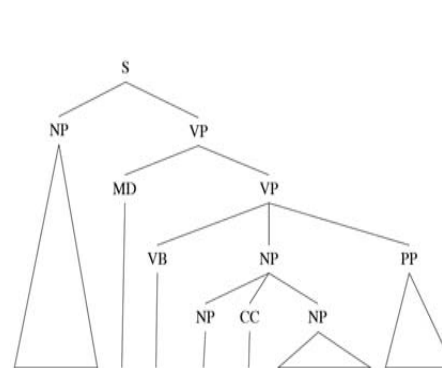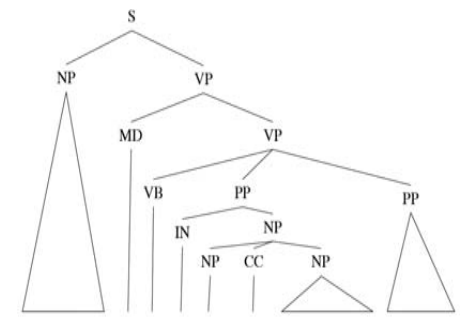- Grammar Induction

# Parsing for Disambiguation

The post office will hold out discounts and service concessions as incentives.
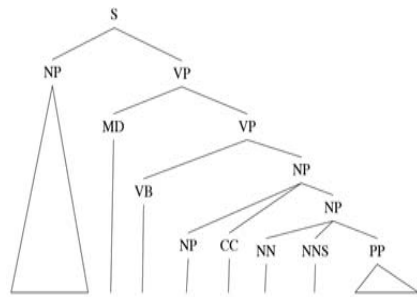
# Parsing for Disambiguation

- There are typically many syntactically possible parses

- Want to find the most likely parses

# Treebanks

- If grammar induction does not work, why not count expansions in many parse trees?

- Penn Treebank

```
( (S
    (NP (NBAR (ADJP (ADJ "Battle-tested/JJ")
                    (ADJ "industrial/JJ"))
              (NPL "managers/NNS")))
    (? (ADV "here/RB"))
    (? (ADV "always/RB"))
    (AUX (TNS *))
    (VP (VPRES "buck/VBP")))
    (? (PP (PREP "up/RP")
           (NP (NBAR (ADJ "nervous/JJ")
                     (NPL "newcomers/NNS")))))
    (? (PP (PREP "with/IN")
           (NP (DART "the/DT")
               (NBAR (N "tale/NN"))
               (PP of/PREP
                   (NP (DART "the/DT")
                       (NBAR (ADJP
                             (ADJ "first/JJ")))))))))
```

# PCFG weaknesses

- ## No Context
  - (immediate prior context, speaker, ...)

- ## No Lexicalization
  - "VP NP NP" more likely if verb is "hand" or "tell"
  - fail to capture lexical dependencies (n-grams do)

- ## No Structural Context
  - How NP expands depends on position

# PCFG weaknesses

| Expansion | % as Subj | % as Obj |
|---|---|---|
| NP ⟶ PRP | 13.7% | 2.1% |
| NP ⟶ NNP | 3.5% | 0.9% |
| NP ⟶ DT NN | 5.6% | 4.6% |
| NP ⟶ NN | 1.4% | 2.8% |
| NP ⟶ NP SBAR | 0.5% | 2.6% |
| NP ⟶ NP PP | 5.6% | 14.1% |

| Expansion | % as 1st Obj | % as 2nd Obj |
|---|---|---|
| NP ⟶ NNS | 7.5% | 0.2% |
| NP ⟶ PRP | 13.4% | 0.9% |
| NP ⟶ NP PP | 12.2% | 14.4% |
| NP ⟶ DT NN | 10.4% | 13.3% |
| NP ⟶ NNP | 4.5% | 5.9% |
| NP ⟶ NN | 3.9% | 9.2% |
| NP ⟶ JJ NN | 1.1% | 10.4% |
| NP ⟶ NP SBAR | 0.3% | 5.1% |

Slide based on "Foundations of Statistical Natural Language Processing" by Christopher Manning and Hinrich Schütze

# Other Approaches

- Challenge: use lexical and structural context, without too many parameters, sparse data

- Other Grammars
  - Probabilistic Left-Corner Grammars
  - Phrase Structure Grammars
  - Dependency Grammars
  - Probabilistic Tree Substitution Grammars
  - History-based Grammars

# Outline

- PCFGs: Inference and Learning
- Parsing English
- Discriminative CFGs
- Grammar Induction

# Generative vs Discriminative

- An HMM (or PCFG) is a *generative* model

$$P(y, w)$$

- Often sufficient is a *discriminative* model

$$P(y|w)$$

- Easier, because does not contain $P(w)$

- Cannot model dependent features in HMM, so one only picks one feature: word's identity

# Generative and Discriminative Models

# Generative and Discriminative Models

# Generative and Discriminative Models

# Discriminative
# Context-Free Grammars

- Terminals $\qquad$ $w^1, w^2, \ldots, w^V$

- Nonterminals $\qquad$ $N^1, N^2, \ldots, N^n$

- Start symbol $\qquad$ $N^1$

- Rules $\qquad$ $N^i \longrightarrow \zeta^j$

  where $\zeta^j$ is a sequence of terminals and nonterminals

- Rule scores

$$S(N^i \longrightarrow \zeta^j, p, q) = \sum_{k=1}^{F} \lambda_k(N^i \longrightarrow \zeta^j) f_k(w_1 w_2 \ldots w_m, p, q, N^i \longrightarrow \zeta^j)$$

# Features

$$S(N^i \longrightarrow \zeta^j, p, q) = \sum_{k=1}^{F} \lambda_k(N^i \longrightarrow \zeta^j) f_k(w_1 w_2 \ldots w_m, p, q, N^i \longrightarrow \zeta^j)$$

- Features can depend on all tokens + span

- Consider feature "AllOnTheSameLine"

```
Mavis Wood
Products
```

```
Mavis Wood Products
```

[compare to linear CRF $\quad f_k(s_t, t_{t-1}, w_1 w_2 \ldots w_m, t)$ ]

# Features

$$S(N^i \longrightarrow \zeta^j, p, q) = \sum_{k=1}^{F} \lambda_k (N^i \longrightarrow \zeta^j) f_k(w_1 w_2 \ldots w_m, p, q, N^i \longrightarrow \zeta^j)$$

- No independence between features necessary
- Can create features based on words, dictionaries, digits, capitalization, …
- Can still do efficient Viterbi inference in O(m³r)

# Example

FIRSTNAME    LASTNAME    ADDRLINE    ADDRLINE    ADDRLINE    CITY    STATE    ZIP
|            |           |           |           |          |       |        |
John         Doe         100         Main        Street     Seattle WA       98195

```
Fred Jones                        Boston College
10 Main St.                       10 Main St.
Cambridge, MA 02146               Cambridge MA 02146
(425) 994-8021                    (425) 994-8021
```

BizContact $\longrightarrow$ BizName Address BizPhone
PersonalContact $\longrightarrow$ BizName Address HomePhone

# Example

# Training

- Train feature weight vector for each rule $\lambda_j(R)$
- Have labels, but not parse trees; efficiently create trees by ignoring leaves

# Collins' Averaged Perceptron

$$\textbf{for } r \leftarrow 1 \ldots numRounds \textbf{ do}$$
$$\quad \textbf{for } i \leftarrow 1 \ldots m \textbf{ do}$$
$$\quad\quad T \leftarrow \text{ optimal parse of } w^i \text{ with current parameters}$$
$$\quad\quad \textbf{if } T \neq T^i \textbf{ then}$$
$$\quad\quad\quad \textbf{for } \text{ each rule } R \text{ used in } T \text{ but not in } T^i \textbf{ do}$$
$$\quad\quad\quad\quad \textbf{if } \text{ feature } f_j \text{ is active in } w^i \textbf{ then}$$
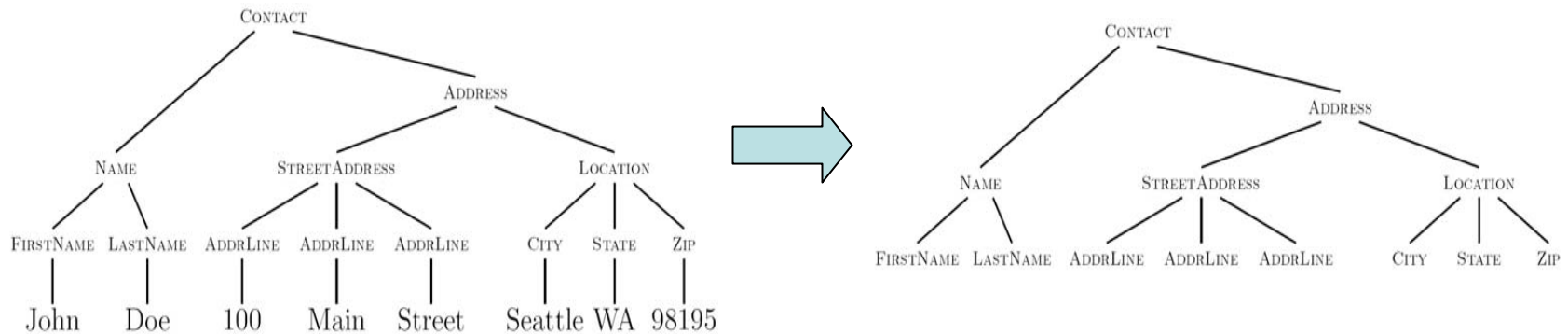$$\quad\quad\quad\quad\quad \lambda_j(R) \leftarrow \lambda_j(R) - 1;$$
$$\quad\quad\quad\quad \textbf{endif}$$
$$\quad\quad\quad \textbf{endfor}$$
$$\quad\quad\quad \textbf{for } \text{ each rule } R \text{ used in } T^j \text{ but not in } T \textbf{ do}$$
$$\quad\quad\quad\quad \textbf{if } \text{ feature } f_j \text{ is active in } w^i \textbf{ then}$$
$$\quad\quad\quad\quad\quad \lambda_j(R) \leftarrow \lambda_j(R) + 1;$$
$$\quad\quad\quad\quad \textbf{endif}$$
$$\quad\quad\quad \textbf{endfor}$$
$$\quad\quad \textbf{endif}$$
$$\quad \textbf{endfor}$$
$$\textbf{endfor}$$

# Results

| | Linear CRF | Discriminative CFG | Improvement |
|---|---|---|---|
| Word Error Rate | 11.57% | 6.29% | 45.63% |
| Record Error Rate | 54.71% | 27.13% | 50.41% |

# Outline

- PCFGs: Inference and Learning
- Parsing English
- Discriminative CFGs
- Grammar Induction

# Gold's Theorem '67

"Any formal language which has hierarchical structure capable of infinite recursion is unlearnable from positive evidence alone."

# Empirical Problems

- Even finite search spaces can be too big
- Noise
- Insufficient data
- Many local optima

# Common Approach

- Minimize total description length
- Simulated Annealing

Initial state

$$\begin{aligned}
D &:= \text{pabikugolatuda...} \\
T &:= T_0 \\
G &:= \begin{cases}
\gamma \to p\ \gamma \\
\gamma \to a\ \gamma \\
\gamma \to b\ \gamma \\
\vdots
\end{cases}
\end{aligned}$$

Repeat:

$$\begin{aligned}
G' &:= random\_neighbor(G) \\
\Delta &:= Energy(G', D) - Energy(G, D) \\
p &:= \begin{cases}
1 & \Delta \leq 0 \\
e^{-\frac{\Delta}{T}} & \Delta > 0
\end{cases} \\
G &:= G' \text{ with probability } p \\
T &:= \alpha T
\end{aligned}$$

Slide based on "Unsupervised grammar induction with Minimum Description Length" by Roni Katzir

# random_neighbor(G)

- Insert:

$$G := \begin{cases} A \to B\ C \\ B \to D\ E \end{cases} \Rightarrow G := \begin{cases} A \to B\ X\ C \\ B \to D\ E \end{cases}$$

- Delete

$$G := \begin{cases} A \to B\ C \\ B \to D\ E \end{cases} \Rightarrow G := \begin{cases} A \to B \\ B \to D\ E \end{cases}$$

- New Rule

$$G := \begin{cases} A \to B\ C \\ B \to D\ E \end{cases} \Rightarrow G := \begin{cases} A \to B\ C \\ B \to D\ E \\ Y \to \end{cases}$$

- Split

$$G := \begin{cases} A \to B\ C \\ B \to D\ E \end{cases} \Rightarrow G := \begin{cases} A \to D\ E\ C \\ B \to D\ E \end{cases}$$

- Substitute

$$G := \begin{cases} A \to B\ C \\ B \to D\ E \end{cases} \Rightarrow G := \begin{cases} A \to Z\ C \\ B \to D\ E \\ Z \to B \end{cases}$$

# Energy

$$Energy(G, D) \quad := \quad |G| + |code(D|G)|$$

Define binary representation for G, code(D|G)

$$G \quad := \quad \begin{cases} A \to & B \; A \\ A \to & B \\ B \to & C \; D \\ \vdots \\ E \to & F \; G \end{cases}$$

$$G \quad := \quad \texttt{ABA\#AB\#BCD\#...\#EFG\#\#}$$

# Experiment 1

- Word segmentation by 8-month old infants
- Vocabulary: `pabiku, golatu, daropi, tibudo`
- Saffran '96: use speech synthesizer, no word breaks, 2 minutes = 180 words
- Infants can distinguish words from non-words
- Now try grammar induction (60 words)

# Experiment 1

Step 37: current temp. = 14.994450998883487 Grammar:
$\gamma \to g\ \gamma; \gamma \to l\ \gamma; \gamma \to a\ \gamma; \gamma \to p\ \gamma; \gamma \to k\ \gamma; \gamma \to .\ \gamma; \gamma \to$
$t\ \gamma; \gamma \to u\ \gamma; \gamma \to o\ \gamma; \gamma \to d\ \gamma; \gamma \to b\ \gamma; \gamma \to i\ \gamma; \gamma \to$
$r\ \gamma; l \to$
Grammar length: 146 Encoding length: 1442

Step 618: current temp. = 14.907585393190937
Grammar:
$k \to k\ g\ 23; i \to p; 23 \to o; \gamma \to; \gamma \to g\ \gamma; \gamma \to$ *p a* $\gamma; \gamma \to$
$a\ \gamma; \gamma \to t\ \gamma; \gamma \to u\ \gamma; \gamma \to o\ \gamma; \gamma \to$ *k u* $\gamma; \gamma \to$ *o l* $\gamma; \gamma \to$
*o p* $\gamma; \gamma \to d\ \gamma; \gamma \to b\ \gamma; \gamma \to i\ \gamma; \gamma \to r\ \gamma; o \to; t \to g$
Grammar length: 200 Encoding length: 1199

Step 5837: current temp. = 14.149508793558308
Grammar: $k \to p; \gamma \to$ *t i b u d o* $\gamma; \gamma \to \gamma; \gamma \to$
*p a b i k u* $\gamma; \gamma \to i\ a\ \gamma\ \gamma\ g; \gamma \to$ *d a r o p i* $\gamma; \gamma \to$
*g o l a t u* $\gamma; b \to u\ d\ o\ o\ \gamma; b \to k\ u\ l; d \to; l \to; t \to \gamma$
Grammar length: 179 Encoding length: 183

# Experiment 2

$$S \rightarrow NP\ VP$$
$$NP \rightarrow Nm \mid D\ N$$
$$Nm \rightarrow max \mid sam \mid kim \mid bill \mid mary$$
$$D \rightarrow the \mid a$$
$$N \rightarrow man \mid dog \mid cat \mid turtle$$
$$VP \rightarrow Vin \mid Vtr\ NP \mid Vtl\ NP\ CP \mid Vsy\ CP$$
$$Vin \rightarrow walks \mid runs$$
$$Vtr \rightarrow kills \mid hits$$

. . .

► Sample sentences:

```
aturtleknowsthatsamtellskimthattheturtlewalksands
amkillsmax.kimknowsthatsamtellsmaxthatkimtellsk
imthatkimhitstheman.kimtellsaturtlethatmaxruns
```

# Experiment 2

Step : 423000 Temperature : 26.2

$1 \rightarrow$ *aman*$1$; $1 \rightarrow$ *hits*$1o$; $m \rightarrow h$; $1 \rightarrow$ *bill*$1m$; $m \rightarrow au$; $m \rightarrow \epsilon$; $1 \rightarrow$ *or*$1$; $1 \rightarrow$ *knowsthat*$1u$; $1 \rightarrow$ *urtle*$1$; $a \rightarrow n$; $1 \rightarrow$ e*man*; $m \rightarrow a$; $1 \rightarrow$ *heui*; $1 \rightarrow$ e*dog*; $1 \rightarrow$ *saysthat*$1$; $t \rightarrow e$; $a \rightarrow x$; $1 \rightarrow$ *wak*$1e$; $a \rightarrow o$; $1 \rightarrow$ *and*$1$; $1 \rightarrow \epsilon$; $1 \rightarrow$ *tells*$1$; $1 \rightarrow$ *walks*$1c$; $1 \rightarrow$ *raac*; $1 \rightarrow$ *runs*$1$; $a \rightarrow uu$; $1 \rightarrow$ *acat*$1cx$; $1 \rightarrow x$; $1 \rightarrow$ *kim*$1ky$; $m \rightarrow o$; $a \rightarrow \epsilon$; $1 \rightarrow$ *th*$11$; $1 \rightarrow$ *at*$1w$; $m \rightarrow r$; $1 \rightarrow$ *ma*$11$; $1 \rightarrow$ *et*; $1 \rightarrow r$; $1 \rightarrow$ *ecat*; $1 \rightarrow$ *adog*$1ty$; $1 \rightarrow$ *inksthat*; $1 \rightarrow ry$; $1 \rightarrow$ *sam*$1os$; $1 \rightarrow$ *kills*$1$; $1 \rightarrow$ *a*

Grammar length: 778 Encoding length: 10620 Energy: 11398

- Accurate segmentation
- Inaccurate structural learning

# Prototype-Driven Grammar Induction

- Semi-supervised approach
- Give only a few dozen prototypical examples (for NP e.g. determiner noun, pronouns, …)
- On English Penn Treebank: F1 = 65.1 (52% reduction over naïve PCFG induction)

Aria Haghighi and Dan Klein.
*Prototype-Driven Grammar Induction*.
ACL 2006

Dan Klein and Chris Manning.
*A Generative Constituent-Context Model for Improved Grammar Induction*.
ACL 2001

# That's it!