

Basic issues on Parsing ₁

- Introduction
- Parsing issues
- Parsing context free grammar (CFG)
- Robust parsing
- Statistical parsing

Basic issues on Parsing ₂

- Parsing goals
 - Syntactic structure
 - Semantic structure
- Syntax/semantic interaction
 - Only syntax
 - Only semantics
 - Performing in sequence
 - Performing in parallel

Basic issues on Parsing ₃

- Parsing as searching in a search space
 - Characterizing the states
 - (if possible) enumerate them
 - Define the initial state (s)
 - Define (if possible) final states or the condition to reach one of them

Basic issues on Parsing ₃

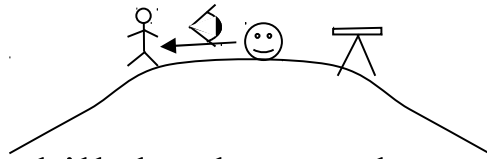
Ambiguity in parsing

A sentence is structurally ambiguous if the grammar assigns it more than a possible parse.

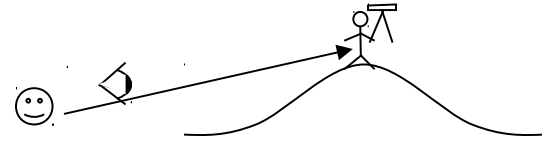
Common kinds of structural ambiguity include:

- PP-attachment
- Coordination ambiguity

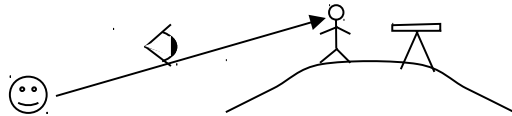
Basic issues on Parsing 5



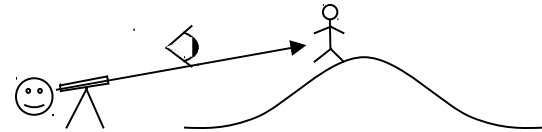
“I was on the hill that has a telescope when I saw a man.”



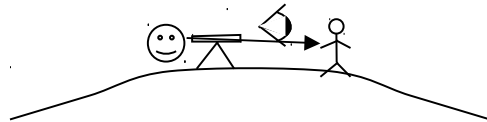
“I saw a man who was on a hill and who had a telescope.”



“I saw a man who was on the hill that has a telescope on it.”



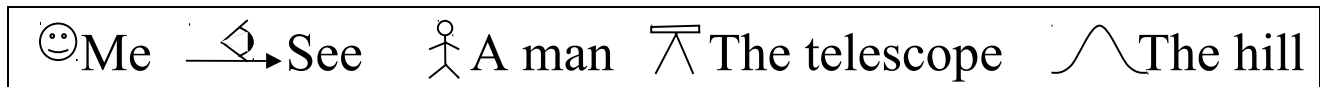
“Using a telescope, I saw a man who was on a hill.”



“I was on the hill when I used the telescope to see a man.”

...

I saw the man on the hill with the telescope



Basic issues on Parsing ₄

Factors in parsing

- Grammar expressivity
- Coverage
- Involved Knowledge Sources
- Parsing strategy
- Parsing direction
- Ambiguity management
- (in)determinism
- Parsing engineering

Basic issues on Parsing ₆

- Parsers today
 - Context free grammars (extended or not)
 - Tabular
 - Charts
 - Others
 - Unification-based
 - Statistical
 - Dependency parsing
 - Robust parsing (shallow, fragmental, chunkers, spotters)

Basic issues on Parsing 7

Parsing strategies

Parsing can be viewed as a search problem

Two common architectural approaches for this search are

Top-down: Starting with the root **S** and growing trees down to the input words

Bottom-up: Starting with the words and growing trees up toward the root **S**.

Basic issues on Parsing ₈

CFG grammar

Non terminal

sentence \rightarrow NP VP

NP \rightarrow det n

NP \rightarrow n

VP \rightarrow vi

VP \rightarrow vt NP

Terminal

det \rightarrow the

n \rightarrow cat

n \rightarrow fish

v \rightarrow eats

the cat eats fish

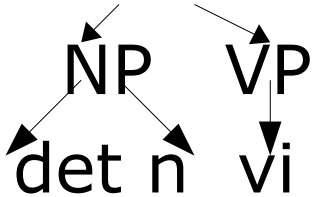
Basic issues on Parsing 9

Top-down
First step

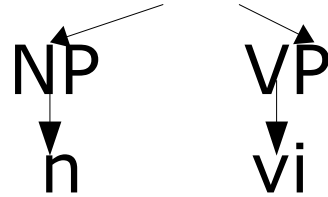


Second step

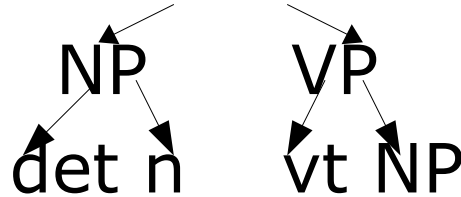
Sentence



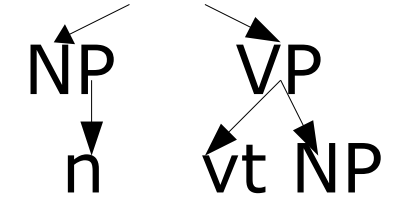
Sentence



Sentence



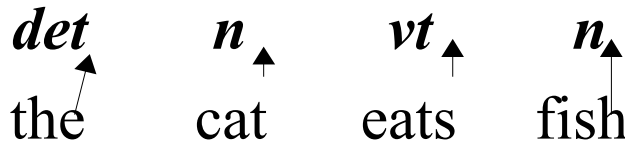
Sentence



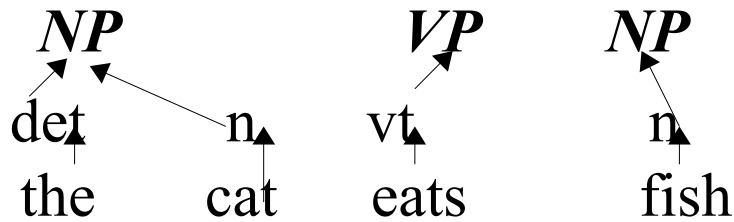
Basic issues on Parsing ₁₀

Bottom-up

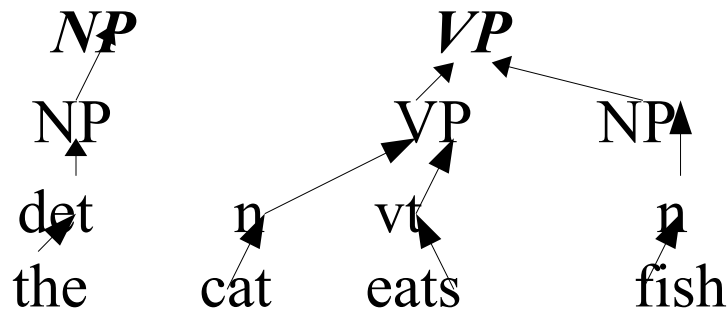
First step



Second step



Third step



Fourth step



Basic issues on Parsing ₁₁

Parsing strategy

- **Top Down**

- Guided by goals
- Starts with a goal (or set of goals) to be built.
- Tries to solve one of the pending goals
- If more than one production can be applied:
 - search problem
- Pending goals can be reordered
- Several search criteria (including heuristics) can be applied
- The process ends when all the goals have been reached

Parsing strategy

- **Bottom up**
 - Data driven
 - Starts from the sequence of words to be parsed (facts)
 - Proceeds bottom up
 - Several search criteria (including heuristics) can be applied
 - The process ends when the list of facts contains the initial symbol of the grammar

Basic issues on Parsing ₁₃

- Problems of ***top down*** strategy
 - Left recursivity
 - Many productions expanding the same non terminal
 - Useless work
 - Search basically guided by the grammar
 - Repeated work
 - Problems of backtracking algorithms

Basic issues on Parsing ₁₄

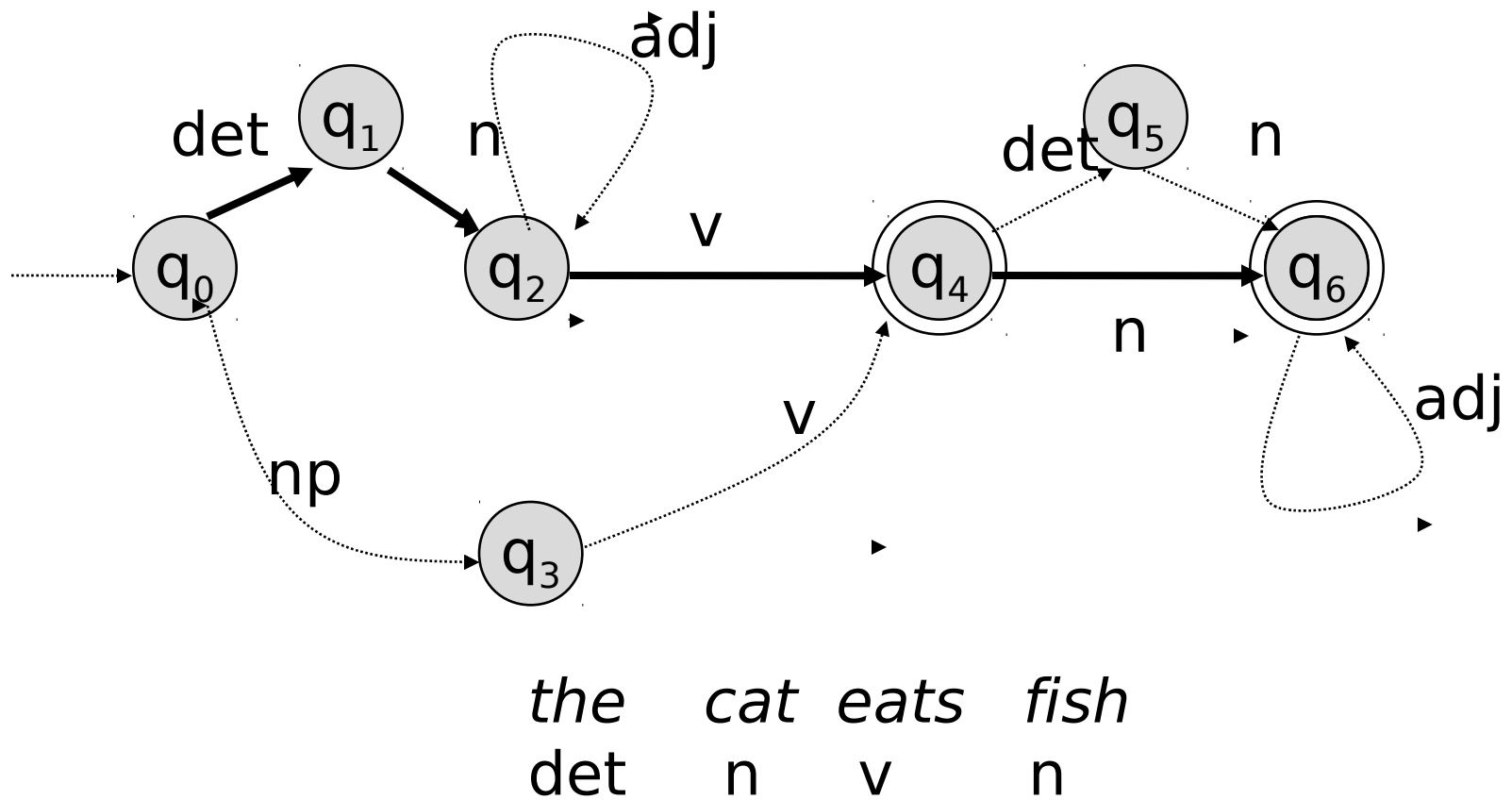
- Problems of ***bottom up*** parsing
 - empty (optional) categories
 - Useless work (locally possible but globally impossible)
 - Inefficient when there is a high lexical ambiguity
 - Repeated work

Transition Networks ₁

Finite state automata -> Transition Network (TN)

- States associated to the positions in the sentence
- Arcs (transitions)
 - Labeled with part of speech (POS)
 - An arc can be traversed if the current word has the same POS as the arc.
- Non determinism
 - More than one initial state
 - Current word with more than 1 POS
 - More than one arc for the same POS

Transition Networks ₂



Transition Networks ₃

Transition networks limitations

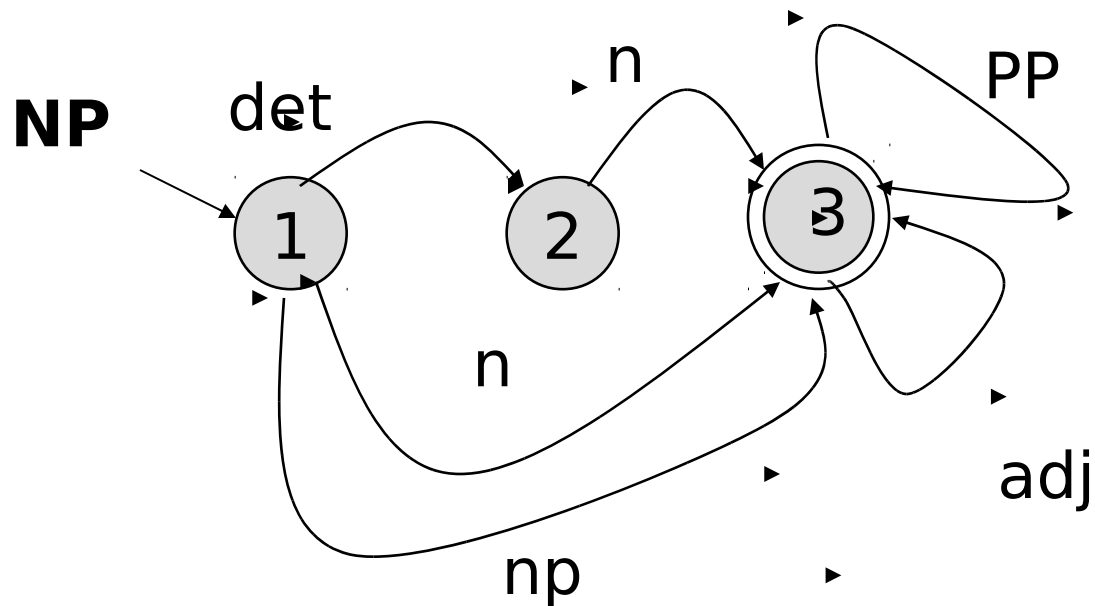
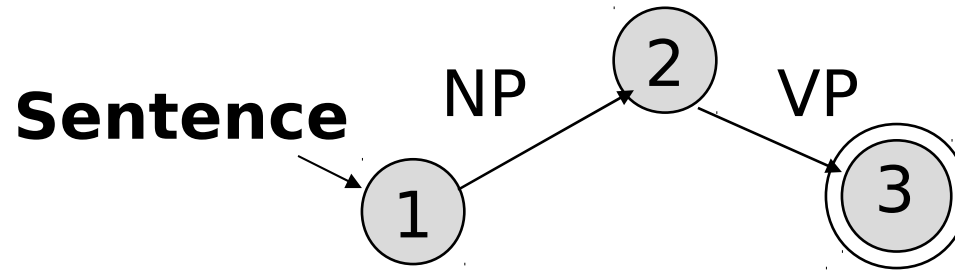
- Only regular grammars
- Only recognition
- Non-determinism -> backtracking
- No separation between grammar and parser
 - grammar -> syntactic model description
 - parser -> control

Transition Networks ₄

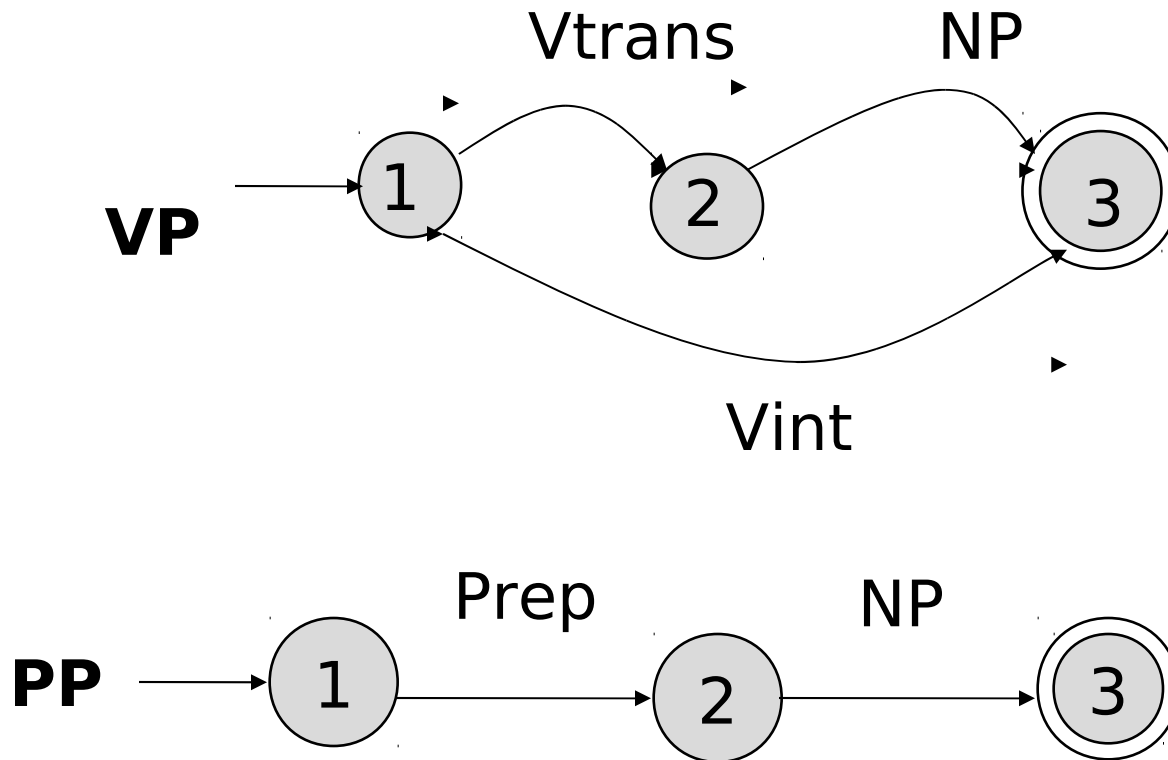
Recursive transition networks (RTN)

- Collection of TNs labeled with a name
 - Arcs
 - Labeled as in TN with POS
 - Terminal labels
 - Labeled with RTN identifiers
 - Non terminal labels
 - Final states in RTN produce coming back to the target state of the arc producing the call
- RTN are weakly equivalent to CFG

Transition Networks ₅



Transition Networks ⁶



Transition Networks ₇

Recursive transition networks (RTN) Limitations

- Transitions depend only on the categories
 - CFG
- Only recognizing
- In fact fixed top-down strategy

Transition Networks ₈

- Woods (1970)
- **Aumented transition networks (ATN)** = RTN with *operations* attached to arcs and use of *registers*

Operations

Conditions

Filter transitions between states

Actions

Building intermediate and output structures.

Initializations

- ATN allow expressing contextual constraints

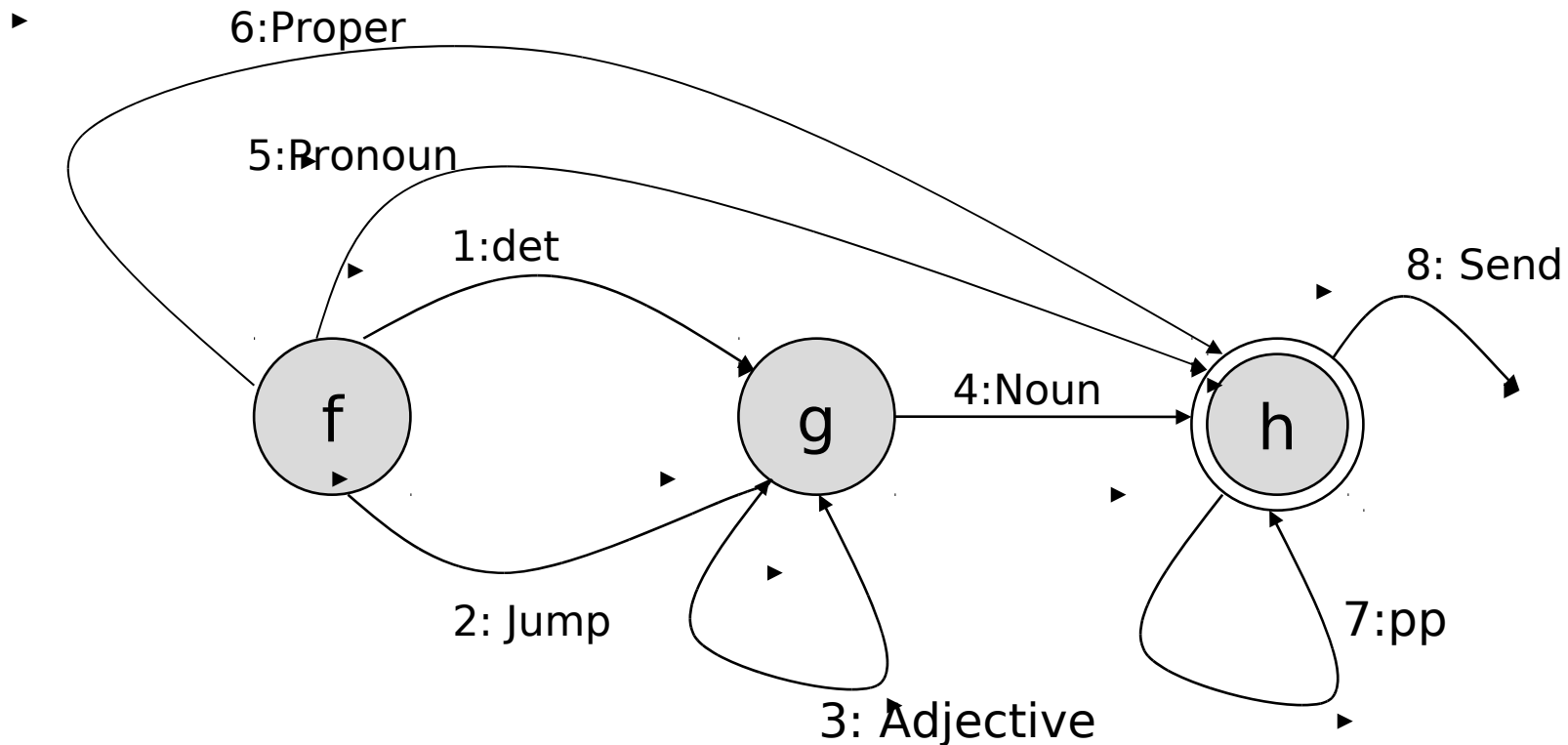
Transition Networks ₉

Features

Number: Singular, Plural Default: empty

Person: 1st, 2nd, 3rd Default: 3rd

Rols: Subject



Transition Networks 10

Initializations, Conditions and Actions

NP-1: f **Determiner** g

A: Set Number to the number of *

NP-4: g **Noun** h

C: Number is empty or number is the number of *

A: Set Number to the number of *
Set Subject to *

NP-5: f **Pronoun** h

A: Set Number to the number of *
Set Person to the Person of *
Set Subject to *

NP-6: f **Proper** h

A: Set Number to the number of *
Set Subject to *

Transition Networks ₁₁

Aumented Transition networks limitations

- Fixed top-down strategy
- Redundancy in backtracking operations
- Problems of notational expressivity:
 - Very difficult to transport

Basic issues on Parsing ₁₁

- Unified mechanism of parser description
 - Sikkel, 1997
- **Parser (schema):**
 - Given a sentence, an initial set of items is build
 - Given a grammar, a set of rules can be used for getting additional items
- **Parser (algorithm):**
 - Parsing schema
 - + data structures
 - + control structures
 - (+ communication structures)

Charts₁

- A *Chart* is a directed graph built dynamically along parsing
- Extension of WFST
- Nodes correspond to the start and end of the sentence and to the positions between words.
- Active arcs (goals or hypothesis) and inactive arcs (facts)
 - Notation active arcs: dotted rules
 - inactive arcs : category

0 1 2 3 4
○ the ○ cat ○ eats ○ fish ○

Charts₂

CFG grammar

Non terminal

sentence \rightarrow NP VP

NP \rightarrow det n

NP \rightarrow n

VP \rightarrow vi

VP \rightarrow vt NP

Terminal

det \rightarrow the

n \rightarrow cat

n \rightarrow fish

v \rightarrow eats

the cat eats fish

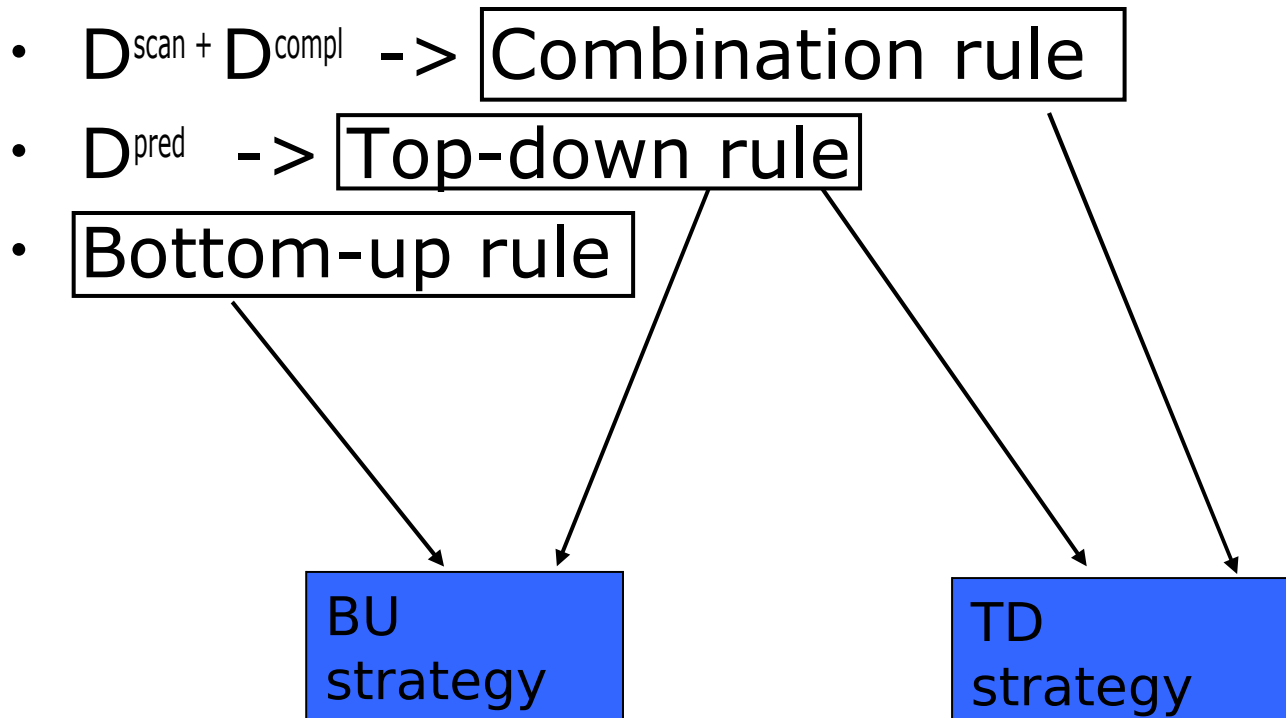
Charts₃

program chart

```
{ initialize the chart with H;  
  initialize the agenda with items which can be deduced without  
  antecedents;  
  while not empty (agenda)  
  {extract current_item from agenda and put it on the chart;  
    foreach item which could be deduced with one step including  
    current_item  
    {if item not in agenda and not in chart  
      then add item to agenda  
    }  
  }  
}
```

Charts₄

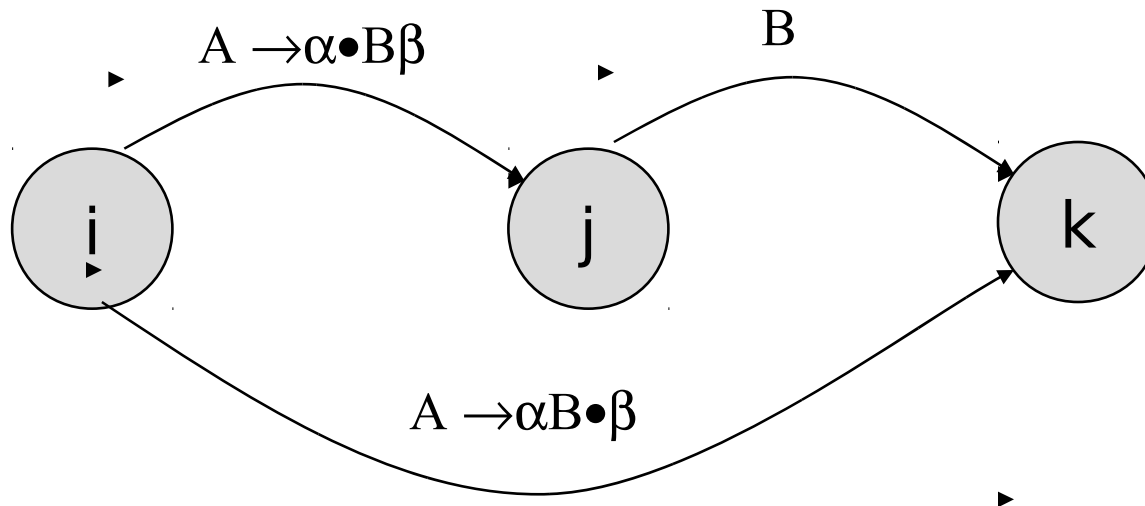
- A concrete chart algorithm should:
 - define the structure of *agenda* and its scheduling criteria
 - define order of performing deductive steps



Charts₅

Combination rule

When an active arc of the Chart reaches a node j and from this node starts an inactive arc labeled with the category the active arc was waiting for, both arcs are combined for building a new arc (active or not) starting in the start node of the active arc and ending in the ending node of the inactive arc.

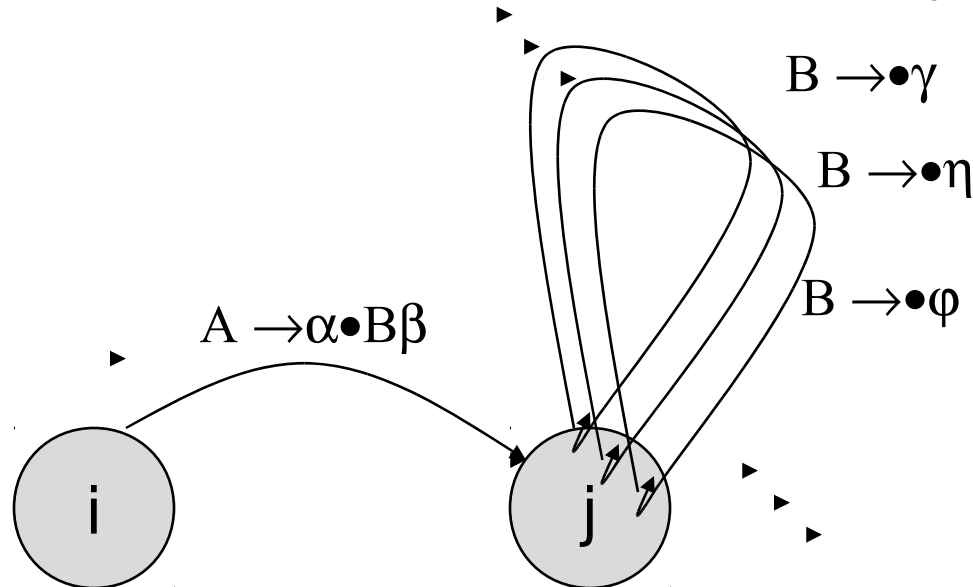


Charts₆

Top-down rule

When an active arc of the Chart reaches a node j , for all the productions of the grammar expanding the category the active arc is waiting for a new active arc is built starting and ending in j corresponding to the dotted rule with dot in the initial position.

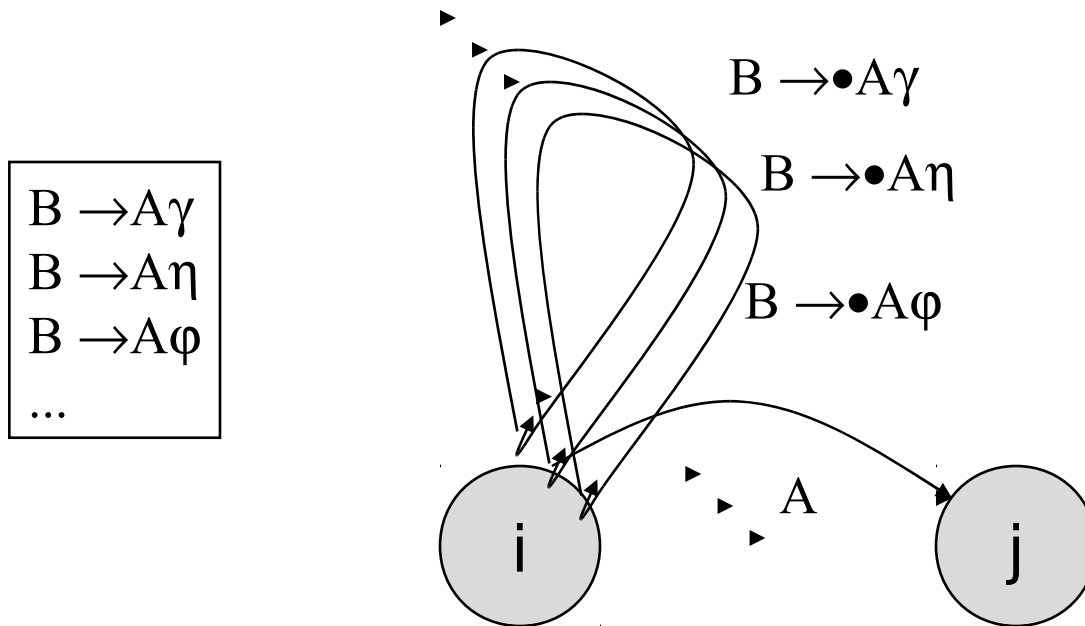
- $B \rightarrow \gamma$
- $B \rightarrow \eta$
- $B \rightarrow \varphi$
- ...



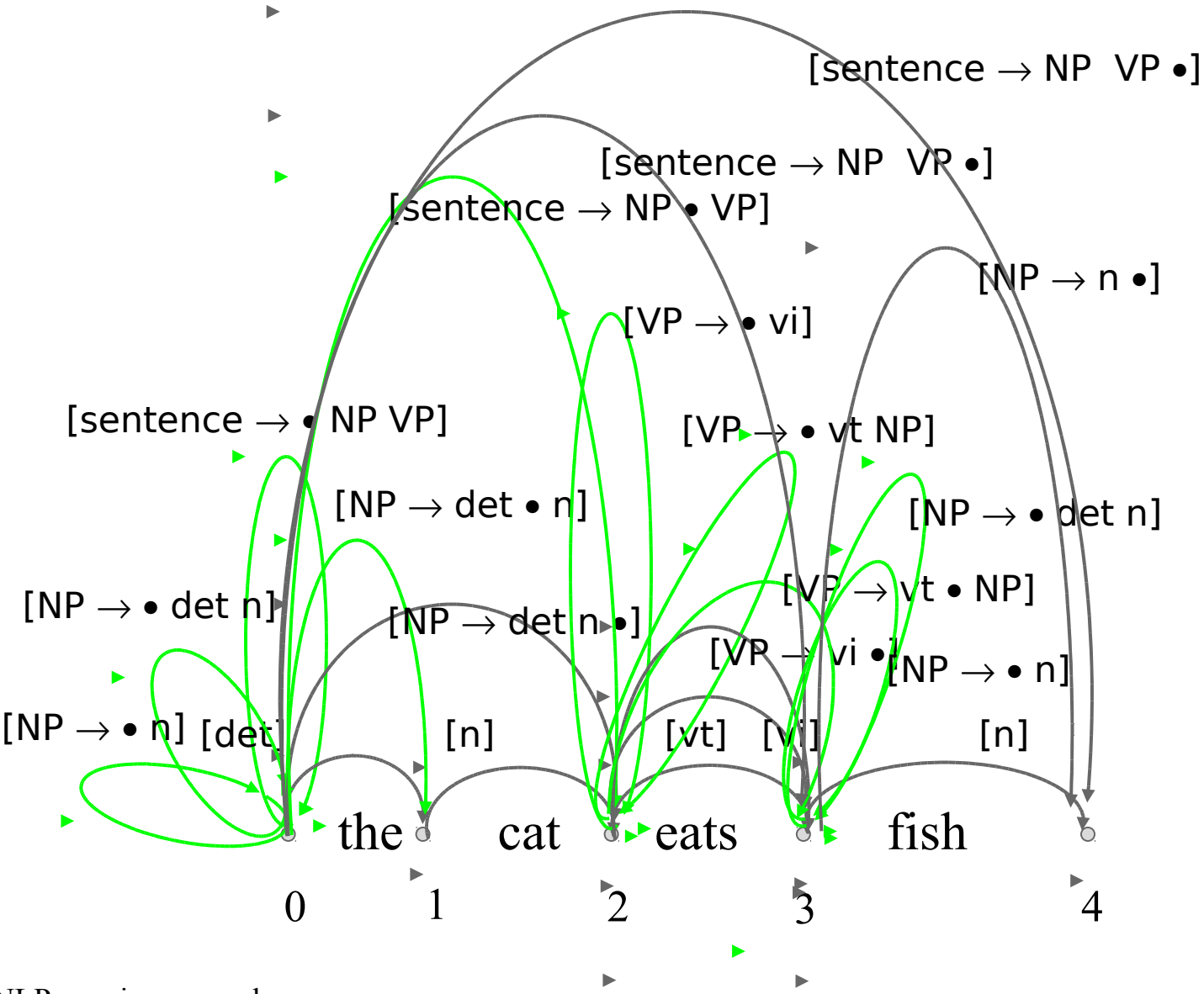
Charts₇

Bottom-up rule

When an inactive arc of the Chart starts in a node i , for each Production of the grammar owning as first copnstituent of the right side the category of the inactive arc a new active arc is built starting and ending in i corresponding to the dotted rule with dot in the initial position



Charts₈



Charts₉

- Problems
 - The size of the chart grows with the size of the grammar making the algorithm difficult to scale up.
 - A lot of useless active and inactive arcs are built.
 - In practice, lacking appropriate knowledge, a fixed bottom-up strategy, eventually corrected with top-down predictions, is used

Charts₁₀

- Ambiguity combined with the repeated parsing of sub-trees are a difficulty for parsing algorithms. Those algorithms use simple backtracking mechanisms.
- There are parsing algorithms that use dynamic programming techniques, such as a table of partial parsers to efficiently parse ambiguous sentences.
- The CKY, Earley and Chart-Parsing algorithms all use dynamic-programming to solve the repeated parsing of subtrees problem.

Robust Parsing₁

- Conventional methods of parsing are insufficient for dealing with non restricted texts

- Adequate segmentation
- Disambiguation
- Coverage

Which are the units to parse

Which is the most likely parse

Parsing beyond the lexical coverage

- What to do

- Not parsing all
- Not parsing in depth

fragmental parsing

shallow parsing

Robust Parsing₂

- Problems when parsing non restricted corpus
- Adaptation of a grammar to a corpus or sublanguage
- Selection of the correct (!?) parse between the ones allowed by the grammar
- Production of good parses for entries outside the coverage of the grammar (Robustness)

Robust Parsing₂

- Partial parsing and chunking are methods for identifying shallow syntactic constituents in a text.
- High accuracy partial parsing can be achieved either through rule-based or machine learning-based methods.

Robust Parsing₃

Partial parsers

phrasal parsers

- chunkers, spotters
- Church,1988

cooccurrence parsers

- Church,Hanks,1989, Brent,1993

fragmental parsers

- Fidditch, Hindle,1994, MITFP, Abney,1991

constraint-based parsers

- Voutilainen,1995

Robust Parsing₄

- Chunking
 - detection of phrases nominal, verbal, adjectival, adverbial basic (without recursion)
 - Finite state techniques
 - Performance of a cascade of transducers
 - Hidden Markov Models
 - Machine Learning

Abney, 1996
Argamon et al, 1998
Cardie, Pierce, 1998
Church, 1988
Ramshaw, Marcus, 1995
Skut, Brants, 1998

Robust Parsing₅

Definition of chunk

With linguistic basis: Abney

Only pragmatic:

- Contiguous sequences of related tokens
 - Not confusing with terms
- e.g. Base NP

Approaches to chunking

Look for (include) information

Remove information

- e.g. Chink

Robust Parsing₆

Representing chunks

Labels

- e.g. tags
- BEGIN, INSIDE, OUTSIDE

Trees

Chunk parser

Looking for non overlapped chunks for reaching a maximal coverage

Robust Parsing₇

Frequently regular expressions over sequences of POS tags

agglomerative (chunk rules) vs divisive (chunk rules)

- Rules for fusion of adjacent chunks
- Rules for splitting a chunk in smaller components

Statistical Parsing₁

Using **statistical models** for

- Guiding parsing
 - Get the most likely parse
 - Ambiguity resolution (pp-attachment)
 - Grammatical induction from corpora

Goal: Parsing of non restricted texts with a reasonable level of accuracy (>90%) and efficiency.

Requirements:

Corpora tagged (with POS): Brown, LOB, Clic-Talp

Corpora analyzed: Penn treebank, Susanne, Ancora

Statistical Parsing₂

- Lexical approaches
 - Context free: unigram
 - Context dependent: N-gram, HMM
- Syntactic approaches
 - Statistical context free grammar (SCFG)
- Hybrid approaches
 - Stochastic lexicalized Tags
 - Computing the most likely (most probable) parse

Statistical Parsing₃

Probabilistic grammars assign a probability to a sentence or string of words.

Usually they capture more general syntactic information than the N-gram grammars.

In a probabilistic context-free grammar (PCFG):

- Associate a probability to each rule
- Associate a probability to each lexical entry

Each PCFG rule is treated as if it were conditionally independent; thus the probability of a sentence is computed by multiplying the probabilities of each rule in the parse of the sentence.

Statistical Parsing₄

S	⇒	NP	VP	1.0
VP	⇒	Vi		0.4
VP	⇒	Vt	NP	0.4
VP	⇒	VP	PP	0.2
NP	⇒	DT	NN	0.3
NP	⇒	NP	PP	0.7
PP	⇒	P	NP	1.0

Vi	⇒	sleeps	1.0
Vt	⇒	saw	1.0
NN	⇒	man	0.7
NN	⇒	woman	0.2
NN	⇒	telescope	0.1
DT	⇒	the	1.0
IN	⇒	with	0.5
IN	⇒	in	0.5

- Probability of a tree t with rules

$$\alpha_1 \rightarrow \beta_1, \alpha_2 \rightarrow \beta_2, \dots, \alpha_n \rightarrow \beta_n$$

is

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

where $q(\alpha \rightarrow \beta)$ is the probability for rule $\alpha \rightarrow \beta$.

Statistical Parsing₅

1. A context-free grammar $G = (N, \Sigma, S, R)$.
2. A parameter

$$q(\alpha \rightarrow \beta)$$

for each rule $\alpha \rightarrow \beta \in R$. The parameter $q(\alpha \rightarrow \beta)$ can be interpreted as the conditional probability of choosing rule $\alpha \rightarrow \beta$ in a left-most derivation, given that the non-terminal being expanded is α . For any $X \in N$, we have the constraint

$$\sum_{\alpha \rightarrow \beta \in R: \alpha = X} q(\alpha \rightarrow \beta) = 1$$

In addition we have $q(\alpha \rightarrow \beta) \geq 0$ for any $\alpha \rightarrow \beta \in R$.

Given a parse-tree $t \in \mathcal{T}_G$ containing rules $\alpha_1 \rightarrow \beta_1, \alpha_2 \rightarrow \beta_2, \dots, \alpha_n \rightarrow \beta_n$, the probability of t under the PCFG is

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

Statistical Parsing₆

- Assigns a probability to each *left-most derivation*, or parse-tree, allowed by the underlying CFG
- Say we have a sentence s , set of derivations for that sentence is $\mathcal{T}(s)$. Then a PCFG assigns a probability $p(t)$ to each member of $\mathcal{T}(s)$. i.e., *we now have a ranking in order of probability.*
- The most likely parse tree for a sentence s is

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

Statistical Parsing₇

- ▶ **Learning.** How to obtain a PCFG from a treebank?
- ▶ **Inference.** Given a PCFG and a sentence s . Denote by $\mathcal{T}(s)$ the set of derivations that yield s .
 - ▶ How to compute the best parse for s ?

$$\operatorname{argmax}_{t \in \mathcal{T}(s)} p(t)$$

- ▶ How to compute the probability of s ?

$$\sum_{t \in \mathcal{T}(s)} p(t)$$

Statistical Parsing₇

Pros and cons of SCFG

Some idea of the probability of a parse but not very good
CFG cannot be learned without negative examples, SCFG
can SCFGs provide a language model for a language

In practice SCFG provide a worse language model than a
3-gram

$P([N [N \text{ toy}] [N [N \text{ coffee}] [N \text{ grinder}]]) =$

$P([N [N [N \text{ cat}] [N \text{ food}]] [N \text{ tin}]))$

$P(NP \rightarrow \text{Pro})$ is $>$ in Subj position than in Obj position.

Statistical Parsing₈

- Robust
- Possibility of combining SCFG with 3-grams
- SCFG assign a lot of probability mass to short sentences (a small tree is more probable than a big one)
- Parameter estimation (probabilities)
- Problem of sparseness
- Volume

Statistical Parsing,

- Association to the rule. Information about the point of application of the rule in the derivation tree is lost.
- Low frequency constructions are penalized
- Probability of a derivation. Contextual independence is assumed
- Possibility of relax conditional independence:
 - Sensitivity to structure
 - Lexicalization

Statistical Parsing₁₀

Sensitivity to structure

Node expansion depends of its position in the tree

Pronoun Lexical

Subject	91%	9%
Object	34%	66%

Enrichment of a node with information of its ancestors

- ^SNP is different of ^VP NP

Pronouns as arguments of ditransitive verbs

- *I gave Charlie the book*
- *I gave the book to Charlie*
- *I gave you the book*
- ? *I gave the book to you*

Statistical Parsing₁₁

(Head) Lexicalization

put takes both an NP and a VP

- *Sue put [the book]_{NP} [on the table]_{PP}*
- * *Sue put [the book]_{NP}*
- * *Sue put [on the table]_{PP}*

like usually takes an NP and not a PP

- *Sue likes [the book]_{NP}*
- * *Sue likes [on the table]_{PP}*

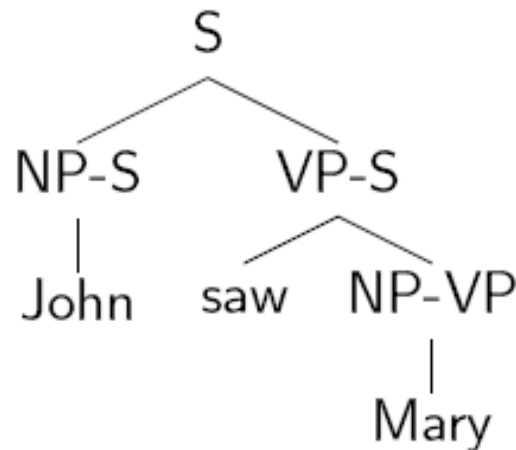
Statistical Parsing₁₂

<i>Local Tree</i>	<i>Come</i>	<i>Take</i>	<i>Think</i>	<i>Want</i>
VP-> V	9.5%	2.6%	4.6%	5.7%
VP-> V NP	1.1%	32.1%	0.2%	13.9%
VP-> V PP	34.5%	3.1%	7.1%	0.3%
VP- V SBAR	6.6%	0.3%	73.0%	0.2%
VP-> V S	2.2%	1.3%	4.8%	70.8%
VP->V NP S	0.1%	5.7%	0.0%	0.3%
VP->V PRT NP	0.3%	5.8%	0.0%	0.0%
VP->V PRT PP	6.1%	1.5%	0.2%	0.0%

Statistical Parsing₁₃

PCFGs with Parent Annotations

(Johnson, 1999)

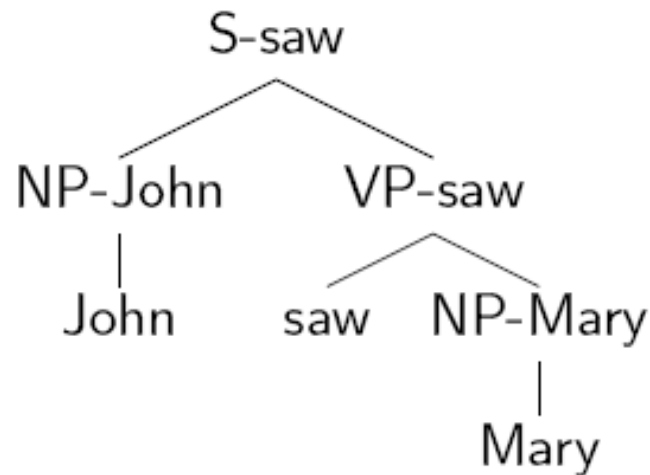


$$P(Tree) = q(S \rightarrow NP-S \ VP-S \mid S) \times \\ q(NP-S \rightarrow John \mid NP-S) \times \\ \dots$$

Statistical Parsing₁₄

Lexicalized PCFGs

(Collins, 1999)

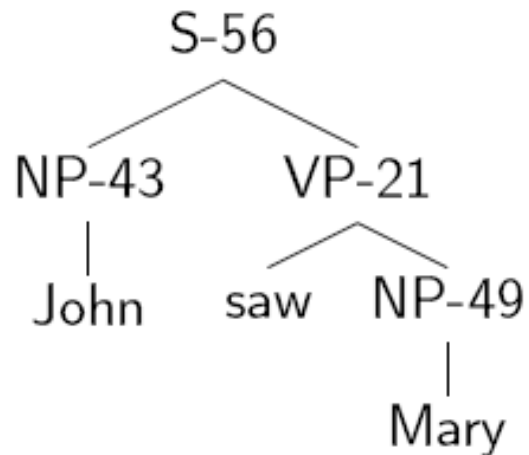


$$P(Tree) = P(S-saw \rightarrow NP-John \ VP-saw \mid S-saw) \times \\ P(NP-John \rightarrow John \mid NP-John) \times \\ \dots$$

Statistical Parsing₁₅

PCFGs with Latent Variables

(e.g., Petrov and Klein, 2007)



- ▶ Each non-terminals (e.g., S) is split into a number of new non-terminals (e.g., S-1, S-2, ..., S-128)
- ▶ Latent annotations learned using EM

Statistical Parsing₁₅

Two models

Conditional/Discriminative model :

The probability of a parse tree is directly estimated

$$P(t | s, G) \text{ con } \sum_t P(t | s, G) = 1$$

Probabilities are conditioned on a concrete sentence.

No sentence distribution probabilities is assumed

Statistical Parsing₁₆

- CFG
- SCFG
 - For each rule of G , $(A \rightarrow \alpha) \in P_G$ we should be able to define a probability $P(A \rightarrow \alpha)$

$$\sum_{(A \rightarrow \alpha) \in P_G} P(A \rightarrow \alpha) = 1$$

- Probability of a tree

$$P(\psi) = \prod_{(A \rightarrow \alpha) \in P_G} P(A \rightarrow \alpha) f(A \rightarrow \alpha; \psi)$$

Statistical Parsing₁₇

$P(t)$ -- Probability of a tree t (product of probabilities of the rules generating it).

$P(w_{ln})$ -- Probability of a sentence is the sum of the probabilities of all the valid parse trees of the sentence

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

Statistical Parsing₁₈

- Positional invariance:

The probability of a subtree is independent of its position in the derivation tree

- Context-free

Independence from ancestors

Statistical Parsing₁₉

Parameter estimation

- Supervised learning
 - From a treebank
- Non supervised learning

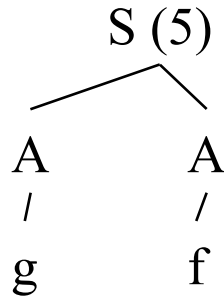
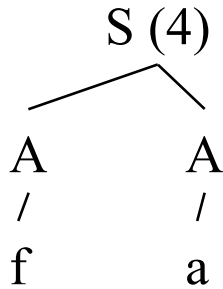
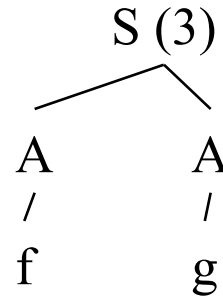
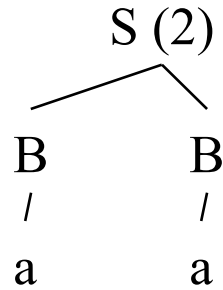
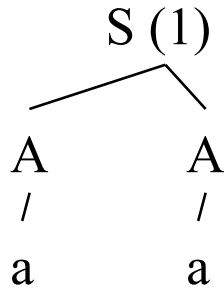
Statistical Parsing₂₀

A Penn Treebank tree (POS tags not shown)

```
( (S (NP-SBJ The move)
    (VP followed
      (NP (NP a round)
        (PP of
          (NP (NP similar increases)
            (PP by
              (NP other lenders))
            (PP against
              (NP Arizona real estate loans))))))
    (S-ADV (NP-SBJ *)
      (VP reflecting
        (NP (NP a continuing decline)
          (PP-LOC in
            (NP that market))))))
  .))
```

Treebank grammars₁

Consider a treebank containing the following trees



Treebank grammars₂

Suppose that (1) occurs 40 times, (2) occurs 10 times, (3) occurs 5 times, (4) occurs 5 times, and (5) occurs once.

We want to induce a SCFG reflexing this treebank.

Parameter estimation

$$\sum_j P(N^i \rightarrow \zeta^j | N^i) = 1$$

Treebank grammars₃

Rules

$$S \rightarrow A A \quad : \quad 40 + 5 + 5 + 1 = 51$$

$$S \rightarrow B B \quad : \quad 10$$

$$A \rightarrow a \quad : \quad 40 + 40 + 5 = 85$$

$$A \rightarrow f \quad : \quad 5 + 5 + 1 = 11$$

$$A \rightarrow g \quad : \quad 5 + 1 = 6$$

$$B \rightarrow a \quad : \quad 10$$

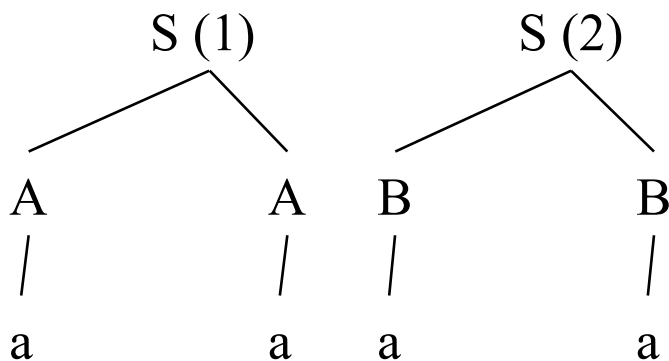
Treebank grammars₄

Parameters maximizing global likelihood of the corpus:

G	Frequency	Total	Probability
$S \rightarrow A A$	51	61	0.836
$S \rightarrow B B$	10	61	0.164
$A \rightarrow a$	85	102	0.833
$A \rightarrow f$	11	102	0.108
$A \rightarrow g$	6	102	0.059
$B \rightarrow a$	10	10	1.0

Treebank grammars₅

Given this parametrization, what is the most likely parse for "a a"?



$$\begin{aligned} P(1) &= P(S \rightarrow A A) * P(A \rightarrow a) * P(A \rightarrow a) \\ &= 0.836 * 0.833 * 0.833 = 0.580 \end{aligned}$$

$$\begin{aligned} P(2) &= P(S \rightarrow B B) * P(B \rightarrow a) * P(B \rightarrow a) \\ &= 0.164 * 1.0 * 1.0 = 0.164 \end{aligned}$$

Treebank grammars₈

- Removing non linguistically valid rules
- Assign probabilities (MLE) to the initial rules
- Remove a rule except in the the probability of the structure built from its application is greater than the probability of building applying simpler rules.
- Thresholding Removing rules occurring $< n$ times

	Full	Simply thresholded	Fully compacted	Linguistically Compacted Grammar 1	Linguistically Compacted Grammar 2
Recall	70.55	70.78	30.93	71.55	70.76
precision	77.89	77.66	19.18	72.19	77.21
Grammar size	15,421	7,278	1,122	4,820	6,417