

### Last Class: Parsing

1. Grammars and parsing
2. Top-down and bottom-up parsing
3. A top-down parser

### Today: Parsing

1. Bottom-up chart parsing
2. Earley algorithm

Slide CS674-1

### General Parsing Strategies

Grammar	Top-Down	Bottom-Up
1. $S \rightarrow NP VP$	$S \rightarrow NP VP$	$\rightarrow \text{NAME ate the cat}$
2. $VP \rightarrow V NP$	$\rightarrow \text{NAME VP}$	$\rightarrow \text{NAME V the cat}$
3. $NP \rightarrow \text{NAME}$	$\rightarrow \text{Beav VP}$	$\rightarrow \text{NAME V ART cat}$
4. $NP \rightarrow \text{ART N}$	$\rightarrow \text{Beav V NP}$	$\rightarrow \text{NAME V ART N}$
5. $\text{NAME} \rightarrow \text{Beavis}$	$\rightarrow \text{Beav ate NP}$	$\rightarrow \text{NP V ART N}$
6. $V \rightarrow \text{ate}$	$\rightarrow \text{Beav ate ART N}$	$\rightarrow \text{NP V NP}$
7. $\text{ART} \rightarrow \text{the}$	$\rightarrow \text{Beav ate the N}$	$\rightarrow \text{NP VP}$
8. $N \rightarrow \text{cat}$	$\rightarrow \text{Beav ate the cat}$	$\rightarrow S$

Slide CS674-2

### Algorithm for a Top-Down Parser

$PSL \leftarrow (((S) 1))$

1. *Check for failure.* If PSL is empty, return NO.
2. *Select the current state, C.*  $C \leftarrow \text{pop}(PSL)$ .
3. *Check for success.* If  $C = (()) <\text{final-position}>$ , YES.
4. *Otherwise, generate the next possible states.*
  - (a)  $s_1 \leftarrow \text{first-symbol}(C)$
  - (b) If  $s_1$  is a *lexical symbol* and next word can be in that class, create new state by removing  $s_1$ , updating the word position, and adding it to PSL.
  - (c) If  $s_1$  is a *non-terminal*, generate a new state for each rule in the grammar that can rewrite  $s_1$ . Add all to PSL.

Slide CS674-3

### Problems with the Top-Down Parser

1. Only judges grammaticality.
2. Stops when it finds a single derivation.
3. No semantic knowledge employed.
4. No way to rank the derivations.
5. Problems with left-recursive rules.
6. Problems with ungrammatical sentences.
7. **Terribly inefficient!!**

*Have the first year Phd students in the computer science department take the Q-exam.*

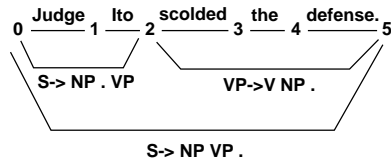
*Have the first year Phd students in the computer science department taken the Q-exam?*

Slide CS674-4

## Chart Parsers

**chart:** data structure that stores partial results of the parsing process in such a way that they can be reused. The chart for an  $n$ -word sentence consists of:

- $n + 1$  **vertices**
- a number of **edges** that connect vertices



Slide CS674-5

## Chart Parsing: The General Idea

The process of parsing an  $n$ -word sentence consists of forming a chart with  $n + 1$  vertices and adding edges to the chart one at a time.

- Goal: To produce a complete edge that spans from vertex 0 to  $n$  and is of category  $S$ .
- There is no backtracking.
- Everything that is put in the chart stays there.
- Chart contains all information needed to create parse tree.

Slide CS674-6

## Bottom-UP Chart Parsing Algorithm

Do until there is no input left:

1. If the agenda is empty, get next word from the input, look up word categories, add to agenda (as constituent spanning two positions).
2. Select a constituent from the agenda: constituent  $C$  from  $p_1$  to  $p_2$ .
3. Insert  $C$  into the chart from position  $p_1$  to  $p_2$ .
4. For each rule in the grammar of form  $X \rightarrow C X_1 \dots X_n$ , add an active edge of form  $X \rightarrow C \circ X_1 \dots X_n$  from  $p_1$  to  $p_2$ .

Slide CS674-7

5. Extend existing edges that are looking for a  $C$ .

- (a) For any active edge of form  $X \rightarrow X_1 \dots \circ C X_n$  from  $p_0$  to  $p_1$ , add a new active edge  $X \rightarrow X_1 \dots C \circ X_n$  from  $p_0$  to  $p_2$ .
- (b) For any active edge of form  $X \rightarrow X_1 \dots X_n \circ C$  from  $p_0$  to  $p_1$ , add a new (completed) constituent of type  $X$  from  $p_0$  to  $p_2$  to the agenda.

Slide CS674-8

## Grammar and Lexicon

**Grammar:**

1. $S \rightarrow NP VP$	3. $NP \rightarrow ART ADJ N$
2. $NP \rightarrow ART N$	4. $VP \rightarrow V NP$

**Lexicon:**

the: ART	man: N, V
old: ADJ, N	boat: N

**Sentence:** <sub>1</sub> The <sub>2</sub> old <sub>3</sub> man <sub>4</sub> the <sub>5</sub> boat <sub>6</sub>

1.  $S \rightarrow NP VP$
2.  $NP \rightarrow ART N$
3.  $NP \rightarrow ART ADJ N$
4.  $VP \rightarrow V NP$

the: ART	man: N, V
old: ADJ, N	boat: N

Slide CS674-9

The diagram illustrates the hierarchical structure of the sentence "The old man the boat." using a syntax tree. The root node is **S (rule 1)**. It branches into **NP2 (rule 3)** and **VP2 (rule 4)**. **NP2 (rule 3)** branches into **NP1 (rule 2)** and **VP1**. **NP1 (rule 2)** branches into **ART1** ("The") and **ADJ1** ("old"). **VP1** branches into **N2** ("man") and **V1** (implied). **VP2 (rule 4)** branches into **NP1 (rule 2)** and **N3** ("boat."). **NP1 (rule 2)** branches into **ART2** ("the") and **N3** ("boat.").

Below the tree, numbered nodes (1-6) are shown, corresponding to the words in the sentence. Arrows indicate the semantic interpretation rules for each node:

- 1 → NP → ART . N
- 2 → NP → ART . ADJ N
- 3 → NP → ART ADJ . N
- 4 → S → NP . VP
- 5 → VP → V . NP
- 6 → S → NP . VP

NP2 (rule 3) VP2 (rule 4)

NP1 (rule 2) YP1 NP1 (rule 2)

1 The 2 old 3 man 4 the 5 boat. 6

ART1 ADJ1 N2 V1 ART2 N3

NP->ART . N

NP->ART . ADJ N

NP -> ART ADJ . N

S -> NP . VP

VP -> V . NP

S -> NP . VP

Slide CS674-12

## Bottom-up Chart Parser

Is it any less naive than the top-down parser?

1. Only judges grammaticality.[fixed]
2. Stops when it finds a single derivation.[fixed]
3. No semantic knowledge employed.
4. No way to rank the derivations.
5. Problems with ungrammatical sentences.[better]
6. Terribly inefficient[???

1. Only judges grammaticality.[fixed]
2. Stops when it finds a single derivation.[fixed]
3. No semantic knowledge employed.
4. No way to rank the derivations.
5. Problems with ungrammatical sentences.[better]
6. Terribly inefficient[???

Slide CS674-11

### Top-Down Chart Parser

For all S rules of the form  $S \rightarrow X_1 \dots X_k$ , add a (top-down) edge from 1 to 1 labeled:  $S \rightarrow \circ X_1 \dots X_k$ .

Do until there is no input left:

1. If the agenda is empty, look up word categories for next word, add to agenda.
2. Select a constituent from the agenda: constituent  $C$  from  $p_1$  to  $p_2$ .
3. Using the (bottom-up) edge extension algorithm, combine  $C$  with every active edge on the chart (adding  $C$  to chart as well). Add any new constituents to the agenda.
4. For any active edges created in Step 3, add them to the chart using the top-down edge introduction algorithm.

Slide CS674–13

*Top-down edge introduction.*

To add an edge  $S \rightarrow C_1 \dots \circ C_i \dots C_n$  ending at position  $j$ :

For each rule in the grammar of form  $C_i \rightarrow X_1 \dots X_k$ ,  
recursively add the new edge  $C_i \rightarrow \circ X_1 \dots X_k$  from  $j$  to  $j$ .

Slide CS674–14

### Grammar and Lexicon

Grammar

Lexicon

- |                               |                |
|-------------------------------|----------------|
| 1. $S \rightarrow NP VP$      | the: ART       |
| 2. $NP \rightarrow ART ADJ N$ | large: ADJ     |
| 3. $NP \rightarrow ART N$     | can: N, AUX, V |
| 4. $NP \rightarrow ADJ N$     | hold: N, V     |
| 5. $VP \rightarrow AUX VP$    | water: N, V    |
| 6. $VP \rightarrow V NP$      |                |

Sentence: <sub>1</sub> The <sub>2</sub> large <sub>3</sub> can <sub>4</sub> can <sub>5</sub> hold <sub>6</sub> water <sub>7</sub>

Slide CS674–15