

Handout 16: How parsers work

Top-down parser

1. Consider the grammar:

$S \rightarrow$ [a] NP VP
 $NP \rightarrow$ [a] Det N N [b] Det N
 $VP \rightarrow$ [a] V [b] V NP

2. Lexicon:

$Det \rightarrow$ [a] the
 $N \rightarrow$ [a] police [b] man [c] boat
 $V \rightarrow$ [a] man [b] boat

3. Parse “the police man the boat”

- a. Start from S.

1 [0] * S

- b. **Expand** category after * in every way possible. First way: $S \rightarrow NP$ VP.

1 [1a] * NP VP
2 [2a] * Det N N VP
3 [3a] * the N N VP

- c. Bottomed out. It **matches** the sentence. Advance the *.

1 [4] the * N N VP

- d. Now back to **expanding**.

1 [5a] the * police N VP
2 [6] the police * N VP
3 [7a] the police * police VP

- e. Does not match. Go to an alternative.

1 [7b] the police * man VP
2 [8] the police man * VP
3 [9a] the police man * V
4 [10a] the police man * man
5 [10b] the police man * boat
6 [9b] the police man * V NP
7 [11a] the police man * man
8 [11b] the police man * boat
9 [7c] the police * boat VP
10 [5b] the * man N VP
11 [5c] the * boat N VP
12 [2b] * Det N VP

f. Finally, we're on the right track

```
1      [12a] * the N VP
2      [13] the * N VP
3      [14a] the * police VP
4      [15] the police * VP
5      [16a] the police * V
6      [17a] the police * man
7      [18] the police man *
```

g. We have generated a complete tree, but the sentence is not done.

```
1      [17b] the police * boat
2      [16b] the police * V NP
3      [19a] the police * man NP
4      [20] the police man * NP
5      [21a] the police man * Det N N
6      [22a] the police man * the N N
7      [23] the police man the * N N
8      [24a] the police man the * police N
9      [24b] the police man the * man N
10     [24c] the police man the * boat N
11     [25] the police man the boat * N
```

h. We consumed all of the input sentence, but we don't have a complete tree.

```
1      [21b] the police man * Det N
2      [26a] the police man * the N
3      [27] the police man the * N
4      [28a] the police man the * police
5      [28b] the police man the * man
6      [28c] the police man the * boat
```

i. Success! But we continue, to see if there are more parses.

```
1      [19b] the police * boat NP
2      [14b] the * man VP
3      [14c] the * boat VP
```

j. All options exhausted.

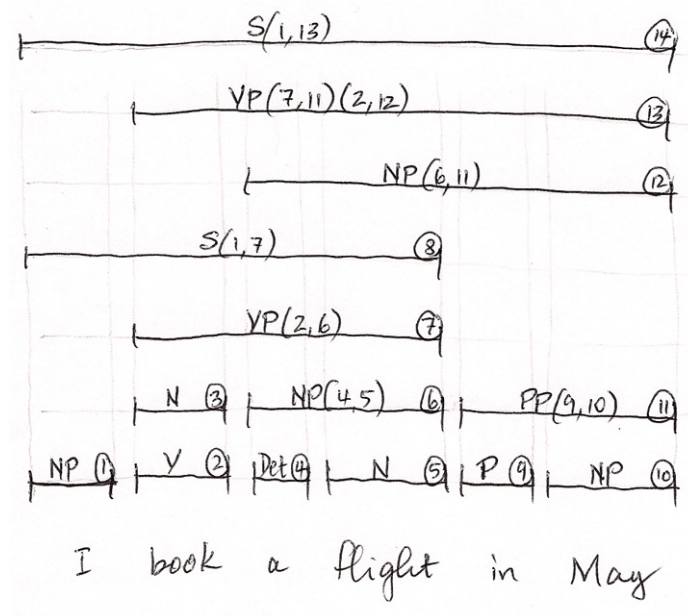
Bottom-up parsing: CKY

4. Grammar

- | | | | |
|---|------------------------|----|------------------------|
| 1 | $S \rightarrow NP VP$ | 7 | $NP \rightarrow I$ |
| 2 | $NP \rightarrow Det N$ | 8 | $N \rightarrow book$ |
| 3 | $NP \rightarrow NP PP$ | 9 | $V \rightarrow book$ |
| 4 | $VP \rightarrow V NP$ | 10 | $Det \rightarrow a$ |
| 5 | $VP \rightarrow VP PP$ | 11 | $N \rightarrow flight$ |
| 6 | $PP \rightarrow P NP$ | 12 | $NP \rightarrow May$ |

5. Chart

a. The chart



b. Positions:

₀ I ₁ book ₂ a ₃ flight ₄ in ₅ May ₆

c. Chart nodes: $_i X_j$. E.g., #1: $_0 NP_1$, #7: $_1 VP_4$.

d. No duplicates: rather, assign multiple child-lists. E.g., #13: $_1 VP_6$.

6. **unwind()**. Read a tree out of the chart: make a non-deterministic choice wherever there is more than one child-list.

```

14          (S
01          (NP I)
13          (VP
07          (VP
02          (V book)
06          (NP
04          (Det a)
05          (N flight))))
11          (PP
09          (P in)
10          (NP May))))

```

```

14          (S
01          (NP I)
13          (VP
02          (V book)
12          (NP
06          (NP
04          (Det a)
05          (N flight)))
11          (PP
09          (P in)
10          (NP May))))))

```

7. Algorithm for **fill_chart()**
- Work from left to right: j ranges from 1 to 6
 - shift(j)**. Let w be the word ending at position j . For each rule $X \rightarrow w$, create a node $_{j-1}X_j$.
 - extend_edges($node$)**. Each time you create a node $_kZ_j$, look for nodes $_iY_k$ ending where $_kZ_j$ starts. If there is a rule $X \rightarrow YZ$, then create a new node $_iX_j$.
 - When you create a new node, record the two children that you used to create it.
 - Do not create duplicates: if you want to create a new node $_iX_j$, but it already exists, just add a new child-list to the old edge instead of creating a new one.
8. CKY parser requires grammar in **Chomsky Normal Form**. Only two kinds of rules allowed:

$$X \rightarrow YZ$$

$$X \rightarrow a$$

Bottom-up: chart parser

9. Example: figure 1
10. Two basic data structures
 - a. Node ${}_iX_j$
 - b. Edge ${}_iX \rightarrow \alpha \bullet_j \beta$
11. Four basic operations (bottom-up parser)
 - a. **shift**(j). For each part of speech X for word[$j - 1$], create a node ${}_{j-1}X_j$.
 - b. **bu_predict**($node$). Invoked after creating node ${}_iY_j$. For each rule $X \rightarrow Y\beta$, create edge ${}_iX \rightarrow Y \bullet_j \beta$.
 - c. **extend_edges**($node$). Invoked after creating node ${}_kZ_j$. For each edge ${}_iX \rightarrow \alpha \bullet_k Z\beta$, create edge ${}_iX \rightarrow \alpha Z \bullet_j \beta$.
 - d. **complete**($edge$). Invoked after creating edge ${}_iX \rightarrow \alpha \bullet_j$. Create node ${}_iX_j$. But if ${}_iX_j$ already exists, just add a new expansion instead.
12. Two supporting operations that stitch everything together:
 - a. **create_node**(i, X, j). Create node ${}_iX_j$ and put it in the chart. Then call **bu_predict** and **extend_edges**.
 - b. **create_edge**(i, X, α, j, β). Create the edge ${}_iX \rightarrow \alpha \bullet_j \beta$ and put it in the chart. If the dot is at the end, call **complete**.
13. Top level
 - a. **fill_chart**($words$). Initialize the chart, then call **shift** at each sentence position, from 1 to the length of the sentence. Each call to **shift** triggers a cascade of actions.
 - b. **unwind**(\cdot). Look for an S node spanning the whole sentence. If found, extract the list of trees from it and return them.

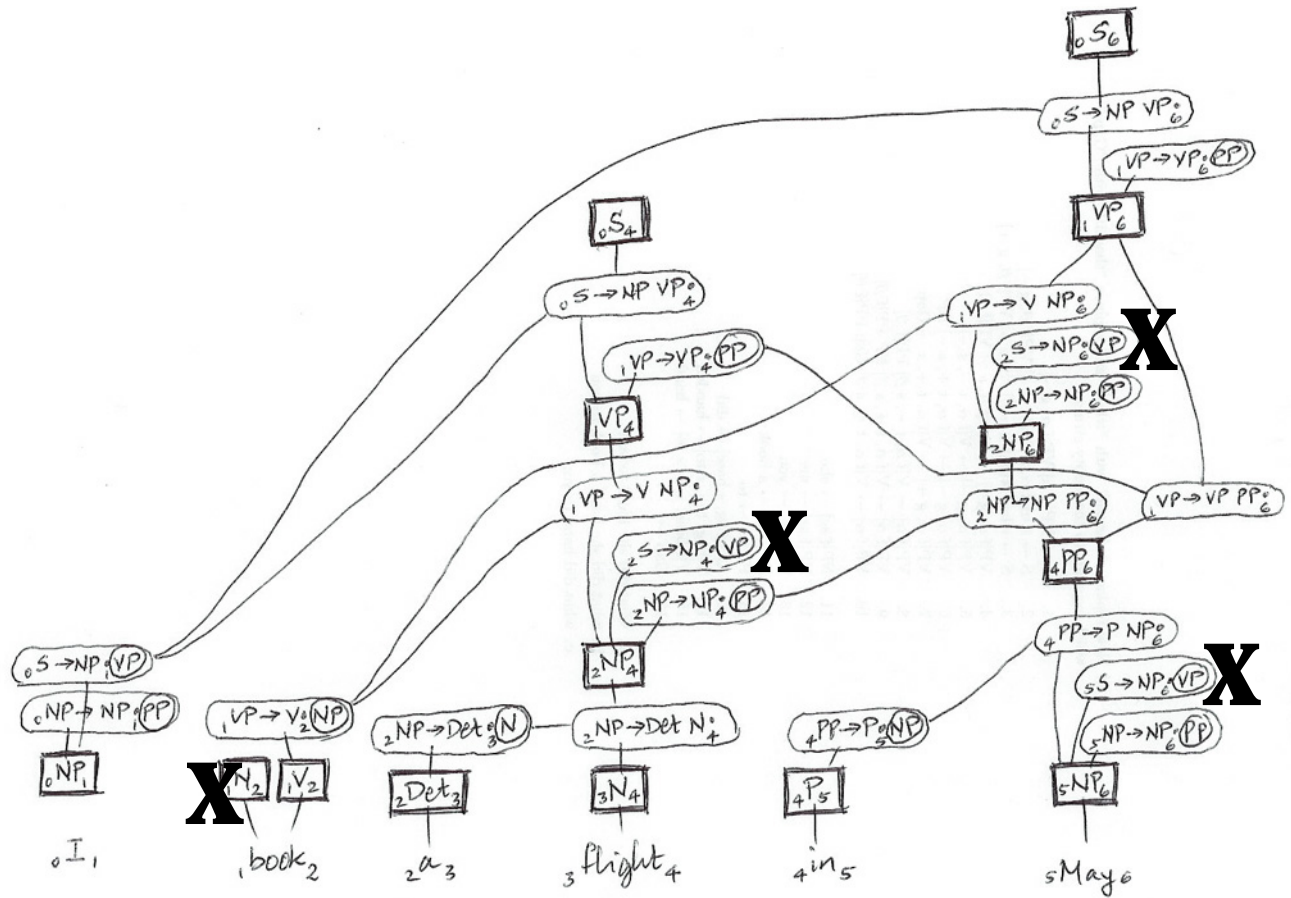


Figure 1: Filled chart for sentence "I book a flight in May." The node and edges marked with "X" are filtered out when topdown prediction is used. (Not shown is an initial prediction $\rightarrow \bullet_0 S$.)