# Bottom-up Chart Parsing: the CKY algorithm
Parsing
ISCL-BA-06

Çağrı Çöltekin
ccoltekin@sfs.uni-tuebingen.de

University of Tübingen
Seminar für Sprachwissenschaft
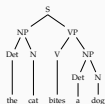
Winter Semester 2020/21

---

## Parsing so far

- Parsing is the task of automatic syntactic analysis
- For most practical purposes, context-free grammars are the most useful formalism for parsing
- We can formulate parsing as
  - Top-down: begin with the start symbol, try to *produce* the input string to be parsed
  - Bottom-up: begin with the input, and try to *reduce* it to the start symbol
- Both strategies can be cast as search with backtracking
- Backtracking parsers are inefficient: they recompute sub-trees multiple times

---

## Bottom-up parsing as search



S    → NP VP
NP   → Det N
VP   → V NP
VP   → V
Det  → a
Det  → the
N    → cat
N    → dog
V    → bites
N    → bites

---

## Dealing with ambiguity

S    → NP VP
NP → Prm N
NP → Prm
VP → V NP
VP → V
VP → V S
N    → duck
V    → duck
V    → saw
Prm → I                      ←
Prm → she
Prm → her

1
I    saw    her    duck

---

## Dealing with ambiguity

S    → NP VP
NP → Prm N
NP → Prm                     ←
VP → V NP
VP → V
VP → V S
N    → duck
V    → duck
V    → saw
Prm → I
Prm → she
Prm → her

Prm

2
I    saw    her    duck

---

## Dealing with ambiguity

S    → NP VP
NP → Prm N
NP → Prm
VP → V NP
VP → V
VP → V S
N    → duck
V    → duck
V    → saw                   ←
Prm → I
Prm → she
Prm → her

NP
Prm

3
I    saw    her    duck

---

## Dealing with ambiguity

S    → NP VP
NP → Prm N
NP → Prm
VP → V NP
VP → V
VP → V S
N    → duck
V    → duck
V    → saw
Prm → I
Prm → she
Prm → her                    ←

NP
Prm    V

4
I    saw    her    duck

---

## Dealing with ambiguity

S    → NP VP
NP → Prm N
NP → Prm                     ←
VP → V NP
VP → V
VP → V S
N    → duck
V    → duck
V    → saw
Prm → I
Prm → she
Prm → her

NP
Prm    V

Prm

5
I    saw    her    duck

---

## Dealing with ambiguity

S    → NP VP
NP → Prm N
NP → Prm
VP → V NP
VP → V
VP → V S
N    → duck                  ←
V    → duck
V    → saw
Prm → I
Prm → she
Prm → her

NP
Prm    V
            NP
            Prm

6
I    saw    her    duck

---

## Dealing with ambiguity

S    → NP VP
NP → Prm N
NP → Prm
VP → V NP
VP → V
VP → V S
N    → duck
V    → duck
V    → saw
Prm → I
Prm → she
Prm → her

NP
Prm    V
            NP
            Prm    V

7
I    saw    her    duck

---

## Dealing with ambiguity

S    → NP VP                 ←
NP → Prm N
NP → Prm
VP → V NP
VP → V
VP → V S
N    → duck
V    → duck
V    → saw
Prm → I
Prm → she
Prm → her

NP
Prm    V
            NP    VP
            Prm    V

8
I    saw    her    duck

---

## Dealing with ambiguity

S    → NP VP
NP → Prm N
NP → Prm
VP → V NP
VP → V
VP → V S                     ←
N    → duck
V    → duck
V    → saw
Prm → I
Prm → she
Prm → her

NP
Prm    V        S
            NP    VP
            Prm    V

9
I    saw    her    duck

## Dealing with ambiguity

Tree for: I saw her duck

```
S → NP VP      ←
NP → Prn N
NP → Prn
VP → V NP
VP → V
VP → V S
N → duck
V → duck
V → saw
Prn → I
Prn → she
Prn → her
```

10

---

## Dealing with ambiguity

Tree for: I saw her duck

```
S → NP VP
NP → Prn N
NP → Prn
VP → V NP
VP → V
VP → V S
N → duck
V → duck
V → saw
Prn → I
Prn → she
Prn → her
```

11

---

## Dealing with ambiguity

Trees for: I saw her duck

```
S → NP VP
NP → Prn N
NP → Prn
VP → V NP
VP → V
VP → V S
N → duck      ←
V → duck      ←
V → saw
Prn → I
Prn → she
Prn → her
```

12

---

## Dealing with ambiguity

Trees for: I saw her duck

```
S → NP VP
NP → Prn N      ←
NP → Prn
VP → V NP
VP → V
VP → V S
N → duck
V → duck
V → saw
Prn → I
Prn → she
Prn → her
```

13

---

## Dealing with ambiguity

Trees for: I saw her duck

```
S → NP VP
NP → Prn N
NP → Prn
VP → V NP      ←
VP → V
VP → V S
N → duck
V → duck
V → saw
Prn → I
Prn → she
Prn → her
```

14

---

## Dealing with ambiguity

Trees for: I saw her duck

```
S → NP VP      ←
NP → Prn N
NP → Prn
VP → V NP
VP → V
VP → V S
N → duck
V → duck
V → saw
Prn → I
Prn → she
Prn → her
```

15

---

## Dealing with ambiguity

Trees for: I saw her duck

```
S → NP VP
NP → Prn N
NP → Prn
VP → V NP
VP → V
VP → V S
N → duck
V → duck
V → saw
Prn → I
Prn → she
Prn → her
```

16

---

## How to represent multiple parses
parse forest grammar

Trees for: I saw her duck

$$S_{0:4} \to NP_{0:1}\ VP_{1:4}$$
$$NP_{0:1} \to Prn_{0:1}$$
$$Prn_{0:1} \to I_{0:1}$$
$$VP_{1:4} \to V_{1:2}\ S_{2:4}$$
$$V_{1:2} \to saw_{1:2}$$
$$S_{2:4} \to Prn_{2:3}\ V_{3:4}$$
$$V_{3:4} \to duck_{3:4}$$
$$VP_{1:4} \to V_{1:2}\ NP_{2:4}$$
$$NP_{2:4} \to Prn_{2:3}\ N_{3:4}$$

---

## CKY algorithm

- The CKY (Cocke-Kasami-Younger) parsing algorithm is a dynamic programming algorithm (Kasami 1965; Younger 1967; Cocke and Schwartz 1970)
- It processes the input *bottom up*, and saves the intermediate results on a *chart*
- Time complexity for *recognition* is $O(n^3)$
- Space complexity is $O(n^2)$
- It requires the CFG to be in *Chomsky normal form* (CNF) (can somewhat be relaxed, but not common)

---

## Chomsky normal form (CNF)

- A CFG is in CNF, if the rewrite rules are one of the following forms
  - A → B C
  - A → a
  where A, B, C are non-terminals and a is a terminal
- Any CFG can be converted to CNF
- Resulting grammar is *weakly equivalent* to the original grammar:
  - it generates/accepts the same language
  - but the derivations are different

---

## Converting to CNF: example

```
S   → NP VP
S   → Aux NP VP
NP  → the N
VP  → V NP
VP  → V
V   → V
N   → cat
N   → dog
V   → bites
N   → bites
```

- S → Aux NP VP
  S → Aux NP VP  ⇒  S → Aux X
                     X → NP VP

- NP → the N
  NP → the N  ⇒  NP → X N
                  X → the

- VP → V
  VP → V  ⇒  VP → bites

---

## Converting to CNF

1. Eliminate the ε rules: if A → ε is in the grammar
   - replace any rule B → α A β with two rules
     B → α β
     B → α A β
   - add A → α for all α (except ε) whose LHS is A
   - repeat the process for newly created ε rules
   - remove the rules with ε on the RHS (except S → ε)
2. Eliminate unit rules: for a rule A → B
   - Replace the rule with $A \to \alpha_1 \mid \ldots \mid \alpha_n$, where $\alpha_1, \ldots, \alpha_n$ are all RHS or rule B
   - Remove the rule A → B
   - Repeat the process until no unit rules remain
3. Binarize all the non-binary rules with non-terminal on the RHS: for a rule $A \to X_1\ X_2\ \ldots X_n$
   - Replace the rule with $A \to A_1\ X_2 \ldots X_n$, and add $A_1 \to X_1\ X_2$
   - Repeat the process until all new rules are binary

# CKY demonstration
*an ambiguous example*

| | |
|---|---|
| S | → NP VP |
| NP | → Prn N |
| VP | → V NP |
| VP | → V S |
| N | → duck |
| VP | → duck \| saw |
| V | → duck \| saw |
| Prn | → I \| she \| her |
| NP | → I \| she \| her |

0  I  1  saw  2  her  3  duck  4

Prn, NP   V, VP   Prn, NP   N, V, VP

S → NP VP

VP → V NP

NP → Prn N
S → NP VP

S → NP VP

# CKY demonstration
an ambiguous example

S → NP VP



S   → NP VP
NP  → Prn N
VP  → V NP
VP  → V S
N   → duck
VP  → duck | saw
V   → duck | saw
Prn → I | she | her
NP  → I | she | her

---

# CKY demonstration
an ambiguous example



S   → NP VP
NP  → Prn N
VP  → V NP
VP  → V S
N   → duck
VP  → duck | saw
V   → duck | saw
Prn → I | she | her
NP  → I | she | her

---

# CKY demonstration
an ambiguous example



S   → NP VP
NP  → Prn N
VP  → V NP
VP  → V S
N   → duck
VP  → duck | saw
V   → duck | saw
Prn → I | she | her
NP  → I | she | her

---

# CKY demonstration: the chart
our chart is a 2D array



Space complexity is $O(n^2)$.

---

# CKY demonstration: the chart
our chart is a 2D array – this is more convenient for programming
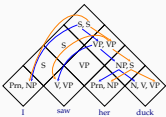


Space complexity is $O(n^2)$.

---

# Parsing vs. recognition

- We went through a recognition example
- Note that the algorithm is not directional: it takes the complete input
- Recognition accepts or rejects a sentence based on a grammar
- For parsing, we want to know the derivations that yielded a correct parse
- To recover parse trees, we
  - we follow the same procedure as recognition
  - add back links to keep track of the derivations

---

# Chart parsing example (CKY parsing)



The chart stores a *parse forest* efficiently.

---

# Summary

- \+ CKY avoids re-computing the analyses by storing the earlier analyses (of sub-spans) in a table
- − It still computes lower level constituents that are not allowed by the grammar
- − CKY requires the grammar to be in CNF
- CKY has $O(n^3)$ recognition complexity
- For parsing we need to keep track of the backlinks
- CKY can efficiently store all possible parses in a chart
- Enumerating all possible parses have exponential complexity (worst case)

Next:

- Top-down chart parsing: Earley algorithm
- Suggested reading: Grune and Jacobs (2007, section 7.2)

---

# Acknowledgments, references, additional reading material

Cocke, John and J. T. Schwartz (1970). *Programming languages and their compilers: preliminary notes.* Tech. rep. Courant Institute of Mathematical Sciences, NYU.

Grune, D. and C.J.H. Jacobs (2007). *Parsing Techniques: A Practical Guide.* second. Monographs in Computer Science. The first edition is available at http://dickgrune.com/Books/PTAPG_1st_Edition/BookBody.pdf. Springer New York. ISBN: 9780387689548.

Kasami, Tadao (1965). *An Efficient Recognition and Syntax-Analysis Algorithm for Context-Free Languages.* Tech. rep. ITTC Document.

Younger, Daniel H (1967). "Recognition and parsing of context-free languages in time $n^3$". In: *Information and control* 10.2, pp. 189–208.