# Bottom-up Chart Parsing: the CKY algorithm
### Data Structures and Algorithms for Computational Linguistics III
### (ISCL-BA-07)

Çağrı Çöltekin
ccoltekin@sfs.uni-tuebingen.de

University of Tübingen
Seminar für Sprachwissenschaft

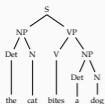Winter Semester 2023/24

---

## 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

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 |

---

## Dealing with ambiguity

Prn

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 |

---

## Dealing with ambiguity

NP
Prn

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 |

---

## Dealing with ambiguity

NP
Prn V

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 ← |

---

## Dealing with ambiguity

NP
Prn V

Prn

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 |

---

## Dealing with ambiguity

NP
Prn V

NP
Prn

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 |

---

## Dealing with ambiguity

NP
Prn V

NP
Prn V

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 |

---

## Dealing with ambiguity

NP
Prn V

NP V
Prn

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 |

---

## Dealing with ambiguity

NP
Prn V

NP VP
Prn V

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 |

---

## Dealing with ambiguity

NP
Prn V S

NP VP
Prn V

I saw her duck

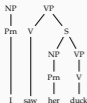| | |
|---|---|
| S | → NP VP |
| NP | → Prn N |
| NP | → Prn |
| VP | → V NP |
| VP | → V |
| VP | → V S ← |
| N | → duck |
| V | → duck |
| V | → saw |
| Prn | → I |
| Prn | → she |
| Prn | → her |

## Dealing with ambiguity

(parse tree for "I saw her duck")

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

---

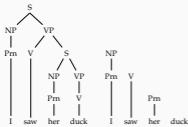## Dealing with ambiguity

(parse 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

---

## Dealing with ambiguity

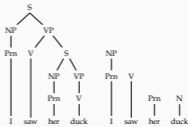(parse tree for "I saw her duck I saw her duck")

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

---

## Dealing with ambiguity

(parse tree for "I saw her duck I saw her duck")
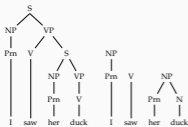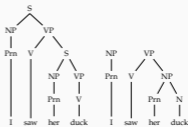
S → NP VP
NP → Prn N ←
NP → Prn
VP → V NP
VP → V
VP → V S
N → duck
V → duck
V → saw
Prn → I
Prn → she
Prn → her

---

## Dealing with ambiguity

(parse tree for "I saw her duck 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

---

## Dealing with ambiguity

(parse tree for "I saw her duck I saw her duck")
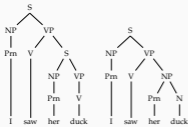
S → NP VP ←
NP → Prn N
NP → Prn
VP → V NP
VP → V
VP → V S
N → duck
V → duck
V → saw
Prn → I
Prn → she
Prn → her

---

## Dealing with ambiguity

(parse trees for "I saw her duck 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

---

## How to represent multiple parses
### parse forest grammar

(parse trees for "I saw her duck 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
- 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 in 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
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
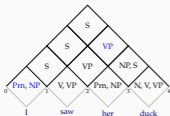
---

## 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 a (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 → $\alpha_1$ | ... | $\alpha_n$, where $\alpha_1, ..., \alpha_n$ 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 → $X_1\ X_2 ... X_n$,
   - Replace the rule with A → $A_1\ X_1...X_n$, and add $A_1 \to X_1\ X_2$
   - Repeat the process until all new rules are binary

# CKY demonstration
an ambiguous example

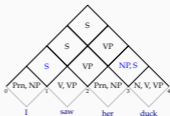| S | → NP VP |
| NP | → Prn N |
| VP | → V NP |
| VP | → V S |
| N | → duck |
| VP | → duck \| saw |
| V | → duck \| saw |
| Prn | → I \| she \| her |
| NP | → I \| she \| her |

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

0   I   1   saw   2   her   3   duck   4

S → NP VP

VP → V NP

NP → Prn N
S → NP VP

S → NP VP

VP → V NP
VP → V S

?

# CKY demonstration
*an ambiguous example*

$$S \quad \to \quad NP\ VP$$



| |
|---|
| 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
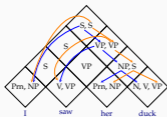*our chart is a 2D array – this is more convenient for programming*



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
  - 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 $O(n^2)$ recognition complexity
- For parsing we need to keep track of backlinks
- CKY can efficiently store all possible parses in a chart
- Enumerating all possible parses have exponential complexity (worst case)
- Suggested reading: Jurafsky and Martin (2009, draft 3rd ed, section 13.2)

Next:
- Top-down chart parsing: Earley algorithm
- Suggested reading:
  - Jurafsky and Martin (2009, section 13.2.4)
  - Grune and Jacobs (2007, section 7.2)

---

# Acknowledgments, references, additional reading material

Grune, Dick and Ceriel 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.

Jurafsky, Daniel and James H. Martin (2009). *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition.* second edition. Pearson Prentice Hall. ISBN: 978-0-13-504196-3.