# CFGs & Parsing

<u>Recall</u>: We showed that CFLs = languages accepted by PDAs

We used Greibach Normal form, where every rule is of the form
$$A \longrightarrow c B_1 \cdots B_k, \text{ where}$$
$$k \geqslant 0, \quad c \in T, \quad A, B_1, \ldots, B_k \in NT.$$

<u>Today</u>: More about normal forms, and why they are useful.

Normal forms give us a uniform "shape" for a structure.

Main use of CFGs is in parsing.

Given a string $\omega$ and a CFG $G$, does $G$ generate $\omega$?

Constructing + operating over PDA is expensive.

Can we somehow directly operate over the grammar?

Normal forms help here!

Greibach Normal form not so much, but

There is another, Chomsky Normal form, where every rule is either of the form $A \rightarrow c$ or $A \rightarrow BC$, where $c \in T$, $A, B, C \in NT$

Thm: For every CFG $G$, there is a CFG $G_1$ in Chomsky Normal form, and a CFG $G_2$ in Greibach Normal form s.t.
$$\mathcal{L}(G_1) = \mathcal{L}(G_2) = \mathcal{L}(G) \setminus \{\varepsilon\}.$$

For any CFL $L$, $L$ can be generated by a grammar where every rule is of the form
$$A \rightarrow c, \text{ or } A \rightarrow BC, \text{ or } (\text{if } \varepsilon \in L) \ S \rightarrow \varepsilon.$$

# Parsing algorithm: [Cocke, Younger, Kasami : 1961]

Main idea: Determine, for each substring $x$ of $\omega$,

the set of all non-terminals that generate $x$.

Needs to do this systematically: ⟵ grammar in Chomsky Normal Form

First check all substrings of length 1, ($A \longrightarrow c$ ?)

then all substrings of length 2 or more

Doing this repeatedly?
Use dynamic programming!

Consider every possible partitioning into two parts and match against $A \longrightarrow BC$, where
B generates the left part, and
C generates the right part.

$$S \rightarrow AB \mid BA \mid SS \mid Ac \mid BD$$

$$A \rightarrow a \qquad B \rightarrow b \qquad C \rightarrow SB \qquad D \rightarrow SA$$

Consider a string $\omega = {}_{|}a_{|}b_{|}b_{|}b_{|}a_{|}a_{|} \quad \rightarrow n = 6 = |\omega|$

$\quad\quad\quad\quad\quad\quad\quad\quad 0 \; 1 \; 2 \; 3 \; 4 \; 5 \; 6$

$\omega_{ij}$ : substring of $\omega$ between markers $i$ and $j$

Build a table with an entry for every $(i,j)$, where $0 \le i < j \le n$.

|   | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 1 |   |   |   |   |   |   |
| 2 |   |   |   |   |   |   |
| 3 |   |   |   |   |   |   |
| 4 |   |   |   |   |   |   |
| 5 |   |   |   |   |   |   |
| 6 |   |   |   |   |   |   |

(columns labeled $i$; rows labeled $j$)

Fill $T(i,j)$ with the non-terminals that generate $\omega_{ij}$

$$S \to AB \mid BA \mid SS \mid AC \mid BD$$

$$A \to a \qquad B \to b$$

$$C \to SB \qquad D \to SA$$

$$\omega = {}_{|}a_{|}b_{|}b_{|}b_{|}a_{|}a_{|}$$
$$\phantom{\omega =}\; 0 \; 1 \; 2 \; 3 \; 4 \; 5 \; 6$$

$\omega_{01} = a$  So $T(0,1) = \{A\}$. Similarly, $\omega_{12} = b$, so $T(1,2) = \{B\}$.

If there are multiple non-terminals which yield $\omega_{ij}$, write them all in $T(i,j)$.

Now look at substrings of length 2.

$$S \rightarrow AB \mid BA \mid SS \mid AC \mid BD$$
$$A \rightarrow a \qquad B \rightarrow b$$
$$C \rightarrow SB \qquad D \rightarrow SA$$

$$\omega = {}_|a_|b_|b_|b_|a_|a_|$$
$$\quad\ 0\ \ 1\ \ 2\ \ 3\ \ 4\ \ 5\ \ 6$$

$\omega_{01}$  $\omega_{12}$

$\omega_{02} = ab$. Break this into two substrings of length 1, $a$, and $b$. Look for all combinations of non-terminals which can yield these substrings, and look for a non-terminal which goes to this pair. $T(0,1) = A$, $T(1,2) = B$, and $S \rightarrow AB$, so $T(0,2) = \{S\}$. Fill in the diagonal below the top one this way.

$$S \rightarrow AB \mid BA \mid SS \mid AC \mid BD$$
$$A \rightarrow a \qquad B \rightarrow b$$
$$C \rightarrow SB \qquad D \rightarrow SA$$

$$\omega = \mid a \mid b \mid b \mid b \mid a \mid a \mid$$
$$\phantom{\omega=\mid}0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6$$

$\omega_{03} = abb$. There are two ways to break this string:

$\omega_{01}$ — a and bb, or ab and b.

$\omega_{01}$, $\omega_{13}$, $\omega_{02}$, $\omega_{23}$

$T(0,1) = A$, $T(1,3) = \emptyset$

$T(0,2) = S$, $T(2,3) = B$

$C \rightarrow SB$, so $T(0,3) = C$

Do this for the diagonal.

$$S \rightarrow AB \mid BA \mid SS \mid Ac \mid BD$$
$$A \rightarrow a \qquad B \rightarrow b$$
$$C \rightarrow SB \qquad D \rightarrow SA$$

$$\omega = {}_|a_|b_|b_|b_|a_|a_|$$
$$\phantom{\omega = } 0\ 1\ 2\ 3\ 4\ 5\ 6$$

| $j \backslash i$ | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 1 | A | | | | | |
| 2 | S | B | | | | |
| 3 | C | ∅ | B | | | |
| 4 | ∅ | ∅ | ∅ | B | | |
| 5 | | ∅ | ∅ | S | A | |
| 6 | | | S | D | ∅ | A |

$\omega_{04} = abbb \qquad a, bbb \quad$ or $\quad ab, bb \quad$ or $\quad abb, b$

None of these has any possible generating non-terminals.

$\omega_{26} = bbaa \qquad$ Only possibility: $b, baa$

$\phantom{\omega_{26} = bbaa \qquad \text{Only possibility: }} B \qquad D \qquad S \rightarrow BD.$ So $T(2,6) = S.$

$$S \rightarrow AB \mid BA \mid SS \mid AC \mid BD$$
$$A \rightarrow a \qquad B \rightarrow b$$
$$C \rightarrow SB \qquad D \rightarrow SA$$

$$\omega = |a|b|b|b|a|a|$$
$$0\ 1\ 2\ 3\ 4\ 5\ 6$$

CYK table (rows $j$, columns $i$):

| $j \backslash i$ | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 1 | A | | | | | |
| 2 | S | B | | | | |
| 3 | C | $\emptyset$ | B | | | |
| 4 | $\emptyset$ | $\emptyset$ | $\emptyset$ | B | | |
| 5 | $\emptyset$ | $\emptyset$ | $\emptyset$ | S | A | |
| 6 | S | $\emptyset$ | S | D | $\emptyset$ | A |

$$\omega_{05} = abbba \qquad\qquad \omega_{06} = \omega : ab, bbaa$$
$$\omega_{16} = bbbaa \qquad\qquad\qquad S' \quad\searrow S \qquad S \rightarrow SS$$

Since $S \in T(0,6)$, this string is generated by the grammar.

```
for i := 0 to n − 1 do                              /* first do substrings of length 1 */
  begin
    T_{i,i+1} := ∅;                                 /* initially assign the empty set */
    for A → a a production of G do
      if a = x_{i,i+1} then T_{i,i+1} := T_{i,i+1} ∪ {A}
  end;
for m := 2 to n do                                  /* for each length m ≥ 2 */
  for i := 0 to n − m do                            /* for each substring of length m */
    begin
      T_{i,i+m} := ∅;                               /* initially assign the empty set */
      for j := i + 1 to i + m − 1 do    /* for all breaks of the string */
        for A → BC a production of G do
          if B ∈ T_{i,j} ∧ C ∈ T_{j,i+m}
            then T_{i,i+m} := T_{i,i+m} ∪ {A}
    end;
```
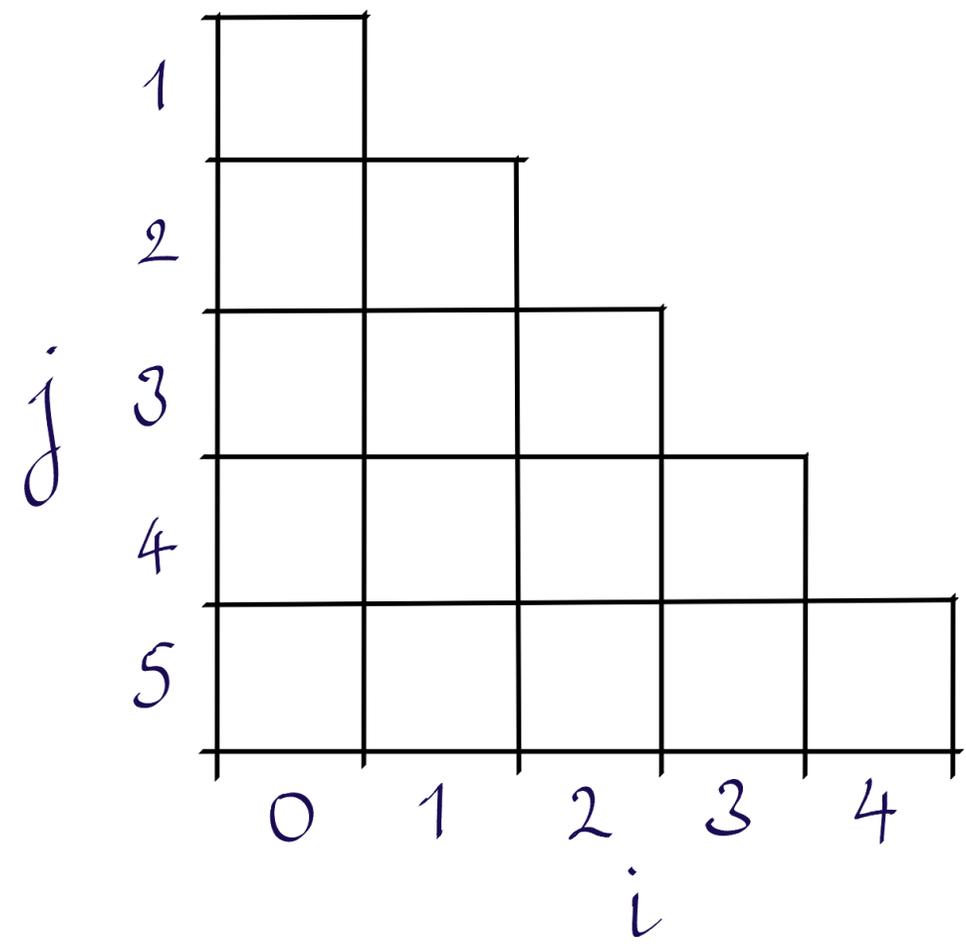
If $S \to \varepsilon$ in $G$, it is used only if $\varepsilon$ is provided as input; NEVER ELSE!

$$S \rightarrow AB \mid BA \mid SS \mid AC \mid BD$$

$$A \rightarrow a \qquad B \rightarrow b$$

$$C \rightarrow SB \qquad D \rightarrow SA$$

$$\omega = {}_{|}a_{|}b_{|}a_{|}b_{|}a_{|}$$
$$\quad\;\, 0 \;\; 1 \;\; 2 \;\; 3 \;\; 4 \; 5$$