# Nondeterministic Top-Down Parsing

# Parsing

- We have seen the relationship between CFGs and syntax trees.
- Now, we have to find an efficient algorithm to obtain the syntax tree for a given input $w \in \Sigma^*$ and the grammar $G$
  - We have to handle ambiguity in the grammar
  - We have to decide which rule to apply for $\alpha \Rightarrow \beta$ if $\alpha$ contains multiple non-terminal symbols
- Like for REs, we will use automata to analyze and understand the parsing process

# Nondeterministic Top-Down Automaton

- The nondetermistic Top-Down parsing automaton $NTA(G)$ of the CFG $G = <\Sigma, N, P, S>$ is defined by
  - Input alphabet $\Sigma$
  - Pushdown alphabet $X = N \cup \Sigma$
  - Output alphabet $U =$ the rule numbers 1,2,3,...
  - States $\Sigma^* \times X^* \times U^*$
  - Two types of transitions for $w \in \Sigma^*, \alpha \in X, z \in U$:
    - **Expansion** of non-terminal symbol $A$ using rule $A \to \beta$ with number $i$:
      $$(w, A\alpha, z) \to (w, \beta\alpha, z\,i)$$
    - **Matching** of terminal symbol $a \in \Sigma$:
      $$(aw, a\alpha, z) \to (w, \alpha, z)$$
  - Initial state $(w, S, \varepsilon)$ for $w \in \Sigma^*$
  - Final state $(\varepsilon, \varepsilon, u)$ where $u \in U^*$

# Example

$$E \rightarrow E + T \mid T \qquad\qquad \text{(1, 2)}$$
$$T \rightarrow T * F \mid F \qquad\qquad \text{(3, 4)}$$
$$F \rightarrow (E) \mid a \mid b \qquad\quad \text{(5, 6, 7)}$$

- Running the NTA on $w = (a) * b$ gives the leftmost analysis of $w$:
  - Initial state $\qquad\qquad\qquad ((a) * b, \qquad E, \ \varepsilon)$
  - Rule 2 to expand $E \qquad\quad ((a) * b, \qquad T, \ 2)$
  - Rule 3 to expand $T \qquad\quad ((a) * b, \quad T * F, \ 23)$
  - Rule 4 to expand $T \qquad\quad ((a) * b, \quad F * F, \ 234)$
  - Rule 5 to expand $E \qquad\quad ((a) * b, \ (E) * F, \ 2345)$
  - Match the "(" $\qquad\qquad ( \ a) * b, \quad E) * F, \ 2345)$
  - Rule 2 to expand $E \qquad\quad ( \ a) * b, \quad T) * F, \ 23452)$
  - Rule 4 to expand $T \qquad\quad ( \ a) * b, \quad F) * F, \ 234524)$
  - Rule 6 to expand $F \qquad\quad ( \ a) * b, \quad a) * F, \ 2345246)$
  - Match "a" $\qquad\qquad\quad ( \quad ) * b, \quad ) * F, \ 2345246)$
  - Match ")" $\qquad\qquad\quad ( \qquad * b, \qquad * F, \ 2345246)$
  - Match "*" $\qquad\qquad\quad ( \qquad b, \qquad F, \ 2345246)$
  - Rule 7 to expand $F \qquad\quad ( \qquad b, \qquad b, \ 23452467)$
  - Match "b" $\qquad\qquad\quad ( \qquad \varepsilon, \qquad \varepsilon, \ 23452467)$

# Nondeterminism

- In the example on the previous slide, we "magically" knew which rule to apply, even in complicated situations

- For a CFG like
  $$S \rightarrow A \mid B \qquad \text{(1, 2)}$$
  $$A \rightarrow b \qquad \text{(3)}$$
  $$B \rightarrow Cd \qquad \text{(4)}$$
  $$C \rightarrow a \qquad \text{(5)}$$

  how should the NTA know that in the initial state $(ad, S, \varepsilon)$ the non-terminal $S$ should be expanded to $B$ and not $A$ ?

- Possible implementation: Backtracking
  1. The parser first tries rule 1.
  2. After rule 3, it sees that it cannot match the input "a"
  3. It goes back to step 1 and tries rule 2.

- Backtracking works, but it requires:
  - Remember the places where a choice was made (here: at rule 1)
  - "Unread" the tokens already read when backtracking