

# Syntax Analysis

# Syntax Analysis

- **Syntax analysis** checks the correctness of a program according to the grammar:
  - gets scanner generated stream of tokens as an input;
  - outputs a syntax-tree corresponding to the program;
  - in the presence of syntactic errors, locates them:
  - ... reports possible causes of the errors;
  - ... tries to recover and continue the analysis (in order to discover more errors).
- Syntax analysis is called **parsing** and the corresponding analyzer **parser**.

# Grammars

- Syntax is usually described by context-free grammars.
- **Grammar** is a quadruple  $G = \langle N, T, P, S \rangle$ , where
  - $N$  is a finite alphabet of **non-terminal symbols**;
  - $T$  is a finite alphabet of **terminal symbols**;
  - $N \cap T = \emptyset$  and  $V = N \cup T$ ;
  - $P \subset \{\alpha \rightarrow \beta \mid \alpha \in V^+, \beta \in V^*\}$  is a finite set of **production rules**;
  - $S \in N$  is a **start symbol**.
- Grammar is **context-free** if production rules are in the form  $A \rightarrow \alpha$ , where  $A \in N$  and  $\alpha \in V^*$ .

## Grammars

- A sequence  $w \in V^*$  is called a **sentential form**.
- The sentential form  $v \in V^*$  is **directly derivable** from the sentential form  $u \in V^*$  (notation  $u \Rightarrow v$ ), if there are  $w_1, w_2, \alpha, \beta \in V^*$  such, that  $u = w_1\alpha w_2$ ,  $v = w_1\beta w_2$  and  $\alpha \rightarrow \beta \in P$ .
- Reflexive transitive closure of the relation  $\Rightarrow$  is called **derivation** (notation  $\Rightarrow^*$ ).
- The grammar  $G = \langle N, T, P, S \rangle$  generates a **language**

$$L(G) = \{w \in T^* \mid S \Rightarrow^* w\}$$

- Grammars  $G_1$  and  $G_2$  are **equivalent** if  $L(G_1) = L(G_2)$ .

# Grammars

## Chomsky hierarchy:

	Productions	Languages	Automata
$L_0$	$\alpha \rightarrow \beta$	Semi-Thue systems	Turing machines
$L_1$	$\alpha A \beta \rightarrow \alpha \gamma \beta$	Context-dependent	Bounded TM-s
$L_2$	$A \rightarrow \alpha$	Context-free	Push-down automata
$L_3$	$A \rightarrow w, A \rightarrow wB$	Regular	Finite automata
$(L_4)$	$A \rightarrow w$	Finite	FA without cycles

where  $A, B \in N$ ,  $\alpha, \beta, \gamma \in V^*$  ja  $w \in T^*$ .

**Lemma:** Chomsky hierarchy is strict; ie.:

$$(L_4) \subset L_3 \subset L_2 \subset L_1 \subset L_0$$

## Context-free Grammars

- From now on we consider only context-free grammars.
- Production rules of context-free grammars are usually described using **Backus-Naur Form** (BNF).
- Example: let  $N = \{\text{Exp}\}$  and  $T = \{+, *, (, ), id\}$ , then

$$\begin{array}{l} \text{Exp} \rightarrow \text{Exp} + \text{Exp} \\ \quad | \quad \text{Exp} * \text{Exp} \\ \quad | \quad ( \text{Exp} ) \\ \quad | \quad id \end{array}$$

describes the set of production rules

$$P = \{ \text{Exp} \rightarrow \text{Exp} + \text{Exp}, \text{Exp} \rightarrow ( \text{Exp} ), \\ \text{Exp} \rightarrow \text{Exp} * \text{Exp}, \text{Exp} \rightarrow id \}.$$

## Context-free Grammars

- Non-terminal  $A$  is **productive** if there exists  $w \in T^*$  such that  $A \Rightarrow^* w$ .
- Non-terminal  $A$  is **reachable** if there exist sentential forms  $u, v \in V^*$  such that  $S \Rightarrow^* uAv$ .
- CF-grammar  $G = \langle N, T, P, S \rangle$  is **reduced** if every non-terminal is productive and reachable.
- **Lemma:** Every CF-grammar can be transformed into an equivalent reduced CF-grammar.

## Context-free Grammars

- A sentential form may have several derivations.
- Canonical derivations:
  - **left-derivation** – on every derivation step the leftmost non-terminal is replaced;
  - **right-derivation** – on every derivation step the rightmost non-terminal is replaced.
- Example:

$$\begin{aligned} \text{Exp} &\Rightarrow_{lm} \text{Exp} + \text{Exp} \\ &\Rightarrow_{lm} id + \text{Exp} \\ &\Rightarrow_{lm} id + \text{Exp} * \text{Exp} \\ &\Rightarrow_{lm} id + id * \text{Exp} \\ &\Rightarrow_{lm} id + id * id \end{aligned}$$
$$\begin{aligned} \text{Exp} &\Rightarrow_{rm} \text{Exp} + \text{Exp} \\ &\Rightarrow_{rm} \text{Exp} + \text{Exp} * \text{Exp} \\ &\Rightarrow_{rm} \text{Exp} + \text{Exp} * id \\ &\Rightarrow_{rm} \text{Exp} + id * id \\ &\Rightarrow_{rm} id + id * id \end{aligned}$$

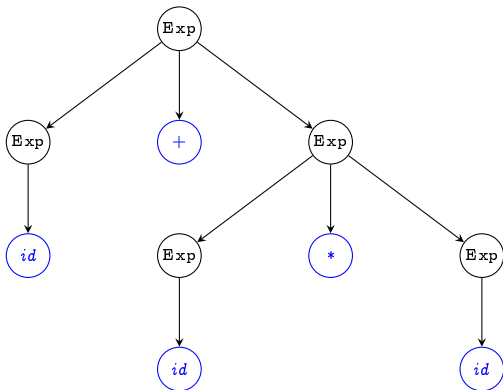


## Context-free Grammars

- Every derivation determines a unique **syntax-tree** (or **parse-tree**) – a tree with ordered nodes, where:
  - the root is labelled by the start symbol  $S$ ;
  - intermediate nodes are labelled by non-terminals;
  - leaves are labelled by terminals or the empty symbol  $\epsilon$ ;
  - when intermediate node is labelled by the non-terminal  $A$  and roots of its subtrees (from left to right)  $t_1, \dots, t_n$  are labelled by  $A_1, \dots, A_n$ , then  $A \rightarrow A_1 \dots A_n \in P$ .
- Labels of leaves (read from left to right) form the derived sentential form.
- Syntax-tree uniquely determines which production rules were used, but not the order of their application.

## Context-free Grammars

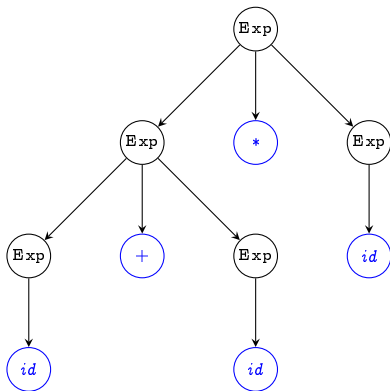
Example: previously given left- and right-derivations correspond to the same syntax-tree



# Context-free Grammars

**NB!** A sentence may have several syntax-trees!

$\text{Exp} \Rightarrow_{lm} \text{Exp} * \text{Exp}$   
 $\Rightarrow_{lm} \text{Exp} + \text{Exp} * \text{Exp}$   
 $\Rightarrow_{lm} id + \text{Exp} * \text{Exp}$   
 $\Rightarrow_{lm} id + id * \text{Exp}$   
 $\Rightarrow_{lm} id + id * id$



## Context-free Grammars

- CF-grammar is **ambiguous** if for the same sentence there are several syntax-trees.
- For every syntax-tree, there is exactly one left- and right-derivation; thus:
  - non-ambiguous sentence has exactly one left- and one right-derivation;
  - ambiguous sentence has at least two left- and right-derivations..
- Different syntax-trees of a sentence usually correspond to different semantic interpretations of the sentence.
- An ambiguous grammar can sometimes (but not always) be transformed to an equivalent non-ambiguous one.

## Context-free Grammars

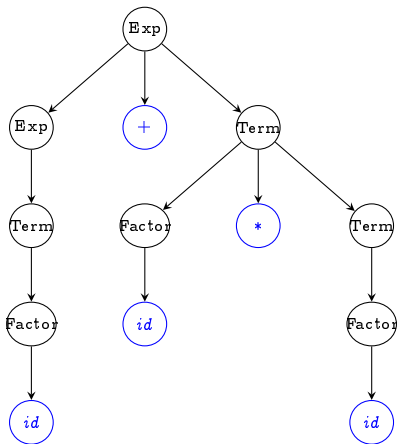
- Elimination of ambiguity – binary operators:
  - every priority level introduces a new non-terminal;
  - left-associative operators use left-recursion;  
right-associative operators right-recursion;
  - rules corresponding to operators of higher priorities are placed "deeper".
- Example:

$$\begin{array}{l} \text{Exp} \quad \rightarrow \quad \text{Exp} + \text{Term} \\ \quad \quad \quad | \quad \text{Term} \\ \text{Term} \quad \rightarrow \quad \text{Factor} * \text{Term} \\ \quad \quad \quad | \quad \text{Factor} \\ \text{Factor} \quad \rightarrow \quad ( \text{Exp} ) \\ \quad \quad \quad | \quad \textit{id} \end{array}$$

# Context-free Grammars

Example:

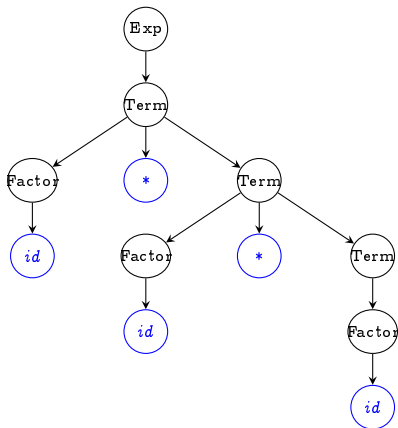
$\text{Exp} \Rightarrow_{lm} \text{Exp} + \text{Term}$   
 $\Rightarrow_{lm} \text{Term} + \text{Term}$   
 $\Rightarrow_{lm} \text{Factor} + \text{Term}$   
 $\Rightarrow_{lm} id + \text{Term}$   
 $\Rightarrow_{lm} id + \text{Factor} * \text{Term}$   
 $\Rightarrow_{lm} id + id * \text{Term}$   
 $\Rightarrow_{lm} id + id * \text{Factor}$   
 $\Rightarrow_{lm} id + id * id$



# Context-free Grammars

Example:

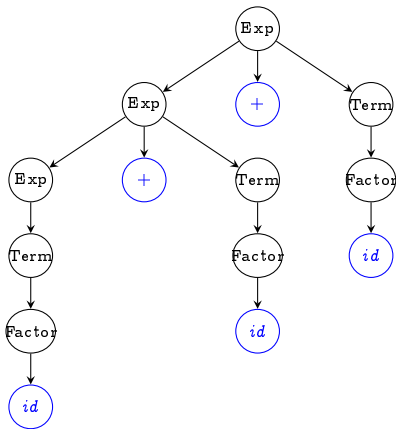
$\text{Exp} \Rightarrow_{lm} \text{Term}$   
 $\Rightarrow_{lm} \text{Factor} * \text{Term}$   
 $\Rightarrow_{lm} id * \text{Term}$   
 $\Rightarrow_{lm} id * \text{Factor} * \text{Term}$   
 $\Rightarrow_{lm} id * id * \text{Term}$   
 $\Rightarrow_{lm} id * id * \text{Factor}$   
 $\Rightarrow_{lm} id * id * id$



# Context-free Grammars

Example:

$\text{Exp} \Rightarrow_{lm} \text{Exp} + \text{Term}$   
 $\Rightarrow_{lm} \text{Exp} + \text{Term} + \text{Term}$   
 $\Rightarrow_{lm} \text{Term} + \text{Term} + \text{Term}$   
 $\Rightarrow_{lm} \text{Factor} + \text{Term} + \text{Term}$   
 $\Rightarrow_{lm} id + \text{Term} + \text{Term}$   
 $\Rightarrow_{lm} id + \text{Factor} + \text{Term}$   
 $\Rightarrow_{lm} id + id + \text{Term}$   
 $\Rightarrow_{lm} id + id + \text{Factor}$   
 $\Rightarrow_{lm} id + id + id$





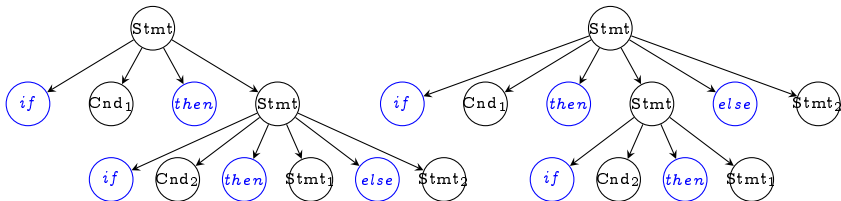
# Context-free Grammars

- Elimination of ambiguity – condition statements:

$$\begin{array}{l} \text{Stmt} \rightarrow \textit{if} \text{ Cnd } \textit{then} \text{ Stmt} \\ \quad \quad | \textit{if} \text{ Cnd } \textit{then} \text{ Stmt } \textit{else} \text{ Stmt} \\ \quad \quad | \text{ Other} \end{array}$$

- The following sentence has two different syntax-trees:

*if* Cnd<sub>1</sub> *then if* Cnd<sub>2</sub> *then* Stmt<sub>1</sub> *else* Stmt<sub>2</sub>

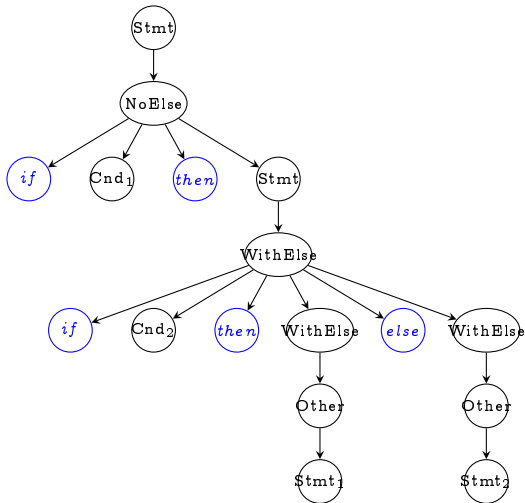


## Context-free Grammars

Usually, the first one is considered to be the correct one; ie. *else* belongs to the innermost conditional sentence:

Stmt	→	WithElse
		NoElse
WithElse	→	<i>if</i> Cnd <i>then</i> WithElse <i>else</i> WithElse
		Other
NoElse	→	<i>if</i> Cnd <i>then</i> Stmt
		<i>if</i> Cnd <i>then</i> WithElse <i>else</i> NoElse

# Context-free Grammars



# Parsing Techniques

## Top-down parsing:

- starts constructing the syntax-tree from the root downward towards leaves;
- on every step selects a production rule and tries to match it with the input string;
- if the rule doesn't match the process backtracks;
- results to the leftmost derivation.

## Bottom-up parsing:

- starts constructing the syntax-tree from leaves working up toward the root;
- applies suitable rules from right to left until reaches the start symbol;
- results to the rightmost derivation.

# Top-down Parsing

General algorithm of top-down parsing:

- construct a root node, label it with the start symbol, and continue construction of the tree towards leaves from left to right;
- if the node under consideration is a non-terminal  $A$ , then choose a rule in the form of  $A \rightarrow \alpha$ , construct nodes corresponding to its RHS, and continue with its leftmost subnode.
- if the node is a terminal which doesn't match with the input symbol, then backtrack to the choice of the production rule which introduced this terminal, and continue from there by choosing another production rule;
- if the node is terminal matching to the input symbol, then continue with the leftmost unexpanded node.

# Top-down Parsing

Example:

$E \rightarrow E + T$   
|  $E - T$   
|  $T$   
 $T \rightarrow T * F$   
|  $T / F$   
|  $F$   
 $F \rightarrow ( E )$   
|  $id$   
|  $num$

$\bullet id - num * id$   
E.1  $\bullet id - num * id$   
E.3  $\bullet id - num * id$   
T.3  $\bullet id - num * id$   
F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
T.1  $id - \bullet num * id$   
T.3  $id - \bullet num * id$   
F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

# Top-down Parsing

Example:

$E \rightarrow E + T$   
|  $E - T$   
|  $T$   
 $T \rightarrow T * F$   
|  $T / F$   
|  $F$   
 $F \rightarrow ( E )$   
|  $id$   
|  $num$

$\bullet id - num * id$   
E.1  $\bullet id - num * id$   
E.3  $\bullet id - num * id$   
T.3  $\bullet id - num * id$   
F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
T.1  $id - \bullet num * id$   
T.3  $id - \bullet num * id$   
F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

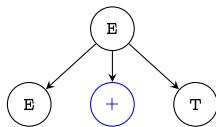
E

# Top-down Parsing

Example:

$E \rightarrow E + T$   
|  $E - T$   
|  $T$   
 $T \rightarrow T * F$   
|  $T / F$   
|  $F$   
 $F \rightarrow ( E )$   
|  $id$   
|  $num$

$\bullet id - num * id$   
E.1  $\bullet id - num * id$   
E.3  $\bullet id - num * id$   
T.3  $\bullet id - num * id$   
F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
T.1  $id - \bullet num * id$   
T.3  $id - \bullet num * id$   
F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
F.2  $id - num * \bullet id$   
 $id - num * id \bullet$



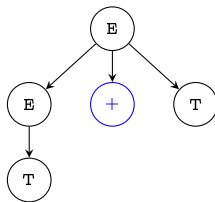


# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.1  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

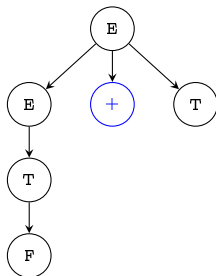


# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.1  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

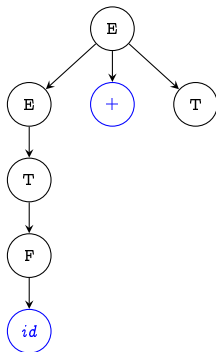


# Top-down Parsing

Example:

$E \rightarrow E + T$   
|  $E - T$   
|  $T$   
 $T \rightarrow T * F$   
|  $T / F$   
|  $F$   
 $F \rightarrow ( E )$   
|  $id$   
|  $num$

•  $id - num * id$   
E.1 •  $id - num * id$   
E.3 •  $id - num * id$   
T.3 •  $id - num * id$   
F.2 •  $id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
T.1  $id - \bullet num * id$   
T.3  $id - \bullet num * id$   
F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

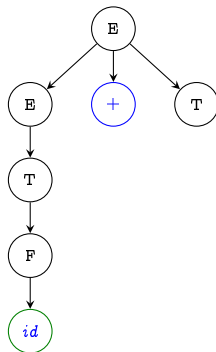


# Top-down Parsing

Example:

$E \rightarrow E + T$   
|  $E - T$   
|  $T$   
 $T \rightarrow T * F$   
|  $T / F$   
|  $F$   
 $F \rightarrow ( E )$   
|  $id$   
|  $num$

•  $id - num * id$   
E.1 •  $id - num * id$   
E.3 •  $id - num * id$   
T.3 •  $id - num * id$   
F.2 •  $id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
T.1  $id - \bullet num * id$   
T.3  $id - \bullet num * id$   
F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

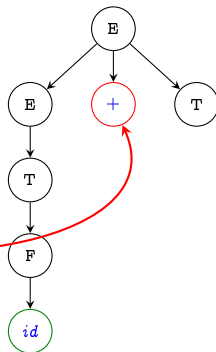


# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.1  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$



# Top-down Parsing

Example:

$E \rightarrow E + T$   
|  $E - T$   
|  $T$   
 $T \rightarrow T * F$   
|  $T / F$   
|  $F$   
 $F \rightarrow ( E )$   
|  $id$   
|  $num$

$\bullet id - num * id$   
E.1  $\bullet id - num * id$   
E.3  $\bullet id - num * id$   
T.3  $\bullet id - num * id$   
F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
T.1  $id - \bullet num * id$   
T.3  $id - \bullet num * id$   
F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

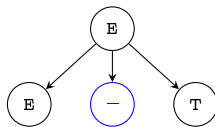
E

# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.2  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

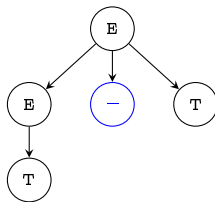


# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.2  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$



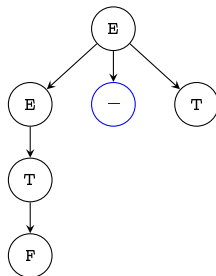


# Top-down Parsing

Example:

$E \rightarrow E + T$   
|  $E - T$   
|  $T$   
 $T \rightarrow T * F$   
|  $T / F$   
|  $F$   
 $F \rightarrow ( E )$   
|  $id$   
|  $num$

•  $id - num * id$   
E.2 •  $id - num * id$   
E.3 •  $id - num * id$   
T.3 •  $id - num * id$   
F.2 •  $id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
T.1  $id - \bullet num * id$   
T.3  $id - \bullet num * id$   
F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

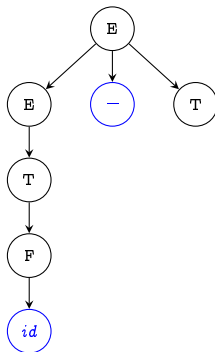


# Top-down Parsing

Example:

$E \rightarrow E + T$   
|  $E - T$   
|  $T$   
 $T \rightarrow T * F$   
|  $T / F$   
|  $F$   
 $F \rightarrow ( E )$   
|  $id$   
|  $num$

•  $id - num * id$   
E.2 •  $id - num * id$   
E.3 •  $id - num * id$   
T.3 •  $id - num * id$   
F.2 •  $id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
T.1  $id - \bullet num * id$   
T.3  $id - \bullet num * id$   
F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

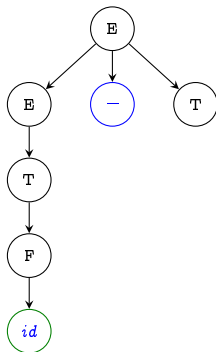


# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.2  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

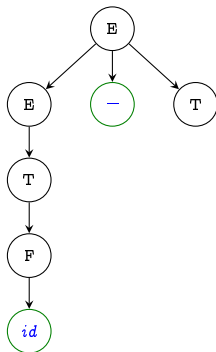


# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.2  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

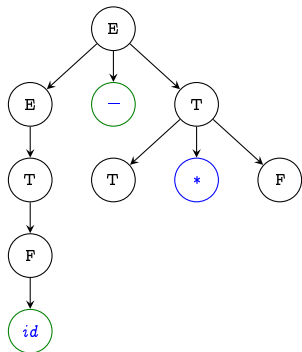


# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.2  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

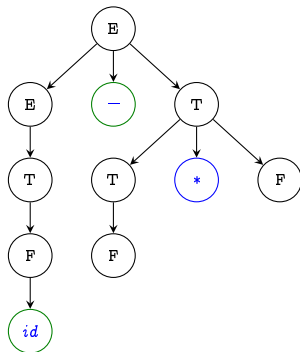


# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.2  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

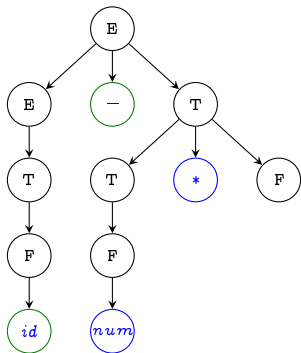


# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.2  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

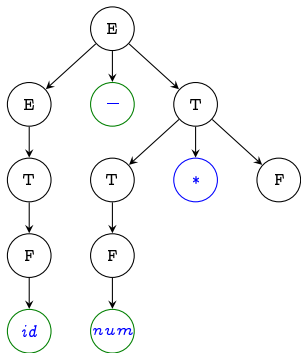


# Top-down Parsing

Example:

$E \rightarrow E + T$   
|  $E - T$   
|  $T$   
 $T \rightarrow T * F$   
|  $T / F$   
|  $F$   
 $F \rightarrow ( E )$   
|  $id$   
|  $num$

•  $id - num * id$   
E.2 •  $id - num * id$   
E.3 •  $id - num * id$   
T.3 •  $id - num * id$   
F.2 •  $id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
T.1  $id - \bullet num * id$   
T.3  $id - \bullet num * id$   
F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
F.2  $id - num * \bullet id$   
 $id - num * id \bullet$



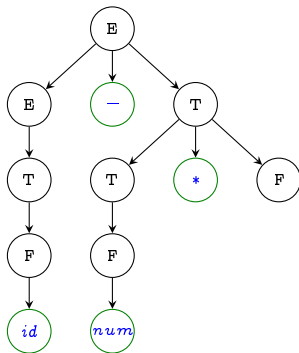


# Top-down Parsing

Example:

$E \rightarrow E + T$   
|  $E - T$   
|  $T$   
 $T \rightarrow T * F$   
|  $T / F$   
|  $F$   
 $F \rightarrow ( E )$   
|  $id$   
|  $num$

•  $id - num * id$   
E.2 •  $id - num * id$   
E.3 •  $id - num * id$   
T.3 •  $id - num * id$   
F.2 •  $id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
T.1  $id - \bullet num * id$   
T.3  $id - \bullet num * id$   
F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

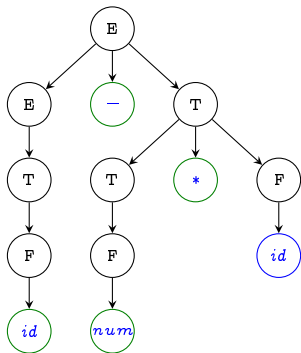


# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.2  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$

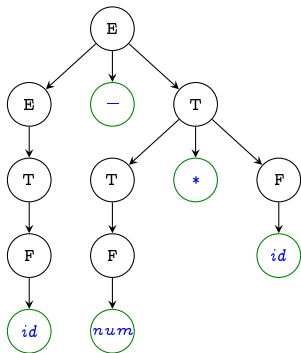


# Top-down Parsing

Example:

$E \rightarrow E + T$   
 $\quad | E - T$   
 $\quad | T$   
 $T \rightarrow T * F$   
 $\quad | T / F$   
 $\quad | F$   
 $F \rightarrow ( E )$   
 $\quad | id$   
 $\quad | num$

$\bullet id - num * id$   
 E.2  $\bullet id - num * id$   
 E.3  $\bullet id - num * id$   
 T.3  $\bullet id - num * id$   
 F.2  $\bullet id - num * id$   
 $id \bullet - num * id$   
 $id - \bullet num * id$   
 T.1  $id - \bullet num * id$   
 T.3  $id - \bullet num * id$   
 F.3  $id - \bullet num * id$   
 $id - num \bullet * id$   
 $id - num * \bullet id$   
 F.2  $id - num * \bullet id$   
 $id - num * id \bullet$



## Top-down Parsing

- Efficiency of parsing strongly depends from the choice of a production rule.
- Choosing a wrong rule causes a later backtracking.
- In the case of the grammar has left-recursive rules, the top-down parsing may **not to terminate**.

- *id - num \* id*

E.2 ● *id - num \* id*

E.2 ● *id - num \* id*

E.2 ● *id - num \* id*

...

# Top-down Parsing

- Efficiency of parsing strongly depends from the choice of a production rule.
- Choosing a wrong rule causes a later backtracking.
- In the case of the grammar has left-recursive rules, the top-down parsing may **not to terminate**.



- *id - num \* id*

E.2 ● *id - num \* id*

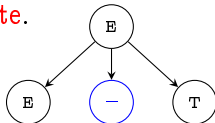
E.2 ● *id - num \* id*

E.2 ● *id - num \* id*

...

# Top-down Parsing

- Efficiency of parsing strongly depends from the choice of a production rule.
- Choosing a wrong rule causes a later backtracking.
- In the case of the grammar has left-recursive rules, the top-down parsing may **not to terminate**.



- *id - num \* id*

E.2 ● *id - num \* id*

E.2 ● *id - num \* id*

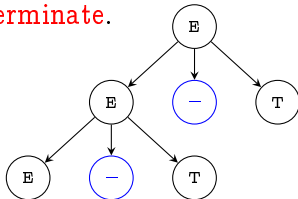
E.2 ● *id - num \* id*

...

# Top-down Parsing

- Efficiency of parsing strongly depends from the choice of a production rule.
- Choosing a wrong rule causes a later backtracking.
- In the case of the grammar has left-recursive rules, the top-down parsing may **not terminate**.

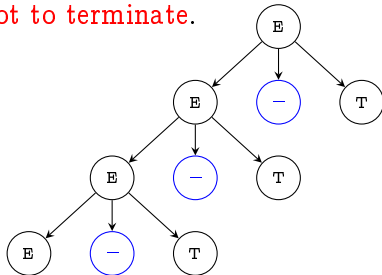
- *id - num \* id*
- E.2 ● *id - num \* id*
- E.2 ● *id - num \* id*
- E.2 ● *id - num \* id*



# Top-down Parsing

- Efficiency of parsing strongly depends from the choice of a production rule.
- Choosing a wrong rule causes a later backtracking.
- In the case of the grammar has left-recursive rules, the top-down parsing may **not terminate**.

- *id - num \* id*
- E.2 ● *id - num \* id*
- E.2 ● *id - num \* id*
- E.2 ● *id - num \* id*



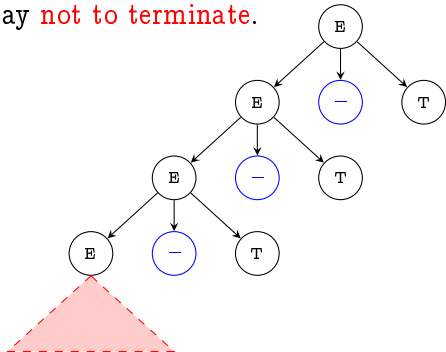


# Top-down Parsing

- Efficiency of parsing strongly depends from the choice of a production rule.
- Choosing a wrong rule causes a later backtracking.
- In the case of the grammar has left-recursive rules, the top-down parsing may **not to terminate**.

- *id - num \* id*
- E.2 ● *id - num \* id*
- E.2 ● *id - num \* id*
- E.2 ● *id - num \* id*

...



## Left-recursion

- A grammar is **left-recursive**, if there is a non-terminal  $A \in N$  such that

$$A \Longrightarrow^+ A\alpha,$$

where  $\alpha \in V^*$ .

- Left-recursion is **direct**, if there is a rule in the form  $A \rightarrow A\alpha$ .
- Otherwise, the left-recursion is **indirect**.

## Left-recursion Elimination

Elimination of the direct left-recursion:

- Introduce a new non-terminal and replace the left-recursion with the right-recursion

$$A \rightarrow A \alpha \mid \beta \quad \longrightarrow \quad \begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \epsilon \end{array}$$

- In general

$$\begin{array}{l} A \rightarrow A \alpha_1 \mid A \alpha_2 \mid \dots \mid \beta_1 \mid \beta_2 \mid \dots \\ \longrightarrow \quad \begin{array}{l} A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \\ A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \epsilon \end{array} \end{array}$$

## Left-recursion Elimination

Example:

$$\begin{array}{l} E \rightarrow E + T \\ \quad | \quad E - T \\ \quad | \quad T \end{array}$$



$$\begin{array}{l} E \rightarrow T E' \\ E' \rightarrow + T E' \\ \quad | \quad - T E' \\ \quad | \quad \epsilon \end{array}$$

$$\begin{array}{l} T \rightarrow T * F \\ \quad | \quad T / F \\ \quad | \quad F \end{array}$$



$$\begin{array}{l} T \rightarrow F T' \\ T' \rightarrow * F T' \\ \quad | \quad / F T' \\ \quad | \quad \epsilon \end{array}$$

- The new grammar generates the same language, but is less intuitive.
- In both grammars the operators are left associative.

## Left-recursion Elimination

Elimination of the indirect left-recursion:

- Example:

$$\begin{aligned}A_1 &\rightarrow A_2 \alpha \mid \beta \\A_2 &\rightarrow A_1 \gamma \mid A_2 \delta\end{aligned}$$

- Transform the indirect left-recursion to the direct one.
- In the right-hand sides of  $A_2$  production rules, replace all occurrences of  $A_1$  with its definition:

$$\begin{aligned}A_1 &\rightarrow A_2 \alpha \mid \beta \\A_2 &\rightarrow A_2 \alpha \gamma \mid \beta \gamma \mid A_2 \delta\end{aligned}$$

- Eliminate the immediate left-recursion:

$$\begin{aligned}A_1 &\rightarrow A_2 \alpha \mid \beta \\A_2 &\rightarrow \beta \gamma A'_2 \\A'_2 &\rightarrow \alpha \gamma A'_2 \mid \delta A'_2 \mid \epsilon\end{aligned}$$

## Left-recursion Elimination

General algorithm for left-recursion elimination:

Assign some order to non-terminals  $A_1, \dots, A_n$

**for**  $i \leftarrow 1$  **to**  $n$

**for**  $j \leftarrow 1$  **to**  $i - 1$

    Replace productions in the form  $A_i \rightarrow A_j \alpha$

    with productions  $A_i \rightarrow \beta_1 \alpha \mid \beta_2 \alpha \mid \dots \mid \beta_k \alpha$ ,

    where  $A_j \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_k$  are all  $A_j$  productions

  Eliminate the immediate left-recursion from productions of  
  the non-terminal  $A_i$

**NB!** Assumes that the original grammar doesn't have neither  $\epsilon$ -productions nor cycles (ie.  $A_i \Longrightarrow^+ A_i$ ).

# Left-recursion Elimination

Example:

$$\begin{aligned} A &\rightarrow C \alpha \mid A \alpha \\ B &\rightarrow A \beta \mid \gamma \\ C &\rightarrow B \delta \mid \varphi \end{aligned}$$



$$\begin{aligned} A &\rightarrow C \alpha \mid A' \\ B &\rightarrow A \beta \mid \gamma \\ C &\rightarrow B \delta \mid \varphi \\ A' &\rightarrow \alpha A' \mid \varepsilon \end{aligned}$$



$$\begin{aligned} A &\rightarrow C \alpha \mid A' \\ B &\rightarrow C \alpha \mid A' \beta \mid \gamma \\ C &\rightarrow B \delta \mid \varphi \\ A' &\rightarrow \alpha A' \mid \varepsilon \end{aligned}$$



$$\begin{aligned} A &\rightarrow C \alpha \mid A' \\ B &\rightarrow C \alpha \mid A' \beta \mid \gamma \\ C &\rightarrow C \alpha \mid A' \beta \delta \mid \gamma \delta \mid \varphi \\ A' &\rightarrow \alpha A' \mid \varepsilon \end{aligned}$$



$$\begin{aligned} A &\rightarrow C \alpha \mid A' \\ B &\rightarrow C \alpha \mid A' \beta \mid \gamma \\ C &\rightarrow \alpha A' \beta \delta \mid C' \\ A' &\rightarrow \alpha A' \mid \varepsilon \\ C' &\rightarrow \gamma \delta \mid C' \mid \varphi \mid C' \mid \varepsilon \end{aligned}$$



$$\begin{aligned} A &\rightarrow C \alpha \mid A' \\ C &\rightarrow \alpha A' \beta \delta \mid C' \\ A' &\rightarrow \alpha A' \mid \varepsilon \\ C' &\rightarrow \gamma \delta \mid C' \mid \varphi \mid C' \mid \varepsilon \end{aligned}$$

# Left-recursion Elimination

Example:

$$\begin{array}{l} A \rightarrow C \alpha \mid A \alpha \\ B \rightarrow A \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \end{array} \quad \longrightarrow \quad \begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow A \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$

$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha \mid A' \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array} \quad \longrightarrow \quad \begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha \mid A' \beta \mid \gamma \\ C \rightarrow C \alpha \mid A' \beta \delta \mid \gamma \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$

$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha \mid A' \beta \mid \gamma \\ C \rightarrow \alpha A' \beta \delta \mid C' \\ A' \rightarrow \alpha A' \mid \varepsilon \\ C' \rightarrow \gamma \delta C' \mid \varphi C' \mid \varepsilon \end{array} \quad \longrightarrow \quad \begin{array}{l} A \rightarrow C \alpha \mid A' \\ C \rightarrow \alpha A' \beta \delta \mid C' \\ A' \rightarrow \alpha A' \mid \varepsilon \\ C' \rightarrow \gamma \delta C' \mid \varphi C' \mid \varepsilon \end{array}$$



# Left-recursion Elimination

Example:

$$\begin{array}{l} A \rightarrow C \alpha \mid A \alpha \\ B \rightarrow A \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \end{array} \quad \longrightarrow \quad \begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow A \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$

$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha A' \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array} \quad \longrightarrow \quad \begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha A' \beta \mid \gamma \\ C \rightarrow C \alpha A' \beta \delta \mid \gamma \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$

$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha A' \beta \mid \gamma \\ C \rightarrow \alpha A' \beta \delta C' \\ A' \rightarrow \alpha A' \mid \varepsilon \\ C' \rightarrow \gamma \delta C' \mid \varphi C' \mid \varepsilon \end{array} \quad \longrightarrow \quad \begin{array}{l} A \rightarrow C \alpha \mid A' \\ C \rightarrow \alpha A' \beta \delta C' \\ A' \rightarrow \alpha A' \mid \varepsilon \\ C' \rightarrow \gamma \delta C' \mid \varphi C' \mid \varepsilon \end{array}$$

# Left-recursion Elimination

Example:

$$\begin{array}{l} A \rightarrow C \alpha \mid A \alpha \\ B \rightarrow A \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow A \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha \mid A' \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha \mid A' \beta \mid \gamma \\ C \rightarrow C \alpha \mid A' \beta \delta \mid \gamma \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha \mid A' \beta \mid \gamma \\ C \rightarrow \alpha A' \beta \delta \mid C' \\ A' \rightarrow \alpha A' \mid \varepsilon \\ C' \rightarrow \gamma \delta \mid C' \mid \varphi \mid C' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ C \rightarrow \alpha A' \beta \delta \mid C' \\ A' \rightarrow \alpha A' \mid \varepsilon \\ C' \rightarrow \gamma \delta \mid C' \mid \varphi \mid C' \mid \varepsilon \end{array}$$

# Left-recursion Elimination

Example:

$$\begin{array}{l} A \rightarrow C \alpha \mid A \alpha \\ B \rightarrow A \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow A \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha \mid A' \beta \mid \gamma \\ C \rightarrow B \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha \mid A' \beta \mid \gamma \\ C \rightarrow C \alpha \mid A' \beta \delta \mid \gamma \delta \mid \varphi \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ B \rightarrow C \alpha \mid A' \beta \mid \gamma \\ C \rightarrow \alpha A' \beta \delta \mid C' \\ A' \rightarrow \alpha A' \mid \varepsilon \\ C' \rightarrow \gamma \delta C' \mid \varphi C' \mid \varepsilon \end{array}$$



$$\begin{array}{l} A \rightarrow C \alpha \mid A' \\ C \rightarrow \alpha A' \beta \delta \mid C' \\ A' \rightarrow \alpha A' \mid \varepsilon \\ C' \rightarrow \gamma \delta C' \mid \varphi C' \mid \varepsilon \end{array}$$

# Left-recursion Elimination

Example:

$$\begin{array}{l}
 A \rightarrow C \alpha \mid A \alpha \\
 B \rightarrow A \beta \mid \gamma \\
 C \rightarrow B \delta \mid \varphi
 \end{array}
 \longrightarrow
 \begin{array}{l}
 A \rightarrow C \alpha \mid A' \\
 B \rightarrow A \beta \mid \gamma \\
 C \rightarrow B \delta \mid \varphi \\
 A' \rightarrow \alpha A' \mid \varepsilon
 \end{array}$$

$$\begin{array}{l}
 A \rightarrow C \alpha \mid A' \\
 B \rightarrow C \alpha \mid A' \beta \mid \gamma \\
 C \rightarrow B \delta \mid \varphi \\
 A' \rightarrow \alpha A' \mid \varepsilon
 \end{array}
 \longrightarrow
 \begin{array}{l}
 A \rightarrow C \alpha \mid A' \\
 B \rightarrow C \alpha \mid A' \beta \mid \gamma \\
 C \rightarrow C \alpha \mid A' \beta \delta \mid \gamma \delta \mid \varphi \\
 A' \rightarrow \alpha A' \mid \varepsilon
 \end{array}$$

$$\begin{array}{l}
 A \rightarrow C \alpha \mid A' \\
 B \rightarrow C \alpha \mid A' \beta \mid \gamma \\
 C \rightarrow \alpha A' \beta \delta \mid C' \\
 A' \rightarrow \alpha A' \mid \varepsilon \\
 C' \rightarrow \gamma \delta \mid C' \mid \varphi \mid C' \mid \varepsilon
 \end{array}
 \longrightarrow
 \begin{array}{l}
 A \rightarrow C \alpha \mid A' \\
 C \rightarrow \alpha A' \beta \delta \mid C' \\
 A' \rightarrow \alpha A' \mid \varepsilon \\
 C' \rightarrow \gamma \delta \mid C' \mid \varphi \mid C' \mid \varepsilon
 \end{array}$$

## Predictive Parsing

- Choosing a wrong rule causes a later backtracking.
- The "rightness" of the rule can often be decided by lookahead of some number of input symbols.
- In general case, one needs unbounded lookahead.
  - Ex.: parsing algorithms by Cocke-Younger-Kasam or Earley.
- **Predictive parsing** is a top-down parsing where it is always possible to choose a correct rule without backtracking.
  - A grammar must be such, that the next input symbol (or some fixed number of symbols) determines uniquely the correct rule.

## Predictive Parsing

- For every sentential form  $\alpha \in (N \cup T)^*$  define a set:

$$\begin{aligned} \text{first}(\alpha) &= \{a \in T \mid \alpha \Longrightarrow^* a\beta\} \\ &\cup \{\varepsilon \mid \alpha \Longrightarrow^* \varepsilon\} \end{aligned}$$

where  $\beta \in (N \cup T)^*$ .

- For every non-terminal  $A \in N$  define a set:

$$\begin{aligned} \text{follow}(A) &= \{a \in T \mid S \Longrightarrow^* \alpha A a \beta\} \\ &\cup \{\$ \mid S \Longrightarrow^* \alpha A\} \end{aligned}$$

where  $\alpha, \beta \in (N \cup T)^*$  and  $\$$  is a special end of input marker.

# Predictive Parsing

Example:

$$S \rightarrow ABC$$

$$A \rightarrow aA \mid \varepsilon$$

$$B \rightarrow b \mid \varepsilon$$

$$C \rightarrow c \mid d$$

$$\text{first}(C) = \{c, d\}$$

$$\text{first}(B) = \{b, \varepsilon\}$$

$$\text{first}(A) = \{a, \varepsilon\}$$

$$\text{first}(S) = \text{first}(ABC)$$

$$= (\text{first}(A) \setminus \{\varepsilon\})$$

$$\cup (\text{first}(B) \setminus \{\varepsilon\})$$

$$\cup \text{first}(C)$$

$$= \{a, b, c, d\}$$

$$\text{follow}(C) = \{\$\}$$

$$\text{follow}(B) = \text{first}(C)$$

$$= \{c, d\}$$

$$\text{follow}(A) = (\text{first}(B) \setminus \{\varepsilon\})$$

$$\cup \text{first}(C)$$

$$= \{b, c, d\}$$

## Predictive Parsing

- If production rules  $A \rightarrow \alpha$  and  $A \rightarrow \beta$  are such that  $first(\alpha) \cap first(\beta) = \emptyset$ , then lookahead of one input symbol is enough to decide which rule to choose.
- **NB!** Holds only when  $\epsilon \notin \{first(\alpha) \cup first(\beta)\}$ .
- Otherwise, the set  $follow(A)$  should also be inspected.
- For each rule  $A \rightarrow \alpha$  define a set:

$$first^+(\alpha) = \begin{cases} (first(\alpha) \setminus \{\epsilon\}) \cup follow(A), & \text{kui } \epsilon \in first(\alpha) \\ first(\alpha), & \text{kui } \epsilon \notin first(\alpha) \end{cases}$$

- A CF-grammar is **LL(1)**, if for all (pairwise different) production rules  $A \rightarrow \alpha$  and  $A \rightarrow \beta$

$$first^+(\alpha) \cap first^+(\beta) = \emptyset$$



# Predictive Parsing

- **NB!** LL(1) grammar cannot be neither left-recursive nor ambiguous!
- Example:

$$S \rightarrow S a \mid \beta$$

- If  $\beta \neq \epsilon$ 
  - Then  $first(\beta) \subseteq first(S) = first(S a)$
  - Thus  $first^+(S a) \cap first^+(\beta) \neq \emptyset$
- If  $\beta = \epsilon$ 
  - Then  $a \in first(S a)$  ja  $a \in follow(S) = \{a, \$\}$
  - Thus  $first^+(S a) \cap first^+(\beta) \neq \emptyset$

## Predictive Parsing

- Often, a grammar can be transformed to LL(1) using:
  - left-recursion elimination;
  - left-factoring;
  - in worst case, one can generalize the grammar a bit (and check for removed restrictions after parsing).
- **Left-factoring** replaces rules with a common prefix with new ones, where the prefix is only in one RHS.
- Example:

$$\begin{array}{l} A \rightarrow B a C D \\ \quad | \quad B a C E \end{array} \quad \longrightarrow \quad \begin{array}{l} A \rightarrow B a C Z \\ Z \rightarrow D \\ \quad | \quad E \end{array}$$

# Predictive Parsing

Left-factoring algorithm:

- 1 For every non-terminal  $A \in N$  find a longest prefix  $\alpha$  which appears in two or more right-hand sides of  $A$  production rules.
- 2 If  $\alpha \neq \epsilon$ , then replace all productions of  $A$

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$$

with productions

$$\begin{aligned} A &\rightarrow \alpha Z \mid \gamma \\ Z &\rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n \end{aligned}$$

where  $Z \in N$  is a new non-terminal.

- 3 Repeat the process until none of right-hand sides have common prefix.

## Predictive Parsing

- **Recursive descent** parsing is a top-down parsing method where:
  - parser consists of a set of (mutually recursive) procedures, one for each non-terminal, recognizing sentential forms derivable from the corresponding non-terminal;
  - depending from the input, each procedure chooses a production rule and calls one after another procedures corresponding to non-terminals in the right-hand side of the production rule.
- Recursive descent parsing is commonly used for hand-written parsers.

## Predictive Parsing

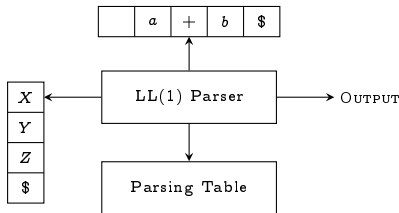
Example: let non-terminal  $A$  have production rules

$A \rightarrow a B \mid b C A \mid D E$

```
parseA() {  
  if token = a then {  
    token := nextWord(); parseB();  
  } else if token = b then {  
    token := nextWord(); parseC(); parseA();  
  } else {  
    parseD(); parseE();  
  }  
}
```

# Predictive Parsing

- Automatically generated LL(1) parsers are usually **table driven**:
  - construct a table  $M$  where rows and columns are indexed by non-terminals and terminals respectively;
  - cells of the table contain production rules to be chosen for the given non-terminal and input symbol.
- Structure of a table driven LL(1) parser:



# Predictive Parsing

LL(1) parsing algorithm:

```
push($); push(S);  
token := nextWord();  
while stack ≠ empty do {  
    A := pop();  
    if  $A \in N$  then {  
        if  $M[A, \textit{token}] = B_1 \dots B_n$  then {  
            push( $B_n$ ); ...; push( $B_1$ );  
        } else error();  
    } else if  $A = \textit{token}$  then {  
        token := nextWord(); pop();  
    } else error();  
}
```

## Shift-reduce Parsing

**Shift-reduce parsing** is a general method for a bottom-up syntax analysis:

- constructing a tree starts from leaves working up toward the root with a goal of "reducing" the input string to the start symbol.
- during parsing there is a forest of trees, which correspond to the different, already recognized, substrings;
- two basic actions:
  - **shift** reads a new input symbol and pushes it to the stack;
  - **reduce** applies production rules in reverse to the top of the stack; ie. it replaces a sequence of symbols in the top of the stack, forming a right-hand side of a rule, with the left-hand symbol of that rule;
- construction corresponds to right derivation.



# Shift-reduce Parsing

Example:

Input String: • *a b c c d e*

*S* → *a A B e*

*A* → *b c A* | *c*

*B* → *d*

shift

shift

shift

shift

reduce *A* → *c*

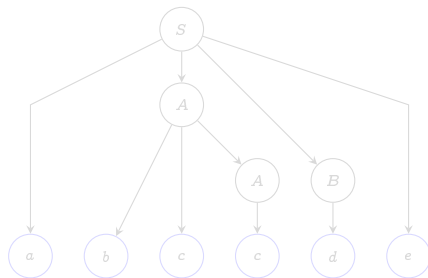
reduce *A* → *b c A*

shift

reduce *B* → *d*

shift

reduce *S* → *a A B e*



$S \Rightarrow_{rm} a A B e \Rightarrow_{rm} a A d e$   
 $\Rightarrow_{rm} a b c A d e \Rightarrow_{rm} a b c c d e$

# Shift-reduce Parsing

Example:

Input String:  $a \bullet b c c d e$

$S \rightarrow a A B e$

$A \rightarrow b c A \mid c$

$B \rightarrow d$

**shift**

shift

shift

shift

reduce  $A \rightarrow c$

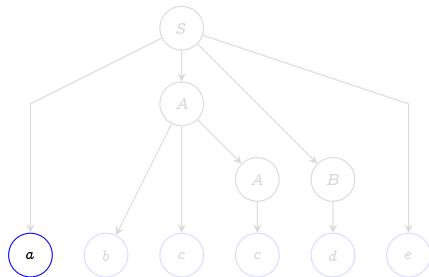
reduce  $A \rightarrow b c A$

shift

reduce  $B \rightarrow d$

shift

reduce  $S \rightarrow a A B e$



$S \xRightarrow{rm} a A B e \xRightarrow{rm} a A d e$   
 $\xRightarrow{rm} a b c A d e \xRightarrow{rm} a b c c d e$

# Shift-reduce Parsing

Example:

Input String:  $a b \bullet c c d e$

$S \rightarrow a A B e$

$A \rightarrow b c A \mid c$

$B \rightarrow d$

**shift**

**shift**

**shift**

**shift**

reduce  $A \rightarrow c$

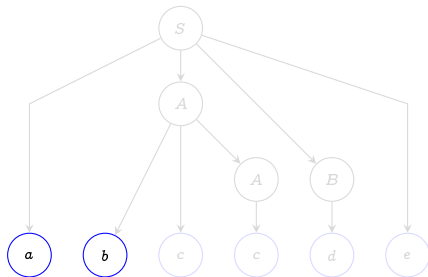
reduce  $A \rightarrow b c A$

**shift**

reduce  $B \rightarrow d$

**shift**

reduce  $S \rightarrow a A B e$



$S \xRightarrow{rm} a A B e \xRightarrow{rm} a A d e$   
 $\xRightarrow{rm} a b c A d e \xRightarrow{rm} a b c c d e$

# Shift-reduce Parsing

Example:

Input String:  $a b c \bullet c d e$

$S \rightarrow a A B e$

$A \rightarrow b c A \mid c$

$B \rightarrow d$

**shift**

**shift**

**shift**

**shift**

reduce  $A \rightarrow c$

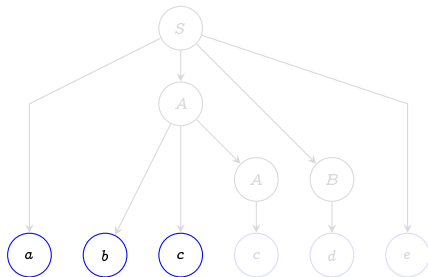
reduce  $A \rightarrow b c A$

**shift**

reduce  $B \rightarrow d$

**shift**

reduce  $S \rightarrow a A B e$



$S \xRightarrow{rm} a A B e \xRightarrow{rm} a A d e$   
 $\xRightarrow{rm} a b c A d e \xRightarrow{rm} a b c c d e$

# Shift-reduce Parsing

Example:

Input String:  $a b c c \bullet d e$

$S \rightarrow a A B e$

$A \rightarrow b c A \mid c$

$B \rightarrow d$

**shift**

**shift**

**shift**

**shift**

reduce  $A \rightarrow c$

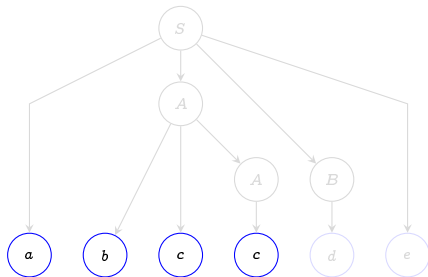
reduce  $A \rightarrow b c A$

**shift**

reduce  $B \rightarrow d$

**shift**

reduce  $S \rightarrow a A B e$



$S \Rightarrow_{rm} a A B e \Rightarrow_{rm} a A d e$   
 $\Rightarrow_{rm} a b c A d e \Rightarrow_{rm} a b c c d e$

# Shift-reduce Parsing

Example:

Input String:  $a b c c \bullet d e$

$S \rightarrow a A B e$

$A \rightarrow b c A \mid c$

$B \rightarrow d$

**shift**

**shift**

**shift**

**shift**

**reduce**  $A \rightarrow c$

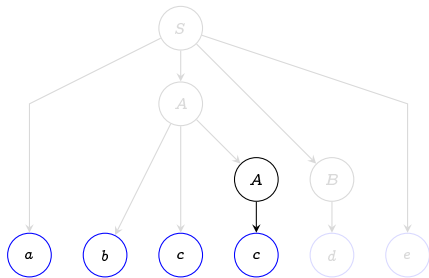
*reduce*  $A \rightarrow b c A$

*shift*

*reduce*  $B \rightarrow d$

*shift*

*reduce*  $S \rightarrow a A B e$



$S \Rightarrow_{rm} a A B e \Rightarrow_{rm} a A d e$   
 $\Rightarrow_{rm} a b c A d e \Rightarrow_{rm} a b c c d e$

# Shift-reduce Parsing

Example:

Input String:  $a b c c \bullet d e$

$S \rightarrow a A B e$

$A \rightarrow b c A \mid c$

$B \rightarrow d$

shift

shift

shift

shift

reduce  $A \rightarrow c$

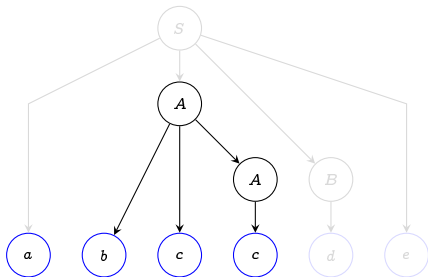
reduce  $A \rightarrow b c A$

shift

reduce  $B \rightarrow d$

shift

reduce  $S \rightarrow a A B e$



$S \xRightarrow{rm} a A B e \xRightarrow{rm} a A d e$   
 $\xRightarrow{rm} a b c A d e \xRightarrow{rm} a b c c d e$

# Shift-reduce Parsing

Example:

Input String:  $a b c c d \bullet e$

$S \rightarrow a A B e$

$A \rightarrow b c A \mid c$

$B \rightarrow d$

shift

shift

shift

shift

reduce  $A \rightarrow c$

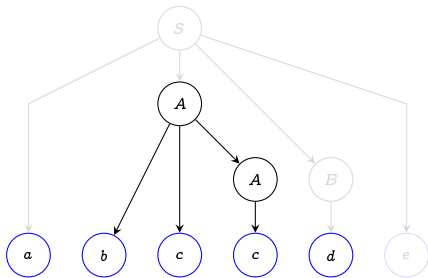
reduce  $A \rightarrow b c A$

shift

reduce  $B \rightarrow d$

shift

reduce  $S \rightarrow a A B e$



$S \xRightarrow{rm} a A B e \xRightarrow{rm} a A d e$   
 $\xRightarrow{rm} a b c A d e \xRightarrow{rm} a b c c d e$



# Shift-reduce Parsing

Example:

Input String:  $a b c c d \bullet e$

$S \rightarrow a A B e$

$A \rightarrow b c A \mid c$

$B \rightarrow d$

shift

shift

shift

shift

reduce  $A \rightarrow c$

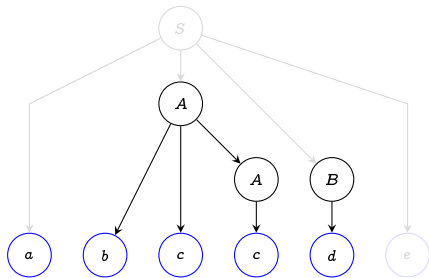
reduce  $A \rightarrow b c A$

shift

reduce  $B \rightarrow d$

shift

reduce  $S \rightarrow a A B e$



$S \xRightarrow{rm} a A B e \xRightarrow{rm} a A d e$   
 $\xRightarrow{rm} a b c A d e \xRightarrow{rm} a b c c d e$

# Shift-reduce Parsing

Example:

Input String:  $a b c c d e \bullet$

$S \rightarrow a A B e$

$A \rightarrow b c A \mid c$

$B \rightarrow d$

**shift**

**shift**

**shift**

**shift**

**reduce**  $A \rightarrow c$

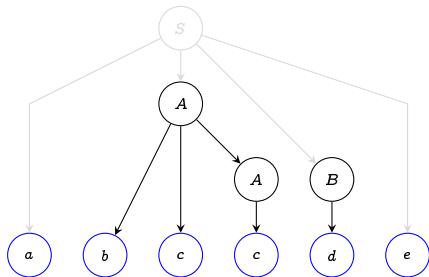
**reduce**  $A \rightarrow b c A$

**shift**

**reduce**  $B \rightarrow d$

**shift**

**reduce**  $S \rightarrow a A B e$



$S \xRightarrow{rm} a A B e \xRightarrow{rm} a A d e$   
 $\xRightarrow{rm} a b c A d e \xRightarrow{rm} a b c c d e$

# Shift-reduce Parsing

Example:

Input String:  $a b c c d e \bullet$

$S \rightarrow a A B e$

$A \rightarrow b c A \mid c$

$B \rightarrow d$

shift

shift

shift

shift

reduce  $A \rightarrow c$

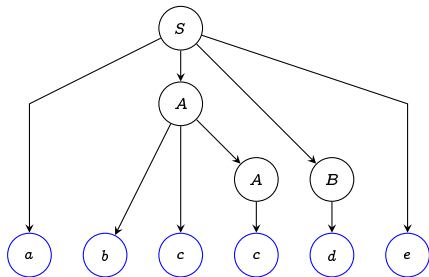
reduce  $A \rightarrow b c A$

shift

reduce  $B \rightarrow d$

shift

reduce  $S \rightarrow a A B e$



$S \Rightarrow_{rm} a A B e \Rightarrow_{rm} a A d e$   
 $\Rightarrow_{rm} a b c A d e \Rightarrow_{rm} a b c c d e$

# Shift-reduce Parsing

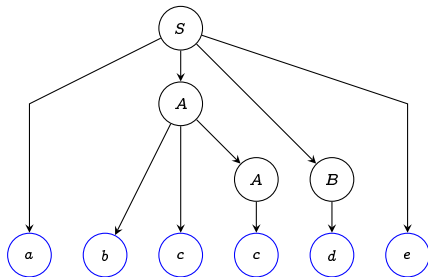
Example:

$S \rightarrow a A B e$   
 $A \rightarrow b c A \mid c$   
 $B \rightarrow d$

shift  
shift  
shift  
shift

reduce  $A \rightarrow c$   
reduce  $A \rightarrow b c A$   
shift  
reduce  $B \rightarrow d$   
shift  
reduce  $S \rightarrow a A B e$

Input String:



$S \xRightarrow{rm} a A B e \xRightarrow{rm} a A d e$   
 $\xRightarrow{rm} a b c A d e \xRightarrow{rm} a b c c d e$

## Shift-reduce Parsing

- Sentential form is called **right-sentential form**, if it occurs in the rightmost derivation of some sentence.
- A **handle** of a right-sentential form  $\gamma$  is a production rule  $A \rightarrow \beta$  and a position of  $\gamma$  where the string  $\beta$  may be found and replaced by  $A$  to produce the previous right-sentential form in a rightmost derivation of  $\gamma$ .
- Equivalently, a handle is a substring  $\beta$ , such that it matches RHS of some production and  $S \xRightarrow{*}_{rm} \delta A w \xRightarrow{rm} \delta \beta w = \gamma$ , where  $\beta, \gamma, \delta \in V^*$  and  $w \in T^*$ .
- The process of discovering a handle and reducing it to the appropriate LHS is called "**handle pruning**".
- **NB!** In the case of an unambiguous grammar, rightmost derivations, and hence handles, are unique.

## Shift-reduce Parsing

Example: given the grammar

$$\begin{aligned} S &\rightarrow a A B e \\ A &\rightarrow b c A \mid c \\ B &\rightarrow d \end{aligned}$$

and the rightmost derivation

$$S \Longrightarrow_{rm} a A B e \Longrightarrow_{rm} a A d e \Longrightarrow_{rm} a b c A d e$$

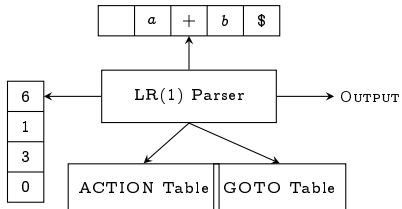
The handle of the right-sentential form  $abcAde$  is  $bcA$ .

## Shift-reduce Parsing

- In a stack based implementation of shift-reduce parser, the handle will always eventually appear on top of the stack.
- A **viable prefix** of a right-sentential form  $\gamma$  is any prefix of  $\gamma$  ending no farther right than the right end of the handle of  $\gamma$ .
- Viable prefixes are possible stacks of the shift-reduce parser!
- **NB!** The "language of viable prefixes" is regular!
- Hence, there is a finite automaton accepting viable prefixes.
- This automaton is the essential ingredient of all the LR parsing technique.

# LR Parsing

Structure of LR parser:





## LR Parsing

Skeleton of LR(1) parser:

```
push(Invalid); push( $s_0$ ); found := false;
token := nextWord();
while found  $\neq$  true do {
    s := top();
    if  $ACTION[s, token] = \text{reduce}(A \rightarrow \beta)$  then
        pop( $2 * |\beta|$ );
        s := top(); push( $A$ ); push( $GOTO[s, A]$ );
    else if  $ACTION[s, token] = \text{shift}(s_i)$  then
        push(token); push( $s_i$ );
        token := nextWord();
    else if  $ACTION[s, token] = \text{accept} \ \& \ token = \$$  then
        found := true;
    else report error;
}
report success;
```

## LR Parsing

- **LR(0)-item** (or simply **item**) is a production rule with a dot in the RHS.
- An item  $[A \rightarrow \alpha \cdot \beta]$  is **valid** for a viable prefix  $\varphi$  if there is a rightmost derivation  $S \xRightarrow{*}_{rm} \delta A w \xRightarrow{rm} \delta \alpha \beta w$  and  $\delta \alpha = \varphi$ .
- Item in the form  $[A \rightarrow \cdot \alpha]$  is an **initial item** and in the form  $[A \rightarrow \alpha \cdot]$  is a **complete item**.

# LR Parsing

- Example: given a grammar

$$\begin{aligned} S &\rightarrow a A c \\ A &\rightarrow A b \mid \varepsilon \end{aligned}$$

Its LR(0)-items are:

$$\begin{array}{ll} [S \rightarrow \cdot a A c] & [A \rightarrow \cdot A b] \\ [S \rightarrow a \cdot A c] & [A \rightarrow A \cdot b] \\ [S \rightarrow a A \cdot c] & [A \rightarrow A b \cdot] \\ [S \rightarrow a A c \cdot] & [A \rightarrow \cdot] \end{array}$$

- **NB!** For each CF-grammar the set of LR(0)-items is finite.

## LR Parsing

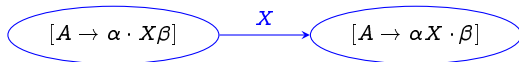
- A valid item  $[A \rightarrow \cdot \beta \gamma]$  means that the input seen so far is consistent with the use of  $A \rightarrow \beta \gamma$  immediately after the symbol on top of the stack.
- A valid item  $[A \rightarrow \beta \cdot \gamma]$  means that the input seen so far is consistent with the use of  $A \rightarrow \beta \gamma$  at this point of the parse, and that the parser has already recognized  $\beta$ .
- A complete valid item  $[A \rightarrow \beta \gamma \cdot]$  means that the parser has seen  $\beta \gamma$ , and that this is consistent with reducing to  $A$ .

## LR Parsing

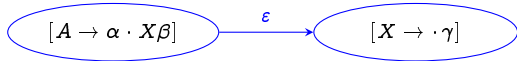
- If  $[A \rightarrow \alpha \cdot]$  is a complete valid item for a viable prefix  $\gamma$ , then  $A \rightarrow \alpha$  **might** have been used to derive  $\gamma w$  from  $\delta A w$ .
- However, in general, it might be also that this was not a case.
  - $[A \rightarrow \alpha \cdot]$  may be valid because of a different rightmost derivation  $S \xRightarrow{*}_{rm} \delta A w' \xRightarrow{rm} \gamma w'$ .
  - There could be several complete valid items for  $\gamma$ .
  - There could be a handle of  $\gamma w$  that includes some symbols of  $w$ .
- A context-free grammar for which knowing a complete valid item is enough to determine the previous right-sentential form is called **LR(0)**.

# LR Parsing

- A nondeterministic LR(0)-automaton is a NFA, where states are items and there are two kinds of transitions:
  - for every pair of items  $[A \rightarrow \alpha \cdot X\beta]$  and  $[A \rightarrow \alpha X \cdot \beta]$ , a transition labelled by a (terminal and nonterminal) symbol  $X$ .

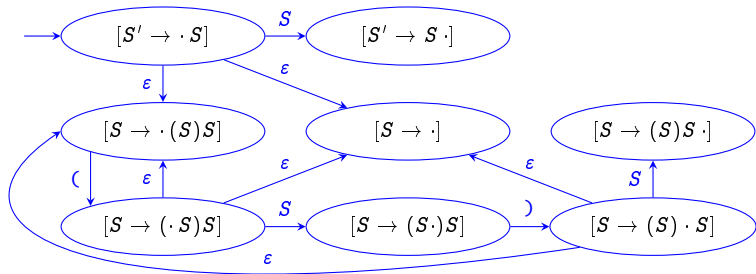


- for every pair of items  $[A \rightarrow \alpha \cdot X\beta]$  and  $[X \rightarrow \cdot \gamma]$ , a transition labelled by the empty symbol  $\epsilon$ .



- The grammar is **augmented** with a new start symbol  $S'$  and a single rule  $S' \rightarrow S$ .
- The state containing item  $[S' \rightarrow \cdot S]$  is the starting state of the automaton.

# LR Parsing



$$S' \rightarrow S$$

$$S \rightarrow (S)S \mid \epsilon$$

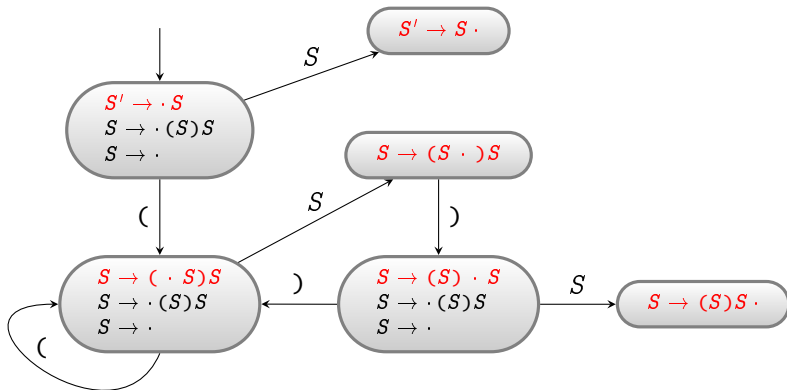
$$\begin{aligned} & [S' \rightarrow \cdot S] \\ & [S' \rightarrow S \cdot] \\ & [S \rightarrow \cdot (S)S] \\ & [S \rightarrow ( \cdot S)S] \\ & [S \rightarrow (S \cdot)S] \\ & [S \rightarrow (S) \cdot S] \\ & [S \rightarrow (S)S \cdot] \\ & [S \rightarrow (S)S \cdot] \end{aligned}$$

## LR Parsing

- DFA can be constructed from NFA by subset construction.
- DFA has sets of items as its states.
- Items in sets are called **kernel items** if they originate as targets of non- $\epsilon$ -transitions.
- Items added to the state during the  $\epsilon$ -closure step are called **closure items**.
- Kernel items uniquely determine the state.



# LR Parsing



# LR Parsing

LR(0) parsing algorithm:

```
while true do {
  state := top();
  if  $\exists[A \rightarrow \alpha \cdot X\beta] \in state \wedge X \in T$  then
    Y := getToken();
    if  $\exists[A \rightarrow \alpha \cdot Y\beta] \in state$  then
      shift(Y); push([A  $\rightarrow \alpha Y \cdot \beta$ ]);
    else error;
  else if  $\exists[A \rightarrow \gamma \cdot] \in state$  then
    if  $A = S' \wedge \gamma = S$  then accept;
    else pop(2 * |\gamma|);
      state := top();
      if  $\exists[B \rightarrow \alpha \cdot A\beta] \in state$  then
        push(A); push([B  $\rightarrow \alpha A \cdot \beta$ ]);
      else error;
}
```

## LR Parsing

- A grammar is **LR(0) grammar** iff any complete item  $[A \rightarrow \alpha \cdot] \in \text{state}$  is the only item of the state.
  - **Shift-reduce conflict**: if  $\exists [A \rightarrow \alpha \cdot X\beta] \in \text{state}$  where  $X \in T$ ;
  - **Reduce-reduce conflict**: if  $\exists [B \rightarrow \beta \cdot] \in \text{state}$ .
- **Reduce states**: states containing a complete item.
- **Shift states**: all other states.

## LR Parsing

- $SLR(1)$  = Simple LR(1) parsing.
- Uses the DFA of sets of LR(0) items.
- Uses the next lookahead token in the input string to direct its actions.
- Similar to LR(0) parsing, except that decision on which token to use is delayed until the last possible moment.
  - Consults the input token before a shift to make sure that an appropriate DFA transition exists.
  - Uses the Follow set of a nonterminal to decide if a reduction should be performed.
- Effective extension to LR(0) parsing that is powerful enough to handle almost all practical languages.

# LR Parsing

SLR(1) parsing algorithm:

```
while true do {  
  state := top(); X := getToken();  
  if  $\exists[A \rightarrow \alpha \cdot X\beta] \in state$  then  
    shift(X); push([A  $\rightarrow \alpha X \cdot \beta$ ]);  
  else if  $\exists[A \rightarrow \gamma \cdot] \wedge X \in Follow(A)$  then  
    if  $A = S' \wedge \gamma = S \wedge X = \$$  then accept;  
    else pop(2 * | $\gamma$ |);  
      state := top();  
      if  $\exists[B \rightarrow \alpha \cdot A\beta]$  then  
        push(A); push([B  $\rightarrow \alpha A \cdot \beta$ ]);  
      else error;  
  else error;  
}
```

## LR Parsing

- A grammar is **SLR(1)** iff it does not have the following two conflicts: for all states
  - **Shift-reduce conflict:**

$$\begin{aligned} \forall [A \rightarrow \alpha \cdot X \beta] \in \text{state} \wedge X \in T \\ \Rightarrow \neg(\exists [B \rightarrow \gamma \cdot] \wedge X \in \text{Follow}(B)) \end{aligned}$$

- **Reduce-reduce conflict:**

$$\forall [A \rightarrow \alpha \cdot] \wedge [B \rightarrow \beta \cdot] \Rightarrow \text{Follow}(A) \cap \text{Follow}(B) = \emptyset$$

## LR Parsing

- **LR(k)-item** is a pair  $[A \rightarrow \alpha \cdot \beta, w]$ , where  $A \rightarrow \alpha\beta$  is a production rule and  $w \in T^*$  is a word of length  $|w| \leq k$ .
- $[A \rightarrow \cdot \beta\gamma, w]$  means that the input seen so far is consistent with the use of  $A \rightarrow \beta\gamma$  immediately after the symbol on top of the stack.
- $[A \rightarrow \beta \cdot \gamma, w]$  means that the input seen so far is consistent with the use of  $A \rightarrow \beta\gamma$  at this point of the parse, and that the parser has already recognized  $\beta$ .
- $[A \rightarrow \beta\gamma \cdot, w]$  means that the parser has seen  $\beta\gamma$ , and that lookahead symbol of  $w$  is consistent with reducing to  $A$ .

## Shift-reduce Parsing

- Example: given a grammar

$$\begin{aligned} S &\rightarrow a A c \\ A &\rightarrow A b \mid \varepsilon \end{aligned}$$

Its LR(k)-items for lookahead string  $w$  are:

$$\begin{array}{ll} [S \rightarrow \cdot a A c, w] & [A \rightarrow \cdot A b, w] \\ [S \rightarrow a \cdot A c, w] & [A \rightarrow A \cdot b, w] \\ [S \rightarrow a A \cdot c, w] & [A \rightarrow A b \cdot, w] \\ [S \rightarrow a A c \cdot, w] & [A \rightarrow \cdot, w] \end{array}$$

- **NB!** For each CF-grammar the set of LR(k)-items is finite.



# LR Parsing

## Finding a closure of LR(1)-items

- $Closure(S)$  adds all the items implied by items already in the set of items  $S$ .
  - An item  $[A \rightarrow \beta \cdot B\delta, a]$  implies  $[B \rightarrow \cdot\tau, x]$  for each production  $B \rightarrow \tau$  and symbol  $x \in first(\delta a)$ .

- Algorithm:

```
 $Closure(S)$  {  
  while ( $S$  is still changing) do {  
     $\forall [A \rightarrow \beta \cdot B\delta, a] \in S$   
     $\forall B \rightarrow \tau \in P$   
     $\forall b \in first(\delta a)$   
    if  $[B \rightarrow \cdot\tau, b] \notin S$  then  
       $S := S \cup \{[B \rightarrow \cdot\tau, b]\}$   
  }  
}
```

- The algorithm terminates, as the set of items is finite.

## LR Parsing

- $Goto(s, X)$  computes the state that the parser would reach if it recognized  $X \in V$  while in state  $s$ :

$$Goto(s, X) = Closure(\{[A \rightarrow \beta X \cdot \delta, a] \mid [A \rightarrow \beta \cdot X \delta, a] \in s\})$$

- Building the Canonical Collection:

$s_0 := Closure(\{[S' \rightarrow \cdot S, \$]\});$

$\mathcal{S} := \{s_0\}; \Delta := \emptyset; k := 1;$

**while** ( $\mathcal{S}$  is still changing) **do** {

$\forall s_j \in \mathcal{S}, x \in V$

$s_k := Goto(s_j, x);$

$\Delta := \Delta \cup \{(s_j, x) \mapsto s_k\};$

**if**  $s_k \notin \mathcal{S}$  **then**

$\mathcal{S} := \mathcal{S} \cup \{s_k\}; k := k + 1;$

}

# LR Parsing

Example:

$$S \rightarrow E$$

$$E \rightarrow T - E \mid T$$

$$T \rightarrow F * T \mid F$$

$$F \rightarrow n$$

Symbol	<i>first</i>
$S$	$\{n\}$
$E$	$\{n\}$
$T$	$\{n\}$
$F$	$\{n\}$
$-$	$\{-\}$
$*$	$\{*\}$
$n$	$\{n\}$

Start state:

$$\begin{aligned} s_0 &= \text{Closure}(\{[S \rightarrow \cdot E, \$]\}) \\ &= \{ [S \rightarrow \cdot E, \$], [E \rightarrow \cdot T - E, \$], [E \rightarrow \cdot T, \$], \\ &\quad [T \rightarrow \cdot F * T, \$], [T \rightarrow \cdot F * T, -], [T \rightarrow \cdot F, \$], \\ &\quad [T \rightarrow \cdot F, -], [F \rightarrow \cdot n, \$], [F \rightarrow \cdot n, -], [F \rightarrow \cdot n, *] \} \end{aligned}$$

# LR Parsing

1st iteration:

$$\begin{aligned} s_1 &= \text{Goto}(s_0, E) \\ &= \{ [S \rightarrow E\cdot, \$] \} \\ s_2 &= \text{Goto}(s_0, T) \\ &= \{ [E \rightarrow T\cdot -E, \$], [E \rightarrow T\cdot, \$] \} \\ s_3 &= \text{Goto}(s_0, F) \\ &= \{ [T \rightarrow F\cdot *T, \$], [T \rightarrow F\cdot *T, -], \\ &\quad [T \rightarrow F\cdot, \$], [T \rightarrow F\cdot, -] \} \\ s_4 &= \text{Goto}(s_0, n) \\ &= \{ [F \rightarrow n\cdot, \$], [F \rightarrow n\cdot, -], [F \rightarrow n\cdot, *] \} \end{aligned}$$

# LR Parsing

2nd iteration:

$$\begin{aligned} s_5 &= \text{Goto}(s_2, -) \\ &= \{ [E \rightarrow T - \cdot E, \$], [E \rightarrow \cdot T - E, \$], [E \rightarrow \cdot T, \$], \\ &\quad [T \rightarrow \cdot F * T, -], [T \rightarrow \cdot F, -], \\ &\quad [T \rightarrow \cdot F * T, \$], [T \rightarrow \cdot F, \$], \\ &\quad [F \rightarrow \cdot n, *], [F \rightarrow \cdot n, -], [F \rightarrow \cdot n, \$] \} \\ s_6 &= \text{Goto}(s_3, *) \\ &= \{ [T \rightarrow F * \cdot T, \$], [T \rightarrow F * \cdot T, -], \\ &\quad [T \rightarrow \cdot F * T, \$], [T \rightarrow \cdot F * T, -], \\ &\quad [T \rightarrow \cdot F, \$], [T \rightarrow \cdot F, -], \\ &\quad [F \rightarrow \cdot n, \$], [F \rightarrow \cdot n, -], [F \rightarrow \cdot n, *] \} \end{aligned}$$

3rd iteration:

$$\begin{aligned} s_7 &= \text{Goto}(s_5, E) \\ &= \{ [E \rightarrow T - E \cdot, \$] \} \\ s_8 &= \text{Goto}(s_6, T) \\ &= \{ [T \rightarrow F * T \cdot, \$], [T \rightarrow F * T \cdot, -] \} \end{aligned}$$

# LR Parsing

Transition relation  $\Delta$

State	$E$	$T$	$F$	$-$	$*$	$n$
0	1	2	3			4
1						
2				5		
3					6	
4						
5	7	2	3			4
6		8	3			4
7						
8						

# LR Parsing

Generation of LR(1)-tables:

$\forall s_x \in \mathcal{S}$

$\forall i \in s_x$

**if**  $i = [A \rightarrow \alpha \cdot a\beta, b], \Delta(s_x, a) = s_k, a \in T$  **then**

$ACTION[x, a] := \text{shift}(k);$

**else if**  $i = [S' \rightarrow S \cdot, \$]$  **then**

$ACTION[x, a] := \text{accept};$

**else if**  $i = [A \rightarrow \beta \cdot, a]$  **then**

$ACTION[x, a] := \text{reduce}(A \rightarrow \beta);$

$\forall A \in N$

**if**  $\Delta(s_x, A) = s_k$  **then**

$GOTO[x, A] := k;$

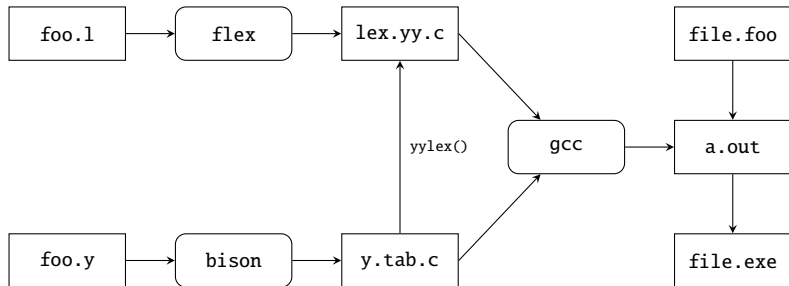
# LR Parsing

LR(1)-tables for the example grammar:

	<i>ACTION</i>				<i>GOTO</i>		
	<i>n</i>	—	*	\$	<i>E</i>	<i>T</i>	<i>F</i>
0	<b>shift(4)</b>				1	2	3
1				<b>accept</b>			
2		<b>shift(5)</b>		<b>reduce(3)</b>			
3		<b>reduce(5)</b>	<b>shift(6)</b>	<b>reduce(5)</b>			
4		<b>reduce(6)</b>	<b>reduce(6)</b>	<b>reduce(6)</b>			
5	<b>shift(4)</b>				7	2	3
6	<b>shift(4)</b>					8	3
7				<b>reduce(2)</b>			
8		<b>reduce(4)</b>		<b>reduce(4)</b>			



# Parser generator Bison



# Parser generator `Bison`

## Format of the input file:

- An input file of `Bison` has three parts:

definitions

%%

rules

%%

user code

- The same general structure as `Flex` has.

# Parser generator Bison

```
%{
#include <stdio.h>
%}
%token INTEGER
%%
program:
    | program expr '\n' { printf("%d\n", $2); }
    ;
expr:
    INTEGER                { $$ = $1; }
    | expr '+' expr        { $$ = $1 + $3; }
    | expr '-' expr        { $$ = $1 - $3; }
    ;
%%
void yyerror(char *s) {
    fprintf(stderr, "%s\n", s);
}
int main(void) {
    yyparse();
    return 0;
}
```

## Parser generator Bison

```
%{  
#include <stdio.h>  
%}  
%token INTEGER  
%%
```

The definitions part consists of:

- C code surrounded by `%{` and `%}`, which is copied verbatim into the generated file;
- **Bison** declarations:
  - `%token` list of terminal symbols (used in production rules, but also in the scanner);
  - `%start` declaration of the start symbol (if absent, the first non-terminal is the start symbol);
  - `%union`, `%left`, `%right`, ...

## Parser generator Bison

```
program:
  | program expr '\n' { printf("%d\n", $2); }
  ;
expr:
  INTEGER           { $$ = $1; }
  | expr '+' expr   { $$ = $1 + $3; }
  | expr '-' expr   { $$ = $1 - $3; }
  ;
```

- The second part consists of production rules.
- Must contain at least one rule.
- RHS of a rule is built of terminal and non-terminal symbols, and may also contain actions.
- Terminal symbols may be the ones declared before or singleton characters.

## Parser generator Bison

```
program:
  | program expr '\n' { printf("%d\n", $2); }
  ;
expr:
  INTEGER           { $$ = $1; }
  | expr '+' expr   { $$ = $1 + $3; }
  | expr '-' expr   { $$ = $1 - $3; }
  ;
```

- Actions are C code fragments surrounded by curly braces.
- They correspond to semantic rules of attribute grammars.
- Can refer to attribute values of grammar symbols appearing in the rule:
  - \$\$ corresponds to the value of LHS symbol;
  - \$1 corresponds to the value of the first symbol in RHS;
  - \$2 corresponds to the value of the second symbol in RHS;
  - ...

## Parser generator Bison

```
program:
  | program expr '\n' { printf("%d\n", $2); }
  ;
expr:
  INTEGER           { $$ = $1; }
  | expr '+' expr   { $$ = $1 + $3; }
  | expr '-' expr   { $$ = $1 - $3; }
  ;
```

- Actions are usually at the end of RHS, and they are executed when the rule is reduced.
- They may appear also in between the symbols, in which case it's equivalent of having an extra non-terminal in the place; this non-terminal has the action as RHS of its rule (and is otherwise empty).
- If an action is missing, then the default action is

$$\{ \$\$ = \$1; \}$$

## Parser generator Bison

```
%%  
void yyerror(char *s) {  
    fprintf(stderr, "%s\n", s);  
}  
int main(void) {  
    yyparse();  
    return 0;  
}
```

The third part of the specification is a C code which will be copied into the generated file verbatim.

- Most important functions:
  - main() calls yyparse();
  - yyerror() reports of syntax errors;
  - yylex() recognizes tokens (usually defined in a **Flex** generated scanner).
- The third part may be absent in what case the separator line may also be missing.



## Parser generator **Bison**

- The function `yyparse()` uses the function `yylex()` to get the next token.
- The interface between scanner and parser is specified in **Bison**:
  - terminal symbols are declared with the command `%token;`
  - single character terminal symbols may be left undeclared;
  - the return value of `yylex()` is either a declared terminal symbol or a single character.
- Scanner should include the header file `"*.tab.h"`
  - can be generated automatically by **Bison** with the argument option `-d`.
- An alternative is to include the file of scanner `"lex.yy.c"` in the third part of the parser specification file.

## Parser generator **Bison**

- Attribute values of non-terminals are in the variable `yylval`.
- Attributes are of type `YYSTYPE`, which by default is `int`.
- If attributes of different non-terminals have different types, one has to:
  - declare all types with the command `%union`

```
%union {  
    type1 name1;  
    type2 name2;  
    ...  
}
```
  - specify the types of terminal and non-terminal symbols

```
%token <name> TOKEN  
%type <name> non-terminal
```
- An attribute value corresponding to a terminal symbol (and assigned by the scanner) should be a field of the union (`yylval.name`).

## Parser generator **Bison**

- Shift-reduce actions are decided using one symbol lookahead.
- Conflicts are resolved by precedence and default rules:
  - commands `%left`, `%right` and `%nonassoc` are used to determine associativity and precedence of symbols;
  - the precedence is determined implicitly by textual ordering (later one have higher precedence);
  - the precedence of a rule is the same as the last nonterminal in RHS (can be changed using the command `%prec`);
  - shift/reduce conflict is solved by preferring shift action;
  - reduce/reduce conflict is solved by using textually the first production rule.