# 3 Parsing

## 3.1   Context-Free Grammars

- Now that we have moved on to the **Parse** phase of our compiler, we can—thanks to lexical analysis—work with tokens (instead of character sequences).

- The **parser** analyses the incoming token stream and

  ① checks that the stream is indeed valid according to the **syntax rules** of the Tiger language, and

  ② builds a **parse tree** for the token stream.

- The parse tree is already suitable to **derive the actual meaning** of the source program. We thus need to be careful to build the correct parse tree for a given program.

- We will use **context-free grammars** (**CFGs**) to formally specify the syntax rules of our source language.

- Why not use REs for this task?

  [ *An RE can describe valid sequences of symbols, namely the RE's language. Let's just use REs operating on tokens rather than single characters.* ]
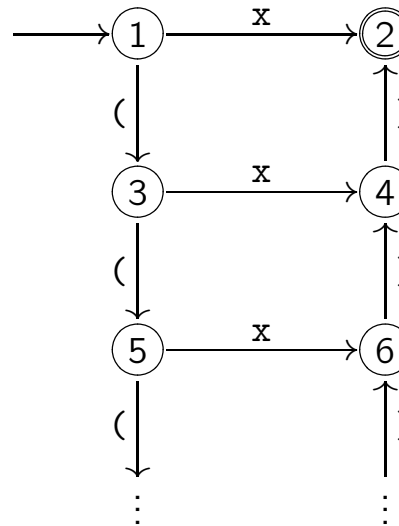
  – REs are simply not powerful enough to represent the **recursive** nature of the syntax of typical source languages.
  – **Example:** consider the simple syntax rule (this is actually a tiny CFG):

$$E \quad \rightarrow \quad x \mid ( \, E \, )$$

  Clearly, our intention is $L(E) = \{\, \texttt{"x"}, \texttt{"(x)"}, \texttt{"((x))"}, \texttt{"(((x)))"}, \dots \,\}$.

– If we can define this language via an RE then there must be an equivalent DFA $A$ with $L(A) = L(E)$.

A would need to look similar to what is shown below:



– No DFA will be "deep enough" to accept all parenthesized inputs ("*A DFA cannot count.*").

• Parsing algorithms for CFGs therefore use a **stack** to remember which input tokens have already been consumed.

- **Definition CFG:**

  - A CFG is defined by a **set of productions** (**rules**) of the form

$$S \quad \rightarrow \quad \underbrace{S\ S\ \cdots\ S}_{\text{zero or more}}$$

  - Each **symbol** $S$ is either a **terminal** ($=$ token) or a **non-terminal**.
  - Only non-terminals may appear on the left-hand sides (i.e., left of the $\rightarrow$) of productions.
  - One non-terminal symbol is marked as the **start symbol**.

  - **Example:** a CFG describing a straight-line programming language (start symbol $S'$):

$$
\begin{array}{lll}
_0\ S' \rightarrow S\ \$ & _4\ E \rightarrow \texttt{id} & _8\ L \rightarrow E \\
_1\ S \rightarrow S\ ;\ S & _5\ E \rightarrow \texttt{num} & _9\ L \rightarrow L\ ,\ E \\
_2\ S \rightarrow \texttt{id} := E & _6\ E \rightarrow E + E & \\
_3\ S \rightarrow \texttt{print}\ (\ L\ ) & _7\ E \rightarrow (\ S\ ,\ E\ ) &
\end{array}
$$

- **N.B.**: $\$$ is the end-of-file (EOF) symbol which marks the end of any input token stream.

# 3.1.1 Derivations (and Parse Trees)

- Just like an RE or DFA/NFA, a CFG describes a language.

- A **derivation** is like a proof that a given token sequence is **accepted** by a CFG $G$:

  ① Select the start symbol $S$ of $G$ as the next non-terminal $E$ to expand.

  ② Replace the selected non-terminal $E$ by *some* right-hand side of an $E$-production in $G$.

  ③ This right-hand side will, in general, contain more non-terminals. Choose of these non-terminals to be the new $E$, *goto* ②.
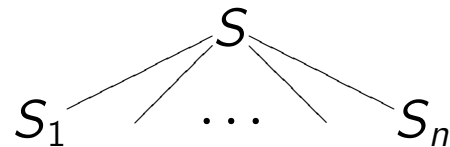  If only terminals remain, *stop*.

  The resulting token sequence is in $L(G)$, i.e., accepted by CFG $G$.

- **Example:** derive a token sequence from the start symbol $S'$ of the CFG shown ealier (the next non-terminal chosen for replacement is underlined)[6]:
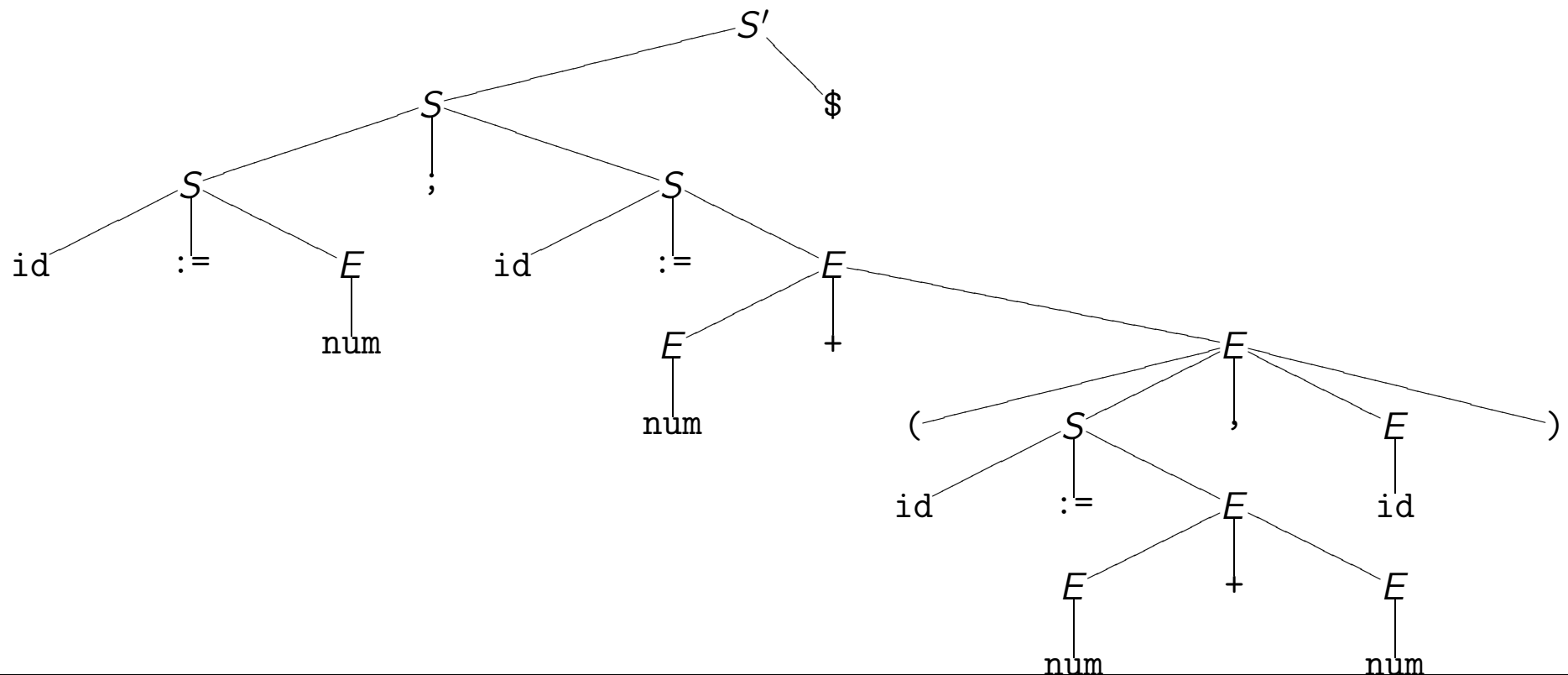
$S'$

$S$ \$

$S$ ; $S$ \$

$S$ ; id := $E$ \$

id := $E$ ; id := $E$ \$

id := num ; id := $E$ \$

id := num ; id := $E$ + $E$ \$

id := num ; id := $E$ + ( $S$ , $E$ ) \$

id := num ; id := id + ( $S$ , $E$ ) \$

id := num ; id := id + ( id := $E$ , $E$ ) \$

id := num ; id := id + ( id := $E$ + $E$ , $E$ ) \$

id := num ; id := id + ( id := $E$ + $E$ , id ) \$

id := num ; id := id + ( id := num + $E$ , id ) \$

id := num ; id := id + ( id := num + num , id ) \$

---

[6]We omit the semantic values associated with token types id and num.

- The derivation of a token sequence yields the **parse tree** associated with the derivation:

  - For each expansion of a production $S \rightarrow S_1 \cdots S_n$ build the parse tree fragment



- **Example:** parse tree associated with the derivation shown on previous slide:

- **Remarks:**

  – In any parse tree, the **start symbol** will always appears at the **root**.

  – The **front** of the parse tree, read from left to right, is the input token sequence.

## 3.1.2   Ambigiuous Grammars
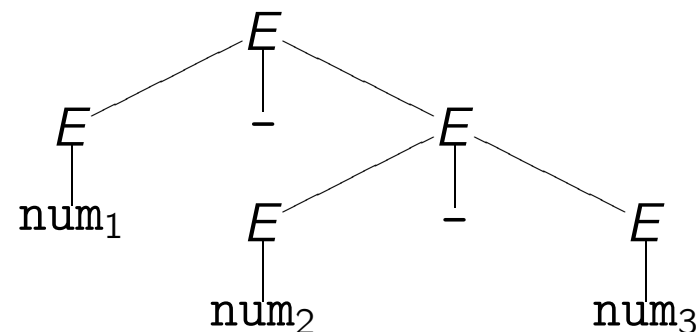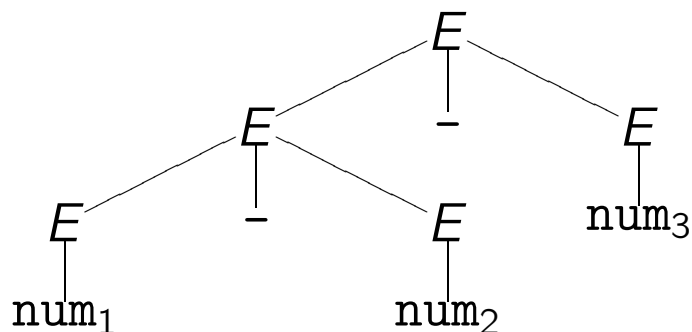
- Compilers use parse tree to **infer the meaning** of a given source program.

  – If we can build more than one parse tree (e.g., have more than one derivation) for a given token sequence, which meaning shall be the correct one?

- Grammars that associate more than one parse tree with a given token sequence are **ambigiuous**. Ambiguous grammars are *not* suitable for parsing.

- **Example:** consider the CFG $G$ for expressions over the usual binary arithmetic operators (with their well-known mathematical semantics):

$$
\begin{array}{rcl}
E & \rightarrow & \texttt{id} \\
E & \rightarrow & \texttt{num} \\
E & \rightarrow & (\ E\ )
\end{array}
\qquad
\begin{array}{rcl}
E & \rightarrow & E * E \\
E & \rightarrow & E\ /\ E \\
E & \rightarrow & E + E \\
E & \rightarrow & E - E
\end{array}
$$

- $G$ is ambiguous since we may derive two parse trees for the token sequence $\texttt{num}_1\texttt{-num}_2\texttt{-num}_3$:



*Shall the binary operators* **associate** *to the left or right?*

– *G* is ambigiuous since we may derive two parse trees for the token sequence $\texttt{num}_1\texttt{+num}_2\texttt{*num}_3$:



- The associated meanings would yield different machine code and thus different results.

- It is often possible to **transform** a CFG such that

  ① the accepted language does *not* change, and
  ② any input token sequence is assigned exactly one parse tree.

– **Example:** transform the expression grammar $G$ into a CFG $G'$ such that the parse trees reflect usual mathematical semantics:

$$
\begin{array}{lcl lcl lcl}
E & \to & E + T & \quad T & \to & T * F & \quad F & \to & \texttt{id} \\
E & \to & E - T & \quad T & \to & T / F & \quad F & \to & \texttt{num} \\
E & \to & T & \quad T & \to & F & \quad F & \to & (\,E\,)
\end{array}
$$

($E, T, F$ stand for *expression*, *term*, *factor*, respectively.)

– **Note:** with $G'$, we will *never* be able to derive parse trees shaped as shown below:



– The "layering" in the grammar expresses operator **priority**.
– The **left recursion** in the $E, T$-productions expresses operator **associativity**.

# 3.2 LR Parsing

- We will now discuss the LR($k$) family of **parsing algorithms** which also lies at the core of modern parser program generators, like `yacc` or `bison`.

- An LR($k$) parser

  **L -** reads the input token sequence from **left** to right,

  **R -** performs a **right-most** derivation during parsing,

  **$k$ -** uses a **lookahead** of $k$ tokens to guide itself.

- Performing a *right-most* derivation when reading the input *left* to right might sound odd at first, but:

  ① LR($k$) decides which production to use *after* it has seen all tokens on the right-hand side of the production ($+k$ tokens more).

  ② Before this decision is made by the parser, it temporarily saves the tokens on a **stack**.

- **Overview LR parsing:**

  – The parser maintains its **input** (token sequence) and a **stack** (of CFG symbols).
  – While the parser reads the input it can perform two **actions**:

  ① **shift**:
  Consume first token of input, push token onto the stack:

  | Stack | Input | Action |
  |---|---|---|
  | $S_1 \cdots S_n$ | id + num $\cdots$ \$ | shift |
  | $S_1 \cdots S_n$ id | + num $\cdots$ \$ | |

  ② **reduce** (production $S \to S_1 \cdots S_n$):
  Pop $S_n, \cdots, S_1$ from the stack, then push $S$ onto the stack:

  | Stack | Input | Action |
  |---|---|---|
  | $S_0\ S_1 \cdots S_n$ | id + num $\cdots$ \$ | reduce $S \to S_1 \cdots S_n$ |
  | $S_0\ S$ | id + num $\cdots$ \$ | |

  – Shifting the \$ token is also called **accepting** (the parser has consumed all tokens).

| Stack | Input | Action |
|---|---|---|
| id := num ; id := id + ( id := num + num , id ) $ | | *shift* |
| id := num ; id := id + ( id := num + num , id ) $ | | *shift* |
| id := num ; id := id + ( id := num + num , id ) $ | | *shift* |
| id := num ; id := id + ( id := num + num , id ) $ | | *reduce E → num* |
| id := E ; id := id + ( id := num + num , id ) $ | | *reduce S → id := E* |
| S ; id := id + ( id := num + num , id ) $ | | *shift* |
| S ; id := id + ( id := num + num , id ) $ | | *shift* |
| S ; id := id + ( id := num + num , id ) $ | | *shift* |
| S ; id := id + ( id := num + num , id ) $ | | *shift* |
| S ; id := id + ( id := num + num , id ) $ | | *reduce E → id* |
| S ; id := E + ( id := num + num , id ) $ | | *shift* |
| S ; id := E + ( id := num + num , id ) $ | | *shift* |
| S ; id := E + ( id := num + num , id ) $ | | *shift* |
| S ; id := E + ( id := num + num , id ) $ | | *shift* |
| S ; id := E + ( id := num + num , id ) $ | | *shift* |
| S ; id := E + ( id := num + num , id ) $ | | *reduce E → num* |
| S ; id := E + ( id := E + num , id ) $ | | *shift* |
| S ; id := E + ( id := E + num , id ) $ | | *shift* |
| S ; id := E + ( id := E + num , id ) $ | | *reduce E → num* |
| S ; id := E + ( id := E + E , id ) $ | | *reduce E → E + E* |
| S ; id := E + ( id := E , id ) $ | | *reduce S → id := E* |
| S ; id := E + ( S , id ) $ | | *shift* |
| S ; id := E + ( S , id ) $ | | *shift* |
| S ; id := E + ( S , id ) $ | | *reduce E → id* |
| S ; id := E + ( S , E ) $ | | *shift* |
| S ; id := E + ( S , E ) $ | | *reduce E → ( S , E )* |
| S ; id := E + E $ | | *reduce E → E + E* |
| S ; id := E $ | | *reduce S → id := E* |
| S ; S $ | | *reduce S → S ; S* |
| S $ | | *accept (shift $)* |

- ## Remarks:

  - The concatenation of stack and remaining input always represents the whole input.
  - Read bottom-up, the LR parser determines a right-most derivation of the input.

- The LR parser needs some control to decide when to shift and when to reduce.

  > **Idea:** We *cannot* use a DFA applied to the input (DFAs/REs are too weak to handle CFGs), but we can **apply the DFA idea to the LR stack**.

  - Use a **parsing table** (transition matrix) $goto[d][S]$ to implement the DFA ($d$: DFA state, $S$: CFG symbol).
  - The DFA edges are labelled with the symbols that may occur on the LR stack, i.e., tokens *and* non-terminals.
  - As an optimization, with each symbol we push onto the stack we also push the current DFA state (later, when we pop the stack, we can easily determine which state we are in after the pops are done).
  - The transition matrix entries indicate which **action** the LR parser has to perform:

| Matrix entry | Action |
| --- | --- |
| $s\, n$ | *shift* (next state is $n$) |
| $r\, S \rightarrow S_1 \cdots S_k$ | *reduce* by given rule |
| $a$ | *accept* |

- **LR parsing algorithm:**

  - **Given:** input token sequence $t_0 \, t_1 \cdots t_n \, \$$, LR parsing table $goto[][]$, empty LR stack[7].
  - **Output:** *accept* or *error*.

```
ip ← 0;          /* pointer into input token sequence */
push(1);         /* we start in state 1 */


while true do
    s ← top();
    switch goto[s][t_ip] do
        case s n:                       /* shift */
            push(t_ip);
            push(n);
            ip ← ip + 1;
        case r S → S_1 ··· S_k:          /* reduce */
            pop() 2 × k times;
            s' ← top();
            push(S);
            push(goto[s'][S]);
        case a:                         /* accept */
            return accept;

        otherwise
            return error;
```

---

[7]With the usual stack operations *push()*, *pop()*, *top()*.

**CFG symbol**

| | id | num | print | ; | , | + | := | ( | ) | $ | | S | E | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | s4 | | s7 | | | | | | | | | 2 | | |
| 2 | | | | s3 | | | | | | a | | | | |
| 3 | s4 | | s7 | | | | | | | | | 5 | | |
| 4 | | | | | | | s6 | | | | | | | |
| 5 | | | | $r\,S \to S;S$ | $r\,S \to S;S$ | | | | | $r\,S \to S;S$ | | | | |
| 6 | s20 | s10 | | | | | | s8 | | | | | 11 | |
| 7 | | | | | | | | s9 | | | | | | |
| 8 | s4 | | s7 | | | | | | | | | 12 | | |
| 9 | s20 | s10 | | | | | | s8 | | | | | 15 | 14 |
| 10 | | | | $r\,E \to \texttt{num}$ | $r\,E \to \texttt{num}$ | $r\,E \to \texttt{num}$ | | | $r\,E \to \texttt{num}$ | $r\,E \to \texttt{num}$ | | | | |
| 11 | | | | $r\,S \to \texttt{id:=}E$ | $r\,S \to \texttt{id:=}E$ | s16 | | | | $r\,S \to \texttt{id:=}E$ | | | | |
| 12 | | | | s3 | s18 | | | | | | | | | |
| 13 | | | | $r\,S \to \texttt{print}(L)$ | $r\,S \to \texttt{print}(L)$ | | | | | $r\,S \to \texttt{print}(L)$ | | | | |
| 14 | | | | | s19 | | | | s13 | | | | | |
| 15 | | | | | $r\,L \to E$ | | | | $r\,L \to E$ | | | | | |
| 16 | s20 | s10 | | | | | | s8 | | | | | 17 | |
| 17 | | | | $r\,E \to E\texttt{+}E$ | $r\,E \to E\texttt{+}E$ | s16 | | | $r\,E \to E\texttt{+}E$ | $r\,E \to E\texttt{+}E$ | | | | |
| 18 | s20 | s10 | | | | | | s8 | | | | | 21 | |
| 19 | s20 | s10 | | | | | | s8 | | | | | 23 | |
| 20 | | | | $r\,E \to \texttt{id}$ | $r\,E \to \texttt{id}$ | $r\,E \to \texttt{id}$ | | | $r\,E \to \texttt{id}$ | $r\,E \to \texttt{id}$ | | | | |
| 21 | | | | | | | | | s22 | | | | | |
| 22 | | | | $r\,E \to (S,E)$ | $r\,E \to (S,E)$ | $r\,E \to (S,E)$ | | | $r\,E \to (S,E)$ | $r\,E \to (S,E)$ | | | | |
| 23 | | | | | $r\,L \to L,E$ | s16 | | | $r\,L \to L,E$ | | | | | |

---

| Stack | Input | Action |
|---|---|---|
| $_1$ | id := num ; id := id + ( id := num + num , id ) \$ | shift |
| $_1$ id $_4$ | := num ; id := id + ( id := num + num , id ) \$ | shift |
| $_1$ id $_4$ := $_6$ | num ; id := id + ( id := num + num , id ) \$ | shift |
| $_1$ id $_4$ := $_6$ num $_{10}$ | ; id := id + ( id := num + num , id ) \$ | reduce $E \to$ num |
| $_1$ id $_4$ := $_6$ $E$ $_{11}$ | ; id := id + ( id := num + num , id ) \$ | reduce $S \to$ id := $E$ |
| $_1$ $S$ $_2$ | ; id := id + ( id := num + num , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ | id := id + ( id := num + num , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ | := id + ( id := num + num , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ | id + ( id := num + num , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ id $_{20}$ | + ( id := num + num , id ) \$ | reduce $E \to$ id |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ | ( id := num + num , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ | ( id := num + num , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ | id := num + num , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ id $_4$ | := num + num , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ id $_4$ := $_6$ | num + num , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ id $_4$ := $_6$ num $_{10}$ | + num , id ) \$ | reduce $E \to$ num |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ id $_4$ := $_6$ $E$ $_{11}$ | + num , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ | num , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ num $_{10}$ | , id ) \$ | reduce $E \to$ num |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ $E$ $_{17}$ | , id ) \$ | reduce $E \to E$ + $E$ |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ id $_4$ := $_6$ $E$ $_{11}$ | , id ) \$ | reduce $S \to$ id := $E$ |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ $S$ $_{12}$ | , id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ $S$ $_{12}$ , $_{18}$ | id ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ $S$ $_{12}$ , $_{18}$ id $_{20}$ | ) \$ | reduce $E \to$ id |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ $S$ $_{12}$ , $_{18}$ $E$ $_{21}$ | ) \$ | shift |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ ( $_8$ $S$ $_{12}$ , $_{18}$ $E$ $_{21}$ ) $_{22}$ | \$ | reduce $E \to$ ( $S$ , $E$ ) |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ + $_{16}$ $E$ $_{17}$ | \$ | reduce $E \to E$ + $E$ |
| $_1$ $S$ $_2$ ; $_3$ id $_4$ := $_6$ $E$ $_{11}$ | \$ | reduce $S \to$ id := $E$ |
| $_1$ $S$ $_2$ ; $_3$ $S$ $_5$ | \$ | reduce $S \to S$ ; $S$ |
| $_1$ $S$ $_2$ | \$ | accept (shift \$) |

# 3.3  LR Parsing Table Generation

- We are still missing a method of deriving an **LR parsing table** from a given CFG. This is what we discuss next.

  – In what follows, we will consider this example CFG (start symbol $S$):

$$
\begin{array}{lll}
{}_1\, S \rightarrow A\; C\; \$ & {}_4\, A \rightarrow \text{a}\; B\; C\; \text{d} & {}_7\, B \rightarrow \text{b}\; B \\
{}_2\, C \rightarrow \text{c} & {}_5\, A \rightarrow B\; Q & {}_8\, B \rightarrow \text{d} \\
{}_3\, C \rightarrow & {}_6\, A \rightarrow & {}_9\, Q \rightarrow \text{q}
\end{array}
$$

- The LR parsing algorithm imitates a DFA. So, given a CFG, how can we characterize the **state** the LR parser is in currently?

  – When parsing starts, the LR parser tries to derive the whole input from the right-hand side of the start symbol. No token has been consumed yet.
  We indicate this state via the **LR item** shown below:

$$
S \quad \rightarrow \quad \bullet\, A\; C\; \$
$$

– The dot ● (i.e., the current parser position) is just before non-terminal $A$. The start of input must be derivable from non-terminal $A$.

We indicate this by **adding the corresponding items** to our state:

$$\begin{array}{rl}
S \to & \bullet\ A\ C\ \$ \\
A \to & \bullet\ \text{a}\ B\ C\ \text{d} \\
A \to & \bullet\ B\ Q \\
A \to & \bullet
\end{array}$$

– As long as ● is just before a non-terminal (here: $B$), keep adding the corresponding items (*closure*):

$$\begin{array}{rl}
S \to & \bullet\ A\ C\ \$ \\
A \to & \bullet\ \text{a}\ B\ C\ \text{d} \\
A \to & \bullet\ B\ Q \\
A \to & \bullet \\
B \to & \bullet\ \text{b}\ B \\
B \to & \bullet\ \text{d}
\end{array}$$

– Once the closure is stable, the LR state is complete.

- We then label each item of the state with an LR action: *shift, reduce, goto,* and *accept* (shifting $):

$$
\begin{bmatrix}
S \rightarrow \ \bullet \ A \ C \ \$ & goto \ 2 \\
\hline
A \rightarrow \ \bullet \ \texttt{a} \ B \ C \ \texttt{d} & shift \ 3 \\
A \rightarrow \ \bullet \ B \ Q & goto \ 4 \\
A \rightarrow \ \bullet & reduce \ A \rightarrow \\
B \rightarrow \ \bullet \ \texttt{b} \ B & shift \ 5 \\
B \rightarrow \ \bullet \ \texttt{d} & shift \ 6
\end{bmatrix}_{①}
$$

- *shift n*: consume token after ●, then goto state *n*
  (in state *n*: ● moved past the token).
- *goto n*: a right-hand side of the non-terminal after ● has just been recognized
  (reduced), goto state *n*
  (in state *n*: move ● past the non-terminal).
- *reduce*: ● is at the end of a right-hand side, the parser is ready to reduce.
- *accept*: ● would move past $.

- Now simply continue to compute LR states:

$$
\begin{bmatrix}
S \rightarrow \quad \bullet \ A\ C\ \$ & \quad goto\ 2 \\
\hline
A \rightarrow \quad \bullet \ \texttt{a}\ B\ C\ \texttt{d} & \quad shift\ 3 \\
A \rightarrow \quad \bullet \ B\ Q & \quad goto\ 4 \\
A \rightarrow \quad \bullet & \quad reduce\ A \rightarrow \\
B \rightarrow \quad \bullet \ \texttt{b}\ B & \quad shift\ 5 \\
B \rightarrow \quad \bullet \ \texttt{d} & \quad shift\ 6
\end{bmatrix}_{①}
$$

$$
\begin{bmatrix}
S \rightarrow A\ \bullet\ C\ \$ & \quad goto\ 7 \\
\hline
C \rightarrow \quad \bullet\ \texttt{c} & \quad shift\ 8 \\
C \rightarrow \quad \bullet & \quad reduce\ C \rightarrow
\end{bmatrix}_{②}
$$

$$
\begin{bmatrix}
A \rightarrow \texttt{a}\ \bullet\ B\ C\ \texttt{d} & \quad goto\ 9 \\
\hline
B \rightarrow \quad \bullet\ \texttt{b}\ B & \quad shift\ 5 \\
B \rightarrow \quad \bullet\ \texttt{d} & \quad shift\ 6
\end{bmatrix}_{③}
$$

  - Now complete the computation for states ④ … ⑮!

# 3.3.1   Shift/Reduce Conflicts

- In some states, it is not obvious for the parser if *shifting* or *reducing* is the appropriate action: the LR parser table would allow both alternatives.

  - Consider this simple expression CFG and its LR states:

$$S \rightarrow E \; ; \; \$ \qquad E \rightarrow T$$
$$E \rightarrow T + E \qquad T \rightarrow \mathrm{x}$$

$$
\left[
\begin{array}{lll}
S \rightarrow \quad \bullet \; E \; ; \; \$ & goto\ 2 \\
\hline
E \rightarrow \quad \bullet \; T + E & goto\ 3 \\
E \rightarrow \quad \bullet \; T & goto\ 3 \\
T \rightarrow \quad \bullet \; \mathrm{x} & shift\ 4
\end{array}
\right]_{①}
\qquad
\left[
\begin{array}{lll}
E \rightarrow T + \; \bullet \; E & goto\ 7 \\
\hline
E \rightarrow \quad \bullet \; T + E & goto\ 3 \\
E \rightarrow \quad \bullet \; T & goto\ 3 \\
T \rightarrow \quad \bullet \; \mathrm{x} & shift\ 4
\end{array}
\right]_{⑥}
$$

$$
\left[\; S \rightarrow E \; \bullet \; ; \; \$ \quad shift\ 5 \;\right]_{②}
\qquad
\left[
\begin{array}{ll}
E \rightarrow T \; \bullet \; + E & shift\ 6 \\
E \rightarrow T \; \bullet & reduce\ E \rightarrow T
\end{array}
\right]_{③}
$$

$$
\left[\; T \rightarrow \mathrm{x} \; \bullet \quad reduce\ T \rightarrow \mathrm{x} \;\right]_{④}
\qquad
\left[\; S \rightarrow E \; ; \; \bullet \; \$ \quad accept \;\right]_{⑤}
$$

$$
\left[\; E \rightarrow T + E \; \bullet \quad reduce\ E \rightarrow T + E \;\right]_{⑦}
$$

- There is an apparent conflict in state ③:

$$\begin{bmatrix} E \rightarrow T \bullet + E & \quad \textit{shift } 6 \\ E \rightarrow T \bullet & \quad \textit{reduce } E \rightarrow T \end{bmatrix}_{③}$$

- Shall we shift the + token or reduce via rule $E \rightarrow T$[8]?
- Clearly, if the situation is as follows

$$2+3+4 \bullet \; ;$$

we want to reduce $E \rightarrow T$, but in a case like this

$$2+3 \bullet +4 \, ;$$

we want to shift the +.

> **SLR Rule:** choose *reduce* if the next token is among the tokens that can **follow** the non-terminal (here: $E$) in the input.

[8]Remember, *reducing* does not consume any token.

- Resolving shift/reduce conflicts with the help of **follow**ing tokens is called **SLR Parsing** (Simple LR).

- The resulting SLR parsing table thus is

**CFG symbol**

|   | x | + | ; | $ | E | T |
|---|---|---|---|---|---|---|
| 1 | s4 | | | | 2 | 3 |
| 2 | | | s5 | | | |
| 3 | | s6 | r $E \rightarrow T$ | | | |
| 4 | r $T \rightarrow$ x | r $T \rightarrow$ x | r $T \rightarrow$ x | r $T \rightarrow$ x | | |
| 5 | | | | a | | |
| 6 | s4 | | | | 7 | 3 |
| 7 | r $E \rightarrow T{+}E$ | r $E \rightarrow T{+}E$ | r $E \rightarrow T{+}E$ | r $E \rightarrow T{+}E$ | | |

- To resolve such conflicts we need information about the *set of tokens which may possibly follow* some given non-terminal $S$. We call this set *FOLLOW(S)*.

  **Note:** in the example above we have $FOLLOW(E) = \{\,;\,\}$.

- To compute $FOLLOW(S)$, we will additionally need information about the *set of tokens which may occur first* when we derive a non-terminal $S$. We call this set $FIRST(S)$[9].

  - Consider the CFG

$$Z \rightarrow z \qquad\qquad Y \rightarrow \qquad\qquad X \rightarrow Y$$
$$Z \rightarrow X\ Y\ Z \qquad Y \rightarrow y \qquad X \rightarrow x$$

  - $FIRST(Y) = \{y\}$.
  - $FIRST(X) = \{\quad ?\quad \}$.
  - To compute $FIRST(Z)$, is it sufficient to compute $FIRST(X)$ and then simply say $FIRST(Z) = FIRST(X)$?

  - Similarly, to compute $FOLLOW(X)$, can we simply say $FOLLOW(X) = FIRST(Y)$?

---

[9]Later, we will need $FIRST(S)$ anyway to build the more sophisticated LR(1) parsing tables.

- Compute *FIRST*(), *FOLLOW*(), *nullable*()[10] for a given CFG *G*.

$FIRST(S) \leftarrow \varnothing$, $FOLLOW(S) \leftarrow \varnothing$    *for all symbols S in G;*
$nullable(S) \leftarrow false$    *for all symbols S in G;*
$FIRST(\mathtt{x}) \leftarrow \{\mathtt{x}\}$    *for all terminals* $\mathtt{x}$ *in G;*

```
repeat
    foreach X → Y₁ ⋯ Yₖ do
        if k = 0 or nullable(Y₁) ∧ ⋯ ∧ nullable(Yₖ) then
            nullable(X) ← true;
```
$FIRST(X) \leftarrow FIRST(X) \cup FIRST(Y_1);$
$FOLLOW(Y_k) \leftarrow FOLLOW(Y_k) \cup FOLLOW(X);$
```
        foreach i ← 2 … k do
```
$FOLLOW(Y_{i-1}) \leftarrow FOLLOW(Y_{i-1}) \cup FIRST(Y_i);$
/* $X \to \varepsilon \cdots \varepsilon Y_i \cdots$ */
```
            if nullable(Y₁) ∧ ⋯ ∧ nullable(Yᵢ₋₁) then
```
$FIRST(X) \leftarrow FIRST(X) \cup FIRST(Y_i);$

/* $X \to \cdots Y_{i-1}\varepsilon \cdots \varepsilon$ */
```
            if nullable(Yᵢ) ∧ ⋯ ∧ nullable(Yₖ) then
```
$FOLLOW(Y_{i-1}) \leftarrow FOLLOW(Y_{i-1}) \cup FOLLOW(X);$

/* $X \to \cdots Y_{i-1}\varepsilon \cdots \varepsilon Y_j \cdots$ */
```
            foreach j ← i + 1 … k do
                if nullable(Yᵢ) ∧ ⋯ ∧ nullable(Yⱼ₋₁) then
```
$FOLLOW(Y_{i-1}) \leftarrow FOLLOW(Y_{i-1}) \cup FIRST(Y_j);$

```
until FIRST(), FOLLOW(), nullable() stable;
```

---

[10] *nullable*(S) = *true* if $\varepsilon$ is derivable from *S*.

- **Example:**

$$Z \rightarrow z \qquad\qquad Y \rightarrow \qquad\qquad X \rightarrow Y$$
$$Z \rightarrow X\ Y\ Z \qquad\quad Y \rightarrow y \qquad\quad X \rightarrow x$$

we start the algorithm with the initial assumption

|  | nullable() | FIRST() | FOLLOW() |
|---|:---:|:---:|:---:|
| $X$ | f | | |
| $Y$ | f | | |
| $Z$ | f | | |

In the first iteration we find

|  | nullable() | FIRST() | FOLLOW() |
|---|:---:|:---:|:---:|
| $X$ | t | x, y | |
| $Y$ | t | y | z |
| $Z$ | f | z | |

and the third iteration already gives the stable result

|  | nullable() | FIRST() | FOLLOW() |
|---|:---:|:---:|:---:|
| $X$ | t | x, y | x, y, z |
| $Y$ | t | y | x, y, z |
| $Z$ | f | x, y, z | |

- **Example:** for the expression CFG $G$

$$S \rightarrow E \ \$ \qquad T \rightarrow T * F$$
$$E \rightarrow E + T \qquad T \rightarrow F$$
$$E \rightarrow T \qquad F \rightarrow (\ E\ )$$
$$F \rightarrow \texttt{num}$$

we get

|   | nullable() | FIRST() | FOLLOW() |
|---|-----------|---------|----------|
| S | f | (, num | |
| E | f | (, num | +, ), $ |
| T | f | (, num | *, +, ), $ |
| F | f | (, num | *, +, ), $ |

which we use to resolve the shift/reduce conflicts in two LR states associated with $G$:

$$\left[\begin{array}{ll} E \rightarrow T \bullet & reduce\ E \rightarrow T \\ E \rightarrow T \bullet * F & shift\ 9 \end{array}\right]_{③} \qquad \left[\begin{array}{ll} E \rightarrow E + T \bullet & reduce\ E \rightarrow E + T \\ E \rightarrow T \qquad \bullet * F & shift\ 9 \end{array}\right]_{⑪}$$

Since $* \notin FOLLOW(E)$ we reduce for any token $\neq *$, but shift for $*$.

- The SLR parsing table algorithm already leads to useful LR parsers for quite a wide range of input CFGs.

- SLR however is not sufficient to parse typical constructs to be found in even moderately complex programming language grammars.

  – Consider the CFG $G$

$$
\begin{array}{ll}
{}_0\; S' \rightarrow S\; \$ & {}_3\; E \rightarrow V \\
{}_1\; S \rightarrow V = E & {}_4\; V \rightarrow \texttt{id} \\
{}_2\; S \rightarrow E & {}_5\; V \rightarrow *\, E
\end{array}
$$

  – $G$ captures the essence of C-style variable assignments in which a dereference operator $*$ may be applied to (pointer-typed) variables. $G$ accepts input token sequences like

$$
\texttt{x = y} \qquad\qquad \texttt{*x = y} \qquad\qquad \texttt{x = **y}
$$

– To construct $G$'s SLR parsing table, we compute LR items (only states ① ... ③ shown here):

$$
\begin{bmatrix}
S' \to & \bullet\ S\ \$ & goto\ 2 \\
\hline
S \to & \bullet\ V = E & goto\ 3 \\
S \to & \bullet\ E & goto\ 4 \\
V \to & \bullet\ \texttt{id} & shift\ 5 \\
V \to & \bullet * E & shift\ 6 \\
E \to & \bullet\ V & goto\ 3
\end{bmatrix}_{①}
\qquad
\begin{bmatrix} S' \to S \bullet \$ & accept \end{bmatrix}_{②}
$$

$$
\begin{bmatrix}
S \to V \bullet = E & shift\ 7 \\
E \to V \bullet & reduce\ E \to V
\end{bmatrix}_{③}
$$

and the associated $FOLLOW()$ sets to resolve any shift/reduce conflicts which might occur:

|   | nullable() | FIRST() | FOLLOW() |
|---|---|---|---|
| $S$ | f | x, * | \$ |
| $E$ | f | x, * | =, \$ |
| $V$ | f | x, * | =, \$ |

– According to the SLR rule, we resolve the shift/reduce conflict in state ③ by choosing $reduce\ E \to V$ when we see tokens in $FOLLOW(E) = \{=, \$\}$ in the input.

- The resulting SLR parsing table thus is

**CFG symbol**

| | id | * | = | $ | S | E | V |
|---|---|---|---|---|---|---|---|
| 1 | $s5$ | $s6$ | | | 2 | 4 | 3 |
| 2 | | | | $a$ | | | |
| 3 | $r\,E \to V$ | $r\,E \to V$ | $r\,E \to V$ | $r\,E \to V$ | | | |
| 4 | $r\,S \to E$ | $r\,S \to E$ | $r\,S \to E$ | $r\,S \to E$ | | | |
| 5 | $r\,V \to \mathtt{id}$ | $r\,V \to \mathtt{id}$ | $r\,V \to \mathtt{id}$ | $r\,V \to \mathtt{id}$ | | | |
| 6 | | | | | | 8 | 3 |
| 7 | | | | | | 9 | 3 |
| 8 | $r\,V \to * E$ | $r\,V \to * E$ | $r\,V \to * E$ | $r\,V \to * E$ | | | |
| 9 | $r\,S \to V = E$ | $r\,S \to V = E$ | $r\,S \to V = E$ | $r\,S \to V = E$ | | | |

– When we feed the input `id = *id` into the resulting LR parser, we observe the following actions:

| Stack | Input | Action |
|---|---|---|
| 1 | id = * id $ | *shift* |
| 1 id 5 | = * id $ | *reduce $V \to$ `id`* |
| 1 V 3 | = * id $ | *reduce $E \to V$* |
| 1 E 4 | = * id $ | *reduce $S \to E$* |
| 1 S 2 | = * id $ | *error* |

**CFG symbol**

| | id | * | = | $ | S | E | V |
|---|---|---|---|---|---|---|---|
| 1 | s5 | s6 | | | 2 | 4 | 3 |
| 2 | | | | a | | | |
| 3 | r $E \to V$ | r $E \to V$ | r $E \to V$ | r $E \to V$ | | | |
| 4 | r $S \to E$ | r $S \to E$ | r $S \to E$ | r $S \to E$ | | | |
| 5 | r $V \to$ id | r $V \to$ id | r $V \to$ id | r $V \to$ id | | | |
| 6 | | | | | | 8 | 3 |
| 7 | | | | | | 9 | 3 |
| 8 | r $V \to *E$ | r $V \to *E$ | r $V \to *E$ | r $V \to *E$ | | | |
| 9 | r $S \to V = E$ | r $S \to V = E$ | r $S \to V = E$ | r $S \to V = E$ | | | |

- The problem here is that

  ① it may make sense to reduce via $E \rightarrow V$ when we see a = token in the input (since
     = $\in$ *FOLLOW*($E$)), but

  ② this might not be the correct choice in *all* situations.


- SLR is completely indifferent to differing parsing situations (i.e., different lookahead
  tokens) since it bases its shift/reduce conflict resolution on the global *FOLLOW*() sets.


- We can improve LR parser table generation considerably if we consider the next
  **lookahead token** while we are computing LR states.

  - The resulting LR parsing tables are called **LR(1)** (*left-to-right parse, right derivation,
    1 lookahead token*).
  - The resulting LR(1) parsing tables contain more states than SLR tables. LR(1) can
    generate parsers for a strictly larger class of CFGs than SLR.

- An **LR(1) item** takes the form

$$(S \to S_1 \cdots S_{i-1} \bullet S_i \; S_{i+1} \cdots S_k, x)$$

- $x$ is the **lookahead token** (the token following the tokens derived from $S_k$). In this configuration, the parser expects to see input derivable from $S_i \cdots S_k \; x$.

- When we compute the closure of an LR(1) item, we add an item for every possible lookahead token.

- In the configuration above, the possible lookahead tokens are the elements of[11]

$$FIRST(S_{i+1} \cdots S_k \; x)$$

---

[11]$FIRST(S_1 \; S_2) = \begin{cases} FIRST(S_1) \cup FIRST(S_2) & \text{if } nullable(S_1) \\ FIRST(S_1) & \text{otherwise.} \end{cases}$

- **Example:** compute closure of LR(1) item $(S' \to \bullet\, S\, \$, ?)$:[12]

$$\begin{bmatrix} \dfrac{(S' \to\ \bullet\ S\, \$ \qquad, ?)}{\begin{array}{llll} S & \to\ \bullet\ V = E & , \$) \\ S & \to\ \bullet\ E & , \$) \\ E & \to\ \bullet\ V & , \$) \\ V & \to\ \bullet\ \mathtt{id} & , \$) \\ V & \to\ \bullet\ *\, E & , \$) \\ V & \to\ \bullet\ \mathtt{id} & , =) \\ V & \to\ \bullet\ *\, E & , =) \end{array}} \end{bmatrix}_{①}$$

  - **Note:** $FIRST(\$\ ?) = \{\$\}$, $FIRST(= E\, \$) = \{=\}$.
  - You can find the set of LR(1) items for this grammar on 🐯 page 66.
    (Typo: replace $T \to \bullet\, *\, E$ by $V \to \bullet\, *\, E$ in state ⑥.)
  - When we determine the LR(1) parsing table entries, we set

$$goto[s][\mathrm{x}] = r\ S \to S_1\ \cdots\ S_k$$

    for the LR(1) item $(S \to S_1\ \cdots\ S_k\ \bullet, \mathrm{x})$ in state $s$. Otherwise, table generation is
    like in the SLR case.

---

[12]For the start symbol $S'$, the lookahead token is immaterial since we never shift (*look beyond*) \$.

- In the LR(1) case, the former problematic state ③ now looks like

$$\left[ \begin{array}{lll} (S \to V \bullet = E & , \$) & shift\ 4 \\ (E \to V \bullet & , \$) & reduce\ E \to V \end{array} \right]_{③}$$

- With the LR(1) parsing table (also see 🐯 page 66), we now observe the following LR parsing trace for the input `id = *id`:

| Stack | Input | Action |
|---|---|---|
| $_1$ | id = * id $ | shift |
| $_1$ id $_8$ | = * id $ | reduce $V \to$ id |
| $_1 V_3$ | = * id $ | shift |
| $_1 V_3 =_4$ | * id $ | shift |
| $_1 V_3 =_4 *_{13}$ | id $ | shift |
| $_1 V_3 =_4 *_{13}$ id $_{11}$ | $ | reduce $V \to$ id |
| $_1 V_3 =_4 *_{13} V_7$ | $ | reduce $E \to V$ |
| $_1 V_3 =_4 *_{13} E_{14}$ | $ | reduce $V \to * E$ |
| $_1 V_3 =_4 V_7$ | $ | reduce $E \to V$ |
| $_1 V_3 =_4 E_9$ | $ | reduce $S \to V = E$ |
| $_1 S_2$ | $ | accept |

# 3.4    Parser Generators

- The construction of an LR(1) (or the more compact LALR(1) variant) parsing table from a given CFG is tedious but simple enough to be done automatically.

- **Parser generators** like `yacc`[13], or the GNU project's `bison`

  ① read the description of a CFG $G$,
  ② internally construct the LALR(1) parsing table for $G$, and then
  ③ output a table-driven C program (with entry point `yyparse()`) which implements an LALR(1) parser for $G$, see slide 89.

- How can `yyparse()` return anything useful but `Parse OK` or `Parse Error`?

  – By attaching **semantic actions** to each CFG rule which are executed whenever a *reduce* action for that rule is executed (see next chapter).

---

[13]Yet Another Compiler Compiler

- GNU `bison` parser skeleton:

---
`bison` **parser skeleton**

| | |
|---|---|
| 1 | *parser declarations* |
| 2 | *%%* |
| 3 | *CFG rules and semantic actions* |
| 4 | *%%* |
| 5 | *program fragments, helper functions, ...* |

---

- **Example:** the expression CFG

$$
\begin{array}{lll}
E \ \rightarrow \ \mathrm{id} & E \ \rightarrow \ E * E & E \ \rightarrow \ E + E \\
E \ \rightarrow \ \mathrm{num} & E \ \rightarrow \ E / E & E \ \rightarrow \ E - E \\
E \ \rightarrow \ ( \ E \ ) & &
\end{array}
$$

may bewritten as shown on the next slide to be able to feed it into `bison`:

```
                    ──── bison expression CFG file expr.y ────
    1  %token ID NUM LPAREN RPAREN TIMES DIV PLUS MINUS

    2

    3  %%

    4

    5  E : ID

    6    | NUM

    7    | LPAREN E RPAREN

    8    | E TIMES E

    9    | E DIV E

   10    | E PLUS E

   11    | E MINUS E

   12    ;

   13

   14  %%
```

– `bison` can write a C header file containing definitions[14] for the declared tokens (`%token`) of the CFG.

– A lexer can include this header file and then return properly encoded tokens to the parser.

---

[14]The written C header file includes a series of #define ID 258, #define NUM 259, ... or a C enum like enum yytokentype { ID = 258, NUM = 259, ... }

- Now let us feed `expr.y` into `bison`:

```
───────────── bison invocation ─────────────
1  $ bison -d expr.y
2  expr.y: warning: 16 shift/reduce conflicts
3  $ ▮
```

- Remember from slide 82 that the expression grammar was highly ambiguous:

  – We neither incorporated knowledge about **operator precedence**
    (∗ binds tighter than +, . . . ),

  – nor made **operator associativity** explicit (arith. operators associate to the left).

- ⚠ LR(k) table generation fails for any ambiguous CFG: some parser states will have conflicting associated actions.

- We can use `bison` itself to analyze the conflicts:

```
───────────── bison conflict analysis ─────────────
1  $ bison --verbose expr.y
2  expr.y: warning: 16 shift/reduce conflicts
3  $ more expr.output
```

```
1   State 12 contains 4 shift/reduce conflicts.
2   State 13 contains 4 shift/reduce conflicts.
3   State 14 contains 4 shift/reduce conflicts.
4   State 15 contains 4 shift/reduce conflicts.
5
6
7   Grammar
8
9       0 $accept: E $end
10
11      1 E: ID
12      2  | NUM
13      3  | LPAREN E RPAREN
14      4  | E TIMES E
15      5  | E DIV E
16      6  | E PLUS E
17      7  | E MINUS E
18
19  [...]
```

```
 1  state 0
 2      0 $accept: . E $end
 3
 4      ID      shift, and go to state 1
 5      NUM     shift, and go to state 2
 6      LPAREN  shift, and go to state 3
 7
 8      E  go to state 4
 9
10  state 1
11      1 E: ID .
12
13      $default  reduce using rule 1 (E)
14
15  state 2
16      2 E: NUM .
17
18      $default  reduce using rule 2 (E)
19
20  state 3
21      3 E: LPAREN . E RPAREN
22
23      ID      shift, and go to state 1
24      NUM     shift, and go to state 2
25      LPAREN  shift, and go to state 3
26
27      E  go to state 5
```

```
             expr.output continued                       expr.output continued
 1   state 12                               20   state 13
 2                                          21
 3      4 E: E . TIMES E                    22      4 E: E . TIMES E
 4      4  | E TIMES E .                    23      5  | E . DIV E
 5      5  | E . DIV E                      24      5  | E DIV E .
 6      6  | E . PLUS E                     25      6  | E . PLUS E
 7      7  | E . MINUS E                    26      7  | E . MINUS E
 8                                          27
 9      TIMES  shift, and go to state 7     28      TIMES  shift, and go to state 7
10      DIV    shift, and go to state 8     29      DIV    shift, and go to state 8
11      PLUS   shift, and go to state 9     30      PLUS   shift, and go to state 9
12      MINUS  shift, and go to state 10    31      MINUS  shift, and go to state 10
13                                          32
14      TIMES     [reduce using rule 4 (E)] 33      TIMES     [reduce using rule 5 (E)]
15      DIV       [reduce using rule 4 (E)] 34      DIV       [reduce using rule 5 (E)]
16      PLUS      [reduce using rule 4 (E)] 35      PLUS      [reduce using rule 5 (E)]
17      MINUS     [reduce using rule 4 (E)] 36      MINUS     [reduce using rule 5 (E)]
18      $default  reduce using rule 4 (E)   37      $default  reduce using rule 5 (E)
19                                          38
```

- We can clearly see the *shift/reduce* conflicts and that `bison` has decided to favor the *shift*. This is the `bison` default conflict resolution for ambiguous grammars.

  – Are we happy with `bison`'s default decision in this case?

– Trace the `bison` generated parser on input `num * num + num`:

| Stack | Input | Action |
|---|---:|---|
| $_0$ | num * num + num $ | _shift_ |
| $_0$ num $_2$ | * num + num $ | _r_ $E \to$ num |
| $_0$ $E$ $_4$ | * num + num $ | _shift_ |
| $_0$ $E$ $_4$ * $_7$ | num + num $ | _shift_ |
| $_0$ $E$ $_4$ * $_7$ num $_2$ | + num $ | _r_ $E \to$ num |
| $_0$ $E$ $_4$ * $_7$ $E$ $_{12}$ | + num $ | _shift_    ↯ |
| $_0$ $E$ $_4$ * $_7$ $E$ $_{12}$ + $_9$ | num $ | _shift_ |
| $_0$ $E$ $_4$ * $_7$ $E$ $_{12}$ + $_9$ num $_2$ | $ | _r_ $E \to$ num |
| $_0$ $E$ $_4$ * $_7$ $E$ $_{12}$ + $_9$ $E$ $_{14}$ | $ | _r_ $E \to E + E$ |

$$
\begin{bmatrix}
(E \to E & \bullet * E & ,?) \\
(E \to E & \bullet / E & ,?) \\
(E \to E & \bullet + E & ,?) \\
(E \to E & \bullet - E & ,?) \\
(E \to E * E \bullet & & ,\$) \\
(E \to E * E \bullet & & ,)) \\
(E \to E * E \bullet & & ,*) \\
(E \to E * E \bullet & & ,/) \\
(E \to E * E \bullet & & ,+) \\
(E \to E * E \bullet & & ,-)
\end{bmatrix}_{\circled{12}}
$$

- For our example input of `num * num + num` the following two LR items of state 12 comprise the *shift/reduce* dilemma:

$$
\begin{bmatrix}
\vdots \\
(E \rightarrow E \quad \bullet + E \quad , ?) \\
(E \rightarrow E * E \bullet \qquad , +) \\
\vdots
\end{bmatrix}_{\,\text{\textcircled{12}}}
$$

  - Here we want to choose the *reduce* because this interprets $*$ as an operator with higher precedence as $+$.
  - So, to sensibly resolve such ambiguity, for an LR item of the form

$$
(E \rightarrow E \bullet \odot E, ?) \\
(E \rightarrow E \odot E \bullet, \otimes)
$$

  choose

  ① *reduce*, if the precedence of $\odot$ is higher than $\otimes$,
  ② *shift*, if the precedence is lower.

- For an LR item of the form

$$(E \to E \bullet \odot E, ?)$$
$$(E \to E \odot E \bullet, \odot)$$

choose

① [          ], if operator $\odot$ is left-associative,

② [          ], if operator $\odot$ is right-associative,

③ [          ], if operator $\odot$ is not associative at all.

- In `bison`, the declarations `%left`, `%right`, and `%nonassoc` assign **associativity** to the specified token.
  The order in which these declarations appear in the `bison` input file specifies **precedence** (lower precedence comes first).

- The precedence of a rule is the precedence of the **rightmost** token appearing in it[15].

---

[15]This may be overridden by `bison`'s `%prec` directive. See 🐯 **p. 74** and the `bison` manual.

- The resulting modified `bison` input file:

```
                    bison expression CFG file expr.y, modified
 1  %token ID NUM LPAREN RPAREN TIMES DIV PLUS MINUS
 2
 3  %left PLUS MINUS
 4  %left TIMES DIV
 5
 6  %%
 7
 8  E : ID
 9    | NUM
10    | LPAREN E RPAREN
11    | E TIMES E
12    | E DIV E
13    | E PLUS E
14    | E MINUS E
15    ;
16
17  %%
```

- `bison` generates an LALR(1) parsing table for this grammar without any ambiguity problems.

# 3.5 Syntax vs. Semantics

- There still remain decisions to be made during program analysis for which *no* parsing technique whatsoever will be sufficiently clever.

  - **Example:** in a small language, we have *boolean* (operators and, or, =) and *arithmetic* expressions (operators +, *). Arithmetic operators bind tighter than boolean operators, such that

  $$x := a + b = c$$

  is parsed as x := (a + b) = c.
  - The bison grammar tries to separate boolean from arithmetic expressions as shown on the next slide. In some sense, the CFG tries to implement a bit of **type checking** in order to prevent acceptance of input like

  $$x := a + b \text{ and } c$$

```
 1  %token ID ASSIGN PLUS TIMES AND OR EQUAL
 2
 3  %left OR
 4  %left AND
 5  %left EQUAL
 6  %left PLUS
 7  %left TIMES
 8
 9  %%
10
11  stm    : ID ASSIGN ae
12         | ID ASSIGN be
13         ;
14
15  ae     : ae PLUS ae
16         | ae TIMES ae
17         | ID
18         ;
19
20  be     : be OR be
21         | be AND be
22         | be EQUAL be
23         | ae EQUAL ae
24         | ID
25         ;
```

– bison reports a *reduce/reduce* conflict for this CFG. Where/Why?

– For input such as

$$x := a$$

there is no way to tell for the parser if the variable `a` is an arithmetic or a boolean variable (reduce `ae` → `ID` or reduce `be` → `ID`?).

- If there is **no syntactical criterion** to support a parsing decision, then either

  ① the grammar is severely malformed (fix the CFG!), or
  ② the decision has to be **postponed until semantical analysis**.

- Remark on ②: make the CFG more permissive, then reject type incorrect input such as

$$x := a + b \ and \ c$$

  during **semantical analysis**.

```
1   %token ID ASSIGN PLUS TIMES AND OR EQUAL
2
3   %left OR
4   %left AND
5   %left EQUAL
6   %left PLUS
7   %left TIMES
8
9   %%
10
11  stm    : ID ASSIGN e
12         ;
13
14  e      : e PLUS e
15         | e TIMES e
16         | e AND e
17         | e OR e
18         | e EQUAL e
19         | ID
20         ;
21
22  %%
```