# Syntax Analyzer --- Parser
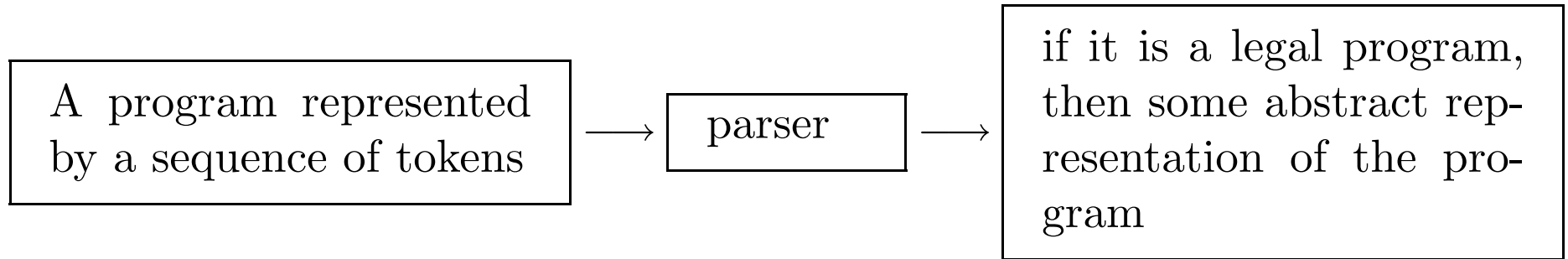
## ASU Textbook Chapter 4.2--4.9 (w/o error handling)

Tsan-sheng Hsu

*tshsu@iis.sinica.edu.tw*

`http://www.iis.sinica.edu.tw/~tshsu`

# main tasks

| A program represented by a sequence of tokens | $\longrightarrow$ | parser | $\longrightarrow$ | if it is a legal program, then some abstract representation of the program |
|---|---|---|---|---|

- **Abstract representations of the input program:**
  - **abstract-syntax tree + symbol table**
  - **intermediate code**
  - **object code**
- **Context free grammar (CFG) is used to specify the structure of legal programs.**

# Context free grammar (CFG)

- **Definitions:** $G = (T, N, P, S)$, **where**

  - $T$: **a set of** terminals **(in lower case letters);**

  - $N$: **a set of** nonterminals **(in upper case letters);**

  - $P$: productions **of the form**
    $A \rightarrow X_1, X_2, \ldots, X_m$, **where** $A \in N$ **and** $X_i \in T \cup N$;
  - $S$: **the starting nonterminal,** $S \in N$.

- **Notations:**
  - **terminals : lower case English strings, e.g.,** $a$, $b$, $c \cdots$
  - **nonterminals: upper case English strings, e.g.,** $A$, $B$, $C \cdots$
  - $\alpha, \beta, \gamma \in (T \cup N)^*$
    - $\triangleright$ $\alpha$, $\beta$, $\gamma$**: alpha, beta and gamma.**
    - $\triangleright$ $\epsilon$**: epsilon.**

  -

$$
\left.
\begin{array}{ccc}
A & \rightarrow & X_1 \\
A & \rightarrow & X_2
\end{array}
\right\} \equiv A \rightarrow X_1 \mid X_2
$$

# How does a CFG define a language?

- **The language defined by the grammar is the set of strings (sequence of terminals) that can be "derived" from the starting nonterminal.**

- **How to "derive" something?**
  - **Start with:**
    **"current sequence" = the starting nonterminal.**
  - **Repeat**
    - ▷ *find a nonterminal $X$ in the current sequence*
    - ▷ *find a production in the grammar with $X$ on the left of the form $X \rightarrow \alpha$, where $\alpha$ is $\epsilon$ or a sequence of terminals and/or nonterminals.*
    - ▷ *create a new "current sequence" in which $\alpha$ replaces $X$*
  - **Until "current sequence" contains no nonterminals.**

- **We derive either $\epsilon$ or a string of terminals. This is how we derive a string of the language.**

# Example

Grammar:

- $E \rightarrow int$
- $E \rightarrow E - E$
- $E \rightarrow E \,/\, E$
- $E \rightarrow (\ E\ )$

$$E$$
$$\Longrightarrow E - E$$
$$\Longrightarrow 1 - E$$
$$\Longrightarrow 1 - E/E$$
$$\Longrightarrow 1 - E/2$$
$$\Longrightarrow 1 - 4/2$$

- **Details:**
  - **The first step was done by choosing the 2nd of the 4 productions.**
  - **The second step was by choosing the first production.**
- **Conventions:**
  - $\Longrightarrow$**: means "derives in one step";**
  - $\overset{+}{\Longrightarrow}$**: means "derives in one or more steps";**
  - $\overset{*}{\Longrightarrow}$**: means "derives in zero or more steps";**
  - **In the above example, we can write** $E \overset{+}{\Longrightarrow} 1 - 4/2$**.**

# Language

- **The  language  defined by a grammar $G$ is**

$$L(G) = \{w \mid S \overset{+}{\Longrightarrow} \omega\},$$

  **where $S$ is the starting nonterminal and $\omega$ is a sequence of terminals or $\epsilon$.**
- **An  element  in a language is $\epsilon$ or a sequence of terminals in the set defined by the language.**
- **More terminology:**
    - $E \Longrightarrow \cdots \Longrightarrow 1 - 4/2$ **is a  derivation  of $1 - 4/2$ from $E$.**
    - **There are several kinds of derivations that are important:**
        - ▷ *The derivation is a  leftmost  one if the leftmost nonterminal always gets to be chosen (if we have a choice) to be replaced.*
        - ▷ *It is a  rightmost  one if the rightmost nonterminal is replaced all the times.*

# A different way to derive

- **Construct a** **derivation** **or** **parse tree** **as follows:**
  - start with the starting nonterminal as a single-node tree
  - REPEAT
    - ▷ *choose a leaf nonterminal $X$*
    - ▷ *choose a production $X \rightarrow \alpha$*
    - ▷ *symbols in $\alpha$ become children of $X$*
  - UNTIL no more leaf nonterminal left
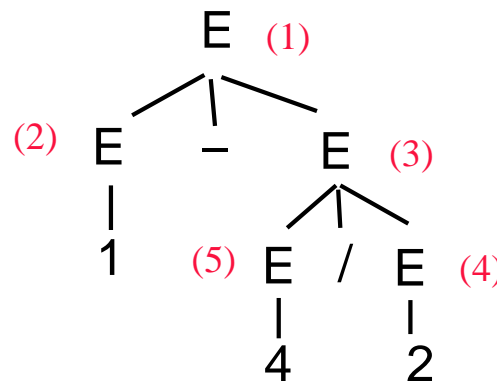- **Need to annotate the order of derivation on the nodes.**

$E$

$\Longrightarrow E - E$

$\Longrightarrow 1 - E$

$\Longrightarrow 1 - E/E$

$\Longrightarrow 1 - E/2$

$\Longrightarrow 1 - 4/2$

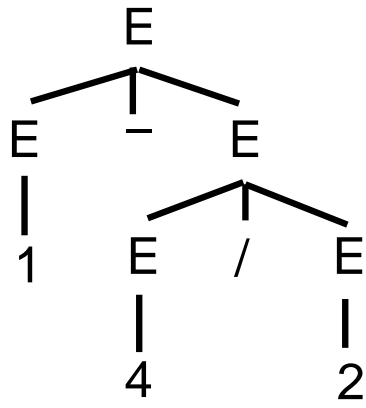# Parse tree examples

- **Example:**

**Grammar:**

$$E \rightarrow int$$

$$E \rightarrow E - E$$

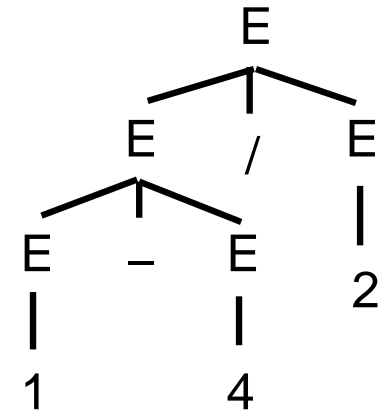$$E \rightarrow E/E$$

$$E \rightarrow (E)$$



leftmost derivation

- **Using $1 - 4/2$ as the input, the left parse tree is derived.**

- **A string is formed by reading the lead nodes from left to right, given $1 - 4/2$.**

- **The string $1 - 4/2$ has another parse tree on the right.**



rightmost derivation

- **Some standard notations:**
  - **Given a parse tree and a fixed order (for example leftmost or rightmost) we can derive the order of derivation.**
  - **For the "semantic" of the parse tree, we normally "interpret" the meaning in a bottom-up fashion. That is, the one that is derived last will be "serviced" first.**

# Ambiguous Grammar

- **If for grammar $G$ and string $S$, there are**
  - more than one leftmost derivation for $S$, or
  - more than one rightmost derivation for $S$, or
  - more than one parse tree for $S$,

  **then $G$ is called** ambiguous **.**
  - Note: the above three conditions are equivalent in that if one is true, then all three are true.

- **Problems with an ambiguous grammar:**
  - Ambiguity can make parsing difficult.
  - Underlying structure is ill-defined: in the example, the precedence is not uniquely defined, e.g., the leftmost parse tree groups $4/2$ while the rightmost parse tree groups $1-4$, resulting in two different semantics.

# Grammar that expresses precedence correctly

- **Use one nonterminal for each precedence level**
- **Start with lower precedence (in our example $-$)**

**Original grammar:**

$E \rightarrow int$

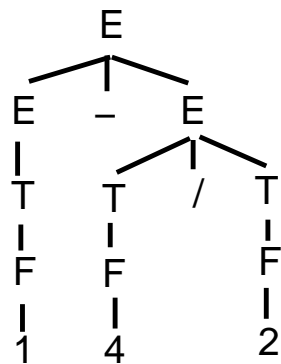$E \rightarrow E - E$

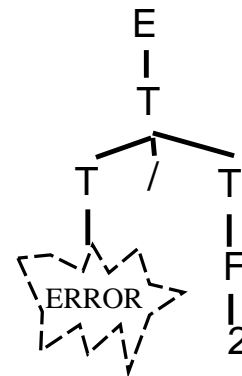$E \rightarrow E/E$

$E \rightarrow (E)$

**Revised grammar:**

$E \rightarrow E - E \mid T$

$T \rightarrow T/T \mid F$

$F \rightarrow int \mid (E)$

```
        E                          E
      / | \                        |
    E   –   E                      T
    |      /|\                    /|\
    T     T / T                  T / T
    |     |   |                  |   |
    F     F   F                  F   F
    |     |   |                ERROR |
    1     4   2                      2

   rightmost derivation
```

# More problems with associativity

- **However, the above grammar is still ambiguous, and parse trees may not express the associative of $-$ and $/$.**
  **Example:** $2-3-4$

**Revised grammar:**

$E \rightarrow E - E \mid T$

$T \rightarrow T/T \mid F$

$F \rightarrow int \mid (E)$

rightmost derivation
value = (2–3)–4 = –5

rightmost derivation
value = 2 – (3–4) = 3

- **Problems with associativity:**
  - **The rule $E \rightarrow E - E$ has $E$ on both sides of "$-$".**
  - **Need to make the second $E$ to some other nonterminal parsed earlier.**
  - **Similarly for the rule $E \rightarrow E/E$.**

# Grammar considering associative rules

Original grammar:

$E \rightarrow int$

$E \rightarrow E - E$

$E \rightarrow E/E$

$E \rightarrow (E)$

Revised grammar:

$E \rightarrow E - E \mid T$
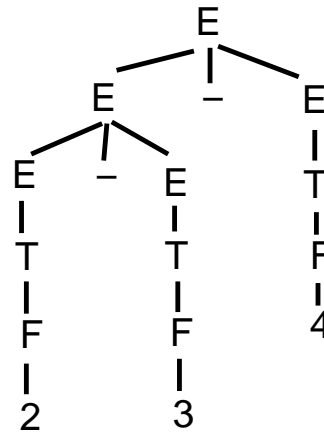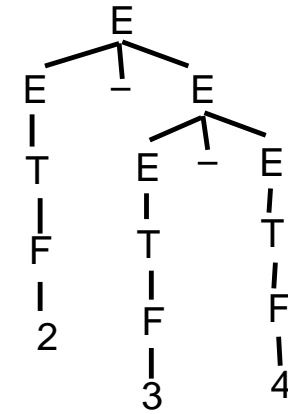
$T \rightarrow T/T \mid F$

$F \rightarrow int \mid (E)$

Final revised grammar:

$E \rightarrow E - T \mid T$

$T \rightarrow T/F \mid F$

$F \rightarrow int \mid (E)$

- **Recursive productions:**
  - $E \rightarrow E - T$ **is called a** left recursive **production.** $A \stackrel{+}{\Longrightarrow} A\alpha$.
  - $E \rightarrow T - E$ **is called a** right recursive **production.** $A \stackrel{+}{\Longrightarrow} \alpha A$.
  - $E \rightarrow E - E$ **is both left and right recursion.**
  - **If one wants left associativity, use left recursion.**
  - **If one wants right associativity, use right recursion.**

# Common grammar problems

- **Expressions: precedence and associativity as discussed above.**
- **Lists: that is, zero or more ID's separated by commas:**
  - **Note it is easy to express one or more ID's:**
    **<idlist>→<idlist>, ID | ID**
  - **For zero or more ID's,**
    - ▷ $<idlist> \rightarrow \epsilon \mid ID \mid <idlist>, <idlist>$
      *won't work due to $\epsilon$; it can generate: $ID, , ID$*
    - ▷ $<idlist> \rightarrow \epsilon \mid <idlist>, ID \mid ID$
      *won't work either because it can generate: $, ID, ID$*
  - **We should separate out the empty list from the general list of one or more ID's.**
    - ▷ $<opt\text{-}idlist> \rightarrow \epsilon \mid <nonEmptyIdlist>$
    - ▷ $<nonEmptyIdlist> \rightarrow <nonEmptyIdlist>, ID \mid ID$

# How to use CFG

- **Breaks down the problem into pieces:**
  - **Think about a C program:**
    - ▷ *Declarations: typedef, struct, variables, . . .*
    - ▷ *Procedures: type-specifier, function name, parameters, function body.*
    - ▷ *function body: various statements.*
  - **Example:**

**<procedure>→<type-def> ID <opt-params><opt-decl> {<opt-statements>}**

    - ▷ *<opt-params>→ (<list-params>)*
    - ▷ *<list-params>→ ε |<nonEmptyParlist>*
    - ▷ *<nonEmptyParlist>→<nonEmptyIdlist>, ID | ID*

- **One of purposes to write a grammar for a language is for others to understand. It will be nice to break things up into different levels in a top-down easily understandable fashion.**

# Useless terms

- **A non-terminal $X$ is useless if either**
  - a sequence includes $X$ cannot be derived from the starting nonterminal, or
  - no string can be derived starting from $X$, where a string means $\epsilon$ or a sequence of terminals.
- **Example 1:**
  - $S \to A\ B$
  - $A \to +\ |\ -\ |\ \epsilon$
  - $B \to digit\ |\ B\ digit$
  - $C \to .\ B$
- **In Example 1:**
  - $C$ is useless and so is the last production.
  - **Any nonterminal not in the right-hand side of any production is useless!**

# More examples for useless terms

- **Example 2: $Y$ is useless.**
  - $S \to X \mid Y$
  - $X \to (\ )$
  - $Y \to (\ Y\ Y\ )$
- $Y$ **derives more and more nonterminals and is useless.**
- **Any recursively defined nonterminal without a production of deriving $\epsilon$ all terminals is useless!**
  - Direct useless.
  - Indirect useless: one can only derive direct useless terms.
- **From now on, we assume a grammar contains no useless nonterminals.**

# Non-context free grammars

- **Some grammar is not CFG, that is, it may be context sensitive.**
- **Expressive power of grammars (in the order of small to large):**
  - **Regular expressions $\equiv$ FA.**
  - **Context-free grammar**
  - **Context-sensitive**
  - $\cdots$
- $\{\omega c\omega \mid \omega$ **is a string of** $a$ **and** $b$**'s**$\}$ **cannot be expressed by CFG.**

# Top-down parsing

- **There are $O(n^3)$-time algorithms to parse a language defined by CFG, where $n$ is the number of input tokens.**
- **For practical purpose, we need faster algorithms. Here we make restrictions to CFG so that we can design $O(n)$-time algorithms.**

- **Recursive-descent parsing** : top-down parsing that allows backtracking.**
  - **Attempt to find a leftmost derivation for an input string.**
  - **Try out all possibilities, that is, do an exhaustive search to find a parse tree that parses the input.**

# Example for recursive-descent parsing

$$S \to cAd$$

$$A \to bc \mid a$$

Input: cad



*error!! backtrack*

- **Problems with the above approach:**
  - **still too slow!**
  - **want to select a derivation without ever causing backtracking!**
  - **trick: use lookahead symbols.**
- **Solution: use $LL(1)$ grammars that can be parsed in $O(n)$ time.**

  - **first $L$: scan the input from left-to-right**
  - **second $L$: find a leftmost derivation**
  - **$(1)$: allow one lookahead token!**

# Predictive parser for $LL(1)$ grammars

- **How a predictive parser works:**
  - **start by pushing the starting nonterminal into the STACK and calling the scanner to get the first token.**
  
    **LOOP: if top-of-STACK is a nonterminal, then**
      - ▷ *use the current token and the PARSING TABLE to choose a production*
      - ▷ *pop the nonterminal from the STACK and push the above production's right-hand-side*
      - ▷ *GOTO LOOP.*
  - **if top-of-STACK is a terminal and matches the current token, then**
      - ▷ *pop STACK and ask scanner to provide the next token*
      - ▷ *GOTO LOOP.*
  - **if STACK is empty and there is no more input, then ACCEPT!**
  - **If none of the above succeed, then FAIL!**
      - ▷ *STACK is empty and there is input left.*
      - ▷ *top-of-STACK is a terminal, but does not match the current token*
      - ▷ *top-of-STACK is a nonterminal, but the corresponding PARSE TABLE entry is ERROR!*

# Example for parsing an $LL(1)$ grammar

- **grammar:** $S \rightarrow \epsilon \mid (S) \mid [S]$      **input:** $([\,])$

| input | stack | action |
|-------|-------|--------|
| (     | S     | pop, push "(S)" |
| (     | (S)   | pop, match with input |
| ([    | S)    | pop, push "[S]" |
| ([    | [S])  | pop, match with input |
| ([ ]  | S])   | pop, push $\epsilon$ |
| ([ ]  | ])    | pop, match with input |
| ([ ]) | )     | pop, match with input |
| ([ ]) |       | accept |

```
          S
        / | \
      (   S   )
        / | \
      [   S   ]
          |
          ε
```
leftmost derivation

- **Use the current input token to decide which production to derive from the top-of-STACK nonterminal.**

# About $LL(1)$

- **It is not always possible to build a predictive parser given a CFG; It works only if the CFG is $LL(1)$!**
- **For example, the following grammar is not $LL(1)$, but is $LL(2)$.**
- **Grammar: $S \rightarrow (S) \mid [S] \mid () \mid [\,]$**
  **Try to parse the input $()$.**

  | input | stack | action |
  |-------|-------|--------|
  | ( | S | pop, but use which production? |

- **In this example, we need 2-token look-ahead.**
  - **If the next token is $)$, push $()$.**
  - **If the next token is $($, push $(S)$.**
- **Two questions:**
  - **How to tell whether a grammar $G$ is $LL(1)$?**
  - **How to build the PARSING TABLE?**

# Properties of non-$LL(1)$ grammars

- **Theorem 1: A CFG grammar is not $LL(1)$ if it is left-recursive.**
- **Definitions:**
  - **recursive grammar: a grammar is recursive if the following is true for a nonterminal $X$ in $G$:**
    $X \stackrel{+}{\Longrightarrow} \alpha X \beta$.
  - $G$ **is** left-recursive **if** $X \stackrel{+}{\Longrightarrow} X\beta$.
  - $G$ **is** immediately left-recursive **if** $X \Longrightarrow X\beta$.

# Example of removing immediate left-recursion

- **Need to remove left-recursion to come out an $LL(1)$ grammar. Example:**
  - **Grammar $G$: $A \rightarrow A\alpha \mid \beta$, where $\beta$ does not start with $A$**
  - **Revised grammar $G'$:**
    - $\triangleright$ $A \rightarrow \beta A'$
    - $\triangleright$ $A' \rightarrow \alpha A' \mid \epsilon$
  - **The above two grammars are equivalent. That is $L(G) \equiv L(G')$.**

- **Example:**

  | **input** $baa$ |
  |:---:|
  | $\beta \equiv b$ |
  | $\alpha \equiv a$ |

  leftmost derivation
  original grammar G

  leftmost derivation
  revised grammar G'

# Rule for removing immediate left-recursion

- **Both grammar recognize the same string, but $G'$ is not left-recursive.**
- **However, $G$ is clear and intuitive.**
- **General rule for removing immediately left-recursion:**
  - **Replace $A \to A\alpha_1 \mid \cdots \mid A\alpha_m \mid \beta_1 \mid \cdots \beta_n$**
  - **with**
    - $\triangleright \ A \to \beta_1 A' \mid \cdots \mid \beta_n A'$
    - $\triangleright \ A' \to \alpha_1 A' \mid \cdots \mid \alpha_m A' \mid \epsilon$

# Algorithm 4.1

- **Algorithm 4.1 systematically eliminates left recursion from a grammar if it is possible to do so.**
  - **Algorithm 4.1 works only if the grammar has no cycles or $\epsilon$-productions.**
  - **It is possible to remove cycles and $\epsilon$-productions using other algorithms.**

---

- **Input: grammar $G$ without cycles and $\epsilon$-productions.**
- **Output: An equivalent grammar without left recursion.**
- **Number the nonterminals in some order $A_1, A_2, \ldots, A_n$**
- **for $i = 1$ to $n$ do**
  - **for $j = 1$ to $i - 1$ do**
    - ▷ *replace $A_i \rightarrow A_j \gamma$*
    - ▷ *with $A_i \rightarrow \delta_1 \gamma \mid \cdots \mid \delta_k \gamma$*
    - ▷ *where $A_j \rightarrow \delta_1 \mid \cdots \mid \delta_k$ are all the current $A_j$-productions.*
    - ▷ *Eliminate immediate left-recursion for $A_i$*

---

After each $i$-loop, only productions of the form $A_i \rightarrow A_k \gamma$, $i < k$ are left.

# Intuition for algorithm 4.1

- **Intuition:** if $A_{i_1} \xRightarrow{+} \alpha_2 A_{i_2} \beta_2 \xRightarrow{+} \alpha_3 A_{i_3} \beta_3 \xRightarrow{+} \cdots$ and $i_1 < i_2 < i_3 < \cdots$, then it is not possible to have recursion.
- **Trace Algorithm 4.1**
  - $i = 1$
    - ▷ *do nothing*
    - ▷ *allow $A_1 \to A_k \alpha$, $\forall k$ before removing immediate left-recursion*
    - ▷ *remove allow $A_1 \to A_1 \alpha$ after removing immediate left-recursion*
  - $i = 2$
    - ▷ *$j = 1$:*
      *replace $A_2 \to A_1 \gamma$ by $A_2 \to A_k \alpha \gamma$*
    - ▷ *$A_k \neq A_1$ since all immediate left-recursion of the form $A_1 \to A_1 \beta$ are removed.*
  - $i = 3$
    - ▷ *$j = 1$:*
      *remove $A_3 \to A_1 \delta_1$*
    - ▷ *$j = 2$:*
      *remove $A_3 \to A_2 \delta_2$*
  - $\cdots$

# Example

- **Original Grammar:**
  - **(1)** $S \rightarrow Aa \mid b$
  - **(2)** $A \rightarrow Ac \mid Sd \mid e$
- **Ordering of nonterminals:** $S \equiv A_1$ **and** $A \equiv A_2$.
- $i = 1$
  - **do nothing as there is no immediate left-recursion for** $S$
- $i = 2$
  - **replace** $A \rightarrow Sd$ **by** $A \rightarrow Aad \mid bd$
  - **hence (2) becomes** $A \rightarrow Ac \mid Aad \mid bd \mid e$
  - **after removing immediate left-recursion:**
    - ▷ $A \rightarrow bdA' \mid eA'$
    - ▷ $A' \rightarrow cA' \mid adA' \mid \epsilon$

# Second property for non-$LL(1)$ grammars

- **Theorem 2: $G$ is not $LL(1)$ if a nonterminal has two productions whose right-hand-sides have a common prefix.**
- **Example:**
  - $S \rightarrow (S) \mid ()$
- **In this example, the common prefix is "(".**

- **This problem can be solved by using the** <span style="background-color:#c0c0c0">left-factoring</span> **trick.**

  - $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$
  - **Transform to:**
    - $\triangleright$ $A \rightarrow \alpha A'$
    - $\triangleright$ $A' \rightarrow \beta_1 \mid \beta_2$

- **Example:**
  - $S \rightarrow (S) \mid ()$
  - **Transform to**
    - $\triangleright$ $S \rightarrow (S'$
    - $\triangleright$ $S' \rightarrow S) \mid )$

# Algorithm for left-factoring

- **Input: context free grammar $G$**
- **Output: equivalent left-factored context-free grammar $G'$**
- **for each nonterminal $A$ do**
  - **find the longest non-$\epsilon$ prefix $\alpha$ that is common to right-hand sides of two or more productions**
  - **replace**

  - $A \rightarrow \alpha\beta_1 \mid \cdots \mid \alpha\beta_n \mid \gamma_1 \mid \cdots \mid \gamma_m$ **with**
    - $\triangleright$ $A \rightarrow \alpha A' \mid \gamma_1 \mid \cdots \mid \gamma_m$
    - $\triangleright$ $A' \rightarrow \beta_1 \mid \cdots \mid \beta_n$

  - **repeat the above process until $A$ has no two productions with a common prefix.**

# Left-factoring and left-recursion removal

- **Original grammar:**
$S \rightarrow (S) \mid SS \mid ()$

- **To remove immediate left-recursion, we have**
  - $S \rightarrow (S)S' \mid ()S'$
  - $S' \rightarrow SS' \mid \epsilon$

- **To do left-factoring, we have**
  - $S \rightarrow (S''$
  - $S'' \rightarrow S)S' \mid )S'$
  - $S' \rightarrow SS' \mid \epsilon$

- **A grammar is not $LL(1)$ if it is**
  - left recursive or
  - not left-factored.

**However, grammars that are not left recursive and are left-factored may still not be $LL(1)$.**

# Definition of $LL(1)$ grammars

- **To see if a grammar is $LL(1)$, we need to compute its FIRST and FOLLOW sets, which are used to build its parsing table.**
- **FIRST sets:**
  - **Definition: let $\alpha$ be a sequence of terminals and/or nonterminals or $\epsilon$**
    - ▷ *$FIRST(\alpha)$ is the set of terminals that begin the strings derivable from $\alpha$*
    - ▷ *if $\alpha$ can derive $\epsilon$, then $\epsilon \in FIRST(\alpha)$*
  - **FIRST$(\alpha) = \{(t \mid t$ is a terminal and $\alpha \stackrel{*}{\Longrightarrow} t\beta)$ or( $t = \epsilon$ and $\alpha \stackrel{*}{\Longrightarrow} \epsilon)\}$**

# How to compute $\textbf{FIRST}(X)$?

- $X$ **is a terminal:**
  - $\textbf{FIRST}(X) = \{X\}$
- $X$ **is** $\epsilon$**:**
  - $\textbf{FIRST}(X) = \{\epsilon\}$
- $X$ **is a nonterminal: must check all productions with** $X$ **on the left-hand side.**
  **That is,** $X \rightarrow Y_1 Y_2 \cdots Y_k$

  - **put** $\textbf{FIRST}(Y_1) - \{\epsilon\}$ **into** $\textbf{FIRST}(X)$
  - **if** $\epsilon \in \textbf{FIRST}(Y_1)$**, then put**
    $\textbf{FIRST}(Y_2) - \{\epsilon\}$ **into** $\textbf{FIRST}(X)$
  - $\cdots$
  - **if** $\epsilon \in \textbf{FIRST}(Y_{k-1})$**, then put**
    $\textbf{FIRST}(Y_k) - \{\epsilon\}$ **into** $\textbf{FIRST}(X)$
  - **if** $\epsilon \in \textbf{FIRST}(Y_i)$ **for each** $1 \leq i \leq k$**, then put** $\epsilon$ **into** $\textbf{FIRST}(X)$

# Example for computing $\textbf{FIRST}(X)$

- **Start with computing FIRST for the last production and walk your way up.**

**Grammar**

$$E \to E'T$$

$$E' \to -TE' \mid \epsilon$$

$$T \to FT'$$

$$T' \to /\ FT' \mid \epsilon$$

$$F \to int \mid (E)$$

$\textbf{FIRST}(F) = \{int, (\}$

$\textbf{FIRST}(T') = \{/, \epsilon\}$

$\textbf{FIRST}(T) = \{int, (\}$,
since $\epsilon \notin \textbf{FIRST}(F)$, **that's all.**

$\textbf{FIRST}(E') = \{-, \epsilon\}$

$\textbf{FIRST}(H) = \{-, int, (\}$

$\textbf{FIRST}(E) = \{-, int, (\}$,
since $\epsilon \in \textbf{FIRST}(E')$.

$\textbf{FIRST}(\epsilon) = \{\epsilon\}$

# How to compute $\textbf{FIRST}(\alpha)$?

- **Given $\textbf{FIRST}(X)$ for each terminal and nonterminal $X$, compute $\textbf{FIRST}(\alpha)$ for $\alpha$ being a sequence of terminals and/or nonterminals**
- **To build a parsing table, we need $\textbf{FIRST}(\alpha)$ for all $\alpha$ such that $X \to \alpha$ is a production in the grammar.**
- **$\alpha = X_1 X_2 \cdots X_n$**
  - **put $\textbf{FIRST}(X_1) - \{\epsilon\}$ into $\textbf{FIRST}(\alpha)$**
  - **if $\epsilon \in \textbf{FIRST}(X_1)$, then put $\textbf{FIRST}(X_2) - \{\epsilon\}$ into $\textbf{FIRST}(\alpha)$**
  - **$\cdots$**
  - **if $\epsilon \in \textbf{FIRST}(X_{n-1})$, then put $\textbf{FIRST}(X_n) - \{\epsilon\}$ into $\textbf{FIRST}(\alpha)$**
  - **if $\epsilon \in \textbf{FIRST}(X_i)$ for each $1 \leq i \leq n$, then put $\{\epsilon\}$ into $\textbf{FIRST}(\alpha)$.**

# Example for computing FIRST$(\alpha)$

Grammar
$E \rightarrow E'T$
$E' \rightarrow -TE' \mid \epsilon$
$T \rightarrow FT'$
$T' \rightarrow /FT' \mid \epsilon$
$F \rightarrow int \mid (E)$

**FIRST**$(F) = \{int, (\}$

**FIRST**$(T') = \{/, \epsilon\}$

**FIRST**$(T) = \{int, (\}$, **since** $\epsilon \notin$ **FIRST**$(F)$, **that's all.**

**FIRST**$(E') = \{-, \epsilon\}$

**FIRST**$(E) = \{-, int, (\}$, **since** $\epsilon \in$ **FIRST**$(E')$.

**FIRST**$(\epsilon) = \{\epsilon\}$

**FIRST**$(E'T) = \{-, int, (\}$

**FIRST**$(-TE') = \{-\}$

**FIRST**$(\epsilon) = \{\epsilon\}$

**FIRST**$(FT') = \{int, )\}$

**FIRST**$(/FT') = \{/\}$

**FIRST**$(\epsilon) = \{\epsilon\}$

**FIRST**$(int) = \{int\}$

**FIRST**$((E)) = \{(\}$

# Why do we need FIRST$(\alpha)$?

- **During parsing, suppose top-of-stack is a nonterminal $A$ and there are several choices**
  - $A \rightarrow \alpha_1$
  - $A \rightarrow \alpha_2$
  - $\cdots$
  - $A \rightarrow \alpha_k$

  **for derivation, and the current lookahead token is $a$**
- **If $a \in$ FIRST$(\alpha_i)$, then pick $A \rightarrow \alpha_i$ for derivation, pop, and then push $\alpha_i$.**
- **If $a$ is in several FIRST$(\alpha_i)$'s, then the grammar is not $LL(1)$.**
- **Question: if $a$ is not in any FIRST$(\alpha_i)$, does this mean the input stream cannot be accepted?**
  - **Maybe not!**
  - **What happen if $\epsilon$ is in some FIRST$(\alpha_i)$?**

# FOLLOW sets

- **Assume there is a special EOF symbol "$" ends every input.**
- **Add a new terminal "$."**
- **Definition: for a nonterminal $X$, FOLLOW$(X)$ is the set of terminals that can appear immediately to the right of $X$ in some partial derivation.**

  **That is, $S \overset{+}{\Longrightarrow} \alpha_1 X t \alpha_2$, where $t$ is a terminal.**
- **If $X$ can be the rightmost symbol in a derivation, then $ is in FOLLOW$(X)$.**
- **FOLLOW$(X) =$**

  $\{t \mid (t$ **is a terminal and** $S \overset{+}{\Longrightarrow} \alpha_1 X t \alpha_2)$ **or (** $t$ **is $ and** $S \overset{+}{\Longrightarrow} \alpha X)\}$.

# How to compute $\textbf{FOLLOW}(X)$

- **If $X$ is the starting nonterminal, put \$ into FOLLOW$(X)$.**
- **Find the productions with $X$ on the right-hand-side.**
  - **for each production of the form $Y \to \alpha X \beta$, put FIRST$(\beta) - \{\epsilon\}$ into FOLLOW$(X)$.**
  - **if $\epsilon \in$ FIRST$(\beta)$, then put FOLLOW$(Y)$ into FOLLOW$(X)$.**
  - **for each production of the form $Y \to \alpha X$, put FOLLOW$(Y)$ into FOLLOW$(X)$.**
- **To see if a given grammar is $LL(1)$ and also to build its parsing table:**
  - **compute FIRST$(\alpha)$ for every production $X \to \alpha$**
  - **compute FOLLOW$(X)$ for all nonterminals $X$**
- **Note that FIRST and FOLLOW sets are always sets of terminals, plus, perhaps, $\epsilon$ for some FIRST sets.**

# A complete example

- **Grammar**
  - $S \to Bc \mid DB$
  - $B \to ab \mid cS$
  - $D \to d \mid \epsilon$

| $\alpha$ | **FIRST**$(\alpha)$ | **FOLLOW**$(\alpha)$ |
|---|---|---|
| $D$ | $\{d, \epsilon\}$ | $\{a, c\}$ |
| $B$ | $\{a, c\}$ | $\{c, \$\}$ |
| $S$ | $\{a, c, d\}$ | $\{c, \$\}$ |
| $Bc$ | $\{a, c\}$ | |
| $DB$ | $\{d, a, c\}$ | |
| $ab$ | $\{a\}$ | |
| $cS$ | $\{c\}$ | |
| $d$ | $\{d\}$ | |
| $\epsilon$ | $\{\epsilon\}$ | |

# Why do we need FOLLOW sets?

- **Note FOLLOW$(S)$ always includes \$!**
- **Situation:**
  - During parsing, the top-of-stack is a nonterminal $X$ and the lookahead symbol is $a$.
  - Assume there are several choices for the nest derivation:
    - ▷ $X \rightarrow \alpha_1$
    - ▷ $\cdots$
    - ▷ $X \rightarrow \alpha_k$
  - If $a \in$ **FIRST**$(\alpha_{g_i})$ for only one $g_i$, then we use that derivation.
  - If $a \in$ **FIRST**$(\alpha_i)$ for two $i$, then this grammar is not $LL(1)$.
  - If $a \notin$ **FIRST**$(\alpha_i)$ for all $i$, then this grammar can still be $LL(1)$!
- **If some $\alpha_{g_i} \overset{*}{\Longrightarrow} \epsilon$ and $a \in$ FOLLOW$(X)$, then we can can use the derivation $X \rightarrow \alpha_{g_i}$.**

# Grammars that are not $LL(1)$

- **A grammar is not $LL(1)$ if any/both of the following is/are true.**
  - **There exists productions**
    - ▷ $A \rightarrow \alpha$
    - ▷ $A \rightarrow \beta$

    **such that $\textbf{FIRST}(\alpha) \cap \textbf{FIRST}(\beta) \neq \emptyset$.**
  - **There exists productions**
    - ▷ $A \rightarrow \alpha$
    - ▷ $A \rightarrow \beta$

    **such that $\epsilon \in \textbf{FIRST}(\alpha)$ and $\textbf{FIRST}(\beta) \cap \textbf{FOLLOW}(A) \neq \emptyset$.**
- **If a grammar is not $LL(1)$, you cannot write a linear-time predictive parser as described above.**

# A complete example (1/2)

- **Grammar:**
  - **<prog_head>→ PROG ID <file_list> SEMICOLON**
  - **<file_list>→ $\epsilon$ | L_PAREN <file_list> SEMICOLON**
- **FIRST and FOLLOW sets:**

| $\alpha$ | FIRST($\alpha$) | FOLLOW($\alpha$) |
|---|---|---|
| <prog_head> | {PROG} | {\$} |
| <file_list> | {$\epsilon$, L_PAREN} | {$\epsilon$, SEMICOLON} |
| PROG ID <file_list> SEMICOLON | {PROG} | |
| $\epsilon$ | {$\epsilon$} | |
| L_PAREN <file_list> SEMICOLON | {LPAREN} | |

# A complete example (2/2)

Input: PROG ID SEMICOLON

| Input | stack | action |
|-------|-------|--------|
| | <prog_head> $ | |
| PROG | <prog_head> $ | pop, push |
| PROG | PROG ID <file_list> SEMICOLON $ | match input |
| ID | ID <file_list> SEMICOLON $ | match input |
| SEMICOLON | <file_list> SEMICOLON $ | WHAT TO DO? |

- **Last actions:**
  - **Two choices:**
    - ▷ *<file_list>→ ε | L_PAREN <file_list> SEMICOLON*
  - **SEMICOLON ∉ FIRST($\epsilon$) and
    SEMICOLON ∉ FIRST(L_PAREN <file_list> SEMICOLON)**
  - **<file_list>$\overset{*}{\Longrightarrow}$ ε**
  - **SEMICOLON ∈ FOLLOW(<file_list>)**
  - **Hence we use the derivation
    <file_list>→ ε**

# $LL(1)$ **Parsing table (1/2)**

Grammar:

- $S \rightarrow XC$
- $X \rightarrow a \mid \epsilon$
- $C \rightarrow a \mid \epsilon$

| $\alpha$ | FIRST$(\alpha)$ | FOLLOW$(\alpha)$ |
|---|---|---|
| $S$ | $\{a, \epsilon\}$ | $\{\$\}$ |
| $X$ | $\{a, \epsilon\}$ | $\{a, \$\}$ |
| $C$ | $\{a, \epsilon\}$ | $\{\$\}$ |
| $\epsilon$ | $\{\epsilon\}$ | |
| $a$ | $\{a\}$ | |
| $XC$ | $\{a, \epsilon\}$ | |

- **Check for possible conflicts in $X \rightarrow a \mid \epsilon$.**
  - **FIRST**$(a) \cap$ **FIRST**$(\epsilon) = \emptyset$
  - $\epsilon \in$ **FIRST**$(\epsilon)$ **and FOLLOW**$(X) \cap$ **FIRST**$(a) = \{a\}$ **Conflict!!**
  - $\epsilon \notin$ **FIRST**$(a)$
- **Check for possible conflicts in $C \rightarrow a \mid \epsilon$.**
  - **FIRST**$(a) \cap$ **FIRST**$(\epsilon) = \emptyset$
  - $\epsilon \in$ **FIRST**$(\epsilon)$ **and FOLLOW**$(C) \cap$ **FIRST**$(a) = \emptyset$
  - $\epsilon \notin$ **FIRST**$(a)$

# $LL(1)$ **Parsing table (2/2)**

- **Parsing table:**

| | $a$ | $\$$ |
|---|---|---|
| $S$ | $S \to XC$ | $S \to XC$ |
| $X$ | **conflict!** | $X \to \epsilon$ |
| $C$ | $C \to a$ | $C \to \epsilon$ |

# Bottom-up parsing (Shift-reduce parsers)

- **Intuition: construct the parse tree from leaves to the root.**

- **Grammar:**
$$S \rightarrow AB$$

- **Example:**
$$A \rightarrow x \mid Y$$
$$B \rightarrow w \mid Z$$

- **Input** $xw$.

$$S \underset{rm}{\Longrightarrow} AB \underset{rm}{\Longrightarrow} Aw \underset{rm}{\Longrightarrow} xw.$$

# Definitions (1/2)

- **Left-most derivation:**
  - $S \underset{rm}{\Longrightarrow} \alpha$: **the rightmost nonterminal is replaced.**

  - $S \underset{rm}{\overset{+}{\Longrightarrow}} \alpha$: $\alpha$ **is derived from** $S$ **using one or more rightmost derivations.**

  - $\alpha$ **is called a** right-sentential form .
- **Define similarly leftmost derivations.**
- handle : **a handle for a right-sentential form** $\gamma$ **is the combining of the following two information:**
  - **a production rule** $A \rightarrow \beta$ **and**
  - **a position in** $\gamma$ **where** $\beta$ **can be found.**

# Definitions (2/2)

**Example:**

$$S \to aABe$$
$$A \to Abc \mid b$$
$$B \to d$$

input: **abbcde**

$\gamma \equiv aAbcde$ **is a right-sentential form**

$A \to Abc$ **and position 2 in** $\gamma$ **is a handle for** $\gamma$

- **reduce** : **replace a handle in a right-sentential form with its left-hand-side. In the above example, replace** $Abc$ **in** $\gamma$ **with** $A$.
- **A right-most derivation in reverse can be obtained by handle reducing.**

# STACK implementation

- **Four possible actions:**
  - shift: shift the input to **STACK**.
  - reduce: perform a reversed rightmost derivation.
  - accept
  - error

| STACK | INPUT | ACTION |
|-------|-------|--------|
| $ | xw$ | shift |
| $x | w$ | reduce by $A \rightarrow x$ |
| $A | w$ | shift |
| $Aw | $ | reduce by $B \rightarrow w$ |
| $AB | $ | reduce by $S \rightarrow AB$ |
| $S | $ | accept |

- **viable prefix** : the set of prefixes of right sentential forms that can appear on the stack.

# Model of a shift-reduce parser



- **Push-down automata!**
- **Current state $S_m$ encodes the symbols that has been shifted and the handles that are currently being matched.**
- $\$S_0 S_1 \cdots S_m a_i a_{i+1} \cdots a_n \$$ **represents a right sentential form.**
- **GOTO table:**
  - when a "reduce" action is taken, which handle to replace;
- **Action table:**
  - when a "shift" action is taken, which state currently in, that is, how to group symbols into handles.
- **The power of context free grammars is equivalent to nondeterministic push down automata.**

# LR parsers

- **By Don Knuth at 1965.**
- $LR(k)$**: see all of what can be derived from the right side with** $k$ **input tokens lookahead.**
  - **first** $L$**: scan the input from left to right**
  - **second** $R$**: reverse rightmost derivation**
  - $(k)$**: with** $k$ **lookahead tokens.**
- **Be able to decide the whereabout of a handle after seeing all of what have been derived so far plus** $k$ **input tokens lookahead.**

$$x_1, x_2, \ldots, \boxed{x_i, x_{i+1}, \ldots, x_{i+j},} \boxed{x_{i+j+1}, \ldots, x_{i+j+k-1},} \ldots$$

  **a handle**      **lookahead tokens**

- **Top-down parsing for** $LL(k)$ **grammars: be able to choose a production by seeing only the first** $k$ **symbols that will be derived from that production.**

# $LR(0)$ parsing

- **Construct a FSA to recognize all possible viable prefixes.**

- **An  $LR(0)$ item  ( item  for short) is a production, with a dot at some position in the RHS (right-hand side). For example:**
  - $A \rightarrow XYZ$
    - ▷ $A \rightarrow \cdot XYZ$
    - ▷ $A \rightarrow X \cdot YZ$
    - ▷ $A \rightarrow XY \cdot Z$
    - ▷ $A \rightarrow XYZ\cdot$

  - $A \rightarrow \epsilon$
    - ▷ $A \rightarrow \cdot$

  **The dot indicates the place of a handle.**

- **Assume  $G$  is  a  grammar  with  the  starting  symbol  $S$.  Augmented grammar  $G'$  is  to  add  a  new  starting  symbol  $S'$ and a new production $S' \rightarrow S$ to $G$. We assume working on the augmented grammar from now on.**

# Closure

- **The closure operation $closure(I)$, where $I$ is a set of items is defined by the following algorithm:**
    - **If $A \rightarrow \alpha \cdot B\beta$ is in $closure(I)$, then**
        - ▷ *at some point in parsing, we might see a substring derivable from $B\beta$ as input;*
        - ▷ *if $B \rightarrow \gamma$ is a production, we also see a substring derivable from $gamma$ at this point.*
        - ▷ *Thus $B \rightarrow \cdot\gamma$ should also be in $closure(I)$.*

- **What does $closure(I)$ means informally:**
    - when $A \rightarrow \alpha \cdot B\beta$ is encountered during parsing, then this means we have seen $\alpha$ so far, and expect to see $B\beta$ later before reducing to $A$.
    - at this point if $B \rightarrow \gamma$ is a production, then we may also want to see $B \rightarrow \cdot\gamma$ in order to reduce to $B$, and then advance to $A \rightarrow \alpha B \cdot \beta$.

- **Using $closure(I)$ to record all possible things that we have seen in the past and expect to see in the future.**

# Example for the closure function

- **Example:**
  - $E' \to E$
  - $E \to E + T \mid T$
  - $T \to T * F \mid F$
  - $F \to (E) \mid id$

$$closure(\{E' \to \cdot E\}) =$$
  - $E' \to \cdot E$
  - $E \to \cdot E + T$
  - $E \to \cdot T$
  - $T \to \cdot T * F$
  - $T \to \cdot F$
  - $F \to \cdot (E)$
  - $F \to \cdot id$

# GOTO table

- $GOTO(I, X)$, **where** $I$ **is a set of items and** $X$ **is a legal symbol is defined as**
  - **If** $A \rightarrow \alpha \cdot X\beta$ **is in** $I$, **then**
  - $closure(\{A \rightarrow \alpha X \cdot \beta\}) \subseteq GOTO(I, X)$
- **Informal meanings:**
  - **currently we have seen** $A \rightarrow \alpha \cdot X\beta$
  - **expect to see** $X$
  - **if we see** $X$,
  - **then we should be in the state** $closure(\{A \rightarrow \alpha X \cdot \beta\})$.
- **Use the GOTO table to denote the state to go to once we are in** $I$ **and have seen** $X$.

# Sets-of-items construction

- **Canonical $LR(0)$ items** : the set of all possible DFA states, where each state is a group of $LR(0)$ items.
- **Algorithm for constructing $LR(0)$ parsing table.**
  - $C \leftarrow \{closure(\{S' \rightarrow \cdot S\}\}$
  - **repeat**
    - ▷ *for each set of items $I$ in $C$ and each grammar symbol $X$ such that $GOTO(I, X) \neq \emptyset$ and not in $C$ do*
    - ▷ *add $GOTO(I, X)$ to $C$*
  - **until no more sets can be added to $C$**
- **Kernel of a state: items**
  - **that is not of the form $X \rightarrow \cdot \beta$ or**
  - $S' \rightarrow S$
- **Given the kernel of a state, all items in the state can be derived.**

# Example of sets of $LR(0)$ items

- **Grammar:**

$$E' \to E$$
$$E \to E + T \mid T$$
$$T \to T * F \mid F$$
$$F \to (E) \mid id$$

$I_0 = closure(\{E' \to \cdot E\})$:

$E' \to \cdot E$

$E \to \cdot E + T$

$E \to \cdot T$

$T \to \cdot T * F$

$T \to \cdot F$

$F \to \cdot (E)$

$F \to \cdot id$

- **Canonical $LR(0)$ items:**
  - $I_1 = GOTO(I_0, E)$
    - ▷ $E' \to E\cdot$
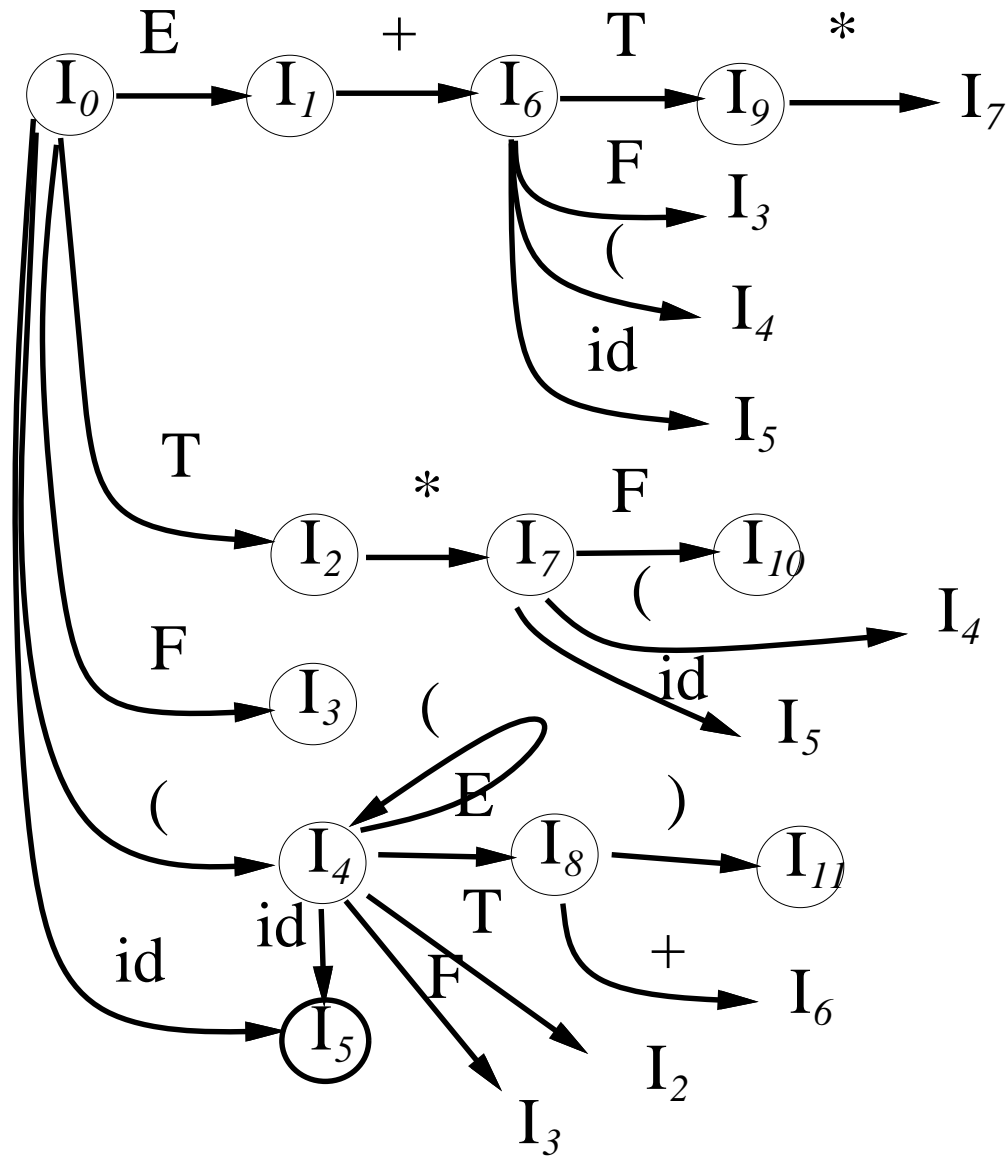    - ▷ $E \to E \cdot + T$
  - $I_2 = GOTO(I_0, T)$
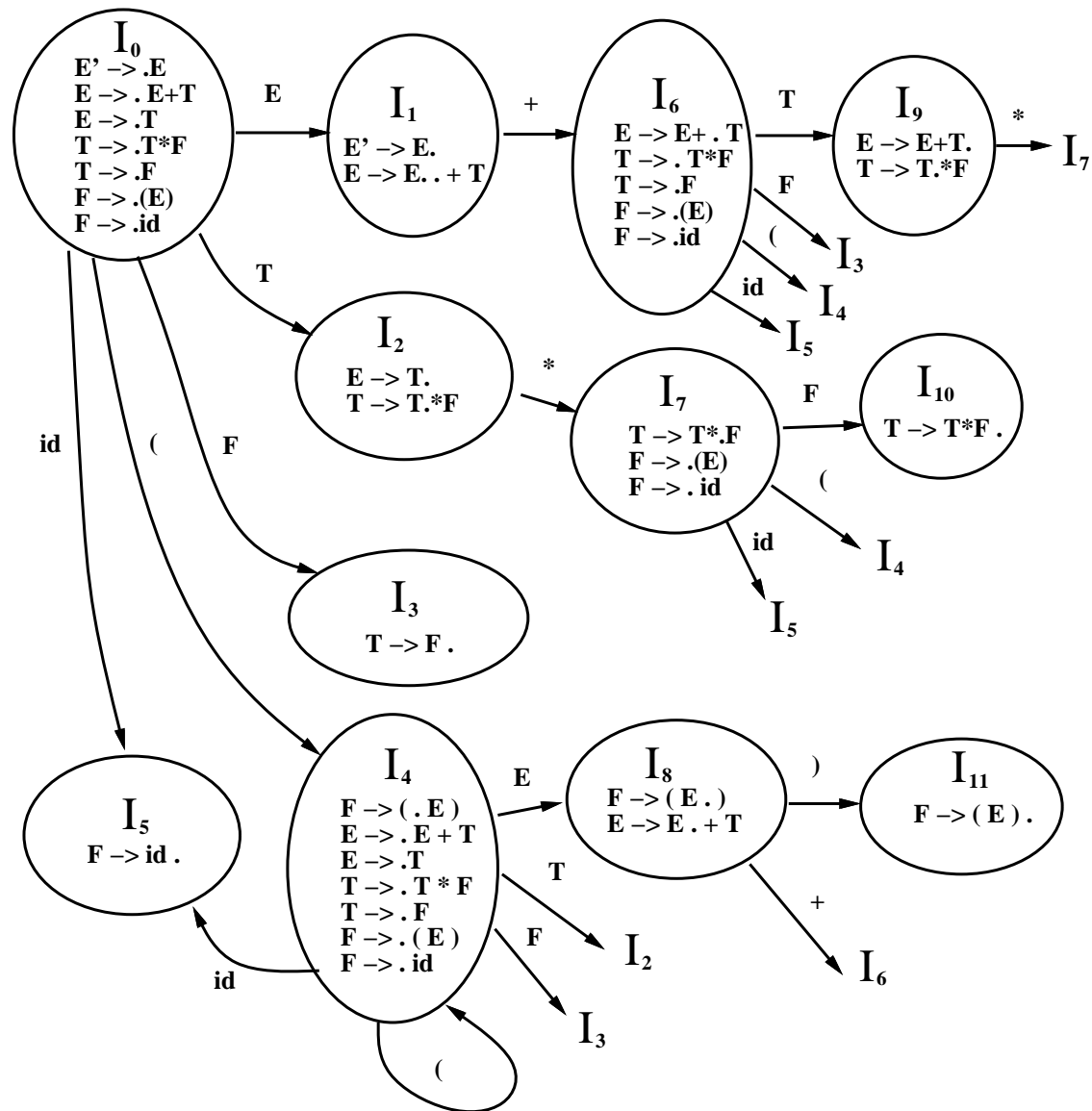    - ▷ $E \to T\cdot$
    - ▷ $T \to T \cdot * F$

# Transition diagram (1/2)

# Transition diagram (2/2)

# Meaning of $LR(0)$ transition diagram

- $E + T*$ is a viable prefix that can happen on the top of the stack while doing parsing.

- after seeing $E + T*$, we are in state $I_7$. $I_7 =$
  - $\{T \to T * \cdot F,$
  - $F \to \cdot (E),$
  - $F \to \cdot id\}$

- We expect to follow one of the following three possible derivations:

$$E' \underset{rm}{\Longrightarrow} E$$
$$\underset{rm}{\Longrightarrow} E + T$$
$$\underset{rm}{\Longrightarrow} E + T * F$$
$$\underset{rm}{\Longrightarrow} E + T * id$$
$$\underset{rm}{\Longrightarrow} E + \underline{T*}F * id$$
$$\ldots$$

$$E' \underset{rm}{\Longrightarrow} E$$
$$\underset{rm}{\Longrightarrow} E + T$$
$$\underset{rm}{\Longrightarrow} E + T * F$$
$$\underset{rm}{\Longrightarrow} \underline{E + T*(E)}$$
$$\ldots$$

$$E' \underset{rm}{\Longrightarrow} E$$
$$\underset{rm}{\Longrightarrow} E + T$$
$$\underset{rm}{\Longrightarrow} E + T * F$$
$$\underset{rm}{\Longrightarrow} \underline{E + T*id}$$
$$\ldots$$

# Definition of $closure(I)$ and $GOTO(I, X)$

- $closure(I)$: a state/configuration during parsing recording all possible things that we are expecting.
- If $A \rightarrow \alpha \cdot B\beta \in I$, then it means
  - in the middle of parsing, $\alpha$ is on the top of the stack;
  - at this point, we are expecting to see $B\beta$;
  - after we saw $B\beta$, we will reduce $\alpha B\beta$ to $A$ and make $A$ top of stack.
- To achieve the goal of seeing $B\beta$, we expect to perform some operations
  - We expect to see $B$ on the top of the stack first.
  - If $B \rightarrow \gamma$ is a production, then it might be the case that we shall see $\gamma$ on the top of the stack.
  - Then we reduce $\gamma$ to $B$.
  - Hence we need to include $B \rightarrow \cdot\gamma$ into $closure(I)$.
- $GOTO(I, X)$: when we are in the state described by $I$, and then a new symbol $X$ is pushed into the stack, If $A \rightarrow \alpha \cdot X\beta$ is in $I$, then $closure(\{A \rightarrow \alpha X \cdot \beta\}) \subseteq GOTO(I, X)$.

# Parsing example

- **Input: id * id + id**

| STACK | input | action |
|---|---|---|
| $\$\ I_0$ | id*id+id$ | |
| $\$\ I_0$ id $I_5$ | * id + id$ | shift 5 |
| $\$\ I_0$ F | * id + id$ | reduce by $F \rightarrow id$ |
| $\$\ I_0$ F $I_3$ | * id + id$ | in $I_0$, saw F, goto $I_3$ |
| $\$\ I_0$ T $I_2$ | * id + id$ | reduce by $T \rightarrow F$ |
| $\$\ I_0$ T $I_2$ * $I_7$ | id + id$ | shift 7 |
| $\$\ I_0$ T $I_2$ * $I_7$ id $I_5$ | + id$ | shift 5 |
| $\$\ I_0$ T $I_2$ * $I_7$ F $I_{10}$ | + id$ | reduce by $F \rightarrow id$ |
| $\$\ I_0$ T $I_2$ | + id$ | reduce by $T \rightarrow F$ |
| $\$\ I_0$ E $I_1$ | + id$ | reduce by $T \rightarrow T * F$ |
| $\$\ I_0$ E $I_1$ + $I_6$ | id$ | shift 6 |
| $\$\ I_0$ E $I_1$ + $I_6$ id $I_5$ | id$ | shift 5 |
| $\$\ I_0$ E $I_1$ + $I_6$ F $I_3$ | id$ | reduce by $F \rightarrow id$ |
| $\cdots$ | $\cdots$ | $\cdots$ |

# $LR(0)$ **parsing**

- $LR$ **parsing without lookahead symbols.**
- **Constructed from DFA for recognizing viable prefixes.**
- **In state $I_i$**
  - if $A \rightarrow \alpha \cdot a\beta$ **is in $I_i$ then perform "shift" while seeing the terminal $a$ in the input, and then go to the state** $closure(\{A \rightarrow \alpha a \cdot \beta\})$
  - if $A \rightarrow \beta \cdot$ **is in $I_i$, then perform "reduce by $A \rightarrow \beta$" and then goto the state $GOTO(I, A)$ where $I$ is the state on the top of the stack after removing $\beta$**
- **Conflicts:**
  - **shift/reduce conflict**
  - **reduce/reduce conflict**
- **Very few grammars are $LR(0)$. For example:**
  - **in $I_2$, you can either perform a reduce or a shift when seeing "*" in the input**
  - **However, it is not possible to have $E$ followed by "*". Thus we should not perform "reduce".**
- **Use FOLLOW$(E)$ as look ahead information to resolve some conflicts.**

# $SLR(1)$ parsing algorithm

- **Using FOLLOW sets to resolve conflicts in constructing $SLR(1)$ parsing table, where the first "S" stands for "simple".**
  - Input: an augmented grammar $G'$
  - Output: The $SLR(1)$ parsing table.
- **Construct $C = \{I_0, I_1, \ldots, I_n\}$ the collection of sets of $LR(0)$ items for $G'$.**
- **The parsing table for state $I_i$ is determined as follows:**
  - if $A \to \alpha \cdot a\beta$ is in $I_i$ and $GOTO(I_i, a) = I_j$, then $action(I_i, a)$ is "shift $j$" for $a$ being a terminal.
  - If $A \to \alpha \cdot$ is in $I_i$, then $action(I_i, a)$ is "reduce by $A \to \alpha$" for all terminal $a \in$ **FOLLOW**$(A)$; here $A \neq S'$
  - if $S' \to S \cdot$ is in $I_i$, then $action(I_i, \$)$ is "accept".
- **If any conflicts are generated by the above algorithm, we say the grammar is not $SLR(1)$.**

# $SLR(1)$ **parsing table**

| state | action | | | | | | GOTO | | |
|---|---|---|---|---|---|---|---|---|---|
| | id | + | * | ( | ) | $ | E | T | F |
| 0 | s5 | | | s4 | | | 1 | 2 | 3 |
| 1 | | s6 | | | | accept | | | |
| 2 | | r2 | s7 | | r2 | r2 | | | |
| 3 | | r4 | r4 | | r4 | r4 | | | |
| 4 | s5 | | | s4 | | | 8 | 2 | 3 |
| 5 | | r6 | r6 | | r6 | r6 | | | |
| 6 | s5 | | | s4 | | | | 9 | 3 |
| 7 | s5 | | | s4 | | | | | 10 |
| 8 | | s6 | | | s11 | | | | |
| 9 | | r1 | s7 | | r1 | r1 | | | |
| 10 | | r3 | r3 | | r3 | r3 | | | |
| 11 | | r5 | r5 | | r5 | r5 | | | |

- **r$i$ means reduce by production numbered $i$.**
- **s$i$ means shift and then go to state $I_i$.**
- **Use FOLLOW$(A)$ to resolve some conflicts.**

# Discussion (1/3)

- **Every $SLR(1)$ grammar is unambiguous, but there are many unambiguous grammars that are not $SLR(1)$.**
- **Example:**
  - $S \rightarrow L = R \mid R$
  - $L \rightarrow *R \mid id$
  - $R \rightarrow L$
- **States:**
  - $I_0$:
    - ▷ $S' \rightarrow \cdot S$
    - ▷ $S \rightarrow \cdot L = R$
    - ▷ $S \rightarrow \cdot R$
    - ▷ $L \rightarrow \cdot * R$
    - ▷ $L \rightarrow \cdot id$
    - ▷ $R \rightarrow \cdot L$
  - $I_1$: $S' \rightarrow S\cdot$
  - $I_2$:
    - ▷ $S \rightarrow L\cdot = R$
    - ▷ $R \rightarrow L\cdot$

# Discussion (2/3)

$I_3$: $S \rightarrow R\cdot$

$I_4$:

   ▷ $L \rightarrow * \cdot R$
   ▷ $R \rightarrow \cdot L$
   ▷ $L \rightarrow \cdot * R$
   ▷ $L \rightarrow \cdot id$

$I_5$: $L \rightarrow id\cdot$
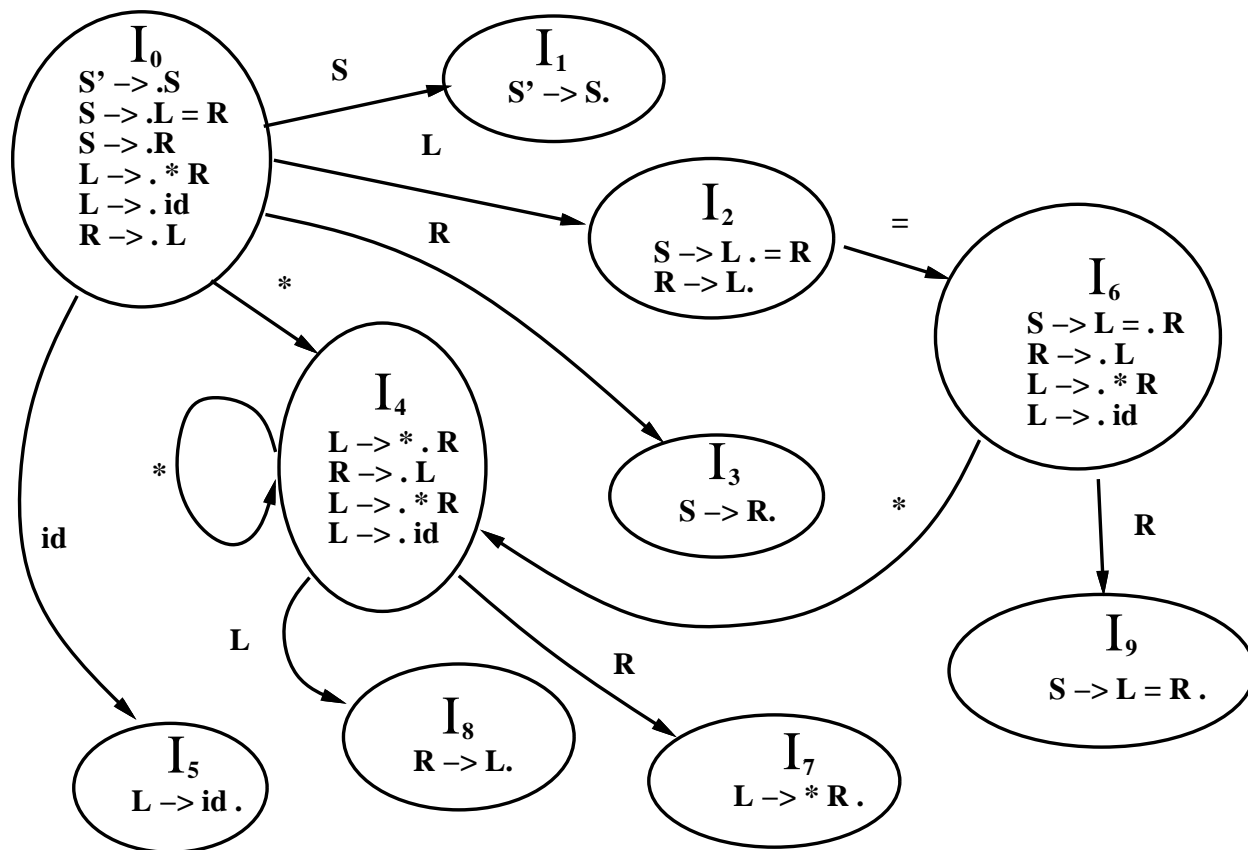
$I_6$:

   ▷ $S \rightarrow L = \cdot R$
   ▷ $R \rightarrow \cdot L$
   ▷ $L \rightarrow \cdot * R$
   ▷ $L \rightarrow \cdot id$

$I_7$: $L \rightarrow *R\cdot$

$I_8$: $R \rightarrow L\cdot$

$I_9$: $S \rightarrow L = R\cdot$

# Discussion (3/3)

- **Suppose the stack has $\$I_0LI_2$ and the input is "=". We can either**
  - shift 6, or
  - reduce by $R \to L$, since $= \in$ **FOLLOW**$(R)$.
- **This grammar is ambiguous for $SLR(1)$ parsing.**
- **However, we should not perform a $R \to L$ reduction.**
  - after performing the reduction, the viable prefix is $\$R$;
  - $= \notin$ **FOLLOW**$(\$R)$
  - $= \in$ **FOLLOW**$(*R)$
  - That is to say, we cannot find a right sentential form with the prefix $R = \cdots$.
  - We can find a right sentential form with $\cdots * R = \cdots$

# Canonical LR — $LR(1)$

- In $SLR(1)$ **parsing, if** $A \rightarrow \alpha\cdot$ **is in state** $I_i$**, and** $a \in$ **FOLLOW**$(A)$**, then we perform the reduction** $A \rightarrow \alpha$**.**
- **However, it is possible that when state** $I_i$ **is on the top of the stack, the viable prefix** $\beta\alpha$ **on the stack is such that** $\beta A$ **cannot be followed by** $a$**.**
- **We can solve the problem by knowing more left context using the technique of** lookahead propagation **.**

# $LR(1)$ **items**

- **An $LR(1)$ item is in the form of $[A \rightarrow \alpha \cdot \beta, a]$, where the first field is an $LR(0)$ item and the second field $a$ is a terminal belonging to a subset of FOLLOW$(A)$.**

- **Intuition: perform a reduction based on an $LR(1)$ item $[A \rightarrow \alpha \cdot, a]$ only when the next symbol is $a$.**

- **Formally: $[A \rightarrow \alpha \cdot \beta, a]$ is valid (or reachable) for a viable prefix $\gamma$ if there exists a derivation**

$$S \underset{rm}{\overset{*}{\Longrightarrow}} \delta A \omega \underset{rm}{\Longrightarrow} \delta \alpha \beta \omega,$$

  **where**
    - $\gamma = \delta \alpha$
    - **either $a \in$ FIRST$(\omega)$ or**
    - $\omega = \epsilon$ **and $a = \$$.**

# $LR(1)$ **parsing example**

- **Grammar:**
  - $S \rightarrow BB$
  - $B \rightarrow aB \mid b$

$$S \underset{rm}{\overset{*}{\Longrightarrow}} aaBab \underset{rm}{\Longrightarrow} aaaBab$$

**viable prefix** $aaa$ **can reach** $[B \rightarrow a \cdot B, a]$

$$S \underset{rm}{\overset{*}{\Longrightarrow}} BaB \underset{rm}{\Longrightarrow} BaaB$$

**viable prefix** $Baa$ **can reach** $[B \rightarrow a \cdot B, \$]$

# Finding all $LR(1)$ items

- **Ideas: redefine the closure function.**
    - **suppose $[A \rightarrow \alpha \cdot B\beta, a]$ is valid for a viable prefix $\gamma \equiv \delta\alpha$**
    - **in other words**

$$S \underset{rm}{\overset{*}{\Longrightarrow}} \delta A a\omega \underset{rm}{\Longrightarrow} \delta\alpha B\beta a\omega$$

 - **Then for each production $B \rightarrow \eta$ assume $\beta a\omega$ derives the sequence of terminals $bc$.**

$$S \underset{rm}{\overset{*}{\Longrightarrow}} \delta\alpha B \boxed{\beta a\omega} \underset{rm}{\overset{*}{\Longrightarrow}} \delta\alpha B \boxed{bc} \underset{rm}{\overset{*}{\Longrightarrow}} \delta\alpha \boxed{\eta} bc$$

 **Thus $[B \rightarrow \cdot\eta, b]$ is also valid for $\gamma$ for each $b \in$ FIRST$(\beta a)$.**
 **Note $a$ is a terminal. So FIRST$(\beta a) =$ FIRST$(\beta a\omega)$.**

- **Lookahead propagation .**

# Algorithm for $LR(1)$ parsing functions

- $closure(I)$
  - **repeat**
    - ▷ **for each item** $[A \rightarrow \alpha \cdot B\beta, a]$ **in** $I$ **do**
    - ▷      **if** $B \rightarrow \cdot\eta$ **is in** $G'$
    - ▷      **then add** $[B \rightarrow \cdot\eta, b]$ **to** $I$ **for each** $b \in \mathbf{FIRST}(\beta a)$
  - **until no more items can be added to** $I$
  - **return** $i$
- $GOTO(I, X)$
  - **let** $J = \{[A \rightarrow \alpha X \cdot \beta, a] \mid [A \rightarrow \alpha \cdot X\beta, a] \in I\}.$
  - **return** $closure(J)$
- $items(G')$
  - $C \leftarrow \{closure(\{[S' \rightarrow \cdot S, \$]\})\}$
  - **repeat**
    - ▷ **for each set of items** $I \in C$ **and each grammar symbol** $X$ **such that** $GOTO(I, X) \neq \emptyset$ **and** $GOTO(I, X) \notin C$ **do**
    - ▷      **add** $GOTO(I, X)$ **to** $C$
  - **until no more sets of items can be added to** $C$

# Example for constructing $LR(1)$ closures

- **Grammar:**
  - $S' \rightarrow S$
  - $S \rightarrow CC$
  - $C \rightarrow cC \mid d$
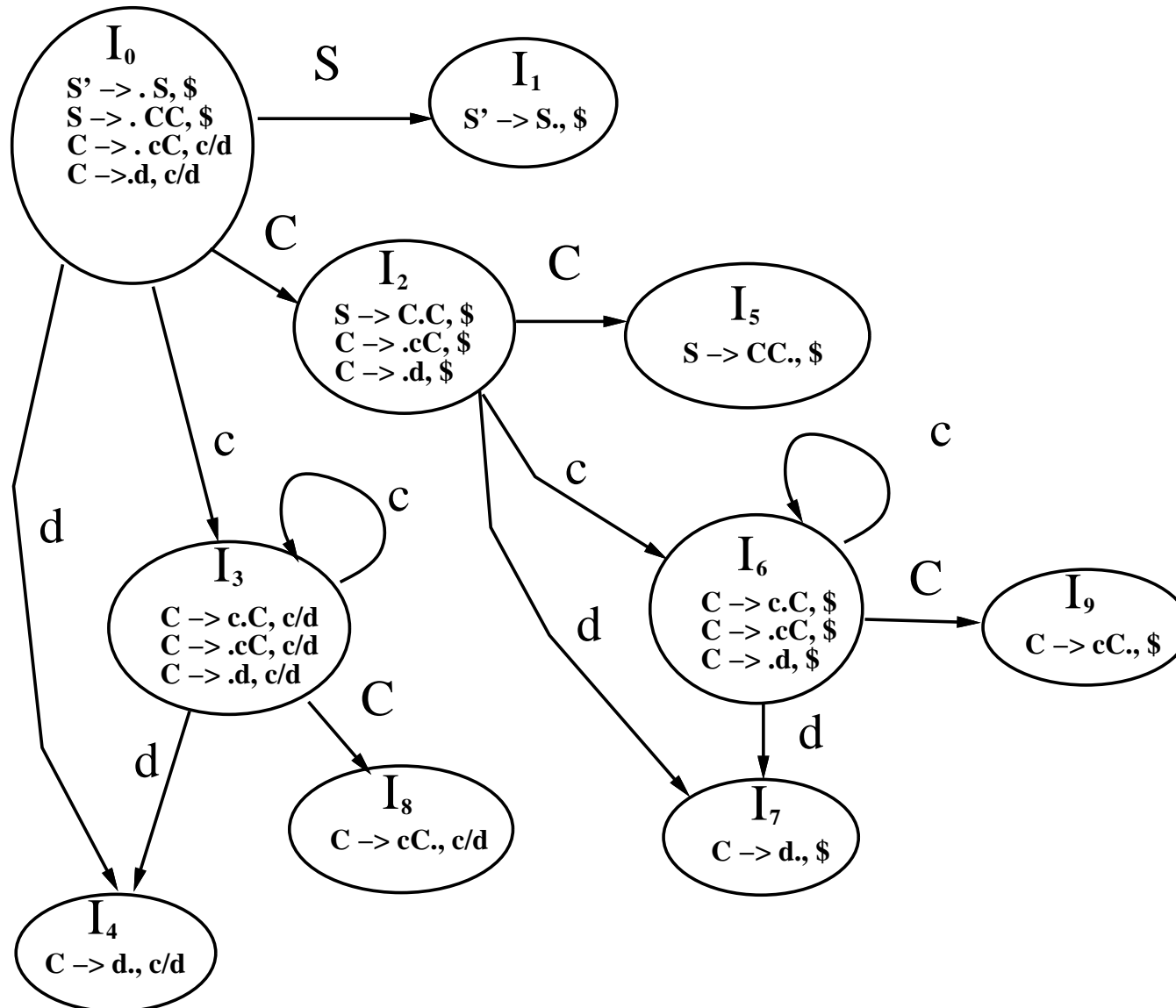- $closure(\{[S' \rightarrow \cdot S, \$]\}) =$
  - $\{[S' \rightarrow \cdot S, \$],$
  - $[S \rightarrow \cdot CC, \$],$
  - $[C \rightarrow \cdot cC, c/d],$
  - $[C \rightarrow \cdot d, c/d]\}$
- **Note:**
  - **FIRST**$(\epsilon\$) = \{\$\}$
  - **FIRST**$(C\$) = \{c, d\}$
  - $[C \rightarrow \cdot cC, c/d]$ **means**
    - ▷ $[C \rightarrow \cdot cC, c]$ **and**
    - ▷ $[C \rightarrow \cdot cC, d]$.

# $LR(1)$ **Transition diagram**

# $LR(1)$ parsing example

- **Input** $cdccd$

| STACK | INPUT | ACTION |
|---|---|---|
| $\$ \ I_0$ | cdccd$\$$ | |
| $\$ \ I_0 \ c \ I_3$ | dccd$\$$ | shift 3 |
| $\$ \ I_0 \ c \ I_3 \ d \ I_4$ | ccd$\$$ | shift 4 |
| $\$ \ I_0 \ c \ I_3 \ C \ I_8$ | ccd$\$$ | reduce by $C \to d$ |
| $\$ \ I_0 \ C \ I_2$ | ccd$\$$ | reduce by $C \to cC$ |
| $\$ \ I_0 \ C \ I_2 \ c \ I_6$ | cd$\$$ | shift 6 |
| $\$ \ I_0 \ C \ I_2 \ c \ I_6 \ c \ I_6$ | d$\$$ | shift 6 |
| $\$ \ I_0 \ C \ I_2 \ c \ I_6 \ c \ I_6$ | d$\$$ | shift 6 |
| $\$ \ I_0 \ C \ I_2 \ c \ I_6 \ c \ I_6 \ d \ I_7$ | $\$$ | shift 7 |
| $\$ \ I_0 \ C \ I_2 \ c \ I_6 \ c \ I_6 \ C \ I_9$ | $\$$ | reduce by $C \to cC$ |
| $\$ \ I_0 \ C \ I_2 \ c \ I_6 \ C \ I_9$ | $\$$ | reduce by $C \to cC$ |
| $\$ \ I_0 \ C \ I_2 \ C \ I_5$ | $\$$ | reduce by $S \to CC$ |
| $\$ \ I_0 \ S \ I_1$ | $\$$ | reduce by $S' \to S$ |
| $\$ I_0 \ S'$ | $\$$ | accept |

# Algorithm for $LR(1)$ parsing table

- **Construction of canonical $LR(1)$ parsing tables.**
    - **Input: an augmented grammar $G'$**
    - **Output: The canonical $LR(1)$ parsing table, i.e., the ACTION table.**
- **Construct $C = \{I_0, I_1, \ldots, I_n\}$ the collection of sets of $LR(1)$ items form $G'$.**
- **Action table is constructed as follows:**
    - **if $[A \rightarrow \alpha \cdot a\beta, b] \in I_i$ and $GOTO(I_i, a) = I_j$, then $action[I_i, a] =$ "shift $j$" for $a$ is a terminal.**
    - **if $[A \rightarrow \alpha\cdot, a] \in I_i$ and $A \neq S'$, then $action[I_i, a] =$ "reduce by $A \rightarrow \alpha$"**
    - **if $[S' \rightarrow S., \$] \in I_i$, then $action[I_i, \$] =$ "accept."**
- **If conflicts result from the above rules, then the grammar is not $LR(1)$.**
- **The initial state of the parser is the one constructed from the set containing the item $[S' \rightarrow \cdot S, \$]$.**

# An example of an $LR(1)$ parsing table

| state | action | | | GOTO | |
|-------|--------|------|--------|------|------|
|       | c      | d    | $      | S    | C    |
| 0     | s3     | s4   |        | 1    | 2    |
| 1     |        |      | accept |      |      |
| 2     | s6     | s7   |        |      | 5    |
| 3     | s3     | s4   |        |      | 8    |
| 4     | r3     | r3   |        |      |      |
| 5     |        |      | r1     |      |      |
| 6     | s6     | s7   |        |      | 9    |
| 7     |        |      | r3     |      |      |
| 8     | r2     | r2   |        |      |      |
| 9     |        |      | r2     |      |      |

- **Canonical $LR(1)$ parser**
  - **too many states and thus occupy too much space**
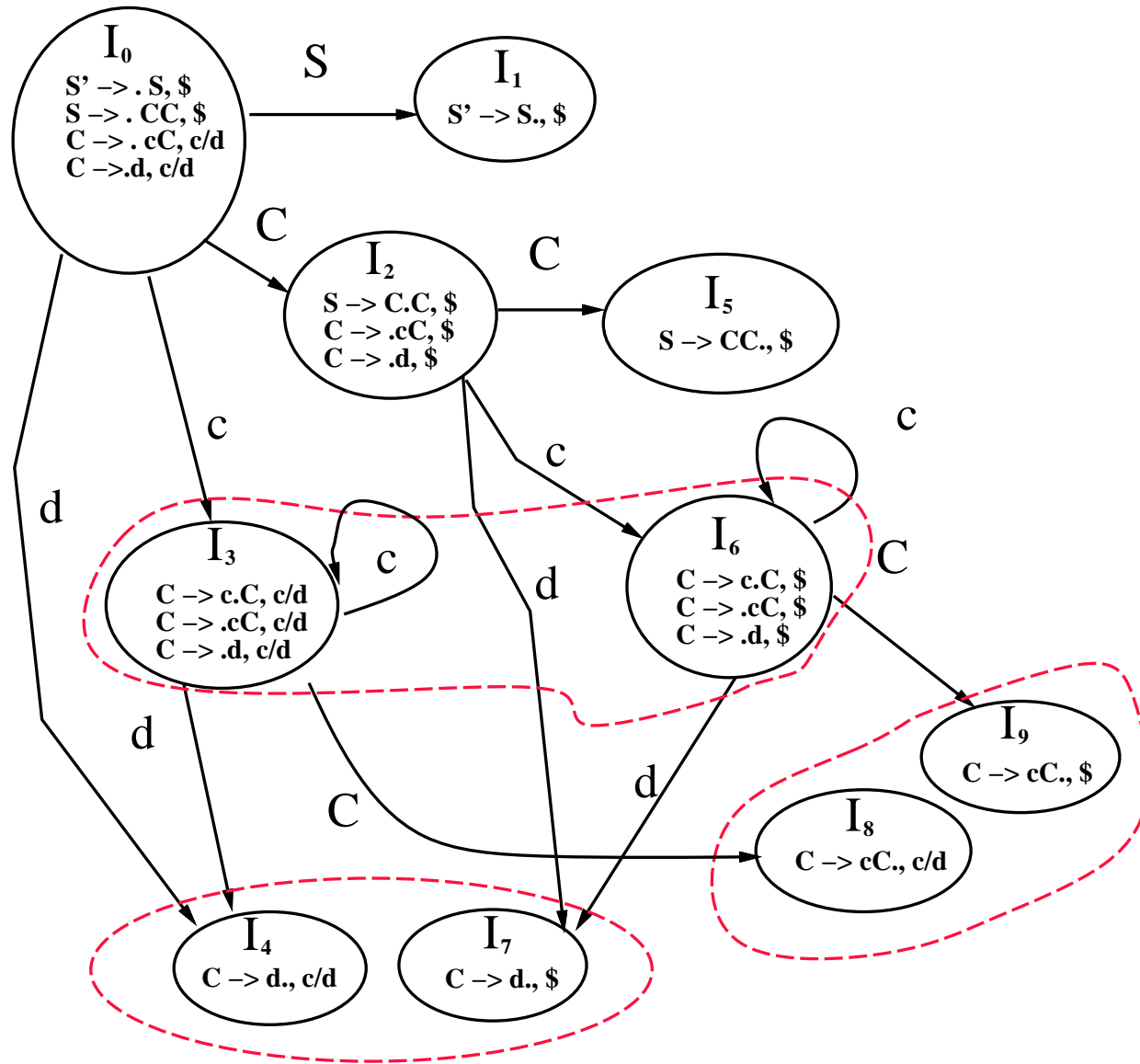  - **most powerful**

# $LALR(1)$ parser — Lookahead LR

- **The method that is often used in practice.**
- **Most common syntactic constructs of programming languages can be expressed conveniently by an $LALR(1)$ grammar.**
- **$SLR(1)$ and $LALR(1)$ always have the same number of states.**
- **Number of states is about 1/10 of that of $LR(1)$.**
- **Simple observation:**
  **an $LR(1)$ item is in the form of $[A \to \alpha \cdot \beta, c]$**

- **We call $A \to \alpha \cdot \beta$ the  first component .**

- **Definition: in an $LR(1)$ state, set of first components is called its  core .**

# Intuition for $LALR(1)$ grammars

- **In $LR(1)$ parser, it is a common thing that several states only differ in lookahead symbol, but have the same core.**
- **To reduce the number of states, we might want to merge states with the same core.**
  - If $I_4$ and $I_7$ are merged, then the new state is called $I_{4,7}$
- **After merging the states, revise the GOTO table accordingly.**
- **merging of states can never produce a shift-reduce conflict that was not present in one of the original states.**
  - $I_1 = \{[A \rightarrow \alpha\cdot, a], \dots\}$
  - $I_2 = \{[B \rightarrow \beta \cdot a\gamma, b], \dots\}$
  - **For $I_1$, we perform a reduce on $a$.**
  - **For $I_2$, we perform a shift on $a$.**
  - **Merging $I_1$ and $I_2$, the new state $I_{1,2}$ has shift-reduce conflicts.**
  - **This is impossible, in the original table since $I_1$ and $I_2$ have the same core.**
  - **So $[A \rightarrow \alpha\cdot, c] \in I_2$ and $[B \rightarrow \beta \cdot a\gamma, d] \in I_1$.**
  - **The shift-reduce conflict already occurs in $I_1$ and $I_2$.**

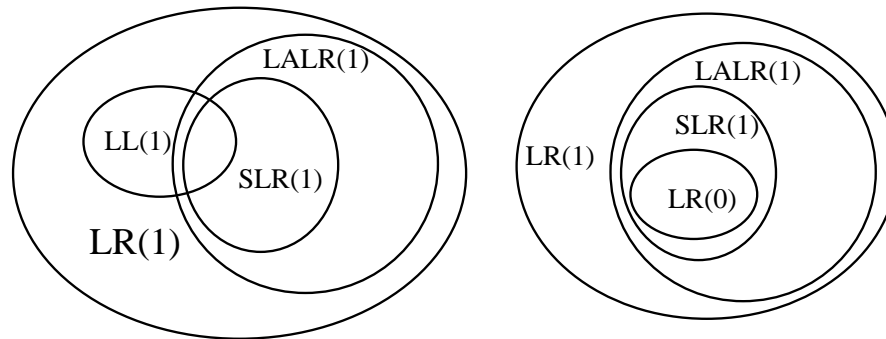# $LALR(1)$ **Transition diagram**

# Possible new conflicts from $LALR(1)$

- **May produce a new reduce-reduce conflict.**
- **For example (textbook page 238), grammar:**
  - $S' \to S$
  - $S \to aAd \mid bBf \mid aBe \mid bAe$
  - $A \to c$
  - $B \to c$
- **The language recognized by this grammar is $\{acd, ace, bcd, bce\}$.**
- **You may check that this grammar is $LR(1)$ by constructing the sets of items.**
- **You will find the set of items $\{[A \to c\cdot, d], [B \to c\cdot, e]\}$ is valid for the viable prefix $ac$, and $\{[A \to c\cdot, e], [B \to c\cdot, d]\}$ is valid for the viable prefix $bc$.**
- **Neither of these sets generates a conflict, and their cores are the same. However, their union, which is**
  - $\{[A \to c\cdot, d/e],$
  - $[B \to c\cdot, d/e]\}$

**generates a reduce-reduce conflict, since reductions by both $A \to c$ and $B \to c$ are called for on inputs $d$ and $e$.**
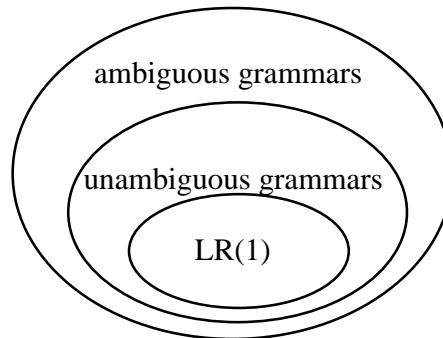
# How to construct $LALR(1)$ parsing table

- **Naive approach:**
  - Construct $LR(1)$ parsing table, which takes lots of intermediate spaces.
  - Merging states.
- **Space efficient methods to construct an $LALR(1)$ parsing table are known.**
  - Construction and merging on the fly.



- **Summary:**
- $LR(1)$ and $LALR(1)$ can almost handle all programming languages, but $LALR(1)$ is easier to write.
- $LL(1)$ is easier to understand.

# Using ambiguous grammars



- **Ambiguous grammars provides a shorter, more natural specification than any equivalent unambiguous grammars.**
- **Sometimes need ambiguous grammars to specify important language constructs.**
- **For example: declare a variable before its usage.**

```
var xyz : integer
begin
    ...
    xyz := 3;
    ...
```

# Ambiguity from precedence and associativity

- **Use precedence and associativity to resolve conflicts.**
- **Example:**
    - $G_1$:
        - $\triangleright$ $E \rightarrow E + E \mid E * E \mid (E) \mid id$
        - $\triangleright$ *ambiguous, but easy to understand!*
    - $G_2$:
        - $\triangleright$ $E \rightarrow E + T \mid T$
        - $\triangleright$ $E \rightarrow T * F \mid F$
        - $\triangleright$ $F \rightarrow (E) \mid id$
        - $\triangleright$ *unambiguous, but it is difficult to change the precedence;*
        - $\triangleright$ *parse tree is much larger for $G_2$, and thus takes more time to parse.*

- **When parsing the following input for $G_1$: $id + id * id$.**
    - **Assume the input parsed so far is $id + id$.**
    - **We now see "*".**
    - **We can either shift or perform "reduce by $E \rightarrow E + E$".**
    - **When there is a conflict, say in $SLR(1)$ parsing, we use precedence and associativity information to resolve conflicts.**

# Dangling-else ambiguity

- **Grammar:**
  - $S \to a \mid if <$**condition**$> then <$**statement**$>$
    $\mid if <$**condition**$> then <$**statement**$> else <$**statement**$>$
- **When seeing**

  ### if c then S else S

  - shift or reduce conflict;
  - always favor a shift.
  - Intuition: favor a longer match.

# Special cases

- **Ambiguity from special-case productions:**
  - **Sometime a very rare happened special case causes ambiguity.**
  - **It's too costly to revise the grammar. We can resolve the conflicts by using special rules.**
  - **Example:**
    - $\triangleright$ $E \rightarrow E \ sub \ E \ sup \ E$
    - $\triangleright$ $E \rightarrow E \ sub \ E$
    - $\triangleright$ $E \rightarrow E \ sup \ E$
    - $\triangleright$ $E \rightarrow \{E\} \mid character$
  - **Meanings:**
    - $\triangleright$ $W \ sub \ U$: $W_U$.
    - $\triangleright$ $W \ sup \ U$: $W^U$.
    - $\triangleright$ $W \ sub \ U \ sup \ V$ **is** $W_U^V$**, not** $W_U{}^V$
  - **Resolve by semantic and special rules.**
  - **Pick the right one when there is a reduce/reduce conflict.**
    - $\triangleright$ *Reduce the production listed earlier.*

# YACC (1/2)

- **Yet Another Compiler Compiler:**
  - **A UNIX utility for generating $LALR(1)$ parsing tables.**
  - **Convert your YACC code into C programs.**

  - **file.y** $\longrightarrow$ | **yacc file.y** | $\longrightarrow$ **y.tab.c**

  - **y.tab.c** $\longrightarrow$ | **cc -ly -ll y.tab.c** | $\longrightarrow$ **a.out**

- **Format:**
  - **declarations**
  - **%%**
  - **translation rules**
    - ▷ *<left side>: <production>*
    - ▷ *{ semantic rules }*
  - **%%**
  - **supporting C-routines.**

# YACC (2/2)

- **Assume the Lexical analyzer routine is** $yylex()$.
- **When there are ambiguities:**
  - **reduce/reduce conflict: favor the one listed first.**
  - **shift/reduce conflict: favor shift. (longer match!)**
- **Error handling:**
  - **Use special error handling productions.**
  - **Example:**
    ```
    lines: error '\n' {...}
    ```
  - **when there is an error, skip until newline.**
  - $error$**: special token.**
  - $yyerror(string)$**: pre-defined routine for printing error messages.**
  - $yyerrok()$**: reset error flags.**

# YACC code example (1/2)

```
%{
#include <stdio.h>
#include <ctype.h>
#include <math.h>
#define YYSTYPE int /* integer type for YACC stack */

%}

%token NUMBER
%left '+' '-'
%left '*' '/'
%left UMINUS

%%
```

# YACC code example (2/2)

```
lines   :   lines expr '\n'         {printf("%d\n", $2);}
        |   lines '\n'
        |   /* empty, i.e., epsilon */
        |   lines error '\n' { yyerror("Please reenter:"); yyerrok; }
        ;
expr    :   expr '+' expr    { $$ = $1 + $3; }
        |   expr '-' expr    { $$ = $1 - $3; }
        |   expr '*' expr    { $$ = $1 * $3; }
        |   expr '/' expr    { $$ = $1 / $3; }
        |   '(' expr ')'     { $$ = $2; }
        |   '-' expr  %prec UMINUS  { $$ = - $2; }
        |   NUMBER       { $$ = atoi(yytext);}
        ;

%%
#include "lex.yy.c"
```

# Included Lex program

```
%{
%}
Digit        [0-9]
IntLit       {Digit}+
%%
[ \t] {/* skip white spaces */}
[\n] {return('\n');}
{IntLit}                                {return(NUMBER);}
"+"                                     {return('+');}
"-"                                     {return('-');}
"*"                                     {return('*');}
"/"                                     {return('/');}
.         {printf("error token <%s>\n",yytext); return(ERROR);}
%%
```

# YACC rules

- **Can assign associativity and precedence.**
  - in increasing precedence
  - left/right or non-associativity
    - ▷ *Dot products of vectors has no associativity.*

- **Semantic rules: every item in the production is associated with a value.**
  - **YYSTYPE**: the type for return values.
  - **$$**: the return value if the production is reduced.
  - **$$i$**: the return value of the $i$th item in the production.

- **Actions can be inserted in the moddle of a production, each such action is treated as a nonterminal.**
  - **Example:**

```
expr    :   expr { $$ = 32;}  '+' expr    { $$ = $1 + $2 + $4; };

is equivalent to

expr  :   expr $ACT '+' expr {$$ = $1 + $2 + $4;};
$ACT : {$$ = 32;};
```

# YACC programming styles

- **Avoid in-production actions.**
  - Replace them by markers.
- **Keep the right hand side of a production short.**
  - Better to be less than 4 symbols.
- **Try to find some unique symbols for each production.**
  - array $\rightarrow$ ID [ elist ]
  - 
    - ▷ *arrary $\rightarrow$ aelist ]*
    - ▷ *aelist $\rightarrow$ aelist, ID*
    - ▷ *aelist $\rightarrow$ ID [ ID | ID*