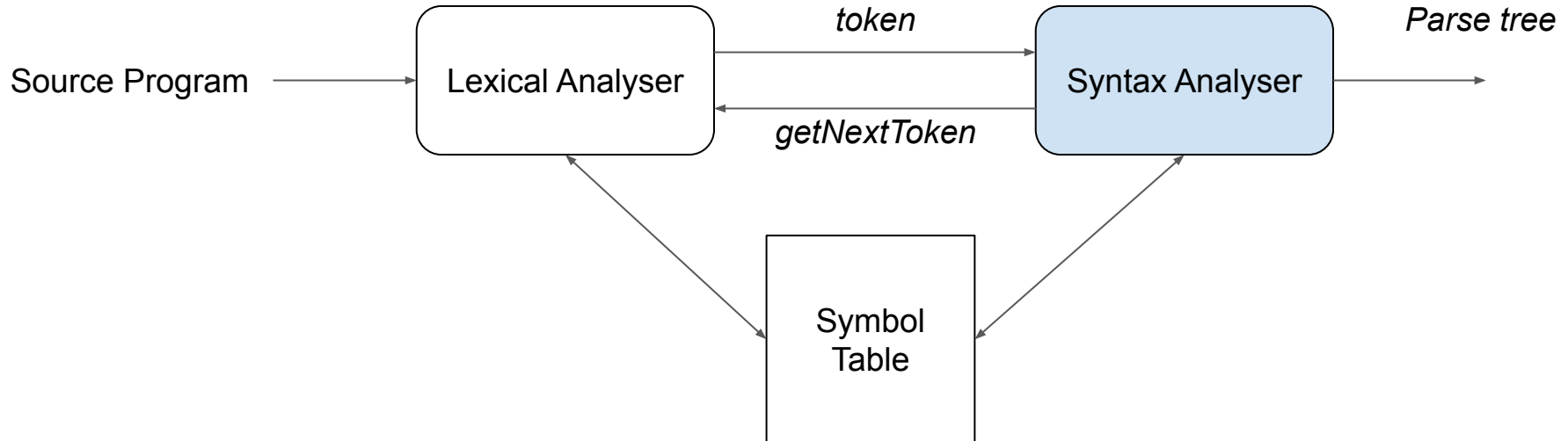


Syntax Analysis

Md Shad Akhtar
Assistant Professor
IIIT Dharwad

Syntax Analyser (Parser)

- Define the syntactic structure for a programming language
- Reads the sequence of tokens from lexical analysis and create|validate the syntactic structure (parse tree) for the sequence of tokens.



Syntactic structure and grammar

- Syntactic structure is defined by the context-free grammar (CFG)
- Steps to create parse tree
 - Parser checks whether a given source program satisfies the rules implied by a CFG or not
 - If it satisfies, the parser creates the parse tree of that program
 - Otherwise, the parser gives the error messages

Grammar (G)

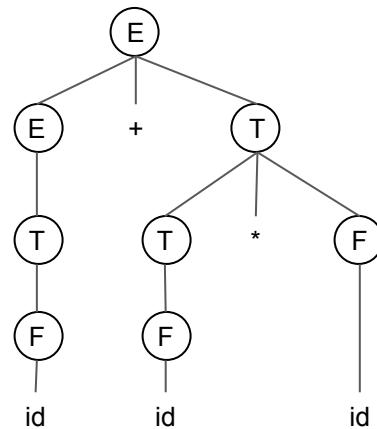
$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow (E) \mid \text{id}$

Token sequence

id + id * id



Syntax Errors

- Role of error handler in parser
 - Report the presence of errors clearly and accurately
 - Recover from each error quickly enough to detect subsequent errors
 - Add minimal overhead to the processing of correct programs.
- Error Detection:
 - Sequence of tokens that can not be accepted by any grammar rule.
 - E.g.:
 - A switch statement without a case statement
 - Missing closing braces
 - Operator without operands $c = a +$
 - Operands without operator $c = a \quad b$

Syntax Error Recovery

- Panic-mode:
 - On discovering an error, discards input symbols one at a time until one of a designated set of synchronizing tokens is found, e.g., semicolon, closing brace, etc.
- Phrase-level recovery:
 - Perform local correction on the remaining input to continue
 - Replace the prefix of the remaining input by some string that allows the parser to continue.
E.g., Replace comma by semicolon, delete|insert an extra|missing semicolon.
- Error Production
 - For common errors, add special production rules to handle such scenario
- Global correction
 - Ideally, we want as few changes as possible to process incorrect inputs.
 - We can design an algorithm for choosing a minimal sequence of changes to obtain a globally least-cost correction.
 - Given incorrect input x and grammar G , find a correct related input y with as less changes as possible.

Types of parsers

- In general three types of parsers
 - Universal
 - Capable to parse any grammar but too complex to use in compiler
 - E.g.: Cocke-Younger-Kasami (CYK) parser, Earley's parser
 - Top-down
 - Build parse tree from root to leaf
 - Bottom-Up
 - Build parse tree from leaf to root

Context-free Grammar (CFG)

- Provides a precise syntactic specification of a programming language
- A CFG $G = \langle N, T, P, S \rangle$
 - **Non-terminals:**
 - A finite set of non-terminals (variables) [usually in capital letters]
 - **Terminals:**
 - A finite set of terminals (input symbols|tokens) [usually in small letters]
 - **Production:**
 - A finite set of productions rules in the following form $A \rightarrow \alpha$ where A is a non-terminal and α is a string of terminals and non-terminals (including the empty string); $|A| \leq |\alpha|$
 - **Start symbol:**
 - One of the non-terminal symbols

CFG: An example

- CFG $G = \langle N, T, P, S \rangle$
 - **Non-terminal** = $\{E\}$
 - **Terminals** = $\{+, -, *, |, (,), \text{id}\}$
 - **Start symbol** = $\{E\}$
 - **Production**

$$E \rightarrow E + E \mid E - E \mid E * E \mid E \mid E \mid - E$$
$$E \rightarrow (E)$$
$$E \rightarrow \text{id}$$

Derivations

- Starting with the start symbol, replace each non-terminals with the body of one of its production rules till all non-terminals are replaced by terminal symbols.

- $E \Rightarrow E+E \Rightarrow id + E \Rightarrow id + id$

- In general a derivation step is

$$\alpha A \beta \Rightarrow \alpha \gamma \beta$$

if there is a production rule $A \rightarrow \gamma$ in our grammar, where α and β are arbitrary strings of terminal and non-terminal symbols.

- $\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$ (α_n derives from α_1 or α_1 derives α_n)
- \Rightarrow drives in one step
- \Rightarrow^* drives in zero or more steps
- \Rightarrow^+ drives in zero or one step

Derivations

- $S \Rightarrow^* \alpha$
 - If α contains non-terminals, it is called as a **sentential form** of G
 - If α does not contain non-terminals, it is called as a **sentence** of G
- **Left-most derivation:** Always chooses the left-most non-terminal in each derivation step

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(id+E) \Rightarrow -(id+id)$$

- **Right-most derivation:** Always chooses the right-most non-terminal in each derivation step

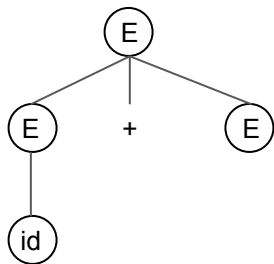
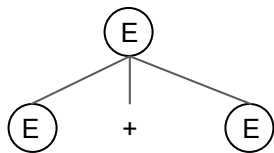
$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(E+id) \Rightarrow -(id+id)$$

- **Top-down parsers:** Finds the left-most derivation of the given source program
- **Bottom-up parsers:** Finds the right-most derivation of the given source program in the reverse order

Parse Tree

- A graphical representation of a derivation
- Intermediate nodes: Inner nodes of a parse tree
- Leaves: Terminal symbols

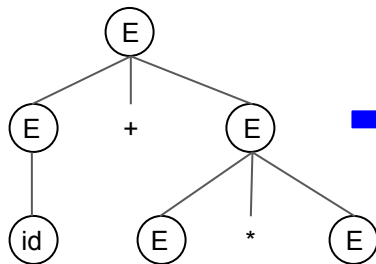
$E \Rightarrow E + E$



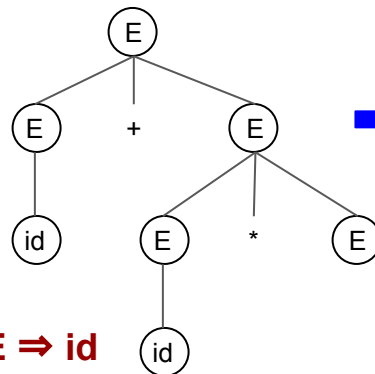
$E \Rightarrow id$



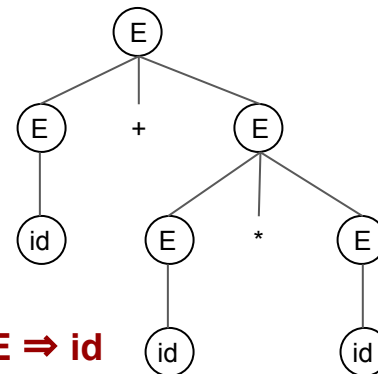
$E \Rightarrow E * E$



$E \Rightarrow id$



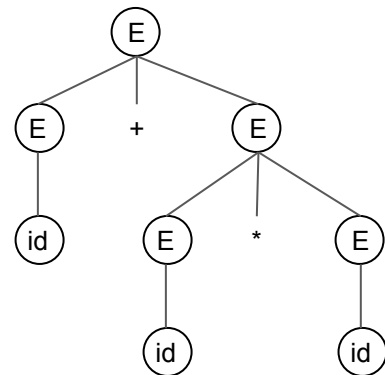
$E \Rightarrow id$



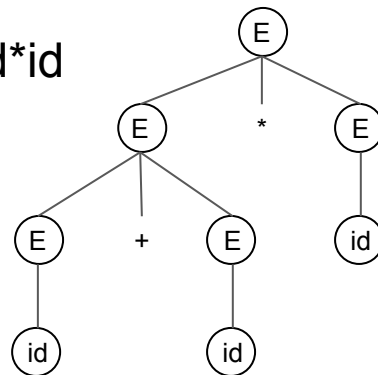
Ambiguity

- A grammar that produces more than one parse tree for a sentence is called as an ambiguous grammar

- $E \Rightarrow E + E \Rightarrow id + E \Rightarrow id + E * E \Rightarrow^* id + id * E \Rightarrow id + id * id$



- $E \Rightarrow E * E \Rightarrow E + E * E \Rightarrow id + E * E \Rightarrow id + id * E \Rightarrow id + id * id$



Ambiguity and Parser

- For the most parsers, the grammar must be unambiguous.
 - unique selection of the parse tree for a sentence
- Disambiguation of an ambiguous grammar
 - Necessary to eliminate the ambiguity in the grammar during the design phase of the compiler
 - Choose one of the parse trees of a sentence to restrict to this choice

Ambiguity disambiguation

- Stmt \rightarrow if Expr then Stmt | if Expr then Stmt else Stmt | other_stmts
- Input string: if E_1 then if E_2 then S_1 else S_2
- **Interpretation 1:** S_2 being executed when E_1 is false (thus attaching the else to the first if)
 - if E_1 then (if E_2 then S_1) else S_2
- **Interpretation 2:** S_2 being executed when E_1 is true and E_2 is false (thus attaching the else to the second if)
 - if E_1 then (if E_2 then S_1 else S_2)

Ambiguity disambiguation

- In general, we prefer the second parse tree (else matches with closest if)
- So, we have to disambiguate our grammar to reflect this choice
- Unambiguous grammar:

Stmt → matchedStmt | unmatchedStmt

matchedStmt → if Expr then matchedStmt else matchedStmt |
Otherstmts

unmatchedStmt → if Expr then Stmt |
if Expr then matchedStmt else unmatchedStmt

Ambiguity disambiguation

- Operator precedence grammar:

$$E \rightarrow E + E \mid E * E \mid E \wedge E \mid \text{id} \mid (E)$$

- Unambiguous grammar

$$E \rightarrow E + T \mid T$$
$$T \rightarrow T * F \mid F$$
$$F \rightarrow G \wedge F \mid G$$
$$G \rightarrow \text{id} \mid (E)$$

Precedence

\wedge (right to left)

$*$ (left to right)

$+$ (left to right)

Left Recursion

- A grammar is left recursive if it has a non-terminal A such that there is a derivation
 - $A \Rightarrow^+ A\alpha$ for some string α
- Top-down parsing techniques **cannot handle** left-recursive grammars
 - Conversion of left-recursive grammar into an equivalent non-recursive grammar is **mandatory**.
- Possible ways of left-recursion
 - It may appear in a single step of the derivation (immediate left-recursion)
 - It may appear in more than one step of the derivation

Removing Left Recursion

In general,

$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \dots \mid \beta_n$ Where $\beta_1 \dots \beta_n$ do
not start with A

\Downarrow

eliminate immediate left recursion

$A \rightarrow \beta_1 A' \mid \dots \mid \beta_n A'$

$A' \rightarrow \alpha_1 A' \mid \dots \mid \alpha_m A' \mid \varepsilon$

an equivalent grammar

Removing Left Recursion: An example

$$E \rightarrow E+T \mid T$$
$$T \rightarrow T*F \mid F$$
$$F \rightarrow \text{id} \mid (E)$$

↓

eliminate immediate left recursion

$$E \rightarrow T E'$$
$$E' \rightarrow +T E' \mid \varepsilon$$
$$T \rightarrow F T'$$
$$T' \rightarrow *F T' \mid \varepsilon$$
$$F \rightarrow \text{id} \mid (E)$$

Why left-recursion is a problem?

- Given

- $A \rightarrow Aa \mid b$

generate a top-down parse tree from input string 'aaaaaa'

- On first input symbol 'a', you apply first production since second production expects first character to be 'b'. [Note that you don't know what's your second input]
 - $A \Rightarrow Aa \Rightarrow Aaa \Rightarrow \dots \Rightarrow Aaaaaaa$
 - We are waiting to reduce 'A' to 'a'
 - After infinite/many steps, we may get to know that the path we chose was not correct.

Non-immediate Left-recursion

- A grammar cannot be immediately left-recursive, but it still can be left-recursive
- Just elimination of the immediate left-recursion does not guarantee a grammar which is not left-recursive

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Sc \mid d$$

- This grammar is not immediately left-recursive, but it is still left-recursive

$$S \Rightarrow Aa \Rightarrow Sca$$

Or

$$A \Rightarrow Sc \Rightarrow Aac$$

Elimination of left-recursion: Algorithm

Input: A grammar G without *e-moves* or *cycle*

Output: An equivalent grammar without left recursion

1. Arrange non-terminals in some order: $A_1 \dots A_n$
2. for $i = 1$ to n
 - a. for $j = 1$ to $i-1$
 - i. replace each production of the form

$$A_i \rightarrow A_j \gamma \quad \Rightarrow \quad A_i \rightarrow \alpha_1 \gamma \mid \alpha_2 \gamma \mid \dots \mid \alpha_k \gamma, \\ A_j \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$$

- b. eliminate the immediate left-recursions among A_i productions

If there are e-moves, the algorithm does not guarantee to work.

Elimination of left-recursion: Example

- Let grammar G:
$$\begin{aligned} S &\rightarrow Aa \mid b \\ A &\rightarrow Ac \mid Sd \mid f \end{aligned}$$
- Order of non-terminals: S, A
- For S: There is no immediate left recursion in S.
- For A: Replace $A \rightarrow Sd$ with $A \rightarrow Aad \mid bd \Rightarrow A \rightarrow Ac \mid Aad \mid bd \mid f$
Eliminate the immediate left-recursion in A

$$A \rightarrow bdA' \mid fA'$$

$$A' \rightarrow cA' \mid adA' \mid \varepsilon$$

- So, the resulting equivalent grammar which is not left-recursive is:

$$S \rightarrow Aa \mid b$$

$$A \rightarrow bdA' \mid fA'$$

$$A' \rightarrow cA' \mid adA' \mid \varepsilon$$

Elimination of left-recursion: Exercises

1. $A \rightarrow ABd \mid Aa \mid a$

$B \rightarrow Be \mid b$

2. $A \rightarrow Ba \mid Aa \mid c$

$B \rightarrow Bb \mid Ab \mid d$

3. $X \rightarrow XSb \mid Sa \mid b$

$S \rightarrow Sb \mid Xa \mid a$

Elimination of left-recursion: Solutions

1.

$$\begin{aligned}A &\rightarrow aA' \\A' &\rightarrow BdA' \mid aA' \mid \varepsilon \\B &\rightarrow bB' \\B' &\rightarrow eB' \mid \varepsilon\end{aligned}$$

2.

$$\begin{aligned}A &\rightarrow BaA' \mid cA' \\A' &\rightarrow aA' \mid \varepsilon \\B &\rightarrow cA'bB' \mid dB' \\B' &\rightarrow bB' \mid aA'bB' \mid \varepsilon\end{aligned}$$

3.

$$\begin{aligned}X &\rightarrow SaX' \mid bX' \\X' &\rightarrow SbX' \mid \varepsilon \\S &\rightarrow bX'aS' \mid aS' \\S' &\rightarrow bS' \mid aX'aS' \mid \varepsilon\end{aligned}$$

Left-factoring

- Top-down parser without backtracking (predictive parser) insists that the grammar must be left left-factored

```
stmt    →  if expr then stmt else stmt |  
          if expr then stmt
```

- After seeing `if`, we cannot decide which production rule to choose to re-write `stmt` in the derivation

Left-factoring

- In general,

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$$

where α is non-empty and the first symbols of β_1 and β_2 (if they have one) are different

- Choice involved when processing α

$$A \text{ to } \alpha\beta_1 \text{ or}$$

$$A \text{ to } \alpha\beta_2$$

- Rewrite the grammar as follows:

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2$$

so, we can immediately expand

$$A \rightarrow \alpha A'$$

Elimination of Left-factoring: Algorithm

- For each non-terminal A with two or more alternatives (production rules) with a common non-empty prefix,

$$A \rightarrow \alpha\beta_1 \mid \dots \mid \alpha\beta_n \mid \gamma_1 \mid \dots \mid \gamma_m$$

Convert it into

$$A \rightarrow \alpha A' \mid \gamma_1 \mid \dots \mid \gamma_m$$

$$A' \rightarrow \beta_1 \mid \dots \mid \beta_n$$

Elimination of Left-factoring: Example

Example 1:

$A \rightarrow abB \mid aB \mid cdg \mid cdeB \mid cdfB$

\Downarrow

$A \rightarrow aA' \mid cdg \mid cdeB \mid cdfB$

$A' \rightarrow bB \mid B$

\Downarrow

$A \rightarrow aA' \mid cdA''$

$A' \rightarrow bB \mid B$

$A'' \rightarrow g \mid eB \mid fB$

Example 2:

$A \rightarrow ad \mid a \mid ab \mid abc \mid b$

\Downarrow

$A \rightarrow aA' \mid b$

$A' \rightarrow d \mid \varepsilon \mid b \mid bc$

\Downarrow

$A \rightarrow aA' \mid b$

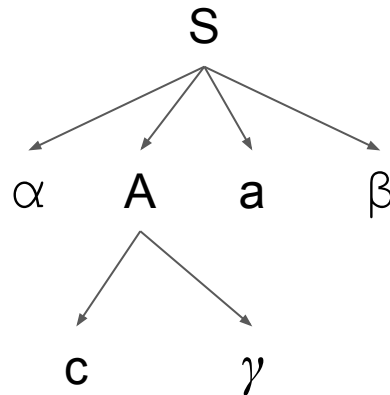
$A' \rightarrow d \mid \varepsilon \mid bA''$

$A'' \rightarrow \varepsilon \mid c$

FIRST() and FOLLOW()

- The construction of top-down and bottom-up parsing is aided by two functions on grammar G
 - $FIRST(\alpha)$: The set of *first character* that can be derived from α
 - $FOLLOW(A)$: The set of *character that can come immediately after* the non-terminal A .

$$S \rightarrow \alpha \ A \ a \ \beta$$
$$A \rightarrow c \ \gamma$$



$FIRST(a) = \{a\}$

$FIRST(A) = FIRST(c) = \{c\}$

$FIRST(S) = FIRST(\alpha) = \{\dots\}$

$FIRST(\beta) = \{\dots\}$

$FIRST(\gamma) = \{\dots\}$

$FOLLOW(A) = \{a\}$

$FOLLOW(S) = \{\$ \}$

\$: A special symbol
for the end marker.

FIRST()

- FIRST(α)

a. If α is a terminal

- $\text{FIRST}(\alpha) = \{\alpha\}$

b. If α is a non-terminal and $\alpha \rightarrow \beta_1 \beta_2 \beta_3 \dots \beta_k$

- $\text{FIRST}(\alpha) = \text{FIRST}(\alpha) \cup \text{FIRST}(\beta_i)$ if $\beta_1 \beta_2 \dots \beta_{i-1} \Rightarrow^* \varepsilon$

c. If $\alpha \rightarrow \varepsilon$

- $\text{FIRST}(\alpha) = \text{FIRST}(\alpha) \cup \varepsilon$

FOLLOW()

- FOLLOW(A)
 - a. If A is the start symbol and \$ is the special end marker
 - $\text{FOLLOW}(A) = \{\$ \}$
 - b. If $A \rightarrow \alpha B \beta$
 - $\text{FOLLOW}(B) = \text{FOLLOW}(B) \cup \{ \text{FIRST}(\beta) - \varepsilon \}$
 - c. If $A \rightarrow \alpha B$ OR $A \rightarrow \alpha B \beta$ with $\text{FIRST}(\beta)$ has ε
 - $\text{FOLLOW}(B) = \text{FOLLOW}(B) \cup \text{FOLLOW}(A)$

FIRST() and FOLLOW()

1. G: $A \rightarrow aBe \mid cBd \mid C$
 $B \rightarrow bB \mid \varepsilon$
 $C \rightarrow f$

$\text{FIRST}(a) = \{a\}$

$\text{FIRST}(b) = \{b\}$

$\text{FIRST}(c) = \{c\}$

$\text{FIRST}(d) = \{d\}$

$\text{FIRST}(e) = \{e\}$

$\text{FIRST}(f) = \{f\}$

$\text{FIRST}(A) = \{a, c, f\}$

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

$\text{FIRST}(C) = \{f\}$

$\text{FOLLOW}(A) = \{\$ \}$

$\text{FOLLOW}(B) = \{e, d\}$

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

2. G: $A \rightarrow aBc$
 $B \rightarrow bC$
 $C \rightarrow c \mid \varepsilon$

$\text{FIRST}(a) = \{a\}$

$\text{FIRST}(b) = \{b\}$

$\text{FIRST}(c) = \{c\}$

$\text{FIRST}(A) = \{a\}$

$\text{FIRST}(B) = \{b\}$

$\text{FIRST}(C) = \{c, \varepsilon\}$

$\text{FOLLOW}(A) = \{\$ \}$

$\text{FOLLOW}(B) = \{c\}$

$\text{FOLLOW}(C) = \{c\}$

FIRST() and FOLLOW()

3. G: $E \rightarrow TE'$
 $E' \rightarrow +TE' \mid \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' \mid \epsilon$
 $F \rightarrow \text{id} \mid (E)$

$\text{FIRST}(+) = \{+\}$, $\text{FIRST}(*) = \{*\}$, $\text{FIRST}(\text{id}) = \{\text{id}\}$, $\text{FIRST}('(') = \{($, $\text{FIRST}(')') = \{)\}$

$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{\text{id}, ($

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

$\text{FIRST}(E') = \{+, \epsilon\}$

$\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{), \$\}$

$\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{+,), \$\}$

$\text{FOLLOW}(F) = \{+, *,), \$\}$

Top-Down Parsing

Top-Down Parsing

- Parse tree are created top to bottom
 - Begin with the start symbol to generate the input string.
- Top-down parser
 - Recursive-Descent Parsing
 - Predictive Parsing
 - Non-recursive Predictive Parsing (LL(1) parsing)

Recursive-Descent Parsing

- Tries to find the left-most derivation
- Recursively applies production rules
- If current production fails, **backtrack**, apply another rule.
- Its simple but not widely used
- Not efficient
 - Cost of backtracking is involved which may be huge.

Designing a recursive-descent parser

- Write a procedure/function for each non-terminal.
- Call the associated function whenever a non-terminal is encountered during derivation.

```
function S()  
{  
    // match the input and/or call non-terminal functions  
    // Backtrack, if it does not apply.  
}
```

Designing a recursive-descent parser

- Design a recursive-descent parser for the following grammar.

$$S \rightarrow aBc$$
$$B \rightarrow bc \mid b$$

(Recursive) Predictive Parser

- A special form of recursive-descent parsing without backtracking.
- Since no backtracking, its efficient
- But needs a special kind of grammar, i.e., LL(1) grammar
- Uniquely choose a production rule by looking at the current symbol in the input string

Production

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$

Input:

... a ...



Current token

Predictive Parser

- **Constraints on grammar**
 - a. Unambiguous
 - b. No left recursion should be there
 - c. Grammar should be left-factored
- Still, no 100% guarantee

Predictive Parser

- Let grammar G:

$A \rightarrow aBe \mid cBd \mid C$

$B \rightarrow bB \mid \epsilon$

$C \rightarrow f$

- Left recursion?

- No left recursion

\Rightarrow This ensures that, given the current token, we don't have to backtrack

- Left factored?

- Yes

\Rightarrow This ensures that, for a input symbol, we no longer have to make a decision

Predictive Parser

- For predictive parser, we need **lookahead symbols**
 - a. Compute **FIRST()** and **FOLLOW()** for the grammar
- For a given non-terminal A and the current input symbol a
 - a. IF **FIRST(A) contains symbol a** ,
 - Apply the production associated with symbol a .
 - b. Else IF **FIRST(A) contains symbol ϵ** ,
 - IF **FOLLOW(A) contains symbol a** ,
 - Apply the production $A \rightarrow \epsilon$ and proceed.
 - c. Else
 - Error

Input string: **a b e**

Predictive Parser

- Grammar G:

$A \rightarrow aBe \mid cBd \mid C$

$B \rightarrow bB \mid \epsilon$

$C \rightarrow f$

$\text{FIRST}(A) = \{a, c, f\}$

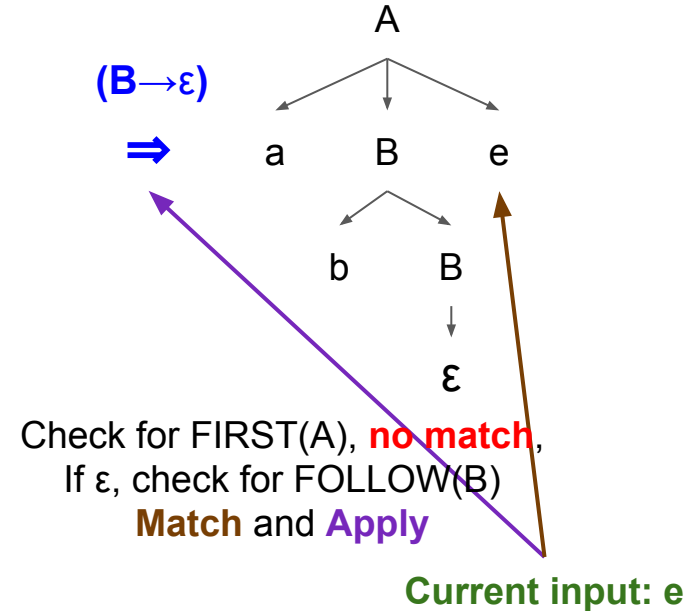
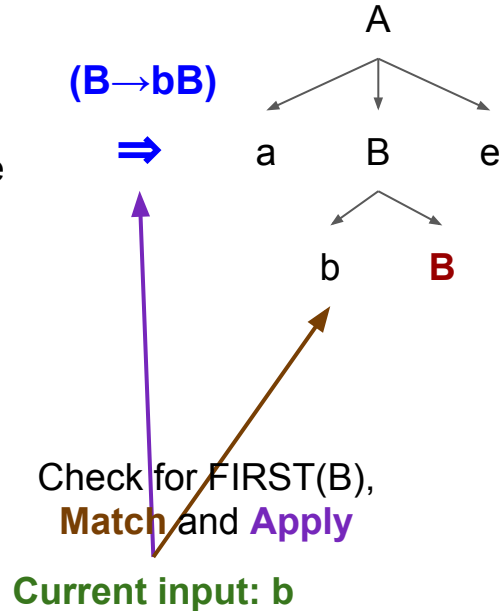
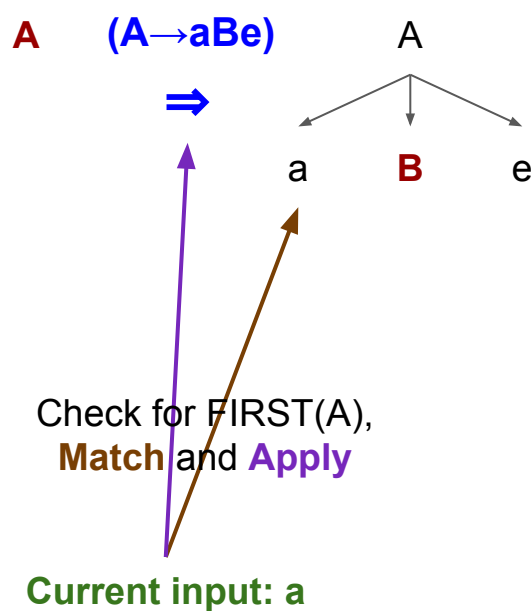
$\text{FIRST}(B) = \{b, \epsilon\}$

$\text{FIRST}(C) = \{f\}$

$\text{FOLLOW}(A) = \{\$ \}$

$\text{FOLLOW}(B) = \{e, d\}$

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



Predictive Parser

- Let grammar G:

$$S \rightarrow aBc$$

$$B \rightarrow bc \mid b$$

- Left recursion?

- No left recursion

- Left factored?

- No

$$S \rightarrow aBc$$

$$B \rightarrow bB'$$

$$B' \rightarrow c \mid \epsilon$$

Predictive Parser

- Let the new grammar G' :

$$S \rightarrow aBc$$

$$B \rightarrow bB'$$

$$B' \rightarrow c \mid \varepsilon$$

- Find FIRST and FOLLOW

$$\text{a. } \text{FIRST}(S) = \{a\} \quad \text{FIRST}(B) = \{b\} \quad \text{FIRST}(B') = \{c, \varepsilon\}$$

$$\text{b. } \text{FOLLOW}(S) = \{\$ \} \quad \text{FOLLOW}(B) = \{c\} \quad \text{FOLLOW}(B') = \{c\}$$

- Input string: **a b c**

Designing a Predictive Parser

- Write a procedure/function for each non-terminal.
- Call the associated function whenever a non-terminal is encountered during derivation.
- Match lookahead with the current input and apply the rule

```
function S(lookahead, current)
{
    // match the input and/or call non-terminal functions
}
```


Designing a Predictive Parser

- Design a recursive-descent parser for the following grammar.

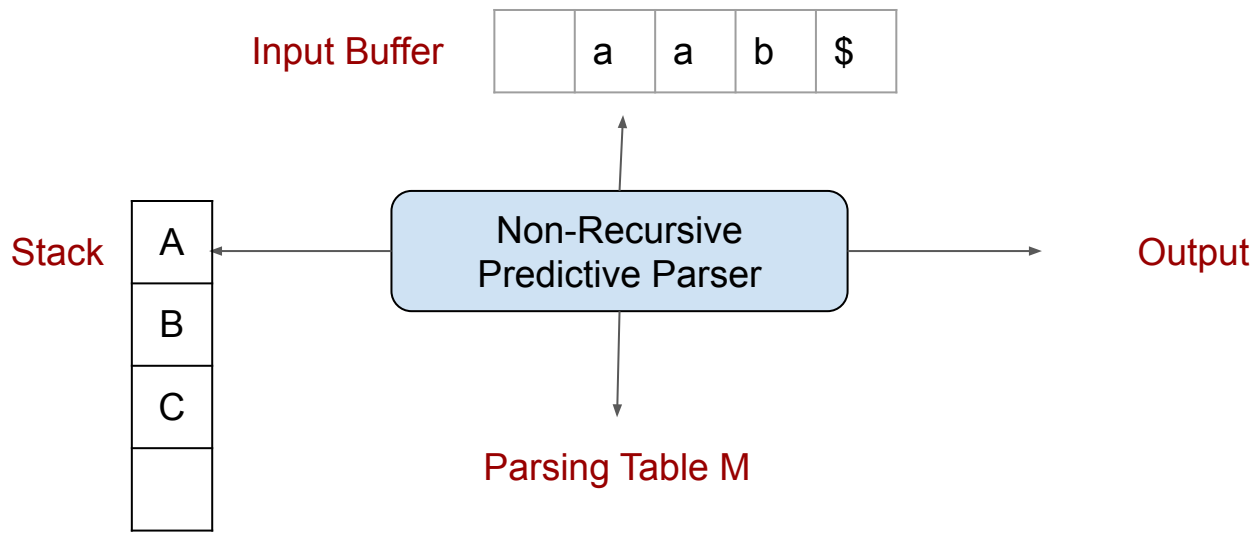
$$A \rightarrow aBe \mid cBd \mid C$$
$$B \rightarrow bB \mid \varepsilon$$
$$C \rightarrow f$$

```
proc A {  
  current token {  
    a:  
      - match the current token with a, and move to the next token;  
      - call B;  
      - match the current token with e, and move to the next token;  
    c:  
      - match the current token with c, and move to the next token;  
      - call B;  
      - match the current token with d, and move to the next token;  
    f:  
      - call C  
  }  
}  
proc B {  
  current token {  
    b:  
      - match the current token with b, and move to the next token;  
      - call B  
    e,d:  
      do nothing                //FOLLOW(B)  
  }  
}  
proc C {  
  f: - match the current token with f, and move to the next token;  
}
```

LL(1) Parser

Non-Recursive Predictive or LL(1) Parser

- Top-down parser
- Table-driven parser
- $LL(k)$, with $k = 1$
 - Left-to-right Left-most-derivation with k lookahead symbols



Non-Recursive Predictive or LL(1) Parser

- **Input buffer**
 - Contains the string to be parsed with end marked with a special symbol \$
- **Output**
 - A production rule representing a step of the derivation sequence (left-most derivation) of the string in the input buffer
- **Stack**
 - Contains the grammar symbols
 - At the bottom of the stack, there is a special end marker symbol \$
 - Initially the stack contains only the symbol \$ and the starting symbol S
 - \$S ← Initial stack
 - Parsing completes when both input and stack becomes empty (i.e., only \$ left in stack)
- **Parsing table**
 - A two-dimensional array $M[A, a]$
 - Each row is a non-terminal symbol
 - Each column is a terminal symbol or the special symbol \$
 - Entries holds a production rule.

LL(1) Grammar

- For any grammar G , if we can build an LL(1), then the grammar is called LL(1) grammar.
 - No *left-recursive*, *non-left-factored* or *ambiguous* grammar can be LL(1)
 - Still, there are some grammar which are *non-left-recursive*, *left-factored* and *unambiguous* but not a LL(1) grammar.
- A grammar G is LL(1) if and only if whenever $A \rightarrow \alpha \mid \beta$ are two distinct productions of G , the following conditions hold:
 - Both α and β cannot derive strings starting with same terminals
 - At most one of α and β can derive to ϵ
 - If β can derive to ϵ , then α cannot derive to any string starting with a terminal in FOLLOW(A) and vice-versa.

Construction of LL(1) parsing table

Input: Grammar G .

Output: Parsing Table M .

1. For each production $A \rightarrow \alpha$ of the grammar,
2. do
 - a. For each terminal a in $\text{FIRST}(\alpha)$
 - i. Add $A \rightarrow \alpha$ to $M[A, a]$
 - b. If ϵ is in $\text{FIRST}(\alpha)$
 - i. For each terminal a in $\text{FOLLOW}(A)$
 1. Add $A \rightarrow \alpha$ to $M[A, a]$
 - c. If ϵ is in $\text{FIRST}(\alpha)$ and $\$$ is in $\text{FOLLOW}(A)$
 - i. Add $A \rightarrow \alpha$ to $M[A, a]$

Construction of LL(1) parsing table

$E \rightarrow TE'$ $E' \rightarrow +TE' \mid \varepsilon$
 $T \rightarrow FT'$ $T' \rightarrow *FT' \mid \varepsilon$
 $F \rightarrow \text{id} \mid (E)$

$\text{FIRST}(T') = \{*, \varepsilon\}$

$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{\text{id}, (\}$

$\text{FIRST}(E') = \{+, \varepsilon\}$

$\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{), \$\}$

$\text{FOLLOW}(F) = \{+, *,), \$\}$

$\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{+,), \$\}$

Non-Term	Input Symbols					
	id	+	*	()	\$
E						
E'						
T						
T'						
F						

Construction of LL(1) parsing table

$E \rightarrow TE'$ $E' \rightarrow +TE' \mid \varepsilon$
 $T \rightarrow FT'$ $T' \rightarrow *FT' \mid \varepsilon$
 $F \rightarrow \text{id} \mid (E)$

$\text{FIRST}(T') = \{*, \varepsilon\}$

$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{\text{id}, (\}$

$\text{FIRST}(E') = \{+, \varepsilon\}$

$\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{), \$\}$

$\text{FOLLOW}(F) = \{+, *,), \$\}$

$\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{+,), \$\}$

Production $E \rightarrow TE' \Rightarrow \text{First}(TE') = \text{First}(T) = \{ (, \text{id} \} \Rightarrow \text{Add } E \rightarrow TE' \text{ to}$

$M[E, \text{id}]$ and $M[E, (]$

Non-Term	Input Symbols					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'						
T						
T'						
F						

Construction of LL(1) parsing table

$E \rightarrow TE'$
 $T \rightarrow FT'$
 $F \rightarrow id \mid (E)$

$E' \rightarrow +TE' \mid \epsilon$
 $T' \rightarrow *FT' \mid \epsilon$

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

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

$FOLLOW(F) = \{+, *,), \$\}$

$FIRST(E) = FIRST(T) = FIRST(F) = \{id, (\}$

$FOLLOW(E) = FOLLOW(E') = \{), \$\}$

$FOLLOW(T) = FOLLOW(T') = \{+,), \$\}$

Production $T \rightarrow FT'$ $\Rightarrow First(FT') = First(F) = \{ (, id \}$ \Rightarrow Add $T \rightarrow FT'$ to

$M[T, id]$ and $M[T, (]$

Non-Term	Input Symbols					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'						
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'						
F						

Construction of LL(1) parsing table

$E \rightarrow TE'$ $E' \rightarrow +TE' \mid \varepsilon$
 $T \rightarrow FT'$ $T' \rightarrow *FT' \mid \varepsilon$
 $F \rightarrow id \mid (E)$

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

$FIRST(E) = FIRST(T) = FIRST(F) = \{id, (\}$

$FIRST(E') = \{+, \varepsilon\}$

$FOLLOW(E) = FOLLOW(E') = \{), \$\}$

$FOLLOW(F) = \{+, *,), \$\}$

$FOLLOW(T) = FOLLOW(T') = \{+,), \$\}$

Production $F \rightarrow id \Rightarrow First(id) = \{id\}$

\Rightarrow Add $F \rightarrow id$ to

$M[F, id]$

Non-Term	Input Symbols					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'						
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'						
F	$F \rightarrow id$					

Construction of LL(1) parsing table

$E \rightarrow TE'$ $E' \rightarrow +TE' \mid \varepsilon$
 $T \rightarrow FT'$ $T' \rightarrow *FT' \mid \varepsilon$
 $F \rightarrow id \mid (E)$

$FIRST(T') = \{*, \varepsilon\}$ $FIRST(E) = FIRST(T) = FIRST(F) = \{id, (\}$
 $FIRST(E') = \{+, \varepsilon\}$ $FOLLOW(E) = FOLLOW(E') = \{), \$\}$
 $FOLLOW(F) = \{+, *,), \$\}$ $FOLLOW(T) = FOLLOW(T') = \{+,), \$\}$

Production $F \rightarrow (E) \Rightarrow First((E)) = \{ (\}$

\Rightarrow Add $F \rightarrow (E)$ to
 $M[F, (]$

Non-Term	Input Symbols					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'						
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'						
F	$F \rightarrow id$			$F \rightarrow (E)$		

Construction of LL(1) parsing table

$E \rightarrow TE'$ $E' \rightarrow +TE' \mid \varepsilon$
 $T \rightarrow FT'$ $T' \rightarrow *FT' \mid \varepsilon$
 $F \rightarrow id \mid (E)$

$FIRST(T') = \{*, \varepsilon\}$ $FIRST(E) = FIRST(T) = FIRST(F) = \{id, (\}$
 $FIRST(E') = \{+, \varepsilon\}$ $FOLLOW(E) = FOLLOW(E') = \{), \$\}$
 $FOLLOW(F) = \{+, *,), \$\}$ $FOLLOW(T) = FOLLOW(T') = \{+,), \$\}$

Production $E' \rightarrow +TE'$ $\Rightarrow First(+TE') = \{ + \}$

\Rightarrow Add $E' \rightarrow +TE'$ to
 $M[E', +]$

Non-Term	Input Symbols					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$				
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'						
F	$F \rightarrow id$			$F \rightarrow (E)$		

Construction of LL(1) parsing table

$E \rightarrow TE'$
 $T \rightarrow FT'$
 $F \rightarrow id \mid (E)$

$E' \rightarrow +TE' \mid \epsilon$
 $T' \rightarrow *FT' \mid \epsilon$

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

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

$FOLLOW(F) = \{+, *,), \$\}$

$FIRST(E) = FIRST(T) = FIRST(F) = \{id, (\}$

$FOLLOW(E) = FOLLOW(E') = \{), \$\}$

$FOLLOW(T) = FOLLOW(T') = \{+,), \$\}$

Production $T' \rightarrow *FT' \Rightarrow First(*FT') = \{ * \}$

\Rightarrow Add $T' \rightarrow *FT'$ to

$M[T', *]$

Non-Term	Input Symbols					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$				
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'			$T' \rightarrow *FT'$			
F	$F \rightarrow id$			$F \rightarrow (E)$		

Construction of LL(1) parsing table

$E \rightarrow TE'$
 $T \rightarrow FT'$
 $F \rightarrow id \mid (E)$

$E' \rightarrow +TE' \mid \varepsilon$
 $T' \rightarrow *FT' \mid \varepsilon$

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

$FIRST(E') = \{+, \varepsilon\}$

$FOLLOW(F) = \{+, *,), \$\}$

$FIRST(E) = FIRST(T) = FIRST(F) = \{id, (\}$

$FOLLOW(E) = FOLLOW(E') = \{), \$\}$

$FOLLOW(T) = FOLLOW(T') = \{+,), \$\}$

Production $E' \rightarrow \varepsilon \Rightarrow Follow(E') = \{), \$\}$

\Rightarrow Add $E' \rightarrow \varepsilon$ to

$M[E',)]$ and $M[E', \$]$

Non-Term	Input Symbols					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \varepsilon$	$E' \rightarrow \varepsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'			$T' \rightarrow *FT'$			
F	$F \rightarrow id$			$F \rightarrow (E)$		

Construction of LL(1) parsing table

$E \rightarrow TE'$
 $T \rightarrow FT'$
 $F \rightarrow id \mid (E)$

$E' \rightarrow +TE' \mid \varepsilon$
 $T' \rightarrow *FT' \mid \varepsilon$

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

$FIRST(E') = \{+, \varepsilon\}$

$FOLLOW(F) = \{+, *,), \$\}$

$FIRST(E) = FIRST(T) = FIRST(F) = \{id, (\}$

$FOLLOW(E) = FOLLOW(E') = \{), \$\}$

$FOLLOW(T) = FOLLOW(T') = \{+,), \$\}$

Production $T' \rightarrow \varepsilon \Rightarrow Follow(T') = \{+,), \$\}$

\Rightarrow Add $T' \rightarrow \varepsilon$ to

$M[T', +], M[T',)]$ and $M[T', \$]$

Non-Term	Input Symbols					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \varepsilon$	$E' \rightarrow \varepsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \varepsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \varepsilon$	$T' \rightarrow \varepsilon$
F	$F \rightarrow id$			$F \rightarrow (E)$		

LL(1) parsing

- Once we have built a parsing table M , verify whether a given string w is part of the language or not.
- **LL(1) parsing algorithm**
 - Look at the symbol at the top of the stack (e.g., X) and the current symbol in the input string (e.g., a)
 - If $X == a == \$$ → Halt with success;
 - If $X == a \neq \$$ → Pop; Move to next input;
 - If $X == \text{terminal OR } M[X, a] \text{ is empty}$ → Halt with error;
 - If X in non-terminal and $M[X, a] == X \rightarrow \alpha_1 \alpha_2 \alpha_3 \dots \alpha_k$
 - Pop
 - Push $\alpha_k \alpha_{k-1} \alpha_{k-2} \dots \alpha_1$ [top of stack = α_1]
 - Output the production $X \rightarrow \alpha_1 \alpha_2 \alpha_3 \dots \alpha_k$

LL(1) parsing

Input: **id + id * id**

Matched

id
id
id
id +
id +
id +
id + id
id + id
id + id *
id + id *
id + id * id
id + id * id
id + id * id

STACK

E \$
T E' \$
F T' E' \$
id T' E' \$
T' E' \$
E' \$
+ T E' \$
T E' \$
F T' E' \$
id T' E' \$
T' E' \$
***** F T' E' \$
F T' E' \$
id T' E' \$
T' E' \$
E' \$
\$

Stack: Empty
Input: Empty
Accept and stop!!

Input

id + id * id \$
id + id * id \$
id + id * id \$
id + id * id \$
+ id * id \$
+ id * id \$
+ id * id \$
id * id \$
id * id \$
id * id \$
* id \$
***** id \$
id \$
id \$
\$
\$
\$

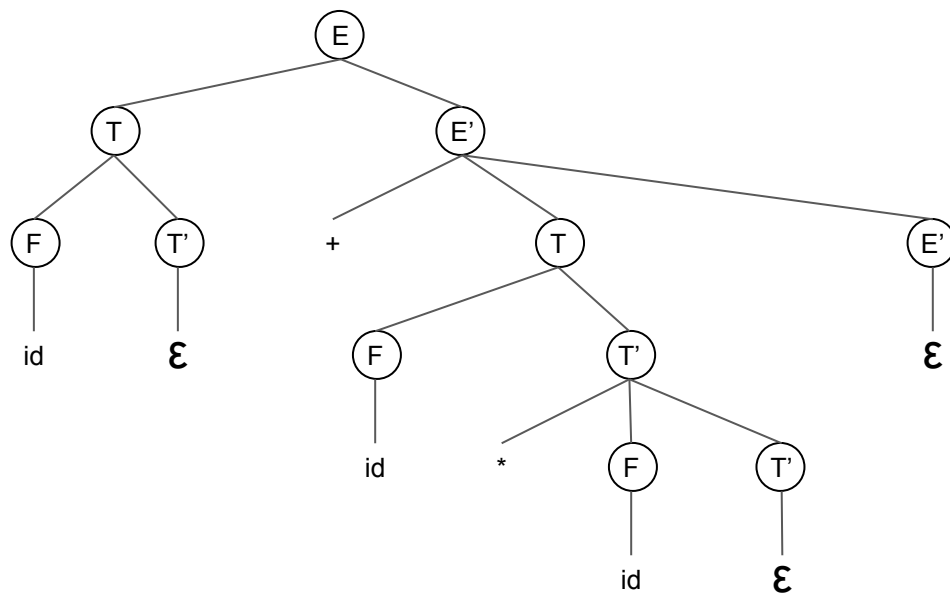
Output

Output: $E \rightarrow T E'$
Output: $T \rightarrow F T'$
Output: $F \rightarrow \text{id}$
Match: **id**
Output: $T' \rightarrow \epsilon$
Output: $E' \rightarrow + T E'$
Match: **+**
Output: $T \rightarrow F T'$
Output: $F \rightarrow \text{id}$
Match: **id**
Output: $T' \rightarrow * F T'$
Match: *****
Output: $F \rightarrow \text{id}$
Match: **id**
Output: $T' \rightarrow \epsilon$
Output: $E' \rightarrow \epsilon$

LL(1) parsing

- Following the output sequence gives you left-most derivation for the input

$E \rightarrow T E'$
 $\rightarrow F T' E'$
 $\rightarrow id T' E'$
 $\rightarrow id \epsilon E'$
 $\rightarrow id + T E'$
 $\rightarrow id + F T' E'$
 $\rightarrow id + id T' E'$
 $\rightarrow id + id * F T' E'$
 $\rightarrow id + id * id T' E'$
 $\rightarrow id + id * id \epsilon E'$
 $\rightarrow id + id * id \epsilon$



LL(1) parsing

- Construct parsing table for the following grammar:

$$S \rightarrow i E t S \mid i E t S e S \mid a$$
$$E \rightarrow b$$

[S' on e] = ??

Non-Term	Input Symbols					
	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow i E t S S'$		
S'			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

Error Recovery in LL(1) Parsing

- An error may occur in the predictive parsing (LL(1) parsing), if
 - The terminal symbol on the top of stack does not match with the current input symbol
 - top of stack is a non-terminal A , the current input symbol is a , and the parsing table entry $M[A, a]$ is empty.

Panic-mode Error Recovery in LL(1) Parsing

- Skip over the symbols on the input until a synchronization [sync] token is found.
- **Synchronization tokens**
 - Place all the symbols in the FOLLOW(A) into the synchronizing token set for the non-terminal A.
- If a non-terminal A can generate ϵ , then $A \rightarrow \epsilon$ can be used as default choice.
- If a terminal on the top of stack cannot be matched, pop the terminal.
- If a non-terminal on the top of stack has an entry sync on a terminal a, skip the terminal.

Modified LL(1) parsing table with “sync” tokens

$E \rightarrow TE'$ $E' \rightarrow +TE' \mid \varepsilon$
 $T \rightarrow FT'$ $T' \rightarrow *FT' \mid \varepsilon$
 $F \rightarrow \text{id} \mid (E)$

$\text{FIRST}(T') = \{*, \varepsilon\}$ $\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{\text{id}, (\}$
 $\text{FIRST}(E') = \{+, \varepsilon\}$ $\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{), \$\}$
 $\text{FOLLOW}(F) = \{+, *,), \$\}$ $\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{+,), \$\}$

Non-Term	Input Symbols					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$	sync	sync
E'		$E' \rightarrow +TE'$			$E' \rightarrow \varepsilon$	$E' \rightarrow \varepsilon$
T	$T \rightarrow FT'$	sync		$T \rightarrow FT'$	sync	sync
T'		$T' \rightarrow \varepsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \varepsilon$	$T' \rightarrow \varepsilon$
F	$F \rightarrow \text{id}$	sync	sync	$F \rightarrow (E)$	sync	sync

Panic-mode Recovery in LL(1) parsing

Input: **) id * + id**

STACK

E \$
E \$
T E' \$
F T' E' \$
id T' E' \$
T' E' \$
* F T' E' \$
F T' E' \$
F T E' \$
id T' E' \$
T' E' \$
E' \$
\$

Input

) id * + id \$
id * + id \$
id * + id \$
id * + id \$
id * + id \$
* + id \$
***** + id \$
+ id \$
id \$
id \$
\$
\$
\$

Remarks

Error, $M[E,)]$ = sync, skip)
id is in FIRST(E)

Error, $M[F, +]$ = sync, skip +
id is in FIRST(F)

Bottom-Up Parsing

Bottom-Up parsing

- Construct a parse tree for an input string beginning at the leaves (the bottom) and working up towards the root (the top)
 - Reducing a string w to the start symbol of a grammar.
 - At each reduction step a particular substring matching the right side of a production is replaced by the symbol on the left of that production.
 - Gives the right-most derivation in the reverse order.

Bottom-Up parsing: An example

G: $S \rightarrow aABe$
 $A \rightarrow Abc \mid b$
 $B \rightarrow d$

Input: **abbcde**

- Procedure
 - Scan the string from left to right looking for a substring that matches the right side of a production: **b and d qualifies**
 - Choose *leftmost* b and apply $A \rightarrow b$, So string becomes aAbcde
 - Scan left to right: **Abc, b and d qualifies**
 - Choose *leftmost* Abc and apply $A \rightarrow Abc$, So string becomes aAde
 - Scan left to right: **d qualifies**
 - Apply $B \rightarrow d$, so the string becomes aABe
 - Scan left to right: **aABe qualifies**
 - Apply $S \rightarrow aABe$
- **abbcde** \Rightarrow **aAbcde** \Rightarrow **aAde** \Rightarrow **aABe** \Rightarrow S

Right-most derivation

Handle

- Handle of a string is a substring that
 - *matches the right side of a production rule*; and
 - whose *reduction* to the nonterminal on the left side of the production represents *one step along the reverse of a rightmost derivation*;
- Therefore, *not every substring (or more specifically, the leftmost substring)* that matches the right side of a production rule is *handle*.

E.g.:

G: $E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow \text{id} \mid (E)$

Input: $\text{id}_1 * \text{id}_2$

$\text{id}_1 * \text{id}_2$

$\Rightarrow F * \text{id}_2$

$\Rightarrow T * \text{id}_2$

matched substring{ id_1, id_2 }

matched substring{ F, id_2 }

matched substring{ T, id_2 }

In the next step, shall we reduce the leftmost substring $E \rightarrow T$ or $F \rightarrow \text{id}_2$?

Shift-Reduce Parsing

- A stack implementation of bottom-up parsing
 - **Shift** → Current input symbol is pushed onto the stack
 - **Reduce** → Right side of a production is replaced by the left side non-terminal in the stack.
- Shift zero or more input symbols onto the stack, until it is ready to reduce a string α to a non-terminal A on top of the stack, if the grammar has production $A \rightarrow \alpha$.
- Repeat the process, until
 - It generate an error signal OR
 - Stack contains the start symbol and input is empty. Accept the input.

Shift-Reduce Parsing

G: $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow \text{id} \mid (E)$

Input: $\text{id}_1 * \text{id}_2$

Stack	Input	Action
\$	$\text{id}_1 * \text{id}_2$ \$	Shift
\$ id_1	$*$ id_2 \$	Reduce by $F \rightarrow \text{id}$
\$ F	$*$ id_2 \$	Reduce by $T \rightarrow F$
\$ T	$*$ id_2 \$	Shift
\$ $T *$	id_2 \$	Shift
\$ $T * \text{id}_2$	\$	Reduce by $F \rightarrow \text{id}$
\$ $T * F$	\$	Reduce by $T \rightarrow T * F$
\$ T	\$	Reduce by $E \rightarrow T$
\$ E	\$	Accept

Observe, handle is always at the top of stack.

Shift-Reduce Parsing: Few key points

- Four primary operations
 - **Shift** → Current input symbol is pushed onto the stack
 - **Reduce** → Right side of a production is replaced by the left side non-terminal on the stack.
 - **Accept** → Announce successful completion of parsing
 - **Error** → Discover a syntax error and call error handling mechanism.
- Handle always appear on top of the stack
- For an unambiguous grammar, for every right-sentential form there is exactly one handle.
 - Remember given $S \Rightarrow^* \alpha$,
 - If α contains non-terminals, it is called as a sentential form of G

Conflicts

- There are grammars for which shift-reduce parsing cannot be used.
- Shift-Reduce parser for such grammars may have a configuration where the parser cannot decide whether to
 - Shift the symbols onto the stack or Reduce the handle to a non-terminal
 - OR
 - Reduce the handle with some non-terminal A or B.
- These situations are called **conflicts**.
 - *Shift/Reduce* conflict
 - *Reduce/Reduce* conflict

Conflicts: Example 1

- Ambiguous grammar can not have shift-reduce parser

- $\text{stmt} \rightarrow \text{if expr then stmt} \mid$
 $\text{if expr then stmt else stmt} \mid$
 other

- Let the configuration of parser is

Stack

\$... if expr then stmt

Reduce

stmt

Shift

Input

else \$

Shift/Reduce
conflict



Conflicts: Example 2

- stmt → id (param_list) | expr = expr
param_list → param_list, param | param
param → id
expr → id (expr_list) | id
expr_list → expr_list, expr | expr

- Let the configuration is

Stack

\$.... id (id



Input

, id) \$

Reduce with param → id

OR

Reduce with expr → id

LR Parsers

LR Parser

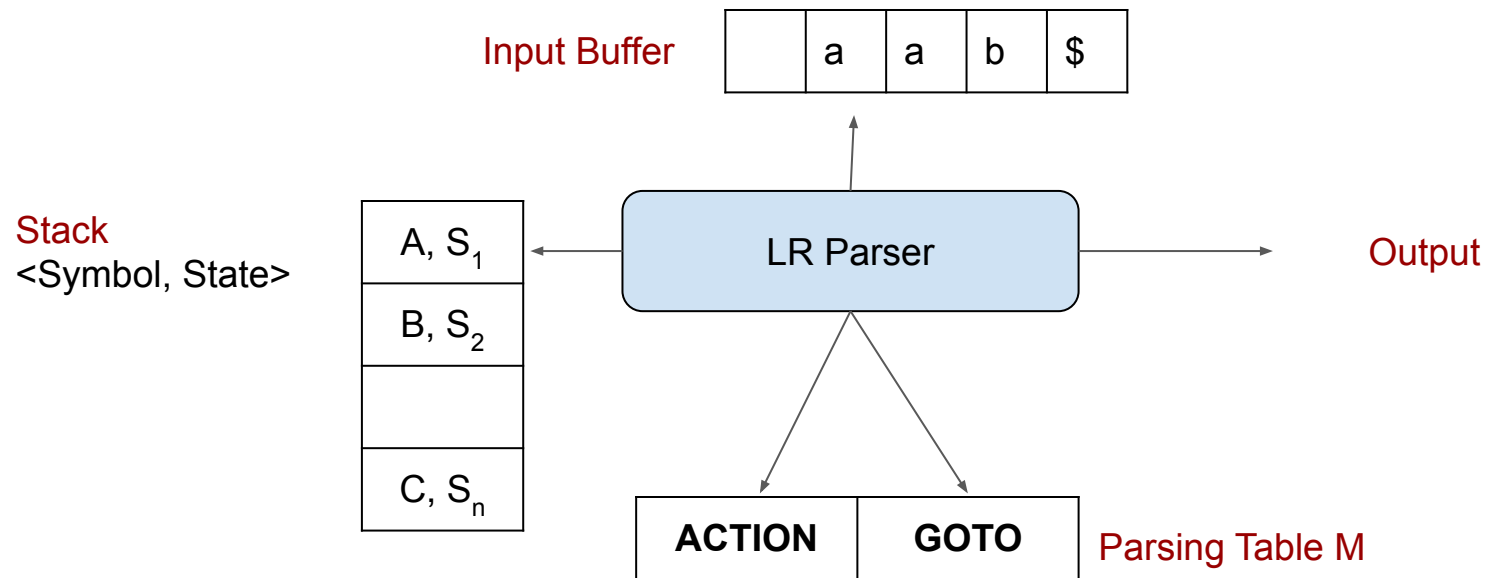
- LR(k) parsers are the most powerful and efficient shift-reduce parser
 - Left-to-right scanning, Right-most derivation (with k lookahead symbols)
 - In general, $k = 1$
 - In both LL(k) and LR(k), if k is omitted, it is assumed LL(1) and LR(1)
- A grammar for which we can construct a LR parser are called LR grammar
- Three main types of parse
 - Simple LR or SLR or LR(0)
 - Canonical LR or LR(1)
 - Look-ahead LR or LALR
- Parsing of all three parsers are similar, only their parsing tables are different

Why LR parsers?

- LR parsers can be constructed to recognize virtually all programming language constructs for which CFGs can be written.
- LR parsers are most general non-backtracking shift-reduce parser and yet its implementation is as efficient as others.
- An LR parser can detect a syntactic error as soon as it is possible to do on a left-to-right scan of the input
- Class of grammars that can be parsed using LR methods is a proper superset of the class of grammars that can be parsed with predictive parsers or LL methods

$LL(1) \text{ grammars} \subset LR(1) \text{ grammars}$

LR parsing



Configuration of LR parsing

- Each symbol on stack has an associated state.
- Initial stack configuration $\$ S_0$ (no symbol is associated with S_0)

$$(\$ S_0 X_1 S_1 \dots X_m \mathbf{S}_m, \quad \mathbf{a}_i a_{i+1} \dots a_n \$)$$

Stack

Input

- \mathbf{S}_m and \mathbf{a}_i decides the next parser action by consulting the parsing table M.

Configuration of LR parsing

- S_m and a_i decides the next parser action by consulting the parsing table M.

- **Shift:**

- Push a_i and its associated state S_i onto the stack

$$(\$S_0X_1S_1 \dots X_mS_m, \quad a_i a_{i+1} \dots a_n \$) \rightarrow (\$S_0X_1S_1 \dots X_mS_m a_i S_i, \quad a_{i+1} \dots a_n \$)$$

- **Reduce:**

- If $A \rightarrow X_{m-r-1}S_{m-r-1} \dots X_mS_m$ is a handle
 - Pop $r = |X_{m-r-1}S_{m-r-1} \dots X_mS_m|$ items from the stack
 - Push A and S onto the stack, where $S = \text{GOTO}[S_{m-r}, A]$

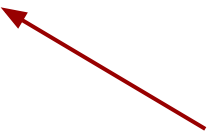
$$(\$S_0X_1S_1 \dots X_mS_m, \quad a_i a_{i+1} \dots a_n \$) \rightarrow (\$S_0X_1S_1 \dots X_{m-r}S_{m-r} A S, \quad a_i a_{i+1} \dots a_n \$)$$

Construction of LR(k) parser

- Building a parser
 1. Build LR(k) automation
 - a. Canonical set of “items”
 2. Build parsing table using LR(k) automation
- Once the parsing table is built, we can parse any given input string using LR(k) parsing algorithm.

Building LR(0) parser: Canonical set of items

Canonical set of “items” for LR(0) automation

- LR parser makes shift-reduce decision based on the states in an automation.
 - Each state contains a set of items that reflects the progress in parsing.
 - Collection of sets of LR(0) items are called canonical LR(0) collection.
 - An LR(0) item (or simply item) of a grammar G is a production with a dot (\cdot) at some position of the right side of the rule.
 - For the production $A \rightarrow XYZ$, we have four items
 - $A \rightarrow \cdot XYZ$
 - $A \rightarrow X \cdot YZ$
 - $A \rightarrow XY \cdot Z$
 - $A \rightarrow XYZ \cdot$
- The position of \cdot indicates the amount of processing completed.
- 
1. Parser has PROCESSED X on a portion of the input; and
 2. HOPE to derive the rest of the input from YZ

(Dot) Closure of items

- To build the LR(0) automation, we need to find the closure of each item set(I)

- Let the grammar G:
$$E \rightarrow E + T \mid T$$
$$T \rightarrow T * F \mid F$$
$$F \rightarrow (E) \mid id$$

- Then, the dot closure of item $E \rightarrow \cdot E + T$ is

$E \rightarrow \cdot E + T$

$E \rightarrow \cdot T$

$T \rightarrow \cdot T * F$

$T \rightarrow \cdot F$

$F \rightarrow \cdot (E)$

$F \rightarrow \cdot id$

Closure($E \rightarrow \cdot E + T$)

(Dot) Closure of items

Closure (I)

1. Add every item in I to **Closure (I)**
2. If $A \rightarrow \alpha \cdot B \beta$ is in **Closure(I)** and $B \rightarrow \gamma$ is a production
 - a. Add item $B \rightarrow \cdot \gamma$ to **Closure (I)**
3. Repeat step 2, until no new items can be added to **Closure (I)**.

Transition function GOTO()

- If $\text{Closure}(I)$ has an item $A \rightarrow \alpha \cdot B \beta$
 - $\text{GOTO}(I, B) = \text{Closure}(A \rightarrow \alpha B \cdot \beta)$
- Let $\text{Closure}(I) = \{[E \rightarrow \cdot T], [E \rightarrow E \cdot + T]\}$
 - $\text{GOTO}(I, +) = \{ [E \rightarrow E + \cdot T],$
 $[T \rightarrow \cdot T * F]$
 $[T \rightarrow \cdot F]$
 $[F \rightarrow \cdot (E)]$
 $[F \rightarrow \cdot id] \}$

LR(0) Automation

- The state of the automation is defined by the $\text{Closure}(I)$ of items
- The $\text{GOTO}(I, X)$ function defines the transition from state I on symbol X
- For every grammar, augment a production $S' \rightarrow S$, if S was the starting symbol.
 - S' becomes new start symbol
 - $S' \rightarrow S.$ signifies the acceptance of the input.

Computation of the canonical LR(0) collection

Items(G')

1. $C = \text{Closure}(\{[S' \rightarrow \cdot S]\})$
2. Repeat
 - a. For each set of items I in C
 - i. For each grammar symbol X
 1. If $\text{GOTO}(I, X)$ is not empty and not in C
 - a. Add $\text{GOTO}(I, X)$ to C
3. Until no new sets of items are added to C

LR(0) Automation

0: $E' \rightarrow E$

1: $E \rightarrow E + T$

2: $E \rightarrow T$

3: $T \rightarrow T * F$

4: $T \rightarrow F$

5: $F \rightarrow (E)$

6: $F \rightarrow id$

$E' \rightarrow \cdot E$

$E \rightarrow \cdot E + T$

$E \rightarrow \cdot T$

$T \rightarrow \cdot T * F$

$T \rightarrow \cdot F$

$F \rightarrow \cdot (E)$

$F \rightarrow \cdot id$

LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

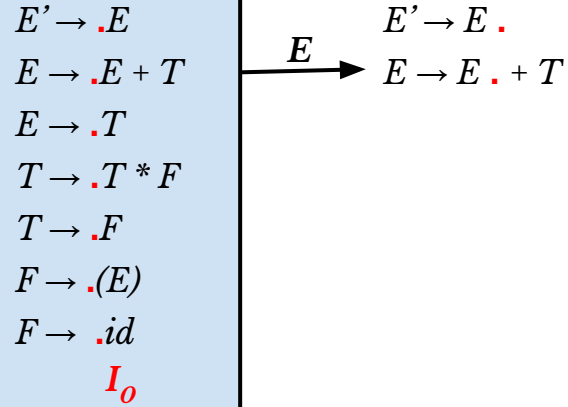
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

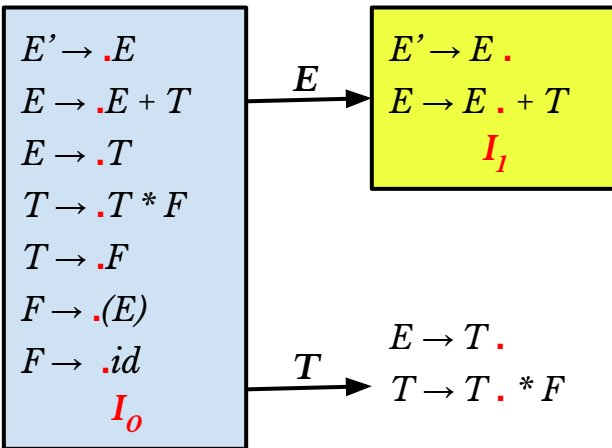
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

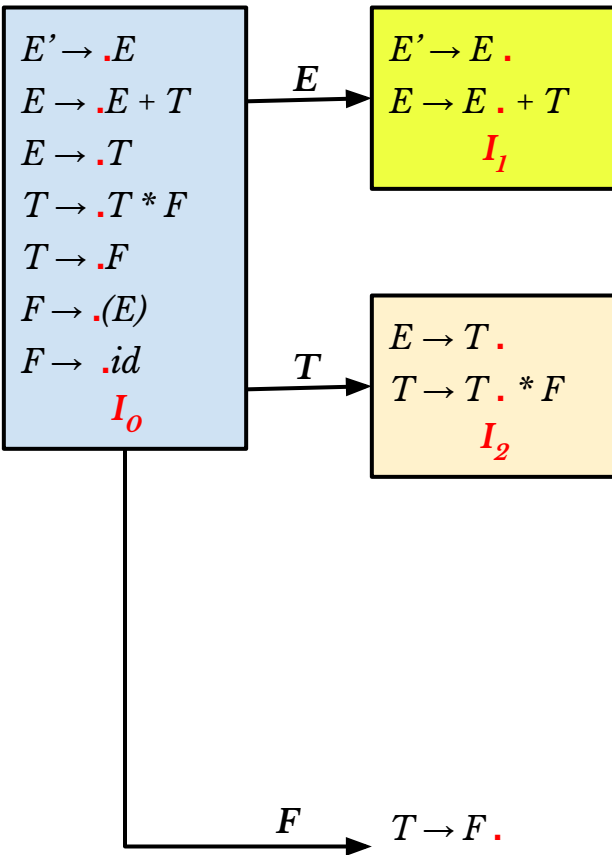
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

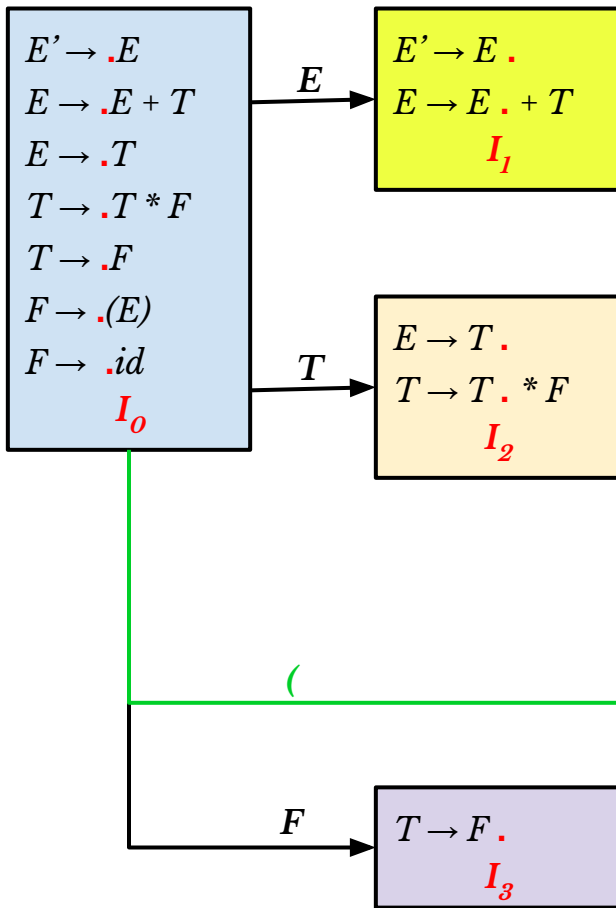
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



$F \rightarrow (\cdot E)$
 $E \rightarrow \cdot E + T$
 $E \rightarrow \cdot T$
 $T \rightarrow \cdot T * F$
 $T \rightarrow \cdot F$
 $F \rightarrow \cdot (E)$
 $F \rightarrow \cdot id$

LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

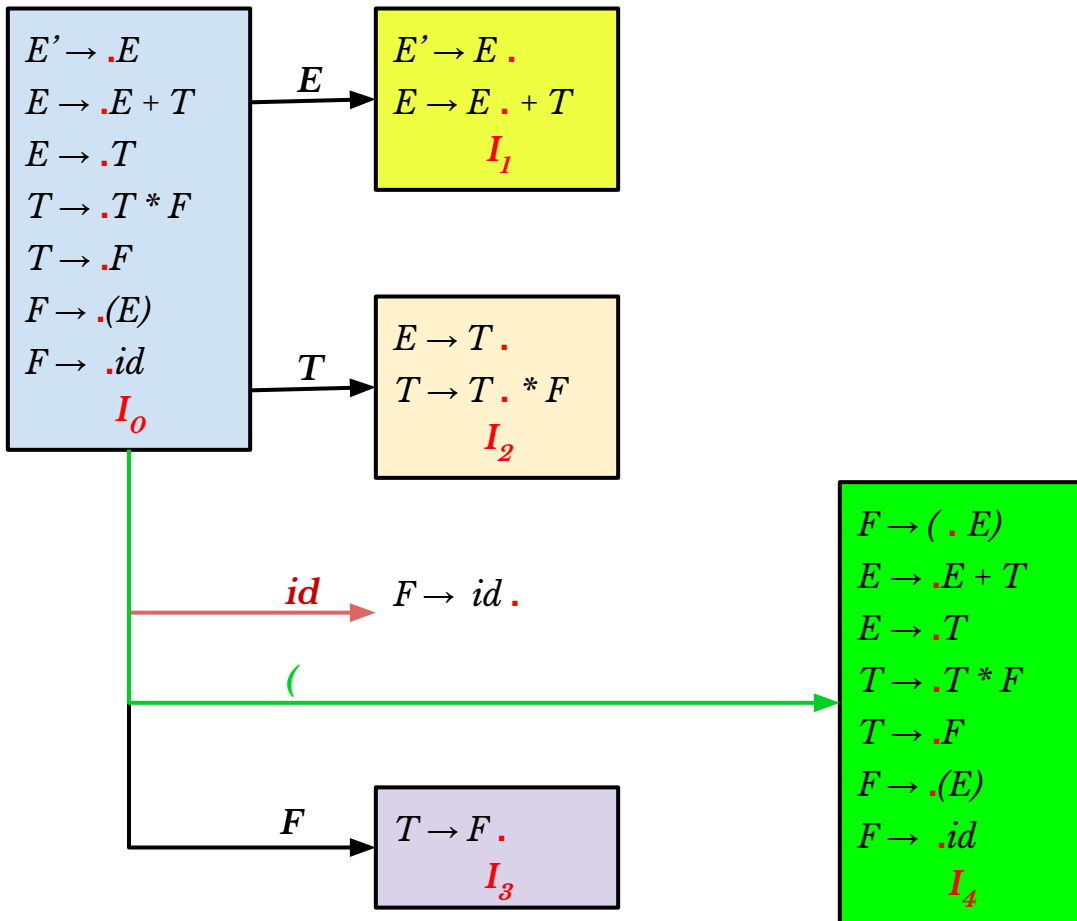
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

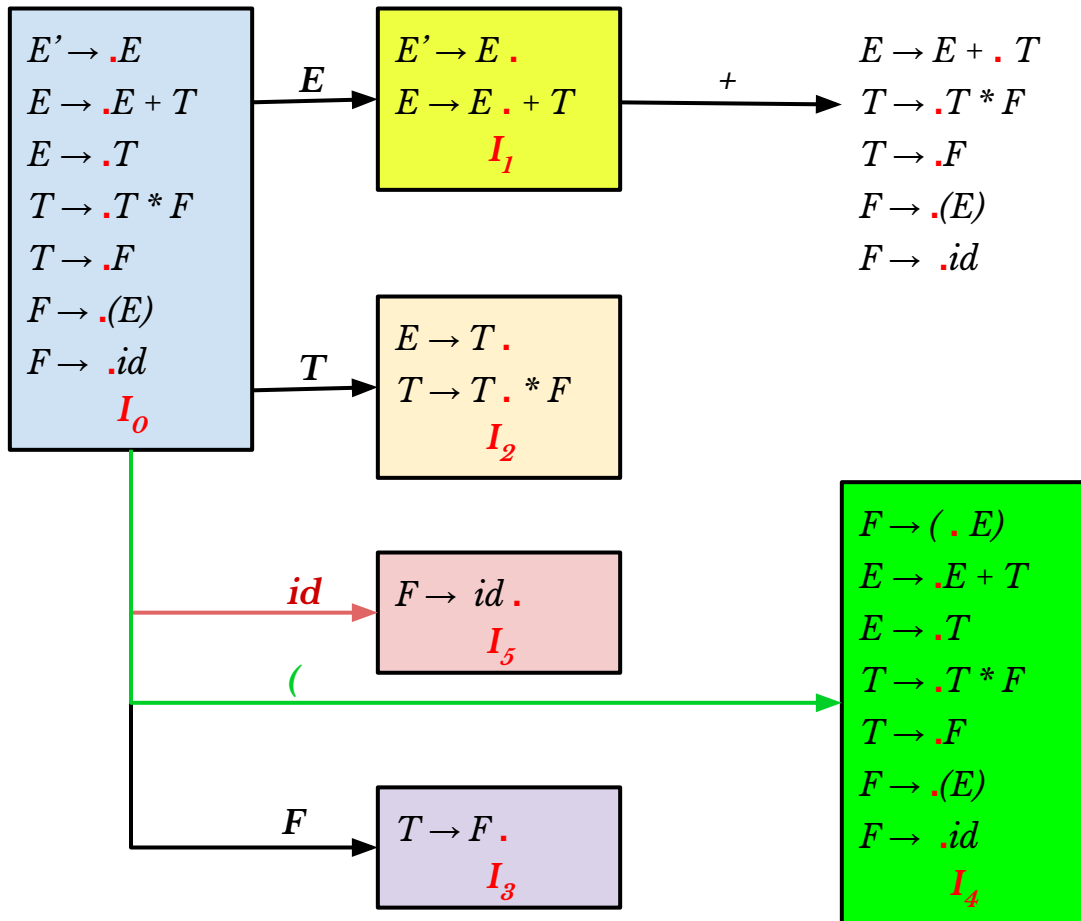
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

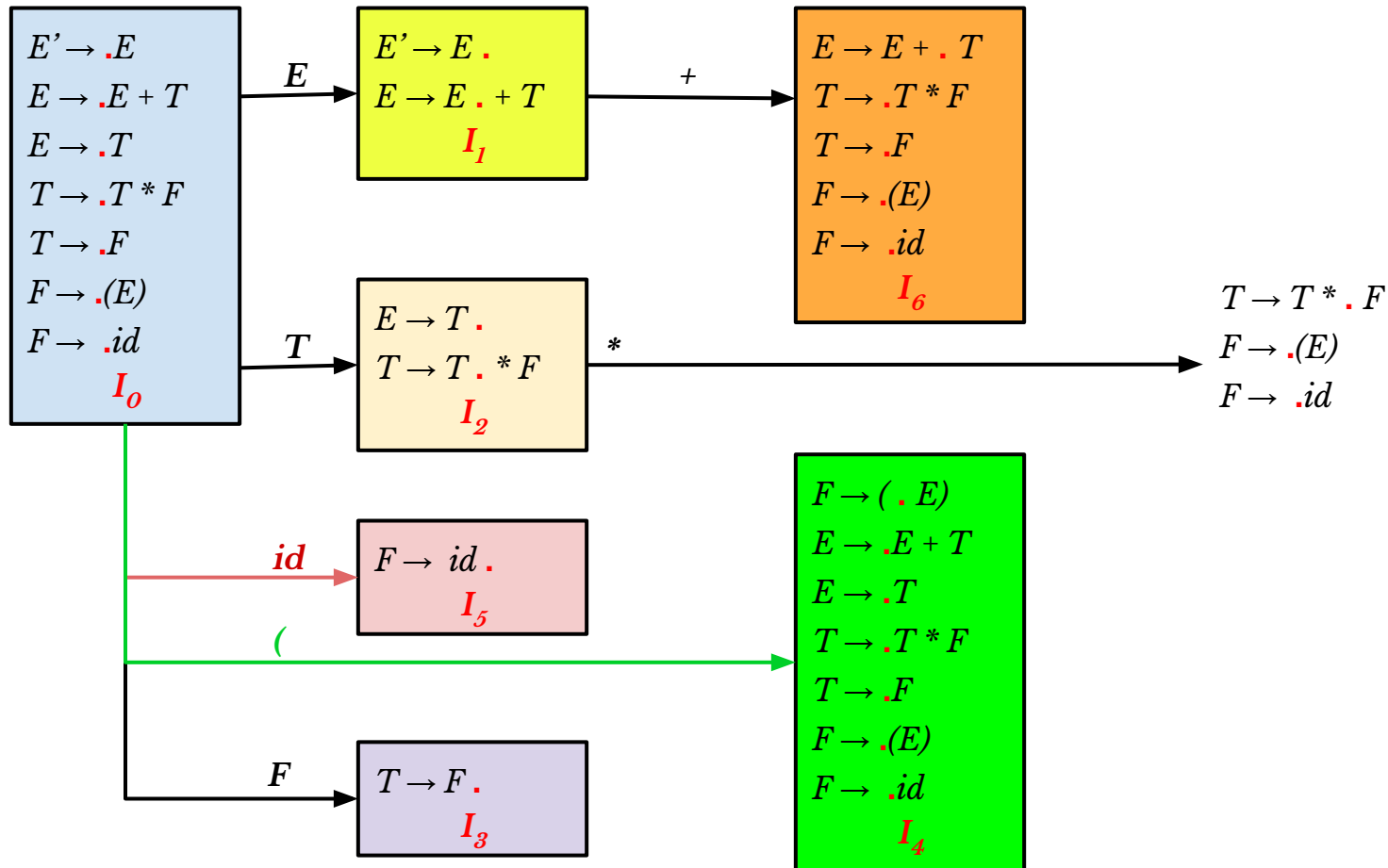
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

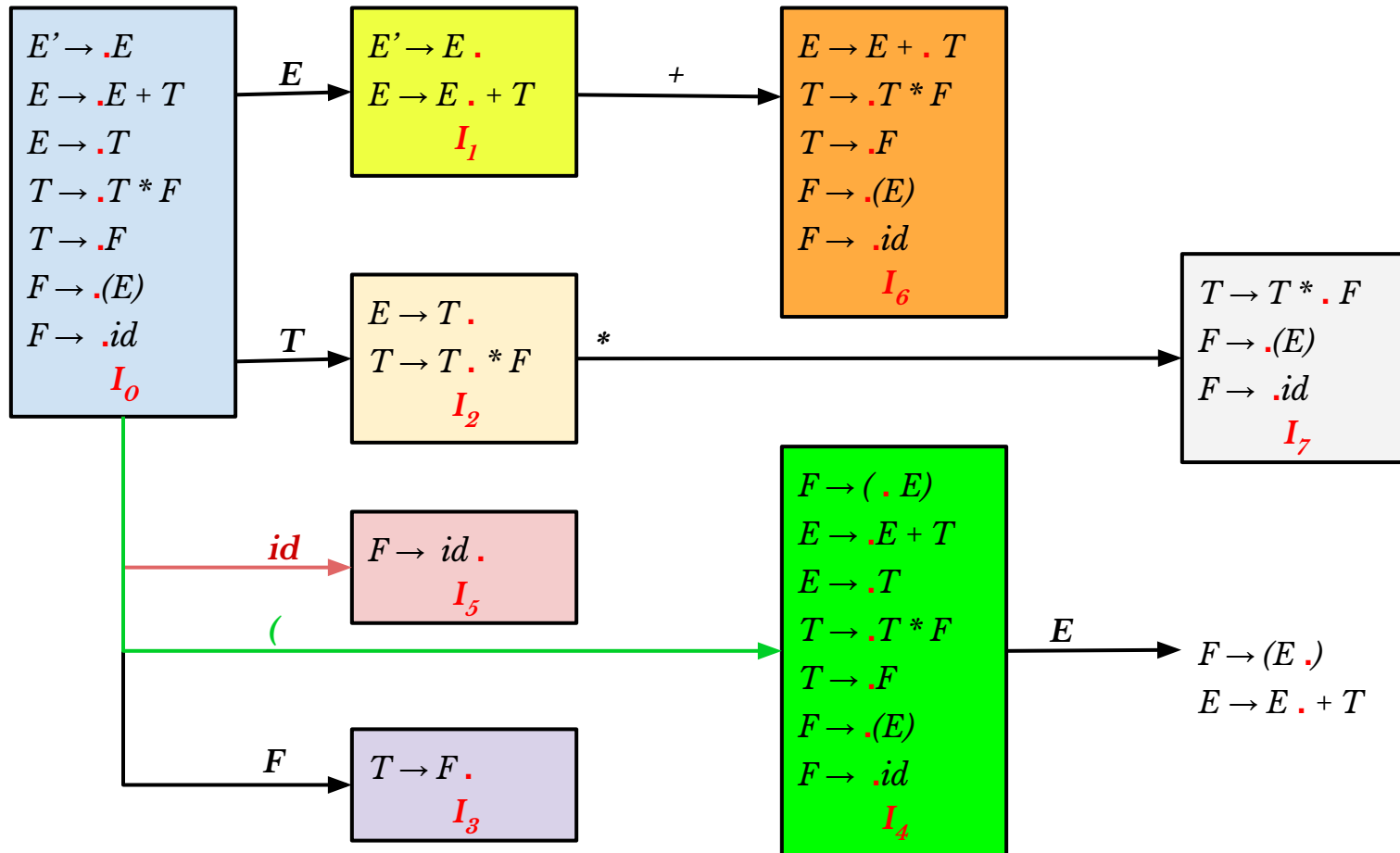
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

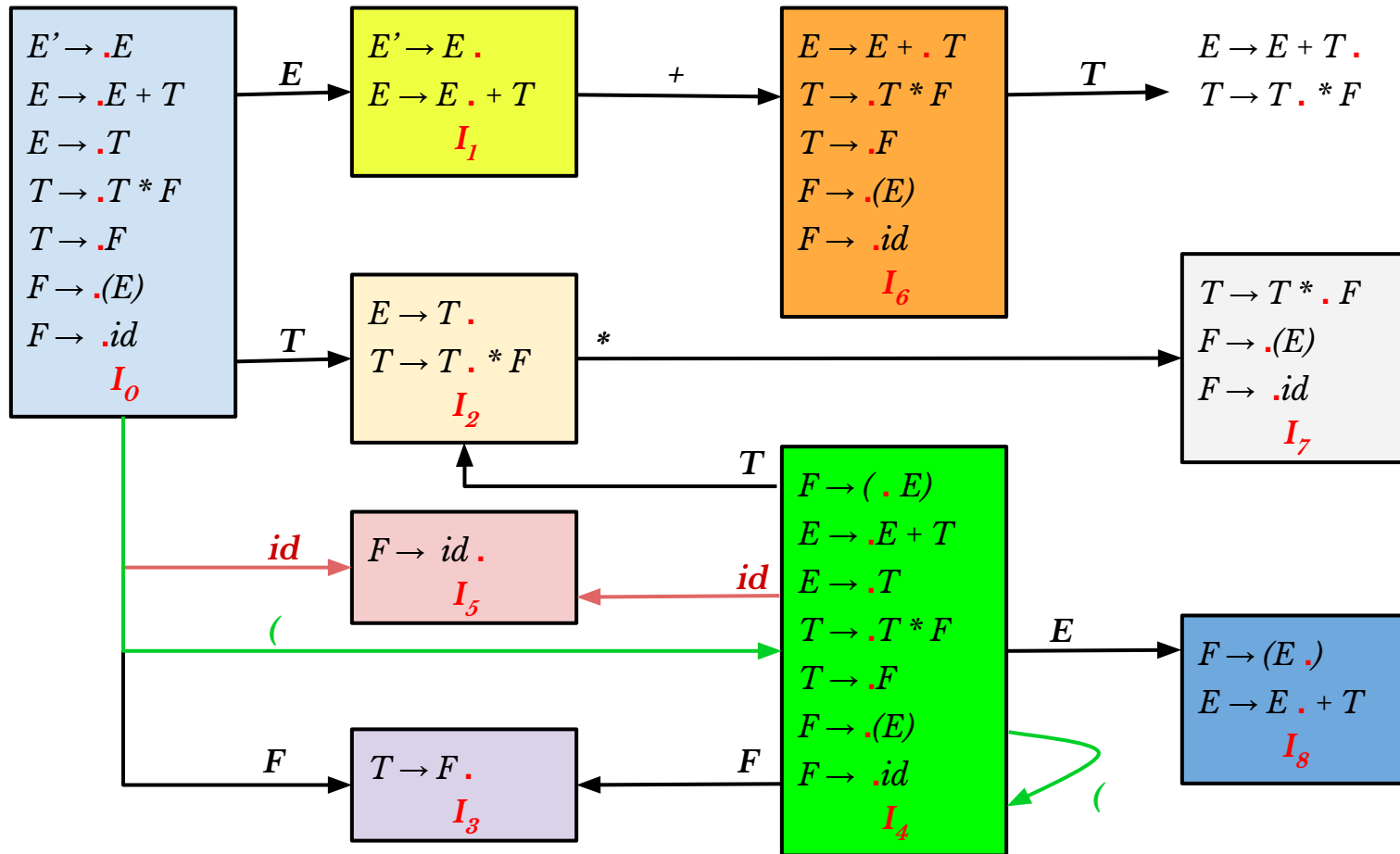
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

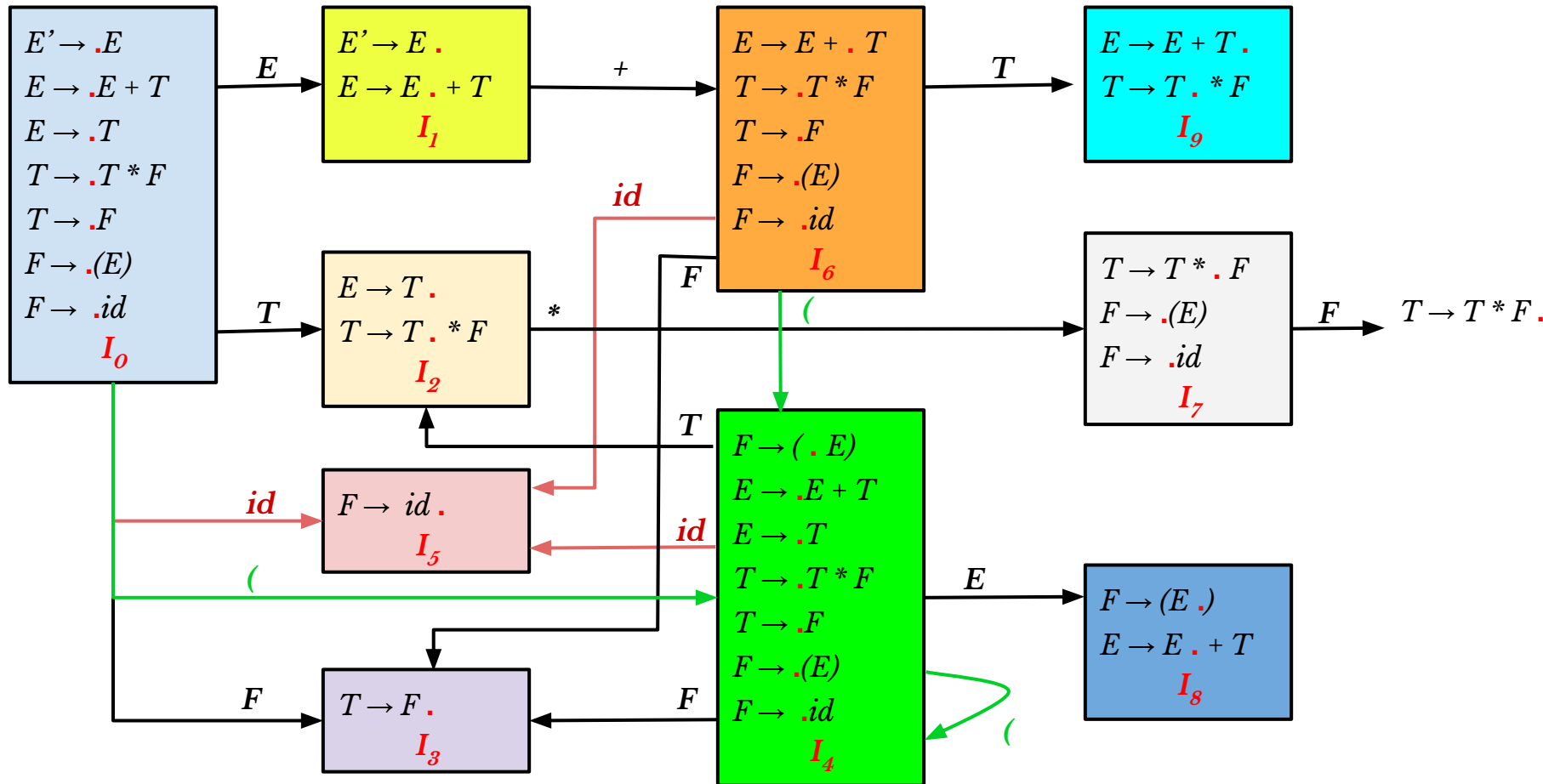
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

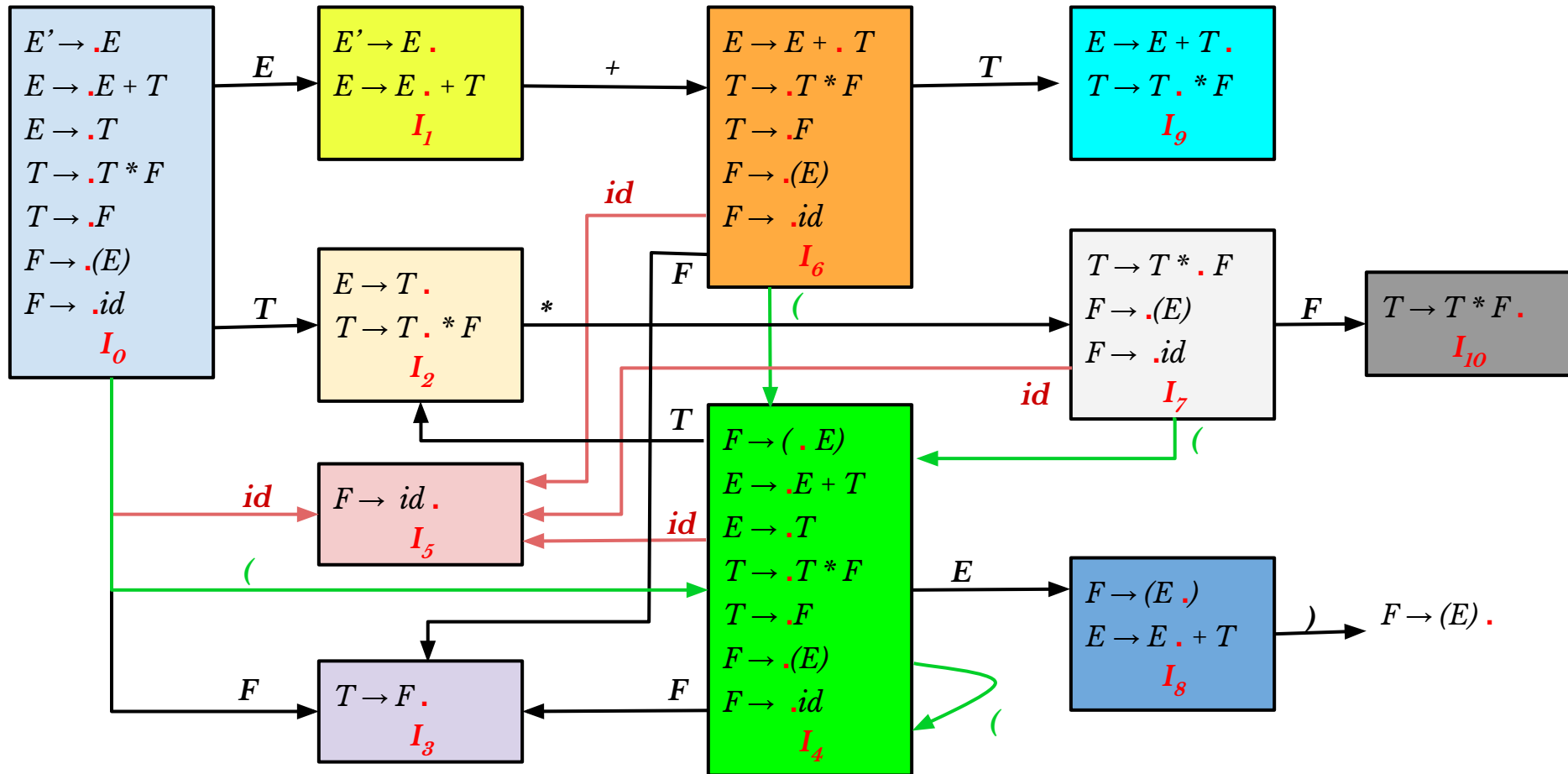
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



LR(0) Automation

0: $E' \rightarrow E$

4: $T \rightarrow F$

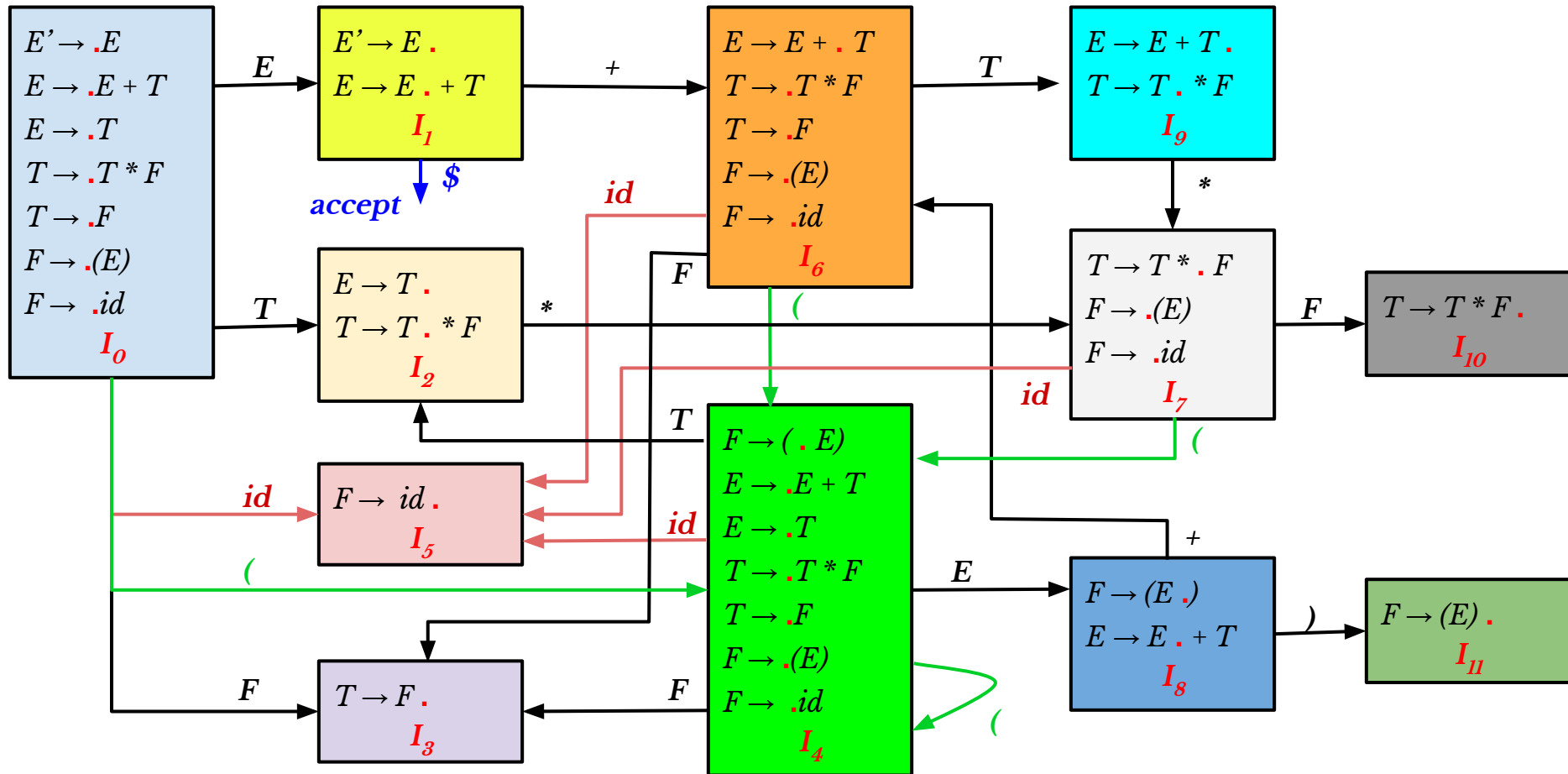
1: $E \rightarrow E + T$

5: $F \rightarrow (E)$

2: $E \rightarrow T$

6: $F \rightarrow id$

3: $T \rightarrow T * F$



Building LR(0) parser: Parsing table

Constructing SLR parsing table

- Remember, LR parsing table has two parts
 - Action: Takes only terminals
 - GOTO: Takes only non-terminals

SLR-Table(G')

1. Construct LR(0) collection for the grammar G'
2. Let I_i represents state S_i , then the parsing action for state i are as follows
 - a. If $[A \rightarrow \alpha \cdot a \beta]$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
 - i. $\text{Action}[i, a] = \text{"shift } j\text{"}$
 - b. If $[A \rightarrow \alpha \cdot]$ is in I_i
 - i. $\text{Action}[i, a] = \text{"reduce } A \rightarrow \alpha\text{"}$ for all $a \in \text{FOLLOW}(A)$
 - c. If $[S' \rightarrow S \cdot]$ is in I_i
 - i. $\text{Action}[i, \$] = \text{"accept"}$
3. For all non-terminals A ,
 - a. if $\text{GOTO}(I_i, A) = I_j$,
 - i. $\text{GOTO}[i, A] = j$

SLR parsing table

GOTO(S, X): Transition from state S to a new state on non-terminal symbol X

For all transitions on non-terminals in state 0

$$\text{GOTO}(0, E) = 1$$

GOTO(0, T) = 2

GOTO(0, F) = 3

State	Action						GOTO		
	id	+	*	()	\$	E	T	F
0							1	2	3
1									
2									
3									
4									
5									
6									
7									
8									
9									
10									
11									

[illegible]

SLR parsing table

GOTO(S, X): Transition from state S to a new state on non-terminal symbol X

For all transitions on non-terminals in state 4

GOTO(4, E) = 8

GOTO(4, T) = 2

GOTO(4, F) = 3

For all transitions on non-terminals in state 6

GOTO(6, T) = 9

GOTO(6, F) = 3

For all transitions on non-terminals in state 7

GOTO(7, F) = 10

[illegible]

SLR parsing table

If $\llbracket A \rightarrow \alpha.a\beta \rrbracket$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
then

$$\text{Action}[i, a] = \text{"shift } j\text{"}$$

Note: a is terminal

For all transitions on terminals in state 0

Action[0, id] = “shift 5” or “s5”

Action[0, () = "s4"

[illegible]

SLR parsing table

If $\llbracket A \rightarrow \alpha.a\beta \rrbracket$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
then

$$\text{Action}[i, a] = \text{"shift } j\text{"}$$

Note: a is terminal

For all transitions on terminals in state 1

Action[1, +] = "s6"

[illegible]

SLR parsing table

If $[A \rightarrow \alpha.a\beta]$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
then

$$\text{Action}[i, a] = \text{"shift } j\text{"}$$

Note: a is terminal

For all transitions on terminals in state 2

```
Action[2, *] = "s7"
```

[illegible]

SLR parsing table

If $\llbracket A \rightarrow \alpha.a\beta \rrbracket$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
then

$$\text{Action}[i, a] = \text{"shift } j\text{"}$$

Note: a is terminal

For all transitions on terminals in state 4

```
Action[4, id] = "s5"
```

```
Action[4, ( ] = "s4"
```

[illegible]

SLR parsing table

If $\llbracket A \rightarrow \alpha.a\beta \rrbracket$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
then

$$\text{Action}[i, a] = \text{"shift } j\text{"}$$

Note: a is terminal

For all transitions on terminals in state 6

```
Action[6, id] = "s5"
```

Action[6, (] = "s4"

[illegible]

SLR parsing table

If $\llbracket A \rightarrow \alpha.a\beta \rrbracket$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
then

$$\text{Action}[i, a] = \text{"shift } j\text{"}$$

Note: a is terminal

For all transitions on terminals in state 7

```
Action[7, id] = "s5"
```

```
Action[7, ( ] = "s4"
```

[illegible]

SLR parsing table

If $[A \rightarrow \alpha_{\bullet}]$ is in I_i .

Action[i, a] = “reduce $\mathcal{A} \rightarrow \alpha$ ” for all $a \in \text{FOLLOW}(\mathcal{A})$

For all terminals in Follow($T \rightarrow F$) in state 3

```
Action[3, +] = "r4"
```

```
Action[3, *] = "r4"
```

```
Action[3, )] = "r4"
```

```
Action[3, $] = "r4"
```

[illegible]

SLR parsing table

If $[A \rightarrow \alpha_\bullet]$ is in I_i

Action[i, a] = “reduce $A \rightarrow \alpha$ ” for all $a \in \text{FOLLOW}(A)$

For all terminals in Follow($F \rightarrow id$) in state 5

```
Action[5, +] = "r6"
```

```
Action[5, *] = "r6"
```

```
Action[5, )] = "r6"
```

```
Action[5, $] = "r6"
```

[illegible]

SLR parsing table

If $[A \rightarrow \alpha.]$ is in I_i .

Action $[i, a]$ = “reduce $A \rightarrow \alpha$ ” for all $a \in \text{FOLLOW}(A)$

For all terminals in $\text{Follow}(T \rightarrow T^*F)$ in state 10

Action[10, +] = "r3"

```
Action[10, *] = "r3"
```

```
Action[10, )] = "r3"
```

```
Action[10, $] = "r3"
```

[illegible]

SLR parsing table

If $[A \rightarrow \alpha \cdot]$ is in I_i

Action $[i, a]$ = “reduce $A \rightarrow \alpha$ ” for all $a \in \text{FOLLOW}(A)$

For all terminals in Follow($F \rightarrow (E)$) in state 11

Action $[11, +]$ = “r5”

Action $[11, *]$ = “r5”

Action $[11,)]$ = “r5”

Action $[11, \$]$ = “r5”

State	Action						GOTO		
	id	+	*	()	\$	E	T	F
0	s5			s4			1	2	3
1		s6							
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			

SLR parsing table

If $[S' \rightarrow S_i]$ is in I_i
 Action $[i, \$]$ = “accept”

Action $[1, \$]$ = “accept”

All empty entries are “error” case.

State	Action						GOTO		
	id	+	*	()	\$	E	T	F
0	s5			s4			1	2	3
1		s6				acc			
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			

LR(0) parser: Parsing an input string

SLR parsing algorithm

1. Let a be the first symbol in $w\$$
2. Repeat
 - a. Let s be the state on top of the stack
 - b. If $\text{Action}[s, a] == s\#t$
 - i. Push t on to the stack
 - ii. Let a be the next symbol
 - c. Else if $\text{Action}[s, a] == \text{reduce } A \rightarrow B$
 - i. Pop $|B|$ symbols off the stack
 - ii. Push $\text{GOTO}[t, A]$ on to the stack
 - iii. Output production $A \rightarrow B$
 - d. Else if $\text{Action}[s, a] == \text{"accept"}$
 - i. Halt
 - e. Else
 - i. Error : Call error handler

SLR parsing

Input: **id * id + id**

Stack

0
0 5
0 3
0 2
0 2 7
0 2 7 5
0 2 7 10
0 2
0 1
0 1 6
0 1 6 5
0 1 6 3
0 1 6 9
0 1

Symbol

id
F
T
T *
T * id
T * F
T
E
E +
E + id
E + F
E + T
E

Input

id * id + id \$
* id + id \$
* id + id \$
* id + id \$
id + id \$
+ id \$
+ id \$
+ id \$
+ id \$
id \$
\$
\$
\$
\$

Action

Shift
Reduce $F \rightarrow id$
Reduce $T \rightarrow F$
Shift
Shift
Reduce $F \rightarrow id$
Reduce $T \rightarrow T * F$
Reduce $E \rightarrow T$
Shift
Shift
Reduce $F \rightarrow id$
Reduce $T \rightarrow F$
Reduce $E \rightarrow E + T$
Accept

LR(0) Automation: Another example

Grammar G: $S \rightarrow L = R \mid R$
 $L \rightarrow *R \mid \text{id}$
 $R \rightarrow L$

Construct SLR parser for the above grammar

$I_0:$	$S' \rightarrow \cdot S$ $S \rightarrow \cdot L = R$ $S \rightarrow \cdot R$ $L \rightarrow \cdot * R$ $L \rightarrow \cdot \text{id}$ $R \rightarrow \cdot L$
$I_1:$	$S' \rightarrow S \cdot$

$I_2:$	$S \rightarrow L \cdot = R$ $R \rightarrow L \cdot$
$I_3:$	$S \rightarrow R \cdot$
$I_4:$	$L \rightarrow * \cdot R$ $R \rightarrow \cdot L$ $L \rightarrow \cdot * R$ $L \rightarrow \cdot \text{id}$

$I_5:$	$L \rightarrow \text{id} \cdot$
$I_6:$	$S \rightarrow L = \cdot R$ $R \rightarrow \cdot L$ $L \rightarrow \cdot * R$ $L \rightarrow \cdot \text{id}$

$I_7:$	$L \rightarrow * R \cdot$
$I_8:$	$R \rightarrow L \cdot$
$I_9:$	$S \rightarrow L = R \cdot$

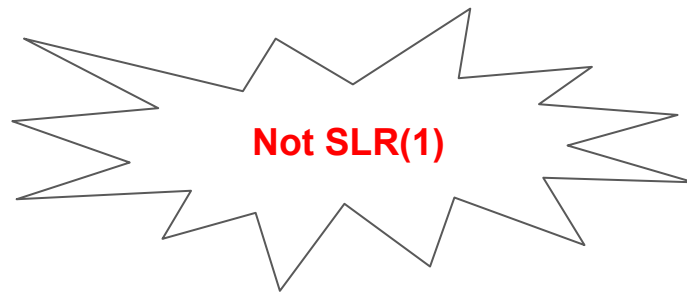
Parsing table entry for state 2

Action[2, =] = “shift 6” or “reduce $R \rightarrow L$ ” ?

SLR(1) grammar

- If any cell in the parsing table has multiple entries, then
 - Grammar is not SLR(1) or LR(0)

Grammar G:

$$S \rightarrow L = R \mid R$$
$$L \rightarrow *R \mid \text{id}$$
$$R \rightarrow L$$


- Every SLR(1) grammar is unambiguous.
- But, there are many unambiguous grammar that are not SLR(1)

LR(0) Automation: Another example - 2

Grammar G: $S \rightarrow AaAb \mid BbBa$

$A \rightarrow \varepsilon$

$B \rightarrow \varepsilon$

Construct SLR parser for the above grammar

$I_0:$	$S' \rightarrow \cdot S$
	$S \rightarrow \cdot AaAb$
	$S \rightarrow \cdot BbBa$
	$A \rightarrow \cdot$
	$B \rightarrow \cdot$

Parsing table entry for state I_0

Action[0, a] = “reduce $A \rightarrow \varepsilon$ ” or “reduce $B \rightarrow \varepsilon$ ” ?

Follow(A) = {a, b}

Follow(B) = {a, b}

More Powerful LR Parsers

LR(1) and LALR parsers

- Canonical LR or LR(1) parser
 - Makes full use of lookahead symbols, i.e., both First() and Follow() symbols
 - Recall, LR(0) uses only Follow() symbols
- Lookahead LR or LALR parser
 - Inclusion of lookahead symbols in LR(0) sets of items
 - Can handle more grammars than SLR method
 - LALR tables are no bigger than the SLR tables
 - Fewer states than typical LR(1) parser

SLR vs LR(1)

- Reduction during parsing in SLR
 - $\text{Stack}(\beta\alpha) \Rightarrow \text{Stack}(\beta A)$ if we had $[A \rightarrow \alpha.]$
- What if, βA is not followed by a in right-sentential form
 - Reduction $A \rightarrow \alpha$ is invalid
- E.g.:
 - Grammar G:
$$\begin{aligned} S &\rightarrow L = R \mid R \\ L &\rightarrow *R \mid \text{id} \\ R &\rightarrow L \end{aligned}$$
 - If we apply “reduce $R \rightarrow L$ ” on Action[2, =]
“ $L = \dots$ ” \Rightarrow “ $R = \dots$ ”
 - However, there is no right-sentential form of the grammar the begins with “ $R = \dots$ ”

Building LR(1) parser: Canonical set of items

Canonical LR(1) Items

- General form of an LR(1) item is $[A \rightarrow \alpha.\beta, a]$
 - $A \rightarrow \alpha\beta$ is a production, and
 - a is a lookahead terminal symbol or endmarker $\$$
- Lookahead symbol has no effect in an item of the form $[A \rightarrow \alpha.\beta, a]$, where β is not ϵ
- Lookahead symbol is required during the reduction only
 - We reduce $[A \rightarrow \alpha., a]$, only if the next input symbol is a

Computation of the canonical LR(1) collection

Items(G')

1. $C = \text{Closure}(\{[S' \rightarrow \cdot S, \$]\})$
2. Repeat
 - a. For each set of items I in C
 - i. For each grammar symbol X
 1. If $\text{GOTO}(I, X)$ is not empty and not in C
 - i. Add $\text{GOTO}(I, X)$ to C
3. Until no new sets of items are added to C

Closure(I)

1. Repeat
 - a. For each items $[A \rightarrow \alpha \cdot B\beta, a]$ in I
 - i. For each production $B \rightarrow \gamma$ in G
 1. For each terminal b in $\text{FIRST}(\beta a)$
 - i. Add $[B \rightarrow \cdot \gamma, b]$
2. Until no new items are added to I

GOTO(I, X)

1. For each items $[A \rightarrow \alpha \cdot X\beta, a]$ in I
 - a. Add item $[A \rightarrow \alpha X \cdot \beta, a]$ to J
2. return Closure(J)

LR(1) Automation

0: $S' \rightarrow S$

1: $S \rightarrow CC$

2: $C \rightarrow c C$

3: $C \rightarrow d$

$S' \rightarrow \cdot S, \$$

$S \rightarrow \cdot CC, \$$

$C \rightarrow \cdot cC, c/d$

$C \rightarrow \cdot d, c/d$

LR(1) Automation

0: $S' \rightarrow S$

1: $S \rightarrow CC$

2: $C \rightarrow c C$

3: $C \rightarrow d$

$S' \rightarrow \cdot S, \$$
 $S \rightarrow \cdot CC, \$$
 $C \rightarrow \cdot cC, c/d$
 $C \rightarrow \cdot d, c/d$
 I_0

$\xrightarrow{S} S' \rightarrow S \cdot, \$$

LR(1) Automation

0: $S' \rightarrow S$

1: $S \rightarrow CC$

2: $C \rightarrow c C$

3: $C \rightarrow d$

$S' \rightarrow \cdot S, \$$
 $S \rightarrow \cdot CC, \$$
 $C \rightarrow \cdot cC, c/d$
 $C \rightarrow \cdot d, c/d$

I_0

S

$S' \rightarrow S \cdot, \$$
 I_1

C

$S \rightarrow C \cdot C, \$$
 $C \rightarrow \cdot cC, \$$
 $C \rightarrow \cdot d, \$$

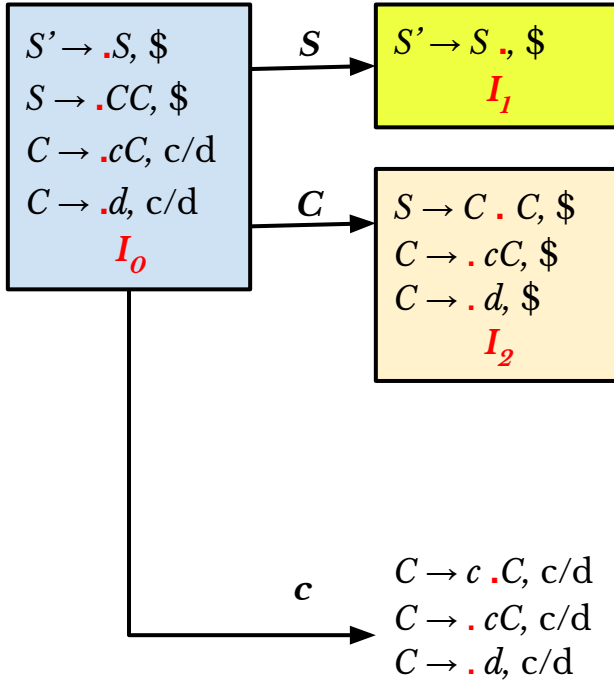
LR(1) Automation

0: $S' \rightarrow S$

1: $S \rightarrow CC$

2: $C \rightarrow c C$

3: $C \rightarrow d$



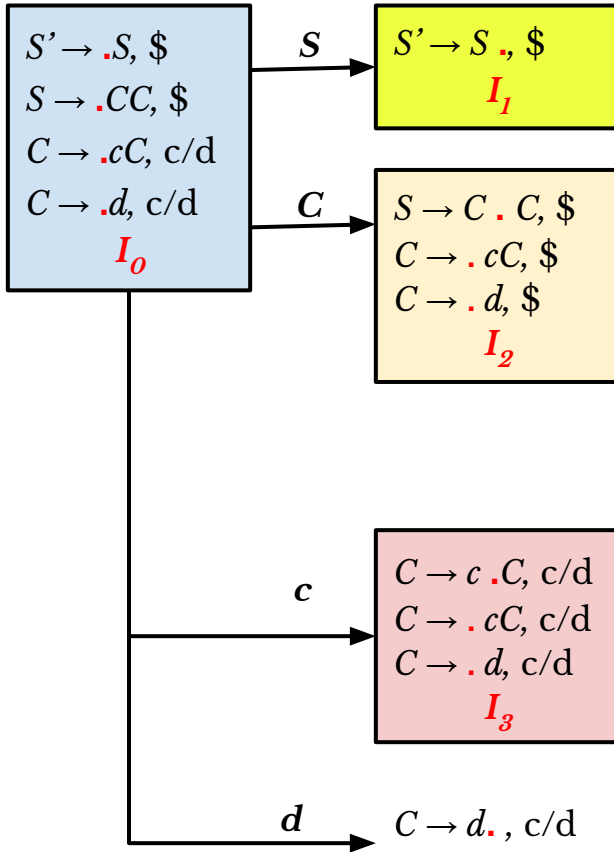
LR(1) Automation

0: $S' \rightarrow S$

1: $S \rightarrow CC$

2: $C \rightarrow c C$

3: $C \rightarrow d$



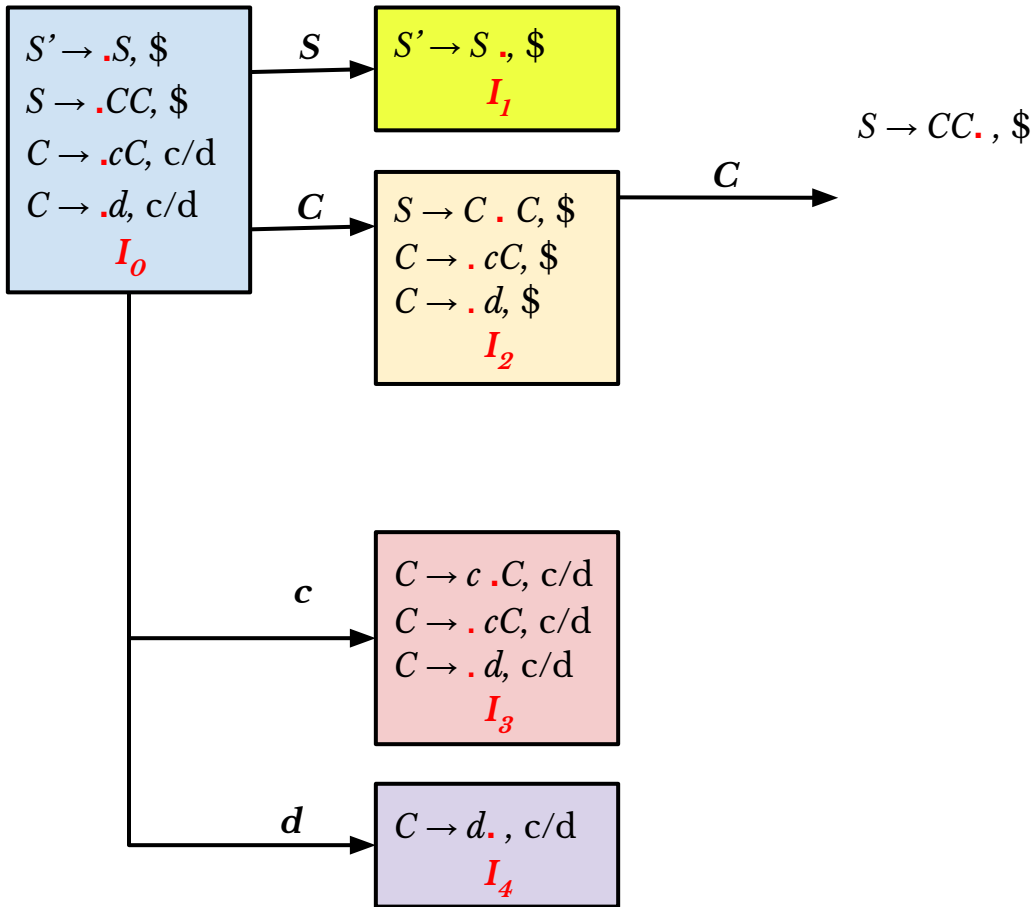
LR(1) Automation

0: $S' \rightarrow S$

1: $S \rightarrow CC$

2: $C \rightarrow c C$

3: $C \rightarrow d$



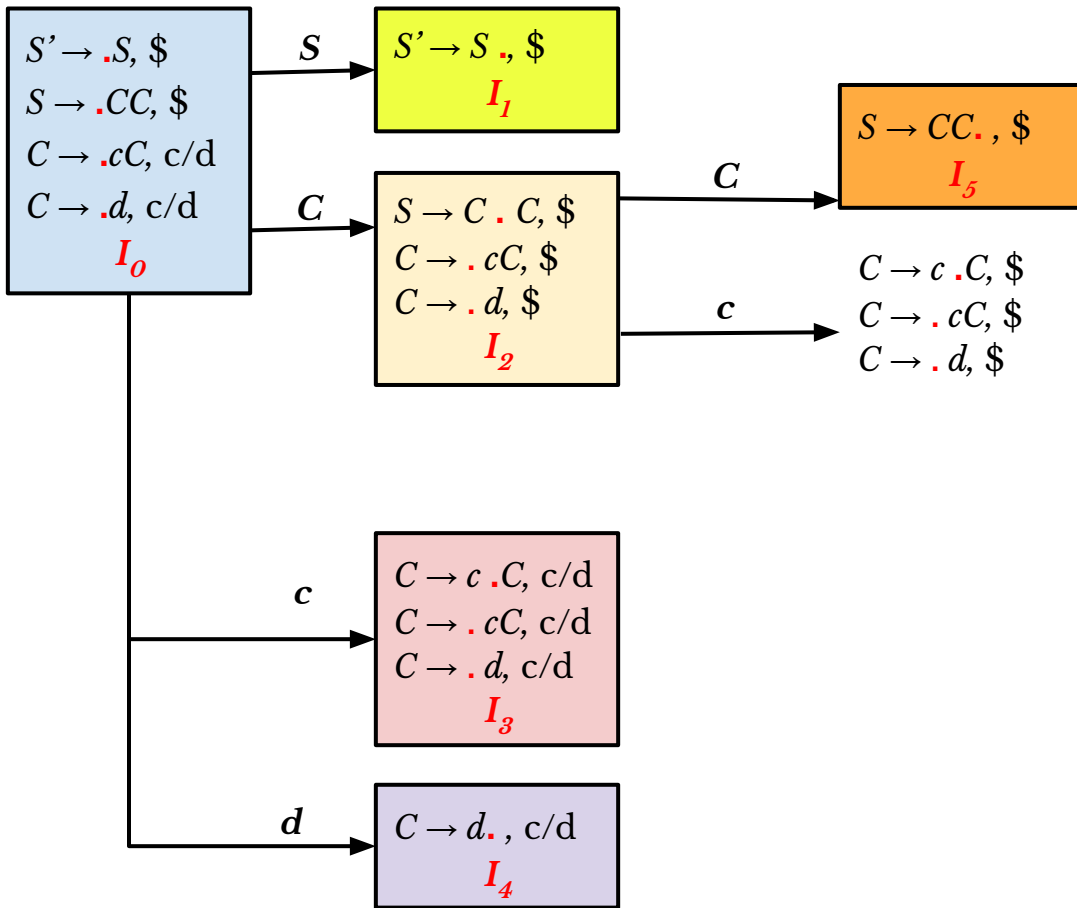
LR(1) Automation

0: $S' \rightarrow S$

1: $S \rightarrow CC$

2: $C \rightarrow cC$

3: $C \rightarrow d$



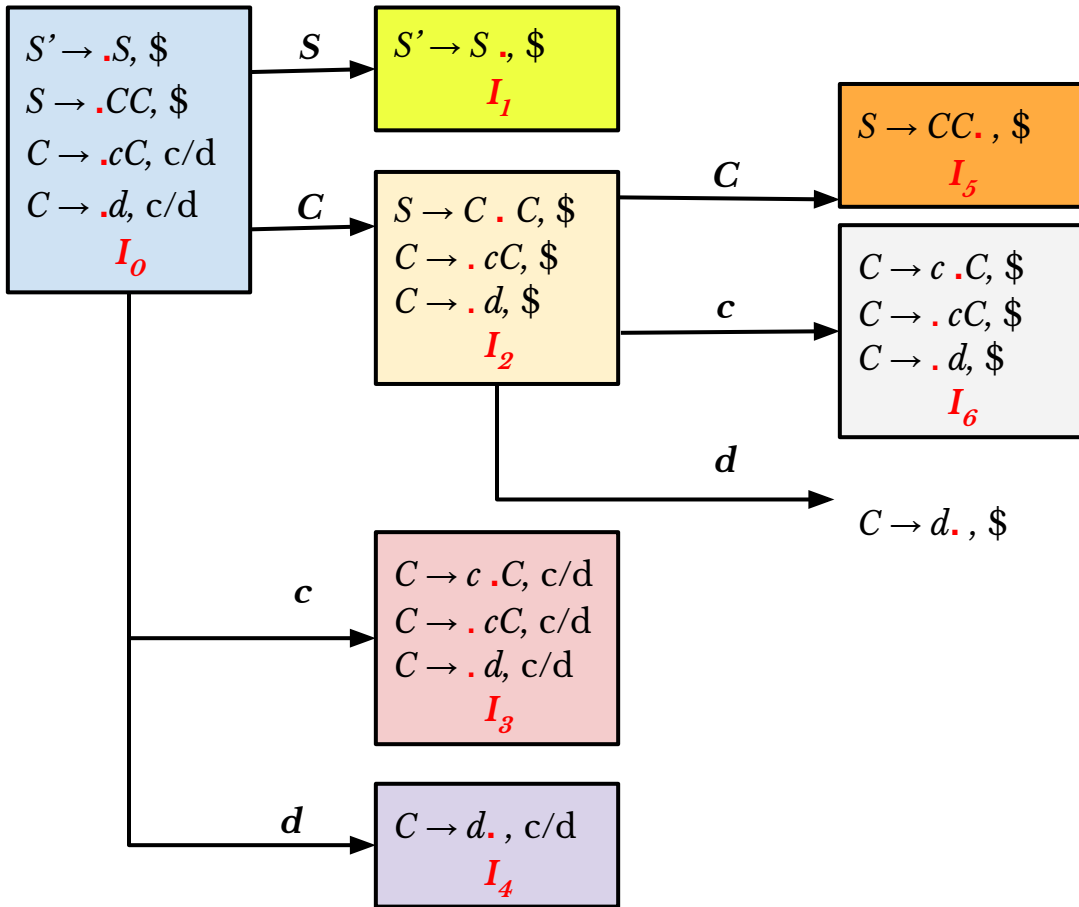
LR(1) Automation

0: $S' \rightarrow S$

1: $S \rightarrow CC$

2: $C \rightarrow cC$

3: $C \rightarrow d$



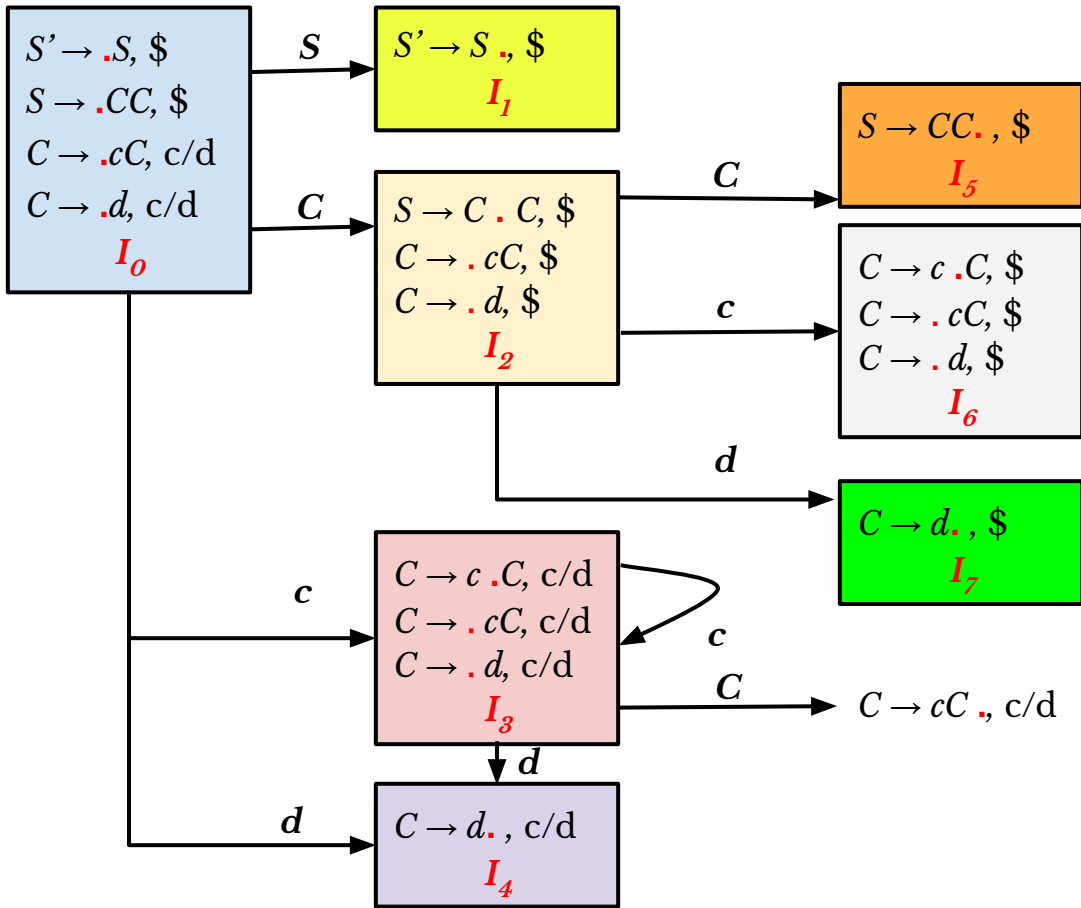
LR(1) Automation

0: $S' \rightarrow S$

1: $S \rightarrow CC$

2: $C \rightarrow cC$

3: $C \rightarrow d$



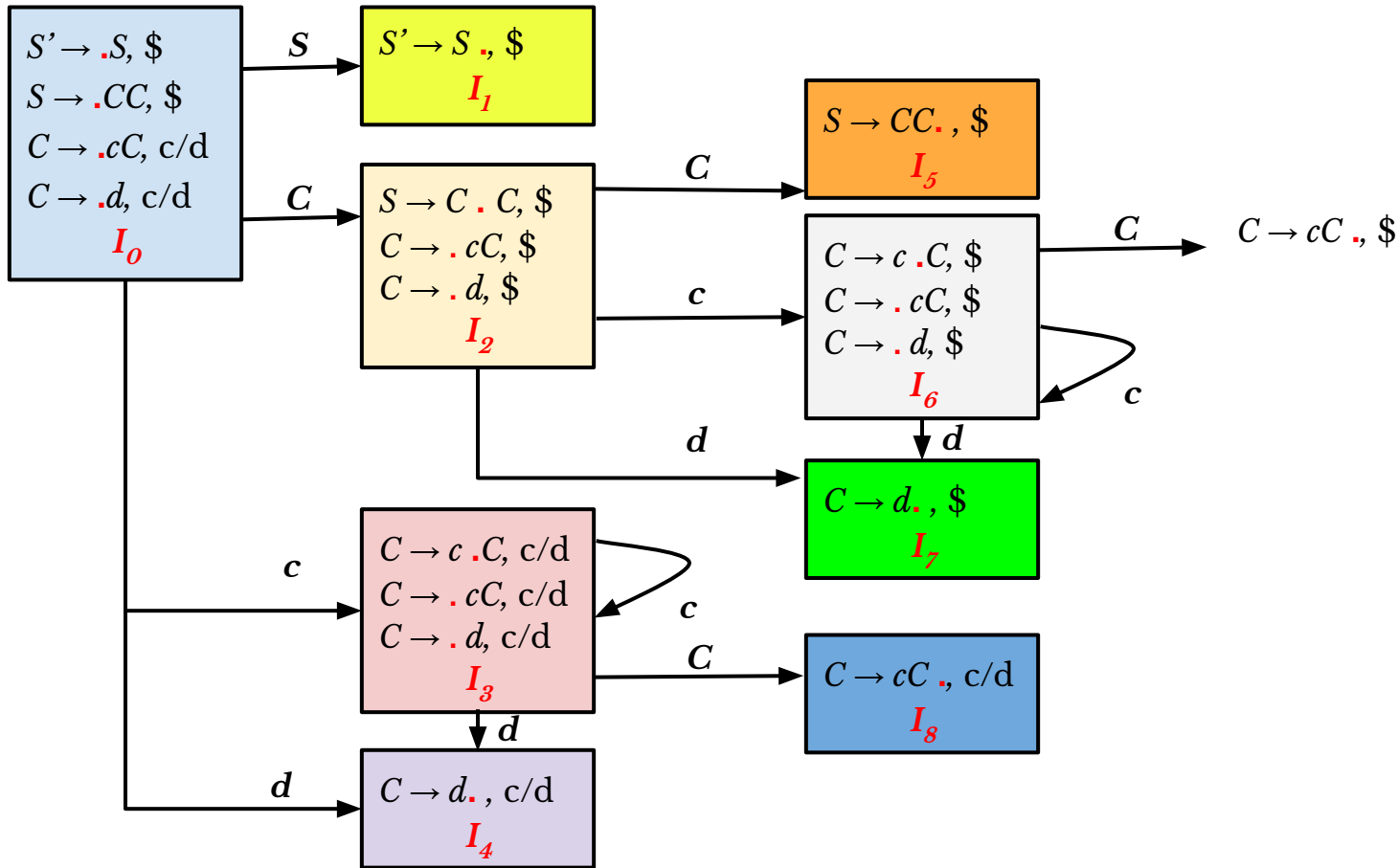
LR(1) Automation

0: $S' \rightarrow S$

1: $S \rightarrow CC$

2: $C \rightarrow cC$

3: $C \rightarrow d$



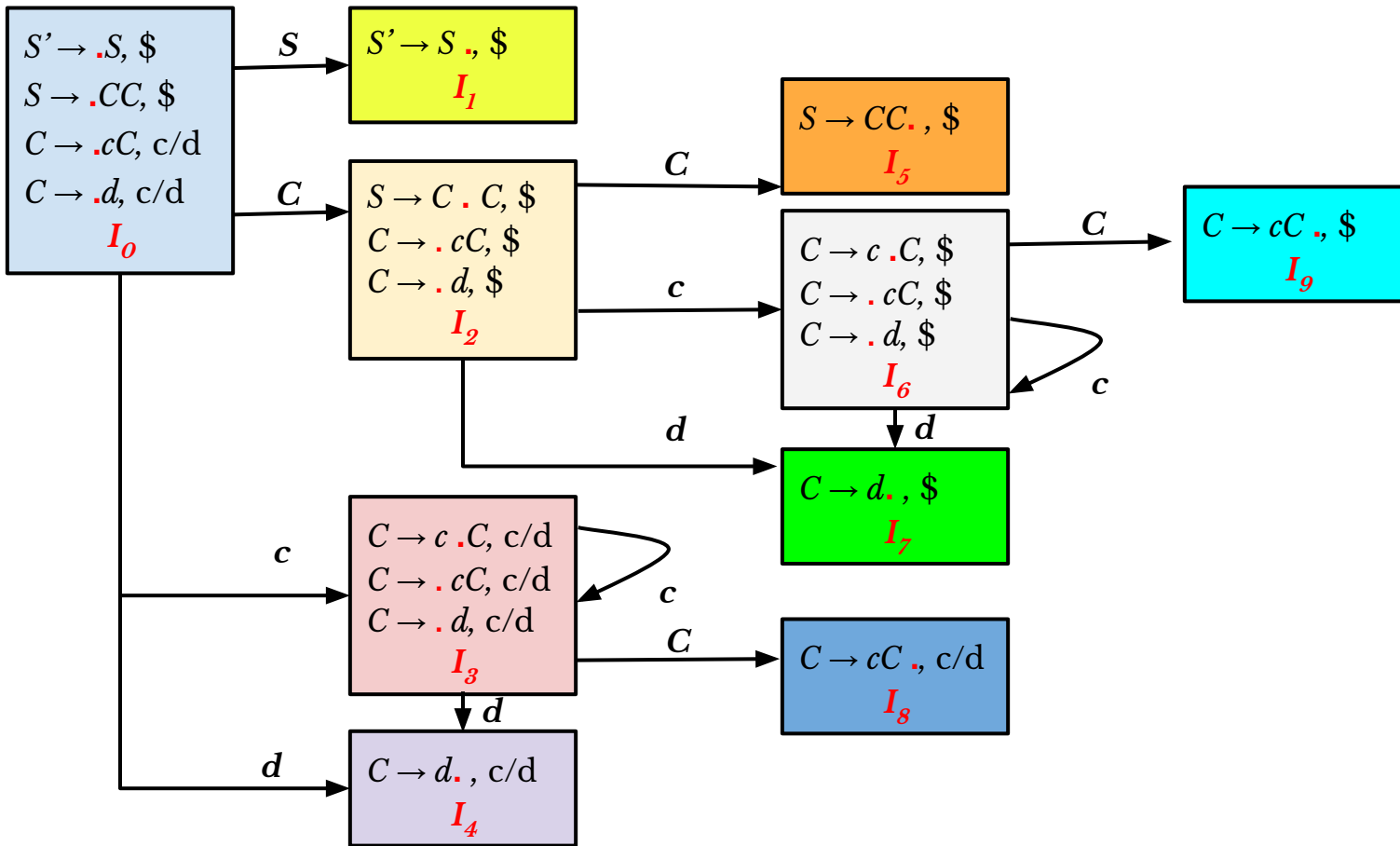
LR(1) Automation

0: $S' \rightarrow S$

1: $S \rightarrow CC$

2: $C \rightarrow cC$

3: $C \rightarrow d$



Building LR(1) parser: Parsing table

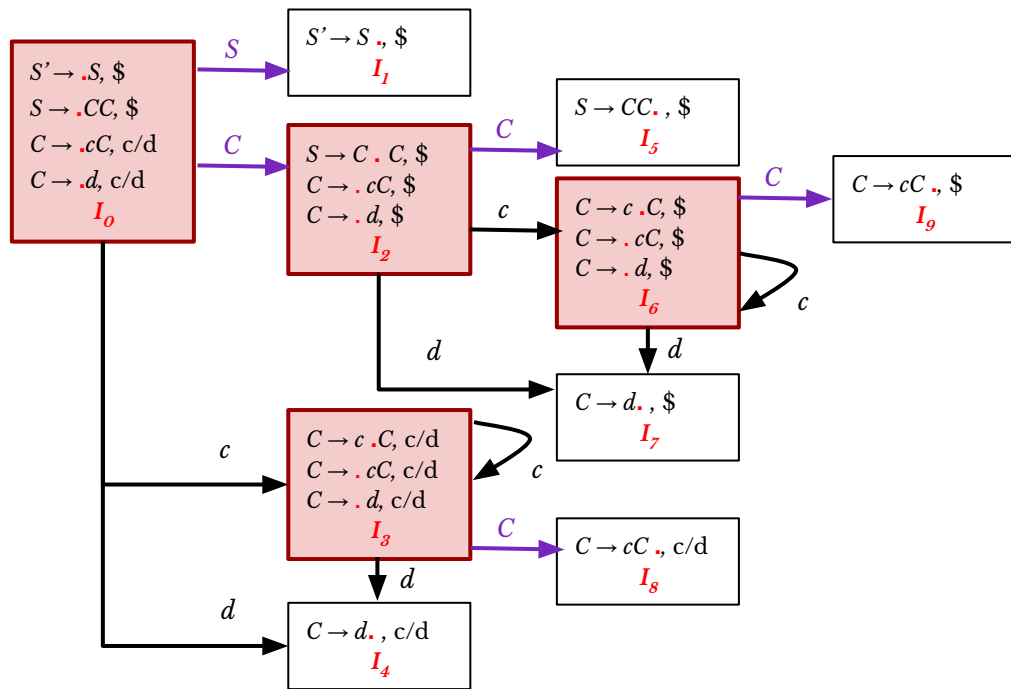
Constructing LR(1) parsing table

LR(1)-Table(G')

1. Construct LR(1) collection for the grammar G'
2. Let I_i represents state S_i , then the parsing action for state i are as follows
 - a. If $[A \rightarrow \alpha \cdot a \beta, b]$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
 - i. $\text{Action}[i, a] = \text{"shift } j\text{"}$
 - b. If $[A \rightarrow \alpha \cdot, b]$ is in I_i
 - i. $\text{Action}[i, b] = \text{"reduce } A \rightarrow \alpha\text{"}$
 - c. If $[S' \rightarrow S \cdot, \$]$ is in I_i
 - i. $\text{Action}[i, \$] = \text{"accept"}$
3. For all non-terminals A ,
 - a. if $\text{GOTO}(I_i, A) = I_j$,
 - i. $\text{GOTO}[i, A] = j$

LR(1) parsing table

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



State	Action			GOTO	
	c	d	\$	S	C
0				1	2
1					
2					5
3					8
4					
5					
6					9
7					
8					
9					

Rule 3: For all non-terminals A ,
if $\text{GOTO}(I_i, A) = I_j$, then $\text{GOTO}[i, A] = j$

$\text{GOTO}[0, S] = 1$

$\text{GOTO}[0, C] = 2$

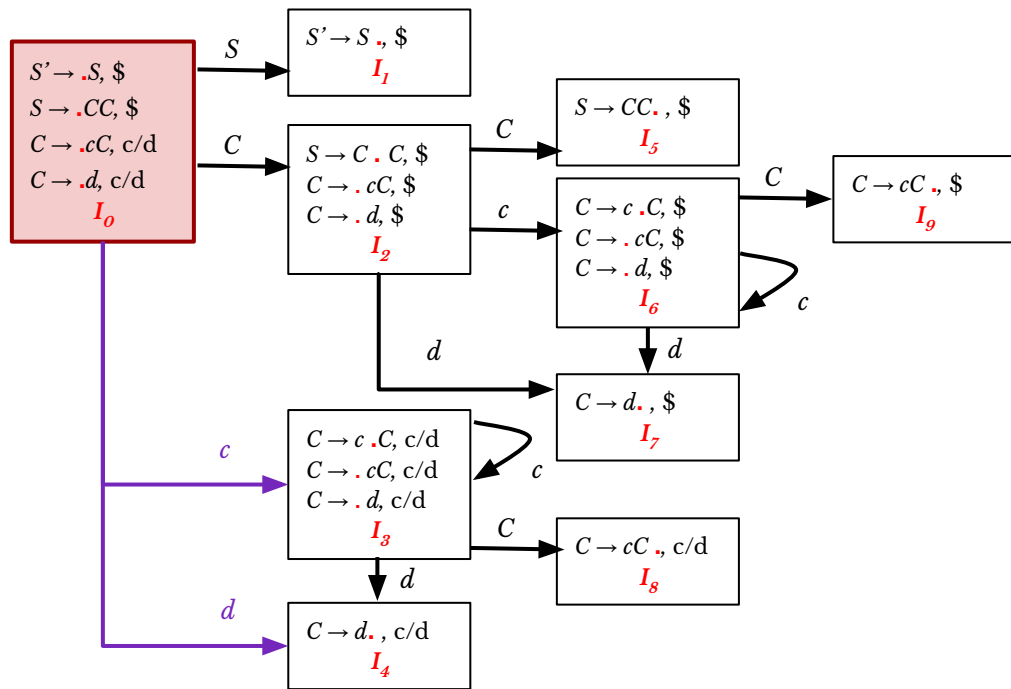
$\text{GOTO}[2, C] = 5$

$\text{GOTO}[3, C] = 8$

$\text{GOTO}[6, C] = 9$

LR(1) parsing table

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



State	Action			GOTO	
	c	d	\$	S	C
0	s3	s4		1	2
1					
2					5
3					8
4					
5					
6					9
7					
8					
9					

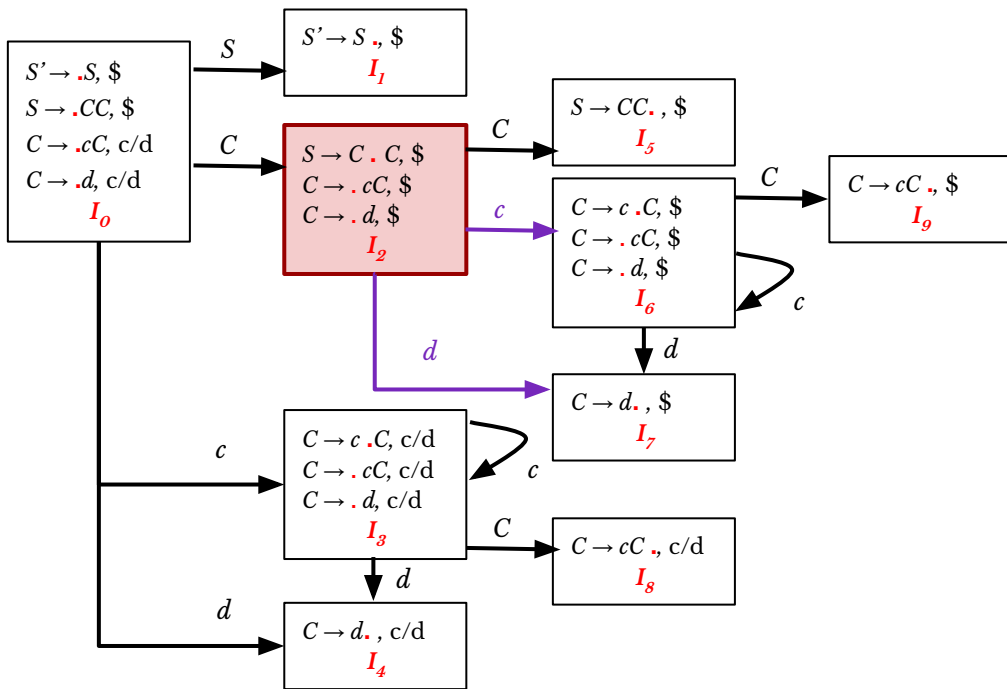
Rule 2a: If $[A \rightarrow \alpha \cdot a \beta, b]$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
 Action $[i, a]$ = "shift j " or "s j "

Action[0, c] = s3

Action[0, d] = s4

LR(1) parsing table

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



State	Action			GOTO	
	c	d	\$	S	C
0	s3	s4		1	2
1					
2	s6	s7			5
3					8
4					
5					
6					9
7					
8					
9					

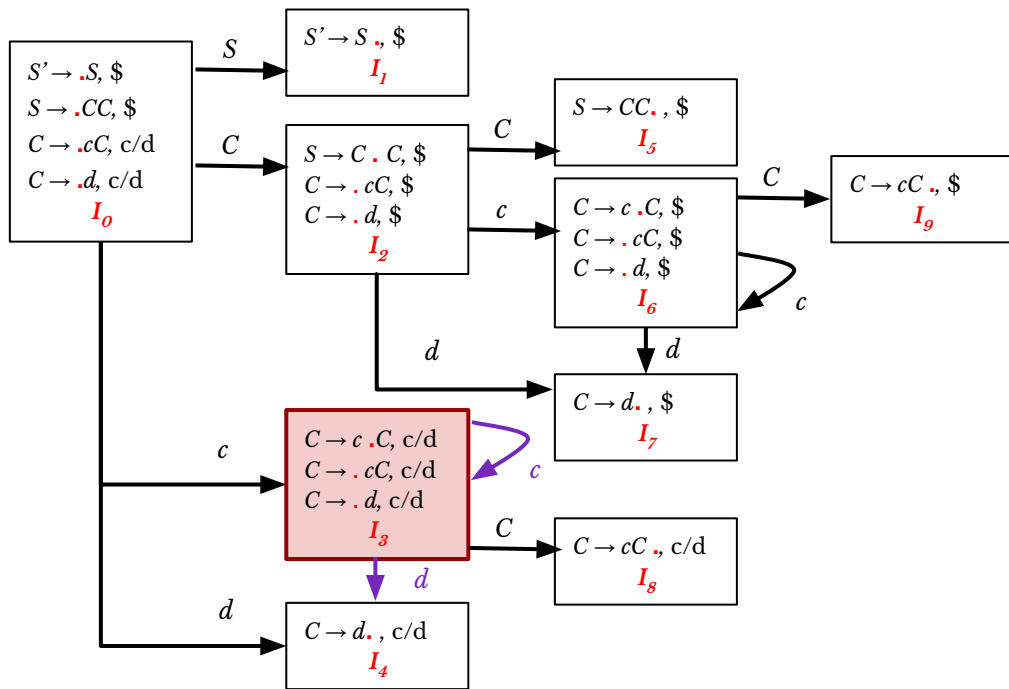
Rule 2a: If $[A \rightarrow \alpha \cdot a \beta, b]$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
 Action $[i, a]$ = "shift j " or "s j "

Action[2, c] = s6

Action[2, d] = s7

LR(1) parsing table

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



State	Action			GOTO	
	c	d	\$	S	C
0	s3	s4		1	2
1					
2	s6	s7			5
3	s3	s4			8
4					
5					
6					9
7					
8					
9					

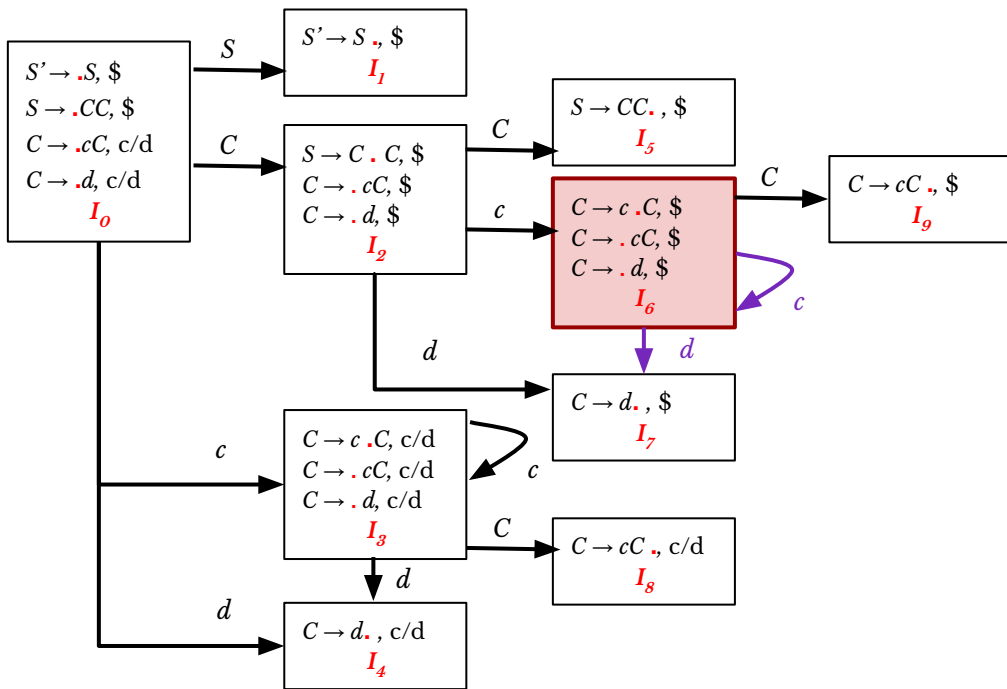
Rule 2a: If $[A \rightarrow \alpha.ab\beta, b]$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
 Action $[i, a]$ = "shift j " or "s j "

Action[3, c] = s3

Action[3, d] = s4

LR(1) parsing table

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



State	Action			GOTO	
	c	d	\$	S	C
0	s3	s4		1	2
1					
2	s6	s7			5
3	s3	s4			8
4					
5					
6	s6	s7			9
7					
8					
9					

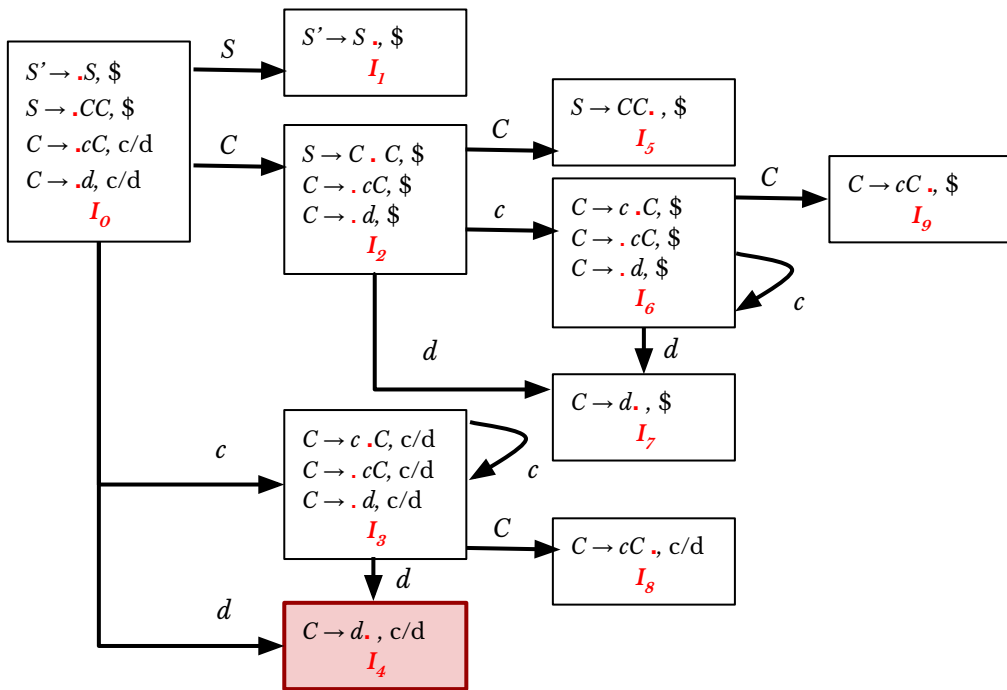
Rule 2a: If $[A \rightarrow \alpha \cdot a \beta, b]$ is in I_i and $\text{GOTO}(I_i, a) = I_j$
 Action $[i, a]$ = "shift j " or "s j "

Action[6, c] = s6

Action[6, d] = s7

LR(1) parsing table

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



State	Action			GOTO	
	c	d	\$	S	C
0	s3	s4		1	2
1					
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5					
6	s6	s7			9
7					
8					
9					

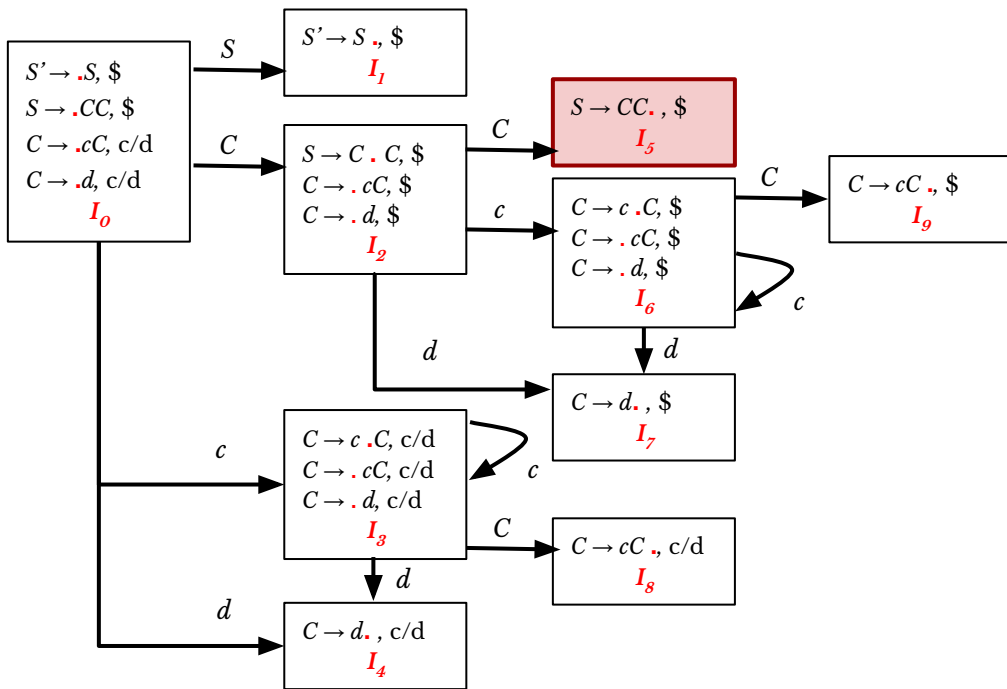
Rule 2b: If $[A \rightarrow \alpha., b]$ is in I_i
 Action $[i, b]$ = "reduce $A \rightarrow \alpha$ "

Action $[4, c]$ = r3

Action $[4, d]$ = r3

LR(1) parsing table

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



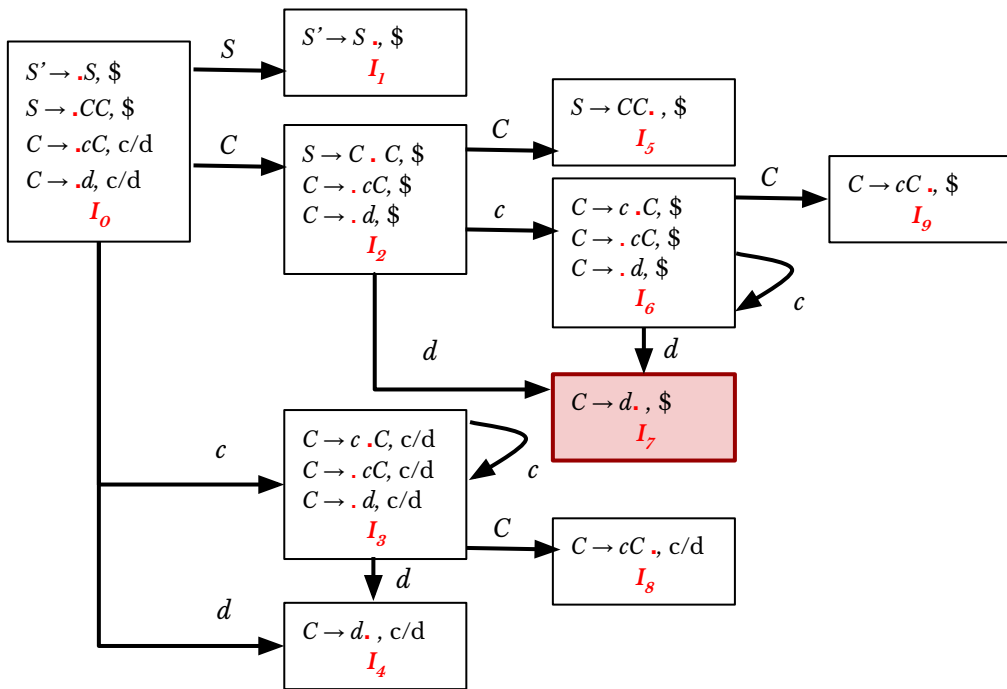
State	Action			GOTO	
	c	d	\$	S	C
0	s3	s4		1	2
1					
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5			r1		
6	s6	s7			9
7					
8					
9					

Rule 2b: If $[A \rightarrow \alpha., b]$ is in I_i
 Action $[i, b]$ = "reduce $A \rightarrow \alpha$ "

Action[5, \$] = r1

LR(1) parsing table

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



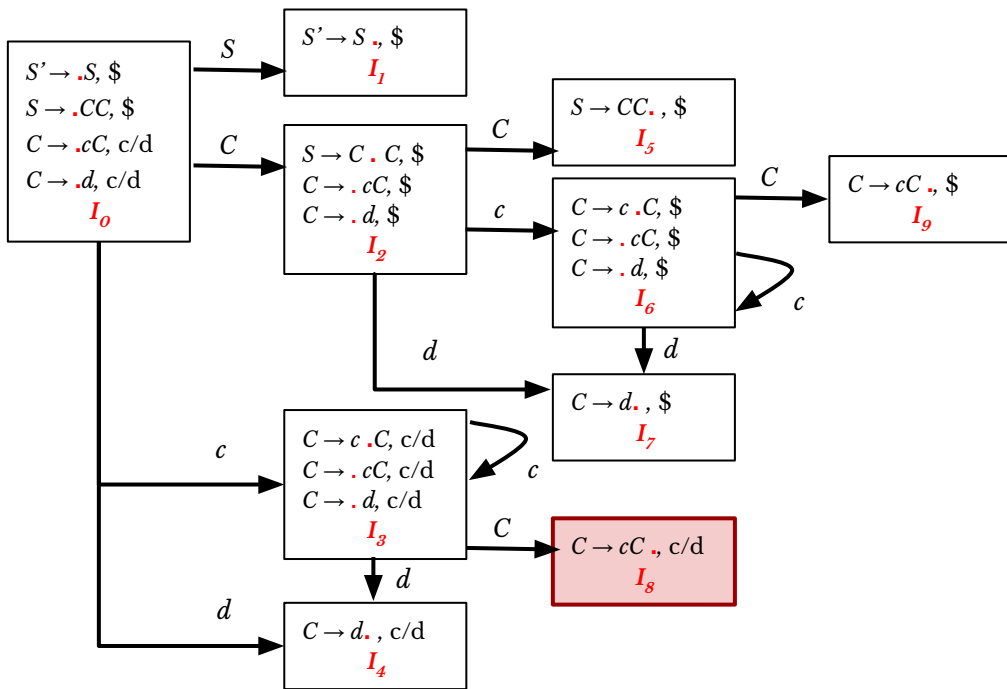
State	Action			GOTO	
	c	d	\$	S	C
0	s3	s4		1	2
1					
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5			r1		
6	s6	s7			9
7			r3		
8					
9					

Rule 2b: If $[A \rightarrow \alpha, b]$ is in I_i
 Action $[i, b]$ = "reduce $A \rightarrow \alpha$ "

Action[7, \$] = r3

LR(1) parsing table

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



State	Action			GOTO	
	c	d	\$	S	C
0	s3	s4		1	2
1					
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5			r1		
6	s6	s7			9
7			r3		
8	r2	r2			
9					

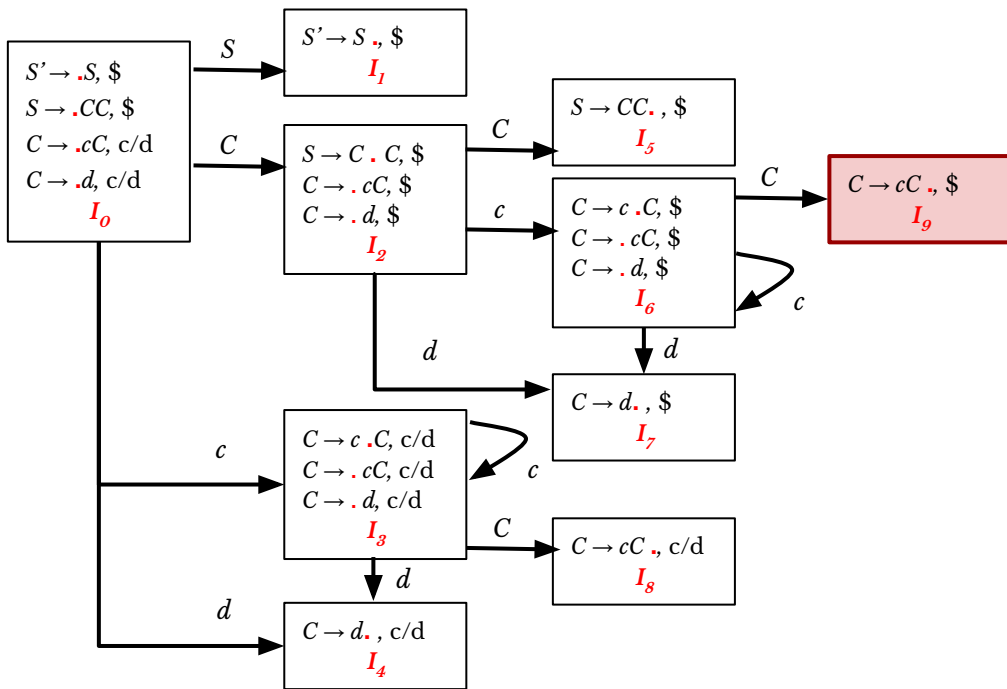
Rule 2b: If $[A \rightarrow \alpha, b]$ is in I_i
 Action $[i, b]$ = "reduce $A \rightarrow \alpha$ "

Action $[8, c]$ = r2

Action $[8, d]$ = r2

LR(1) parsing table

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



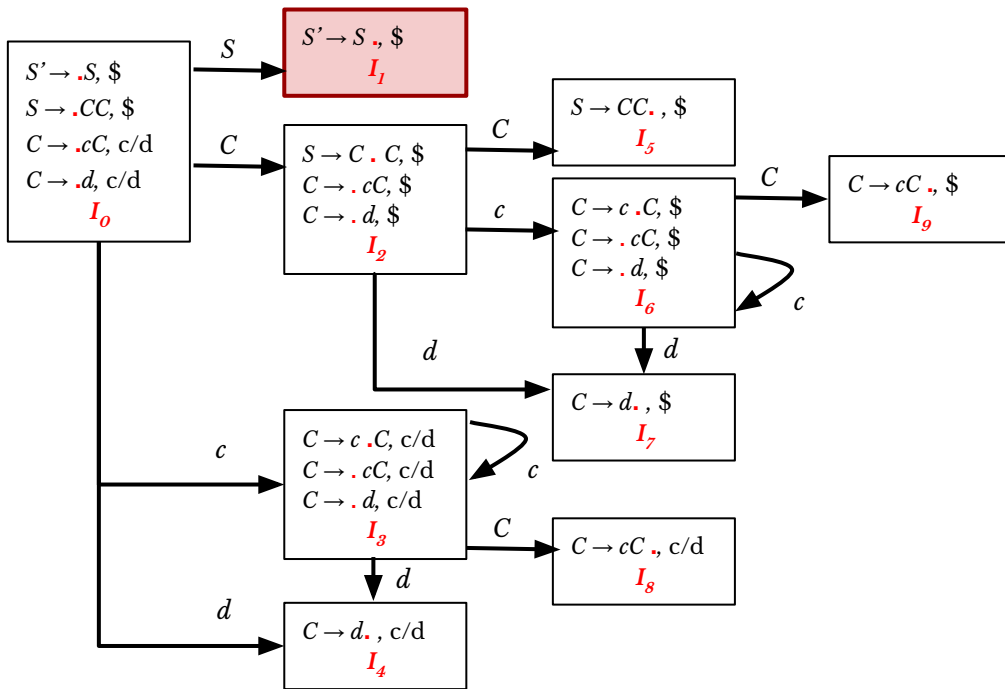
State	Action			GOTO	
	c	d	\$	S	C
0	s3	s4		1	2
1					
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5			r1		
6	s6	s7			9
7			r3		
8	r2	r2			
9			r2		

Rule 2b: If $[A \rightarrow \alpha., b]$ is in I_i
 Action $[i, b]$ = "reduce $A \rightarrow \alpha$ "

Action $[9, \$]$ = r2

LR(1) parsing table

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



Rule 2c: If $[S' \rightarrow S., \$]$ is in I_i
 Action $[i, \$]$ = "accept"

Action $[1, \$]$ = accept

State	Action			GOTO	
	c	d	\$	S	C
0	s3	s4		1	2
1			acc		
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5			r1		
6	s6	s7			9
7			r3		
8	r2	r2			
9			r2		

All empty entries are "error" case.

LR(1) parser: Parsing

LR(1) parsing

- Parsing algorithm for all LR parsers are same.
- Exercise:
 - Input: “cccdddd”

LALR parsing

- Often used in practice
- Most common syntactic constructs of programming languages can be expressed conveniently by an LALR.
- SLR and LALR tables always have the same number of states
 - Roughly, several hundred states for the C language
- Table size is considerably small than canonical LR or LR(1)
 - Canonical LR has, roughly, several thousand states for the same language.

For item $[A \rightarrow \alpha \cdot a \beta, b]$

Core: $A \rightarrow \alpha \cdot a \beta$

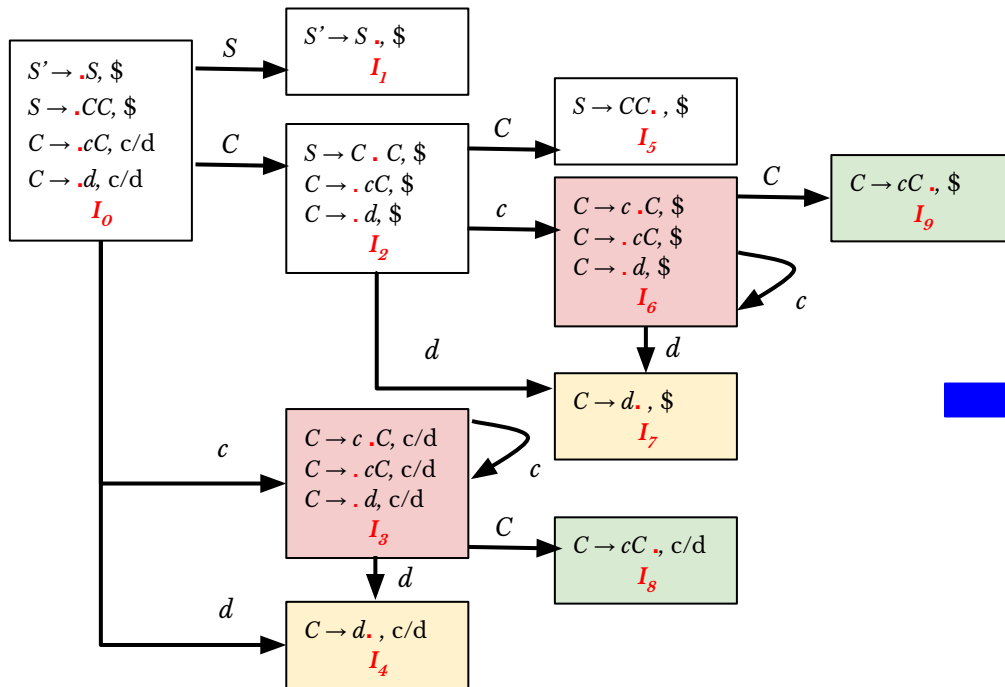
Lookahead: b

Constructing LALR parsers

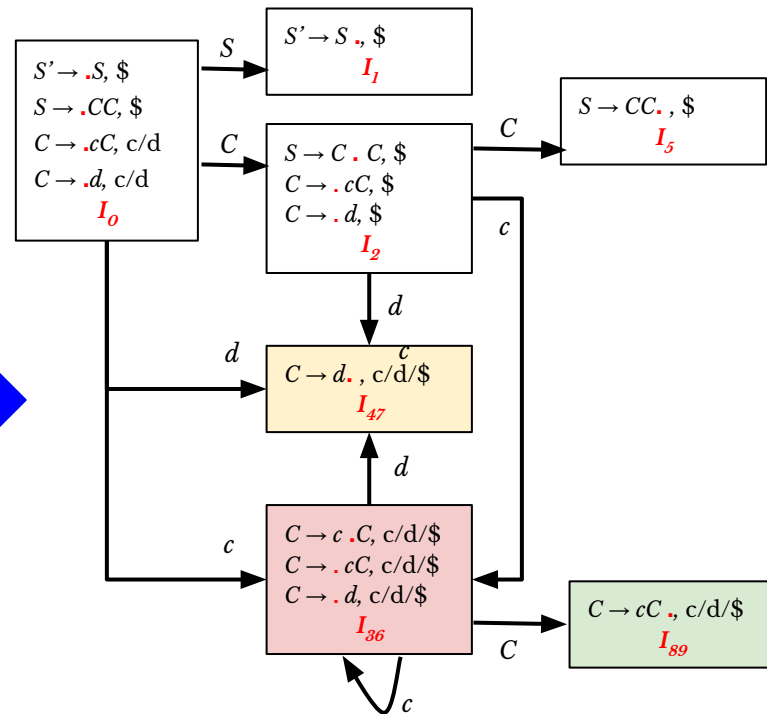
- Canonical set of items for LALR automation
 - Look for the item sets I_i and I_j in LR(1) automation, such that
 - Cores of $\text{Item}(I_i) == \text{Cores of Item}(I_j)$, with different lookahead symbols
 - $I_4 = [C \rightarrow d \cdot, c/d]$
 - $I_7 = [C \rightarrow d \cdot, \$]$
 - Merge the item sets I_i and I_j into I_{ij} , such that
 - Item set I_{ij} contains all items, with lookahead symbols merged
 - $I_{47} = [C \rightarrow d \cdot, c/d/\$]$

LALR Automation

0: $S' \rightarrow S$ 1: $S \rightarrow CC$ 2: $C \rightarrow cC$ 3: $C \rightarrow d$



LR Automation



LALR Automation

LALR parsing table

0: $S' \rightarrow S$ **1:** $S \rightarrow CC$ **2:** $C \rightarrow c C$ **3:** $C \rightarrow d$

State	Action			GOTO	
	c	d	\$	S	C
0	s3	s4		1	2
1			acc		
2	s36	s47			5
36	s36	s47			89
47	r3	r3	r3		
5			r1		
89	r2	r2	r2		

LALR grammar

- Its possible to introduce reduce/reduce conflicts during merger

$$I_1 = \{[A \rightarrow \alpha., a], [B \rightarrow \beta., b]\} \quad I_2 = \{[A \rightarrow \alpha., b], [B \rightarrow \beta., c]\}$$

$$\Rightarrow I_{12} = \{[A \rightarrow \alpha., a/b], [B \rightarrow \beta., b/c]\}$$

Action(12, b) = reduce with A or B??

- Cannot introduce a shift/reduce conflict

- Suppose the merged item introduced a shift/reduce conflict, e.g., on symbol a

$$I_{34} = \{[A \rightarrow \alpha., a/b], [B \rightarrow \beta.a \gamma, b/c]\} \quad \text{Action}(34, a) = \text{shift / reduce?}$$

- This means that we had two items in the LR(1) set as

$$I_3 = \{[A \rightarrow \alpha., a], [B \rightarrow \beta.a \gamma, b]\} \quad I_4 = \{[A \rightarrow \alpha., b], [B \rightarrow \beta.a \gamma, c]\}$$

- Observe, there is a shift/reduce conflict prior to the merge operation, i.e., the original grammar was not LR(1).

- Grammar is LALR(1), if no conflicts are introduced

Parser and Ambiguous grammar

LR parsers for Ambiguous grammars

- Grammars for the construction of LR-parsing tables must be unambiguous
- *Can we create LR-parsing tables for ambiguous grammars?*
 - Yes, but they will have conflicts
 - What if, we can resolve these conflicts in favor of one of them to disambiguate the grammar?
 - At the end, we will have again an unambiguous grammar
- *Why we want to use an ambiguous grammar?*
 - Some of the ambiguous grammars are much natural, and a corresponding unambiguous grammar can be very complex
 - Usage of an ambiguous grammar may eliminate unnecessary reductions
 - $E \Rightarrow T \Rightarrow F \Rightarrow \text{id}$ (using unambiguous expression-grammar)
 - $E \Rightarrow \text{id}$ (using ambiguous expression-grammar)

SLR(1) Automation

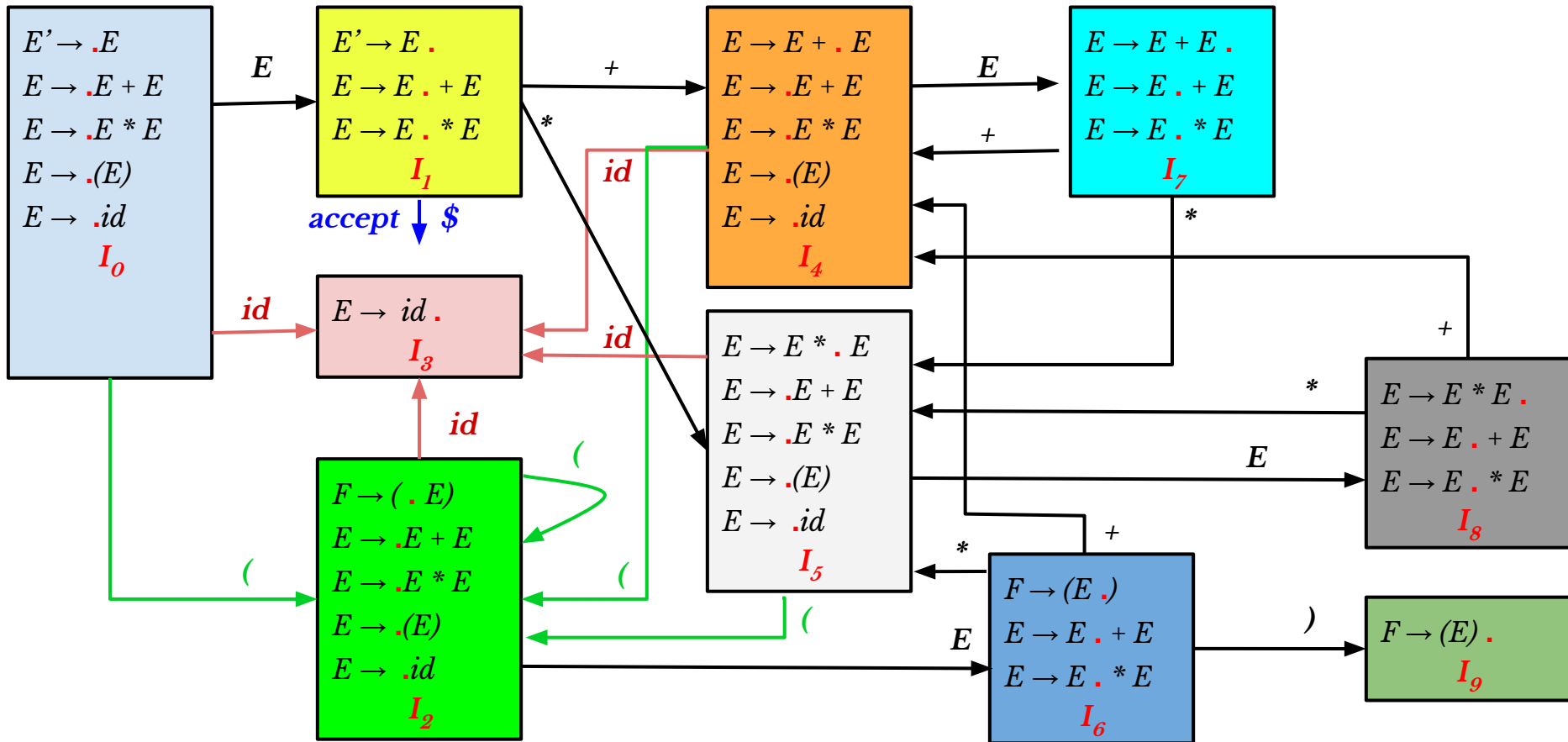
0: $E' \rightarrow E$

1: $E \rightarrow E + E$

2: $E \rightarrow E * E$

3: $E \rightarrow (E)$

4: $E \rightarrow id$



SLR parsing table

State	Action						GOTO
	id	+	*	()	\$	E
0	s3			s2			1
1		s4	s5			acc	
2	s3			s2			6
3		r4	r4		r4	r4	
4	s3			s2			7
5	s3			s2			8
6		s4	s5		s9		
7		s4/r1	s5/r1		r1	r1	
8		s4/r2	s5/r2		r2	r2	
9		r3	r3		r3	r3	

LR parsers for Ambiguous grammars

Why ambiguity is a problem?

We have a decision to make and not sure which parse tree to pick?

- Ambiguous grammars G: $E \rightarrow E + E \mid E * E \mid (E) \mid id$
- We can have two parse trees for an input

- Input 1: $id + id + id$

P1: $E \Rightarrow E + E \Rightarrow id + E \Rightarrow id + E + E \Rightarrow id + id + E \Rightarrow id + id + id$

P2: $E \Rightarrow E + E \Rightarrow E + E + E \Rightarrow id + E + E \Rightarrow id + id + E \Rightarrow id + id + id$

- Input 2: $id * id * id$

P1: $E \Rightarrow E * E \Rightarrow id * E \Rightarrow id * E * E \Rightarrow id * id * E \Rightarrow id * id * id$

P2: $E \Rightarrow E * E \Rightarrow E * E * E \Rightarrow id * E * E \Rightarrow id * id * E \Rightarrow id * id * id$

- Input 3: $id + id * id$

P1: $E \Rightarrow E + E \Rightarrow id + E \Rightarrow id + E * E \Rightarrow id + id * E \Rightarrow id + id * id$

P2: $E \Rightarrow E * E \Rightarrow E + E * E \Rightarrow id + E * E \Rightarrow id + id * E \Rightarrow id + id * id$

- Input 4: $id * id + id$

P1: $E \Rightarrow E * E \Rightarrow id * E \Rightarrow id * E + E \Rightarrow id * id + E \Rightarrow id * id + id$

P2: $E \Rightarrow E + E \Rightarrow E * E + E \Rightarrow id * E + E \Rightarrow id * id + E \Rightarrow id * id + id$

Computation-wise, which parse tree is correct?

Input 1:

Either P1 or P2

Input 2:

Either P1 or P2

Input 3:

P1

Input 4:

P2

Reason?

Operation '*' has precedence over operator '+'

LR parsers for Ambiguous grammars

- In the parsing-table of an ambiguous grammars, if we can explicitly resolve the conflicts, then the processing of parsing remains unambiguous
 - E.g.,
 - We can look for the precedence of operators for the conflict resolution, 'OR'
 - We can look for the associativity of operators for the conflict resolution

In state 7, we have shift/reduce conflicts for symbols + and *

$$I_0 \xRightarrow{E} I_1 \xRightarrow{+} I_4 \xRightarrow{E} I_7$$

If current input symbol is +

Shift → if + right associative

Reduce → if + left associative

If current input symbol is *

Shift → if * has higher precedence over +

Reduce → if + has higher precedence over *

In state 8, we have shift/reduce conflicts for symbols + and *

$$I_0 \xRightarrow{E} I_1 \xRightarrow{*} I_4 \xRightarrow{E} I_7$$

If current input symbol is *

Shift → if * right associative

Reduce → if * left associative

If current input symbol is +

Shift → if + has higher precedence over *

Reduce → if * has higher precedence over +

SLR parsing table

Conflict-free parsing table for ambiguous grammar.

State	Action						GOTO
	id	+	*	()	\$	E
0	s3			s2			1
1		s4	s5			acc	
2	s3			s2			6
3		r4	r4		r4	r4	
4	s3			s2			7
5	s3			s2			8
6		s4	s5		s9		
7		r1	s5		r1	r1	
8		r2	r2		r2	r2	
9		r3	r3		r3	r3	