

Compiler course

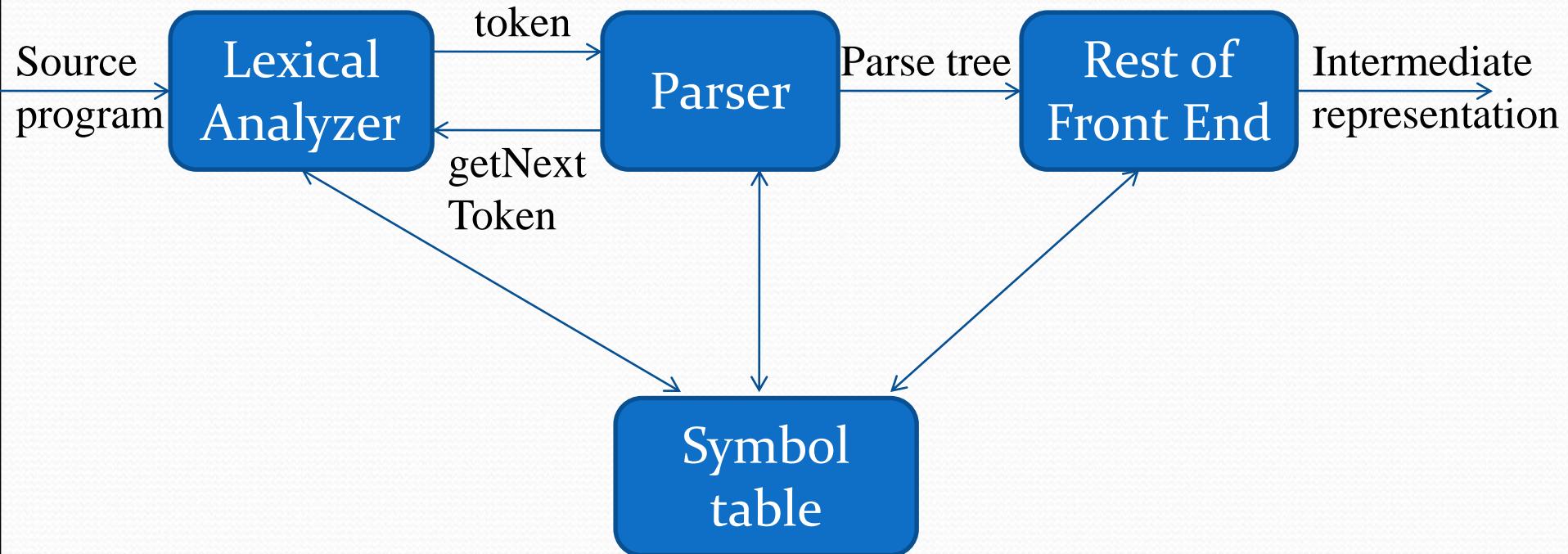
Chapter 4
Syntax Analysis

Outline

- Role of parser
- Context free grammars
- Top down parsing
- Bottom up parsing
- Parser generators

- Checking the structure of source program .
- Known as parsing
- 2nd phase of compiler
- Checks the syntax and produced IC code for the representation of source program.
- Why CFG used for syntax checking?

The role of parser



Uses of grammars

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

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

$E \rightarrow TE'$

$E' \rightarrow +TE' \mid \epsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' \mid \epsilon$

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

Error handling

- Common programming errors
 - Lexical errors
 - Syntactic errors
 - Semantic errors
- Error handler goals
 - 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 recovery:

Panic Mode Recovery

- In this method, successive characters from input are removed one at a time until a designated set of synchronizing tokens is found. Synchronizing tokens are deli-meters such as ; or }
- Advantage is that its easy to implement and guarantees not to go to infinite loop
- Disadvantage is that a considerable amount of input is skipped without checking it for additional errors

Statement Mode recovery

- In this method, when a parser encounters an error, it performs necessary correction on remaining input so that the rest of input statement allow the parser to parse ahead.
- The correction can be deletion of extra semicolons, replacing comma by semicolon or inserting missing semicolon.
- While performing correction, atmost care should be taken for not going in infinite loop.

Error production

- If user has knowledge of common errors that can be encountered then, these errors can be incorporated by augmenting the grammar with error productions that generate erroneous constructs.
- If this is used then, during parsing appropriate error messages can be generated and parsing can be continued.
- Disadvantage is that its difficult to maintain.

Global correction

- The parser examines the whole program and tries to find out the closest match for it which is error free.
- The closest match program has less number of insertions, deletions and changes of tokens to recover from erroneous input.
- Due to high time and space complexity, this method is not implemented practically.

- For Given incorrect input X and grammar G , it will find parse tree for related string y such that no of insertions, deletions and changes of tokens required to transform x into y is as small as possible

Context free grammars

- Terminals
- Nonterminals
- Start symbol
- productions

expression \rightarrow expression + term

expression \rightarrow expression – term

expression \rightarrow term

term \rightarrow term * factor

term \rightarrow term / factor

term \rightarrow factor

factor \rightarrow (expression)

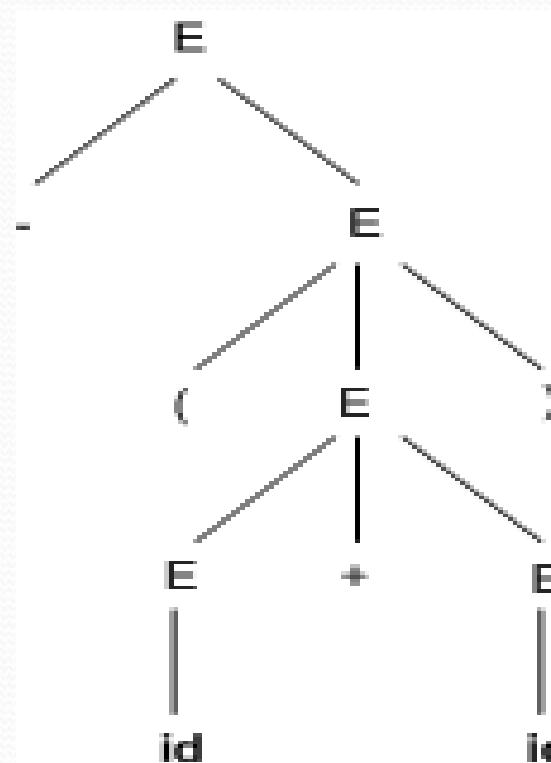
factor \rightarrow **id**

Derivations

- Productions are treated as rewriting rules to generate a string
- Rightmost and leftmost derivations
 - $E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid \mathbf{id}$
 - Derivations for $-(\mathbf{id}+\mathbf{id})$
 - $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(E+E+E) \Rightarrow -(E+E+E+E) = -(id+id)$

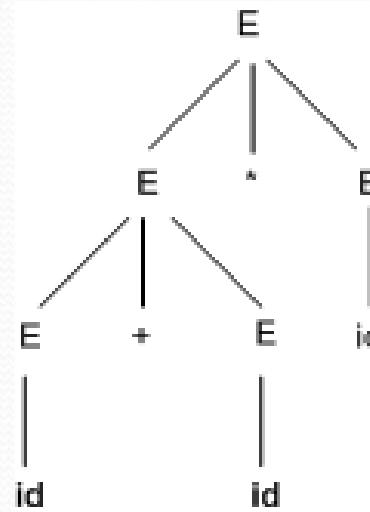
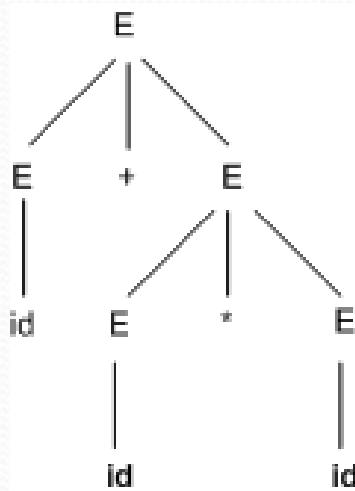
Parse trees

- $-(\mathbf{id}+\mathbf{id})$
- $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(\mathbf{id}+E) \Rightarrow -(\mathbf{id}+\mathbf{id})$



Ambiguity

- For some strings there exist more than one parse tree
- Or more than one leftmost derivation
- Or more than one rightmost derivation
- Example: $\text{id} + \text{id} * \text{id}$



Elimination of left recursion

- A grammar is left recursive if it has a non-terminal A such that there is a derivation $A \Rightarrow^+ A \alpha$
- Top down parsing methods can't handle left-recursive grammars
- A simple rule for direct left recursion elimination:
 - For a rule like:
 - $A \rightarrow A \alpha | \beta$
 - We may replace it with
 - $A \rightarrow \beta A'$
 - $A' \rightarrow \alpha A' | \varepsilon$

Eliminate left recursion

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

$B \rightarrow Be / b$

2. $S \rightarrow (L) / a$

$L \rightarrow L, S / S$

3. $A \rightarrow Ba / Aa / c$

$B \rightarrow Bb / Ab / d$

Left factoring

- Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive or top-down parsing.
- Consider following grammar:
 - Stmt \rightarrow if expr **then** stmt **else** stmt
 - | if expr **then** stmt
- On seeing input **if** it is not clear for the parser which production to use
- We can easily perform left factoring:
 - If we have $A \rightarrow \alpha \beta 1 \mid \alpha \beta 2$ then we replace it with
 - $A \rightarrow \alpha A'$
 - $A' \rightarrow \beta 1 \mid \beta 2$

Left factoring (cont.)

- Example:

1. $S \rightarrow I E t S \mid i E t S e S \mid a$

$E \rightarrow b$

2. $A \rightarrow aAB \mid aBc \mid aAc$

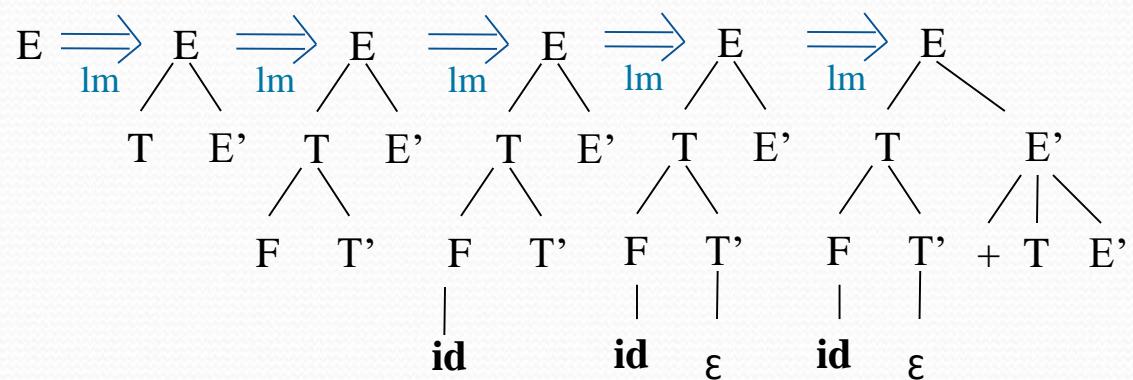
3. $S \rightarrow bSSaaS \mid bSSaSb \mid bSb \mid a$

Top Down Parsing

Introduction

- A Top-down parser tries to create a parse tree from the root towards the leafs scanning input from left to right
- It can be also viewed as finding a leftmost derivation for an input string
- Example: $\text{id} + \text{id} * \text{id}$

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



Recursive descent parsing

- Consists of a set of procedures, one for each nonterminal
- Execution begins with the procedure for start symbol
- A typical procedure for a non-terminal

```
void A() {  
    choose an A-production, A->X1X2..Xk  
    for (i=1 to k) {  
        if (Xi is a nonterminal)  
            call procedure Xi();  
        else if (Xi equals the current input symbol a)  
            advance the input to the next symbol;  
        else /* an error has occurred */  
    }  
}
```

Recursive descent parsing (cont)

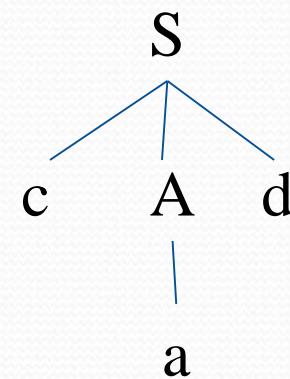
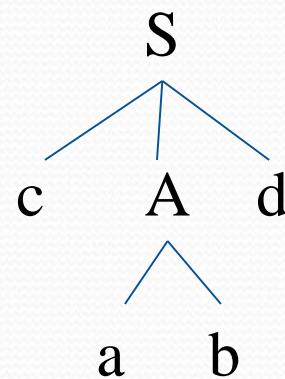
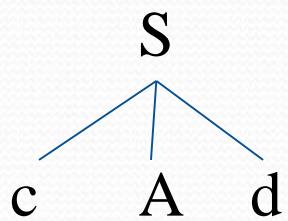
- General recursive descent may require backtracking
- The previous code needs to be modified to allow backtracking
- In general form it can't choose an A-production easily.
- So we need to try all alternatives
- If one failed the input pointer needs to be reset and another alternative should be tried
- Recursive descent parsers can't be used for left-recursive grammars

Example

$S \rightarrow cAd$

$A \rightarrow ab \mid a$

Input: cad



1. $E \rightarrow iE'$

$E' \rightarrow +iE' | \text{null}$

2. $S \rightarrow AB$

$A \rightarrow c | aB$

$B \rightarrow d.$

To understand first and follow

- Each time a predictive parser makes a decision, it needs to determine which production rule to apply to the leftmost non-terminal in an intermediate form, based on the next terminal

- **FIRST?**
- We saw the need of backtrack in the [previous article](#) of on Introduction to Syntax Analysis, which is really a complex process to implement. There can be easier way to sort out this problem:
- If the compiler would have come to know in advance, that what is the “first character of the string produced when a production rule is applied”, and comparing it to the current character or token in the input string it sees, it can wisely take decision on which production rule to apply.

- $S \rightarrow aAb$
- $A \rightarrow a \mid \langle\epsilon\rangle$

- **FOLLOW?**

The parser faces one more problem. Let us consider below grammar to understand this problem.

- $A \rightarrow aBb$
- $B \rightarrow c \mid \epsilon$

First and Follow

- First() is set of terminals that begins strings derived from
- If $\alpha \stackrel{*}{\Rightarrow} \varepsilon$ then ε is also in First(ε)
- In predictive parsing when we have $A \rightarrow \alpha | \beta$, if First(α) and First(β) are disjoint sets then we can select appropriate A-production by looking at the next input
- Follow(A), for any nonterminal A, is set of terminals a that can appear immediately after A in some sentential form
 - If we have $S \stackrel{*}{\Rightarrow} \alpha A a \beta$ for some α and β then a is in Follow(A)
- If A can be the rightmost symbol in some sentential form, then \$ is in Follow(A)

Computing First

- To compute $\text{First}(X)$ for all grammar symbols X , apply following rules until no more terminals or ε can be added to any First set:
 1. If X is a terminal then $\text{First}(X) = \{X\}$.
 2. If X is a nonterminal and $X \rightarrow Y_1 Y_2 \dots Y_k$ is a production for some $k \geq 1$, then place a in $\text{First}(X)$ if for some i , a is in $\text{First}(Y_i)$ and ε is in all of $\text{First}(Y_1), \dots, \text{First}(Y_{i-1})$ that is $Y_1 \dots Y_{i-1} \xrightarrow{*} \varepsilon$. if ε is in $\text{First}(Y_j)$ for $j=1, \dots, k$ then add ε to $\text{First}(X)$.
 3. If $X \rightarrow \varepsilon$ is a production then add ε to $\text{First}(X)$
- Example! *

Rules for First Sets

1. If X is a terminal then $\text{First}(X)$ is just X !
2. If there is a Production $X \rightarrow \epsilon$ then add ϵ to $\text{first}(X)$
3. If there is a Production $X \rightarrow Y_1 Y_2 .. Y_k$ then add $\text{first}(Y_1 Y_2 .. Y_k)$ to $\text{first}(X)$
4. $\text{First}(Y_1 Y_2 .. Y_k)$ is either
 1. $\text{First}(Y_1)$ (if $\text{First}(Y_1)$ doesn't contain ϵ)
 2. OR (if $\text{First}(Y_1)$ does contain ϵ) then $\text{First}(Y_1 Y_2 .. Y_k)$ is Everything in $\text{First}(Y_1)$ <except for ϵ > as well as everything in $\text{First}(Y_2 .. Y_k)$
 3. If $\text{First}(Y_1) \text{First}(Y_2) .. \text{First}(Y_k)$ all contain ϵ then add ϵ to $\text{First}(Y_1 Y_2 .. Y_k)$ as well.

Computing follow

- To compute $\text{First}(A)$ for all nonterminals A, apply following rules until nothing can be added to any follow set:
 1. Place \$ in $\text{Follow}(S)$ where S is the start symbol
 2. If there is a production $A \rightarrow \alpha B \beta$ then everything in $\text{First}(\beta)$ except ϵ is in $\text{Follow}(B)$.
 3. If there is a production $A \rightarrow B$ or a production $A \rightarrow \alpha B \beta$ where $\text{First}(\beta)$ contains ϵ , then everything in $\text{Follow}(A)$ is in $\text{Follow}(B)$
- Example!

Rules for Follow Sets

1. First put \$ (the end of input marker) in $\text{Follow}(S)$ (S is the start symbol)
2. If there is a production $A \rightarrow aBb$, (where a can be a whole string) then everything in $\text{FIRST}(b)$ except for ϵ is placed in $\text{FOLLOW}(B)$.
3. If there is a production $A \rightarrow aB$, then everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(B)$
4. If there is a production $A \rightarrow aBb$, where $\text{FIRST}(b)$ contains ϵ , then everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(B)$

Examples

1)

$S \rightarrow aAB|bA|NULL$

$A \rightarrow aAb|NULL$

$B \rightarrow bB|c$

2)

$S \rightarrow iCtSS'|a$

$S' \rightarrow eS|NULL$

$C \rightarrow b$

3.

$S \rightarrow aBDh$

$B \rightarrow cC$

$C \rightarrow bC \mid \text{NULL}$

$D \rightarrow EF$

$E \rightarrow g \mid \text{NULL}$

$F \rightarrow f \mid \text{NULL}$

4.

$S \rightarrow A$

$A \rightarrow aB \mid Ad$

$B \rightarrow b$

$C \rightarrow g$

LL(1) Grammars

- Predictive parsers are those recursive descent parsers needing no backtracking
- Grammars for which we can create predictive parsers are called LL(1)
 - The first L means scanning input from left to right
 - The second L means leftmost derivation
 - And 1 stands for using one input symbol for lookahead
- 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:
 - For no terminal a do α and β both derive strings beginning with a
 - At most one of α or β can derive empty string
 - If $\alpha \Rightarrow^* \varepsilon$ then β does not derive any string beginning with a terminal in $\text{Follow}(A)$.

Construction of predictive parsing table

- For each production $A \rightarrow \alpha$ in grammar do the following:
 1. For each terminal a in $\text{First}(\alpha)$ add $A \rightarrow a$ in $M[A,a]$
 2. If ϵ is in $\text{First}(\alpha)$, then for each terminal b in $\text{Follow}(A)$ add $A \rightarrow \epsilon$ to $M[A,b]$. If ϵ is in $\text{First}(\alpha)$ and $\$$ is in $\text{Follow}(A)$, add $A \rightarrow \epsilon$ to $M[A,\$]$ as well
- If after performing the above, there is no production in $M[A,a]$ then set $M[A,a]$ to error

Example

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

	First	Follow
F	$\{ (, id \}$	$\{ +, *,), \$ \}$
T	$\{ (, id \}$	$\{ +,), \$ \}$
E	$\{ (, id \}$	$\{), \$ \}$
E'	$\{ +, \epsilon \}$	$\{), \$ \}$
T'	$\{ *, \epsilon \}$	$\{ +,), \$ \}$

Non-terminal	Input Symbol					
	id	+	*	()	\$
E	$E \rightarrow TE'$				$E \rightarrow TE'$	
E'		$E' \rightarrow +TE'$				$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$				$T \rightarrow FT'$	
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$			$T' \rightarrow \epsilon$
F	$F \rightarrow id$				$F \rightarrow (E)$	

Another example

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow eS \mid \epsilon$$

$$E \rightarrow b$$

Non-terminal	Input Symbol					
	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow \epsilon$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

Check the following grammar is LL(1) or not?

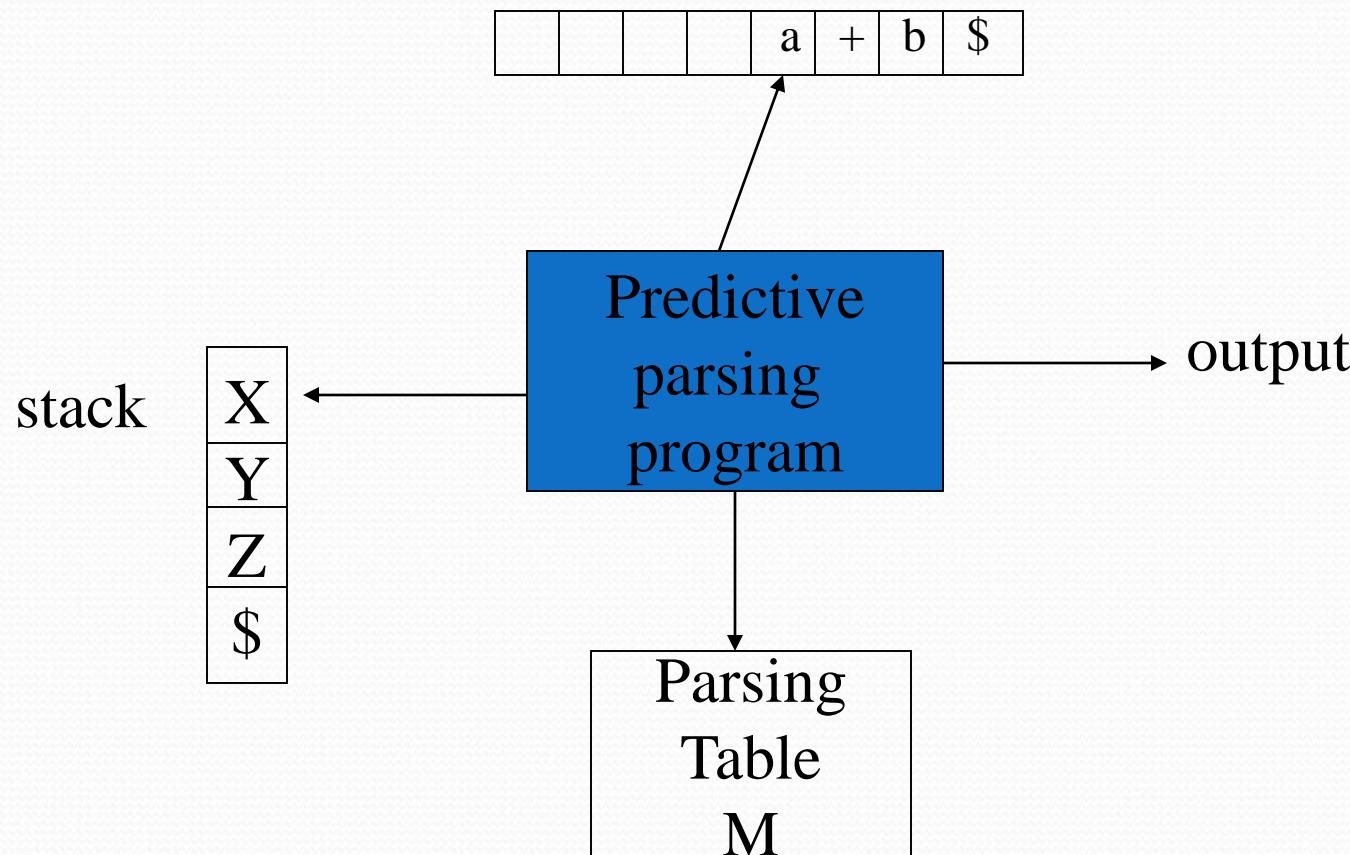
If grammar is LL1 then parse the given string: aaaaab#

$S \rightarrow S\# \mid aA \mid b \mid cB \mid d$

$A \rightarrow aA \mid b$

$B \rightarrow cB \mid d$

Non-recursive predicting parsing



Predictive parsing algorithm

Set ip point to the first symbol of w;

Set X to the top stack symbol;

While ($X \neq \$$) { /* stack is not empty */

 if (X is a) pop the stack and advance ip;

 else if (X is a terminal) error();

 else if ($M[X,a]$ is an error entry) error();

 else if ($M[X,a] = X \rightarrow Y_1 Y_2 \dots Y_k$) {

 output the production $X \rightarrow Y_1 Y_2 \dots Y_k$;

 pop the stack;

 push Y_k, \dots, Y_2, Y_1 on to the stack with Y_1 on top;

 }

 set X to the top stack symbol;

}

Example

- $\text{id} + \text{id} * \text{id} \$$

Matched	Stack	Input	Action
	E\$	$\text{id} + \text{id} * \text{id} \$$	

Error recovery in predictive parsing

- Panic mode
 - Place all symbols in $\text{Follow}(A)$ into synchronization set for nonterminal A: skip tokens until an element of $\text{Follow}(A)$ is seen and pop A from stack.
 - Add to the synchronization set of lower level construct the symbols that begin higher level constructs
 - Add symbols in $\text{First}(A)$ to the synchronization set of nonterminal A
 - If a nonterminal can generate the empty string then the production deriving can be used as a default
 - If a terminal on top of the stack cannot be matched, pop the terminal, issue a message saying that the terminal was inserted

Example

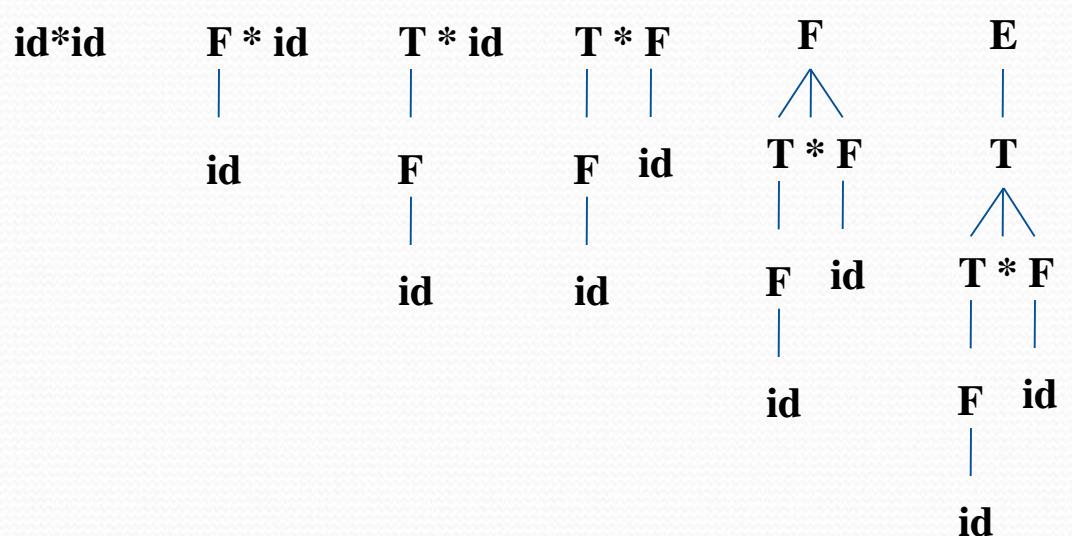
Non-terminal	Input Symbol					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$	synch	synch
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$	synch		$T \rightarrow FT'$	synch	synch
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$	synch	synch	$F \rightarrow (E)$	synch	synch

Stack	Input	Action
E\$	+id*+id\$	Error, Skip +
E\$	id*+id\$	id is in First(E)
TE'\$	id*+id\$	
FT'E'\$	id*+id\$	
idT'E'\$	id*+id\$	
T'E'\$	*+id\$	
*FT'E'\$	*+id\$	
FT'E'\$	+id\$	Error, M[F,+]=synch
T'E'\$	+id\$	F has been popped

Bottom-up Parsing

Introduction

- Constructs parse tree for an input string beginning at the leaves (the bottom) and working towards the root (the top)
- Example: id^*id

$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow (E) \mid \text{id} \end{array}$$


Shift-Reduce Parsing

- Shift-reduce parsing uses two unique steps for bottom-up parsing. These steps are known as shift-step and reduce-step.
- **Shift step:** The shift step refers to the advancement of the input pointer to the next input symbol, which is called the shifted symbol. This symbol is pushed onto the stack. The shifted symbol is treated as a single node of the parse tree.
- **Reduce step :** When the parser finds a complete grammar rule (RHS) and replaces it to (LHS), it is known as reduce-step. This occurs when the top of the stack contains a handle. To reduce, a POP function is performed on the stack which pops off the handle and replaces it with LHS non-terminal symbol.

Shift-reduce parser

- The general idea is to shift some symbols of input to the stack until a reduction can be applied
- At each reduction step, a specific substring matching the body of a production is replaced by the nonterminal at the head of the production
- The key decisions during bottom-up parsing are about when to reduce and about what production to apply
- A reduction is a reverse of a step in a derivation
- The goal of a bottom-up parser is to construct a derivation in reverse:

$S \rightarrow aABe$

$A \rightarrow Abc \mid b$

$B \rightarrow d$

The sentence abbcde can be reduced to S by ?

- $E \Rightarrow aAB \Rightarrow aAde \Rightarrow aAbcde \Rightarrow abbcde$
- Handles?
- Handle of a string is a substring that matches the right side of a production and whose reduction of the non terminal on the left side of the production represents one step along the reverse of the right most derivation.
- **In many cases the leftmost substring β that matches the right side of some production $A \rightarrow \beta$ is not a handle.**
- Why?

Consider below grammar

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

$E \rightarrow id$

Handle pruning

- A Handle is a substring that matches the body of a production and whose reduction represents one step along the reverse of a rightmost derivation

Right sentential form	Handle	Reducing production
$\text{id1} + \text{id2} * \text{id3}$	id1	$E \rightarrow \text{id}$
$\text{E} + \text{id2} * \text{id3}$	id2	$E \rightarrow \text{id}$
$\text{E} + \text{E} * \text{id3}$	id3	$E \rightarrow \text{id}$
$\text{E} + \text{E} * \text{E}$	$\text{E} * \text{E}$	$E \rightarrow \text{E} * \text{E}$
$\text{E} + \text{E}$	$\text{E} + \text{E}$	$E \rightarrow \text{E} + \text{E}$
E		

Shift reduce parsing

- A stack is used to hold grammar symbols
- Handle always appear on top of the stack
- Initial configuration:

Stack	Input
\$	w\$

- Acceptance configuration

Stack	Input
\$S	\$

Stack Implementation on Shift reduce parsing

- Basic operations:

- Shift
- Reduce
- Accept
- Error

- Example: id+id*id

- STACK

INPUT

ACTION`

Conflicts during shit reduce parsing

- Two kind of conflicts
 - Shift/reduce conflict
 - Reduce/reduce conflict
- Example:

```
stmt → If expr then stmt  
      | If expr then stmt else stmt  
      | other
```

Stack
... if expr then stmt

Input
else ...\$

Reduce/reduce conflict

stmt -> id(parameter_list)

stmt -> expr:=expr

parameter_list->parameter_list, parameter

parameter_list->parameter

parameter->id

expr->id(expr_list)

expr->id

expr_list->expr_list, expr

expr_list->expr

Stack

... id(id

Input
,id) ...\$

OPERATOR PRECEDENCE PARSING

- One of the bottom up parser
- It is mainly use to define the mathematical operators.
- What is operator grammar?
- A grammar which is use to define generally mathematical operations it is called as operator grammar.
- With some restrictions.?

Operator Grammar

- No ϵ -transition.
- No two adjacent non-terminals.

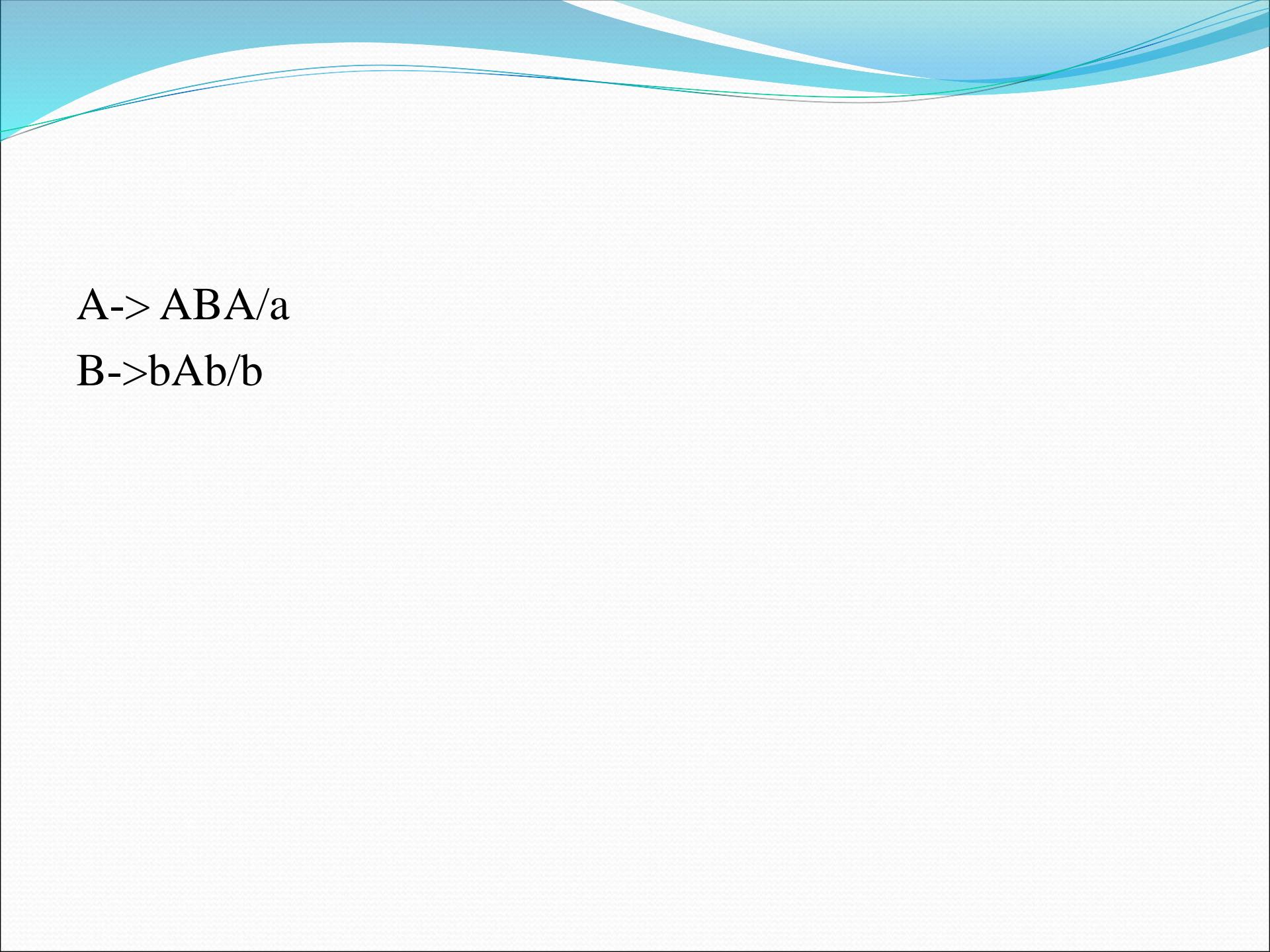
Eg.

$$E \rightarrow E \text{ op } E \mid id$$

$$\text{op} \rightarrow + \mid *$$

The above grammar is not an operator grammar but:

$$E \rightarrow E + E \mid E^* E \mid id$$



A->ABA/a

B->bAb/b

Operator Precedence

- If a has higher precedence over b; $a .> b$
- If a has lower precedence over b; $a <. b$
- If a and b have equal precedence; $a =. b$

Note:

- id has higher precedence than any other symbol
- \$ has lowest precedence.
- if two operators have equal precedence, then we check the **Associativity** of that particular operator.

Precedence Table

	id	+	*	\$
id		.>	.>	.>
+	<.	.>	<.	.>
*	<.	.>	.>	.>
\$	<.	<.	<.	.>

Example: w= \$id + id * id\$
 \$<.id.>+<.id.>*<.id.>\$

BASIC PRINCIPLE

- Scan input string left to right, try to detect .> and put a pointer on its location.
- Now scan backwards till reaching <.
- String between <. And .> is our handle.
- Replace handle by the head of the respective production.
- REPEAT until reaching start symbol.

Operator-Precedence Parsing Algorithm

set p to point to the first symbol of w\$;
repeat forever
if (\$ is on top of the stack and p points to \$) then return
else {
let a be the topmost terminal symbol on the stack and let b be the symbol
pointed to by p;
if (a <. b or a =. b) then { /* SHIFT */
push b onto the stack;
advance p to the next input symbol;
}
else if (a .> b) then /* REDUCE */
repeat pop stack
until (the top of stack terminal is related by <.
to the terminal most recently popped);
else error();
}

EXAMPLE

STACK	INPUT	ACTION/REMARK
\$	id + id * id\$	\$ <. Id
\$ id	+ id * id\$	id >. +
\$	+ id * id\$	\$ <. +
\$ +	id * id\$	+ <. Id
\$ + id	* id\$	id .> *
\$ +	* id\$	+ <. *
\$ + *	id\$	* <. Id
\$ + * id	\$	id .> \$
\$ + *	\$	* .> \$
\$ +	\$	+ .> \$
\$	\$	accept

PRECEDENCE FUNCTIONS

- Operator precedence parsers use **precedence functions** that map terminal symbols to integers.

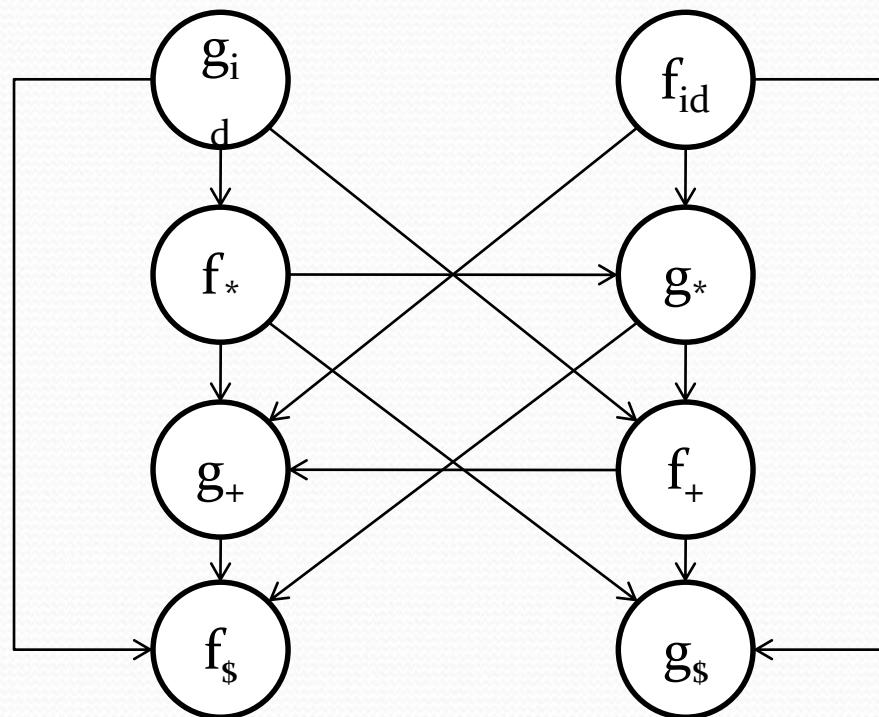
Algorithm for Constructing Precedence Functions

1. Create functions f_a for each grammar terminal a and for the end of string symbol.
2. Partition the symbols in groups so that f_a and g_b are in the same group if $a \cdot\cdot b$ (there can be symbols in the same group even if they are not connected by this relation).
3. Create a directed graph whose nodes are in the groups, next for each symbols a and b do: place an edge from the group of g_b to the group of f_a if $a < \cdot b$, otherwise if $a \cdot > b$ place an edge from the group of f_a to that of g_b .
4. If the constructed graph has a cycle then no precedence functions exist. When there are no cycles collect the length of the longest paths from the groups of f_a and g_b respectively.

- Consider the following table:

	id	+	*	\$
id		.>	.>	.>
+	<.	.>	<.	.>
*	<.	.>	.>	.>
\$	<.	<.	<.	.>

- Resulting graph:



- From the previous graph we extract the following precedence functions:

	id	+	*	\$
f	4	2	4	0
id	5	1	3	0

LR Parsing

- The most prevalent type of bottom-up parsers
- LR(k), mostly interested on parsers with $k \leq 1$
- Why LR parsers?
 - Table driven
 - Can be constructed to recognize all programming language constructs
 - Most general non-backtracking shift-reduce parsing method
 - Can detect a syntactic error as soon as it is possible to do so
 - Class of grammars for which we can construct LR parsers are superset of those which we can construct LL parsers

States of an LR parser

- States represent set of items
- An LR(0) item of G is a production of G with the dot at some position of the body:
 - For $A \rightarrow XYZ$ we have following items
 - $A \rightarrow .XYZ$
 - $A \rightarrow X.YZ$
 - $A \rightarrow XY.Z$
 - $A \rightarrow XYZ.$
 - In a state having $A \rightarrow .XYZ$ we hope to see a string derivable from XYZ next on the input.
 - What about $A \rightarrow X.YZ$?

Constructing canonical LR(0) item sets

- Augmented grammar:
 - G with addition of a production: $S' \rightarrow S$
- Closure of item sets:
 - If I is a set of items, $\text{closure}(I)$ is a set of items constructed from I by the following rules:
 - Add every item in I to $\text{closure}(I)$
 - If $A \rightarrow \alpha.B\beta$ is in $\text{closure}(I)$ and $B \rightarrow \gamma$ is a production then add the item $B \rightarrow .\gamma$ to $\text{closure}(I)$.

- Example:

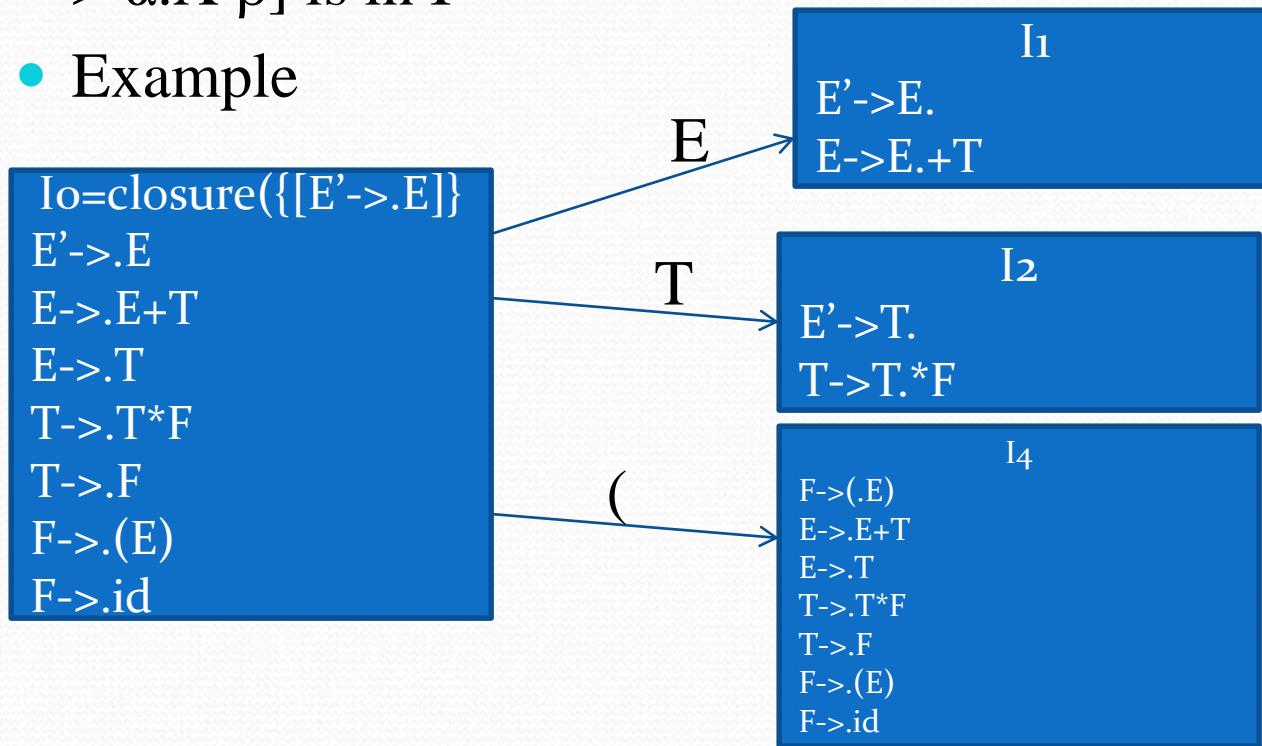
$E' \rightarrow E$
 $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid \mathbf{id}$

$I_0 = \text{closure}(\{[E' \rightarrow .E]\})$

$E' \rightarrow .E$
 $E \rightarrow .E + T$
 $E \rightarrow .T$
 $T \rightarrow .T * F$
 $T \rightarrow .F$
 $F \rightarrow .(E)$
 $F \rightarrow .\mathbf{id}$

Constructing canonical LR(0) item sets (cont.)

- Goto (I, X) where I is an item set and X is a grammar symbol is closure of set of all items $[A \rightarrow \alpha X. \beta]$ where $[A \rightarrow \alpha X \beta]$ is in I
- Example



Closure algorithm

SetOfItems CLOSURE(I) {

 J=I;

 repeat

 for (each item A-> α.Bβ in J)

 for (each production B->γ of G)

 if (B->.γ is not in J)

 add B->.γ to J;

 until no more items are added to J on one round;

 return J;

GOTO algorithm

```
SetOfItems GOTO(I,X) {  
    J=empty;  
    if (A-> α.X β is in I)  
        add CLOSURE(A-> αX. β ) to J;  
    return J;  
}
```

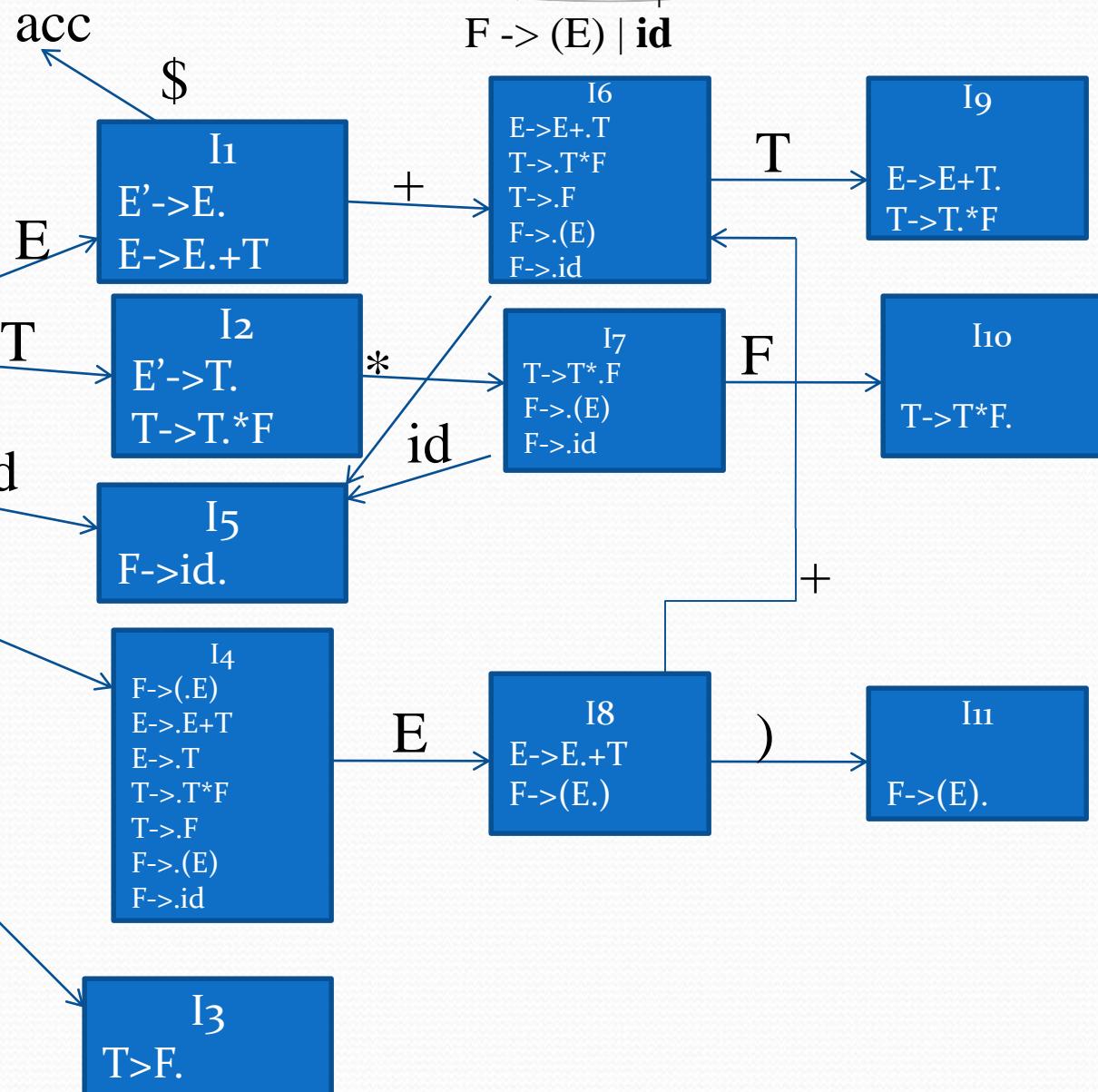
Canonical LR(0) items

```
Void items(G') {  
    C= CLOSURE({[S'->.S]});  
    repeat  
        for (each set of items I in C)  
            for (each grammar symbol X)  
                if (GOTO(I,X) is not empty and not in C)  
                    add GOTO(I,X) to C;  
    until no new set of items are added to C on a round;  
}
```

Example

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

$I_0 = \text{closure}(\{[E' \rightarrow .E]\})$
 $E' \rightarrow .E$
 $E \rightarrow .E + T$
 $E \rightarrow .T$
 $T \rightarrow .T * F$
 $T \rightarrow .F$
 $F \rightarrow .(E)$
 $F \rightarrow .id$

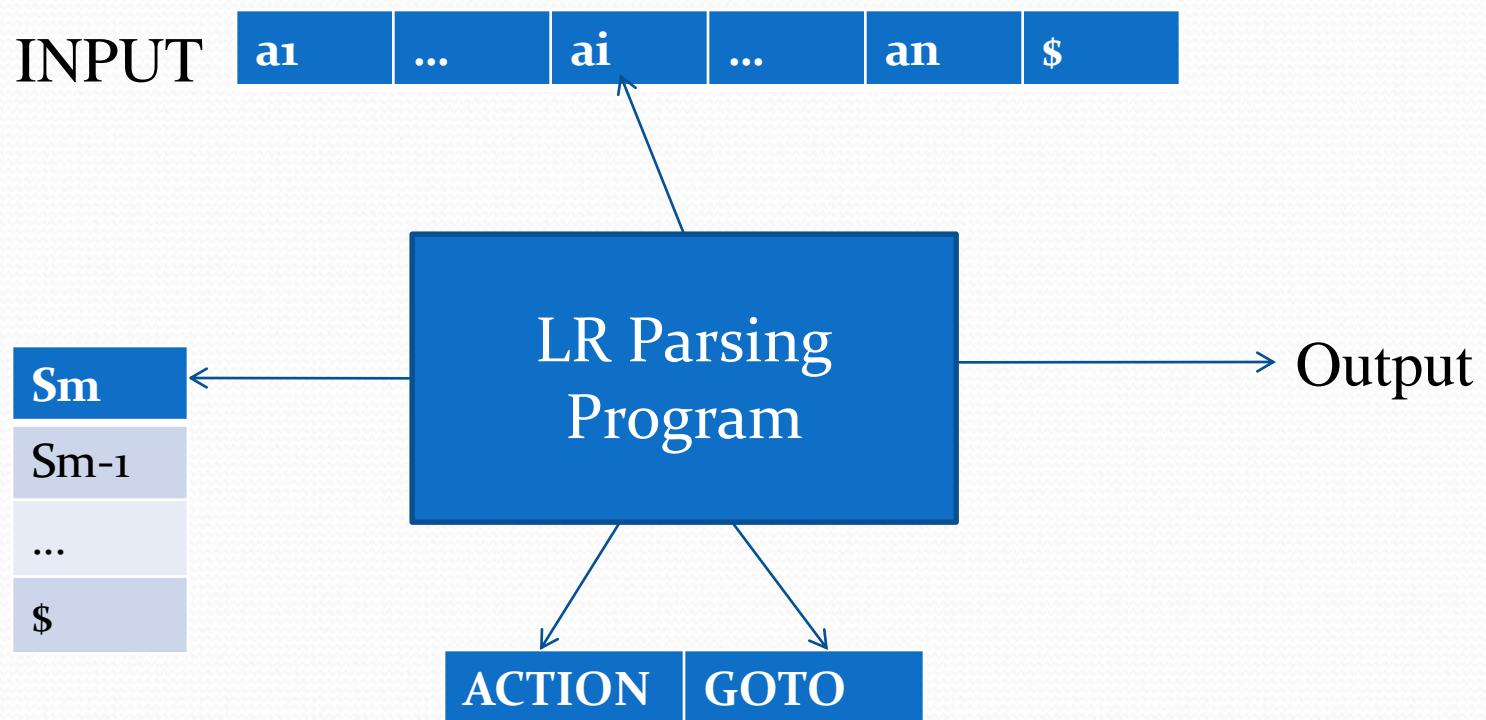


Use of LR(0) automaton

- Example: $\text{id}^* \text{id}$

Line	Stack	Symbols	Input	Action
(1)	o	\$	$\text{id}^* \text{id} \$$	Shift to 5
(2)	o5	\$id	$^* \text{id} \$$	Reduce by F->id
(3)	o3	\$F	$^* \text{id} \$$	Reduce by T->F
(4)	o2	\$T	$^* \text{id} \$$	Shift to 7
(5)	o27	\$T*	id\$	Shift to 5
(6)	o275	\$T*id	\$	Reduce by F->id
(7)	o2710	\$T*T	\$	Reduce by T->T*T
(8)	o2	\$T	\$	Reduce by E->T
(9)	o1	\$E	\$	accept

LR-Parsing model



LR parsing algorithm

```
let a be the first symbol of w$;  
while(1) { /*repeat forever */  
    let s be the state on top of the stack;  
    if (ACTION[s,a] = shift t) {  
        push t onto the stack;  
        let a be the next input symbol;  
    } else if (ACTION[s,a] = reduce A->β) {  
        pop |β| symbols of the stack;  
        let state t now be on top of the stack;  
        push GOTO[t,A] onto the stack;  
        output the production A->β;  
    } else if (ACTION[s,a]=accept) break; /* parsing is done */  
    else call error-recovery routine;  
}
```

Example

STATE	ACTION						GOTO		
	id	+	*	()	\$	E	T	F
0	S ₅			S ₄			1	2	3
1		S ₆				Acc			
2		R ₂	S ₇		R ₂	R ₂			
3		R ₄	R ₇		R ₄	R ₄			
4	S ₅			S ₄			8	2	3
5		R ₆	R ₆		R ₆	R ₆			
6	S ₅			S ₄			9		3
7	S ₅			S ₄					10
8		S ₆			S ₁₁				
9		R ₁	S ₇		R ₁	R ₁			
10		R ₃	R ₃		R ₃	R ₃			
11		R ₅	R ₅		R ₅	R ₅			

- (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) $F \rightarrow id$
- $id * id + id ?$

Line	Stack	Symbol	Input	Action
(1)	o		id*id+id\$	Shift to 5
(2)	o5	id	*id+id\$	Reduce by F->id
(3)	o3	F	*id+id\$	Reduce by T->F
(4)	o2	T	*id+id\$	Shift to 7
(5)	o27	T*	id+id\$	Shift to 5
(6)	o275	T*id	+id\$	Reduce by F->id
(7)	o2710	T*T	+id\$	Reduce by T->T*T
(8)	o2	T	+id\$	Reduce by E->T
(9)	o1	E	+id\$	Shift
(10)	o16	E+	id\$	Shift
(11)	o165	E+id	\$	Reduce by F->id
(12)	o163	E+F	\$	Reduce by T->F
(13)	o169	E+T'	\$	Reduce by E->E+T
(14)	o1	E	\$	accept

Constructing SLR parsing table

- Method
 - Construct $C = \{I_0, I_1, \dots, I_n\}$, the collection of LR(0) items for G'
 - State i is constructed from state I_i :
 - If $[A \rightarrow \alpha.a\beta]$ is in I_i and $\text{Goto}(I_i, a) = I_j$, then set $\text{ACTION}[i, a]$ to “shift j ”
 - If $[A \rightarrow \alpha.]$ is in I_i , then set $\text{ACTION}[i, a]$ to “reduce $A \rightarrow \alpha$ ” for all a in $\text{follow}(A)$
 - If $\{S' \rightarrow .S\}$ is in I_i , then set $\text{ACTION}[i, \$]$ to “Accept”
 - If any conflicts appears then we say that the grammar is not SLR(1).
 - If $\text{GOTO}(I_i, A) = I_j$ then $\text{GOTO}[i, A] = j$
 - All entries not defined by above rules are made “error”
 - The initial state of the parser is the one constructed from the set of items containing $[S' \rightarrow .S]$

Example grammar which is not SLR(1)

SLR(1)

$S \rightarrow L=R \mid R$

$L \rightarrow *R \mid id$

$R \rightarrow L$

I0

$S' \rightarrow .S$

$S \rightarrow .L=R$

$S \rightarrow .R$

$L \rightarrow .*R \mid$

$L \rightarrow .id$

$R \rightarrow .L$

I1

$S' \rightarrow S.$

I2

$S \rightarrow L.=R$

$R \rightarrow L.$

I3

$S \rightarrow R.$

I4

$L \rightarrow *.R$

$R \rightarrow .L$

$L \rightarrow .*R$

$L \rightarrow .id$

I5

$L \rightarrow id.$

I6

$S \rightarrow L.=R$

$R \rightarrow .L$

$L \rightarrow .*R$

$L \rightarrow .id$

I7

$L \rightarrow *R.$

I8

$R \rightarrow L.$

I9

$S \rightarrow L=R.$

2

Action

=

Shift 6
Reduce $R \rightarrow L$

More powerful LR parsers

- Canonical-LR or just LR method
 - Use lookahead symbols for items: LR(1) items
 - Results in a large collection of items
- LALR: lookaheads are introduced in LR(0) items

Canonical LR(1) items

- In LR(1) items each item is in the form: $[A \rightarrow \alpha.\beta, a]$
- An LR(1) item $[A \rightarrow \alpha.\beta, a]$ is valid for a viable prefix γ if there is a derivation $S \stackrel{*}{\Rightarrow}_{rm} \delta A w \Rightarrow_{rm} \delta \alpha \beta w$, where
 - $\Gamma = \delta \alpha$
 - Either a is the first symbol of w , or w is ϵ and a is $\$$
- Example:
 - $S \rightarrow BB$
 - $B \rightarrow aB | b$
$$S \stackrel{*}{\Rightarrow}_{rm} aaBab \Rightarrow_{rm} aaaBab$$

Item $[B \rightarrow a.B, a]$ is valid for $\gamma = aaa$ and $w = ab$

Constructing LR(1) sets of items

```
SetOfItems Closure(I) {
    repeat
        for (each item [A-> $\alpha$ .B $\beta$ ,a] in I)
            for (each production B-> $\gamma$  in G')
                for (each terminal b in First( $\beta$ a))
                    add [B->. $\gamma$ , b] to set I;
}
```

until no more items are added to I;

return I;

}

```
SetOfItems Goto(I,X) {
    initialize J to be the empty set;
    for (each item [A-> $\alpha$ .X $\beta$ ,a] in I)
        add item [A-> $\alpha$ X. $\beta$ ,a] to set J;
    return closure(J);
}
```

```
void items(G') {
    initialize C to Closure({[S'->. $S$ , $]}));
    repeat
        for (each set of items I in C)
            for (each grammar symbol X)
                if (Goto(I,X) is not empty and not in C)
                    add Goto(I,X) to C;
    until no new sets of items are added to C;
}
```

Example

$S' \rightarrow S$

$S \rightarrow CC$

$C \rightarrow cC$

$C \rightarrow d$

Canonical LR(1) parsing table

- Method
 - Construct $C = \{I_0, I_1, \dots, I_n\}$, the collection of LR(1) items for G'
 - State i is constructed from state I_i :
 - If $[A \rightarrow \alpha.a\beta, b]$ is in I_i and $\text{Goto}(I_i, a) = I_j$, then set $\text{ACTION}[i, a]$ to “shift j ”
 - If $[A \rightarrow \alpha., a]$ is in I_i , then set $\text{ACTION}[i, a]$ to “reduce $A \rightarrow \alpha$ ”
 - If $\{S' \rightarrow .S, \$\}$ is in I_i , then set $\text{ACTION}[I, \$]$ to “Accept”
 - If any conflicts appears then we say that the grammar is not LR(1).
 - If $\text{GOTO}(I_i, A) = I_j$ then $\text{GOTO}[i, A] = j$
 - All entries not defined by above rules are made “error”
 - The initial state of the parser is the one constructed from the set of items containing $[S' \rightarrow .S, \$]$

Example

$S' \rightarrow S$

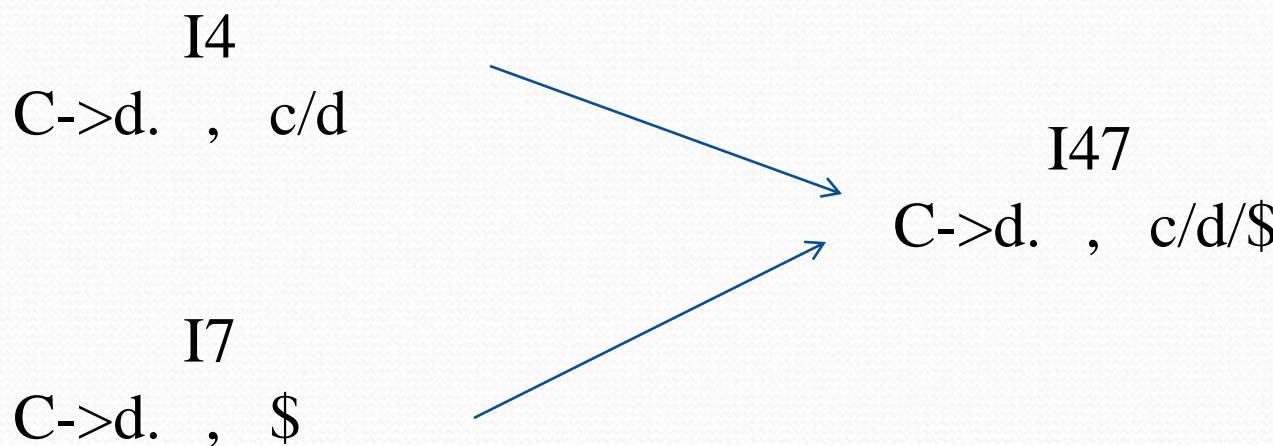
$S \rightarrow CC$

$C \rightarrow cC$

$C \rightarrow d$

LALR Parsing Table

- For the previous example we had:



- State merges can't produce Shift-Reduce conflicts.
Why?
- But it may produce reduce-reduce conflict

Example of RR conflict in state merging

$S' \rightarrow S$

$S \rightarrow aAd \mid bBd \mid aBe \mid bAe$

$A \rightarrow c$

$B \rightarrow c$

An easy but space-consuming LALR table construction

- Method:
 1. Construct $C = \{I_0, I_1, \dots, I_n\}$ the collection of LR(1) items.
 2. For each core among the set of LR(1) items, find all sets having that core, and replace these sets by their union.
 3. Let $C' = \{J_0, J_1, \dots, J_m\}$ be the resulting sets. The parsing actions for state i , is constructed from J_i as before. If there is a conflict grammar is not LALR(1).
 4. If J is the union of one or more sets of LR(1) items, that is $J = I_1 \cup I_2 \dots \cup I_k$ then the cores of $\text{Goto}(I_1, X), \dots, \text{Goto}(I_k, X)$ are the same and is a state like K , then we set $\text{Goto}(J, X) = k$.
- This method is not efficient, a more efficient one is discussed in the book

Compaction of LR parsing table

- Many rows of action tables are identical
 - Store those rows separately and have pointers to them from different states
 - Make lists of (terminal-symbol, action) for each state
 - Implement Goto table by having a link list for each nonterminal in the form (current state, next state)