

Chapter 4

Syntax Analysis

Sr. Nisha C D

Assistant Professor, Dept. of Computer Science

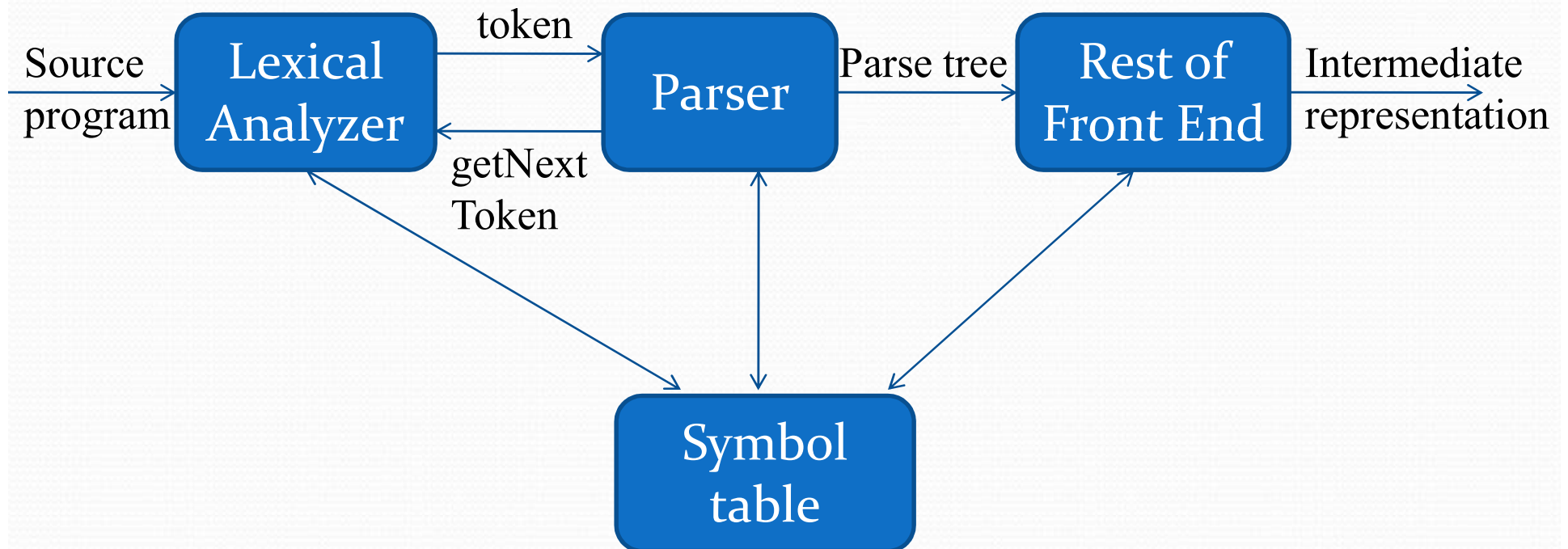
Little Flower College, Guruvayoor



Outline

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

The role of parser



Uses of grammars

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



Error handling

- Common programming errors
 - Lexical errors
 - Syntactic errors
 - Semantic errors
 - Lexical 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-recover strategies

- Panic mode recovery
 - Discard input symbol one at a time until one of designated set of synchronization tokens is found
- Phrase level recovery
 - Replacing a prefix of remaining input by some string that allows the parser to continue
- Error productions
 - Augment the grammar with productions that generate the erroneous constructs
- Global correction
 - Choosing minimal sequence of changes to obtain a globally least-cost correction

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)

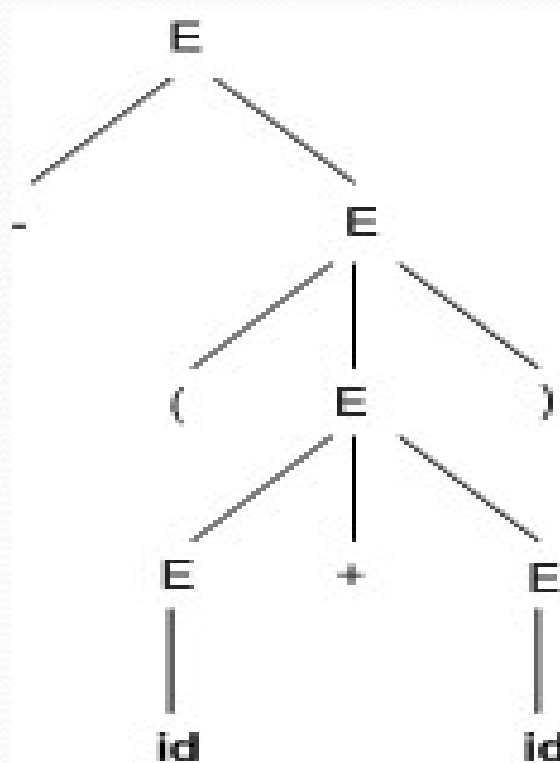
factor \rightarrow **id**

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

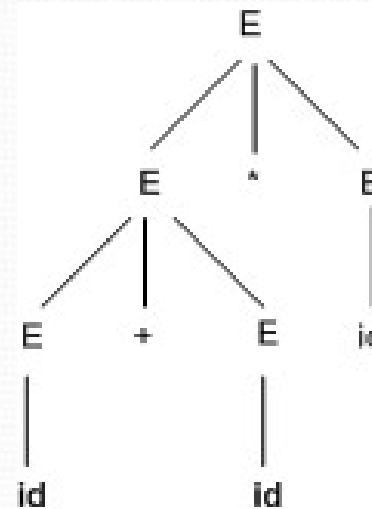
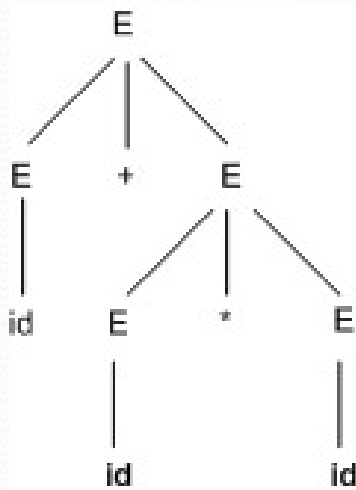
Parse trees

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



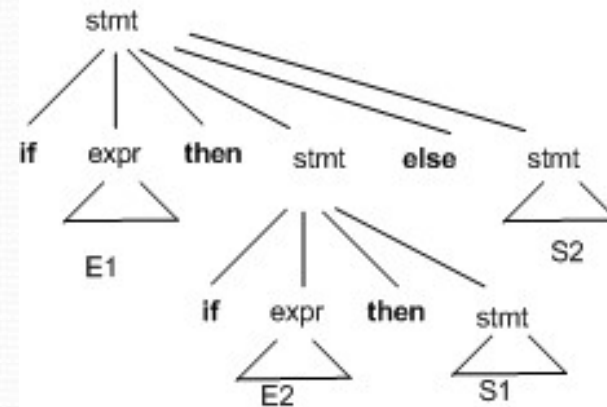
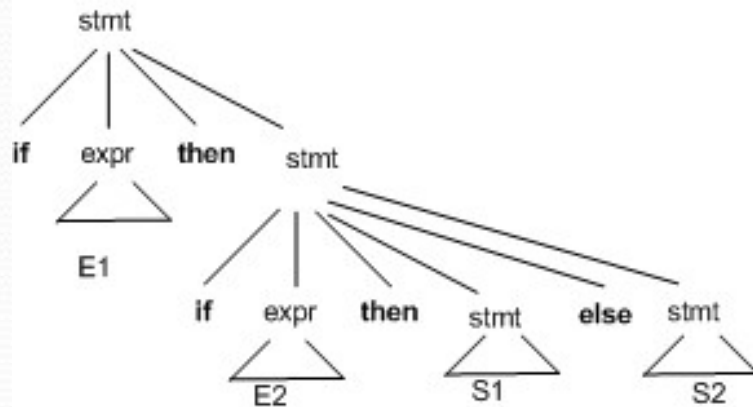
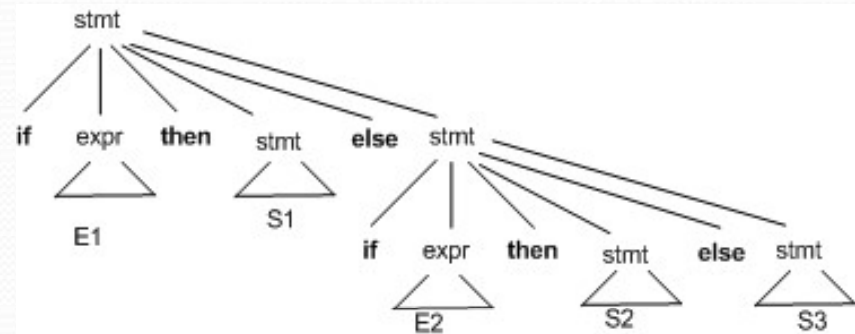
Ambiguity

- 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 ambiguity

stmt \rightarrow If expr then stmt
| If expr then stmt else stmt
| other



Elimination of ambiguity (cont.)

- Idea:
 - A statement appearing between a **then** and an **else** must be matched

```
stmt  →  matched_stmt
      |  open_stmt

matched_stmt → If expr then matched_stmt else matched_stmt
             | other

open_stmt  →  If expr then stmt
             |  If expr then matched_stmt else open_stmt
```

Elimination of left recursion

- A grammar is left recursive if it has a non-terminal A such that there is a derivation $A \xRightarrow{+} A \alpha$
- Top down parsing methods cant handle left-recursive grammars
- A simple rule for direct left recursion elimination:
 - For a rule like:
 - $A \rightarrow A \alpha \mid \beta$
 - We may replace it with
 - $A \rightarrow \beta A'$
 - $A' \rightarrow \alpha A' \mid \varepsilon$

Left recursion elimination (cont.)

- There are cases like following
 - $S \rightarrow Aa \mid b$
 - $A \rightarrow Ac \mid Sd \mid \varepsilon$
- Left recursion elimination algorithm:
 - Arrange the nonterminals in some order A_1, A_2, \dots, A_n .
 - For (each i from 1 to n) {
 - For (each j from 1 to $i-1$) {
 - Replace each production of the form $A_i \rightarrow A_j \gamma$ by the production $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$ where $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ are all current A_j productions
 - }
 - Eliminate left recursion among the A_i -productions
 - }

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.)

- Algorithm

- For each non-terminal A , find the longest prefix α common to two or more of its alternatives. If $\alpha \neq \epsilon$, then replace all of A -productions $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$ by

- $A \rightarrow \alpha A' \mid \gamma$

- $A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

- Example:

- $S \rightarrow I E t S \mid i E t S e S \mid a$

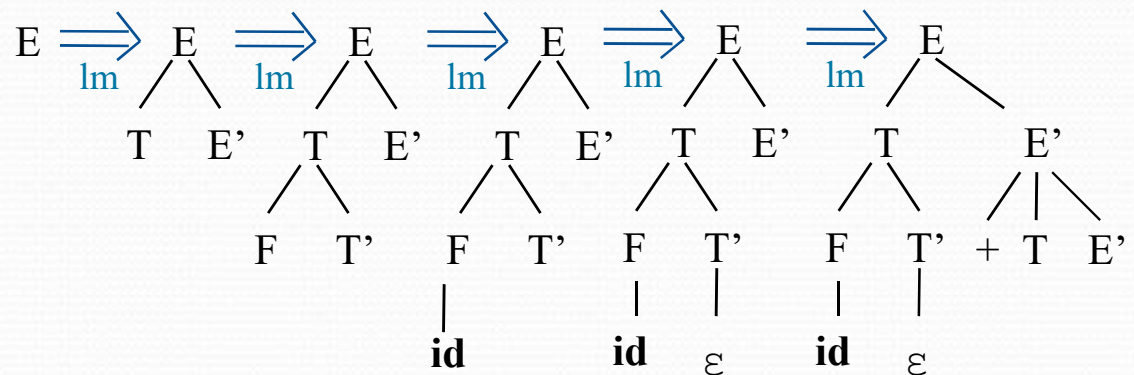
- $E \rightarrow b$

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: $id+id*id$

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



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,  $A \rightarrow X_1 X_2 \dots X_k$   
    for (i=1 to k) {  
        if ( $X_i$  is a nonterminal  
            call procedure  $X_i()$ ;  
        else if ( $X_i$  equals the current input symbol a)  
            advance the input to the next symbol;  
        else /* an error has occurred */  
    }  
}
```



Recursive descent parsing (cont)

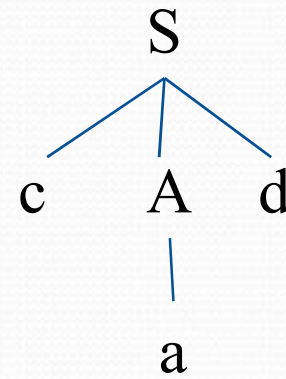
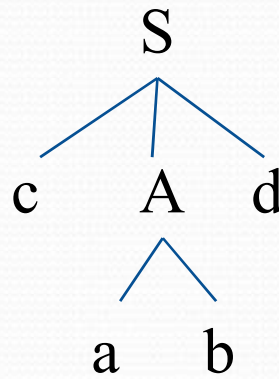
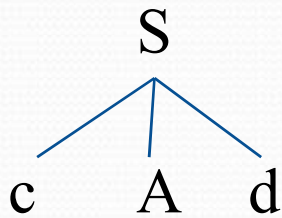
- General recursive descent may require backtracking
- The previous code needs to be modified to allow backtracking
- In general form it can't choose an A -production easily.
- So we need to try all alternatives
- If one failed the input pointer needs to be reset and another alternative should be tried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

Example

$S \rightarrow cAd$

$A \rightarrow ab \mid a$

Input: cad





First and Follow

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

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 $Y_1 \dots Y_{i-1} \xRightarrow{*} \epsilon$. if ϵ is in $\text{First}(Y_j)$ for $j=1, \dots, k$ then add ϵ to $\text{First}(X)$.
 3. If $X \rightarrow \epsilon$ is a production then add ϵ to $\text{First}(X)$
- Example!

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!



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 $A \rightarrow$ 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

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

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

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

Another example

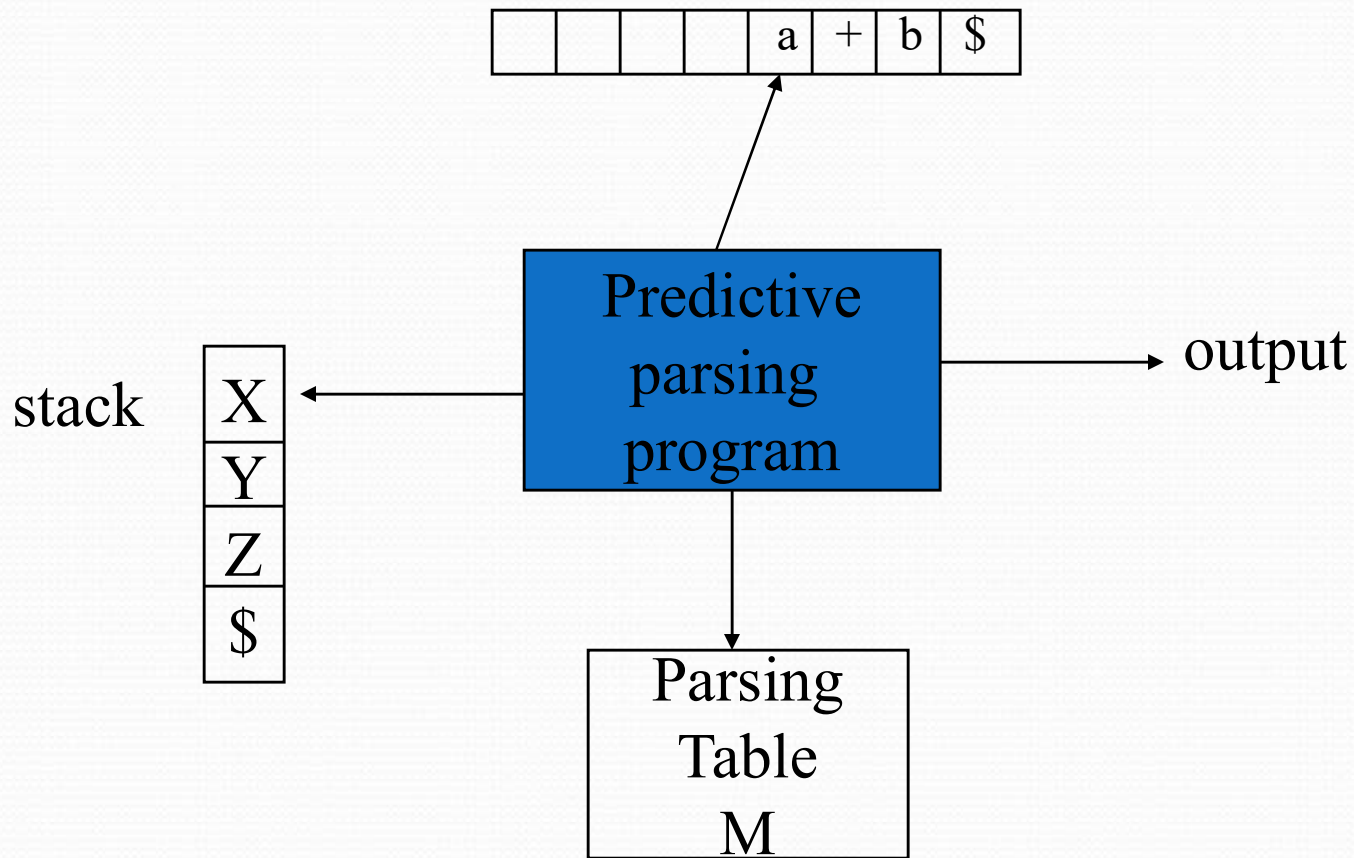
$S \rightarrow iEtSS' \mid a$

$S' \rightarrow eS \mid \epsilon$

$E \rightarrow b$

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

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

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

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



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 → TE'			E → TE'	synch	synch
E'		E' → +TE'			E' → ε	E' → ε
T	T → FT'	synch		T → FT'	synch	synch
T'		T' → ε	T' → *FT'		T' → ε	T' → ε
F	F → id	synch	synch	F → (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: $\text{id}*\text{id}$

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

$\text{id}*\text{id}$

$F * \text{id}$

$T * \text{id}$

$T * F$

F

E

id

F

F

id

$T * F$

F

id

id

F

id

$T * F$

F

id

id



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:
 - $E \Rightarrow T \Rightarrow T * F \Rightarrow T * id \Rightarrow F * id \Rightarrow id * id$

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
id*id	id	F->id
F*id	F	T->F
T*id	id	F->id
T*F	T*F	E->T*F

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	\$

Shift reduce parsing (cont.)

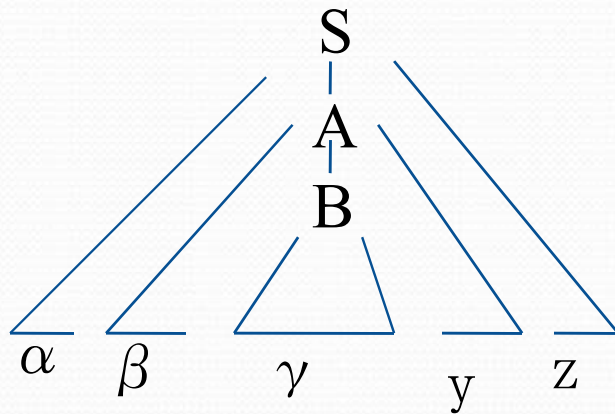
- Basic operations:

- Shift
- Reduce
- Accept
- Error

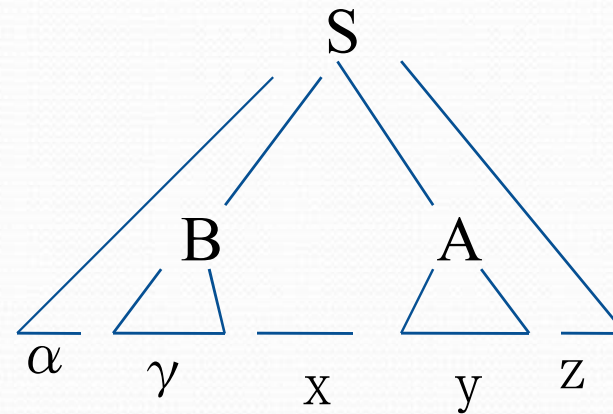
- Example: $id*id$

Stack	Input	Action
\$	$id*id\$$	shift
$\$id$	$*id\$$	reduce by $F \rightarrow id$
$\$F$	$*id\$$	reduce by $T \rightarrow F$
$\$T$	$*id\$$	shift
$\$T*$	$id\$$	shift
$\$T*id$	$\$$	reduce by $F \rightarrow id$
$\$T*F$	$\$$	reduce by $T \rightarrow T*F$
$\$T$	$\$$	reduce by $E \rightarrow T$
$\$E$	$\$$	accept

Handle will appear on top of the stack



Stack	Input
\$ α β γ	yz\$
\$ α β B	yz\$
\$ α β By	z\$



Stack	Input
\$ α γ	xyz\$
\$ α Bxy	z\$

Conflicts during shift 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) ...\$



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(o) 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 $\text{closure}(I)$.

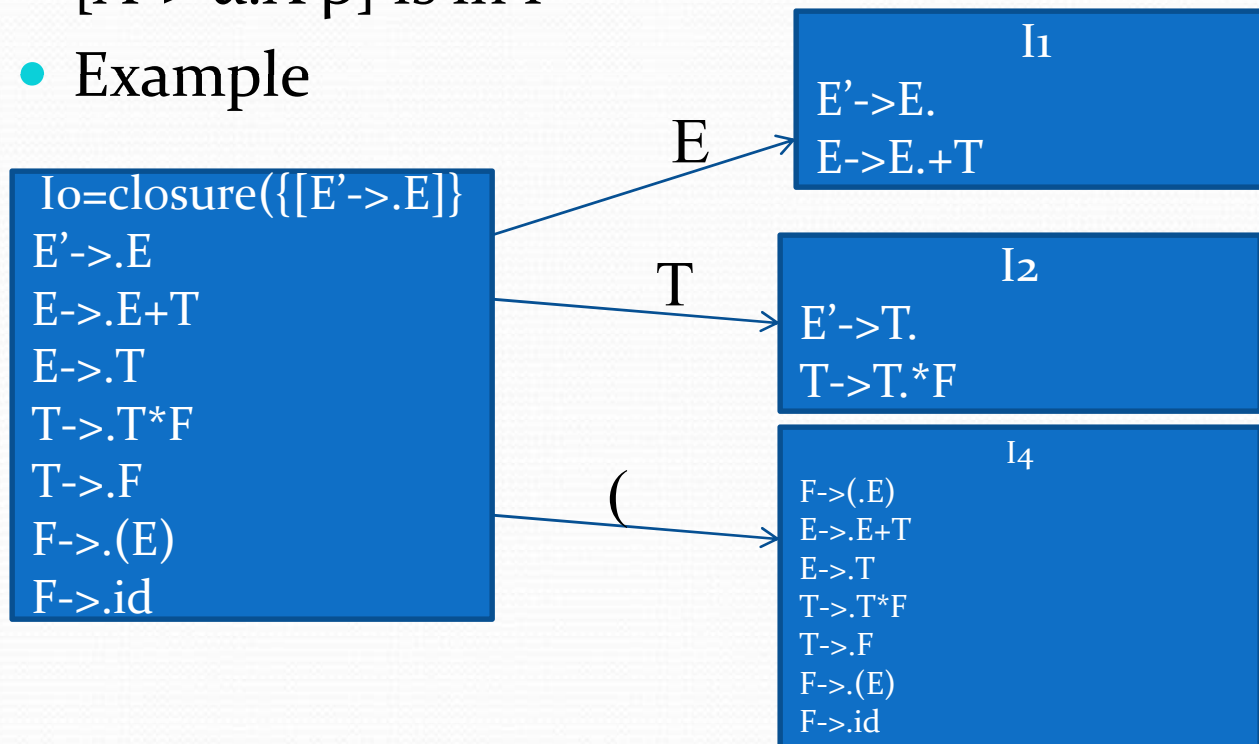
- Example:

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

```
Io=closure({[E'-.E]})
E'-.E
E-.-E+T
E-.-T
T-.-T*F
T-.-F
F-.-(E)
F-.-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 $A \rightarrow \alpha.B\beta$ in J)

 for (each production $B \rightarrow \gamma$ of G)

 if ($B \rightarrow \cdot \gamma$ is not in J)

 add $B \rightarrow \cdot \gamma$ 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- \rightarrow α .X β is in I)

 add CLOSURE(A- \rightarrow α X. β) to J;

 return J;

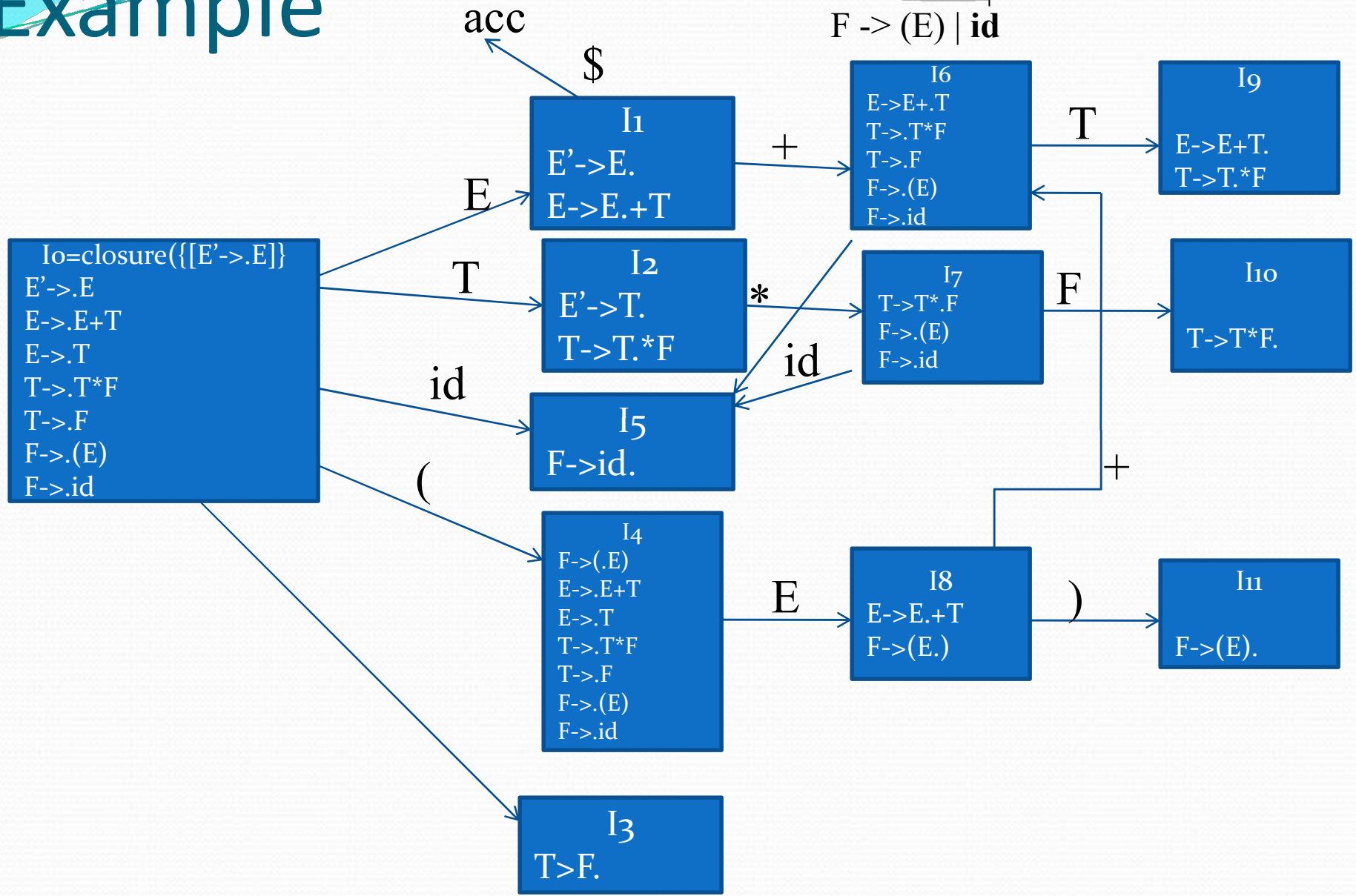
}

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

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

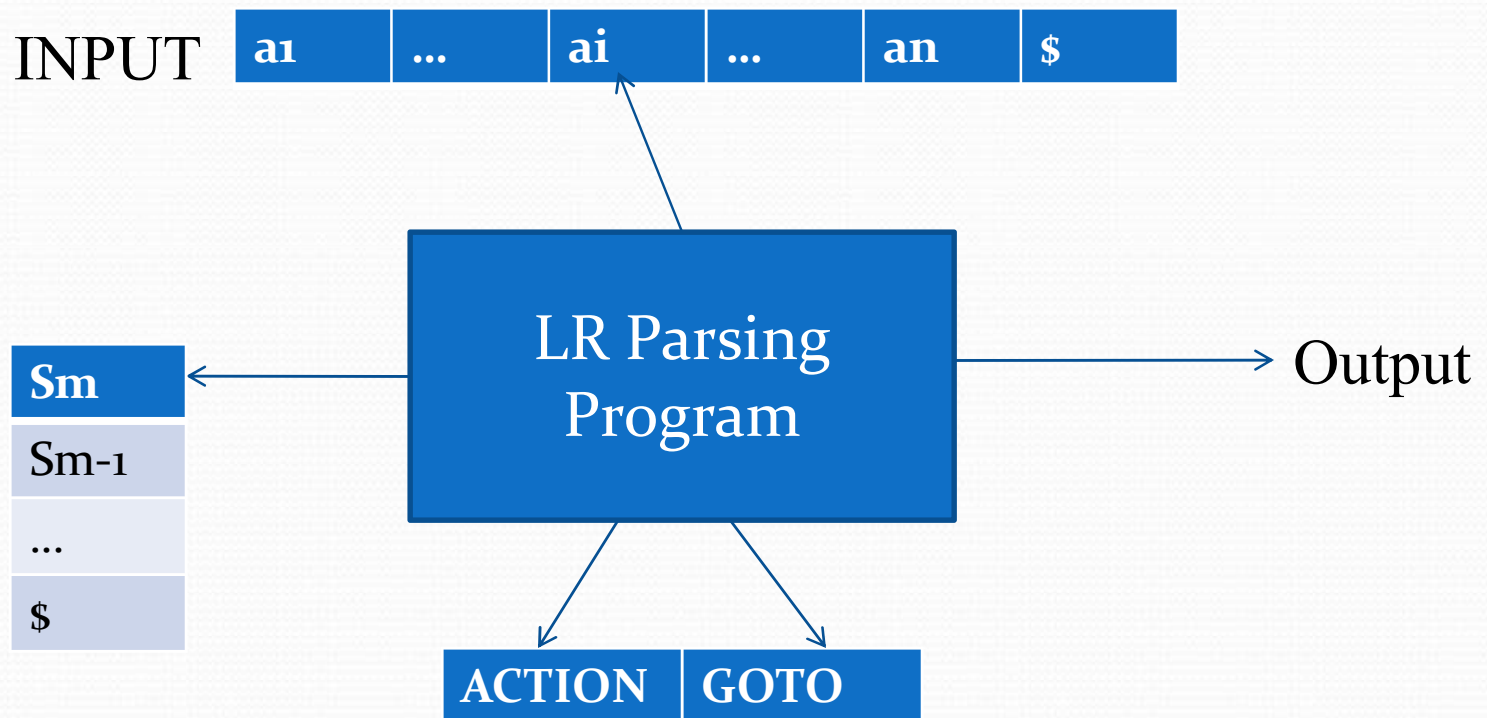


Use of LR(0) automaton

- Example: $id*id$

Line	Stack	Symbols	Input	Action
(1)	0	\$	id*id\$	Shift to 5
(2)	05	\$id	*id\$	Reduce by $F \rightarrow id$
(3)	03	\$F	*id\$	Reduce by $T \rightarrow F$
(4)	02	\$T	*id\$	Shift to 7
(5)	027	\$T*	id\$	Shift to 5
(6)	0275	\$T*id	\$	Reduce by $F \rightarrow id$
(7)	02710	\$T*F	\$	Reduce by $T \rightarrow T*F$
(8)	02	\$T	\$	Reduce by $E \rightarrow T$
(9)	01	\$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-> $\beta$ ) {
        pop  $|\beta|$  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-> $\beta$ ;
    } else if (ACTION[s,a]=accept) break; /* parsing is done */
    else call error-recovery routine;
}
```

Example

- (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?

STATE	ACTION						GOTO		
	id	+	*	()	\$	E	T	F
0	S5			S4			1	2	3
1		S6				Acc			
2		R2	S7		R2	R2			
3		R4	R7		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			

Line	Stack	Symbols	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*F	+id\$	Reduce by T->T*F
(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 $\text{SLR}(1)$.
 - 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)

$S \rightarrow L=R \mid R$

$L \rightarrow *R \mid id$

$R \rightarrow L$

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

$S \rightarrow L=R.$

Action

=

2

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: $[A \rightarrow \alpha \cdot \beta, a]$
- An LR(1) item $[A \rightarrow \alpha \cdot \beta, a]$ is valid for a viable prefix γ if there is a derivation $S \xRightarrow{*} \delta A w \xRightarrow{rm} \delta \alpha \beta w$, where
 - $\Gamma = \delta \alpha$
 - Either a is the first symbol of w , or w is ϵ and a is $\$$
- Example:

- $S \rightarrow BB$

$$S \xRightarrow{*} aaBab \xRightarrow{rm} aaaBab$$

- $B \rightarrow aB \mid b$

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

Constructing LR(1) sets of items

```
SetOfItems Closure(I) {
  repeat
    for (each item  $[A \rightarrow \alpha.B\beta, a]$  in I)
      for (each production  $B \rightarrow \gamma$  in  $G'$ )
        for (each terminal  $b$  in  $\text{First}(\beta a)$ )
          add  $[B \rightarrow \cdot\gamma, b]$  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  $[A \rightarrow \alpha.X\beta, a]$  in I)
    add item  $[A \rightarrow \alpha X \cdot \beta, a]$  to set J;
  return closure(J);
}
```

```
void items( $G'$ ){
  initialize C to  $\text{Closure}(\{[S' \rightarrow \cdot S, \$]\})$ ;
  repeat
    for (each set of items I in C)
      for (each grammar symbol X)
        if ( $\text{Goto}(I, X)$  is not empty and not in C)
          add  $\text{Goto}(I, X)$  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

$S' \rightarrow S$

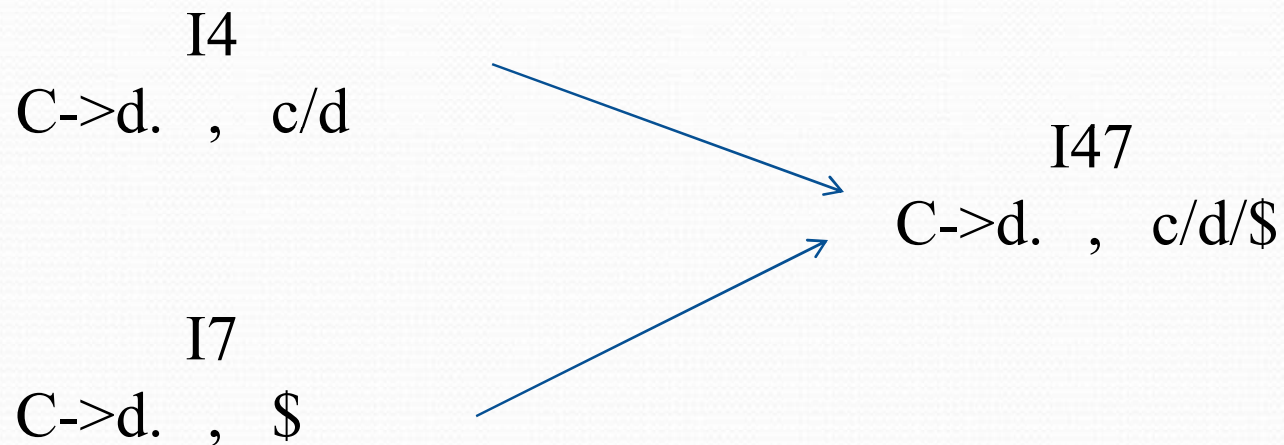
$S \rightarrow CC$

$C \rightarrow cC$

$C \rightarrow d$

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 $\text{Goto}(I_1, X)$, ..., $\text{Goto}(I_k, X)$ 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)

Using ambiguous grammars

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

$E \rightarrow id$

I0: $E' \rightarrow .E$

$E \rightarrow .E + E$

$E \rightarrow .E * E$

$E \rightarrow .(E)$

$E \rightarrow .id$

I1: $E' \rightarrow E.$

$E \rightarrow E. + E$

$E \rightarrow E. * E$

I4: $E \rightarrow E + .E$

$E \rightarrow E. + E$

$E \rightarrow E. * E$

$E \rightarrow E. (E)$

$E \rightarrow E. id$

I2: $E \rightarrow (.E)$

$E \rightarrow .E + E$

$E \rightarrow .E * E$

$E \rightarrow .(E)$

$E \rightarrow .id$

I5: $E \rightarrow E * .E$

$E \rightarrow (.E)$

$E \rightarrow .E + E$

$E \rightarrow .E * E$

$E \rightarrow .(E)$

$E \rightarrow .id$

I6: $E \rightarrow (E.)$

$E \rightarrow E. + E$

$E \rightarrow E. * E$

I8: $E \rightarrow E * E.$

$E \rightarrow E. + E$

$E \rightarrow E. * E$

I7: $E \rightarrow E + E.$

$E \rightarrow E. + E$

$E \rightarrow E. * E$

I9: $E \rightarrow (E).$

STATE	ACTON						GO TO
	id	+	*	()	\$	E
0	S ₃			S ₂			1
1		S ₄	S ₅			Acc	
2	S ₃		S ₂				6
3		R ₄	R ₄		R ₄	R ₄	
4	S ₃			S ₂			7
5	S ₃			S ₂			8
6		S ₄	S ₅				
7		R ₁	S ₅		R ₁	R ₁	
8		R ₂	R ₂		R ₂	R ₂	
9		R ₃	R ₃		R ₃	R ₃	



Readings

- Chapter 4 of the book