

Principles of Programming Languages

<http://www.di.unipi.it/~andrea/Didattica/PLP-15/>

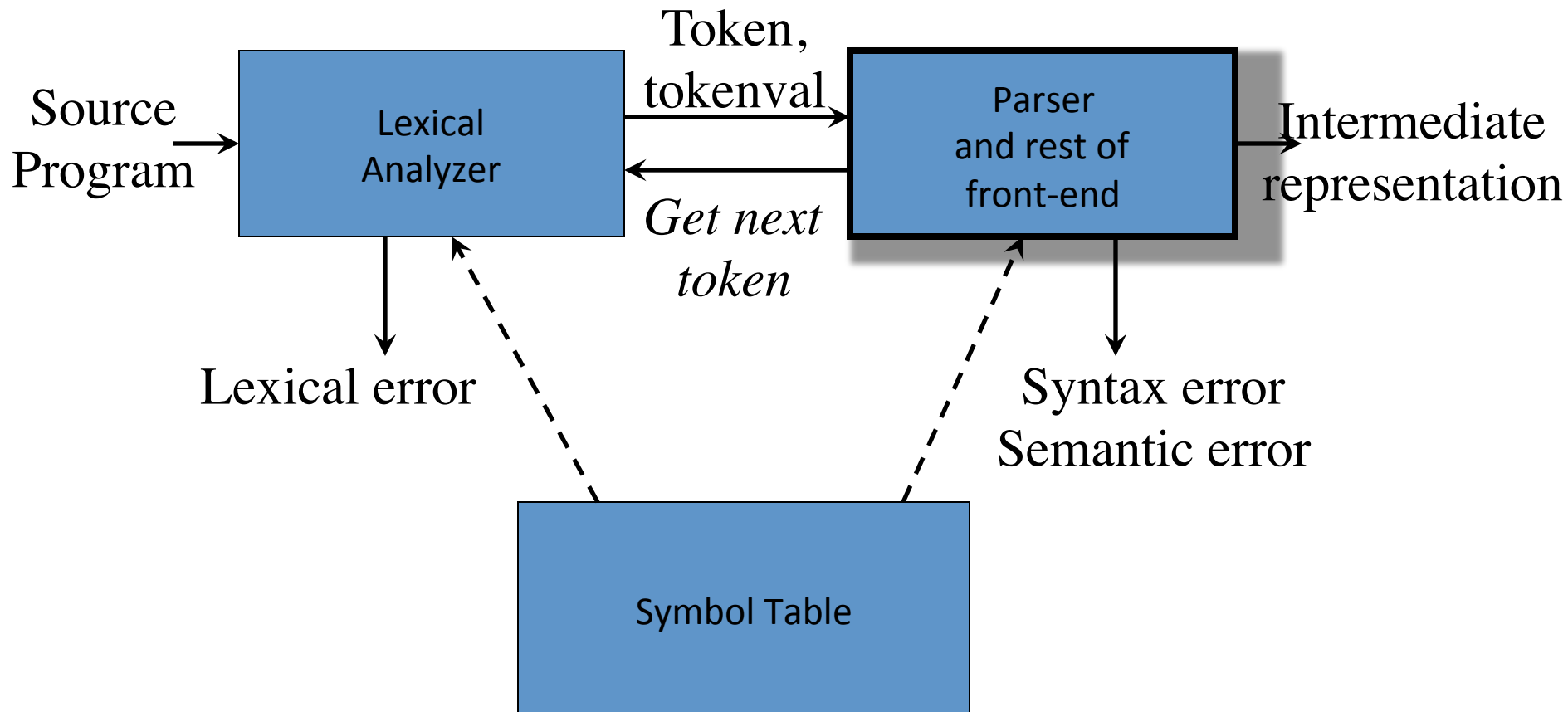
Prof. Andrea Corradini

Department of Computer Science, Pisa

Lesson 8

- Top-down parsing
- Introduction to bottom-up parsing

Position of a Parser in the Compiler Model



The syntax of programming languages

- The syntax of a programming language is typically defined by two grammars
 - Lexical grammar
 - Regular, often presented as regular expressions
 - Terminal symbols are characters
 - Defines tokens
 - Syntax grammar
 - Context-free, often presented in Backus-Naur form
 - Terminal symbols are tokens
 - Defines constructs of the language, not expressible with REs
 - Note: there are non-context free syntact constructs
 - Variables are declared before use $\rightarrow \{wcw \mid w \in (a \mid b)^*\}$
 - Number of actual/formal parameters $\rightarrow \{a^n b^m c^n d^m \mid n > 0, m > 0\}$

Towards parsing

- A parser implements a Context-Free grammar as a recognizer of strings
 - It checks that the input string (of tokens) is generated by the syntax grammar
 - Possibly generates the parse tree
 - Reports syntax errors accurately
 - *Invokes semantic actions*
 - *For static semantics checking, e.g. type checking of expressions, functions, etc.*
 - *For syntax-directed translation of the source code to an intermediate representation*

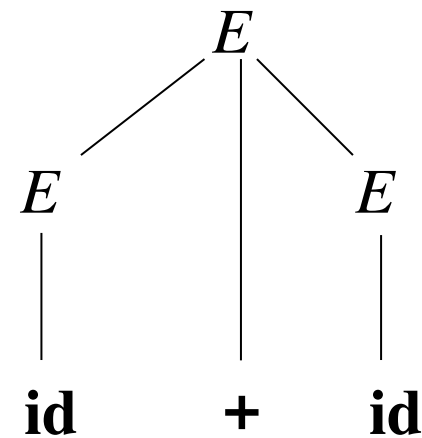
Parse trees and derivations

- A parse tree may correspond to several derivations
- A parse tree has a unique *rightmost* (*leftmost*) derivation

$$P = E \rightarrow E + E \mid \text{id}$$

$$E \Rightarrow_{rm} E + E \Rightarrow_{rm} E + \text{id} \Rightarrow_{rm} \text{id} + \text{id}$$

$$E \Rightarrow_{lm} E + E \Rightarrow_{lm} \text{id} + E \Rightarrow_{lm} \text{id} + \text{id}$$



Parsing algorithms

- *Universal* (any C-F grammar)
 - Cocke-Younger-Kasimi, Earley
 - Based on dynamic programming, $O(n^3)$
- *Top-down* (C-F grammar with restrictions)
 - Recursive descent (predictive parsing)
 - LL (Left-to-right, Leftmost derivation) methods
 - Linear on certain grammars; easier to do manually
- *Bottom-up* (C-F grammar with restrictions)
 - Operator precedence parsing
 - LR (Left-to-right, Rightmost derivation) methods
 - SLR, canonical LR, LALR
 - Linear on certain grammars; typically generated by tools

Top-Down Parsing

- LL methods (Left-to-right, Leftmost derivation) and recursive-descent parsing

Grammar:

$E \rightarrow T + T$

$T \rightarrow (E)$

$T \rightarrow - E$

$T \rightarrow \text{id}$

String

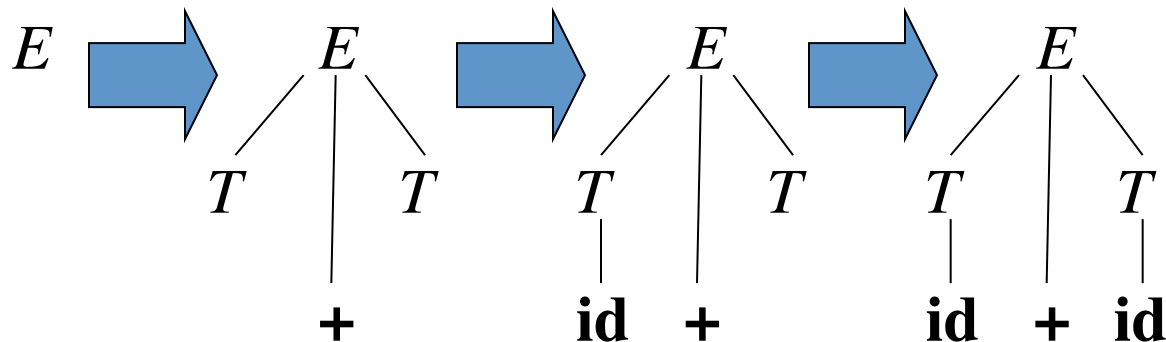
id + id

Leftmost derivation:

$E \Rightarrow_{lm} T + T$

$\Rightarrow_{lm} \text{id} + T$

$\Rightarrow_{lm} \text{id} + \text{id}$



LL(k) parsing

- Top-down parsing is efficient if the grammar satisfies certain conditions
- Whenever we have to expand a non-terminal, the next k tokens should determine the production to use (*lookahead*)
- In this case the grammar is LL(k)
- Most constructs are LL(1), and we will focus on this class of grammars

Left Recursion

- A grammar is *left-recursive* if there is a non-terminal A such that $A \Rightarrow^+ A\eta$ for some string η
 - Example of immediate left-recursion:
 $A \rightarrow A\alpha \mid A\beta \mid \gamma \mid \delta$
 - Left recursion can be indirect
- If the grammar is left-recursive, it cannot be $LL(k)$: a top-down parser loops forever on certain inputs
- Immediate left recursion elimination:

$$A \rightarrow \gamma A_R \mid \delta A_R \qquad A_R \rightarrow \alpha A_R \mid \beta A_R \mid \varepsilon$$

A General Left Recursion Elimination Method

- *Input: Grammar G with no cycles or ε -productions*
- *Arrange the nonterminals in some order A_1, A_2, \dots, A_n*

for $i = 1, \dots, n$ **do**

for $j = 1, \dots, i-1$ **do**

 replace each

$$A_i \rightarrow A_j \gamma$$

 with

$$A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$$

 where

$$A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$$

enddo

 eliminate the *immediate left recursion* in A_i

enddo

Example of left-recursion elimination

$$\left. \begin{array}{l} A \rightarrow B C | a \\ B \rightarrow C A | A b \\ C \rightarrow A B | C C | a \end{array} \right\} \text{Choose arrangement: } A, B, C$$

$i = 1$: nothing to do

$i = 2, j = 1$: $B \rightarrow C A | \underline{A} b$

$\Rightarrow B \rightarrow C A | \underline{B C} b | \underline{a} b$

$\Rightarrow_{(\text{imm})} B \rightarrow C A B_R | a b B_R$

$B_R \rightarrow C b B_R | \varepsilon$

$i = 3, j = 1$: $C \rightarrow \underline{A} B | C C | a$

$\Rightarrow C \rightarrow \underline{B C} B | \underline{a} B | C C | a$

$i = 3, j = 2$: $C \rightarrow \underline{B} C B | a B | C C | a$

$\Rightarrow C \rightarrow \underline{C A B_R} C B | \underline{a b B_R} C B | a B | C C | a$

$\Rightarrow_{(\text{imm})} C \rightarrow a b B_R C B C_R | a B C_R | a C_R$

$C_R \rightarrow A B_R C B C_R | C C_R | \varepsilon$

Example of left-recursion elimination: Grammars for expressions

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned}$$

Grammar after left recursion elimination

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

Left Factoring

- If a nonterminal has two or more productions whose right-hand sides start with the same symbol, the grammar is not LL(1)

- Example:

– $stmt ::= \mathbf{if\ expr\ then\ stmt\ else\ stmt}$
 | $\mathbf{if\ expr\ then\ stmt}$

- Solution: replace productions

$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$

with

$A \rightarrow \alpha A_R \mid \gamma$

$A_R \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

- Example:

– $stmt ::= \mathbf{if\ expr\ then\ stmt\ stmt'}$

– $stmt' ::= \mathbf{else\ stmt} \mid \epsilon$

Predictive Parsing

- Eliminate left recursion from grammar
- Left factor the grammar
- Compute FIRST and FOLLOW, and check if the grammar is LL(1)
- FIRST and FOLLOW are used in the parsing algorithm
- Two variants:
 - Recursive (recursive-descent parsing)
 - Non-recursive (table-driven parsing)

FIRST (Revisited)

- $\text{FIRST}(\alpha) = \{ \text{the set of terminals that begin all strings derived from } \alpha \}$
- $\text{FIRST}(a) = \{a\}$ if $a \in T$
- $\text{FIRST}(\varepsilon) = \{\varepsilon\}$
- $\text{FIRST}(A) = \bigcup_{A \rightarrow \alpha} \text{FIRST}(\alpha)$ for $A \rightarrow \alpha \in P$
- $\text{FIRST}(X_1 X_2 \dots X_k) =$
 - if for all $j = 1, \dots, i-1 : \varepsilon \in \text{FIRST}(X_j)$ then**
add non- ε in $\text{FIRST}(X_i)$ to $\text{FIRST}(X_1 X_2 \dots X_k)$
 - if for all $j = 1, \dots, k : \varepsilon \in \text{FIRST}(X_j)$ then**
add ε to $\text{FIRST}(X_1 X_2 \dots X_k)$

FOLLOW

- $\text{FOLLOW}(A) = \{ \text{the set of terminals that can immediately follow nonterminal } A \}$
- $\text{FOLLOW}(A) =$
 - for** all $(B \rightarrow \alpha A \beta) \in P$ **do**
 - add $\text{FIRST}(\beta) \setminus \{\epsilon\}$ to $\text{FOLLOW}(A)$
 - for** all $(B \rightarrow \alpha A \beta) \in P$ and $\epsilon \in \text{FIRST}(\beta)$ **do**
 - add $\text{FOLLOW}(B)$ to $\text{FOLLOW}(A)$
 - for** all $(B \rightarrow \alpha A) \in P$ **do**
 - add $\text{FOLLOW}(B)$ to $\text{FOLLOW}(A)$
 - if** A is the start symbol S **then**
 - add $\$$ to $\text{FOLLOW}(A)$

LL(1) Grammar

- A grammar G is LL(1) if it is not left recursive and for each collection of productions

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

for nonterminal A the following holds:

1. $\text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \emptyset$ for all $i \neq j$
2. if $\alpha_j \Rightarrow^* \varepsilon$ then
 - 2.a. $\alpha_i \not\Rightarrow^* \varepsilon$ for all $i \neq j$
 - 2.b. $\text{FIRST}(\alpha_j) \cap \text{FOLLOW}(A) = \emptyset$ for all $i \neq j$

Non-LL(1) Examples

<i>Grammar</i>	<i>Not LL(1) because:</i>
$S \rightarrow S a \mid a$	Left recursive
$S \rightarrow a S \mid a$	$\text{FIRST}(a S) \cap \text{FIRST}(a) \neq \emptyset$
$S \rightarrow a R \mid \varepsilon$ $R \rightarrow S \mid \varepsilon$	For R : $S \Rightarrow^* \varepsilon$ and $\varepsilon \Rightarrow^* \varepsilon$
$S \rightarrow a R a$ $R \rightarrow S \mid \varepsilon$	For R : $\text{FIRST}(S) \cap \text{FOLLOW}(R) \neq \emptyset$

Recursive-Descent Parsing

- Grammar must be LL(1)
- Every nonterminal has one (recursive) procedure responsible for parsing the nonterminal's syntactic category of input tokens
- When a nonterminal has multiple productions, each production is implemented in a branch of a selection statement based on input look-ahead information

Using FIRST and FOLLOW in a Recursive-Descent Parser

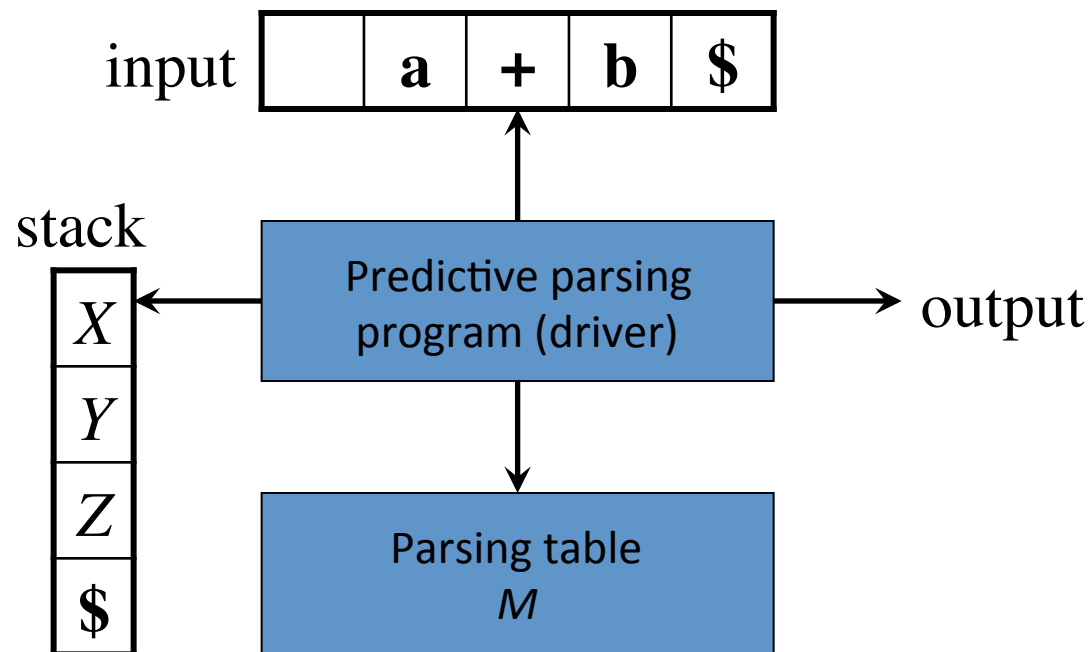
$expr \rightarrow term\ rest$
 $rest \rightarrow +\ term\ rest$
 $| -\ term\ rest$
 $| \epsilon$
 $term \rightarrow id$

```
procedure rest();  
begin  
  if lookahead in FIRST(+ term rest) then  
    match( '+' ); term(); rest()  
  else if lookahead in FIRST(- term rest) then  
    match( '-' ); term(); rest()  
  else if lookahead in FOLLOW(rest) then  
    return  
  else error()  
end;
```

where $FIRST(+\ term\ rest) = \{ + \}$
 $FIRST(-\ term\ rest) = \{ - \}$
 $FOLLOW(rest) = \{ \$ \}$

Non-Recursive Predictive Parsing: Table-Driven Parsing

- Given an LL(1) grammar $G = (N, T, P, S)$ construct a table M and use a *driver program* with a *stack*
- The stack replaces the runtime stack of the recursive algorithm. It will contain symbols of the grammar.



Constructing an LL(1) Predictive Parsing Table

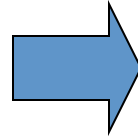
- Table M has one entry $M[A, a]$ for each $A \in N$ and $a \in T$
- Entry $M[A, a]$ contains the production to apply when A has to be reduced and a is the lookahead

```
for each production  $A \rightarrow \alpha$  do
  for each  $a \in \text{FIRST}(\alpha)$  do
    add production  $A \rightarrow \alpha$  to  $M[A, a]$ 
  enddo
  if  $\epsilon \in \text{FIRST}(\alpha)$  then
    for each  $b \in \text{FOLLOW}(A)$  do
      add  $A \rightarrow \alpha$  to  $M[A, b]$ 
    enddo
  endif
enddo
```

- Mark each undefined entry in M error
- **Note:** The grammar is LL(1) iff $M[A, a]$ contains at most one production for each $A \in N$ and $a \in T$

Example Table

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



$A \rightarrow \alpha$	FIRST(α)	FOLLOW(A)
$E \rightarrow T E_R$	(id	\$)
$E_R \rightarrow + T E_R$	+	\$)
$E_R \rightarrow \epsilon$	ϵ	
$T \rightarrow F T_R$	(id	+ \$)
$T_R \rightarrow * F T_R$	*	+ \$)
$T_R \rightarrow \epsilon$	ϵ	
$F \rightarrow (E)$	(* + \$)
$F \rightarrow \mathbf{id}$	id	* + \$)

	id	+	*	()	\$
E	$E \rightarrow T E_R$			$E \rightarrow T E_R$		
E_R		$E_R \rightarrow + T E_R$			$E_R \rightarrow \epsilon$	$E_R \rightarrow \epsilon$
T	$T \rightarrow F T_R$			$T \rightarrow F T_R$		
T_R		$T_R \rightarrow \epsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \epsilon$	$T_R \rightarrow \epsilon$
F	$F \rightarrow \mathbf{id}$			$F \rightarrow (E)$		

Predictive Parsing Program (Driver)

```
push($)  
push(S)  
a := lookahead  
repeat  
     $X := \text{pop}()$   
    if  $X$  is a terminal or  $X = \$$  then  
        match(X) // moves to next token and a := lookahead  
    else if  $M[X,a] = X \rightarrow Y_1 Y_2 \dots Y_k$  then  
        push( $Y_k, Y_{k-1}, \dots, Y_2, Y_1$ ) // such that  $Y_1$  is on top  
        ... invoke actions and/or produce IR output ...  
    else    error()  
    endif  
until  $X = \$$ 
```


Example Table-Driven Parsing

Stack	Input	Production applied
$\$E$	<u>id</u> +id*id\$	$E \rightarrow T E_R$
$\$E_R T$	<u>id</u> +id*id\$	$T \rightarrow F T_R$
$\$E_R T_R F$	<u>id</u> +id*id\$	$F \rightarrow \text{id}$
$\$E_R T_R \text{id}$	<u>id</u> +id*id\$	
$\$E_R T_R$	<u>+</u> id*id\$	$T_R \rightarrow \epsilon$
$\$E_R$	<u>+</u> id*id\$	$E_R \rightarrow + T E_R$
$\$E_R T +$	<u>+</u> id*id\$	
$\$E_R T$	<u>id</u> *id\$	$T \rightarrow F T_R$
$\$E_R T_R F$	<u>id</u> *id\$	$F \rightarrow \text{id}$

Stack	Input	Prod. applied
$\$E_R T_R \text{id}$	<u>id</u> *id\$	
$\$E_R T_R$	<u>*</u> id\$	$T_R \rightarrow * F T_R$
$\$E_R T_R F *$	<u>*</u> id\$	
$\$E_R T_R F$	<u>id</u> \$	$F \rightarrow \text{id}$
$\$E_R T_R \text{id}$	<u>id</u> \$	
$\$E_R T_R$	<u>\$</u>	$T_R \rightarrow \epsilon$
$\$E_R$	<u>\$</u>	$E_R \rightarrow \epsilon$
$\$$	<u>\$</u>	

	id	+	*	()	\$
E	$E \rightarrow T E_R$			$E \rightarrow T E_R$		
E_R		$E_R \rightarrow + T E_R$			$E_R \rightarrow \epsilon$	$E_R \rightarrow \epsilon$
T	$T \rightarrow F T_R$			$T \rightarrow F T_R$		
T_R		$T_R \rightarrow \epsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \epsilon$	$T_R \rightarrow \epsilon$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

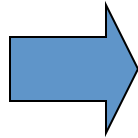
LL(1) Grammars are Unambiguous

Ambiguous grammar

$$S \rightarrow i E t S S_R \mid a$$

$$S_R \rightarrow e S \mid \epsilon$$

$$E \rightarrow b$$



$A \rightarrow \alpha$	FIRST(α)	FOLLOW(A)
$S \rightarrow i E t S S_R$	i	e \$
$S \rightarrow a$	a	
$S_R \rightarrow e S$	e	e \$
$S_R \rightarrow \epsilon$	ϵ	
$E \rightarrow b$	b	t

Error: duplicate table entry

	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow i E t S S_R$		
S_R			$S_R \rightarrow \epsilon$ $S_R \rightarrow e S$			$S_R \rightarrow \epsilon$
E		$E \rightarrow b$				

Error Handling

- A good compiler should assist in identifying and locating errors
 - *Lexical errors*: compiler can easily recover and continue (e.g. misspelled identifiers)
 - *Syntax errors*: can almost always recover (e.g. missing ';' or '{', misplaced **case**)
 - *Static semantic errors*: can sometimes recover (e.g. type mismatches, variable used before declaration)
 - *Dynamic semantic errors*: hard or impossible to detect at compile time, runtime checks are required (e.g. null pointer, division by zero, invalid array access)
 - *Logical errors*: hard or impossible to detect (e.g. `if (b = true) ...`)

Error Recovery Strategies

- *Panic mode*
 - Discard input until a token in a set of designated “synchronizing tokens” is found (e.g. “}”, “;”)
- *Phrase-level recovery*
 - Perform local correction on the input to repair the error
- *Error productions*
 - Augment grammar with productions for erroneous constructs
- *Global correction*
 - Choose a minimal sequence of changes to obtain a global least-cost correction

Panic Mode Recovery

Add synchronizing actions to undefined entries based on FOLLOW

$\text{FOLLOW}(E) = \{) \$ \}$
 $\text{FOLLOW}(E_R) = \{) \$ \}$
 $\text{FOLLOW}(T) = \{ +) \$ \}$
 $\text{FOLLOW}(T_R) = \{ +) \$ \}$
 $\text{FOLLOW}(F) = \{ + *) \$ \}$

Pro: Can be automated
 Cons: Error messages are needed

	id	+	*	()	\$
E	$E \rightarrow T E_R$			$E \rightarrow T E_R$	<i>synch</i>	<i>synch</i>
E_R		$E_R \rightarrow + T E_R$			$E_R \rightarrow \epsilon$	$E_R \rightarrow \epsilon$
T	$T \rightarrow F T_R$	<i>synch</i>		$T \rightarrow F T_R$	<i>synch</i>	<i>synch</i>
T_R		$T_R \rightarrow \epsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \epsilon$	$T_R \rightarrow \epsilon$
F	$F \rightarrow \text{id}$	<i>synch</i>	<i>synch</i>	$F \rightarrow (E)$	<i>synch</i>	<i>synch</i>

synch: the driver pops current nonterminal A and skips input till *synch* token or skips input until one of $\text{FIRST}(A)$ is found

Phrase-Level Recovery

Change input stream by inserting missing tokens

For example: **id id** is changed into **id * id**

Pro: Can be fully automated

Cons: Recovery not always intuitive

Can then continue here

	id	+	*	()	\$
E	$E \rightarrow T E_R$			$E \rightarrow T E_R$	<i>synch</i>	<i>synch</i>
E_R		$E_R \rightarrow + T E_R$			$E_R \rightarrow \epsilon$	$E_R \rightarrow \epsilon$
T	$T \rightarrow F T_R$	<i>synch</i>		$T \rightarrow F T_R$	<i>synch</i>	<i>synch</i>
T_R	<i>insert *</i>	$T_R \rightarrow \epsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \epsilon$	$T_R \rightarrow \epsilon$
F	$F \rightarrow \mathbf{id}$	<i>synch</i>	<i>synch</i>	$F \rightarrow (E)$	<i>synch</i>	<i>synch</i>

*insert **: driver inserts missing * and retries the production

Error Productions

$$\begin{aligned}
 E &\rightarrow T E_R \\
 E_R &\rightarrow + T E_R \mid \varepsilon \\
 T &\rightarrow F T_R \\
 T_R &\rightarrow * F T_R \mid \varepsilon \\
 F &\rightarrow (E) \mid \mathbf{id}
 \end{aligned}$$

Add “*error production*”:

$$T_R \rightarrow F T_R$$

to ignore missing *, e.g.: **id id**

Pro: Powerful recovery method

Cons: Manual addition of productions

	id	+	*	()	\$
E	$E \rightarrow T E_R$			$E \rightarrow T E_R$	<i>synch</i>	<i>synch</i>
E_R		$E_R \rightarrow + T E_R$			$E_R \rightarrow \varepsilon$	$E_R \rightarrow \varepsilon$
T	$T \rightarrow F T_R$	<i>synch</i>		$T \rightarrow F T_R$	<i>synch</i>	<i>synch</i>
T_R	$T_R \rightarrow F T_R$	$T_R \rightarrow \varepsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \varepsilon$	$T_R \rightarrow \varepsilon$
F	$F \rightarrow \mathbf{id}$	<i>synch</i>	<i>synch</i>	$F \rightarrow (E)$	<i>synch</i>	<i>synch</i>

Shift-Reduce Parsing

Grammar:

$S \rightarrow a A B e$

$A \rightarrow A b c \mid b$

$B \rightarrow d$

Reducing a sentence:

$a \underline{b} b c d e$

$a A \underline{b} c d e$

$a A \underline{d} e$

$\underline{a A B e}$

S

Shift-reduce corresponds to a rightmost derivation:

$S \Rightarrow_{rm} a A B e$

$\Rightarrow_{rm} a A d e$

$\Rightarrow_{rm} a A b c d e$

$\Rightarrow_{rm} a b b c d e$

