

Compilerconstructie

najaar 2013

<http://www.1iacs.nl/home/rvv11et/coco/>

Rudy van Vliet

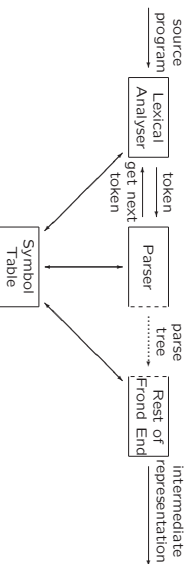
kamer 124 Snellius, tel. 071-527 5777
rvv11et(at)1iacs(dot)nl

college 3, dinsdag 17 september 2013

Syntax Analysis (1)

1

4.1 Parser's Position in a Compiler



- Obtain string of tokens
- Verify that string can be generated by the grammar
- Report and recover from syntax errors

3

4.2 Context-Free Grammars

Context-free grammar is a 4-tuple with

- A set of *nonterminals* (syntactic variables)
- A set of tokens (*terminal* symbols)
- A designated *start* symbol (nonterminal)
- A set of *productions*: rules how to decompose nonterminals

Example: CFG for simple arithmetic expressions:

$G = (\{expr, term, factor\}, \{id, +, -, *, /, (,)\}, expr, P)$

with productions P :

$expr \rightarrow expr + term \mid expr - term \mid term$
 $term \rightarrow term * factor \mid term / factor \mid factor$
 $factor \rightarrow (expr) \mid id$

5

Notational Conventions (Example)

CFG for simple arithmetic expressions:

$G = (\{expr, term, factor\}, \{id, +, -, *, /, (,)\}, expr, P)$

with productions P :

$expr \rightarrow expr + term \mid expr - term \mid term$
 $term \rightarrow term * factor \mid term / factor \mid factor$
 $factor \rightarrow (expr) \mid id$

Can be rewritten concisely as:

$E \rightarrow E + T \mid E - T \mid T$
 $T \rightarrow T * F \mid T / F \mid F$
 $F \rightarrow (E) \mid id$

7

4 Syntax Analysis

- Every language has rules prescribing the syntactic structure of the programs:
 - functions, made up of declarations and statements
 - statements made up of expressions
 - expressions made up of tokens
- Syntax of programming-language constructs can be described by CFG
 - Precise syntactic specification
 - Automatic construction of parsers for certain classes of grammars
 - Structure imparted to language by grammar is useful for translating source programs into object code
 - New language constructs can be added easily
- Syntax analysis is performed by parser

2

Parsing

Finding parse tree for given string

- Universal (any CFG)
 - Cocke-Younger-Kasami
 - Earley
- Top-down (CFG with restrictions)
 - Predictive parsing
 - LL (Left-to-right, Leftmost derivation) methods
 - LL(1): LL parser, needs only one token to look ahead
- Bottom-up (CFG with restrictions)

Today: top-down parsing

Next week: bottom-up parsing

4

Notational Conventions

1. Terminals:
 a, b, c, \dots ; specific terminals: $+, *, (,), 0, 1, id, if, \dots$
2. Nonterminals:
 A, B, C, \dots ; specific nonterminals: $S, expr, stmt, \dots, E, \dots$
3. Grammar symbols: X, Y, Z
4. Strings of terminals: u, v, w, x, y, z
5. Strings of grammar symbols: $\alpha, \beta, \gamma, \dots$
Hence, generic production: $A \rightarrow \alpha$
6. A-productions:
 $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k \quad \Rightarrow \quad A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$
Alternatives for A
7. By default, head of first production is start symbol

6

Derivations

Example grammar:

$E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid id$

- In each step, a nonterminal is replaced by body of one of its productions, e.g.,

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

- One-step derivation:
 $\alpha A \beta \Rightarrow \alpha \gamma \beta$, where $A \rightarrow \gamma$ is production in grammar
- Derivation in zero or more steps: $\xrightarrow{*}$
- Derivation in one or more steps: $\xrightarrow{+}$

8

Derivations

- If $S \xrightarrow{*} \alpha$, then α is **sentential form** of G
- If $S \xrightarrow{*} \alpha$ and α has no nonterminals, then α is **sentence** of G
- **Language generated** by G is $L(G) = \{w \mid w \text{ is sentence of } G\}$
- **Leftmost derivation**: $w \xrightarrow{lm} \gamma \delta w$
- If $S \xrightarrow{*} \alpha$, then α is **left sentential form** of G
- **Rightmost derivation**: $\gamma \delta w \xrightarrow{rm} \gamma \delta w$

Example of leftmost derivation:

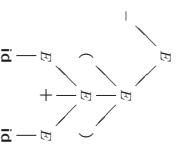
$$E \xrightarrow{lm} -E \xrightarrow{lm} -(E) \xrightarrow{lm} -(E + E) \xrightarrow{lm} -(\text{id} + E) \xrightarrow{lm} -(\text{id} + \text{id})$$

9

Parse Trees and Derivations

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

$$E \xrightarrow{lm} -E \xrightarrow{lm} -(E) \xrightarrow{lm} -(E + E) \xrightarrow{lm} -(\text{id} + E) \xrightarrow{lm} -(\text{id} + \text{id})$$



Many-to-one relationship between derivations and parse trees...

11

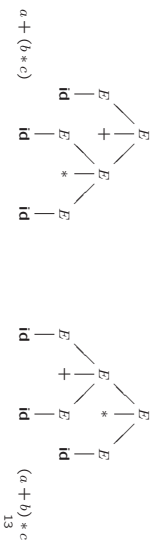
Ambiguity

More than one leftmost/rightmost derivation for same sentence

Example:

$$a + b * c$$

$$\begin{aligned} E &\Rightarrow E + E & E &\Rightarrow E * E \\ &\Rightarrow \text{id} + E & &\Rightarrow E + E * E \\ &\Rightarrow \text{id} + E * E & &\Rightarrow \text{id} + E * E \\ &\Rightarrow \text{id} + \text{id} * E & &\Rightarrow \text{id} + \text{id} * E \\ &\Rightarrow \text{id} + \text{id} * \text{id} & &\Rightarrow \text{id} + \text{id} * \text{id} \end{aligned}$$



13

Eliminating ambiguity

Example: ambiguous "dangling-else"-grammar

$$\begin{aligned} \text{stmt} &\rightarrow \text{if } \text{expr} \text{ then } \text{stmt} \\ &\mid \text{if } \text{expr} \text{ then } \text{stmt} \text{ else } \text{stmt} \\ &\mid \text{other} \end{aligned}$$

Only matched statements between **then** and **else**...

15

Parse Tree

(from college 1)

(derivation tree in F12)

- The root of the tree is labelled by the start symbol
- Each leaf of the tree is labelled by a terminal (\equiv token) or ϵ (\equiv empty)
- Each interior node is labelled by a nonterminal
- If node A has children X_1, X_2, \dots, X_n , then there must be a production $A \rightarrow X_1 X_2 \dots X_n$

Yield of the parse tree: the sequence of leafs (left to right)

10

4.3.1 Why Regular Expressions For Lexical Syntax?

- Convenient way to modularize front end \approx simplifies design
- Regular expressions powerful enough for lexical syntax
- Regular expressions easier to understand than grammars
- More efficient lexical analysers can be constructed automatically from regular expressions than from arbitrary grammars

12

Eliminating ambiguity

- Sometimes ambiguity can be eliminated
- Example: "dangling-else"-grammar

$$\begin{aligned} \text{stmt} &\rightarrow \text{if } \text{expr} \text{ then } \text{stmt} \\ &\mid \text{if } \text{expr} \text{ then } \text{stmt} \text{ else } \text{stmt} \\ &\mid \text{other} \end{aligned}$$

Here, **other** is any other statement

if E_1 then if E_2 then S_1 else S_2



14

Eliminating ambiguity

Example: ambiguous "dangling-else"-grammar

$$\begin{aligned} \text{stmt} &\rightarrow \text{if } \text{expr} \text{ then } \text{stmt} \\ &\mid \text{if } \text{expr} \text{ then } \text{stmt} \text{ else } \text{stmt} \\ &\mid \text{other} \end{aligned}$$

Equivalent unambiguous grammar

$$\begin{aligned} \text{stmt} &\rightarrow \text{matchedstmt} \\ &\mid \text{openstmt} \\ \text{matchedstmt} &\rightarrow \text{if } \text{expr} \text{ then } \text{matchedstmt} \text{ else } \text{matchedstmt} \\ &\mid \text{other} \\ \text{openstmt} &\rightarrow \text{if } \text{expr} \text{ then } \text{stmt} \\ &\mid \text{if } \text{expr} \text{ then } \text{matchedstmt} \text{ else } \text{openstmt} \end{aligned}$$

Only one parse tree for **if E_1 then if E_2 then S_1 else S_2**
Associates each **else** with closest previous unmatched **then**

16

2.4 Parsing (Top-Down Example)

from college 1

```
stmt → expr ;
      | if (expr )stmt
      | for (optexpr ; optexpr )stmt
      | other
optexpr → ε
      | expr
```

How to determine parse tree for

```
for (; expr ; expr )other
```

Use lookahead: current terminal in input

17

Recursive Descent Parsing

Recursive procedure for each nonterminal

```
void A()
1) { Choose an A-production,  $A \rightarrow X_1X_2 \dots X_k$ ;
2)   for ( $i = 1$  to  $k$ )
3)   { if ( $X_i$  is nonterminal)
4)     call procedure  $X_i()$ ;
5)   } else if ( $X_i$  equals current input symbol  $a$ )
6)     advance input to next symbol;
7)   } else /* an error has occurred */;
}
```

Pseudocode is nondeterministic

19

Predictive Parsing

from college 1

- Recursive-descent parsing . . .
- Predictive parsing is a special form of recursive-descent parsing:
 - The lookahead symbol unambiguously determines the production for each nonterminal

Simple example:

```
stmt → expr ;
      | if (expr )stmt
      | for (optexpr ; optexpr )stmt
      | other
```

21

Using FIRST

from college 1

- Let α be string of grammar symbols
- $\text{FIRST}(\alpha)$ is the set of terminals that appear as first symbols of strings generated from α

Simple example:

```
stmt → expr ;
      | if (expr )stmt
      | for (optexpr ; optexpr )stmt
      | other
```

Right-hand side may start with nonterminal . . .

23

Predictive Parsing

from college 1

- Recursive-descent parsing is a top-down parsing method:
 - Executes a set of recursive procedures to process the input
 - Every nonterminal has one (recursive) procedure parsing the nonterminal's syntactic category of input tokens
- Predictive parsing . . .

18

Recursive Descent

- One may use backtracking:
 - Try each A -production in some order
 - In case of failure at line 7 (or call in line 4), return to line 1 and try another A -production
 - Input pointer must then be reset, so store initial value input pointer in local variable

- Example in book

- Backtracking is rarely needed: predictive parsing

20

Predictive Parsing (Example)

from college 1

```
void stmt()
{ switch (lookahead)
  { case expr:
    match(expr); match(';'); break;
    case ff:
    match(ff); match('('); match(expr); match(')'); stmt();
    break;
    case for:
    match(for); match('(');
    optexpr(); match(';'); optexpr(); match(';'); optexpr();
    match(')'); stmt(); break;
    case other:
    match(other); break;
    default:
    report("syntax error");
  }
}

void match(terminal t)
{ if (lookahead==t) lookahead = nextTerminal;
  else report("syntax error");
}
```

22

Using FIRST

from college 1

- Let α be string of grammar symbols
- $\text{FIRST}(\alpha)$ is the set of terminals that appear as first symbols of strings generated from α
- When a nonterminal has multiple productions, e.g.,

$A \rightarrow \alpha \mid \beta$

then $\text{FIRST}(\alpha)$ and $\text{FIRST}(\beta)$ must be disjoint in order for predictive parsing to work

24

Left Recursion

- Productions of the form $A \rightarrow A\alpha \mid \beta$ are left-recursive
 - β does not start with A
 - Example: $E \rightarrow E + T \mid T$
- Top-down parser may loop forever if grammar has left-recursive productions
- Left-recursive productions can be eliminated by rewriting productions

25

Left Recursion Elimination

General left recursion

- Left recursion involving two or more steps

$$\begin{aligned} S &\rightarrow Ba \mid b \\ B &\rightarrow AA \mid a \\ A &\rightarrow Ac \mid Sd \end{aligned}$$

- S is left-recursive because

$$S \Rightarrow Ba \Rightarrow AaA \Rightarrow SdAa \quad (\text{not immediately left-recursive})$$

27

General Left Recursion Elimination

Algorithm for G with **no cycles or ϵ -productions**

- 1) arrange nonterminals in some order A_1, A_2, \dots, A_n
- 2) **for** ($i = 1$ to n)
- 3) **for** ($j = 1$ to $i - 1$)
- 4) { replace each production of form $A_i \rightarrow A_j\gamma$ by the productions $A_i \rightarrow \delta_1\gamma \mid \delta_2\gamma \mid \dots \mid \delta_k\gamma$, where $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ are all current A_j -productions
- 5) }
- 6) eliminate immediate left recursion among A_i -productions
- 7) }

Example with $A \rightarrow \epsilon$

29

Left Factoring (Example)

- Which production to choose when input token is **if**?

$$\begin{array}{l} \text{stmt} \rightarrow \text{if expr then stmt} \\ \quad \mid \text{if expr then stmt else stmt} \\ \quad \mid \text{other} \\ \text{expr} \rightarrow b \end{array}$$
- Or abstract:

$$\begin{array}{l} S \rightarrow iEiS \mid iEiScS \mid a \\ E \rightarrow b \end{array}$$
- Left-factored: ...

31

Left Recursion Elimination

Immediate left recursion

- Productions of the form $A \rightarrow A\alpha \mid \beta$
 - Can be eliminated by replacing the productions by

$$\begin{array}{ll} A &\rightarrow \beta A' & (A' \text{ is new nonterminal}) \\ A' &\rightarrow \alpha A' \mid \epsilon & (A' \rightarrow \alpha A' \text{ is right recursive}) \end{array}$$

- Procedure:

1. Group A -productions as

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

2. Replace A -productions by

$$\begin{array}{l} A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A' \\ A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon \end{array}$$

26

General Left Recursion Elimination

$$\begin{array}{l} S \rightarrow Ba \mid b \\ B \rightarrow AA \mid a \\ A \rightarrow Ac \mid Sd \end{array}$$

- We order nonterminals: S, B, A ($n = 3$)
- Variables may only 'point forward'
- $i = 1$ and $i = 2$: nothing to do
- $i = 3$:
 - substitute $A \rightarrow Sd$
 - substitute $A \rightarrow Bad$
 - eliminate immediate left-recursion in A -productions

28

Left Factoring

Another transformation to produce grammar suitable for predictive parsing

- If $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$ and input begins with nonempty string derived from α
 - How to expand A ? To $\alpha\beta_1$ or to $\alpha\beta_2$?

- Solution: left-factoring
 - Replace two A -productions by

$$\begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta_1 \mid \beta_2 \end{array}$$

30

Left Factoring (Example)

What is result of left factoring for

$$S \rightarrow abs \mid abcA \mid aaaS \mid aab \mid aA$$

32

Non-Context-Free Language Constructs

- Declaration of identifiers before their use

$$L_1 = \{wcu \mid w \in \{a, b\}^*\}$$

- Number of formal parameters in function declaration equals number of actual parameters in function call
Function call may be specified by

$$stmt \rightarrow \mathbf{id} (expr_list)$$

$$expr_list \rightarrow expr_list, expr \mid expr$$

$$L_2 = \{a^n b^m c^n d^n \mid m, n \geq 1\}$$

Such checks are performed during semantic-analysis phase

33

4.4 Top-Down Parsing

- Construct parse tree,
 - starting from the root
 - creating nodes in preorder

Corresponds to finding leftmost derivation

34

Top-Down Parsing (Example)

-

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

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

- Non-left-recursive variant: ...

35

Top-Down Parsing (Example)

-

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

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

- Non-left-recursive variant:

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

- Top-down parse for input **id + id * id** ...
- At each step: determine production to be applied

36

Top-Down Parsing

- Recursive-descent parsing
 - Predictive parsing
 - Eliminate left-recursion from grammar
 - Left-factor the grammar
 - Compute FIRST and FOLLOW
 - Two variants:
 - * Recursive (recursive calls)
 - * Non-recursive (explicit stack)

37

Computing FIRST

Compute $\text{FIRST}(X)$ for all grammar symbols X :

- If X is terminal, then $\text{FIRST}(X) = \{X\}$
- If $X \rightarrow \epsilon$ is production, then add ϵ to $\text{FIRST}(X)$

- Repeat adding symbols to $\text{FIRST}(X)$ by looking at productions

$$X \rightarrow Y_1 Y_2 \dots Y_k$$

(see book) until all FIRST sets are stable

39

FIRST

- Let α be string of grammar symbols
- $\text{FIRST}(\alpha) =$ set of terminals/tokens which begin strings derived from α
- If $\alpha \xrightarrow{*} \epsilon$, then $\epsilon \in \text{FIRST}(\alpha)$
- Example

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

$$\text{FIRST}(FT') = \{(\mathbf{id})\}$$

- When nonterminal has multiple productions, e.g.,

$$A \rightarrow \alpha \mid \beta$$

and $\text{FIRST}(\alpha)$ and $\text{FIRST}(\beta)$ are disjoint, we can choose between these A -productions by looking at next input symbol

38

FIRST (Example)

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

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

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

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

40

FOLLOW

- Let A be nonterminal
- $\text{FOLLOW}(A)$ is set of terminals/tokens that can appear immediately to the right of A in sentential form:

$$\text{FOLLOW}(A) = \{a \mid S \xRightarrow{*} \alpha A \alpha \beta\}$$
- Compute $\text{FOLLOW}(A)$ for all nonterminals A
See book

41

FIRST and FOLLOW (Example)

$$\begin{aligned}
 E &\rightarrow TE^i \\
 E^i &\rightarrow +TE^i \mid \epsilon \\
 T &\rightarrow FT^i \\
 T^i &\rightarrow *FT^i \mid \epsilon \\
 F &\rightarrow (E) \mid \text{id}
 \end{aligned}$$

$$\begin{aligned}
 \text{FIRST}(E) &= \text{FIRST}(T) = \text{FIRST}(F) = \{(\text{id})\} \\
 \text{FIRST}(E^i) &= \{+, \epsilon\} \\
 \text{FIRST}(T^i) &= \{*, \epsilon\} \\
 \text{FOLLOW}(E) &= \text{FOLLOW}(E^i) = \{), \$\} \\
 \text{FOLLOW}(T) &= \text{FOLLOW}(T^i) = \{+, \epsilon\} \\
 \text{FOLLOW}(F) &= \{*, +, \epsilon, \$\}
 \end{aligned}$$

42

Parsing Tables

When next input symbol is a (terminal or input endmarker $\$$), we may choose $A \rightarrow \alpha$

- if $a \in \text{FIRST}(\alpha)$
- if $(\alpha = \epsilon \text{ or } \alpha \xRightarrow{*} \epsilon)$ and $a \in \text{FOLLOW}(A)$

Algorithm to construct parsing table $M[A, a]$

```

for (each production  $A \rightarrow \alpha$ )
  for (each  $a \in \text{FIRST}(\alpha)$ )
    add  $A \rightarrow \alpha$  to  $M[A, a]$ ;
  if ( $\epsilon \in \text{FIRST}(\alpha)$ )
    for (each  $b \in \text{FOLLOW}(A)$ )
      add  $A \rightarrow \alpha$  to  $M[A, b]$ ;
}
if  $M[A, a]$  is empty, set  $M[A, a]$  to error.
    
```

43

LL(1) Grammars

- LL(1)
 - Left-to-right scanning of input, Leftmost derivation, 1 token to look ahead suffices for predictive parsing
- Grammar G is LL(1), if and only if for two distinct productions $A \rightarrow \alpha \mid \beta$,
 - α and β do not both derive strings beginning with same terminal a
 - at most one of α and β can derive ϵ
 - if $\beta \xRightarrow{*} \epsilon$, then α does not derive strings beginning with terminal $a \in \text{FOLLOW}(A)$
- In other words, ...
- Grammar G is LL(1), if and only if parsing table uniquely identifies production or signals error

45

LL(1) Grammars (Example)

- Not LL(1):

$$\begin{aligned}
 E &\rightarrow E + T \mid T \\
 T &\rightarrow T * F \mid F \\
 F &\rightarrow (E) \mid \text{id}
 \end{aligned}$$
- Non-left-recursive variant, LL(1):

$$\begin{aligned}
 E &\rightarrow TE^i \\
 E^i &\rightarrow +TE^i \mid \epsilon \\
 T &\rightarrow FT^i \\
 T^i &\rightarrow *FT^i \mid \epsilon \\
 F &\rightarrow (E) \mid \text{id}
 \end{aligned}$$

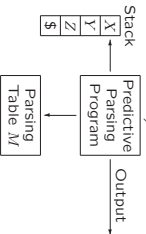
46

Nonrecursive Predictive Parsing

push \$ onto stack;
 push S onto stack;
 let a be first symbol of input w ;
 while ($X \neq \$$) /* stack is not empty */
 { if ($X = a$)
 { pop stack;
 let a be the next symbol of w ;
 }
 else if (X is terminal)
 error (O);
 else if ($M[X, a]$ is error entry)
 error (O);
 else if ($M[X, a] = X \rightarrow Y_1Y_2 \dots Y_k$)
 { output production $X \rightarrow Y_1Y_2 \dots Y_k$;
 pop stack;
 push Y_k, Y_{k-1}, \dots, Y_1 onto stack, with Y_1 on top;
 }
 }
 let X be top stack symbol;

47

Cf. top-down PDA from FIT2



48

Nonrec. Predictive Parsing (Example)

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

Matched	Stack	Input	Action
	ES	$\text{id} + \text{id} * \text{id} \$$	output $E \rightarrow TE'$
	$TE'S$	$\text{id} + \text{id} * \text{id} \$$	output $T \rightarrow FT'$
	$FTE'S$	$\text{id} + \text{id} * \text{id} \$$	output $F \rightarrow \text{id}$
	$idTE'S$	$\text{id} + \text{id} * \text{id} \$$	match id
	$idT'E'S$	$\text{id} + \text{id} * \text{id} \$$	match id
	$idT'E'S$	$+ \text{id} * \text{id} \$$	output $T' \rightarrow \epsilon$
	$idE'S$	$+ \text{id} * \text{id} \$$	output $E' \rightarrow +TE'$
	$+TE'S$	$+ \text{id} * \text{id} \$$	match $+$
	$TE'S$	$\text{id} * \text{id} \$$	output $T \rightarrow FT'$
	$...$	$...$	$...$

Note shift up of last column

49

Error Recovery in Predictive Parsing

Phrase-level recovery

- Local correction on remaining input that allows parser to continue
- Pointer to error routines in blank table entries
 - Change symbols
 - Insert symbols
 - Delete symbols
 - Print appropriate message
- Make sure that we do not enter infinite loop

51

4.1.3 Syntax Error Handling

- Good compiler should assist in identifying and locating errors
 - **Lexical errors:** compiler can easily detect and continue
 - **Syntax errors:** compiler can detect and often recover
 - **Semantic errors:** compiler can sometimes detect
 - **Logical errors:** hard to detect
- Three goals. The error handler should
 - Report errors clearly and accurately
 - Recover quickly to detect subsequent errors
 - Add minimal overhead to processing of correct programs

53

Error-Recovery Strategies

- Continue after error detection, restore to state where processing may continue, but...
 - No universally acceptable strategy.
 - **Panic-mode recovery:** discard input until token in designated set of *synchronizing* tokens is found
 - **Phrase-level recovery:** perform local correction on the input to repair error, e.g., insert missing semicolon
 - Has actually been used
 - **Error productions:** augment grammar with productions for erroneous constructs
 - **Global correction:** choose minimal sequence of changes to obtain correct string
- Costly, but yardstick for evaluating other strategies

55

Error Recovery in Predictive Parsing

Panic-mode recovery

- Discard input until token in set of designated synchronizing tokens is found
- Heuristics
 - Put all symbols in FOLLOW(A) into synchronizing set for A (and remove A from stack)
 - Add symbols based on hierarchical structure of language constructs
 - Add symbols in FIRST(A)
 - If $A \xrightarrow{*} \epsilon$, use production deriving ϵ as default
 - Add tokens to synchronizing sets of all other tokens

50

Predictive Parsing Issues

- What to do in case of multiply-defined entries?
 - Transform grammar
 - * Left-recursion elimination
 - * Left factoring
 - Not always applicable
- Designing grammar suitable for top-down parsing is hard
 - Left-recursion elimination and left factoring make grammar hard to read and to use in translation

Therefore: try to use automatic parser generators

52

Error Detection and Reporting

- **Viable-prefix property** of LL/LR parsers allow detection of syntax errors as soon as possible, i.e., as soon as prefix of input does not match prefix of any string in language (valid program)
- Reporting an error.
 - At least report line number and position
 - Print diagnostic message, e.g.,
 “**semicolon missing at this position**”

54

Compiler constructie

college 3
 Syntax Analysis (1)
 Chapters for reading: 4.1–4.4

Next week: also werkcollege

56