

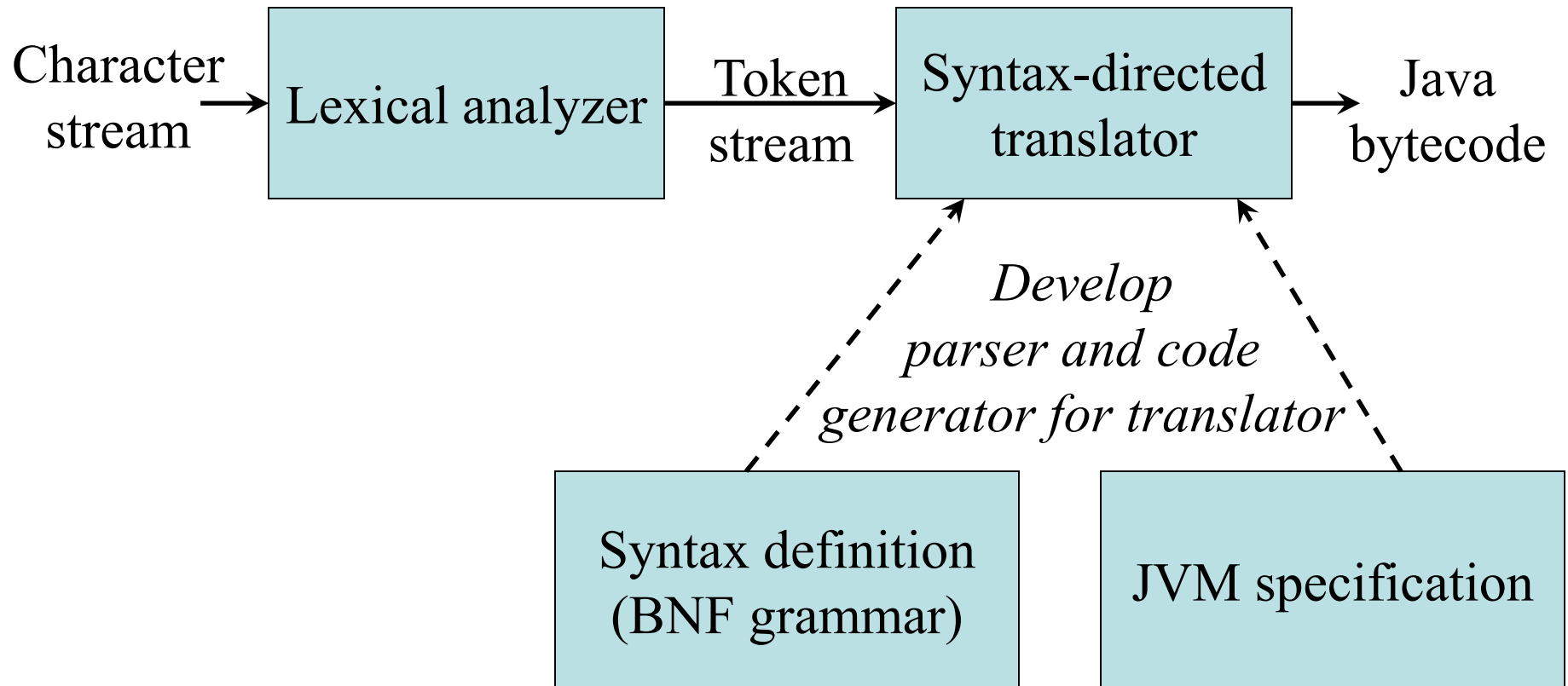
A Simple Syntax-Directed Translator

Chapter 2

Overview

- This chapter contains introductory material to Chapters 3 to 8
- Building a simple compiler
 - Syntax definition
 - Syntax directed translation
 - Predictive parsing

The Structure of our Compiler



Syntax Definition

- Context-free grammar is a 4-tuple with
 - A set of tokens (*terminal* symbols)
 - A set of *nonterminals*
 - A set of *productions*
 - A designated *start symbol*

Example Grammar

Context-free grammar for simple expressions:

$$G = \langle \{list, digit\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, list \rangle$$

with productions $P =$

$$list \rightarrow list + digit$$

$$list \rightarrow list - digit$$

$$list \rightarrow digit$$

$$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

Derivation

- Given a CF grammar we can determine the set of all *strings* (sequences of tokens) generated by the grammar using *derivation*
 - We begin with the start symbol
 - In each step, we replace one nonterminal in the current *sentential form* with one of the right-hand sides of a production for that nonterminal

Derivation for the Example Grammar

$$\begin{aligned}
 &\underline{list} \\
 \Rightarrow &\underline{list} + digit \\
 \Rightarrow &\underline{list} - digit + digit \\
 \Rightarrow &\underline{digit} - digit + digit \\
 \Rightarrow &\mathbf{9} - \underline{digit} + digit \\
 \Rightarrow &\mathbf{9} - \mathbf{5} + \underline{digit} \\
 \Rightarrow &\mathbf{9} - \mathbf{5} + \mathbf{2}
 \end{aligned}$$

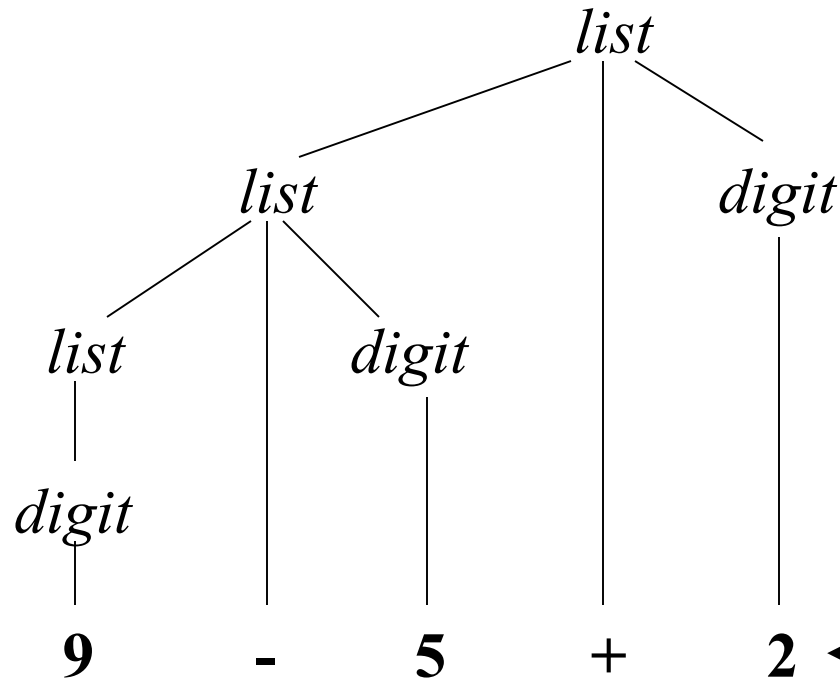
This is an example *leftmost derivation*, because we replaced the leftmost nonterminal (underlined) in each step

Parse Trees

- The root of the tree is labeled by the start symbol
- Each leaf of the tree is labeled by a terminal (=token) or ε
- Each interior node is labeled by a nonterminal
- If $A \rightarrow X_1 X_2 \dots X_n$ is a production, then node A has children X_1, X_2, \dots, X_n where X_i is a (non)terminal or ε (ε denotes the *empty string*)

Parse Tree for the Example Grammar

Parse tree of the string **9-5+2** using grammar G



The sequence of
leaves is called the
yield of the parse tree

Ambiguity

Consider the following context-free grammar:

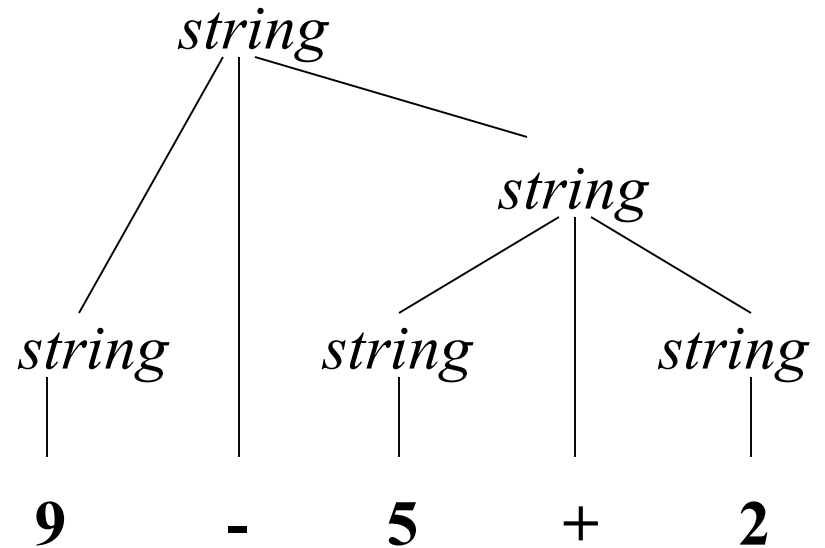
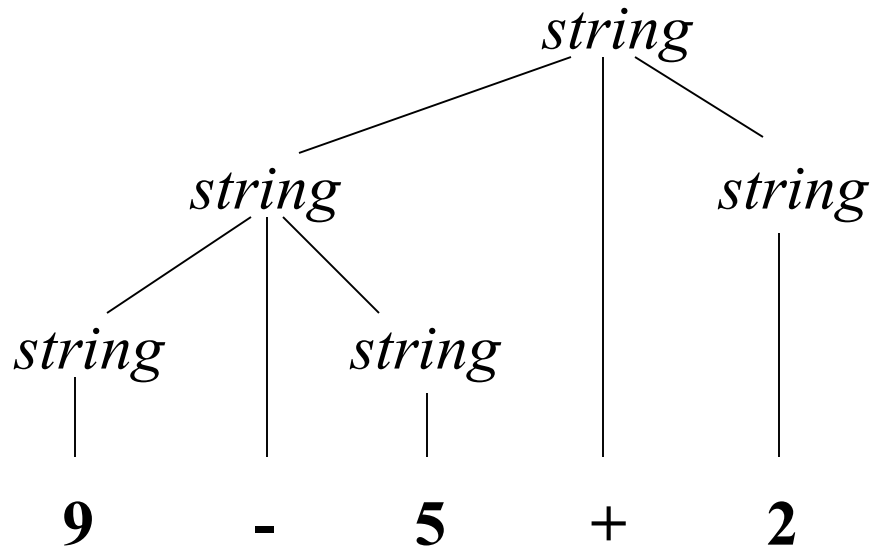
$$G = \langle \{string\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, string \rangle$$

with production $P =$

$$string \rightarrow string + string \mid string - string \mid \mathbf{0} \mid \mathbf{1} \mid \dots \mid \mathbf{9}$$

This grammar is *ambiguous*, because more than one parse tree generates the string **9-5+2**

Ambiguity (cont'd)



Associativity of Operators

Left-associative operators have *left-recursive* productions

$$left \rightarrow left + term \mid term$$

String **a+b+c** has the same meaning as **(a+b)+c**

Right-associative operators have *right-recursive* productions

$$right \rightarrow term = right \mid term$$

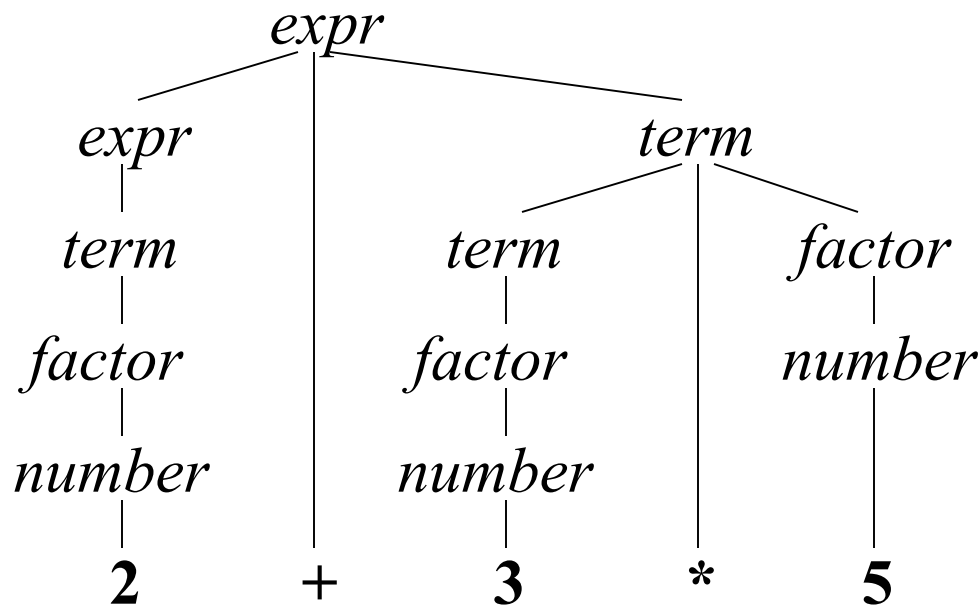
String **a=b=c** has the same meaning as **a=(b=c)**

Precedence of Operators

Operators with higher precedence “bind more tightly”

$$expr \rightarrow expr + term \mid term$$
$$term \rightarrow term * factor \mid factor$$
$$factor \rightarrow number \mid (expr)$$

String **2+3*5** has the same meaning as **2+(3*5)**



Syntax of Statements

$$\begin{aligned} stmt &\rightarrow \mathbf{id} := expr \\ &\quad | \mathbf{if} \, expr \, \mathbf{then} \, stmt \\ &\quad | \mathbf{if} \, expr \, \mathbf{then} \, stmt \, \mathbf{else} \, stmt \\ &\quad | \mathbf{while} \, expr \, \mathbf{do} \, stmt \\ &\quad | \mathbf{begin} \, opt_stmts \, \mathbf{end} \\ opt_stmts &\rightarrow stmt \, ; \, opt_stmts \\ &\quad | \varepsilon \end{aligned}$$

Syntax-Directed Translation

- Uses a CF grammar to specify the syntactic structure of the language
- AND associates a set of *attributes* with (non)terminals
- AND associates with each production a set of *semantic rules* for computing values for the attributes
- The attributes contain the translated form of the input after the computations are completed

Synthesized and Inherited Attributes

- An attribute is said to be ...
 - *synthesized* if its value at a parse-tree node is determined from the attribute values at the children of the node
 - *inherited* if its value at a parse-tree node is determined by the parent (by enforcing the parent's semantic rules)

Example Attribute Grammar

Production

$expr \rightarrow expr_1 + term$

$expr \rightarrow expr_1 - term$

$expr \rightarrow term$

$term \rightarrow 0$

$term \rightarrow 1$

...

$term \rightarrow 9$

Semantic Rule

$expr.t := expr_1.t \parallel term.t \parallel "+"$

$expr.t := expr_1.t \parallel term.t \parallel "-"$

$expr.t := term.t$

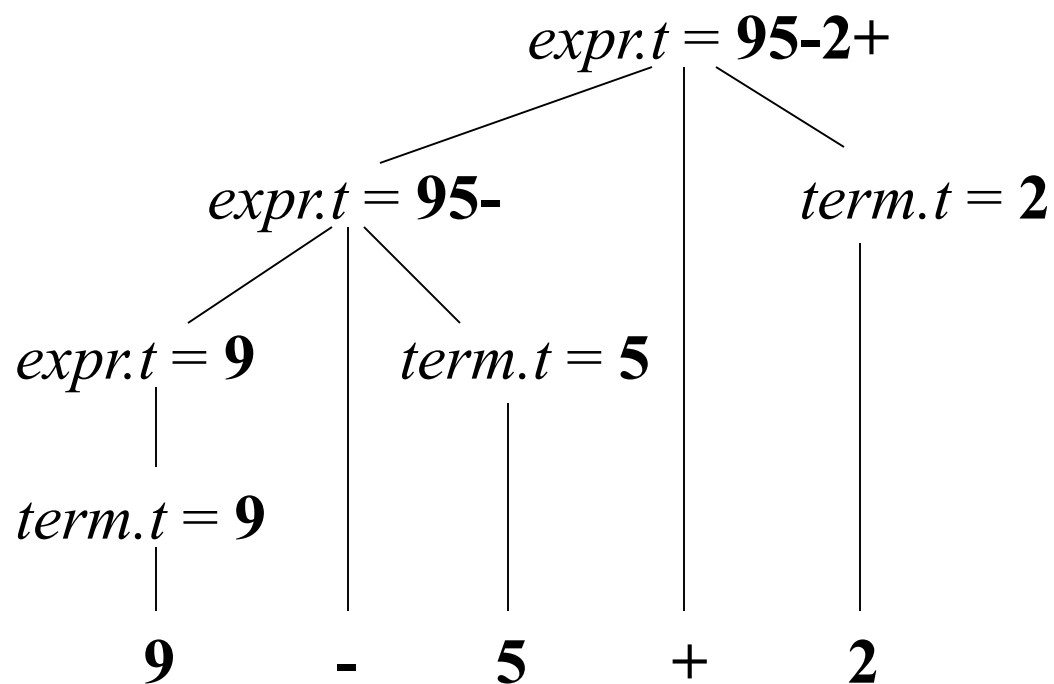
$term.t := "0"$

$term.t := "1"$

...

$term.t := "9"$

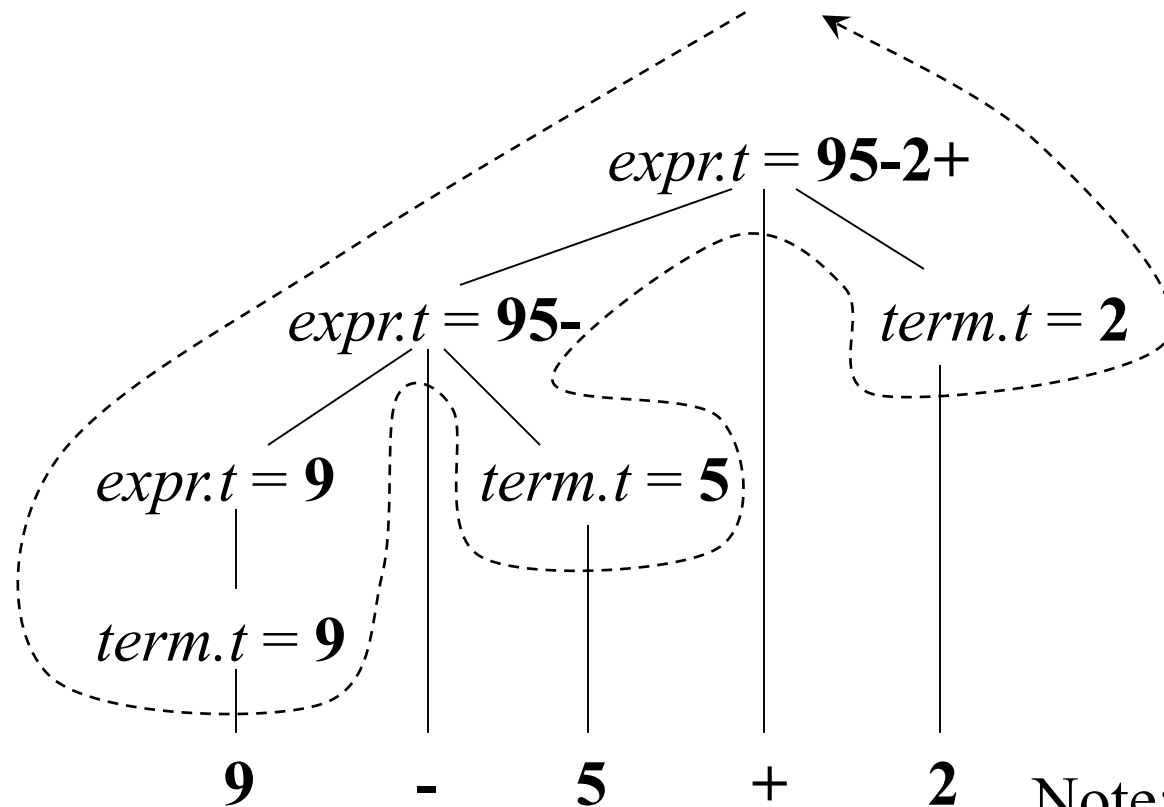
Example Annotated Parse Tree



Depth-First Traversals

```
procedure visit(n : node);  
begin  
    for each child m of n, from left to right do  
        visit(m);  
    evaluate semantic rules at node n  
end
```

Depth-First Traversals (Example)




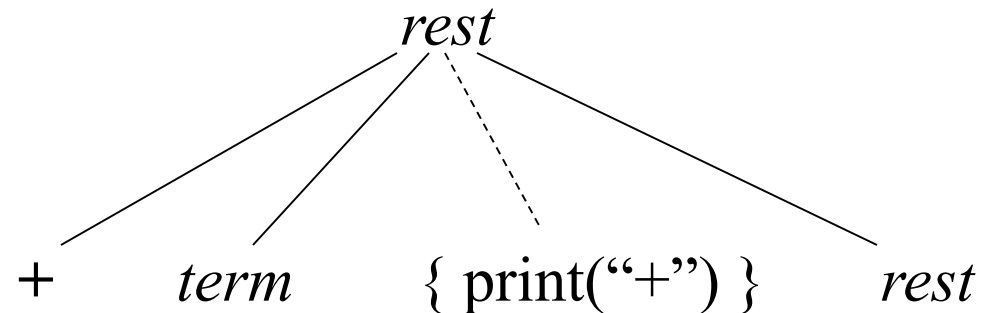
Note: all attributes are of the synthesized type

Translation Schemes

- A *translation scheme* is a CF grammar embedded with *semantic actions*

$$rest \rightarrow + term \{ \text{print}(\text{"+"}) \} rest$$

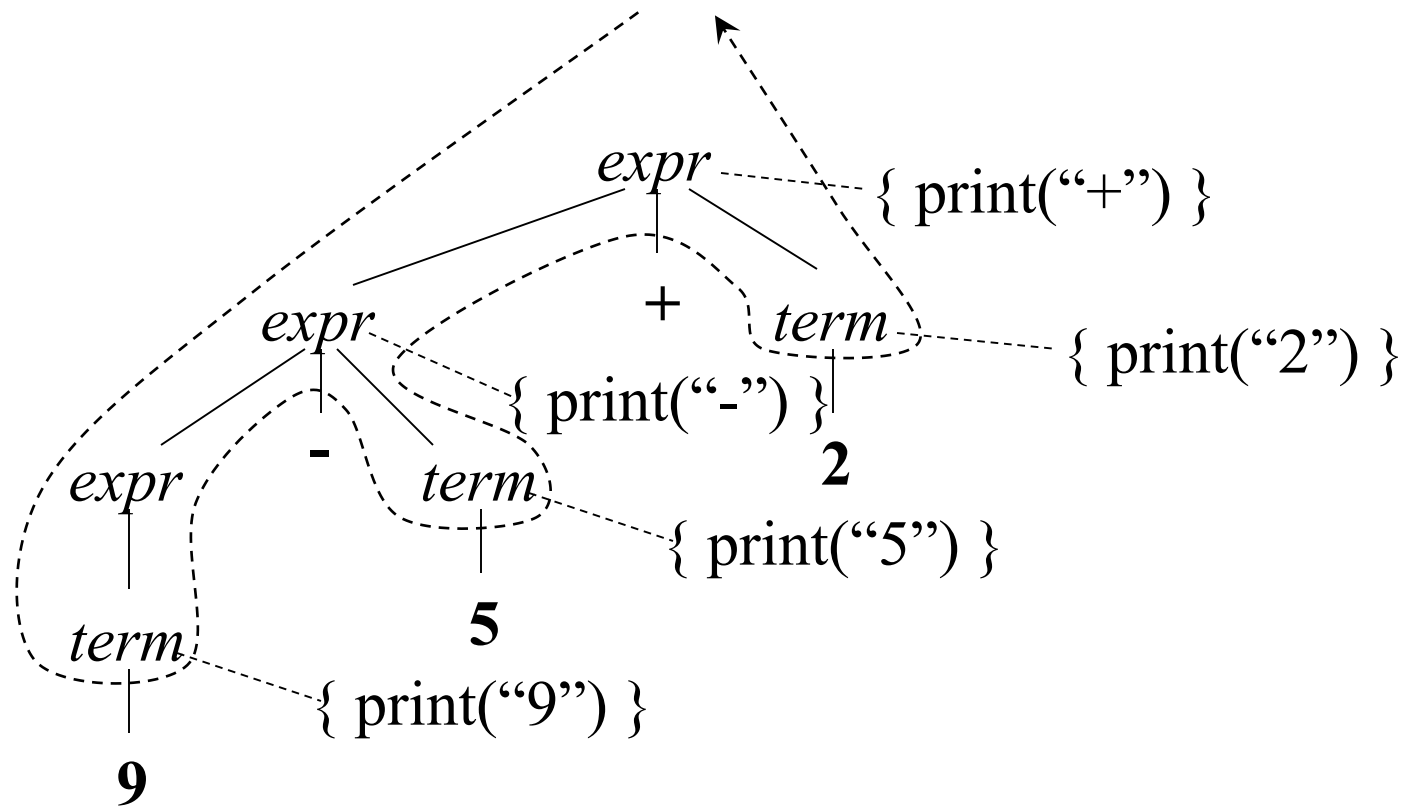

 Embedded
semantic action



Example Translation Scheme

$expr \rightarrow expr + term$	{ print(“+”) }
$expr \rightarrow expr - term$	{ print(“-”) }
$expr \rightarrow term$	
$term \rightarrow 0$	{ print(“0”) }
$term \rightarrow 1$	{ print(“1”) }
...	...
$term \rightarrow 9$	{ print(“9”) }

Example Translation Scheme (cont'd)



Translates **9-5+2** into postfix **95-2+**

Parsing

- Parsing = *process of determining if a string of tokens can be generated by a grammar*
- For any CF grammar there is a parser that takes at most $O(n^3)$ time to parse a string of n tokens
- Linear algorithms suffice for parsing programming language
- *Top-down parsing* “constructs” parse tree from root to leaves
- *Bottom-up parsing* “constructs” parse tree from leaves to root

Predictive Parsing

- *Recursive descent parsing* is a top-down parsing method
 - 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
- *Predictive parsing* is a special form of recursive descent parsing where we use one lookahead token to unambiguously determine the parse operations

Example Predictive Parser (Grammar)

$type \rightarrow simple$
 | **^ id**
 | **array [simple] of type**
 $simple \rightarrow \mathbf{integer}$
 | **char**
 | **num dotdot num**

Example Predictive Parser (Program Code)

```
procedure match(t : token);  
begin  
    if lookahead = t then  
        lookahead := nexttoken()  
    else error()  
end;
```

```
procedure type();  
begin  
    if lookahead in { 'integer', 'char', 'num' } then  
        simple()  
    else if lookahead = '^' then  
        match('^'); match(id)  
    else if lookahead = 'array' then  
        match('array'); match('['); simple();  
        match(']'); match('of'); type()  
    else error()  
end;
```

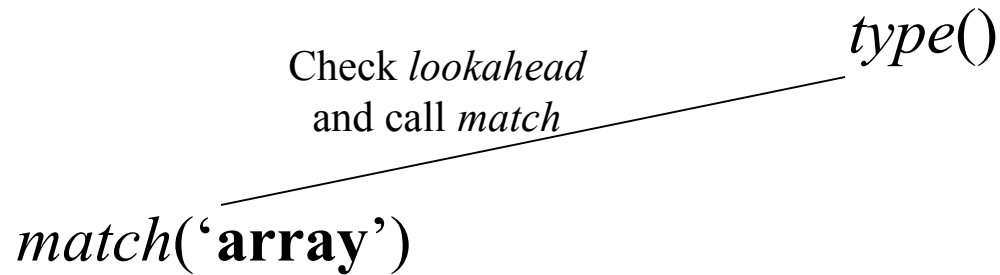
```
procedure simple();  
begin  
    if lookahead = 'integer' then  
        match('integer')  
    else if lookahead = 'char' then  
        match('char')  
    else if lookahead = 'num' then  
        match('num');  
        match('dotdot');  
        match('num')  
    else error()  
end;
```

Example Predictive Parser (Execution Step 1)

match('array')

Check *lookahead*
and call *match*

type()

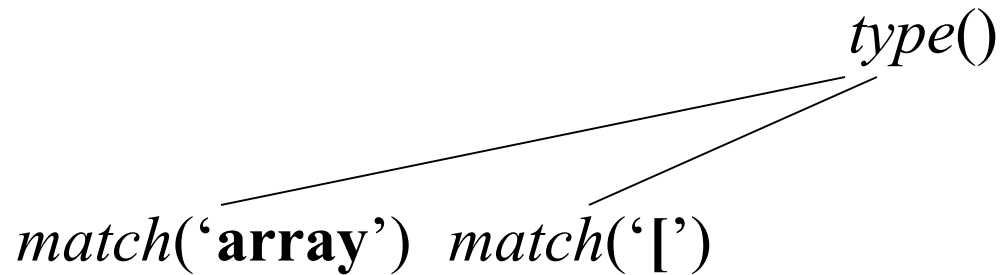


Input: **array** [**num** **dotdot** **num**] **of** **integer**

 ↑
 lookahead



Example Predictive Parser (Execution Step 2)

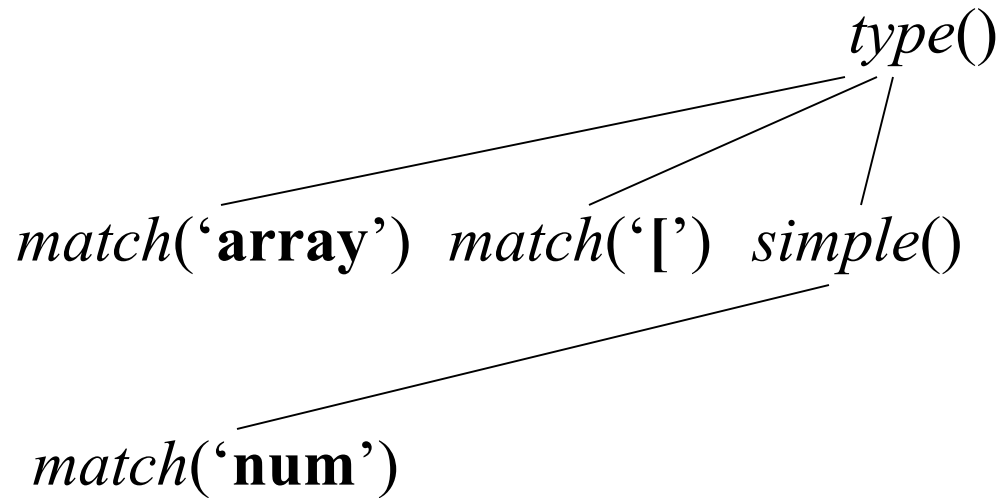


Input: **array** **|** **num** **dotdot** **num** **|** **of** **integer**

\uparrow

lookahead

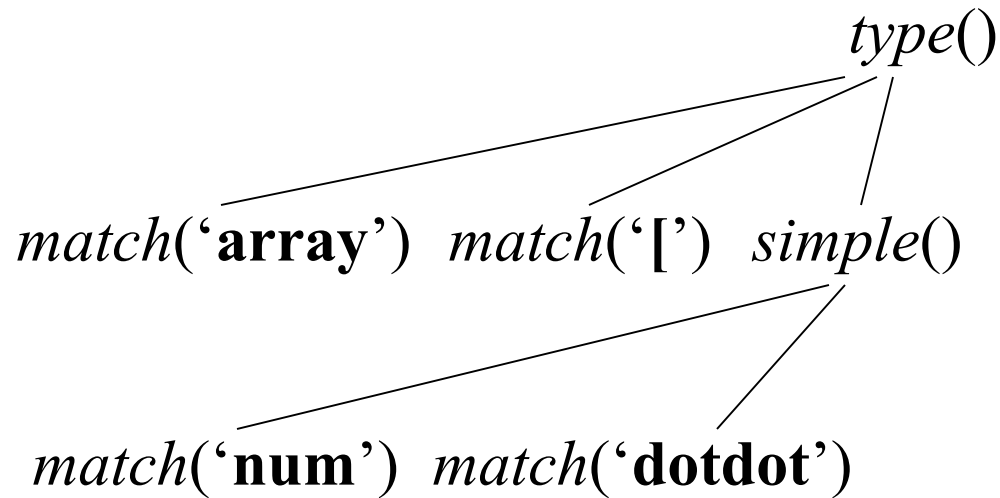
Example Predictive Parser (Execution Step 3)



Input: array [num dotdot num] of integer

 ↑
lookahead

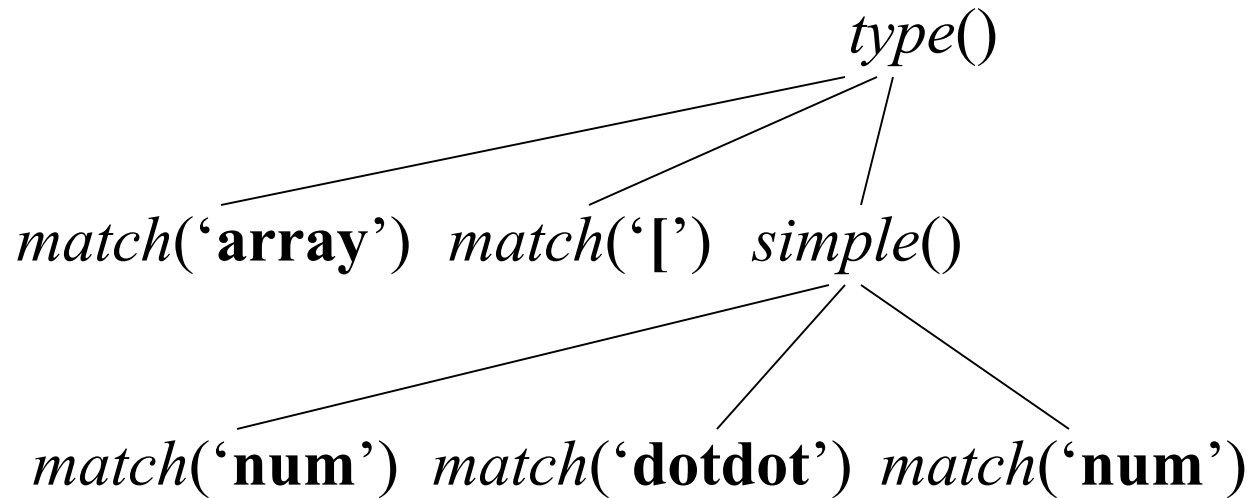
Example Predictive Parser (Execution Step 4)



Input: array [num dotdot num] of integer

 ↑
lookahead

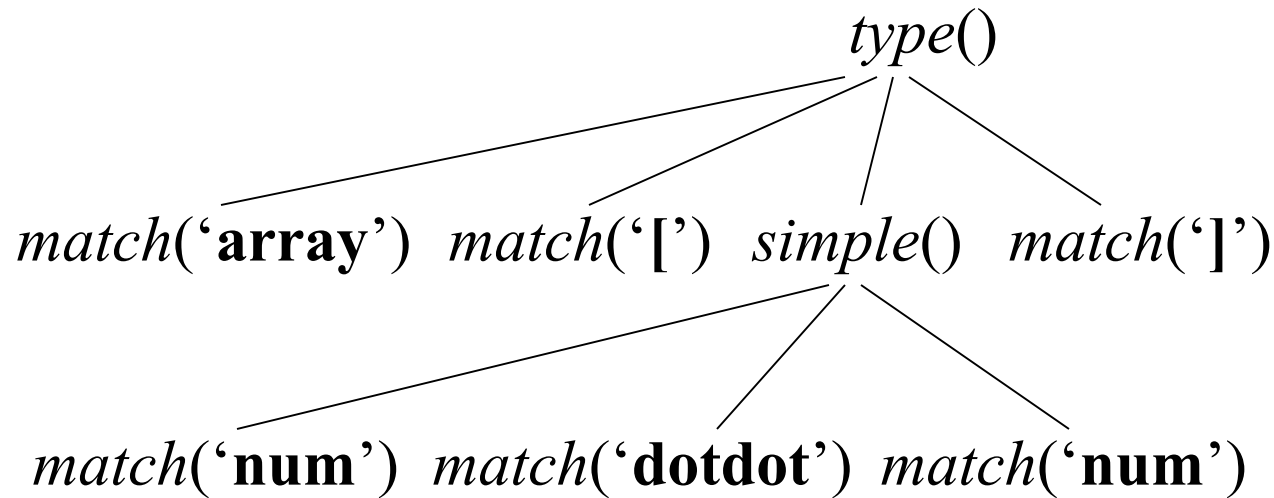
Example Predictive Parser (Execution Step 5)



Input: array [num dotdot num] of integer

 ↑
lookahead

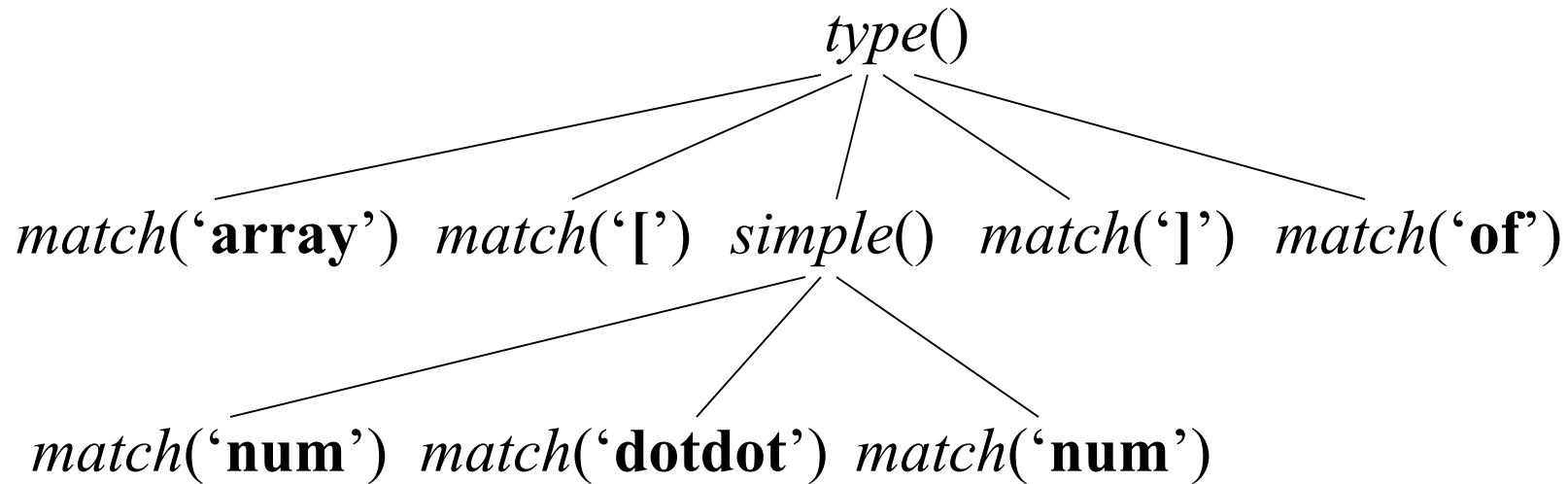
Example Predictive Parser (Execution Step 6)



Input: **array** **[** **num** **dotdot** **num** **]** **of** **integer**

↑
lookahead

Example Predictive Parser (Execution Step 7)



Input: array [num dotdot num] of integer

↑
lookahead

FIRST

$\text{FIRST}(\alpha)$ is the set of terminals that appear as the first symbols of one or more strings generated from α

$$\begin{array}{l}
 \textit{type} \rightarrow \textit{simple} \\
 \quad | \textbf{^ id} \\
 \quad | \textbf{array [simple] of type} \\
 \textit{simple} \rightarrow \textbf{integer} \\
 \quad | \textbf{char} \\
 \quad | \textbf{num dotdot num}
 \end{array}$$

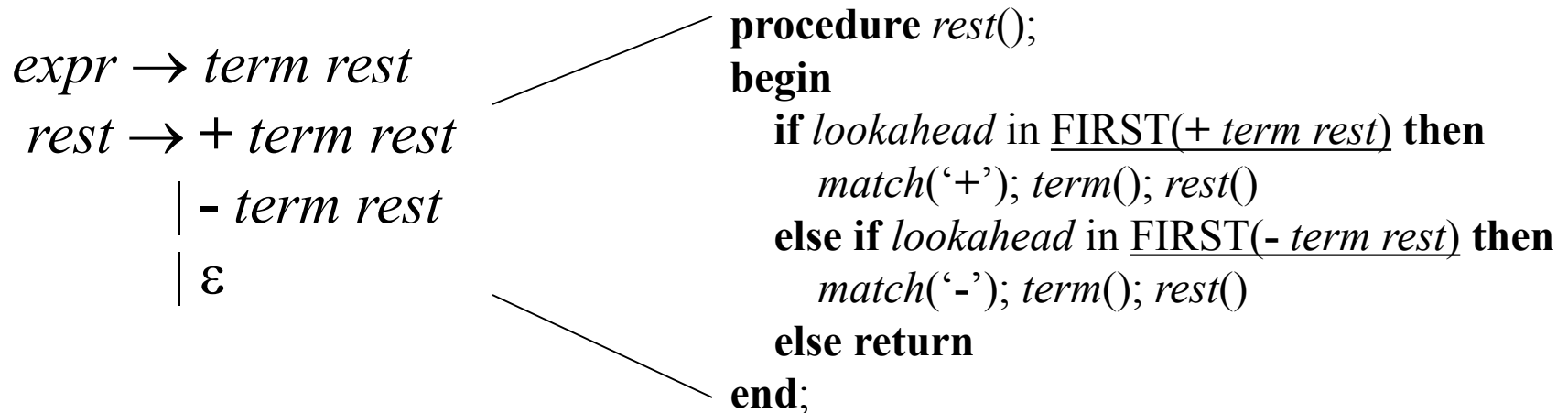
$\text{FIRST}(\textit{simple}) = \{ \textbf{integer, char, num} \}$

$\text{FIRST}(\textbf{^ id}) = \{ \textbf{^} \}$

$\text{FIRST}(\textit{type}) = \{ \textbf{integer, char, num, ^, array} \}$

Using FIRST

We use FIRST to write a predictive parser as follows



When a nonterminal A has two (or more) productions as in

$$\begin{array}{c}
 A \rightarrow \alpha \\
 | \beta
 \end{array}$$

Then FIRST(α) and FIRST(β) must be disjoint for predictive parsing to work

Left Factoring

When more than one production for nonterminal A starts with the same symbols, the FIRST sets are not disjoint

$$\begin{aligned} stmt &\rightarrow \mathbf{if\ expr\ then\ stmt} \\ &\quad | \mathbf{if\ expr\ then\ stmt\ else\ stmt} \end{aligned}$$

We can use *left factoring* to fix the problem

$$\begin{aligned} stmt &\rightarrow \mathbf{if\ expr\ then\ stmt\ opt_else} \\ opt_else &\rightarrow \mathbf{else\ stmt} \\ &\quad | \epsilon \end{aligned}$$

Left Recursion

When a production for nonterminal A starts with a self reference then a predictive parser loops forever

$$\begin{array}{c} A \rightarrow A \alpha \\ | \beta \\ | \gamma \end{array}$$

We can eliminate *left recursive productions* by systematically rewriting the grammar using *right recursive productions*

$$\begin{array}{c} A \rightarrow \beta R \\ | \gamma R \\ R \rightarrow \alpha R \\ | \varepsilon \end{array}$$

A Translator for Simple Expressions

$expr \rightarrow expr + term \quad \{ \text{print}("+") \}$

$expr \rightarrow expr - term \quad \{ \text{print}("-") \}$

$expr \rightarrow term$

$term \rightarrow 0 \quad \{ \text{print}("0") \}$

$term \rightarrow 1 \quad \{ \text{print}("1") \}$

...

$term \rightarrow 9 \quad \{ \text{print}("9") \}$

After left recursion elimination:

$expr \rightarrow term \ rest$

$rest \rightarrow + term \{ \text{print}("+") \} rest \mid - term \{ \text{print}("-") \} rest \mid \varepsilon$

$term \rightarrow 0 \{ \text{print}("0") \}$

$term \rightarrow 1 \{ \text{print}("1") \}$

...

$term \rightarrow 9 \{ \text{print}("9") \}$

$expr \rightarrow term\ rest$
 $rest \rightarrow +\ term\ \{ \text{print}("+") \}\ rest$
 $\quad | -\ term\ \{ \text{print}("-") \}\ rest$
 $\quad | \epsilon$

```

main()
{   lookahead = getchar();
    expr();
}

expr()
{   term();
    while (1) /* optimized by inlining rest()
                and removing recursive calls */
    {   if (lookahead == '+')
        {   match('+'); term(); putchar('+');
        }
        else if (lookahead == '-')
        {   match('-'); term(); putchar('-');
        }
        else break;
    }
}

term()
{   if (isdigit(lookahead))
    {   putchar(lookahead); match(lookahead);
    }
    else error();
}

match(int t)
{   if (lookahead == t)
        lookahead = getchar();
    else error();
}

error()
{   printf("Syntax error\n");
    exit(1);
}

```

Adding a Lexical Analyzer

- Typical tasks of the lexical analyzer:
 - Remove white space and comments
 - Encode constants as tokens
 - Recognize keywords
 - Recognize identifiers

The Lexical Analyzer

y := 31 + 28*x

Lexical analyzer
lexan()

<id, “y”> <assign, > <num, 31> <+, > <num, 28> <*, > <id, “x”>

token

tokenval
(token attribute)

Parser
parse()

Token Attributes

factor \rightarrow (*expr*)
| **num** { print(**num.value**) }

```
#define NUM 256 /* token returned by lexan */
```

```
factor()  
{  
    if (lookahead == '(')  
    {  
        match('('); expr(); match(')');  
    }  
    else if (lookahead == NUM)  
    {  
        printf(" %d ", tokenval); match(NUM);  
    }  
    else error();  
}
```

Symbol Table

The symbol table is globally accessible (to all phases of the compiler)

Each entry in the symbol table contains a string and a token value:

```
struct entry
{
    char *lexptr; /* lexeme (string) */
    int token;
};
struct entry symtable[];
```

insert(s, t): returns array index to new entry for string **s** token **t**

lookup(s): returns array index to entry for string **s** or 0

Possible implementations:

- simple C code (see textbook)
- hashtables

Identifiers

$$\textit{factor} \rightarrow (\textit{expr})$$

| **id** { print(**id**.string) }

```
#define ID 259 /* token returned by lexan() */

factor()
{
    if (lookahead == '(')
    {
        match('('); expr(); match(')');
    }
    else if (lookahead == ID)
    {
        printf(" %s ", symtable[tokenval].lexptr);
        match(NUM);
    }
    else error();
}
```

Handling Reserved Keywords

We simply initialize
the global symbol
table with the set of
keywords

```
/* global.h */
#define DIV 257 /* token */
#define MOD 258 /* token */
#define ID  259 /* token */

/* init.c */
insert("div", DIV);
insert("mod", MOD);

/* lexer.c */
int lexan()
{
    ...
    tokenval = lookup(lexbuf);
    if (tokenval == 0)
        tokenval = insert(lexbuf, ID);
    return symtable[p].token;
}
```

Handling Reserved Keywords (cont'd)

morefactors → **div** *factor* { print('DIV') } *morefactors*
 | **mod** *factor* { print('MOD') } *morefactors*
 | ...

```
/* parser.c */
morefactors()
{
    if (lookahead == DIV)
    {
        match(DIV); factor(); printf("DIV"); morefactors();
    }
    else if (lookahead == MOD)
    {
        match(MOD); factor(); printf("MOD"); morefactors();
    }
    else
        ...
}
```