

# Chapter 2

## A Simple Syntax-Directed Translator

# Overview



- This chapter contains introductory material to chapters 3 to 6 of the Dragon book
- Developing a working Java program that translates representative programming language statements into three-address code, an intermediate representation
- Emphasis is on the front end of a compiler (lexical analysis, parsing, intermediate code generation)

# Building a Simple Compiler



- Creating a syntax-directed translator that maps infix arithmetic expressions into postfix expressions
- Then, extending this translator to map code fragments into three-address code
- The working Java translator appears in Appendix A

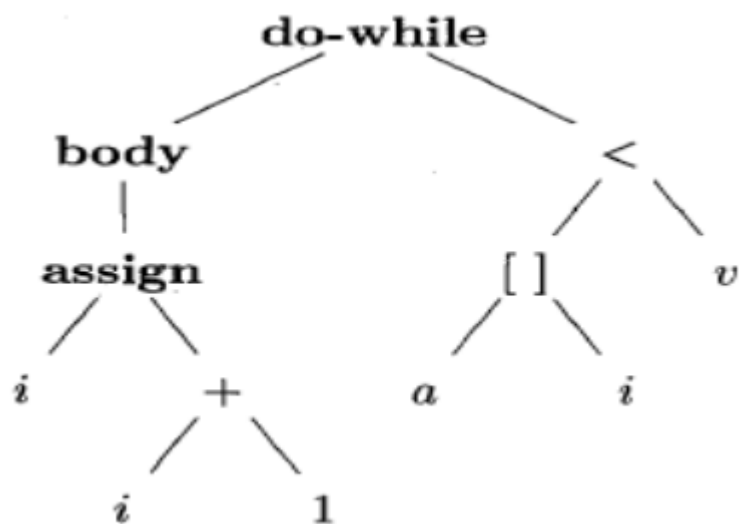
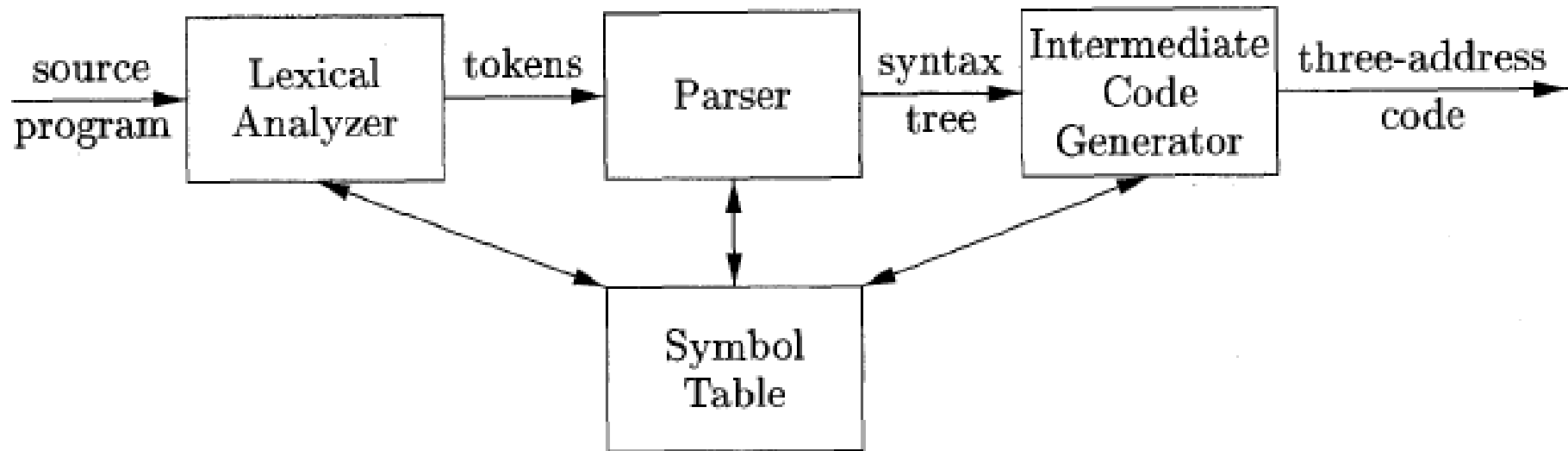
# Output of a Simple Compiler

```
{
int i; int j; float[100] a; float v; float x;

while ( true ) {
    do i = i+1; while ( a[i] < v );
    do j = j-1; while ( a[j] > v );
    if ( i >= j ) break;
    x = a[i]; a[i] = a[j]; a[j] = x;
}
}
```

```
1:  i = i + 1
2:  t1 = a [ i ]
3:  if t1 < v goto 1
4:  j = j - 1
5:  t2 = a [ j ]
6:  if t2 > v goto 4
7:  ifFalse i >= j goto 9
8:  goto 14
9:  x = a [ i ]
10: t3 = a [ j ]
11: a [ i ] = t3
12: a [ j ] = x
13: goto 1
14:
```

# A Model of a Compiler Front-end



`do i = i+1; while ( a[i] < v );`

```
1: i = i + 1
2: t1 = a [ i ]
3: if t1 < v goto 1
```

# Definition of Grammars



- Context-free grammar (CFG) 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



CFG for simple expressions:

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

with productions  $P$ :

$$list \rightarrow list + digit$$

$$list \rightarrow list - digit$$

$$list \rightarrow digit$$

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

$$list \rightarrow list + digit \mid list - digit \mid digit$$

# Syntax Definition using CFG



- If statement in C:  
**if** (expression) statement **else** statement

$$stmt \rightarrow \mathbf{if} ( expr ) stmt \mathbf{else} stmt$$

- List of parameters in a function call

$$call \rightarrow \mathbf{id} ( optparams )$$
$$optparams \rightarrow params / \epsilon$$
$$params \rightarrow params , param / param$$



# Derivation



- Given a CFG, 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

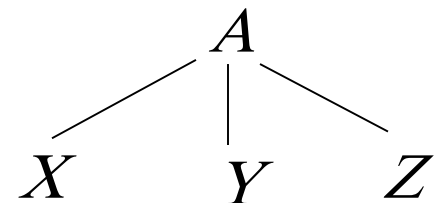
list  
 $\Rightarrow$  list + digit  
 $\Rightarrow$  list - digit + digit  
 $\Rightarrow$  digit - digit + digit  
 $\Rightarrow$  9 - digit + digit  
 $\Rightarrow$  9 - 5 + digit  
 $\Rightarrow$  9 - 5 + 2

- This is an example **leftmost derivation**, because we replaced the leftmost nonterminal (underlined) in each step
- Likewise, a **rightmost derivation** replaces the rightmost nonterminal in each step

# Parse Trees

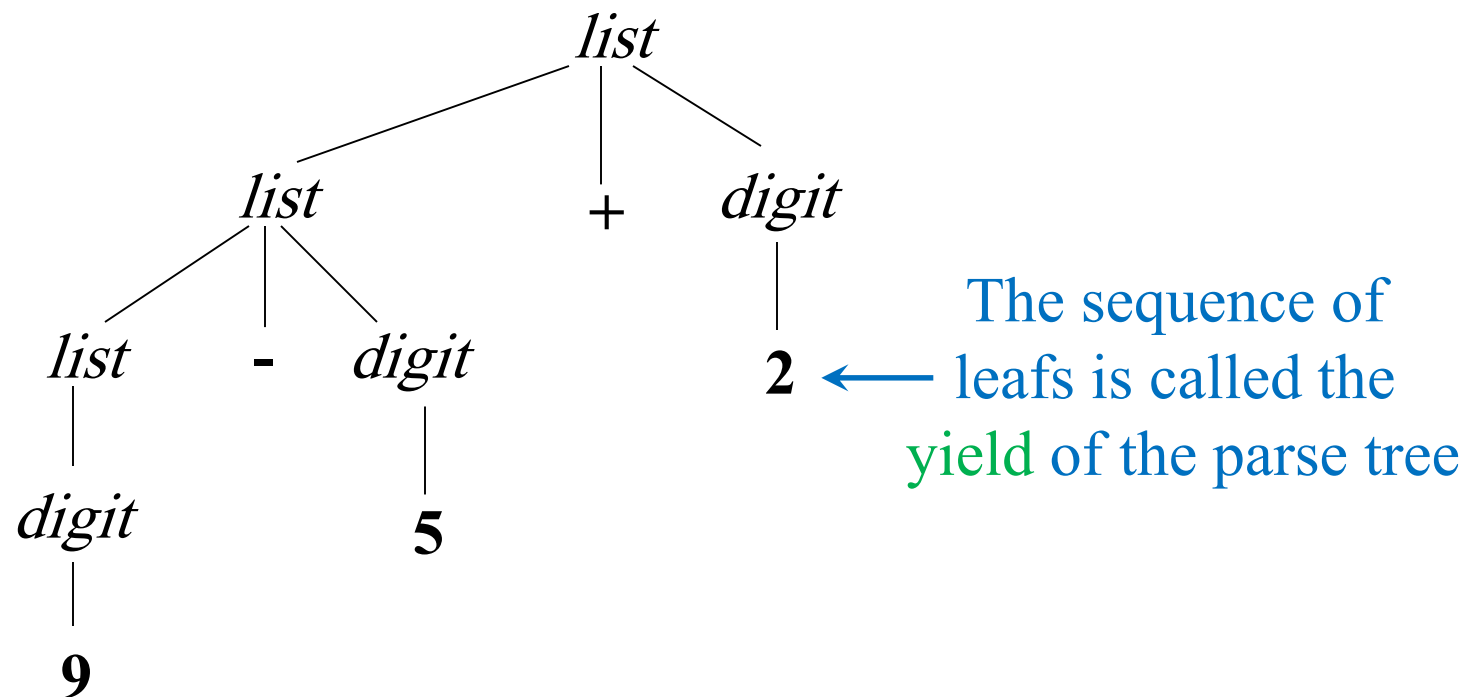
- The **root** of the tree is labeled by the start symbol
- Each **leaf** of the tree is labeled by a terminal (token) or  $\varepsilon$  (**empty string**)
- Each **interior node** is labeled by a nonterminal
- If  $A \rightarrow X_1 X_2 \dots X_n$  is a production, then node  $A$  has immediate **children**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

$$A \rightarrow X Y Z$$



# Parse Tree for the Example Grammar

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



# Ambiguity



Consider the following context-free grammar:

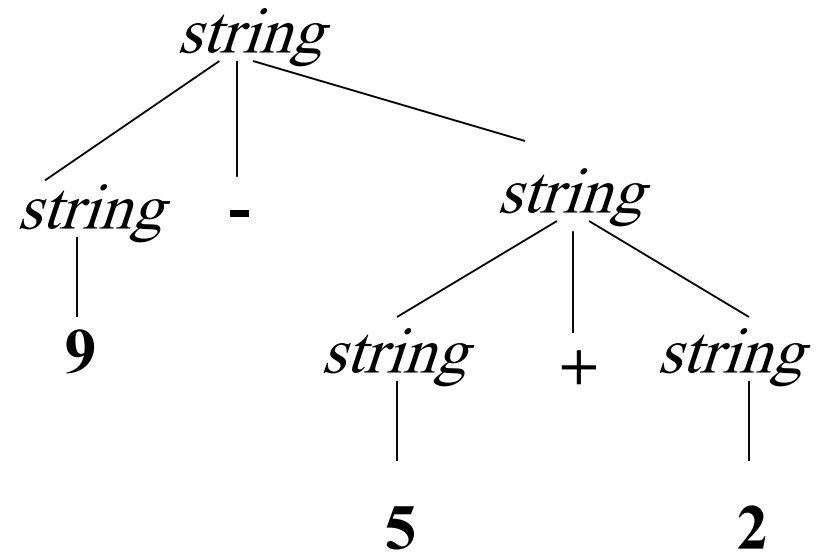
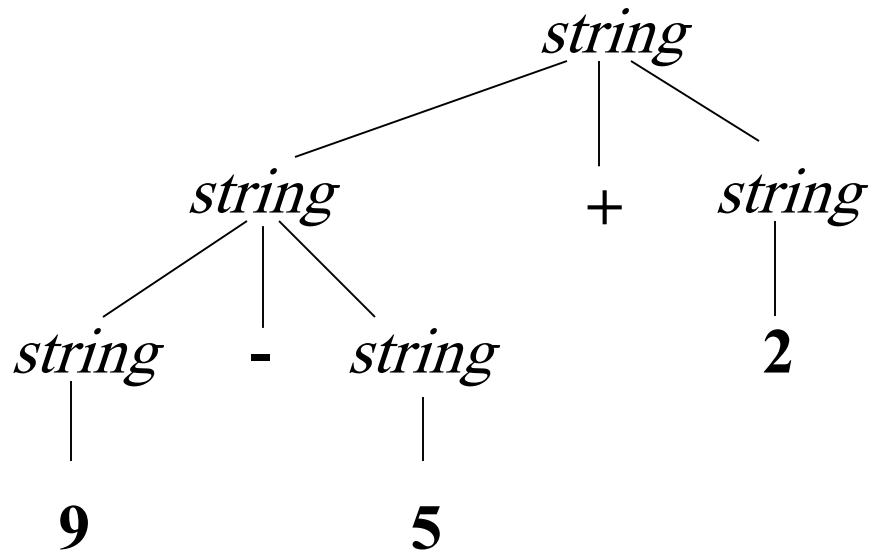
$$G = \langle \{string\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, string \rangle$$

with production  $P =$

$$string \rightarrow string + string \mid string - string \mid 0 \mid 1 \mid \dots \mid 9$$

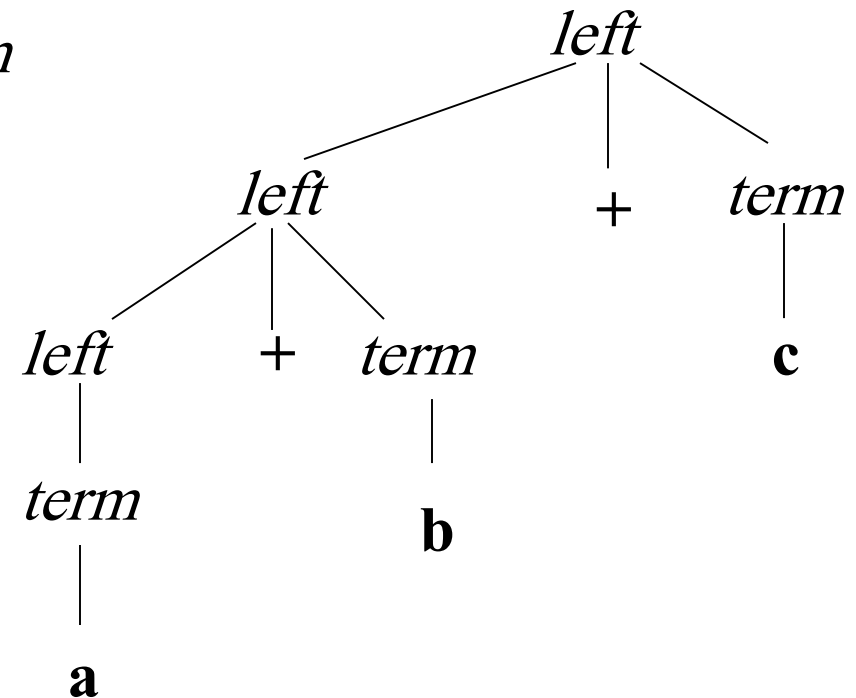
This grammar is **ambiguous**, because more than one parse tree represents the string **9-5+2**

# Ambiguity



# Associativity of Operators

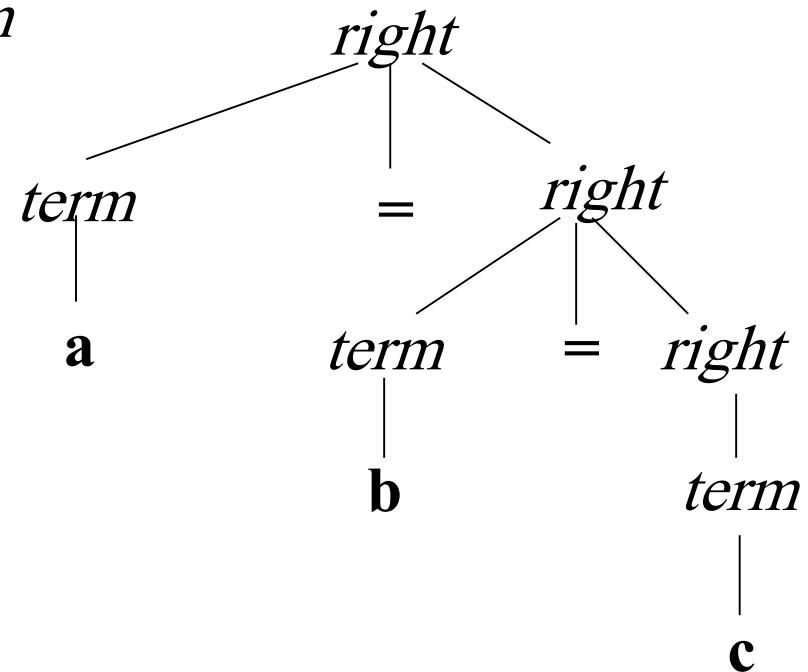
Left-associative operators have left-recursive productions

$$left \rightarrow left + term \mid term$$
$$term \rightarrow a \mid b \mid c$$


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

# Associativity of Operators

Right-associative operators have right-recursive productions

$$right \rightarrow term = right \mid term$$
$$term \rightarrow a \mid b \mid c$$


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



# Precedence of Operators



Table of precedence and associativity of operators:

operators (left-associative):  $+$  ,  $-$       $::$  *expr*

Operators (left-associative):  $*$  ,  $/$       $::$  *term*

basic units (digits and parenthesized expressions)  $::$  *factor*

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

Operators with higher precedence “bind more tightly”

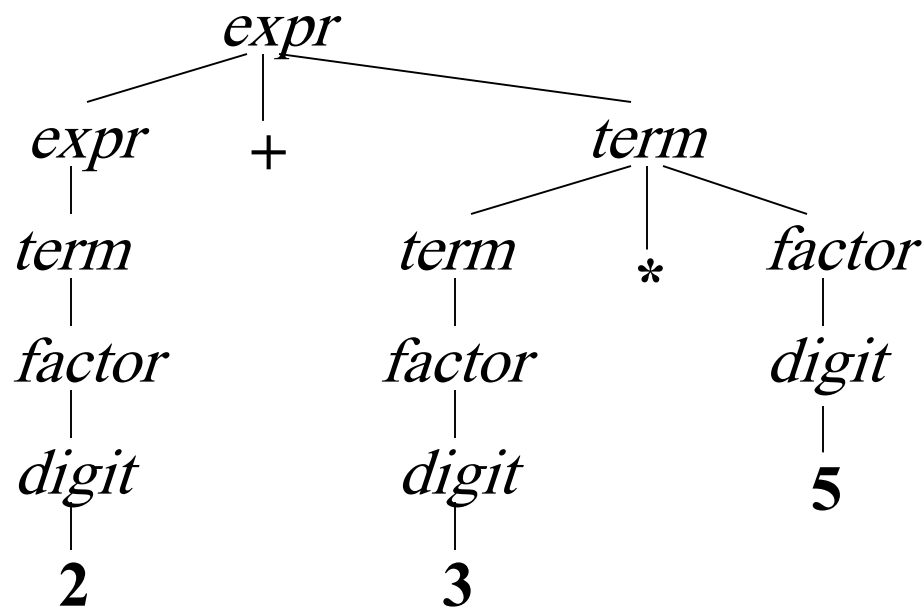
# Precedence of Operators

$expr \rightarrow expr + term \mid expr - term \mid term$

$term \rightarrow term * factor \mid term / factor \mid factor$

$factor \rightarrow digit \mid ( expr )$

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



# Syntax of Statements

Keywords help to recognize statements  
(except assignments and procedure calls):

$$\begin{aligned} stmt \rightarrow & \text{ id = expr } \\ & | \text{ if ( expr ) stmt } \\ & | \text{ if ( expr ) stmt else stmt } \\ & | \text{ while ( expr ) stmt } \\ & | \text{ do stmt while ( expr ) ; } \\ & | \{ stmts \} \end{aligned}$$
$$\begin{aligned} stmts \rightarrow & stmts stmt \\ & | \epsilon \end{aligned}$$

# Parsing



- Parsing: *process of determining if a string of tokens can be generated by a grammar*
- For any CFG, 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 source code
- *Top-down parsing* constructs a parse tree from root to leaves (efficient parsers can be constructed more easily by hand)
- *Bottom-up parsing* constructs a parse tree from leaves to root (can handle a larger class of grammars)

# Top-down Construction of Parse Tree



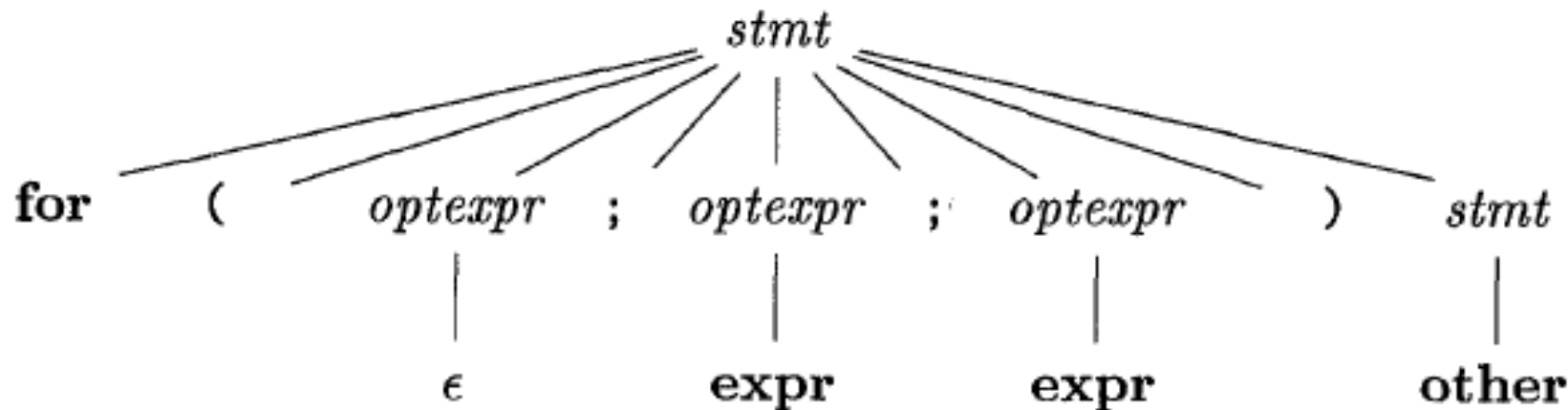
- Starting with the root, labeled with the starting nonterminal, and repeatedly performing the following two steps:
  1. At node  $N$ , labeled with nonterminal  $A$ , select one of the productions for  $A$  and construct children at  $N$  for the symbols in the production body
    - Selecting the appropriate productions for  $A$  depends on the current terminal being scanned in the input, the *lookahead* symbol
  2. Find the next node at which a subtree is to be constructed, typically the leftmost unexpanded nonterminal of the tree

# Top-down Construction of Parse Tree

$stmt \rightarrow \mathbf{expr} ;$   
 $\quad / \mathbf{if} ( \mathbf{expr} ) \mathbf{stmt}$   
 $\quad / \mathbf{for} ( \mathbf{optexpr} ; \mathbf{optexpr} ; \mathbf{optexpr} ) \mathbf{stmt}$   
 $\quad / \mathbf{other}$

$\mathbf{optexpr} \rightarrow \mathbf{expr}$   
 $\quad / \epsilon$

Parse tree of: **for ( ; expr ; expr ) other**



# Top-down Parsing



- For some grammars, the construction steps of *parse tree* can be implemented during a single **left-to-right** scan of the input string
- In general, the selection of a production for a nonterminal may involve trial and error (**try a production and backtrack to try another production if the first is found to be unsuitable**)
  - A production is unsuitable if, after using the production, the tree cannot be completed to match the input

# 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*

$A \rightarrow XbZ$	<code>void A()</code>	<code>void X()</code>	<code>void Z()</code>
	<code>{ ... }</code>	<code>{ ... }</code>	<code>{ ... }</code>

- When a nonterminal has multiple productions, each production is implemented in a branch of a selection statement based on input *lookahead* information

$A \rightarrow bXy$	<code>void A()</code>
$  cZ$	<code>{ if (<i>lookahead</i>=='b')</code>
	<code>...</code>
	<code>elseif (<i>lookahead</i>=='c')</code>
	<code>... }</code>



# Predictive Parsing



- Is a special form of **recursive descent** parsing where one **lookahead** token is sufficient to unambiguously determine the parse operations
- The sequence of procedure calls, during the analysis of an input string, implicitly defines a **parse tree** for the input

*stmt*  $\rightarrow$  **for** ( *optexpr* ; *optexpr* ; *optexpr* ) *stmt*

```
match(for); match('(');  
optexpr(); match(';'); optexpr(); match(';'); optexpr();  
match(')'); stmt();
```

# Example Predictive Parsing (CFG)

$stmt \rightarrow \mathbf{expr};$   
 $\quad / \mathbf{if ( expr ) stmt}$   
 $\quad / \mathbf{for ( optexpr; optexpr; optexpr ) stmt}$   
 $\quad / \mathbf{other}$

$optexpr \rightarrow \mathbf{expr}$   
 $\quad / \epsilon$

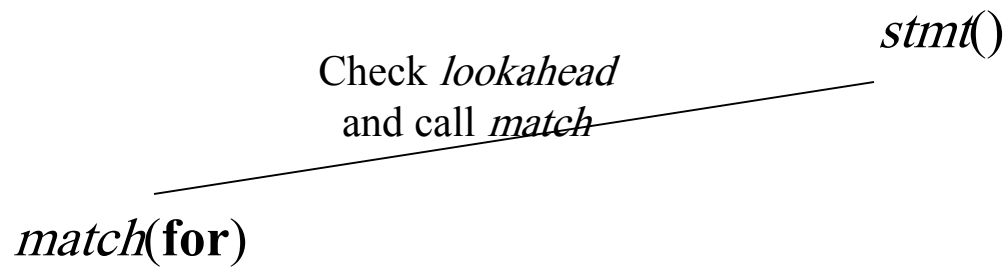
```
void match(terminal t) {  
    if ( lookahead == t ) lookahead = nextTerminal;  
    else report("syntax error");  
}
```

# Example Predictive Parsing (Code)

```
void stmt() {
    switch ( lookahead ) {
    case expr:
        match(expr); match(';'); break;
    case if:
        match(if); 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 optexpr() {
    if ( lookahead == expr ) match(expr);
}
```

# Example Predictive Parser (Execution Step 1)



Input:

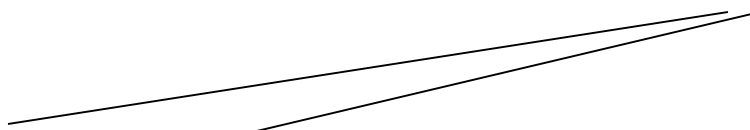
**for ( ; expr ; expr ) other**

↑  
*lookahead*

# Example Predictive Parser (Execution Step 2)

*stmt()*

*match(for) match('(')*



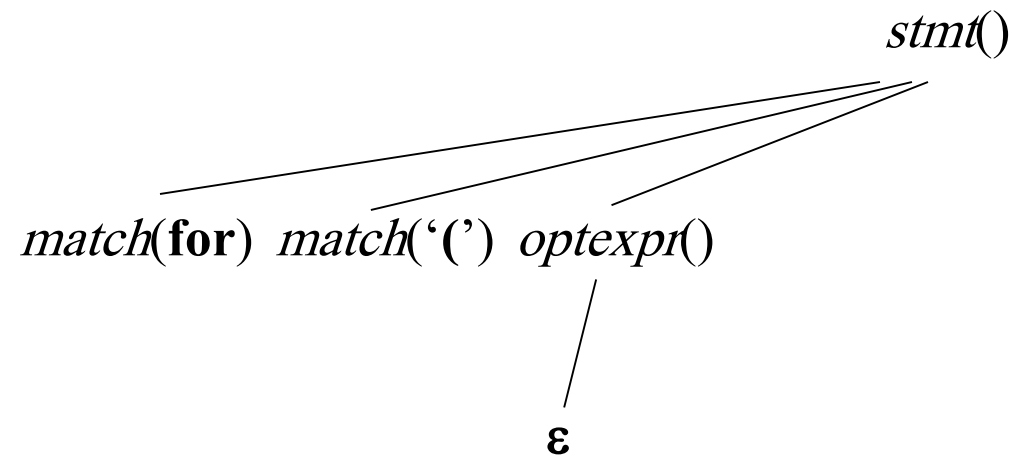
Input:

**for ( ; expr ; expr ) other**



*lookahead*

# Example Predictive Parser (Execution Step 3)



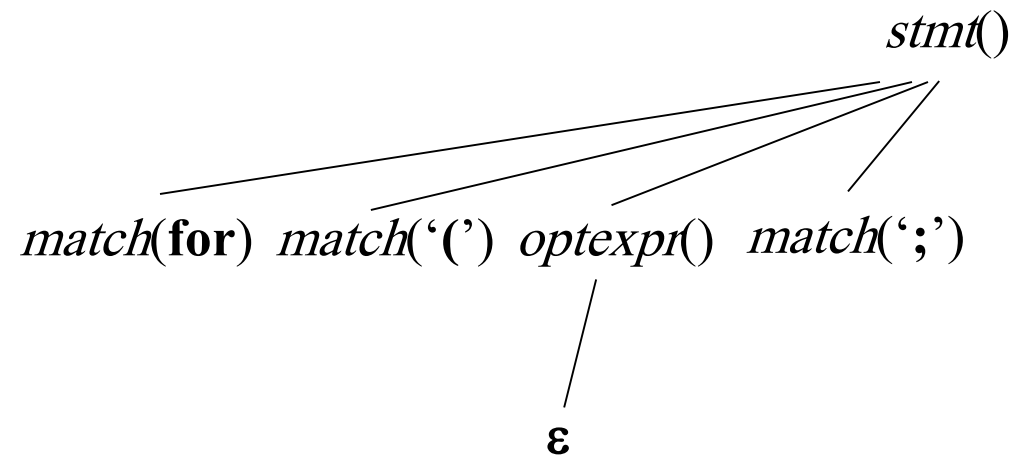
Input:

**for ( ; expr ; expr ) other**



*lookahead*

# Example Predictive Parser (Execution Step 4)



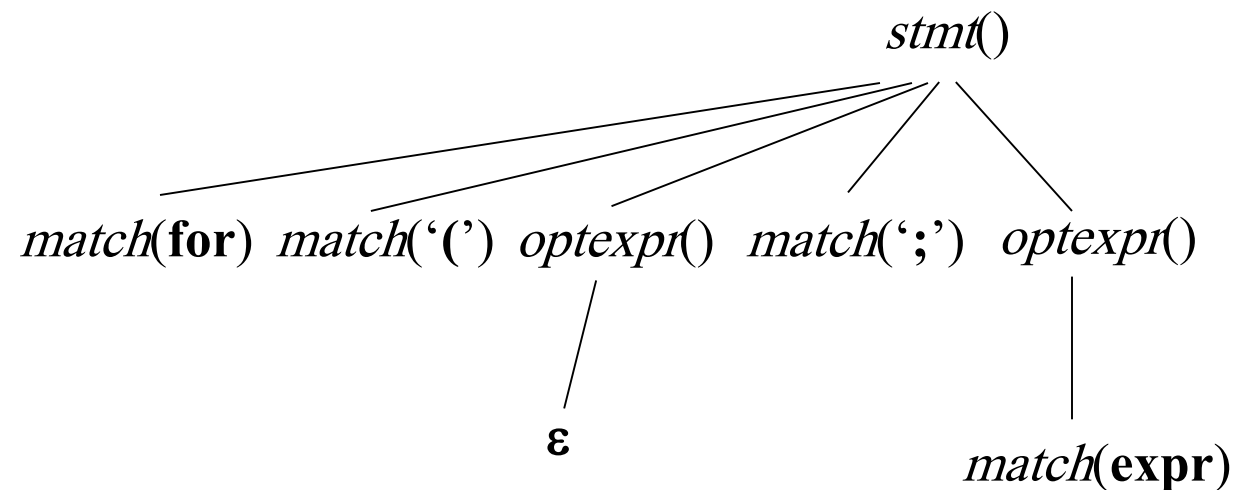
Input:

**for ( ; expr ; expr ) other**



*lookahead*

# Example Predictive Parser (Execution Step 5)



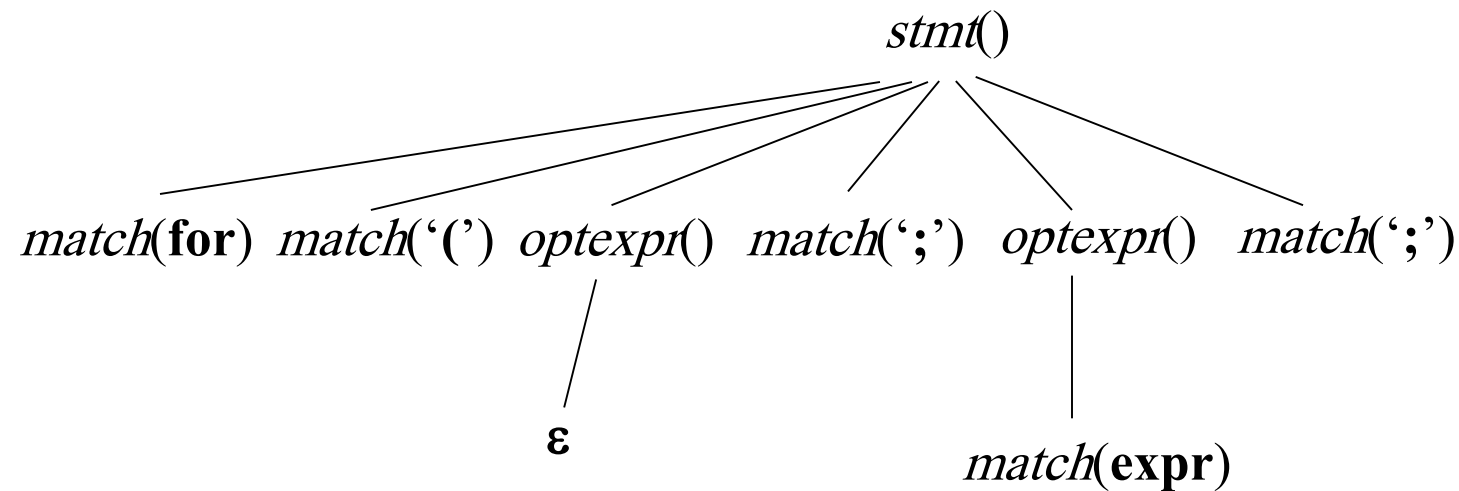
Input:

**for ( ; expr ; expr ) other**

↑  
*lookahead*



# Example Predictive Parser (Execution Step 6)



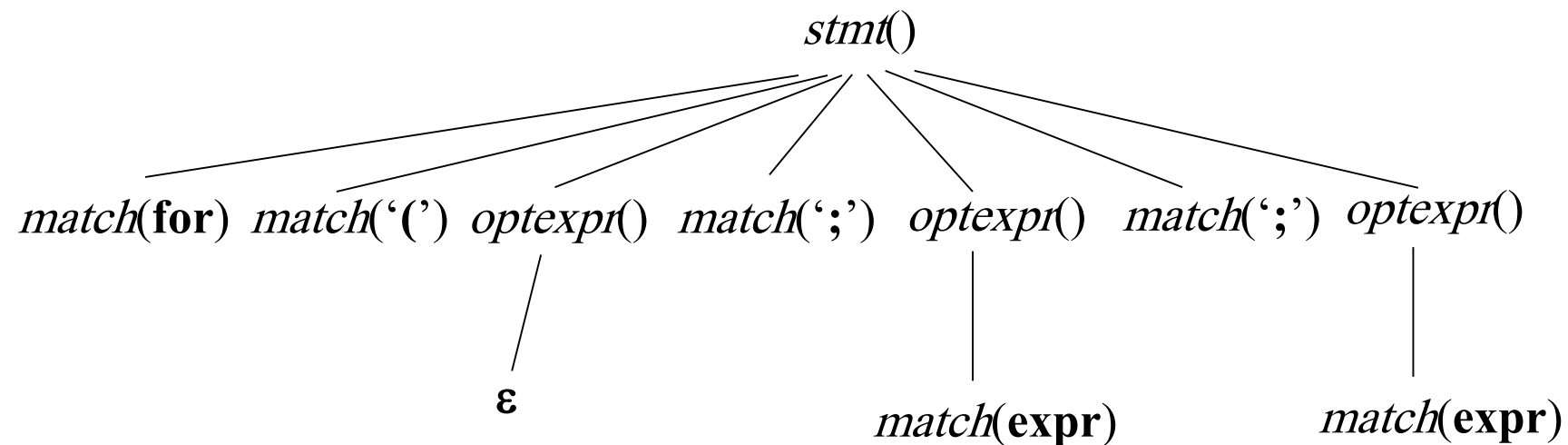
Input:

**for ( ; expr ; expr ) other**

*lookahead*



# Example Predictive Parser (Execution Step 7)



Input:

**for ( ; expr ; expr ) other**



*lookahead*

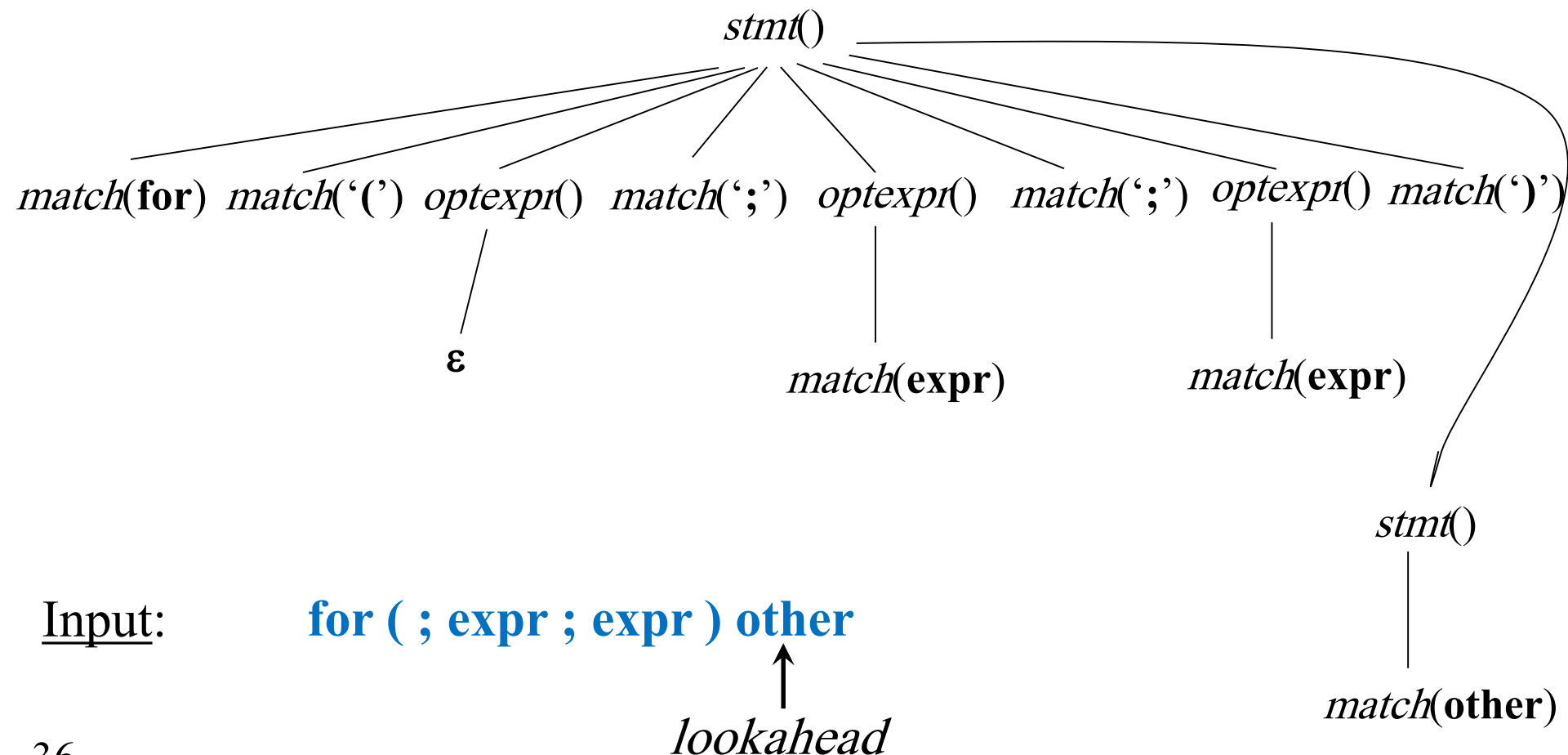
mansoori@shirazu.ac.ir



## for ( ; expr ; expr ) other



# Example Predictive Parser (Execution Step 9)



# FIRST



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

$\text{stmt} \rightarrow \mathbf{expr} ;$   
           $/ \mathbf{if} ( \mathbf{expr} ) \text{stmt}$   
           $/ \mathbf{for} ( \text{optexpr} ; \text{optexpr} ; \text{optexpr} ) \text{stmt}$   
           $/ \mathbf{other}$

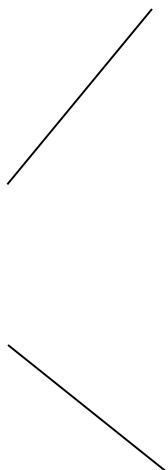
$\text{optexpr} \rightarrow \mathbf{expr}$   
           $/ \epsilon$

$\text{FIRST}(\text{stmt}) = \{ \mathbf{expr}, \mathbf{if}, \mathbf{for}, \mathbf{other} \}$

$\text{FIRST}(\text{optexpr}) = \{ \mathbf{expr}, \epsilon \}$

# Use FIRST to Design Predictive Parser

How to use FIRST set in order to design a predictive parser:

$expr \rightarrow term\ rest$		<pre>void rest() {</pre>
$rest \rightarrow +\ term\ rest$		<pre>  if (lookahead in <u>FIRST(+ term rest)</u> ) {</pre>
$\quad   -\ term\ rest$		<pre>    match('+'); term(); rest() }</pre>
$\quad   \epsilon$		<pre>  else if (lookahead in <u>FIRST(- term rest)</u> ) {</pre>
		<pre>    match('-'); term(); rest() }</pre>
		<pre>  else return</pre>
		<pre>}</pre>

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

$$A \rightarrow \alpha / \beta$$

Then,  $FIRST(\alpha)$  and  $FIRST(\beta)$  must be disjoint for predictive parsing to work

# Left Recursion in CFG

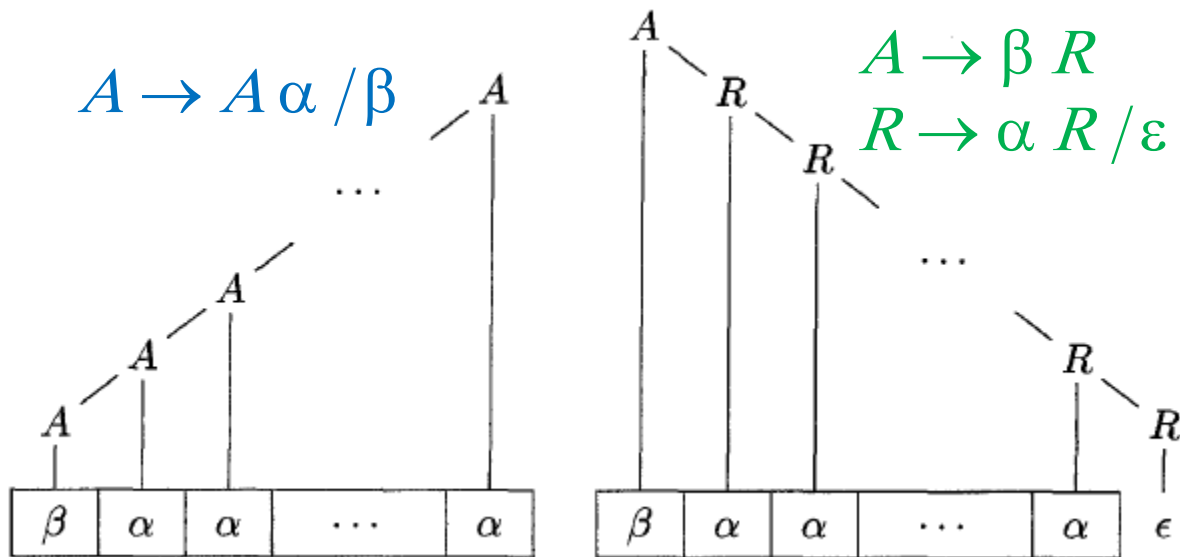
When a production for nonterminal  $A$  starts with a self reference then a predictive parser loops forever

$$A \rightarrow A \alpha / \beta$$

Each **left-recursive** production can be eliminated by rewriting the grammar using **right-recursive** productions

$$\begin{aligned} \text{expr} &\rightarrow \text{expr} + \text{term} \\ &/ \text{expr} - \text{term} \\ &/ \text{term} \end{aligned}$$

$$\begin{aligned} \text{expr} &\rightarrow \text{term rest} \\ \text{rest} &\rightarrow + \text{term rest} \\ &/ - \text{term rest} \\ &/ \epsilon \end{aligned}$$



# Syntax-Directed Translation



- Is done by attaching rules or program fragments to productions in a grammar

*expr*  $\rightarrow$  *expr*<sub>1</sub> + *term*

- In translation of infix expressions into postfix notation:

To translate *expr*:

translate *expr*<sub>1</sub>;

translate *term*;

handle +;



# Syntax-Directed Definition



- Uses a **CFG** to specify the syntactic structure of the language
- AND associates a set of **attributes** with the terminals and nonterminals of the grammar
- AND associates with each production a set of **semantic rules** to compute values of attributes
- A parse tree is traversed and **semantic rules** applied: after the computations are completed, the **attributes** contain the translated form of the input

# Attribute



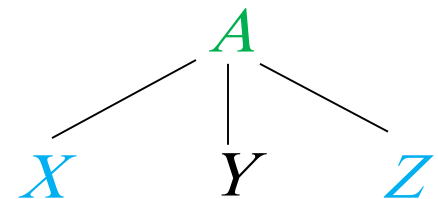
- Is any quantity associated with a programming construct
  - Values
  - Data types of expressions
  - The number of instructions in the generated code
  - The location of the first instruction in the generated code for a construct
- Since nonterminals and terminals represent programming constructs, **attributes** are associated with nonterminals and terminals: *X.x* , *A.a* , *term.value*

# 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

$$A \rightarrow XYZ$$

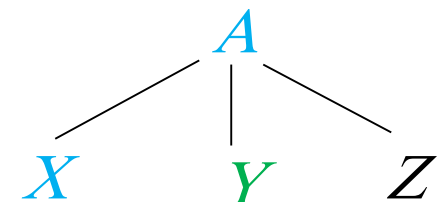
$$A.a = f(X.x, Z.z)$$



- **Inherited**: if its value at a parse-tree node is determined by the parent (by enforcing the parent's semantic rules)

$$A \rightarrow XYZ$$

$$Y.y = g(A.a, X.x)$$



# Semantic Rules



- Rules attached to the **productions** of a grammar
- These rules describe how the **attributes** are computed at the nodes of parse tree
- For a given input string  **$x$** , construct a parse tree for  **$x$**
- Then, apply the **semantic rules** to evaluate **attributes** at each node in the parse tree

# Syntax-Directed Definition



A syntax-directed definition associates

1. With each grammar symbol, a set of attributes, and
2. With each production, a set of semantic rules for computing the values of the attributes associated with the symbols appearing in the production

# Syntax-Directed Definition Example

Infix to postfix translation:

## Production

$expr \rightarrow expr_1 + term$

$expr \rightarrow expr_1 - term$

$expr \rightarrow term$

$term \rightarrow 0$

$term \rightarrow 1$

...

$term \rightarrow 9$

## Semantic rules

$expr.t = expr_1.t \parallel term.t \parallel '+'$

$expr.t = expr_1.t \parallel term.t \parallel '-'$

$expr.t = term.t$

$term.t = '0'$

$term.t = '1'$

...

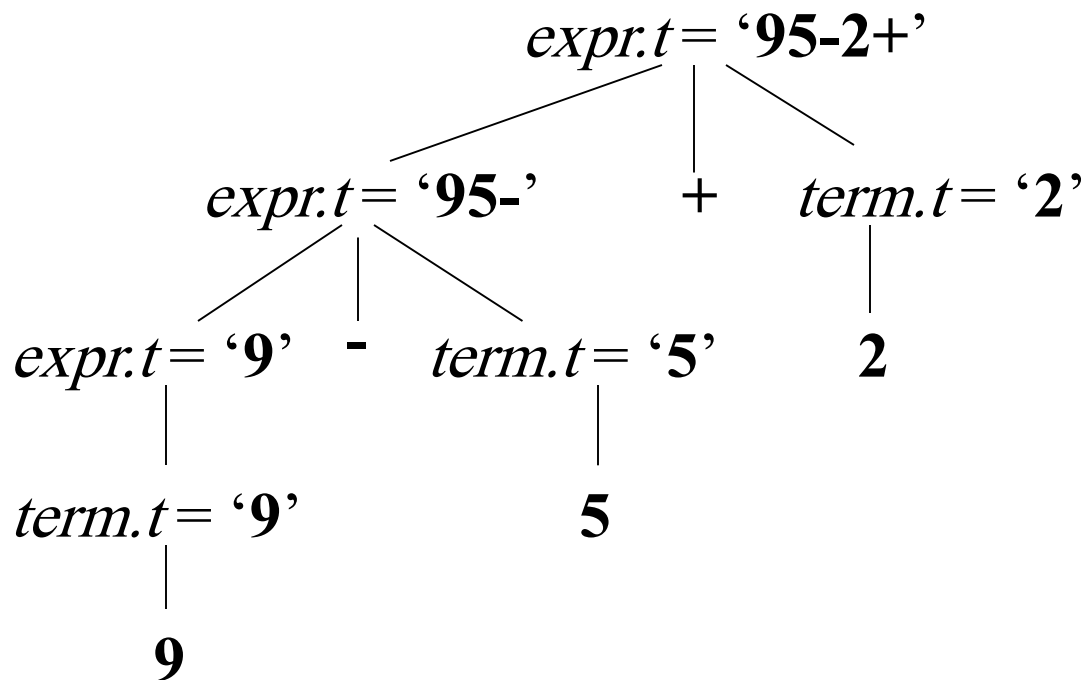
$term.t = '9'$

String concat operator



# Annotated Parse Tree

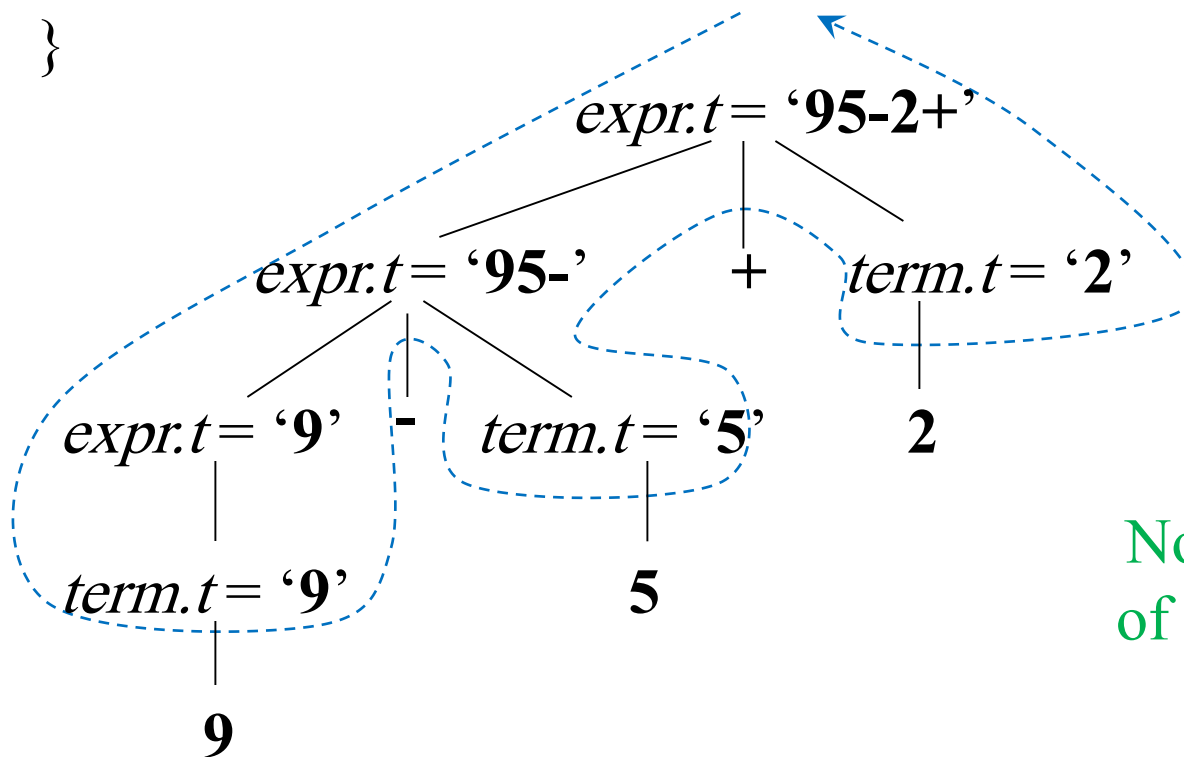
A parse tree showing the **attribute** values at each node is called an **annotated** parse tree



# Depth-first Postorder Traversals

```
procedure visit(node N)  
{  
  for ( each child C of N, from left to right)  
    visit(C);  
  evaluate semantic rules at node N;  
}
```

Bottom-up traversal



Note: all attributes are of the synthesized type



# Translation Schemes

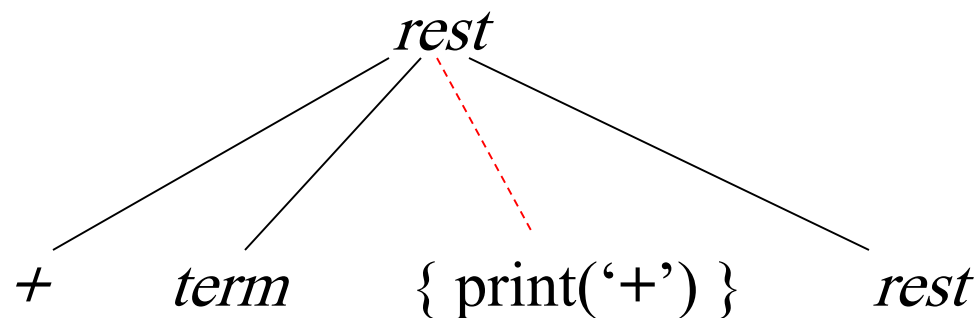


- In **syntax-directed definition**,
  - a translation is built up by evaluating the **attributes** of the nodes in the parse tree
- In **translation scheme**,
  - the same translation is produced incrementally, by executing program fragments (**semantic actions**)
  - is like a syntax-directed definition, except that the **order** of evaluation of the semantic rules is explicitly specified

# Translation Schemes

- A **translation scheme** is a CFG embedded with **semantic actions**

$$rest \rightarrow + term \underbrace{\{ \text{print}('+') \}}_{\substack{\text{Embedded} \\ \text{semantic action}}} rest_1$$



# Example Translation Scheme

$expr \rightarrow expr_1 + term \quad \{ \text{print('+' )} \}$

$expr \rightarrow expr_1 - term \quad \{ \text{print('-' )} \}$

$expr \rightarrow term$

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

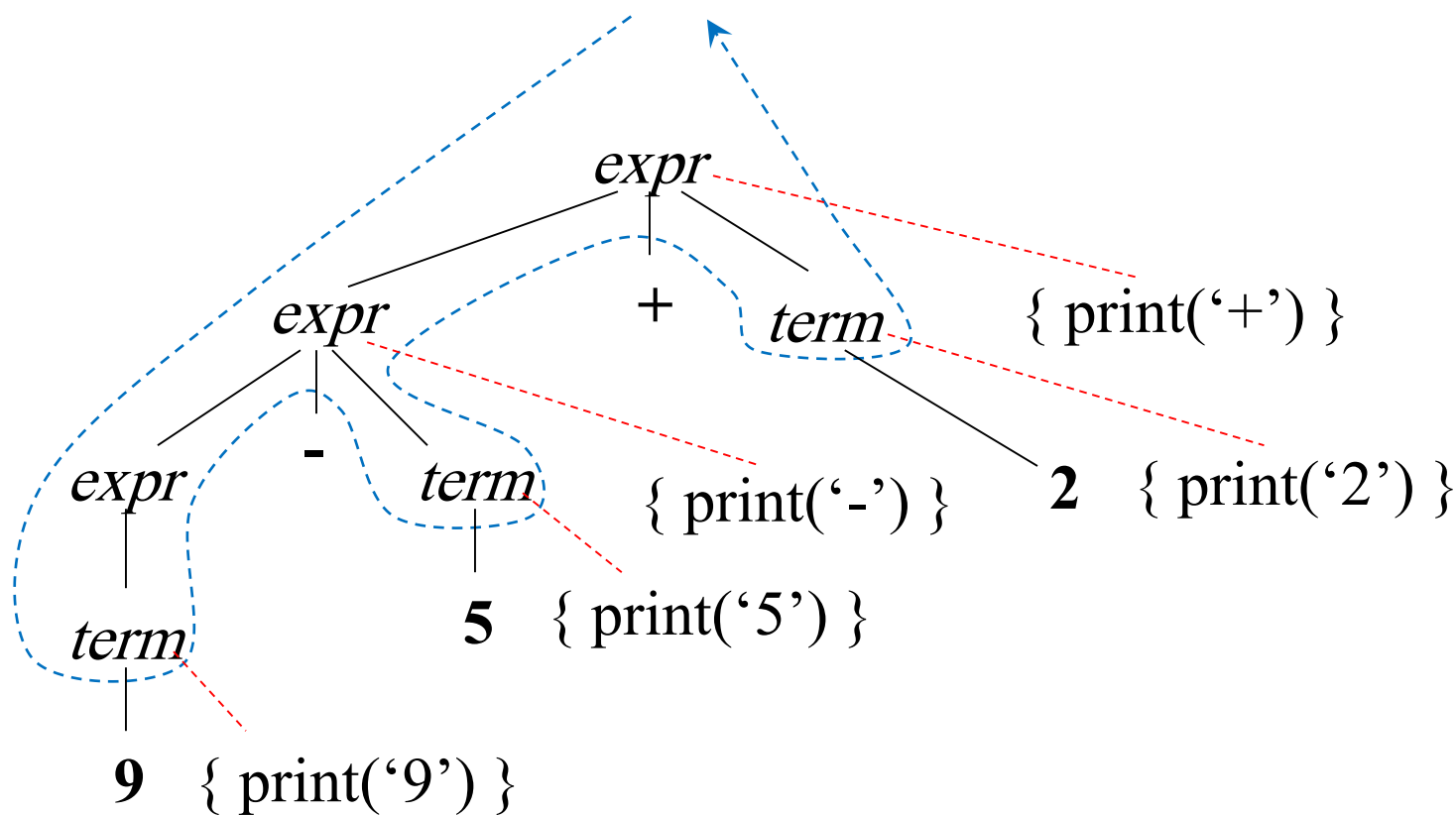
$term \rightarrow 1 \quad \{ \text{print('1' )} \}$

...

...

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

# Example Translation Scheme



Traversing annotated parse tree translates **9-5+2** into postfix **95-2+**

# From Predictive Parser to Translation Scheme



1. Construct a predictive parser, ignoring the actions in productions
2. Copy the actions from the translation scheme into the parser
  - If an action appears after grammar symbol  $X$  in production  $p$ , then it is copied after the implementation of  $X$  in the code for  $p$
  - Otherwise, if it appears at the beginning of the production, then it is copied just before the code for the production body

# Left Recursion in Translation Scheme



When eliminating **left-recursive** from translation scheme, the **semantic actions** are handled as terminals

$$A \rightarrow A \alpha \{a\} / \beta$$

$$A \rightarrow \beta R$$

$$R \rightarrow \alpha \{a\} R / \varepsilon$$

$$\begin{aligned} \text{expr} &\rightarrow \text{expr} + \text{term} \{ \text{print}('+') \} \\ & / \text{term} \end{aligned}$$

$$\text{expr} \rightarrow \text{term rest}$$

$$\begin{aligned} \text{rest} &\rightarrow + \text{term} \{ \text{print}('+') \} \text{rest} \\ & | \varepsilon \end{aligned}$$

# 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:

$expr \rightarrow term \text{ rest}$

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

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

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

...

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

# Code of Translator

$expr \rightarrow term\ rest$

```
void expr() {  
    term(); rest();  
}
```

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

```
void rest() {  
    if ( lookahead == '+' ) {  
        match('+'); term(); print('+'); rest();  
    }  
    else if ( lookahead == '-' ) {  
        match('-'); term(); print('-'); rest();  
    }  
    else { } /* do nothing with the input */ ;  
}
```

$term \rightarrow 0\ \{ \text{print('0' )} \}$   
 $term \rightarrow 1\ \{ \text{print('1' )} \}$   
...  
 $term \rightarrow 9\ \{ \text{print('9' )} \}$

```
void term() {  
    if ( lookahead is a digit ) {  
        t = lookahead; match(lookahead); print(t);  
    }  
    else report("syntax error");  
}
```



# Code of Translator



*expr* → *term rest*

*rest* → + *term* { print('+') } *rest*  
 | - *term* { print('-') } *rest*  
 | ε

*term* → 0 { print('0') }  
*term* → 1 { print('1') }  
 ...  
*term* → 9 { print('9') }

```
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);
}
```

# Complete Code in Java

```
import java.io.*;
class Parser {
    static int lookahead;

    public Parser() throws IOException {
        lookahead = System.in.read();
    }

    void expr() throws IOException {
        term();
        while(true) {
            if( lookahead == '+' ) {
                match('+'); term(); System.out.write('+');
            }
            else if( lookahead == '-' ) {
                match('-'); term(); System.out.write('-');
            }
            else return;
        }
    }

    void term() throws IOException {
        if( Character.isDigit((char)lookahead) ) {
            System.out.write((char)lookahead); match(lookahead);
        }
        else throw new Error("syntax error");
    }

    void match(int t) throws IOException {
        if( lookahead == t ) lookahead = System.in.read();
        else throw new Error("syntax error");
    }
}

public class Postfix {
    public static void main(String[] args) throws IOException {
        Parser parse = new Parser();
        parse.expr(); System.out.write('\n');
    }
}
```

# Tokens **num** and **id** in Expressions

<i>expr</i>	→	<i>expr</i> + <i>term</i>	{ print('+' ) }
		<i>expr</i> - <i>term</i>	{ print('-' ) }
		<i>term</i>	

<i>term</i>	→	<i>term</i> * <i>factor</i>	{ print('*' ) }
		<i>term</i> / <i>factor</i>	{ print('/' ) }
		<i>factor</i>	

<i>factor</i>	→	( <i>expr</i> )	
		<b>num</b>	{ print( <b>num.value</b> ) }
		<b>id</b>	{ print( <b>id.lexeme</b> ) }

# Adding a Lexical Analyzer



- Typical tasks of the lexical analyzer:
  - Read characters from input and group into **tokens**
  - Remove **white space** and **comments**
  - Encode **numbers** and **constants** as tokens
  - Recognize **keywords**
  - Recognize **identifiers** and store identifier names in a global **symbol table**
    - A sequence of input characters that comprises a single token is called a **lexeme**

# Removal of White Space and Comments

1. Lexical analyzer eliminates
  - White space: blank, tab, newline
  - Comments

```
for ( ; ; peek = next input character ) {  
    if ( peek is a blank or a tab ) do nothing;  
    else if ( peek is a newline ) line = line+1;  
    else break;  
}
```

2. Modifying the grammar to incorporate white space

# Reading Ahead



- Lexical analyzer may need to read ahead some characters before it can decide on the token
  - In C or Java, read ahead after ‘>’
    - If next character = , token: ‘greater than or equal to’
    - Otherwise, token: ‘greater than’
- Lexical analyzer needs an input **buffer** from which it can read and push back characters

# Recognizing Constants

- Lexical analyzer
  - Collects the characters into **integers** and computes their **numerical value**
  - Passes to the parser a **token** consisting of the terminal **num** along with an integer-valued **attribute** computed from the digits: **<num,25>**

```
if ( peek holds a digit ) {  
    v = 0;  
    do {  
        v = v * 10 + integer value of digit peek;  
        peek = next input character;  
    } while ( peek holds a digit );  
    return token <num, v>;  
}
```

# Recognizing Keywords and Identifiers



- **Keywords**: fixed character strings used to identify constructs: **for**, **do**, **if**, ...
- **Identifiers (id)**: character strings used to name **variables, arrays, functions**, ...
- Grammars routinely treat identifiers as terminals to **simplify** the **parser**

count = count + increment;                  **id** = **id** + **id** ;

- The token for **id** has an attribute that holds the **lexeme**  
 $\langle \mathbf{id}, \text{"count"} \rangle \Rightarrow \langle \mathbf{id}, \text{"count"} \rangle \langle + \rangle \langle \mathbf{id}, \text{"increment"} \rangle \langle ; \rangle$



# Recognizing Keywords and Identifiers



- **Keywords** generally satisfy the rules for forming **identifiers**, so lexical analyzer should decide when
  - A lexeme forms a **keyword**
  - A lexeme forms an **identifier**
- If **keywords** are **reserved**
  - They **cannot** be used as identifiers
  - Then, a character string forms an identifier only if it **is not** a keyword

# Using Symbol Table

- Single representation
  - Symbol table can **insulate** the compiler from representation of strings, so compiler phases can work with **references** to the table
  - References can be manipulated more **efficiently** than the strings themselves
- Reserved words
  - Reserved words are implemented by **initializing** the symbol table with them
  - When the lexical analyzer reads a **lexeme** of identifier, it first checks whether the **lexeme** is in the table
    - If so, it returns the **token** from the table
    - Otherwise, it returns a **token** of **id**

# Using Symbol Table

- **Java**: a symbol table is implemented as a hash table:

*Hashtable words = new Hashtable();*

- Sets up words as a default **hash table** that maps keys (**lexemes**) to values (**tokens**)

```
if ( peek holds a letter ) {  
    collect letters or digits into a buffer b;  
    s = string formed from the characters in b;  
    w = token returned by words.get(s);  
    if ( w is not null ) return w;  
    else {  
        Enter the key-value pair (s, <id, s>) into words  
        return token <id, s>;  
    }  
}
```

# The Lexical Analyzer

$y = 31 + 28 * x$

Lexical analyzer  
**lexan()**

**<id, "y">** **<=, >** **<num, 31>** **<+, >** **<num, 28>** **<\*, >** **<id, "x">**

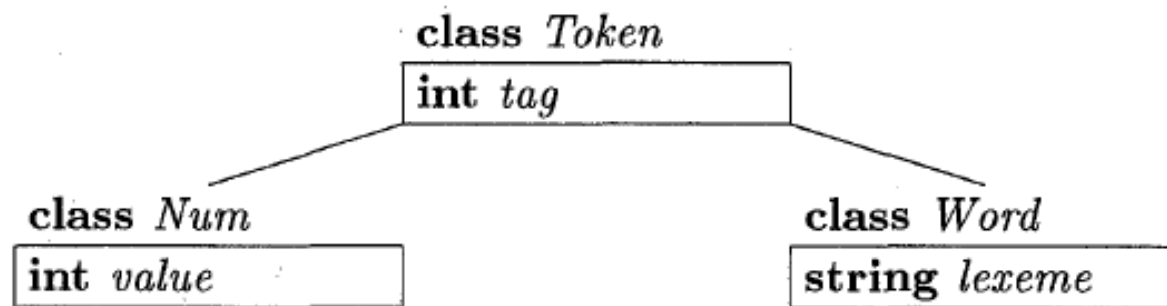
token

attribute

Syntax analyzer  
**parse()**

# Java Package for Lexical Analysis

- The package `lexer` has
  - classes for tokens
  - a class `lexer` containing function `scan`



- Class `Token` has a field `tag` that is used for parsing decisions
  - Subclass `Num` adds a field `value` for an **integer value**
  - Subclass `Word` adds a field `lexeme` for **reserved words** and **identifiers**

# Java Package for Lexical Analysis



```
1) package lexer;                                // File Token.java
2) public class Token {
3)     public final int tag;
4)     public Token(int t) { tag = t; }
5) }
```

- The constructor Token is used to create **token** objects

```
new Token ( ' + ' )
```

- Where the pseudocode has terminals like **num** and **id**, the Java code uses **integer constants**

# Java Package for Lexical Analysis

```
1) package lexer;                                // File Tag.java
2) public class Tag {
3)     public final static int
4)         NUM = 256, ID = 257, TRUE = 258, FALSE = 259;
5) }
```

```
1) package lexer;                                // File Num.java
2) public class Num extends Token {
3)     public final int value;
4)     public Num(int v) { super(Tag.NUM); value = v; }
5) }
```

```
1) package lexer;                                // File Word.java
2) public class Word extends Token {
3)     public final String lexeme;
4)     public Word(int t, String s) {
5)         super(t); lexeme = new String(s);
6)     }
7) }
```

# Java Package for Lexical Analysis

```
1) package lexer;                                // File Lexer.java
2) import java.io.*; import java.util.*;
3) public class Lexer {
4)     public int line = 1;
5)     private char peek = ' ';
6)     private Hashtable words = new Hashtable();
7)     void reserve(Word t) { words.put(t.lexeme, t); }
8)     public Lexer() {
9)         reserve( new Word(Tag.TRUE,  "true")  );
10)        reserve( new Word(Tag.FALSE, "false") );
11)    }
12)    public Token scan() throws IOException {
13)        for( ; ; peek = (char)System.in.read() ) {
14)            if( peek == ' ' || peek == '\t' ) continue;
15)            else if( peek == '\n' ) line = line + 1;
16)            else break;
17)        }
```



# Java Package for Lexical Analysis

```
18)         if( Character.isDigit(peek) ) {
19)             int v = 0;
20)             do {
21)                 v = 10*v + Character.digit(peek, 10);
22)                 peek = (char)System.in.read();
23)             } while( Character.isDigit(peek) );
24)             return new Num(v);
25)         }
26)         if( Character.isLetter(peek) ) {
27)             StringBuffer b = new StringBuffer();
28)             do {
29)                 b.append(peek);
30)                 peek = (char)System.in.read();
31)             } while( Character.isLetterOrDigit(peek) );
32)             String s = b.toString();
33)             Word w = (Word)words.get(s);
34)             if( w != null ) return w;
35)             w = new Word(Tag.ID, s);
36)             words.put(s, w);
37)             return w;
38)         }
39)         Token t = new Token(peek);
40)         peek = ' ';
41)         return t;
42)     }
43) }
```

# Symbol Tables



- The information in the symbol table
  - is **collected** by the **analysis** phases
  - and **used** by the **synthesis** phases to generate the **target code**
- Entries in symbol table contain information about an **identifier**
  - its **lexeme**
  - its **type**
  - its position in **storage**
  - ...
- Symbol tables need to support **multiple** declarations for an **identifier**

# Symbol Tables



- To implement **scopes**, a **separate symbol table** for each scope is set up
  - a **program block** with declarations
  - a **class** with an entry for each field and method
- Example of just key constructs that touch symbol tables: *blocks, declarations, factors*
- A program: *blocks* with optional *declarations* and *statements* of single identifiers (represent use of identifiers)

**Input:** `{ int x; char y; { bool y; x; y; } x; y; }`

**Output:** `{ { x:int; y:bool; } x:int; y:char; }`

# Symbol Table per Scope

- If blocks can be **nested**, several declarations of the same identifier can appear within a single block

*block*  $\rightarrow$  '{' *decls stmts* '}'

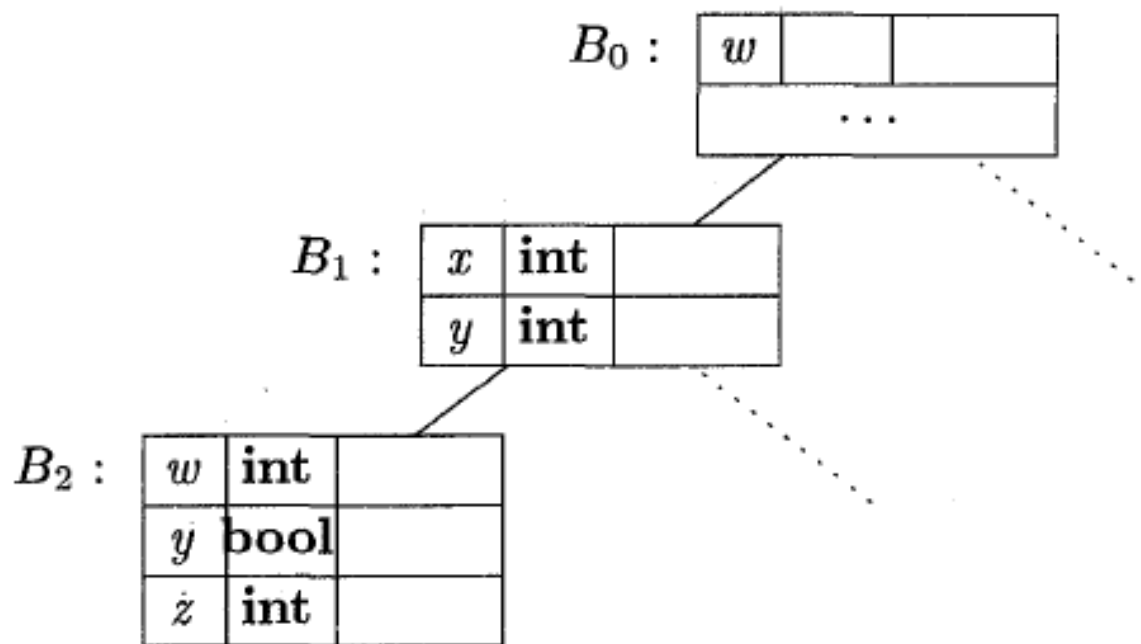
*stmts* can generate a block

- Most-closely nested** rule for blocks: identifier *x* is in the scope of the most-closely nested declaration of *x*

```
1)  {   int x1; int y1;  
2)    {   int w2; bool y2; int z2;  
3)        ... w2 ...; ... x1 ...; ... y2 ...; ... z2 ...;  
4)    }  
5)        ... w0 ...; ... x1 ...; ... y1 ...;  
6) }
```

# Chained Symbol Tables

- The most-closely nested rule for blocks can be implemented by **chaining** symbol tables
  - Table for a **nested** block points to the table for its enclosing block



# Java Imp. of Chained Symbol Tables



```
1) package symbols;                                // File Env.java
2) import java.util.*;
3) public class Env {
4)     private Hashtable table;
5)     protected Env prev;

6)     public Env(Env p) {
7)         table = new Hashtable(); prev = p;
8)     }

9)     public void put(String s, Symbol sym) {
10)         table.put(s, sym);
11)     }

12)     public Symbol get(String s) {
13)         for( Env e = this; e != null; e = e.prev ) {
14)             Symbol found = (Symbol)(e.table.get(s));
15)             if( found != null ) return found;
16)         }
17)         return null;
18)     }
19) }
```

# The Use of Symbol Tables

*program*  $\rightarrow$  *block*

*block*  $\rightarrow$  '{'  
*decls stmts* '}'

*decls*  $\rightarrow$  *decls decl* |  $\epsilon$

*decl*  $\rightarrow$  **type id ;**

*stmts*  $\rightarrow$  *stmts stmt* |  $\epsilon$

*stmt*  $\rightarrow$  *block*  
| *factor ;*

*factor*  $\rightarrow$  **id**

mansoori@shirazu.ac.ir

$$program \rightarrow \{ top = \mathbf{null}; \} block$$

```

block    →  '{ { saved = top; top = new Env(top); print("{ "); }
           decls stmts '}' { top = saved; print("} "); }

```

$$decls \rightarrow decls \ decl \mid \epsilon$$

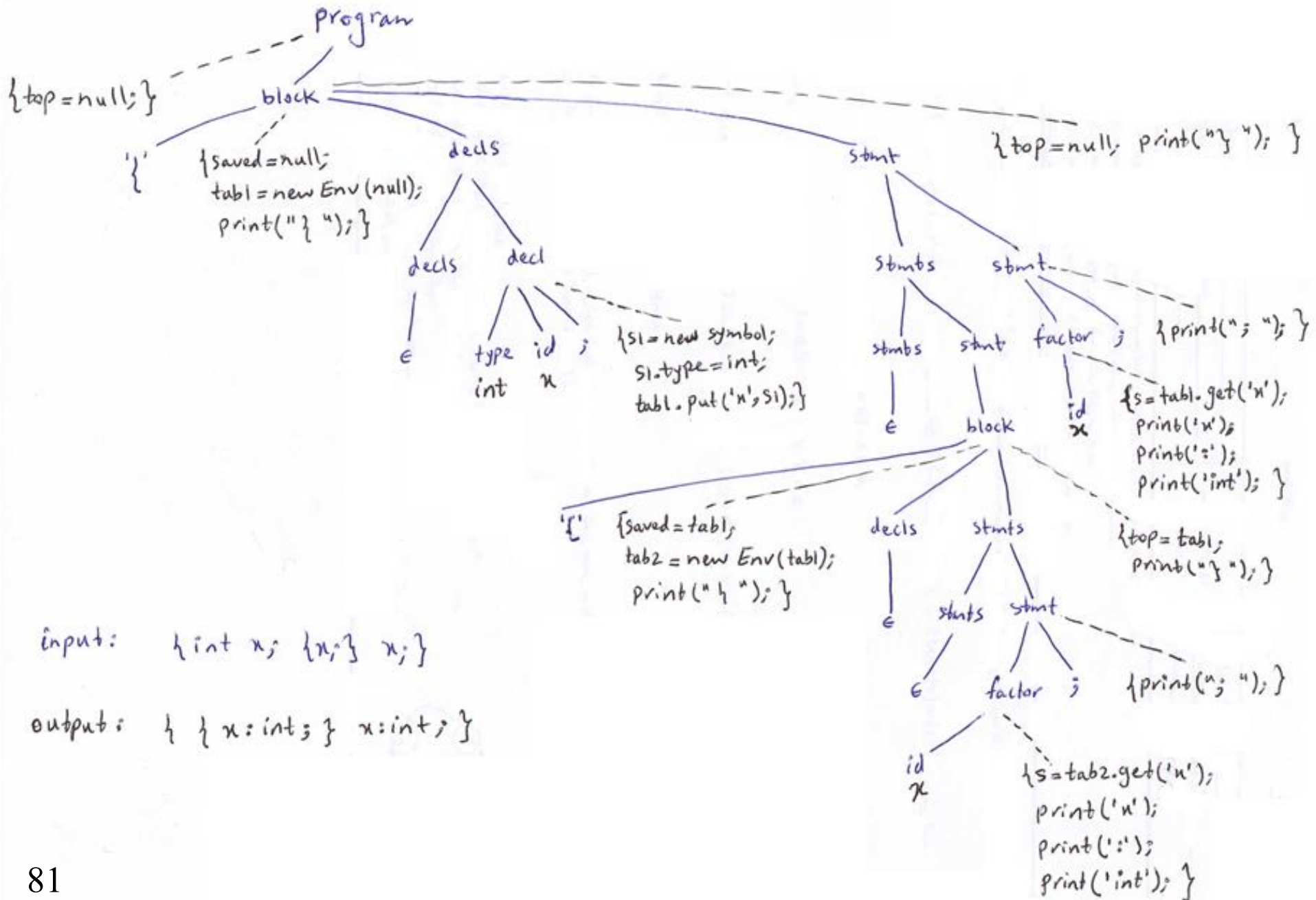
```
decl → type id ; { s = new Symbol; s.type = type.lexeme  
top.put(id.lexeme, s); }
```

$$stmts \rightarrow stmts \; stmt \mid \epsilon$$
$$\begin{array}{lcl} stmt & \rightarrow & block \\ & | & factor ; \quad \{ \text{print}("; "); \} \end{array}$$

```
factor → id { s = top.get(id.lexeme); print(id.lexeme);
              print(":"); print(s.type); }
```



# The Use of Symbol Tables



# Intermediate Code Generation



- The front-end of a compiler constructs an **intermediate representation** of the **source program** from which the back-end generates the **target program**
- Intermediate representations for:
  - **Expressions**
  - **Statements**

# Two Intermediate Representations



- Tree representations: **syntax tree**
  - Syntax-tree nodes represent significant programming **constructs**
  - Information is added to the **nodes** in the form of attributes associated with the nodes
- Linear representations: **three-address code**
  - A **sequence** of elementary program steps
  - Unlike the tree, there is **no hierarchical** structure
  - Three-address code can be **optimized**

# Construction of Syntax Trees for Statements

$program \rightarrow block$

$block \rightarrow \{ stmts \}$

$stmts \rightarrow stmts_1 stmt$   
 $\quad \quad \quad | \epsilon$

$stmt \rightarrow expr ;$   
 $\quad \quad | \text{ if } ( expr ) stmt_1$   
 $\quad \quad | \text{ while } ( expr ) stmt_1$   
 $\quad \quad | \text{ do } stmt_1 \text{ while } ( expr ) ;$   
 $\quad \quad | block$

# Construction of Syntax Trees for Statements

*program*  $\rightarrow$  *block*                      { return *block.n*; }

*block*  $\rightarrow$  '{' *stmts* '}'              { *block.n* = *stmts.n*; }

*stmts*  $\rightarrow$  *stmts*<sub>1</sub> *stmt*              { *stmts.n* = **new** *Seq*(*stmts*<sub>1</sub>.*n*, *stmt.n*); }

          |  $\epsilon$                               { *stmts.n* = **null**; }

*stmt*  $\rightarrow$  *expr* ;                      { *stmt.n* = **new** *Eval*(*expr.n*); }

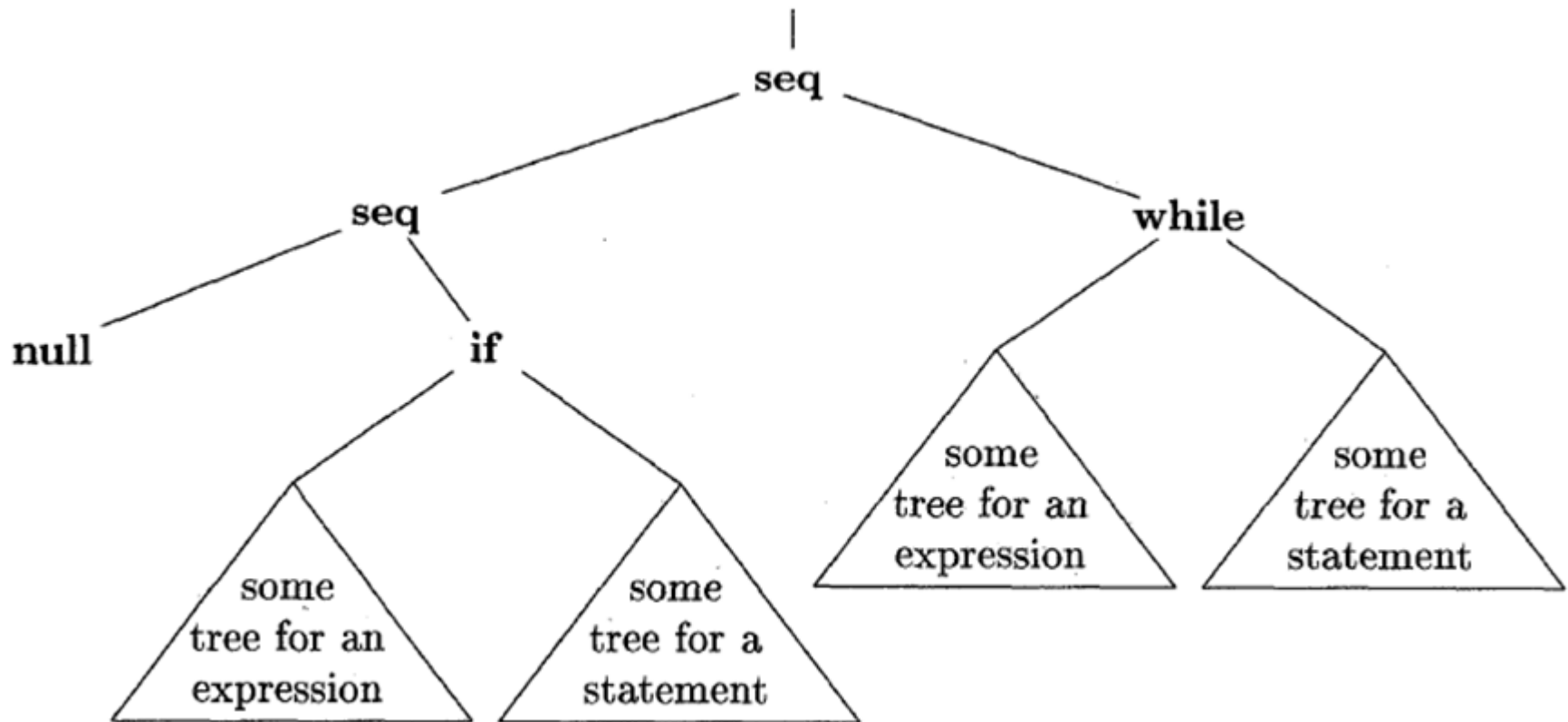
          | **if** ( *expr* ) *stmt*<sub>1</sub>              { *stmt.n* = **new** *If*(*expr.n*, *stmt*<sub>1</sub>.*n*); }

          | **while** ( *expr* ) *stmt*<sub>1</sub>              { *stmt.n* = **new** *While*(*expr.n*, *stmt*<sub>1</sub>.*n*); }

          | **do** *stmt*<sub>1</sub> **while** ( *expr* );              { *stmt.n* = **new** *Do*(*stmt*<sub>1</sub>.*n*, *expr.n*); }

          | *block*                          { *stmt.n* = *block.n*; }

# Construction of Syntax Trees for Statements



# Construction of Syntax Trees for Expressions

$$\begin{array}{l} \text{expr} \rightarrow \text{rel} = \text{expr}_1 \\ \quad | \quad \text{rel} \end{array}$$
$$\begin{array}{l} \text{rel} \rightarrow \text{rel}_1 < \text{add} \\ \quad | \quad \text{rel}_1 \leq \text{add} \\ \quad | \quad \text{add} \end{array}$$
$$\begin{array}{l} \text{add} \rightarrow \text{add}_1 + \text{term} \\ \quad | \quad \text{term} \end{array}$$
$$\begin{array}{l} \text{term} \rightarrow \text{term}_1 * \text{factor} \\ \quad | \quad \text{factor} \end{array}$$
$$\begin{array}{l} \text{factor} \rightarrow ( \text{expr} ) \\ \quad | \quad \mathbf{num} \\ \quad | \quad \mathbf{id} \end{array}$$



# Construction of Syntax Trees for Expressions

$expr \rightarrow rel = expr_1$	$\{ expr.n = \mathbf{new} Assign('=', rel.n, expr_1.n); \}$
$\quad \quad   \quad rel$	$\{ expr.n = rel.n; \}$

$rel \rightarrow rel_1 < add$	$\{ rel.n = \mathbf{new} Rel('<', rel_1.n, add.n); \}$
$\quad \quad   \quad rel_1 \leq add$	$\{ rel.n = \mathbf{new} Rel('<=', rel_1.n, add.n); \}$
$\quad \quad   \quad add$	$\{ rel.n = add.n; \}$

$add \rightarrow add_1 + term$	$\{ add.n = \mathbf{new} Op('+', add_1.n, term.n); \}$
$\quad \quad   \quad term$	$\{ add.n = term.n; \}$

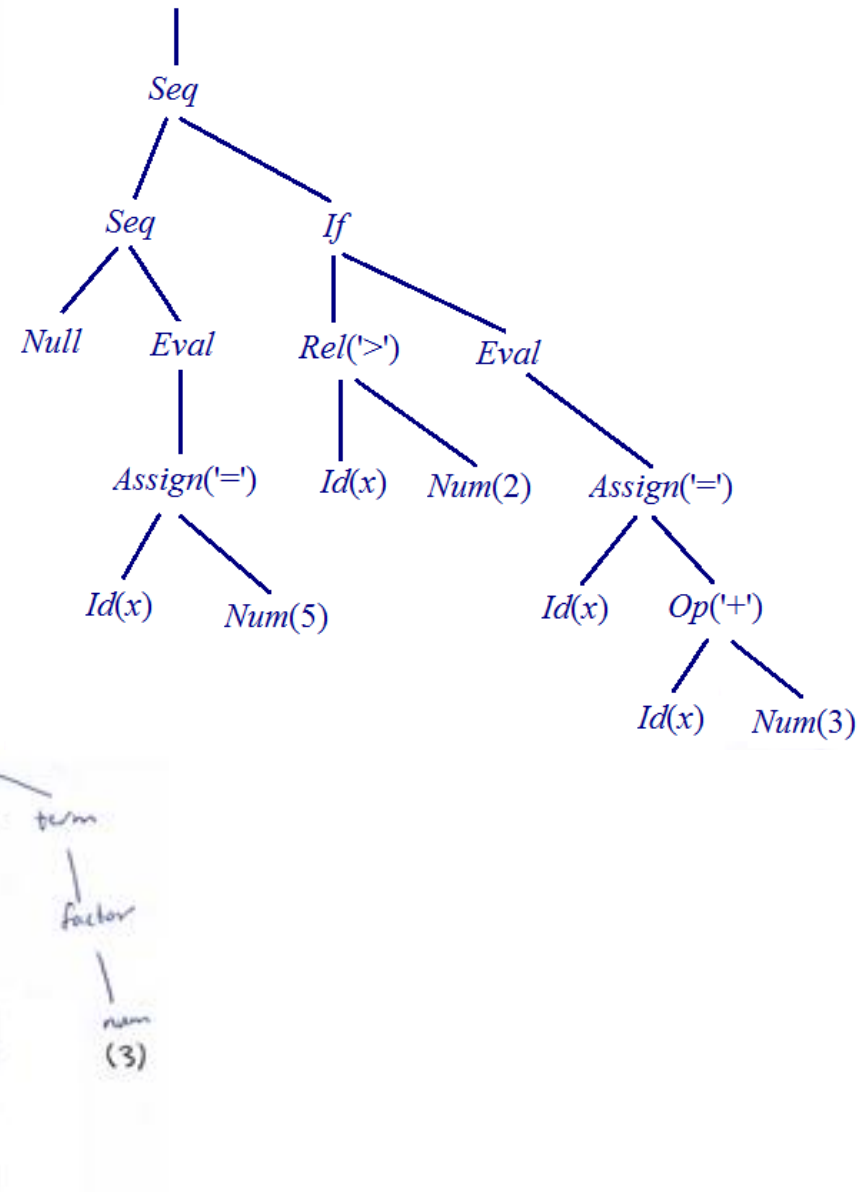
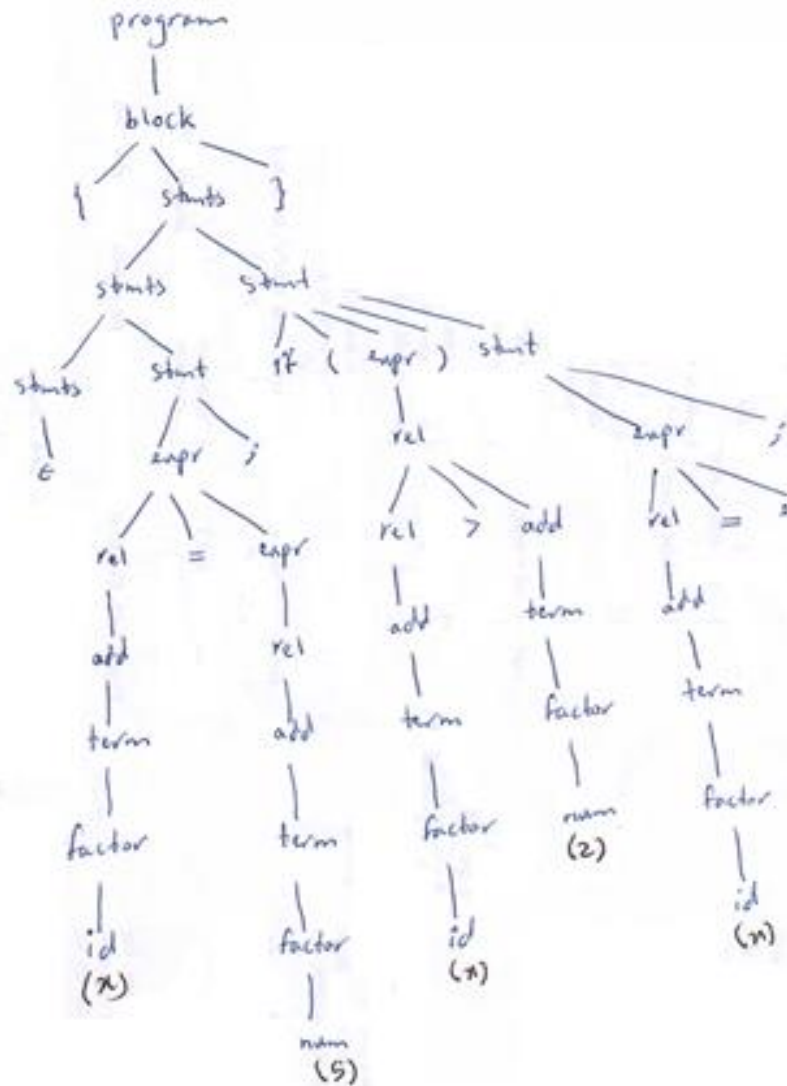
$term \rightarrow term_1 * factor$	$\{ term.n = \mathbf{new} Op('*', term_1.n, factor.n); \}$
$\quad \quad   \quad factor$	$\{ term.n = factor.n; \}$

$factor \rightarrow ( expr )$	$\{ factor.n = expr.n; \}$
$\quad \quad   \quad \mathbf{num}$	$\{ factor.n = \mathbf{new} Num(\mathbf{num.value}); \}$
$\quad \quad   \quad \mathbf{id}$	$\{ factor.n = \mathbf{new} Id(\mathbf{id.lexeme}); \}$



# Construction of Syntax Trees

{ x=5; if (x>2) x=x+3; }



# Static Checking



- Consistency **checks** during compilation
  - To assure that a program can be compiled successfully
  - To catch programming errors before run
- 1. **Syntactic checking**: (not encoded in grammar)
  - An identifier being declared at most once in a scope
  - A break must have an enclosing loop or switch statement
- 2. **L-values** (locations) and **R-values** (values)
  - Left side of an assignment must be an **l-value** (identifier  $i$ , array access  $a[2]$ )
  - A constant has **r-value**, so not on left side of an assignment

# Static Checking



## 3. Type checking:

- Right **number** and **type** of operands in an operator or function

$stmt \rightarrow \text{if} (expr) stmt$

$expr \rightarrow expr_1 \text{ relop } expr_2 \quad \{ \text{if} (expr_1.type == expr_2.type) \\ expr.type = \text{Boolean}; \\ \text{else} \\ type\_error(); \}$

- **Coercion**: type of an operand is **automatically** converted to the type expected by the operator
- **Overloading**: meaning of an operator is determined by considering the known types of its **operands** and **results**

# Three-Address Code



- Walking **syntax trees** to generate **three-address code**
- By evaluating **attributes** and **executing code** fragments at nodes in the tree
- Via writing functions that **process** the syntax tree and **emit** the necessary three-address code

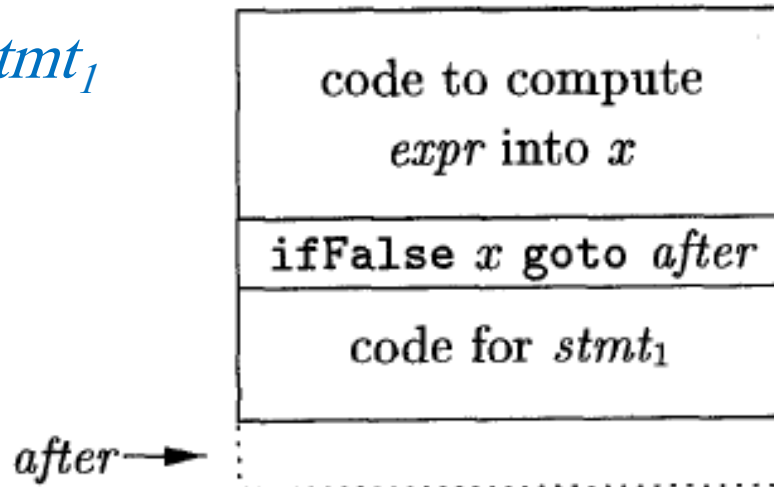
# Three-Address Instructions

- **Binary** operation:  $x = y \text{ op } z$   
 $x, y, z$ : names, constants, compiler-generated temporaries
- **Copy** instruction:  $x = y$
- **Array** access:  $x [ y ] = z$   
 $x = y [ z ]$
- Conditional/Unconditional **jump**:

<code>ifFalse <math>x</math> goto L</code>	if $x$ is false, next execute the instruction labeled L
<code>ifTrue <math>x</math> goto L</code>	if $x$ is true, next execute the instruction labeled L
<code>goto L</code>	next execute the instruction labeled L

# Three-Address Code for Statements

*stmt*  $\rightarrow$  **if** ( *expr* ) *stmt*<sub>1</sub>



```
class If extends Stmt {
    Expr E; Stmt S;
    public If(Expr x, Stmt y) { E = x; S = y; after = newlabel(); }
    public void gen() {
        Expr n = rvalue(E);
        emit( "ifFalse " + n.toString() + " goto " + after );
        S.gen();
        emit( after + ":" );
    }
}
```

94 }

# Three-Address Code for Expressions



- Expressions containing **binary** operators **op**, **array** accesses, **assignments**, **constants**, **identifiers**
- Generating one **three-address** instruction for each **operator** node  $x$  of class *Expr*
  - The computed value is stored into a compiler generated **temporary name**, say  $t$

**$b + c * d$**

**$t1 = c * d$**

**$t2 = b + t1$**

# Three-Address Code for Expressions

```

Expr lvalue(x : Expr) {
    if ( x is an Id node ) return x;
    else error;
}

Expr rvalue(x : Expr) {
    if ( x is an Id or a Constant node ) return x;
    else if ( x is an Op(op, y, z) or a Rel(op, y, z) node ) {
        t = new temporary;
        emit string for t = rvalue(y) op rvalue(z);
        return a new node for t;
    }
    else if ( x is an Assign(y, z) node ) {
        z' = rvalue(z);
        emit string for lvalue(y) = z';
        return z';
    }
}

```

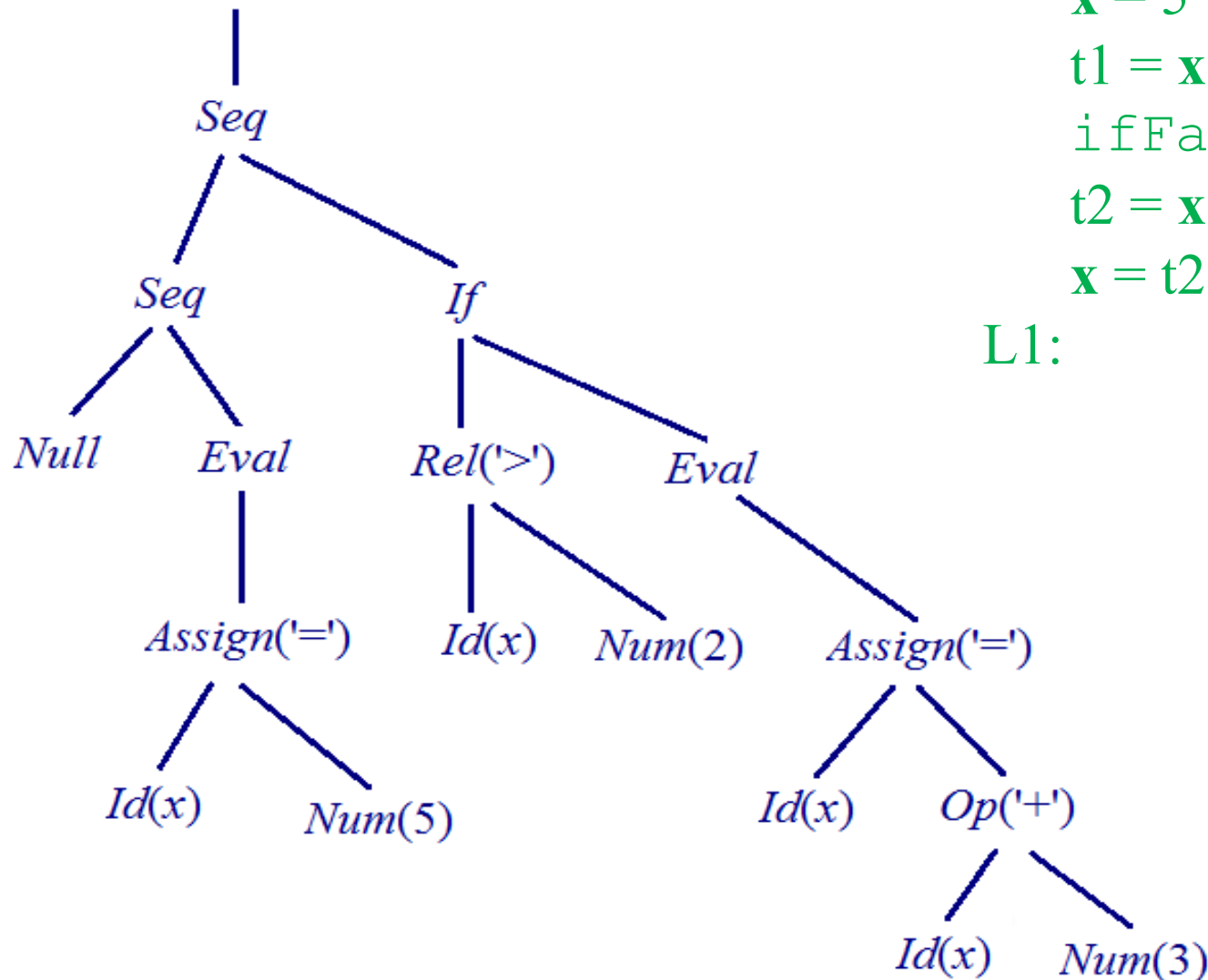


# Example Three-Address Code

```
{ x = 5; if ( x > 2 ) x = x + 3; }
```

```
x = 5
t1 = x > 2
ifFalse t1 goto L1
t2 = x + 3
x = t2
```

L1:



# Three-Address Code for Expressions

- If node  $x$  represents the **array** access:

**a [ b + c ]**                      **t1 = b + c**  
   **a [ t1 ]**

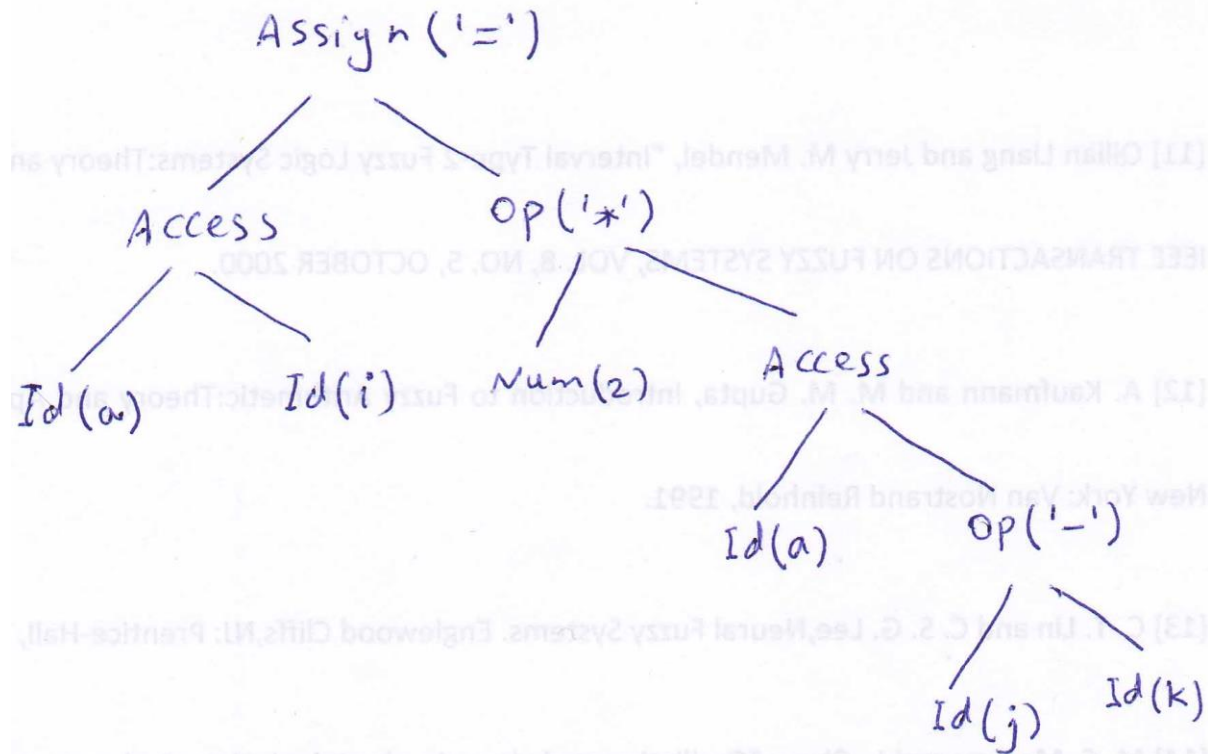
```
Expr lvalue( $x : Expr$ ) {  
    if (  $x$  is an Id node ) return  $x$ ;  
    else if (  $x$  is an Access( $y, z$ ) node and  $y$  is an Id node ) {  
        return new Access( $y, rvalue(z)$ );  
    }  
    else error;  
}
```

# Three-Address Code for Expressions

```
Expr rvalue(x : Expr) {  
    if ( x is an Id or a Constant node ) return x;  
    else if ( x is an Op(op, y, z) or a Rel(op, y, z) node ) {  
        t = new temporary;  
        emit string for t = rvalue(y) op rvalue(z);  
        return a new node for t;  
    }  
    else if ( x is an Access(y, z) node ) {  
        t = new temporary;  
        call lvalue(x), which returns Access(y, z');  
        emit string for t = Access(y, z');  
        return a new node for t;  
    }  
    else if ( x is an Assign(y, z) node ) {  
        z' = rvalue(z);  
        emit string for lvalue(y) = z';  
        return z';  
    }  
}
```

# Three-Address Code for Expressions

`a[i] = 2*a[j-k]`



`t3 = j - k`  
`t2 = a [ t3 ]`  
`t1 = 2 * t2`  
`a [ i ] = t1`