

A Simple One-Pass Compiler

Building a Simple Compiler

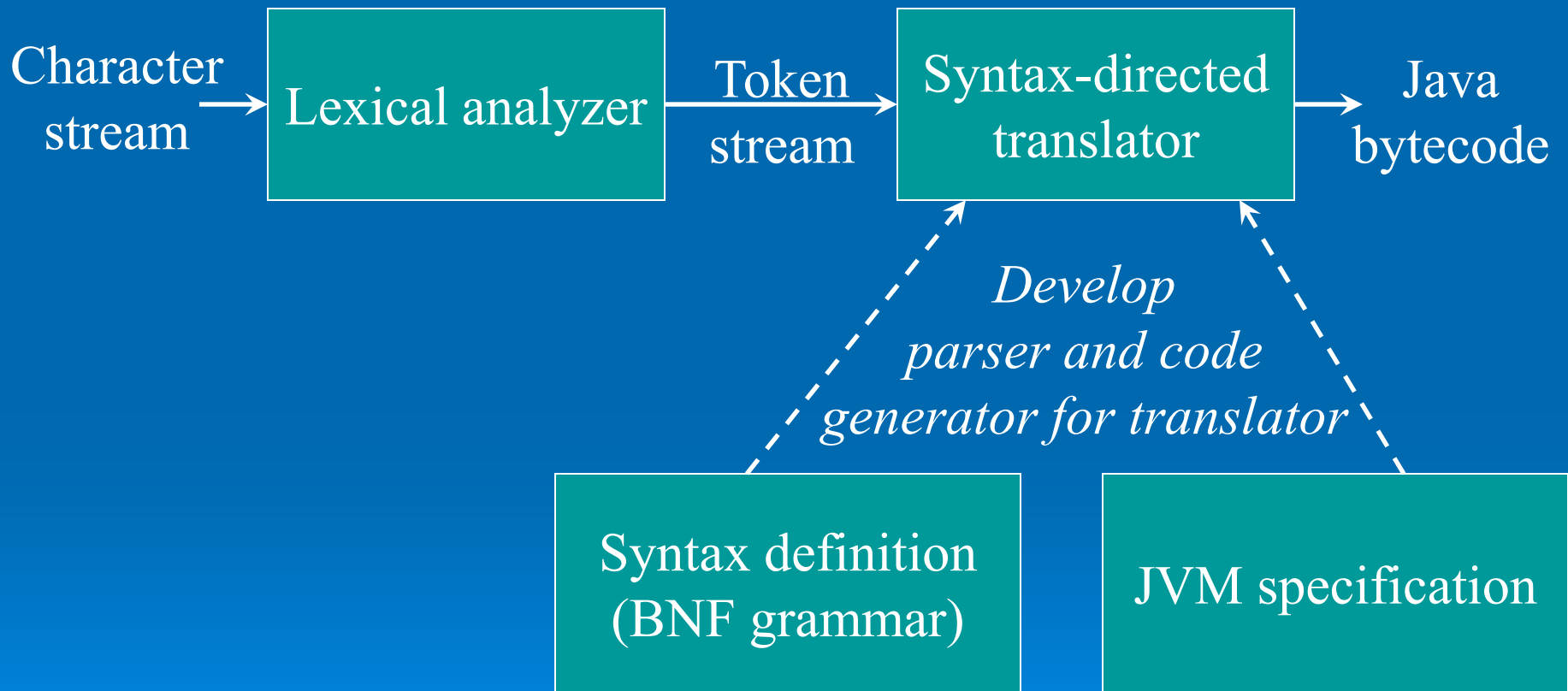
- Building our compiler involves:
 - Defining the *syntax* of a programming language
 - Develop a source code parser: for our compiler we will use *predictive parsing*
 - Implementing *syntax directed translation* to generate intermediate code

Overview

Programming Language can be defined by describing

1. The **syntax** of the language
 1. *How the program look like*
 2. *We use BNF (Backus Naur Form) for defining the structure of the syntax*
2. The **semantics** of the language
 1. *What does the program mean*
 2. *Difficult to describe*
 3. *Use informal descriptions and suggestive examples*

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

list

\Rightarrow list + digit

\Rightarrow list - digit + digit

\Rightarrow digit - digit + digit

\Rightarrow 9 - digit + digit

\Rightarrow 9 - 5 + digit

\Rightarrow 9 - 5 + 2

$list \rightarrow list + digit$

$list \rightarrow list - digit$

$list \rightarrow digit$

$digit \rightarrow 0 | 1 | 2 | 3 |$

$4 | 5 | 6 | 7 | 8 | 9$

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

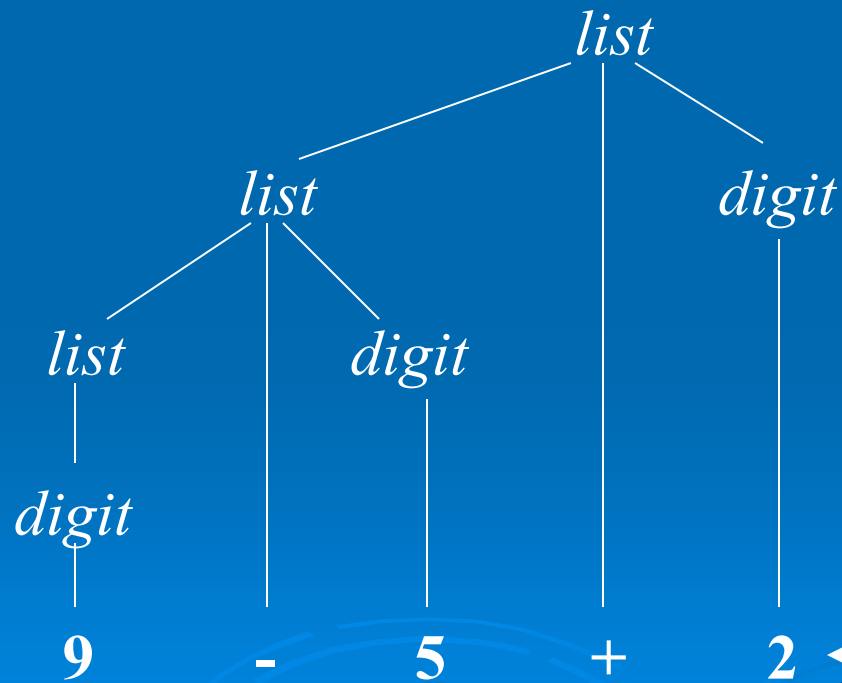
Likewise, a *rightmost derivation* replaces the rightmost nonterminal in each step

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 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 immediate *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:

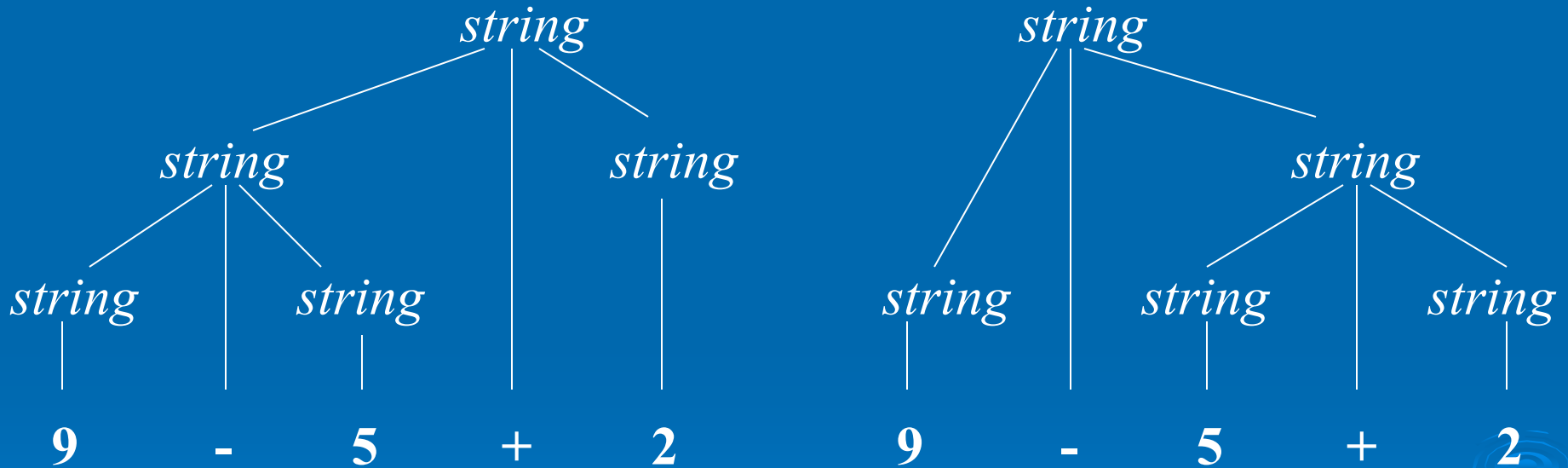
$$G = \langle \{string\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, string \rangle$$

with production $P =$

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

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

Ambiguity (cont'd)



Associativity of Operators

Left-associative operators have *left-recursive* productions

$$\textit{left} \rightarrow \textit{left} + \textit{term} \mid \textit{term}$$

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

Right-associative operators have *right-recursive* productions

$$\textit{right} \rightarrow \textit{term} = \textit{right} \mid \textit{term}$$

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

Precedence of Operators

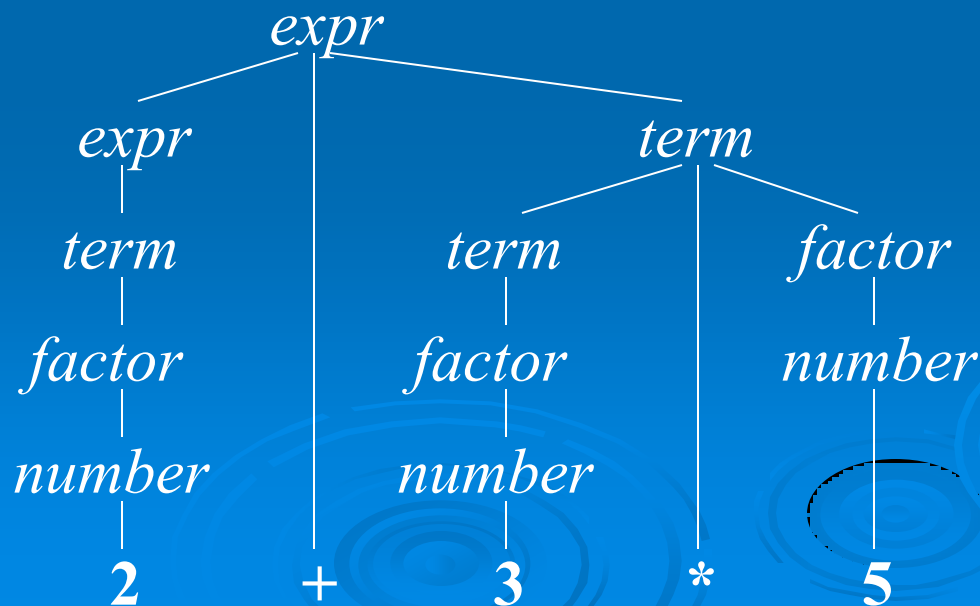
Operators with higher precedence “bind more tightly”

$expr \rightarrow expr + term \mid term$

$term \rightarrow term * factor \mid factor$

$factor \rightarrow number \mid (expr)$

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



Other Important Concepts

Operator Precedence

What does
 $9 + 5 * 2$
mean?

Typically

$\left\{ \begin{array}{l} () \\ * / \\ + - \end{array} \right.$ is precedence order

This can be
incorporated
into a grammar
via rules:

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

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

$factor \rightarrow digit \mid (expr)$

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$

Precedence Achieved by:

expr & term for each precedence level

Rules for each are **left recursive** or **associate to the left**

Syntax of Statements

$stmt \rightarrow id := expr$
| **if** $expr$ **then** $stmt$
| **if** $expr$ **then** $stmt$ **else** $stmt$
| **while** $expr$ **do** $stmt$
| **begin** opt_stmts **end**

$opt_stmts \rightarrow stmt ; opt_stmts$
| ϵ

Syntax-Directed Translation

- Uses a CF grammar 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

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

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

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

$expr.t := term.t$

$term.t := "0"$

$term.t := "1"$

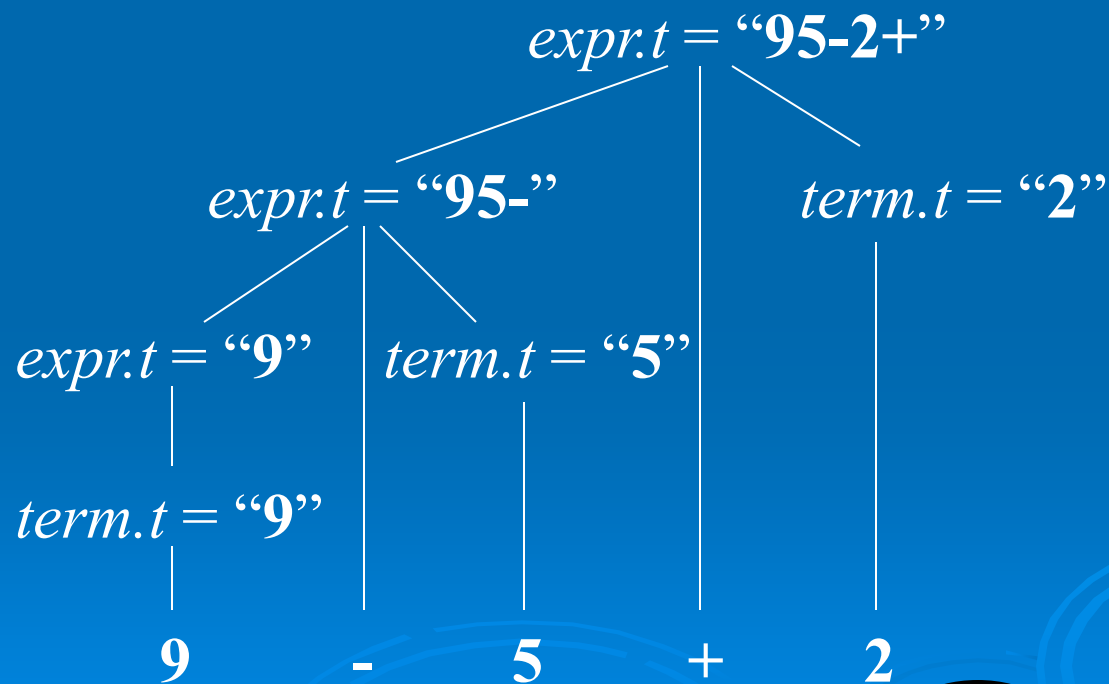
...

$term.t := "9"$

String concat operator



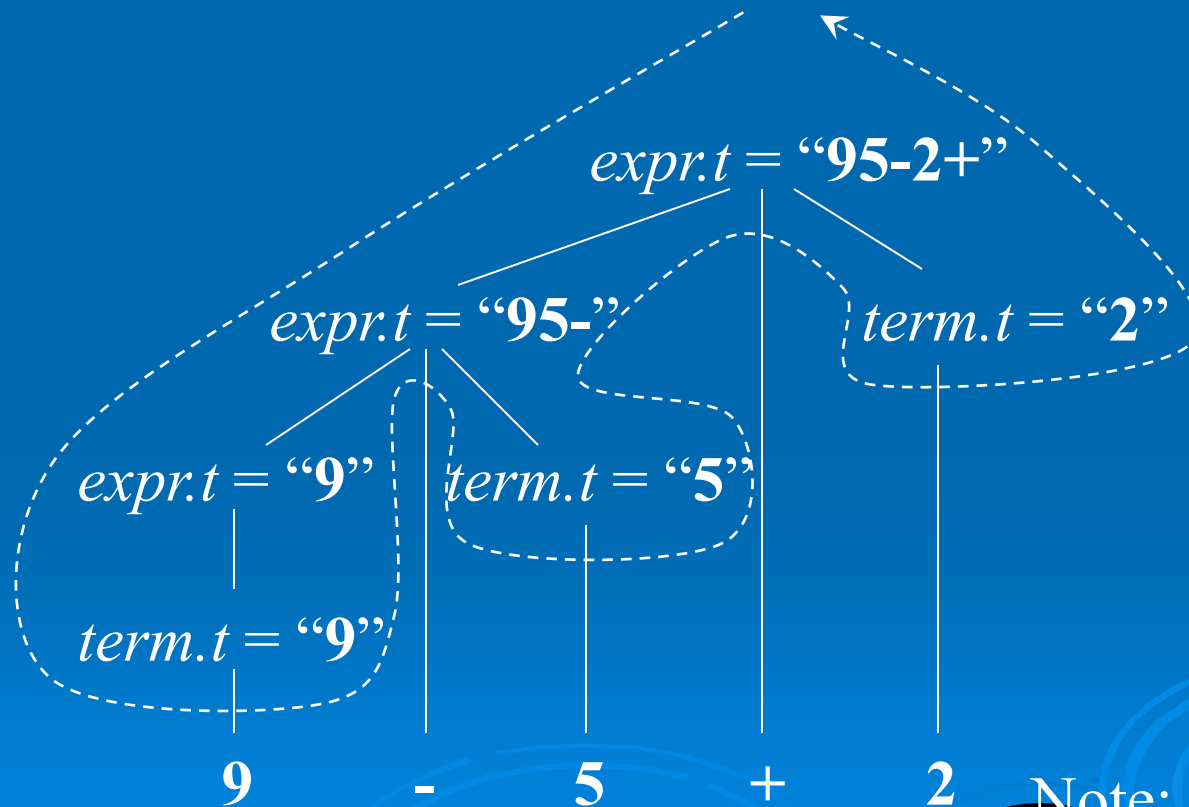
Example Annotated Parse Tree



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)



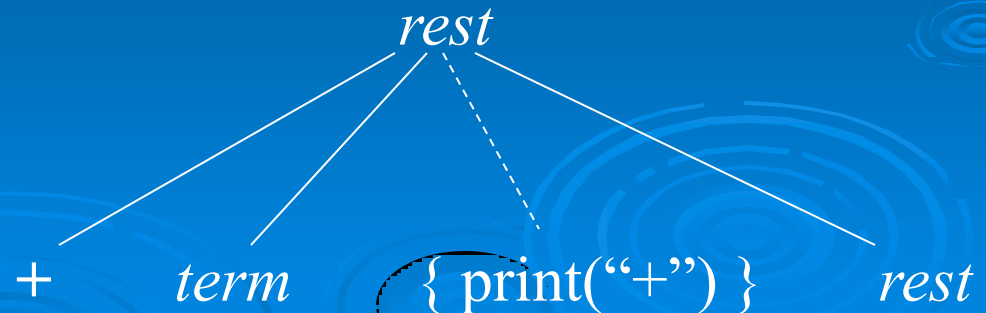
Note: all attributes are of the synthesized type

Translation Schemes

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

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

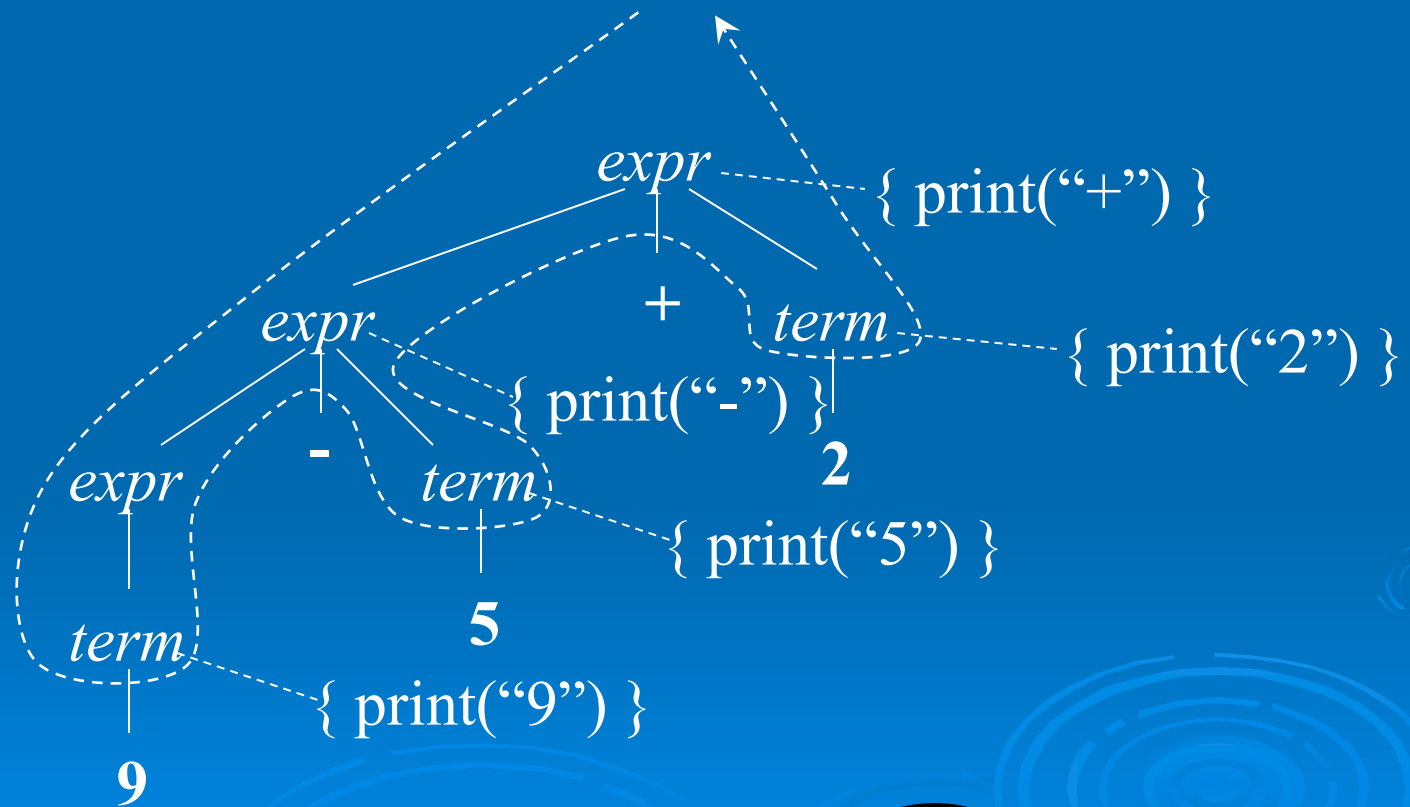
Embedded
semantic action



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 source code
- *Top-down parsing* “constructs” a parse tree from root to leaves
- *Bottom-up parsing* “constructs” a 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 → *simple*
| ^ **id**
| **array** [*simple*] **of** *type*
simple → **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;
```

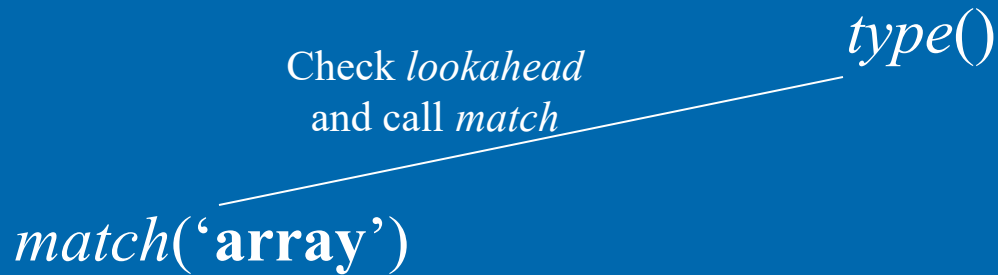
```
procedure simple();  
begin  
  if lookahead = 'integer' then  
    match('integer')  
  else if lookahead = 'char' then  
    match('char')  
  else if lookahead = 'num' then  
    match('num');  
    match('dotdot');  
    match('num')  
  else error()  
end;
```

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)

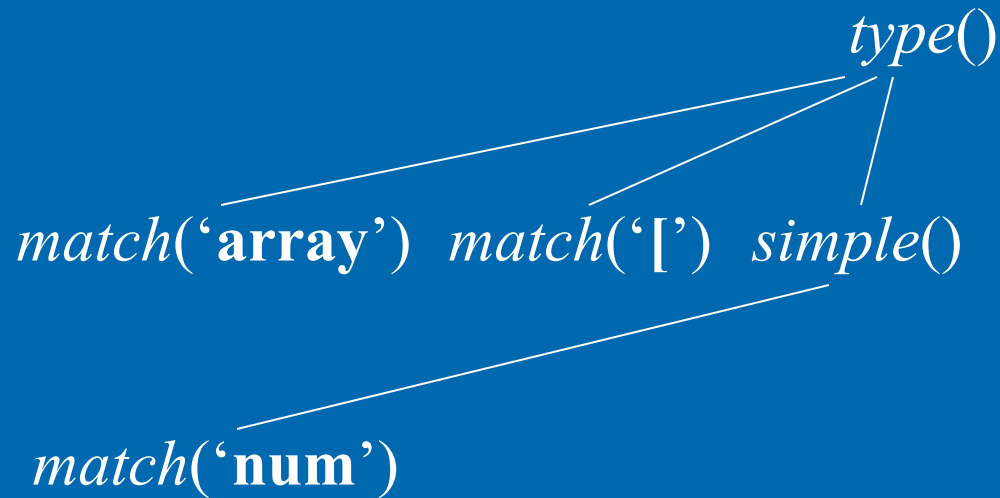
match('array') *match('[')* *type()*

Input:

array num dotdot num] of integer

 ↑
 |
 ↑
lookahead

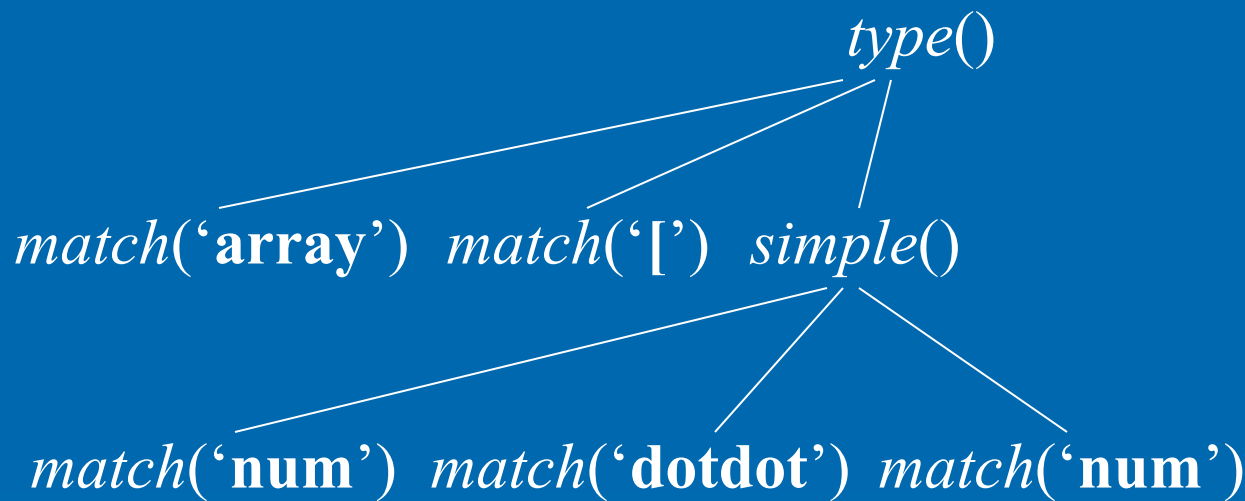
Example Predictive Parser (Execution Step 3)



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

↑
lookahead

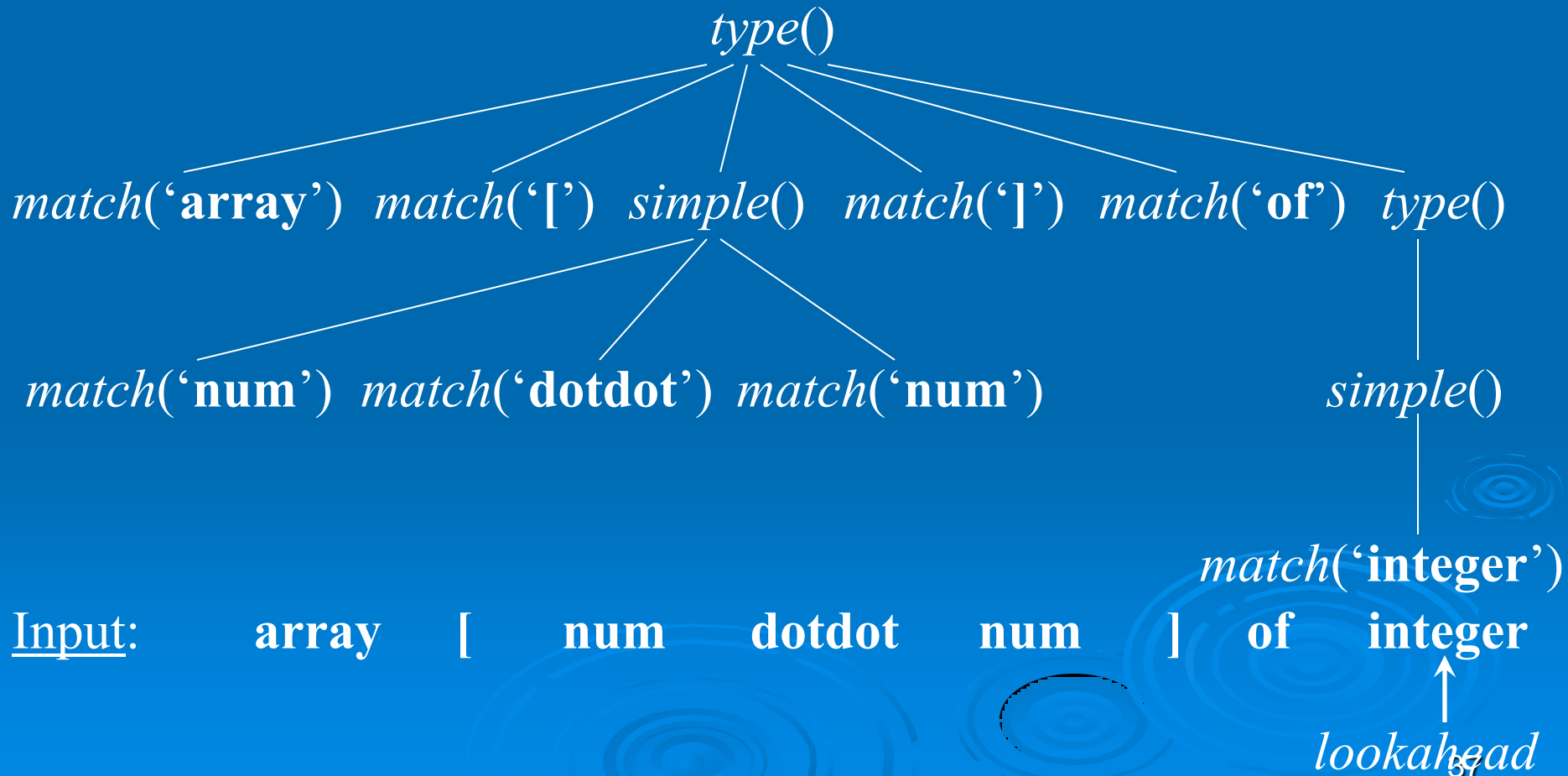
Example Predictive Parser (Execution Step 5)



Input: array [num dotdot num] of integer

↑
lookahead

Example Predictive Parser (Execution Step 8)



FIRST

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

$$\begin{aligned} \textit{type} &\rightarrow \textit{simple} \\ &| \wedge \textbf{id} \\ &| \textbf{array} [\textit{simple}] \textbf{of type} \\ \textit{simple} &\rightarrow \textbf{integer} \\ &| \textbf{char} \\ &| \textbf{num dotdot num} \end{aligned}$$
$$\text{FIRST}(\textit{simple}) = \{ \textbf{integer}, \textbf{char}, \textbf{num} \}$$
$$\text{FIRST}(\wedge \textbf{id}) = \{ \wedge \}$$
$$\text{FIRST}(\textit{type}) = \{ \textbf{integer}, \textbf{char}, \textbf{num}, \wedge, \textbf{array} \}$$

How to use FIRST

We use FIRST to write a predictive parser as follows

| | | |
|----------------------------------|--|---|
| $expr \rightarrow term\ rest$ | | procedure <i>rest</i> (); |
| $rest \rightarrow +\ term\ rest$ | | begin |
| $- \ term\ rest$ | | if <i>lookahead</i> in <u>FIRST(+ <i>term rest</i>)</u> then |
| ϵ | | <i>match</i> ('+'); <i>term</i> (); <i>rest</i> () |
| | | else if <i>lookahead</i> in <u>FIRST(- <i>term rest</i>)</u> then |
| | | <i>match</i> (' - '); <i>term</i> (); <i>rest</i> () |
| | | else return |
| | | end; |

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

$$A \rightarrow \alpha$$
$$| \beta$$

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

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

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

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

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:

$expr \rightarrow term \text{ rest}$

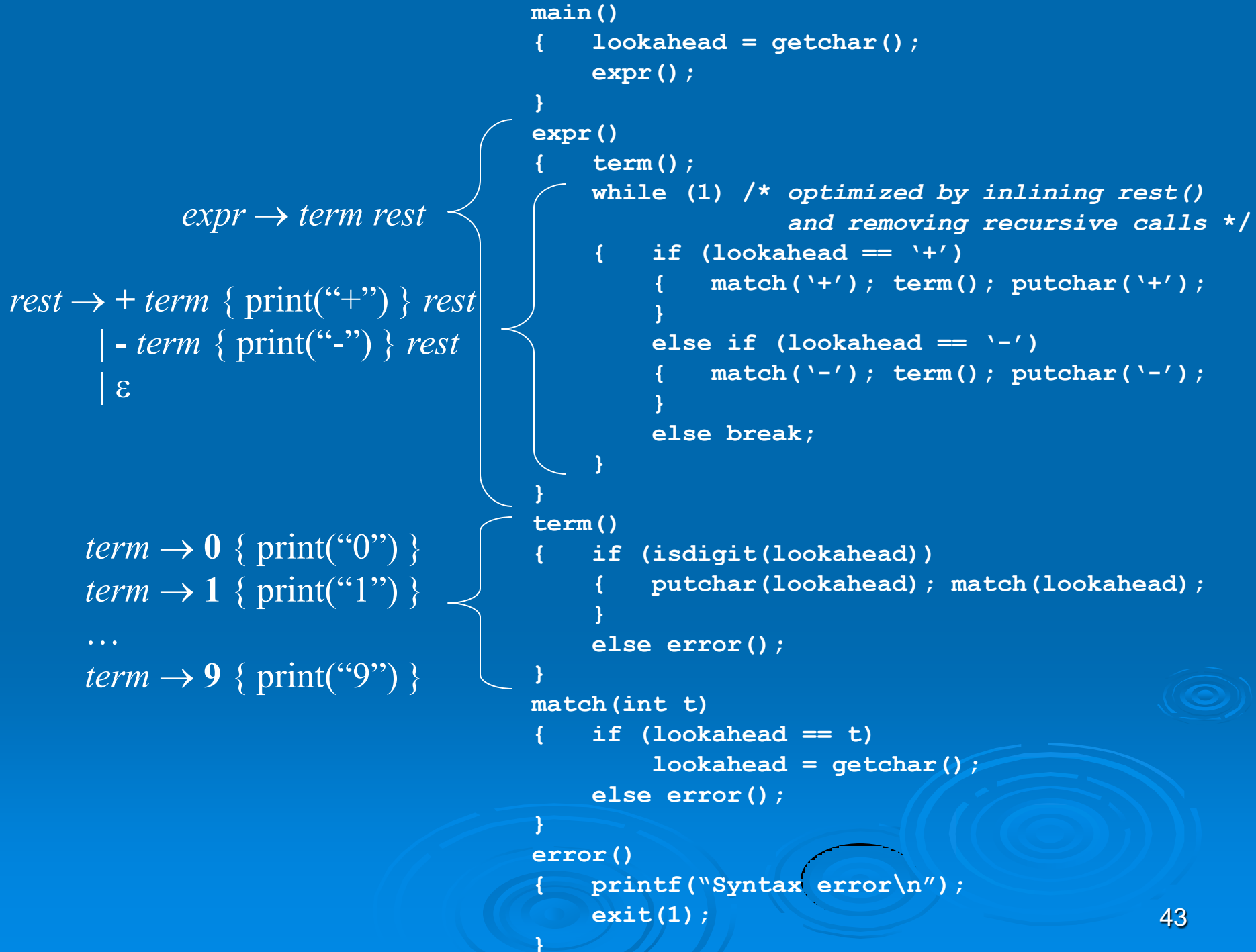
$rest \rightarrow + term \{ \text{print}("+") \} rest \mid - term \{ \text{print}("-") \} rest \mid \varepsilon$

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

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

...

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



Adding a Lexical Analyzer

- Typical tasks of the lexical analyzer:
 - Remove white space and comments
 - Encode constants as tokens
 - Recognize keywords
 - Recognize identifiers and store identifier names in a global symbol table

Lexical Analysis

A lexical analyzer reads and converts the input into a stream of tokens to be analyzed by the parser.

A sequence of input characters that comprises a single token is called a **lexeme**.

Functional Responsibilities

1. White Space and Comments Are Filtered Out

blanks, new lines, tabs are removed

comments can be ignored

Functional Responsibilities (2)

Constants

- The job of collecting digits into integers is generally given to a lexical analyzer because numbers can be treated as single units during translation.
- **num** be the token representing an integer.
- The value of the integer will be passed along as an attribute of the token **num**
- **Example:**

31 + 28 + 59

<num, 31> **<+, >** **<num, 28>** **<+, >** **<num, 31>**

Functional Responsibilities (3)

Recognizing Identifiers and Keywords

Compilers use identifiers as names of

- *Variables*
- *Arrays*
- *Functions*

A grammar for a language treats an **identifier** as **token**

Example:

```
credit = asset + goodwill;
```

Lexical analyzer would convert it like

```
id = id + id ;
```

Functional Responsibilities (3)

Recognizing Identifiers and Keywords (2)

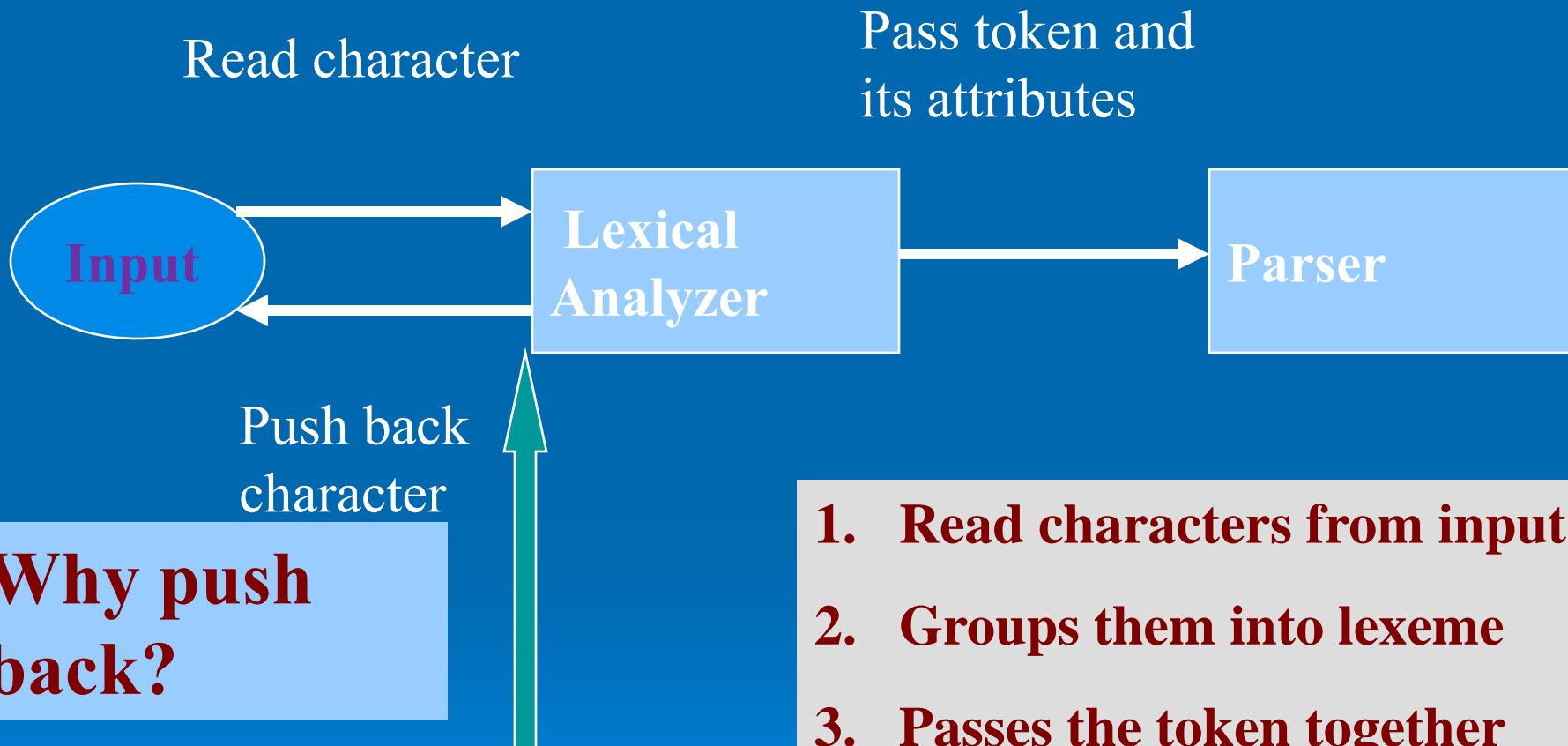
Languages use fixed character strings (if, while, extern) to identify certain construct. We call them *keywords*.

A mechanism is needed fir deciding when a lexeme forms a keyword and when it forms an identifier.

Solution

- 1. Keywords are reserved.**
- 2. The character string forms an identifier only if it is not a keyword.**

Interface to the Lexical Analyzer



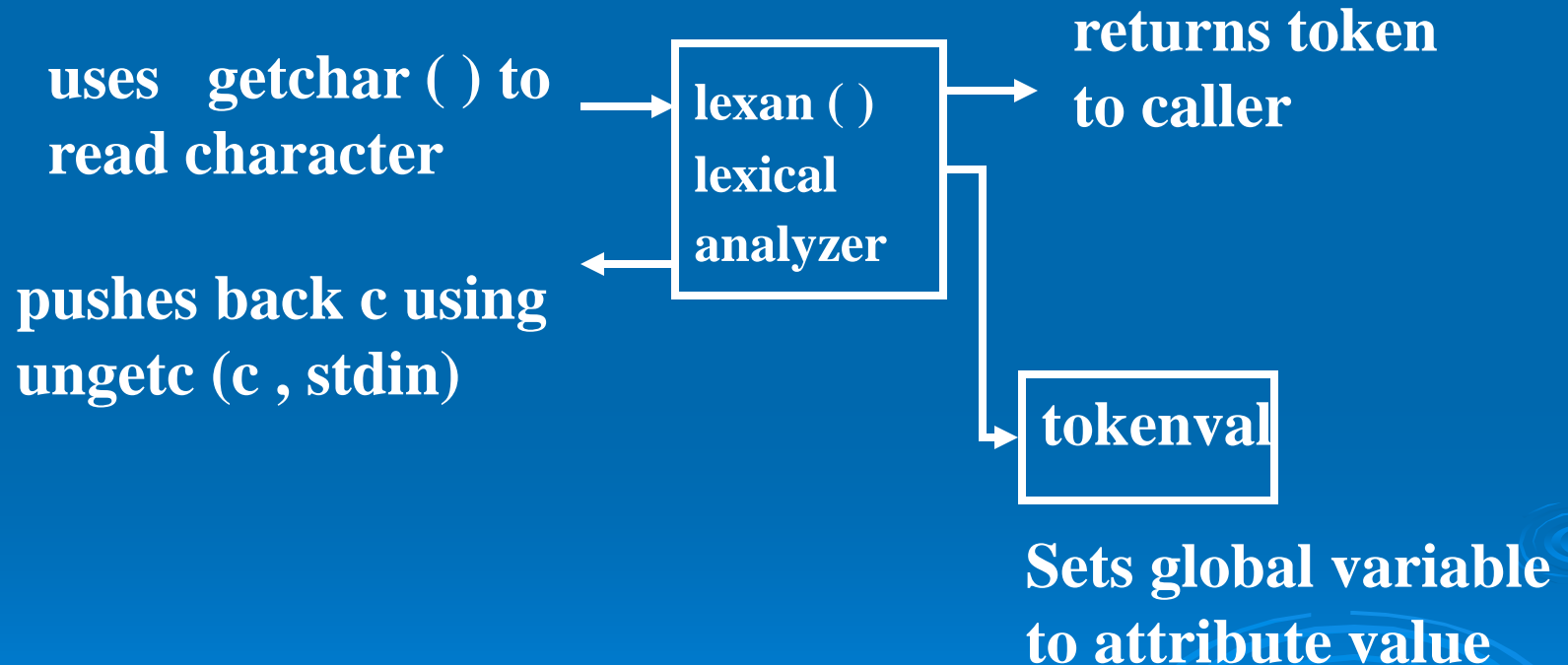
Why push back?

This part is implemented with a buffer

- 1. Read characters from input**
- 2. Groups them into lexeme**
- 3. Passes the token together with attribute values to the later stage**

The Lexical Analysis Process

A Graphical Depiction



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

Symbol Table Considerations

OPERATIONS: Insert (string, token_ID)

Lookup (string)

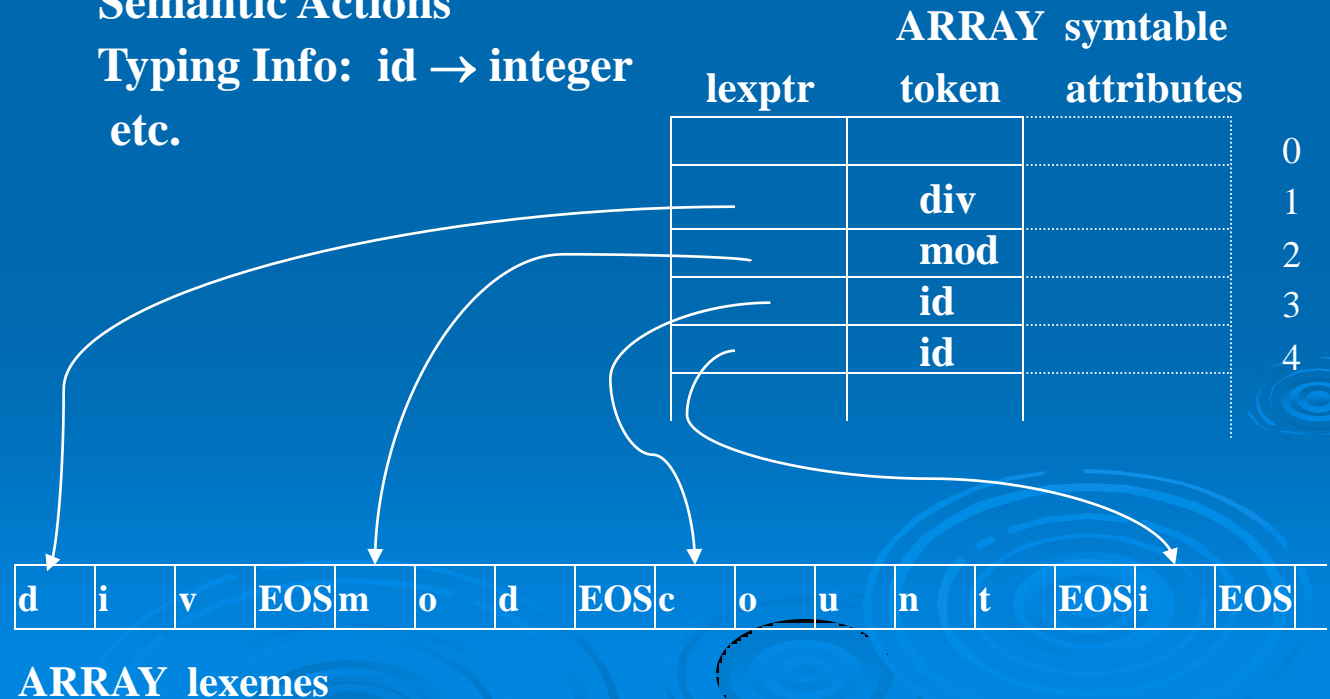
NOTICE: Reserved words are placed into symbol table for easy lookup

Attributes may be associated with each entry, i.e.,

Semantic Actions

Typing Info: id → integer

etc.



Identifiers

factor \rightarrow (*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)

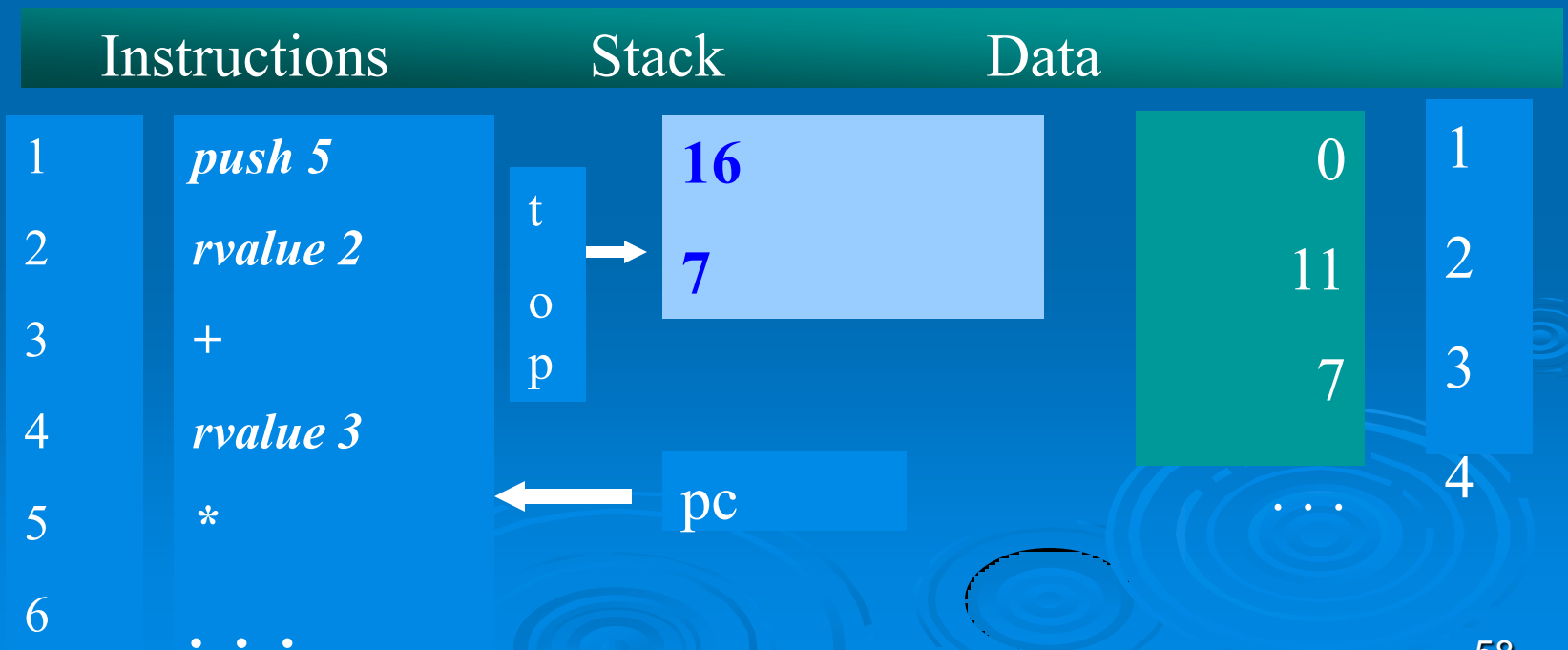
morefactors → **div** *factor* { print('DIV') } *morefactors*
 | **mod** *factor* { print('MOD') } *morefactors*
 | ...

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

Abstract Stack Machines

Instructions fall into three classes.

1. Integer arithmetic
2. Stack manipulation
3. Control flow



L-value and R-value

What is the difference between left and right side identifier?

L-value Vs. R-value of an identifier

$I := 5 ;$ L - Location

$I := I + 1 ;$ R - Contents

The right side specifies an integer value, while left side specifies where the value is to be stored.

Usually,

r-values are what we think as values

l-values are locations.

Generic Instructions for Stack Manipulation

| | |
|----------------------------------|---|
| push <i>v</i> | push constant value <i>v</i> onto the stack |
| rvalue <i>l</i> | push contents of data location <i>l</i> |
| lvalue <i>l</i> | push address of data location <i>l</i> |
| pop | discard value on top of the stack |
| := | the r-value on top is placed in the l-value below it and both are popped |
| copy | push a copy of the top value on the stack |
| + | add value on top with value below it pop both and push result |
| - | subtract value on top from value below it pop both and push result |
| * , / , ... | ditto for other arithmetic operations |
| < , & , ... | ditto for relational and logical operations |

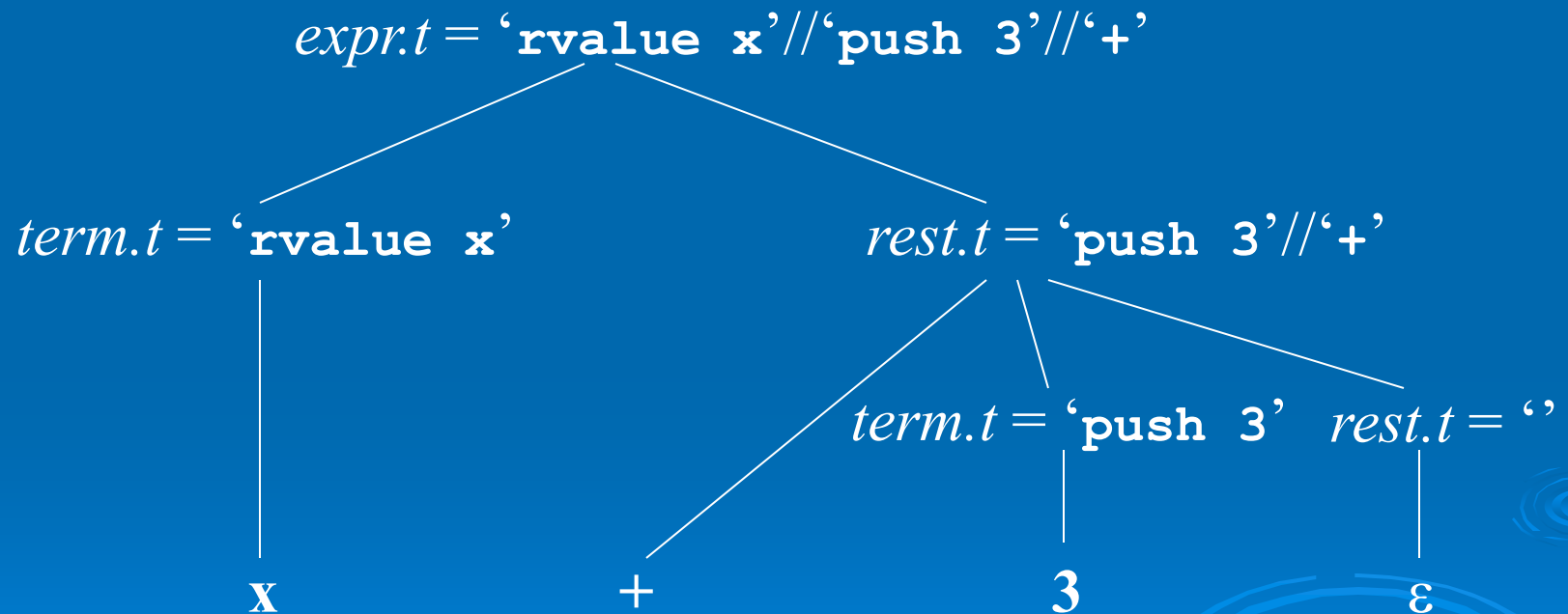
Generic Control Flow Instructions

| | |
|-------------------------|---|
| label <i>l</i> | label instruction with <i>l</i> |
| goto <i>l</i> | jump to instruction labeled <i>l</i> |
| gofalse <i>l</i> | pop the top value, if zero then jump to <i>l</i> |
| gotrue <i>l</i> | pop the top value, if nonzero then jump to <i>l</i> |
| halt | stop execution |
| jsr <i>l</i> | jump to subroutine labeled <i>l</i> , push return address |
| return | pop return address and return to caller |

Syntax-Directed Translation of Expressions

$expr \rightarrow term\ rest \{ expr.t := term.t \ //\ rest.t \}$
 $rest \rightarrow +\ term\ rest_1 \{ rest.t := term.t \ //\ '+' \ //\ rest_1.t \}$
 $rest \rightarrow -\ term\ rest_1 \{ rest.t := term.t \ //\ '-' \ //\ rest_1.t \}$
 $rest \rightarrow \varepsilon \{ rest.t := '' \}$
 $term \rightarrow \mathbf{num} \{ term.t := \mathbf{'push'} \ //\ \mathbf{num.value} \}$
 $term \rightarrow \mathbf{id} \{ term.t := \mathbf{'rvalue'} \ //\ \mathbf{id.lexeme} \}$

Syntax-Directed Translation of Expressions (cont'd)



Translation Scheme to Generate Abstract Machine Code

expr → *term moreterms*

moreterms → + *term* { print('+') } *moreterms*

moreterms → - *term* { print('-') } *moreterms*

moreterms → ε

term → *factor morefactors*

morefactors → * *factor* { print('*') } *morefactors*

morefactors → **div** *factor* { print('DIV') } *morefactors*

morefactors → **mod** *factor* { print('MOD') } *morefactors*

morefactors → ε

factor → (*expr*)

factor → **num** { print('push ' // **num.value**) }

factor → **id** { print('rvalue ' // **id.lexeme**) }

Translation Scheme to Generate Abstract Machine Code (cont'd)

$stmt \rightarrow id := \{ \text{print}('lvalue' // id.lexeme) \} expr \{ \text{print}(':=') \}$

| |
|-------------------------------|
| <code>lvalue id.lexeme</code> |
| code for <code>expr</code> |
| <code>:=</code> |

Translation Scheme to Generate Abstract Machine Code (cont'd)

$stmt \rightarrow \mathbf{if} \ expr \ \{ \ out := \mathbf{newlabel}(); \ \mathbf{print}(\mathbf{'gofalse' // out}) \}$
 $\mathbf{then} \ stmt \ \{ \ \mathbf{print}(\mathbf{'label' // out}) \}$

| |
|---------------------------|
| <i>code for expr</i> |
| gofalse <i>out</i> |
| <i>code for stmt</i> |
| label <i>out</i> |

Translation Scheme to Generate Abstract Machine Code (cont'd)

stmt → **while** { *test* := newlabel(); print('label ' // *test*) }
expr { *out* := newlabel(); print('gofalse ' // *out*) }
do *stmt* { print('goto ' // *test* // 'label ' // *out*) }

| |
|---------------------------|
| label <i>test</i> |
| code for <i>expr</i> |
| gofalse <i>out</i> |
| code for <i>stmt</i> |
| goto <i>test</i> |
| label <i>out</i> |

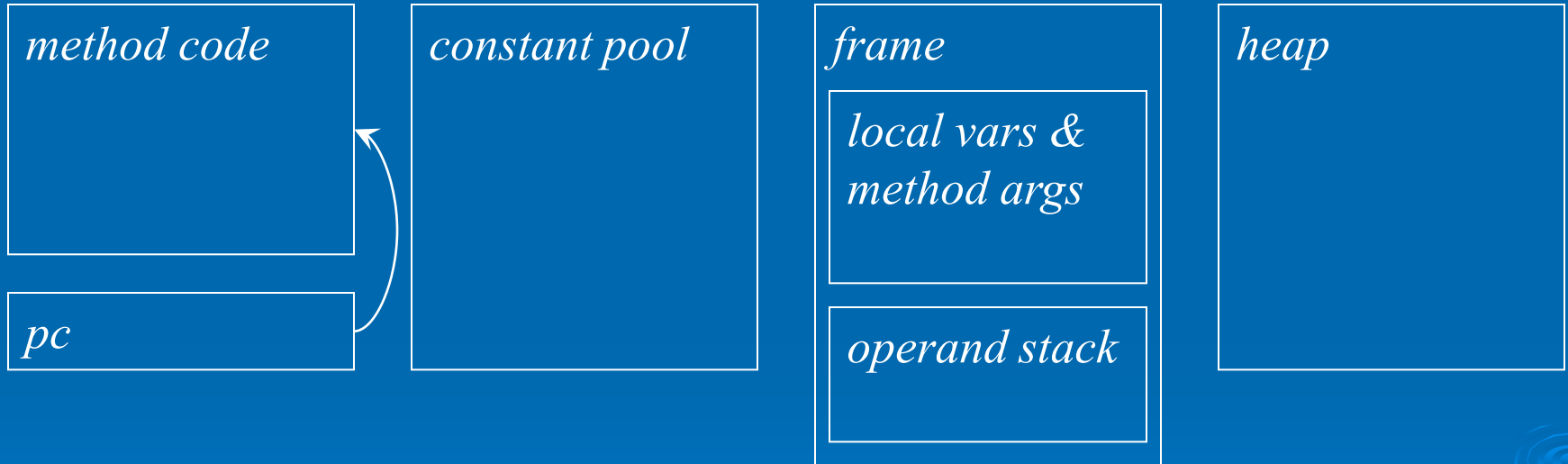
Translation Scheme to Generate Abstract Machine Code (cont'd)

$start \rightarrow stmt \{ \text{print}(\text{'halt'}) \}$
 $stmt \rightarrow \mathbf{begin} \text{ opt_stmts } \mathbf{end}$
 $opt_stmts \rightarrow stmt ; \text{ opt_stmts } \mid \varepsilon$

The JVM

- Abstract stack machine architecture
 - Emulated in software with JVM interpreter
 - *Just-In-Time* (JIT) compilers
 - Hardware implementations available
- *Java bytecode*
 - Platform independent
 - Small
 - Safe
- The Java™ Virtual Machine Specification, 2nd ed.
<http://java.sun.com/docs/books/vmspec>

Runtime Data Areas



Constant Pool

- Serves a function similar to that of a symbol table
- Contains several kinds of constants
- Method and field references, strings, float constants, and integer constants larger than 16 bit cannot be used as operands of bytecode instructions and must be loaded on the operand stack from the constant pool
- Java *bytecode verification* is a pre-execution process that checks the consistency of the bytecode instructions and constant pool

Frames

- A new *frame* (also known as *activation record*) is created each time a method is invoked
- A frame is destroyed when its method invocation completes
- Each frame contains an array of variables known as its *local variables* indexed from 0
 - Local variable 0 is “*this*” (unless the method is static)
 - Followed by method parameters
 - Followed by the local variables of blocks
- Each frame contains an *operand stack*

Data Types

| | |
|----------------------|--|
| byte | a 8-bit signed two's complement integer |
| short | a 16-bit signed two's complement integer |
| int | a 32-bit signed two's complement integer |
| long | a 64-bit signed two's complement integer |
| char | a 16-bit Unicode character |
| float | a 32-bit IEEE 754 single-precision float value |
| double | a 64-bit IEEE 754 double-precision float value |
| boolean | a virtual type only, int is used to represent true (1) false (0) |
| returnAddress | the location of the <i>pc</i> after method invocation |
| reference | a 32-bit address reference to an object of <i>class type</i> , <i>array type</i> , or <i>interface type</i> (value can be NULL) |

Operand stack has 32-bit slots, thus **long** and **double** occupy two slots

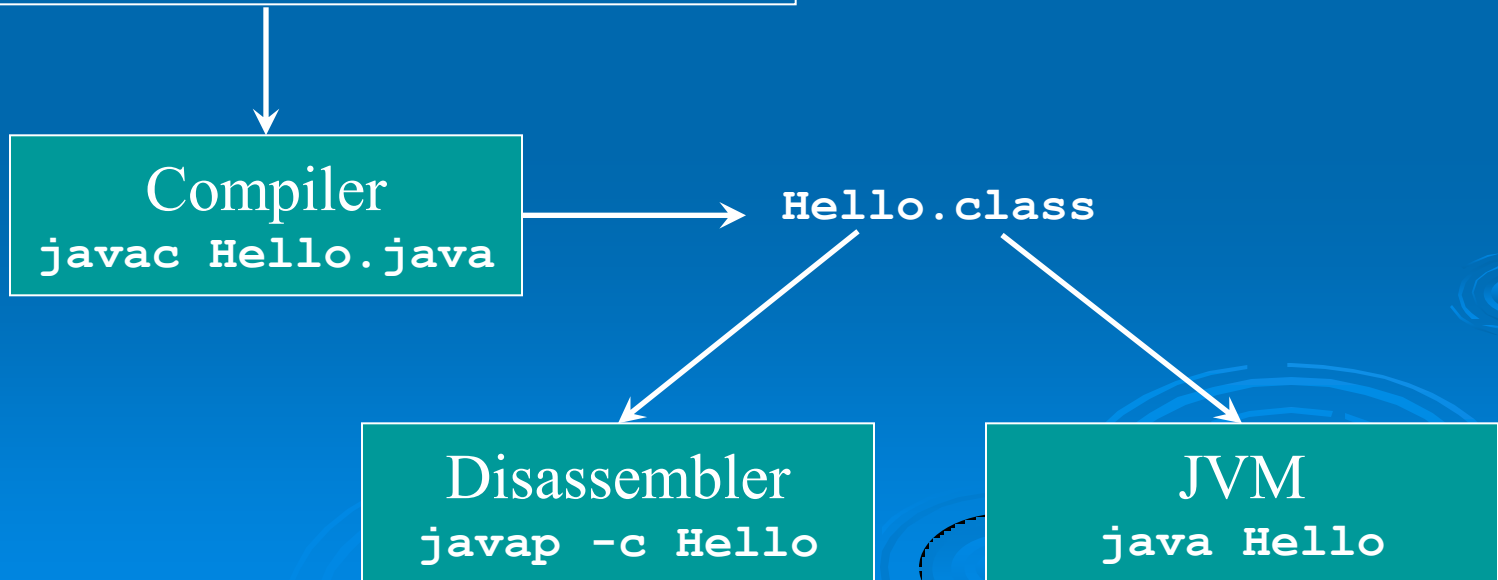
The Class File Format

- A *class file* consists of a stream of 8-bit bytes
- 16-, 32-, and 64-bit quantities are stored in 2, 4, and 8 consecutive bytes in *big-endian* order
- Contains several components, including:
 - Magic number `0xCAFEBABE`
 - Version info
 - Constant pool
 - This and super class references (index into pool)
 - Class fields
 - Class methods

javac, javap, java

Hello.java

```
import java.lang.*;
public class Hello
{ public static void main(String[] arg)
  { System.out.println("Hello World!");
  }
}
```



javap -c Hello

Local variable 0 = "this"

Index into constant pool

Method descriptor

Compiled from "Hello.java"

```
public class Hello extends java.lang.Object{  
public Hello();
```

Code:

```
0:  aload_0  
1:  invokespecial  #1; //Method java/lang/Object."<init>":()V  
4:  return
```

```
public static void main(java.lang.String[]);
```

Code:

```
0:  getstatic      #2; //Field java/lang/System.out:Ljava/io/PrintStream;  
3:  ldc           #3; //String Hello World!  
5:  invokevirtual #4; //Method java/io/PrintStream.println:(Ljava/lang/String;)V  
8:  return
```

```
}
```

String literal

Field descriptor

Generating Code for the JVM

expr → *term* *moreterms*

moreterms → + *term* { emit(**iadd**) } *moreterms*

moreterms → - *term* { emit(**isub**) } *moreterms*

moreterms → ε

term → *factor* *morefactors*

morefactors → * *factor* { emit(**imul**) } *morefactors*

morefactors → **div** *factor* { emit(**idiv**) } *morefactors*

morefactors → **mod** *factor* { emit(**irem**) } *morefactors*

morefactors → ε

factor → (*expr*)

factor → **int8** { emit2(**bipush**, **int8.value**) }

factor → **int16** { emit3(**sipush**, **int16.value**) }

factor → **int32** { *idx* := newpoolint(**int32.value**);
emit2(**ldc**, *idx*) }

factor → **id** { emit2(**iload**, **id.index**) }

Generating Code for the JVM (cont'd)

$stmt \rightarrow id := expr \{ emit2(istore, id.index) \}$

| |
|-------------------------------|
| code for <i>expr</i> |
| istore <i>id.index</i> |

$stmt \rightarrow \text{if } expr \{ emit(iconst_0); loc := pc; emit3(if_icmpeq, 0) \}$
 $\text{then } stmt \{ backpatch(loc, pc - loc) \}$

| |
|--|
| code for <i>expr</i> |
| iconst_0 |
| if_icmpeq <i>off₁</i> <i>off₂</i> |
| code for <i>stmt</i> |
| <i>pc:</i> |

loc: 
pc:

backpatch() sets the offsets of the relative branch when the target *pc* value is known

Thank You