

# *CS201 - PL'S*

## Lexical & Syntactic Analysis

Aravind Mohan

Allegheny College

September 6, 2021



# Most Important Steps in Compilation

- Optional Preprocessing
- Lexical analysis (scanning)
- Syntax analysis (parsing)
- Semantic analysis
- Intermediate code generation
- Optimization (usually machine-independent)
- Final code generation
- Optional final optimization

# Lexical Analysis

0	<unsigned int>
1	'('
2	')'
3	'+'
4	'-'
5	"for"
6	"while"
7	<identifier> (not reserved)

8	';'
9	'='
10	"=="
11	"<="
12	">="
13	<string literal>
14	<
...	... etc. ...

A token is any component of a program that is generally treated as an indivisible piece, e.g., a variable name, an operator such as <=, a punctuation mark such as a semicolon, a string constant, etc.

For each token type, give a description:

- either a literal string (e.g., “ $\leq$ ” or “while” to describe an operator or reserved word),
- or a  $< rule >$  (e.g., the rule  $< \text{unsigned int} >$  might stand for “a sequence of one or more digits”; the rule  $< \text{identifier} >$  might stand for “a letter followed by a sequence of zero or more letters or digits.”)

**Lexical analysis** produces a “token stream” in which the program is reduced to a sequence of token types, each with its identifying number and the actual string (in the program) corresponding to it.

# Lexical Analysis

## Lexical Analysis

```
// see if 3 occurs
while x <= 10
    a = x+1
    while (a == 3)
        found = 1
        a = f(x)
```

Program

6, "while" → 10, "=="  
7, "x" 0, "3"  
11, "<=" 2, ")"  
0, "10" 7, "found"  
7, "a" 9, "="  
9, "=" 0, "1"  
7, "x" 7, "a"  
3, "+" 9, "="  
0, "1" 7, "f"  
6, "while" 1, "("  
1, "(" 7, "x"  
7, "a" 2, ")")

Stream of Tokens

- The syntax of a language is described by a **grammar** that specifies the legal combinations of tokens.
- Grammars are often specified in BNF notation (“Backus Naur Form”):

```
<item1> ::= valid replacements for  
<item1>  
<item2> ::= valid replacements for  
<item2>
```

# Syntactic Analysis

**Example:** an expression can be either a simple variable identifier; an integer; or an expression, followed by an operator, followed by another expression:

$\langle \text{expr} \rangle ::= \langle \text{id} \rangle \mid \langle \text{int} \rangle \mid \langle \text{expr} \rangle \text{ op } \langle \text{expr} \rangle$

Classic BNF notation

**Alternative notations:**

$\text{expr} \rightarrow \text{id} \mid \text{int} \mid \text{expr op expr}$

The book uses this notation  
(but as three separate rules)

$\text{expr} ::= \text{id} \mid \text{int} \mid \text{expr } \{\text{op expr}\}^*$

The symbol " $|$ " means "or"

The " $\{\dots\}^*$ " means "zero or more  
repetitions of the items in  $\{\dots\}$ "

# Grammars (Context-free Grammars)

- Collection of VARIABLES (things that can be replaced by other things), also called NON-TERMINALS.
- Collection of TERMINALS (“constants”, strings that can't be replaced)
- One special variable called the START SYMBOL.
- Collection of RULES, also called PRODUCTIONS.

variable → rule1 | rule2 | rule3 | ...

You can also write each rule on a separate line (as in the book)

## Grammar

A, B, and C are non-terminals.

0, 1, and 2 are terminals.

The start symbol is A.

The rules are:

- $A \rightarrow 0A|1C|2B|0$
- $B \rightarrow 0B|1A|2C|1$
- $C \rightarrow 0C|1B|2A|2$

Can this be parsed?

2011020

## Grammar

A, B, and C are non-terminals.

0, 1, and 2 are terminals.

The start symbol is A, the rules are:

- $A \rightarrow 0A|1C|2B|0$
- $B \rightarrow 0B|1A|2C|1$
- $C \rightarrow 0C|1B|2A|2$

Can this be parsed?

1112202

00102

2120

## Syntactic Analysis

$prog \rightarrow \{ statement \}^+$

The "...+" means "one or more repetitions of the items in..."

$statement \rightarrow assignment \mid loop \mid io$

In this example,  
"=", "while",  
"(", and ")"  
are terminals

$assignment \rightarrow id = expression$

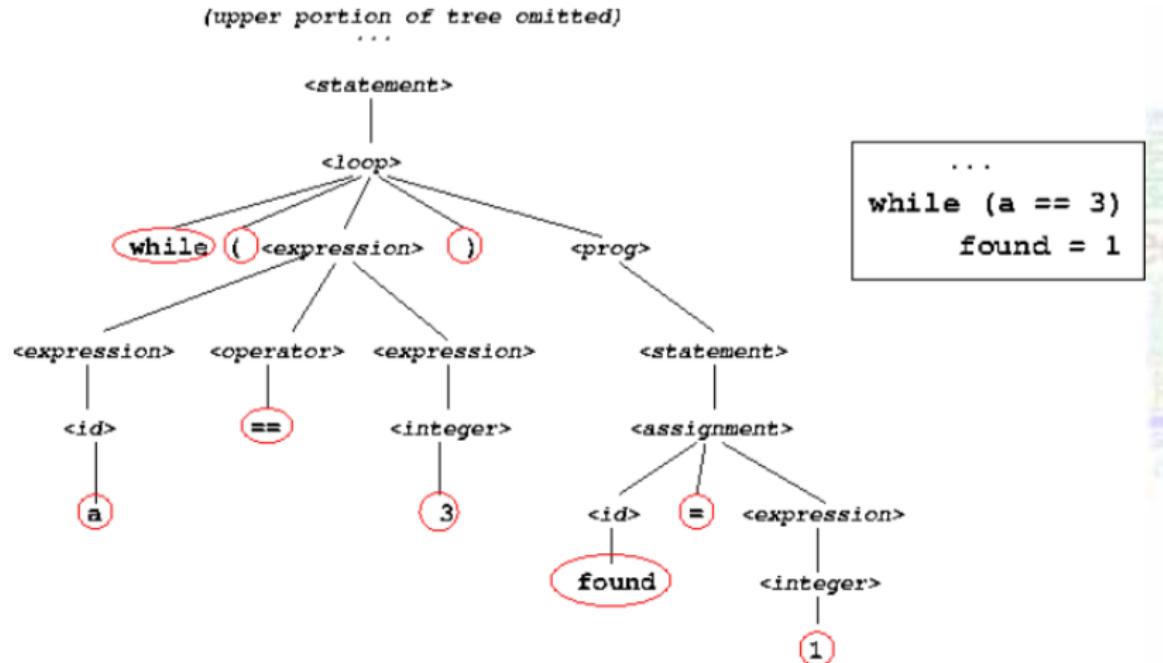
$loop \rightarrow while ( expression ) prog$

"A program is one or more statements."

"A statement is an assignment, a loop, or an input/output command."

"An assignment is an identifier, followed by "=", followed by an expression."

# Sample Parse Tree (portion)



# Sample Parse Tree (failed)

(upper portion of tree omitted)  
...



...

```
while x <= 10
      a = x+1
```

Parse halts after "while" -- unable to match the "(" in  
the rule with a "(" in the input. An error is reported by  
the compiler.

- Overview of notation used:

<https://docs.oracle.com/javase/specs/jls/se8/html/jls-2.html>

- The full syntax grammar:

<https://docs.oracle.com/javase/specs/jls/se8/html/jls-19.html>

So far, we have looked at:

- the scanner (lexical analysis)–tokenizes input
- the parser (syntactic analysis)–validates structure
- Next, we do **semantic analysis**: is the program meaningful?

## EXAMPLE:

In Java,

```
int i, i, i;
```

has the right structure for a declaration, but it's not legal to redeclare *i* within the same block of code.

- During lexical analysis and parsing, as we process tokens we gather the user-defined names into a **symbol table**.
- Symbol table contains information such as:
  - where the symbol first appeared (usually in a declaration)
  - whether it has an initial value (parsing will tell us this)
  - what its type is (parsing tells us), etc.

# Symbol Table

<u>Line#</u>	<u>Program source</u>
...	...
20	int i = 10;
21	double j = 3;
22	int k;
...	...

Entry	Symbol	Declared at	Type	Init Value
...	...	...	...	...
5	int	<library>.103	--	(TYPE)
6	double	<library>.142	--	(TYPE)
7	i	<Lab3>.20	5	10
8	j	<Lab3>.21	6	3.0
9	k	<Lab3>.22	5	uninit
...	...	...	...	...

As the parser encounters names, it looks them up to see if they are already declared; if not, it creates a table entry. (Some names are pre-declared as part of the language.)

# Symbol Table

Line#    Program source

```
...  
20 int i = 10;  
21 double j = 3;  
22 int k;  
23 k = i*j;
```

Entry	Symbol	Declared at	Type	Init Value
...	...	...	...	...
5	int	<library>.103	--	(TYPE)
6	double	<library>.142	--	(TYPE)
7	i	<Lab3>.20	5	10
8	j	<Lab3>.21	6	3.0
9	k	<Lab3>.22	5	uninit
...	...	...	...	...

```
test.java:23: error: possible loss of precision:  
k = i*j;  
          ^  
required: int, found: double
```

# Symbol Table

Line#	Program source
...	...
20	int i = 10;
21	double j = 3;
22	int k;
23	k = i*j;
24	<b>int j;</b>
	<b>test.java:24: error: variable j is already defined in method main (String[])</b>
	int j;
	^

Entry	Symbol	Declared at	Type	Init Value
...	...	...	...	...
5	int	<library>.103	--	(TYPE)
6	double	<library>.142	--	(TYPE)
7	i	<Lab3>.20	5	10
8	j	<Lab3>.21	6	3.0
9	k	<Lab3>.22	5	uninit
...	...	...	...	...

# Symbol Table

Line#	Program source
20	int i = 10;
21	double j = 3;
22	int k;
23	k = i*j;
<b>24</b>	<b>int j;</b>
<b>25</b>	<b>m = 17;</b>

test.java:25: error: cannot find symbol

m = 17;

^

symbol: variable m

Entry	Symbol	Declared at	Type	Init Value
...	...	...	...	...
5	int	<library>.103	--	(TYPE)
6	double	<library>.142	--	(TYPE)
7	i	<Lab3>.20	5	10
8	j	<Lab3>.21	6	3.0
9	k	<Lab3>.22	5	uninit
...	...	...	...	...

- During this phase, the parsed program is converted into a simpler, step-by-step description in some **intermediate language**.
- Intermediate language may exist only as an internal representation within the compiler—it does not need be an “actual” language).
- A simple example is something called “three-address code.”

In “Three-address code”, everything is reduced to simple operations on two values, with the result stored in a (possibly temporary) variable.

## Source code:

```
double f = 3 + i * 7;
```

## Three-address code:

```
temp1 = i
```

```
temp2 = 7
```

```
temp1 = temp1 * temp2
```

```
temp2 = 3
```

```
temp1 = temp2 + temp1
```

```
f = int-to-float(temp1)
```

# Why Intermediate Code Generation?

Simple representation, easy to manipulate,  
close to assembly/machine language. E.g.,

temp1 = x

temp2 = 10

temp1 = temp1 + temp2

**Three-address code**

lw \$t0, x  
li \$t1, 10  
add \$t0, \$t0, \$t1

**MIPS assembly language**

## Many levels, from very basic:

`temp1 = temp1 + 1`

`temp1 = temp1 + x`

`temp1 = temp1 + 5`



`temp1 = temp1 + x`

`temp1 = temp1 + 6`

**“constant folding”**

... to more complex:

**temp1 = 1**

temp2 = x

temp3 = temp2+9

temp2 = temp2\*temp3

temp1 = 8



temp2 = x

temp3 = temp2+9

temp2 = temp2\*temp3

temp1 = 8

**“dead code removal”**

Here is a list of the many, many kinds of optimizations:

<http://www.compileroptimizations.com/>

(The examples show the effects on source code, but the optimizations are usually made on the intermediate code.)

- It is usually very straightforward to generate machine instructions from intermediate code, since the intermediate code is simple.
- Some further machine-specific optimizations may take place during or after this stage.

# Pipelining

- The steps in compilation don't need to be done in whole phases, one after the other, but can be "**pipelined**":
- `int count = 1;`  
create tokens, pass to parser, generate some intermediate code
- `j = j + count;`  
create tokens, pass to parser, generate some intermediate code  
... etc. ...

# Reading Assignment

PLP Chapter 01, Section 1.4, 1.6