

Lexer

COP-3402 Systems Software

Paul Gazzillo

1	position	...
2	initial	...
3	rate	...

SYMBOL TABLE

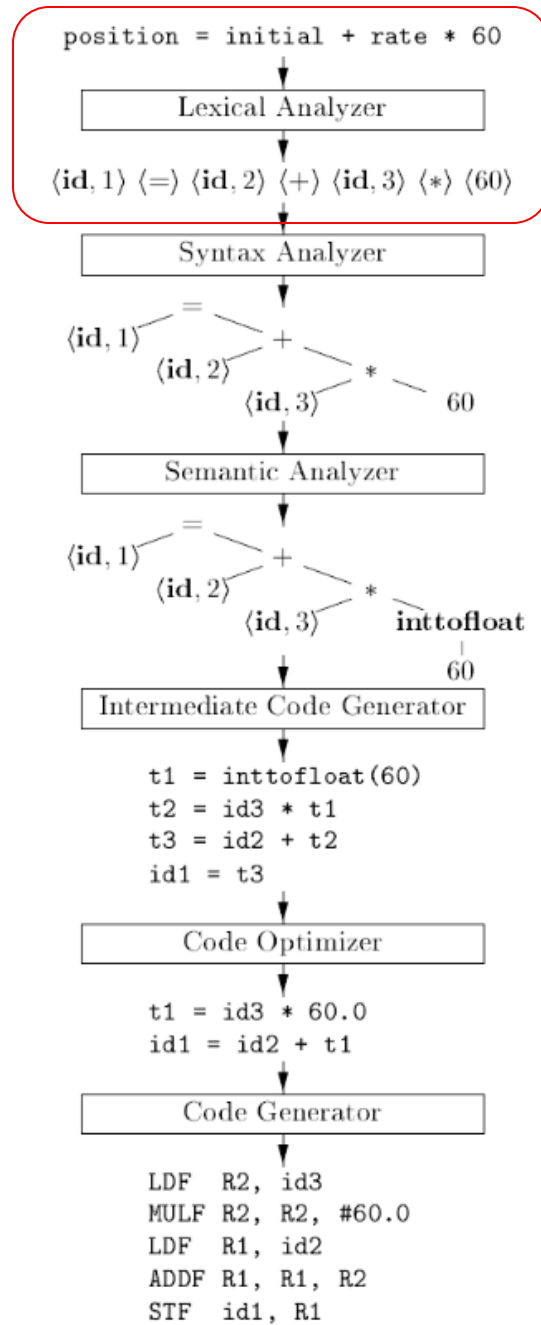


Figure 1.7: Translation of an assignment statement

Core Idea: Convert Characters to Tokens

- Regular expressions *specify* tokens
 - As sets of strings
- Finite automata *recognize* tokens
 - Via state transitions
- A lexer *produces* a stream of tokens
 - From a stream of individual characters

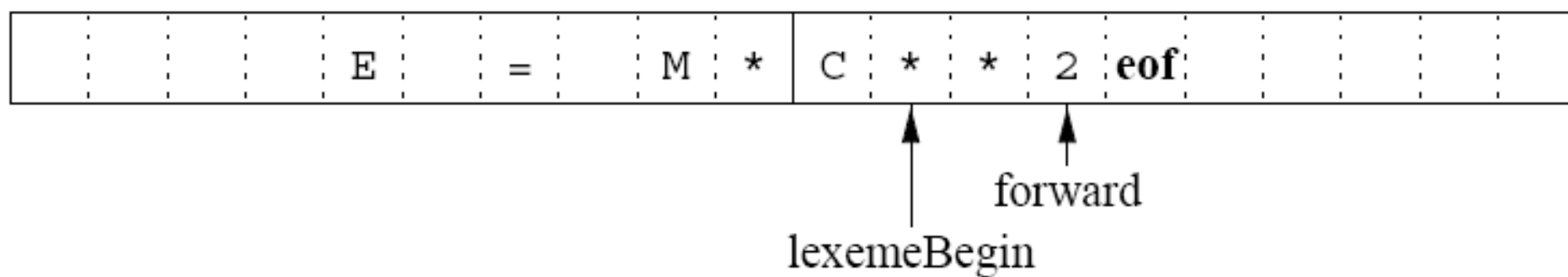


Figure 3.3: Using a pair of input buffers

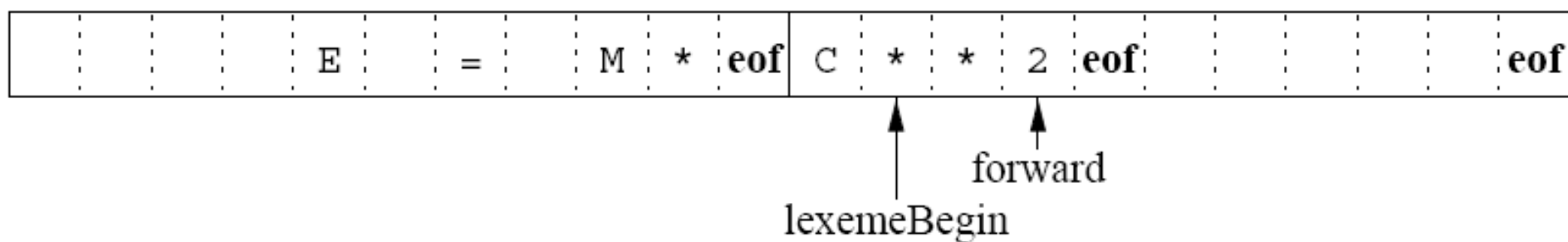


Figure 3.4: Sentinels at the end of each buffer

TOKEN	INFORMAL DESCRIPTION	SAMPLE LEXEMES
if	characters <code>i</code> , <code>f</code>	<code>if</code>
else	characters <code>e</code> , <code>l</code> , <code>s</code> , <code>e</code>	<code>else</code>
comparison	<code><</code> or <code>></code> or <code><=</code> or <code>>=</code> or <code>==</code> or <code>!=</code>	<code><=</code> , <code>!=</code>
id	letter followed by letters and digits	<code>pi</code> , <code>score</code> , <code>D2</code>
number	any numeric constant	<code>3.14159</code> , <code>0</code> , <code>6.02e23</code>
literal	anything but <code>"</code> , surrounded by <code>"</code> 's	<code>"core dumped"</code>

Figure 3.2: Examples of tokens

LEXEMES	TOKEN NAME	ATTRIBUTE VALUE
Any <i>ws</i>	–	–
if	if	–
then	then	–
else	else	–
Any <i>id</i>	id	Pointer to table entry
Any <i>number</i>	number	Pointer to table entry
<	relop	LT
<=	relop	LE
=	relop	EQ
<>	relop	NE
>	relop	GT
>=	relop	GE

Figure 3.12: Tokens, their patterns, and attribute values

Implementing a Lexer

- Regex compiler
 - Regex -> NFA -> DFA -> state table
- Hand-coded lexer
 - Lexer reads one character at a time, e.g., fgetc
 - Concatenation is sequence of statements
 - Union is a conditional
 - Closure is a while loop

Concatenation is sequence of statements, e.g., ab

```
char c;  
c = fgetc(lexerin);  
assert 'a' == c  
c = fgetc(lexerin);  
assert 'b' == c
```

Union is a conditional, e.g., $a|b$

```
char c;  
c = fgetc(lexerin);  
if ('a' == c) {  
    // ...  
} else if ('b' == c) {  
    // ...  
} else {  
    error(...)  
}
```

Closure is a while loop, e.g., a^*

```
char c;  
c = fgetc(lexerin);  
while('a' == c) {  
    c = fgetc(lexerin);  
}  
// c is now the next character
```


Demo

Lexing PL/0

EBNF Standardizes Language Specifications

- EBNF = Extended Backus-Naur Format
- Different notation for regular expressions
 - a^* $\rightarrow \{ a \}$
 - $a|\epsilon$ $\rightarrow [a]$
 - ab $\rightarrow a b$
 - $a | b$ $\rightarrow a | b$

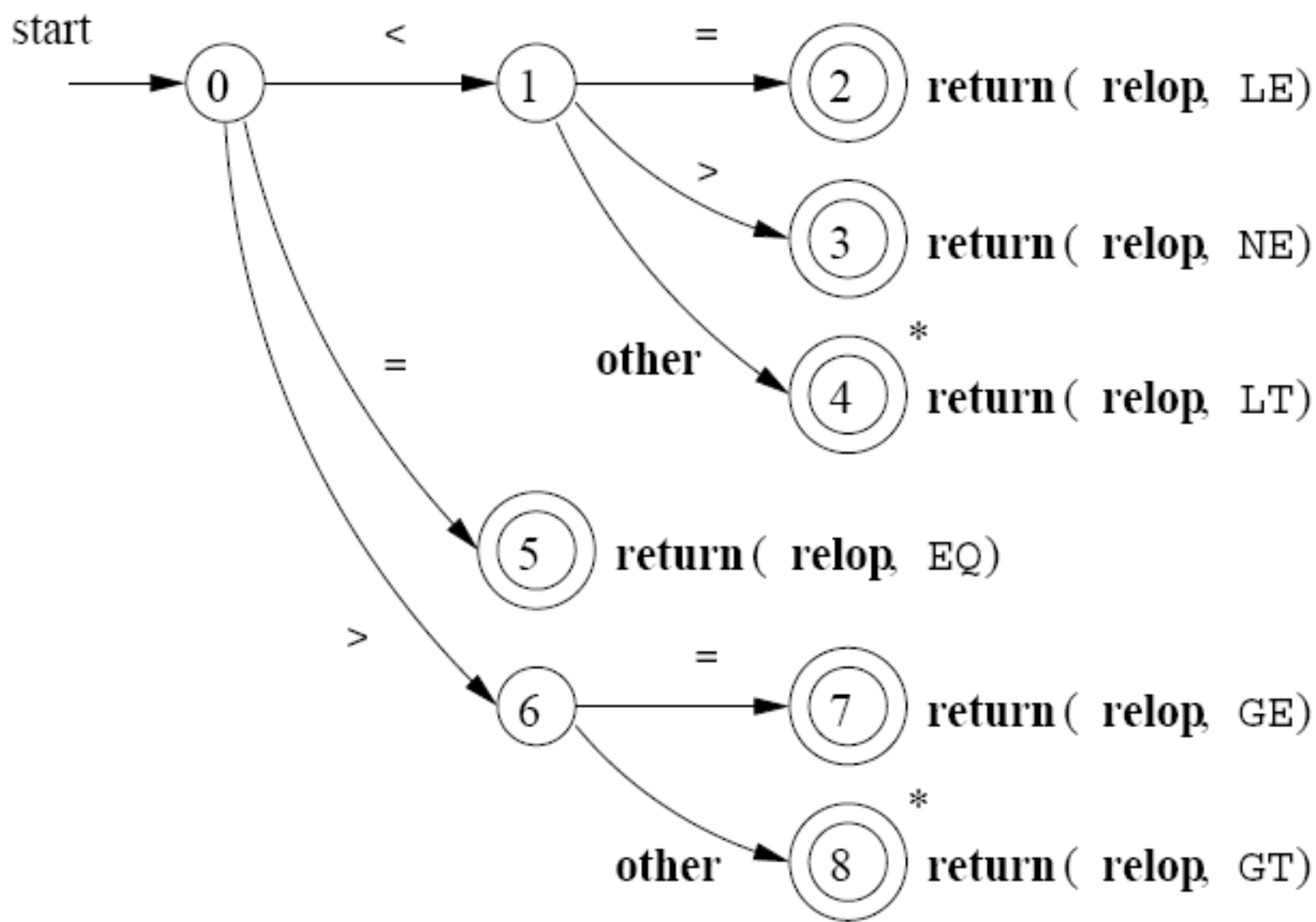


Figure 3.13: Transition diagram for **relop**