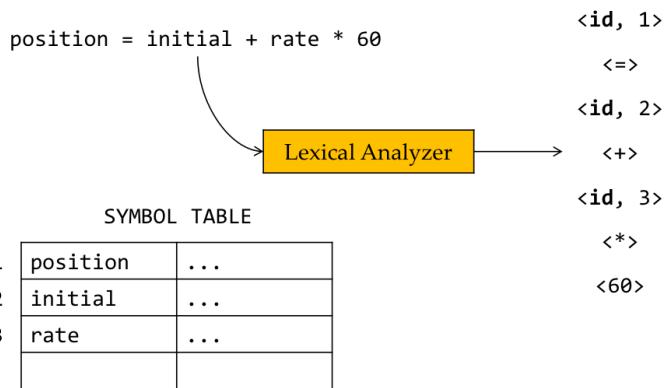
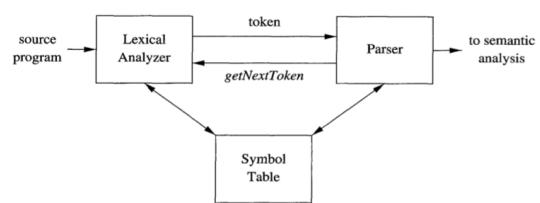


The Role of Lexical Analyzer

- Read the input characters of the source program, group them into lexemes, and produces a sequence of tokens
- Add lexemes into the symbol table when necessary

读入源程序的输入字符，生成 token 序。



Tokens, Patterns, and Lexemes

- A **lexeme** is a string of characters that is a lowest-level syntactic unit in programming languages
- A **token** is a syntactic category representing a class of lexemes. Formally, it is a pair **<token name, attribute value>**
 - **Token name**: an abstract symbol representing the kind of the token
 - **Attribute value** (optional) points to the symbol table
- Each token has a particular **pattern**: a description of the form that the lexemes of the token may take

词素是最小的语法单元

token是一类词素

每个token都有特定的模式，以描述其类中的 lexeme

Examples

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

Consider the C statement: printf("Total = %d\n", score);

Lexeme	printf	score	"Total = %d\n"	(...
Token	id	id	literal	left_parenthesis	...

Attributes for Tokens

- When more than one lexeme match a pattern, the lexical analyzer must provide additional information, named **attribute values**, to the subsequent compiler phases
 - Token names influence parsing decisions
 - Attribute values influence semantic analysis, code generation etc.
- For example, an **id** token is often associated with: (1) its lexeme, (2) type, and (3) the location at which it is first found. Token attributes are stored in the **symbol table**.

`A = B * 2` →

- <id, pointer to symbol-table entry for A>
- <assign_op>
- <id, pointer to symbol-table entry for B>
- <mult_op> <number, integer value 2>

token → parsing
attribute value → 语义分析, 代码生成

Lexical Errors

- When none of the patterns for tokens match any prefix of the remaining input
- Example: `int 3a = a * 3;`

Lexical errors and syntax errors in Java (learn by yourself):

- <https://www.javatpoint.com/lexical-error>
- <https://www.javatpoint.com/syntax-error>

当没有匹配上任何 token 的模式时会报语法错误

Specification of Tokens

- Regular expression** (正则表达式, **regexp for short**) is an important notation for specifying lexeme patterns
- Content of this part
 - Strings and Languages (串和语言)
 - Operations on Languages (语言上的运算)
 - Regular Expressions
 - Regular Definitions (正则定义)
 - Extensions of Regular Expressions

Strings and Languages

- Alphabet** (字母表): any **finite** set of symbols
 - Examples of symbols: **letters**, **digits**, and **punctuations**
 - Examples of alphabets: {1, 0}, ASCII, Unicode
- A **string** (串) over an alphabet is a **finite** sequence of symbols drawn from the alphabet
 - The length of a string s , denoted $|s|$, is the number of symbols in s (i.e., cardinality)
 - Empty string** (空串): the string of length 0, ϵ

Terms : banana

- **Prefix (前缀)** of string s : any string obtained by removing 0 or more symbols from the end of s ([ban](#), [banana](#), ϵ)
- **Proper prefix (真前缀)**: a prefix that is not ϵ and not s itself ([ban](#))
- **Suffix (后缀)**: any string obtained by removing 0 or more symbols from the beginning of s ([nana](#), [banana](#), ϵ).
- **Proper suffix (真后缀)**: a suffix that is not ϵ and not equal to s itself ([nana](#))
- **Substring (子串) of s** : any string obtained by removing any prefix and any suffix from s ([banana](#), [nan](#), ϵ)
- **Proper substring (真子串)**: a substring that is not ϵ and not equal to s itself ([nan](#))
- **Subsequence (子序列)**: any string formed by removing 0 or more not necessarily consecutive symbols from s ([bnn](#))



How many substrings & subsequences does [banana](#) have?

(Two substrings are different as long as they have different start/end index)

String Operations

- **Concatenation (连接)**: the concatenation of two strings x and y , denoted xy , is the string formed by appending y to x
 - $x = \text{dog}$, $y = \text{house}$, $xy = \text{doghouse}$
- **Exponentiation (幂/指数运算)**: $s^0 = \epsilon$, $s^1 = s$, $s^i = s^{i-1}s$
 - $x = \text{dog}$, $x^0 = \epsilon$, $x^1 = \text{dog}$, $x^3 = \text{dogdogdog}$

Language

- A **language** is any **countable set**¹ of strings over some fixed alphabet
 - The set containing only the empty string, that is $\{\epsilon\}$, is a language, denoted \emptyset
 - The set of all [grammatically correct English sentences](#)
 - The set of all [syntactically well-formed C programs](#)

¹In mathematics, a countable set is a set with the same cardinality (number of elements) as some subset of the set of natural numbers. A countable set is either a finite set or a countably infinite set.

Operations on Language

- 并, 连接, Kleene闭包, 正闭包



对串 s , 若 $|s| = n$, 则有
前(后)缀 $n+1$ 个, 其前(后)缀 $n-1$ 个
子串 2^n 个
子序列 $\sum_{i=0}^n C_n^i = 2^n$

子序列可以不连续移除

字符表固定、串可数。

Kleene闭包中含有空串而正闭包不一定有。

Examples

- $L = \{A, B, \dots, Z, a, b, \dots, z\}$
- $D = \{0, 1, \dots, 9\}$

$L \cup D$	$\{A, B, \dots, Z, a, b, \dots, z, 0, 1, \dots, 9\}$
LD	the set of 520 strings of length two, each consisting of one letter followed by one digit
L^4	the set of all 4-letter strings
L^*	the set of all strings of letters, including ϵ
$L(L \cup D)^*$?
D^+	?

Note: L, D might seem to be the alphabets of letters and digits. We define them to be languages: all strings happen to be of length one.

Regular Expressions : For Describing Languages/Patterns

Rules that define regexps over an alphabet Σ :

- **BASIS:** two rules form the basis:
 - ϵ is a regexp, $L(\epsilon) = \{\epsilon\}$
 - If a is a symbol in Σ , then a is a regexp, and $L(a) = \{a\}$
- **INDUCTION:** Suppose r and s are regexps denoting the languages $L(r)$ and $L(s)$
 - $(r)s$ is a regexp denoting the language $L(r) \cup L(s)$
 - $(r)s$ is a regexp denoting the language $L(r)L(s)$
 - $(r)^*$ is a regexp denoting $(L(r))^*$
 - (r) is a regexp denoting $L(r)$. Additional parentheses do not change the language an expression denotes.
- Following the rules, regexps often contain **unnecessary pairs of parentheses**. We may drop some if we adopt the conventions:
 - **Precedence:** closure $*$ > concatenation > union |
 - **Associativity:** All three operators are **left associative**, meaning that operations are grouped from the left, e.g., $a \mid b \mid c$ would be interpreted as $(a \mid b) \mid c$
- Example: $(a) \mid ((b)^*(c)) = a \mid b^*c$

Examples

- Examples: Let $\Sigma = \{a, b\}$
 - $a \mid b$ denotes the language $\{a, b\}$
 - $(a \mid b)(a \mid b)$ denotes $\{aa, ab, ba, bb\}$
 - a^* denotes $\{\epsilon, a, aa, aaa, \dots\}$
 - $(a \mid b)^*$ denotes the set of all strings consisting of 0 or more a 's or b 's: $\{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$
 - $a \mid a^*b$ denotes the string a and all strings consisting of 0 or more a 's and ending in b : $\{a, b, ab, aab, aaab, \dots\}$

多余的括号不会影响表达式的语言。

Regular Language

- A **regular language** is a language that can be defined by a regexp
- If two regexps r and s denote the same language, they are **equivalent**, written as $r = s$
- Each **algebraic law** below asserts that expressions of two different forms are equivalent

LAW	DESCRIPTION
$r s = s r$	is commutative
$r (s t) = (r s) t$	is associative
$r(st) = (rs)t$	Concatenation is associative
$r(s t) = rs rt; (st)r = sr tr$	Concatenation distributes over
$\epsilon r = r\epsilon = r$	ϵ is the identity for concatenation
$r^* = (r \epsilon)^*$	ϵ is guaranteed in a closure
$r^{**} = r^*$	* is idempotent

Is $(a|b)(a|b) = aa|ab|ba|bb$ true?

| can be viewed as + in arithmetics, concatenation can be viewed as \times , * can be viewed as the power operator.

两个正则表达式等价 \Leftrightarrow 两个正则表达式表示同一种语言

I 可视作加法
II 可视作乘法
* 可视作幂

Regular Definitions

- For **notational convenience**, we can give names to certain regexps and use those names in subsequent expressions

If Σ is an alphabet of basic symbols, then a **regular definition** is a sequence of definitions of the form:

$$\begin{aligned} d_1 &\rightarrow r_1 \\ d_2 &\rightarrow r_2 \\ \dots \\ d_n &\rightarrow r_n \end{aligned}$$

where:

- Each d_i is a new symbol not in Σ and not the same as the other d 's
- Each r_i is a regexp over the alphabet $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$

Each new symbol denotes a regular language. The second rule means that you may reuse previously-defined symbols.

Examples

- Regular definition for C identifiers

$$\begin{aligned} \text{letter_} &\rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z \mid - && \text{hello valid?} \\ \text{digit} &\rightarrow 0 \mid 1 \mid \dots \mid 9 \\ \text{id} &\rightarrow \text{letter_} (\text{letter_} \mid \text{digit})^* && 3\text{times valid?} \end{aligned}$$

- Regexp for C identifiers

$$(A|B|\dots|Z|a|b|\dots|z|_-)((A|B|\dots|Z|a|b|\dots|z|_-)|(0|1|\dots|9))^*$$

Extensions of Regular Expressions

- Basic operators:** union |, concatenation, and Kleene closure (proposed by Kleene in 1950s)
- A few **notational extensions:**
 - One or more instances:** the unary, postfix operator *
 - $r^+ = rr^*, r^* = r^+ \mid \epsilon$
 - Zero or one instance:** the unary postfix operator ?
 - $r? = r \mid \epsilon$
 - Character classes:** shorthand for a logical sequence
 - $[a_1a_2\dots a_n] = a_1 \mid a_2 \mid \dots \mid a_n$
 - $[a-e] = a \mid b \mid c \mid d \mid e$
- The extensions are **only for notational convenience**, they do not change the descriptive power of regexps

Recognition of Tokens

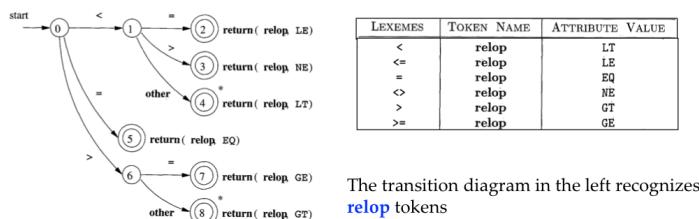
- Lexical analyzer examines the input string and finds a prefix that matches one of the tokens
- The first thing when building a lexical analyzer is to define the patterns of tokens using regular definitions
- A special token:** $\text{ws} \rightarrow (\text{blank} \mid \text{tab} \mid \text{newline})^*$
 - When the lexical analyzer recognizes a **whitespace token**, it does not return it to the parser, but restart from the next character

Example : Patterns and Tokens

$\begin{array}{l} \text{digit} \rightarrow [0-9] \\ \text{digits} \rightarrow \text{digit}^+ \\ \text{number} \rightarrow \text{digits} (. \text{ digits})? (\text{ E } [+ -]?) \text{ digits })? \\ \text{letter} \rightarrow [\text{A-Za-z}] \\ \text{id} \rightarrow \text{letter} (\text{letter} \mid \text{digit})^* \\ \text{if} \rightarrow \text{if} \\ \text{then} \rightarrow \text{then} \\ \text{else} \rightarrow \text{else} \\ \text{rellop} \rightarrow < \mid > \mid \leq \mid \geq \mid = \mid \neq \end{array}$ <p>Patterns for tokens</p>	<table border="1"> <thead> <tr> <th>LEXEMES</th><th>TOKEN NAME</th><th>ATTRIBUTE VALUE</th></tr> </thead> <tbody> <tr> <td>Any ws</td><td>-</td><td>-</td></tr> <tr> <td>if</td><td>if</td><td>-</td></tr> <tr> <td>then</td><td>then</td><td>-</td></tr> <tr> <td>else</td><td>else</td><td>-</td></tr> <tr> <td>Any id</td><td>id</td><td>Pointer to table entry</td></tr> <tr> <td>Any number</td><td>number</td><td>Pointer to table entry</td></tr> <tr> <td><</td><td>rellop</td><td>LT</td></tr> <tr> <td>\leq</td><td>rellop</td><td>LE</td></tr> <tr> <td>=</td><td>rellop</td><td>EQ</td></tr> <tr> <td>\neq</td><td>rellop</td><td>NE</td></tr> <tr> <td>></td><td>rellop</td><td>GT</td></tr> <tr> <td>\geq</td><td>rellop</td><td>GE</td></tr> </tbody> </table> <p>Lexemes, tokens, and attribute values</p>	LEXEMES	TOKEN NAME	ATTRIBUTE VALUE	Any ws	-	-	if	if	-	then	then	-	else	else	-	Any id	id	Pointer to table entry	Any number	number	Pointer to table entry	<	rellop	LT	\leq	rellop	LE	=	rellop	EQ	\neq	rellop	NE	>	rellop	GT	\geq	rellop	GE
LEXEMES	TOKEN NAME	ATTRIBUTE VALUE																																						
Any ws	-	-																																						
if	if	-																																						
then	then	-																																						
else	else	-																																						
Any id	id	Pointer to table entry																																						
Any number	number	Pointer to table entry																																						
<	rellop	LT																																						
\leq	rellop	LE																																						
=	rellop	EQ																																						
\neq	rellop	NE																																						
>	rellop	GT																																						
\geq	rellop	GE																																						

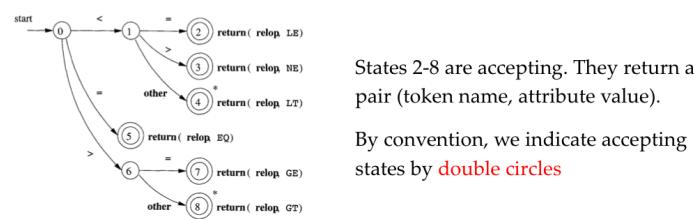
Transition Diagrams

- An important step in constructing a lexical analyzer is to convert patterns into “**transition diagrams**”
- Transition diagrams have a collection of nodes, called **states** (状态) and **edges** (边) directed from one node to another



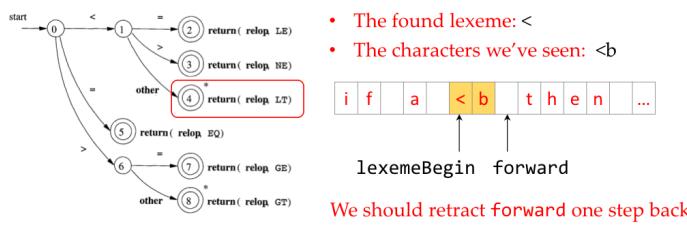
States

- Represent conditions that could occur during the process of scanning (i.e., what characters we have seen)
- The **start state** (开始状态), or **initial state**, is indicated by an edge labeled “**start**”, which enters from nowhere
- Certain states are said to be **accepting** (接受状态), or **final**, indicating that a lexeme has been found



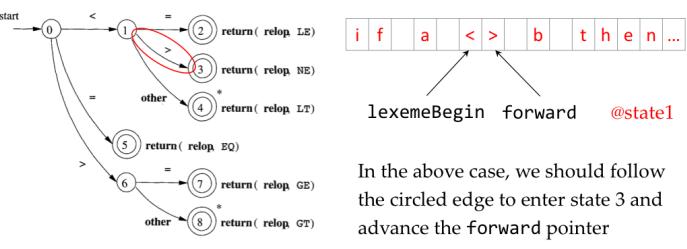
The Retract Action

- At certain accepting states, the found lexeme may not contain all characters that we have seen from the start state (such states are annotated with *)
- When entering * states, it is necessary to **retract** (回退) the forward pointer, which points to the next char in the input string



Edges

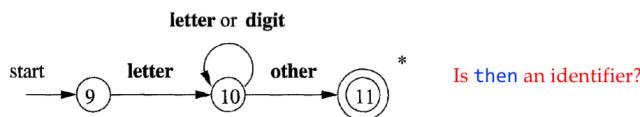
- Edges** are directed from one state to another
- Each edge is labeled by a symbol or set of symbols



Recognition of Reserved Words and Identifiers

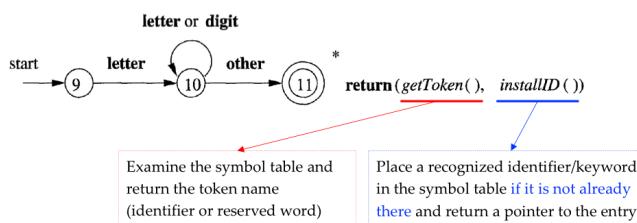
保留字和标识符需要区分

- In many languages, **reserved words** or **keywords** (e.g., then) also match the pattern of identifiers
- Problem:** the transition diagram that searches for identifiers can also recognize reserved words

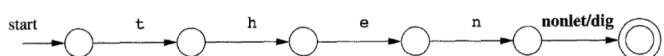


Handling Reserved Words

- Strategy 1:** Preinstall the reserved words in the symbol table. Put a field in the symbol-table entries to indicate that these strings are not ordinary identifiers (预先存表方案)



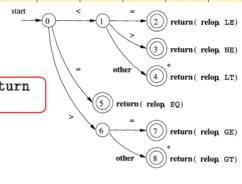
- Strategy 2:** Create a separate transition diagram with a **high priority** for each keyword (多状态转移图方案)



Building a Lexical Analyzer from transition Diagrams

```
TOKEN getRelop()
{
    TOKEN retToken = new(RELOP);
    while(1) { /* repeat character processing until a return
        or failure occurs */
        switch(state) {
            case 0: c = nextChar();
                if (c == '<') state = 1;
                else if (c == '=') state = 5;
                else if (c == '>') state = 6;
                else fail(); /* lexeme is not a relop */
                break;
            case 1: ...
            ...
            case 8: retract();
                retToken.attribute = GT;
                return(retToken);
        }
    }
}
```

Sketch implementation of relop transition diagram



TOKEN getRelop()

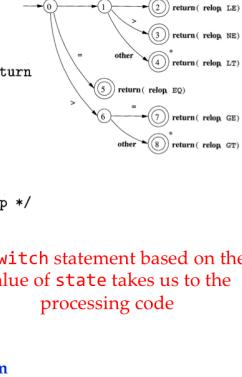
```
TOKEN retToken = new(RELOP);
while(1) { /* repeat character processing until a return
        or failure occurs */
    switch(state) {
        case 0: c = nextChar();
            if (c == '<') state = 1;
            else if (c == '=') state = 5;
            else if (c == '>') state = 6;
            else fail(); /* lexeme is not a relop */
            break;
        case 1: ...
        ...
        case 8: retract();
            retToken.attribute = GT;
            return(retToken);
    }
}
```

Use a variable **state** to record the current state

Sketch implementation of relop transition diagram

```
TOKEN getRelop()
{
    TOKEN retToken = new(RELOP);
    while(1) { /* repeat character processing until a return
        or failure occurs */
        switch(state) {
            case 0: c = nextChar();
                if (c == '<') state = 1;
                else if (c == '=') state = 5;
                else if (c == '>') state = 6;
                else fail(); /* lexeme is not a relop */
                break;
            case 1: ...
            ...
            case 8: retract();
                retToken.attribute = GT;
                return(retToken);
        }
    }
}
```

A **switch statement based on the value of state takes us to the processing code**



TOKEN getRelop()

```
TOKEN retToken = new(RELOP);
while(1) { /* repeat character processing until a return
        or failure occurs */
    switch(state) {
        case 0: c = nextChar();
            if (c == '<') state = 1;
            else if (c == '=') state = 5;
            else if (c == '>') state = 6;
            else fail(); /* lexeme is not a relop */
            break;
        case 1: ...
        ...
        case 8: retract();
            retToken.attribute = GT;
            return(retToken);
    }
}
```

The code of a normal state:
 1. Read the next character
 2. Determine the next state
 3. If step 2 fails, do error recovery

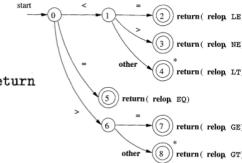
Sketch implementation of relop transition diagram

```
TOKEN getRelop()
{
    TOKEN retToken = new(RELOP);
    while(1) { /* repeat character processing until a return
        or failure occurs */
        switch(state) {
            case 0: c = nextChar();
                if (c == '<') state = 1;
                else if (c == '=') state = 5;
                else if (c == '>') state = 6;
                else fail(); /* lexeme is not a relop */
                break;
            case 1: ...
            ...
            case 8: retract();
                retToken.attribute = GT;
                return(retToken);
        }
    }
}
```

The code of an accepting state:

1. Perform retraction if the state has *
2. Set token attribute values
3. Return the token to parser

Sketch implementation of relop transition diagram



Building the Entire Lexical Analyzer

- **Strategy 1:** Try the transition diagram for each type of token sequentially
 - `fail()` resets the pointer forward and tries the next diagram
- **Problem:** Not efficient
 - May need to try many irrelevant diagrams whose first edge does not match the first character in the input stream
- **Strategy 2:** Run transition diagrams in parallel
 - Need to resolve the case where one diagram finds a lexeme and others are still able to process input.
 - **Solution:** take the longest prefix of the input that matches any pattern
- **Problem:** Requires special hardware for parallel simulation, may degenerate into the sequential strategy on certain machines

- **Strategy 3:** Combining all transition diagrams into one
 - Allow the transition diagram to read input until there is no possible next state
 - Take the longest lexeme that matched any pattern
- This is a commonly-adopted strategy in real-world compiler implementation (efficient & requires no special hardware)



How? Be patient ☺, we will talk about this later.

The Lexical Analyzer Generator Lex

- Lex, or a more recent tool Flex, allows one to specify a lexical analyzer by specifying regexps to describe patterns for tokens
- Often used with Yacc/Bison to create the frontend of compiler

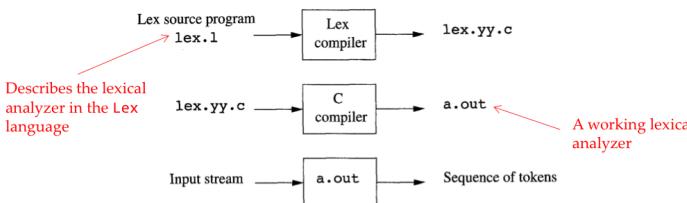


Figure 3.22: Creating a lexical analyzer with Lex

Structure of Lex Programs

- A Lex program has three sections separated by %%
 - Declaration (声明)
 - Variables, constants (e.g., token names)
 - Regular definitions
 - Translation rules (转换规则) in the form “Pattern {Action}”
 - Each pattern (模式) is a regexp (may use the regular definitions of the declaration section)
 - Actions (动作) are fragments of code, typically in C, which are executed when the pattern is matched
 - Auxiliary functions section (辅助函数)
 - Additional functions that can be used in the actions

Example

```
/*
 * definitions of manifest constants
 LT, LE, EQ, NE, GT, GE,
 IF, THEN, ELSE, ID, NUMBER, RELOP */
%
```

```
/* regular definitions */
delim  [\t\n]
ws     [delim]+
letter [A-Za-z]
digit  [0-9]
id     {[letter]({digit})*|{digit})*
number {digit}(\.{digit})+?(E[+-]?{digit})?
```

```
%%
```

Anything in between %{ and }% is copied directly to lex.yy.c.
In the example, there is only a comment, not real C code to define manifest constants

Regular definitions that can be used in translation rules

Section separator

```

{ws}           /* no action and no return */
if            {return(IF);}
then          {return(THEN);}
else          {return(ELSE);}
{id}           {yyval = (int) installID(); return(ID);}
{number}       {yyval = (int) installNum(); return(NUMBER);}

"<"           {yyval = LT; return(RELOP);}
"<="          {yyval = LE; return(RELOP);}
"="           {yyval = EQ; return(RELOP);}
">="           {yyval = NE; return(RELOP);}
">"            {yyval = GT; return(RELOP);}
">="          {yyval = GE; return(RELOP;)}

%%
```

Continue to recognize other tokens

Return token name to the parser

Place the lexeme found in the symbol table

A global variable that stores a pointer to the symbol table entry for the lexeme. Can be used by the parser or a later component of the compiler.

* The characters inside have no special meaning (even if it is a special one such as *).

- Everything in the auxiliary function section is copied directly to the file `lex.yy.c`
- Auxiliary functions may be used in actions in the translation rules

```
int installID() /* function to install the lexeme, whose  
    first character is pointed to by yytext,  
    and whose length is yylen, into the  
    symbol table and return a pointer  
    thereto */  
}  
                                Variables defined and set automatically  
                                by the lexical analyzer Lex generates  
int installNum() /* similar to installID, but puts numer-  
    ical constants into a separate table */  
}
```

Conflict Resolution

- When the generated lexical analyzer runs, it analyzes the input looking for **prefixes that match any of its patterns.***
- **Rule 1:** If it finds multiple such prefixes, it takes the **longest** one
 - The analyzer will treat `<=` as a single lexeme, rather than `<` as one lexeme and `=` as the next
- **Rule 2:** If it finds a prefix matching different patterns, **the pattern listed first** in the Lex program is chosen.
 - If the keyword patterns are listed before identifier pattern, the lexical analyzer will not recognize keywords as identifiers

找匹配到最长的

找第一个匹配到的