

廈門大學



软件学院

《编译技术》实验报告

题 目 标识符和实数的识别

姓 名 陈澄

学 号 32420212202930

班 级 软工三班

实验时间 2024/03/19

2024 年 03 月 19 日

1 实验目的

基本掌握计算机语言的词法分析程序的开发方法。

2 实验内容

编制一个能够分析三种整数、标识符、主要运算符和主要关键字的词法分析程序。

3 实验要求

1. 根据以下的正规式，编制正规文法，画出状态图；

基本要求：

标识符 $\langle \text{字母} \rangle (\langle \text{字母} \rangle | \langle \text{数字字符} \rangle)^*$

十进制整数 $(1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^* | 0$

八进制整数 $0o(0|1|2|3|4|5|6|7)(0|1|2|3|4|5|6|7)^*$

十六进制整数

$0x(0|1|2|3|4|5|6|7|8|9|a|b|c|d|e|f)(0|1|2|3|4|5|6|7|8|9|a|b|c|d|e|f)^*$

运算符和分隔符 $+ \ - \ * \ / \ > \ < \ = \ (\) \ ;$

关键字 `if then else while do`

附加要求：

标识符 $\langle \text{字母} \rangle (\langle \text{字母} \rangle | \langle \text{数字字符} \rangle | \epsilon | _ | .)^* (\langle \text{字母} \rangle | \langle \text{数字字符} \rangle)$

十进制整数 $(1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^* | 0 (\epsilon | .$
 $(0|1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^*)$

八进制整数 $0o(0|1|2|3|4|5|6|7)(0|1|2|3|4|5|6|7)^* (\epsilon | .$
 $(0|1|2|3|4|5|6|7)(0|1|2|3|4|5|6|7)^*)$

十六进制整数

$0x(0|1|2|3|4|5|6|7|8|9|a|b|c|d|e|f)(0|1|2|3|4|5|6|7|8|9|a|b|c|$
 $d|e|f)^* (\epsilon | .$
 $(0|1|2|3|4|5|6|7|8|9|a|b|c|d|e|f)(0|1|2|3|4|5|6|7|8|9|a|b|c|d|$
 $e|f)^*)$

2. 根据状态图，设计词法分析函数 `int scan()`，完成以下功能：
 - 1) 从键盘读入数据，分析出一个单词。
 - 2) 返回单词种别（用整数表示），
 - 3) 返回单词属性（不同的属性可以放在不同的全局变量中）。
3. 编写测试程序，反复调用函数 `scan()`，输出单词种别和属性。

4 实验环境

PC 微机

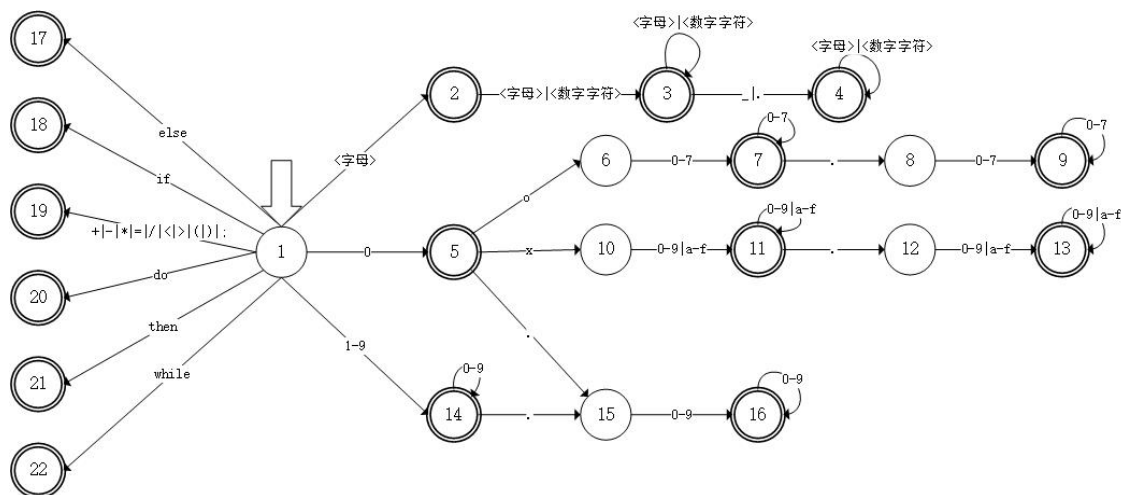
Windows 操作系统

Turbo C 程序集成环境或 Visual C++ 程序集成环境

5 实验步骤

1. 根据状态图，设计词法分析算法。

状态图如下：



设计算法：

采用状态转移算法，大致思路：检查当前状态->根据当前状态将第 i 个字符唯一转化为下一状态->字符串指针+1

伪代码：

```

While(字符串未结束){

    Switch(状态):

    Case 状态 1:{

        Switch(字符):

        Case a:状态 1 转化为状态 2;break;

        Case b:状态 2 转化为状态 3;break;

        .....

    }

    Case 状态 2:{

        .....

    }

    .....

    字符串指针后移;

}

```

2. 采用 C 语言，设计函数 `scan()`，实现该算法。

由于 `switch...case` 内部嵌套的字符选择逻辑较简单，改用 `if...else` 结构进行选择

```

19 int scan(char s[]) {
20     int status = 1, i = 0;
21     //初始状态为1, 不合法状态为0, i为指针
22     while (s[i] && status) {
23         switch (status)
24         {
25             case 1: {
26                 /*char s_copy[10];
27                 strncpy(s_copy, s, n);
28                 s_copy[n] = '\0';*/
29                 if (strcmp(s, "else") == 0) {
30                     status = 17; i += 3; break;
31                 }
32                 if (strcmp(s, "if") == 0) {
33                     status = 18; i += 1; break;
34                 }
35                 if (strcmp(s, "do") == 0) {
36                     status = 20; i += 1; break;
37                 }
38                 if (strcmp(s, "then") == 0) {
39                     status = 21; i += 3; break;
40                 }
41                 if (strcmp(s, "while") == 0) {
42                     status = 22; i += 4; break;
43                 }
44                 if (s[i] == '0') status = 5;
45                 else if (isAlpha(s[i])) status = 2;
46                 else if (isPunct(s[i])) status = 19;
47                 else if (s[i] >= '1' && s[i] <= '9') status = 14;
48                 else status = 0;
49                 break;
50             }
51             case 2: {
52                 if (isAlpha(s[i]) || isNum(s[i])) status = 3;
53                 else status = 0;
54                 break;
55             }
56             case 3: {
57                 if (isAlpha(s[i]) || isNum(s[i])) status = 3;
58                 else if (s[i] == '.' || s[i] == '_') status = 4;
59                 else status = 0;
60                 break;
61             }
62             case 4: {
63                 if (isAlpha(s[i]) || isNum(s[i])) status = 3;
64                 else status = 0;
65                 break;
66             }
67             case 5: {
68                 if (s[i] == 'e' || s[i] == '0') status = 6;
69                 else if (s[i] == 'x' || s[i] == 'X') status = 10;
70                 else if (s[i] == '.') status = 16;
71                 else status = 0;
72                 break;
73             }
74             case 6: {
75                 if (isInt8(s[i])) status = 7;
76                 else status = 0;
77                 break;
78             }
79             case 7: {
80                 if (isInt8(s[i])) status = 7;
81                 else if (s[i] == '.') status = 8;
82                 else status = 0;
83                 break;
84             }
85             case 8: {
86                 if (isInt8(s[i])) status = 9;
87                 else status = 0;
88                 break;
89             }
90             case 9: {
91                 if (isInt8(s[i])) status = 9;
92                 else status = 0;
93                 break;
94             }
95             case 10: {
96                 if (isInt16(s[i])) status = 11;
97                 else status = 0;
98                 break;
99             }
100             case 11: {
101                 if (isInt16(s[i])) status = 11;
102                 else if (s[i] == '.') status = 12;
103                 else status = 0;
104                 break;
105             }
106             case 12: {
107                 if (isInt16(s[i])) status = 13;
108                 else status = 0;
109                 break;
110             }
111             case 13: {
112                 if (isInt16(s[i])) status = 13;
113                 else status = 0;
114                 break;
115             }
116             case 14: {
117                 if (isNum(s[i])) status = 14;
118                 else if (s[i] == '.') status = 15;
119                 else status = 0;
120                 break;
121             }
122             case 15: {
123                 if (isNum(s[i])) status = 16;
124                 else status = 0;
125                 break;
126             }
127             default: status = 0;
128         }
129         i++;
130     }
131     if (status == 1 || status == 6 || status == 8
132         || status == 10 || status == 12 || status == 15) status = 0;
133     //除去非终结态
134     return status;
135     //将状态返回
136 } //状态转移算法
137 string progressStatus(int status) {

```

3. 编制测试程序（主函数 main）。

主要思路：

Step1:识别前 k 个字符，若满足前 k 个字符可以被识别，而前 k+1 个字符无法被识别，则说明前 k 个字符为一个完整的终结符。

Step2:通过 scan 方法识别该终结符，并输出。

Step3:指针后移，将已经识别过的终结符去除，回到 Step1 重复。

```
//循环读取并识别
int i = 0, j = 0;
char temp[100];
while (s[i..j]) {
    if (!(s[i..j] == ' ' || isPunct(s[i..j]))) {
        temp[j] = s[i..j];
        j++;
        continue;
    }
    temp[j] = '\0';
    cout << temp << "\t" << processStatus(scan(temp)) << endl;
    char punct = 0;
    while (s[i..j] == ' ' || isPunct(s[i..j])) {
        if (isPunct(s[i..j]))punct = s[i..j];
        j++;
    }
    if (punct)cout << punct << "\t_" << endl;
    i += j;
    j = 0;
}
```

4. 调试程序：输入一组单词，检查输出结果。

输入数据 1: 0 92+data> 0x3f 0o0 while a+acc>xx do x=x-1;

```

0  92+data>  0x3f  0o0  while a+acc>xx do x=x-1;
0          INT10
92          INT10
+          -
data       IDN
>          -
0x3f       INT16
0o0        INT8
while      WHILE
a          IDN
+          -
acc        IDN
>          -
xx         IDN
do         DO
x          IDN
=          -
x          IDN
-          -
1          INT10
;          -

```

输入数据 2: $a = 6.2 + a * 0X88.80$;

```

a=6.2+a*0X88.80;
a          IDN
=          -
6.2        INT10
+          -
a          IDN
*          -
0X88.80    INT16
;          -

```

输入数据 3: if $a > b$ then $a = b$ else $a = b - 1 + c$;

```

if a>b then a=b else a=b-1+c;
if          IF
a          IDN
>          -
b          IDN
then       THEN
a          IDN
=          -
b          IDN
else      ELSE
a          IDN
=          -
b          IDN
-          -
1          INT10
+          -
c          IDN
;          -

```

6 实验报告要求

1. 词法的正规式描述

实验要求中已经给出

2. 变换后的正规文法

标识符：

文法 $G = (\{S, A, B, C\}, \{<\text{字母}>, <\text{数字}>, ,, _ \}, P, S)$ 其中 P 由下列产生式生成：

$$S \rightarrow <\text{字母}>A$$
$$A \rightarrow \epsilon$$
$$A \rightarrow <\text{字母}>B | <\text{数字}>B$$
$$B \rightarrow <\text{字母}>B | <\text{数字}>B$$
$$B \rightarrow \epsilon$$
$$B \rightarrow .C | _C$$
$$C \rightarrow <\text{字母}>C | <\text{数字}>C$$
$$C \rightarrow \epsilon$$

十进制整数：

文法 $G = (\{S, A, B, C, D\}, \{<\text{数字}>, ,, \}, P, S)$ 其中 P 由下列产生式生成：

$$S \rightarrow 0A$$
$$A \rightarrow \epsilon$$
$$A \rightarrow .C$$
$$S \rightarrow 1B | 2B | 3B | 4B | 5B | 6B | 7B | 8B | 9B$$
$$B \rightarrow <\text{数字}>B$$

$B \rightarrow \epsilon$

$B \rightarrow .C$

$C \rightarrow \langle \text{数字} \rangle D$

$D \rightarrow \langle \text{数字} \rangle D$

$D \rightarrow \epsilon$

八进制整数:

文法 $G = (\{S, A, B, C, D\}, \{0, 1, 2, 3, 4, 5, 6, 7, \dots\}, P, S)$ 其中 P 由下列产生式生成:

$S \rightarrow 0A$

$A \rightarrow 0B | 1B | 2B | 3B | 4B | 5B | 6B | 7B$

$B \rightarrow 0B | 1B | 2B | 3B | 4B | 5B | 6B | 7B$

$B \rightarrow \epsilon$

$B \rightarrow .C$

$C \rightarrow 0D | 1D | 2D | 3D | 4D | 5D | 6D | 7D$

$D \rightarrow 0D | 1D | 2D | 3D | 4D | 5D | 6D | 7D$

$D \rightarrow \epsilon$

十六进制整数:

文法 $G = (\{S, A, B, C, D\}, \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, a, b, c, d, e, f, \dots\}, P, S)$ 其中 P 由下列产生式生成:

$S \rightarrow 0xA$

$A \rightarrow 0B | 1B | 2B | 3B | 4B | 5B | 6B | 7B | 8B | 9B | aB | bB | cB | dB | eB | fB$

$B \rightarrow 0B | 1B | 2B | 3B | 4B | 5B | 6B | 7B | 8B | 9B | aB | bB | cB | dB | eB | fB$

$B \rightarrow \epsilon$

$B \rightarrow .C$

$C \rightarrow 0D|1D|2D|3D|4D|5D|6D|7D|8B|9B|aB|bB|cB|dB|eB|fB$

$D \rightarrow 0D|1D|2D|3D|4D|5D|6D|7D|8B|9B|aB|bB|cB|dB|eB|fB$

$D \rightarrow \epsilon$

运算符和分隔符：

文法 $G = (\{S\}, \{+, -, *, /, <, >, (,), =, ;\}, P, S)$ 其中 P 由下列产生式生成：

$S \rightarrow +|-|*|/|<|>|()|=|;$

关键字：

文法 $G = (\{S\}, \{\text{while}, \text{do}, \text{if}, \text{else}, \text{then}\}, P, S)$ 其中 P 由下列产生式生成：

$S \rightarrow \text{while}|\text{do}|\text{if}|\text{else}|\text{then}$

3. 状态图

实验步骤中已经给出

4. 词法分析程序的数据结构与算法

实验步骤中已经给出

7 思考题

1. 词法分析能否采用空格来区分单词？

答：不能，单词与单词间不仅可以通过空格也可以通过运算符，换行符，制表符等等连接，因此不能仅仅通过空格来区分单词。

2. 程序设计中哪些环节影响词法分析的效率？如何提高效率？

答：

(1) 数据结构的选择：选择合适的数据结构来存储和组织文本数据对词法分析效率至关重要。例如，使用合适的数据结构来存储单词、标记或者字符流，能够提高词

法分析的效率。常见的数据结构如哈希表、树结构等可以被用来加快查找和匹配过程。

(2)算法的设计：采用高效的算法能够显著提高词法分析的效率。例如，使用适当的字符串匹配算法来识别和提取单词，或者利用自动机等数据结构来进行状态转换和匹配。

(3)输入/输出：在词法分析中，输入输出操作可能成为性能瓶颈。因此，优化输入输出操作，减少文件读写次数，或者采用缓冲技术来减少 IO 开销，都能提高词法分析的效率。

(4)并行处理：采用并行处理技术，将词法分析任务拆分成多个子任务并行处理，可以加速整个词法分析过程。