

词法分析程序的设计与实现——工具生成

任飞 2021210724

目录

1 概述	1
1.1 实验内容与要求	1
1.2 实现的内容	1
2 程序设计说明	2
2.1 设计规格	2
2.2 程序结构	3
2.3 词法分析器实现	3
2.4 与手写版本的比较	3
3 测试	3
3.1 基本正常程序测试	3
3.2 复杂数字测试	5
3.3 复杂字符常量、字符串字面量与通用字符名测试	8
4 总结	11

1 概述

1.1 实验内容与要求

1. 选定源语言，比如：C、Pascal、Python、Java 等，任何一种语言均可；
2. 可以识别出用源语言编写的源程序中的每个单词符号，并以记号的形式输出每个单词符号。
3. 可以识别并跳过源程序中的注释。
4. 可以统计源程序中的语句行数、各类单词的个数、以及字符总数，并输出统计结果。
5. 检查源程序中存在的词法错误，并报告错误所在的位置。
6. 对源程序中出现的错误进行适当的恢复，使词法分析可以继续进行，对源程序进行一次扫描，即可检查并报告源程序中存在的所有词法错误。
7. 编写 LEX 源程序，利用 LEX 编译程序自动生成词法分析程序。

1.2 实现的内容

使用 lex 工具构建了一个 C 语言的词法分析程序，支持除了预处理器和替用记号外的 C17 标准中的所有词法元素，具有统计功能。

2 程序设计说明

2.1 设计规格

本程序参照 C17 标准文档编写，尽可能完整地实现了 C17 标准的所有内容，仅考虑实际情况做了以下简化：

1. 不支持与预处理器相关的内容。C 语言的基于文本替换的预处理器是一个比较特殊的设计。C 语言的 token 实际上分为预处理前的 preprocessing-token 和预处理后的 token，尤其是对数字的定义方面两者存在较大的差异。预处理时拼接行尾反斜杠的特性也破坏了词法的正则性。为了降低复杂性，我们只识别预处理后的 token，不考虑与预处理器相关的部分，也就是认为：
 - 所有的宏定义都已经被展开
 - 不存在行尾的反斜杠
 - 不存在三标符
 - 不存在不是合法整数或合法浮点数的预处理数字
 - # 直到行尾将被如同注释地忽略
2. 不支持替用记号，即 <% %> 等记号。

不同于手工编写的版本，本程序只接受 ascii 编码的源文件。

程序使用命令行交互，可通过参数指定输入文件（若不指定则为标准输入），token 流和统计信息输出到标准输出流。

程序将所有的 token 分为以下 7 类：

- 关键字 keyword
- 标识符 identifier
- 标点 punctuator
- 整数常量 integer constant
- 浮点数常量 floating constant
- 字符常量 character constant
- 字符串字面量 string literal

出现错误时会输出为一个类型为 error 的 token。

其中关键字包括：

auto	break	case	char	const	continue
default	do	double	else	enum	extern
float	for	goto	if	inline	int
long	register	restrict	return	short	signed
sizeof	static	struct	switch	typedef	union
unsigned	void	volatile	while	_Alignas	_Alignof
_Atomic	_Bool	_Complex	_Generic	_Imaginary	_Noreturn
_Static_assert	_Thread_local				

标点包括：

{	}	[]	()	;	:	...	?	.	->
~	!	+	-	*	/	%	^	&		=	+=
--	*=	/=	%=	^=	&=	=	==	!=	<	>	<=
>=	&&		<<	>>	<<=	>>=	++	--	,		

2.2 程序结构

lex 源文件为 `lex.l` 经过 lex 程序编译后生成 `lex.yy.c`，与 `main.cpp` 一起编译链接生成可执行文件 `lexer`。

`main` 函数中每次调用 `yylex` 函数即读出一个 `token`。最后输出统计信息。

2.3 词法分析器实现

参照标准中的词法定义编写 lex 正则表达式即可。

2.4 与手写版本的比较

如果输入的源文件在词法上是完全正确的，那么 lex 版本和手写版本的效果是一致的。但如果输入的源文件在词法上有误，两者的表现会有所不同。

lex 生成的词法分析器每次都会遵循最长匹配和优先匹配的原则匹配一个 `token`，有时不会如预期一般给出报错。而手写的词法分析程序会根据具体的情况给出具体的报错。

例如 `123abc`，手写版本认为是一个不合法的整数常量，给出“`abc` 不是合法的整数后缀”的报错。但是 lex 版本会将其识别为一个整数常量 `123` 和一个标识符 `abc`。

相比之下，手写版本的错误处理更加友好。当然使用 lex 也能够编写提供友好错误处理的词法分析程序，但需要手动加入对各种典型错误场景的正则表达式。

3 测试

由于篇幅限制，仅包含较为简短的测试用例。

3.1 基本正常程序测试

输入

```
1 #include <stdio.h>
2 int main(void) {
3     int sum = 0;
4     for (int i = 0; i < 100; i++) {
5         sum += i;
6     }
7     printf("sum: %d\n", sum);
8 }
```

输出

```

1 2:1: <keyword, int>
2 2:5: <identifier, main>
3 2:9: <punctuator, (>
4 2:10: <keyword, void>
5 2:14: <punctuator, )>
6 2:16: <punctuator, {>
7 3:5: <keyword, int>
8 3:9: <identifier, sum>
9 3:13: <punctuator, =>
10 3:15: <error, 0>
11 3:16: <punctuator, ;>
12 4:5: <keyword, for>
13 4:9: <punctuator, (>
14 4:10: <keyword, int>
15 4:14: <identifier, i>
16 4:16: <punctuator, =>
17 4:18: <error, 0>
18 4:19: <punctuator, ;>
19 4:21: <identifier, i>
20 4:23: <punctuator, <>
21 4:25: <integer constant, 100>
22 4:28: <punctuator, ;>
23 4:30: <identifier, i>
24 4:31: <punctuator, ++>
25 4:33: <punctuator, )>
26 4:35: <punctuator, {>
27 5:9: <identifier, sum>
28 5:13: <punctuator, +=>
29 5:16: <identifier, i>
30 5:17: <punctuator, ;>
31 6:5: <punctuator, }>
32 7:5: <identifier, printf>
33 7:11: <punctuator, (>
34 7:12: <string literal, "sum: %d\n">
35 7:23: <punctuator, ,>
36 7:25: <identifier, sum>
37 7:28: <punctuator, )>
38 7:29: <punctuator, ;>
39 8:1: <punctuator, }>
40
41 stat:

```

```

42 total: 144 characters, 8 lines
43 tokens:
44   error:
45       0                      2
46   (total)                    2
47   identifier:
48       i                      4
49   main                      1
50   printf                    1
51   sum                      3
52   (total)                    9
53   integer constant:
54       100                    1
55   (total)                    1
56   keyword:
57       for                    1
58       int                    3
59       void                   1
60   (total)                    5
61   punctuator:
62       (                      3
63       )                      3
64       ++                     1
65       +=                     1
66       ,                      1
67       ;                      5
68       <                      1
69       =                      2
70       {                      2
71       }                      2
72   (total)                    21
73   string literal:
74       "sum: %d\n"            1
75   (total)                    1
76   (total)                    39

```

3.2 复杂数字测试

输入

```

1 int main(void) {
2     -011+0x123u*0123;
3     0xabcdefghijk/0912%0912.0;

```

```

4      1.e+2+0e-2-0x0.4p-2;
5      3.411+-322lu+3L-3lL;
6      12.4qwq+0x12.3P43;
7      0x124+000-111e;
8  }
```

输出

```

1 1:1: <keyword, int>
2 1:5: <identifier, main>
3 1:9: <punctuator, (>
4 1:10: <keyword, void>
5 1:14: <punctuator, )>
6 1:16: <punctuator, {>
7 2:5: <punctuator, ->
8 2:6: <error, 0>
9 2:7: <identifier, ll>
10 2:9: <punctuator, +>
11 2:10: <integer constant, 0x123u>
12 2:16: <punctuator, *>
13 2:17: <integer constant, 0123>
14 2:21: <punctuator, ;>
15 3:5: <integer constant, 0xabcdef>
16 3:13: <identifier, ghijk>
17 3:18: <punctuator, />
18 3:19: <error, 0>
19 3:20: <integer constant, 912>
20 3:23: <punctuator, %>
21 3:24: <floating constant, 0912.0>
22 3:30: <punctuator, ;>
23 4:5: <floating constant, 1.e+2>
24 4:10: <punctuator, +>
25 4:11: <floating constant, 0e-2>
26 4:15: <punctuator, ->
27 4:16: <floating constant, 0x0.4p-2>
28 4:24: <punctuator, ;>
29 5:5: <floating constant, 3.41>
30 5:9: <identifier, l>
31 5:10: <punctuator, +>
32 5:11: <punctuator, ->
33 5:12: <integer constant, 322lu>
34 5:17: <identifier, l>
```


77	qwq	1
78	(total)	8
79	integer constant:	
80	000	1
81	0123	1
82	0x123u	1
83	0x124	1
84	0xabcdef	1
85	111	1
86	322lu	1
87	3L	1
88	3l	1
89	912	1
90	(total)	10
91	keyword:	
92	int	1
93	void	1
94	(total)	2
95	punctuator:	
96	%	1
97	(1
98)	1
99	*	1
100	+	6
101	-	5
102	/	1
103	;	6
104	{	1
105	}	1
106	(total)	24
107	(total)	53

3.3 复杂字符常量、字符串字面量与通用字符名测试

输入

```

1 int main(void) {
2     char *\u1234 = u8"\a\b\f\n\r\t\v\012\xab\U12345678";
3     mm\u123 = L"\qabc";
4     U'3' == '';
5     "this is a
6     '\t
7 }
```


输出

```
1 1:1: <keyword, int>
2 1:5: <identifier, main>
3 1:9: <punctuator, (>
4 1:10: <keyword, void>
5 1:14: <punctuator, )>
6 1:16: <punctuator, {>
7 2:5: <keyword, char>
8 2:10: <punctuator, *>
9 2:11: <identifier, \u1234>
10 2:18: <punctuator, =>
11 2:20: <identifier, u8>
12 2:22: <error, ">
13 2:23: <error, \>
14 2:24: <identifier, a>
15 2:25: <error, \>
16 2:26: <identifier, b>
17 2:27: <error, \>
18 2:28: <identifier, f>
19 2:29: <error, \>
20 2:30: <identifier, n>
21 2:31: <error, \>
22 2:32: <identifier, r>
23 2:33: <error, \>
24 2:34: <identifier, t>
25 2:35: <error, \>
26 2:36: <identifier, v>
27 2:37: <error, \>
28 2:38: <integer constant, 012>
29 2:41: <error, \>
30 2:42: <identifier, xab\U12345678>
31 2:55: <error, ">
32 2:56: <punctuator, ;>
33 3:5: <identifier, mm>
34 3:7: <error, \>
35 3:8: <identifier, u123>
36 3:13: <punctuator, =>
37 3:15: <identifier, L>
38 3:16: <error, ">
39 3:17: <error, \>
```

```

40 3:18: <identifier, qabc>
41 3:22: <error, ">
42 3:23: <punctuator, ;>
43 4:5: <char constant, U'3'>
44 4:10: <punctuator, ==>
45 4:13: <error, '>
46 4:14: <error, '>
47 4:15: <punctuator, ;>
48 5:5: <error, ">
49 5:6: <identifier, this>
50 5:11: <identifier, is>
51 5:14: <identifier, a>
52 6:5: <error, '>
53 6:6: <error, \>
54 6:7: <identifier, t>
55 7:1: <punctuator, }>
56
57 stat:
58 total: 138 characters, 7 lines
59 tokens:
60   char constant:
61       U'3'                1
62       (total)             1
63   error:
64       "                   5
65       '                   3
66       \                   12
67       (total)             20
68   identifier:
69       L                   1
70       \u1234              1
71       a                   2
72       b                   1
73       f                   1
74       is                  1
75       main                1
76       mm                  1
77       n                   1
78       qabc                1
79       r                   1
80       t                   2
81       this                1

```

82	u123	1
83	u8	1
84	v	1
85	xab\U12345678	1
86	(total)	19
87	integer constant:	
88	012	1
89	(total)	1
90	keyword:	
91	char	1
92	int	1
93	void	1
94	(total)	3
95	punctuator:	
96	(1
97)	1
98	*	1
99	;	3
100	=	2
101	==	1
102	{	1
103	}	1
104	(total)	11
105	(total)	55

4 总结

这次实验让我初步掌握了 lex 这个工具。