

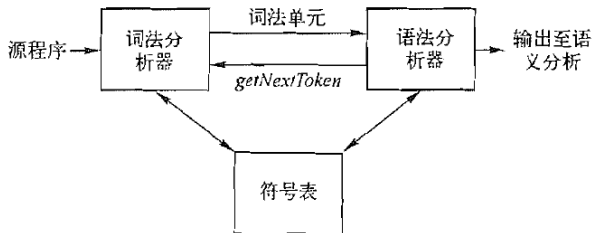
(1. ANTLR v4)

hfwei@nju.edu.cn

20221109 ()

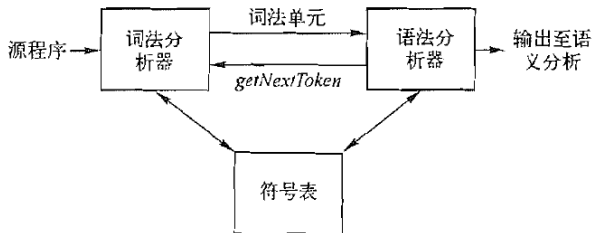


\therefore / s



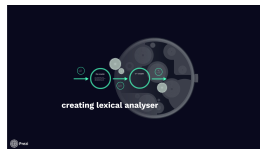
\therefore

$\therefore / s + (\text{token})$



\therefore

()



(gcc)



master ▾

[gcc / gcc / c-family / c-lex.c](#)

iains Objective-C/C++ : Improve '@' keyword locations. ...

19 contributors +7

1435 lines (1278 sloc) | 38.8 KB

master ▾

[gcc / libcpp / lex.c](#)

urnathan cpplib: EOF in pragmas ...

25 contributors +13

4364 lines (3825 sloc) | 119 KB

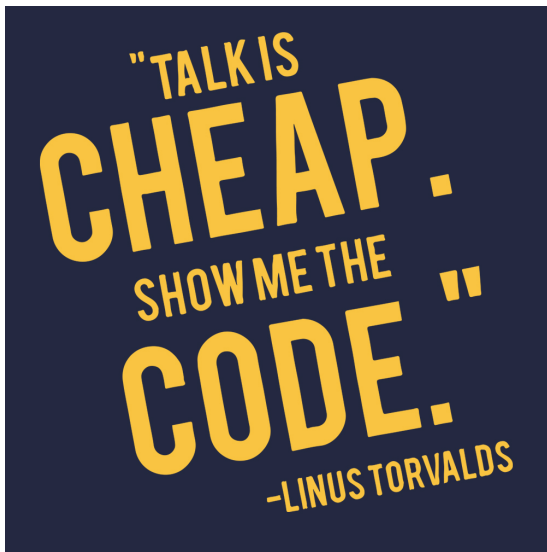


:

SimpleExpr.g4

:

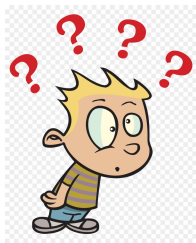
- ▶ SimpleExpr.tokens
- ▶ SimpleExprLexer.java



ANTLR v4

: *vs.*

: `>=, ifhappy, thenext, 1.23`



词法单元	非正式描述	词素示例
if	字符 i, f	if
else	字符 e, l, s, e	else
comparison	< 或 > 或 <= 或 >= 或 == 或 !=	<=, !=
id	字母开头的字母 / 数字串	pi, score, D2
number	任何数字常量	3.14159, 0, 6.02e23
literal	在两个 "之间, 除" 以外的任何字符	"core dumped"

id: /

id , (Language)

id: /

id , (Language)

, (Alphabet)

id: /

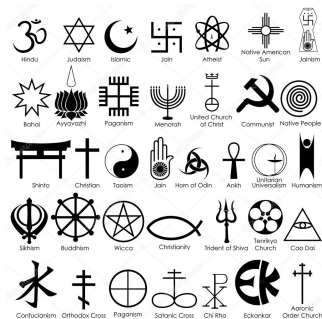
id , (Language)

, (Alphabet)

(,) (String)

Definition ()

Σ



Definition ()

$\Sigma(s) \Sigma$

\mathcal{E}

$: |\epsilon| = 0$

Definition (“”)

$$x = \textit{dog}, y = \textit{house} \quad xy = \textit{doghouse}$$

$$s\epsilon = \epsilon s = s$$

Definition (“”)

$$x = \text{dog}, y = \text{house} \quad xy = \text{doghouse}$$

$$s\epsilon = \epsilon s = s$$

Definition (“”)

$$s^0 \triangleq \epsilon$$

$$s^i \triangleq ss^{i-1}, i > 0$$

Definition ()

Σ

\emptyset

$\{\epsilon\}$

Definition ()

Σ

\emptyset

$\{\epsilon\}$

id : $\{a, b, c, a1, a2, \dots\}$

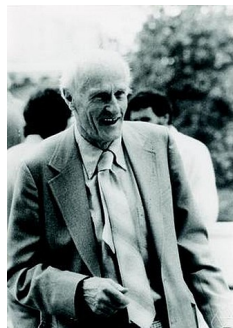
ws : $\{\text{blank}, \text{tab}, \text{newline}\}$

if : $\{\text{if}\}$

,

运算	定义和表示
L 和 M 的并	$L \cup M \doteq \{s \mid s \text{ 属于 } L \text{ 或者 } s \text{ 属于 } M\}$
L 和 M 的连接	$LM = \{st \mid s \text{ 属于 } L \text{ 且 } t \text{ 属于 } M\}$
L 的 Kleene 闭包	$L^* = \bigcup_{i=0}^{\infty} L^i$
L 的正闭包	$L^+ = \bigcup_{i=1}^{\infty} L^i$

$$L^* \ (L^+)$$



Stephen Kleene
(1909 ~ 1994)

$$L = \{A, B, \dots, Z, a, b, \dots, z\}$$

$$D = \{0, 1, \dots, 9\}$$

运算	定义和表示
L 和 M 的并	$L \cup M \doteq \{s \mid s \text{ 属于 } L \text{ 或者 } s \text{ 属于 } M\}$
L 和 M 的连接	$LM = \{st \mid s \text{ 属于 } L \text{ 且 } t \text{ 属于 } M\}$
L 的Kleene 闭包	$L^* = \bigcup_{i=0}^{\infty} L^i$
L 的正闭包	$L^+ = \bigcup_{i=1}^{\infty} L^i$

$$L = \{A, B, \dots, Z, a, b, \dots, z\}$$

$$D = \{0, 1, \dots, 9\}$$

运算	定义和表示
L 和 M 的并	$L \cup M = \{s \mid s \text{ 属于 } L \text{ 或者 } s \text{ 属于 } M\}$
L 和 M 的连接	$LM = \{st \mid s \text{ 属于 } L \text{ 且 } t \text{ 属于 } M\}$
L 的 Kleene 闭包	$L^* = \bigcup_{i=0}^{\infty} L^i$
L 的正闭包	$L^+ = \bigcup_{i=1}^{\infty} L^i$

$$L \cup D \quad LD \quad L^4 \quad L^* \quad D^+ \quad L(L \cup D)^*$$

id : $L(L \cup D)^*$

ANTLR v4 : **id** ?

id : $L(L \cup D)^*$

ANTLR v4 : **id** ?



$$r \models L(r)$$

Syntax

Semantics

,

Definition ()

$\Sigma, \Sigma :$

(1) $\epsilon ;$

(2) $\forall a \in \Sigma, a ;$

(3) $r , (r) ;$

(4) $r \mid s , r|s, rs, r^*$

$: () \prec * \prec \prec |$

$$(a)|((b)^*(c)) \equiv a|b^*c$$

$$r \quad L(r)$$

Definition ()

$$L(\epsilon) = \{\epsilon\} \quad (1)$$

$$L(a) = \{a\}, \forall a \in \Sigma \quad (2)$$

$$L((r)) = L(r) \quad (3)$$

$$L(r|s) = L(r) \cup L(s) \quad L(rs) = L(r)L(s) \quad L(r^*) = (L(r))^* \quad (4)$$

$$\Sigma = \{a, b\}$$

$$L(a|b) = \{a, b\}$$

$$\Sigma = \{a, b\}$$

$$L(a|b) = \{a, b\}$$

$$L((a|b)(a|b))$$

$$\Sigma = \{a, b\}$$

$$L(a|b) = \{a, b\}$$

$$L((a|b)(a|b))$$

$$L(a^*)$$

$$\Sigma = \{a, b\}$$

$$L(a|b) = \{a, b\}$$

$$L((a|b)(a|b))$$

$$L(a^*)$$

$$L((a|b)^*)$$

$$\Sigma = \{a, b\}$$

$$L(a|b) = \{a, b\}$$

$$L((a|b)(a|b))$$

$$L(a^*)$$

$$L((a|b)^*)$$

$$L(a|a^*b)$$

So/
EASY

The text 'So/EASY' is rendered in a hand-drawn, chalk-like style. 'So' is in green and 'EASY' is in blue. A pink pencil is positioned diagonally, pointing towards the 'o' in 'So' and the 'Y' in 'EASY'. Several short, pink, radiating lines surround the text, giving it a sense of motion or emphasis.

表达式	匹配	例子
<code>c</code>	单个非运算符字符 c	<code>a</code>
<code>\c</code>	字符 c 的字面值	<code>*</code>
<code>"s"</code>	串 s 的字面值	<code>"**"</code>
<code>.</code>	除换行符以外的任何字符	<code>a.*b</code>
<code>^</code>	一行的开始	<code>^abc</code>
<code>\$</code>	行的结尾	<code>abc\$</code>
<code>[s]</code>	字符串 s 中的任何一个字符	<code>[abc]</code>
<code>[^s]</code>	不在串 s 中的任何一个字符	<code>[^abc]</code>
<code>r*</code>	和 r 匹配的零个或多个串连接成的串	<code>a*</code>
<code>r+</code>	和 r 匹配的一个或多个串连接成的串	<code>a+</code>
<code>r?</code>	零个或一个 r	<code>a?</code>
<code>r{m,n}</code>	最少 m 个, 最多 n 个 r 的重复出现	<code>a{1,5}</code>
<code>r₁r₂</code>	r_1 后加上 r_2	<code>ab</code>
<code>r₁ r₂</code>	r_1 或 r_2	<code>a b</code>
<code>(r)</code>	与 r 相同	<code>(a b)</code>
<code>r₁/r₂</code>	后面跟有 r_2 时的 r_1	<code>abc/123</code>

表达式	匹配	例子
<code>c</code>	单个非运算符字符 c	<code>a</code>
<code>\c</code>	字符 c 的字面值	<code>*</code>
<code>"s"</code>	串 s 的字面值	<code>"**"</code>
<code>.</code>	除换行符以外的任何字符	<code>a.*b</code>
<code>^</code>	一行的开始	<code>^abc</code>
<code>\$</code>	行的结尾	<code>abc\$</code>
<code>[s]</code>	字符串 s 中的任何一个字符	<code>[abc]</code>
<code>[^s]</code>	不在串 s 中的任何一个字符	<code>[^abc]</code>
<code>r*</code>	和 r 匹配的零个或多个串连接成的串	<code>a*</code>
<code>r+</code>	和 r 匹配的一个或多个串连接成的串	<code>a+</code>
<code>r?</code>	零个或一个 r	<code>a?</code>
<code>r{m,n}</code>	最少 m 个, 最多 n 个 r 的重复出现	<code>a{1,5}</code>
<code>r₁r₂</code>	r_1 后加上 r_2	<code>ab</code>
<code>r₁ r₂</code>	r_1 或 r_2	<code>a b</code>
<code>(r)</code>	与 r 相同	<code>(a b)</code>
<code>r₁/r₂</code>	后面跟有 r_2 时的 r_1	<code>abc/123</code>

`[0 - 9] [a - zA - Z]`

Vim ↕	Java ↕	ASCII ↕	Description ↕
	<code>\p{ASCII}</code>	<code>[\x00-\x7F]</code>	ASCII characters
	<code>\p{Alnum}</code>	<code>[A-Za-z0-9_]</code>	Alphanumeric characters
<code>\w</code>	<code>\w</code>	<code>[A-Za-z0-9_]</code>	Alphanumeric characters plus "_"
<code>\W</code>	<code>\W</code>	<code>[^A-Za-z0-9_]</code>	Non-word characters
<code>\a</code>	<code>\p{Alpha}</code>	<code>[A-Za-z]</code>	Alphabetic characters
<code>\s</code>	<code>\p{Blank}</code>	<code>[\t]</code>	Space and tab
<code>\< \></code>	<code>\b</code>	<code>(?<=\W) (?=\w) (?<=\w) (?=\W)</code>	Word boundaries
	<code>\B</code>	<code>(?<=\W) (?=\W) (?<=\w) (?=\w)</code>	Non-word boundaries
	<code>\p{Cntrl}</code>	<code>[\x00-\x1F\x7F]</code>	Control characters
<code>\d</code>	<code>\p{Digit}</code> or <code>\d</code>	<code>[0-9]</code>	Digits
<code>\D</code>	<code>\D</code>	<code>[^0-9]</code>	Non-digits
	<code>\p{Graph}</code>	<code>[\x21-\x7E]</code>	Visible characters
<code>\l</code>	<code>\p{Lower}</code>	<code>[a-z]</code>	Lowercase letters
<code>\p</code>	<code>\p{Print}</code>	<code>[\x20-\x7E]</code>	Visible characters and the space character
	<code>\p{Punct}</code>	<code>[!\"#\$%&'()*+,-./:;<=>?@^_`{ }~.]</code>	Punctuation characters
<code>_s</code>	<code>\p{Space}</code> or <code>\s</code>	<code>[\t\r\n\v\f]</code>	Whitespace characters
<code>\S</code>	<code>\S</code>	<code>[^ \t\r\n\v\f]</code>	Non-whitespace characters
<code>\u</code>	<code>\p{Upper}</code>	<code>[A-Z]</code>	Uppercase letters
<code>\x</code>	<code>\p{XDigit}</code>	<code>[A-Fa-f0-9]</code>	Hexadecimal digits

$$\left(0|(1(01^*0)^*1)\right)^*$$



<https://regex101.com/r/ED4qgC/1>

REGULAR EXPRESSION v1 ✓

$$/ \wedge (0 | (1 (0 1 * 0) * 1)) * 5$$

TEST STRING

Ⓢ

1

10

11

100

101

110

111

1000

1001

1010

1011

1100

1101

1110

1111

10000

10001

10010

10011

10100

10101

1 0 1 1 0

REGULAR EXPRESSION v1 ✓

$$\frac{1}{(0 + (1(01^*0)^*1))^*5}$$

TEST STRING

①

1

10

11

100

101

110

111

1000

1001

1010

1011

1100

1101

1110

1 1 1 1

10000

10001

10010

10011

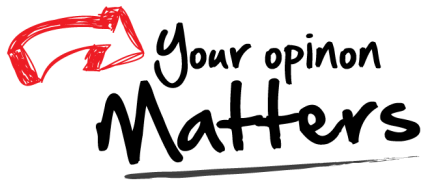
10100

10101

10110

3 ()

Thank
You!



Office 926

hfwei@nju.edu.cn