



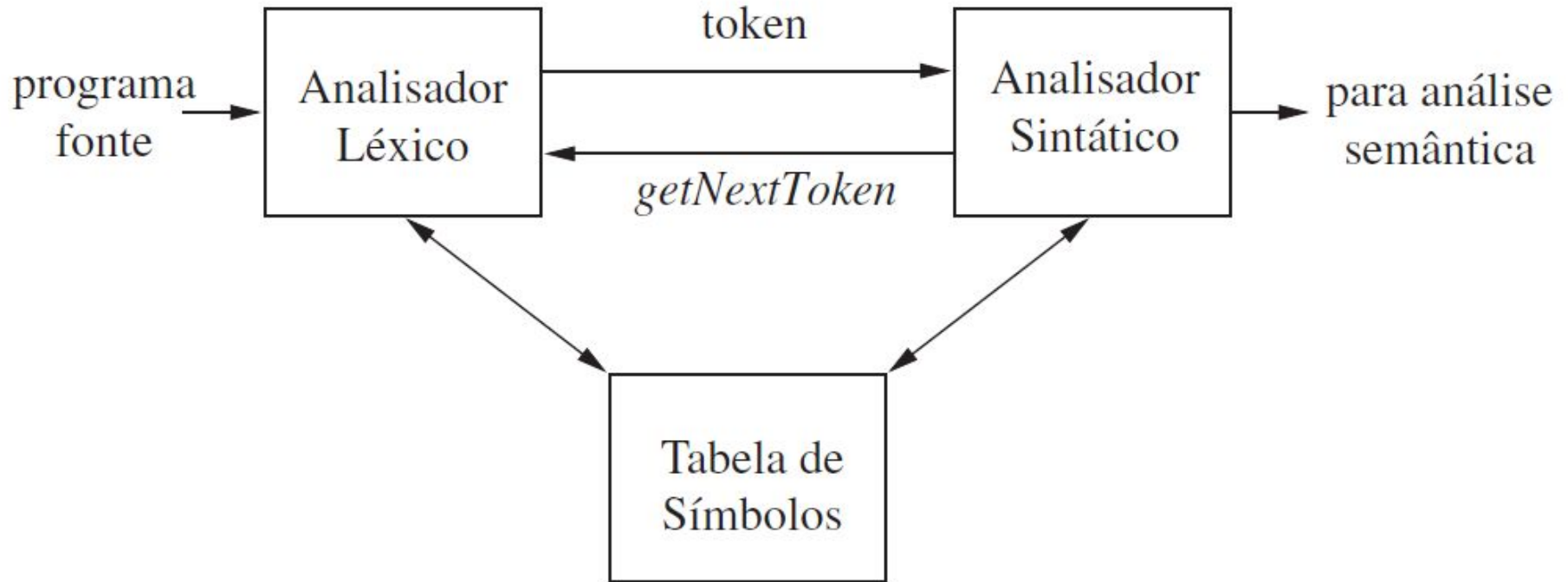
INSTITUTO FEDERAL DE
EDUCAÇÃO, CIÊNCIA E TECNOLOGIA
CEARÁ
CAMPUS ARACATI

Bacharelado em Ciência da Computação

Disciplina: Compiladores
Aula 02 - Criando um analisador Léxico

Professor: Diego Rocha Lima

Analizador Léxico → Analizador Sintático



Tokens X Lexema

TOKEN	DESCRIÇÃO INFORMAL	EXEMPLOS DE LEXEMAS
if	caracteres i, f	if
else	caracteres e, l, s, e	else
comparison	< or > ou <= ou >= ou == ou !=	<=, !=
id	letra seguida por letras e dígitos	pi, score, D2
number	qualquer constante numérica	3.14159, 0, 6.02e23
literal	qualquer caractere diferente de ", cercado por "s	"core dumped"

1. Um token para cada palavra-chave. O padrão para uma palavra-chave é o mesmo que a própria palavra-chave.
2. Tokens para os operadores, seja individualmente ou em classes, como o token `comparison`
3. Um token representando todos os identificadores.
4. Um ou mais tokens representando constantes, como números e cadeias literais.
5. Tokens para cada símbolo de pontuação, como parênteses esquerdo e direito, vírgula e ponto-e-vírgula.

Exemplos de Tokens

E = M * C ** 2

são escritos a seguir como uma seqüência de pares.

< **id**, apontador para entrada da tabela de símbolos de E>

< **assign_op**>

< **id**, apontador para entrada da tabela de símbolos de M>

< **mult_op**>

< **id**, apontador para entrada da tabela de símbolos de C>

< **exp_op**>

< **number**, valor inteiro 2>

Operações

OPERAÇÃO	DEFINIÇÃO E NOTAÇÃO
<i>União de L e M</i>	$L \cup M = \{s \mid s \text{ está em } L \text{ ou } s \text{ está em } M\}$
<i>Concatenação de L e M</i>	$LM = \{st \mid s \text{ está em } L \text{ e } t \text{ está em } M\}$
<i>Fecho Kleene de L</i>	$L^* = \bigcup_{i=0}^{\infty} L^i$

- A união é a conhecida operação sobre conjuntos.
- A concatenação de linguagens são todas as cadeias formadas a partir de uma cadeia da primeira linguagem e uma cadeia da segunda linguagem, em todas as formas possíveis, e concatenando-as.
- O fecho (Kleene) de uma linguagem L , indicado por L^* , é o conjunto de cadeias obtidas concatenando L zero ou mais vezes.

Exemplos

Considere que L seja o conjunto de letras $\{A, B, \dots, Z, a, b, z\}$ e D seja o conjunto de dígitos $\{0, 1, \dots, 9\}$.

1. $L \cup D$ é o conjunto de letras e dígitos — estritamente falando, a linguagem com 62 cadeias de tamanho um, cada uma tendo uma letra ou um dígito.
2. LD é o conjunto de 520 cadeias de tamanho dois, cada uma consistindo em uma letra seguida por um dígito.
3. L^3 é o conjunto de todas as cadeias de 3 letras.
4. L^* é o conjunto de todas as cadeias de letras, incluindo , a cadeia vazia.
5. $L(L \cup D)^*$ é o conjunto de todas as cadeias de letras e dígitos começando com uma letra.
6. D^+ é o conjunto de todas as cadeias de um ou mais dígitos.

Expressões Regulares

Nesta notação, se estabelecermos que *letra_* significa qualquer letra ou o sublinhado, e que *dígito_* significa qualquer dígito, então poderíamos descrever os identificadores da linguagem C por:

$$\textit{letra_} (\textit{letra_} \mid \textit{dígito_})^*$$

Existem quatro partes na indução, por meio das quais expressões regulares maiores são construídas a partir das menores. Suponha que r e s sejam expressões regulares denotando as linguagens $L(r)$ e $L(s)$, respectivamente.

1. $(r)l(s)$ é uma expressão regular denotando a linguagem $L(r) \cup L(s)$.
2. $(r)(s)$ é uma expressão regular denotando a linguagem $L(r)L(s)$.
3. $(r)^*$ é uma expressão regular denotando $(L(r))^*$.
4. (r) é uma expressão regular denotando $L(r)$. Essa última regra diz que podemos acrescentar pares de parênteses adicionais em torno das expressões sem alterar a linguagem que elas denotam.

Relembrando Leis Algébricas

LEI	DESCRIÇÃO
$rls = slr$ $rl(st) = (rls)lt$ $r(st) = (rs)t$ $r(st) = rslrt; (slt)r = srltr$ $\epsilon r = r\epsilon = r$ $r^* = (r\epsilon)^*$ $r^{**} = r^*$	$ $ é comutativo $ $ é associativo A concatenação é associativa A concatenação distribui entre $ $ ϵ é o elemento identidade para concatenação ϵ é garantido em um fechamento $*$ é igual potência

Construindo definições regulares

Os identificadores de C são cadeias de letras, dígitos e sublinhados. Aqui está uma definição regular para os identificadores da linguagem C. Por convenção, vamos usar itálico para os símbolos definidos em definições regulares.

```
letter_ → A | B | ... | Z | a | b | ... | z | _  
digit   → 0 | 1 | ... | 9  
id     → letter_ ( letter_ | digit )*
```

Construindo definições regulares

Números sem sinal (inteiros ou ponto flutuante) são cadeias como 5280, 0.01234, 6.336E4 ou 1.89E-4. A definição regular

$digit \rightarrow 0 \mid 1 \mid \dots \mid 9$

$digits \rightarrow digit\,digit^*$

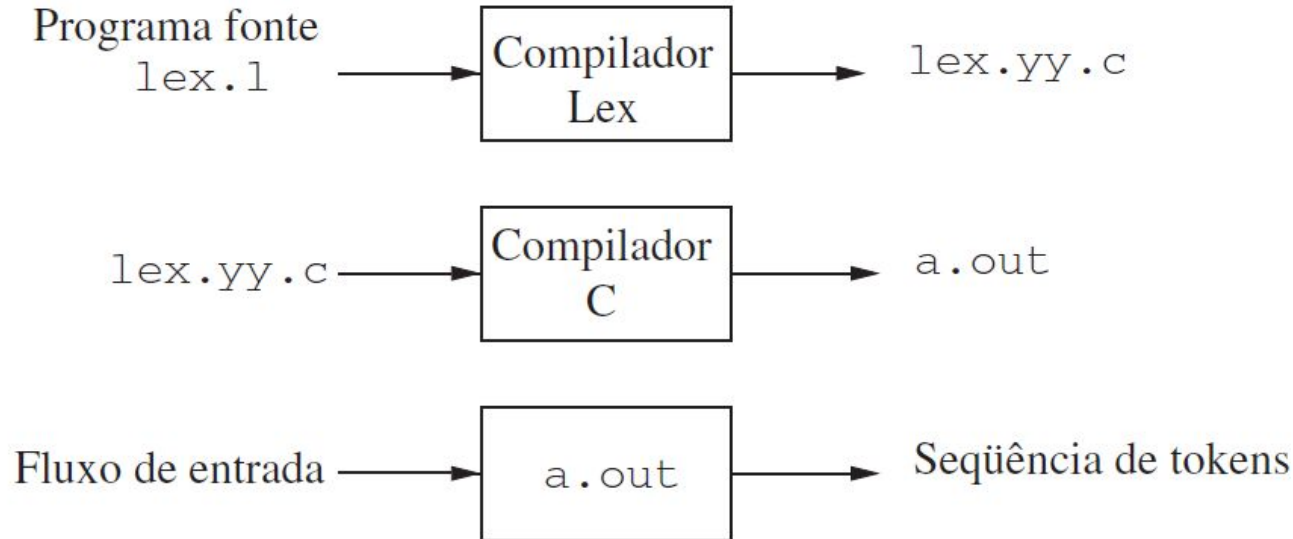
$optionalFraction \rightarrow .\,digits \mid \epsilon$

$optionalExponent \rightarrow (E\, (+|-|\epsilon)\,digits) \mid \epsilon$

$number \rightarrow digits\,optionalFraction\,optionalExponent$

Utilizando o LEX para criar analisadores léxicos

Lex é um programa de computador que gera analisadores léxicos, que são ferramentas que mapeiam expressões regulares em blocos de código. Ele é usado para escrever programas em linguagem C ou C++. A seguir algumas expressões regulares com o LEX



EXPRESSÃO	CASA COM	EXEMPLO
c	o único caractere não operador c	a
$\backslash c$	o caractere C literalmente	$\backslash *$
$"s"$	a cadeia s literalmente	$"**"$
$.$	qualquer caractere menos quebra de linha	$a.*b$
$^$	o início de uma linha	abc
$\$$	o fim de uma linha	$abc\$$
$[s]$	qualquer um dos caracteres na cadeia s	$[abc]$
$[^s]$	qualquer caractere não presente na cadeia s	$[^abc]$
r^*	zero ou mais cadeias casando com r	a^*
r^+	uma ou mais cadeias casando com r	a^+
$r?$	zero ou um r	$a?$
$r\{m,n\}$	entre m e n ocorrências de r	$a[1,5]$
r_1r_2	um r_1 seguido por um r_2	ab
$r_1 \mid r_2$	um r_1 ou um r_2	$a \mid b$
(r)	o mesmo que r	$(a \mid b)$
r_1/r_2	R_1 quando seguido por r_2	$abc/123$

Exemplos de Tokens com o LEX

digit → [0-9]

digits → *digit*⁺

number → *digits* (. *digits*) ? (E [+ -] ? *digits*) ?

lettter → [A-Za-z]

id → *letter* (*letter* | *digit*)^{*}

if → if

then → then

else → else

relop → < | > | <= | >= | = | <>

Valores de atributos

LEXEMAS	NOME DO TOKEN	VALOR DO ATRIBUTO
Qualquer <i>ws</i>	–	–
<i>if</i>	if	–
<i>then</i>	then	–
<i>else</i>	else	–
Qualquer <i>id</i>	id	Apontador para entrada de tabela
Qualquer <i>number</i>	number	Apontador para entrada de tabela
<	relop	LT
<=	relop	LE
=	relop	EQ
<>	relop	NE
>	relop	GT
>=	relop	GE

Estrutura do LEX (arquivo.l)

A seção de declarações inclui declarações de variáveis, constantes manifestas (identificadores declarados para significar uma constante, por exemplo, o nome de um token) e definições regulares.

declarações

%%

regras de tradução

%%

funções auxiliares

Já as regras de produção tem o seguinte formato:

Padrão {Ação}

A terceira seção contém quaisquer funções adicionais usadas nas ações. Alternativamente, essas funções podem ser compiladas separadamente e carregadas com o analisador léxico

Primeiro programa em LEX (1/3)

```
%{  
    /* definições de constantes manifestas  
    LT, LE, EQ, NE, GT, GE,  
    IF, THEN, ELSE, ID, NUMBER, RELOP */  
%}  
  
/* definições regulares */  
delim      [ \t\n]  
ws         {delim}+  
letter     [A-Za-z]  
digit      [0-9]  
id         {letter}({letter}|{digit})*  
number     {digit}+(\.{digit}+)?(E[+-]?{digit}+)?
```


Primeiro programa em LEX (2/3)

%%

```
{ws}      { /* nenhuma ação e nenhum retorno */ }
if         { return(IF); }
then       { return(THEN); }
else       { return(ELSE); }
{id}       { yylval = (int) installID(); return(ID); }
{number}   { yylval = (int) installNum(); return(NUMBER); }
"<"        { yylval = LT; return(RELOP); }
"<="       { yylval = LE; return(RELOP); }
"="        { yylval = EQ; return(RELOP); }
"<>"       { yylval = NE; return(RELOP); }
">"        { yylval = GT; return(RELOP); }
">="       { yylval = GE; return(RELOP); }
```

%%

Primeiro programa em LEX (3/3)

```
int installID() { /* função para instalar o lexema, cujo
                  primeiro caractere é apontado por yytext,
                  e cujo tamanho é yyleng, na tabela de
                  símbolos, e retorna um apontador para lá */
}

int installNum() { /* semelhante a installID, mas coloca constantes
                   numéricas em uma tabela separada */
}
```

Na seção de função auxiliar, vemos duas dessas funções, `installID()` e `installNum()`. Assim como a parte da seção de declaração que aparece entre `%{...%}`, tudo na seção auxiliar é copiado diretamente para o arquivo `lex.yy.c`, mas pode ser usado nas ações.

Variáveis padrão do lex

`yylval` → variável inteira contendo o lexema

`yytext` → é um apontador para o início do lexema, ou seja, uma string que contém o valor atual, lido e armazenado no buffer

`yylen` → é o tamanho do lexema encontrado

Vamos a um
exemplo

aula1.l

```
1  all: aula1.l
2      clear
3      flex -i aula1.l
4      gcc lex.yy.c -o aula1 -lfl
5      ./aula1
```

```
1  %{
2      //definições,funções... (código C)
3  %}
4  /*expressões regulares*/
5  NUM [0-9]+
6  STR [A-Za-z]+
7  SIMB [-+*/=]
8  DESC [ \n\t]
9  %%
10 | /*tokens*/
11 "FIM" {printf("Fim do programa\n");}
12 {NUM} {printf("Número: %s\n",yytext);}
13 {SIMB} {printf("Operação: %s\n",yytext);}
14 {STR} {printf("String: %s\n",yytext);}
15 {DESC} {printf("Caractere ignorado\n");}
16 %%
17 /*definições*/
18 int main(){
19     yyin=fopen("codigo.d","r"); //abre o código fonte
20     yylex(); //chama o analizador léxico
21     fclose(yyin); //fecha o código fonte
22     return 0;
23 }
```

Exercício

- 1 - Criar um analisador léxico utilizando FLEX para reconhecer operações em uma calculadora simples
- 2 - Dar início ao seu projeto para criação de uma linguagem de programação