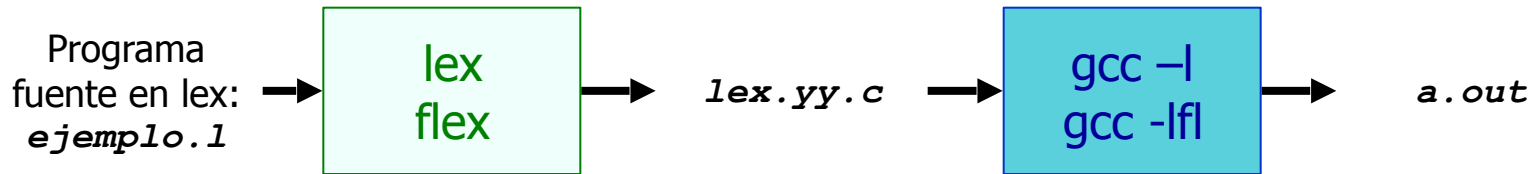


Generador de analizadores léxicos.

Un analizador léxico (scanner) es un programa que reconoce **patrones léxicos** en un texto.

# Funcionamiento de lex

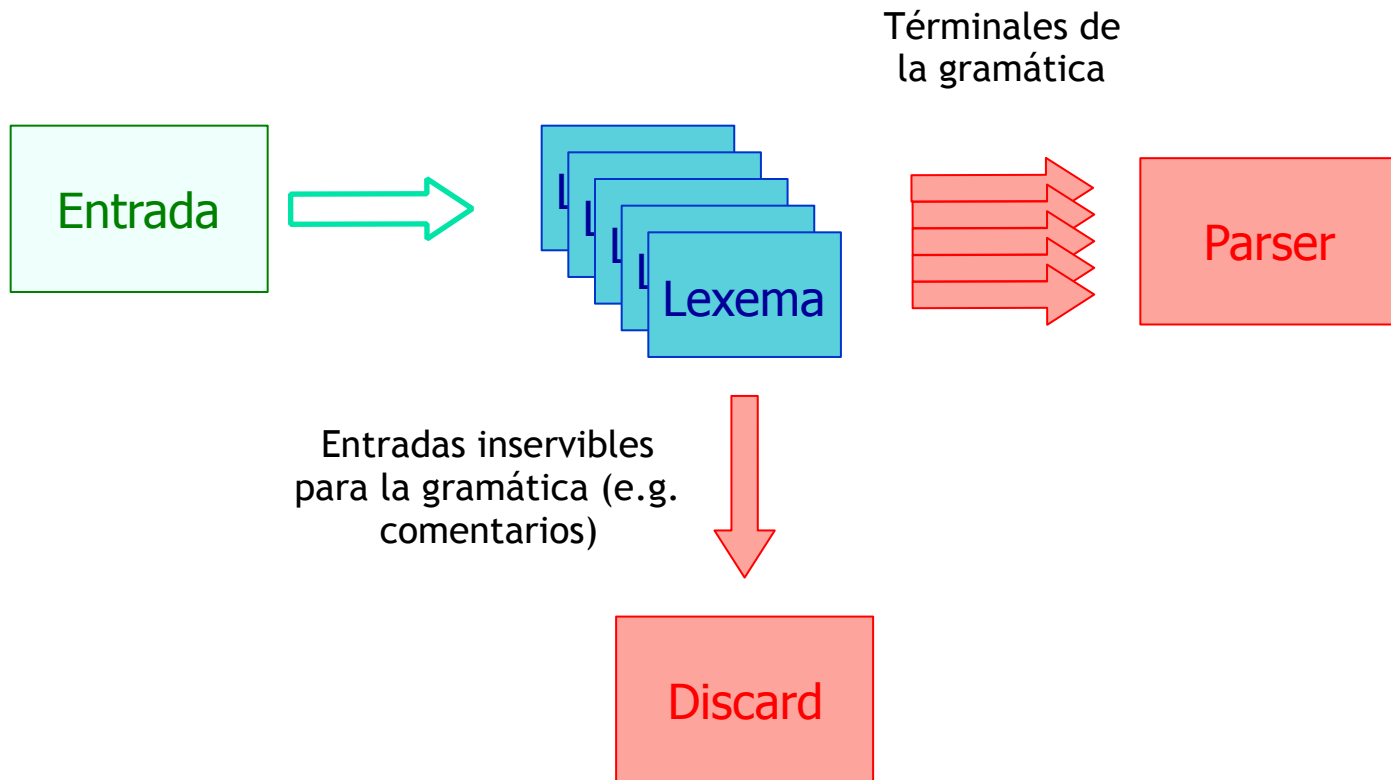


```
flex ejemplo.1 && gcc -oejemplo lex.yy.c -lfl(*)
```



(\*) En MacOS, la librería se referencia con -ll. En Linux instalen flex 2.6.1 e instalen flex-devel

# Tokens / Lexemas





# Formato del archivo fuente

---

```
%{  
    Declaraciones en C  
}%  
{sección de definiciones}  
%%  
{sección de reglas}  
%%  
{sección de código de usuario}
```



**NO PUEDE  
FALTAR!!!**

En la sección de reglas, a cada patrón se asocia una acción.  
Cuando el patrón se encuentra en el texto, se ejecuta la acción asociada.

Si para un patrón no hay acción asociada, no se hace nada.

Regla por default: copiar todo carácter de la entrada en la salida.

# Ejemplos de archivos fuente

```
%%  
.      ECHO;  
\n     ECHO;  
%%  
int yywrap(void) {  
    return 1;  
}  
int main(void) {  
    yylex();  
    return 0;  
}
```



Copia toda la entrada en la salida

El comportamiento del programa generado es el siguiente:

- Imprimir en la salida estándar los lexemas que no se adapten a ningún patrón (ECHO)
- Realizar la acción indicada para los lexemas que se ajustan a un patrón.



# Sección de Reglas

---

`<patron> <reglas>`

Comenzando en primer columna

Después de espacio, una acción de C (o varias encerradas en {})

Si lo primero que aparece en una línea no está en la primer columna, se copia textual en el archivo C generado.



# Patrones

<i>Expresión</i>	<i>Significado</i>
.	Cualquier carácter excepto \n (en algunas versiones de FLEX si incluye el \n)
\	Para usar un carácter operador literalmente
R*	Cero o más apariciones de la expresión R
R+	Una o más apariciones de la expresión R
R?	Cero o una aparición de la expresión R
^R	Si R se da al comienzo de la línea
R\$	Si R se da al final de la línea
R S	Si se da la expresión R o la expresión S
R/S	Si se da la expresión R seguida de S
(RS)+	Una o más apariciones de la concatenación
"R+S"	Literalmente el texto R+S
[...]	Todo carácter que esté comprendido dentro de la clase
[^a-z]	Todo carácter que no sea letra minúscula
R{n,m}	De n a m apariciones de R

# Ejemplos de archivos fuente

```
%%  
[ \t]+ printf(" ");
```

Si la entrada coincide con uno o más espacios o tabulaciones seguidos, muestra uno solo

```
%%
```

Copia toda la entrada en la salida

```
%%  
.      ECHO;  
\n      ECHO;
```

Copia toda la entrada en la salida.  
El '.' excluye el '\n'

```
%%  
[0-9]+ printf("digito");  
.
```

Si la entrada coincide con uno o más dígitos seguidos, escribe "digito". Lo demás no lo muestra





# Coincidencia con un patrón

---

- Cuando el scanner se ejecuta, analiza la entrada en búsqueda de cadenas que coincidan con algún patrón.
  - Si una cadena coincide con más de un patrón, se toma la regla que coincida con el token más largo.
  - Si los dos patrones son de igual longitud, se toma el primero que aparezca en la sección
  - Si un patrón no tiene regla asociada, se elimina
  - Si un patrón tiene asociado una "|", se ejecuta la regla del patrón que sigue (sirve para asociar varios patrones a una única acción)



# Ejemplos:

```
%%  
"alfa"          putchar('a');  
"alfabeto"      printf("de la 'a' a la 'z'");
```

```
%{  
#include <stdio.h>  
int lineno;  
%}  
  
%%  
begin|end      {printf("%s:%d\n",yytext,lineno);}  
\n             lineno++;  
.*             ;
```

```
%%  
" "           |  
"\t"          |  
"\n"          ;
```

```
%%  
[ \t\n] ;
```



Para las tres reglas, se ejecuta la misma acción definida en la última regla (la del \n).



# REJECT

---

- Indica que hay que ejecutar la siguiente acción de coincidencia que pueda corresponderse con la entrada.

```
%{  
int s = 0;  
int h = 0;  
%}  
%%  
she      {s++;REJECT;}  
he       {h++;REJECT;}  
\n      |  
      .      ;
```



# Sección de Código de Usuario

---

- Código que se copia textual en lex.yy.c

- Ej:

```
int
main(int argc, char **argv)
{
    if(argc > 0)
        yyin = fopen(argv[1], "r");
    else
        yyin = stdin;
    yylex();
    return 0;
}
```



# Sección de Código de Usuario

---

- El programa por default es:

```
int yywrap(void) {  
    return 1;  
}  
int main(void) {  
    yylex() ;  
    return 0;  
}
```

**yylex()** es la función que contiene el automata que se define a partir de las reglas.

**yywrap()** es la función que se invoca cuando se agota la entrada. Retorna 0 si se necesita procesar más cosas, como por ejemplo abrir un nuevo archivo de entrada (usando yyin) para continuar con el procesamiento



# Variables de lex

---

<i>Nombre</i>	<i>Significado</i>
char * <b>yytext</b>	Puntero al inicio del texto que se ha hecho coincidir con un patrón (puntero al inicio del token) También puede ser un arreglo (char [])
<b>yytext</b>	Longitud de yytext
FILE * <b>yyout</b>	Output file (stdout por default)
FILE * <b>yyin</b>	Input file (stdin por default)
YY_START	Indicador de la condición de arranque actual.



# Funciones de lex

Nombre	Significado
<b>yymore()</b>	La proxima vez que se compruebe una coincidencia, el token se anexará al actual yytext.
<b>yylless(n)</b>	Retorna todo excepto los n primeros caracteres del token
<b>unput(c)</b>	El caracter c se vuelve a colocar en el input. Será el próximo carácter a procesar por el scanner.
<b>input()</b>	Lee el siguiente carácter del input stream
<b>Yyrestart(FILE*)</b>	Toma el puntero al archivo indicado como nueva entrada.

```
%%  
mega-    ECHO; yymore();  
kludge   ECHO;
```

Si la entrada es mega-kludge muestra **mega-mega-kludge**. (El ECHO de la primera linea imprime el primer **mega-** en tanto que en el segundo se concatena con **kludge**).

```
%%  
foobar   ECHO; ylless(3);  
[a-z]+   ECHO;
```

Si la entrada es foobar muestra **foobarbar**



# Sección de Definiciones

---

- Contiene declaraciones de definiciones más simples de nombres (sinónimos)
  - Ej: `digit [0-9]`
- Contiene declaraciones de condiciones de inicio (start conditions)(%s o %x)
  - Ej: `%s upper`
- Contiene código que se copia textual en el archivo .c generado.

- Ej: 

```
%{  
#include <ctype.h>  
%}
```





# Ejemplos:

---

```
DIGITO  [0-9]
%%
{DIGITO}+      printf("entero");
```

```
%{
#include <ctype.h>
%}
%x      PALABRA
MINUSCULA      [a-z]
MAYUSCULA      [A-Z]
%%
<PALABRA>[ \t\n]      ECHO;BEGIN(INITIAL);
<INITIAL>{MAYUSCULA}  ECHO;BEGIN(PALABRA);
<INITIAL>{MINUSCULA} putchar(toupper(*yytext));BEGIN(PALABRA);
```



# Condiciones de Arranque

---

- **%S CONDICION**

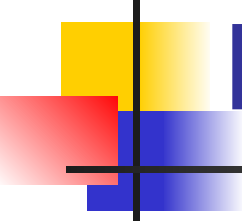
- Se definen en la sección de definiciones con %Start CONDICION
- Una condición de arranque se activa utilizando la acción "BEGIN(CONDICION)".
- El estado inicial es siempre por default INTIAL.
- Las reglas se condicionan con: <CONDICION> { regla }
- Hasta que se ejecute la próxima acción BEGIN, las reglas con la condición de arranque dada estarán activas y las reglas con otras condiciones de arranque estarán inactivas.



# Diferencia entre %x y %s

---

- %s Declara condiciones de arranque inclusivas
  - %X Declara condiciones de arranque exclusivas.
- 
- Si la condición de arranque es inclusiva, entonces las reglas sin condiciones de arranque también estarán activas.
  - Si es exclusiva, entonces solamente las reglas calificadas con la condición de arranque estarán activas.



# Ejemplo %x vs %s

---

```
%s ejemplo  
%% <ejemplo>foo hacer_algo();  
bar algo_mas();
```



```
%x ejemplo  
%% <ejemplo>foo hacer_algo();  
<INITIAL,ejemplo>bar algo_mas();
```

Sin el calificador `<INITIAL,ejemplo>',  
el patrón `bar' en el segundo ejemplo no estará activo



# Ejemplos:

---

```
%{
int line_num = 1;
%}
%x comment
%%

"/*"          BEGIN(comment);

<comment>[^*\n]*      /* eat anything that's not a '*' */
<comment>"*" + [^*/\n]* /* eat up '*'s not followed by '/'s */
<INITIAL>\n           ++line_num;
<comment>"*" + "/"    BEGIN(INITIAL);
%%
int main(int argc, char *argv[])
{
    yylex();
    printf("# of lines = %d\n", line_num);
}
```



Cambien <INITIAL> x <comment>



# Tricky stuff

---

- Lex se llama “flex” ahora en muchas distribuciones. Pueden chequear la versiones con `flex --version`.
- En pampero está instalada la versión 2.6.1. En algunas versiones de ubuntu tienen que compilarlo con:
  - `“gcc yy.tab.c -ly -lfl -ocompiler.o”`
- En algunas instalaciones la función `yywrap` no viene en las librerías por defecto y hay que proporcionarla dentro de nuestros archivos `.l`, o pueden poner `“%option noyywrap”` para que el lex no la utilice.
- Desde bash pueden crear rápido un archivo de texto de entrada con `“cat >> file.in”`. Para cerrar el archivo tienen que enviar EOF que es `CTRL+D`
- Al ejecutable lo pueden arrancar con `“./compiler.o < file.in”`, y va a procesar la entrada que ustedes le pasen.
- Pueden redireccionar la salida del programa con `“./compiler.o < file.in > file.out”` y les va a dejar la salida en “file.out”



# Referencias

---

- The Dragon's Book
- <http://flex.sourceforge.net/manual/>
- *Compiler Implementation in C*, Appel, Ginsburg
- *Engineering a Compiler*, Cooper, Torczon
- Lex & Yacc, Levin and Mason, 1990
- <https://github.com/faturita/YetAnotherCompilerClass>