

آموزش FLEX

با تشکر از آقای مصطفی شاکری و خانم عاطفه عظیمی

Shakery.mostafa@gmail.com

بازنگری و آماده سازی مجدد : رضا فهیمی

rezafahimi@yahoo.com

✓ Flex چیست ؟

اصول و روش های ساخت **تحلیلگر لغوي** برای زبان های مختلف یکسان است در نتیجه ابزارهایی ایجاد شده است که با اخذ مشخصات لغوي زبان مبدا می توانند تحلیلگر لغوي آن را تولید کنند . این گونه نرم افزار ها را تولید کننده ی تحلیلگر لغوي می نامند . یکی از مهمترین این ابزار ها نرم افزار **Flex (Fast lexical analyzer)** است .

معمولا نرم افزار **Flex** به همراه نرم افزار دیگری به نام **Bison** استفاده می شود .

ابزار **Bison** يك مولد **تحليلگر نحوي (Syntax Analyzer)** است که در صورت استفاده در کنار **Flex** ورودی خود را از آن گرفته و سپس عملیات تحلیل معنایی را بر اساس قوانین تعریف شده ی گرامر انجام می دهد.

علاوه بر نرم افزار های **Flex** و **Bison** ، **lex** و **yacc** نیز وجود دارند که در واقع **Flex** و **Bison** به ترتیب مشتق شده از **Lex** و **Yacc** هستند . **Lex** و **Yacc** در ابتدا با سیستم عامل **UNIX** کار می کردند . **Flex** و **Bison** دوباره نوشته شده از **Lex** و **Yacc** برای **DOS** و به تازگی برای **Windows32** می باشند .

✓ مراحل استفاده از Flex جهت تولید تحلیل گر لغوي :

برنامه به زبان Flex

کامپایل به وسیله Flex

برنامه به زبان C یا Pascal

کامپایل به وسیله C

تحلیلگر لغوي

✓ نحوه ي اجراي Flex :

براي اجراي **Flex** ، دو راه وجود داره كه راه اول را به هيچ وجه توصيه نمي كنيم ولي جهت آشنائي بيان مي شود!

(1) در روش اول كافيست بطور مستقيم فايل **Flex.exe** را از طريق خط فرمان (**command prompt**) اجرا كنيم و سپس گرامر را تا انتها بنويسيم ، و در پايان كليد **Ctrl+Z** (نشانگر پايان فايل) را بزنيم تا از ويرايشگر فرمان خارج شده و در صورت صحت گرامر فايل خروجي ايجاد شود .

(2) در روش دوم ابتدا یک فایل (با پسوند **.L**) ایجاد کرده و سپس فایل را از طریق خط فرمان به کامپایلر **Flex** پاس می دهیم و در نهایت در صورتی که فایل مبداء وجود داشته و گرامر صحیح باشد فایل خروجی ایجاد خواهد شد .

بعد از انجام یکی از روش های بالا و در صورت نداشتن خطا در متن برنامه **Flex** فایلی با نام **Yy.lex.c** به عنوان خروجی و به زبان **C** در همان مسیری که برنامه ی **Flex** وجود دارد ایجاد می کند . که برای ساخت فایل اجرایی باید این فایل را با کامپایلر **C (Borlanc C)** باز کرده و کامپایل کرد .

با اجرای فایل **Yy.lex.c** در **Borlanc C** ممکن است با خطاهایی مواجه شویم که با رفع این خطاها و اجرا موفقیت آمیز فایل اجرایی که همان برنامه ی تحلیلگر لغوی می باشد ساخته می شود .

حال برای استفاده از این برنامه دو راه وجود دارد :

1) اجرای فایل اجرایی در خط فرمان و نوشتن متن در خط فرمان به عنوان ورودی و سپس دیدن جواب در همان خط فرمان .

2 نوشتن متن ورودی در یک فایل متنی و اجرای فایل
اجرائی در خط فرمان به همراه نام فایل ورودی و گرفتن یک
فایل متنی به عنوان خروجی .

در قسمت بعد انجام این مراحل به صورت عملی بیان می شود .

✓ کار با خط فرمان :

• اجرای Flex :

فایل **flex.exe** و فایل متنی که به زبان **flex** نوشته ایم و پسوند آن **.l** است را در مسیری مانند **d:** قرار می دهیم و در **command prompt** (خط فرمان ویندوز) به روش زیر عمل می کنیم :

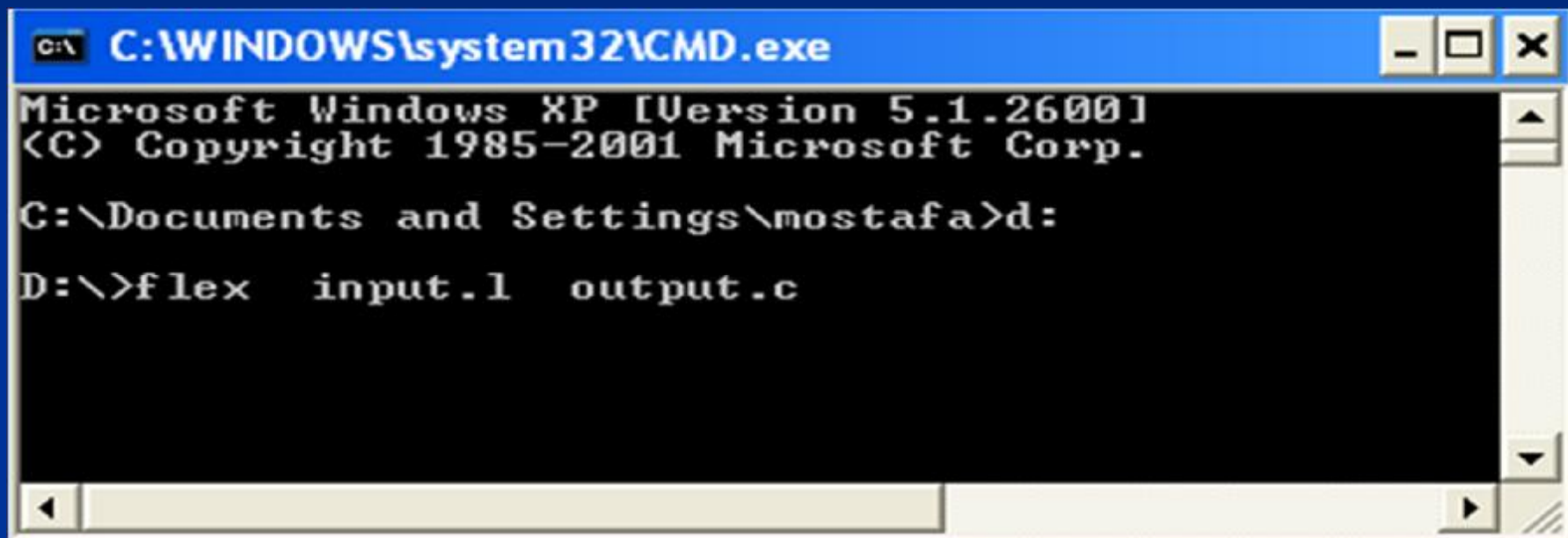
```
D:\flex ykǎ2ΘŸϕŸ kθ t2ΘŸϕŸ
```

ykǎ2ΘŸϕŸ : فایل متنی به زبان **Flex** (با پسوند **.l**)

kθ t2ΘŸϕŸ : فایل متنی به زبان **C** (با پسوند **.c**)

البته در شرایطی هم نمی توان نام فایل مقصد را وارد کرد و در این حالت فایل مقصد با نام **Yy.lex.c** تولید خواهد شد !!!

به عنوان مثال :



```
C:\WINDOWS\system32\CMD.exe
Microsoft Windows XP [Version 5.1.2600]
(C) Copyright 1985-2001 Microsoft Corp.

C:\Documents and Settings\mostafa>d:
D:\>flex input.l output.c
```

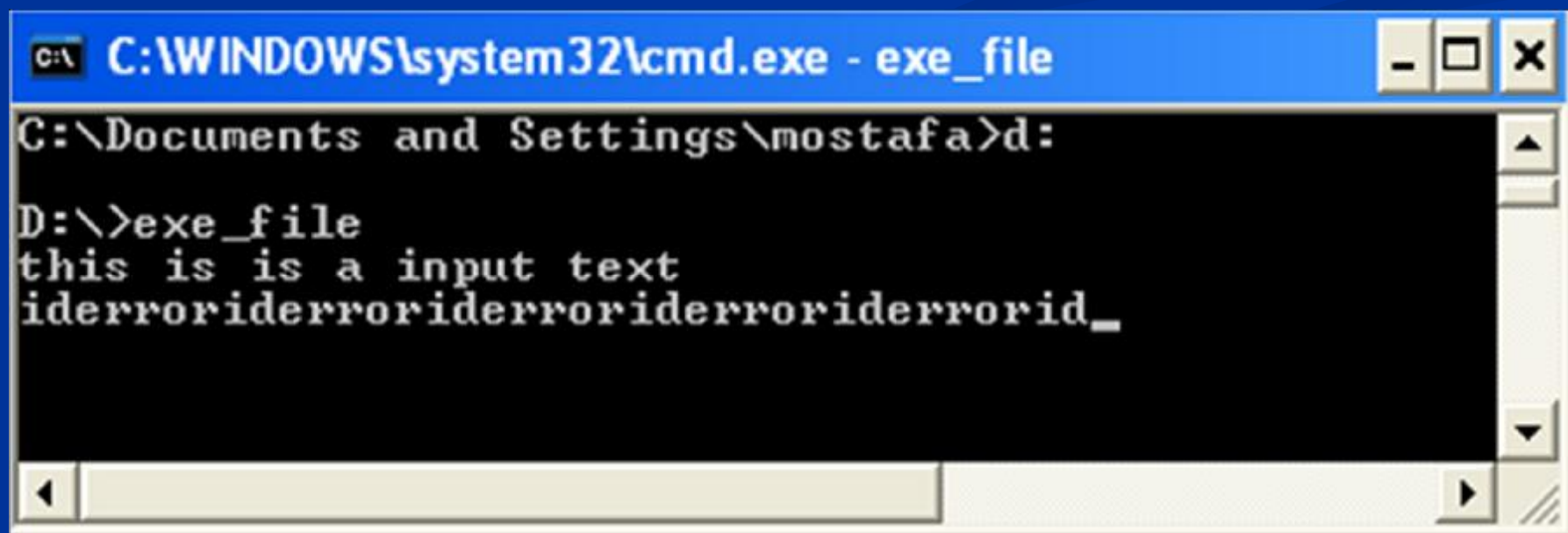
Flex برنامه ی مبدا (**input.l**) را خوانده و برنامه ی تحلیلگر لغوی را به زبان **c** تولید کرده و در فایل مقصد (**output.c**) می ریزد.

• اجرای تحلیلگر لغوی :

روش اول (نوشتن متن ورودی در خط فرمان) :

```
D:\3YfσH TËPŭ  
T2YHă3kYHYΣi2 2Y  
T2YHă2YHYK 2Y
```

به عنوان مثال :



```
C:\WINDOWS\system32\cmd.exe - exe_file  
C:\Documents and Settings\mostafa>d:  
D:\>exe_file  
this is is a input text  
iderroriderroriderroriderrorid_
```

Exe_file : فایل اجرایی (با پسوند **.exe**).

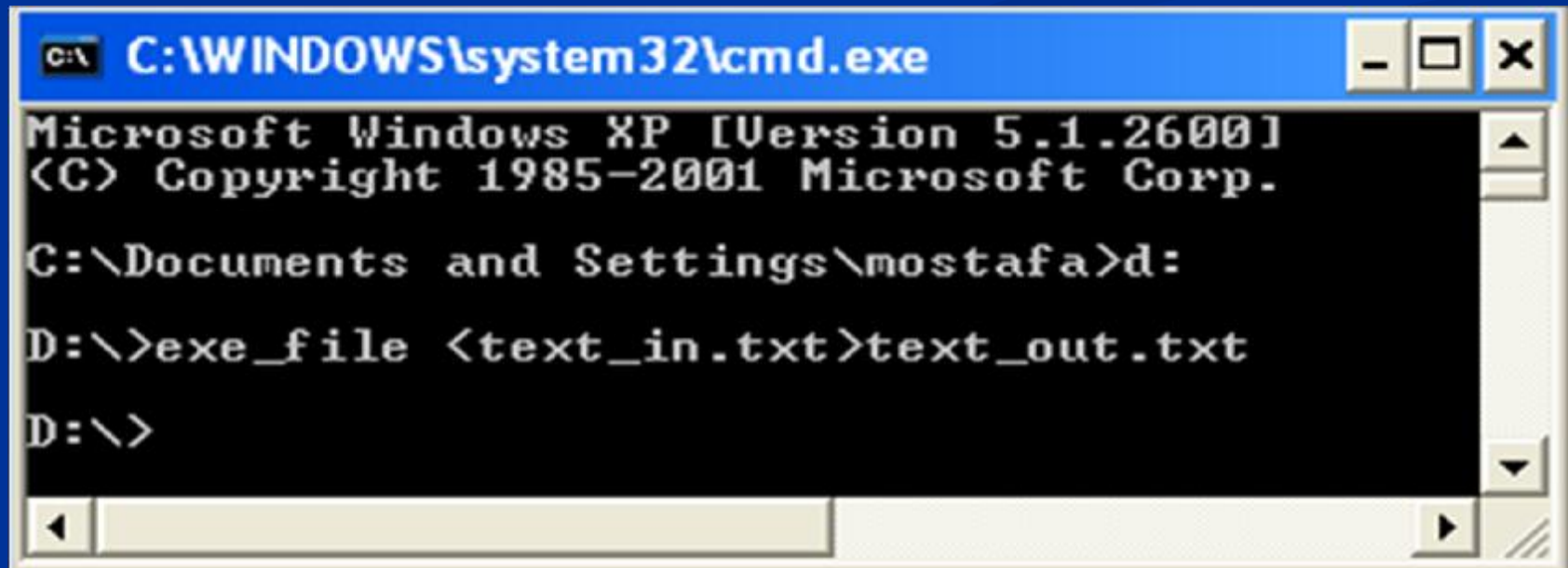
This is ... : متن به عنوان ورودی

Iderror ... : خروجی تحلیلگر لغوی

روش دوم (نوشتن متن ورودی در فایل متنی) :

D:\3Yf0H TĖPŭ <3kYHΥΘỸϕỸ> ΣÜYHκ ΘỸϕỸ

به عنوان مثال :



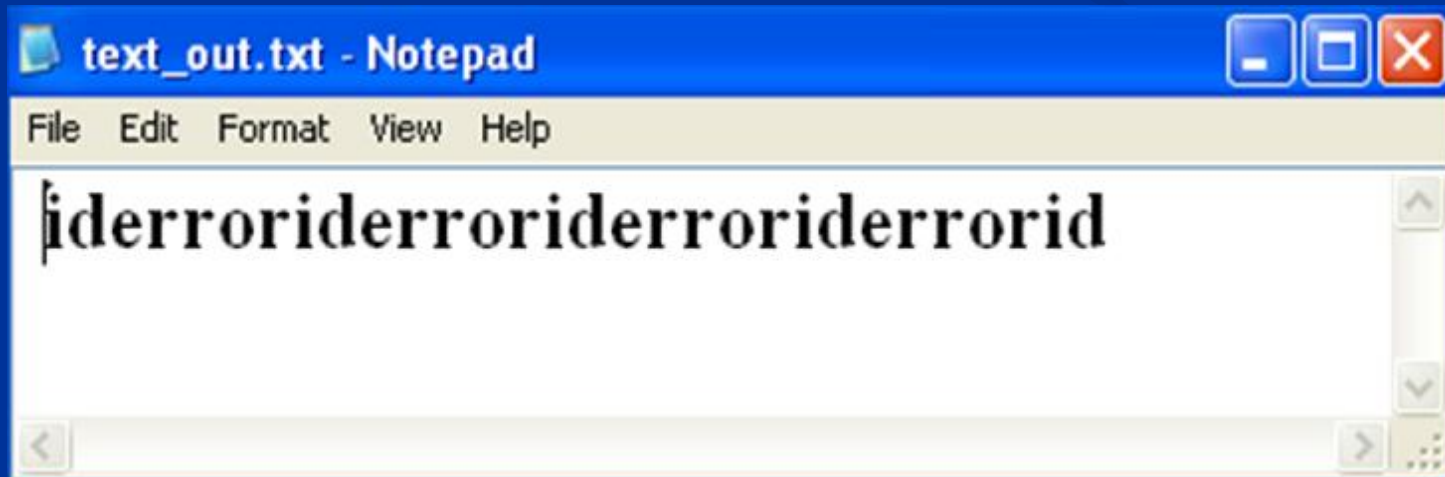
```
C:\WINDOWS\system32\cmd.exe
Microsoft Windows XP [Version 5.1.2600]
(C) Copyright 1985-2001 Microsoft Corp.

C:\Documents and Settings\mostafa>d:
D:\>exe_file <text_in.txt>text_out.txt
D:\>
```

Text_in.txt : فایل ورودی که در آن همان متن روش اول نوشته شده است .



Text_out.txt : فایل خروجی که در آن همان خروجی خط فرمان از روش اول نوشته شده است .



✓ آموزش زبان Flex :

متن درون فایل ورودی باید به زبان flex باشد تا نرم افزار flex آن را بشناسد و خروجی ای به زبان c بسازد . Flex عبارات با قاعده را می گیرد و تحلیلگر لغوی مطابق آن را می سازد . عبارت با قاعده را می توان به صورت رشته به flex معرفی کرد.

برای مثال :

“if”

عبارت با قاعده ی if

” while ”

عبارت با قاعده ی while

ساختار یک فایل *flex* :

(definition) تعاریف

%%

(translation) ترجمه ها

%%

(function) توابع

• تعاریف :

در قسمت تعاریف دو نوع اطلاعات قرار می گیرد :

1)

عبارت با قاعده نام

2)

%{

تعریف متغیرها

اعلان توابع

تعریف توابع

%}

چند مثال :

Lower [a-z]

Upper [A-Z]

Letter lower | upper

این قسمت به صورت کامل در فایلی که
به زبان C تولید می شود کپی می شود

```
%{  
#include<stdio.h>  
Int nchar,nline;  
%}
```

Digit [0-9]

Lower [a-z]

Upper [A-Z]

Letter lower | upper

Var { letter } | ({ letter } | { digit }) *

چند نکته :

این بخش اختیاری است و می توان برنامه ای نوشت که قسمت تعریف در آن نباشد. ولی استفاده از آن می تواند باعث خوانا تر شدن برنامه شود برای مثال می توان عبارت با قاعده ای را تعریف کرده و به آن نامی اختصاص داد و در قسمت های دیگر از آن نام استفاده کرد.

نحوه ی بیان عبارات با قاعده به زبان Flex در
جدول صفحه ی بعد توضیح داده شده :

نشان دهنده ي هر کاراکتر به جز کاراکتر سر خط است	.
رشته هاي يك کاراکتری که با کاراکتر منطبق هستند x	X
رشته هاي يك کاراکتری که با هر يك از کاراکتر هاي مشخص منطبق شوند	[مجموعه ي کاراکتر ها]
يکي از کاراکتر هاي مبدا تا مقصد	[کاراکتر مقصد – کاراکتر مبدا]
يکي الي چند تکرار از r	$r +$
صفر الي يك تکرار از r	$r ?$
صفر الي چند تکرار از r	$r *$
حد اقل m و حداکثر n تکرار از r	$r \{ m, n \}$

$x y$	يعني يا x يا y
r^*	صفر الي چند تکرار از r
xy	يعني اول x مي آید و بعد y
x/y	رشته هايي که قسمت اول آن با x و قسمت دوم آن با y منطبق باشد
[کاراکتر هاي مورد نظر [^]]	هر کاراکتر به جز کاراکتر هاي مشخص شده
$\backslash n$	براي نشان دادن سر خط
$\backslash b$	براي نشان دادن جاي خالي
$\backslash t$	براي نشان دادن tab

(عملگر ها و کاراکتر ها)

اولویت با عبارات داخل پرانتز می باشد و علاوه بر برای ترکیب عملگر ها استفاده کرد مانند : $(a|b)^*$

کاراکترهای مورد نظر $^{\wedge}$

کاراکتر های مورد نظری که حتما در اول خرد آمده باشند

{ }

برای استفاده ی کلماتی که در قسمت تعاریف آمده اند مانند:

Digit [0-9]

Int ({ digit })+

• ترجمه (این بخش اجباری است)

این بخش مشخص می کند هنگام کشف يك لغت مطابق یکی از عبارات با قاعده چه عملی باید انجام شود.

ساختار ترجمه :

{ عملیات } الگویی نشانه (عبارت با قاعده)

الگویی نشانه : عبارت با قاعده و یا یکی از اسامی تعریف شده در بخش تعاریف .

عملیات : دستورالعمل هایی به زبان c را نشان می دهد که هنگام یافتن دنباله ای از کاراکتر ها مطابق الگویی نشانه باید اجرا شوند.

مثال :

تعاریف	Digit	[0-9]
	Lower	[a-z]
	Upper	[A-Z]
	Letter	lower upper
	Var	{ letter } ({ letter } { digit }) *
	Ws	[\n\t] ⁺
ترجمه	%%	
	Ws	{ }
	“if”	{ printf (“ I found “if” keyword ”) ; }
	“else”	{ printf (“ I found “else” keyword) ; }
	var	{ printf (“I found variable ”) ; }
	%%	

حتماً باید آخرین خط نوشته شود (چرا؟)

ورودي: If temp else if id34

برنامه ي صفحه ي قبل

خروجي : I found "if" keyword
I found variable
I found "else" keyword
I found "if" keyword
I found variable

همان طور که ملاحظه گردید **تحلیلگر های لغوي** ساخته شده با کشف هر نوع لغت پیغام مناسب **چاپ** می کند . اما در یک کامپایلر واقعی **تحلیلگر لغوي** به جای چاپ یک پیغام باید یک نشانه مشخص را برای **تحلیلگر نحوي ارسال** کند . به این منظور به جای چاپ یک پیغام می توان از دستور **return** استفاده کرد . در این شرایط **تحلیلگر لغوي** با کشف هر لغت نشانه مناسب را باز می گرداند .

مثال :

```
%{  
# define    ASSIGNMENT    300  
# define    DIGIT    301  
%}  
%%  
":="      return(ASSIGNMENT );  
[0-9]+    return(DIGIT );  
"a"       return( 65 );
```

• توابع (این قسمت اختیاری می باشد)

Flex تحلیلگر لغوی (بخش ترجمه) را به یک تابع به نام **yylex()** تبدیل می کند. پس برای اجرای تحلیلگر لغوی باید این تابع را فراخوانی کرد.

مثال 1) برنامه ای به زبان **flex** بنویسید که تعداد کاراکتر ها و خطوط ورودی را شمارش کرده و نتیجه را چاپ کند.

(صفحه ی بعد)

```
%option noyywrap
```

به flex مي گويد كه فقط يك
فایل ورودی وجود دارد

```
%{
```

```
Int nchar , nline ;
```

```
%}
```

```
%%
```

```
[\n] { nline++ ; }
```

```
. { nchar++ ; }
```

```
%%
```

```
Int main ( void )
```

يعني هر کاراکتر به جز سرخط است

```
{
```

```
    yylex();
```

```
    printf( "%d %d" , nchar , nline+1 );
```

```
    return (0) ;
```

```
}
```

چرا با يك جمع كرديم؟

مثال 2) برنامه ی زیر اغلب لغات برنامه به زبان پاسکال را شناسایی کرده و بر حسب مورد پیغام مناسب را چاپ می کند.

Nquote [^']

%%

[a-zA-Z] ([a-zA-Z0-9])*	printf ("ID ") ;
“:=”	printf ("assignment ") ;
‘({ nquote } “)’	printf ("charakter_string ") ;
“.”	printf ("colon ") ;
“,”	printf ("comma ") ;
[0-9]+	printf ("digseq ") ;
“.”	printf ("dot ") ;
“=”	printf ("equal ") ;

“(”

“<”

“<>”

[0-9]+”.”[0-9]+

“)”

“^”

[\n\t\f]+

%%

Int main()

{

 yylex() ;

Return 0 ;

}

printf (“lparen ”) ;

printf (“lt ”) ;

printf (“notequal ”) ;

 printf (“realnumber ”) ;

printf (“rparen ”) ;

printf (“uparrow ”) ;

;

نکات مهم این مثال :

(1) در این برنامه کلمات کلیدی توسط عبارت با قاعده ي ***([a-zA-Z0-9] ([a-zA-Z]** پذیرفته مي شوند.

(2) در برخي موارد ممکن است **ناسازگاري هاي** بروز کند به عنوان مثال ممکن است بخش هاي مختلف يك دنباله توسط عبارات با قاعده مختلف پذیرفته شود به عنوان مثال فرض کنید تحلیلگر رشته ي "<>" را مي خواهد پردازش کند آیا وقتی کاراکتر اول را مي خواند ('<') پیام **it** چاپ مي کند يا هر دو کاراکتر را در کنار هم ('<>') پردازش کرده و پیام **notequal** چاپ مي شود ؟

“<” printf (“lt ”) ;

“<>” printf (“notequal ”) ;

دراين گونه موارد تحليلگر لغوي طولاني ترين رشته يعني <> را با توجه به عبارت با قاعده ي دوم مي پذيرد و **notequal** چاپ مي کند. به عبارتي مي توان گفت کل کاراکتر هاي بين دو علامت جداکننده (که معمولا همان جاي خالي است) به عنوان يك رشته انتخاب مي شود و سپس پردازش مي شود.

علاوه بر **yylex()** توابعي ديگر هم به صورت پيش فرض در **flex** وجود دارد که از آن جمله مي توان به موارد زير اشاره کرد :

Yytext() : رشته ي پذيرفته شده از ورودي در اين تابع قرار مي گيرد.

Yyleng() : طول رشته اي را که در **yytext** قرار دارد ، نگه مي دارد.

input() : کاراکتر بعدي در متن را برمي گرداند .

output(c) : کاراکتر **c** را در خروجي مي نويسد .

unput(c) : کاراکتر **c** را به متن پس مي فرستد تا بتواند توسط **input()** بعدي، خوانده شود .

yywrap() : هنگامي که کار تحليل به پايان فايل رسيد ، توسط

Flex فراخواني مي شود. اين تابع هميشه مقدار 1 را برمي گرداند

✓ کلمات کلیدی :

برای معرفی کلمات کلیدی به تحلیلگر لغوی از طریق **flex** دو روش وجود دارد :

روش اول) هر کلمه ی کلیدی را به عنوان یک نوع لغت در نظر گرفته و عبارت با قاعده ی مناسب را در لیست عبارات با قاعده درج می کنیم .

...

%%

“program”

“begin”

“end”

[a-zA-Z] ([a-zA-Z0-9])*

...

```
printf ( “program” ) ;
```

```
printf ( “begin” ) ;
```

```
printf ( “end” ) ;
```

```
printf ( “ id ” ) ;
```

باید بعد از کلمات کلیدی نوشته شود (چرا؟)

دقت : مکان تعریف عبارات با قاعده در برنامه مهم است . تحلیلگر لغوی تولیدی **flex** , مقایسه ی دنباله ی ورودی را با عبارات با قاعده به ترتیب از **بالا به پایین** انجام می دهد . بنا بر این عبارات با قاعده ی کلمات کلیدی باید قبل از عبارت با قاعده ی شناسه باشند در غیر این صورت تمام کلمات کلیدی شناسه در نظر گرفته می شوند.

روش دوم) کلمات کلیدی را در يك جدول (**جدول نماد ها**) به عنوان مقدار اولیه وارد می کنیم . هر گاه يك شناسه کشف می شود این شناسه در جدول جستجو می شود اگر این شناسه با کلمات کلیدی جدول یکسان باشد شناسه يك کلمه ی کلیدی است .

مثال :

```
%{  
#include<string.h>  
#include<iostream.h>  
char c;  
void toupper(char k[])  
{  
int i;  
for(i=0;i<=strlen(k);++i)  
if(k[i]<='z' && k[i]>='a')  
k[i]- =32;  
}  
int is_keyword(char id[])  
{  
char keyword [ 40 ][ 20 ] =  
{ "AND","ARRAY","BEGIN","CASE","CONST",  
"IF","FOR","END","WHILE","THEN" } ;
```

خروجي اين تابع رشته با
حروف بزرگ مي باشد

```

int i;
For ( I = 0 ; I < 40 ; i++ )
    if ( strcmp ( id , keyword[ I ] ) == 0 )
        return i;
return -1;
}
}%
%%
[a-zA-Z]([a-zA-Z0-9])* { toupper( yytext );
                        If ( is_keyword( yytext ) != -1 )
                            cout << yytext ;
                        else
                            cout << "ID" ;
                        }
[0-9]+ "." [0-9]+
                        cout << "REALNUMBER " ;

```

الغوريتم تشخيص
كلمات كليدي

"(*"

```
do {  
    c=input() ;  
    if ( c == '*' )  
    {  
        c = input() ;  
        if ( c == ')' )  
            break;  
        else  
            unput(c) ;  
    }  
} while ( 1 ) ;
```

"{"
[\n\t]

.

%%

int main()

{

yylex() ;

return 0 ;

}

```
while ( ( c = input() ) != '}' ) ;  
;  
Printf ( "" ) '\c' : illegal character at  
line %d '\n' "" , yytext[0] , line_no);
```

مثال : تحلیلگری که به حروف کوچک و بزرگ حساس نیست.

```
%{  
#include <stdio.h>  
int line_no=1;  
%}  
A [aA]  
B [bB]  
C [cC]  
.  
.  
.  
X [xX]  
Y [yY]  
Z [zZ]  
%%  
{A}{N}{D}  
{D}{O}  
[a-zA-Z]([a-zA-Z0-9])+
```

```
return ( AND );  
return ( DO );  
return ( IDENTIFIER );
```

"(*" | "{"

نوع register باعث
مي شود كه يك كپي
از اين متغير در
cache قرار بگيرد
تا سرعت عمليات بالا
رود كه براي متغير
هاي پر کاربرد مفيد
است

```
{ register int c;  
while ( (c = input() ) )  
{  
    if ( c == '}' )  
        break;  
    else if ( c == '*' )  
    {  
        if ( ( c = input() ) == '}' )  
            break;  
        else  
            unput(C) ;  
    }  
    else if ( c == '\n' )  
        line_no ++ ;  
    else if ( c == 0 )  
        commenteof() ;  
}  
};
```

کاراکتر را در نوع
int ریخته ایم اما
برنامه error
نمی دهد (چرا؟)

[\\t\\f]

```
\n                line_no++;  
%%  
commentof()  
{  
    printf ( %d\n , line_no ) ;  
    exit (1) ;  
}  
yywrap()  
{  
    Return (1) ;  
}
```


مثال : تعداد کاراکتر و خط

```
%option noyywrap
int num_lines=0;num_chars=0;
%%
\n  ++num_lines;++num_chars;
.  ++num_chars;
%%
Main ( int argc , char **argv ) {
  + + argv , - - argc;
If ( argc == 0 )
  yyin = fopen ( argv [ 0 ] , "r" ) ;
else
  yyin = stdin ;
yylex() ;
printf("# of lines=%d,# of chars = %d\n" , num_lines ,num_chars ) ;
}
```

تعداد پارامتر های ورودی

اشاره گر به اول آرایه ی
پارامتر های ورودی

در صورت نبود پارامتر
ورودی اطلاعات از روی
فایل خوانده می شود

در صورت وجود پارامتر ورودی اطلاعات از صفحه کلید خوانده می شود

مثال :

```
%{  
int mylineno=0;  
void handle comments(void);  
%}
```

```
string  
char  
prepro
```

```
ws  
alpha  
dig
```

```
name  
int_num  
num1  
num2  
Flt_num
```

```
\"[^\\n"]+\\"
```

```
\'[^\\n\']+\'
```

```
^#.*
```



```
[\\t]+
```

```
[A-Za-z]
```

```
[0-9]
```

```
{alpha}({alpha}|{dig})*
```

```
[-+]?{dig}
```

```
[-+]?{dig}+\\.?([eE][-+]?{dig}+)?
```

```
[-+]?{dig}*\\. {dig}+([eE][-+]?{dig}+)?
```

```
{ num1 } | { num2 }
```

یعنی اگر اول خطی با # شروع شده بعد از آن هرچه کاراکتر آمده بود عملیاتی انجام نده (مورد استفاده برای دستورات #include (<...>))

یعنی وقتی به // رسیدی بعد از آن هر کاراکتری و
به هر تعداد بود (به جز کاراکتر خط جدید) برای آن
هیچ عملیاتی انجام نده

%%

{ws}

{prepro}

"//".*

"/**"

\n

{int_num}

{flt_num}

{name}

{string}

.

{ handle_comments() ; }

{ mylineno++ ; }

{ printf("integer number[%s]\n" , yytext) ; }

{ printf("floatin-point number[%s]\n" , yytext) ; }

{ printf ("name[%s]\n" , yytext) ; }

{ printf ("string[%s]\n" , yytext) ; }

{ printf ("don't recognise[%s]\n" , yytext) ; }

%%

```
const char *keywords[ ] = {  
"auto","break","char","const","switch","sizeof","  
static","long" };
```

```
const int keywords_length = sizeof(keywords) / sizeof(keywords [ 0 ] );
```

```
int main()
```

```
{
```

```
    printf ( "**STARTING LEXICAL ANALYSIS**\N" );
```

```
    while ( yylex() != 0 );
```

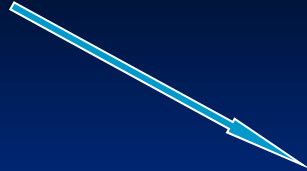
```
    printf ( "**FINISHED LEXICAL ANALYSIS**\N" );
```

```
    return 0 ;
```

```
}
```

```
void handle_comments(void)
```

```
{  
int c;  
While ( ( c = input() ) != 0 )  
{  
    if ( c == '\n' )  
        ++ mylineno ;  
    else if ( c == '*' )  
    {  
        if ( ( c = input() ) == '/' )  
            break ;  
        else  
            unput(c) ;  
    }  
}  
}
```



این تابع در صورت ورود به
قسمت توضیحات تا پایان آن
پیش می رود و پردازش
خاصی انجام نمی دهد فقط
اگر خط جدیدی ایجاد شود به
شمارنده ی خط يك واحد
اضافه می کند