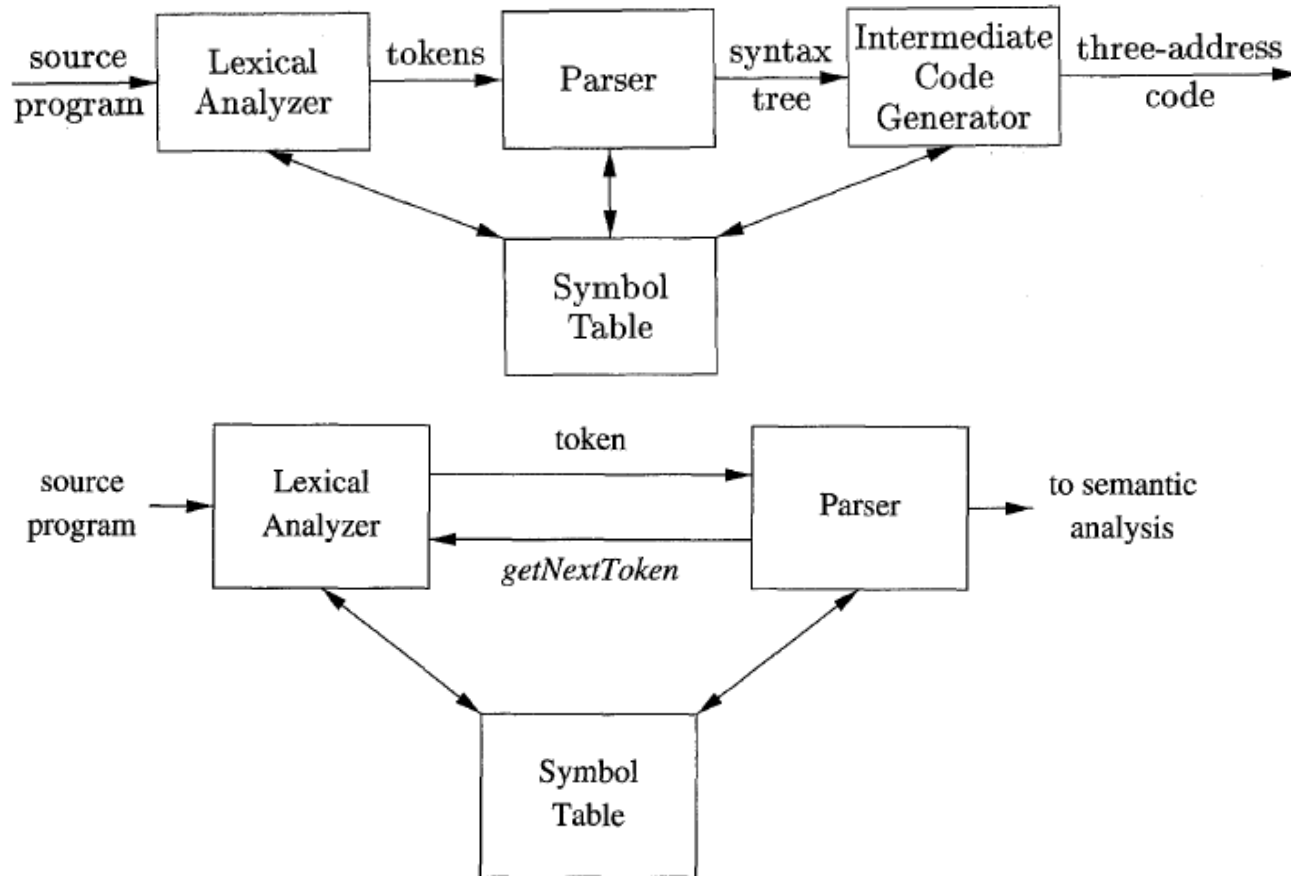

محلل المفردات المحلل اللفظي

Lexical Analyzer

طرق بناء محلل المفردات

- الطريقة اليدوية: يجب كتابة كود للتعرف على كل مفردة من اللغة وإعادة معلومات عنها.
- توليد المحلل اللفظي: من خلال توصيف لنماذج مفردات اللغة. يوجد العديد من التطبيقات التي تعتمد على اكتشاف النماذج
 - أسهل لتعديل المحلل.
 - أسرع كون توصيف مفردات اللغة يتم في مستوى تجريد عال بينما تفاصيل الكود تقع على عاتق مولد الكود.

المحلل اللفظي في سلسلة الترجمة



دور محلل المفردات

□ Scanninig: حذف التعليقات والفراغات.

□ Lexical Analysis: إنتاج سلسلة من مفردات اللغة الموجودة في نص البرنامج.

■ Token: زوج من اسم المفردة (نوع الوحدة اللفظية) وقيمتها (إختيارية).

■ Pattern: كون توصيف مفردات اللغة يتم في مستوى تجريد عال بينما تفاصيل الكود تقع على عاتق مولد الكود.

□ من أجل الكلمات المحجوزة: النموذج هو سلسلة المحارف التي تشكل الكلمة.

□ من أجل المعارف: النموذج هو بنية أعقد تتطابق مع مجموعة كبيرة من الكلمات.

■ Lexeme: سلسلة محارف من الكود المصدري تطابق نموذج مفردة، يكتشفها المحلل اللفظي كغرض من تلك المفردة.

Examples of tokens

Token	Sample Lexemes	Description
const	const	const
if	if	if
else	else	else
relation	<, <=, ==, >, >=	< or <= or == or
id	pi, count, D5	letter followed by letters and digits
unm	3.14, 0, 6.02E23	any numeric constant
literal	"hello world"	any characters between " and " except "

□ Ex: const pi = 3.14;

Token = id

Lexeme = pi

Tokens

Token	pattern
Keywords	the keyword itself
Operators	relation operators themselves.
Identifiers	regular expression
Constants: Numbers & Literals	regular expression
Punctuation symbols as parentheses, commas, semicolons	Symbols themselves

Regular Expressions

- **Ex1**: identifiers are strings of letters, digits, and underscores:
 - $letter_ \rightarrow A/B/.../Z/a/b/.../z/_$
 - $digit \rightarrow 0/1/2/.../9$
 - $Id \rightarrow letter_ (letter_ / digit)^*$
- **Ex2**: Unsigned numbers (integer or floating point) are strings such as 5280, 0.01234, 6.336E4, or 1.89E-4.
 - $digit \rightarrow 0/1/2/.../9$
 - $digits \rightarrow digit digit^*$
 - $optionalFraction \rightarrow . digits / \epsilon$
 - $optionalExponent \rightarrow (E (+ / - / \epsilon) digits) / \epsilon$
 - $number \rightarrow digits optionalFraction optionalExponent$

Extensions of Regular Expressions

Lex (Unix)

- *Operator $^+$* : one or more instances.
- *Operator $^?$* : zero or one instance.
- *Character classes*:
 - $[abc] = a|b|c$
 - $[a-z] = a|b|\dots|z$
 - **Ex:** $Id \rightarrow [A-Za-z][A-Za-z0-9]^*$
 - **Ex:** $Id \rightarrow [A-Za-z_][A-Za-z0-9_]^*$

Extensions of Regular Expressions

- **Ex1**: identifiers are strings of letters, digits, and underscores:
 - $letter_ \rightarrow [A-Za-z_]$
 - $digit \rightarrow [0-9]$
 - $Id \rightarrow letter_ (letter_ / digit)^*$
- **Ex2**: Unsigned numbers (integer or floating point) are strings such as 5280, 0.01234, 6.336E4, or 1.89E-4.
 - $digit \rightarrow [0-9]$
 - $digits \rightarrow digit^+$
 - $number \rightarrow digit^+ (. Digits)? (E[+-]? digits)?$

Extensions (Flex) of Regular Expressions

EXPRESSION	MATCHES	EXAMPLE
<i>c</i>	the one non-operator character <i>c</i>	a
<i>\c</i>	character <i>c</i> literally	*
<i>"s"</i>	string <i>s</i> literally	"**"
<i>.</i>	any character but newline	a.*b
<i>^</i>	beginning of a line	^abc
<i>\$</i>	end of a line	abc\$
<i>[s]</i>	any one of the characters in string <i>s</i>	[abc]
<i>[^s]</i>	any one character not in string <i>s</i>	[^abc]
<i>r*</i>	zero or more strings matching <i>r</i>	a*
<i>r+</i>	one or more strings matching <i>r</i>	a+
<i>r?</i>	zero or one <i>r</i>	a?
<i>r{m,n}</i>	between <i>m</i> and <i>n</i> occurrences of <i>r</i>	a[1,5]
<i>r₁r₂</i>	an <i>r₁</i> followed by an <i>r₂</i>	ab
<i>r₁ r₂</i>	an <i>r₁</i> or an <i>r₂</i>	a b
<i>(r)</i>	same as <i>r</i>	(a b)
<i>r₁/r₂</i>	<i>r₁</i> when followed by <i>r₂</i>	abc/123

Recognition of Tokens

□ **Language Rules**

- $stmt \rightarrow if\ expr\ then\ stmt$
| $if\ expr\ then\ stmt\ else\ stmt$
| ϵ
- $expr \rightarrow term\ relop\ term$
| $term$
- $term \rightarrow id$
| $number$

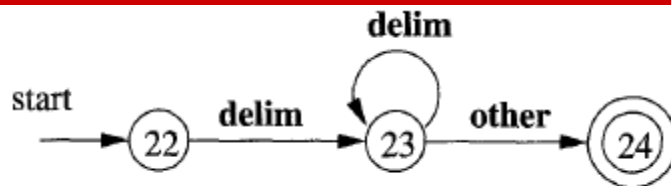
□ **Patterns ...**

- $ws \rightarrow (blank\ | \ tab\ | \ newline)^+$
- $digit \rightarrow [0-9]$
- $digits \rightarrow digit^+$
- $number \rightarrow digit^+ (. \ Digits)^? (E[+-] ? digits)^?$
- $letter \rightarrow [A-Za-z]$
- $id \rightarrow letter(letter/digit)^*$
- $if \rightarrow if$
- $then \rightarrow then$
- $else \rightarrow else$
- $relop \rightarrow < / > / > = / < = / = / < >$

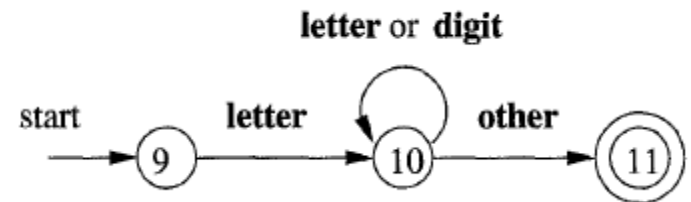
Recognition of Tokens

Lexemes	Tokens	Token attributes
Any ws	-	-
if	if	-
then	then	-
else	else	-
Any id	id	Pointer to entry table
<	relop	LT
<=	relop	LE
>	relop	GT
>=	relop	GE
=	relop	EQ
<>	relop	NE

Recognition of Tokens

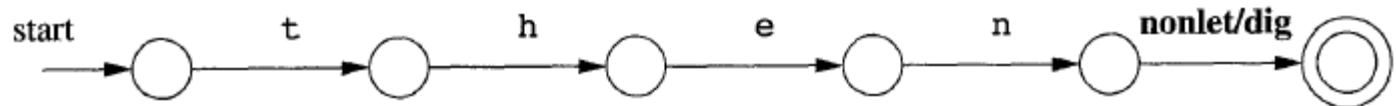


Recognition of white spaces

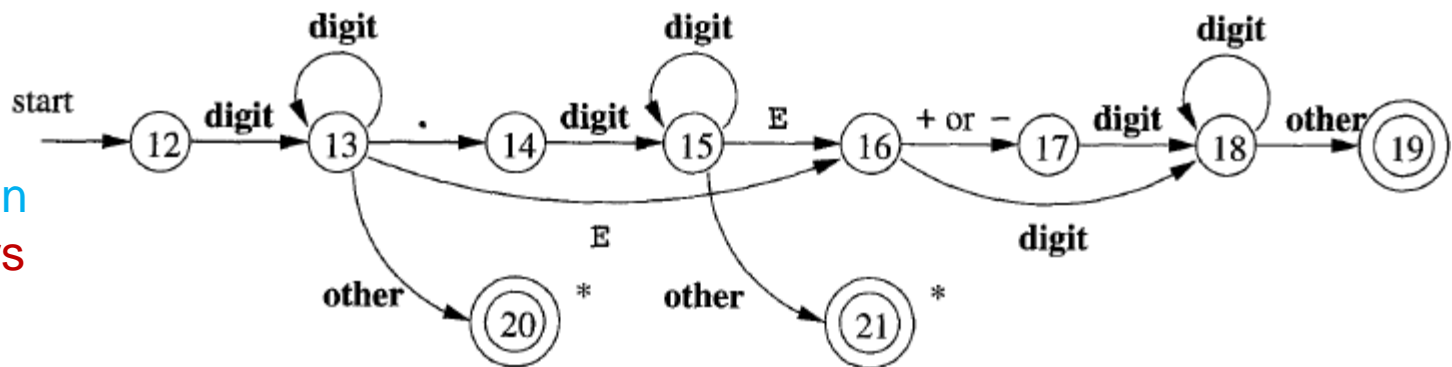


Recognition of identifiers

Recognition
Of then



Recognition
Of numbers



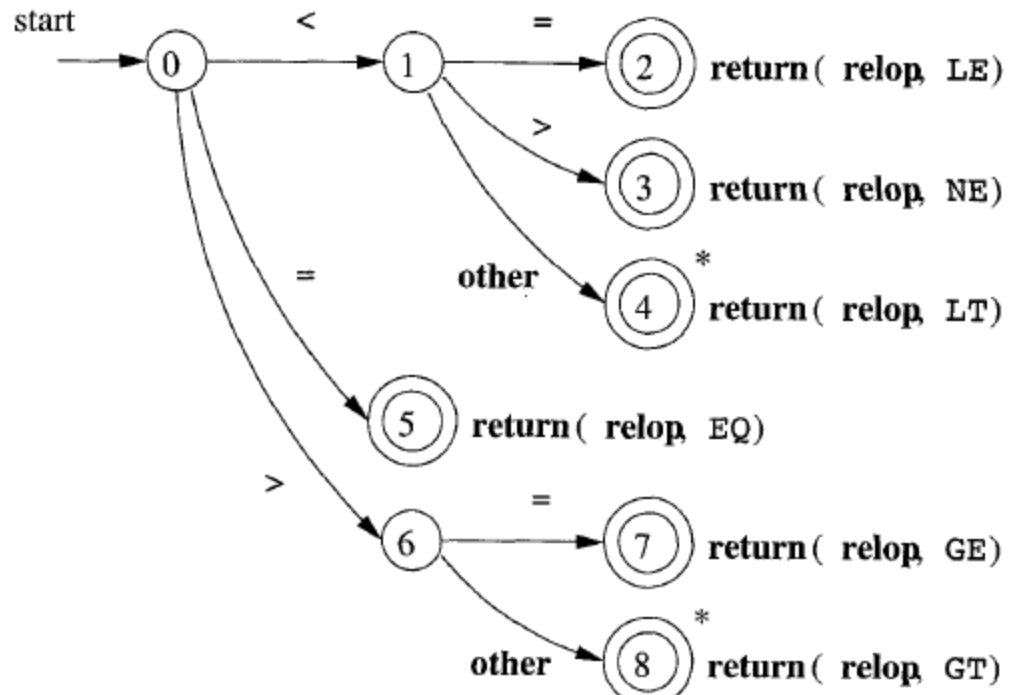
Recognition of Tokens

```

TOKEN getRelop() {
    TOKEN retToken = new(RELOP);
    while(1) {
        switch(state) {
            case 0: c = nextchar();
                if ( c == '<' ) state = 1;
                else if ( c == '=' ) state = 5;
                else if ( c == '>' ) state = 6;
                else fail();
                break;

            case 1: ...

            case 8: retract();
                    retToken.attribute = GT;
                    return(retToken);
                    break;
        }
    }
}
    
```



Recognition of relation

Next start state

```
state  $\leftarrow$  0; start  $\leftarrow$  0;  
lexcal_value;  
int fail() {  
    forward = token_beginning;  
    switch(start) {  
        case 0: start  $\leftarrow$  9; break;  
        case 9: start  $\leftarrow$  12; break;  
        case 12: start  $\leftarrow$  20; break;  
        case 20: start  $\leftarrow$  25; break;  
        case 25: recover;  
    }  
}
```

Code of Lexical Analyzer ...

```
token nexttoken() {
    while(1) {
        switch(state) {
            case 0 : c ← nextchar();
                if (c== blank || c == tab || c == newline) {
                    state ← 0
                    lexeme_beginning++;
                }
            elseif (c == '<') state ← 1;
            elseif (c == '=') state ← 5;
            elseif (c == '>') state ← 6;
            else state = fail();
            break;
            //....
            case 9: c← nextchar();
                //...
```

□ في حال وجود عدة حالات تطابق:

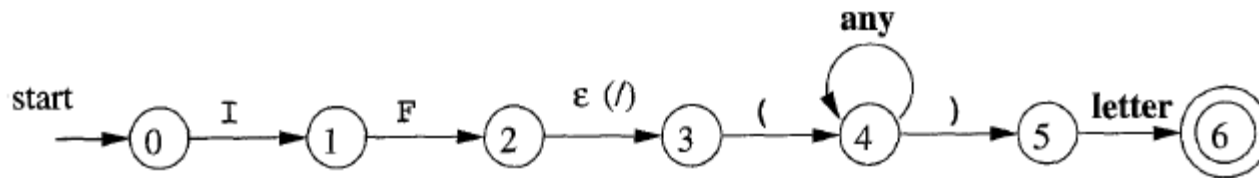
□ نختار الـ lexeme

الأطول.

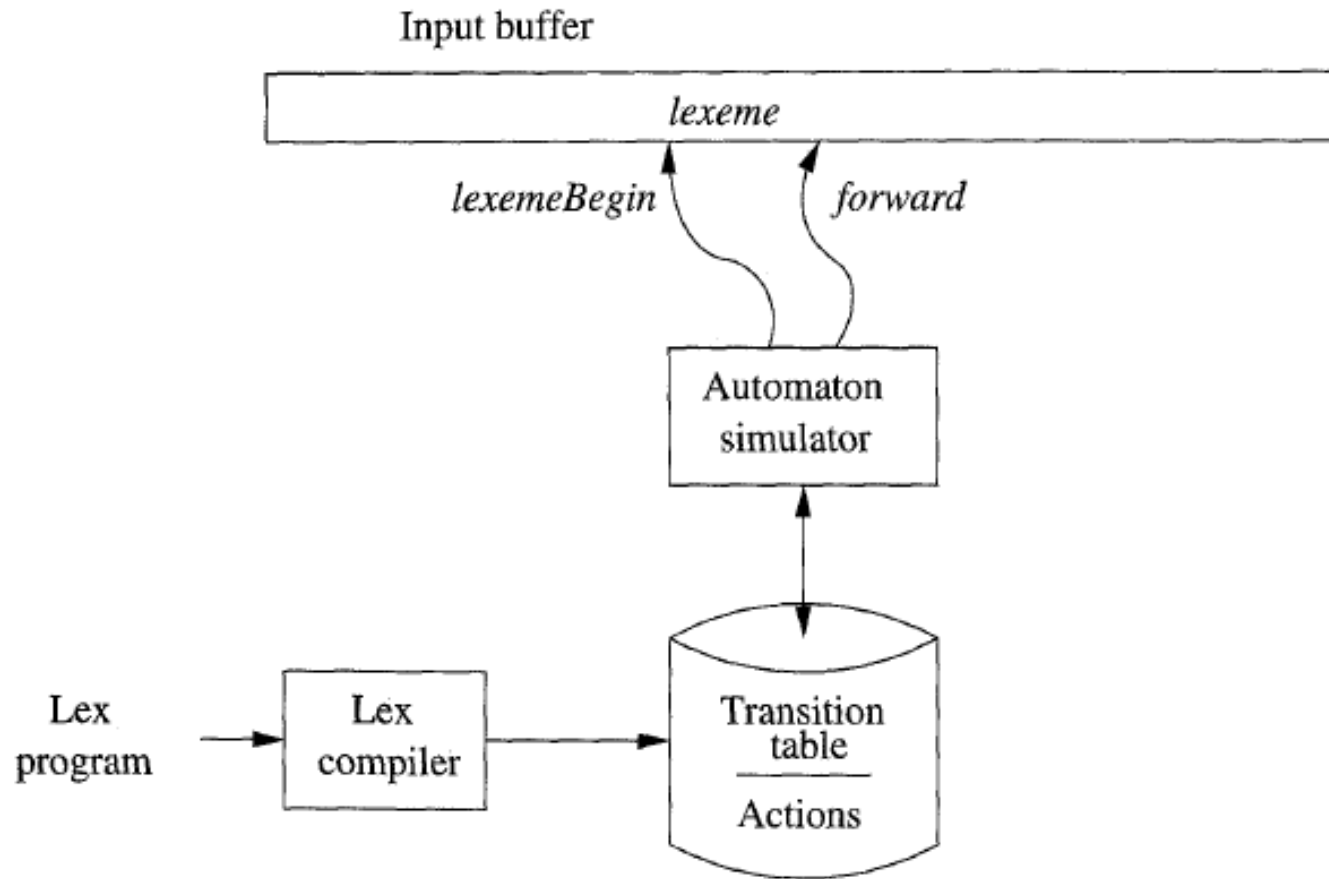
□ وعند تساوي طول أكثر من

lexeme نختار المعرف
أولاً...

Conditional recognition of lexeme





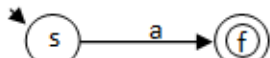
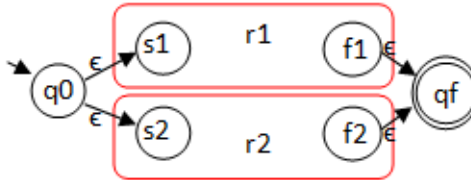
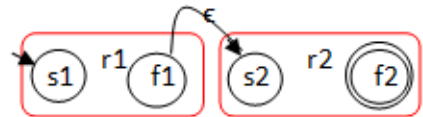
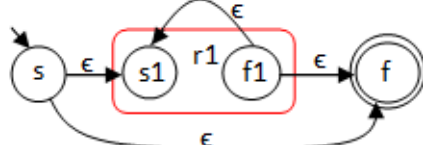
Generation of Lexical Analyzer



Transformation

Regular expressions to ϵ -NFA

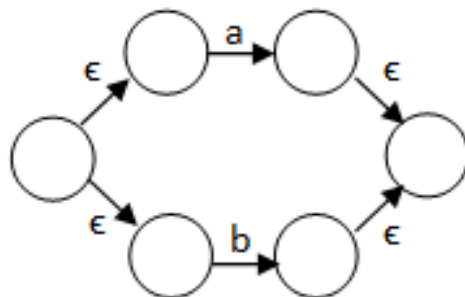
□ RE operations

Regular Expression	Language	Finite automata
Φ	$\{\}$	
ϵ	$\{\epsilon\}$	
a	$\{a\}$	
$r1 + r2$	$L(r1+r2) = L(r1) + L(r2)$	
$r1.r2$	$L(r1.r2) = L(r1).L(r2)$	
$r1^*$	$L(r1^*) = L(r1)^*$	

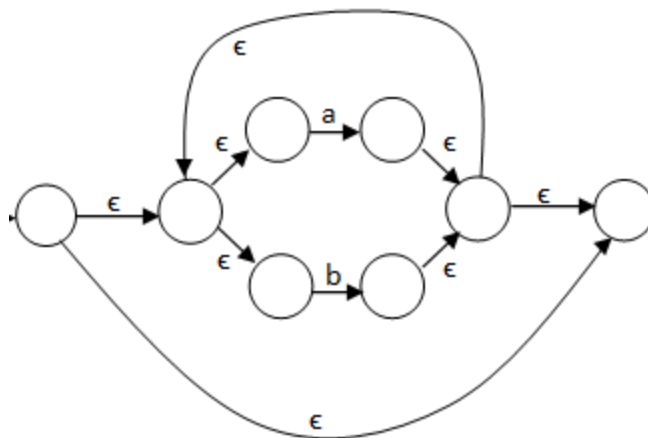
RE \rightarrow ϵ -NFA

□ Transform $a.(a+b)^*.b.b$ to ϵ -NFA

□ $(a+b)$



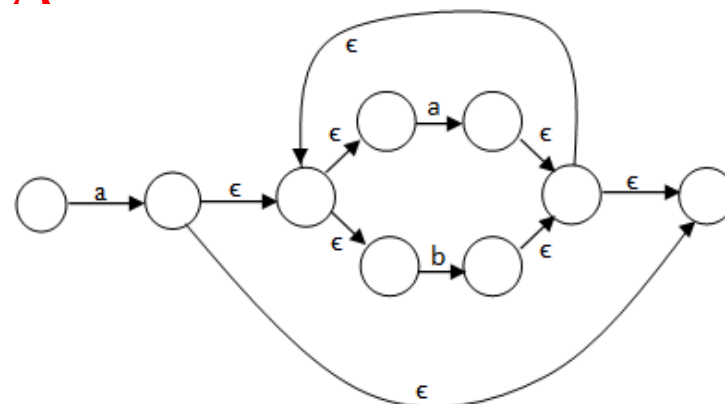
□ $(a+b)^*$



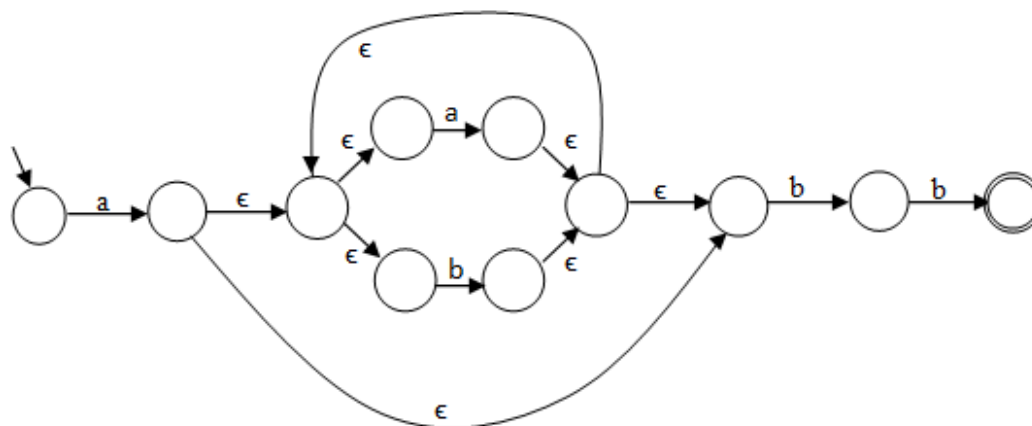
RE \rightarrow ϵ -NFA

□ Transform $a.(a+b)^*.b.b$ to ϵ -NFA

□ $a.(a+b)^*$



□ $a.(a+b)^*.b.b$



Simulation ϵ -NFA ...

```
s  $\leftarrow$   $\epsilon$ -closure({s0})  
a  $\leftarrow$  nextchar  
while(a  $\neq$  eof) do  
    s  $\leftarrow$   $\epsilon$ -closure (move(s, c))  
    c  $\leftarrow$  nextchar  
end while  
if s in finite states then  
    return 'yes'  
else  
    return 'no'
```

move	a	b	ϵ
0	{0, 1}	{0}	Φ
1	Φ	{1}	Φ
2	Φ	Φ	Φ
3	Φ	Φ	Φ

Transition Table

Simulation NFA ... (ϵ -closure)

push all states of T onto stack

initialize ***E-closure(T) to T***

while (stack is not empty) {

 pop t , the top element of stack;

for (each state u with an edge from t to u labeled ϵ){

if (u is not in *e-closure(T)*) {

 add u to *e-closure(T)*;

 push u onto stack

 }

 }

}

Transform ϵ -NFA to DFA

initially, $e\text{-closure}(s_0)$ is the only state in $Dstates$, and it is unmarked;

while (there is an unmarked state T in $Dstates$) {

 mark T ;

for (each input symbol a) {

$U = E\text{-closure}(\text{move}(T, a));$

if (U is not in $Dstates$)

 add U as an unmarked state to $Dstates$;

$Dtran[T, a] = U$;

 }

}

simulation DFA ...

```
s ← s0
c ← nextchar
while(c != eof) do
    s ← move(s, c)
    c ← nextchar
end while

if s in finite states then
    return 'yes'
else
    return 'no'
```

Thanks