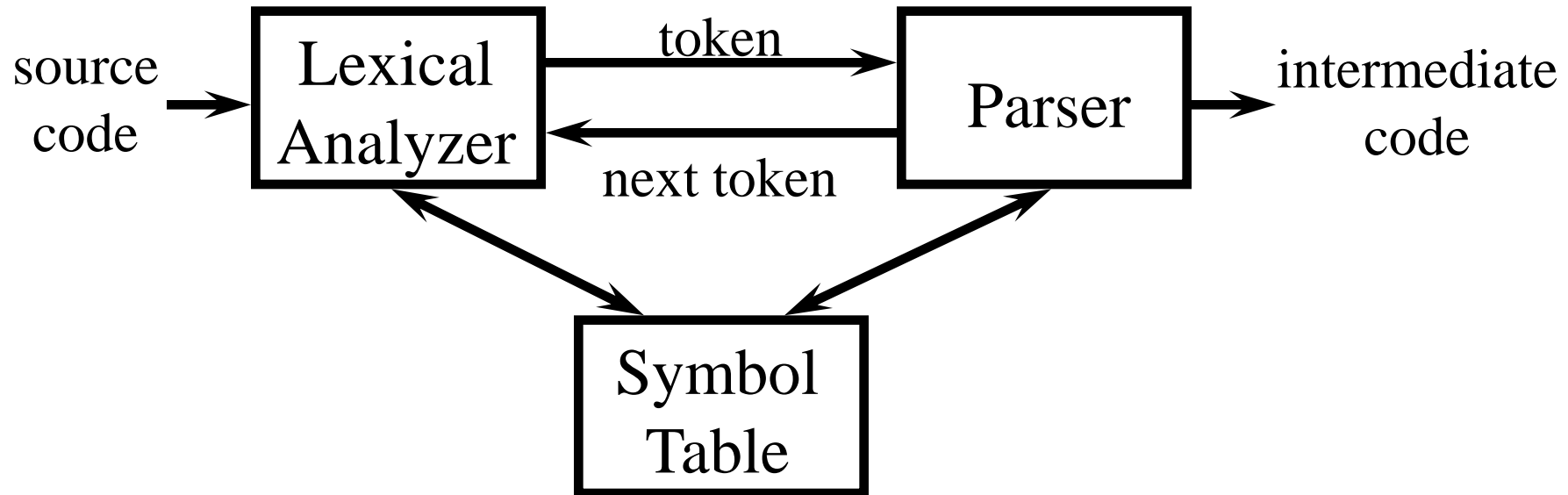


Lexical Analysis

Contents

- ♠ Introduction to lexical analyzer
- ♠ Tokens
- ♠ Regular expressions (RE)
- ♠ Finite automata (FA)
 - deterministic and nondeterministic finite automata (DFA and NFA)
 - from RE to NFA
 - from NFA to DFA
- ♠ Flex - a lexical analyzer generator

Introduction to Lexical Analyzer



Tokens

♠ **Token** (language): a **set of strings**

- if, identifier, relop

♠ **Pattern** (grammar): a **rule** defining a token

- if: if

- identifier: letter followed by letters and digits

- relop: < or <= or = or <> or >= or >

♠ **Lexeme** (sentence): a **string** matched by the pattern of a token

- if, Pi, count, <, <=

Attributes of Tokens

♠ **Attributes** are used to distinguish different lexemes in a token

- < if, >
- < identifier, **pointer to symbol table entry** >
- < relop, '=' >
- < number, **value** >

♠ Tokens affect **syntax analysis** and attributes affect **semantic analysis**

Regular Expressions

- ♠ ϵ is a RE denoting $\{\epsilon\}$
- ♠ If $a \in \text{alphabet}$, then a is a RE denoting $\{a\}$
- ♠ Suppose r and s are RE denoting $L(r)$ and $L(s)$
 - $(r) \mid (s)$ is a RE denoting $L(r) \cup L(s)$
 - $(r) (s)$ is a RE denoting $L(r)L(s)$
 - $(r)^*$ is a RE denoting $(L(r))^*$
 - (r) is a RE denoting $L(r)$

Examples

- ♠ a / b $\{a, b\}$
- ♠ $(a / b)(a / b)$ $\{aa, ab, ba, bb\}$
- ♠ a^* $\{\epsilon, a, aa, aaa, \dots\}$
- ♠ $(a / b)^*$ the set of all strings of a 's and b 's
- ♠ a / a^*b the set containing the string a and all strings consisting of zero or more a 's followed by a b

Regular Definitions

♠ **Names** for regular expressions

$$d_1 \rightarrow r_1$$
$$d_2 \rightarrow r_2$$

...

$$d_n \rightarrow r_n$$

where r_i over alphabet $\cup \{d_1, d_2, \dots, d_{i-1}\}$

♠ **Examples:**

$$\text{letter} \rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z$$
$$\text{digit} \rightarrow 0 \mid 1 \mid \dots \mid 9$$
$$\text{identifier} \rightarrow \{\text{letter}\} (\{\text{letter}\} \mid \{\text{digit}\})^*$$

Notational Shorthands

♠ One or more instances

$(r)^+$ denoting $(L(r))^+$

$$r^* = r^+ \mid \varepsilon$$

$$r^+ = r r^*$$

♠ Zero or one instance

$$r^? = r \mid \varepsilon$$

♠ Character classes

$$[abc] = a \mid b \mid c$$

$$[a-z] = a \mid b \mid \dots \mid z$$

$$[^a a-z] = \text{any character except } [a-z]$$

Examples

delim → [\t\n]

ws → {delim}+

letter → [A-Za-z]

digit → [0-9]

id → {letter}({letter}|{digit})*

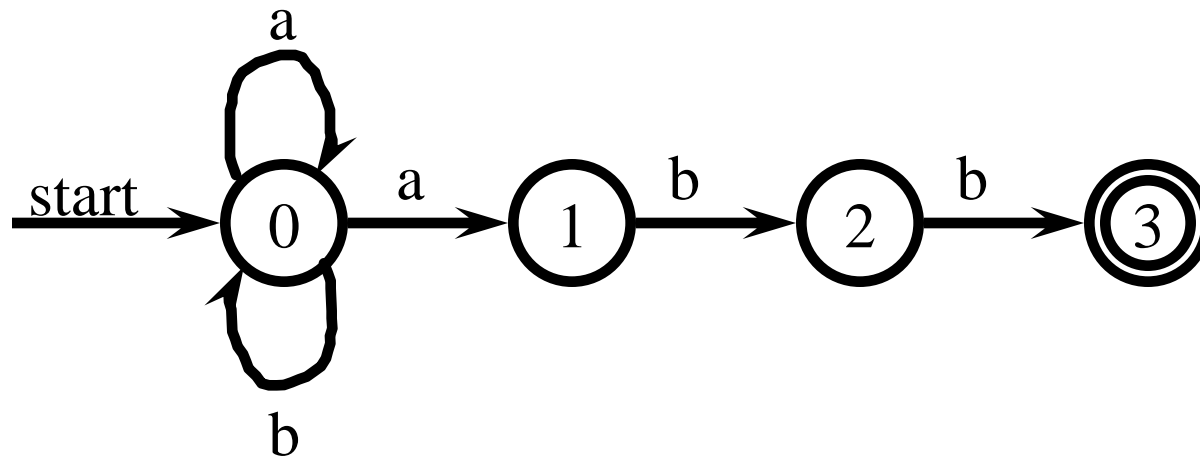
number → {digit}+(.{digit}+)?(E[+\-]?{digit}+)?

Nondeterministic Finite Automata

- ♠ An NFA consists of
 - A finite set of *states*
 - A finite set of *input symbols*
 - A *transition function* (or *transition table*) that maps (state, symbol) pairs to sets of states
 - A state distinguished as *start state*
 - A set of states distinguished as *final states*

Transition Diagram

$(a \mid b)^*abb$



An Example

♠ RE: $(a \mid b)^*abb$

♠ States: $\{0, 1, 2, 3\}$

♠ Input symbols: $\{a, b\}$

♠ Transition function:

$(0,a) = \{0,1\}, \quad (0,b) = \{0\}$

$(1,b) = \{2\}, \quad (2,b) = \{3\}$

♠ Start state: 0

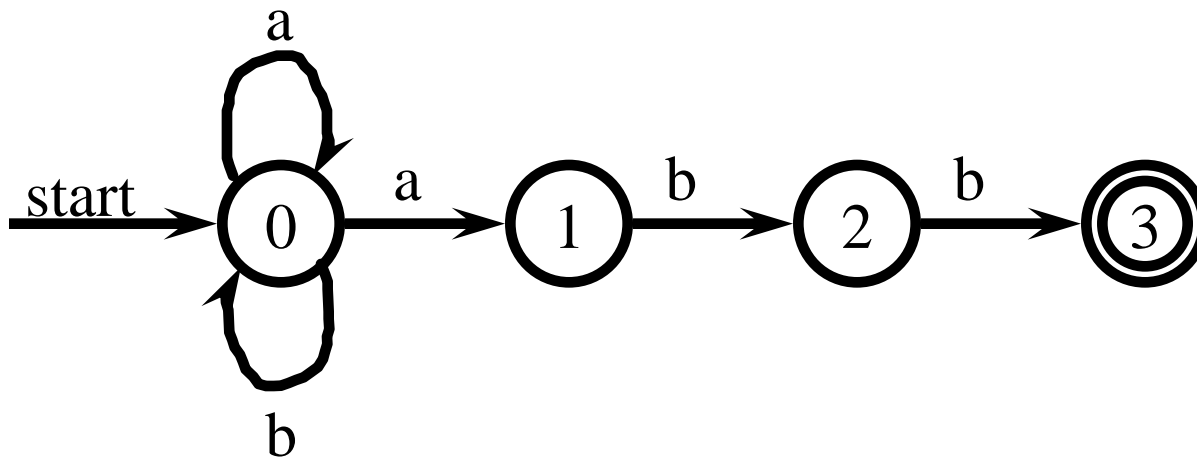
♠ Final states: $\{3\}$

Acceptance of NFA

- ♠ An NFA *accepts* an input string s *iff* there is *some* path in the transition diagram from the *start state* to some *final state* such that the edge labels along this path spell out s

An Example

$(a \mid b)^*abb$



$abb: \{0\} \xrightarrow{a} \{0, 1\} \xrightarrow{b} \{0, 2\} \xrightarrow{b} \{0, 3\}$

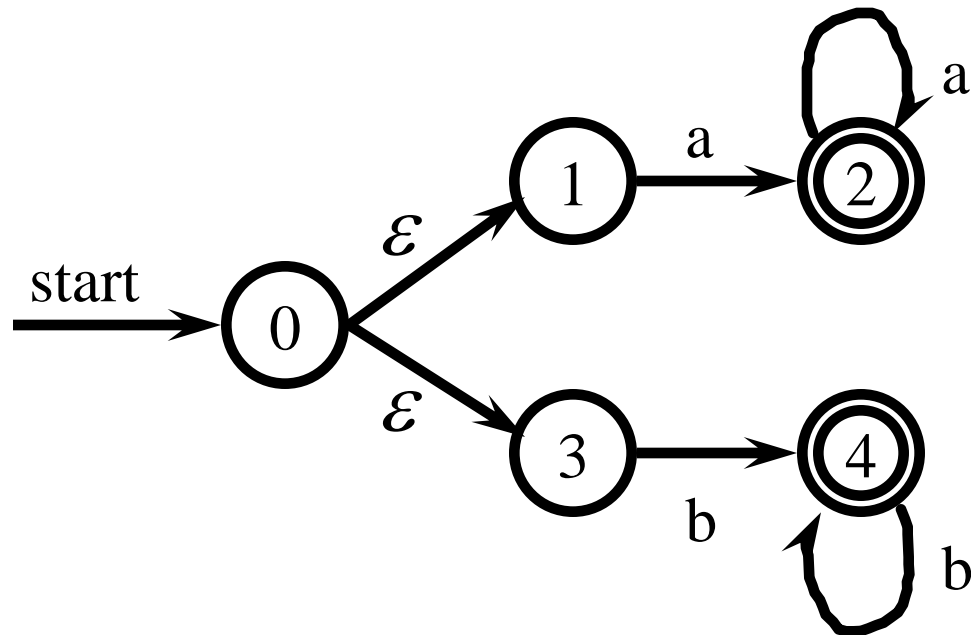
$aabb: \{0\} \xrightarrow{a} \{0, 1\} \xrightarrow{a} \{0, 1\} \xrightarrow{b} \{0, 2\} \xrightarrow{b} \{0, 3\}$

abb
 $aabb$
 $babb$
 $aaabb$
 $ababb$
 $baabb$
 $bbabb$

...

Transition Diagram

$aa^* \mid bb^*$



Another Example

♠ RE: $aa^* \mid bb^*$

♠ States: $\{0, 1, 2, 3, 4\}$

♠ Input symbols: $\{a, b\}$

♠ Transition function:

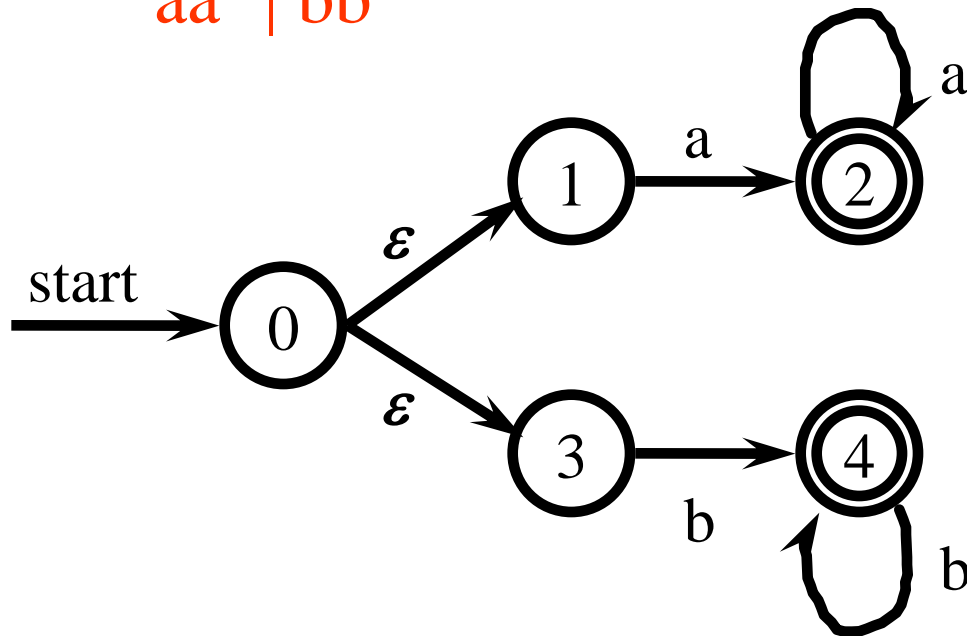
$$(0, \varepsilon) = \{1, 3\}, \quad (1, a) = \{2\}, \quad (2, a) = \{2\} \\ (3, b) = \{4\}, \quad (4, b) = \{4\}$$

♠ Start state: 0

♠ Final states: $\{2, 4\}$

Another Example

$aa^* \mid bb^*$



$aaa: \{0\} \xrightarrow{\epsilon} \{0, 1, 3\} \xrightarrow{a} \{2\} \xrightarrow{\epsilon} \{2\} \xrightarrow{a} \{2\} \xrightarrow{\epsilon} \{2\} \xrightarrow{a} \{2\} \xrightarrow{\epsilon} \{2\}$

Simulating an NFA

Input. An input string ended with **eof** and an NFA with start state s_0 and final states F .

Output. The answer “yes” if accepts, “no” otherwise.

begin

$S := \varepsilon\text{-closure}(\{s_0\});$

$c := \text{nextchar};$

while $c \neq \text{eof}$ **do begin**

$S := \varepsilon\text{-closure}(\text{move}(S, c));$

$c := \text{nextchar}$

end;

if $S \cap F \neq \emptyset$ **then return** “yes”

else return “no”

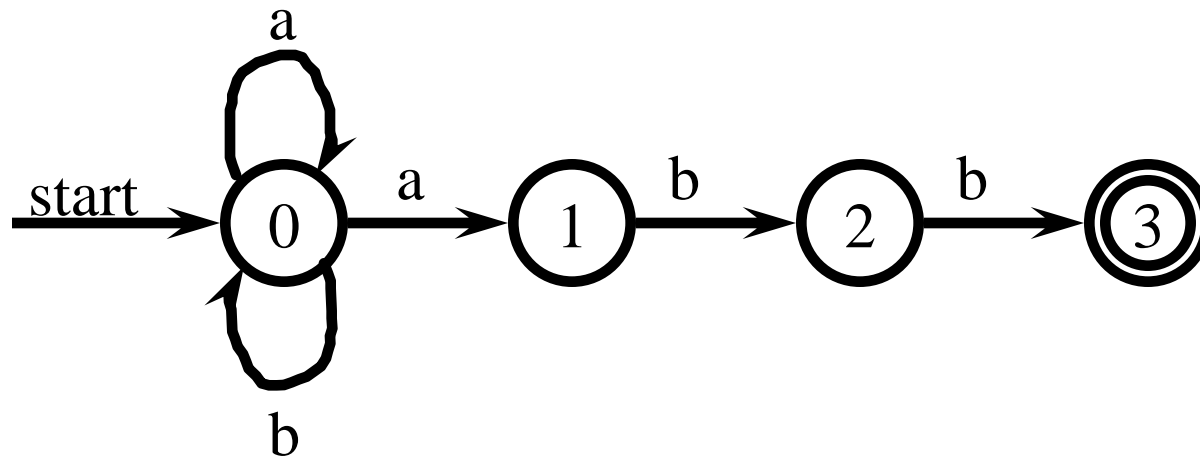
end.

Operations on NFA states

- ♠ $move(s, c)$: set of NFA states reachable from NFA state s on input symbol c
- ♠ $move(S, c)$: set of NFA states reachable from some NFA state s in S on input symbol c
- ♠ $\varepsilon\text{-closure}(s)$: set of NFA states reachable from NFA state s on $\varepsilon\text{-transitions}$ alone
- ♠ $\varepsilon\text{-closure}(S)$: set of NFA states reachable from some NFA state s in S on $\varepsilon\text{-transitions}$ alone

Transition Diagram

$(a \mid b)^*abb$



An Example

$(a \mid b)^*abb$

bbababb

$S = \{0\}$
 $S = \text{move}(\{0\}, b) = \{0\}$
 $S = \text{move}(\{0\}, b) = \{0\}$
 $S = \text{move}(\{0\}, a) = \{0, 1\}$
 $S = \text{move}(\{0, 1\}, b) = \{0, 2\}$
 $S = \text{move}(\{0, 2\}, a) = \{0, 1\}$
 $S = \text{move}(\{0, 1\}, b) = \{0, 2\}$
 $S = \text{move}(\{0, 2\}, b) = \{0, 3\}$
 $S \cap \{3\} \subsetneq \emptyset$

bbabab

$S = \{0\}$
 $S = \text{move}(\{0\}, b) = \{0\}$
 $S = \text{move}(\{0\}, b) = \{0\}$
 $S = \text{move}(\{0\}, a) = \{0, 1\}$
 $S = \text{move}(\{0, 1\}, b) = \{0, 2\}$
 $S = \text{move}(\{0, 2\}, a) = \{0, 1\}$
 $S = \text{move}(\{0, 1\}, b) = \{0, 2\}$
 $S \cap \{3\} = \emptyset$

Computation of ε -closure

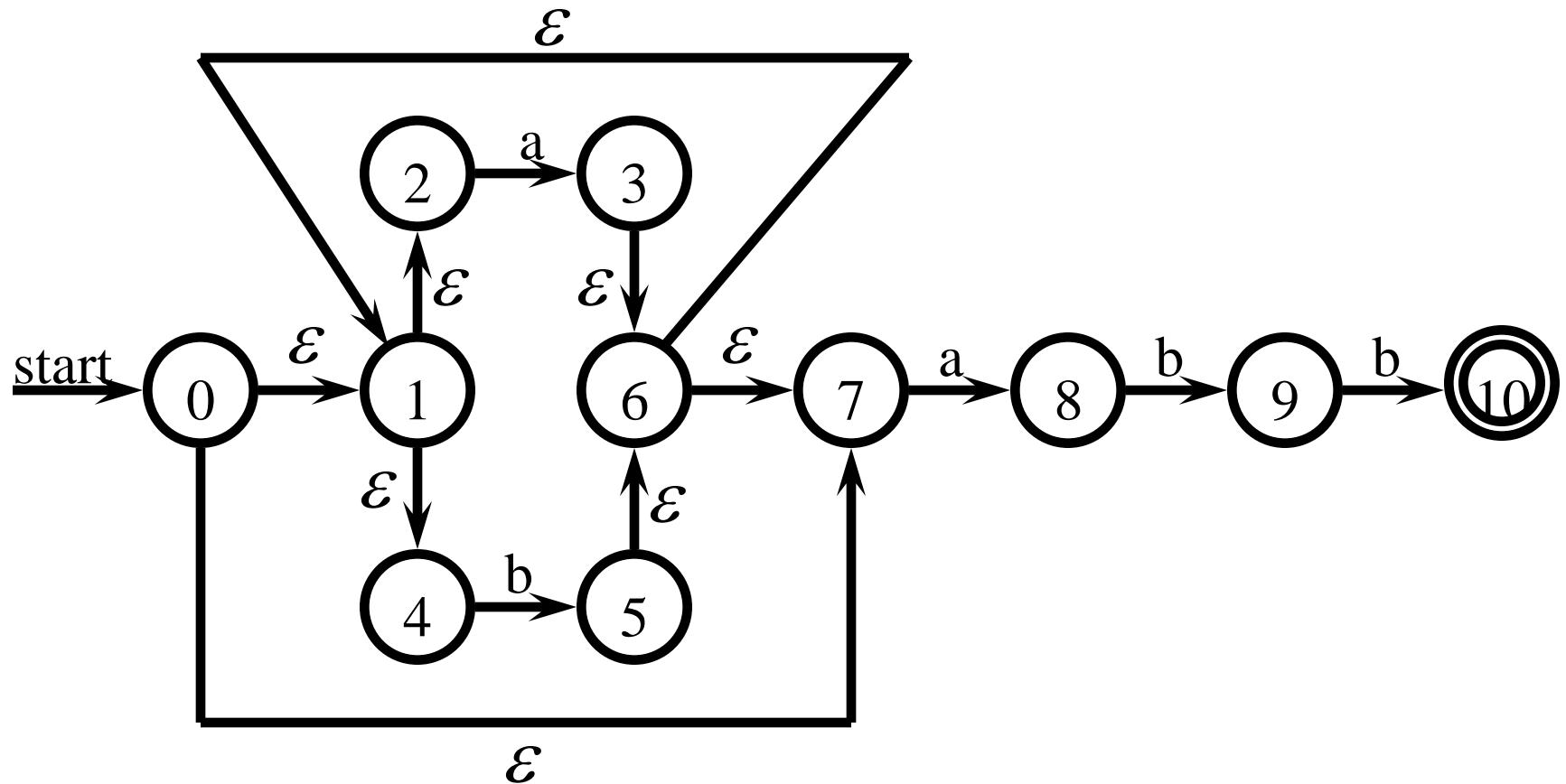
Input. An NFA and a set of NFA states S .

Output. $T = \varepsilon\text{-closure}(S)$.

```
begin  /* A DFT along the  $\varepsilon$ -transitions */  
  push all states in  $S$  onto stack;  $T := S$ ;  
  while stack is not empty do begin  
    pop  $t$ , the top element, off of stack;  
    for each state  $u$  with an edge from  $t$  to  $u$  labeled  $\varepsilon$  do  
      if  $u$  is not in  $T$  do begin  
        add  $u$  to  $T$ ; push  $u$  onto stack  
      end  
    end;  
  return  $T$   
end.
```

An Example

$(a \mid b)^*abb$



An Example

bbabb

$$S = \varepsilon\text{-closure}(\{0\}) = \{0,1,2,4,7\}$$

$$\begin{aligned} S &= \varepsilon\text{-closure}(\text{move}(\{0,1,2,4,7\}, b)) \\ &= \varepsilon\text{-closure}(\{5\}) = \{1,2,4,5,6,7\} \end{aligned}$$

$$\begin{aligned} S &= \varepsilon\text{-closure}(\text{move}(\{1,2,4,5,6,7\}, b)) \\ &= \varepsilon\text{-closure}(\{5\}) = \{1,2,4,5,6,7\} \end{aligned}$$

$$\begin{aligned} S &= \varepsilon\text{-closure}(\text{move}(\{1,2,4,5,6,7\}, a)) \\ &= \varepsilon\text{-closure}(\{3,8\}) = \{1,2,3,4,6,7,8\} \end{aligned}$$

$$\begin{aligned} S &= \varepsilon\text{-closure}(\text{move}(\{1,2,3,4,6,7,8\}, b)) \\ &= \varepsilon\text{-closure}(\{5,9\}) = \{1,2,4,5,6,7,9\} \end{aligned}$$

$$\begin{aligned} S &= \varepsilon\text{-closure}(\text{move}(\{1,2,4,5,6,7,9\}, b)) \\ &= \varepsilon\text{-closure}(\{5,10\}) = \{1,2,4,5,6,7,10\} \end{aligned}$$

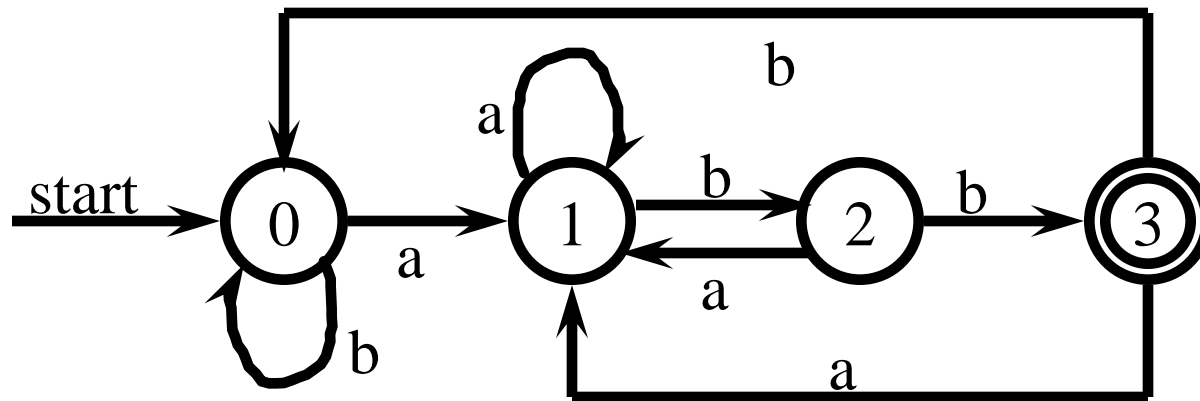
$$S \cap \{10\} \triangleleft \emptyset$$

Deterministic Finite Automata

- ♠ A DFA is a special case of an NFA in which
 - *no* state has an ϵ -transition
 - for each state s and input symbol a , there is *at most one* edge labeled a leaving s

Transition Diagram

$(a \mid b)^*abb$



An Example

♠ RE: $(a \mid b)^*abb$

♠ States: $\{0, 1, 2, 3\}$

♠ Input symbols: $\{a, b\}$

♠ Transition function:

$(0,a) = 1, (1,a) = 1, (2,a) = 1, (3,a) = 1$

$(0,b) = 0, (1,b) = 2, (2,b) = 3, (3,b) = 0$

♠ Start state: 0

♠ Final states: $\{3\}$

Simulating a DFA

Input. An input string ended with **eof** and a DFA with start state s_0 and final states F .

Output. The answer “yes” if accepts, “no” otherwise.

begin

$s := s_0;$

$c := nextchar;$

while $c \neq eof$ **do begin**

$s := move(s, c);$

$c := nextchar$

end;

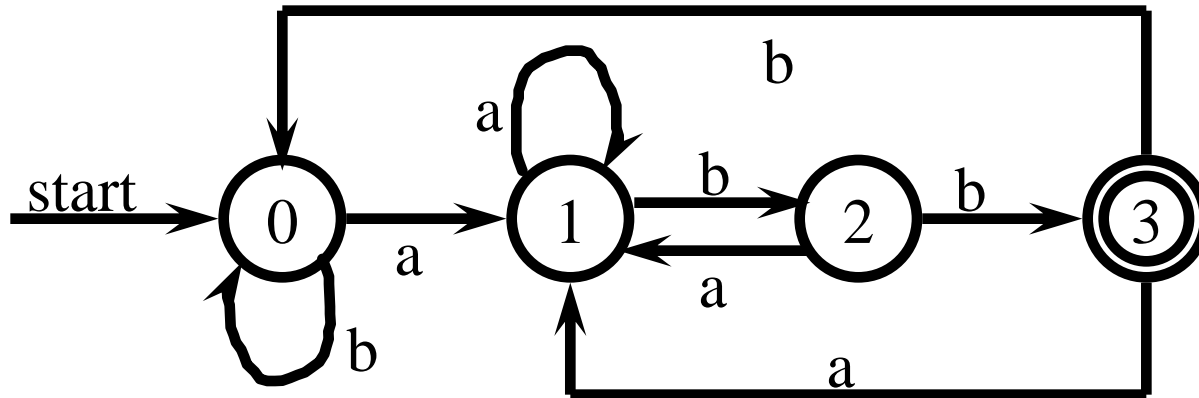
if s is in F **then return** “yes”

else return “no”

end.

An Example

$(a \mid b)^*abb$



$abb: 0 \xrightarrow{a} 1 \xrightarrow{b} 2 \xrightarrow{b} 3$

$aabb: 0 \xrightarrow{a} 1 \xrightarrow{a} 1 \xrightarrow{b} 2 \xrightarrow{b} 3$

共勉

子貢曰：貧而無諂，富而無驕，何如。

子曰：可也，未若貧而樂，富而好禮者也。

子貢曰：詩云：「如切如磋，如琢如磨。」

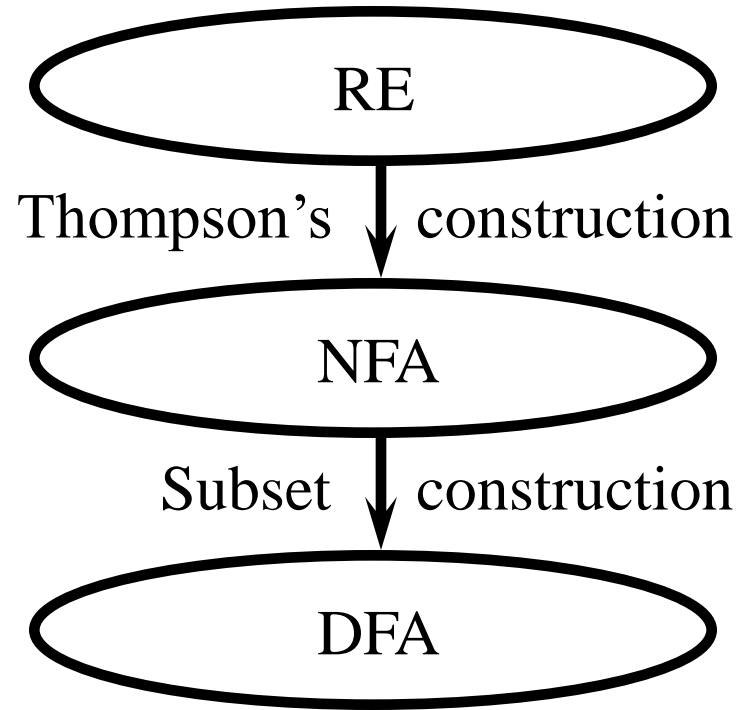
其斯之謂與。

子曰：賜也，始可與言詩已矣；

告諸往而知來者。

-- 論語

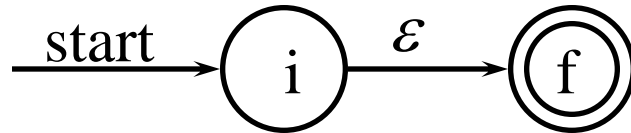
Lexical Analyzer Generator



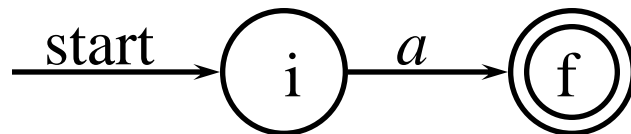
From a RE to an NFA

♠ Thompson's construction algorithm

- For ϵ , construct

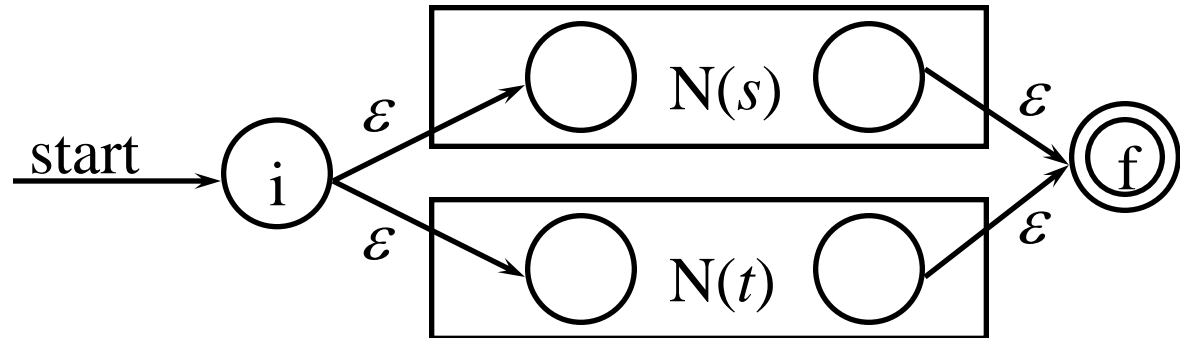


- For a in alphabet, construct

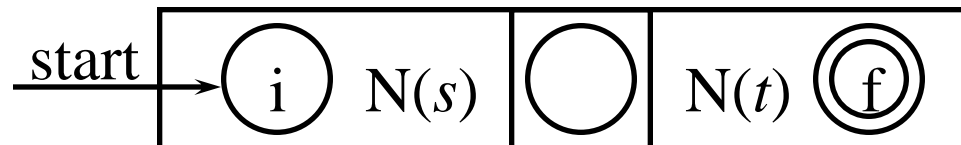


From a RE to an NFA

- Suppose $N(s)$ and $N(t)$ are NFA for RE s and t
 - for $s \mid t$, construct

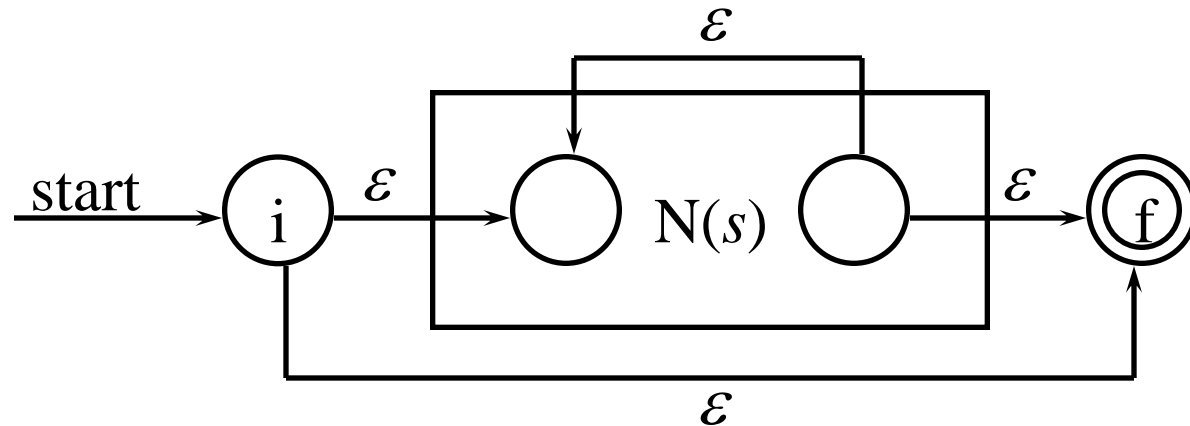


- for st , construct



From a RE to an NFA

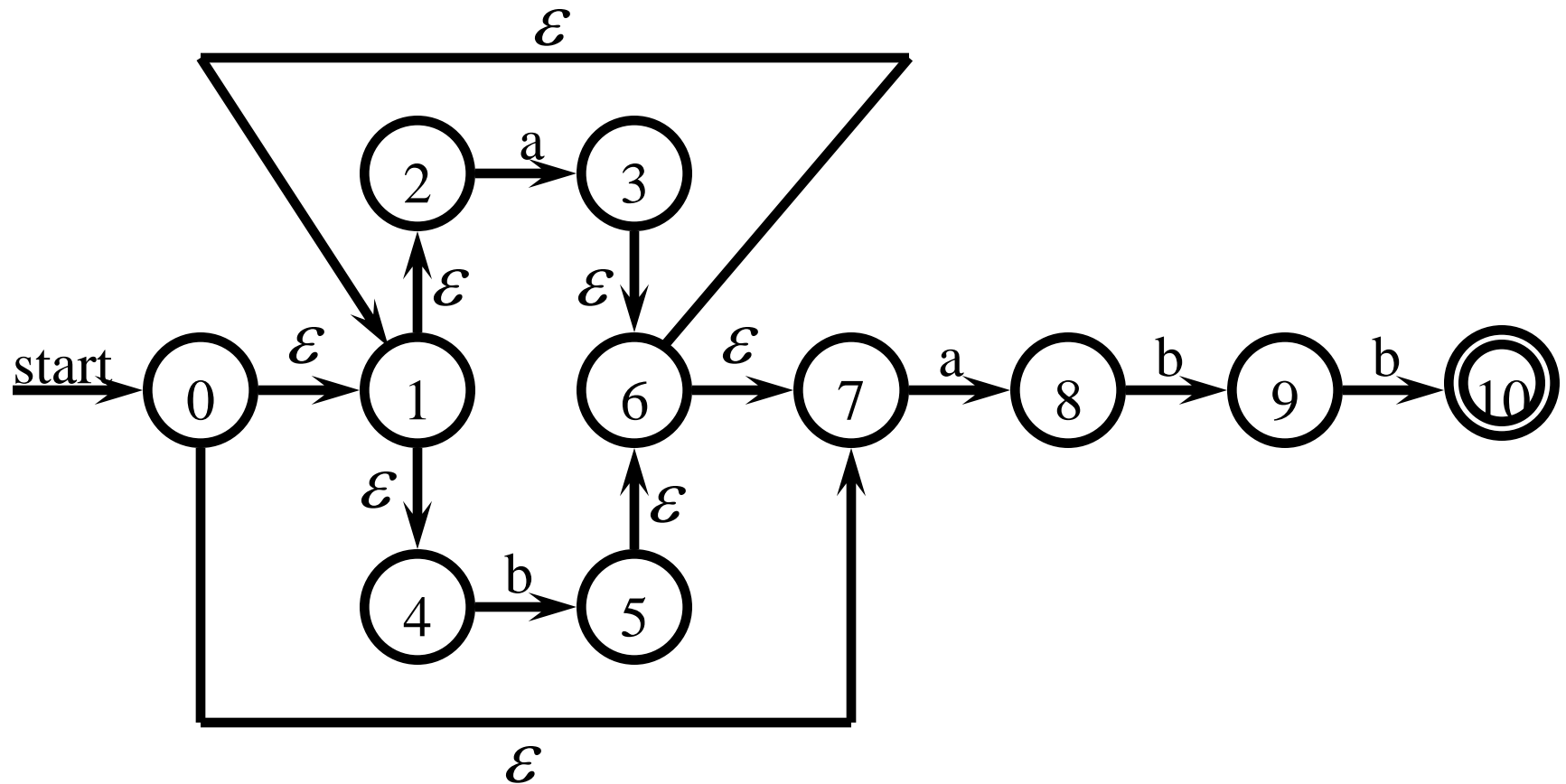
- for s^* , construct



- for (s) , use $N(s)$

An Example

$(a \mid b)^*abb$



From an NFA to a DFA

a set of NFA states \equiv a DFA state

- Find the initial state of the DFA
- Find all the states in the DFA
- Construct the transition table
- Find the final states of the DFA

Subset Construction Algorithm

Input. An NFA N .

Output. A DFA D with states $Dstates$ and transition table $Dtran$.

begin

add ε -closure(s_0) as an unmarked state to $Dstates$;

while there is an unmarked state T in $Dstates$ **do begin**

mark T ;

for each input symbol a **do begin**

$U := \varepsilon$ -closure(move(T, a));

if U is not in $Dstates$ **then**

add U as an unmarked state to $Dstates$;

$Dtran[T, a] := U$

end

end.

An Example

$$\varepsilon\text{-closure}(\{0\}) = \{0,1,2,4,7\} = \textcolor{red}{A}$$

$$\varepsilon\text{-closure}(\text{move}(A, a)) = \varepsilon\text{-closure}(\{3,8\}) = \{1,2,3,4,6,7,8\} = \textcolor{red}{B}$$

$$\varepsilon\text{-closure}(\text{move}(A, b)) = \varepsilon\text{-closure}(\{5\}) = \{1,2,4,5,6,7\} = \textcolor{red}{C}$$

$$\varepsilon\text{-closure}(\text{move}(B, a)) = \varepsilon\text{-closure}(\{3,8\}) = B$$

$$\varepsilon\text{-closure}(\text{move}(B, b)) = \varepsilon\text{-closure}(\{5,9\}) = \{1,2,4,5,6,7,9\} = \textcolor{red}{D}$$

$$\varepsilon\text{-closure}(\text{move}(C, a)) = \varepsilon\text{-closure}(\{3,8\}) = B$$

$$\varepsilon\text{-closure}(\text{move}(C, b)) = \varepsilon\text{-closure}(\{5\}) = C$$

$$\varepsilon\text{-closure}(\text{move}(D, a)) = \varepsilon\text{-closure}(\{3,8\}) = B$$

$$\varepsilon\text{-closure}(\text{move}(D, b)) = \varepsilon\text{-closure}(\{5,10\}) = \{1,2,4,5,6,7,10\} = \textcolor{red}{E}$$

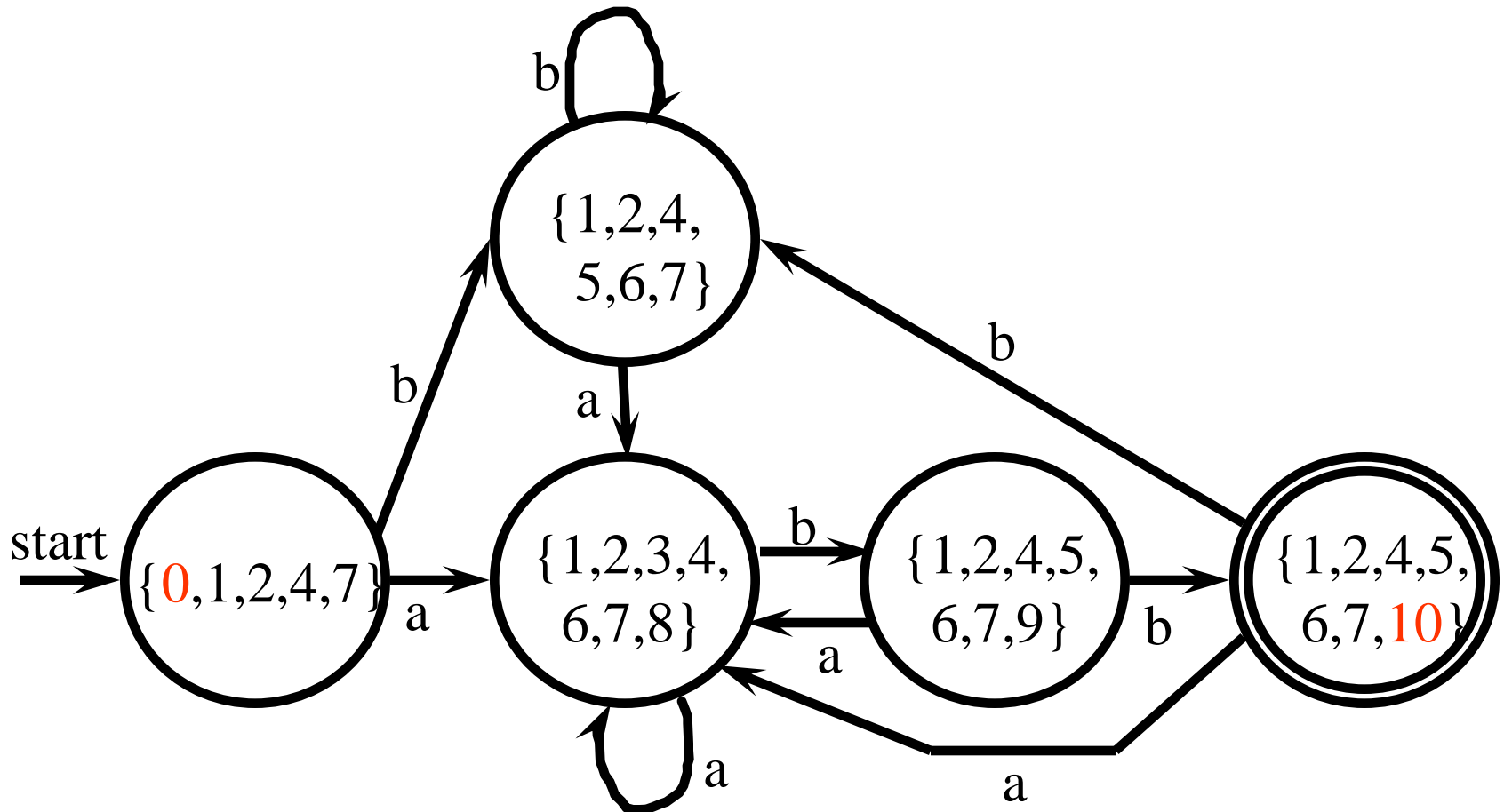
$$\varepsilon\text{-closure}(\text{move}(E, a)) = \varepsilon\text{-closure}(\{3,8\}) = B$$

$$\varepsilon\text{-closure}(\text{move}(E, b)) = \varepsilon\text{-closure}(\{5\}) = C$$

An Example

State	Input Symbol	
	<i>a</i>	<i>b</i>
$A = \{0,1,2,4,7\}$	B	C
$B = \{1,2,3,4,6,7,8\}$	B	D
$C = \{1,2,4,5,6,7\}$	B	C
$D = \{1,2,4,5,6,7,9\}$	B	E
$E = \{1,2,4,5,6,7,10\}$	B	C

An Example

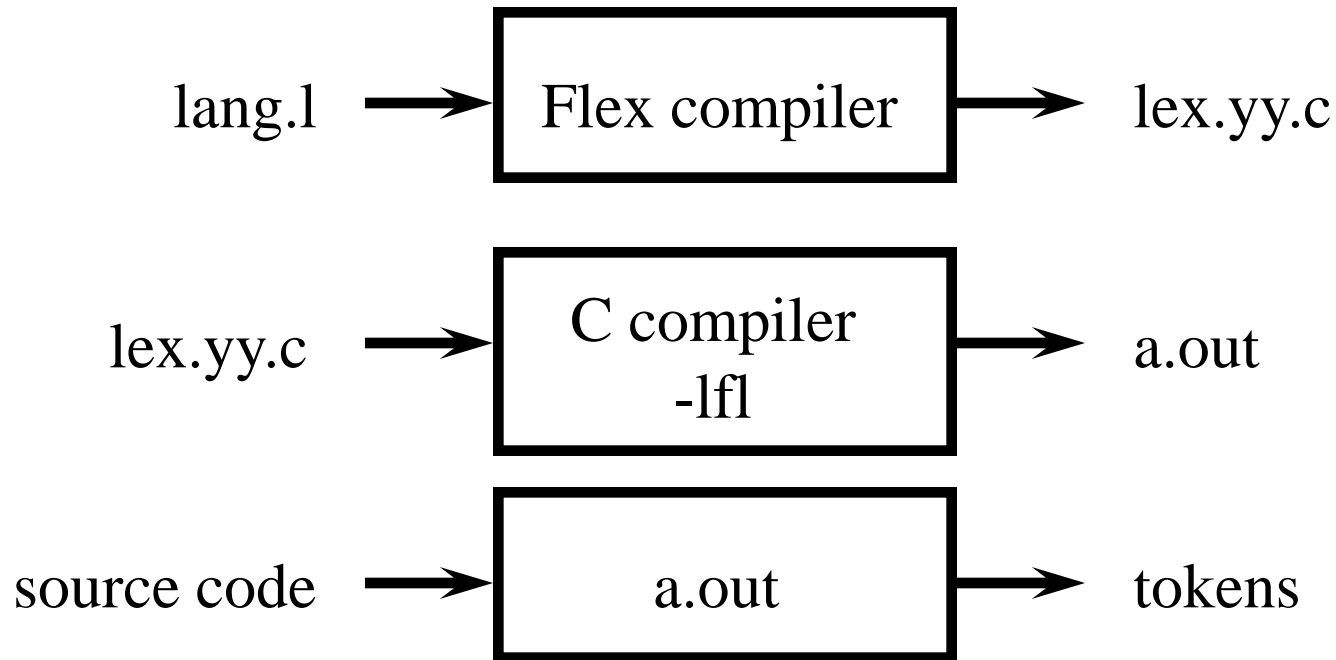


Time-Space Tradeoffs

- ♠ RE to NFA, **simulate NFA**
 - time: $O(|r| * |x|)$, space: $O(|r|)$
- ♠ RE to NFA, NFA to DFA, **simulate DFA**
 - time: $O(|x|)$, space: $O(2^{|r|})$
- ♠ Lazy transition evaluation
 - transitions are computed as needed at run time; computed transitions are stored in cache for later use

Flex – Lexical Analyzer Generator

A language for specifying lexical analyzers



Flex Programs

% {

auxiliary declarations

% }

regular definitions

% %

translation rules

% %

auxiliary procedures

Translation Rules

P_1	action_1
P_2	action_2
\dots	
P_n	action_n

where P_i are **regular expressions** and
 action_i are **C program segments**

An Example

%%

```
username  printf( “%s”, getlogin() );
```

By default, any text not matched by a flex lexical analyzer is copied to the output. This lexical analyzer copies its input file to its output with each occurrence of “username” being replaced with the user’s login name.

An Example

```
% {  
    int num_lines = 0, num_chars = 0;  
% }  
%%  
\n    ++num_lines; ++num_chars;  
.    ++num_chars;    /* all characters except \n */  
%%  
main() {  
    yylex();  
    printf("lines = %d, chars = %d\n",  
        num_lines, num_chars);  
}
```

An Example

```
% {  
#define EOF      0  
#define LE       25  
#define EQ       26  
...  
% }  
delim      [ \t\n]  
ws         { delim }+  
letter     [A-Za-z]  
digit      [0-9]  
id         { letter } ( { letter } | { digit } ) *  
number     { digit } + ( \. { digit } + ) ? ( E [ + \ - ] ? { digit } + ) ?  
% %
```


An Example

```
{ws}      { /* no action and no return */ }
if        {return (IF);}
else      {return (ELSE);}
{id}      {yylval=install_id(); return (ID);}
{number}  {yylval=install_num(); return (NUMBER);}
"<="      {yylval=LE; return (RELOP);}
"=="      {yylval=EQ; return (RELOP);}

...
<<EOF>>   {return(EOF);}
%%
install_id() { ... }
install_num() { ... }
```

Functions and Variables

yylex()

a function implementing the lexical analyzer and returning the token matched

yytext

a global pointer variable pointing to the lexeme matched

yylen

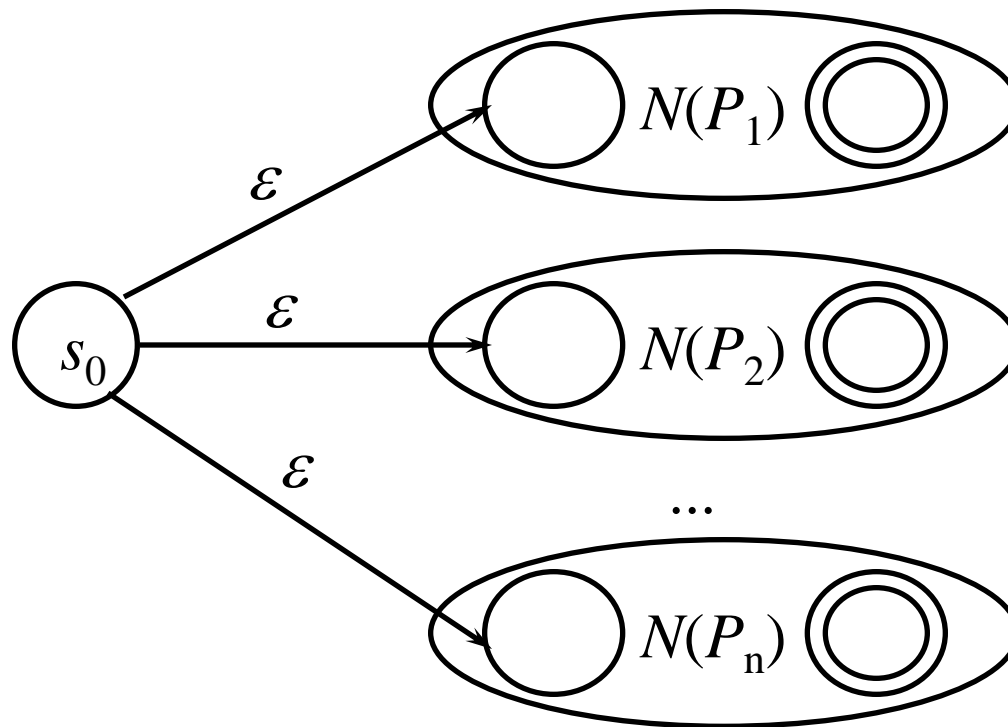
a global variable giving the length of the lexeme matched

yyval

an external global variable storing the attribute of the token

NFA from Flex Programs

$$P_1 \mid P_2 \mid \dots \mid P_n$$



Rules

- ♠ Look for the longest lexeme
 - number
- ♠ Look for the first-listed pattern that matches the longest lexeme
 - keywords and identifiers
- ♠ List frequently occurring patterns first
 - white space

Rules

- ♠ View keywords as exceptions to the rule of identifiers
 - construct a keyword table
- ♠ Lookahead operator: r_1/r_2 - match a string in r_1 only if followed by a string in r_2
 - DO 5 I = 1. 25
DO 5 I = 1, 25
DO/({letter}||{digit})* = ({letter}||{digit})*,

Rules

- **Start condition:** $\langle s \rangle r$ – match r only in start condition s
 $\langle \text{str} \rangle [^"]^* \quad \{ /* \text{eat up string body} */ \}$
- Start conditions are declared in the first section using either **%s** or **%x**
%s str
- A start condition is activated using the **BEGIN** action
 $\backslash'' \quad \text{BEGIN}(\text{str});$
- The default start condition is **INITIAL**

Lexical Error Recovery

- ♠ Error: none of patterns matches a prefix of the remaining input
- ♠ Panic mode error recovery
 - delete successive characters from the remaining input until the pattern-matching can continue
- ♠ Error repair:
 - delete an extraneous character
 - insert a missing character
 - replace an incorrect character
 - transpose two adjacent characters

Maintaining Line Number

- Flex allows to maintain the number of the current line in the global variable **yylineno** using the following option mechanism

%option yylineno

in the first section

共勉

子曰：學而不思則罔，思而不學則殆。

子曰：溫故而知新，可以為師矣。

-- 論語