

Lexical Analysis and Lexical Analyzer Generators

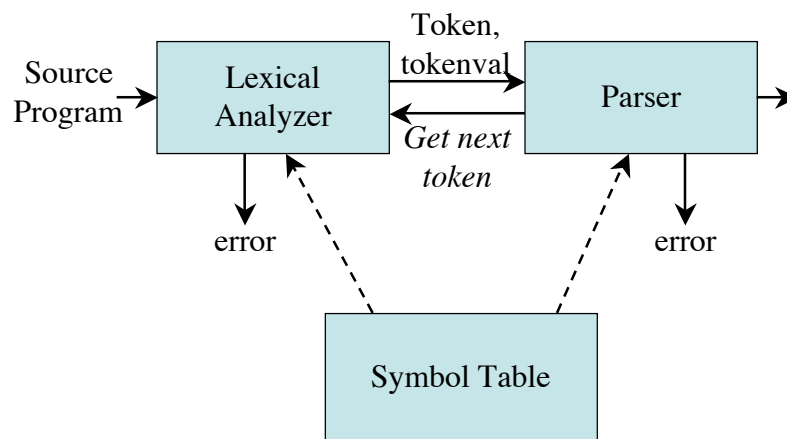
Chapter 3

COP5621 Compiler Construction
Copyright Robert van Engelen, Florida State University, 2005

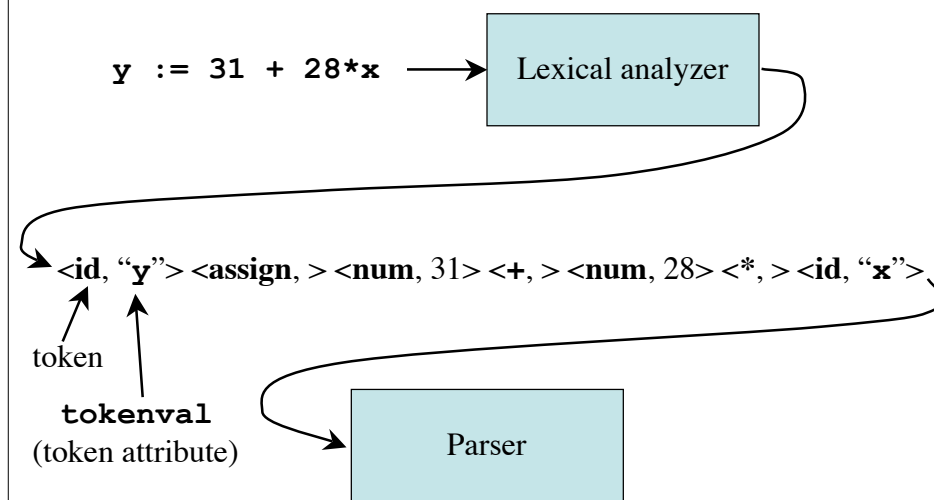
The Reason Why Lexical Analysis is a Separate Phase

- Simplifies the design of the compiler
 - LL(1) or LR(1) with 1 lookahead would not be possible
- Provides efficient implementation
 - Systematic techniques to implement lexical analyzers by hand or automatically
 - Stream buffering methods to scan input
- Improves portability
 - Non-standard symbols and alternate character encodings can be more easily translated

Interaction of the Lexical Analyzer with the Parser



Attributes of Tokens



Tokens, Patterns, and Lexemes

- A *token* is a classification of lexical units
 - For example: **id** and **num**
- *Lexemes* are the specific character strings that make up a token
 - For example: **abc** and **123**
- *Patterns* are rules describing the set of lexemes belonging to a token
 - For example: “*letter followed by letters and digits*” and “*non-empty sequence of digits*”

Specification of Patterns for Tokens: Terminology

- An *alphabet* Σ is a finite set of symbols (characters)
- A *string* s is a finite sequence of symbols from Σ
 - $|s|$ denotes the length of string s
 - ϵ denotes the empty string, thus $|\epsilon| = 0$
- A *language* is a specific set of strings over some fixed alphabet Σ

Specification of Patterns for Tokens: String Operations

- The *concatenation* of two strings x and y is denoted by xy
- The *exponentiation* of a string s is defined by

$$s^0 = \epsilon$$

$$s^i = s^{i-1}s \text{ for } i > 0$$

(note that $s\epsilon = \epsilon s = s$)

Specification of Patterns for Tokens: Language Operations

- *Union*

$$L \cup M = \{s \mid s \in L \text{ or } s \in M\}$$
- *Concatenation*

$$LM = \{xy \mid x \in L \text{ and } y \in M\}$$
- *Exponentiation*

$$L^0 = \{\epsilon\}; L^i = L^{i-1}L$$
- *Kleene closure*

$$L^* = \bigcup_{i=0, \dots, \infty} L^i$$
- *Positive closure*

$$L^+ = \bigcup_{i=1, \dots, \infty} L^i$$

Specification of Patterns for Tokens: Regular Expressions

- Basis symbols:
 - ϵ is a regular expression denoting language $\{\epsilon\}$
 - $a \in \Sigma$ is a regular expression denoting $\{a\}$
- If r and s are regular expressions denoting languages $L(r)$ and $M(s)$ respectively, then
 - $r \mid s$ is a regular expression denoting $L(r) \cup M(s)$
 - rs is a regular expression denoting $L(r)M(s)$
 - r^* is a regular expression denoting $L(r)^*$
 - (r) is a regular expression denoting $L(r)$
- A language defined by a regular expression is called a *regular set*

Specification of Patterns for Tokens: Regular Definitions

- Naming convention for regular expressions:

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

$$\dots$$

$$d_n \rightarrow r_n$$
 where r_i is a regular expression over $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$
- Each d_j in r_i is textually substituted in r_i

Specification of Patterns for Tokens: Regular Definitions

- Example:

letter \rightarrow **A** | **B** | ... | **Z** | **a** | **b** | ... | **z**

digit \rightarrow **0** | **1** | ... | **9**

id \rightarrow **letter** (**letter** | **digit**)^{*}

- Cannot use recursion, this is illegal:

digits \rightarrow **digit digits** | **digit**

Specification of Patterns for Tokens: Notational Shorthands

- We frequently use the following shorthands:

$$r^+ = rr^*$$

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

$$[\mathbf{a-z}] = \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \dots \mid \mathbf{z}$$

- For example:

digit \rightarrow [0-9]

num \rightarrow **digit**⁺ (. **digit**⁺)[?] (**E** (+ | -)[?] **digit**⁺)[?]

Regular Definitions and Grammars

Grammar

$stmt \rightarrow \mathbf{if\ expr\ then\ stmt}$
 $\mid \mathbf{if\ expr\ then\ stmt\ else\ stmt}$
 $\mid \epsilon$

$expr \rightarrow \mathbf{term\ relop\ term}$
 $\mid \mathbf{term}$

$term \rightarrow \mathbf{id}$
 $\mid \mathbf{num}$

Regular definitions

$\mathbf{if} \rightarrow \mathbf{if}$

$\mathbf{then} \rightarrow \mathbf{then}$

$\mathbf{else} \rightarrow \mathbf{else}$

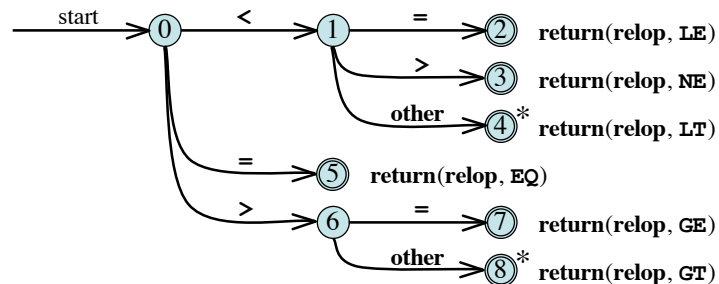
$\mathbf{relop} \rightarrow < \mid <= \mid <> \mid > \mid >= \mid =$

$\mathbf{id} \rightarrow \mathbf{letter\ (letter\ \mid\ digit)^*}$

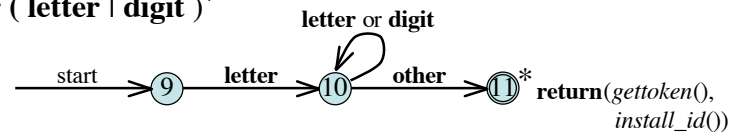
$\mathbf{num} \rightarrow \mathbf{digit^+ \ (. \ digit^+)^? \ (\ E \ (+|-)^? \ digit^+ \)^?}$

Implementing a Scanner Using Transition Diagrams

$\mathbf{relop} \rightarrow < \mid <= \mid <> \mid > \mid >= \mid =$



$\mathbf{id} \rightarrow \mathbf{letter\ (letter\ \mid\ digit)^*}$



Implementing a Scanner Using Transition Diagrams (Code)

15

```
token nexttoken()
{ while (1) {
    switch (state) {
        case 0: c = nextchar();
            if (c==blank || c==tab || c==newline) {
                state = 0;
                lexeme_beginning++;
            }
            else if (c=='<') state = 1;
            else if (c=='=') state = 5;
            else if (c=='>') state = 6;
            else state = fail();
            break;
        case 1:
            ...
        case 9: c = nextchar();
            if (isletter(c)) state = 10;
            else state = fail();
            break;
        case 10: c = nextchar();
            if (isletter(c)) state = 10;
            else if (isdigit(c)) state = 10;
            else state = 11;
            break;
        ...
    }
}
```

Decides what
other start state
is applicable



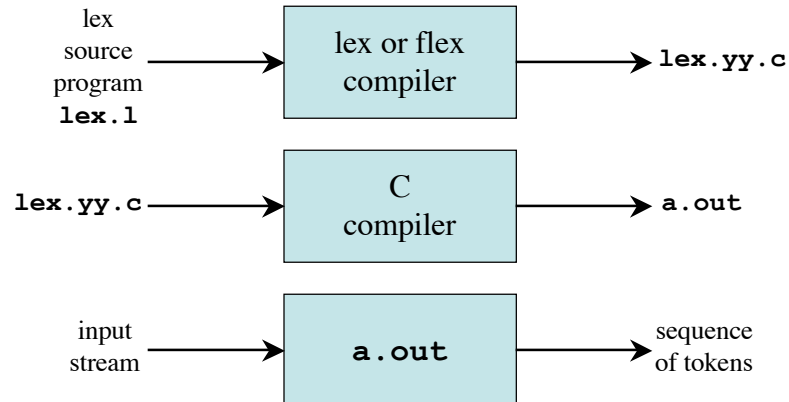
```
int fail()
{ forward = token_beginning;
  switch (start) {
      case 0: start = 9; break;
      case 9: start = 12; break;
      case 12: start = 20; break;
      case 20: start = 25; break;
      case 25: recover(); break;
      default: /* error */
  }
  return start;
}
```

The Lex and Flex Scanner Generators

16

- *Lex* and its newer cousin *flex* are scanner generators
- Systematically translate regular definitions into C source code for efficient scanning
- Generated code is easy to integrate in C applications

Creating a Lexical Analyzer with Lex and Flex



Lex Specification

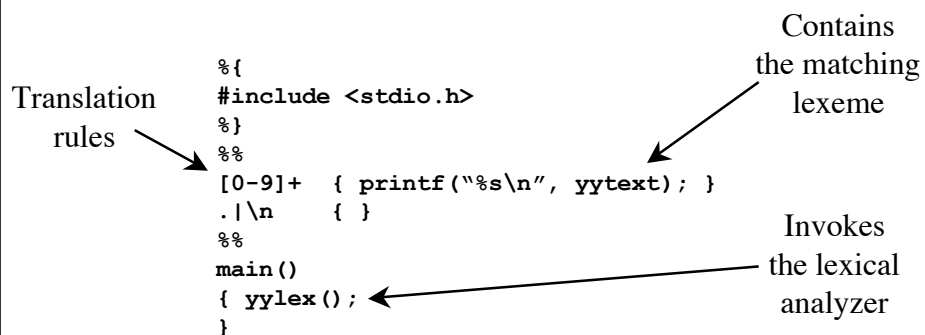
- A *lex specification* consists of three parts:
regular definitions, C declarations in % { % }
% %
translation rules
% %
user-defined auxiliary procedures
- The *translation rules* are of the form:

p_1	$\{ action_1 \}$
p_2	$\{ action_2 \}$
...	
p_n	$\{ action_n \}$

Regular Expressions in Lex

x	match the character x
\.	match the character .
"string"	match contents of string of characters
.	match any character except newline
^	match beginning of a line
\$	match the end of a line
[xyz]	match one character x , y , or z (use \ to escape -)
[^xyz]	match any character except x , y , and z
[a-z]	match one of a to z
r*	closure (match zero or more occurrences)
r+	positive closure (match one or more occurrences)
r?	optional (match zero or one occurrence)
r₁r₂	match r₁ then r₂ (concatenation)
r₁ r₂	match r₁ or r₂ (union)
(r)	grouping
r₁\r₂	match r₁ when followed by r₂
{d}	match the regular expression defined by d

Example Lex Specification 1



```

lex spec.1
gcc lex.yy.c -ll
./a.out < spec.1

```

Example Lex Specification 2

Translation rules →

```

%{
#include <stdio.h>
int ch = 0, wd = 0, nl = 0;
%}
delim      [ \t]+
%%
\n           { ch++; wd++; nl++; }
^{delim}     { ch+=yy leng; }
{delim}      { ch+=yy leng; wd++; }
.            { ch++; }
%%
main()
{ yy lex();
  printf("%8d%8d%8d\n", nl, wd, ch);
}

```

Regular definition →

Example Lex Specification 3

Translation rules →

```

%{
#include <stdio.h>
%}
digit        [0-9]
letter       [A-Za-z]
id           {letter}({letter}|{digit})*
%%
{digit}+     { printf("number: %s\n", yytext); }
{id}         { printf("ident: %s\n", yytext); }
.            { printf("other: %s\n", yytext); }
%%
main()
{ yy lex();
}

```

Regular definitions →

Example Lex Specification 4

```

%{ /* definitions of manifest constants */
#define LT (256)
...
%}
delim      [ \t\n]
ws         {delim}+
letter     [A-Za-z]
digit      [0-9]
id         {letter}({letter}|{digit})*
number     {digit}+(\.{digit})?(E[+-]?{digit})?
%%
{ws}       { }
if         {return IF;}
then       {return THEN;}
else       {return ELSE;}
{id}       {yylval = install_id(); return ID;}
{number}   {yylval = install_num(); return NUMBER;}
"<"       {yylval = LT; return RELOP;}
"<="      {yylval = LE; return RELOP;}
"="        {yylval = EQ; return RELOP;}
">"       {yylval = NE; return RELOP;}
">="      {yylval = GE; return RELOP;}
%%
int install_id()
...

```

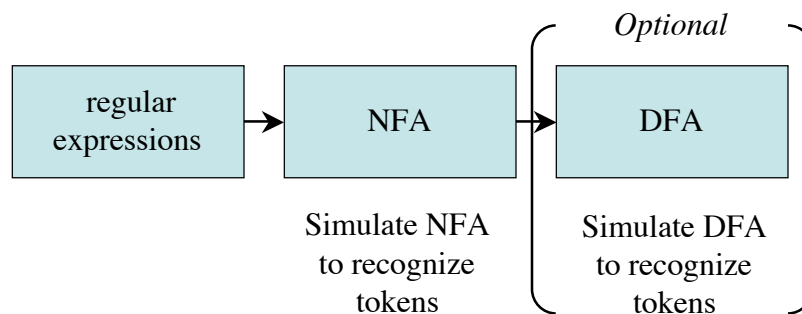
Return token to parser

Token attribute

Install **yytext** as identifier in symbol table

Design of a Lexical Analyzer Generator

- Translate regular expressions to NFA
- Translate NFA to an efficient DFA



Nondeterministic Finite Automata

- Definition: an NFA is a 5-tuple $(S, \Sigma, \delta, s_0, F)$ where

S is a finite set of *states*

Σ is a finite set of *input symbol alphabet*

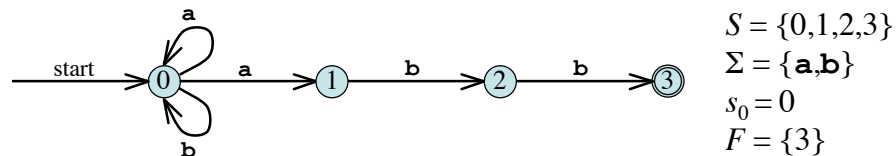
δ is a *mapping* from $S \times \Sigma$ to a set of states

$s_0 \in S$ is the *start state*

$F \subseteq S$ is the set of *accepting (or final) states*

Transition Graph

- An NFA can be diagrammatically represented by a labeled directed graph called a *transition graph*



Transition Table

- The mapping δ of an NFA can be represented in a *transition table*

$$\delta(0, \mathbf{a}) = \{0, 1\}$$

$$\delta(0, \mathbf{b}) = \{0\}$$

$$\delta(1, \mathbf{b}) = \{2\}$$

$$\delta(2, \mathbf{b}) = \{3\}$$

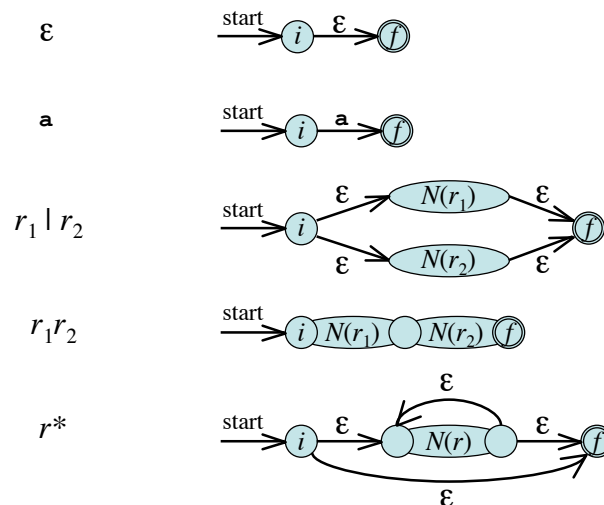
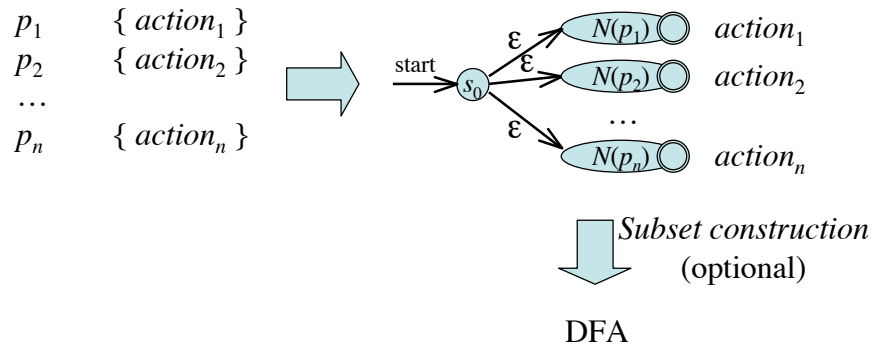


<i>State</i>	<i>Input a</i>	<i>Input b</i>
0	{0, 1}	{0}
1		{2}
2		{3}

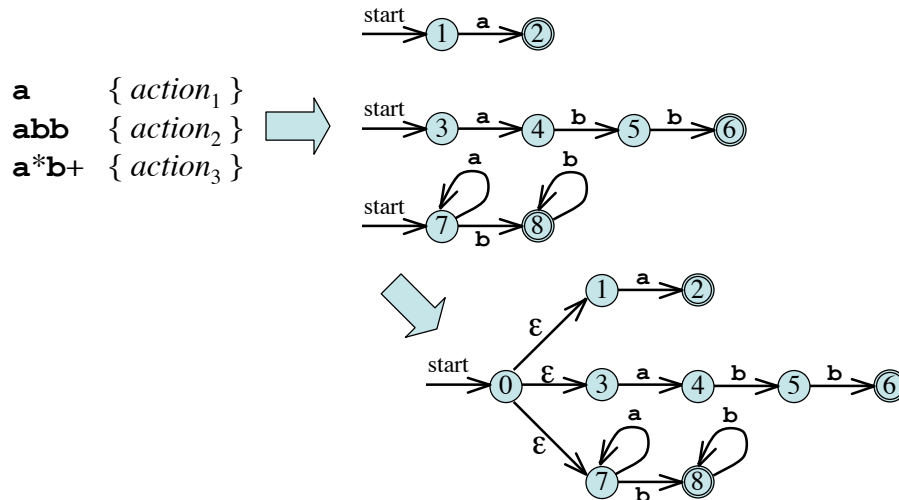
The Language Defined by an NFA

- An NFA *accepts* an input string x **iff** there is some path with edges labeled with symbols from x in sequence from the start state to some accepting state in the transition graph
- A state transition from one state to another on the path is called a *move*
- The *language defined by* an NFA is the set of input strings it accepts, such as $(\mathbf{a|b})^*\mathbf{abb}$ for the example NFA

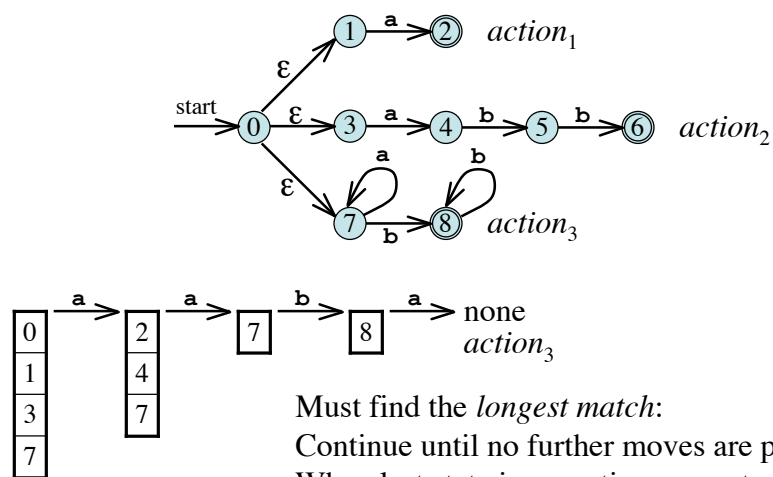
NFA



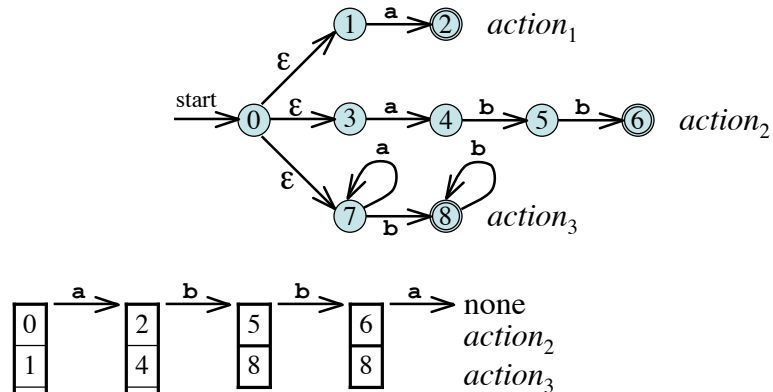
Combining the NFAs of a Set of Regular Expressions



Simulating the Combined NFA Example 1



Simulating the Combined NFA Example 2

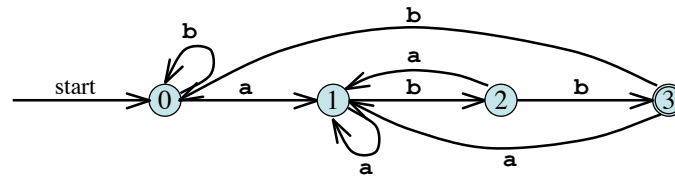


Deterministic Finite Automata

- A *deterministic finite automaton* is a special case of an NFA
 - No state has an ϵ -transition
 - For each state s and input symbol a there is at most one edge labeled a leaving s
- Each entry in the transition table is a single state
 - At most one path exists to accept a string
 - Simulation algorithm is simple

Example DFA

A DFA that accepts $(\mathbf{a|b})^*\mathbf{abb}$



Conversion of an NFA into a DFA

- The *subset construction algorithm* converts an NFA into a DFA using:

$$\varepsilon\text{-closure}(s) = \{s\} \cup \{t \mid s \xrightarrow{\varepsilon} \dots \xrightarrow{\varepsilon} t\}$$

$$\varepsilon\text{-closure}(T) = \bigcup_{s \in T} \varepsilon\text{-closure}(s)$$

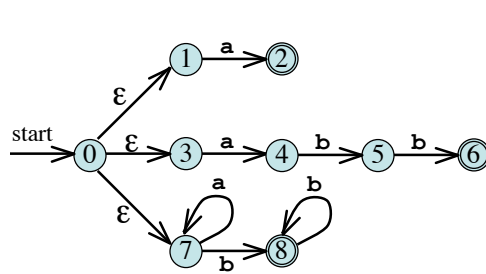
$$\text{move}(T, a) = \{t \mid s \xrightarrow{a} t \text{ and } s \in T\}$$

- The algorithm produces:

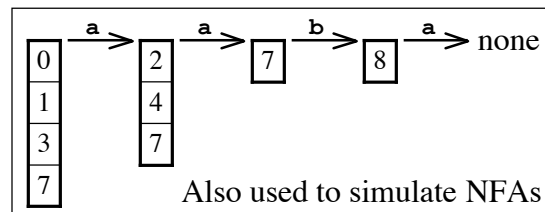
Dstates is the set of states of the new DFA
consisting of sets of states of the NFA

Dtran is the transition table of the new DFA

ϵ -closure and move Examples



$\epsilon\text{-closure}(\{0\}) = \{0, 1, 3, 7\}$
 $\text{move}(\{0, 1, 3, 7\}, \mathbf{a}) = \{2, 4, 7\}$
 $\epsilon\text{-closure}(\{2, 4, 7\}) = \{2, 4, 7\}$
 $\text{move}(\{2, 4, 7\}, \mathbf{a}) = \{7\}$
 $\epsilon\text{-closure}(\{7\}) = \{7\}$
 $\text{move}(\{7\}, \mathbf{b}) = \{8\}$
 $\epsilon\text{-closure}(\{8\}) = \{8\}$
 $\text{move}(\{8\}, \mathbf{a}) = \emptyset$



Simulating an NFA using ϵ -closure and move

```

 $S := \epsilon\text{-closure}(\{s_0\})$ 
 $S_{prev} := \emptyset$ 
 $a := \text{nextchar}()$ 
while  $S \neq \emptyset$  do
     $S_{prev} := S$ 
     $S := \epsilon\text{-closure}(\text{move}(S, a))$ 
     $a := \text{nextchar}()$ 
end do
if  $S_{prev} \cap F \neq \emptyset$  then
    execute action in  $S_{prev}$ 
    return "yes"
else
    return "no"
  
```

The Subset Construction Algorithm

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

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

 mark T

for each input symbol $a \in \Sigma$ **do**

$U := \epsilon\text{-closure}(\text{move}(T, a))$

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

 add U as an unmarked state to $Dstates$

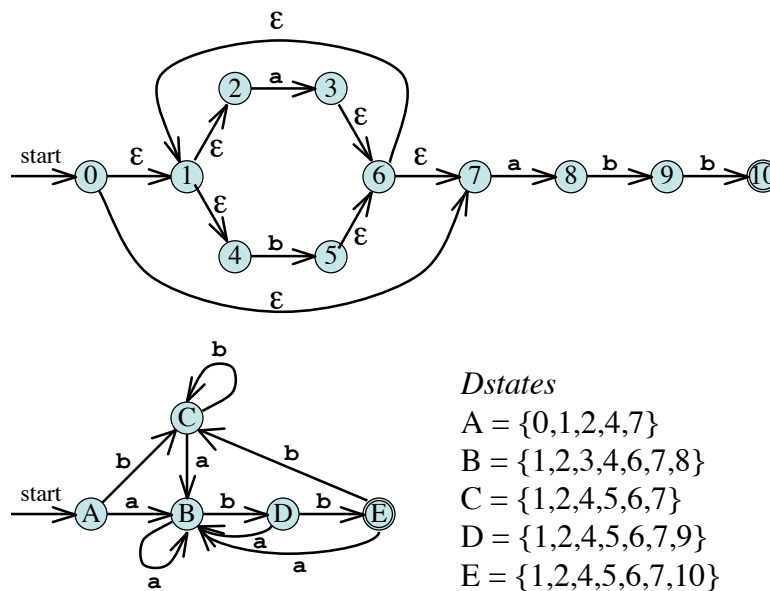
end if

$Dtran[T, a] := U$

end do

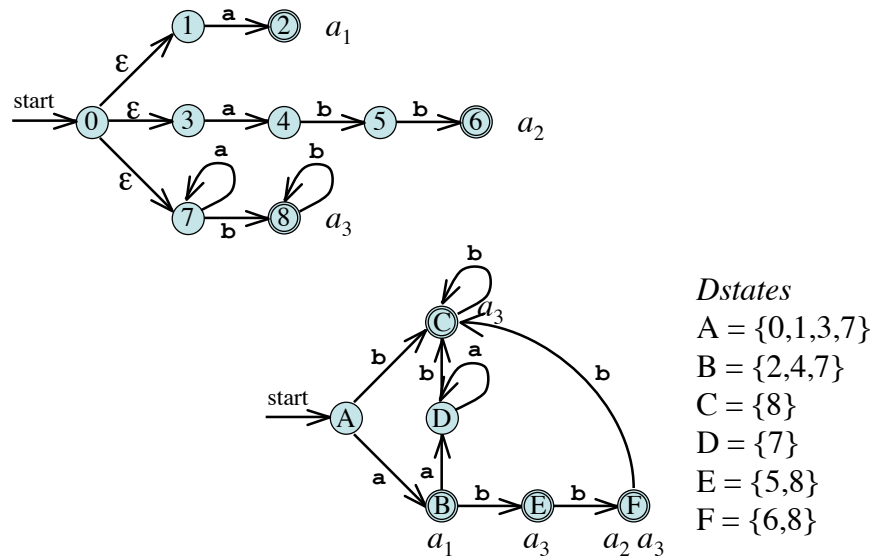
end do

Subset Construction Example 1



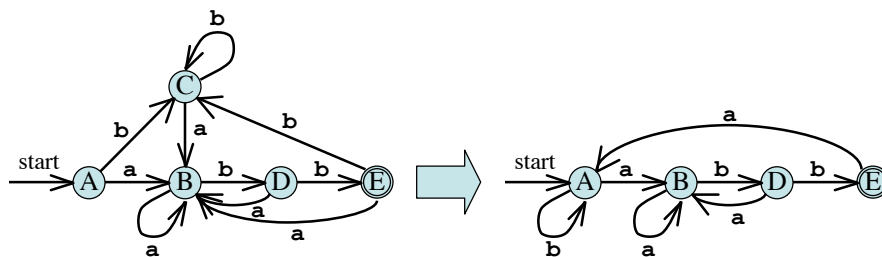
Subset Construction Example 2

41



Minimizing the Number of States of a DFA

42



From Regular Expression to DFA Directly

- The *important states* of an NFA are those without an ϵ -transition, that is if $move(\{s\}, a) \neq \emptyset$ for some a then s is an important state
- The subset construction algorithm uses only the important states when it determines ϵ -closure($move(T, a)$)

From Regular Expression to DFA Directly (Algorithm)

- Augment the regular expression r with a special end symbol $\#$ to make accepting states important: the new expression is $r\#$
- Construct a syntax tree for $r\#$
- Traverse the tree to construct functions *nullable*, *firstpos*, *lastpos*, and *followpos*

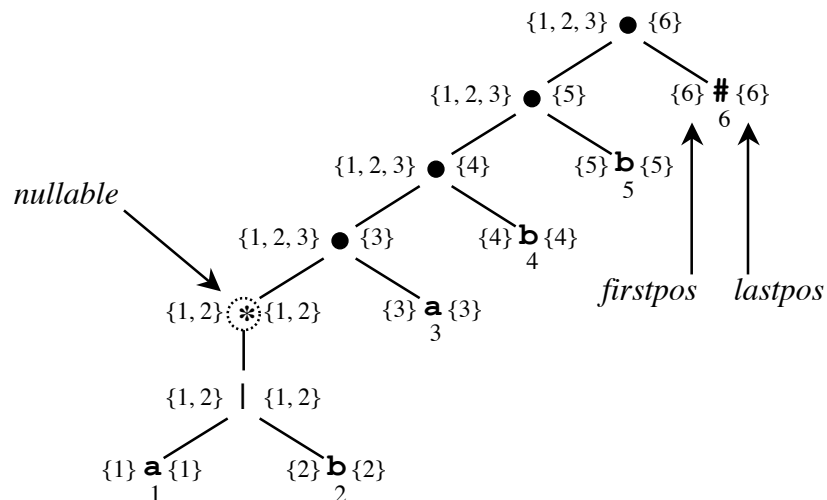


- 23

From Regular Expression to DFA Directly: Annotating the Tree

Node n	$nullable(n)$	$firstpos(n)$	$lastpos(n)$
Leaf ϵ	true	\emptyset	\emptyset
Leaf i	false	$\{i\}$	$\{i\}$
$\begin{array}{c} \\ / \quad \backslash \\ c_1 \quad c_2 \end{array}$	$nullable(c_1)$ or $nullable(c_2)$	$firstpos(c_1) \cup firstpos(c_2)$	$lastpos(c_1) \cup lastpos(c_2)$
$\begin{array}{c} \bullet \\ / \quad \backslash \\ c_1 \quad c_2 \end{array}$	$nullable(c_1)$ and $nullable(c_2)$	if $nullable(c_1)$ then $firstpos(c_1) \cup firstpos(c_2)$ else $firstpos(c_1)$	if $nullable(c_2)$ then $lastpos(c_1) \cup lastpos(c_2)$ else $lastpos(c_2)$
$\begin{array}{c} * \\ \\ c_1 \end{array}$	true	$firstpos(c_1)$	$lastpos(c_1)$

From Regular Expression to DFA Directly: Syntax Tree of $(ab)^*abb\#$



From Regular Expression to DFA Directly: *followpos*

```

for each node  $n$  in the tree do
    if  $n$  is a cat-node with left child  $c_1$  and right child  $c_2$  then
        for each  $i$  in  $lastpos(c_1)$  do
             $followpos(i) := followpos(i) \cup firstpos(c_2)$ 
        end do
    else if  $n$  is a star-node
        for each  $i$  in  $lastpos(n)$  do
             $followpos(i) := followpos(i) \cup firstpos(n)$ 
        end do
    end if
end do

```

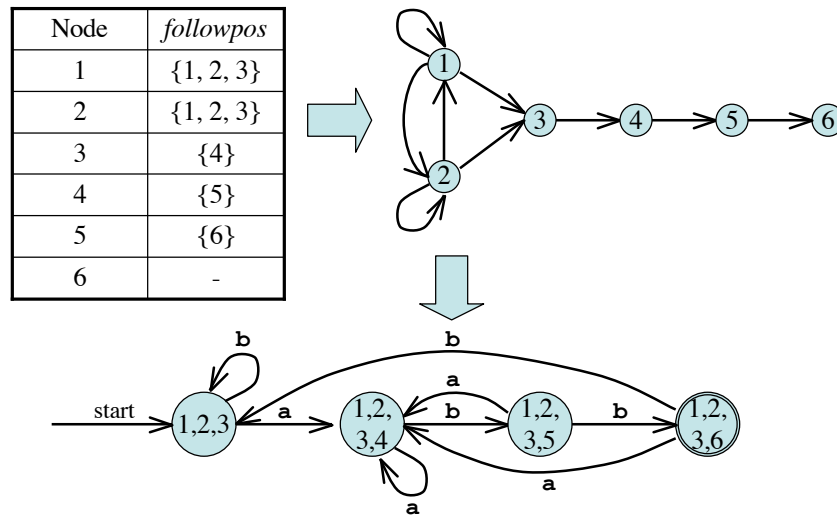
From Regular Expression to DFA Directly: Algorithm

```

 $s_0 := firstpos(root)$  where  $root$  is the root of the syntax tree
 $Dstates := \{s_0\}$  and is unmarked
while there is an unmarked state  $T$  in  $Dstates$  do
    mark  $T$ 
    for each input symbol  $a \in \Sigma$  do
        let  $U$  be the set of positions that are in  $followpos(p)$ 
        for some position  $p$  in  $T$ ,
        such that the symbol at position  $p$  is  $a$ 
        if  $U$  is not empty and not in  $Dstates$  then
            add  $U$  as an unmarked state to  $Dstates$ 
        end if
         $Dtran[T, a] := U$ 
    end do
end do

```

From Regular Expression to DFA Directly: Example



Time-Space Tradeoffs

<i>Automaton</i>	<i>Space (worst case)</i>	<i>Time (worst case)</i>
NFA	$O(r)$	$O(r \times x)$
DFA	$O(2^{ r })$	$O(x)$