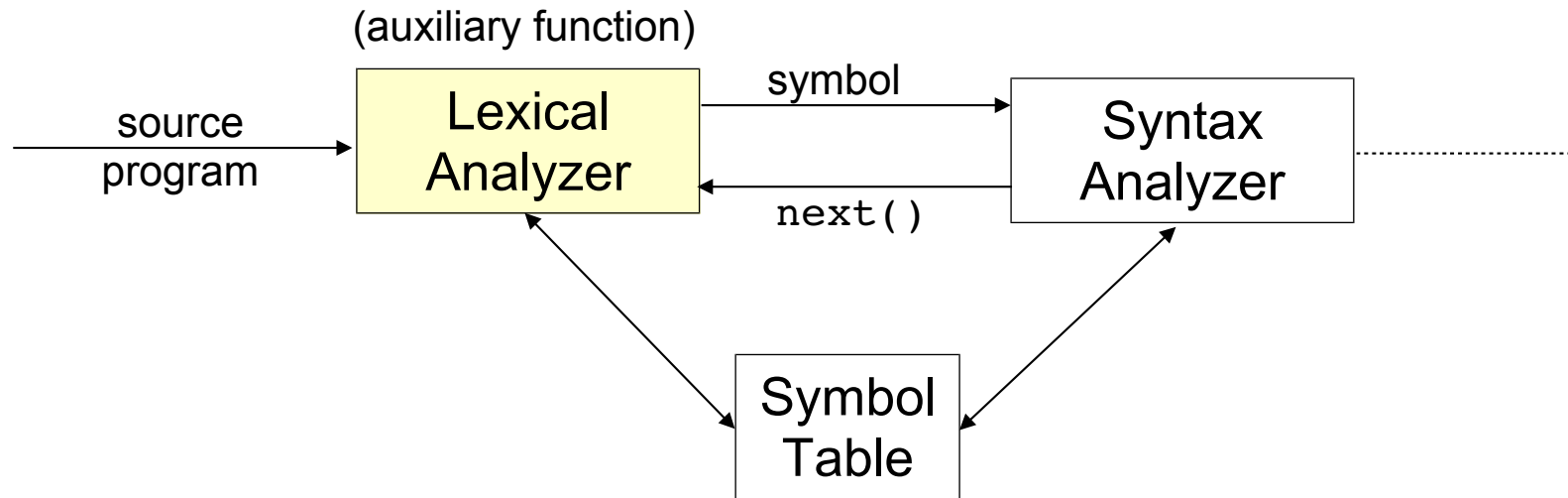


Lexical Analysis

- **Primary** role = abstraction process: $[characters] \rightarrow [symbols]$



- **Secondary** role $\left\{ \begin{array}{l} \text{removal of spacing} \\ \text{removal of comments} \end{array} \right\}$ pseudo-symbols
linking of errors to source lines (supports clarity of error messages)
- Advantages in separating syntax analysis from lexical analysis:
 1. Simplicity of design (syntax do not care about white spaces, comments, ...)
 2. Efficiency
 3. Portability

Symbol Specification

- **Lexical string**: significant sequence of characters (*lexeme*)
- **Symbol**: abstraction of a class of lexical strings
- **Pattern**: rule to describe the extension of a symbol (\approx grammar of lexical symbol)
- **Lexical attribute**: when the symbol is the abstraction of several lexical strings

$\text{alpha} = (\text{beta} * 3)$ \Rightarrow $\langle \text{id}, \uparrow \text{alpha} \rangle \langle \text{assign}, \rangle \langle \text{left}, \rangle \langle \text{id}, \uparrow \text{beta} \rangle \langle \text{times}, \rangle \langle \text{num}, 3 \rangle \langle \text{right}, \rangle$

Regular Expressions

- Definition of regular expression (on alphabet Σ):
 1. ε is a regular expression denoting the language $\{ \varepsilon \}$
 2. If $a \in \Sigma$, then a is a regular expression denoting the language $\{ a \}$
 3. If x, y are regular expressions denoting languages $L(x), L(y)$, then:
 - (x) is a regular expression denoting $L(x)$
 - $(x) \mid (y)$ is a regular expression denoting $L(x) \cup L(y)$
 - $(x)(y)$ is a regular expression denoting $\{ xy \mid x \in L(x), y \in L(y) \}$
 - $(x)^*$ is a regular expression (repetition zero or more times of strings in $L(x)$)
- **Extended regular expressions:** $(x)^+$, \bullet , $[a-z]$, $\sim(a|b)$, $(x)?$
- **Regular definitions:** associations of names with regular expressions

letter $\rightarrow [A-Za-z]$

digit $\rightarrow [0-9]$

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

Regular Expressions for Tokens in Prog. Languages

- Classification of lexical elements:

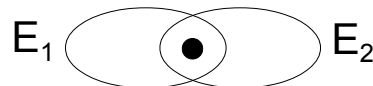
- **Keywords** = { if, then, while, do, begin, end, procedure, function, case, repeat ... }
- **Special symbols** = { +, -, *, /, =, >, >=, !=, ++, --, &&, | |, ... }
- **Identifiers** = { alphanumeric strings starting with a letter }
- **Constants** = { 25, .34, 2.7E-3, "alpha", 'a', ... }
- **Pseudo-symbols** (spacing, comments)

- **Ambiguity**: when the same string of characters is compatible with several regular expressions

Examples: $\left\{ \begin{array}{l} \text{while} \left\{ \begin{array}{l} \text{id} \\ \text{while} \end{array} \right. \\ \text{<>} \left\{ \begin{array}{l} \text{noteq} \\ \text{less greater} \end{array} \right. \end{array} \right.$

PL \rightarrow disambiguating rules $\left\{ \begin{array}{l} \text{keywords = reserved words} \\ \text{principle of longest substring (maximal munch)} \end{array} \right.$

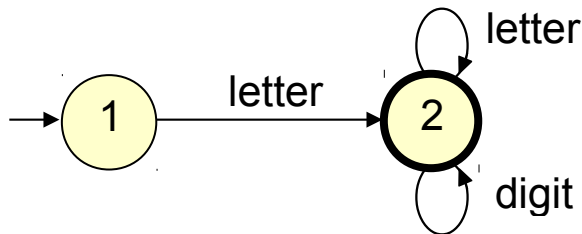
In general: insufficient



Finite Automata

- Mathematical formalism to describe certain types of algorithms (“machines”)
- In particular: to check the membership of a string in a regular language
(matching string – regexpr)
- Strong correlation between $\begin{cases} \text{finite automata} \\ \text{regular expressions} \end{cases}$

`id` \rightarrow letter (letter | digit)* \Rightarrow recognition process of an `id` by means of a diagram



- **States**: what was recognized
- **Transitions**: change of state caused by the matching of a character

- Recognition process of an `id` \rightarrow defined by the sequence of involved transitions

`a l p h a`
 $\rightarrow 1 \rightarrow 2 \rightarrow 2 \rightarrow 2 \rightarrow 2 \rightarrow 2$ $\Rightarrow \exists$ a path generating the string

- Recognized regular language = { paths within automaton }

Deterministic Finite Automata

- **Def:** A DFA is a 5-tuple $M = (\Sigma, S, T, s_0, F)$, where:

- Σ = alphabet
- S = set of states
- $T: S \times \Sigma \rightarrow S$ = transition function (deterministic)
- s_0 = initial state (unique)
- $F \subseteq S$ = set of final states (nonempty set)

- L recognized by $M \equiv L(M) = \{ x \mid x = c_1c_2...c_n, c_i \in \Sigma, \text{ and } \exists$

$$s_1 = T(s_0, c_1),$$

$$s_2 = T(s_1, c_2),$$

...

$$s_n = T(s_{n-1}, c_n), s_n \in F \}$$

$$\rightarrow s_0 \xrightarrow{c_1} s_1 \xrightarrow{c_2} s_2 \rightarrow \dots \rightarrow s_{n-1} \xrightarrow{c_n} s_n$$

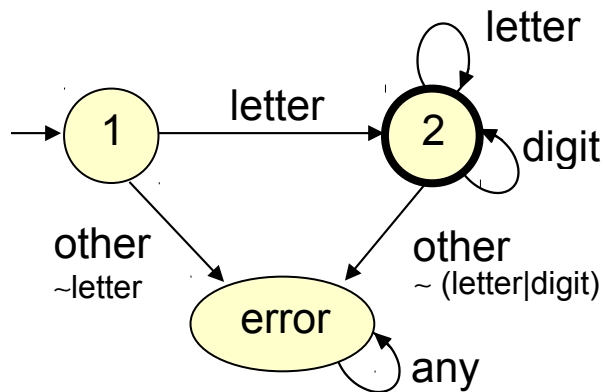
Deterministic Finite Automata (ii)

Notes comparing formal definition of DFA and diagrammatic example (**id**)

1. States identifiable by any type of labels (numbers, letters, strings, ...)
2. Factorization (abstraction) of n transitions, one \forall letter/digit, by a single transition **letter/digit**
3. Def $\Rightarrow T: S \times \Sigma \rightarrow S = \text{function} \Rightarrow T(s, c)$ must have a value $\forall (s, c) \Rightarrow$ missing transitions!



missing transitions \equiv **error transitions!**



possible input characters

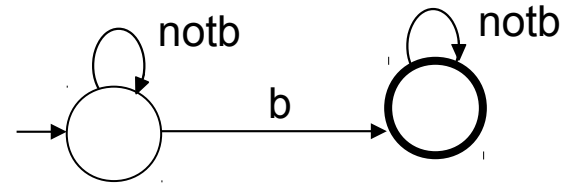
other $\equiv c \in (\Sigma - \{ \text{characters involved in the other transitions exiting the same state} \})$

context-sensitive meaning

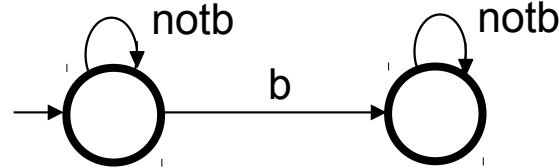
- Error state: double semantics $\left\{ \begin{array}{l} \text{it is not an identifier (from non-final state)} \\ \text{identifier followed by separator (from final state)} \rightarrow \text{maximal munch} \end{array} \right.$

Examples of DFAs in Diagrammatic Form

1. $L = \{ x \mid x \text{ includes exactly one } b \}$



2. $L = \{ x \mid x \text{ includes at most one } b \}$



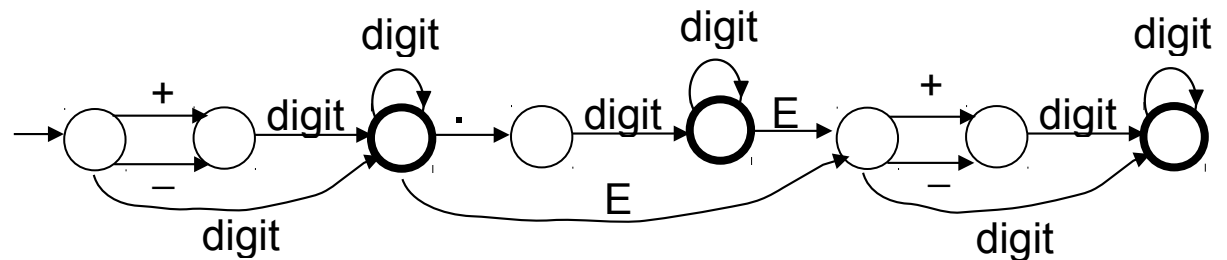
3. $L = \{ x \mid x = \text{numeric constant in scientific notation} \}$

digit $\rightarrow [0-9]$

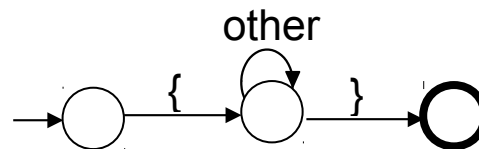
nat $\rightarrow \text{digit}^+$

snat $\rightarrow (+|-) ? \text{ nat}$

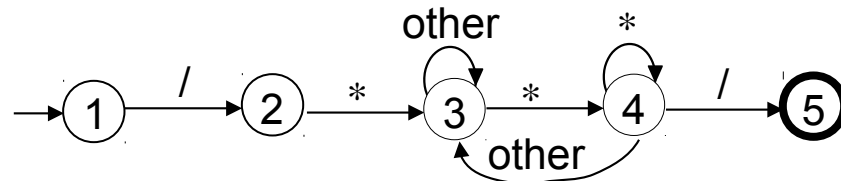
num $\rightarrow \text{snat} ("." \text{ nat}) ? (E \text{ snat}) ?$



4. $L = \{ x \mid x = \text{Pascal comment} \}$

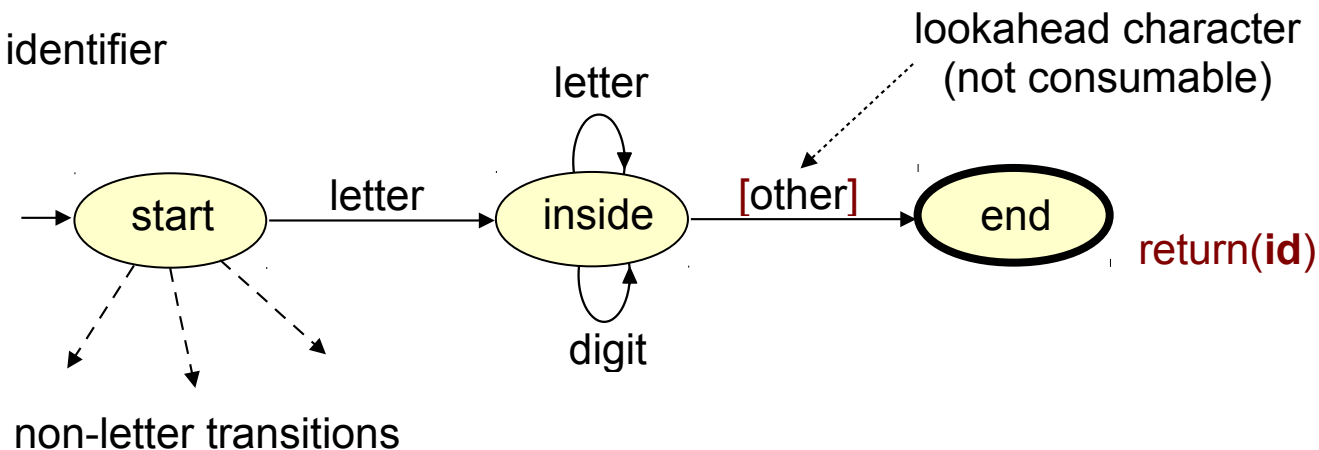


5. $L = \{ x \mid x = \text{C comment} \}$



Augmentation of Semantics in DFAs

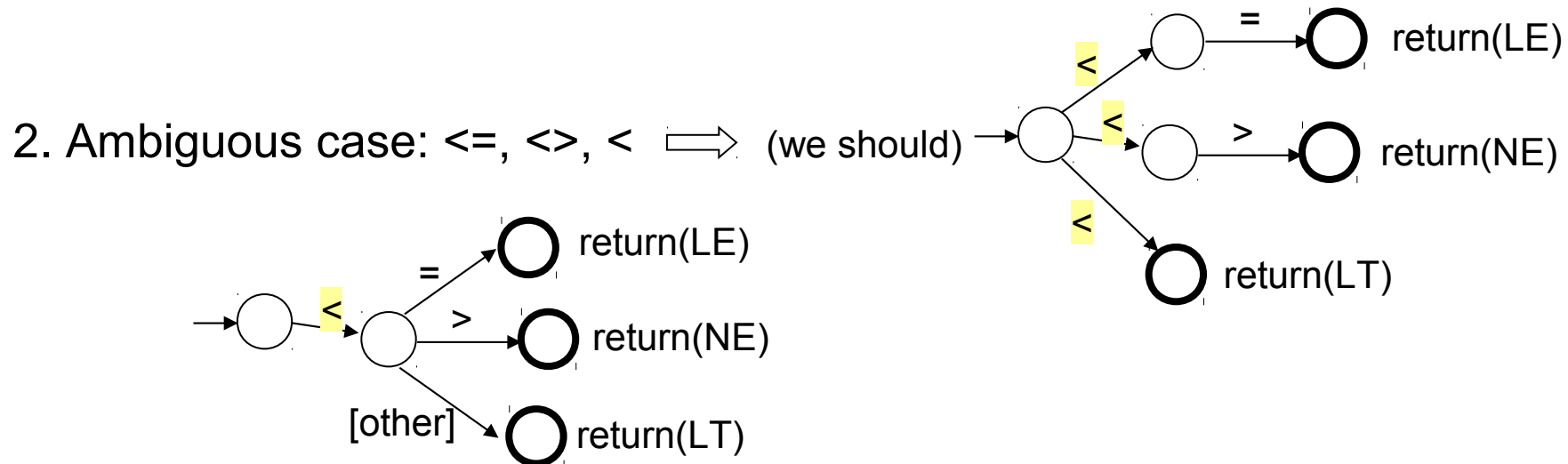
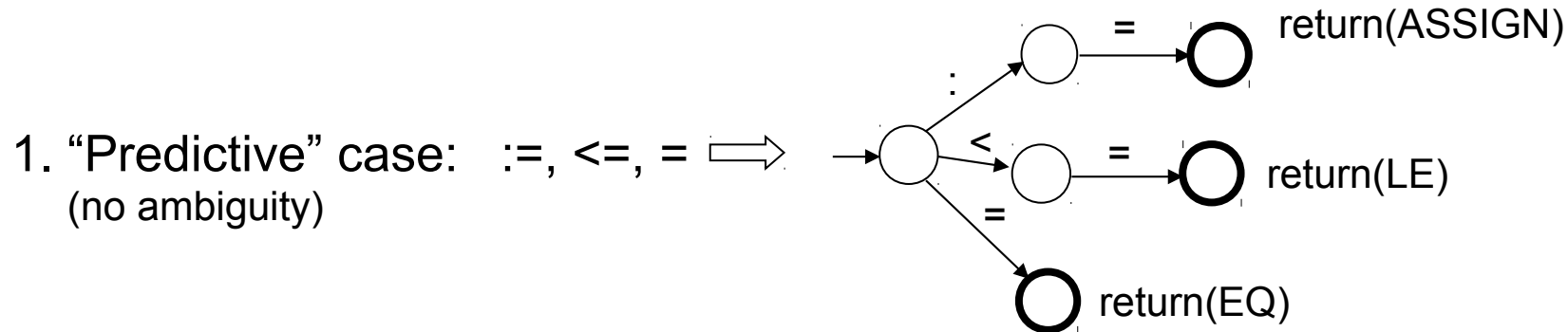
- Representation $\left\{ \begin{array}{l} \text{diagrammatic} \\ \text{formal} \end{array} \right\}$ does not describe every algorithmic aspect of the recognition process
 - Error → backtracking or return(ERROR) ?
 - Final state → ?
 - Matching of character → ? (accumulation in lexeme)
 - Maximal munch ?
- Example: Pascal identifier



Furthermore: it expresses maximal munch

Problem of Initial State

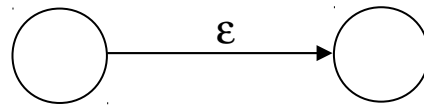
- Since \exists different tokens in a PL \rightarrow combination (fusion) of different automata



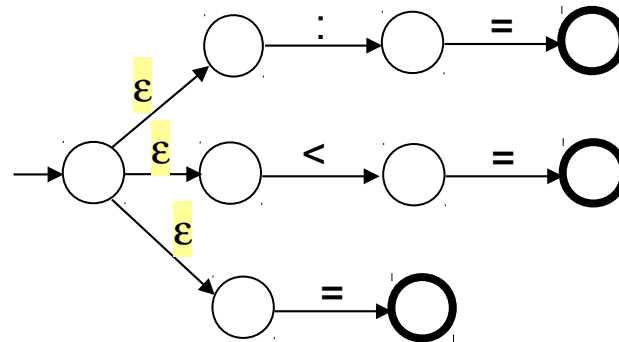
- In theory: manual combination of all tokens in one large DFA \Rightarrow complex!
Better: extending DFA to **NFA** + algorithm: **NFA** \rightarrow equivalent DFA (modular approach)

Empty Transition

- ϵ -transition $\left\{ \begin{array}{l} \text{formally: matching of empty string } (\epsilon) \\ \text{intuitively: performed "spontaneously" (without character consumption)} \end{array} \right.$



- **Advantage:** simple and systematic combination of symbol recognizers ("glue" for NFAs)



Nondeterministic Finite Automata

- **Def:** An NFA is a 5-tuple $M = (\Sigma, S, T, s_0, F)$, where:

- Σ = alphabet
- S = set of states
- $T: S \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^S$ = transition function $\rightarrow 2$ causes for nondeterminism
- s_0 = initial state
- $F \subseteq S$ = set of final states

ϵ -transitions
same c for several trans.

- L recognized by $M \equiv L(M) = \{ x \mid x = c_1c_2\dots c_n, c_i \in (\Sigma \cup \{\epsilon\}), \text{ and } \exists$

$$s_1 \in T(s_0, c_1),$$

$$s_2 \in T(s_1, c_2),$$

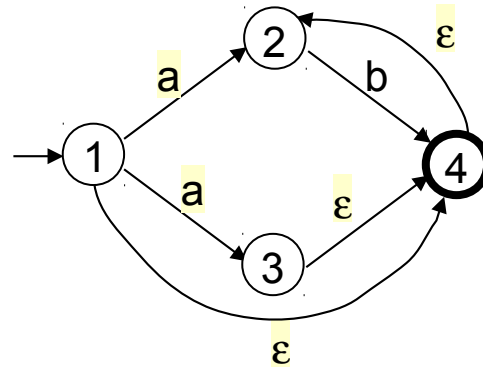
...

$$s_n \in T(s_{n-1}, c_n), s_n \in F \}$$

- Notes:

1. $x = c_1c_2\dots c_n$, actually $|x| \leq n$
2. Nondeterminism when choosing $s_i \in T(s_{i-1}, c_i)$
3. Formally, NFA does not represent an algorithm, yet can be simulated by an algorithm with backtracking (\rightarrow efficiency problem)

Examples of NFAs in Diagrammatic Form



- Recognition of abb

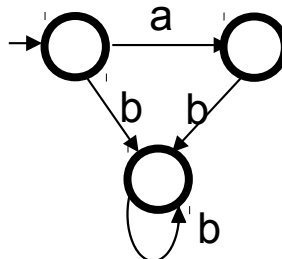
$1 \xrightarrow{a} 2 \xrightarrow{b} 4 \xrightarrow{\epsilon} 2 \xrightarrow{b} 4$

$1 \xrightarrow{a} 3 \xrightarrow{\epsilon} 4 \xrightarrow{\epsilon} 2 \xrightarrow{b} 4 \xrightarrow{\epsilon} 2 \xrightarrow{b} 4$

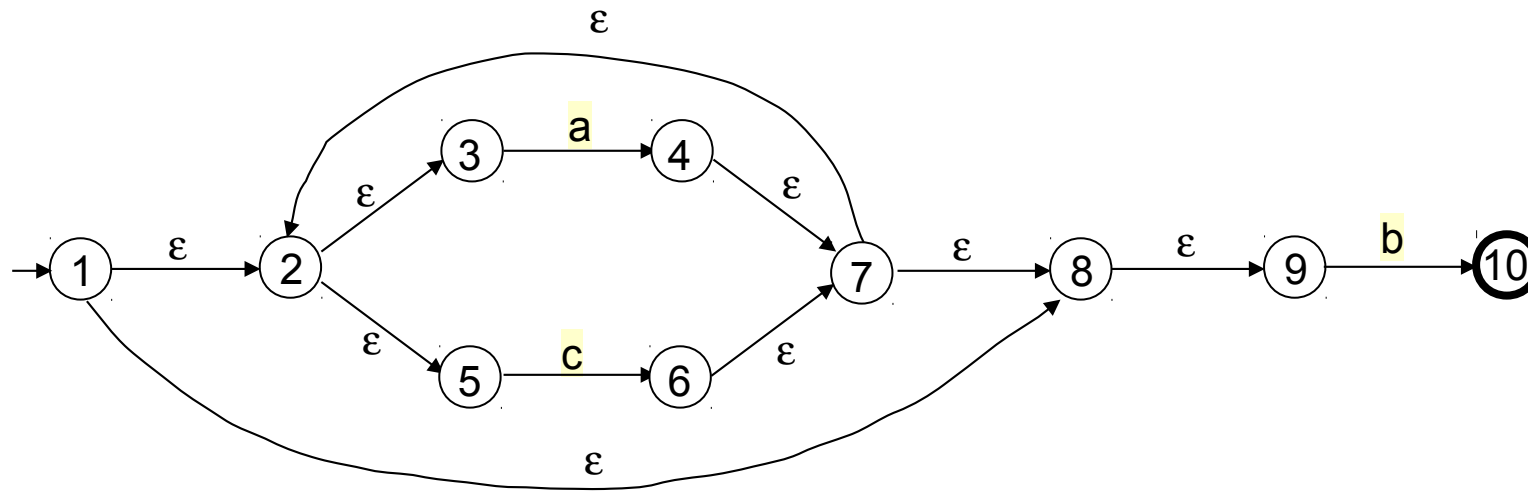
} several paths for the same string

- Recognition of $L((a \mid \epsilon)b^*) \Rightarrow$ in particular: ϵ , a , b

- Equivalent DFA:

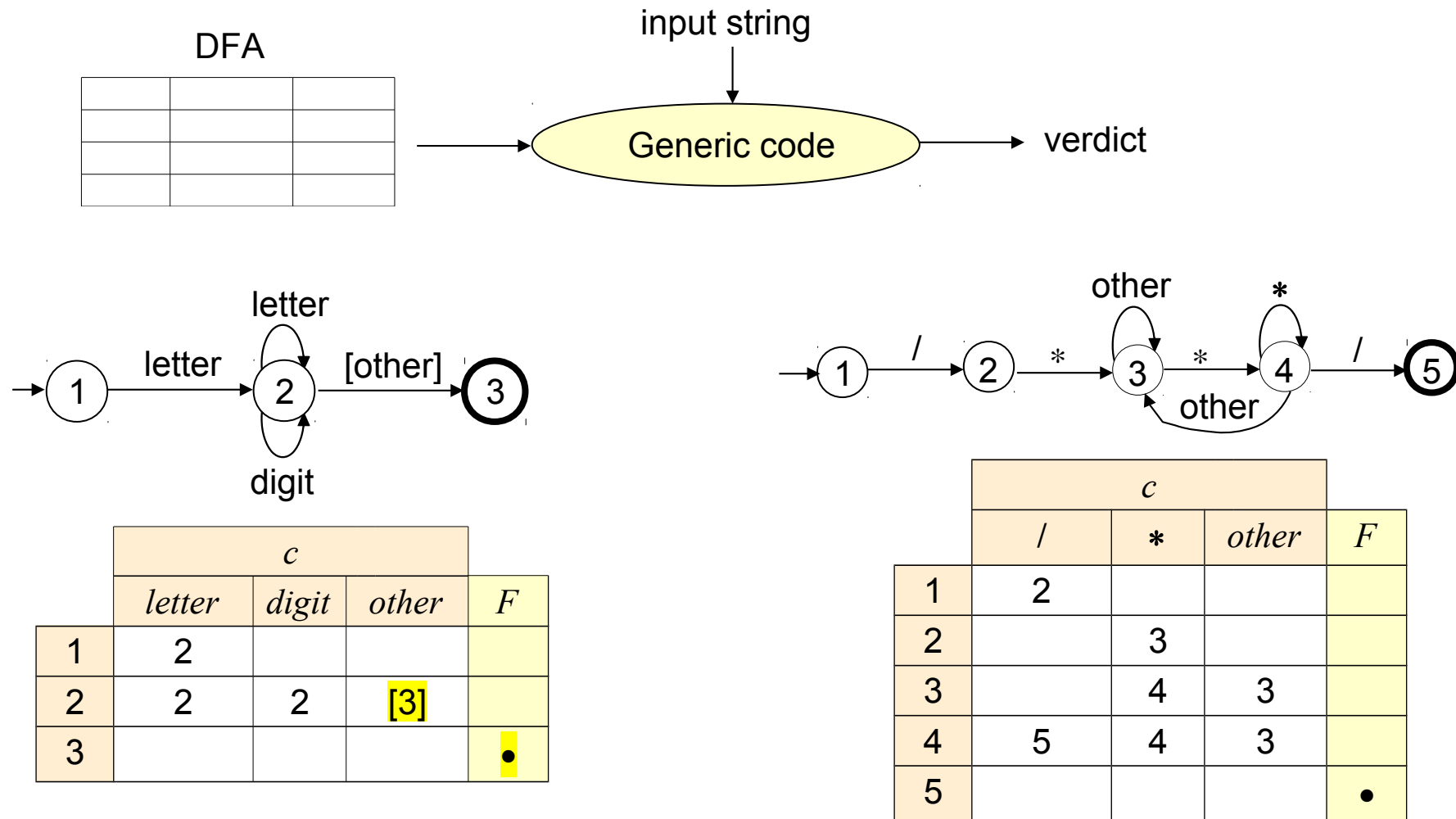


Examples of NFAs in Diagrammatic Form (ii)



$acab \in (a | c)^*b$

Automata Specified by Transition Tables



- Assumptions:
 - Blank boxes → unspecified transitions
 - s_0 = first in list
 - Information extension $\left\{ \begin{array}{l} \text{states} \in F \\ \text{consumption of } c \end{array} \right.$

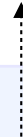
Automata Specified by Transition Tables (ii)

- Generic code \rightarrow 3 data structures (simple, yet universal) $\left\{ \begin{array}{l} T[state, c] \equiv \text{transitions} \rightarrow \text{integers} \\ Cons[state, c] \equiv \text{Booleans} \rightarrow \text{true: go forward (consumption)} \\ F[state] \equiv \text{Booleans} \rightarrow \text{true: final} \end{array} \right.$

- Table-driven algorithm:

```
state := 1; c := next();
while not F[state] and not error(state) do
  begin
    new_state := T[state, c];
    if Cons[state, c] then
      c := next();
      state := new_state
    end;
  if F[state] then accept.
```

exmp: state = -1



Automata Specified by Transition Tables (iii)

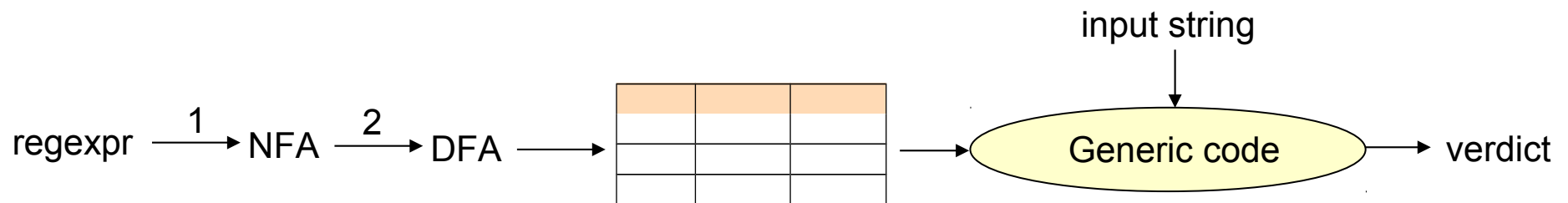
- Notes on table-driven algorithms :

1. Pros: {
Generic code
Reduced size of code
Easy code maintenance

2. Cons: waste of memory (sparse tables \rightarrow compression)

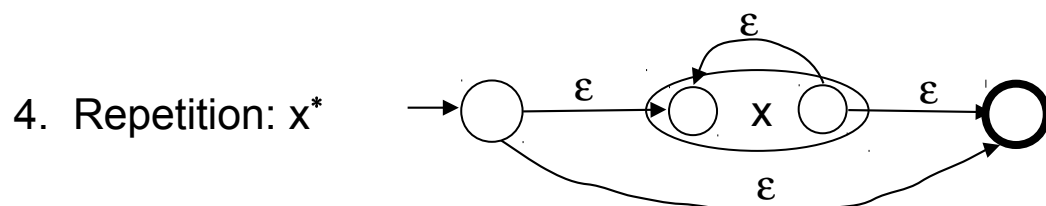
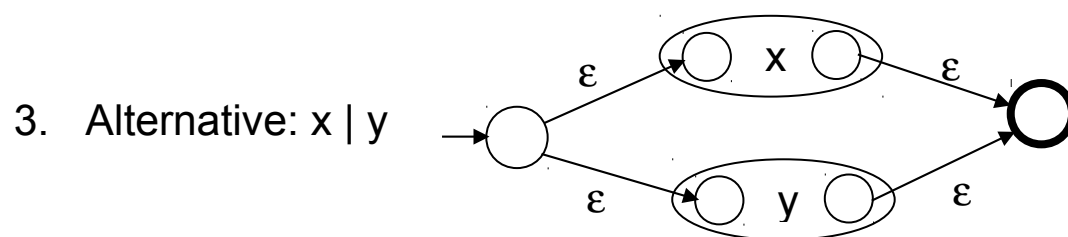
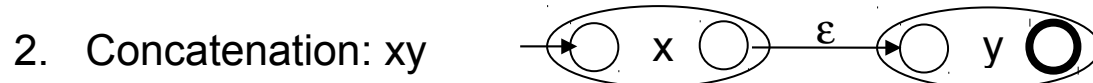
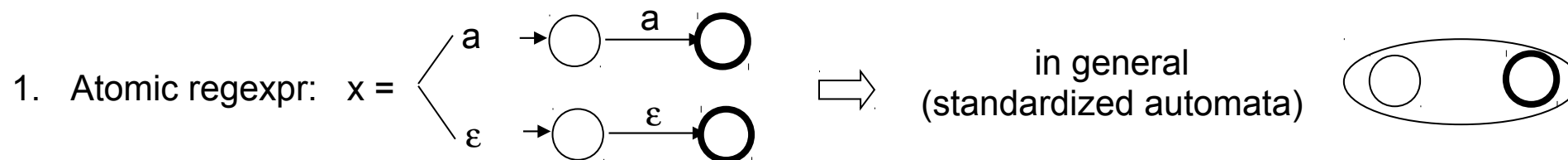
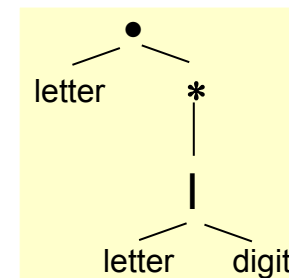
3. Technique extensible to NFA $\left\{ \begin{array}{l} T[\text{state}, c] = \{ \text{states} \} \\ \text{backtracking} \Rightarrow \text{inefficient: better transforming NFA} \rightarrow \text{DFA} \end{array} \right.$

- Underlying problem: **Regexpr \rightarrow DFA** \Rightarrow in two steps (automatic within generators):



Regexpr \rightarrow NFA (Construction of Thompson)

- Method: {
 - Incremental construction starting from automata (a, ϵ)
 - ϵ -transitions: to glue together sub-automata
 - Final automaton = recognizer of $L(r)$

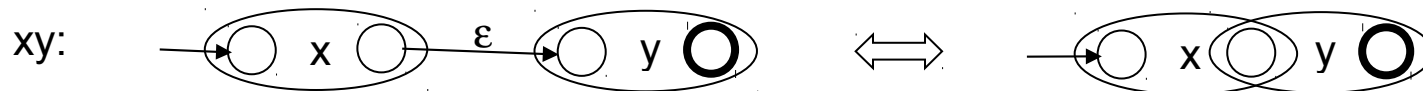


5. Parentheses: $(x) \Rightarrow$ automaton unchanged

Regexpr \rightarrow NFA (Construction of Thompson) (ii)

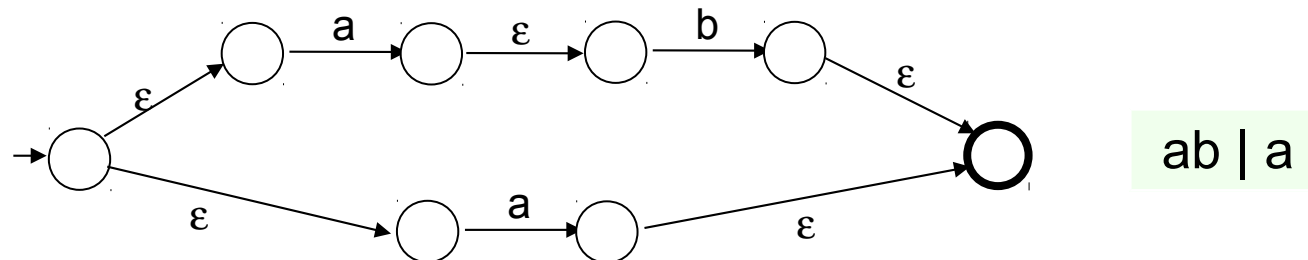
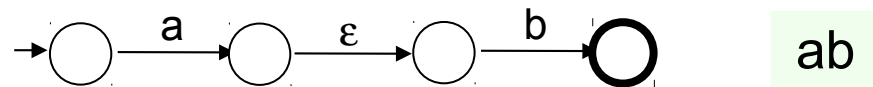
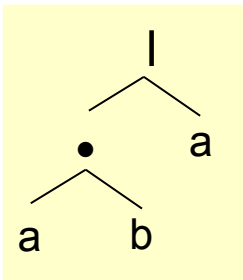
- Notes:

- \forall new state \rightarrow new name
- Properties of the final NFA (by construction):
 1. NFA recognizes $L(r)$
 2. $\text{Total-states(NFA)} \leq 2 * (\text{symbols in } r)$
 3. NFA has $\begin{cases} 1 \text{ initial state (no entering transitions)} \\ 1 \text{ final state (no exiting transitions)} \end{cases}$
 4. $\forall \text{ state} \in \text{NFA}$: possible one of two cases $\begin{cases} 1 \text{ exiting transition marked by } c \in \Sigma \\ \text{at most 2 exiting transitions marked by } \epsilon \end{cases}$
- Not unique construction (simplification):



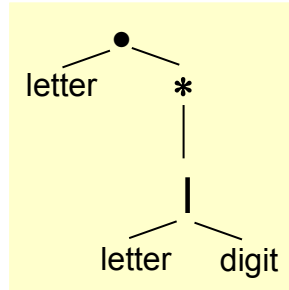
Construction of Thompson: Examples

$ab \mid a$

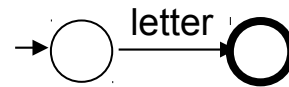


Construction of Thompson: Examples (ii)

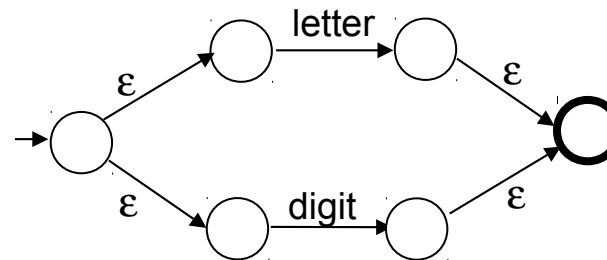
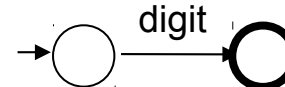
letter (letter | digit)*



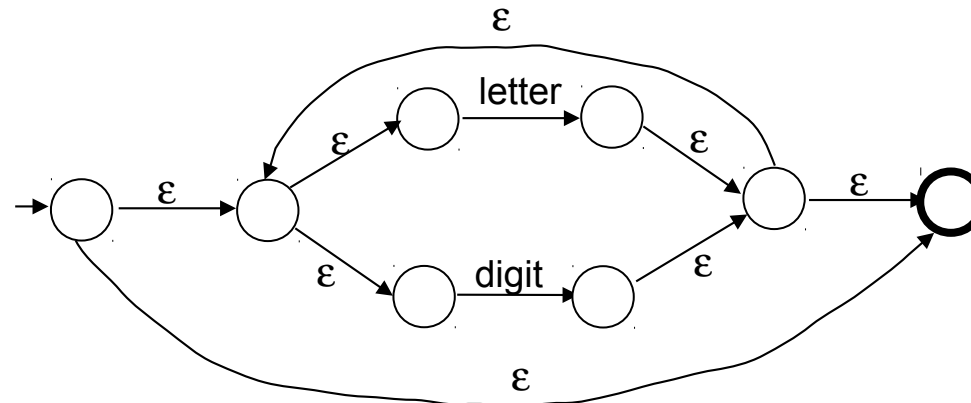
letter



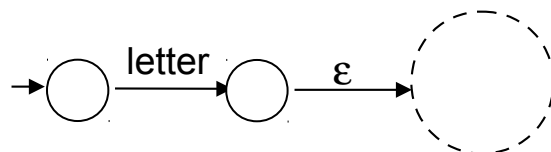
digit



letter | digit



(letter | digit)*



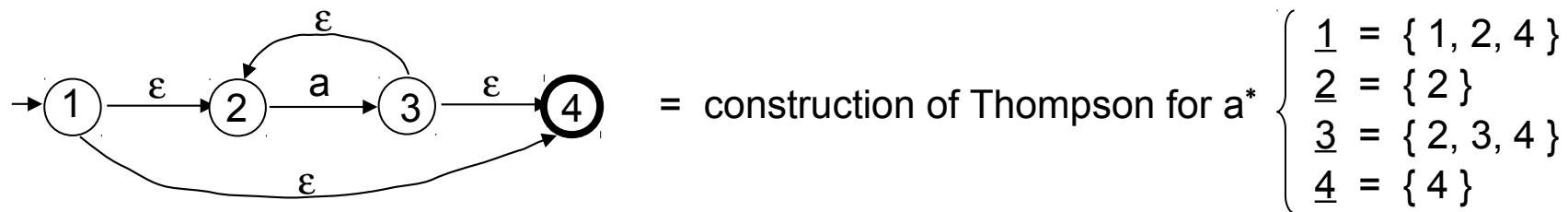
letter (letter | digit)*

NFA \rightarrow DFA

- Removal of transitions (causes for nondeterminism) $\left\{ \begin{array}{l} \text{empty } (\epsilon\text{-transitions}) \Rightarrow \epsilon\text{-closure} = \{ \text{states reachable by } \epsilon\text{-transitions} \} \\ \text{multiple (on single } c) \Rightarrow \{ \text{states reachable with same } c \} \end{array} \right.$

- Def of ϵ -closure:

1. **Single state:** $\underline{s} \equiv \epsilon\text{-closure}(s) \equiv \{ s' \mid \exists s \xrightarrow{\epsilon} \dots \xrightarrow{\epsilon} s' \text{ in } M \} \cup \{ s \}$



2. **Set of states:** $\underline{S} \equiv \epsilon\text{-closure}(S) \equiv \bigcup_{s \in S} \underline{s} \Rightarrow \underline{\{1,3\}} = \underline{1} \cup \underline{3} = \{1,2,4\} \cup \{2,3,4\} = \{1,2,3,4\}$

Subset Construction

$$\begin{array}{l} \text{NFA} \quad \text{DFA} \\ M \rightarrow M' \end{array} \left\{ \begin{array}{l} M = (\Sigma, S, T, s_0, F) \\ M' = (\Sigma, S', T', s'_0, F'), \text{ where} \end{array} \right. \left\{ \begin{array}{l} S' \subseteq 2^S \\ T': S' \times \Sigma \rightarrow S' \\ s'_0 = s_0 \\ F' \subseteq S' \end{array} \right.$$

- Computation of S' , T' , F' :

$S' := \{ s'_0 \}; T' := \emptyset;$

repeat

Choose an unmarked node $A \in S'$;

for each $c \in \Sigma$ marking a transition of M exiting a node in A **do**

begin

$A'_c := \{ z \mid s \in A, s \xrightarrow{c} z \in T \};$

$\underline{A}'_c := \varepsilon\text{-closure}(A'_c);$

if $\underline{A}'_c \notin S'$ **then** $S' := S' \cup \{ \underline{A}'_c \};$

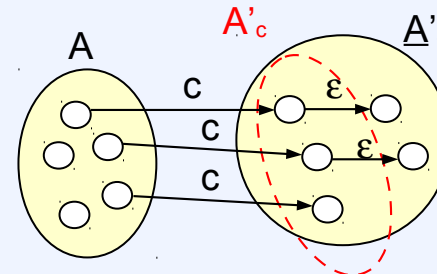
$T' := T' \cup \{ A \xrightarrow{c} \underline{A}'_c \}$

end;

Mark A

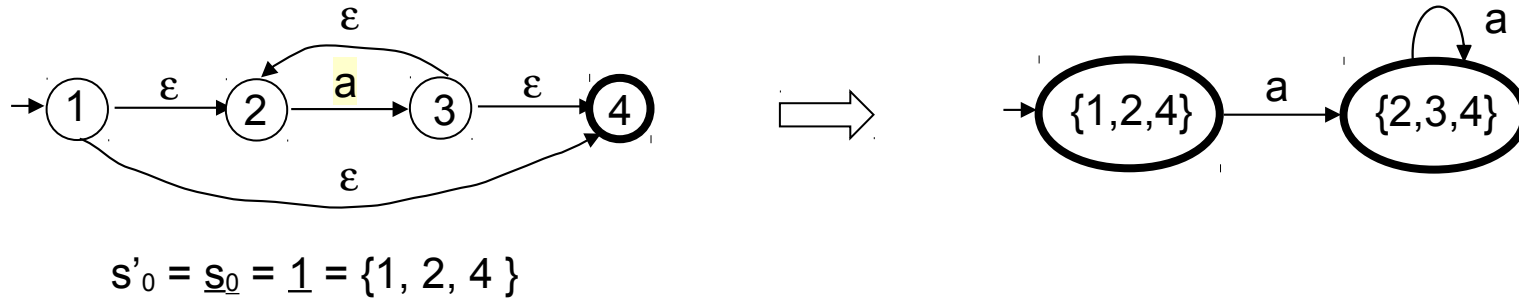
until all elements in S' are marked;

$F' := \{ A \mid A \in S', A \cap F \neq \emptyset \}.$

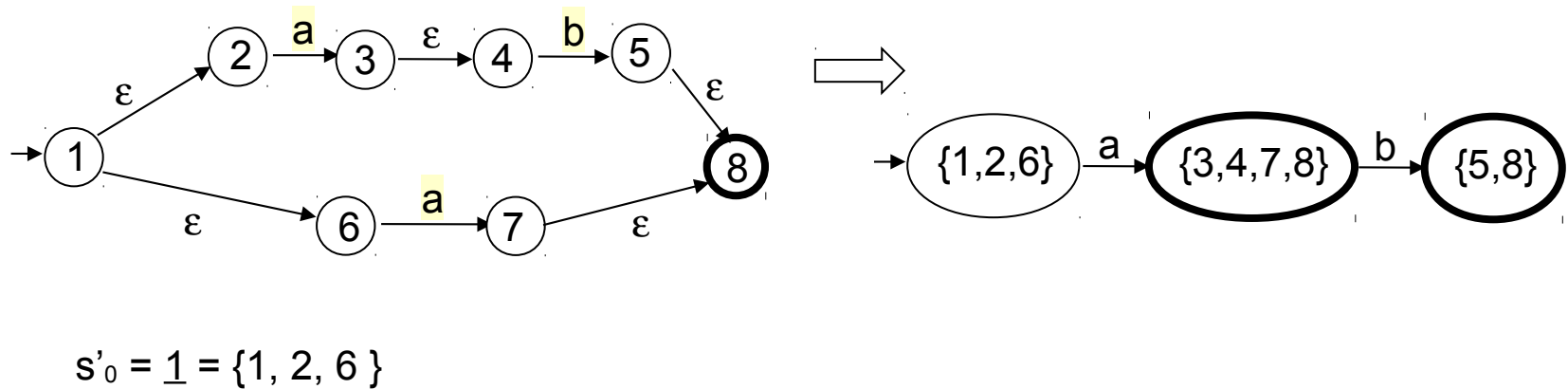


Examples NFA \rightarrow DFA

1. a^*

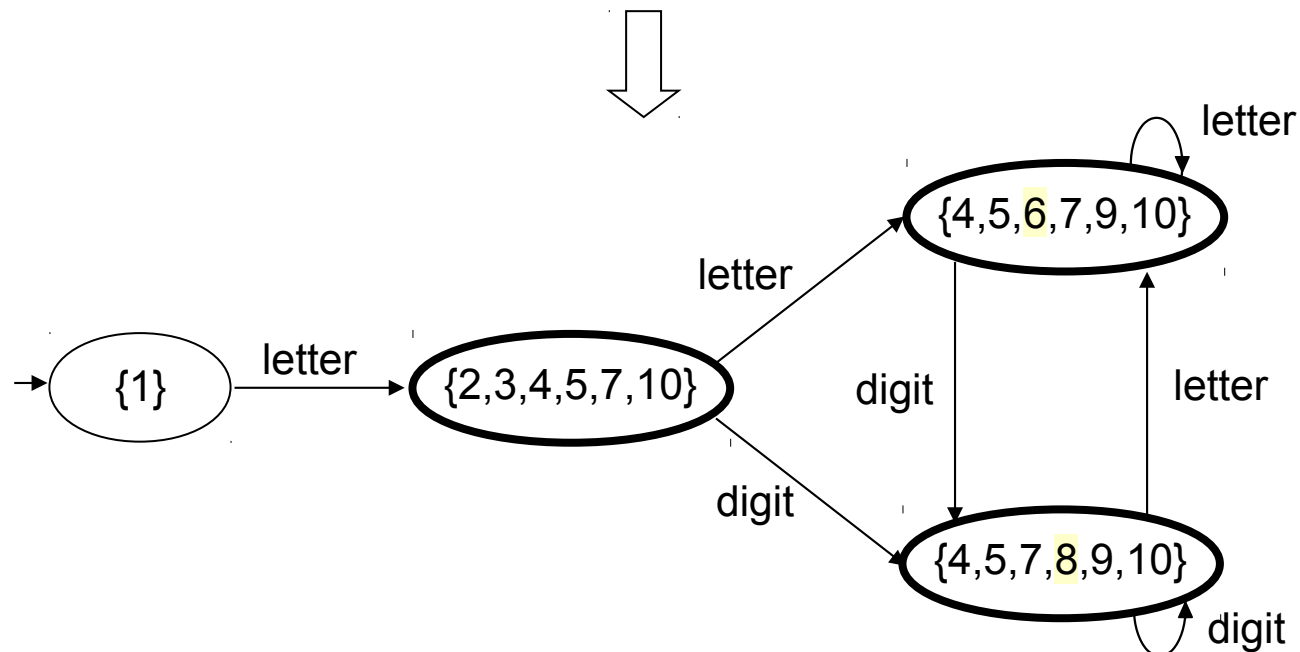
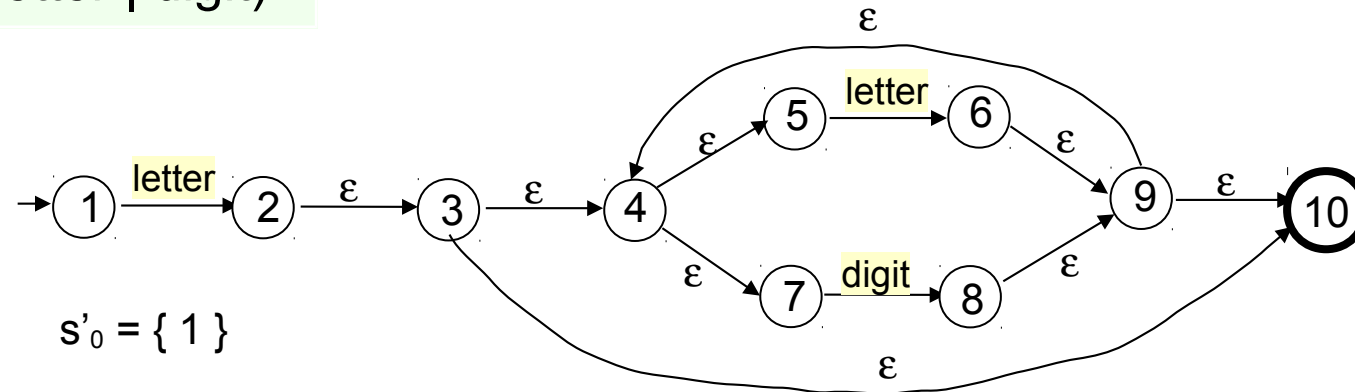


2. $ab \mid a$



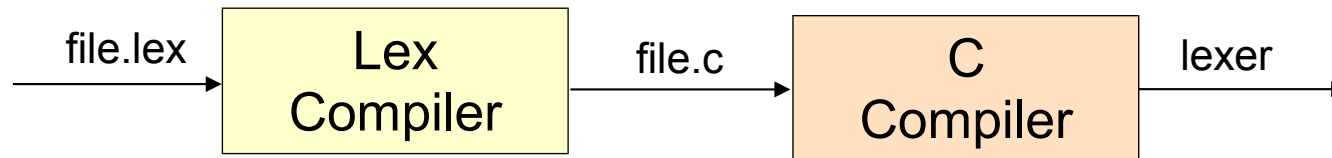
Examples NFA → DFA (ii)

3. letter (letter | digit)*

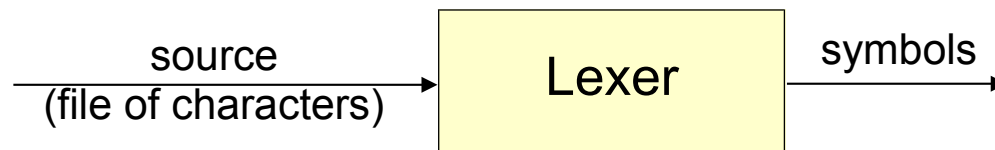


Lex (lexical analyzer generator)

- Broad applicability: sw tool recognizing patterns of characters defined by regexp
→ then ?
- Lex $\left\{ \begin{array}{l} \text{language (regexp, actions)} \\ \text{compiler} \rightarrow \text{C code (source-to-source)} \end{array} \right.$



- Lexical analyzer:



Lex (ii)

Declarations

%%

Translation rules

%%

Auxiliary functions

- Lex specification (program): 3 sections

- Declarations
 - black box (auxiliary definitions): **%{ #include, constants, variables %}**
 - white box (regular definitions): *name expr*

- Translation rules

<i>Regexpr</i>	<i>C fragment</i>
r_1	{ action ₁ }
r_2	{ action ₂ }
r_n	{ action _n }

- Auxiliary functions: necessary to actions → possibly compiled separately

- Flex (GNU) → **yyllex()** = C function in **file.c** that
 - recognizes the symbol
 - executes the action

Lex (iii)

- Notation for the specification of regular expressions

- Matching of strings of characters: `if` or `"if"`

- Neutralization of the meta-character effect $\begin{matrix} \text{"(}" \\ \backslash(\end{matrix} \longrightarrow \text{acts on a single character}$

- Uniformity with C to denote space characters $\begin{matrix} \backslash t \\ \backslash n \end{matrix}$

- Meta-character `•` \equiv any character $\neq \backslash n$

- Canonical interpretation of meta-characters `* + () | ?` \longrightarrow textual form

Example: " Set of strings of *a*, *b*, which start with *aa* or *bb*, and end with an optional *c* "

$$\Sigma = \{a,b,c\} \quad (aa \mid bb) (a \mid b)^* c?$$

- Complement set: `^` = first character within the range `[^0-9abc]` $= \Sigma - \{0,1,\dots,9,a,b,c\}$

Example: "Numbers in scientific notation" `("+" | "-")? [0-9]+ ("." [0-9]+)? (E ("+" | "-")? [0-9]+)?`

- Reference to names of regexpr by `{name}`

<code>nat</code>	<code>[0-9]+</code>	
<code>snat</code>	<code>("+" "-")? {nat}</code>	\longrightarrow only when referenced

- Further meta-characters $\begin{matrix} \text{(expr)\{n\}, (expr)\{n,m\}} \\ \text{(expr)_1/(expr)_2} \end{matrix}$ (*m* missing \rightarrow infinite)

 \searrow `(expr)1/(expr)2`: matching of `expr1` only if followed by `expr2` (contextual)

Lex (ex1.lex)

Generation of lines preceded by their position number

```
%{
#include <stdio.h>
int l = 1;
%}
%option noyywrap .....► termination option
line    .*\\n .....► regular definition
%%
{line}  {printf("%d %s", l++, yytext);} .....► translation rule
%%
main()
{
    yylex();
}
```

.....► lexical string (lexeme)

Lex (ex2.lex)

Generation of lines in odd position

```
%{
#include <stdio.h>
int l = 1;
%}
%option noyywrap
line    .*\\n
%%
{line}  { if(l++%2)
          printf("%s", yytext);
        }

%%
main()
{
    yylex();
}
```

Lex (ex3.lex)

Substitution of numbers from decimal to hexadecimal notation + printing of number of actual substitutions

```
%{
#include <stdio.h>
#include <stdlib.h>
int cont = 0;
}%
%option noyywrap
digit      [0-9]
num        {digit}+
%%
{num} {    int n = atoi(yytext);
           printf("%x", n);
           if (n > 9) cont++; }
%%
main()
{
    yylex();
    fprintf(stderr, "Tot substitutions = %d\n", cont);
}
```

Note: Default action: when a character (or a string of characters) is not part of any symbol



ECHO on output

Lex (ex4.lex)

Substitution of numbers with value ≥ 10 from decimal to hexadecimal notation + printing of number of substitutions

```
%{
#include <stdio.h>
#include <stdlib.h>
int cont = 0;
}%
%option noyywrap
digit      [0-9]
num        {digit}{digit}+
%%
{num} {    int n = atoi(yytext);
           printf("%x", n);
           cont++; }

%%
main()
{
    yylex();
    fprintf(stderr, "Tot substitutions = %d\n", cont);
}
```


Lex (ex5.lex)

*Selection of lines which either start or end with character **a***

```
%{
#include <stdio.h>
%}
%option noyywrap
a_line  a.*\n
line_a  .*a\n
%%
{a_line} ECHO;
{line_a} ECHO;
.*\n    ;
%%
main()
{
    yylex();
}
```

Note: Ambiguous set of rules (a string may correspond to different regexpr, e.g: **a**)



Priority rules (built-in)

1. Principle of maximal munch.
2. If \exists several rules matching the string \rightarrow select the rule specified first.

output of yytext
empty action

.*\n	;
{a_line}	ECHO;
{line_a}	ECHO;

\Rightarrow empty output!

{a_line}	ECHO;
{line_a}	ECHO;

\Rightarrow output = input!

Lex (ex6.lex)

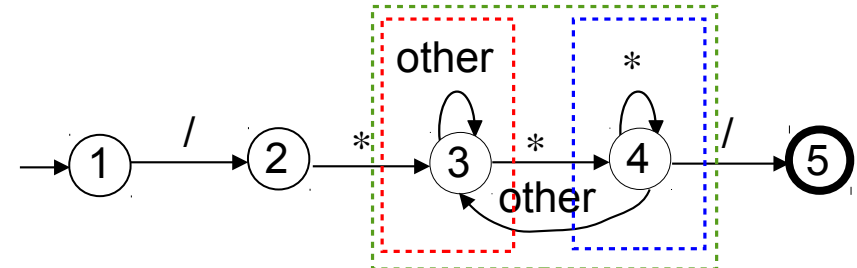
*Selection of lines which both start and end with character **a** and contain other characters, all different from **a***

```
%{
#include <stdio.h>
%}
%option noyywrap
a_line_a    a[^a\n]+a\n
%%
{a_line_a} ECHO;
.*\n      ;
%%
main()
{
    yylex();
}
```

Lex (ex7.lex)

Substitution of uppercase with lowercase letters, except those in C comments

```
%{
#include <stdio.h>
#define FALSE 0
#define TRUE 1
}%
%option noyywrap
%%
[A-Z]    {putchar(tolower(yytext[0]));}
"/*"    {char c; int end = FALSE;
        ECHO;
        do { while ((c=input()) != '*')
                putchar(c);
              while ((c=input()) == '*')
                putchar(c);
              if(c == '/')
                end = TRUE;
            } while(!end);
        }
%%
main(){ yylex(); }
```



Internal Lex identifiers:

Identifier	Description
yylex()	Function of lexical analysis
yytext	Lexical string (<i>lexeme</i>)
yylen	strlen(yytext)
yyin	Input file (default: stdin)
yyout	Output file (default: stdout)
input()	Input of one character ← yyin
ECHO	Default action: prints yytext on yyout

Lex (ex8.lex)

Substitution of uppercase with lowercase letters, except those in Pascal comments

```
%{  
#include <stdio.h>  
%}  
%option noyywrap  
%%  
[A-Z]      {putchar(tolower(yytext[0]));}  
\{[^\}]*\} ECHO;  
%%  
main()  
{  
    yylex();  
}
```

```
{ (~)* }
```

Lex (ex9.lex)

Counting of chars, words and lines (wc), where word = list of non-blank chars

```
%{
#include <stdio.h>
int nc=0, nw=0, nl=0;
%}
%option noyywrap
word    [^ \t\n]+
eol     \n
%%
{word}  {nw++; nc+=yyleng;}
{eol}   {nl++; nc++;}
.       {nc++;}
%%
main()
{
    yylex();
    printf("%d %d %d\n", nl, nw, nc);
}
```

.....> either space or tab

Lex: Lexical Analysis

Recognition of lexical symbols in a programming language

```
%{
#include <stdlib.h>
#include "def.h" /* IF, THEN, ELSE, ID, NUM, RELOP, LT, LE, EQ, NE, GT, GE */
int lexval;
%}
delimiter    [ \t\n]
spacing      {delimiter}+
letter       [A-Za-z]
digit        [0-9]
id           {letter}({letter}|{digit})*
num          {digit}+
%%
{spacing}    ;
if           {return(IF);}
then         {return(THEN);}
else         {return(ELSE);}
{id}         {lexval = store_id(); return(ID);}
{num}        {lexval = atoi(yytext); return(NUM);}
"<"         {lexval = LT; return(RELOP);}
"<="        {lexval = LE; return(RELOP);}
"="          {lexval = EQ; return(RELOP);}
"<>"        {lexval = NE; return(RELOP);}
">"         {lexval = GT; return(RELOP);}
">="        {lexval = GE; return(RELOP);}
%%
int store_id() /* symbol table without keywords */
{ int line;
  if((line = lookup(yytext)) == 0) line = insert(yytext);
  return(line);
}
```