

Lexical Analysis



Matthew Might
University of Utah
matt.might.net

Join the google group!

Plan

- What is lexical analysis?
- Formal languages
- Regular languages
- Deterministic automata
- Nondeterministic automata
- Lexer-generation: `flex`

```
function id(x)
{
    return x ;    // comment
}
```

FUNCTION

IDENT(id)

LPAR

IDENT(x)

RPAR

LBRACE

RETURN

IDENT(x)

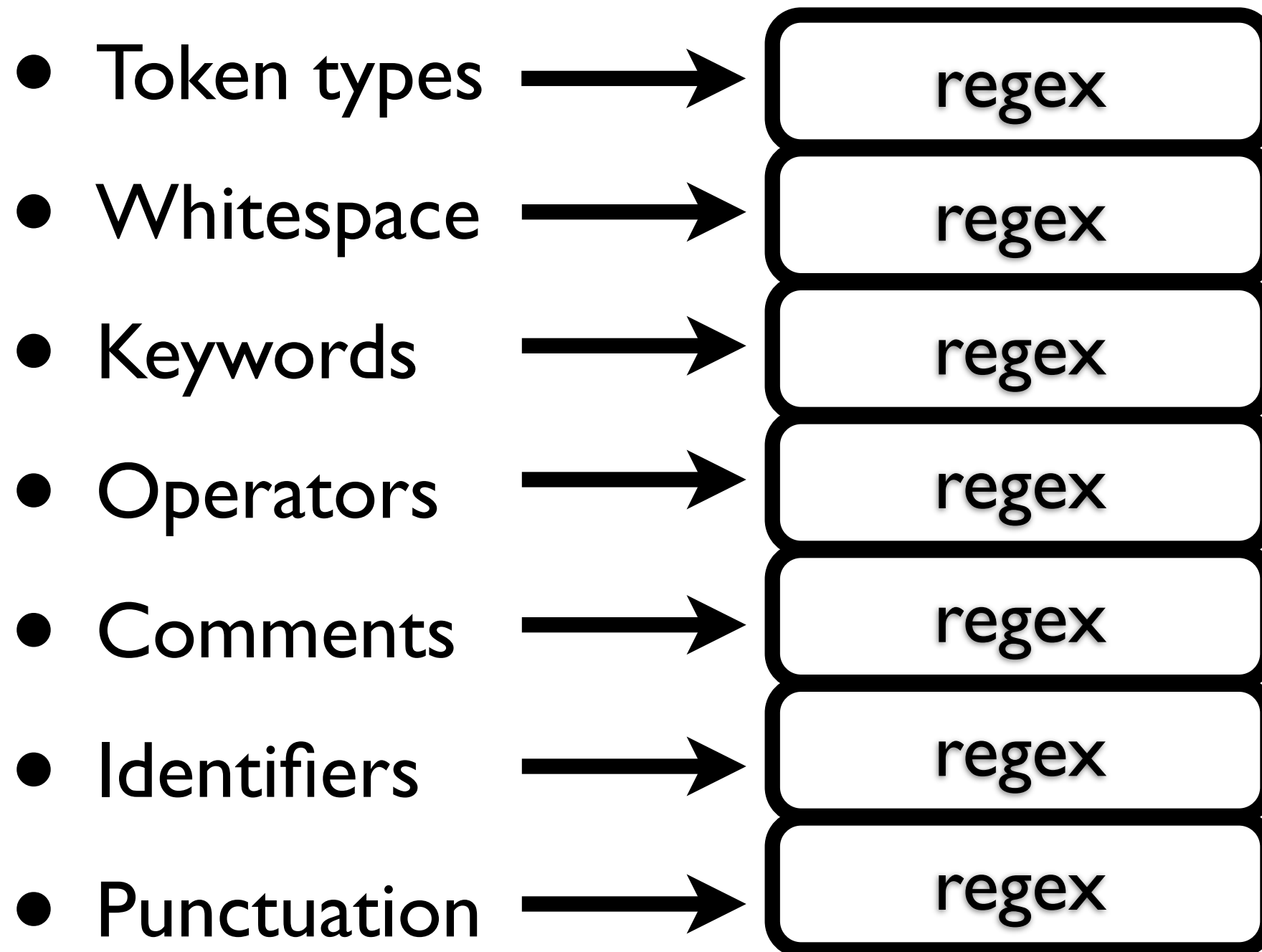
SEMI

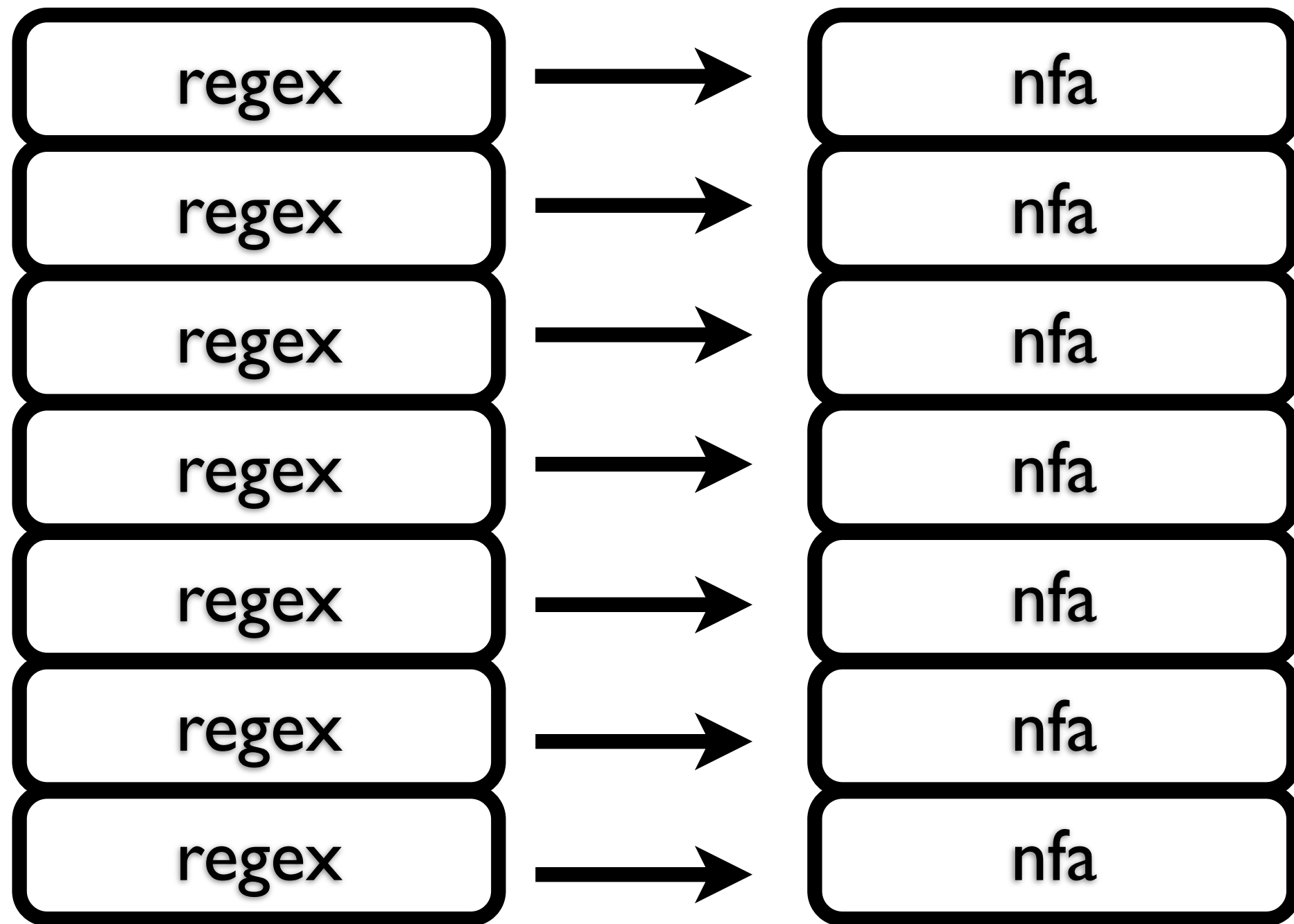
RBRACE

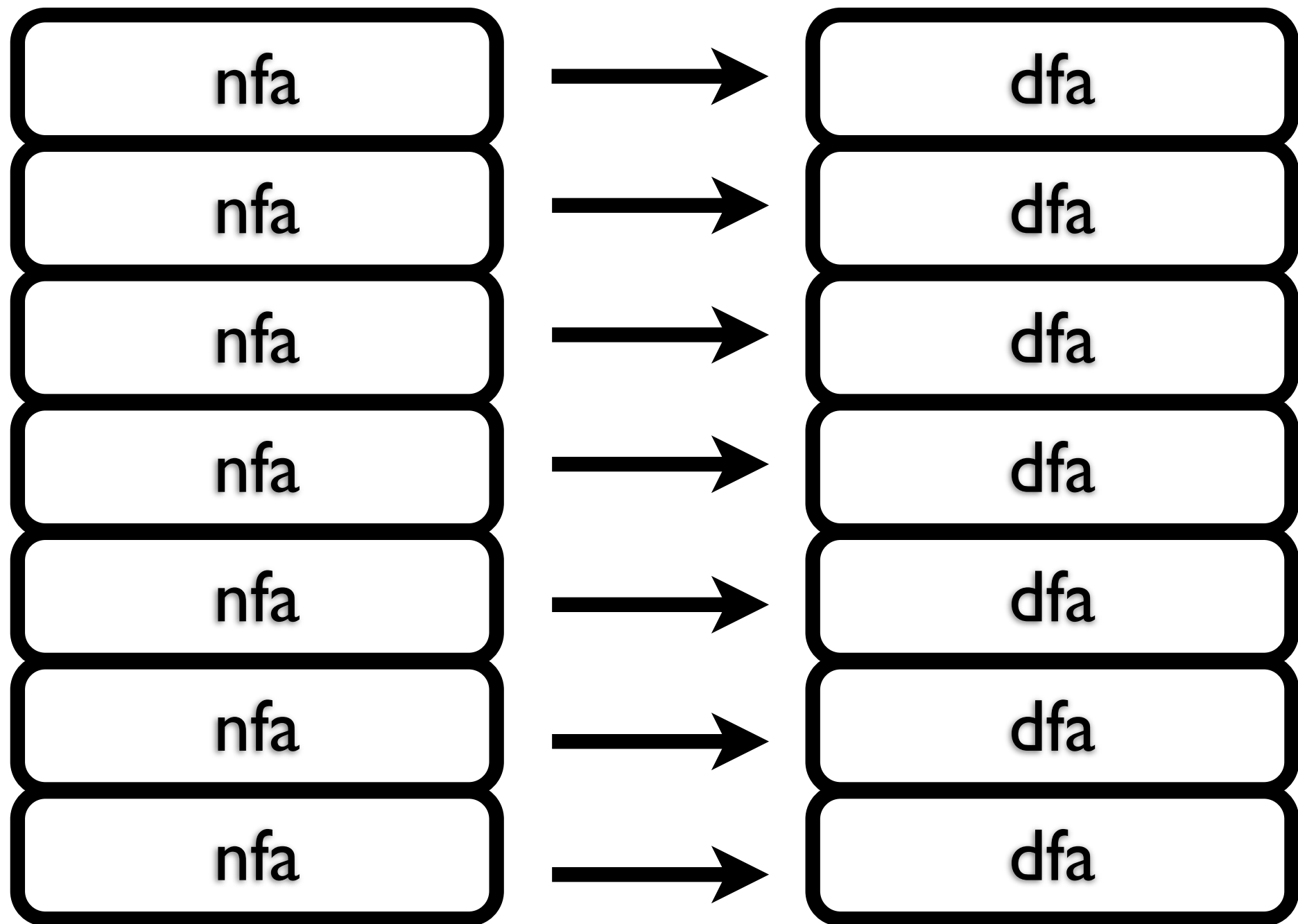
Lexical specification

- Token types
- Whitespace
- Keywords
- Operators
- Comments
- Identifiers
- Punctuation

- Token types
- Whitespace
- Keywords
- Operators
- Comments
- Identifiers
- Punctuation







dfa

dfa

dfa

dfa

dfa

dfa

dfa

fsm



$(x + 3)$



fsm



$(x + 3)$



fsm



LPAR ID(x) ...



Formal languages

A formal language is a
set of strings over an alphabet.

A **string** is a sequence
of characters from an alphabet.

Notation

A, Σ – Alphabets

L – Language

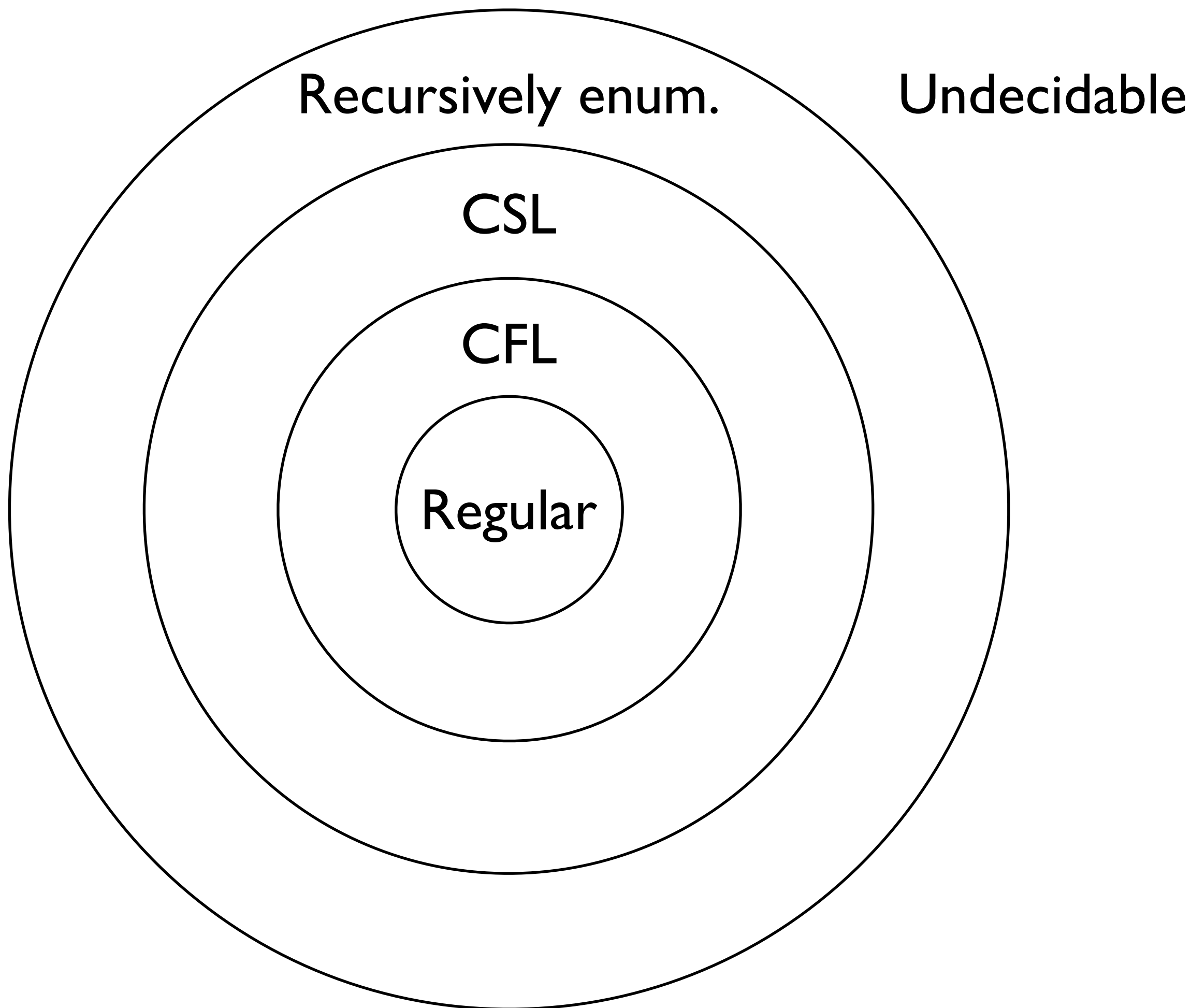
$L \subseteq A^*$ – Language over alphabet A

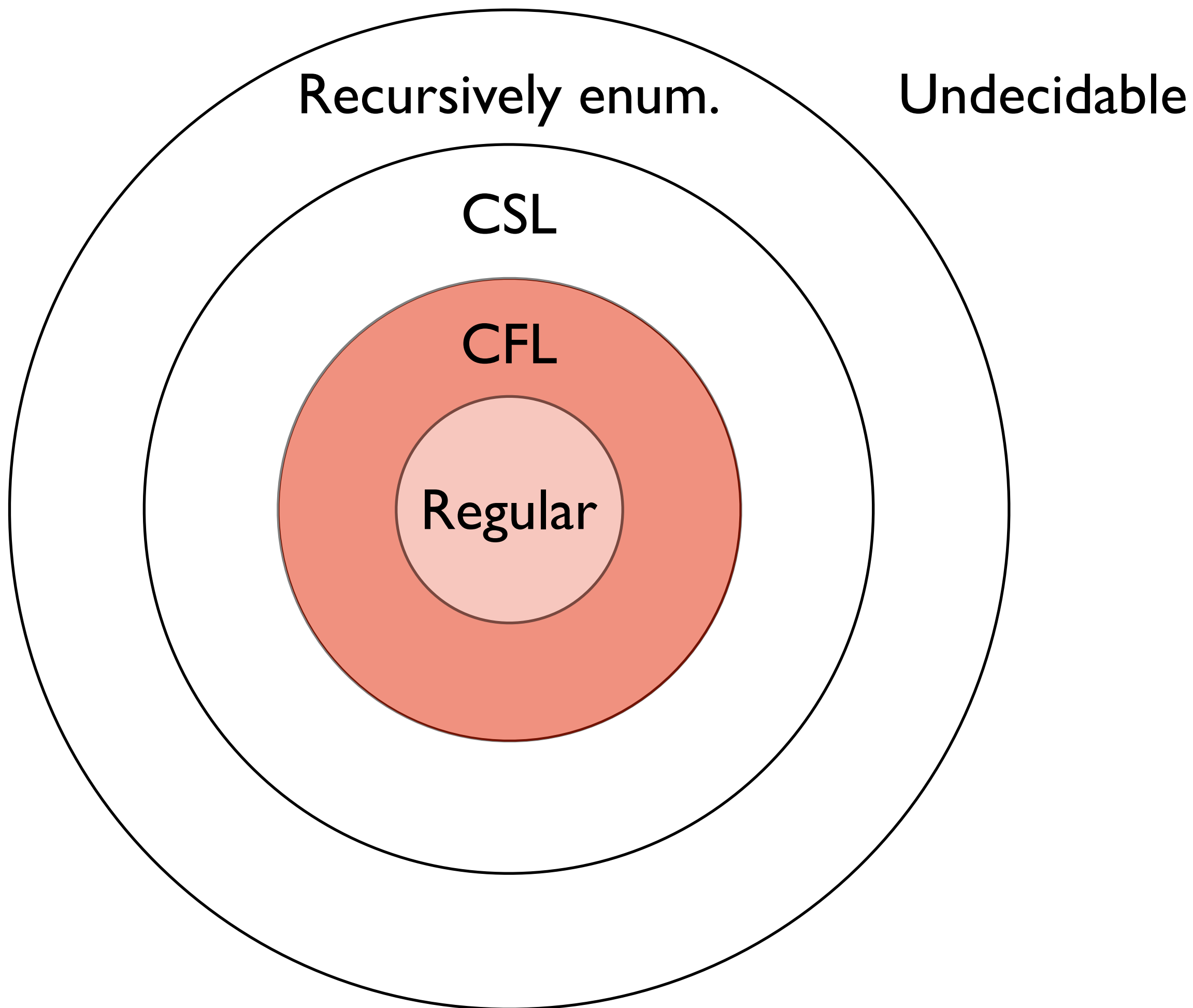
Examples

$$A = \{a, b\}$$

$$L_1 = \{bab, abba\}$$

$$L_2 = \{a, aa, aaa, \dots\}$$





Atomic languages

\emptyset

€



$$\epsilon = \{ \text{" " " " } \}$$

$$\epsilon = ""$$

Primitive languages

c

$$c = \{ " c " \}$$

Operations

Concatenation

$$L_1 \cdot L_2 = \{w_1w_2 : w_1 \in L_1 \text{ and } w_2 \in L_2\}$$

Exercise

$$L_1 = \{a, b\}$$

$$L_2 = \{c, d\}$$

Union

$$L_1 \cup L_2 = \{w : w \in L_1 \text{ or } w \in L_2\}$$

Exercise

$$L_1 = \{a, b\}$$

$$L_2 = \{c, d\}$$

Exponentiation

$$L^n = \{w_1 \dots w_n : w_i \in L\}$$

$$L^0 = \epsilon$$

Exercise

$$\{a, b\}^3$$

Kleene star

$$L^{\star} = \bigcup_{n=0}^{\infty} L^n$$

Exercise

$$\{a, b\}^*$$

Concatenation,
union,
Kleene star.

Regular Languages

More regular operations!

Option

$$L^? \equiv L \cup \{\epsilon\}$$

Kleene plus

$$L^+ = \bigcup_{n=1}^{\infty} L^n$$

Intersection

$$L_1 \cap L_2 = \{w : w \in L_1 \text{ and } w \in L_2\}$$

Difference

$$L_1 - L_2 = \{w : w \in L_1 \text{ and } w \notin L_2\}$$

Complement

$$\overline{L} = \{w : w \notin L\}$$

Reversal

$$L^R = \{ \langle a_n, \dots, a_1 \rangle : \langle a_1, \dots, a_n \rangle \in L \}$$

Prefix

$$L^{\leq} = \{w_1 : \text{there exists } w_2 \text{ such that } w_1 w_2 \in L\}$$

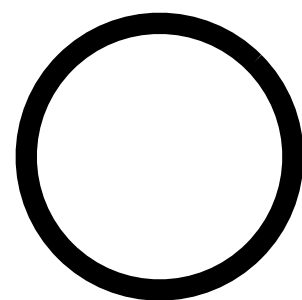
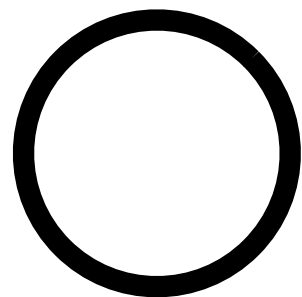
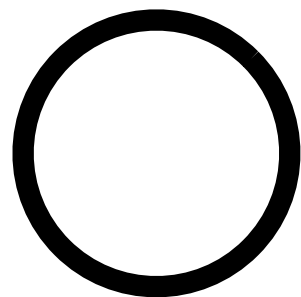
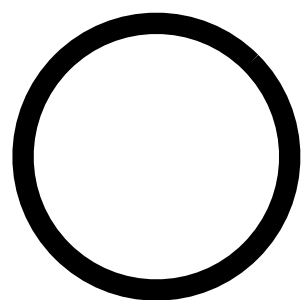
$$w \in L?$$

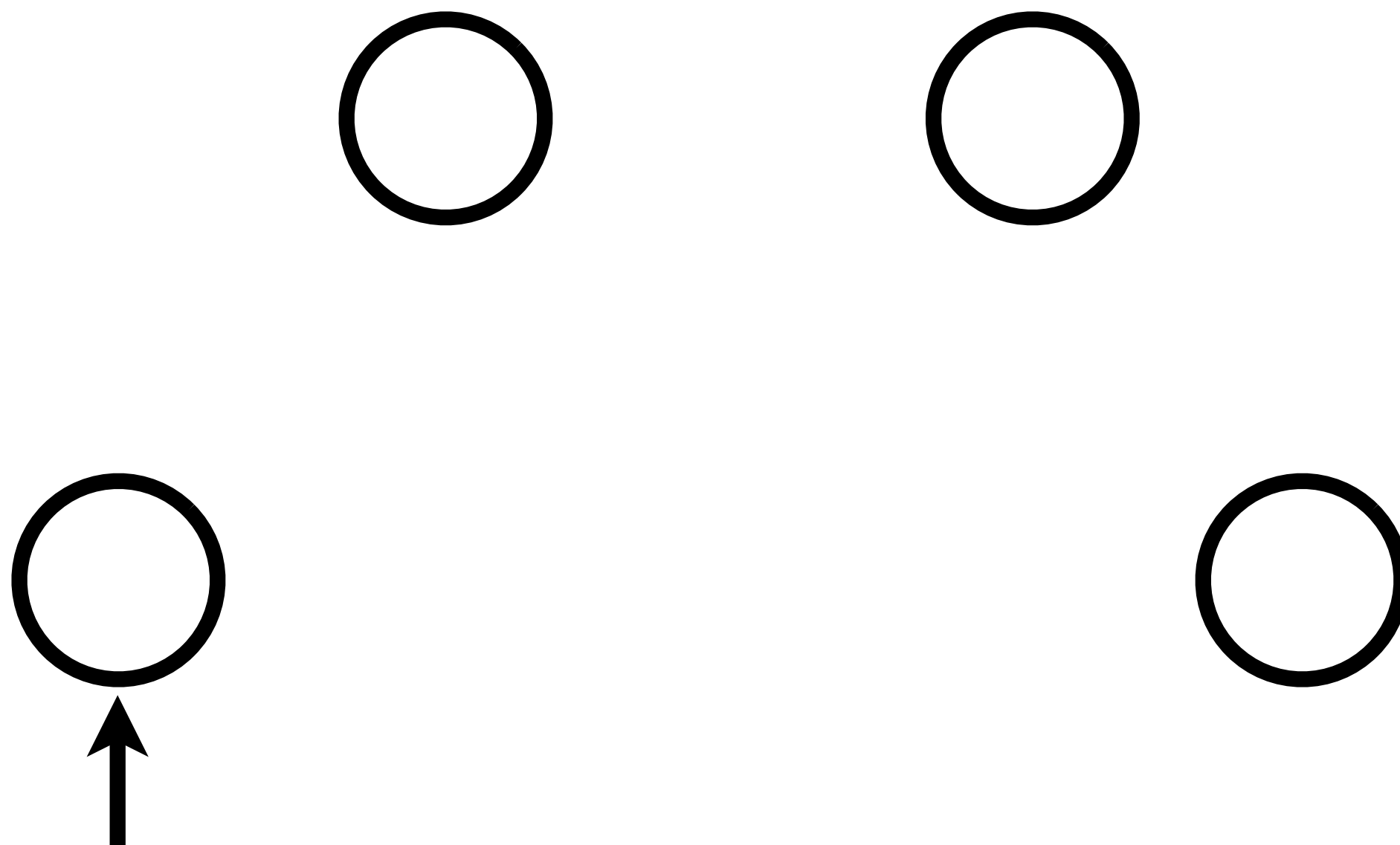
RegEx \Rightarrow NFA \Rightarrow DFA

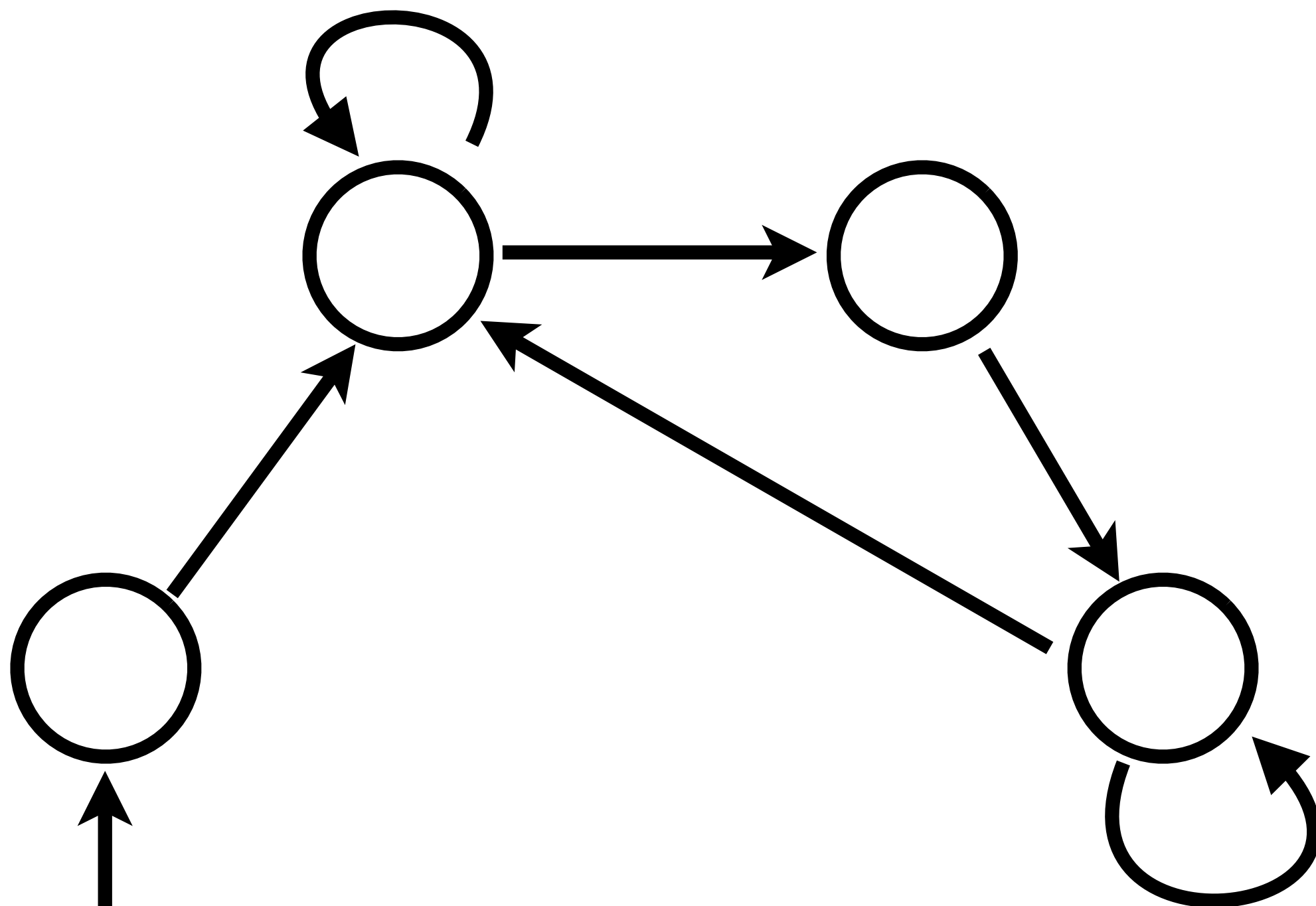
An **automaton** is a state machine
that recognizes a formal language.

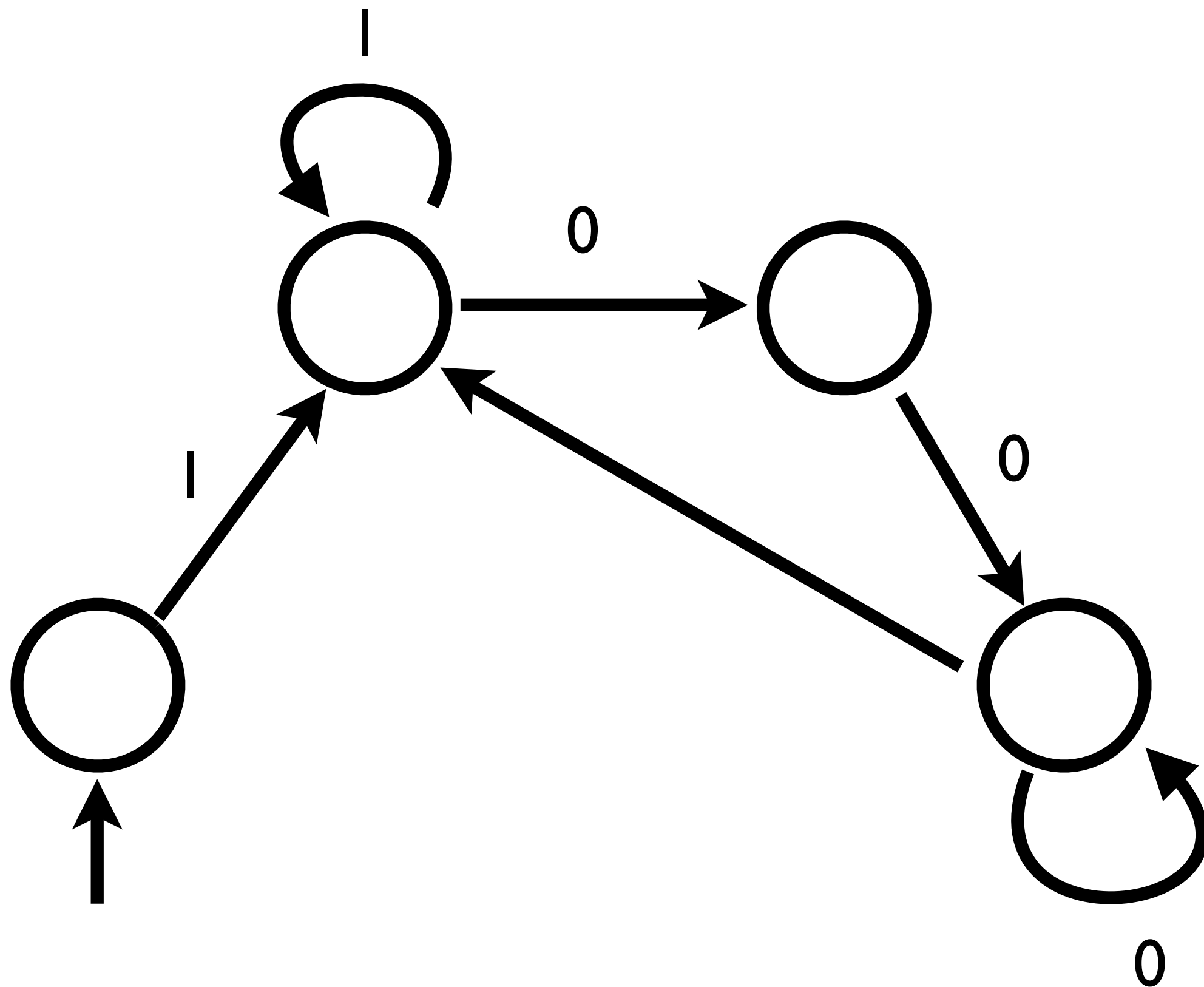
Formally, a **deterministic automaton** M is a 5-tuple (A, Q, q_0, δ, F) where:

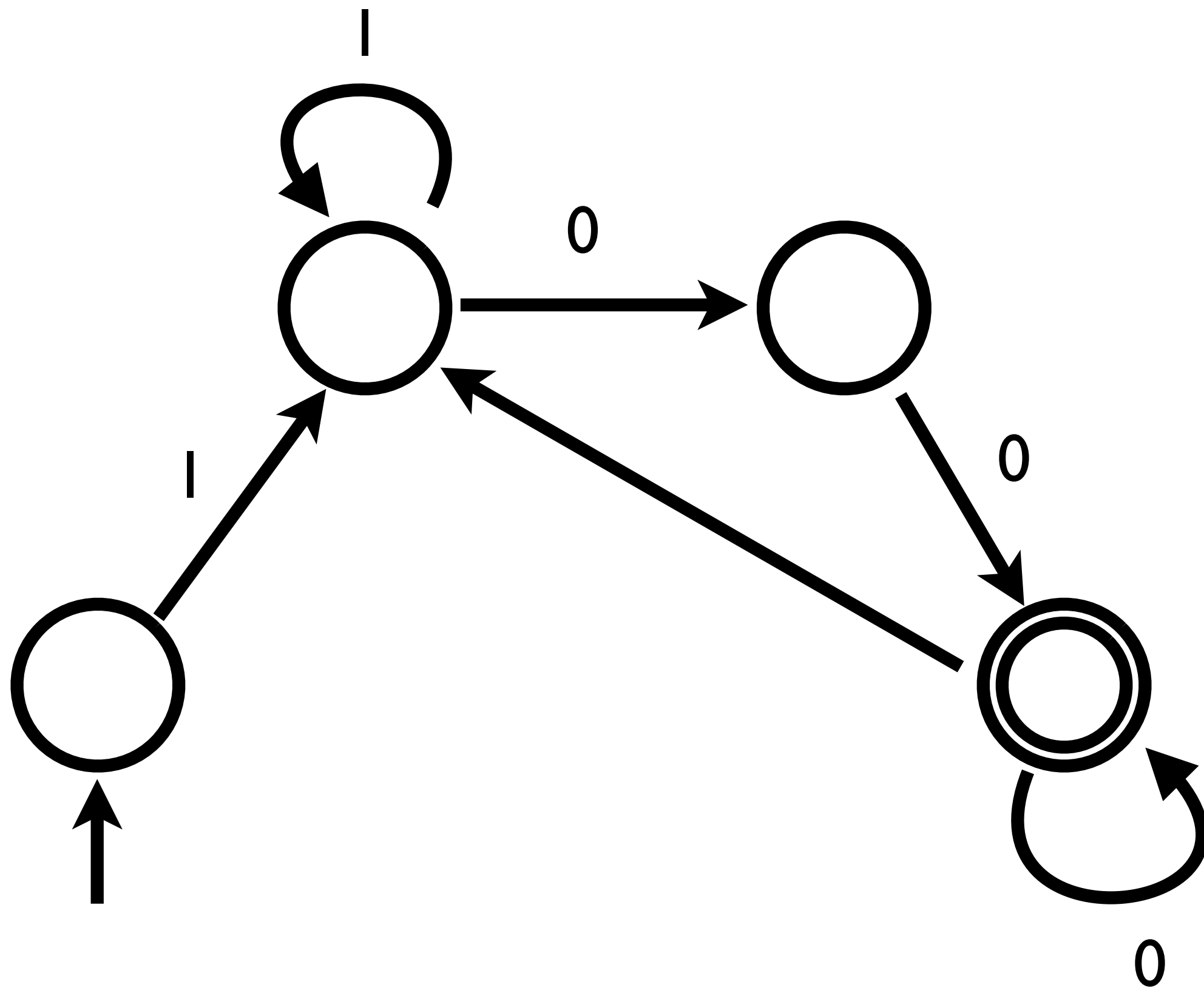
- The set A is an alphabet.
- The set Q is a set of control states.
- The control state q_0 is the initial control state.
- The function $\delta : Q \times A \rightarrow Q$ determines the next state.
- The set F determines the final (accepting) states.











$$\mathcal{L}(A, Q, q_0, \delta, F) \subseteq A^*$$

If P_1 and \dots and P_n , then Q .

$$P_1 \quad \cdot \quad \cdot \quad \cdot \quad P_n$$

$$Q$$

$$\frac{P_1 \quad \cdot \quad \cdot \quad \cdot \quad P_n}{Q}$$

If P_1 and \dots and P_n , then Q .

$$\frac{q_0 \in F}{\epsilon \in \mathcal{L}(A, Q, q_0, \delta, F)}$$

$$\frac{w \in \mathcal{L}(A, Q, \delta(q_0, c), \delta, F)}{cw \in \mathcal{L}(A, Q, q_0, \delta, F)}$$

Example 4.2. The automaton (A, Q, q_0, δ, F) , where:

$$A = \{0\}$$

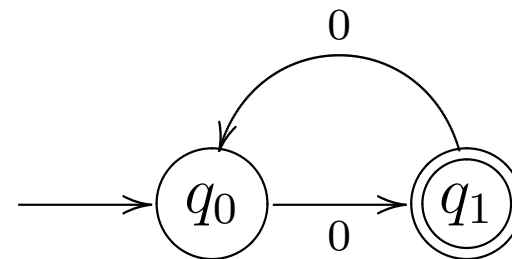
$$Q = \{q_0, q_1\}$$

$$F = \{q_1\},$$

$$\delta(q_0, 0) = q_1$$

$$\delta(q_1, 0) = q_0$$

is depicted graphically as:



and it accepts exactly the set of strings with odd length— $\{0, 000, 00000, \dots\}$.

Formally, a **nondeterministic automaton** M is also a 5-tuple (A, Q, q_0, δ, F) where:

- The set A is an alphabet.
- The set Q is a set of control states.
- The control state q_0 is the initial control state.
- The function $\delta : Q \times (A \cup \{\epsilon\}) \rightarrow \mathcal{P}(Q)$ determines the next state.
- The set F determines the final (accepting) states.

$$\delta : Q \times A \rightarrow Q$$

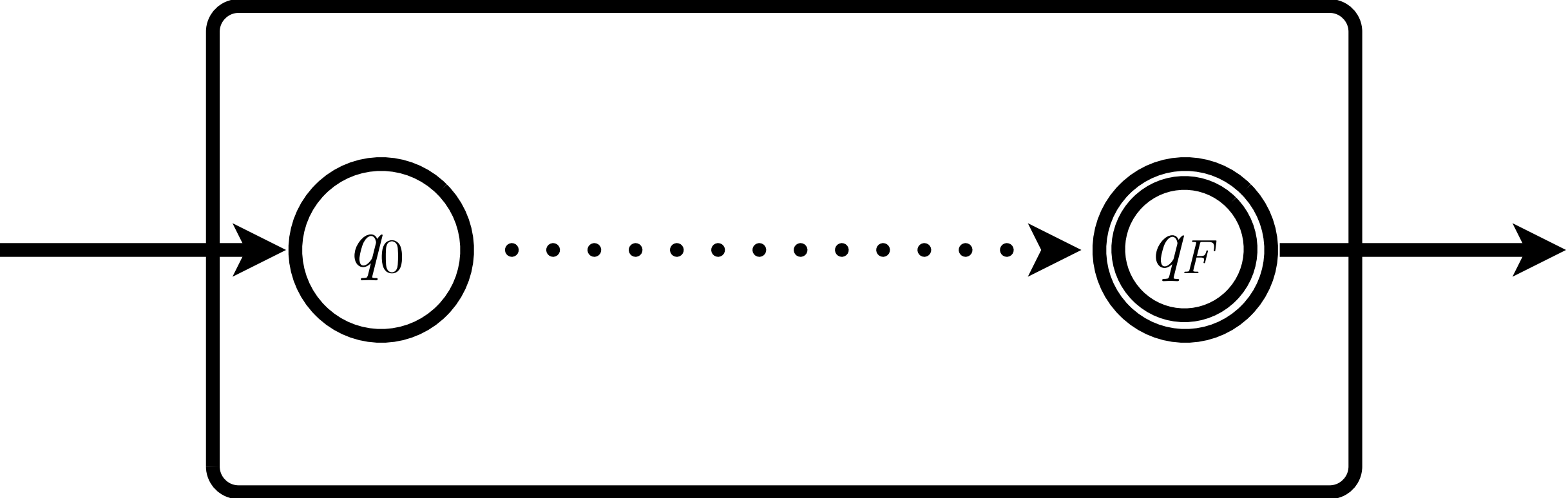
$$\delta : Q \times (A \cup \{\epsilon\}) \rightarrow \mathcal{P}(Q)$$

$$\frac{q_0 \in F}{\epsilon \in \mathcal{L}(A, Q, q_0, \delta, F)}$$

$$\frac{w \in \mathcal{L}(A, Q, q', \delta, F) \quad q' \in \delta(q_0, c)}{cw \in \mathcal{L}(A, Q, q_0, \delta, F)}$$

$$\frac{w \in \mathcal{L}(A, Q, q', \delta, F) \quad q' \in \delta(q_0, \epsilon)}{w \in \mathcal{L}(A, Q, q_0, \delta, F)}$$

RegEx \Rightarrow NFA



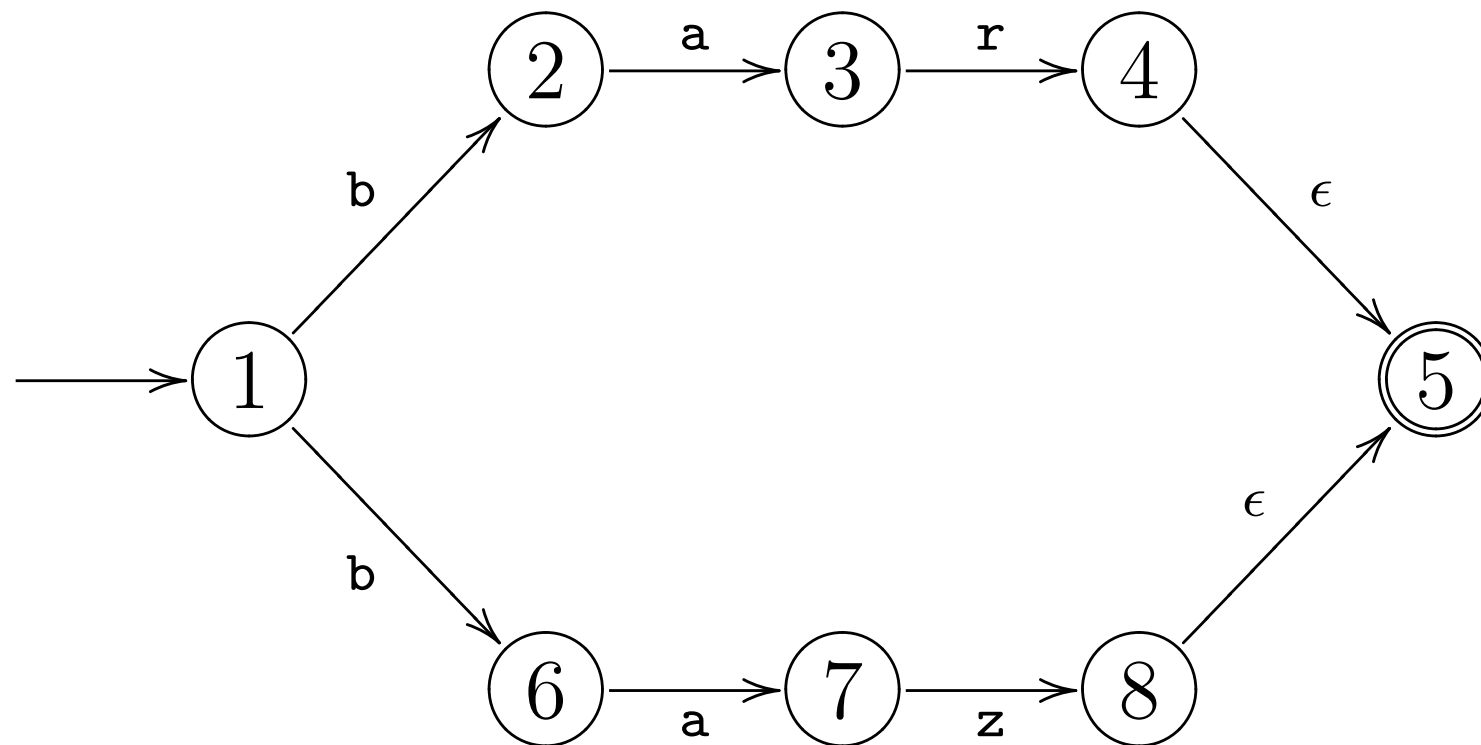
RegEx \Rightarrow NFA

- Empty language
- Empty-string singleton
- One-character singleton
- Concatenation
- Union
- Repetition

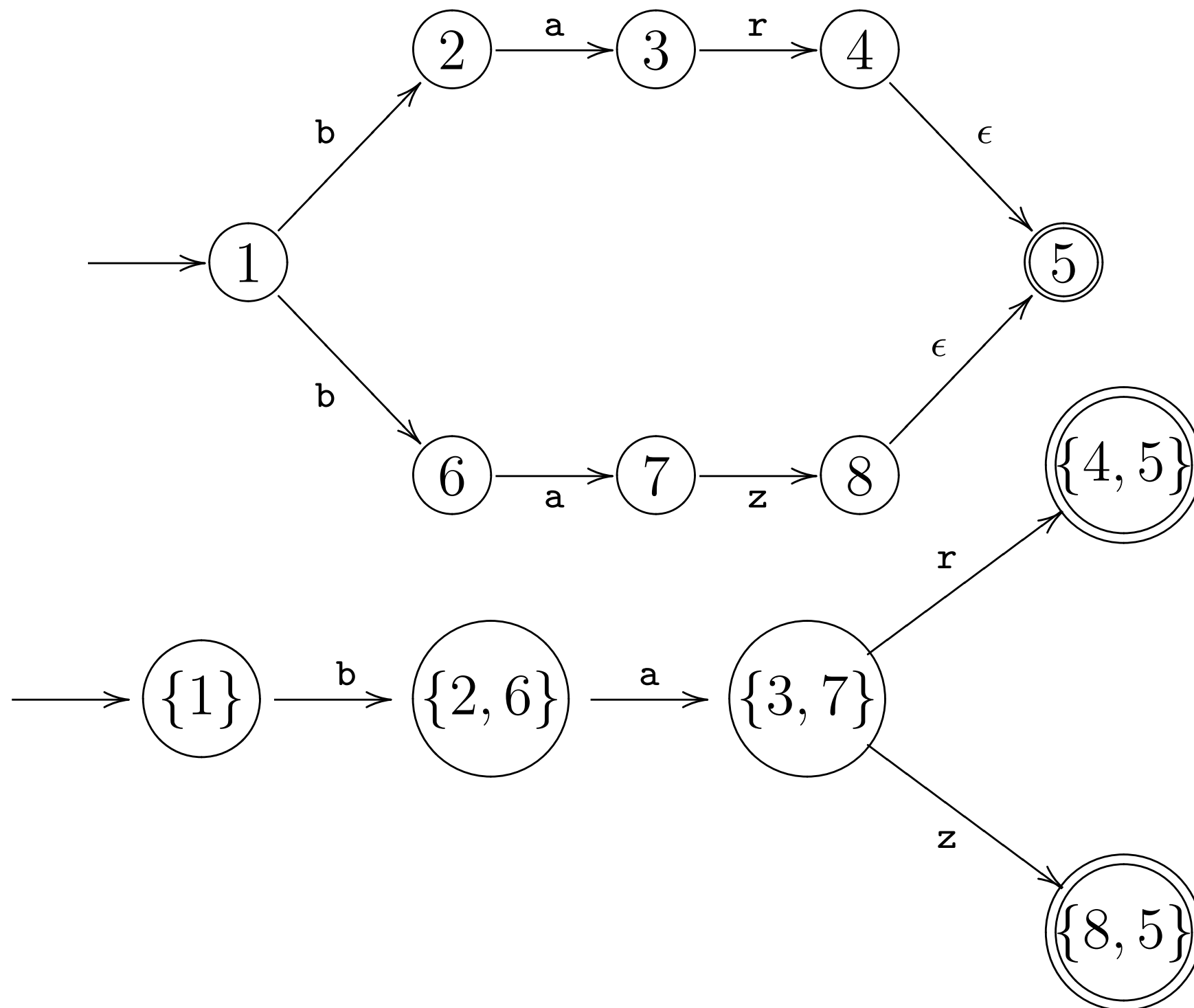
NFA \Rightarrow DFA

Thompson subset construction.

Example: NFA \Rightarrow DFA



Example: NFA \Rightarrow DFA



RegEx \Rightarrow NFA \Rightarrow DFA

lex

flex

lexical spec \Rightarrow lexer

- Identifiers match $[A-Za-z_][A-Za-z0-9_]*$
- Delimiters match $[(\);]$
- Operators match $[-*+/\wedge=]$
- Integers match $-?([1-9][0-9]*|0)$
- Whitespace is ignored
- Comments match $^\#[^\backslash n]*$

How lex works

foo.1

%%

[0-9]+ return 1;

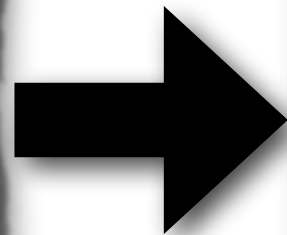
[A-Z]+ return 2;

[\t\n] {}

%%

foo.l

```
%%  
[0-9]+    return 1;  
[A-Z]+    return 2;  
[ \t\n]   {}  
%%
```



lex.yy.c

```
int yylex() {  
    ...  
}
```

defs, opts, decs

%%

rules

%%

C code

Definitions

name regex

{name}

digit [0-9]

int {digit}+

alpha [A-Za-z_]

alnum {alpha}|{digit}

id {alpha}{alnum}*

Options

%option case-insensitive

%option yylineno

%option noyywrap

Declarations

%{

C code

%}

`%s` *state*₁ *state*₂ . . .

`%x` *state*₁ *state*₂ . . .

Rules

pattern { action }

<state> pattern { action }

If in state *state*,
and *pattern* matches the longest prefix,
then perform *action*.

%s state

%%

pat action

<state> pat action

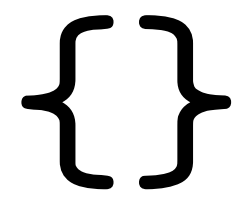
$\%x$ *state*

$\%\%$

pat action

$\langle state \rangle$ *pat action*

Actions



```
return token_type ;
```

BEGIN(*state*);

ECHO ;

REJECT;

Example

- `<MAIN> “(“ { return(LPAR) ; }`
- `<MAIN> “)” { return(RPAR) ; }`
- `<MAIN> “/*” { BEGIN(COMMENT) ; }`
- `<COMMENT> “*/” { BEGIN(MAIN) ; }`
- `<COMMENT> . { }`

Matching length

- Suppose the pattern: $f(oo)^*$
- How many ways can match:
fooooooooooooooooooooo

Options

- Shortest match
- Longest match

Longest match

Given a string w and a set of regular expressions R ,
which regex can match the longest prefix of w ?

Algorithm

```
NaïveRemoveLongestMatch( $w \in A^*, R \subseteq 2^{A^*}$ )  
  suffix  $\leftarrow w$   
  while ( $w \neq \varepsilon$ )  
    if  $\exists L \in R : \varepsilon \in L$   
      suffix  $\leftarrow w$   
     $C:W \leftarrow w$   
     $R \leftarrow D_C.R$   
  return suffix
```

Algorithm

RemoveLongestMatch($w \in A^*, R \subseteq 2^{A^*}$)

 suffix $\leftarrow w$

 while ($R \neq \emptyset$ or $w \neq \varepsilon$)

 if $\exists L \in R : \varepsilon \in L$

 suffix $\leftarrow w$

$C:W \leftarrow w$

$R \leftarrow D_C.R - \{\emptyset\}$

 return suffix

Example

- RegEx: fo*
- RegEx: foobar
- RegEx: foob
- Input string: foobarbaz

Examples

```
%{  
    int num_lines = 0, num_chars = 0;  
}%  
  
%%  
\n        ++num_lines; ++num_chars;  
.        ++num_chars;  
%%  
int main(int argc, char* argv[]) {  
    yylex();  
    printf("# of lines = %d, # of chars = %d\n",  
           num_lines, num_chars );  
}
```

Questions?