

Compiler Construction

Lent Term 2021

Lecture 2 : Lexical analysis

- Recall regular expressions
- Recall Finite Automata
- Recall NFA to DFA transformation
- What is the “lexing problem”?
- How DFAs are used to solve the lexing problem?

Timothy G. Griffin

tgg22@cam.ac.uk

Computer Laboratory
University of Cambridge

What problem are we solving?

Translate a sequence of characters

```
if m = 0 then 1 else if m = 1 then 1 else fib (m - 1) + fib (m - 2)
```

into a sequence of **tokens**

```
IF, IDENT "m", EQUAL, INT 0, THEN, INT 1, ELSE, IF,  
IDENT "m", EQUAL, INT 1, THEN, INT 1, ELSE, IDENT "fib",  
LPAREN, IDENT "m", SUB, INT 1, RPAREN, ADD,  
IDENT "fib", LPAREN, IDENT "m", SUB, INT 2, RPAREN
```

implemented with some data type

```
type token =  
  | INT of int | IDENT of string | LPAREN | RPAREN  
  | ADD | SUB | EQUAL | IF | THEN | ELSE  
  | ...
```

Regular expressions e over alphabet Σ

$$e \rightarrow \phi \mid \varepsilon \mid a \mid e + e \mid ee \mid e^* \quad (a \in \Sigma)$$

$$M(e) \subseteq \Sigma^*$$

$$M(\phi) = \{\}$$

$$M(\varepsilon) = \{\varepsilon\}$$

$$M(a) = \{a\}$$

$$M(e_1 + e_2) = M(e_1) \cup M(e_2)$$

$$M(e_1 e_2) = \{w_1 w_2 \mid w_1 \in M(e_1), w_2 \in M(e_2)\}$$

$$M(e^0) = \{\varepsilon\}$$

$$M(e^{n+1}) = M(ee^n)$$

$$M(e^*) = \bigcup_{n \geq 0} M(e^n)$$

Regular Expression (RE) Examples

$$M((a + b)^* abb) =$$

$$\{abb, aabb, baabb, aaabb, ababb, \\ baabb, bbabb, aaaabb, \dots\}$$

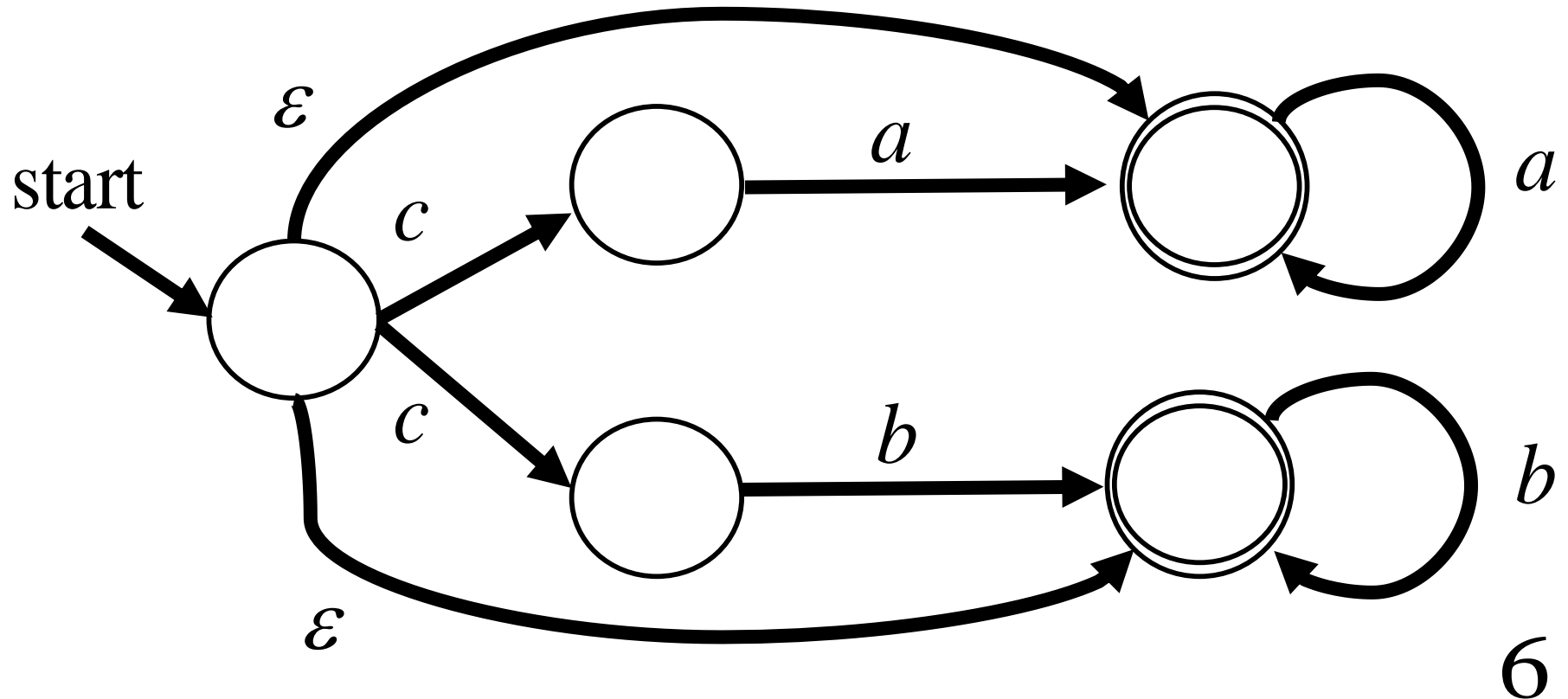
$$M((\square + \otimes)^* \square \otimes \otimes) =$$

$$\{\square \otimes \otimes, \square \square \otimes \otimes, \otimes \square \square \otimes \otimes, \\ \square \square \square \otimes \otimes, \square \otimes \square \otimes \otimes, \otimes \square \square \otimes \otimes, \\ \otimes \otimes \square \otimes \otimes, \square \square \square \square \otimes \otimes, \dots\}$$

NFA Example

An NFA accepting

$$M(a^* + b^* + caa^* + cbb^*)$$



A bit of notation

$$q \xrightarrow{\varepsilon} q$$

For deterministic FA.

$$q_1 \xrightarrow{aw} q_3 \quad \text{if } \delta(q_1, a) = q_2 \text{ and } q_2 \xrightarrow{w} q_3$$

$$L(M) = \{w \mid \exists q \in F, q_0 \xrightarrow{w} q\}$$

$$q \xrightarrow{\varepsilon} q$$

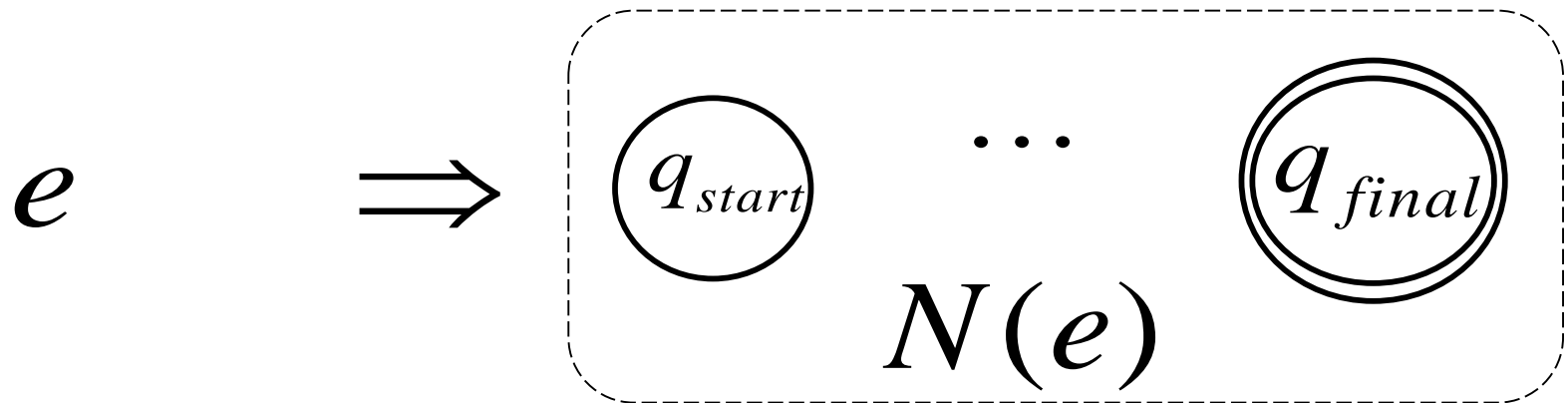
For nondeterministic FA.

$$q_1 \xrightarrow{w} q_3 \quad \text{if } q_2 \in \delta(q_1, \varepsilon) \quad \text{and } q_2 \xrightarrow{w} q_3$$

$$q_1 \xrightarrow{aw} q_3 \quad \text{if } q_2 \in \delta(q_1, a) \quad \text{and } q_2 \xrightarrow{w} q_3$$

$$L(M) = \{w \mid \exists q \in F, q_0 \xrightarrow{w} q\}$$

Review of RE \rightarrow NFA



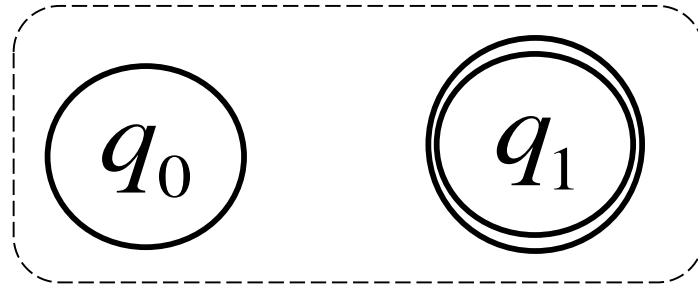
A regular expression.

A nondeterministic FA accepting $M(e)$ with a single final state.

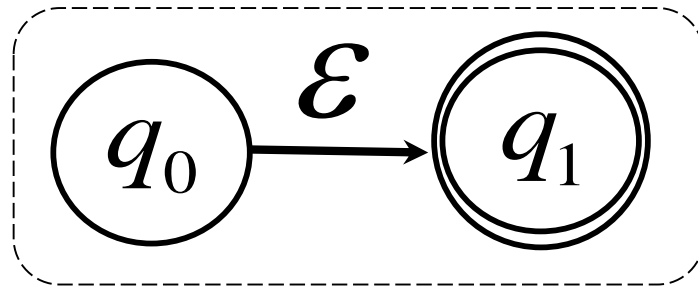
The construction is done by induction on the structure of e .

Review of RE \rightarrow NFA

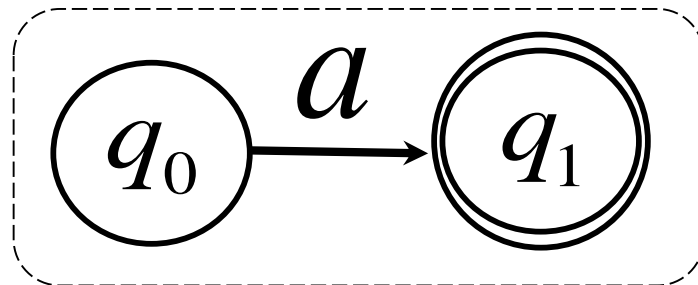
$$N(\phi) =$$



$$N(\varepsilon) =$$

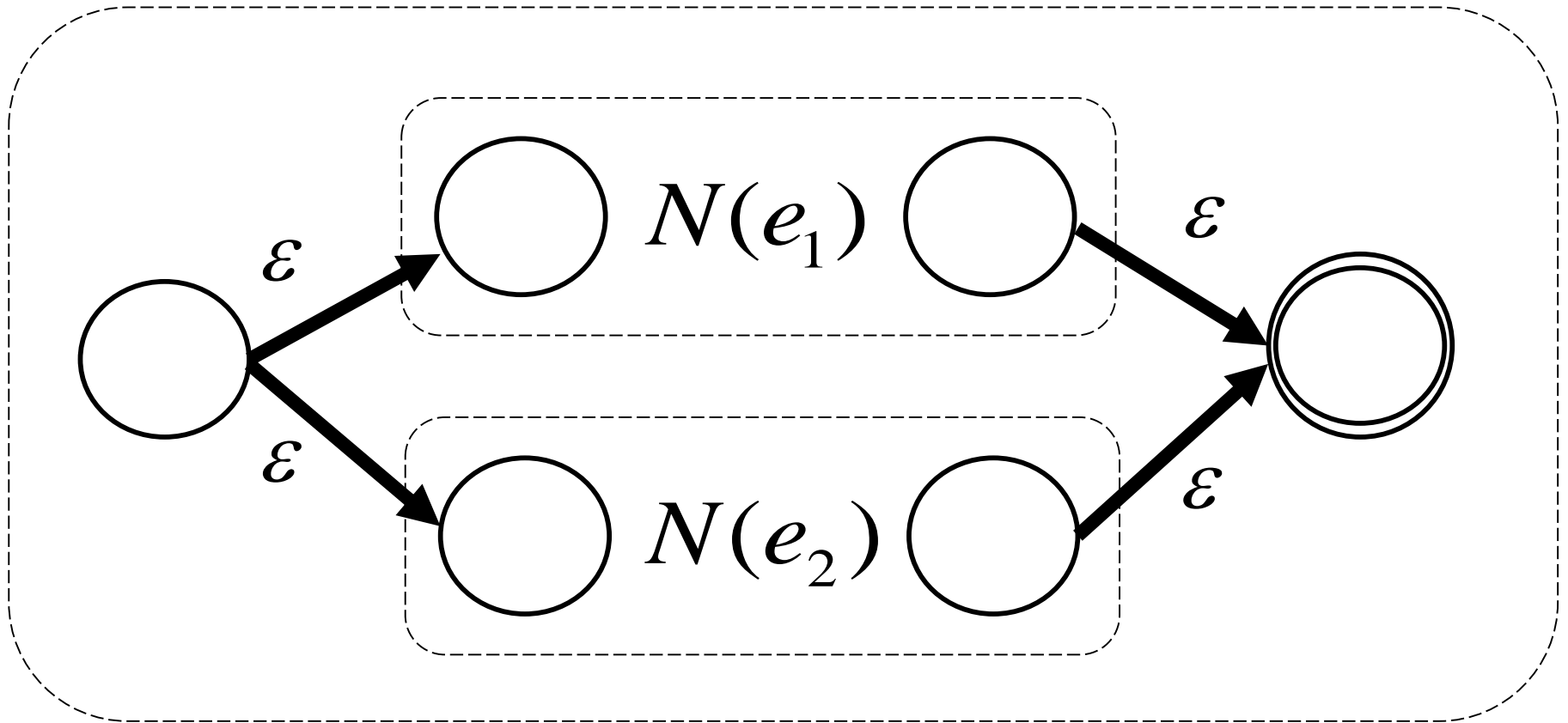


$$N(a) =$$



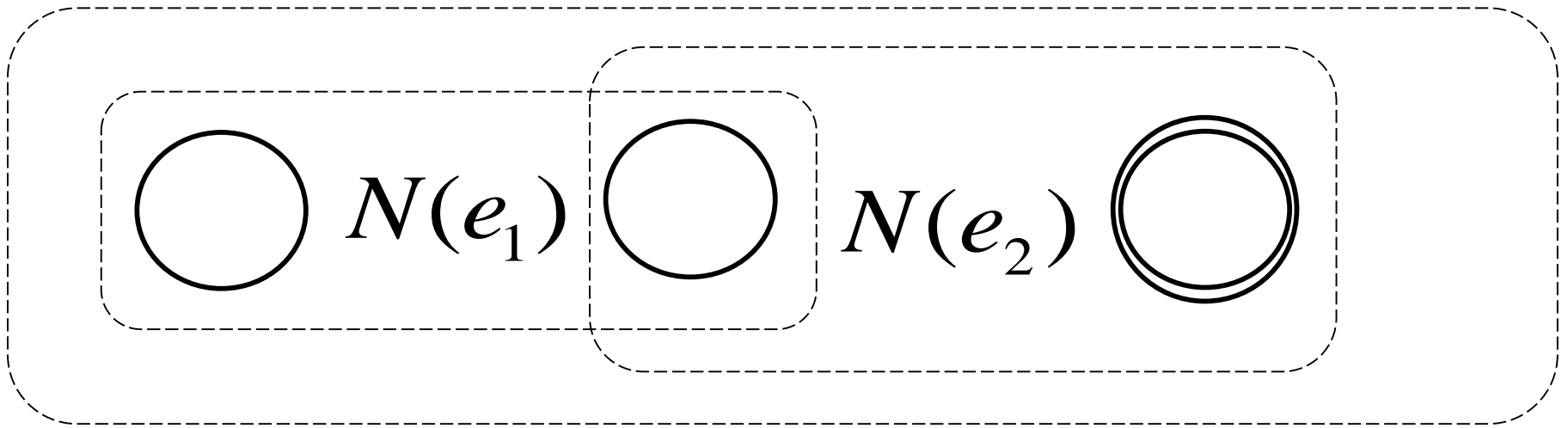
Review of RE \rightarrow NFA

$$N(e_1 + e_2) =$$



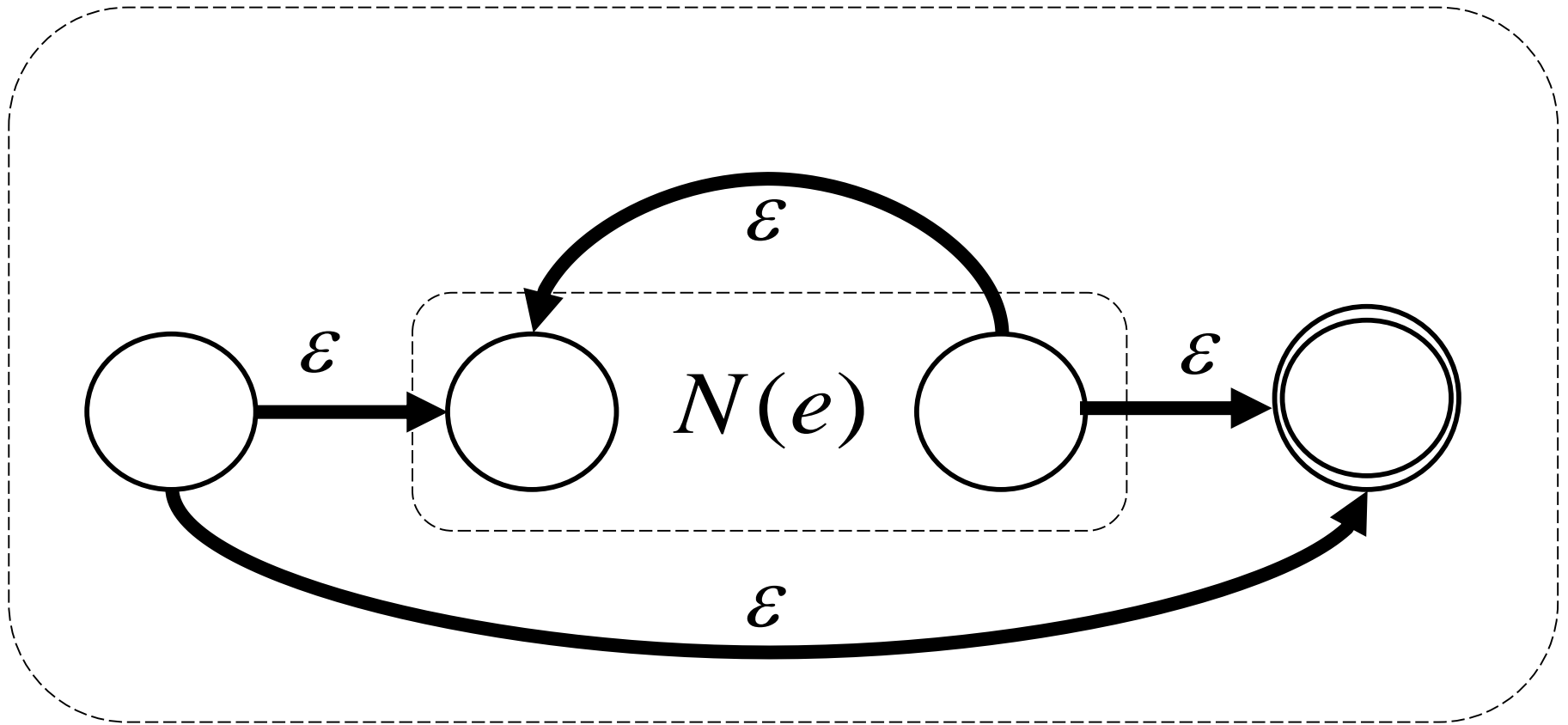
Review of RE \rightarrow NFA

$$N(e_1e_2) =$$

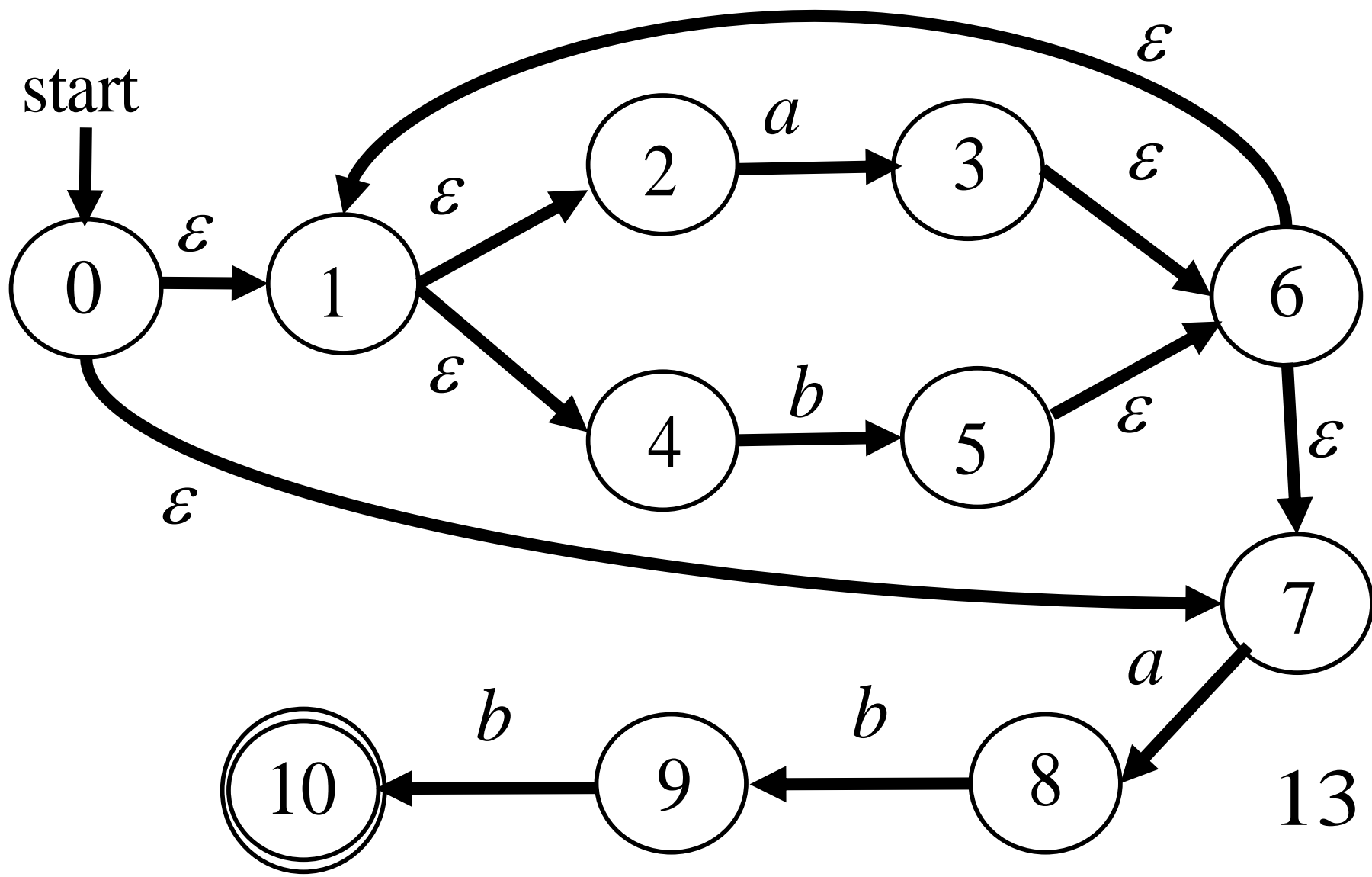


Review of RE \rightarrow NFA

$$N(e^*) =$$



$N((a+b)^*abb)$



Review of NFA -> DFA

$$M = (Q, \Sigma, \delta, q_0, F)$$

$$M' = (Q', \Sigma, \delta', q'_0, F')$$

$$Q' = \{S \mid S \subseteq Q\}$$

$$\varepsilon\text{-closure}(S) = \{q' \in Q \mid \exists q \in S, q \xrightarrow{\varepsilon} q'\}$$

$$\delta'(S, a) = \varepsilon\text{-closure}(\{q' \in \delta(q, a) \mid q \in S\})$$

$$q'_0 = \varepsilon\text{-closure}\{q_0\}$$

$$F' = \{S \subseteq Q \mid S \cap F \neq \emptyset\}$$

How do we compute ε -closure(S)?

ε -closure(S):

push all elements of S onto a stack

result := S

while stack not empty

pop q off the stack

for each $u \in \delta(q, \varepsilon)$

if $u \notin$ result

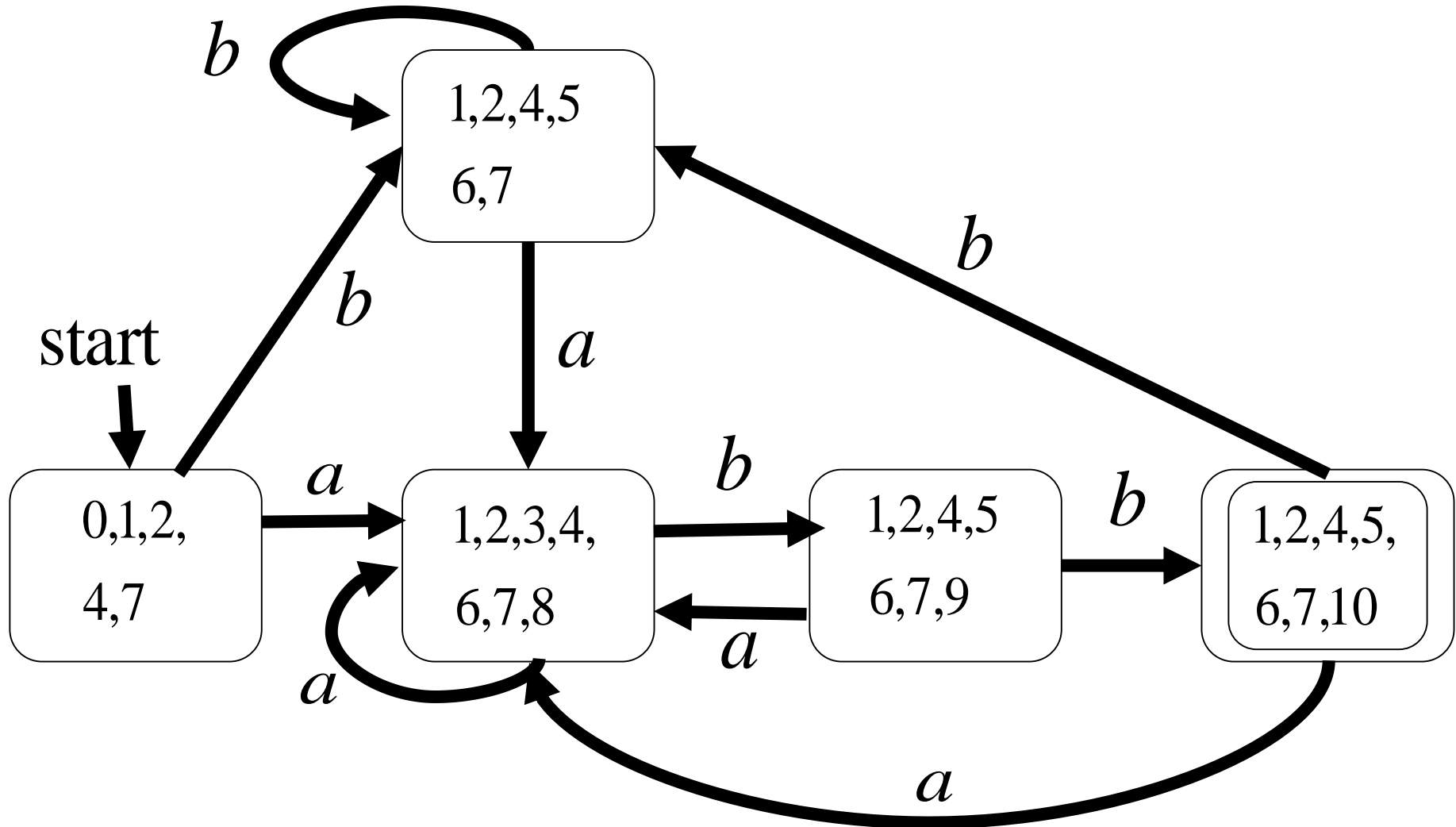
then result := $\{u\} \cup$ result

push u on stack

return result

Look familiar?
It's just a version of
transitive closure!

$DFA(N((a+b)^*abb))$



Traditional Regular Language Problem

Given e and w , is $w \in L(e)$?

Solution : construct NFA from e , then DFA, then run the DFA on w .

But is this a solution to the “lexing problem?”

No!

Something closer to the “lexing problem”

Given $e_1, e_2 \dots, e_k$ and W

find $(i_1, w_1), (i_2, w_2), \dots, (i_n, w_n)$ **so that**

$W = w_1 w_2 \dots w_n$ **and** $\forall i \exists j w_i \in L(e_{i_j})$

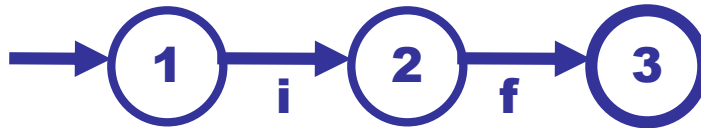
and what else?

The expressions are **ordered** by priority. Why?
Is “if” a variable or a keyword? Need priority to resolve ambiguity (so “if” matched keyword RE before identifier RE).

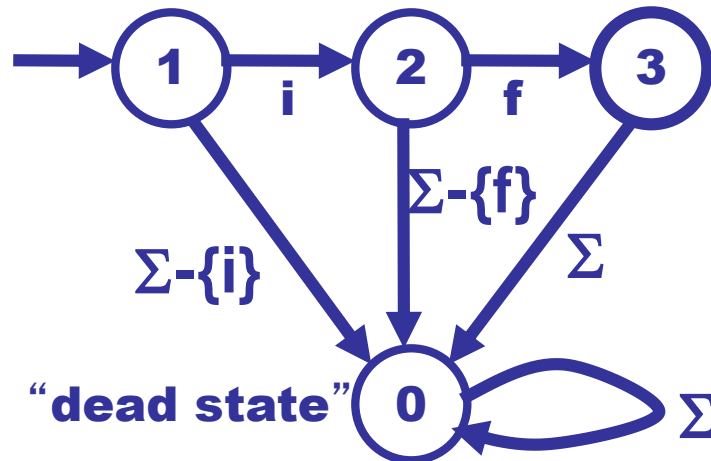
We need to do a **longest** match. Why?
Is “ifif” a variable or two “if” keywords?

Define Tokens with Regular Expressions (Finite Automata)

Keyword: if



This FA is really shorthand for:



Define Tokens with Regular Expressions (Finite Automata)

Regular Expression	Finite Automata	Token
Keyword: if	<p>A finite automata with three states: 1, 2, and 3. State 1 is the start state, indicated by an incoming arrow from the left. Transitions are: 1 to 2 on 'i', and 2 to 3 on 'f'. State 3 is the final state, indicated by a double circle.</p>	KEY(IF)
Keyword: then	<p>A finite automata with five states: 1, 2, 3, 4, and 5. State 1 is the start state, indicated by an incoming arrow from the left. Transitions are: 1 to 2 on 't', 2 to 3 on 'h', 3 to 4 on 'e', and 4 to 5 on 'n'. State 5 is the final state, indicated by a double circle.</p>	KEY(then)
Identifier: [a-zA-Z][a-zA-Z0-9]*	<p>A finite automata with two states: 1 and 2. State 1 is the start state, indicated by an incoming arrow from the left. Transitions are: 1 to 2 on '[a-zA-Z]', and a self-loop on state 2 labeled '[a-zA-Z0-9]*'. State 2 is the final state, indicated by a double circle.</p>	ID(s)

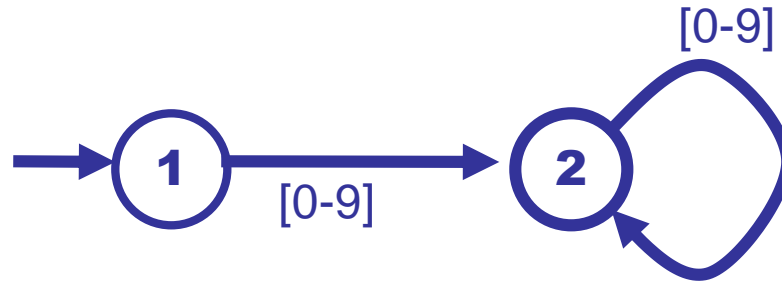
Define Tokens with Regular Expressions (Finite Automata)

Regular Expression

Finite Automata

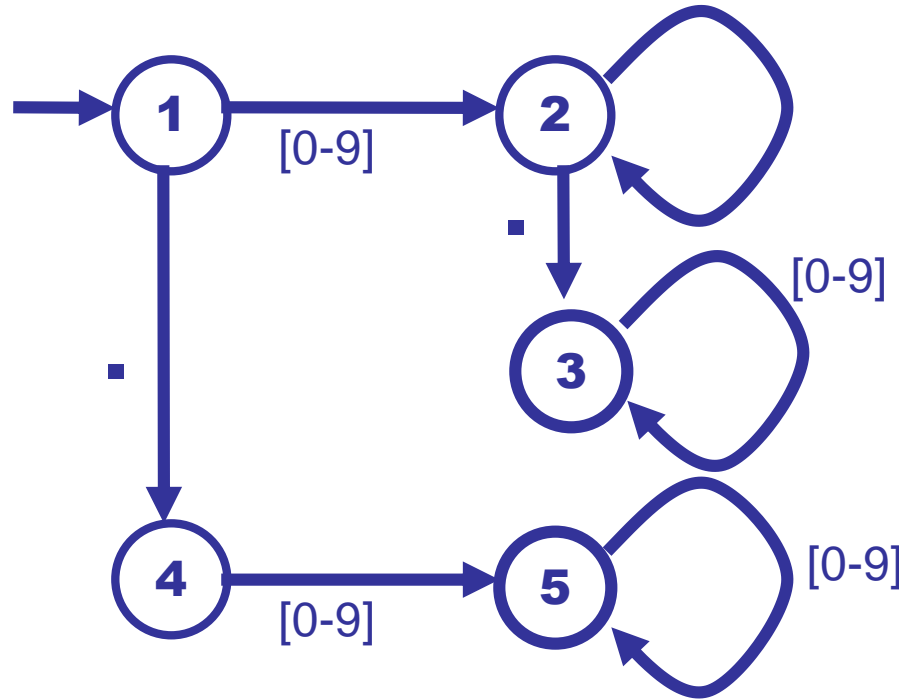
Token

number:
[0-9][0-9]*



NUM(n)

real:
([0-9]+ '.' [0-9]*)
| ([0-9]* '.' [0-9]+)



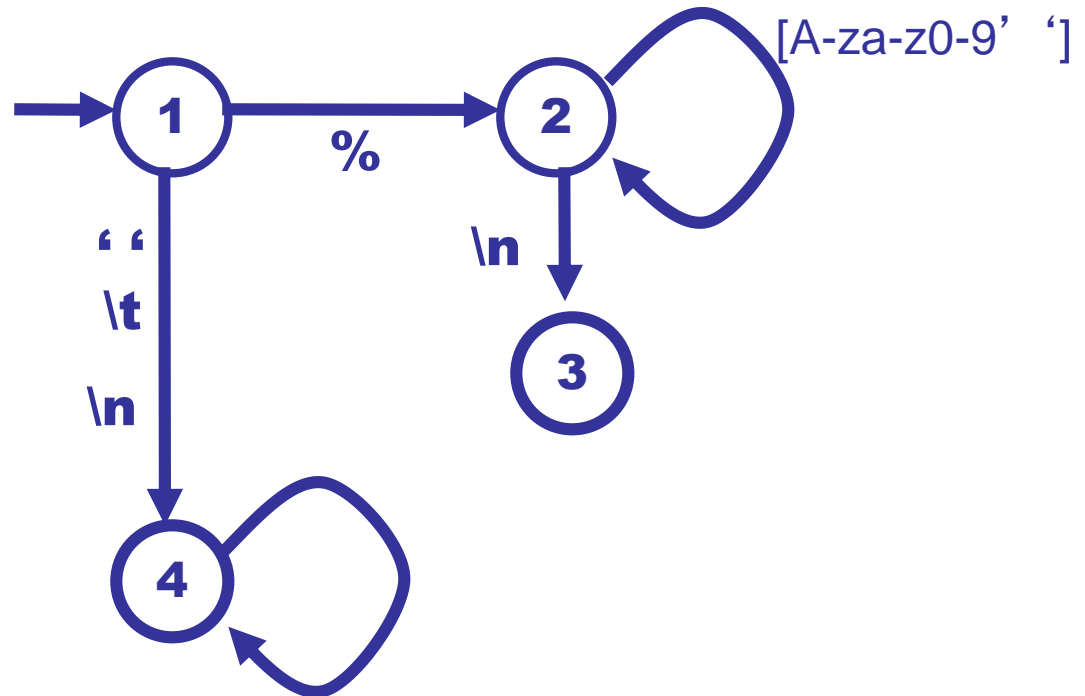
NUM(n)

[0-9]

21

No Tokens for “White-Space”

**White-space with one line
comments starting with %**



Constructing a Lexer

INPUT: $e_1, e_2 \dots, e_k$

an **ordered** list of regular expressions
Highest priority first, lowest last



NFA for $e = e_1 + e_2 + \dots + e_k$



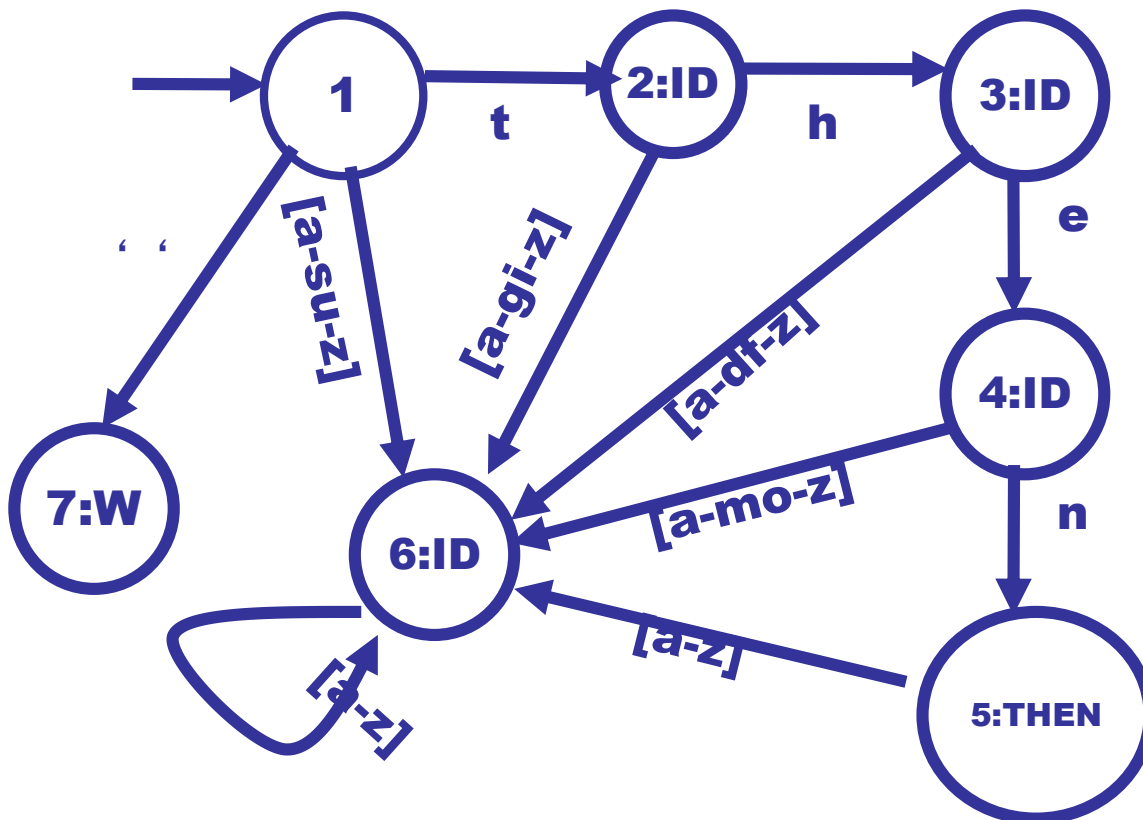
DFA with each final state
associated with the e_i of
highest priority.

Constructing a Lexer

(1) Keyword : then

(2) Ident : [a-z][a-z]*

(2) White-space: ' '



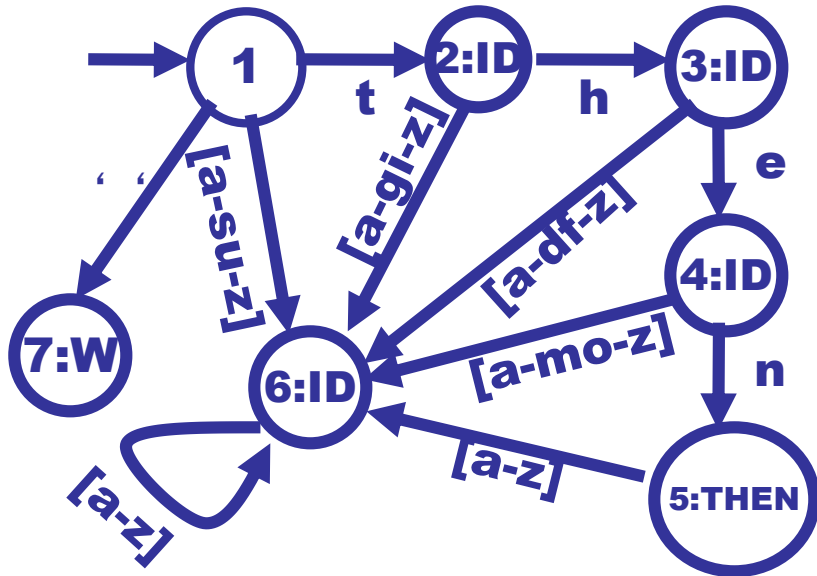
State 5 could accept either an ID or the keyword “then”. The priority rules eliminates this ambiguity and associates state 5 with the keyword.

What about longest match?

Start in initial state,

Repeat:

- (1) read input until dead state is reached. Emit token associated with last accepting state.
- (2) reset state to start state



| = current position, \$ = EOF

Input	current state	last accepting state	
then thenx\$	1	0	
t hen thenx\$	2	2	
th en thenx\$	3	3	
the n thenx\$	4	4	
then thenx\$	5	5	
then thenx\$	0	5	EMIT KEY(THEN)
then thenx\$	1	0	RESET
then thenx\$	7	7	
then t henx\$	0	7	EMIT WHITE(' ')
then thenx\$	1	0	RESET
then t henx\$	2	2	
then th enx\$	3	3	
then the nx\$	4	4	
then then x\$	5	5	
then thenx \$	6	6	
then thenx \$	0	6	EMIT ID(thenx)