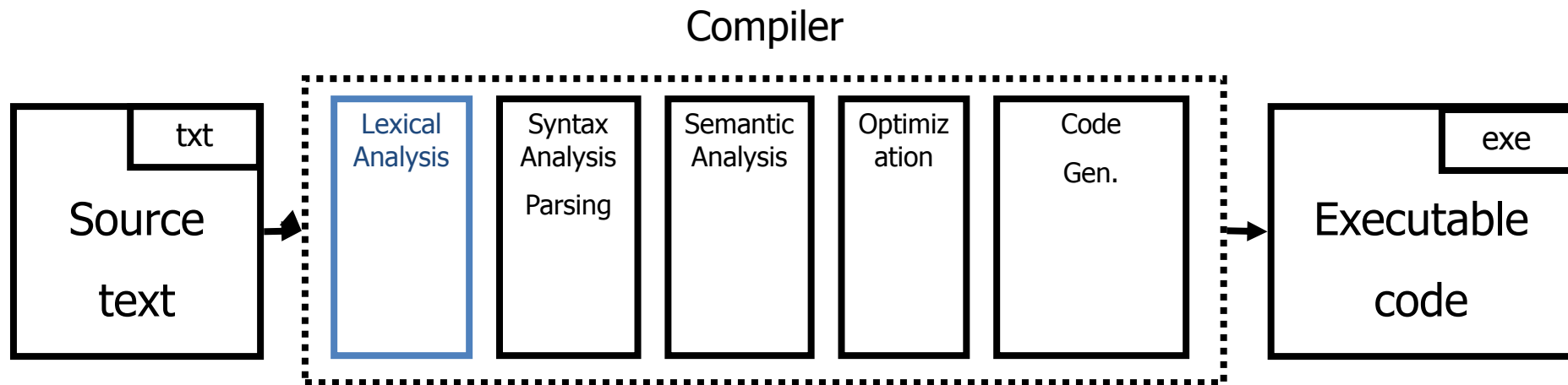


Introduction to Compilation

Lexical Analysis

You are here



Role of lexical analyzer

- Recognize tokens and ignore white spaces, comments

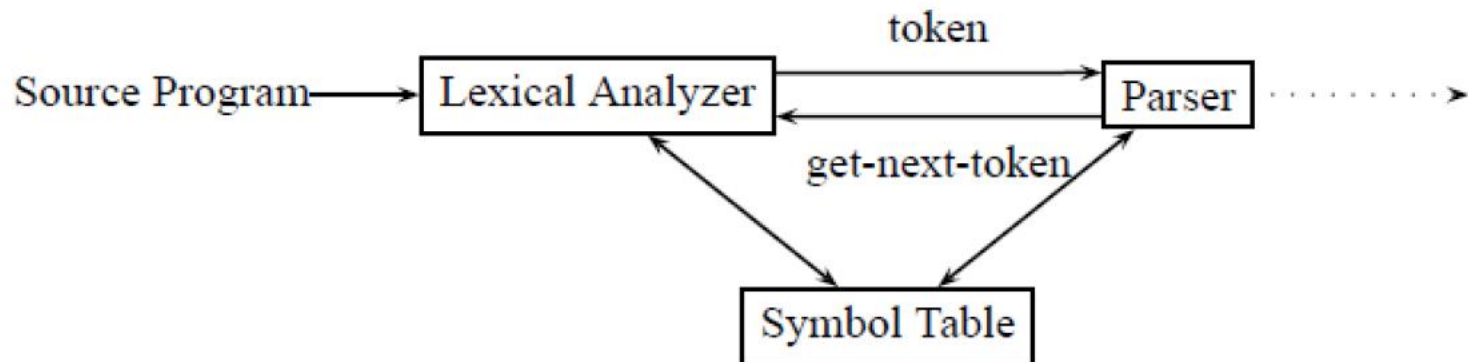
i	f		(x	1		*	x	2	<	1	.	0)	{
---	---	--	---	---	---	--	---	---	---	---	---	---	---	---	---

- Generates token stream

if	(x1	*	x2	<	1.0)	{
----	---	----	---	----	---	-----	---	---

- Error reporting

Diagram is here



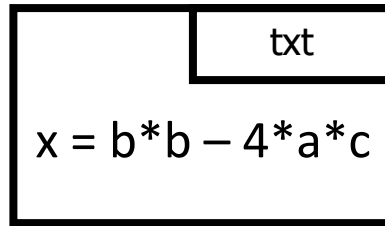
From characters to tokens

- What is a token?
 - Roughly – a “word” in the source language
 - Identifiers
 - Values
 - Language keywords
 - **Really - anything that should appear in the input to syntax analysis**
- Technically
 - Usually a pair of (type,value)

Tokens vs. lexemes

Token type	Example lexemes
Identifier	x, y, z, foo, bar
NUM	42
FLOATNUM	3.141592654
STRING	“so long, and thanks for all the fish”
LPAREN	(
RPAREN)
IF	if
...	

From characters to tokens



Token
Stream

<ID,"x"> <EQ> <ID,"b"> <MULT> <ID,"b"> <MINUS> <INT,4> <MULT> <ID,"a"> <MULT> <ID,"c">

Tokens – using symbol table

Example. Let us consider the following assignment statement:

$$E := M * C * 2$$

then the following pairs $\langle \text{token}, \text{attribute} \rangle$ are passed to the Parser:

$\langle \text{id}, \text{pointer to symbol-table entry for } E \rangle$

$\langle \text{assign-op}, \rangle$

$\langle \text{id}, \text{pointer to symbol-table entry for } M \rangle$

$\langle \text{mult-op}, \rangle$

$\langle \text{id}, \text{pointer to symbol-table entry for } C \rangle$

$\langle \text{exp-op}, \rangle$

$\langle \text{num}, \text{integer value } 2 \rangle$.

- Some Tokens have a null attribute: the Token is sufficient to identify the Lexeme.
- From an implementation point of view, each token is encoded as an integer number.

Errors in lexical analysis

	txt
pi = 3.141.562	



Illegal token

	txt
pi = 3oranges	



Illegal token

	txt
pi = oranges3	



<ID,"pi">, <EQ>, <ID,"oranges3">

Error Handling

- Many errors cannot be identified at this stage.
- Example: “fi (a==f(x))”. Should “fi” be “if”? Or is it a routine name?
 - We will discover this later in the analysis.
 - At this point, we just create an identifier token.

How can we define tokens?

- Keywords – easy!
 - **if, then, else, for, while, ...**
- Identifiers?
- Numerical Values?
- Strings?
- Characterize **unbounded sets of values** using a **bounded description**?

Regular Expressions

Basic Patterns	Matching
x	A single letter 'x' from the alphabet
.	Any character, usually except a new line
[xyz]	Any of the characters x,y,z
Repetition Operators	
R?	An R or nothing (ϵ R)
R*	Zero or more occurrences of R
R+	One or more occurrences of R (RR^*)
Composition Operators	
R1R2	An R1 followed by R2
R1 R2	Either an R1 or R2
Grouping	
(R)	R itself

Examples

- $ab^* | cd? =$
- $(a | b)^+ =$
- $(0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9)^* =$

Escape characters

- What is the expression for one or more + symbols?
 - `(+)+` won't work
 - `(\+)+` will
- backslash `\` before an operator turns it to standard character
- `*`, `\?`, `\+`, ...

Shorthands

- Use names for expressions
 - letter = a | b | ... | z | A | B | ... | Z
 - letter_ = letter | _
 - digit = 0 | 1 | 2 | ... | 9
 - id = letter_ (letter_ | digit)*
- Use “-” to denote a range
 - letter = a-z | A-Z
 - digit = 0-9

Examples

- `if = if`
- `then = then`
- `relop = < | > | <= | >= | = | <>`
- `digit = 0-9`
- `digits = digit+`

Example

- A number is:

$$\begin{aligned} \text{number} = & (0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)^+ \\ & (\varepsilon \mid . (0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)^+) \\ & (\varepsilon \mid E (0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)^+) \end{aligned}$$

- Using shorthands it can be written as (this time with negative exponent as well):

$$\begin{aligned} \text{digit} & \rightarrow 0 \mid 1 \mid \dots \mid 9 \\ \text{digits} & \rightarrow \text{digit}^+ \\ \text{optional-fraction} & \rightarrow (. \text{digits})? \\ \text{optional-exponent} & \rightarrow (E(+ \mid -)? \text{digits})? \\ \text{num} & \rightarrow \text{digits optional-fraction optional-exponent} \end{aligned}$$

Ambiguity

- if = if
 - id = letter_ (letter_ | digit)*
- “if” is a valid word in the language of identifiers... so what should it be?
- How about the identifier “iffy”?
- Solution
 - **Always find longest matching token**
 - Break ties using **order of definitions**... first definition wins (= > list rules for keywords before identifiers)

Creating a lexical analyzer

- Input
 - List of token definitions (pattern name, regex)
 - String to be analyzed
- Output
 - List of tokens
- How do we build an analyzer?

Character classification

```
#define is_end_of_input(ch) ((ch) == '\0');  
#define is_uc_letter(ch) ('A' <= (ch) && (ch) <= 'Z')  
#define is_lc_letter(ch) ('a' <= (ch) && (ch) <= 'z')  
#define is_letter(ch) (is_uc_letter(ch) || is_lc_letter(ch))  
#define is_digit(ch) ('0' <= (ch) && (ch) <= '9')  
...
```

Main reading routine

```
Token get_next_token() {  
do {  
    char c = getchar();  
    switch(c) {  
        case is_letter(c) : return recognize_identifier(c);  
        case is_digit(c) : return recognize_number(c);  
        ...  
    } while (c != EOF);  
}
```

But we have a much better way!

- Generate a lexical analyzer **automatically** from token definitions
- Main idea
 - Use finite-state automata to match regular expressions

Overview

- Construct a nondeterministic finite-state automaton (NFA) from regular expression (automatically)
- Determinize the NFA into a deterministic finite-state automaton (DFA)
- DFA can be directly used to identify tokens

Reminder: Finite-State Automaton

- **Deterministic** automaton
- $M = (\Sigma, Q, \delta, q_0, F)$
 - Σ - alphabet
 - Q – finite set of state
 - $q_0 \in Q$ – initial state
 - $F \subseteq Q$ – final states
 - $\delta : Q \times \Sigma \rightarrow Q$ - transition function

Reminder: Finite-State Automaton

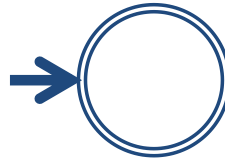
- **Non-Deterministic** automaton
- $M = (\Sigma, Q, \delta, q_0, F)$
 - Σ - alphabet
 - Q – finite set of state
 - $q_0 \in Q$ – initial state
 - $F \subseteq Q$ – final states
 - $\delta : Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^Q$ - transition function
- Possible ε -transitions
- For a word w , M can reach a number of states or get stuck. If **some state** reached is final, M accepts w .

From regular expressions to NFA

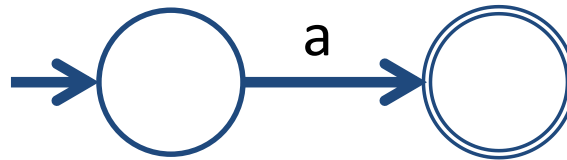
- Step 1: assign expression names and obtain pure regular expressions $R_1 \dots R_m$
- Step 2: construct an NFA M_i for each regular expression R_i
- Step 3: combine all M_i into a single NFA
- Ambiguity resolution: prefer longest accepting word

Basic constructs

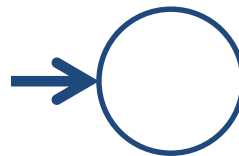
$R = \varepsilon$



$R = a$

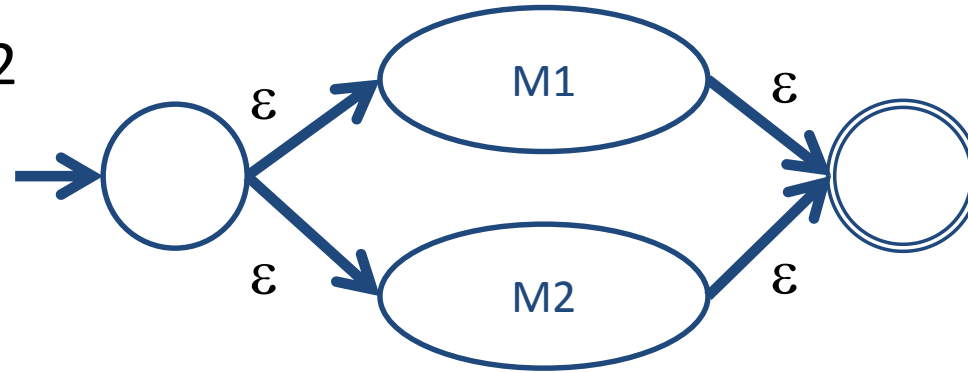


$R = \phi$

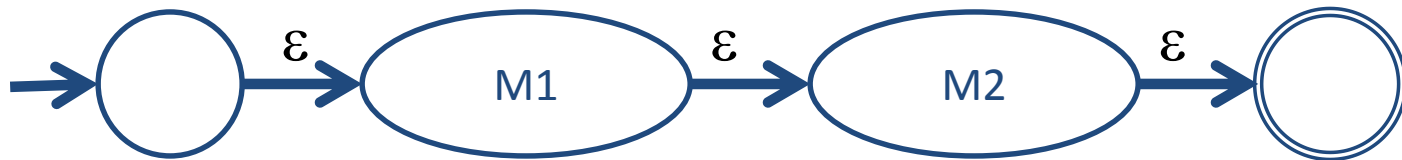


Composition

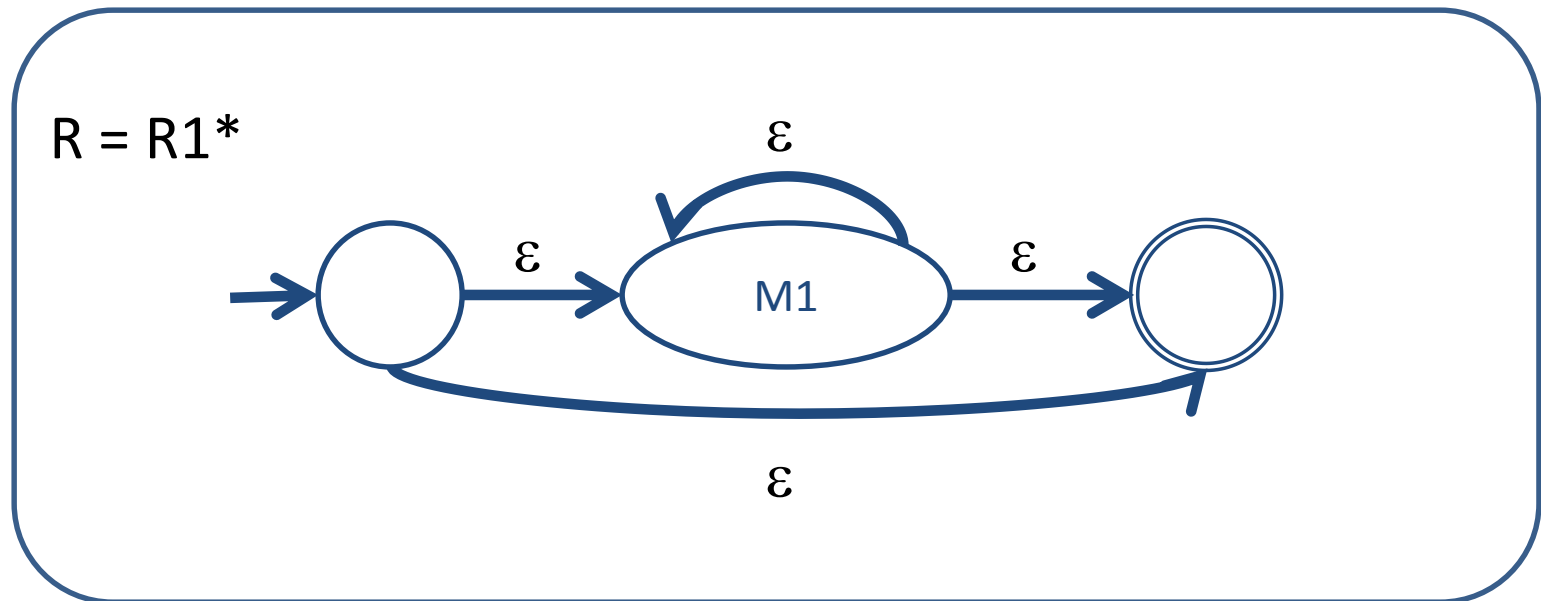
$R = R1 \mid R2$



$R = R1R2$



Repetition



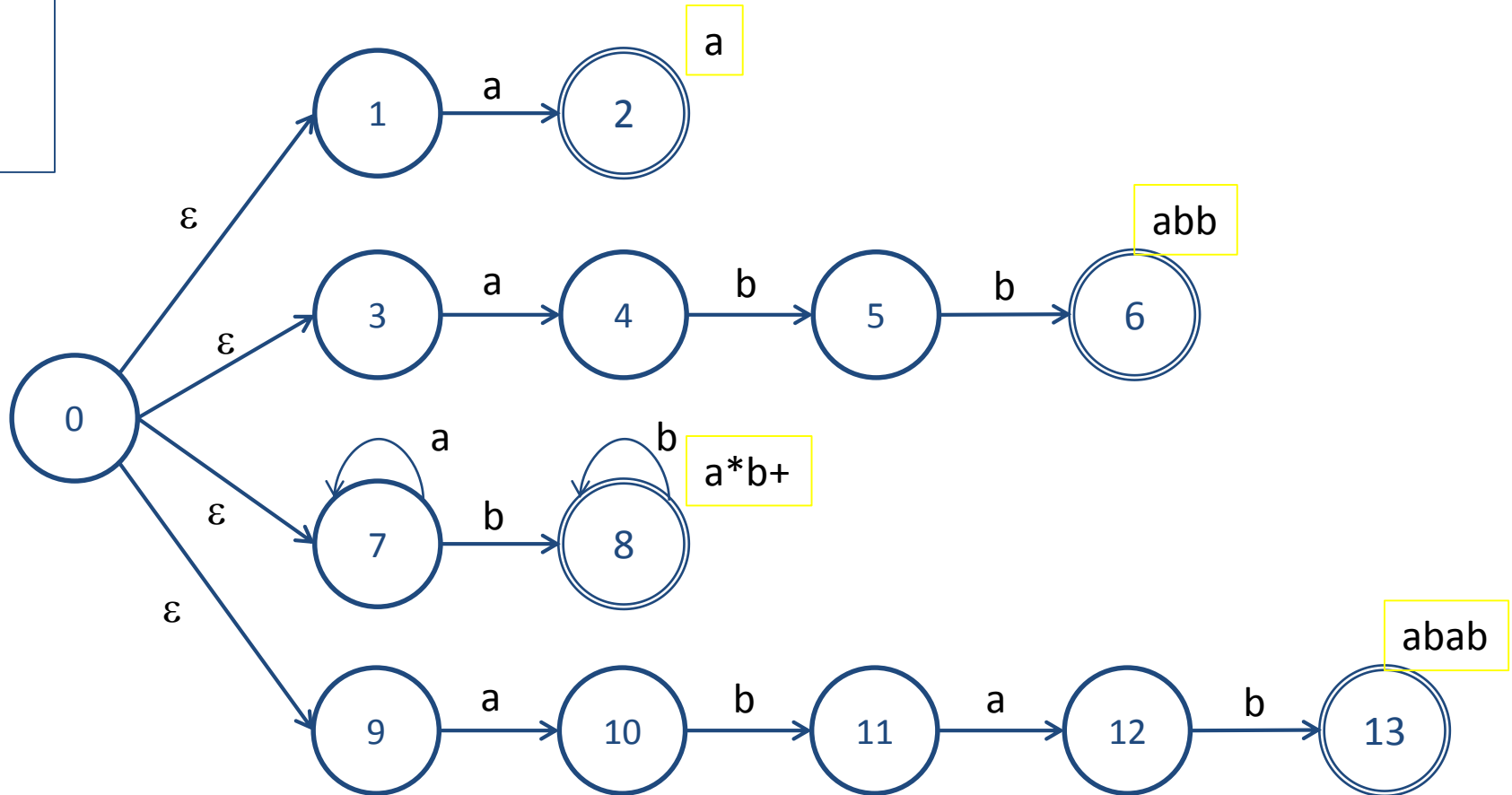
What now?

- Naïve approach: try each automaton separately
- Given a word w :
 - Try $M_1(w)$
 - Try $M_2(w)$
 - ...
 - Try $M_n(w)$
- Requires resetting after every attempt

Combine automata

combines

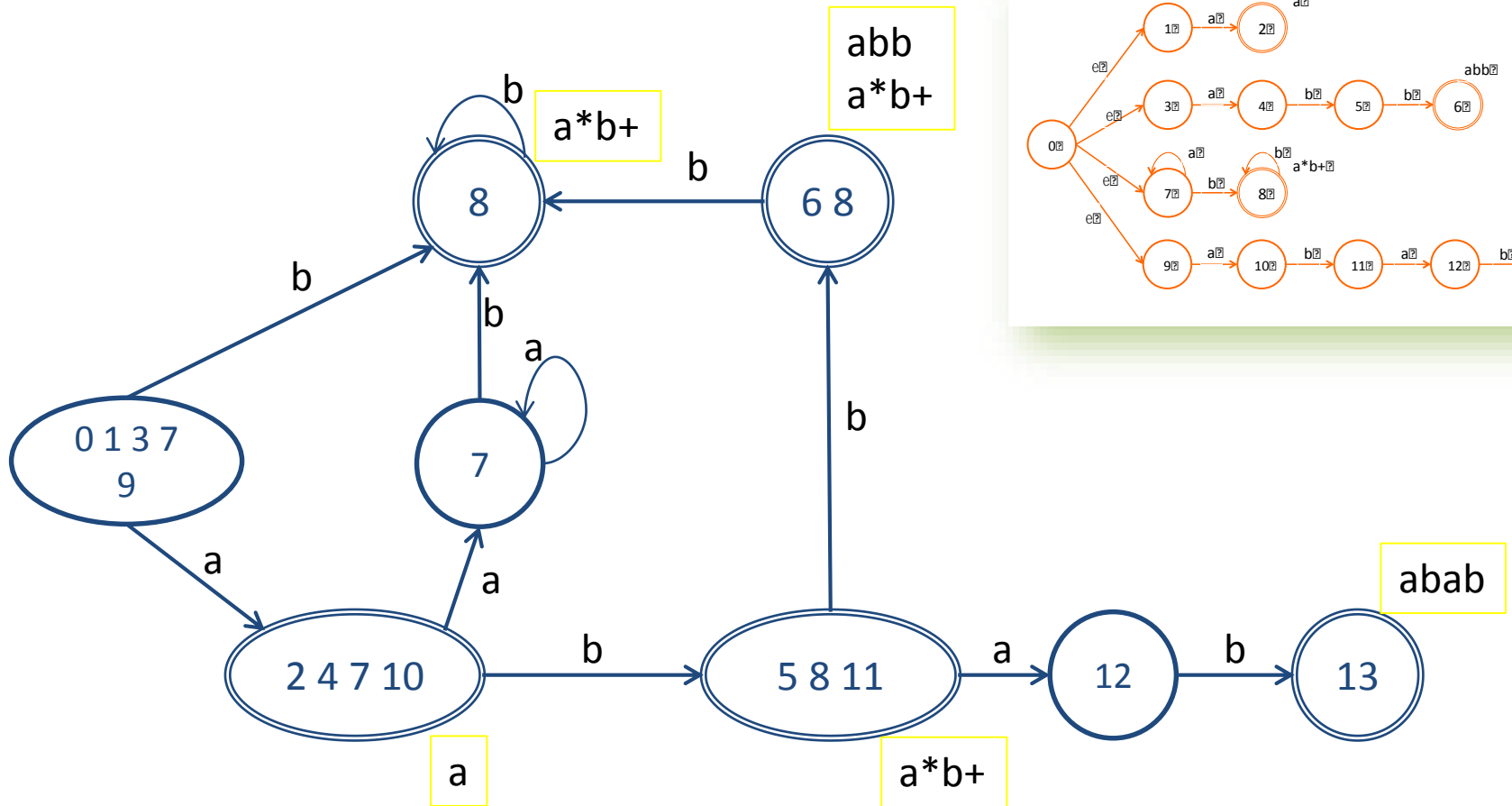
a
abb
 a^+b^+
abab



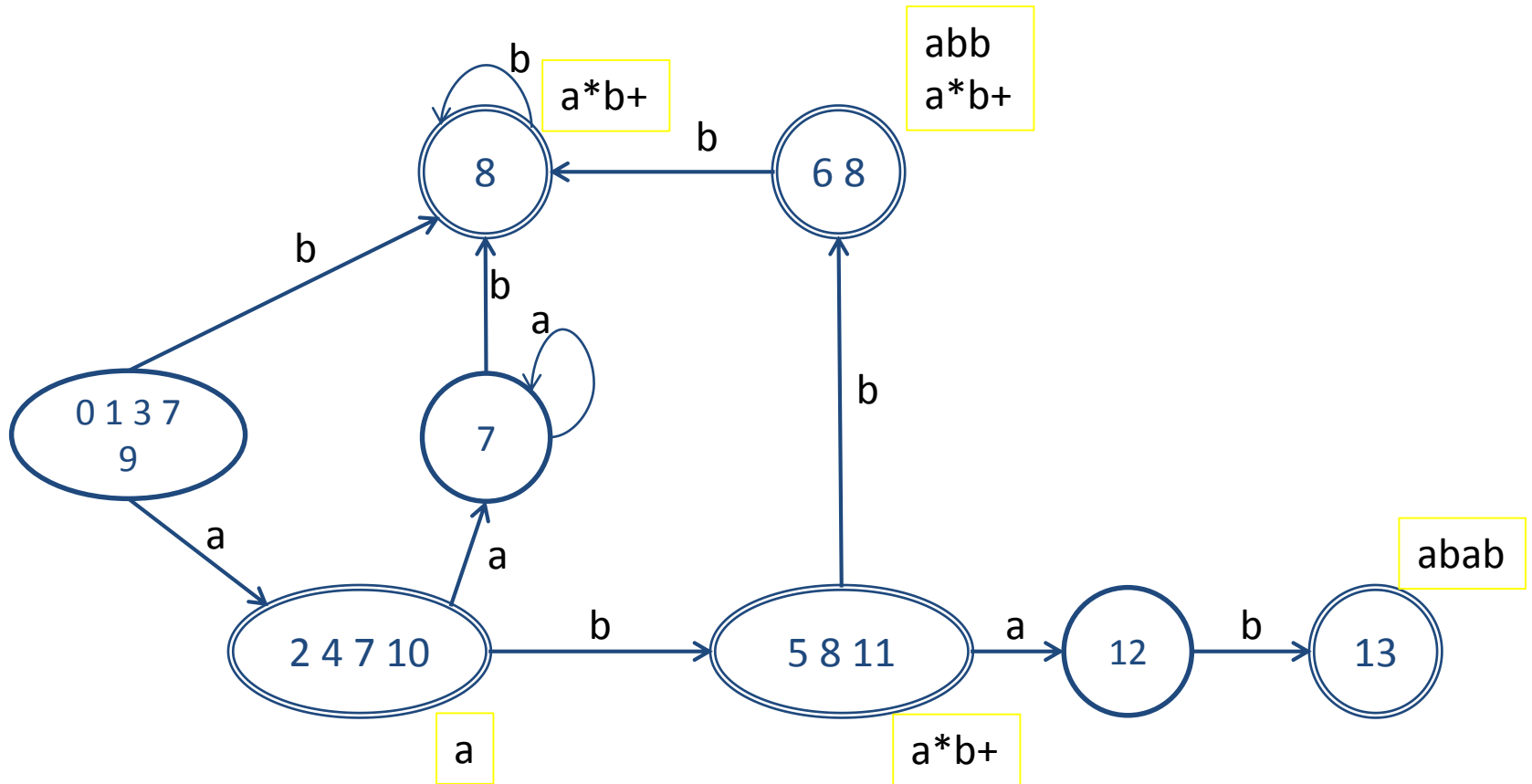
Ambiguity resolution

- Recall...
- Longest word
- Tie-breaker based on **order of rules** when words have same length
- Recipe
 - Turn NFA to DFA
 - **Run until stuck, remember last accepting state, this is the token to be returned**

Corresponding DFA



Examples

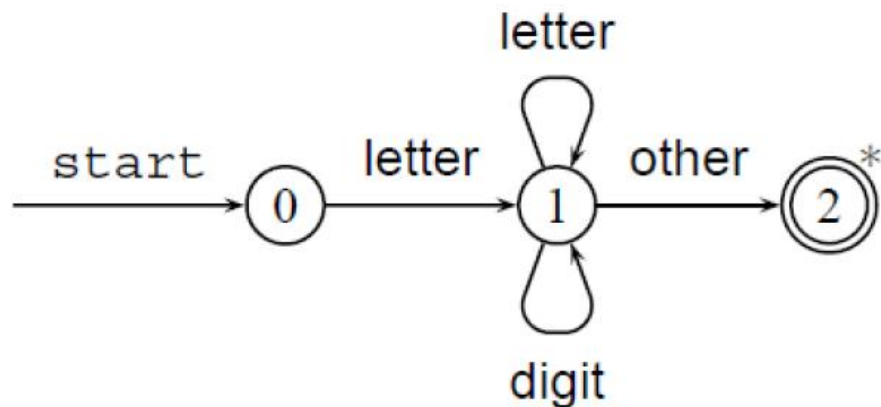


abaa: gets stuck after aba in state 12, backs up to state (5 8 11) pattern is a^*b^+ , token is ab
abba: stops after second b in (6 8), token is abb because it comes first in spec

Maximal lexemes strategy

Solution 1. **Use Automata with *Lookahead*.**

- *Example:* Automaton for id with lookahead.



- The label **other** refers to any character not indicated by any other edges leaving the node;
- The * indicates that we read an extra character and we must retract the *forward* input pointer by one character.

From NFA to DFA

- NFA are hard to simulate with a computer program.
 1. There are many possible paths for a given input string caused by the nondeterminism;
 2. The acceptance condition says that there must be *at least one* path ending with a final state;
 3. We may need to find all the paths before accepting/excluding a string.
- The algorithm to map an NFA to a DFA is called the *Subset Construction*.
- *Main Idea*: Each DFA state corresponds to a set of NFA states, thus encoding all the possible states an NFA could reach after reading an input symbol.
- The algorithm makes use of the operation ϵ -closure: Given a set of states, T , ϵ -closure(T) returns the set of NFA states reachable from some $s \in T$ on ϵ -transition only.

ϵ -closure algorithm

ϵ -closure(T)

CL_T := T; */* Current Closure of T */*

NOT_EX := CL_T; */* Not Examined States */*

repeat */* Updates the Closure of T */*

 NEW_CL := \emptyset ;

for *each* t **in** NOT_EX **do**

$Q := \delta(t, \epsilon)$;

if $Q \not\subseteq \text{CL_T}$ **then**

 NEW_CL := NEW_CL \cup ($Q - \text{CL_T}$);

end

 CL_T := CL_T \cup NEW_CL;

 NOT_EX := NEW_CL;

until NEW_CL == \emptyset

return CL_T */* ϵ -Closure of T */*

NFA to DFA definition

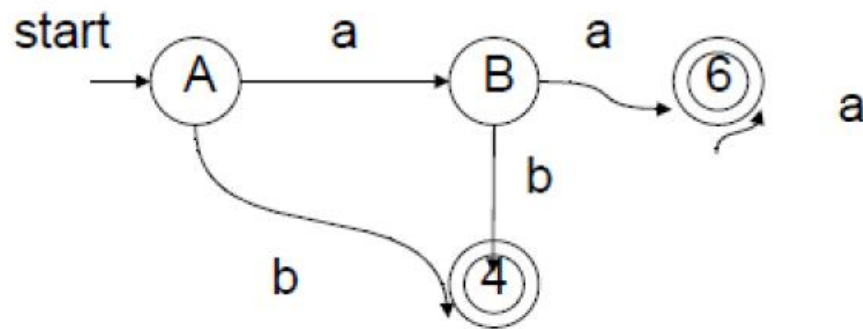
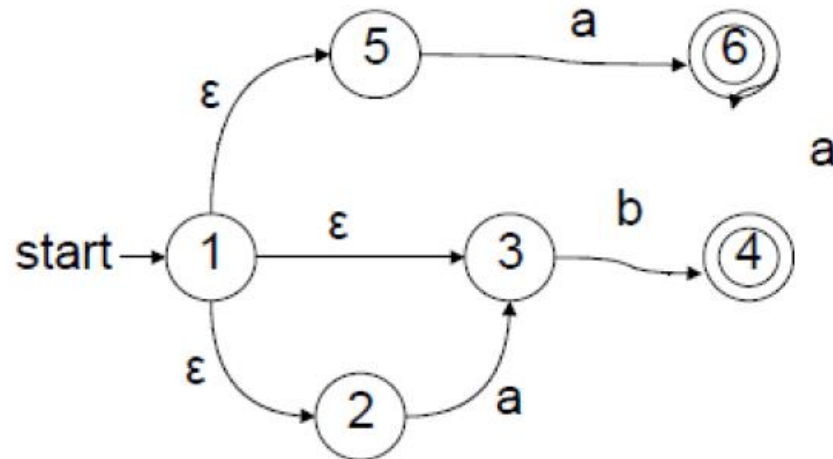
Let $\text{NFA} = (S, V, \delta, s_0, F)$ then the equivalent DFA $= (S', V, \delta', s'_0, F')$ where:

- $S' \subseteq 2^S$. With $[s_1, \dots, s_n]$ we denote an element in S' , which stands for the set of states $\{s_1, \dots, s_n\}$.
- $s'_0 = [\epsilon\text{-closure}(\{s_0\})]$;
- F' is the set of states in S' containing some element in F ;
- $\delta'([s_1, \dots, s_n], a) = [q_1, \dots, q_m]$ iff $\{q_1, \dots, q_m\} = \epsilon\text{-closure}(\delta(\{s_1, \dots, s_n\}, a))$, where $\delta(\{s_1, \dots, s_n\}, a) = \bigcup_{i=1}^n \delta(s_i, a)$.

NFA to DFA – Subset construction algorithm

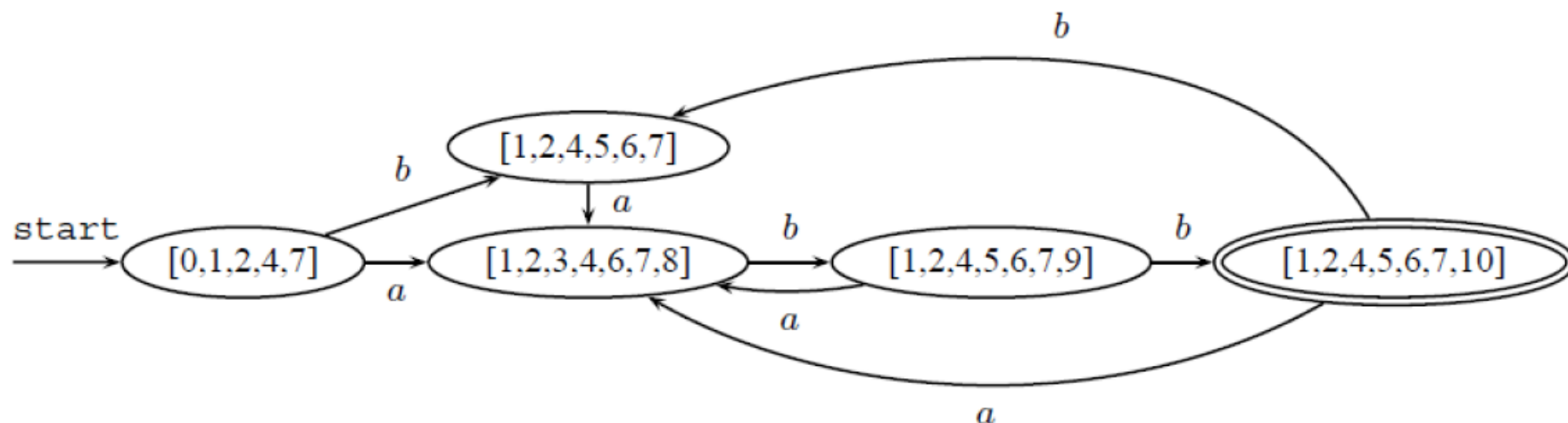
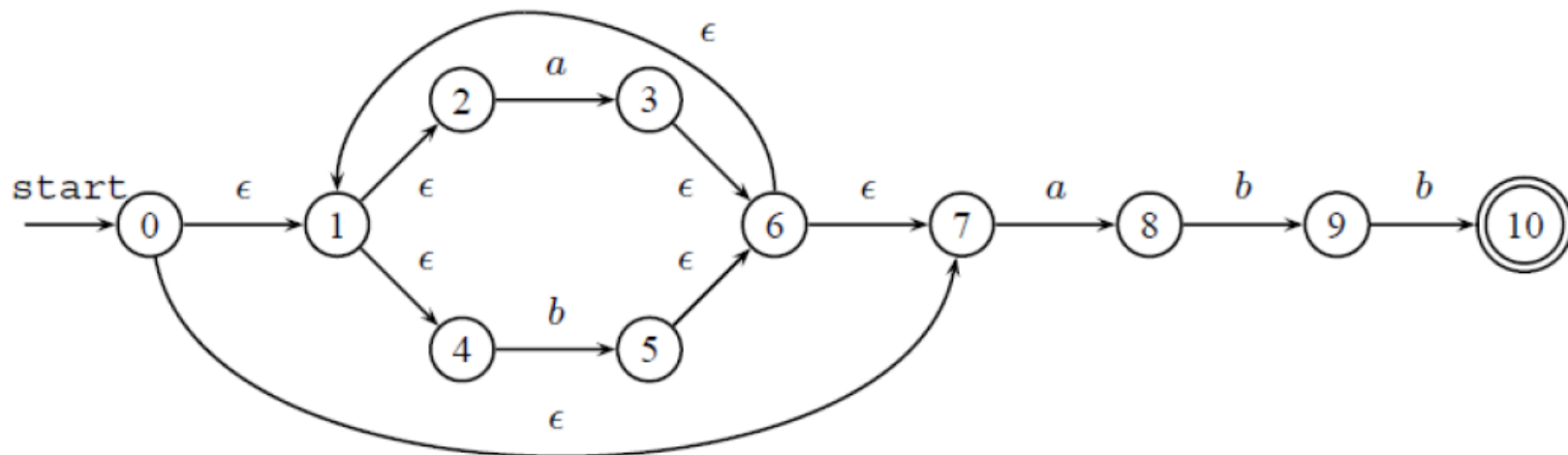
```
Subset_Construction()           /* Subset Construction algorithm */
    DS := { $\epsilon$ -closure( $\{s_0\}$ )};    /* Current Deterministic states */
    NOT_EX := DS;                /* Not Examined States */
    repeat                      /* DFA Construction */
        NEW_DS :=  $\emptyset$ ;
        for each  $T$  in NOT_EX do
            for each symbol  $a \in V$  do
                 $Q := \epsilon$ -closure( $\delta(T, a)$ );
                if  $Q \notin DS$  then NEW_DS := NEW_DS  $\cup$   $Q$ ;
                 $\delta'(T, a) := Q$ ;      /* DFA's Transition Function update */
            end
        end
        DS := DS  $\cup$  NEW_DS;    /* DFA's States update */
        NOT_EX := NEW_DS;
    until NEW_DS ==  $\emptyset$ 
    return DS,  $\delta'$            /* DFA's States and Transition Function */
```


NFA to DFA example



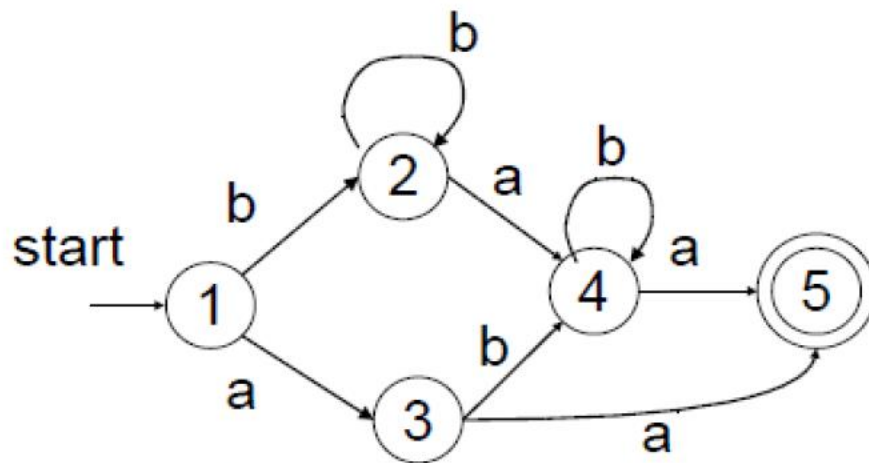
- ϵ -closure(1) = {1, 2, 3, 5}
- Create a new state A = {1, 2, 3, 5} and examine transitions out of it
- $\text{move}(A, a) = \{3, 6\}$
- Call this a new subset state = B = {3, 6}
- $\text{move}(A, b) = \{4\}$
- $\text{move}(B, a) = \{6\}$
- $\text{move}(B, b) = \{4\}$
- Complete by checking $\text{move}(4, a)$; $\text{move}(4, b)$; $\text{move}(6, a)$; $\text{move}(6, b)$

NFA to DFA example 2

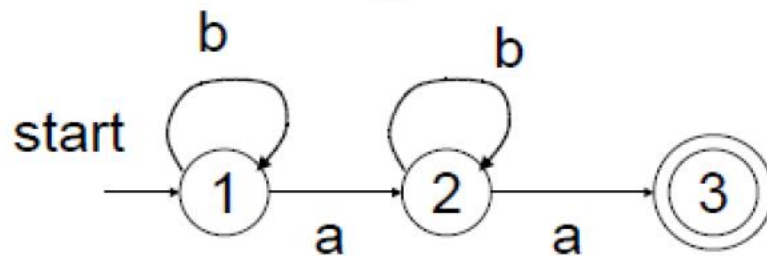


State minimization

- Resulting DFA can be quite large
 - Contains redundant or equivalent states

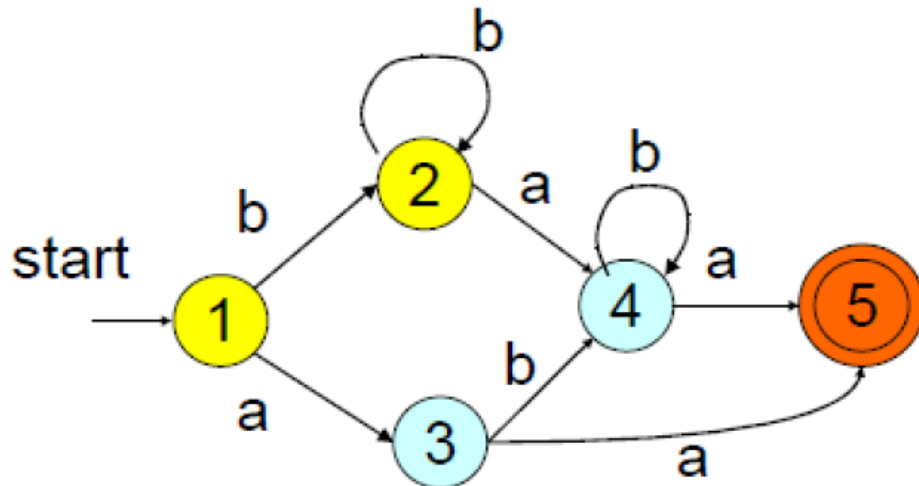


Both DFAs accept b^*ab^*a



State minimization - strategy

- Idea – find groups of equivalent states and merge them
 - All transitions from states in group G1 go to states in another group G2
 - Construct minimized DFA such that there is 1 state for each group of states



Basic strategy: identify distinguishing transitions

State minimization - algorithm

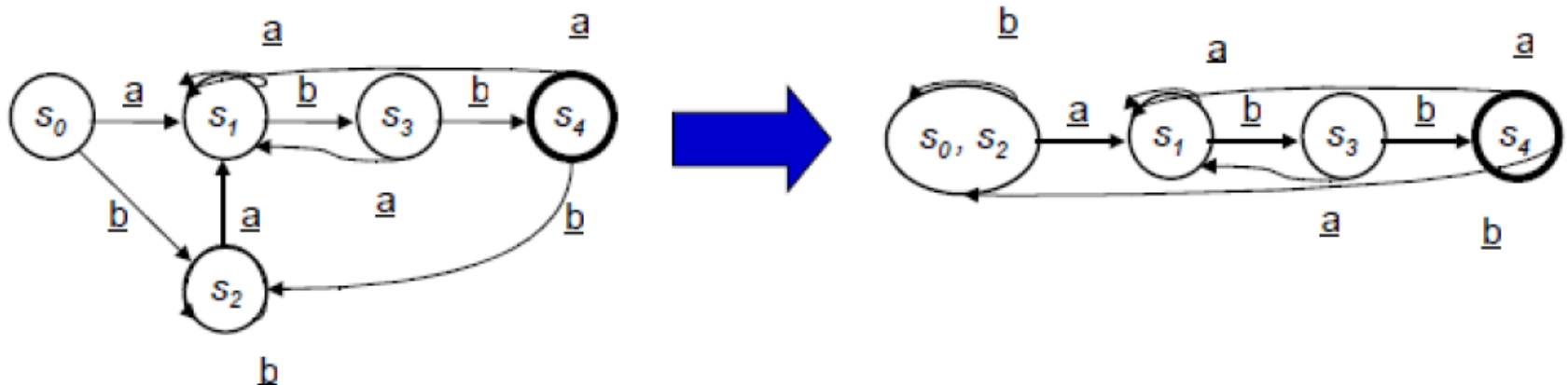
```
 $P \leftarrow \{ F, \{Q-F\} \}$   
while (  $P$  is still changing )  
   $T \leftarrow \{ \}$   
  for each set  $s \in P$   
    for each  $\alpha \in \Sigma$   
      partition  $s$  by  $\alpha$   
        into  $s_1, s_2, \dots, s_k$   
       $T \leftarrow T \cup s_1, s_2, \dots, s_k$   
  if  $T \neq P$  then  
     $P \leftarrow T$ 
```

State minimization - example

example: $(\underline{a} \mid \underline{b})^* \underline{a} \underline{b} \underline{b}$

	<i>Current Partition</i>	<i>Split on <u>a</u></i>	<i>Split on <u>b</u></i>
P_0	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	none	$\{s_0, s_1, s_2\} \{s_3\}$
P_1	$\{s_4\} \{s_3\} \{s_0, s_1, s_2\}$	none	$\{s_0, s_2\} \{s_1\}$
P_2	$\{s_4\} \{s_3\} \{s_1\} \{s_0, s_2\}$	none	None

final state

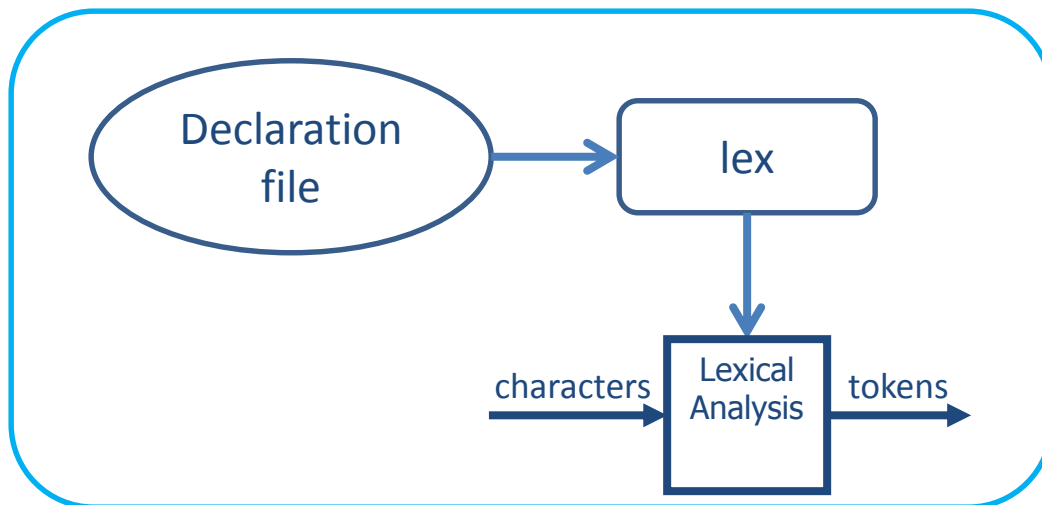


Summary of Construction

- Describe tokens as regular expressions.
- Regular expressions are turned into NFA.
- NFA is turned to DFA via subset construction.
- State minimization of DFA.
- Lexical analyzer simulates the run of an automata with the given transition table on any input string.

Good News

- Construction is done automatically by common tools
- lex is your friend
 - Automatically generates a lexical analyzer from declaration file
- Advantages: short declaration file, easily checked, easily modified and maintained



Intuitively:

- Lex builds DFA table
- Analyzer simulates (runs) the DFA on a given input