



中山大學
SUN YAT-SEN UNIVERSITY



国家超级计算广州中心
NATIONAL SUPERCOMPUTER CENTER IN GUANGZHOU

Compilation Principle 编译原理

第2讲：词法分析(2)

张献伟

xianweiz.github.io

DCS290, Spring 2021



中山大學
SUN YAT-SEN UNIVERSITY



Review Questions

Q1: input and output of lexical analysis?

character stream \rightarrow tokens

Q2: how to denote a token?

$\langle \text{class, lexeme} \rangle$

Q3: atomic and compound REs?

atomic: ϵ , $\{a\}$

compound: $R1 | R2$, $R1R2$, $R1^*$

Q4: $(+|-)?([0-9])^*(0|2|4|6|8)$

even numbers

Q5: RE of identifiers in C language?

$(_letter)(_letter | digit)^*$

Alphabet Operations[字母表运算]

- **Product**[乘积]: $\Sigma_1 \Sigma_2 = \{ab \mid a \in \Sigma_1, b \in \Sigma_2\}$
 - E.g., $\{0, 1\}\{a, b\} = \{0a, 0b, 1a, 1b\}$
- **Power**[幂]: $\Sigma^n = \Sigma^{n-1} \Sigma, n \geq 1; \Sigma^0 = \{\epsilon\}$
 - Set of strings of length n
 - $\{0, 1\}^3 = \{0, 1\}\{0, 1\}\{0, 1\} = \{000, 001, 010, 011, 100, 101, 110, 111\}$
- **Positive Closure**[正闭包]: $\Sigma^+ = \Sigma \cup \Sigma^2 \cup \Sigma^3 \cup \dots$
 - $\{a, b, c\}^+ = \{a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, aaa, aab, \dots\}$
- **Kleene Closure**[闭包]: $\Sigma^* = \Sigma^0 \cup \Sigma^+$

Regular Expressions

- Atomic[原子]

- ϵ is a RE: $L(\epsilon) = \{\epsilon\}$
- If $a \in \Sigma$, then a is a RE: $L(a) = \{a\}$

- Compound[組合]

- If both r and s are REs, corr. to languages $L(r)$ and $L(s)$, then:
- $r|s$ is a RE: $L(r|s) = L(r) \cup L(s)$
- rs is a RE: $L(rs) = L(r)L(s)$
- r^* is a RE: $L(r^*) = (L(r))^*$
- (r) is a RE: $L((r)) = L(r)$

Different REs of the Same Language

- $(a|b)^* = ?$

- $L((a|b)^*) = (L(a|b))^* = (L(a) \cup L(b))^* = (\{a\} \cup \{b\})^* = \{a, b\}^*$
- $= \{a, b\}^0 + \{a, b\}^1 + \{a, b\}^2 + \dots$
- $= \{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$

- $(a^*b^*)^* = ?$

- $L((a^*b^*)^*) = (L(a^*b^*))^* = (L(a^*)L(b^*))^*$
- $= L(\{\epsilon, a, aa, \dots\}\{\epsilon, b, bb, \dots\})^*$
- $= L(\{\epsilon, a, b, aa, ab, bb, \dots\})^*$
- $= \epsilon + \{\epsilon, a, b, aa, ab, bb, \dots\} + \{\epsilon, a, b, aa, ab, bb, \dots\}^2 + \{\epsilon, a, b, aa, ab, bb, \dots\}^3 + \dots$

Lexical Specification of a Language

- S0: write a regex for the lexemes of each token class
 - Numbers = digit^+
 - Keywords = 'if' + 'else' + ...
 - Identifiers = $\text{letter}(\text{letter} + \text{digit})^*$
- S1: construct R, matching all lexemes for all tokens
 - $R = \text{numbers} + \text{keywords} + \text{identifiers} + \dots = R1 + R2 + R3 + \dots$
- S2: let input be $x_1 \dots x_n$, for $1 \leq i \leq n$, check $x_1 \dots x_i \in L(R)$
- S3: if successful, then we know $x_1 \dots x_i \in L(R_j)$ for some j
- S4: remove $x_1 \dots x_i$ from input and go to step S2

Lexical Specification of a Language

- How much input is used?
 - $x_1 \dots x_i \in L(R)$, $x_1 \dots x_j \in L(R)$, $i \neq j$
 - Which one do we want? (e.g., '==' or '=')
 - Maximal match: always choose the longer one
- Which token is used if more than one matches?
 - $x_1 \dots x_i \in L(R)$ where $R = R_1 + R_2 + \dots + R_n$
 - $x_1 \dots x_i \in L(R_m)$, $x_1 \dots x_i \in L(R_n)$, $m \neq n$
 - E.g., keywords = 'if', identifier = letter(letter+digit)*
 - Keyword has higher priority
 - Rule of thumb: choose the one listed first
- What if no rule matches?
 - $x_1 \dots x_i \notin L(R) \rightarrow \text{Error}$

Summary: RE

- We have learnt how to specify tokens for lexical analysis
 - Regular expressions
 - Concise notations for the string patterns
- Used in lexical analysis with some extensions
 - To resolve ambiguities
 - To handle errors
- REs is only a language specification
 - An implementation is still needed
 - Next: to construct a token recognizer for languages given by regular expressions – by using **finite automata**

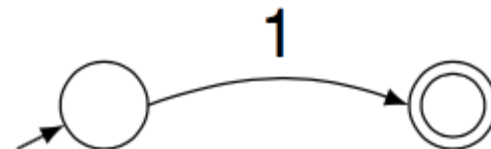
Implementation of Lexical Analyzer

- How do we go from specification to implementation?
 - RE \rightarrow finite automata
- **Solution 1:** to implement using a tool — Lex (for C), Flex (for C++), Jlex (for java)
 - Programmer specifies tokens using REs
 - The tool generates the source code from the given REs
 - The Lex tool essentially does the following translation: REs (Specification) \Rightarrow FAs (Implementation)
- **Solution 2:** to write the code yourself
 - More freedom; even tokens not expressible through REs
 - But difficult to verify; not self-documenting; not portable; usually not efficient
 - Generally not encouraged

Transition Diagram[转换图]

- REs → transition diagrams

- By hand
- Automatic



- Node: state

- Each state represents a condition that may occur in the process
- Initial state (Start): only one, circle marked with 'start →'
- Final state (Accepting): may have multiple, double circle

- Edge: directed, labeled with symbol(s)

- From one state to another on the input

Finite Automata[有穷自动机]

- **Regular Expression** = specification
- **Finite Automata** = implementation
- Automaton (pl. automata): a machine or program
- Finite automaton (FA): a program with a finite number of states
- Finite Automata are similar to transition diagrams
 - they have states and labelled edges
 - there are one unique start state and one or more than one final states

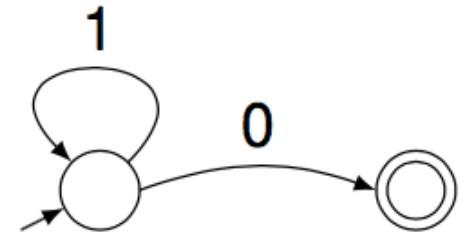
FA: Language

- An FA is a program for classifying strings (accept, reject)
 - In other words, a program for recognizing a language
 - The Lex tool essentially does the following translation: REs (Specification) \Rightarrow FAs (Implementation)
 - For a given string 'x', if there is transition sequence for 'x' to move from start state to certain accepting state, then we say 'x' is accepted by the FA
- Language of FA = set of strings accepted by that FA
 - $L(\text{FA}) \equiv L(\text{RE})$

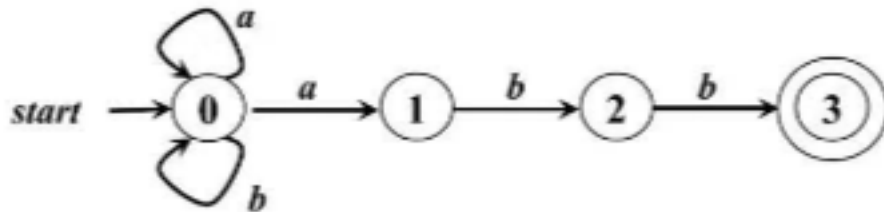
Example

- Are the following strings acceptable?

- 0 ✓
- 1 ✗
- 11110 ✓
- 11101 ✗
- 11100 ✗
- 1111110 ✓



- What language does the state graph recognize? $\Sigma = \{0, 1\}$
Any number of '1's followed by a single 0



L(FA): all strings of $\Sigma \{a, b\}$, ending with 'abb'

L(RE) = $(a|b)^*abb$

DFA and NFA

- Deterministic Finite Automata (**DFA**): the machine can exist in only one state at any given time
 - One transition per input per state
 - No ϵ -moves
 - Takes only one path through the state graph
- Nondeterministic Finite Automata (**NFA**): the machine can exist in multiple states at the same time
 - Can have multiple transitions for one input in a given state
 - Can have ϵ -moves
 - Can choose which path to take
 - An NFA accepts if some of these paths lead to accepting state at the end of input

State Graph

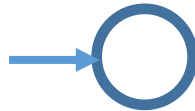
- 5 components $(\Sigma, S, n, F, \delta)$

- An input alphabet Σ

- A set of states S



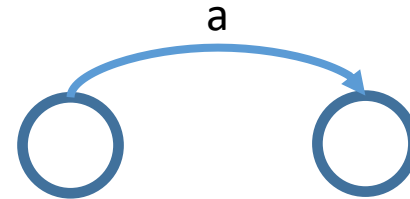
- A start state $n \in S$



- A set of accepting states $F \subseteq S$



- A set of transitions $\delta: S_a \xrightarrow{\text{input}} S_b$

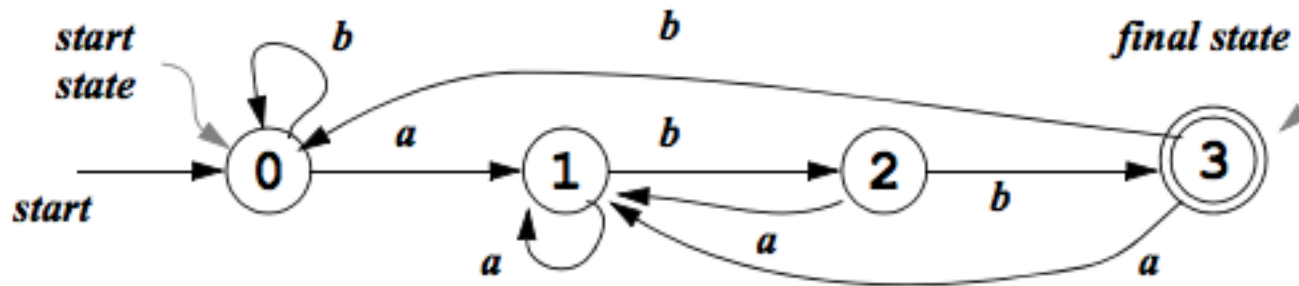


Example: DFA

- There is **only one** possible sequence of moves --- either lead to a final state and accept or the input string is rejected

– Input string: **aabb**

– Successful sequence: $0 \xrightarrow{a} 1 \xrightarrow{a} 1 \xrightarrow{b} 2 \xrightarrow{b} 3$



A DFA accepts $(a|b)^*abb$

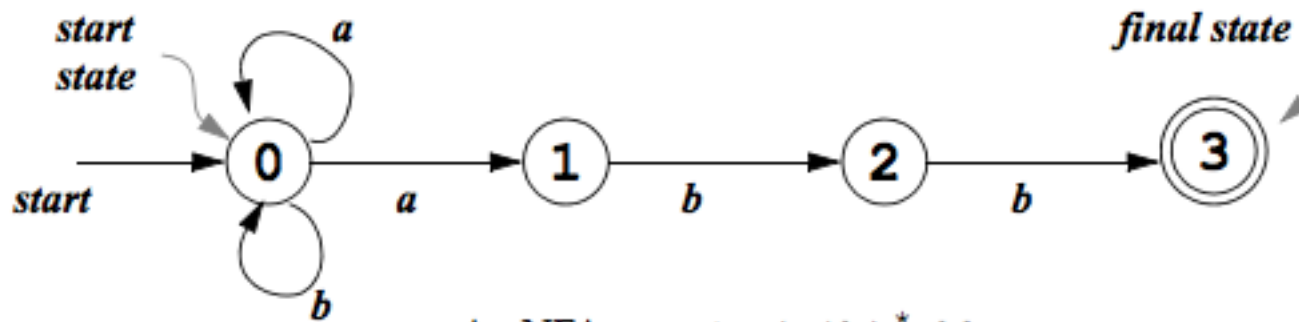
Example: NFA

- There are **many possible** moves --- to accept a string, we only need one sequence of moves that lead to a final state

- Input string: **aabb**

- Successful sequence: $0 \xrightarrow{a} 0 \xrightarrow{a} 1 \xrightarrow{b} 2 \xrightarrow{b} 3$

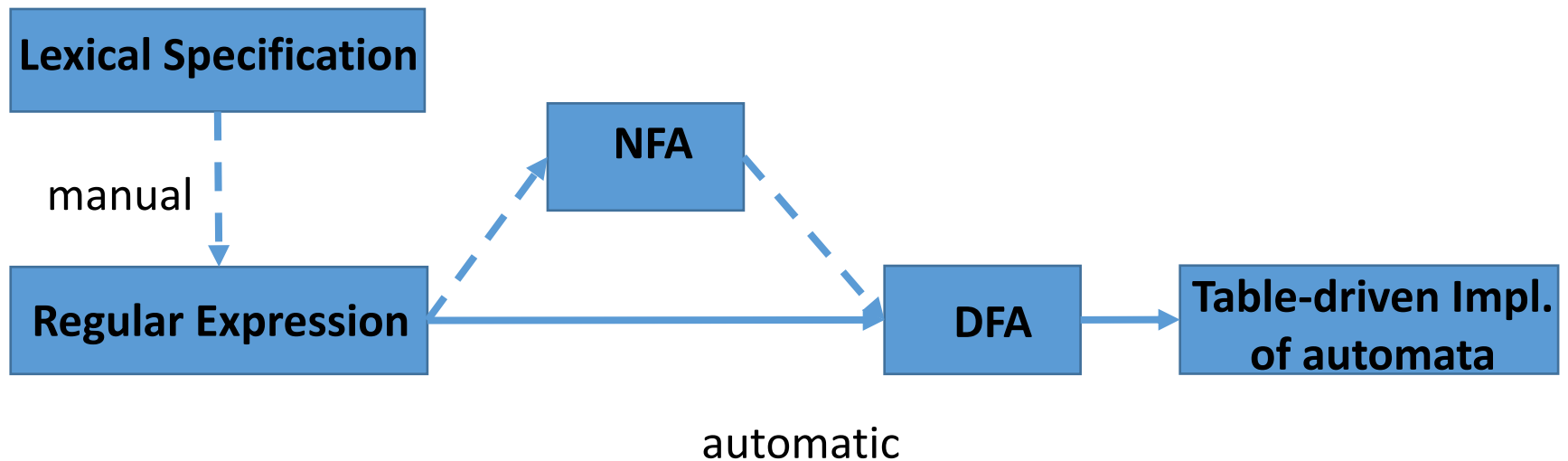
- Unsuccessful sequence: $0 \xrightarrow{a} 0 \xrightarrow{a} 0 \xrightarrow{b} 0 \xrightarrow{b} 0$



An NFA accepts $(a|b)^*abb$

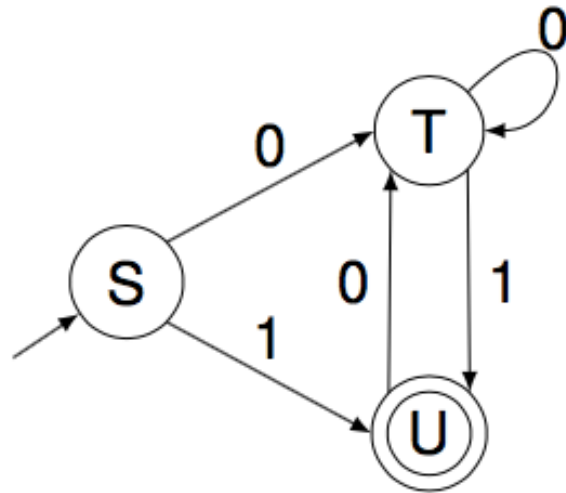
Conversion Flow

- Outline: RE \rightarrow NFA \rightarrow DFA \rightarrow Table-driven Implementation
 - Converting DFAs to table-driven implementations
 - Converting REs to NFAs
 - Converting NFAs to DFAs



DFA \rightarrow Table

- FA can also be represented using transition table



alphabet			
state		0	1
	S	T	U
	T	T	U
	U	T	x

Table-driven Code:

```
DFA() {  
    state = "S";  
    while (!done) {  
        ch = fetch_input();  
        state = Table[state][ch];  
        if (state == "x")  
            print("reject");  
    }  
    if (state ∈ F)  
        printf("accept");  
    else  
        printf("reject");  
}
```

Q: which is/are accepted?

111

000

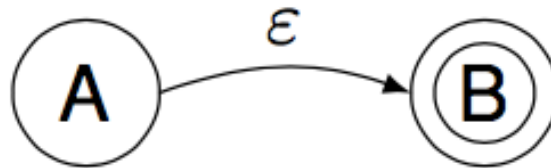
001

Discussion

- Implementation is efficient
 - Table can be automatically generated
 - Need finite memory $O(S \times \Sigma)$
 - Size of transition table
 - Need finite time $O(\text{input length})$
 - Number of state transitions
- Pros and cons of table:
 - Pro: can easily find the transitions on a given state and input
 - Con: takes a lot of space, when the input alphabet is large, yet most states do not have any moves on most of the input symbols

RE \rightarrow NFA

- NFA can have ϵ -moves
 - Edges labelled with ϵ
 - move from state A to state B without reading any input



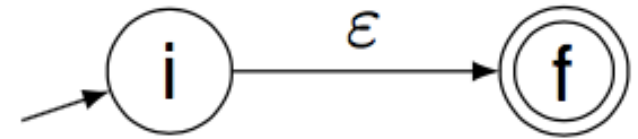
- M-Y-T algorithm to convert any RE to an NFA that defines the same language
 - Input: RE r over alphabet Σ
 - Output: NFA accepting $L(r)$

RE \rightarrow NFA (cont.)

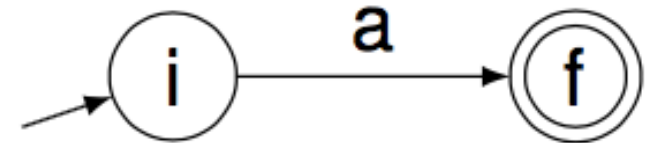
- Step 1: processing atomic REs

- ϵ expression

- i is new state, the start state of NFA
 - f is another new state, the accepting state of NFA

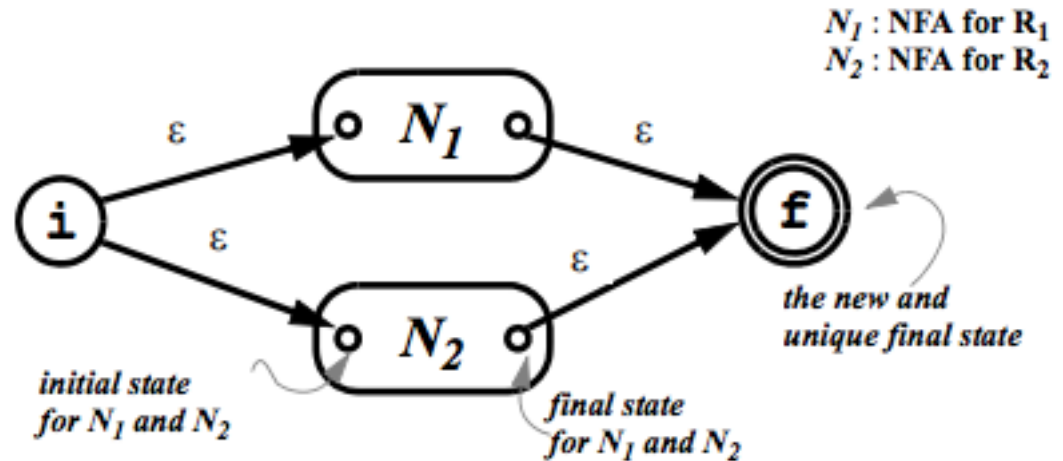


- Single character RE a

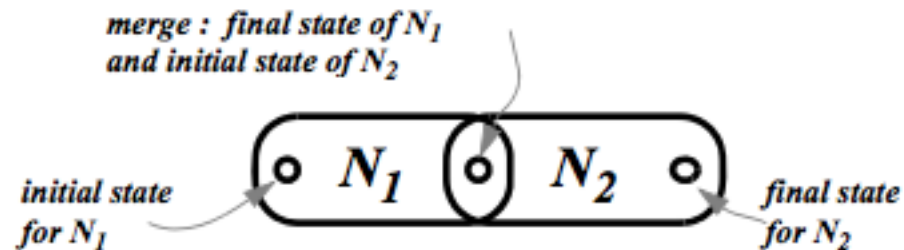


RE \rightarrow NFA (cont.)

- Step 2: processing compound REs
 - $R = R_1 \mid R_2$

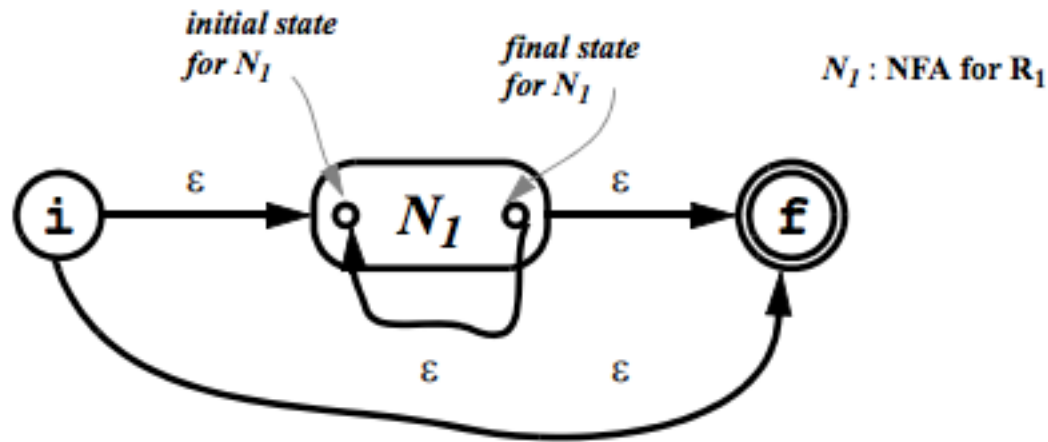


- $R = R_1R_2$



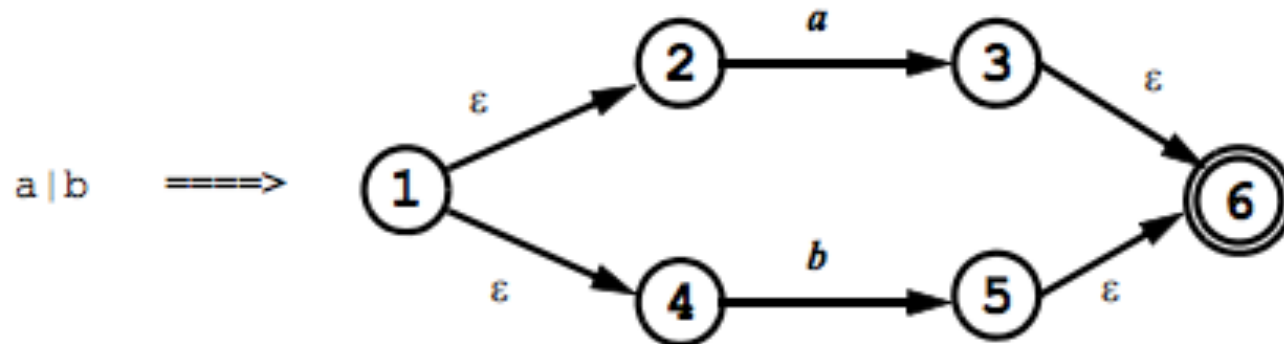
RE \rightarrow NFA (cont.)

- Step 2: processing compound REs
 - $R = R_1^*$



Example

- Convert “(a|b)*abb” to NFA



Example (cont.)

- Convert “(a|b)*abb” to NFA

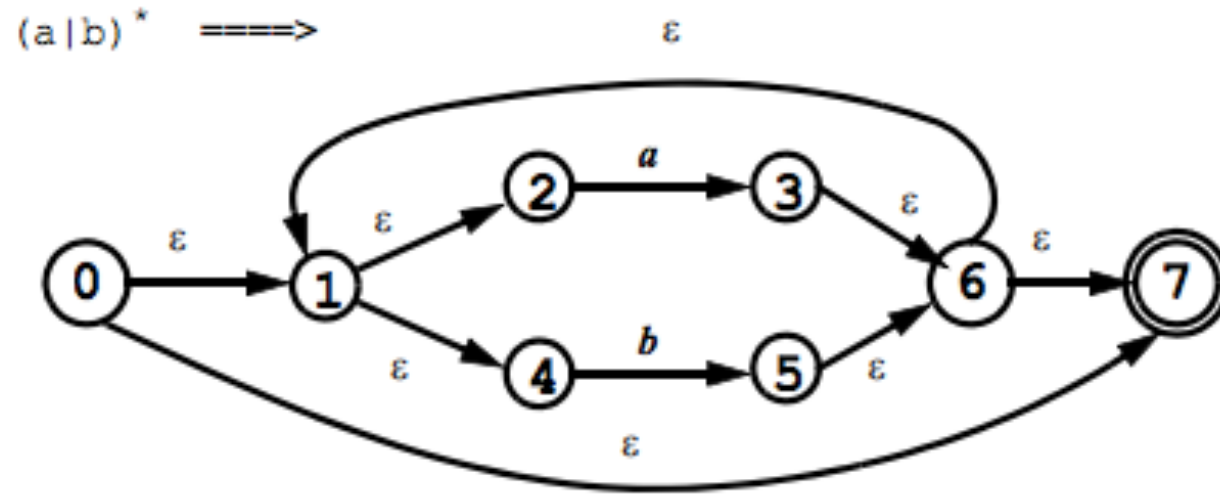
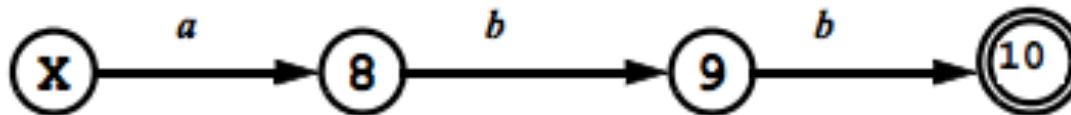
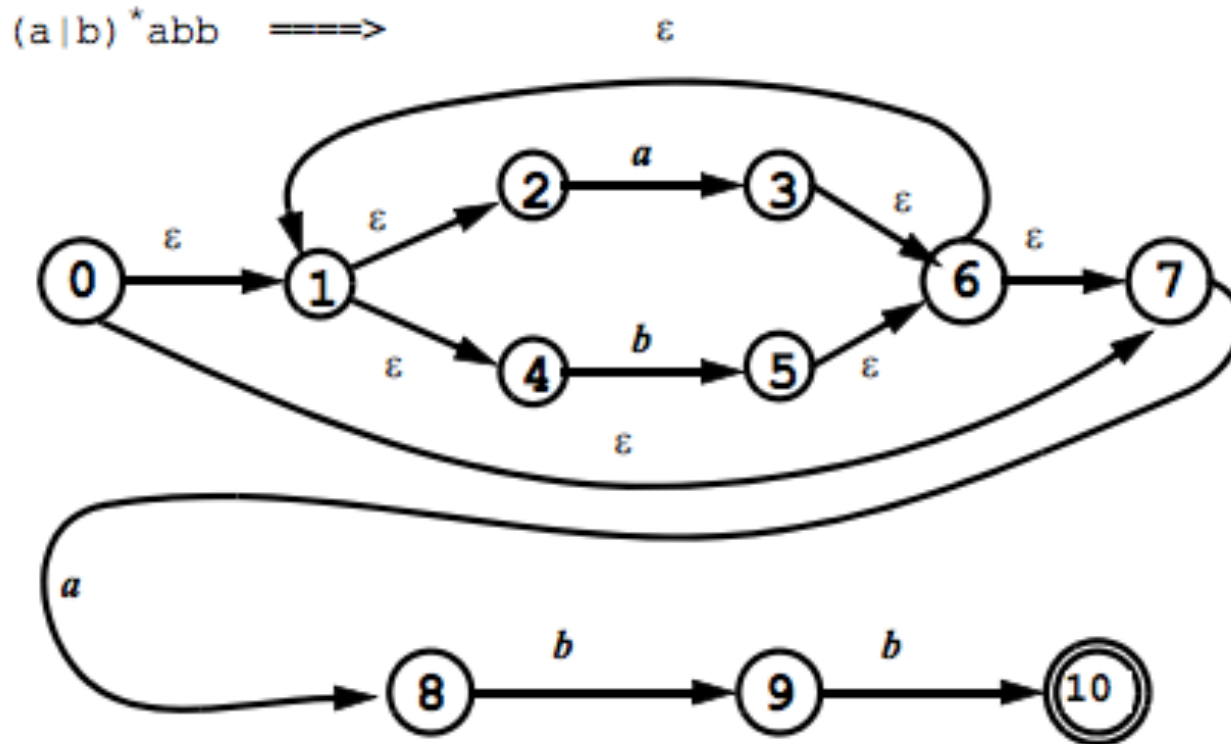


abb \implies (several steps are omitted)



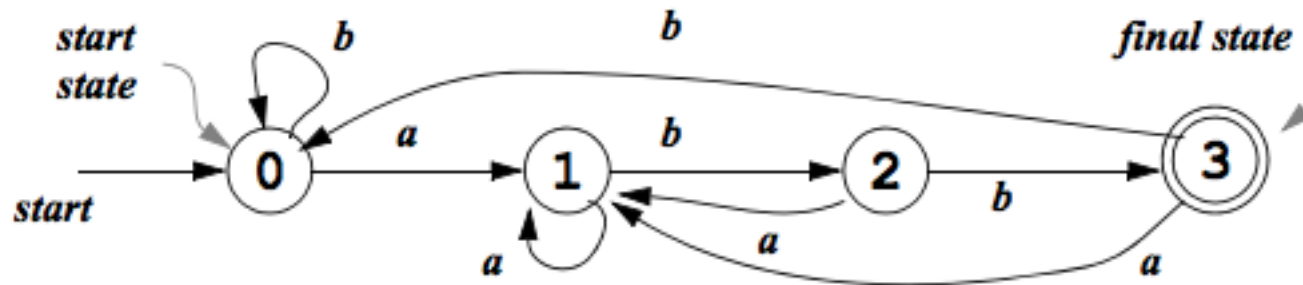
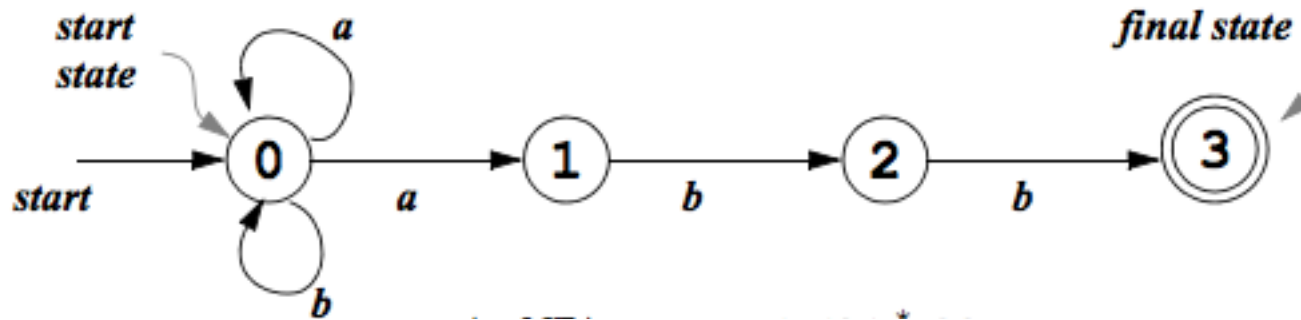
Example (cont.)

- Convert “(a|b)*abb” to NFA



NFA \rightarrow DFA: Same

- NFA and DFA are equivalent



NFA \rightarrow DFA: Theory

- Question: is $L(\text{NFA}) \subseteq L(\text{DFA})$
 - Otherwise, conversion would be futile
- Theorem: $L(\text{NFA}) \equiv L(\text{DFA})$
 - Both recognize regular languages $L(\text{RE})$
 - Will show $L(\text{NFA}) \subseteq L(\text{DFA})$ by construction ($\text{NFA} \rightarrow \text{DFA}$)
 - Since $L(\text{DFA}) \subseteq L(\text{NFA})$, $L(\text{NFA}) \equiv L(\text{DFA})$
- Resulting DFA consumes more memory than NFA
 - Potentially larger transition table as shown later
- But DFAs are faster to execute
 - For DFAs, number of transitions == length of input
 - For NFAs, number of potential transitions can be larger
- $\text{NFA} \rightarrow \text{DFA}$ conversion is done because the speed of DFA far outweighs its extra memory consumption