



中山大學
SUN YAT-SEN UNIVERSITY



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

DCS²⁹⁰

Compilation Principle 编译原理

第三章 词法分析

郑馥丹

zhengfd5@mail.sysu.edu.cn

CONTENTS

目录

01

概述

Introduction

02

词法规范

Lexical
Specification

03

有穷自动机

Finite
Automata

04

转换和等价

Transformation
and Equivalence

05

词法分析实践

Lexical Analysis
in Practice

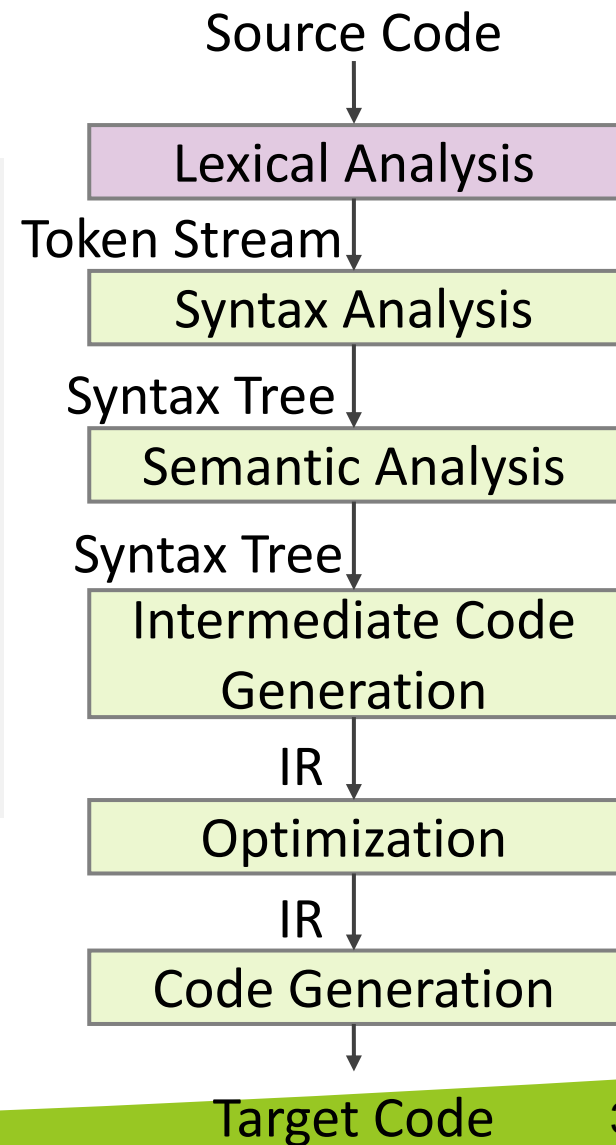
1. 什么是词法分析[Lexical Analysis]?

- 扫描源程序字符流，识别并分解出有词法意义的单词或符号 (token)

- 输入：源程序，输出：token序列
- token表示：<类别，属性值>
 - ✓ 保留字、标识符、常量、运算符等
- token是否符合词法规则？
 - ✓ 0var, \$num

```
void main()  
{  
    int arr[10], i, x = 1;  
    for (i = 0; i < 10; i++)  
        arr[i] = x * 5;  
}
```

keyword(for)	id(arr)
symbol((symbol([
id(i)	id(i)
symbol(=)	symbol(]
num(0)	symbol(=)
symbol(;	id(x)
id(i)	symbol(*)
symbol(<)	num(5)
num(10)	symbol(;
symbol(;	id(i)
id(i)	symbol(++)
symbol(++)	symbol())
symbol())	



1. 什么是词法分析[Lexical Analysis]?

- 扫描源程序字符流，识别并分解出有词法意义的单词或符号 (token)

- 例：

- 输入：

- ✓ 字符串“if (i == j)\n\tz = 0; \nelse\n\tz = 1; \n”

- 目标：将字符串划分为一组**tokens**

- 步骤：

- ① 移除注释：/* simple example */

- ② 识别tokens：‘if’ ‘(’ ‘i’ ‘==’ ‘j’

- ③ 识别tokens所属的类别：(**keyword**, ‘if’), (**LPAR**, ‘(’), (**id**, ‘i’)

- 输出：(**keyword**, ‘if’), (**LPAR**, ‘(’), (**id**, ‘i’)

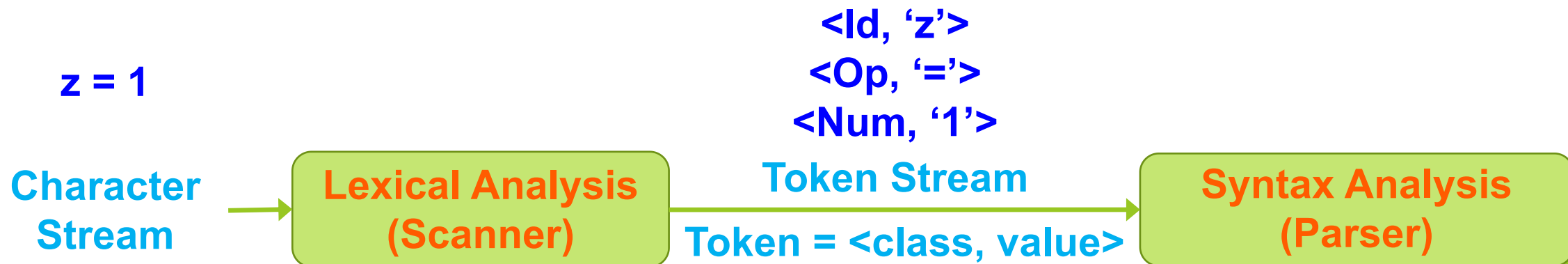
```
/* simple example */  
if (i == j)  
    z = 0;  
else  
    z = 1;
```

2. 什么是词[token]?

- Token: 词, **最小意义单元**
 - 英语中的token类别:
 - ✓ 名词noun, 动词verb, 形容词adjective, ...
 - 编程语言中的token类别:
 - ✓ 数字number, 关键词keyword, 空白whitespace, 标识符identifier, ...
- 每个token类别对应一个字符串集合[a set of strings]
 - number: 非空的数字字符串
 - keyword: 一组固定的**保留字**(“for”, “if”, “else”, ...)
 - whitespace: 由空格blank、制表符tab、换行符newline组成的非空序列
 - identifier: 用户自定义的要标识的实体名称

3. 词法分析：分词[Tokenization/Scanner]

- 词法分析器也称**标记化Tokenization**或**扫描器Scanner**
 - 将输入字符串划分为token序列
 - 对token进行分类——token类别
 - ✓ **词素[lexeme]**：token类的实例，例如：‘z’，‘=’，‘1’
 - ✓ 事先定义好token类别：例如keyword, identifier, whitespace, integer等
- 得到的token会被送入**语法分析器[Syntax Analyzer]**（也称**Parser**）
 - parser依赖token类来识别角色（例如，关键字与标识符的处理方式不同）。



3. 词法分析：设计[Design]

- 定义token类别

- 描述所有感兴趣的项
- 依赖于语言、parser的设计

✓ `"if (i == j)\n\tz = 0; \nelse\n\tz = 1; \n"`

Keyword, identifier, whitespace, integer

- 识别哪些字符串属于哪个token类别

```
if (i == j)
    z = 0;
else
    z = 1;
```

'==' or '='?

keyword or identifier?

5种最常见的编程语言token类别:

- Identifiers: `getBalance`, `weight`,....
- Reserved words: `if`, `else`, `while`,...
- Constants: `10`, `3.14`, `"abc"`, `'a'`,...
- Operators: `+`, `-`, `*`, `/`, `<<`,...
- Punctuation: `(`, `)`, `;`, `:`,...

3. 词法分析：实现[Implementation]

- 实现时需完成2件事情：
 - 识别字符串的token分类
 - 返回token的值或词素
- 一个token是一个二元组(class, lexeme)
- 丢弃无意义词
 - 例如：whitespace, comments等
- 如果token类别是无二义性的，则可以在对输入字符串进行一次从左到右的扫描中识别标记
- 但token类别可能有歧义[ambiguous]

3. 词法分析：实现[Implementation]

- 歧义举例

- C++ 模板语法

- ✓ Foo<Bar>

- C++ 流语法

- ✓ cin >> var

- 二义性

- ✓ Foo<Bar<Bar>>>

- ✓ cin >>> var

✓ 疑问：‘>>’ 应当成是流操作符，还是两个连续的括号呢？

```
Template <typename T>
T getMax(T x, T y) {
    return (x > y) ? x : y;
}
```

```
int main (int argc, char* argv[]) {
    getMax<int>(3, 7);
    getMax<double>(3.0, 2.0);
    getMax<char>('g', 'e');

    return 0;
}
```

3. 词法分析：实现[Implementation]

- “**向前看[look ahead]**”可以展望消除歧义
 - 要查看更大的上下文或结构 `cin >> var`
 - 有时词法分析需要parser的反馈
- 如果token没有重叠
 - tokenizing可以在没有parser反馈的情况下一次完成
 - 可明确区分词法和语法分析

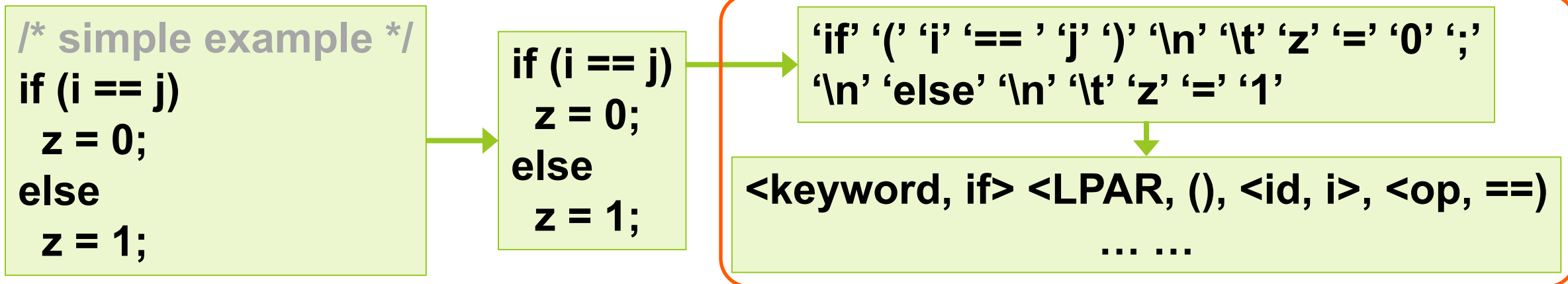
3. 词法分析：总结

- 词法分析

- 将输入符号串分解成token
- 识别每个token的类别

- 从左到右扫描

- 有时候需要 “向前看[look ahead]”
- 应该尽量减少 “向前看[look ahead]” 的数量



CONTENTS

目 录

01

概述

Introduction

02

词法规范

Lexical
Specification

03

有穷自动机

Finite
Automata

04

转换和等价

Transformation
and Equivalence

05

词法分析实践

Lexical Analysis
in Practice

1. 词法规范

- 三种词法规范

- 表达式[Expressions]

- ✓ 正则表达式[Regular expression]

- ✓ 正则定义[Regular definition]

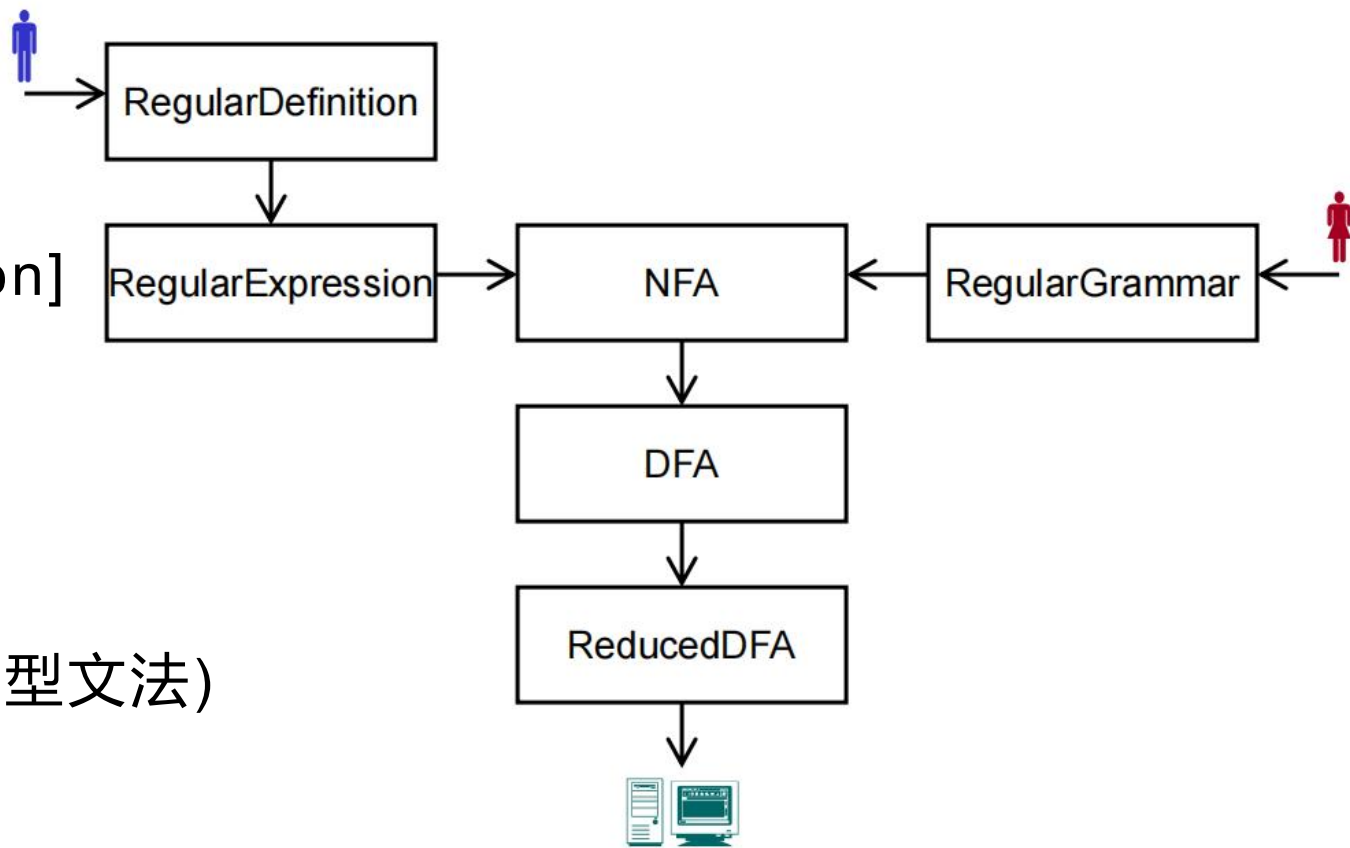
- 语法

- ✓ 正规文法[Regular grammar](3型文法)

- 有穷自动机

- ✓ 确定的有穷自动机[Deterministic Finite Automata (DFA)]

- ✓ 不确定的有穷自动机[Nondeterministic Finite Automata (NFA)]



2. 正则表达式[Regular expressions, REs]

- 正则表达式可用于定义token class
 - **简单但功能强大**（能够表达模式）
 - 可以从规范**自动生成**Tokenizer实现（使用翻译工具）
 - 最终实现是**有效的**

2. 正则表达式[Regular expressions, REs]

• 定义

- ① ε 和 Φ 是 Σ 上的正规式，正规集 $L(\varepsilon)=\{\varepsilon\}$ ， $L(\Phi)=\Phi$ ；
- ② 对于任何 $a \in \Sigma$ ， a 是 Σ 上的正规式，它所表示的正规集为 $L(a)=\{a\}$ ；
- ③ 假定 r 和 s 都是 Σ 上的正规式，他们所表示的正规集分别为 $L(r)$ 和 $L(s)$ ，那么，以下也都是正规式和它们所表示的正规集；

正规式	正规集
(r)	$L(r) = L((r)) = L(r)$
$r \mid s$	$L(r \mid s) = L(r) \cup L(s)$
$r \cdot s$	$L(r \cdot s) = L(r) L(s)$
r^*	$L(r^*) = (L(r))^*$

算符的优先顺序: $'()' > '*' > '\cdot' > '\mid'$
 $'\cdot'$ 和 $'\mid'$ 都是左结合

$ab^*c \mid d$

$a(b^*)c \mid d$

$(a(b^*))c \mid d$

$((a(b^*))c) \mid d$

- ④ 仅由有限次使用上述三步定义的表达式才是 Σ 上的正规式，仅由这些正规式所表示的字集才是 Σ 上的正规集。

2. 正则表达式[Regular expressions, REs]

- 例子

- 令 $\Sigma = \{a, b\}$, Σ 上的正规式和相应的正规集有:

正规式

正规集

a

{a}

a | b

{a,b}

ab

{ab}

(a | b)(a | b)

{aa,ab,ba,bb}

a *

{ ϵ , a, aa, ... 任意个a串}

(a | b)*

{ ϵ , a, b, aa, ab 所有由a和b组成的串}

2. 正则表达式[Regular expressions, REs]

• 正则表达式的代数性质

$$- r \mid s = s \mid r$$

“ \mid ” 满足交换律

$$- r \mid (s \mid t) = (r \mid s) \mid t$$

“ \mid ” 的结合律

$$- (r s) t = r (s t)$$

“ \cdot (连接)” 的结合律

$$- r(s \mid t) = r s \mid r t$$

$$- (r \mid s) t = r t \mid s t$$

分配律

$$- \varepsilon r = r \varepsilon = r$$

ε 是 “ \cdot ” 的恒等元素

$$- r \mid r = r$$

“ \mid ” 的抽取律 $r^* = \varepsilon \mid r \mid rr \mid \dots$

2. 正则表达式[Regular expressions, REs]

- 常用表达

- At least one: $A^+ \equiv AA^*$
- Option: $A? \equiv A \mid \epsilon$
- Characters: $[a_1a_2...a_n] \equiv a_1|a_2|...|a_n$
- Range: $'a' + 'b' + ... + 'z' \equiv [a-z]$
- Excluded range: **complement of** $[a-z] \equiv [^a-z]$

2. 正则表达式[Regular expressions, REs]

• 例子

Regular Expression	Explanation
a^*	0 or more a's (ϵ , a, aa, aaa, aaaa, ...)
a^+	1 or more a's (a, aa, aaa, aaaa, ...)
$(a b)(a b)$	(aa, ab, ba, bb)
$(a b)^*$	all strings of a's and b's (including ϵ)
$(aa ab ba bb)^*$	all strings of a's and b's of even偶数 length
$[a-zA-Z]$	shorthand简写 for "a b ...z A B ... Z"
$[0-9]$	shorthand for "0 1 2 ... 9"
$0([0-9])^*0$	numbers that start and end with 0
$1^*(0 \epsilon)1^*$	binary strings that contain at most one zero
$(0 1)^*00(0 1)^*$	all binary strings that contain '00' as substring

2. 正则表达式[Regular expressions, REs]

- 思考: $(a|b)^*$ 和 $(a^*b^*)^*$ 是否等价? **等价**

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

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

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

$$\begin{aligned} (L(a^*b^*))^* &= (L(a^*)L(b^*))^* \\ &= (\{\epsilon, a, aa, \dots\}\{\epsilon, b, bb, \dots\})^* \\ &= L(\{\epsilon, a, b, aa, ab, bb, \dots\})^* \\ &= \epsilon + \{\epsilon, a, b, aa, ab, bb, \dots\} + \{\epsilon, a, b, aa, ab, bb, \dots\}^2 + \{\epsilon, a, b, aa, ab, \\ &\quad bb, \dots\}^3 + \dots \end{aligned}$$

2. 正则表达式[Regular expressions, REs]

- 更多例子

- Keywords: 'if' 或 'else' 或 'then' 或 'for' ...
 - ✓ RE = 'i' 'f' + 'e' 'l' 's' 'e' + ... = 'if' + 'else' + 'then' + ...
- Numbers: 非空数字字符串
 - ✓ digit = '0' + '1' + '2' + '3' + '4' + '5' + '6' + '7' + '8' + '9'
 - ✓ 若定义integer = digit digit*, 请问'000' 是否为integer?
- Identifier: 以字母开头的字母或数字的字符串
 - ✓ letter = 'a' + 'b' + ... 'z' + 'A' + 'B' + ... + 'Z' = [a-zA-Z]
 - ✓ 若定义RE = letter(letter + digit)*, 请问RE符合C语言中标识符的定义吗?
- Whitespace: 由空格、换行符、制表符组成的非空序列
 - ✓ (' ' + '\n' + '\t')+

2. 正则表达式[Regular expressions, REs]

- 编程语言中的REs <https://docs.python.org/3/howto/regex.html>

Symbol	Meaning		
\d	Any decimal digit, i.e. [0-9]		
\D	Any non-digit char, i.e., [^0-9]		
\s	Any whitespace char, i.e., [\t\n\r\f\v]		
\S	Any non-whitespace char, i.e., [^ \t\n\r\f\v]		
\w	Any alphanumeric char, i.e., [a-zA-Z0-9_]		
\W	Any non-alphanumeric char, i.e., [^a-zA-Z0-9_]		
.	Any char	\.	Matching “.”
[a-f]	Char range	[^a-f]	Exclude range
^	Matching string start	\$	Matching string end
(...)	Capture matches		

2. 正则表达式[Regular expressions, REs]

- 正则表达式vs.正规文法[Regular grammar]

	正则表达式	正规文法
标识符集合	$l(l d)^*$ 其中: l为a-z的字母, d为0-9的数字	<ul style="list-style-type: none">$\langle \text{标识符} \rangle \rightarrow l l \langle \text{字母数字} \rangle$$\langle \text{字母数字} \rangle \rightarrow l d l \langle \text{字母数字} \rangle d \langle \text{字母数字} \rangle$
无符号整数	dd^*	$\langle \text{无符号整数} \rangle \rightarrow d d \langle \text{无符号整数} \rangle$

正规式比正规文法更容易让人理解单词是按怎样的规律构成的, 且可以从某个正规式**自动地**构造识别程序。

随堂练习(1)

- 写出C语言中标识符的正则表达式

(_|letter)(_|letter|digit)*

- 写出C语言中十进制数字的正则表达式

0|([1-9][0-9]*)

- 写出2的倍数的二进制表示的正则表达式

(0|1)*0

CONTENTS

目录

01

概述

Introduction

02

词法规范

Lexical
Specification

03

有穷自动机

Finite
Automata

04

转换和等价

Transformation
and Equivalence

05

词法分析实践

Lexical Analysis
in Practice

1. 转换图[Transition Diagram]

• 结点[Node]: 状态

- 词法分析器在扫描输入串的过程中寻找和某个模式匹配的词素
- 转换图的每一个状态代表一个在此过程中可能发生的情况
- **初始状态**(Start/Initial State): **只有一个**, 一般由没有出发结点的箭头表示
- **最终状态**(Accepting/Final States): **可以有多个, 用双圈表示**

• 边[Edge]: 有向[directed], 标记有符号[symbol(s)]



- 从一个状态指向另一个状态
- 如果处于某个状态S, 且下一个输入符号是a, 则会寻找一条从S状态离开且标号为a的边

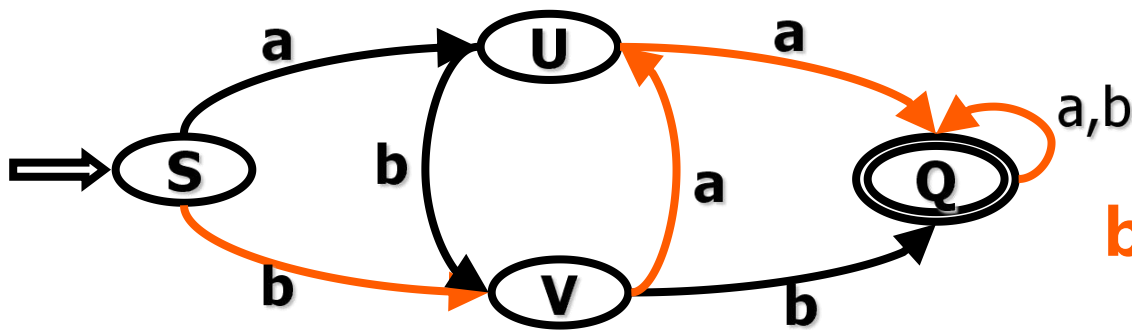
2. 有穷自动机[Finite Automata, FA]

- 正则表达是定义[Regular Expression = **specification**]
- 自动机是实现[Finite Automata = **implementation**]
- 自动机[Automata]: 一个机器或一个程序
- **有穷**自动机[Finite Automata, FA]: 具有**有穷个状态**的程序
- 有穷自动机基于转换图
 - **有状态和标记的边**
 - **有一个唯一的开始状态和一个或多个最终状态**

2. 有穷自动机[Finite Automata, FA]

- 有穷自动机FA是一个对字符串进行**分类（接受、拒绝）**的程序
 - 是一个识别语言的程序
 - Lex tool: 将REs(specification)转换成FAs(implementation)
 - 对于给定的字符串x, 如果存在一条从FA的**初始结点**到某一个最终结点的通路, 且该**通路上所有边上的符号连接成的字符串等于x**, 则称**x可被FA识别**, 或称**x被FA接受**

✓ 否则, 称**x被FA拒绝**



baab被FA接受

abcd被FA拒绝

- **FA所能表达的语言即为该FA接受的字符串集合**

- $L(FA) \equiv L(RE)$

2. 有穷自动机[Finite Automata, FA]

- 例：有右图FA

(1) 判断以下字符串是否被FA接受

✓ 0 ✓

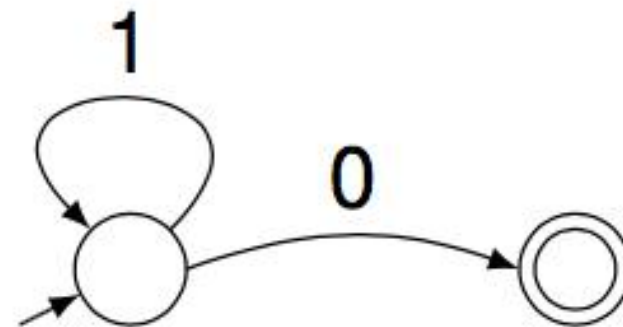
✓ 1 ✗

✓ 11110 ✓

✓ 11101 ✗

✓ 11100 ✗

✓ 1111110 ✓



(2) 在字母表 $\Sigma = \{0, 1\}$ 上，这个FA的语言是什么？

✓ 任意多个 '1' 后跟一个 '0'

2. 有穷自动机[Finite Automata, FA]

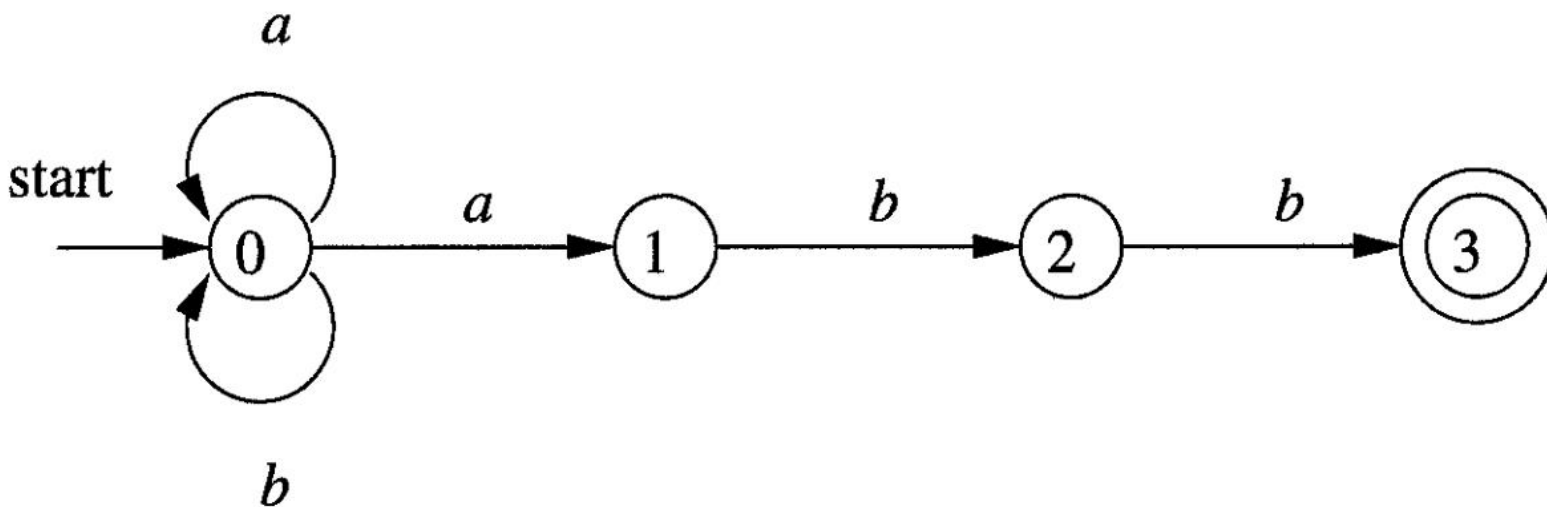
• 例：有右图FA

(1) 在字母表 $\Sigma = \{a, b\}$ 上，这个FA的语言是什么？

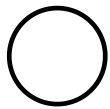
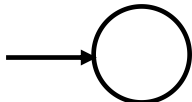
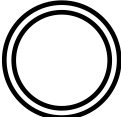
✓ 所有以若干个a和b的组成的串后跟 'abb' 的字符串

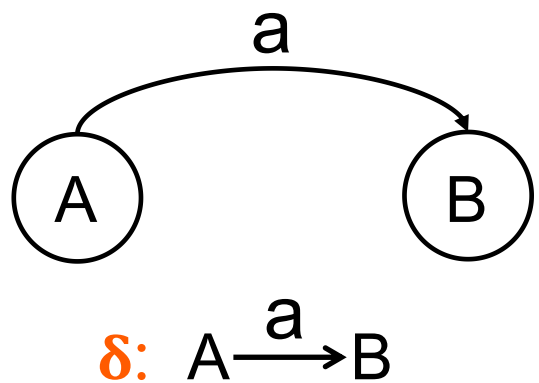
(2) 尝试用正则表达式表达这个语言

✓ $L(RE) = (a|b)^*abb$



2. 有穷自动机[Finite Automata, FA]

- 一个有穷自动机是一个五元组： $M = (S, \Sigma, \delta, s_0, F)$ ，其中：
 - S ：有穷状态集合 
 - Σ ：输入字母表
 - δ ：一个**转换函数[transition function]**，它为每个状态和 $\Sigma \cup \{\varepsilon\}$ 中的每个符号都给出了相应的后继状态集合
 - $s_0 \in S$ ：开始状态/初始状态 
 - $F \subseteq S$ ：接受状态/终止状态的**集合** 



3. DFA和NFA

- **DFA**(Deterministic Finite Automata): 机器在任何给定时间只能以**一种**状态存在(**确定**)
 - 每个状态对于每个输入只对应**一个转换**
 - **没有 ϵ 转移**[ϵ -moves]
 - 只通过状态图的一条**路径**
- **NFA**(Nondeterministic Finite Automata): 机器可以同时以**多种**状态存在(**非确定**)
 - 在给定状态下, 对一个输入可以有**多个转换**
 - 可以有 **ϵ 转移**[ϵ -moves]
 - **多条路径**: 可以选择采取哪条路径
 - ✓ 如果其中**某些**路径在输入结束时导致接受状态, 则NFA接受

对DFA: $\delta: S \times \Sigma \rightarrow S$

对NFA: $\delta: S \times \Sigma \rightarrow 2^S$

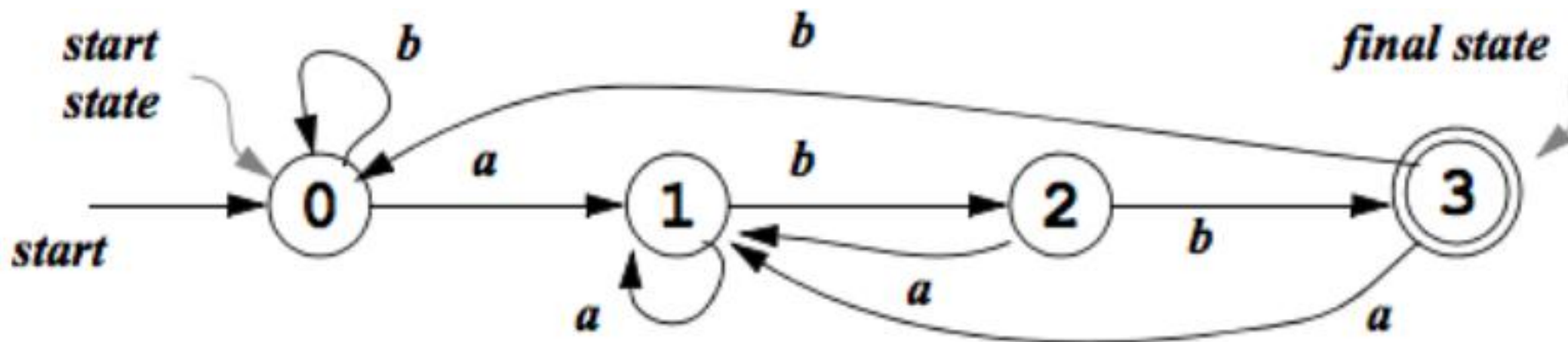
对 ϵ -NFA: $\delta: S \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^S$

- DFA示例

- **只有一种**可能的转移[moves]序列：**要么**导致最终状态并**接受**，**要么拒绝**输入字符串

- 输入字符串：aabb

- 成功序列：**0** \xrightarrow{a} **1** \xrightarrow{a} **1** \xrightarrow{b} **2** \xrightarrow{b} **3** **接受**



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

3. DFA和NFA

• NFA示例

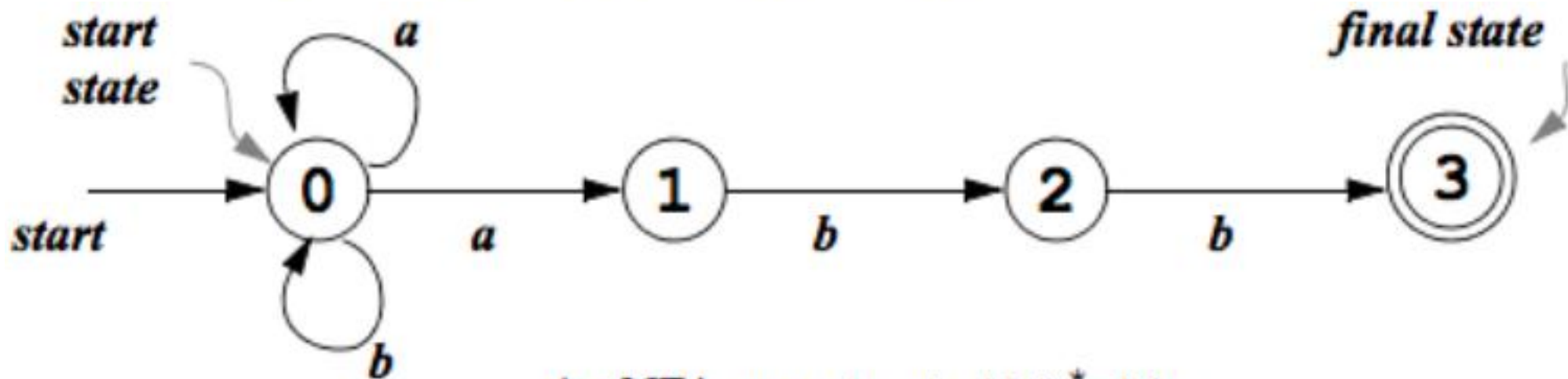
– 有**多个**可能的moves: **只要有一个moves序列导致最终状态, 即认为接受**

– 输入字符串: aabb

– 成功序列: $0 \xrightarrow{a} 0 \xrightarrow{a} 1 \xrightarrow{b} 2 \xrightarrow{b} 3$

– 不成功序列: $0 \xrightarrow{a} 0 \xrightarrow{a} 0 \xrightarrow{b} 0 \xrightarrow{b} 0$

接受



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

CONTENTS

目 录

01

概述

Introduction

02

词法规范

Lexical
Specification

03

有穷自动机

Finite
Automata

04

转换和等价

Transformation
and Equivalence

05

词法分析实践

Lexical Analysis
in Practice

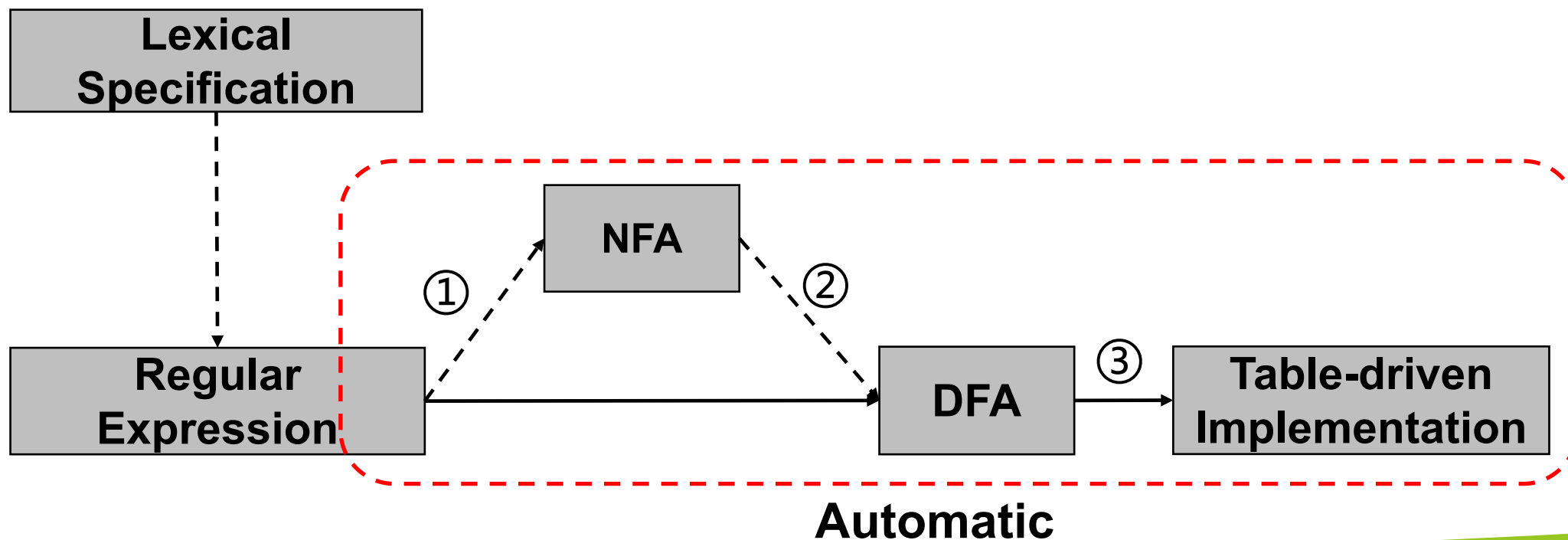
1. 转换流程[Conversion Flow]

• 流程：RE→NFA→DFA→Table-drive Implementation

① RE→NFA

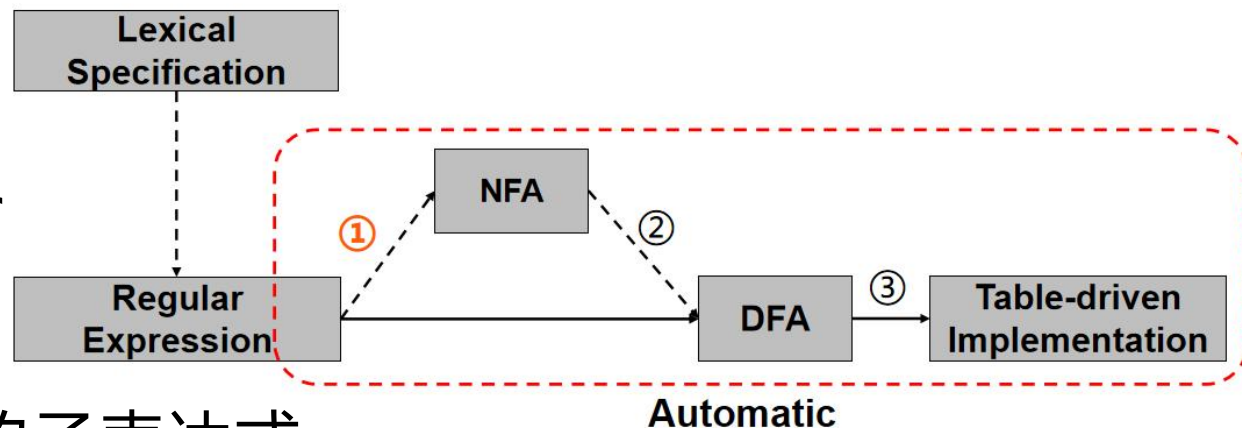
② NFA→DFA

③ DFA→Table-drive Implementation



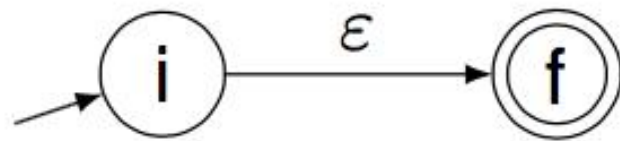
2. $RE \rightarrow NFA$ • $RE \rightarrow NFA$ 采用 Thompson 构造算法

- 输入：字母表 Σ 上的一个正则表达式 r
- 输出：一个接受 $L(r)$ 的 NFA N
- 方法：对 r 进行分析，分解出组成它的子表达式



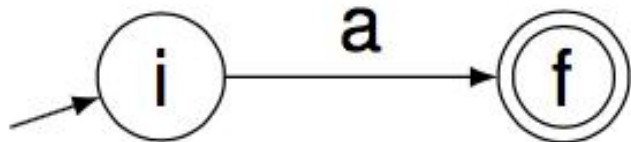
• Step 1: 处理原子 REs (基本规则)

- 对于表达式 ϵ ：
 - 增加新状态 i ，作为 NFA 的初始状态
 - 增加另一个新状态 f ，作为 NFA 的接受状态



- 对于表达式 a ：

- 同上



Thompson

Ken Thompson



Kenneth Lane Thompson (born February 4, 1943) is an American pioneer of [computer science](#) & Computer Chess Development. Thompson worked at [Bell Labs](#) for most of his career where he designed and implemented the original [Unix operating system](#). He also invented the [B programming language](#), the direct predecessor to the [C programming language](#), and was one of the creators and early developers of the [Plan 9](#) operating system. Since 2006, Thompson has worked at [Google](#), where he co-developed the [Go programming language](#).

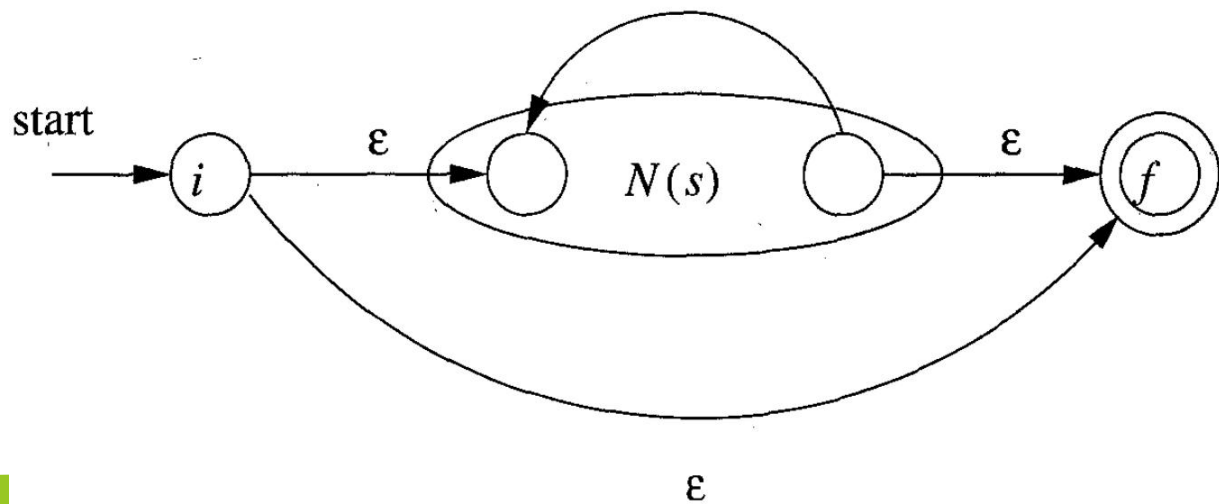
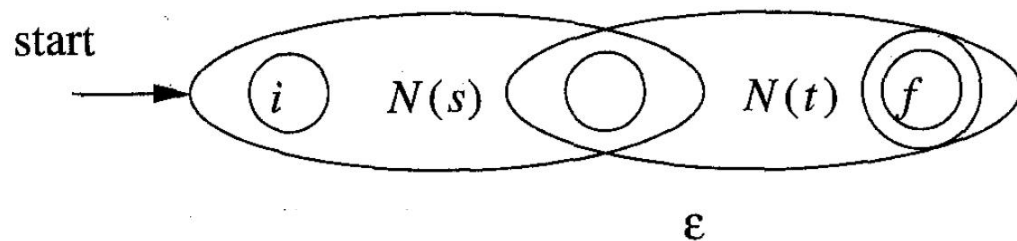
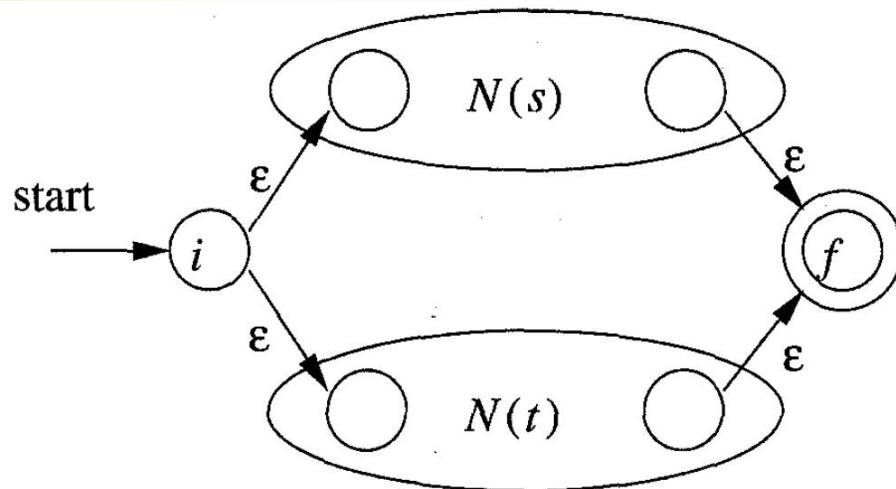
Other notable contributions included his work on [regular expressions](#) and early computer text editors [QED](#) and [ed](#), the definition of the [UTF-8](#) encoding, and his work on computer chess that included the creation of [endgame tablebases](#) and the chess machine [Belle](#). He [won the Turing Award in 1983](#) with his long-term colleague [Dennis Ritchie](#).

In the 1960s, Thompson also began work on [regular expressions](#). Thompson had developed the [CTSS](#) version of the editor [QED](#), which included regular expressions for searching text. QED and Thompson's later editor [ed](#) (the standard text editor on Unix) contributed greatly to the eventual popularity of regular expressions, and regular expressions became pervasive in Unix text processing programs. Almost all programs that work with regular expressions today use some variant of Thompson's notation. He also invented [Thompson's construction algorithm](#) used for converting regular expressions into [nondeterministic finite automata](#) in order to make expression matching faster.^[12]

2. $RE \rightarrow NFA$

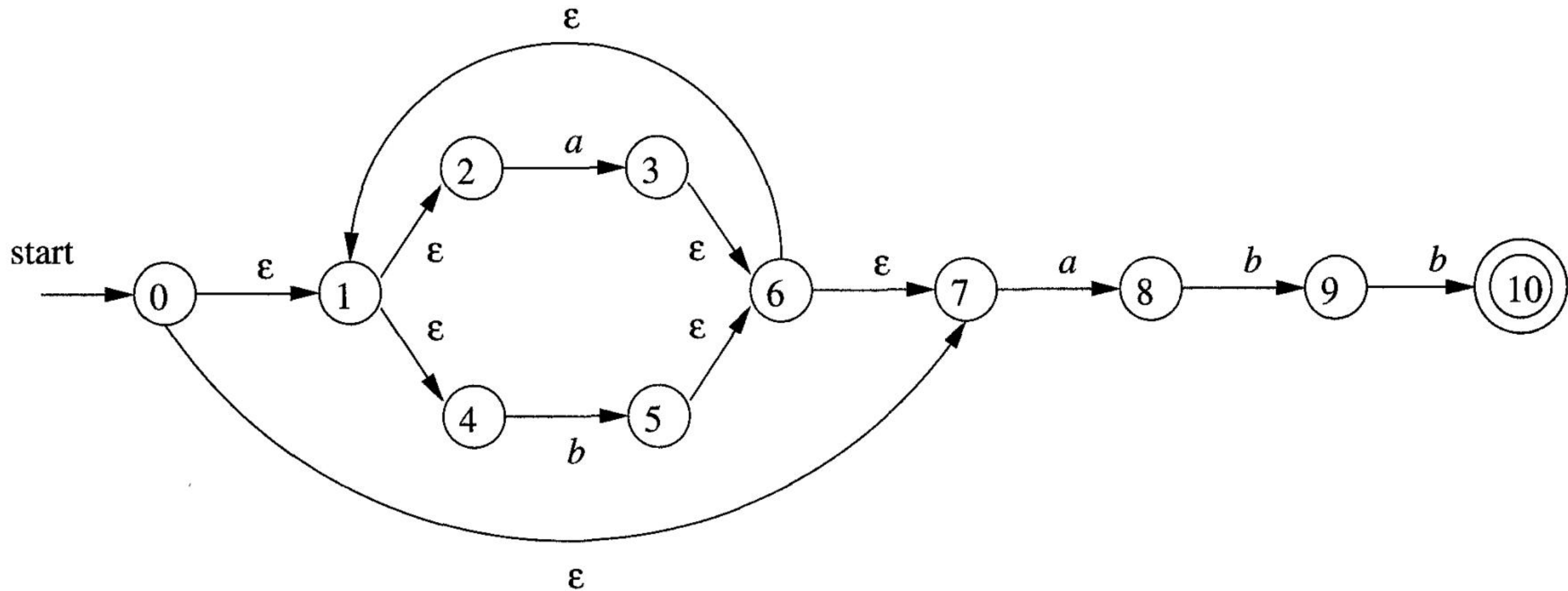
• Step 2: 处理组合REs (归纳规则)

- 对于表达式 $r = s|t$
 - 增加2个新状态, 4个新的 ϵ 转移
- 对于表达式 $r = st$
 - 无需增加新状态或转移
- 对于表达式 $r = s^*$
 - 增加2个新状态, 4个新的 ϵ 转移



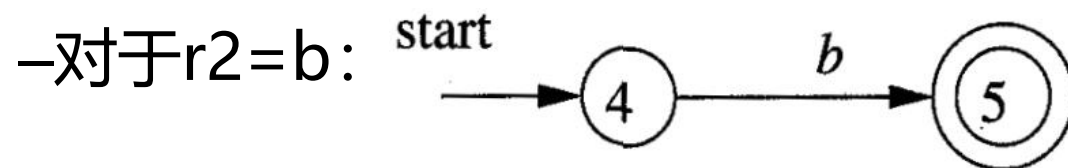
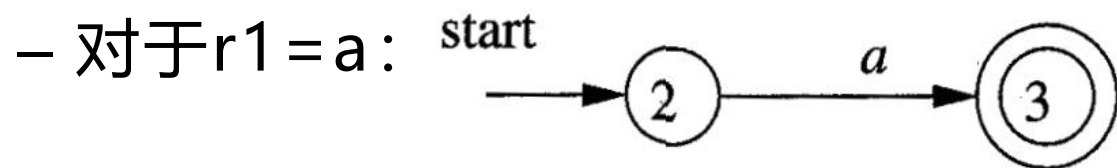
2. $RE \rightarrow NFA$

- 例：将正则表达式 $(a|b)^*abb$ 转成NFA
- 最终结果：

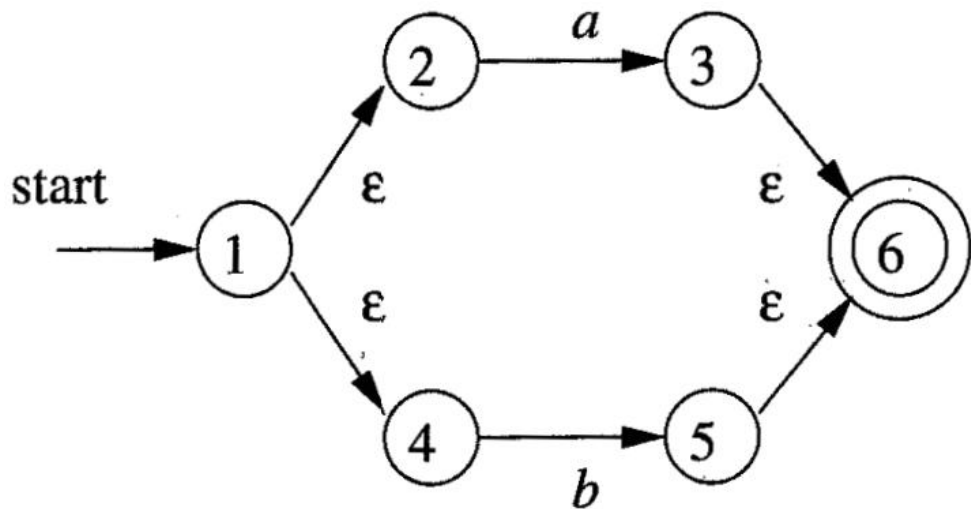


2. $RE \rightarrow NFA$

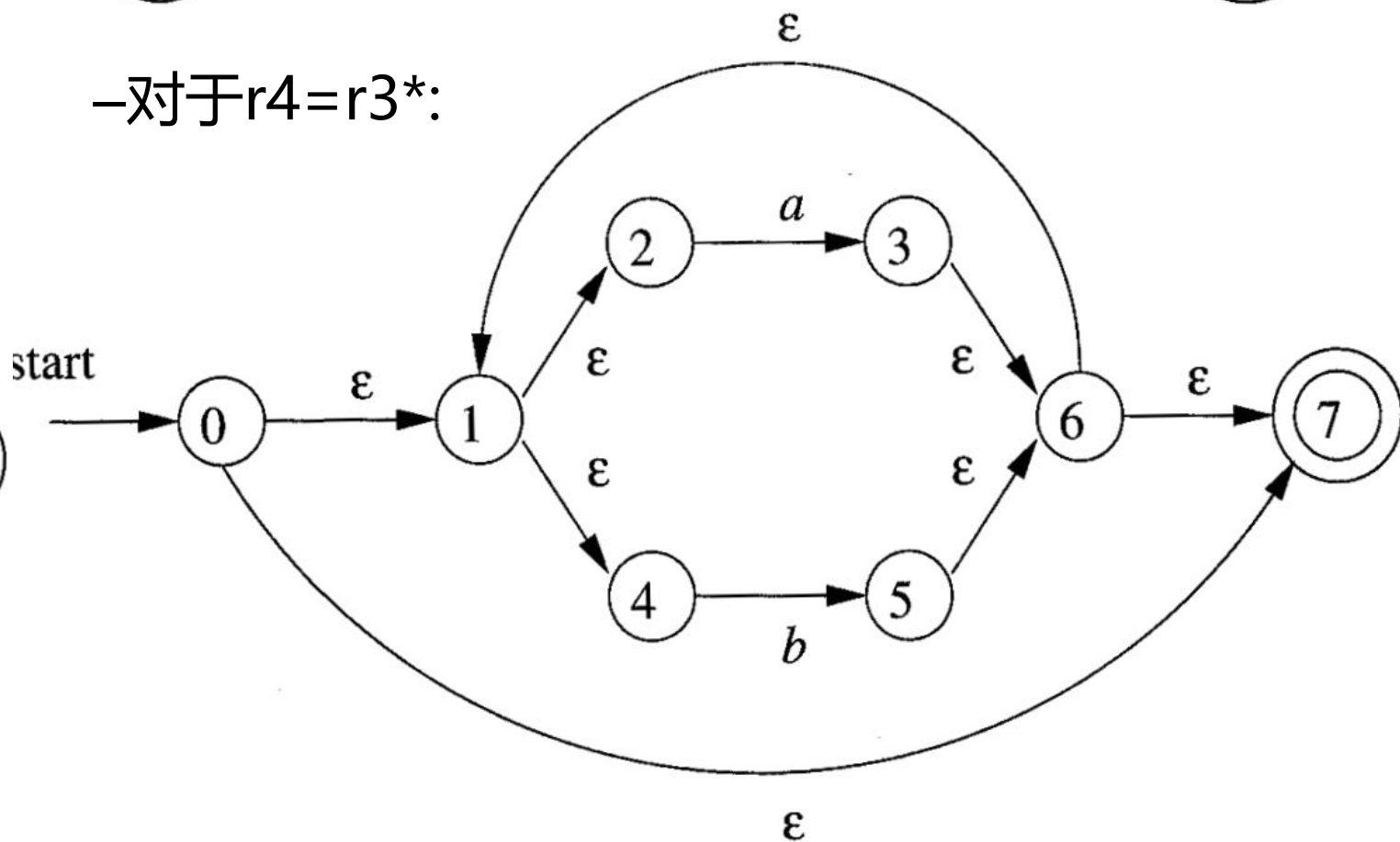
- 例：将正则表达式 $(a|b)^*abb$ 转成 NFA



- 对于 $r_3 = a|b$:



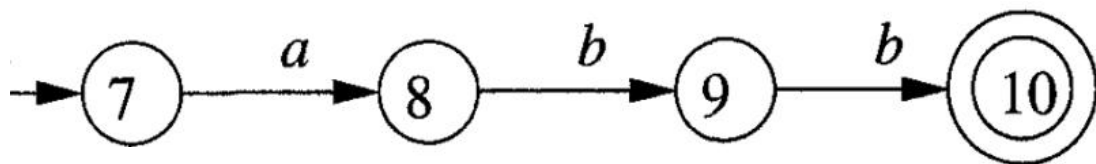
- 对于 $r_4 = r_3^*$:



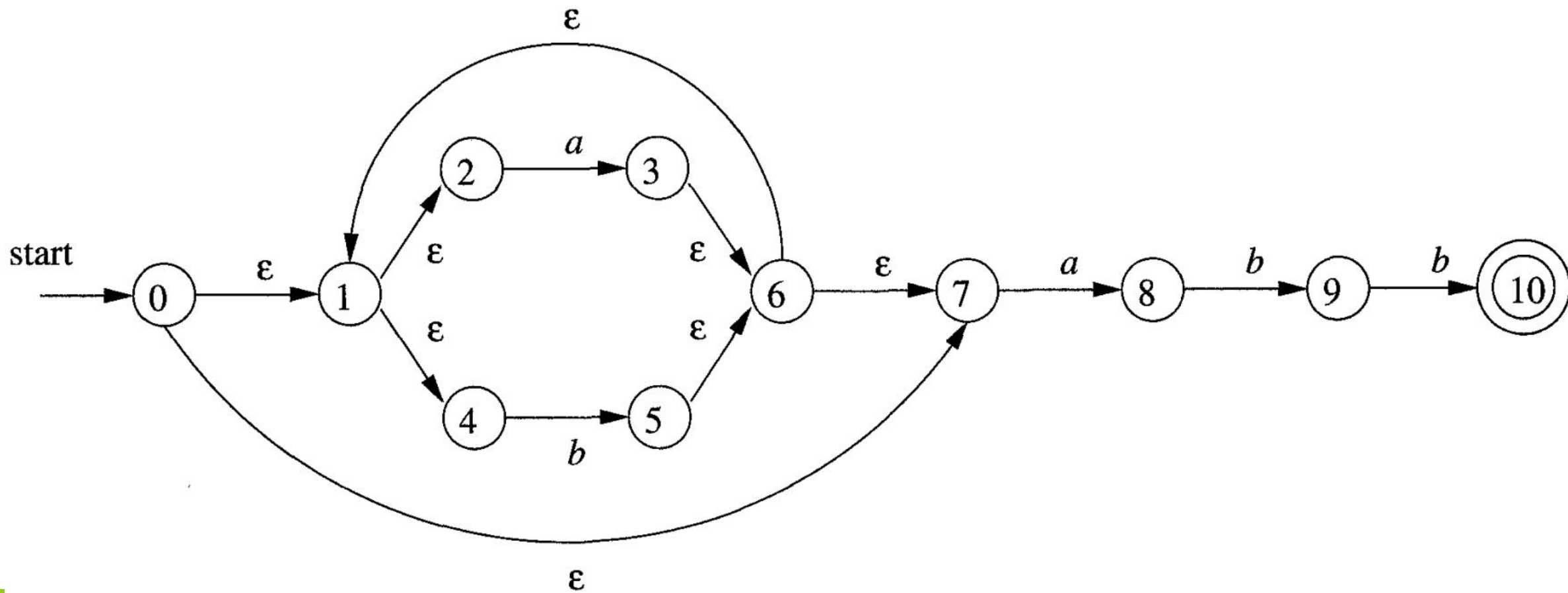
2. $RE \rightarrow NFA$

- 例：将正则表达式 $(a|b)^*abb$ 转成 NFA

- 对于 abb :



- 最终:

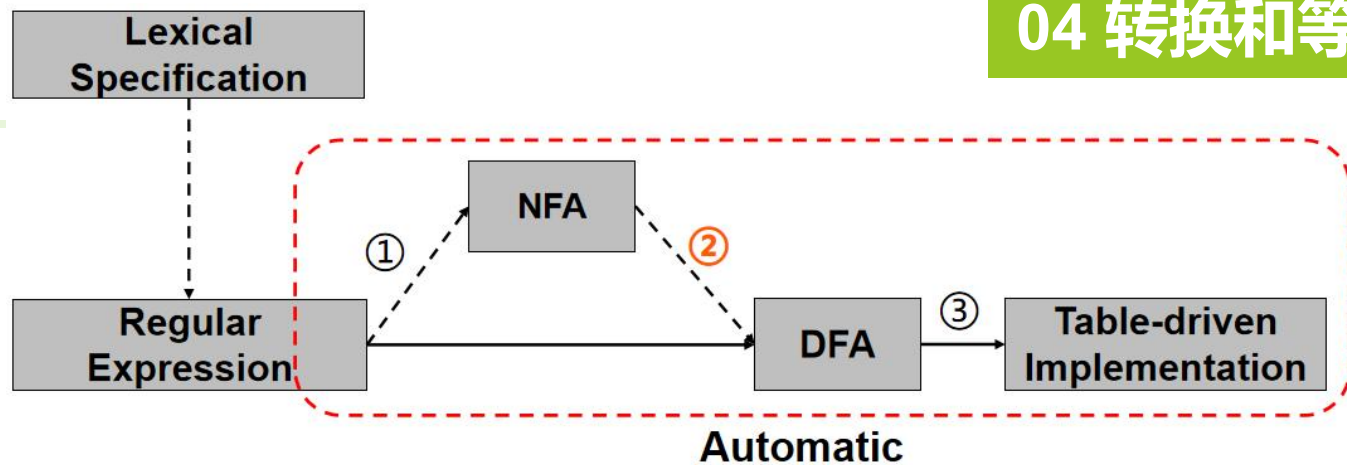


随堂练习(2)

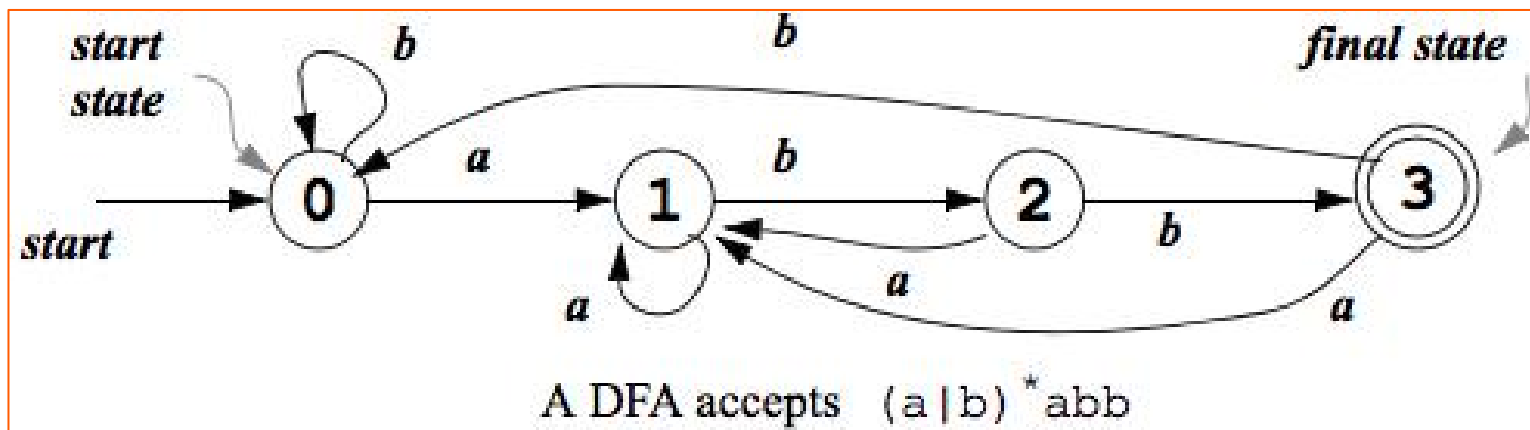
- 将正则表达式 $1^*(0|01)^*$ 转成NFA

3. NFA \rightarrow DFA

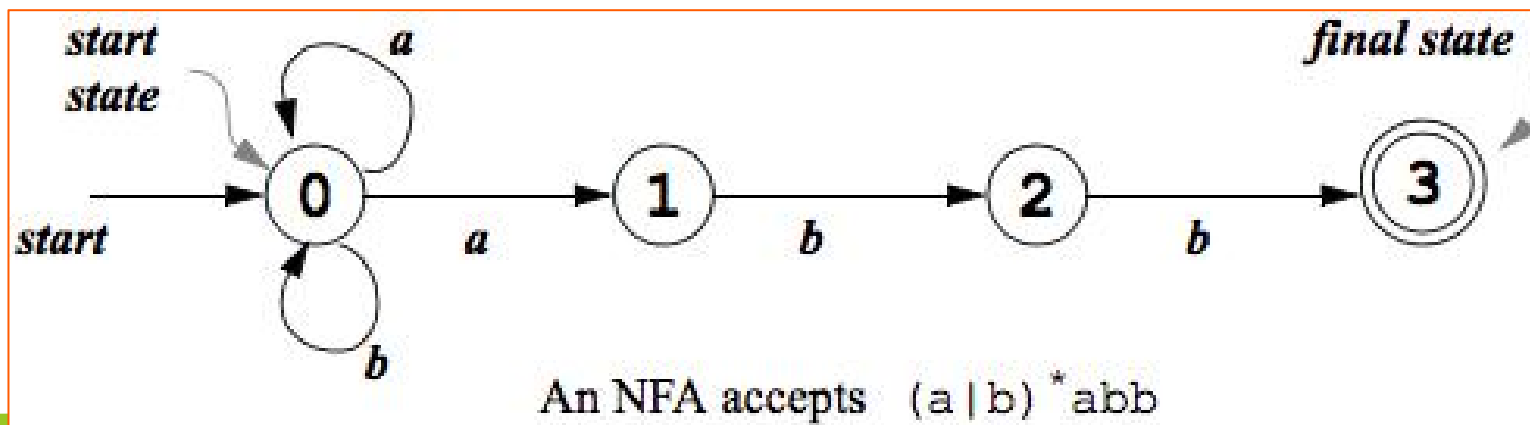
- 对于每个NFA M , 存在一个与其**等价**的DFA M'



- 但**NFA**的转移是不确定的, 更难以用机器模拟

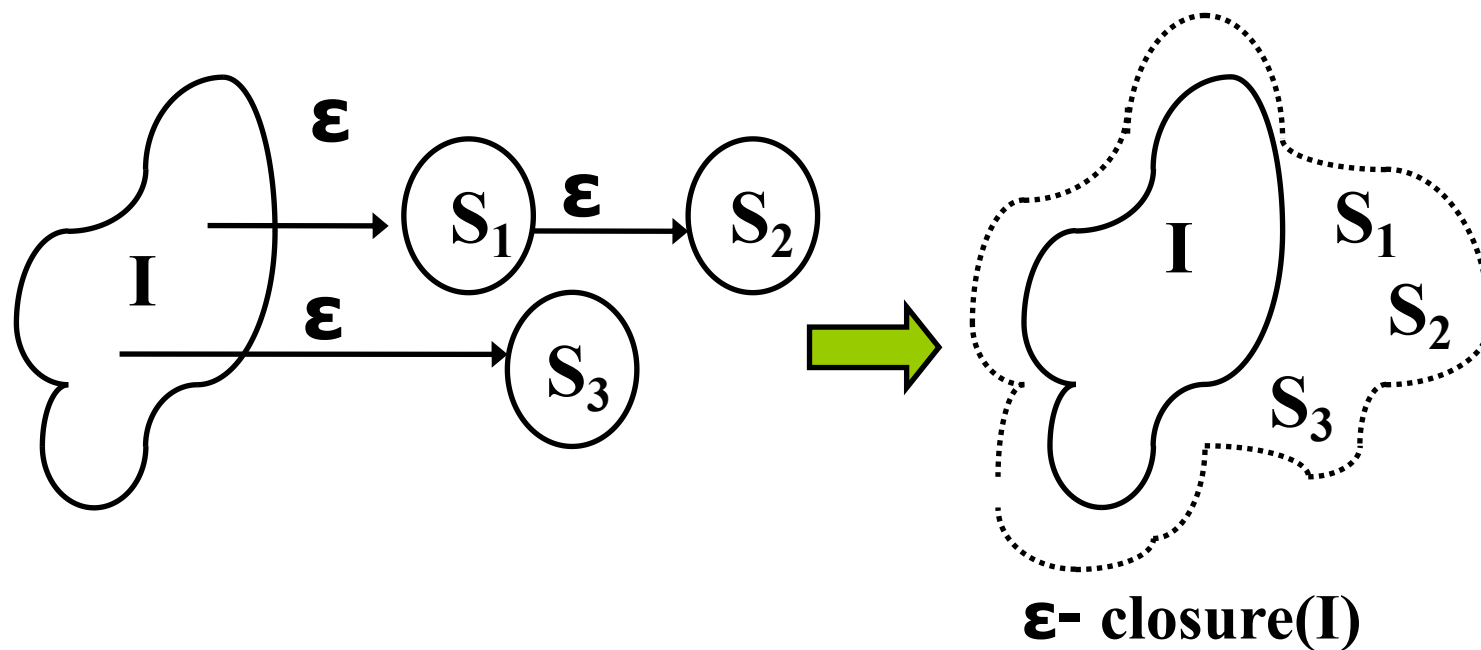


- 故须: NFA \rightarrow DFA



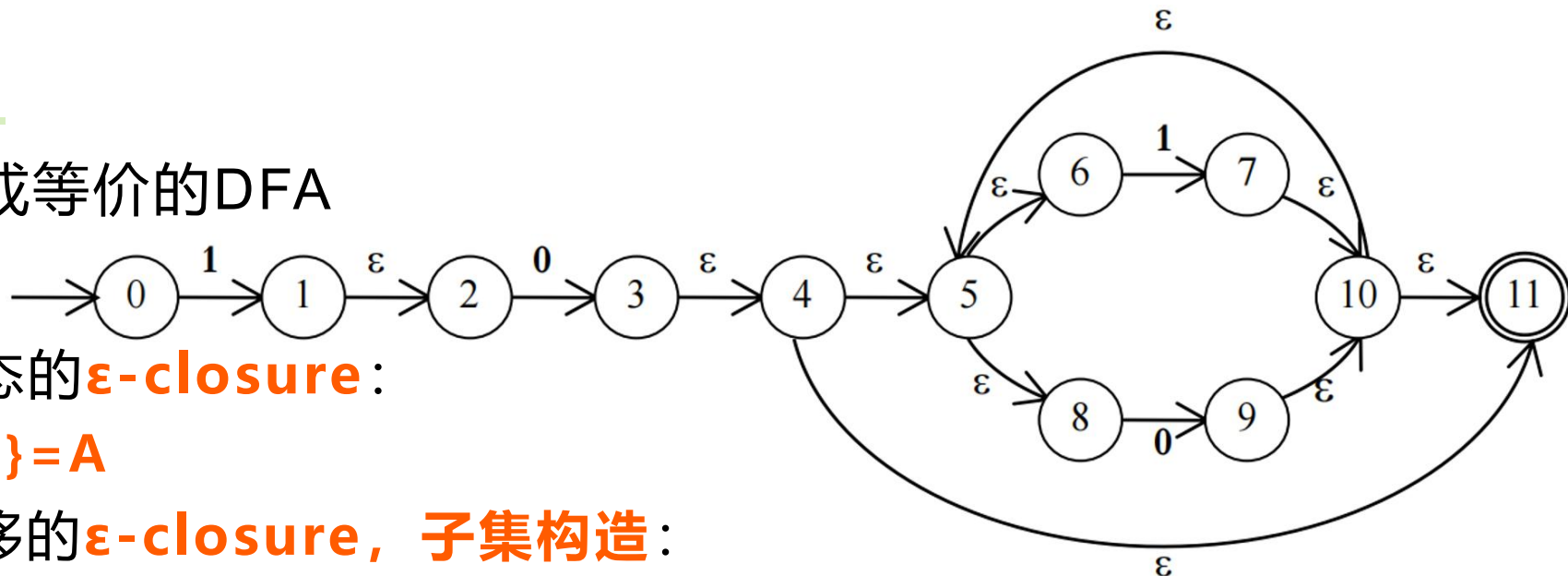
3. NFA \rightarrow DFA• ϵ -NFA \rightarrow DFA– ϵ 闭包构造 [ϵ -closure construction] 算法

- ✓ 状态集合 I 的 ϵ 闭包 ϵ -closure(I) 是状态集 I 中的所有状态 S 以及经任意条 ϵ 弧而能到达的状态的集合。



3. NFA \rightarrow DFA

- 例1：将右图 ϵ -NFA 转换成等价的 DFA



- **Step 1**：计算初始状态的 ϵ -closure：

✓ $\epsilon\text{-closure}(\{0\}) = \{0\} = A$

- **Step 2**：计算状态转移的 ϵ -closure，子集构造：

✓ $(A, 1) = \epsilon\text{-closure}(\{1\}) = \{1, 2\} = B$

✓ $(B, 0) = \epsilon\text{-closure}(\{3\}) = \{3, 4, 5, 6, 8, 11\} = C$

✓ $(C, 0) = \epsilon\text{-closure}(\{9\}) = \{5, 6, 8, 9, 10, 11\} = D$

✓ $(C, 1) = \epsilon\text{-closure}(\{7\}) = \{5, 6, 7, 8, 10, 11\} = E$

✓ $(D, 0) = \epsilon\text{-closure}(\{9\}) = D$

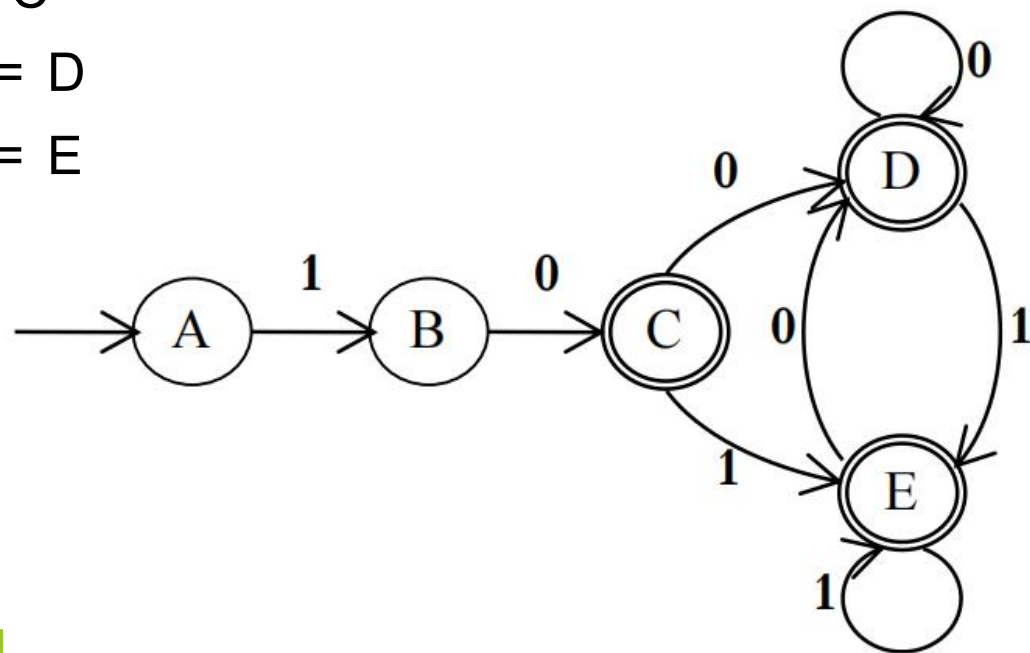
✓ $(D, 1) = \epsilon\text{-closure}(\{7\}) = E$

✓ $(E, 0) = \epsilon\text{-closure}(\{9\}) = D$

✓ $(E, 1) = \epsilon\text{-closure}(\{7\}) = E$

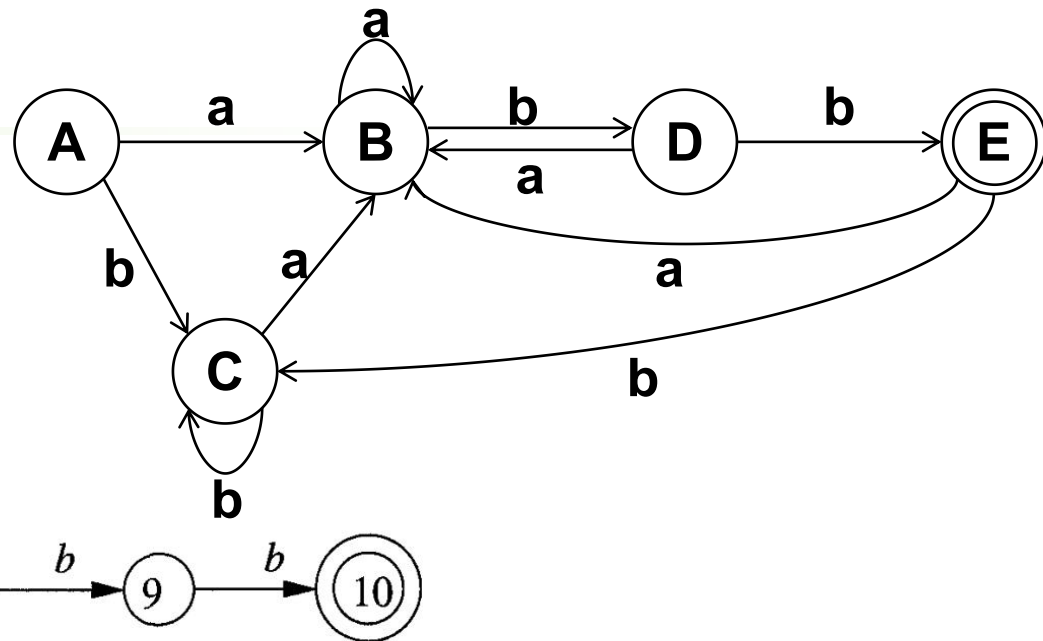
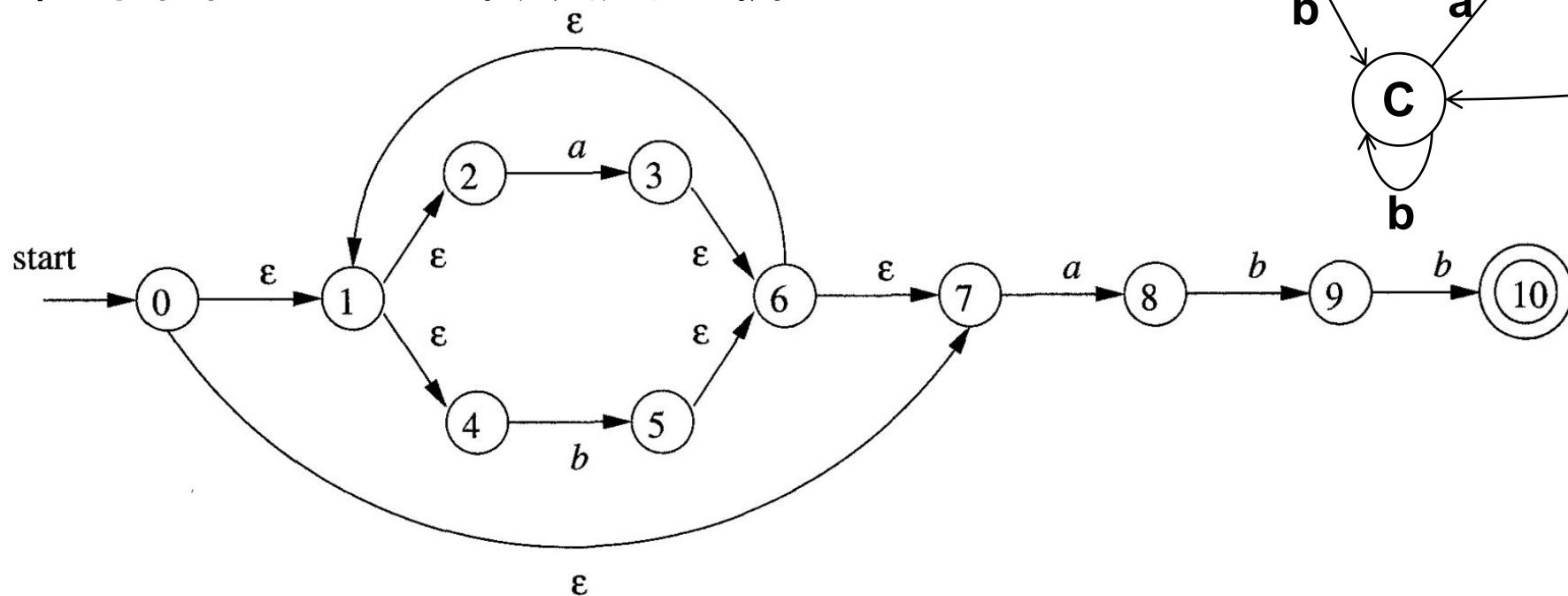
没有新的集合诞生 \rightarrow Step2 结束

- **Step 3**：根据新的状态集合画出 DFA



随堂练习(3)

- 将下图的 ϵ -NFA转换成等价的DFA



- ϵ -closure($\{0\}$) = $\{0, 1, 2, 4, 7\}$ = A
- (A, a) = ϵ -closure($\{3, 8\}$) = $\{1, 2, 3, 4, 6, 7, 8\}$ = B
- (A, b) = ϵ -closure($\{5\}$) = $\{1, 2, 4, 5, 6, 7\}$ = C
- (B, a) = ϵ -closure($\{3, 8\}$) = B
- (B, b) = ϵ -closure($\{5, 9\}$) = $\{1, 2, 4, 5, 6, 7, 9\}$ = D
- (C, a) = ϵ -closure($\{3, 8\}$) = B

- (C, b) = ϵ -closure($\{5\}$) = C
- (D, a) = ϵ -closure($\{3, 8\}$) = B
- (D, b) = ϵ -closure($\{5, 10\}$) = $\{1, 2, 4, 5, 6, 7, 10\}$ = E
- (E, a) = ϵ -closure($\{3, 8\}$) = B
- (E, b) = ϵ -closure($\{5\}$) = C
- 没有新的集合诞生-->结束**

3. NFA \rightarrow DFA

- NFA(不含 ϵ 转移) \rightarrow DFA

- 子集构造[subset construction]算法

- ✓ 让DFA的每个状态对应NFA的一个状态集合
 - ✓ 即DFA用它的一个状态记录在NFA读入一个输入符号后可能达到的所有状态

3. NFA \rightarrow DFA

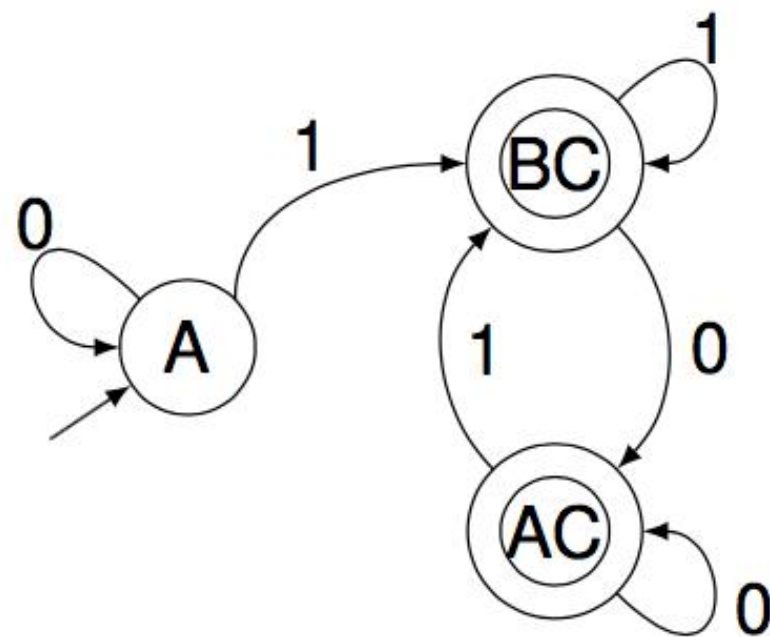
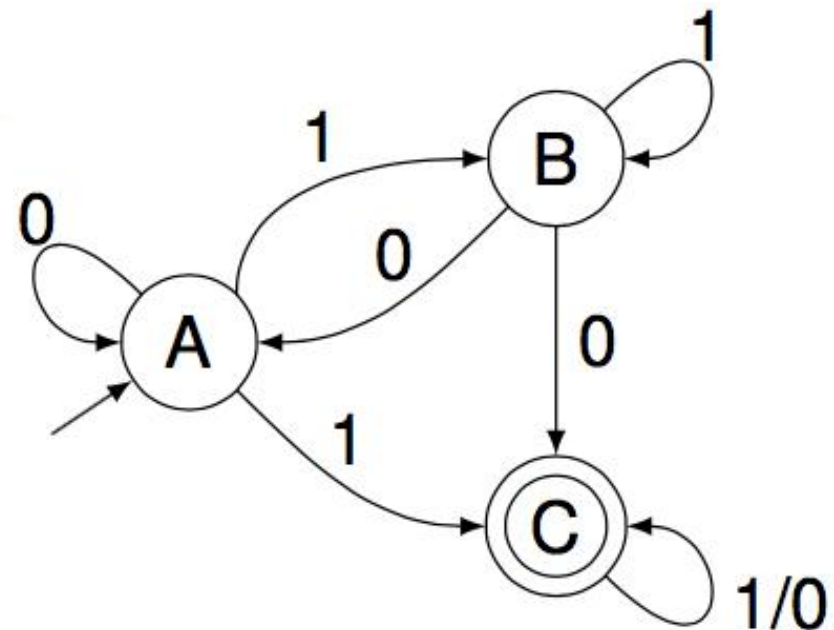
- 例2：将右图NFA(不含 ϵ 转移)转换成等价的DFA

- 可借助**状态转换矩阵**来进行
- Step 1: **子集构造**

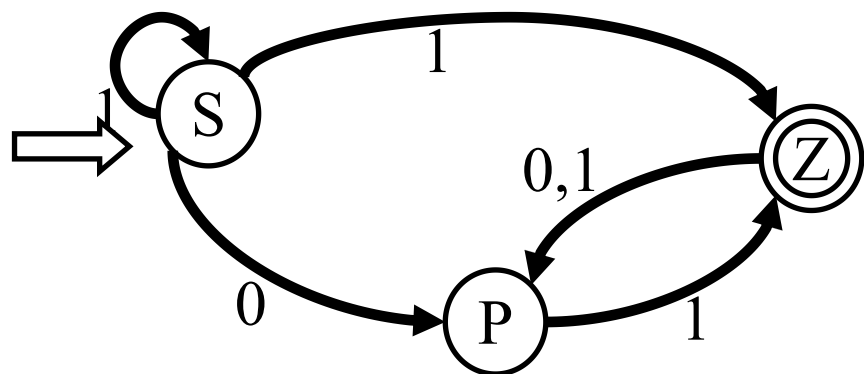
		alphabet	
state		0	1
	A	A	BC
	BC	AC	BC
	AC	AC	BC

没有新的集合诞生 \rightarrow Step1 结束

- Step 2: 根据新的状态集合画出DFA



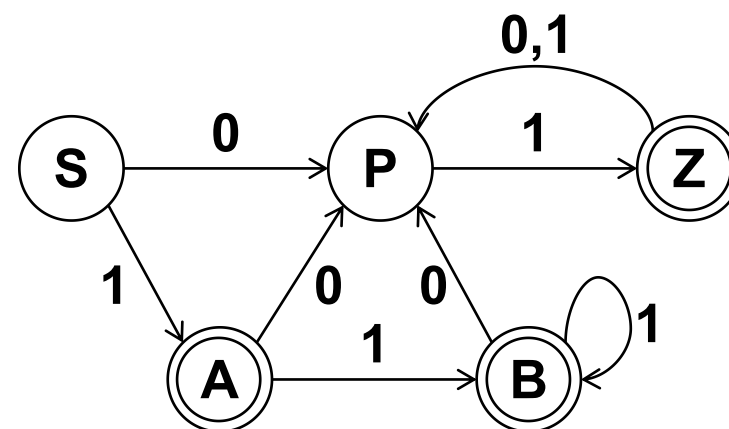
- 将下图的NFA(不含 ϵ 转移)转换成等价的DFA



令 $A=SZ$, $B=SZP$, 则有:

	0	1
S	P	SZ
P	-	Z
SZ	P	SZP
Z	P	P
SZP	P	SZP

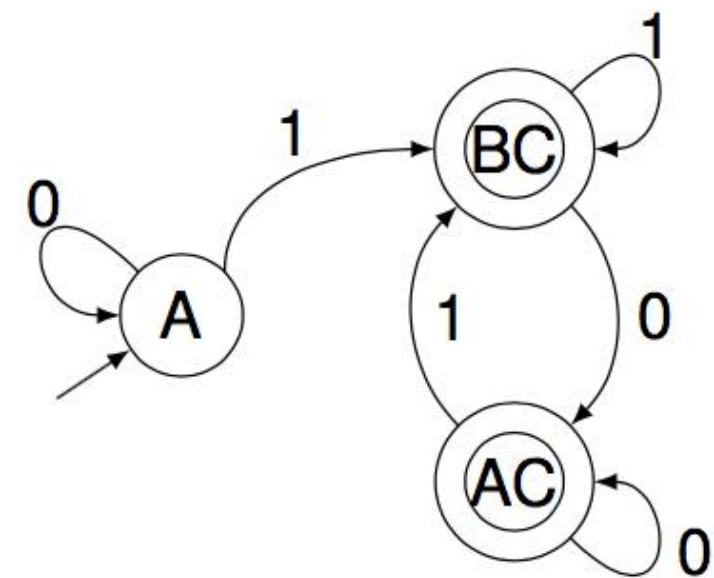
	0	1
S	P	A
P	-	Z
A	P	B
Z	P	P
B	P	B



4. DFA化简[Reduction]

• 任务：去掉多余状态，合并等价状态

- 多余状态：从开始状态出发无法到达的状态。
- 等价状态：两个状态s和t等价的条件是：
 - ✓ 1. **一致性条件**—状态s和t必须同为可接受状态或不可接受状态。
 - ✓ 2. **蔓延性条件**—对于所有输入符号，状态s和状态t必须转换到等价的状态里。

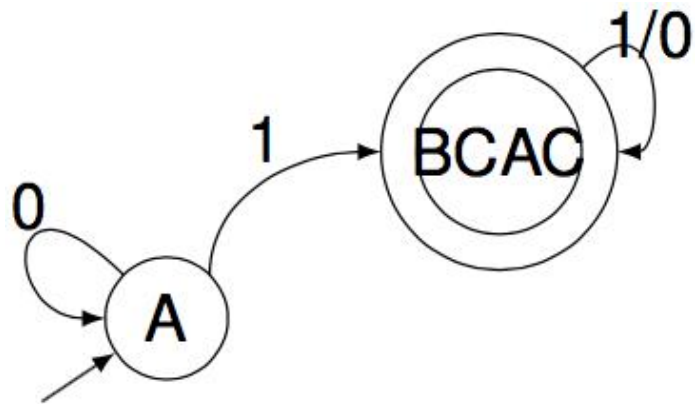


1. BC和AC**同为可接受状态**

2. BC和AC对所有输入符号**都转移到等价的状态里**

- BC on '0' → AC, AC on '0' → AC
- BC on '1' → BC, AC on '1' → BC

因此：BC和AC是**等价状态，可合并。**



4. DFA化简[Reduction]

• 例1：将右图DFA化简

– Step 1: 将状态分成非终态和终态集

✓ {S,A,B} {C,D,E,F}

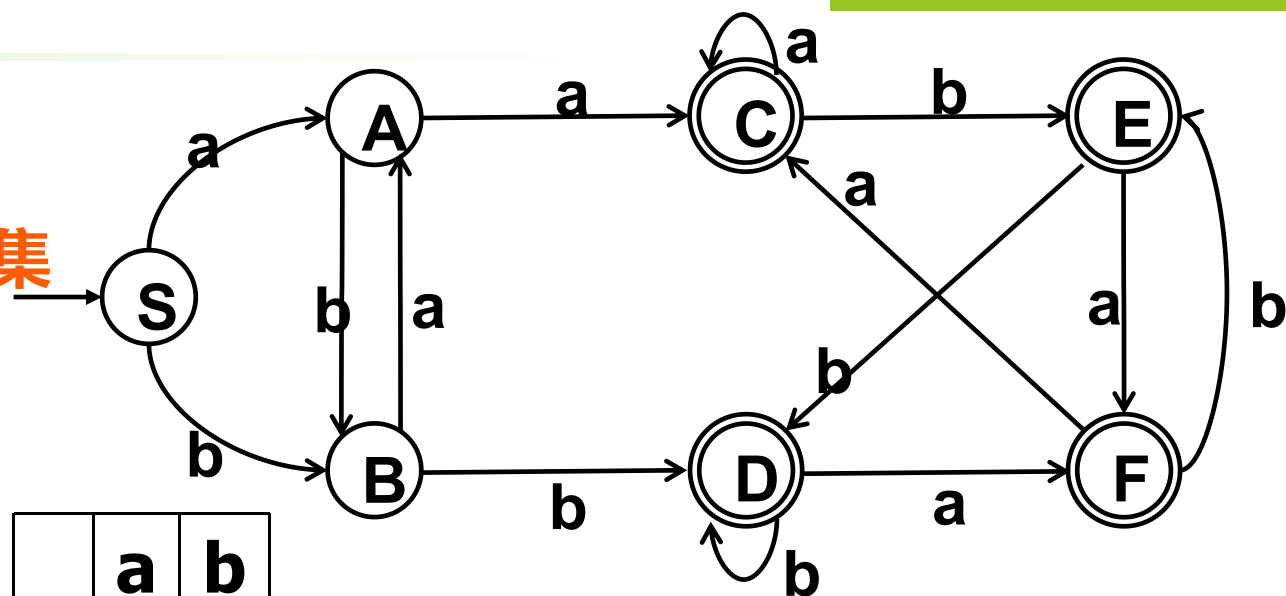
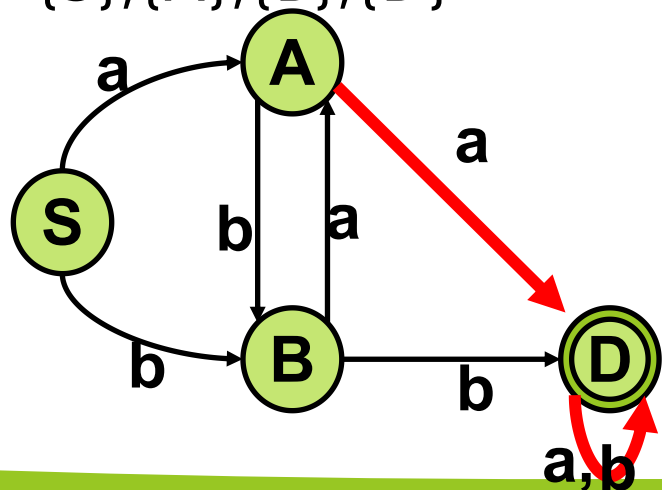
– Step 2: 寻找子集中不等价状态

✓ {S,A,B} => {S}{A,B} => {S}{A}{B}

✓ {C,D,E,F}

– Step 3: 令D代表{C,D,E,F}

✓ {S},{A},{B},{D}



	a	b
S	A	B
A	C	B
B	A	D
C	C	E
D	F	D
E	F	D
F	C	E

	a	b
S	A	B
A	C	B
B	A	D
CF	C	E
DE	F	D

	a	b
S	A	B
A	C	B
B	A	D
CFDE	CF	DE

4. DFA化简[Reduction]

- 例2：将右图DFA化简

- Step 1：将状态分成非终态和终态集

- ✓ $\{1,2,3,4\}$ $\{5,6,7\}$

- Step 2：寻找子集中不等价状态

- ✓ $\{1,2,3,4\} \Rightarrow \{1,2\}\{3,4\}$

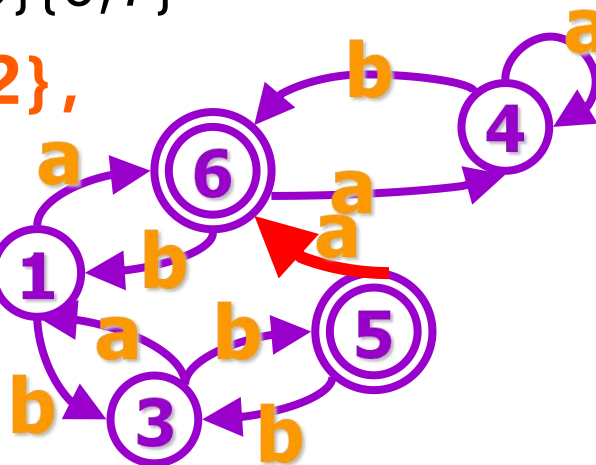
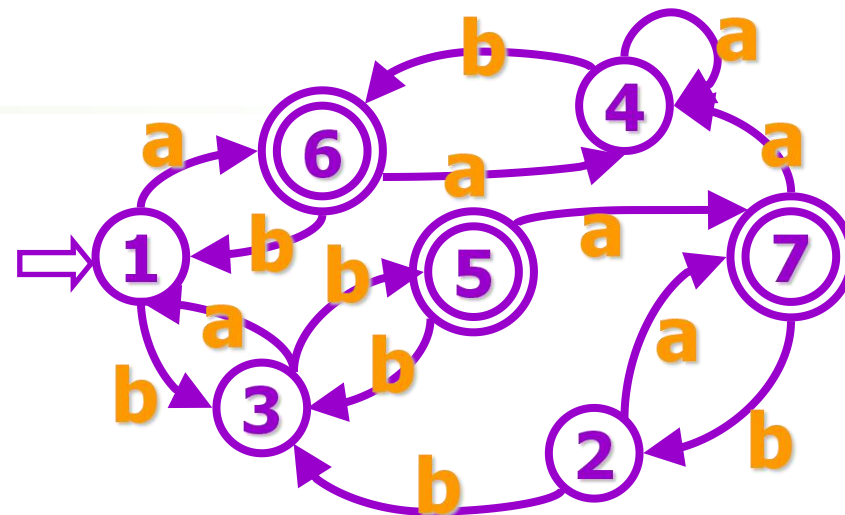
- ✓ $\{3,4\} \Rightarrow \{3\}\{4\}$

- ✓ $\{5,6,7\} \Rightarrow \{5\}\{6,7\}$

- ✓ 最终： $\{1,2\}\{3\}\{4\}\{5\}\{6,7\}$

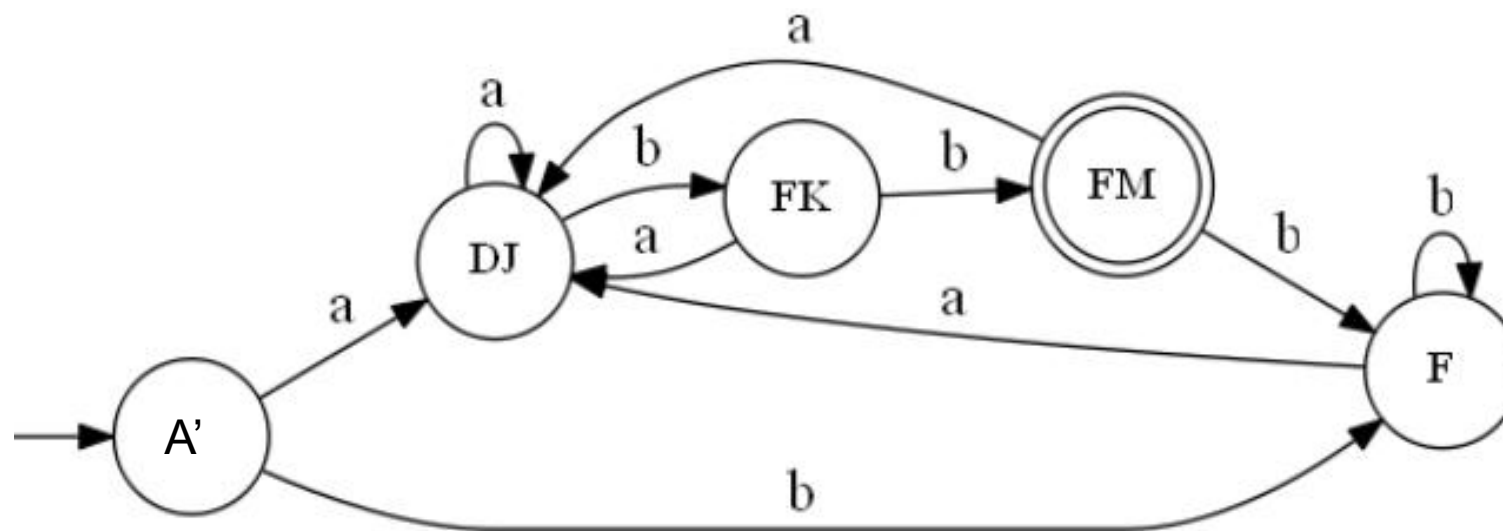
- Step 3：令1代表 $\{1,2\}$,
6代表 $\{6,7\}$

- ✓ $\{1\}, \{3\}, \{4\}, \{5\}, \{6\}$



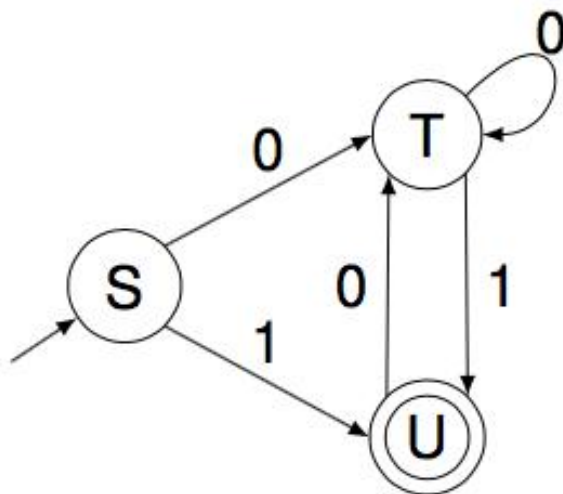
	a	b		a	b		a	b		a	b
1	6	3	1	6	3	1	6	3	1	6	3
2	7	3	2	7	3	2	7	3	2	7	3
3	1	5	3	1	5	3	1	5	3	1	5
4	4	6	4	4	6	4	4	6	4	4	6
5	7	3	5	7	3	5	7	3	5	7	3
6	4	1	6	4	1	6	4	1	6	4	1
7	4	2	7	4	2	7	4	2	7	4	2

- 将下图DFA化简



5. DFA → Table-drive Implementation

- DFA可转成Table-drive Code



alphabet →

state ↓	0	1
S	T	U
T	T	U
U	T	x

Lexical Specification

Regular Expression

①

NFA

②

DFA

③

Table-driven Implementation

Automatic

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

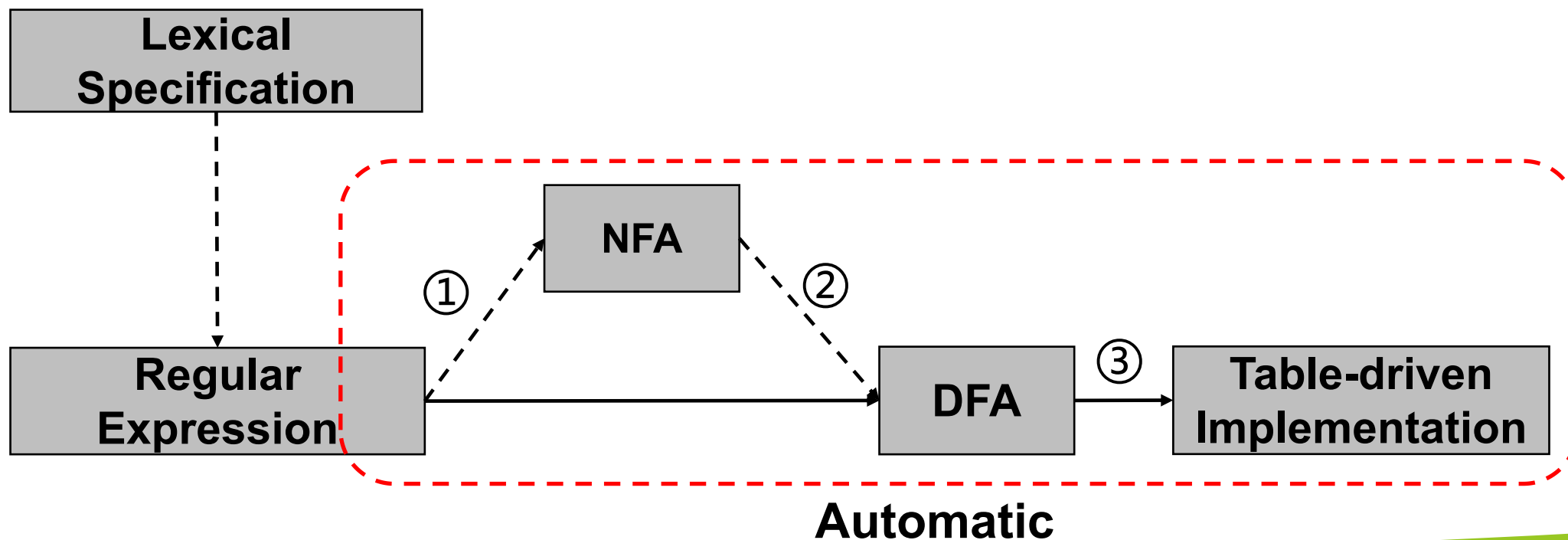
转换流程[Conversion Flow]

• 流程：RE→NFA→DFA→Table-drive Implementation

① RE→NFA

② NFA→DFA

③ DFA→Table-drive Implementation



6. 一些讨论

- 关于Table-drive Implementation

- 表格是一种高效实现[efficient]

- ✓ 仅需**有穷空间** $O(S \times \Sigma)$

- 转换表的Size

- ✓ 仅需**有穷时间** $O(\text{input length})$

- 状态转换的个数

- 表格实现的优劣

- ✓ pros: 在给定的状态和输入下能轻松找到转换

- ✓ cons: 当输入字母很大时, 会占用大量空间, 但大多数状态对大多数输入符号没有任何移动

6. 一些讨论

- 关于**空间复杂度[Space Complexity]**
 - NFA在任何时刻可能有多个状态
 - DFA:
 - ✓ 如果NFA有N个状态，则DFA一定在这N个状态的某个子集中
 - ✓ 非空子集： $2^N - 1$
 - ✓ 空间复杂度： $O(2^N)$ ，其中N为NFA中的原始状态个数

6. 一些讨论

- 关于**时间复杂度[Time Complexity]**

- DFA

- ✓ 需要 **$O(|X|)$** 步，其中 $|X|$ 是输入长度
 - ✓ **每一步花费 $O(1)$** 常数时间
 - 若当前状态是 S 且输入为 c ，则读表 $T[S,c]$ ，更新当前状态为 $T[S,c]$
 - ✓ 时间复杂度为 **$O(|X|)$**

- NFA

- ✓ 需要 **$O(|X|)$** 步，其中 $|X|$ 是输入长度
 - ✓ **每一步花费 $O(N^2)$** 时间，其中 N 为状态的个数
 - 当前状态是一组潜在状态，最多 N 个
 - 对于输入 c ，必须合并所有 $T[S_{\text{potential}}, c]$ ，最多 N 次，每次合并操作花费 $O(N)$ 时间
 - ✓ 时间复杂度为 **$O(|X| * N^2)$**

CONTENTS

目 录

01

概述

Introduction

02

词法规范

Lexical
Specification

03

有穷自动机

Finite
Automata

04

转换和等价

Transformation
and Equivalence

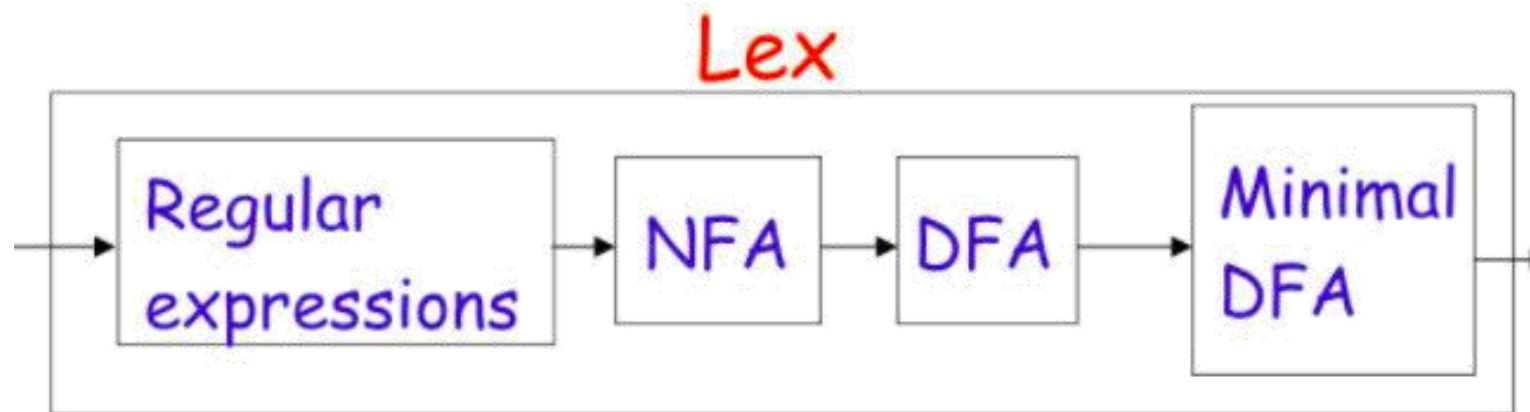
05

词法分析实践

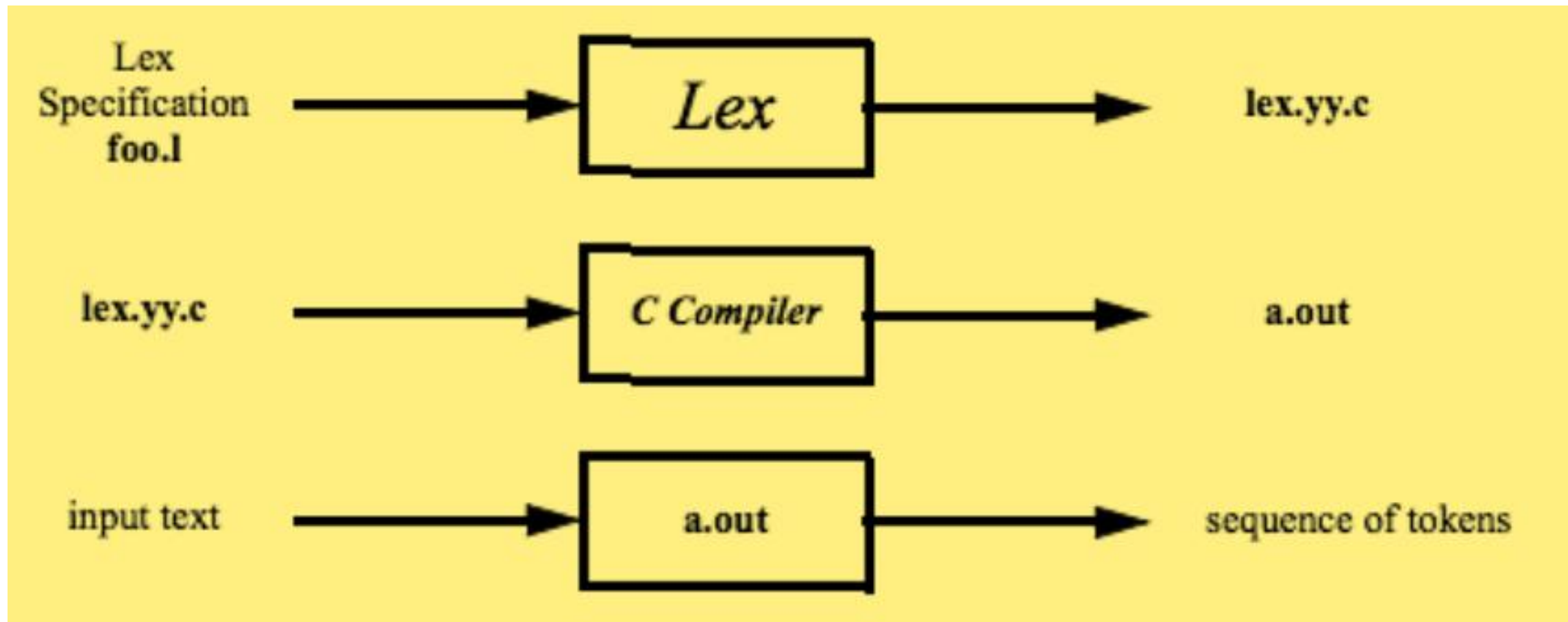
Lexical Analysis
in Practice

1. 实际实现

- Lex: RE \rightarrow NFA \rightarrow DFA \rightarrow Table



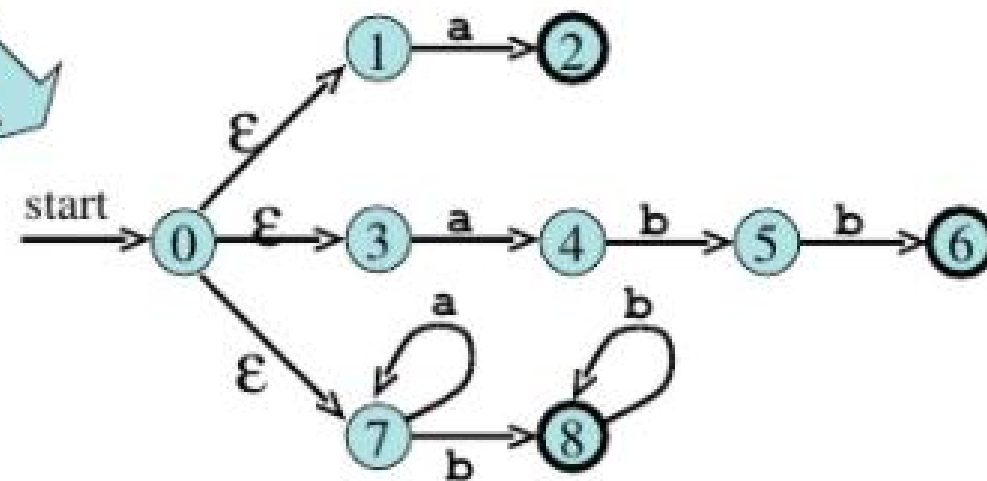
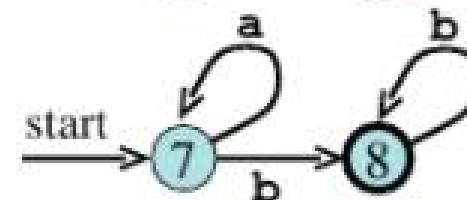
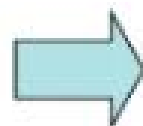
- 大多数其他自动词法分析器也选择DFA而不是NFA
 - 用空间换取速度[Trade off space for speed]



2. Lex

- 例：假设现有3种模式，对应3个NFA
 - 将3个NFA组合成1个NFA
 - 添加开始状态和 ϵ 转移

a $\{ action_1 \}$
abb $\{ action_2 \}$
a*b+ $\{ action_3 \}$



2. Lex

```
ptn1    a
ptn2    abb
ptn3    a*b+
```

```
%%
```

```
{ptn1} { printf("\n<%s, %s>", "ptn1", yytext); }
{ptn2} { printf("\n<%s, %s>", "ptn2", yytext); }
{ptn3} { printf("\n<%s, %s>", "ptn3", yytext); }
```

```
%%
```

```
int main(){
    yylex();
    return 0;
}
```

\$flex lex.l

\$clang lex.yy.c -o mylex -ll

```
[root@aa51dde06c76:~/test# echo "aaba" | ./mylex
```

```
<ptn3, aab>
<ptn1, a>
```

```
[root@aa51dde06c76:~/test# echo "abba" | ./mylex
```

```
<ptn2, abb>
<ptn1, a>
```


• NFA

– 输入: **aaba**

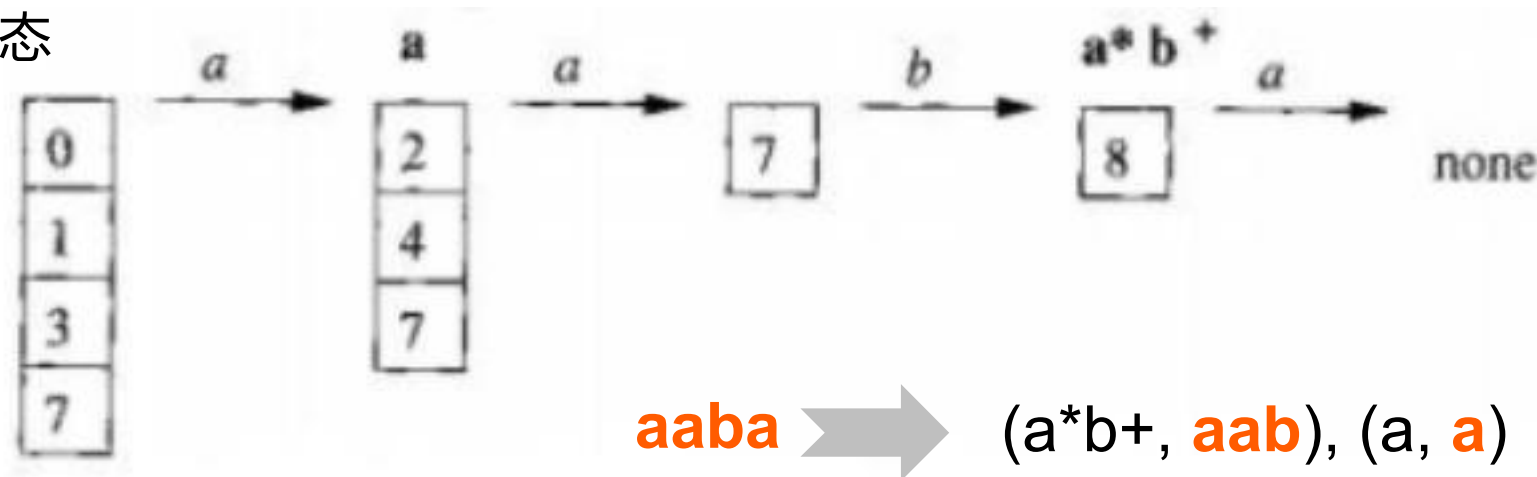
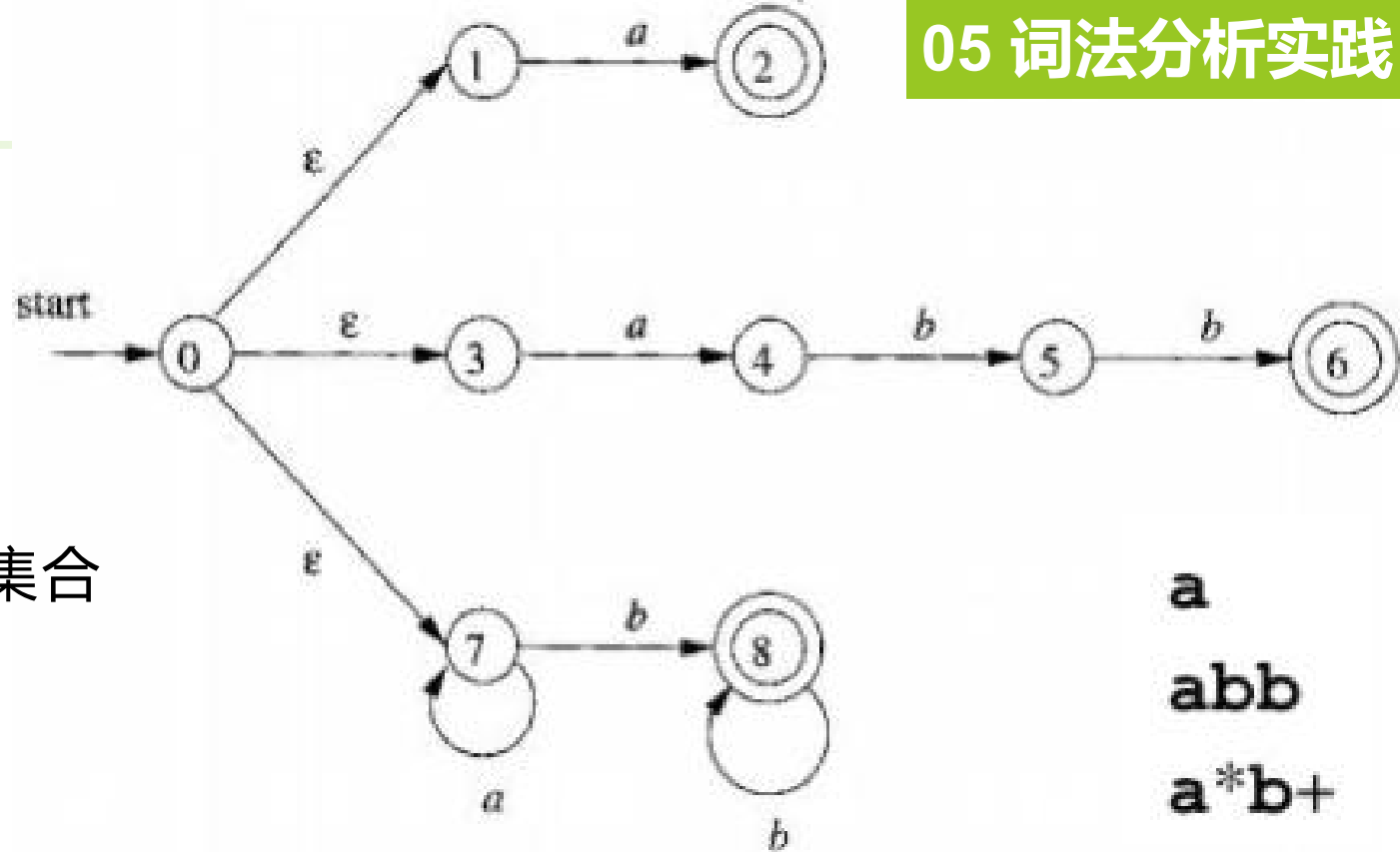
✓ ϵ -closure(0) = {0, 1, 3, 7}

✓ 寻找一组**包含接受状态**的状态集合

✓ 状态8: a^*b^+ 匹配成功

- **aab**即为词素

- 继续寻找接受状态



2. Lex

- DFA

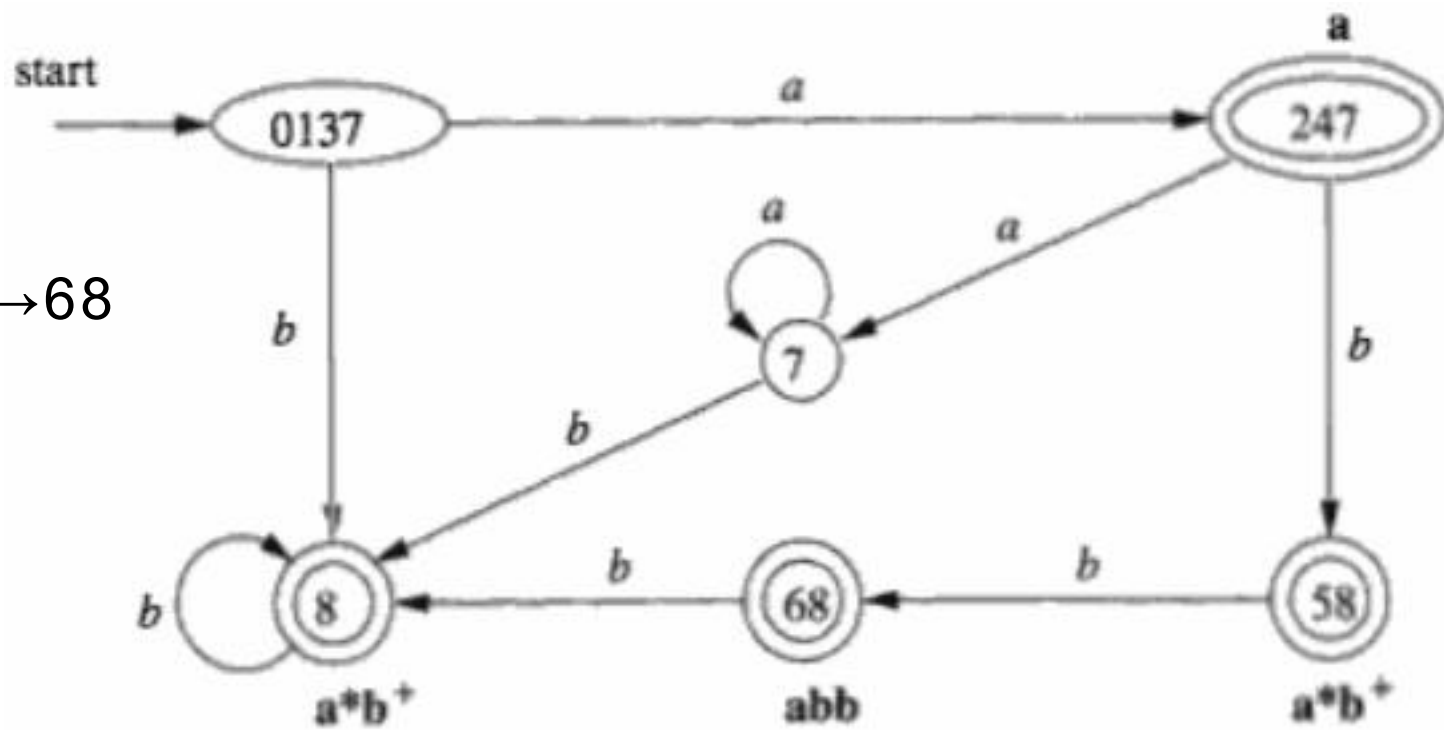
- 输入: **abba**

- ✓ 状态序列: $0137 \rightarrow 247 \rightarrow 58 \rightarrow 68$

- ✓ 状态68: **abb**匹配成功

- **abb**即为词素

- 继续寻找接受状态



a

abb

a^*b^+

2. Lex

- 若有多种匹配的可能性?

- 寻找最长匹配

✓ 例：输入 **aabbb**，符合pattern2和pattern3，取最长匹配——pattern3

- 先出现的规则优先级更高

✓ 例：输入 **abb**，符合pattern2和pattern3，取先出现的pattern2

a { *action*₁ }

abb { *action*₂ }

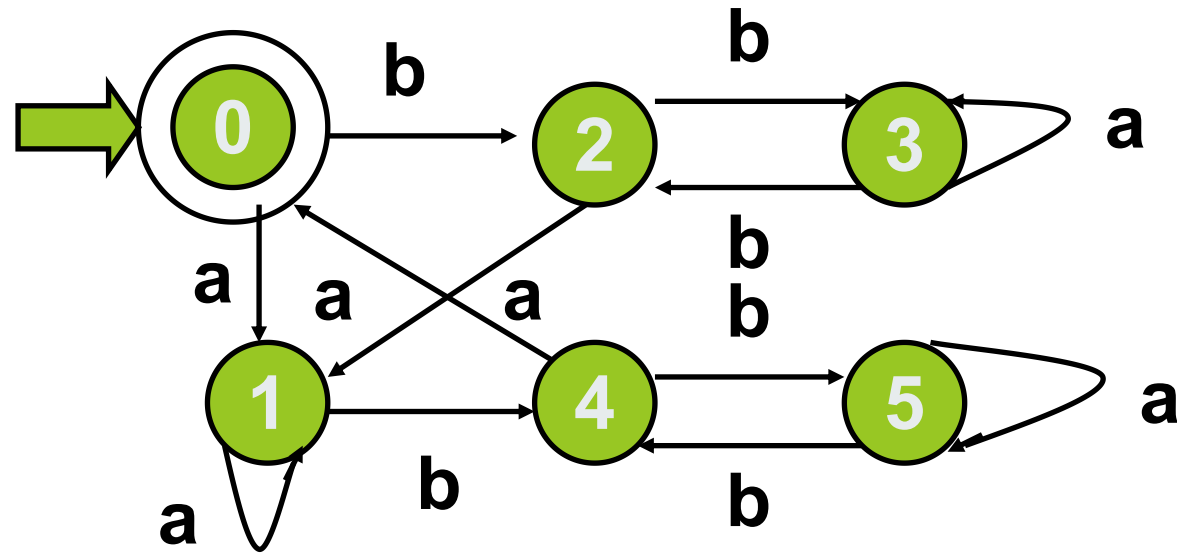
a*b+ { *action*₃ }

2. Lex

- 如何匹配关键字[keywords]?
 - 方法1：为关键字创建正则表达式，并将它们放在标识符的正则表达式之前，让他们具有更高的优先级
 - ✓ 会导致更臃肿的FA
 - 方法2：使用相同的正则表达式识别关键字和标识符，但使用特殊的**关键字表[keyword table]**进行区分
 - ✓ 会导致更精简的FA
 - ✓ 但是需要额外的表查找
 - **方法2更常用**

第三章课后作业

- 将下图DFA最小化



第三章课后作业

- 提交要求：
 - 文件命名：学号-姓名-第三章作业；
 - 文件格式：.pdf文件；
 - 手写版、电子版均可；若为手写版，则拍照后转成pdf提交，但**须注意将照片旋转为正常角度，且去除照片中的多余信息**；电子版如word等转成pdf提交；
 - 提交到超算习堂（第三章作业）处；
 - 提交ddl：**3月20日晚上12:00**；
 - **重要提示：不得抄袭！**