



重庆大学
CHONGQING UNIVERSITY

第二讲 形式语言与自动机基础

重庆大学 计算机学院 张敏





本讲要点

- 语言和文法
- 有限自动机
- 自动机、文法、表达式等价性



讨论

要开发一个编译器，第一步应该做什么？

定义程序设计语言



思考

如何定义程序设计语言？

如何定义自然语言？

- 汉语--所有符合汉语语法的句子的全体

程序设计语言--所有该语言的程序的全体

语言中哪些部分是有穷的？哪些部分是无穷的？

- 程序是有穷的吗？

如何表示一种语言中的每个成分？

- 枚举？
- 以有限表示无限？



2.1 语言 and 文法

如果语言是有穷的，可以将句子逐一列出来表示

如果语言是无穷的，找出语言的有穷表示：

- ✓ 生成方式：文法
- ✓ 识别方式：自动机



2.1.1 语言定义相关概念

每种语言具有两个可识别的特性：

- 语言的形式
- 与形式相关联的意义

语言的实例相关联的意义：

- 该句子的创立者所想要表示的意义
- 接收者所检验到的意义

程序设计语言

- 每个程序构成的规律
- 每个程序的含义
- 每个程序和使用者的关系



2.1.2 形式语言

语法意义下的语言称作形式语言。
形式语言抽象地定义为一个数学系统。



2.1.2.1 字母表和符号串

字母表

符号的非空有限集合

典型的符号是字母、数字、各种标点和运算符等。

英语字母表: $\Sigma = \{a, b, c, \dots, x, y, z, A, B, \dots, Y, Z\}$

机器语言字母表: $\Sigma = \{0, 1\}$

符号串

定义在某一字母表上, 由该字母表中的符号组成的有限符号序列

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

$\epsilon, a, b, aa, ab, aabba \dots$ 都是 Σ 上的符号串

1. 空符号串 ϵ (没有符号的符号串) 是 Σ 上的符号串
2. 若 α 是 Σ 上的符号串, a 是 Σ 的元素, 则 αa 是 Σ 上的符号串
3. y 是 Σ 上的符号串, 当且仅当它可以由 1 和 2 导出



2.1.2.1 字母表和符号串

长度：符号串 α 的长度是指 α 中出现的符号的个数，记作 $|\alpha|$ 。空串的长度为0，常用 ε 表示。

前缀：符号串 α 的前缀是指从符号串 α 的末尾删除0个或多个符号后得到的符号串。如：univ 是 university 的前缀。

后缀：符号串 α 的后缀是指从符号串 α 的开头删除0个或多个符号后得到的符号串。如：sity 是 university 的后缀。

子串：符号串 α 的子串是指删除了 α 的前缀和/或后缀后得到的符号串。如：ver 是 university 的子串。

真前缀、真后缀、真子串：如果非空符号串 β 是 α 的前缀、后缀或子串，并且 $\beta \neq \alpha$ ，则称 β 是 α 的真前缀、真后缀、或真子串。

子序列：符号串 α 的子序列是指从 α 中删除0个或多个符号(这些符号可以是不连续的)后得到的符号串。如：nvst是university中的子序列



2.1.2.2 语言

语言：某一确定字母表上的符号串的集合。空集 ϕ ，集合 $\{\epsilon\}$ 也是符合此定义的语言。



2.1.2.3 语言的运算

假设 L 和 M 表示两个语言

L和M的并记作 $L \cup M$: $L \cup M = \{s | s \in L \text{ 或 } s \in M\}$

L和M的连接记作 LM : $LM = \{st | s \in L \text{ 并且 } t \in M\}$

L的闭包记作 L^* : 即L的0次或若干次连接。

$$L^* = \bigcup_{i=0}^{\infty} L^i = L^0 \cup L^1 \cup L^2 \cup L^3 \cup \dots$$

L的正闭包记作 L^+ : 即L的1次或若干次连接。

$$L^+ = \bigcup_{i=1}^{\infty} L^i = L^1 \cup L^2 \cup L^3 \cup L^4 \cup \dots$$



2.1.2.3 语言的运算

$L = \{A, B, \dots, Z, a, b, \dots, z\}$, $D = \{0, 1, \dots, 9\}$

可以把L和D看作是字母表

可以把L和D看作是语言

语言	描述
$L \cup D$	全部字母和数字的集合
LD	由一个字母后跟一个数字组成的所有符号串的集合
L^4	由4个字母组成的所有符号串的集合
L^*	由字母组成的所有符号串（包括 ϵ ）的集合
$L(L \cup D)^*$	以字母开头，后跟字母、数字组成的所有符号串的集合
D^+	由一个或若干个数字组成的所有符号串的集合



2.1.2.3 语言的运算

例: $\Sigma = \{a, b\}$

$\Sigma^* = \{\epsilon, a, b, aa, ab, ba, bb, aaa, aab, \dots\}$

$\Sigma^+ = \{a, b, aa, ab, ba, bb, aaa, aab, \dots\}$

➤ 集合 $\{ab, aabb, aaabbb, \dots, a^n b^n, \dots\}$ 或 $\{w | w \in \Sigma^* \text{ 且 } w = a^n b^n, n \geq 1\}$ 为字母表 Σ 上的一个语言。

➤ 集合 $\{a, aa, aaa, \dots\}$ 或 $\{w | w \in \Sigma^* \text{ 且 } w = a^n, n \geq 1\}$ 为字母表 Σ 上的一个语言。

➤ $\{\epsilon\}$ 是字母表 Σ 上的一个语言。 Φ 即 $\{\}$ 是字母表 Σ 上的一个语言。



2.1.3 文法及其形式定义

- **文法**：所谓文法就是描述语言的语法结构的形式规则。
- 任何一个文法都可以表示为一个**四元组** $G=(V_T, V_N, S, P)$
 - V_T 是一个非空的有限集合，它的每个元素称为**终结符号**。
 - V_N 是一个非空的有限集合，它的每个元素称为**非终结符号**。
$$V_T \cap V_N = \phi$$
 - S 是一个特殊的非终结符号，称为文法的**开始符号**。
 - P 是一个非空的有限集合，它的每个元素称为**产生式**。



2.1.3.1 产生式的定义

规则 (产生式或生成式), 是形如 $\alpha \rightarrow \beta$ 或 $\alpha ::= \beta$ 的 (α, β) 有序对, 且 $\alpha \in V^+$, $\beta \in V^*$ 。

➤ α 称为规则的**左部**(或生成式的左部)。

➤ β 称为规则的**右部**(或生成式的右部)。

左部相同的产生式 $\alpha \rightarrow \beta_1$ 、 $\alpha \rightarrow \beta_2$ 、.....、 $\alpha \rightarrow \beta_n$ 可以缩写为 $\alpha \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$

"|" 表示 “或”, 每个 $\beta_i (i=1, 2, \dots, n)$ 称为 α 的一个候选式

BNF (Backus-Naur Form)

$::=$ 表示 “定义为” 或 “由.....组成”

$\langle \dots \rangle$ 表示非终结符号

| 表示 “或”

算术表达式文法的BNF表示:

$\langle \text{表达式} \rangle ::= \langle \text{表达式} \rangle + \langle \text{项} \rangle \mid \langle \text{表达式} \rangle - \langle \text{项} \rangle \mid \langle \text{项} \rangle$

$\langle \text{项} \rangle ::= \langle \text{项} \rangle * \langle \text{因子} \rangle \mid \langle \text{项} \rangle / \langle \text{因子} \rangle \mid \langle \text{因子} \rangle$

$\langle \text{因子} \rangle ::= (\langle \text{表达式} \rangle) \mid i$



2.1.3.2 产生式书写约定

非终结符号：次序靠后的大写字母，如：X、Y、Z

终结符号串：次序靠后的小写字母，如：u、v、...、z

文法符号串：小写的希腊字母，如： α 、 β 、 γ 、 δ

可以直接用产生式的集合代替四元组来描述文法，**第一个产生式的左部符号是文法的开始符号。**

示例：

1、文法 $G = (V_N, V_T, P, S)$ ， $V_N = \{ S \}$ ， $V_T = \{ 0, 1 \}$ ， $P = \{ S \rightarrow 0S1, S \rightarrow 01 \}$ ，S为开始符号

2、文法 $G = (V_N, V_T, P, S)$ ， $V_N = \{ \text{标识符}, \text{字母}, \text{数字} \}$ ， $V_T = \{ a, b, c, \dots, x, y, z, 0, 1, \dots, 9 \}$

$P = \{ \langle \text{标识符} \rangle \rightarrow \langle \text{字母} \rangle$

$\langle \text{标识符} \rangle \rightarrow \langle \text{标识符} \rangle \langle \text{字母} \rangle$

$\langle \text{标识符} \rangle \rightarrow \langle \text{标识符} \rangle \langle \text{数字} \rangle$

$\langle \text{字母} \rangle \rightarrow a, \dots, \langle \text{字母} \rangle \rightarrow z$

$\langle \text{数字} \rangle \rightarrow 0, \dots, \langle \text{数字} \rangle \rightarrow 9 \}$

$S = \langle \text{标识符} \rangle$



2.1.4 文法推导

文法所产生的语言

从文法的开始符号出发，推导得到的所有字符串的集合，为该文法定义的语言。

推导

假定 $A \rightarrow \gamma$ 是一个产生式， α 和 β 是任意的文法符号串，则： $\alpha A \beta \Rightarrow \alpha \gamma \beta$ 称 $\alpha A \beta$ 直接推导出 $\alpha \gamma \beta$ ，也可以说 $\alpha \gamma \beta$ 是 $\alpha A \beta$ 的直接推导 或说 $\alpha \gamma \beta$ 直接归约到 $\alpha A \beta$

如果有直接推导序列： $\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$ 则说 α_1 推导出 α_n ，记作： $\alpha_1 \Rightarrow^* \alpha_n$

称这个序列是从 α_1 到 α_n 的长度为 n 的推导

示例：

$\langle \text{句子} \rangle \rightarrow \langle \text{名词} \rangle \langle \text{动词} \rangle \langle \text{名词} \rangle$

$\langle \text{名词} \rangle \rightarrow \text{狼} \mid \text{羊} \mid \text{草} \mid \text{人} \mid \text{水}$

$\langle \text{动词} \rangle \rightarrow \text{吃} \mid \text{喝}$



2.1.4 文法推导

示例: $G = (\{+, *, (,), i\}, \{E, T, F\}, E, \varphi)$

$\varphi:$ $E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow (E) \mid i$

$\alpha A \beta$	$\alpha \gamma \beta$	α	β	所用产生式	从E到 $\alpha \gamma \beta$ 的推导长度
E	E+T	ϵ	ϵ	$E \rightarrow E+T$	1
E+T	T+T	ϵ	+T	$E \rightarrow T$	2
T+T	F+T	ϵ	+T	$T \rightarrow F$	3
F+T	i+T	ϵ	+T	$F \rightarrow i$	4
i+T	i+F	i+	ϵ	$T \rightarrow F$	5
i+F	i+i	i+	ϵ	$F \rightarrow i$	6

$E \Rightarrow E+T \Rightarrow T+T \Rightarrow F+T \Rightarrow i+T \Rightarrow i+F \Rightarrow i+i$



2.1.4 文法推导

最左推导

如果 $\alpha \xRightarrow{*} \beta$, 并且在每“一步推导”中, 都替换 α 中最左边的非终结符号, 则称这样的推导为最左推导。记作:

$$\alpha \xRightarrow{*}_{lm} \beta$$

$$E \Rightarrow E+T \Rightarrow T+T \Rightarrow F+T \Rightarrow i+T \Rightarrow i+F \Rightarrow i+i$$

最右推导

如果 $\alpha \xRightarrow{*} \beta$, 并且在每“一步推导”中, 都替换 α 中最右边的非终结符号, 则称这样的推导为最右推导。记作:

$$\alpha \xRightarrow{*}_{rm} \beta$$

最右推导也称为规范推导

$$E \Rightarrow E+T \Rightarrow E+F \Rightarrow E+i \Rightarrow T+i \Rightarrow F+i \Rightarrow i+i$$



根据文法对某个句型/句子进行最左推导或者最右推导，推导过程是唯一的吗？什么情况下不是唯一的？（例如表达式 $1+2-3$ ）



2.1.5 句子、句型、语言

句型

对于文法 $G=(V_T, V_N, S, \varphi)$ ，如果 $S \xRightarrow{*} \alpha$ ，则称 α 是当前文法的一个句型。

若 $S \xRightarrow[lm]{*} \alpha$ ，则 α 是当前文法的一个左句型，

若 $S \xRightarrow[rm]{*} \alpha$ ，则 α 是当前文法的一个右句型。

句子

仅含有终结符号的句型是文法的一个句子。

语言

文法 G 产生的所有句子组成的集合是文法 G 所定义的语言，记作 $L(G)$ 。



2.1.6 短语、直接短语和句柄

- 对于文法 $G=(V_T, V_N, S, \varphi)$, 假定 $\alpha\beta\delta$ 是文法 G 的一个句型。如果存在:

$$S \Rightarrow^* \alpha A \delta, \text{ 并且 } A \overset{+}{\Rightarrow} \beta$$

则称 β 是句型 $\alpha\beta\delta$ 关于非终结符号 A 的**短语**。如果存在:

$$S \Rightarrow^* \alpha A \delta, \text{ 并且 } A \Rightarrow \beta$$

则称 β 是句型 $\alpha\beta\delta$ 关于非终结符号 A 的**直接短语**。

一个句型的最左直接短语称为该句型的**句柄**。

例如:

$$\begin{array}{cccccccccccc} \underline{E} \Rightarrow \underline{T} \Rightarrow \underline{T * F} \Rightarrow \underline{T * (E)} \Rightarrow \underline{F * (E)} \Rightarrow \underline{i * (E)} \Rightarrow \underline{i * (E + T)} \Rightarrow \underline{i * (T + T)} \Rightarrow \underline{i * (F + T)} \Rightarrow \underline{i * (i + T)} \\ \textcircled{1} \quad \textcircled{2} \quad \textcircled{3} \quad \quad \textcircled{4} \quad \quad \textcircled{5} \quad \quad \textcircled{6} \quad \quad \textcircled{7} \quad \quad \textcircled{8} \quad \quad \textcircled{9} \quad \quad \textcircled{10} \end{array}$$



2.1.7 语言的定义

由文法 G 生成的语言记为 $L(G)$, 是文法 G 的一切句子的集合:
 $L(G) = \{x | S \Rightarrow x, \text{ 其中 } S \text{ 为文法的开始符号, 且 } x \in V_T^*\}$

例1: $G[S]: S \rightarrow 0S1, S \rightarrow 01$
 $L(G) = \{0^n 1^n | n \geq 1\}$

例2: 文法 $G[S]$:

- (1) $S \rightarrow aSBE$
- (2) $S \rightarrow aBE$
- (3) $EB \rightarrow BE$
- (4) $aB \rightarrow ab$
- (5) $bB \rightarrow bb$
- (6) $bE \rightarrow be$
- (7) $eE \rightarrow ee$

$$L(G) = \{ a^n b^n e^n \mid n \geq 1 \}$$



2.1.8 文法的等价

若 $L(G_1) = L(G_2)$, 则称文法 G_1 和 G_2 是等价的。

如文法 $G_1[A]$: $A \rightarrow 0R$ 与 $G_2[S]$: $S \rightarrow 0S1$ 等价
 $A \rightarrow 01$ $S \rightarrow 01$
 $R \rightarrow A1$



2.1.9 分析树及二义性

分析树：推导的图形表示即分析树，又称推导树。

分析树的每一个结点都有标记。

- 根结点由文法的开始符号标记；
- 每个内部结点由非终结符号标记，它的子结点由这个非终结符号的这次推导所用产生式的右部各符号从左到右依次标记；
- 叶结点由非终结符号或终结符号标记，它们从左到右排列起来，构成句型

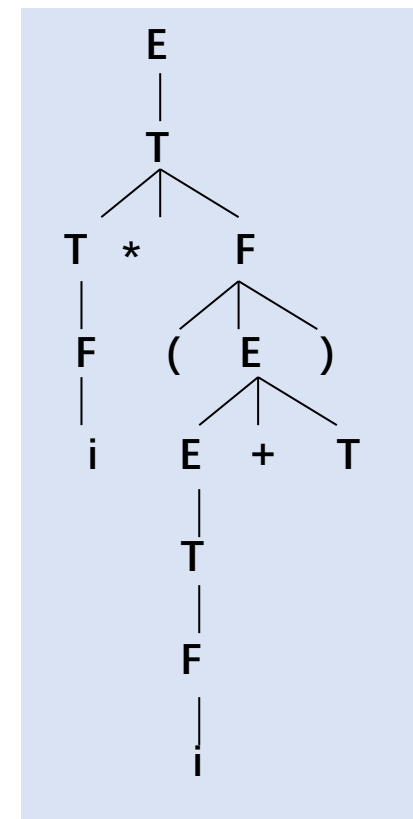
2.1.9 分析树及二义性

$G[E]: E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid i$

句型: $i*(i+T)$

推导过程:

$E \Rightarrow T \Rightarrow T * F \Rightarrow T * (E) \Rightarrow F * (E) \Rightarrow i * (E)$
 $\Rightarrow i * (E + T) \Rightarrow i * (T + T) \Rightarrow i * (F + T) \Rightarrow i * (i + T)$

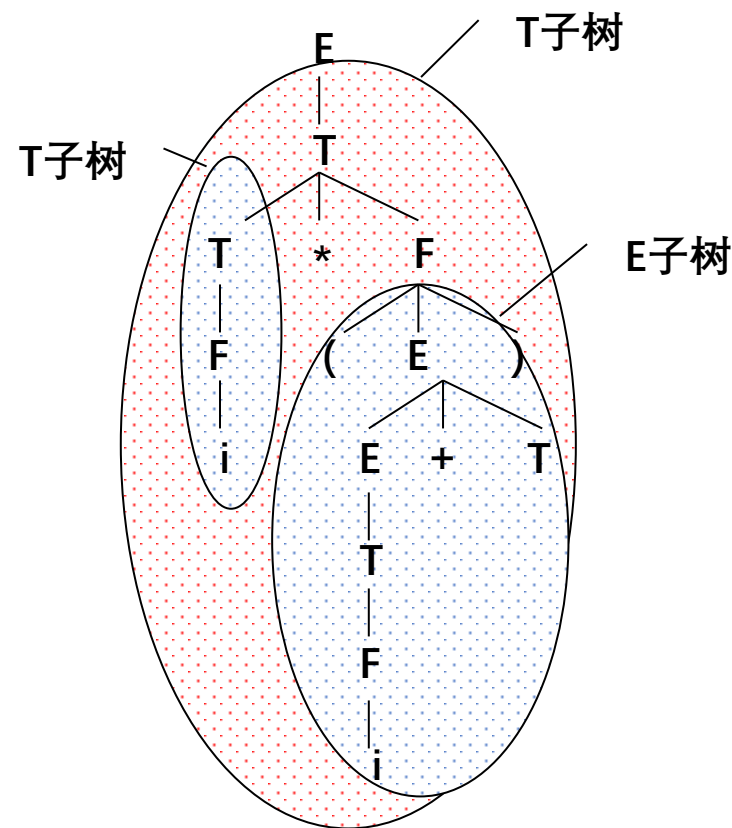


2.1.9 分析树及二义性

子树

分析树中一个特有的**结点**、连同它的**全部后裔结点**、连接这些结点的**边**、以及这些结点的**标记**。

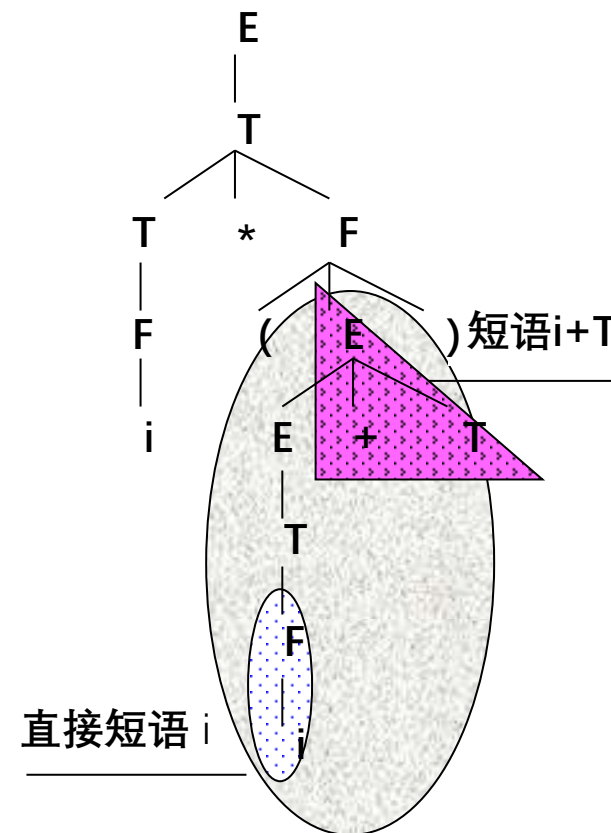
如果子树的根结点标记为非终结符号A，则可称该子树为**A-子树**。



2.1.9 分析树与二义性

子树与短语的关系

一棵子树的所有叶结点自左至右排列起来，形成此句型相对于该子树根的短语；
分析树中只有父子两代的子树的所有叶结点自左至右排列起来，形成此句型相对于该子树根的直接短语；
分析树中最左边的那棵只有父子两代的子树的所有叶结点自左至右排列起来，就是该句型的句柄。





讨论

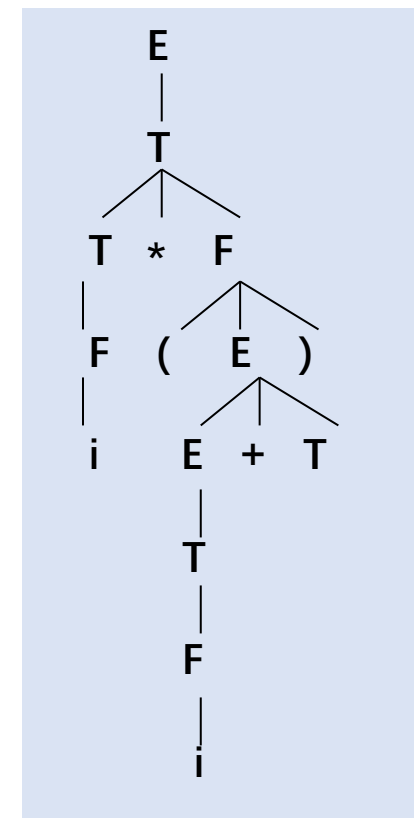
一个句子对应的分析树是唯一的吗?

$G[E]: E \rightarrow E+T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid I$

句型: $i*(i+T)$

$E \Rightarrow T \Rightarrow T * F \Rightarrow T * (E) \Rightarrow F * (E) \Rightarrow i * (E) \Rightarrow i * (E+T) \Rightarrow i * (T+T) \Rightarrow i * (F+T) \Rightarrow i * (i+T)$

$E \Rightarrow T \Rightarrow T * F \Rightarrow F * F \Rightarrow i * F \Rightarrow i * (E) \Rightarrow i * (E+T) \Rightarrow i * (T+T) \Rightarrow i * (F+T) \Rightarrow i * (i+T)$





2.1.9 分析树与二义性

文法的二义性

如果一个文法的某个句子有不只一棵分析树，则这个句子是**二义性的**。
含有二义性句子的文法是**二义性的文法**。

例：考虑文法 $G = (\{+, *, (,), i\}, \{E\}, E, \varphi)$

$\varphi: E \rightarrow E + E \mid E * E \mid (E) \mid id$

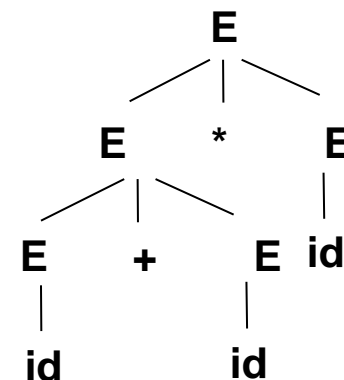
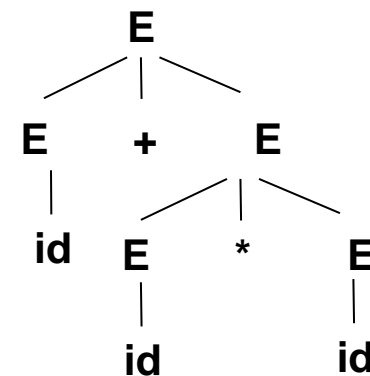
句子 $id + id * id$ 存在两个不同的最左推导：

$E \Rightarrow E + E \Rightarrow id + E \Rightarrow id + E * E \Rightarrow id + id * E \Rightarrow id + id * id$

$E \Rightarrow E * E \Rightarrow E + E * E \Rightarrow id + E * E \Rightarrow id + id * E \Rightarrow id + id * id$

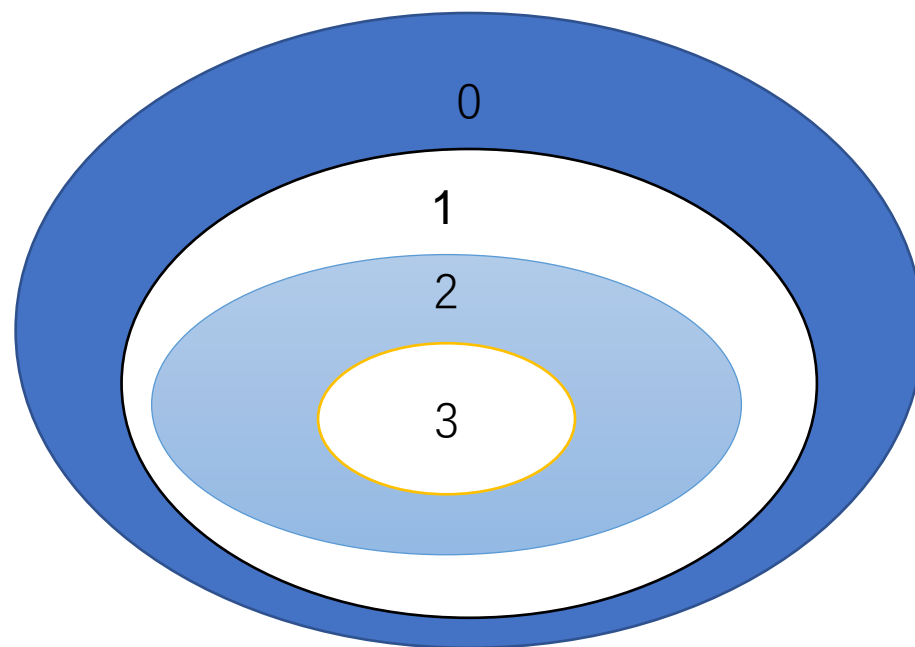
语言的二义性

有些语言，根本就不存在无二义性的文法，这样的语言称为**二义性的语言**。
有时，一个二义性的文法可以变换为一个等价的、无二义性的文法。
二义性问题是**不可判定的**





2.2 Chomsky语法体系





2.2 Chomsky语法体系

四类文法和相应的四种形式语言类

文法类型	产生式形式的限制	文法产生的语言类
0型文法	$\alpha \rightarrow \beta$ 其中 $\alpha, \beta \in (V_T \cup V_N)^*$ $ \alpha \neq 0$	0型语言
1型文法, 即 上下文有关文法	$\alpha \rightarrow \beta$ 其中 $\alpha, \beta \in (V_T \cup V_N)^*$ $ \alpha \leq \beta $	1型语言, 即 上下文有关语言
2型文法, 即 上下文无关文法	$A \rightarrow \beta$ 其中 $A \in V_N, \beta \in (V_T \cup V_N)^*$	2型语言, 即 上下文无关语言
3型文法, 即 正规文法 (线性文法)	$A \rightarrow a$ 或 $A \rightarrow aB$ (右线性), 或 $A \rightarrow a$ 或 $A \rightarrow Ba$ (左线性) 其中 $A, B \in V_N, a \in V_T \cup \{\epsilon\}$	3型语言, 即 正规语言



2.2 Chomsky文法体系

1型（上下文有关）文法

文法G[S]:

$$\begin{aligned} S &\rightarrow aSBE \\ S &\rightarrow aBE \\ EB &\rightarrow BE \\ aB &\rightarrow ab \\ bB &\rightarrow bb \\ bE &\rightarrow be \\ eE &\rightarrow ee \end{aligned}$$

2型（上下文无关）文法

文法G[S]:

$$\begin{aligned} S &\rightarrow aB \mid bA \\ A &\rightarrow a \mid aS \mid bAA \\ B &\rightarrow b \mid bS \mid aBB \end{aligned}$$

3型（正则）文法

文法G[S]:

$$\begin{aligned} S &\rightarrow 0A \mid 1B \mid 0 \\ A &\rightarrow 0A \mid 1B \mid 0S \\ B &\rightarrow 1B \mid 1 \mid 0 \end{aligned}$$

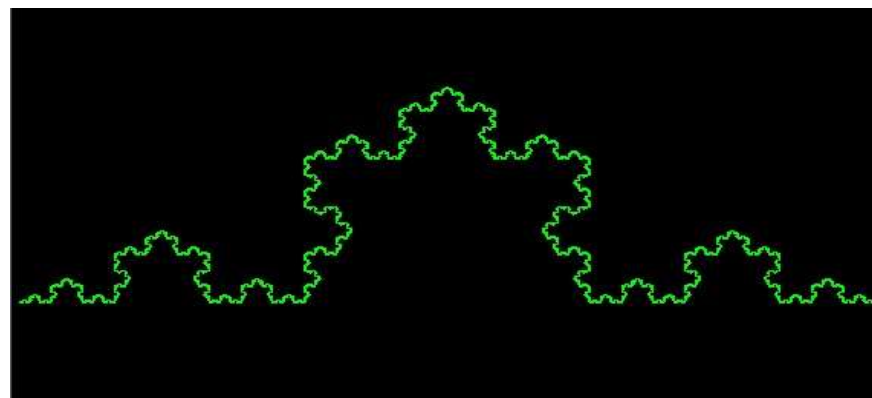


比较几种类型的文法，程序设计语言中不同组成成分，各自适合用哪种文法进行描述？为什么？



拓展：文法应用-分形

分形 (fractals)



<http://paulbourke.net/fractals/>



拓展：文法应用-分形

序列1: 0

序列2: 0 1

序列3: 0 1 1 2

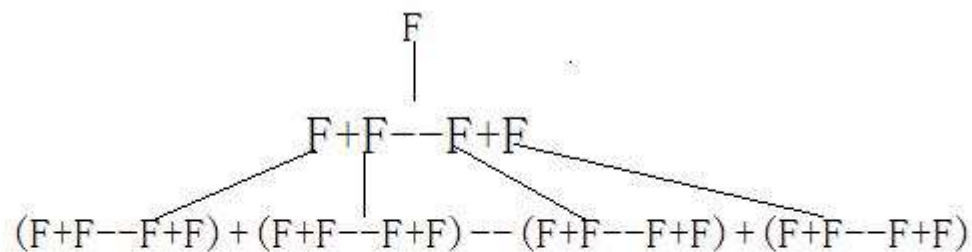
序列4: 0 1 1 2 1 2 2 3

序列5: 0 1 1 2 1 2 2 3 1 2 2 3 2 3 3 4

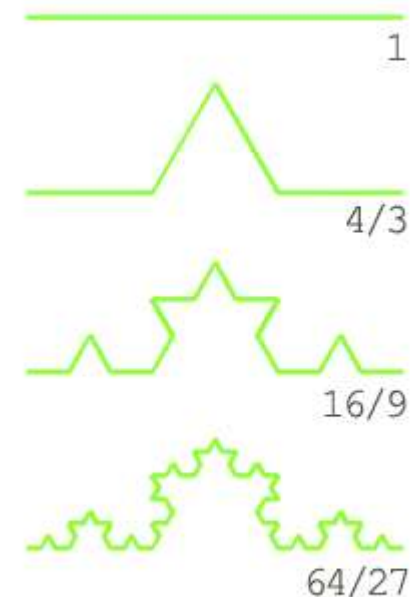
初始元：

第一级变换：

第二级变换：



·
·
·





拓展：文法应用-分形



重庆大学
CHONGQING UNIVERSITY

```
Koch.py > ...
4
5 # 定义初始变量
6 step = 10 # 步长
7 n = 5 # 迭代次数
8 # 初始条件
9 origin = 'F'
10 angel = 60
11 # 迭代规则
12 rule = { 'F' : "F+F--F+F" }
13 # 初始坐标
14 x = -300
15 y = 0
16 # 动画速度
17 speed = 'normal'
18
19
20 if __name__ == "__main__":
21     ls = GetLs(list(origin).rule.n)
```

OUTPUT TERMINAL DEBUG CONSOLE PROBLEMS 48

1: powershell

```
File "C:\Python38\lib\turtle.py", line 1292, in _incrementdc
    raise Terminator
turtle.Terminator
PS C:\Python38> python Koch.py
PS C:\Python38> python Koch.py
Traceback (most recent call last):
  File "Koch.py", line 22, in <module>
    Draw(ls,step,angel,x,y,speed)
  File "C:\Python38\drawLS.py", line 52, in draw
    drawMinus(angel)
  File "C:\Python38\drawLS.py", line 31, in drawMinus
    right(a)
  File "<string>", line 5, in right
turtle.Terminator
PS C:\Python38> LS文法
```

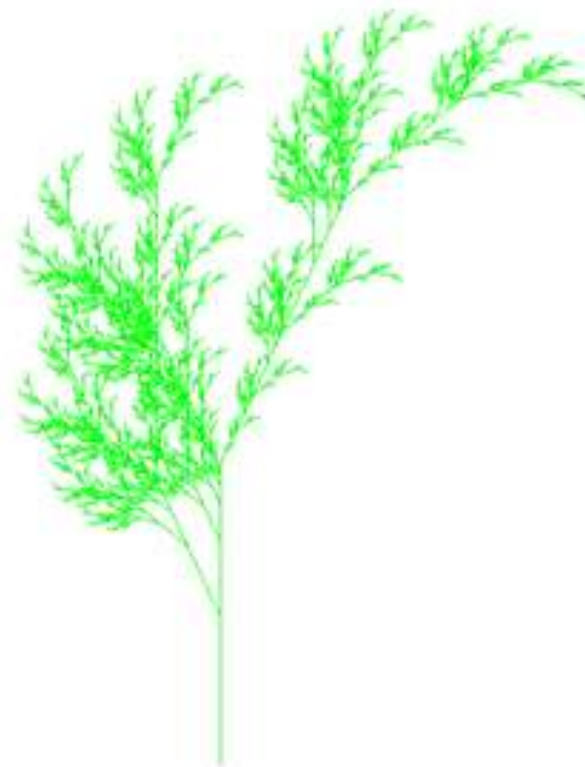
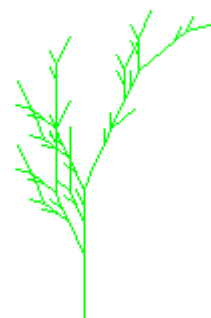


拓展:



重庆大学
CHONGQING UNIVERSITY

L系统—植物生长建模





拓展:



```
3
4 # 定义初始变量
5 step = 10 # 步长
6 n = 3 # 迭代次数
7 # 初始条件
8 origin = 'F'
9
10 angel = 20
11 # 迭代规则
12 rule = { 'F' : "F[-F]F[+F]F" }
13 # 初始坐标
14 x = 0
15 y = -250
16 # 动画速度s
17 speed = 'fast'
18 # 初始角度s
19 head = 90
20 # 渲染速度
```

OUTPUT TERMINAL DEBUG CONSOLE PROBLEMS 1: powershell + []

Windows PowerShell
版权所有 (C) Microsoft Corporation。保留所有权利。

PS C:\Python38> []



2.3 有限自动机

有限自动机(finite automata)

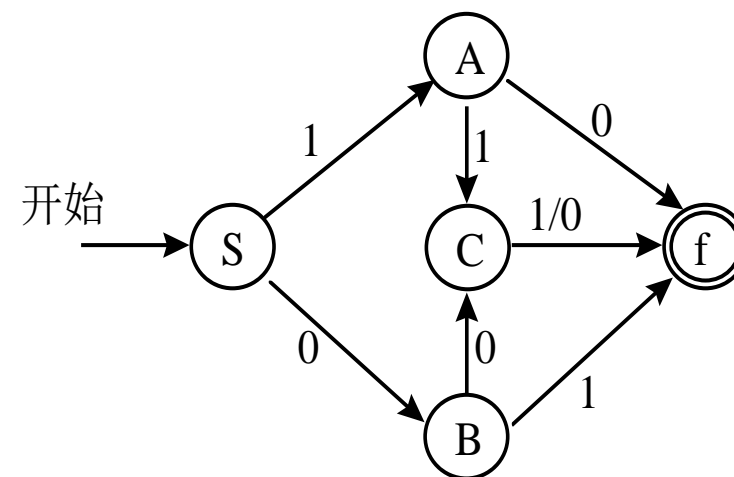
离散数字系统的抽象数学模型，具有离散的输入和输出

有限自动机分类

- 确定的有限自动机(Deterministic Finite Automata,DFA)
- 不确定的有限自动机(Nondeterministic Finite Automata,NFA)

2.3.1 确定的有限自动机 (DFA)

- 一张有限的方向图
- 图中结点代表状态，用圆圈表示
- 只含有限个状态，有一个初始状态，可以有若干个终结状态，终态用双圆圈表示。
- 状态之间用有向边连接
- 边上的标记表示在射出结点状态下可能出现的输入符号

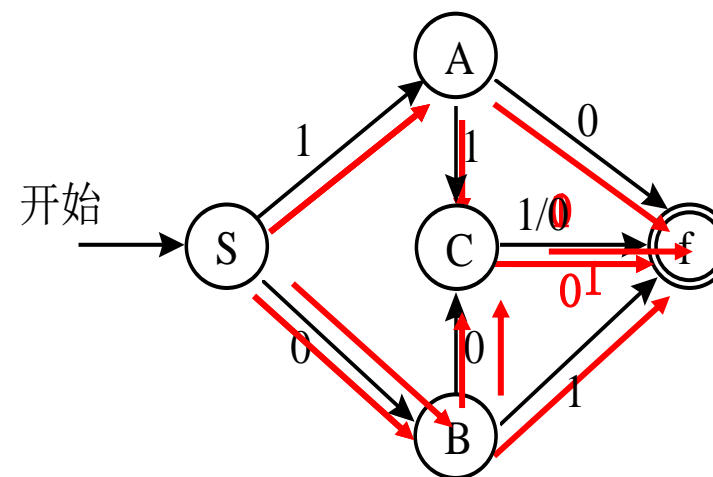


状态转换图

利用状态转换图识别符号串

- 识别方法：
 - (1) 起点为初态S，从w的最左符号开始，重复步骤(2)，直到达到w的最右符号为止。
 - (2) 扫描w的下一个符号，在当前状态的所有射出边中找出标记为该字符的边，沿此边过度到下一个状态。
- 状态转换图所能识别的符号串的全体称为该状态转换图所识别的语言。

$$L(M) = \{10, 110, 111, 01, 000, 001\}$$





确定的有限自动机的定义

一个确定的有限自动机M(记作: **DFA M**)是一个五元组: $M=(\Sigma, Q, q_0, F, \delta)$

其中 Σ : 是一个**字母表**, 它的每个元素称为一个输入符号

Q : 是一个有限的**状态集合**

$q_0 \in Q$: q_0 称为**初始状态**

$F \subseteq Q$: F 称为**终结状态集合**

δ : 是一个从 $Q \times \Sigma$ 到 Q 的单值**映射**

转换函数 $\delta(q, a)=q'$ (其中 $q, q' \in Q, a \in \Sigma$) 表示当前状态为 q , 输入符号为 a 时, 自动机将转换到下一个状态 q' , q' 称为 q 的一个**后继**。

若 $Q=\{q_1, q_2, \dots, q_n\}$, $\Sigma=\{a_1, a_2, \dots, a_m\}$, 则 $Q \times \Sigma = (\delta(q_i, a_j))_{n \times m}$ 是一个 n 行 m 列的矩阵, 它称为DFA M的**状态转换矩阵**, 也称为**转换表**



DFAM与状态转换图

DFAM可用一张状态转换图来表示:

- 若 q, q' 若DFA M 有 n 个状态, m 个输入符号, 则状态转换图有 n 个状态结点, 每个结点最多有 m 条射出边。
- $q \in Q, a \in \Sigma$, 并且 $\delta(q, a) = q'$, 则从 q 到 q' 有一条标记为 a 的有向边。
- 整个图含有唯一的一个初态。

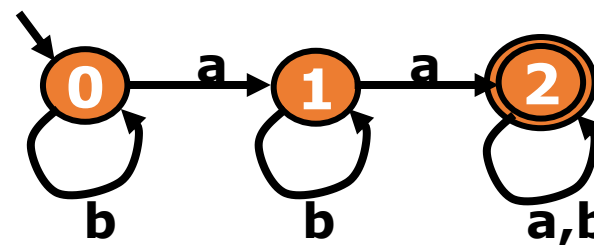


思考



图中自动机可接受什么样的字符串?

转移函数:

$$\{ (q_0, a) \rightarrow q_1, (q_0, b) \rightarrow q_0, \\ (q_1, a) \rightarrow q_2, (q_1, b) \rightarrow q_1, \\ (q_2, a) \rightarrow q_2, (q_2, b) \rightarrow q_2 \}$$




2.3.2 DFA M所识别的语言

对 Σ 上的任何符号串 $\omega \in \Sigma^*$ ，若存在一条从初态结点到终态结点的路径，该路径上每条边的标记连接成的符号串恰好是 ω ，则称 ω 为DFA M所识别。

DFA M所能识别的符号串的全体记为 $L(M)$ ，称为 DFA M 所识别的语言。

如果我们对所有 $\omega \in \Sigma^*$ ，递归地扩张 δ 的定义：对任何 $a \in \Sigma$ ， $q \in Q$ 定义

$$\delta(q, \epsilon) = q$$

$$\delta(q, \omega a) = \delta(\delta(q, \omega), a)$$

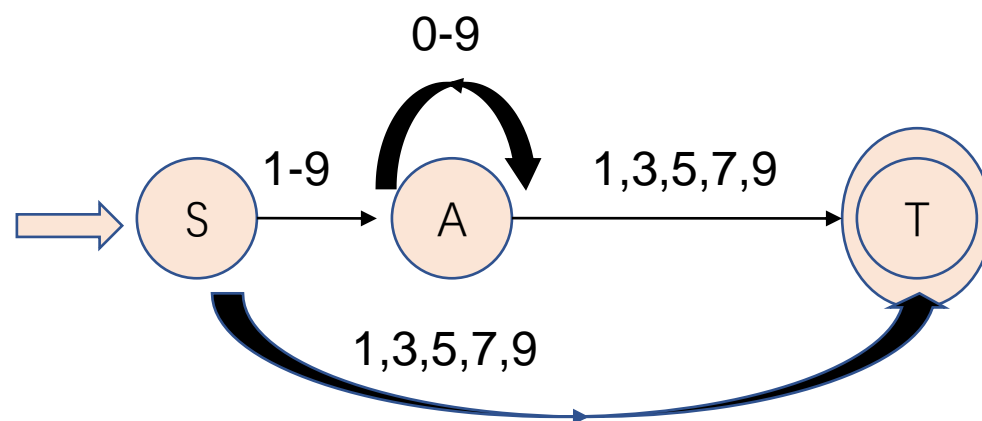
$$L(M) = \{ \omega \mid \omega \in \Sigma^*, \text{ 并且存在 } q \in F, \text{ 使 } \delta(q_0, \omega) = q \}$$



思考



如下所示的FA，所能识别的语言是什么？





DFA M所识别的语言

对 Σ 上的任何符号串 $\omega \in \Sigma^*$ ，若存在一条从初态结点到终态结点的路径，该路径上每条边的标记连接成的符号串恰好是 ω ，则称 ω 为DFA M所识别。

DFA M所能识别的符号串的全体记为 $L(M)$ ，称为 DFA M 所识别的语言。

如果我们对所有 $\omega \in \Sigma^*$ ，递归地扩张 δ 的定义：对任何 $a \in \Sigma$ ， $q \in Q$ 定义

$$\delta(q, \epsilon) = q$$

$$\delta(q, \omega a) = \delta(\delta(q, \omega), a)$$

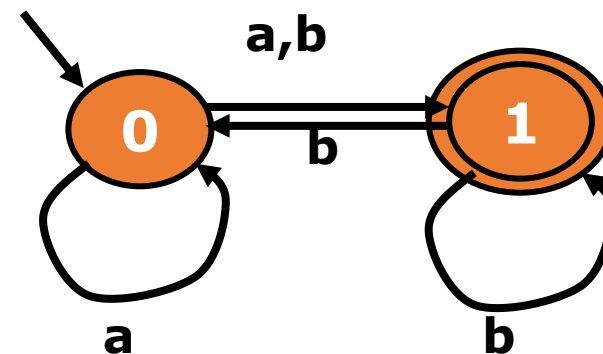
$$L(M) = \{ \omega \mid \omega \in \Sigma^*, \text{ 并且存在 } q \in F, \text{ 使 } \delta(q_0, \omega) = q \}$$

NFA的状态转换图及识别的语言

- 状态转换图
 - NFA M 含有 n 个状态, m 个输入符号
 - 图中含有 n 个状态结点, 每个结点可射出若干条边
 - 图中有唯一的一个初态结点
- 对 Σ 上的任何符号串 $\omega \in \Sigma^*$, 若存在一条从初态结点到终态结点的路径, 该路径上每条边的标记连接成的符号串恰好是 ω , 则称 ω 为 NFA M 所识别
- NFA M 所能识别的符号串的全体记为 $L(M)$, 称为 NFA M 所识别的语言

转换函数:

- $\{(q_0, a) \rightarrow \{q_0, q_1\},$
 $(q_0, b) \rightarrow \{q_1\},$
 $(q_1, a) \rightarrow \emptyset,$
 $(q_1, b) \rightarrow \{q_0, q_1\}\}$





2.3.3 NFA转化DFA

定理 2.1:

对任何一个NFA M ，都存在一个与之等价的DFA D ，即 $L(M)=L(D)$ 。

例：构造与下面的 NFA M 等价的 DFA D

NFA $M = (\{a,b\}, \{A,B\}, A, \{B\}, \delta)$

其中 δ : $\delta(A,a)=\{A,B\}$ $\delta(A,b)=\{B\}$ $\delta(B,b)=\{A,B\}$



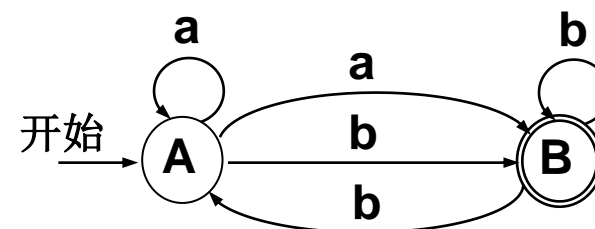
1、转换函数构造法:

首先, 画出该NFA M的状态转换图

假设 DFA $D = (\{a,b\}, Q', q_0', F', \delta')$, 则 Q' : Q 的所有子集组成的集合, 即 $Q' = 2^Q$, $Q' = \{\varnothing, \{A\}, \{B\}, \{A, B\}\}$

$q_0' = \{q_0\}$, 此DFAM中 $q_0' = \{A\}$

F' : 所有含有原NFA M终态的 Q 的子集组成的集合 $F' = \{\{B\}, \{A, B\}\}$



δ' 的构成

$$\delta'(\{q_1, q_2, \dots, q_k\}, a) = \delta(q_1, a) \cup \delta(q_2, a) \cup \dots \cup \delta(q_k, a)$$

$$\delta'(\varphi, a) = \varphi$$

$$\delta'(\varphi, b) = \varphi$$

$$\delta'(\{A\}, a) = \delta(A, a) = \{A, B\} \quad \delta'(\{A\}, b) = \delta(A, b) = \{B\}$$

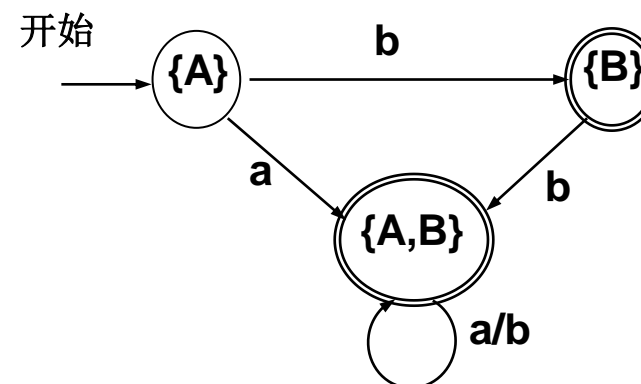
$$\delta'(\{B\}, a) = \delta(B, a) = \varphi \quad \delta'(\{B\}, b) = \delta(B, b) = \{A, B\}$$

$$\delta'(\{A, B\}, a) = \delta(A, a) \cup \delta(B, a) = \{A, B\} \cup \varphi = \{A, B\}$$

$$\delta'(\{A, B\}, b) = \delta(A, b) \cup \delta(B, b) = \{B\} \cup \{A, B\} = \{A, B\}$$

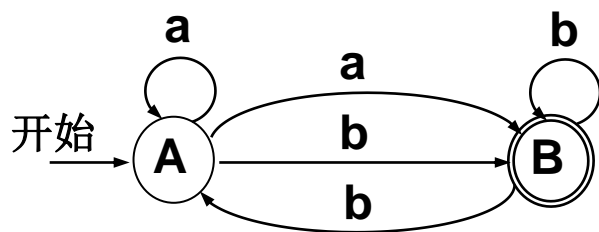
DFA D的状态转换矩阵和状态转换图

	a	b
{A}	{A, B}	{B}
{B}	—	{A, B}
{A, B}	{A, B}	{A, B}



2、子集构造法：

列出NFA M 的每个子集及该子集相对于每个输入符号的后继子集，对所有子集重新命名，得到DFA D 的状态转换矩阵。

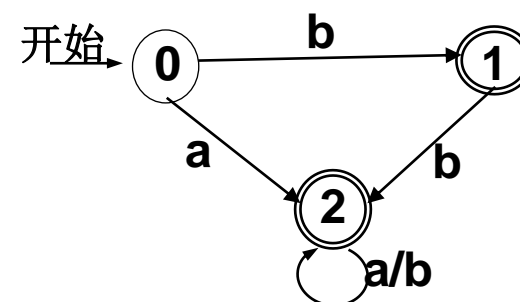


NFA M 的状态转换矩阵

状态子集 \ 输入	a	b
{A}	{A, B}	{B}
{B}	—	{A, B}
{A, B}	{A, B}	{A, B}

DFA D 的状态转换矩阵

状态子集 \ 输入	a	b
0	2	1
1	—	2
2	2	2



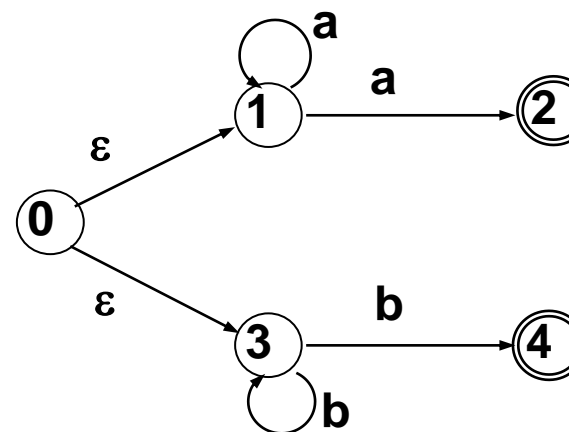
示例

有NFAM= $(\{a,b\},\{0,1,2,3,4\},0,\{2,4\},\delta)$ 其中 $\delta(0,\varepsilon)=\{1,3\}$, $\delta(1,a)=\{1,2\}$, $\delta(3,b)=\{3,4\}$

- NFA M的状态转换矩阵

	ε	a	b
0	$\{1,3\}$	—	—
1	—	$\{1,2\}$	—
2	—	—	—
3	—	—	$\{3,4\}$
4	—	—	—

- NFA M的状态转换图



NFA M所识别的语言为 $L(M)=\{ a^+|b^+ \}$ 。



具有 ε -转移的非确定有限自动机

定义：一个具有 ε -转移的非确定有限自动机M(记作：**NFA M**)，是一个五元组

$M = (\Sigma, Q, q_0, F, \delta)$ ，其中

Σ ：是一个**字母表**，它的每个元素称为一个输入符号

Q ：是一个**有限的状态集合**

$q_0 \in Q$ ： q_0 称为**初始状态**

$F \subseteq Q$ ：F称为**终结状态集合**

δ ：是一个从 $Q \times (\Sigma \cup \{\varepsilon\})$ 到 Q 的子集的**映射**，即 $\delta: Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^Q$

对任何 $q \in Q$ 及 $a \in (\Sigma \cup \{\varepsilon\})$ ，转移函数 δ 的值具有如下的形式

$\delta(q, a) = \{q_1, q_2, \dots, q_k\}$ 其中 $q_i \in Q$ ($i = 1, 2, \dots, k$)

NFA 的状态转换图

–图中可能有标记为 ε 的边

–当 $\delta(q, \varepsilon) = \{q_1, q_2, \dots, q_k\}$ 时，从 q 出发有 k 条标记为 ε 的边分别指向 q_1, q_2, \dots, q_k

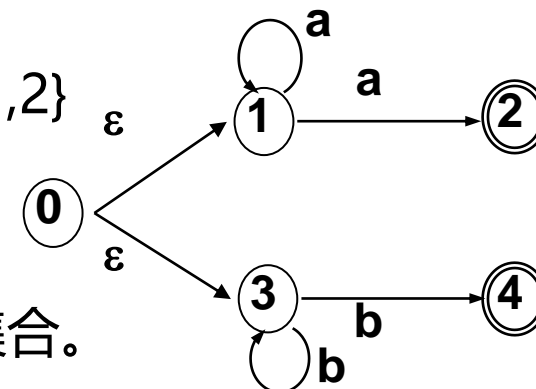


定理2.2

对任何一个具有 ε -转移的NFA M ，都存在一个等价的不具有 ε -转移的NFA N ，即 $L(M)=L(N)$ 。

例：构造与NFA M 等价的不具有 ε -转移的NFA

设NFA $M=(\{a,b\},\{0,1,2,3,4\},0,\{2,4\},\delta)$,其中 $\delta(0,\varepsilon)=\{1,3\}$ $\delta(1,a)=\{1,2\}$
 $\delta(3,b)=\{3,4\}$,状态转换图如右：



假设 NFA $N=(\{a,b\},\{0,1,2,3,4\},0,F',\delta')$

$\varepsilon_closure(q)$: 从状态 q 出发，经过 ε -道路可以到达的所有状态的集合。

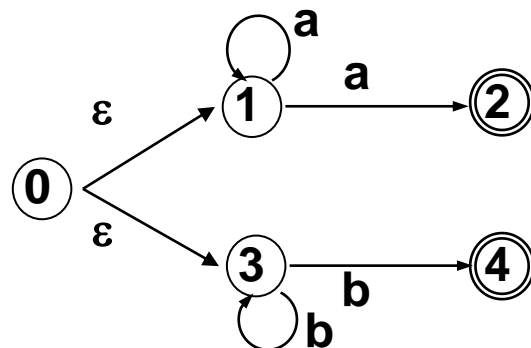
F' 的构成: $F' = \begin{cases} F \cup \{q_0\} & \text{如果 } \varepsilon_closure(q_0) \text{ 中包含 } F \text{ 的一个终态} \\ F & \text{否则} \end{cases}$

由于 $\varepsilon_closure(0)=\{0,1,3\}$ ，不包含NFA M 的终态，因此 $F'=F=\{2, 4\}$



δ' 的构成:

$\delta'(q,a)=\{q' \mid q' \text{ 为从 } q \text{ 出发, 经过标记为 } a \text{ 的道路所能到达的状态} \}$



$$\delta'(0, a) = \{1, 2\}$$

$$\delta'(0, b) = \{3, 4\}$$

$$\delta'(1, a) = \{1, 2\}$$

$$\delta'(1, b) = \phi$$

$$\delta'(2, a) = \phi$$

$$\delta'(2, b) = \phi$$

$$\delta'(3, a) = \phi$$

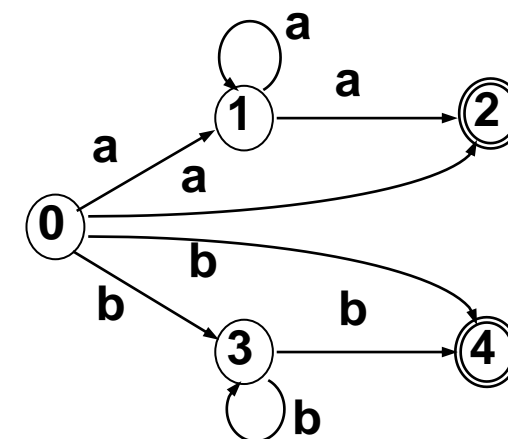
$$\delta'(3, b) = \{3, 4\}$$

$$\delta'(4, a) = \phi$$

$$\delta'(4, b) = \phi$$

NFA N的状态转换矩阵和状态转换图:

	a	b
0	$\{1, 2\}$	$\{3, 4\}$
1	$\{1, 2\}$	—
2	—	—
3	—	$\{3, 4\}$
4	—	—





推论:

对于任何一个具有 ε -转移的NFA M , 都存在一个与之等价的DFA D , 即 $L(M)=L(D)$ 。

DFA D 的每个状态对应NFA M 的一个状态子集。

对给定的输入符号串, 为了使 D “并行地” 模拟 M 所能产生的所有可能的转换, 令 q 为NFA M 的状态, T 为 NFA M 的状态子集, 引入以下操作:

(1) $\varepsilon_closure(q)=\{q' \mid \text{从} q \text{ 出发, 经过} \varepsilon\text{-道路可以到达状态} q'\}$

(2) $\varepsilon_closure(T)=\bigcup \varepsilon_closure(q_i)$ 其中 $q_i \in T$,

从 T 中任一状态出发, 经过 ε -道路后可以到达的状态集合。

(3) $move(T,a)=\{q \mid \delta(q_i,a)=q, \text{ 其中} q_i \in T\}$

从某个状态 $q_i \in T$ 出发, 经过输入符号 a 之后可到达的状态集合



推论:

对于任何一个具有 ε -转移的NFA M , 都存在一个与之等价的DFA D , 即 $L(M)=L(D)$ 。

DFA D 的每个状态对应NFA M 的一个状态子集。

对给定的输入符号串, 为了使 D “并行地”模拟 M 所能产生的所有可能的转换, 令 q 为NFA M 的状态, T 为 NFA M 的状态子集, 引入以下操作:

(1) $\varepsilon_closure(q)=\{q' \mid \text{从} q \text{出发, 经过} \varepsilon\text{-道路可以到达状态} q'\}$

(2) $\varepsilon_closure(T)=\bigcup \varepsilon_closure(q_i)$ 其中 $q_i \in T$,

从 T 中任一状态出发, 经过 ε -道路后可以到达的状态集合。

(3) $move(T,a)=\{q \mid \delta(q_i,a)=q, \text{ 其中} q_i \in T\}$

从某个状态 $q_i \in T$ 出发, 经过输入符号 a 之后可到达的状态集合



算法：由NFA构造等价的DFA

输入：一个NFA M

输出：一个与NFA M 等价（即接受同样语言）的DFA D

方法：为DFA D 构造状态转换矩阵DTT

初态： $\varepsilon_closure(q_0)$ 是DQ中唯一的状态，且未标记。

while (DQ中存在一个未标记的状态T)

{

 标记 T;

 for (each $a \in \Sigma$)

 {

$U = \varepsilon_closure(move(T, a));$

 if ($U \notin DQ$)

 把U做为一个未标记的状态加入DQ;

$DTT[T, a] = U;$

 }

}

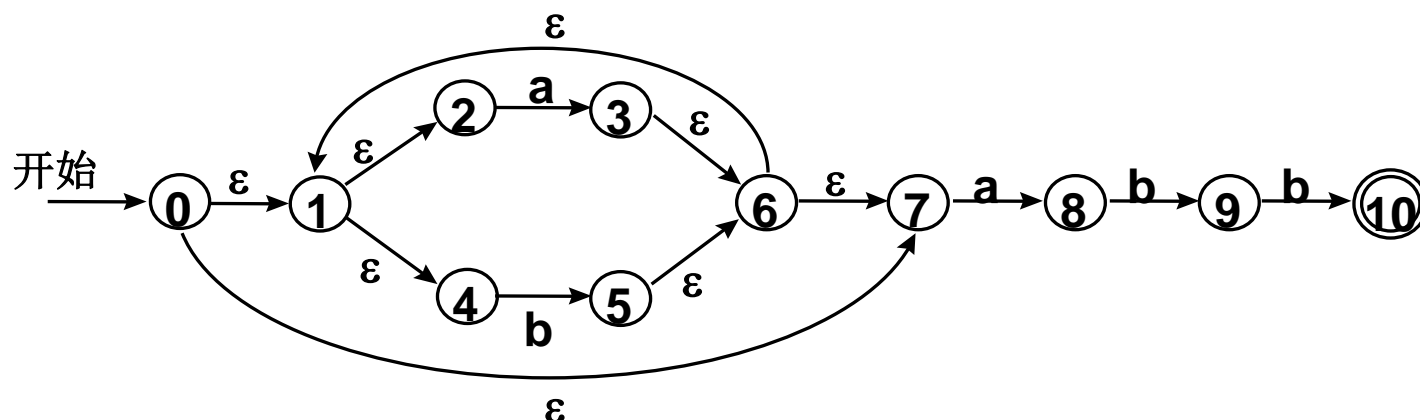


算法：计算 ε _closure(T)

```
把T中所有状态压入栈;  
 $\varepsilon$ _closure(T)的初值置为T;  
while 栈不空  
{  
    弹出栈顶元素t;  
    for (each  $q \in \varepsilon$ _closure(t))  
        if ( $q \notin \varepsilon$ _closure(T)) {  
            把q加入 $\varepsilon$ _closure(T);  
            把q压入栈;  
        }  
}
```



示例：构造与下面的NFA M等价的DFA D



字母表 $\Sigma = \{a, b\}$

初态为A, $A = \epsilon\text{-closure}(0) = \{0, 1, 2, 4, 7\}$

$DTT[A, a] = \epsilon\text{-closure}(\text{move}(A, a))$

$= \epsilon\text{-closure}(\text{move}(0, a) \cup \text{move}(1, a) \cup \text{move}(2, a) \cup \text{move}(4, a) \cup \text{move}(7, a))$

$= \epsilon\text{-closure}(\{3, 8\}) = \epsilon\text{-closure}(3) \cup \epsilon\text{-closure}(8) = \{1, 2, 3, 4, 6, 7, 8\} = B$

$DTT[A, b] = \epsilon\text{-closure}(\text{move}(A, b)) = \epsilon\text{-closure}(5) = \{1, 2, 4, 5, 6, 7\} = C$

$DTT[B, a] = \epsilon\text{-closure}(\text{move}(B, a)) = \epsilon\text{-closure}(\{3, 8\}) = B$

$DTT[B, b] = \epsilon\text{-closure}(\text{move}(B, b)) = \epsilon\text{-closure}(\{5, 9\}) = \{1, 2, 4, 5, 6, 7, 9\} = D$



示例 (续)

$$DTT[C, a] = \varepsilon\text{-closure}(\text{move}(C, a)) = \varepsilon\text{-closure}(\{3, 8\}) = B$$

$$DTT[C, b] = \varepsilon\text{-closure}(\text{move}(C, b)) = \varepsilon\text{-closure}(5) = C$$

$$DTT[D, a] = \varepsilon\text{-closure}(\text{move}(D, a)) = \varepsilon\text{-closure}(\{3, 8\}) = B$$

$$\begin{aligned} DTT[D, b] &= \varepsilon\text{-closure}(\text{move}(D, b)) = \varepsilon\text{-closure}(\{5, 10\}) \\ &= \{1, 2, 4, 5, 6, 7, 10\} = E \end{aligned}$$

$$DTT[E, a] = \varepsilon\text{-closure}(\text{move}(E, a)) = \varepsilon\text{-closure}(\{3, 8\}) = B$$

$$DTT[E, b] = \varepsilon\text{-closure}(\text{move}(E, b)) = \varepsilon\text{-closure}(5) = C$$

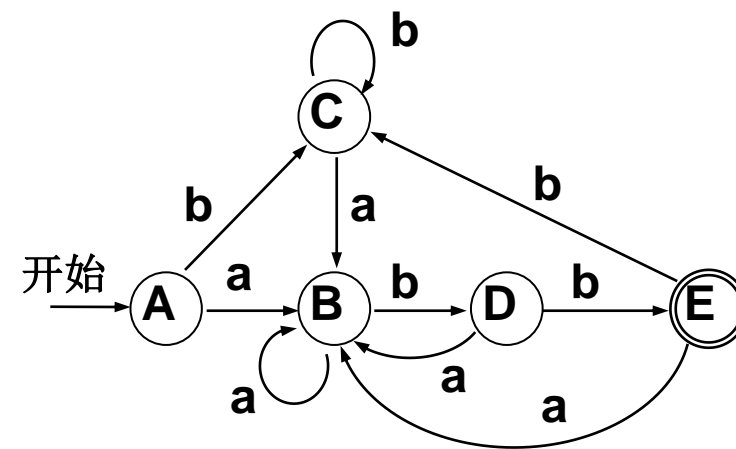
DFA D有5个状态，即A、B、C、D、E，

其中A为初态

E为终态，因为E的状态集合中包括原NFA M的终态10。

DFA D的状态转换矩阵和状态转换图

	a	
b		
A	B	C
B	B	D
C	B	C
D	B	E
E	B	C





算法：计算 ε _closure(T)

```
把T中所有状态压入栈;  
 $\varepsilon$ _closure(T)的初值置为T;  
while 栈不空  
{  
    弹出栈顶元素t;  
    for (each  $q \in \varepsilon$ _closure(t))  
        if ( $q \notin \varepsilon$ _closure(T)) {  
            把q加入 $\varepsilon$ _closure(T);  
            把q压入栈;  
        }  
}
```




思考



重庆大学
CHONGQING UNIVERSITY

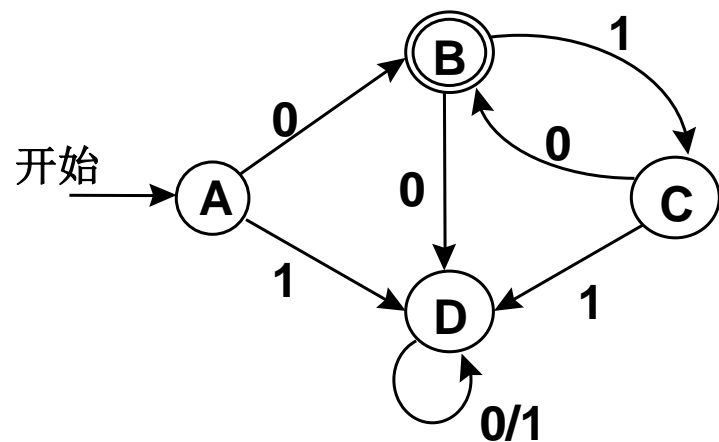
DFA/NFA 如何实现?

2.3.4 DFA的化简

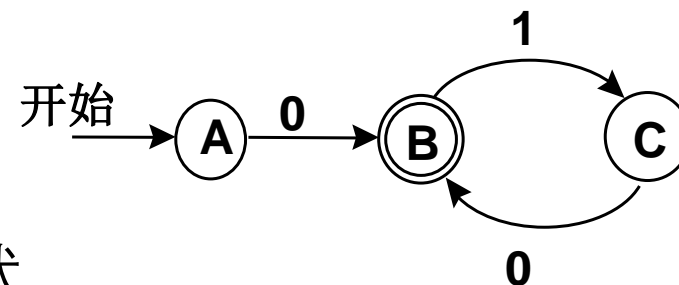
对于任何一个含有 n 个状态的DFA，都存在含有 m ($m > n$) 个状态的DFA与之等价。

DFA D的化简，指寻找一个状态数比较少的DFA D' ，使 $L(D) = L(D')$ 。

可以证明，存在一个最少状态的DFA D'' ，使 $L(D) = L(D'')$ ，并且这个 D'' 是唯一的



- 状态D是一个“死状态”，是一个无用的状态。
- 去掉状态D及与之相连接的边，可以得到等价的状态转换图





2.3.4 DFA的化简

状态的区分

定义： 设 $s, t \in Q$ ，若对任何 $\omega \in \Sigma^*$ ， $\delta(s, \omega) \in F$ 当且仅当 $\delta(t, \omega) \in F$ ，则称状态 s 和 t 是等价的。否则称状态 s 和 t 是可区分的。

首先，把 D 的状态集合分割成一些互不相交的子集，使每个子集中的任何两个状态是等价的，而任何两个属于不同子集的状态是可区分的。

然后，在每个子集中任取一个状态作“代表”，删去该子集中其余的状态，并把射向其它结点的边改为射向“代表”结点。

最后，如果得到的 DFA 中有死状态、或从初态无法到达的状态，则把它们删除。



Hopcroft算法

- 把状态集合Q划分成两个子集：终态子集F和非终态子集G。
- 对每个子集进行划分：
 - 取某个子集 $S = \{s_1, s_2, \dots, s_k\}$
 - 取某个输入符号c，检查S中的每个状态对该输入符号的转换。
 - 如果S中的状态相对于c，转换到不同子集中的状态，则要对S进行划分。使S中能够转换到同一子集的状态作为一个新的子集。
 - 重复上述过程，直到每个子集都不能再划分为止。

split(S)

foreach (character c)

if (c can split S)

split S into T1, ..., Tk

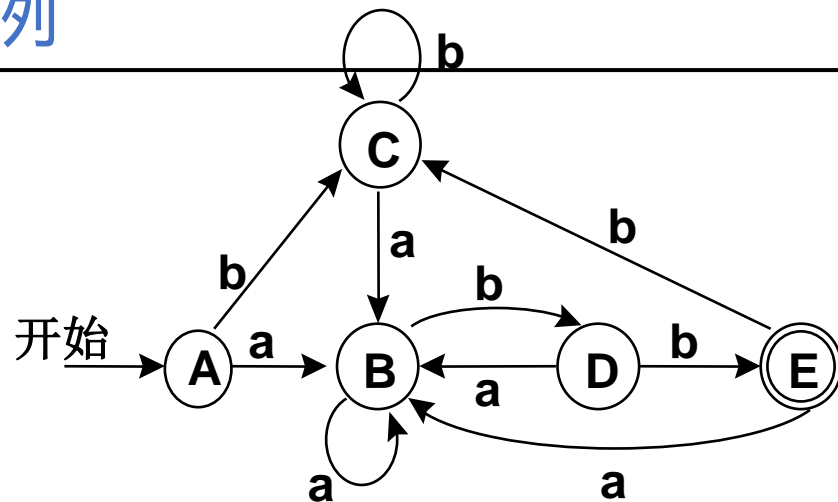
hopcroft ()

split all nodes into F, G

while (set is still changes) split(S)



示例



第一步：把DFA D的状态集合划分为子集，使每个子集中的状态相互等价，不同子集中的状态可区分。

把D的状态集合划分为两个子集： $\{A, B, C, D\}$ 和 $\{E\}$

考察非终态子集 $\{A, B, C, D\}$

- 对于a，状态A,B,C,D都转换到状态B，所以对输入符号a而言，该子集不能再划分。
- 对于b，状态A,B,C都转换到子集 $\{A, B, C, D\}$ 中的状态，而状态D则转换到子集 $\{E\}$ 中的状态。
- 应把子集 $\{A, B, C, D\}$ 划分成两个新的子集 $\{A, B, C\}$ 和 $\{D\}$ 。

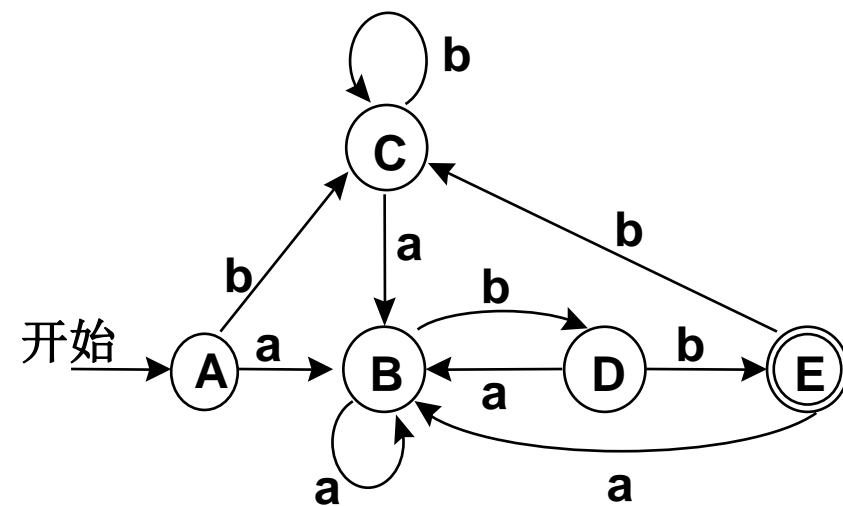
示例

- 考察子集 $\{A, B, C\}$
 - 对于a, 状态A,B,C都转换到状态B, 所以对输入符号a而言, 该子集不能再划分。
 - 对于b, 状态A,C转换到C, 状态B转换到D。状态C和D分属于不同的子集。
 - 应把子集 $\{A, B, C\}$ 划分成两个新的子集 $\{A, C\}$ 和 $\{B\}$ 。

D的状态集合被划分为: $\{A, C\}$ 、 $\{B\}$ 、 $\{D\}$ 和 $\{E\}$

- 考察子集 $\{A, C\}$
 - 对于a, 状态A,C都转换到状态B。
 - 对于b, 状态A,C都转换到状态C。
 - 该子集不可再划分。

D的状态集合最终被划分为: $\{A, C\}$ 、 $\{B\}$ 、 $\{D\}$ 和 $\{E\}$

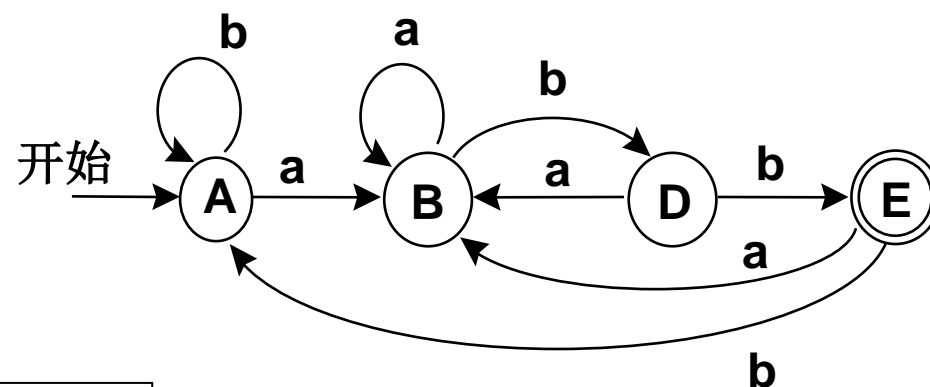


示例

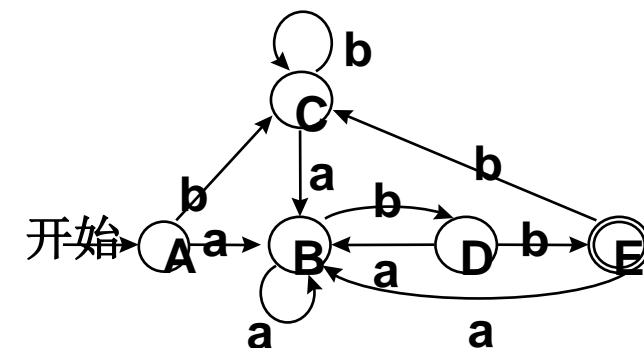
- **第二步：** 为每个子集选择一个代表状态。

- 选择A为子集 {A,C} 的代表状态

- D'的状态转换图



状态	输入符号	
	a	b
A	B	A
B	B	D
D	B	E
E	B	A





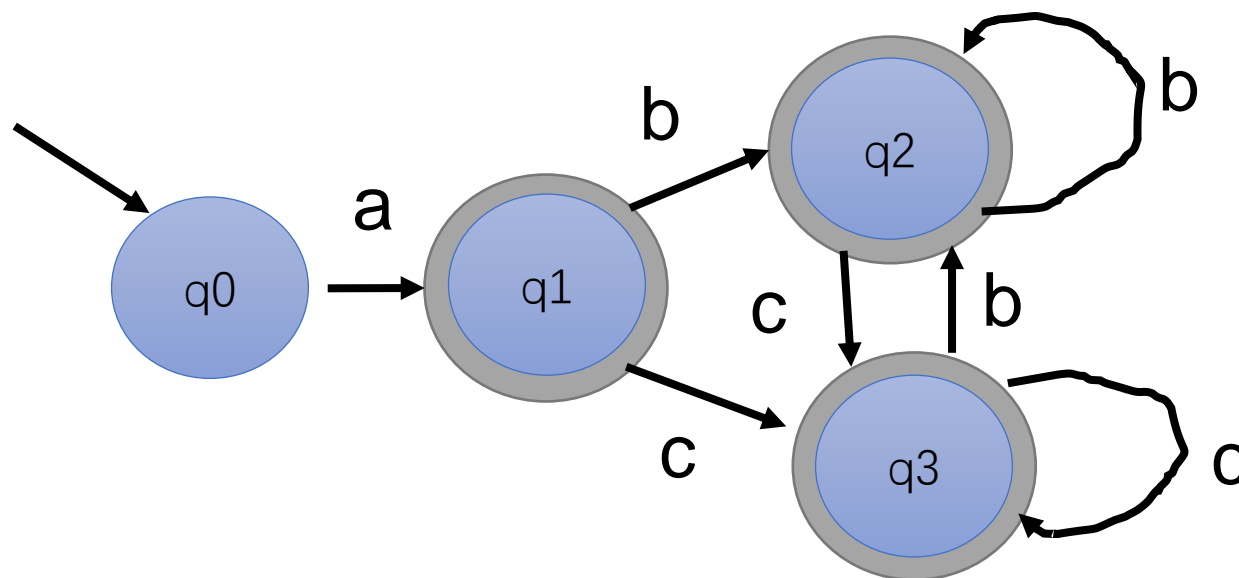
课堂练习

请对图示的DFA进行最小化



重庆大学
CHONGQING UNIVERSITY

请将图示中的DFA最小化





2.4 文法与自动机

- **0型文法**（短语文法）的能力相当于**图灵机**，可以表征任何递归可枚举集，而且任何0型语言都是递归可枚举的
- **1型文法**（上下文有关文法）：产生式的形式为 $\alpha_1 A \alpha_2 \rightarrow \alpha_1 \beta \alpha_2$ ，即只有A出现在 α_1 和 α_2 的上下文中时，才允许 β 取代A。其识别系统是**线性界限自动机**。
- **2型文法**（上下文无关文法、CFG）：产生式的形式为 $A \rightarrow \beta$ ， β 取代A时与A的上下文无关。其识别系统是**不确定的下推自动机**。
- **3型文法**（正规文法、右线性文法）：产生的语言是**有穷自动机（FA）**所接受的集合



正规文法与有限自动机的等价性

如果对于某个正规文法 G 和某个有限自动机 M , 有 $L(G)=L(M)$, 则称 G 和 M 是等价的。

定理: 对每一个右线性文法 G 或左线性文法 G , 都存在一个等价的有限自动机 M 。

定理: 对每一个DFA M , 都存在一个等价的右线性文法 G 和一个等价的左线性文法 G' 。

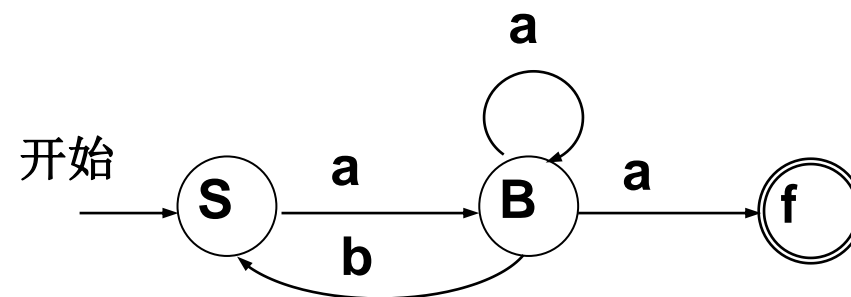
正规文法与有限自动机的等价性

示例:

设有右线性文法 $G = (\{a, b\}, \{S, B\}, S, \varphi)$,

其中 $\varphi: S \rightarrow aB, B \rightarrow aB \mid bS \mid a$, 试构造与 G 等价的有限自动机 M 。

- 设 FA $M = (\Sigma, Q, q_0, F, \delta)$
- $\Sigma = \{a, b\}$ $q_0 = S$ $F = \{f\}$ $Q = \{S, B, f\}$
- 转换函数 δ :
 - 对于产生式 $S \rightarrow aB$, 有 $\delta(S, a) = \{B\}$
 - 对于产生式 $B \rightarrow aB$, 有 $\delta(B, a) = \{B\}$
 - 对于产生式 $B \rightarrow bS$, 有 $\delta(B, b) = \{S\}$
 - 对于产生式 $B \rightarrow a$, 有 $\delta(B, a) = \{f\}$
- FA M 的状态转换图:



正规文法与有限自动机的等价性

示例：设有DFA $M = (\{a, b\}, \{q_0, q_1, q_2, q_3\}, q_0, \{q_3\}, \delta)$

其中转换函数 δ 如下：

$\delta(q_0, a) = q_1$, $\delta(q_1, a) = q_3$, $\delta(q_2, a) = q_2$

$\delta(q_0, b) = q_2$, $\delta(q_1, b) = q_1$, $\delta(q_2, b) = q_3$

试构造与之等价的右线性文法 G 。

构造右线性文法 $G = (V_T, V_N, S, \varphi)$

$V_T = \{a, b\}$ $V_N = \{q_0, q_1, q_2, q_3\}$ $S = q_0$

产生式集合 φ

$\delta(q_0, a) = q_1$, $\therefore q_0 \rightarrow aq_1$

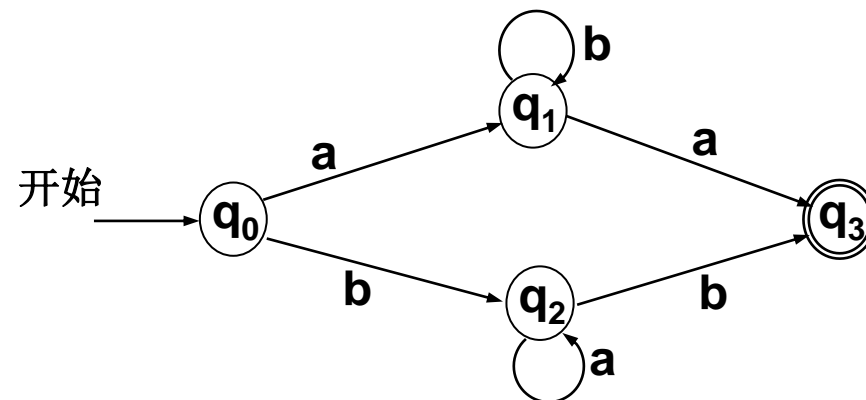
$\delta(q_0, b) = q_2$, $\therefore q_0 \rightarrow bq_2$

$\delta(q_1, a) = q_3$, $q_3 \in F$, $\therefore q_1 \rightarrow a \mid aq_3$

$\delta(q_1, b) = q_1$, $\therefore q_1 \rightarrow bq_1$

$\delta(q_2, a) = q_2$, $\therefore q_2 \rightarrow aq_2$

$\delta(q_2, b) = q_3$, $q_3 \in F$, $\therefore q_2 \rightarrow b \mid bq_3$



DFA M的状态转换图

构造的文法 G ：

$G = (\{a, b\}, \{q_0, q_1, q_2\}, q_0, \varphi)$

φ : $q_0 \rightarrow aq_0 \mid bq_2$

$q_1 \rightarrow a \mid bq_1$

$q_2 \rightarrow aq_2 \mid b$



2.5 正规表达式与有限自动机的等价性

2.5.1 正规表达式

用正规表达式可以精确地定义集合,

定义Pascal语言标识符的正规表达式: $\text{letter}(\text{letter}|\text{digit})^*$

定义: 字母表 Σ 上的正规表达式

正规表达式表示的语言叫做正规集。

- (1) ϵ 是正规表达式, 它表示的语言是 $\{\epsilon\}$
- (2) 如果 $a \in \Sigma$, 则 a 是正规表达式, 它表示的语言是 $\{a\}$
- (3) 如果 r 和 s 都是正规表达式, 分别表示语言 $L(r)$ 和 $L(s)$, 则:
 - $(r) | (s)$ 是正规表达式, 表示的语言是 $L(r) \cup L(s)$
 - $(r)(s)$ 是正规表达式, 表示的语言是 $L(r)L(s)$
 - $(r)^*$ 是正规表达式, 表示的语言是 $(L(r))^*$



2.5 正规表达式与有限自动机的等价性

2.5.1 正规表达式

用正规表达式可以精确地定义集合,

定义Pascal语言标识符的正规表达式: $\text{letter}(\text{letter}|\text{digit})^*$

定义: 字母表 Σ 上的正规表达式

正规表达式表示的语言叫做正规集。

- (1) ϵ 是正规表达式, 它表示的语言是 $\{\epsilon\}$
- (2) 如果 $a \in \Sigma$, 则 a 是正规表达式, 它表示的语言是 $\{a\}$
- (3) 如果 r 和 s 都是正规表达式, 分别表示语言 $L(r)$ 和 $L(s)$, 则:
 - $(r) | (s)$ 是正规表达式, 表示的语言是 $L(r) \cup L(s)$
 - $(r)(s)$ 是正规表达式, 表示的语言是 $L(r)L(s)$
 - $(r)^*$ 是正规表达式, 表示的语言是 $(L(r))^*$



2.5.1 正规表达式的书写约定

- 一元闭包 $'^*$ 具有最高优先级，并且遵从左结合
- 连接运算的优先级次之，遵从左结合
- 并运算 $'|'$ 的优先级最低，遵从左结合

例：如果 $\Sigma = \{a, b\}$ ，则有：

正规表达式 $a|b$ 表示集合 $\{a, b\}$

$(a|b)(a|b)$ 表示： $\{aa, ab, ba, bb\}$

a^* 表示：由0个或多个a组成的所有符号串的集合

$a|a^*b$ 表示：a和0个或多个a后跟一个b的所有符号串的集合

$(a|b)^*$ 表示：由a和b构成的所有符号串的集合

如果两个正规表达式r和s表示同样的语言，即 $L(r) = L(s)$ ，则称r和s等价，写作 $r = s$ 。

如： $(a|b) = (b|a)$



2.5.2 正规表达式遵从的代数定律

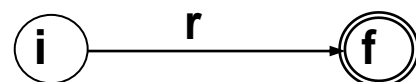
定律	说明
$r \mid s = s \mid r$	“并”运算是可交换的
$r \mid (s \mid t) = (r \mid s) \mid t$	“并”运算是可结合的
$(rs)t = r(st)$	连接运算是可结合的
$r(s \mid t) = rs \mid rt$ $(s \mid t)r = sr \mid tr$	连接运算对并运算的分配
$\epsilon r = r, r\epsilon = r$	对连接运算而言, ϵ 是单位元素
$r^* = (r \mid \epsilon)^*$	$*$ 和 ϵ 之间的关系
$r^{**} = r^*$	$*$ 是等幂的
$r^* = r^+ \mid \epsilon, r^+ = rr^*$	$+$ 和 $*$ 之间的关系



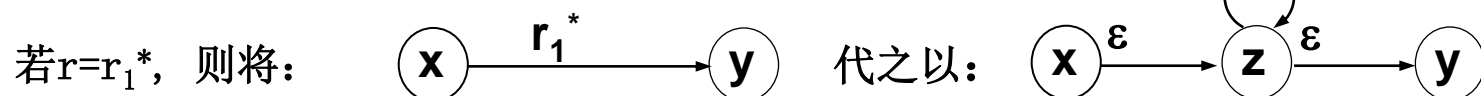
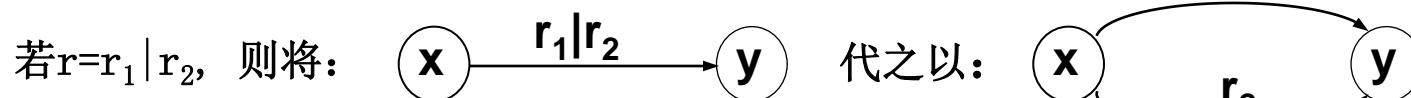
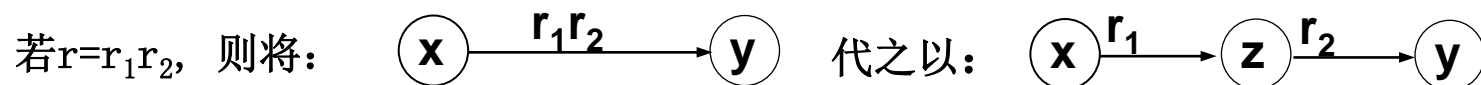
2.5.3 正规式与有限自动机

定理：对任何一个正规表达式 r ，都存在一个FA M ，使 $L(r)=L(M)$ ，反之亦然。

- 证1：设 r 是 Σ 上的一个正规表达式，则存在一个具有 ε -转移的NFA M 接受 $L(r)$ 。首先，为正规表达式 r 构造如下图所示的拓广转换图。



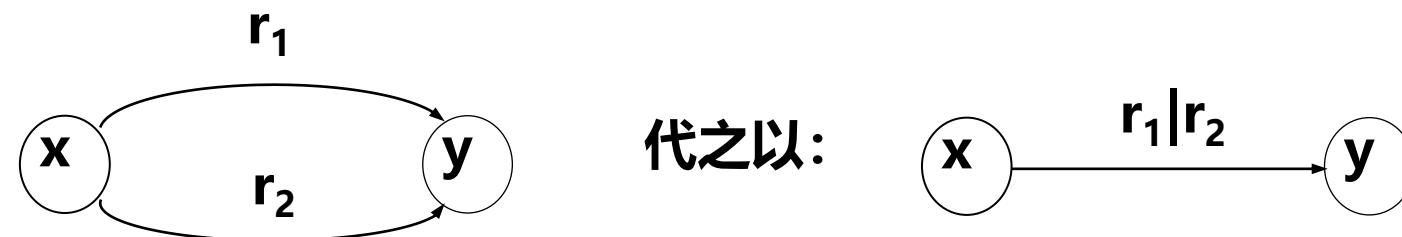
然后，按照下面的转换规则，对正规表达式 r 进行分裂、加入新的结点，直到每条边的标记都为基本符号为止。



2.5.3 正规式与有限自动机

证2: 设有FAM, 则存在一个正规表达式R, 它表示的语言即该FAM所识别的语言。

- 首先, 在FAM的转换图中增加两个结点i和f, 并且增加 ϵ 边, 将i连接到M的所有初态结点, 并将M的所有终态结点连接到f。形成一个新的与M等价的NFA N。
- 然后, 反复利用下面的替换规则, 逐步消去N中的中间结点, 直到只剩下结点i和f为止。



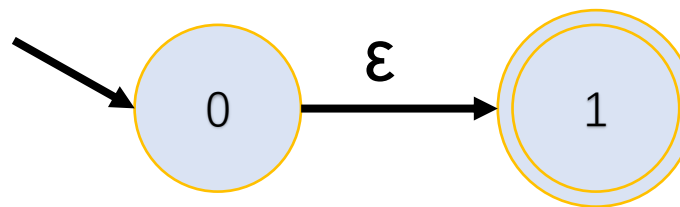
2.5.4 RE到NFA——Thompson算法

- 对RE的结构进行归纳
 - 对基本的RE直接构造
 - 对复合的RE递归构造

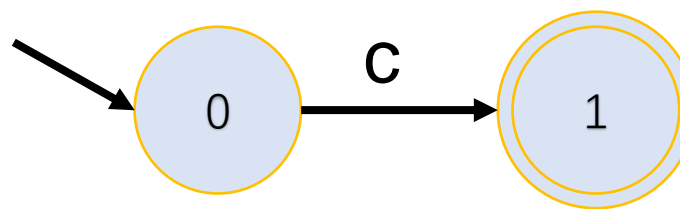


2.5.4 RE到NFA——Thompson算法

$e \rightarrow \varepsilon$



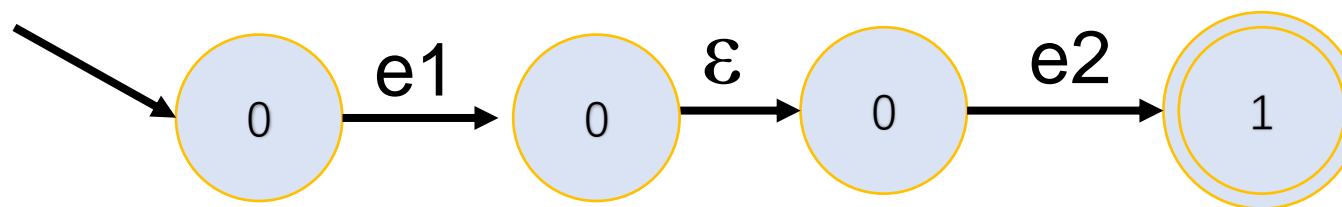
$e \rightarrow c$



$e \rightarrow e1 e2$

$e \rightarrow e1 / e2$

$e \rightarrow e1^*$



2.5.4 RE到NFA——Thompson算法

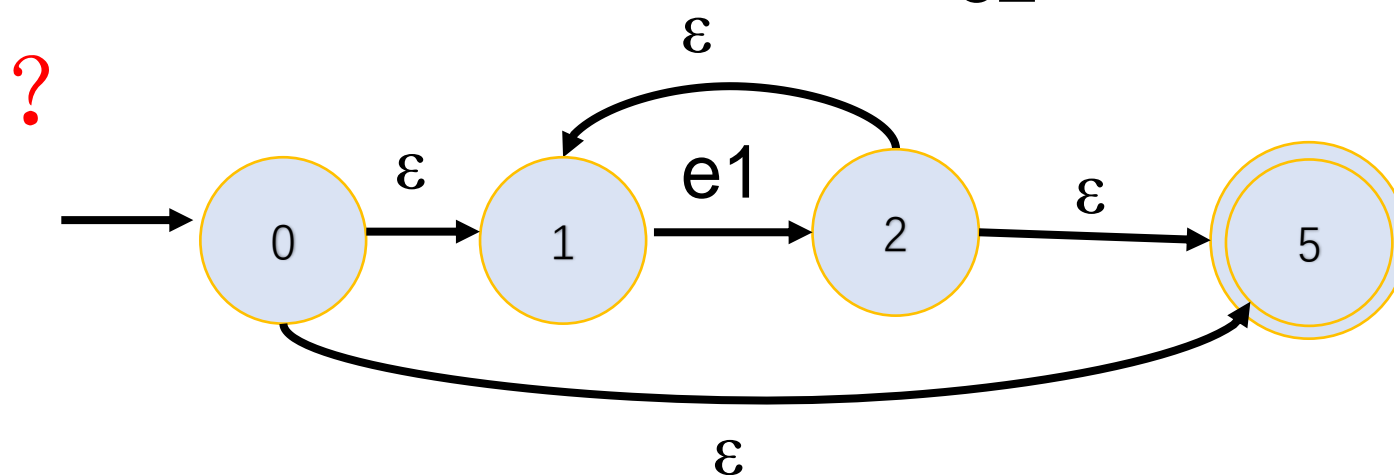
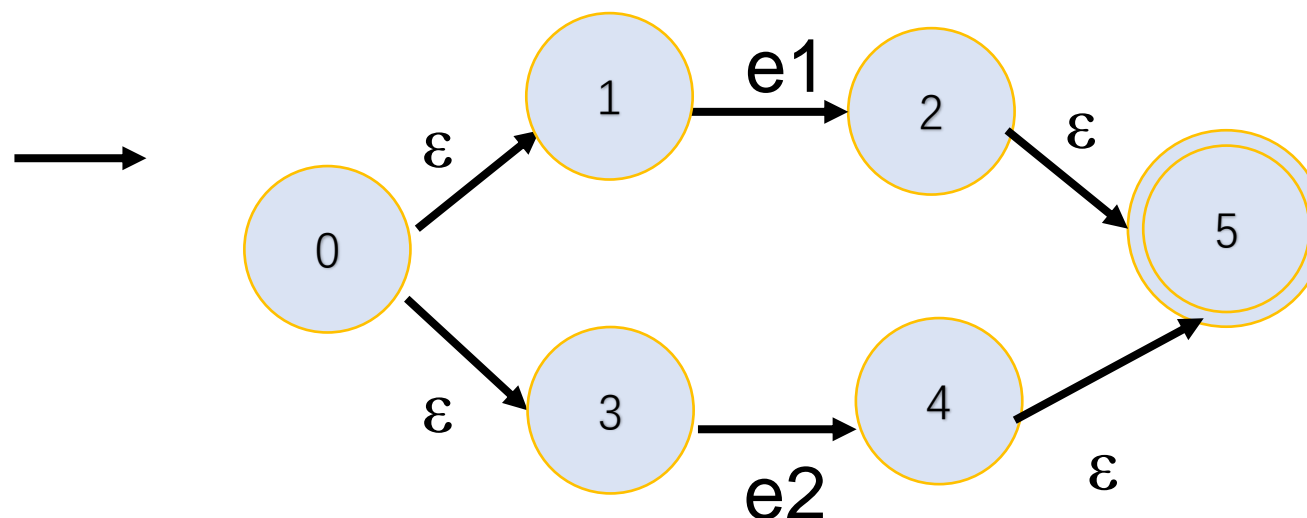
$e \rightarrow \epsilon$

$e \rightarrow c$

$e \rightarrow e_1 e_2$

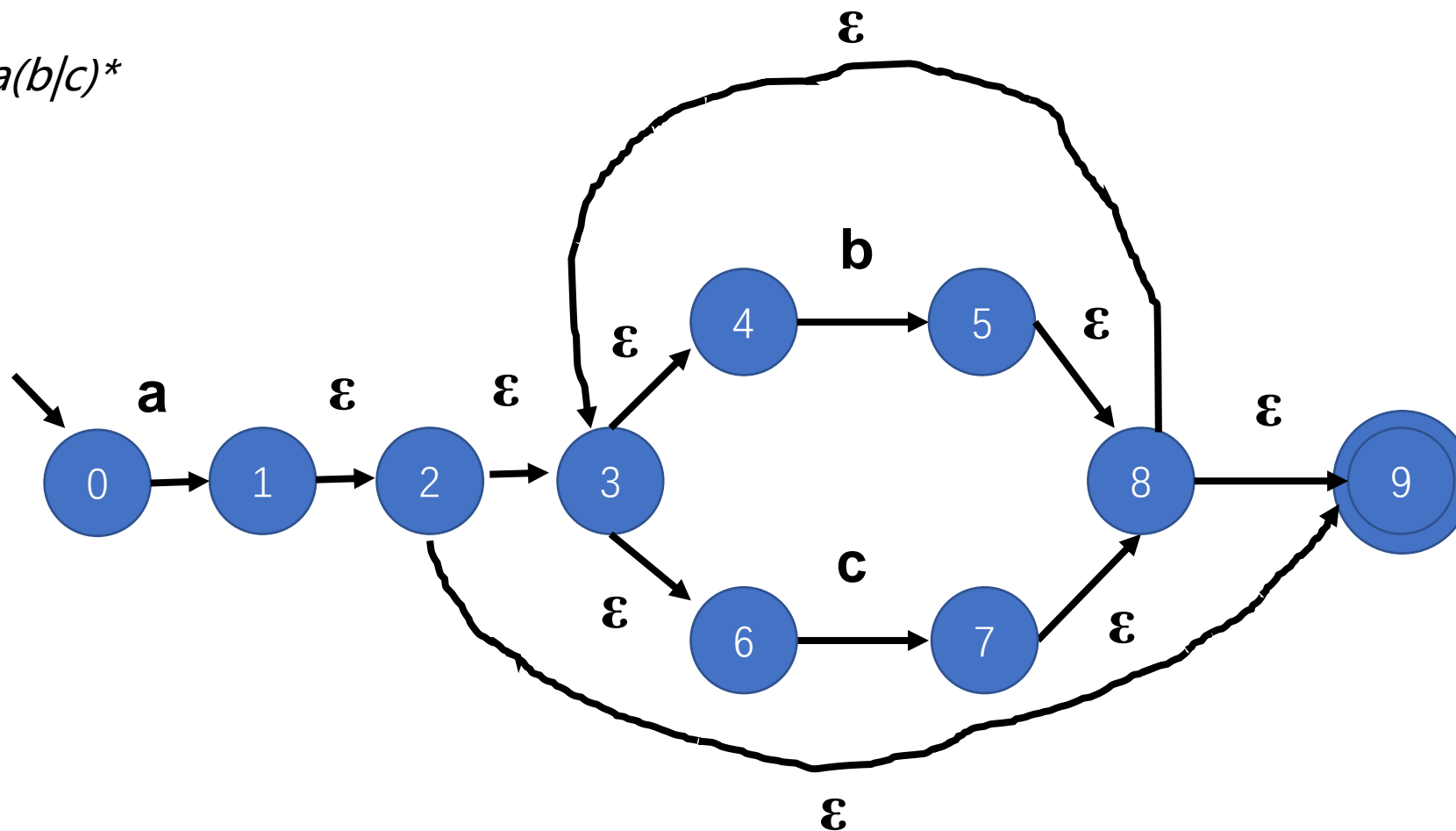
$e \rightarrow e_1 \mid e_2$

$e \rightarrow e_1^*$

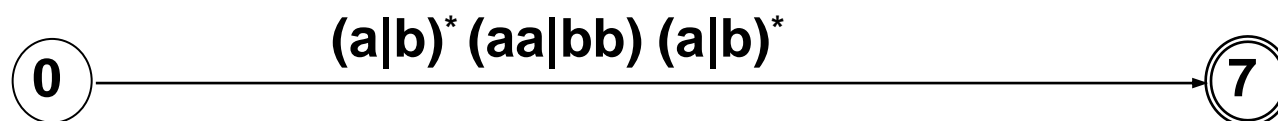
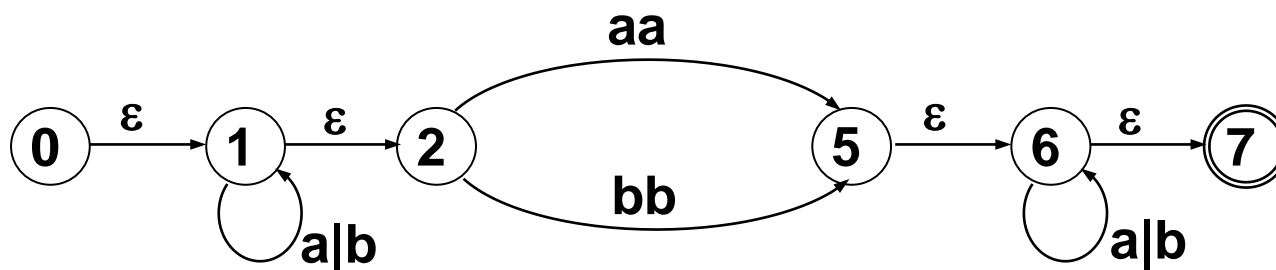
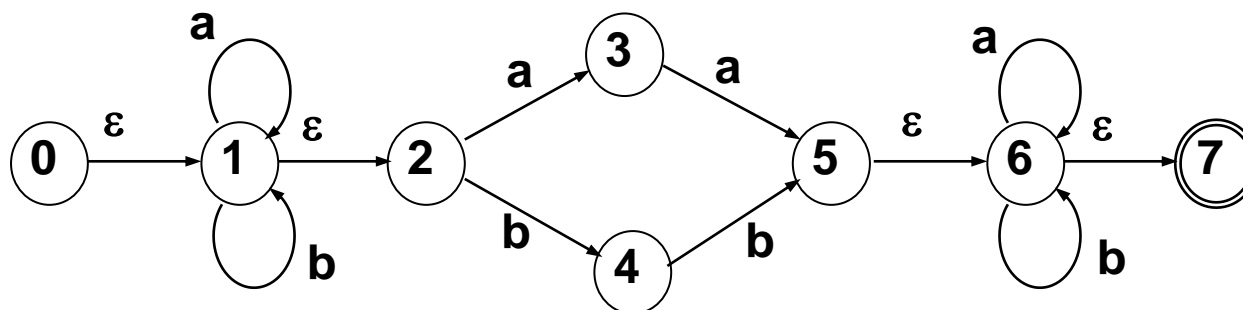


2.5.4 RE到NFA——Thompson算法

$a(b/c)^*$



示例：构造与如下的NFAM等价的正规表达式





2.5.5 正规表达式与正规文法的等价性

- 正规表达式与正规文法具有**同样的表达能力**
- 对任何一个正规表达式都可以找到一个正规文法，使这个正规文法所产生的语言（即正规语言）恰好是该正规表达式所表示的语言（即正规集），反之亦然。
- 正规表达式和正规文法都可以用来描述程序设计语言中单词符号的结构
 - 用正规表达式描述，清晰而简洁；
 - 用正规文法描述，易于识别。



2.5.5.1 正规定义式

定义：令 Σ 是字母表，正规定义式是如下形式的定义序列：

$$d_1 \rightarrow r_1$$
$$d_2 \rightarrow r_2$$
$$\dots$$
$$d_n \rightarrow r_n$$

其中 d_i 是不同的名字， r_i 是 $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$ 上的正规表达式。

例：Pascal语言的无符号数可用如下的正规表达式来描述：

$\text{digit}^+(\text{.digit}^+|\epsilon)(\text{E}(+|-|\epsilon)\text{digit}^+|\epsilon)$

正规定义式：

$$\text{digit} \rightarrow 0|1|\dots|9$$
$$\text{digits} \rightarrow \text{digit digit}^*$$



2.5.5.1 正规定义式

表示的缩写

引入正闭包运算符 ‘+’

$$r^* = r^+ | \epsilon, \quad r^+ = rr^*$$

$$\text{digits} \rightarrow \text{digit}^+$$

引入可选运算符 ‘?’

$$r? = r | \epsilon$$

$$\text{optional_fraction} \rightarrow (. \text{digits})?$$

$$\text{optional_exponent} \rightarrow (E(+|-)? \text{digits})?$$

引入表示 ‘[...]’

字符组[abc], 表示正规表达式a|b|c

$$\text{digit} \rightarrow [0-9]$$

标识符的正规表达式: $[A-Za-z][A-Za-z0-9]^*$



2.5.5.1 正规定义式

正规文法的产生式和正规定义式中的正规定义

两个不同的概念，具有不同的含义。

产生式：左部是一个非终结符号，右部是一个符合特定形式的文法符号串 α ， α 中的非终结符号可以与该产生式左部的非终结符号相同，即允许非终结符号的递归出现。

正规定义：左部是一个名字，右部是一个正规表达式，表达式中出现的名字是有限制的，即只能是此定义之前已经定义过的名字。



2.5.5.1 正规定义式

正规表达式转换为等价的正规文法

例：Pascal语言标识符的正规表达式：

$\text{letter}(\text{letter}|\text{digit})^*$

正规定义式：

$\text{letter} \rightarrow A|B|\dots|Z|a|b|\dots|z$

$\text{digit} \rightarrow 0|1|\dots|9$

$\text{id} \rightarrow \text{letter}(\text{letter}|\text{digit})^*$



2.5.5.1 正规定义式

分析:

为子表达式 $(\text{letter}|\text{digit})^*$ 取一个名字 rid , 展开第三个正规定义

$$\begin{aligned}(\text{letter}|\text{digit})^* &= \epsilon | (\text{letter}|\text{digit})^+ \\&= \epsilon | (\text{letter}|\text{digit}) (\text{letter}|\text{digit})^* \\&= \epsilon | \text{letter}(\text{letter}|\text{digit})^* | \text{digit}(\text{letter}|\text{digit})^* \\&= \epsilon | (A|B|\cdots|Z|a|b|\cdots|z)(\text{letter}|\text{digit})^* | (0|1|\cdots|9)(\text{letter}|\text{digit})^* \\&= \epsilon | A(\text{letter}|\text{digit})^* | B(\text{letter}|\text{digit})^* | \cdots | Z(\text{letter}|\text{digit})^* | a(\text{letter}|\text{digit})^* | b(\text{letter}|\text{digit})^* \\&\quad | \cdots | z(\text{letter}|\text{digit})^* | 0(\text{letter}|\text{digit})^* | 1(\text{letter}|\text{digit})^* | \cdots | 9(\text{letter}|\text{digit})^*\end{aligned}$$

id 和 rid 看成是文法的非终结符号, 产生式:

$$\begin{aligned}\text{id} &\rightarrow A \text{ rid} | B \text{ rid} | \cdots | Z \text{ rid} | a \text{ rid} | b \text{ rid} | \cdots | z \text{ rid} \\ \text{rid} &\rightarrow \epsilon | A \text{ rid} | B \text{ rid} | \cdots | Z \text{ rid} | a \text{ rid} | b \text{ rid} | \cdots | z \text{ rid} \\ &\quad | 0 \text{ rid} | 1 \text{ rid} | \cdots | 9 \text{ rid}\end{aligned}$$



2.6 DFAM的实现及应用

2.6.1 DFA表示方式

- 状态转换表
- 状态转换图
- 映射函数集合
- 状态转换矩阵



2.6 DFAM的实现及应用

2.6.2 DFA实现方式

实际中，有不同的DFA代码表示

转移表（类似于邻接矩阵）

哈希表

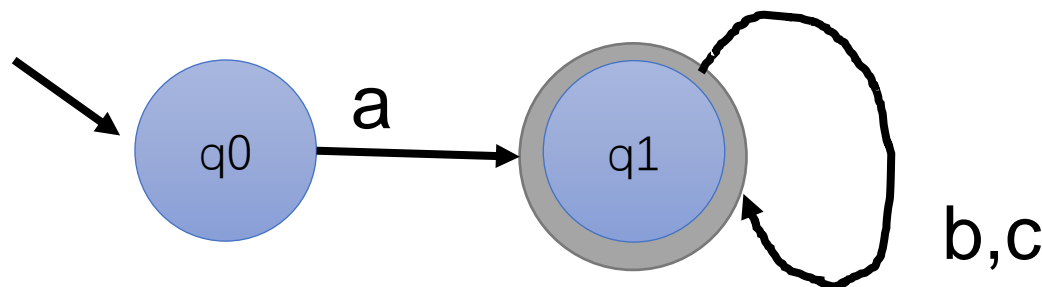
跳转表

...

取决于在实际实现中，对时间空间的权衡

2.6 DFAM的实现及应用

转移表实现



状态/ 字符	a	b	c
0	1		
1		1	1

```
char table[M][N];
```

```
table[0]['a']=1;
```

```
table[1]['b']=1;
```

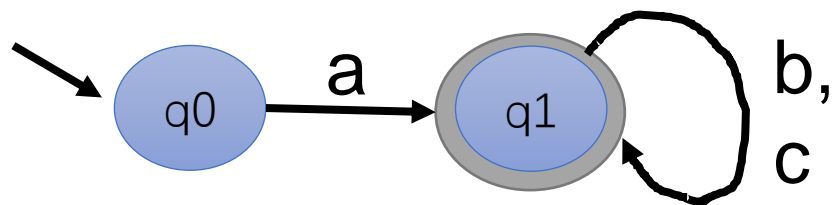
```
table[1]['c']=1;
```

```
// other table entries
```




2.6 DFAM的实现及应用

跳转表实现



状态/字符	a	b	c
0	1		
1		1	1

```
nextToken()  
  state = 0  
  stack = []  
  goto q0
```

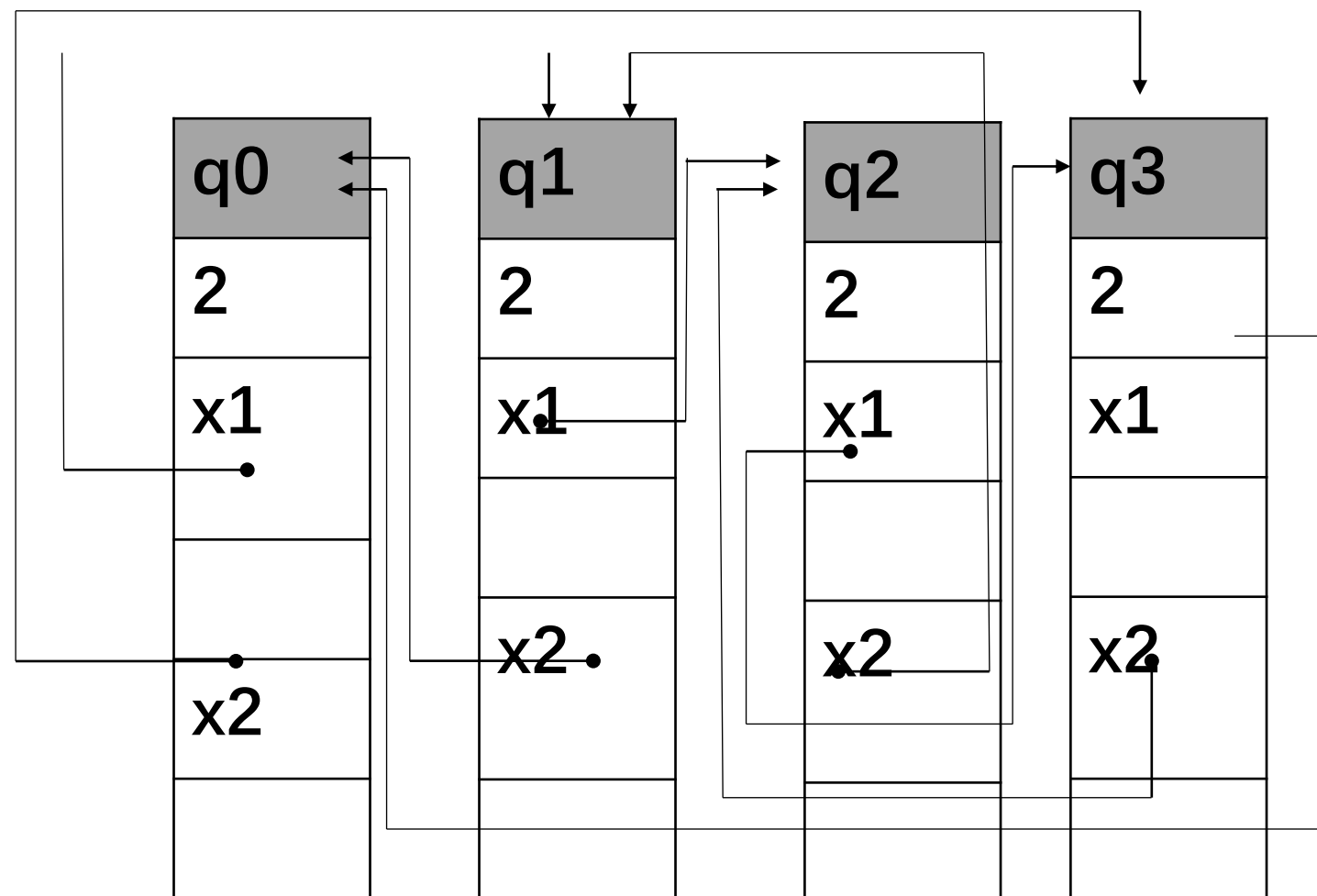
```
q0:  
  c = getChar()  
  if (state is ACCEPT)  
    clear (stack)  
  push (state)  
  if (c=='a') goto q1:
```

```
q1:  
  c = getChar()  
  if (state is ACCEPT)  
    clear (stack)  
  push (state)  
  if (c=='b' || c=='c')  
    goto q1
```



2.6 DFAM的实现及应用

链表实现





2.6.3 DFA应用

示例

某公共汽车始发站，在始发站有乘客的条件下，每隔一固定时间段发出一班车。车站发出一班车后，如果在下一发车时刻车站没有乘客，则停一班次，当再等到下一发车时刻时，不管始发站有没有乘客，则车站必发出一班车。

用有限自动机刻画这一过程：设 x_0 ， x_1 分别表示始发站没有乘客和有乘客，为自动机的输入信息；自动机的状态有两个，分别为上一发车时刻没有发车和发车两个状态，我们分别用 q_0 ， q_1 表示这两个状态， q_0 为初始状态。

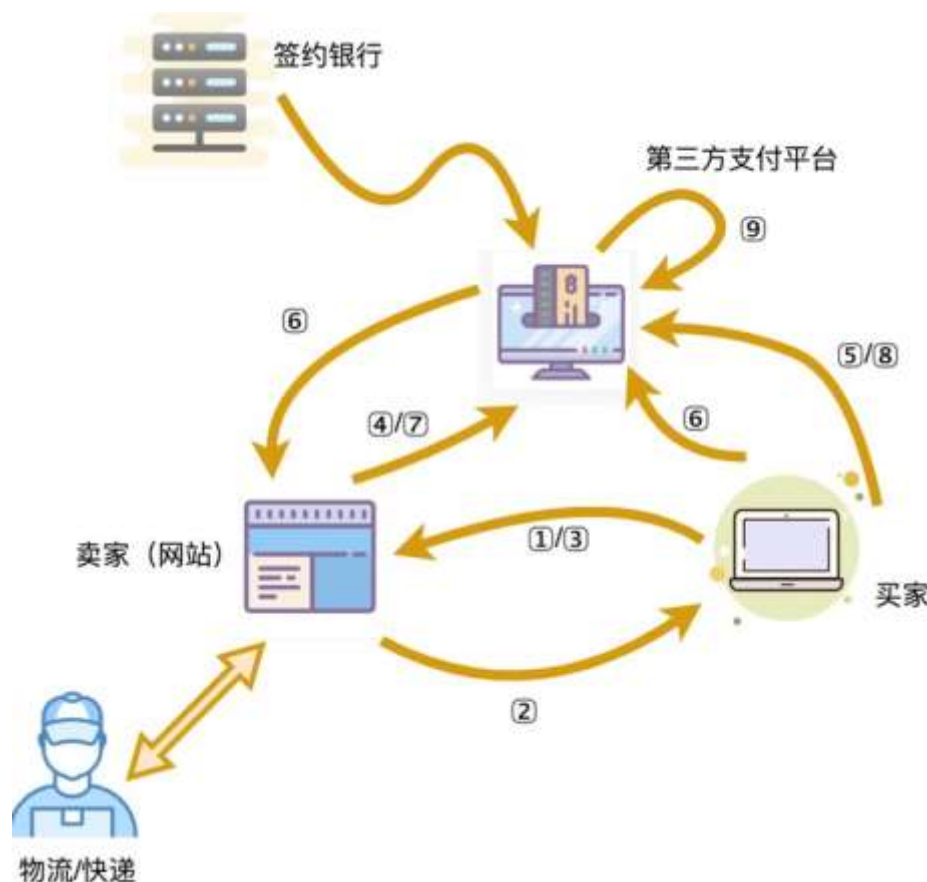
则有限自动机 $M = (Q, \Sigma, \delta, q_0, F)$ 可表述如下：

$$Q = \{q_0, q_1\}, \Sigma = \{x_0, x_1\}, F \subseteq Q$$

$$\delta = \{ \delta(q_0, x_0) = q_1, \delta(q_0, x_1) = q_1, \delta(q_1, x_0) = q_0, \delta(q_1, x_1) = q_1 \}$$

2.6.3 DFA应用

DFA在网上购物平台中的应用



交易具体步骤:

1. 商品浏览
2. 商品推送
3. 订单确认
4. 订单提交
5. 货款支付
6. 支付响应
7. 发货确认
8. 收货确认
9. 交易完成



课堂练习

请给出描述上述网上购物流程的DFA。

操作	动作描述
e1	买家下单
e2	买家付款
e3	买家取消订单
e4	买家确认收货
e5	买家退货
e6	卖家发货
e7	卖家确认退款
e8	卖家取消订单
e9	其他操作

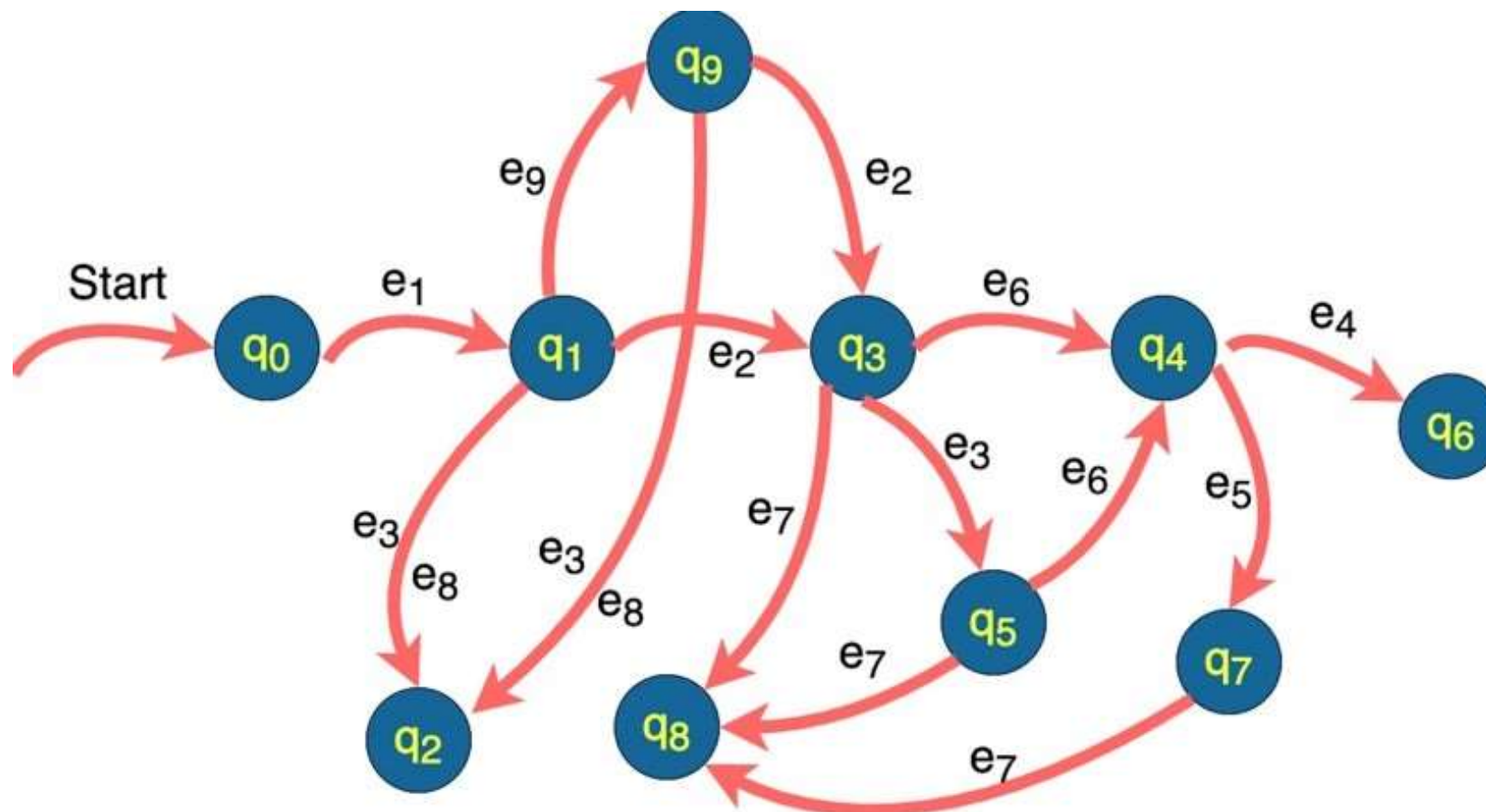
交易操作描述/输入

状态	状态描述
q0	开始状态，等待订单
q1	已下单状态
q2	终结状态1， 订单取消
q3	已付款状态
q4	已发货状态
q5	申请退款状态
q6	终结状态2， 已收货状态
q7	已退货状态
q8	终结状态3， 确认退款状态
q9	其他状态

DFA状态描述

2.6.3 DFA应用

购物流程DFA状态图





2.6.3 DFA应用

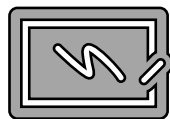
实际应用

设备管理：车辆管理

业务流程管理：航班调度

嵌入式菜单：嵌入式HTML解析器

自然语言识别：微博舆情监控



作业 1

1、已知语言描述，写出产生该语言的文法。

语言由0、1符号串组成，串中0和1的个数相同，构造描述该语言的文法。

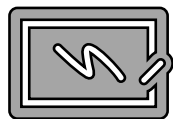
2、已知文法，写出所产生语言的描述。

G: $S \rightarrow a$

$S \rightarrow aB$

$B \rightarrow aS$





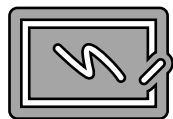
作业 2

根据语法规则，扩展编程作业3，实现正整数运算（加、减、乘、除）的解释器。

输入：正整数四则运算表达式

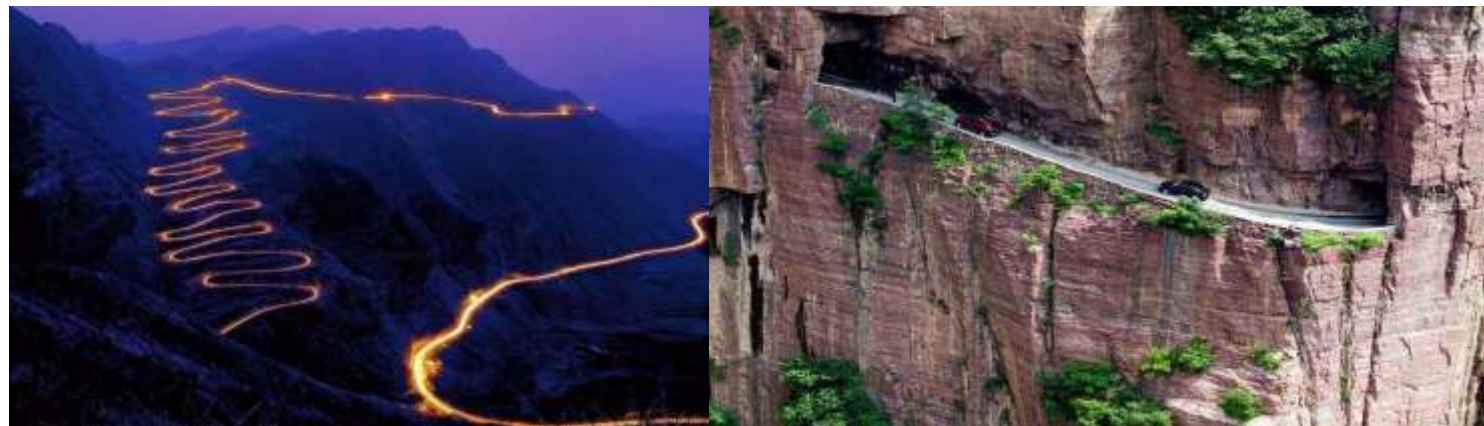
输出：运算结果

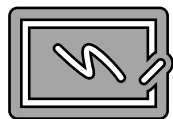




作业 3

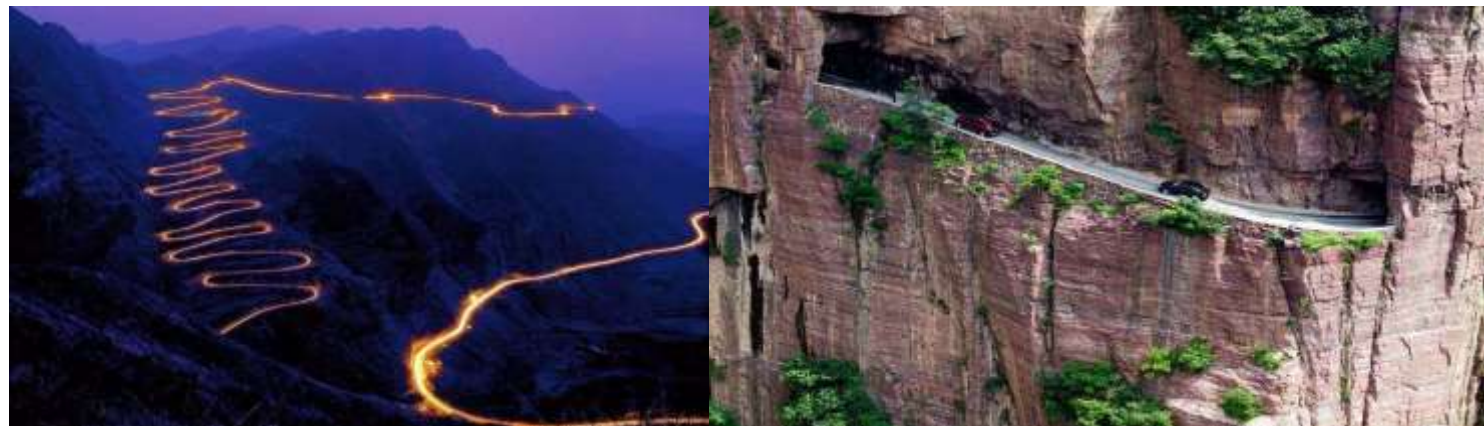
构造识别输入串 $(a|b)^*a(a|b)^*$ 的DFA

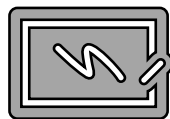




作业 4

构造识别输入串 $(a|b)^*a(a|b)^*$ 的DFA





作业 5

美团共享电动车

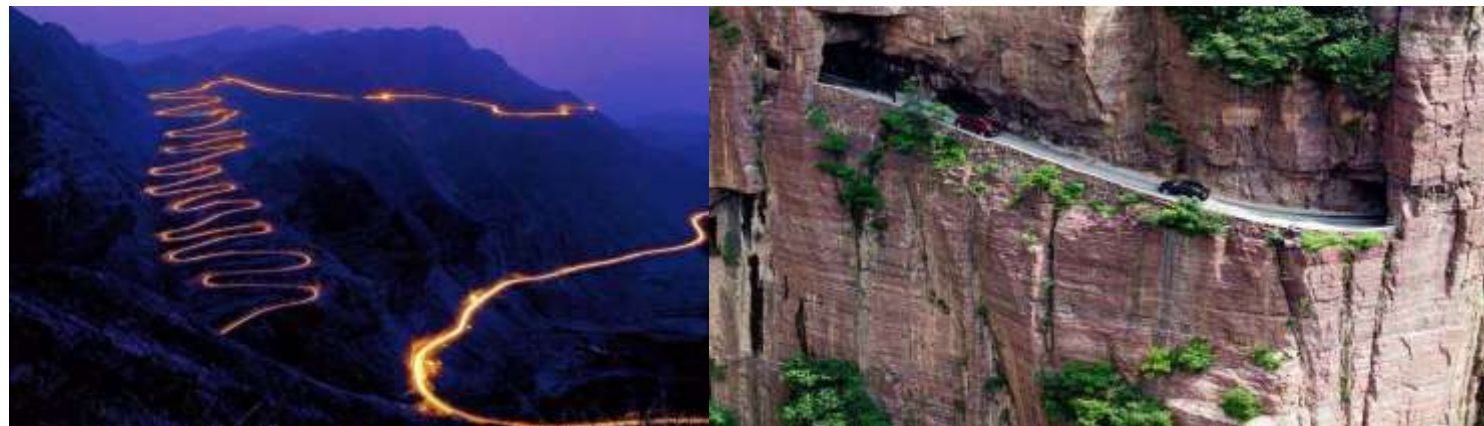
快递收发柜

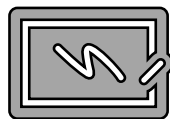
扫地机器人

无人驾驶汽车

...

请给出一个应用实例，并用DFA进行描述。





课后扩展 1

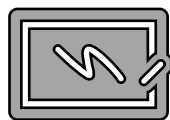
ABNF文法是扩展的BNF，在互联网领域应用广泛，请阅读相关资料。

<http://oss.org.cn/man/develop/rfc/RFC2234.txt>).

[科大讯飞ABNF文法规范](#)

[科大讯飞翻译器](#)

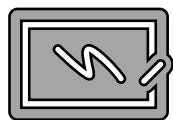




课后扩展 2

课后请查阅资料并思考：文法还能应用于哪些领域？你能给出一些实例吗？





课后扩展 3

课后请阅读：文法实验解读。

