

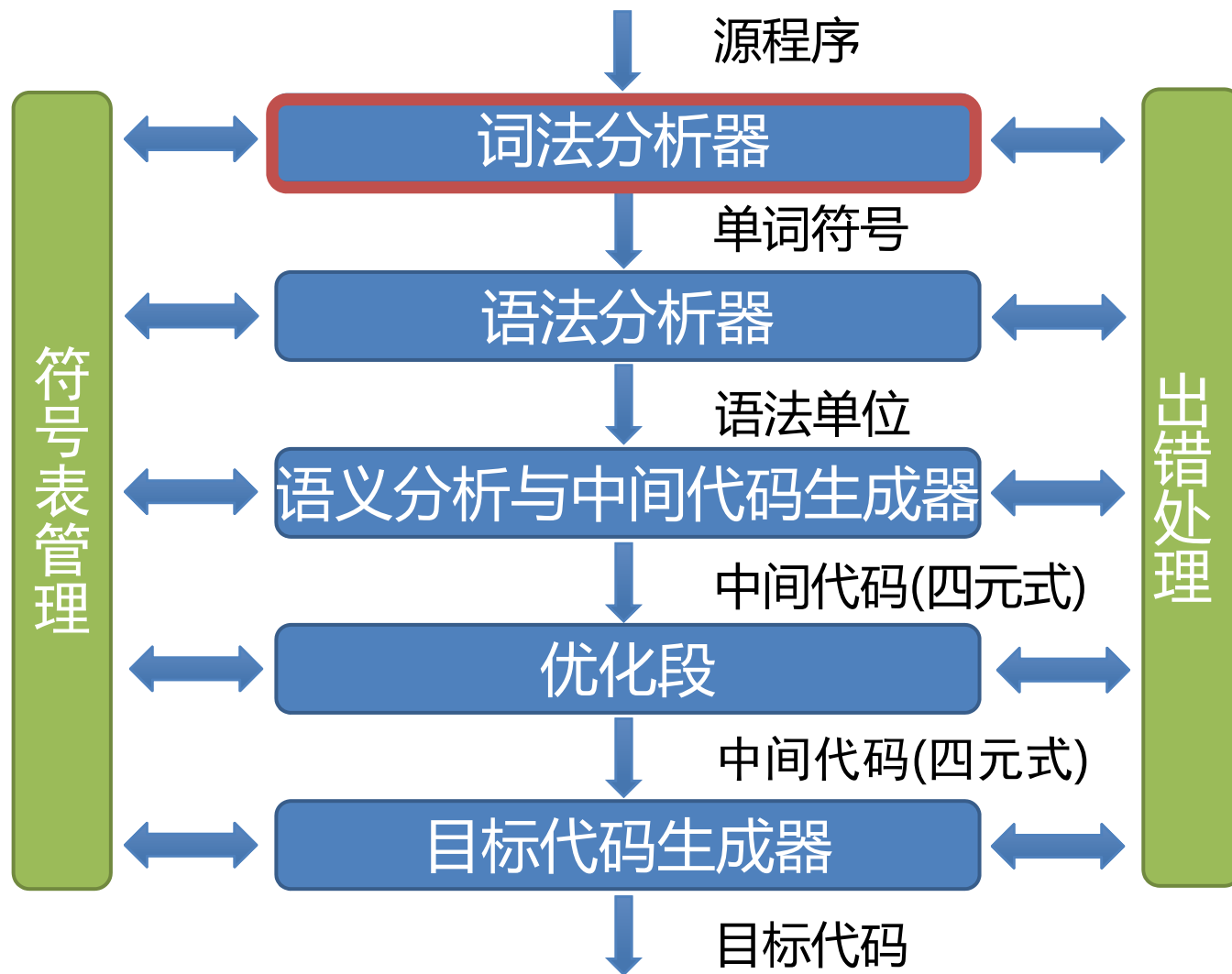
编译原理

词法规则的形式化

编译原理

词法规则的形式化 ——回顾

编译程序总框



词法分析

- ▶ 词法分析器的设计
- ▶ 正规表达式与有限自动机
- ▶ 词法分析器的自动产生--LEX

词法分析器的设计

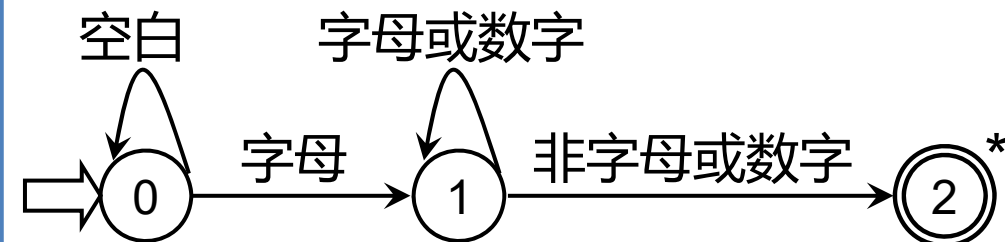
- ▶ 词法分析器的功能

- ▶ 输入源程序、输出单词符号

- ▶ 词法分析器的设计

- ▶ 给出程序设计语言的单词规范——单词表
 - ▶ 对照单词表设计识别该语言所有单词的状态转换图
 - ▶ 根据状态转换图编写词法分析程序

将状态图的代码



- ▶ 变量curState用于保存现有的状态
- ▶ 用二维数组表示状态图：stateTrans[state][ch]

```
curState = 初态
GetChar();
while( stateTrans[curState][ch]有定义) {
    //存在后继状态, 读入、拼接
    Concat();
    //转换入下一状态, 读入下一字符
    curState= stateTrans[curState][ch];
    if curState是终态 then 返回strToken中的单词
    GetChar();
}
```

词法分析

- ▶ 词法分析器的设计
- ▶ 正规表达式与有限自动机
- ▶ 词法分析器的自动产生--LEX

语法描述的几个基本概念

- ▶ **字母表**：一个有穷字符集，记为 Σ
- ▶ **字母表**中每个元素称为**字符**
- ▶ Σ 上的**字**(也叫**字符串**) 是指由 Σ 中的字符所构成的一个有穷序列
- ▶ 不包含任何字符的序列称为**空字**，记为 ε
- ▶ 用 Σ^* 表示 Σ 上的所有**字的全体**，包含空字 ε
- ▶ 例如：设 $\Sigma = \{a, b\}$ ，则

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

语法描述的几个基本概念

- ▶ Σ^* 的子集U和V的**连接** (**积**) 定义为

$$UV = \{ \alpha\beta \mid \alpha \in U \ \& \ \beta \in V \}$$

- ▶ V自身的 n次积记为

$$\underbrace{V^n}_{n\uparrow} = V \ V \dots V$$

- ▶ $V^0 = \{\varepsilon\}$

- ▶ V^* 是V的**闭包**: $V^* = V^0 \cup V^1 \cup V^2 \cup V^3 \cup \dots$

- ▶ V^+ 是V的**正规闭包**: $V^+ = V \ V^*$

编译原理

词法规则的形式化 ——正规集和正规式

单词集合、正规集和正规式

- ▶ 程序设计语言的单词符号都是一些特殊的字符串
- ▶ 用**正规集**和**正规表达式**（简称**正规式**）来描述

单词符号	种别编码	助忆符	内码值
DIM	1	\$DIM	-
IF	2	\$IF	-
DO	3	\$DO	-
STOP	4	\$STOP	-
END	5	\$END	-
标识符	6	\$ID	内部字符串
常数(数)	7	\$INT	标准二进制形式
=	8	\$ASSIGN	-
+	9	\$PLUS	-
*	10	\$STAR	-
**	11	\$POWER	-
,	12	\$COMMA	-
(13	\$LPAR	-
)	14	\$RPAR	-

正规式和正规集

- ▶ 正规集可以用正规式表示
- ▶ 正规式是表示正规集一种方法
- ▶ 一个字集合是正规集当且仅当它能用正规式表示

正规式和正规集的递归定义

► 对给定的字母表 Σ

- 1) ϵ 和 \emptyset 都是 Σ 上的正规式, 它们所表示的正规集为 $\{\epsilon\}$ 和 \emptyset ;

测试： ϵ 是什么？

A. 字符

B. 字 \checkmark

C. 正规式 \checkmark

测试： \emptyset 是什么？

A. 集合 \checkmark

B. 字

C. 正规式 \checkmark

正规式和正规集的递归定义

► 对给定的字母表 Σ

- 1) ε 和 \emptyset 都是 Σ 上的正规式，它们所表示的正规集为 $\{\varepsilon\}$ 和 \emptyset ;
- 2) 任何 $a \in \Sigma$ ， a 是 Σ 上的正规式，它所表示的正规集为 $\{a\}$;

测试: a ($a \in \Sigma$)是什么?

A. 字符 ✓

B. 字 ✓

C. 正规式 ✓

正规式和正规集的递归定义

► 对给定的字母表 Σ

- 1) ε 和 \emptyset 都是 Σ 上的正规式，它们所表示的正规集为 $\{\varepsilon\}$ 和 \emptyset ;
- 2) 任何 $a \in \Sigma$ ， a 是 Σ 上的正规式，它所表示的正规集为 $\{a\}$;

正规式和正规集的递归定义 (续)

3) 假定 e_1 和 e_2 都是 Σ 上的正规式, 它们所表示的正规集为 $L(e_1)$ 和 $L(e_2)$, 则

i) $(e_1|e_2)$ 为正规式, 它所表示的正规集为
 $L(e_1) \cup L(e_2)$

ii) $(e_1.e_2)$ 为正规式, 它所表示的正规集为
 $L(e_1)L(e_2)$

iii) $(e_1)^*$ 为正规式, 它所表示的正规集为
 $(L(e_1))^*$

仅由有限次使用上述三步骤而定义的表达式才是 Σ 上的正规式, 仅由这些正规式表示的字集才是 Σ 上的正规集。

正规式的等价性

- 若两个正规式所表示的正规集相同，则称这两个正规式**等价**。如

$$b(ab)^* = (ba)^*b$$

$$\begin{aligned} & L(b(ab)^*) \\ &= L(b)L((ab)^*) \\ &= L(b) (L(ab))^* \\ &= L(b) (L(a)L(b))^* \\ &= \{b\} \{ab\}^* \\ &= \{b\} \{\varepsilon, ab, abab, ababab, \dots\} \\ &= \{b, bab, babab, bababab, \dots\} \end{aligned}$$

$$\begin{aligned} & L((ba)^*b) \\ &= L((ba)^*) L(b) \\ &= (L(ba))^* L(b) \\ &= (L(b)L(a))^* L(b) \\ &= \{ba\}^* \{b\} \\ &= \{\varepsilon, ba, baba, bababa, \dots\} \{b\} \\ &= \{b, bab, babab, bababab, \dots\} \end{aligned}$$

$$\therefore L(b(ab)^*) = L((ba)^*b)$$

$$\therefore b(ab)^* = (ba)^*b$$

作业

- ▶ 利用正规式与正规集的对应关系, 证明 $(a^*b^*)^* = (a|b)^*$

正规式的性质

▶ 对正规式，下列等价成立

▶ $e_1|e_2 = e_2|e_1$ 交换律

▶ $e_1|(e_2|e_3) = (e_1|e_2)|e_3$ 结合律

▶ $e_1(e_2e_3) = (e_1e_2)e_3$ 结合律

▶ $e_1(e_2|e_3) = e_1e_2|e_1e_3$ 分配律

▶ $(e_2|e_3)e_1 = e_2e_1|e_3e_1$ 分配律

▶ $e\varepsilon = \varepsilon e = e$ $e_1e_2 <> e_2e_1$

$$\begin{aligned} L(e_1|e_2) &= L(e_1) \cup L(e_2) \\ &= L(e_2) \cup L(e_1) \\ &= L(e_2|e_1) \end{aligned}$$

小结：正规式和正规集

单词符号	种别编码	助忆符	内码值
DIM	1	\$DIM	-
IF	2	\$IF	-
DO	3	\$DO	-
STOP	4	\$STOP	-
END	5	\$END	-
标识符	6	\$ID	内部字符串
常数(数)	7	\$INT	标准二进制形式
=	8	\$ASSIGN	-
+	9	\$PLUS	-
*	10	\$STAR	-
**	11	\$POWER	-
,	12	\$COMM A	-
(13	\$LPAR	-
)	14	\$RPAR	-

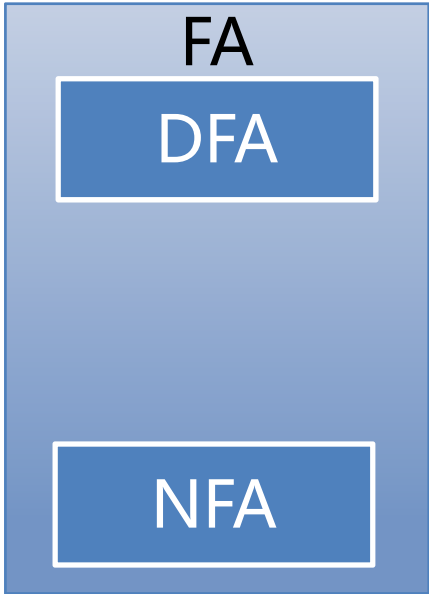
DIM,IF, DO,STOP,END
number, name, age
125, 2169
...

DIM
IF
DO
STOP
END
letter(letter|digit)*
digit(digit)*

正规集



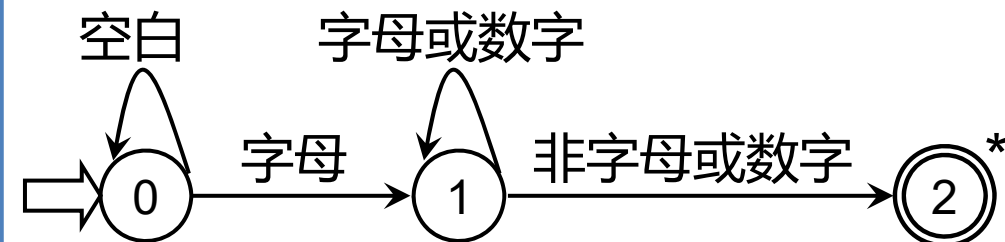
正规式



编译原理

词法规则的形式化
——确定有限自动机(DFA)

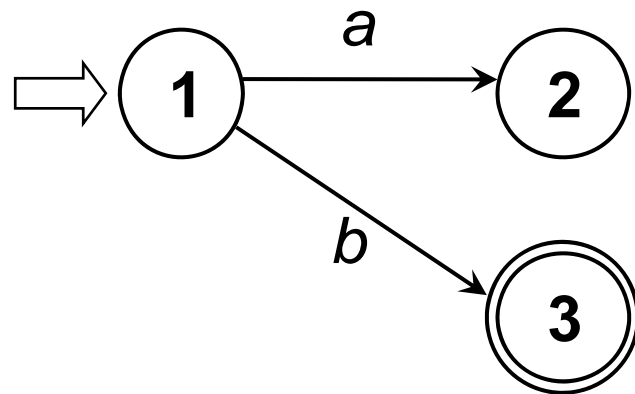
将状态图的代码



- ▶ 变量curState用于保存现有的状态
- ▶ 用二维数组表示状态图：stateTrans[state][ch]

```
curState = 初态
GetChar();
while( stateTrans[curState][ch]有定义) {
    //存在后继状态, 读入、拼接
    Concat();
    //转换入下一状态, 读入下一字符
    curState= stateTrans[curState][ch];
    if curState是终态 then 返回strToken中的单词
    GetChar();
}
```

确定有限自动机



- ▶ 对状态图进行形式化定义
- ▶ **确定有限自动机**(Deterministic Finite Automata, DFA) M 是一个五元式 $M=(S, \Sigma, f, S_0, F)$, 其中:
 1. S : 有穷**状态集**
 2. Σ : 输入**字母表**(有穷)
 3. f : **状态转换函数**, 为 $S \times \Sigma \rightarrow S$ 的**单值部分映射**, $f(s, a)=s'$ 表示: 当现行状态为 s , 输入字符为 a 时, 将状态转换到下一状态 s' , s' 称为 s 的一个后继状态
 4. $S_0 \in S$ 是唯一的一个**初态**
 5. $F \subseteq S$: **终态集**(可空)

确定有限自动机

► DFA $M = (\{0, 1, 2, 3\}, \{a, b\}, f, 0, \{3\})$,
其中 f 定义如下:

$f(0, a) = 1$ $f(0, b) = 2$

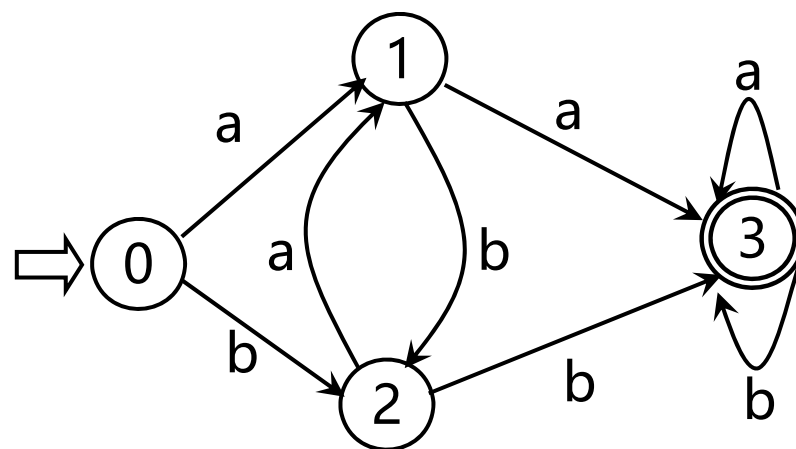
$f(1, a) = 3$ $f(1, b) = 2$

$f(2, a) = 1$ $f(2, b) = 3$

$f(3, a) = 3$ $f(3, b) = 3$

状态
转换
矩阵

	a	b
0	1	2
1	3	2
2	1	3
3	3	3



状态转换图

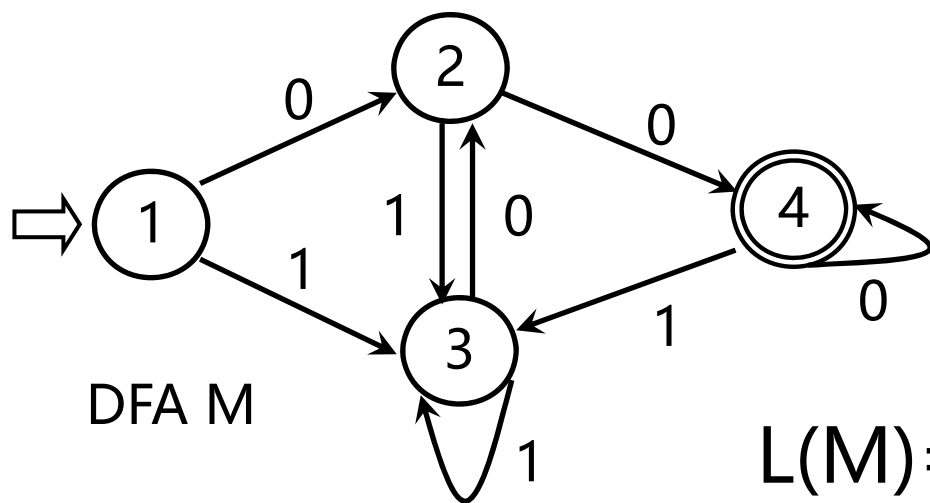
确定有限自动机

► DFA表示为状态转换图

- 假定DFA M 含有 m 个状态和 n 个输入字符
- 对应的状态转换图含有 m 个状态结点，每个结点顶多含有 n 条箭弧射出，且每条箭弧用 Σ 上的不同的输入字符来作标记

确定有限自动机

- ▶ 对于 Σ^* 中的任何字 α ，若存在一条从初态到某一终态的道路，且这条路上所有弧上的标记符连接成的字等于 α ，则称 α 为DFA M 所识别(接收)
- ▶ DFA M 所识别的字的全体记为 $L(M)$



$L(M) = \{\text{以00结尾的串}\}$

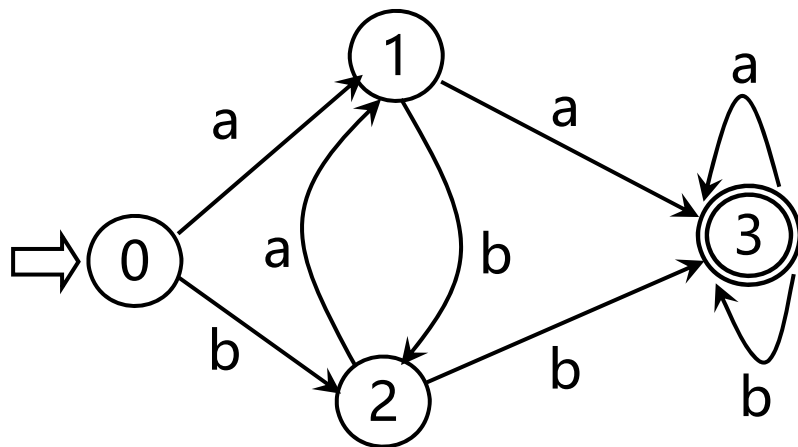
测试

► 图中DFA M识别的 $L(M)$ 是什么?

A. $L(M) = \{\text{以aa或bb开头的字}\}$

B. $L(M) = \{\text{含aa或bb的字}\}$ ✓

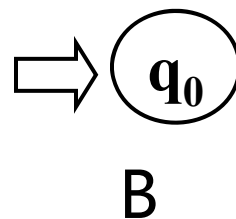
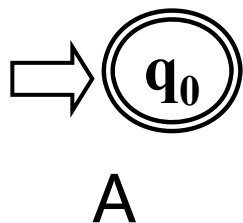
C. $L(M) = \{\text{以aa或bb结尾的字}\}$



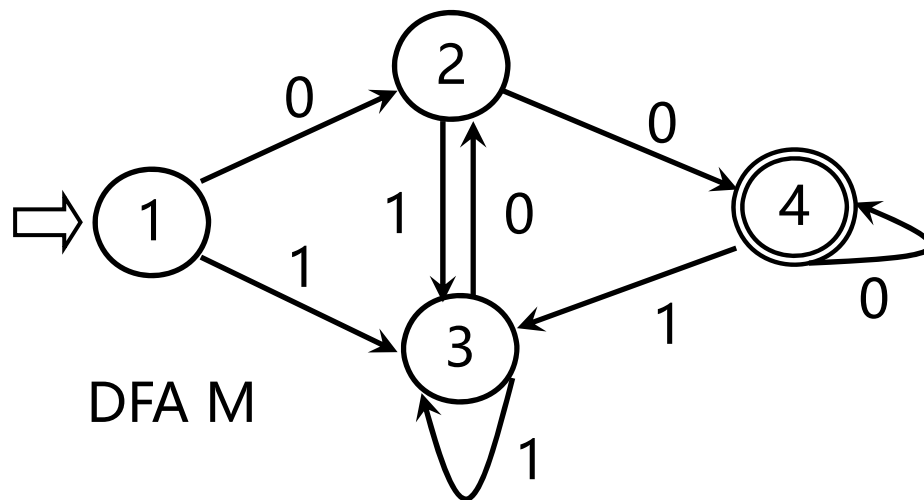
DFA M

测试

► 哪个DFA识别 $\{\epsilon\}$?



DFA的程序实现



```
curState = 初态
GetChar();
while( stateTrans[curState][ch]有定义) {
    //存在后继状态, 读入、拼接
    Concat();
    //转换入下一状态, 读入下一字符
    curState= stateTrans[curState][ch];
    if curState是终态 then 返回strToken中的单词
    GetChar();
}
```


单词符号	种别编码	助忆符	内码值
DIM	1	\$DIM	-
IF	2	\$IF	-
DO	3	\$DO	-

正规式、正规集

DIM,IF, DO,STOP,E
number, name, age
125, 2169
...

```
curState = 初态
GetChar();
while( stateTrans[curState][ch]有定义){
    //存在后继状态, 读入、拼接
    Concat();
    //转换入下一状态, 读入下一字符
    curState= stateTrans[curState][ch];
    if curState是终态 then 返回strToken中的单词
    GetChar();
}
```

DIM
IF
DO
STOP
END
letter(letter|digit)*
digit(digit)*

正规集

正规式

FA

DFA

NFA

编译原理

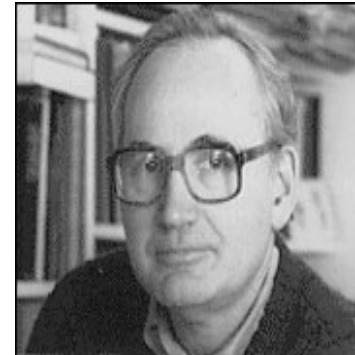
词法规则的形式化 ——非确定有限自动机(NFA)

非确定有限自动机

- 1976年图灵奖: For their joint paper "**Finite Automata and Their Decision Problem**," which introduced the idea of nondeterministic machines, which has proved to be an enormously valuable concept. Their (Scott & Rabin) classic paper has been a continuous source of inspiration for subsequent work in this field.



Michael O. Rabin



Dana S. Scott

非确定有限自动机

- ▶ 定义：一个非确定有限自动机
(Nondeterministic Finite Automata, NFA)
 M 是一个五元式 $M = (S, \Sigma, f, S_0, F)$ ，其中：
 1. S : 有穷状态集
 2. Σ : 输入字母表(有穷)
 3. f : 状态转换函数，为 $S \times \Sigma^* \rightarrow 2^S$ 的部分映射
 4. $S_0 \subseteq S$ 是非空的初态集
 5. $F \subseteq S$: 终态集(可空)

非确定有限自动机

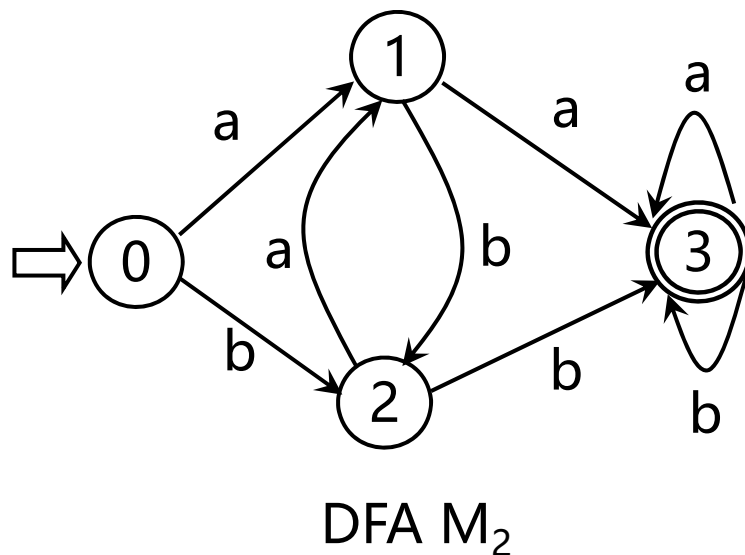
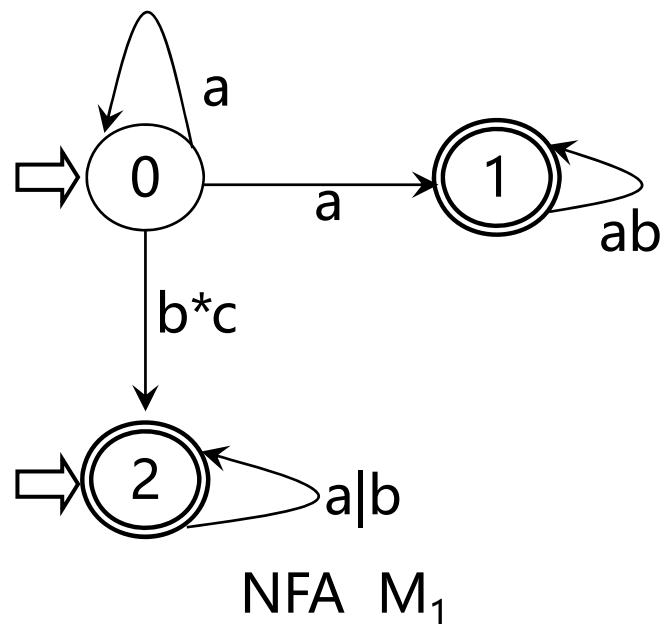
▶ 从状态图看NFA 和DFA的区别

- ▶ NFA可以有多个初态

- ▶ 弧上的标记可以是 Σ^* 中的一个字 (甚至可以是一个正规式), 而不一定是单个字符

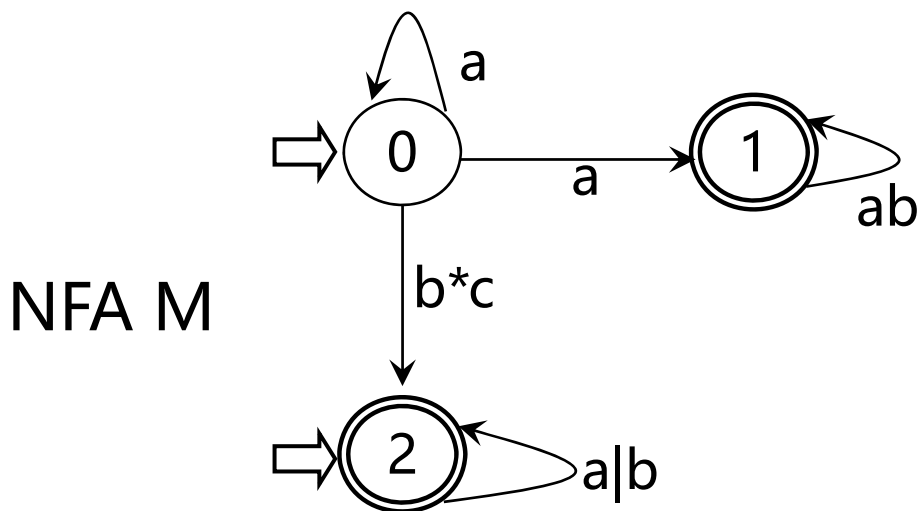
- ▶ 同一个字可能出现在同状态射出的多条弧上

▶ DFA是NFA的特例



非确定有限自动机

- ▶ 对于 Σ^* 中的任何字 α ，若存在一条从初态到某一终态的道路，且这条路上所有弧上的标记字连接成的字等于 α （忽略那些标记为 ε 的弧），则称 α 为NFA M 所**识别(接收)**
- ▶ NFA M 所识别的字的全体记为 $L(M)$



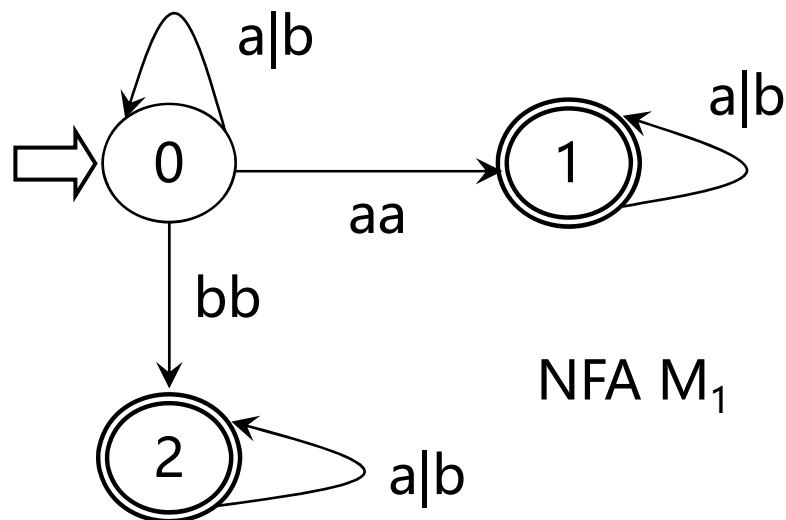
测试

► 图中NFA M_1 识别的 $L(M_1)$ 是什么?

A. $L(M_1) = \{\text{以aa或bb开头的字}\}$

B. $L(M_1) = \{\text{含aa或bb的字}\}$

C. $L(M_1) = \{\text{以aa或bb结尾的字}\}$



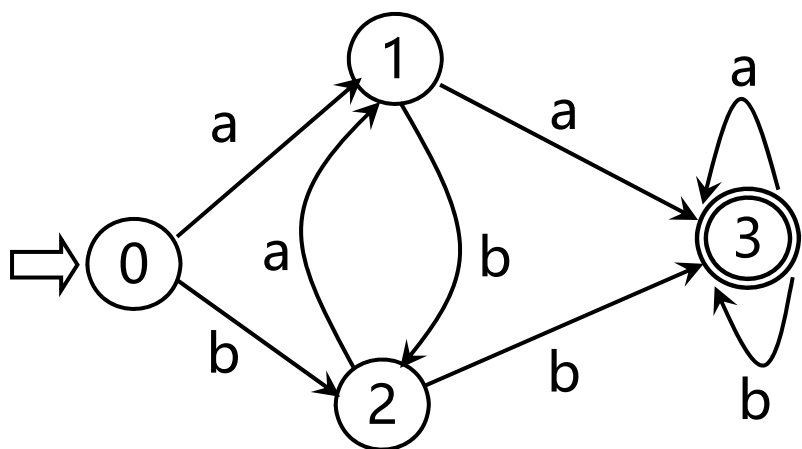
测试

► 图中NFA M_1 识别的 $L(M_1)$ 是什么?

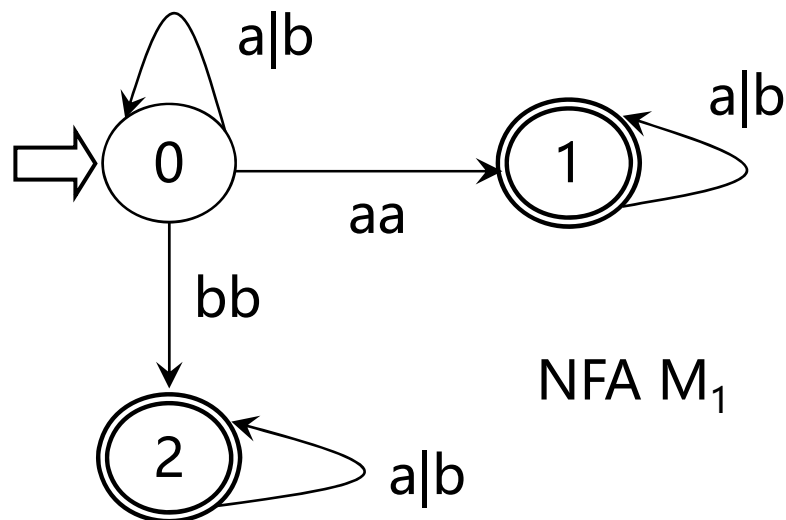
A. $L(M_1) = \{\text{以aa或bb开头的字}\}$

B. $L(M_1) = \{\text{含aa或bb的字}\}$

C. $L(M_1) = \{\text{以aa或bb结尾的字}\}$



DFA M



NFA M_1

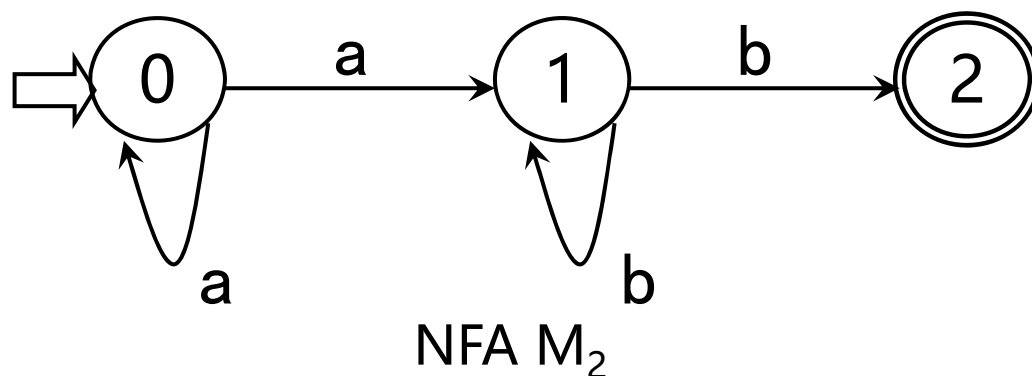
测试

► 图中NFA M_2 识别的 $L(M_2)$ 是什么?

A. $L(M_2) = \{ab^n \mid n \geq 1\}$

B. $L(M_2) = \{a^n b^n \mid n \geq 1\}$

C. $L(M_2) = \{a^m b^n \mid m, n \geq 1\}$



DFA和NFA

- ▶ 定义：对于任何两个有限自动机 M 和 M' ，如果 $L(M)=L(M')$ ，则称 M 与 M' 等价
- ▶ 自动机理论中一个重要的结论：判定两个自动机等价性的算法是存在的
- ▶ 对于每个NFA M 存在一个DFA M' ，使得 $L(M)=L(M')$
- ▶ DFA与NFA识别能力相同!

单词符号	种别编码	助忆符	内码值
DIM	1	\$DIM	-
IF	2	\$IF	-

正规式、正规集

DIM,IF, DO,STOP,END
number, name, age
125, 2169
 ...

DIM
IF
DO
STOP
END
letter(letter|digit)*
digit(digit)*

```

curState = 初态
GetChar();
while( stateTrans[curState][ch]有定义){
    //存在后继状态, 读入、拼接
    Concat();
    //转换入下一状态, 读入下一字符
    curState= stateTrans[curState][ch];
    if curState是终态 then 返回strToken中的单词
    GetChar( );
}
  
```

正规集

正规式

FA

DFA

NFA

易于人工设计

小结

- ▶ 正规式和正规集
- ▶ 确定有限自动机和非确定有限自动机