

程序

```
int main() '\n' { int count=read(); '\n'
//if number of entries read is greater
than 1 '\n' //then sort() and compact()
'\n' if (count > 1) { sort();
compact(); } '\n' if (count ==0) '\n'
count << "no sales for this month\n";
'\n' else write(); '\n' return; '\n' }
```

Excellence in  
BUAA SEI

北京航空航天大学计算机学院

程序

```
int main()
{ int count=read();
//if number of entries read is greater than 1
//then sort() and compact()
if (count > 1) { sort(); compact(); }
if (count ==0)
count << "no sales for this month\n";
else write();
return;
}
```

Excellence in  
BUAA SEI

北京航空航天大学计算机学院

内 容

- 3.1 词法分析程序的功能及实现方案
- 3.2 单词的种类及词法分析程序的输出形式
- 3.3 正则文法和状态图**
- 3.4 词法分析程序的设计与实现
- 11.1-2 正则表达式与有穷自动机
- 11.3 词法分析程序的自动生成器

Excellence in  
BUAA SEI

北京航空航天大学计算机学院

3.3 正则文法和状态图

• 状态图的画法（根据文法画出状态图）

例如：正则文法

$$Z ::= U0 | V1$$

$$U ::= Z1 | 1$$

$$V ::= Z0 | 0$$

左线性文法。该文法所定义的语言为：

$$L(G[Z]) = \{ B^n \mid n > 0 \}, \text{ 其中 } B = \{01, 10\}$$

Excellence in  
BUAA SEI

北京航空航天大学计算机学院

例：正则文法

$$Z ::= U0 | V1$$

$$U ::= Z1 | 1$$

$$V ::= Z0 | 0$$

左线性文法的状态图的画法：

1. 令G的每个非终结符都是一个状态；
2. 设一个开始状态S；
3. 若  $Q ::= T, Q \in V_n, T \in V_t$ , 则:
4. 若  $Q ::= RT, Q, R \in V_n, T \in V_t$ , 则:
5. 按自动机方法, 可加上开始状态和终止状态标志。

Excellence in  
BUAA SEI

北京航空航天大学计算机学院

例：正则文法

$$Z ::= U0 | V1$$

$$U ::= Z1 | 1$$

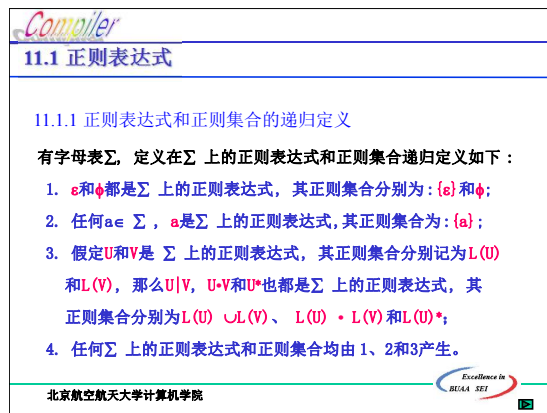
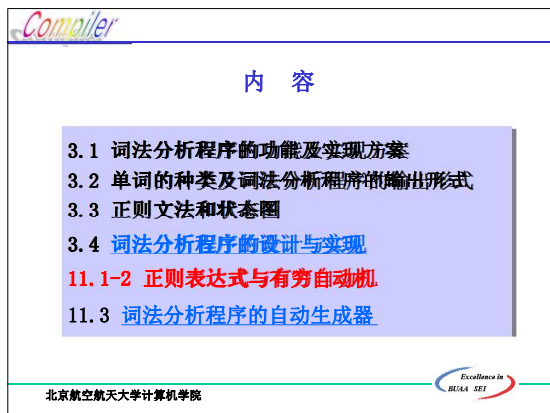
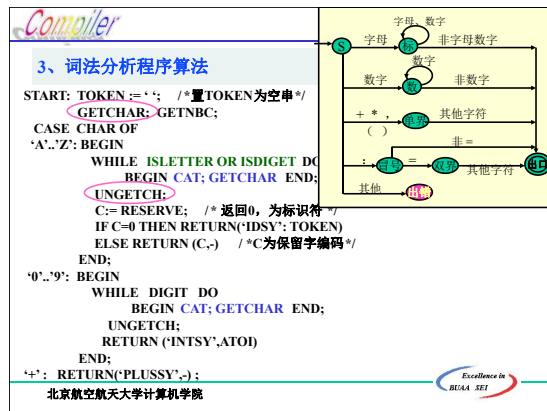
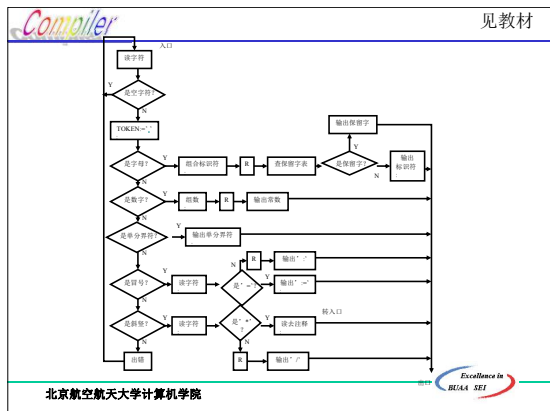
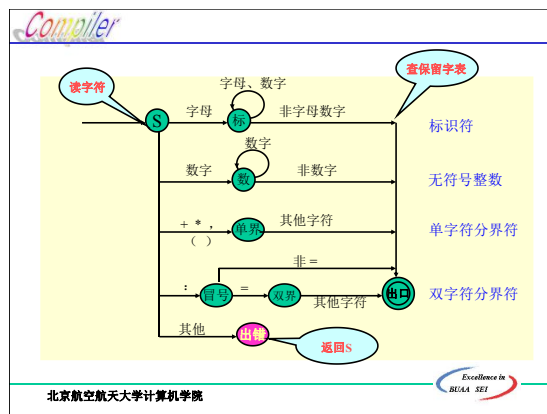
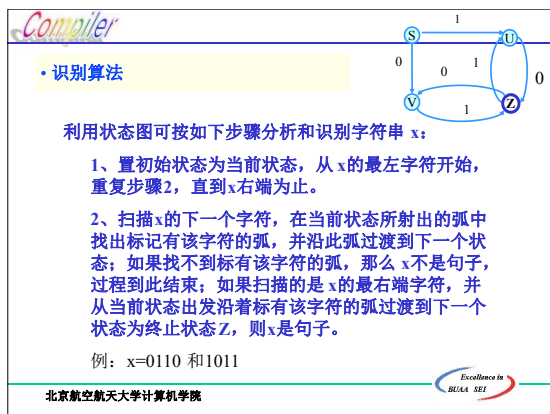
$$V ::= Z0 | 0$$

其状态图为：

1. 每个非终结符设一个状态；
2. 设一个开始状态S；
3. 若  $Q ::= T, Q \in V_n, T \in V_t$ ,
4. 若  $Q ::= RT, Q, R \in V_n, T \in V_t$ ,
5. 加上开始状态和终止状态标志

Excellence in  
BUAA SEI

北京航空航天大学计算机学院



**正则表达式中的运算符：**

	——或（选择）	.	——连接
*	或 { } ——重复	()	——括号

与集合的闭包运算有区别  
这里  $a^*$  表示由任意个  $a$  组成的串，  
而  $\{a, b\}^* = \{e, a, b, aa, ab, ba, bb, \dots\}$

**运算符的优先级：**

先 $.$ ，后 $*$ ，最后 $|$   
在正则表达式中可以省略。

**正则表达式相等  $\Leftrightarrow$  这两个正则表达式表示的语言相等**

如： $b\{ab\} = \{ba\}b$   
 $\{a\}b = \{a\}\{b\} = (a^*b^*)^*$

北京航空航天大学计算机学院

**例：设  $\Sigma = \{a, b\}$ ，下面是定义在  $\Sigma$  上的正则表达式和正则集合**

正则表达式	正则集合
$ba^*$	以 $b$ 为首，后跟 0 个和多个 $a$ 的符号串
$a(a b)^*$	$\Sigma$ 上以 $a$ 为首的所有符号串
$(a b)^*(aa bb)(a b)^*$	$\Sigma$ 上含有 $aa$ 或 $bb$ 的所有符号串

北京航空航天大学计算机学院

**正则表达式的性质：**

设  $e_1, e_2$  和  $e_3$  均是某字母表上的正则表达式，则有：

单位正则表达式： $\varepsilon$        $\varepsilon e = e\varepsilon = e$

交换律： $e_1 | e_2 = e_2 | e_1$

结合律： $e_1 | (e_2 | e_3) = (e_1 | e_2) | e_3$   
 $e_1 (e_2 e_3) = (e_1 e_2) e_3$

分配律： $e_1 (e_2 | e_3) = e_1 e_2 | e_1 e_3$   
 $(e_1 | e_2) e_3 = e_1 e_3 | e_2 e_3$

此外： $r^* = (r | \varepsilon)^*$        $r^{**} = r^*$   
 $(r | s)^* = (r^* s^*)^*$

北京航空航天大学计算机学院

**正则表达式与 3 型文法等价**

例如：

正则表达式： $ba^*$	$a(a b)^*$
3 型文法： $Z ::= Za b$	$Z ::= Za Zb a$

**例：**

3 型文法	正则表达式
$S ::= aS aB$	$aS aba^*a \rightarrow a^*aba^*a$
$B ::= bC$	$ba^*a$
$C ::= aC a$	$a^*a$

北京航空航天大学计算机学院

**11.2.1 确定的有穷自动机 (DFA) — 状态图的形式化 (Deterministic Finite Automata)**

一个确定的有穷自动机 (DFA)  $M$  是一个五元式：

$$M = (S, \Sigma, \delta, s_0, Z)$$

其中：

- $S$  — 有穷状态集
- $\Sigma$  — 输入字母表
- $\delta$  — 映射函数 (也称状态转换函数)  
 $S \times \Sigma \rightarrow S$   
 $\delta(s, a) = s', s, s' \in S, a \in \Sigma$
- $s_0$  — 初始状态       $s_0 \in S$
- $Z$  — 终止状态集       $Z \subseteq S$

北京航空航天大学计算机学院

$M = (S, \Sigma, \delta, s_0, Z)$

例如： $M = (\{0, 1, 2, 3\}, \{a, b\}, \delta, 0, \{3\})$

$\delta(0, a) = 1$	$\delta(0, b) = 2$
$\delta(1, a) = 3$	$\delta(1, b) = 2$
$\delta(2, a) = 1$	$\delta(2, b) = 3$
$\delta(3, a) = 3$	$\delta(3, b) = 3$

**状态转换函数  $\delta$  可用一矩阵来表示：**

输入 字符	a	b
状态		
0	1	2
1	3	2
2	1	3
3	3	3

所谓确定的状态机，其确定性表现在状态转换函数是单值函数！

北京航空航天大学计算机学院

Compiler

DFA 也可以用一状态转换图表示:

输入 字符	a	b
状态		
0	1	2
1	3	2
2	1	3
3	3	3

DFA 的状态图表示:

北京航空航天大学计算机学院

Compiler

DFA M 所接受的符号串:

令  $\alpha = a_1 a_2 \dots a_n$ ,  $\alpha \in \Sigma^*$ , 若  $\delta(\delta(s_0, a_1), a_2) \dots a_n$ ,  $\delta(s_{n-1}, a_n) = s_n$ , 且  $s_n \in Z$ , 则可以写成  $\delta(s_0, \alpha) = s_n$ , 我们称  $\alpha$  可为 M 所接受。

$\delta(s_0, a_1) = s_1$   
 $\delta(s_1, a_2) = s_2$   
 $\dots$   
 $\delta(s_{n-2}, a_{n-1}) = s_{n-1}$   
 $\delta(s_{n-1}, a_n) = s_n$

换言之: 若存在一条初始状态到某一终止状态的路径, 且这条路径上能有弧的标记符号连接成符号串  $\alpha$ , 则称  $\alpha$  为 DFA M (接受) 识别。

DFA M 所接受的语言为:  $L(M) = \{ \alpha \mid \delta(s_0, \alpha) = s_n, s_n \in Z \}$

北京航空航天大学计算机学院

Compiler

### 11.2.2 不确定的有穷自动机(NFA) (Nondeterministic Finite Automata)

若  $\delta$  是一个多值函数, 且输入可允许为  $\epsilon$ , 则有穷自动机是不确定的, 即在某个状态下, 对于某个输入字符存在多个后继状态。

从同一状态出发, 有以同一字符标记的多条边, 或者有以  $\epsilon$  标记的特殊边的自动机。

北京航空航天大学计算机学院

Compiler

在始态, 输入字母 a 时, 自动机既可以

如果向右, 则可接受由 a 组成长度为偶数的字符串。

如果向左, 则可接受由 a 组成长度为 3 的倍数的字符串。

因此, 该 NFA 所能接受的语言是所有由 a 组成, 长度为 2 和 3 的倍数的字符串的集合。在第一次状态转换中, 自动机需要选择要走的路径。只要有任何路径可匹配输入字符串, 该串就必须被接受, 因此 NFA 必须正确“猜测”所需的路径。

北京航空航天大学计算机学院

Compiler

经过以  $\epsilon$  标记的边无须任何字符输入。这里是接受同一语言的另一个 NFA:

图示自动机需要选择沿哪一条标记有  $\epsilon$  的边前进。如果一个状态同时引出以  $\epsilon$  标记的边和以其它字符标记的边, 则自动机可以选择处理一个输入字符并沿其对应的边前进, 或者仅沿  $\epsilon$  边前进。

北京航空航天大学计算机学院

Compiler

NFA 的形式定义为:

一个非确定的有穷自动机 NFA M' 是一个五元式:

NFA M' = (S,  $\Sigma \cup \{ \epsilon \}$ ,  $\delta$ ,  $s_0$ , Z)

其中

- S — 有穷状态集
- $\Sigma \cup \{ \epsilon \}$  — 输入符号加上  $\epsilon$ , 即自动机的每个结点所射出的弧可以是  $\Sigma$  中的一个字符或是  $\epsilon$
- $s_0$  — 初态  $s_0 \in S$
- Z — 终态集  $Z \subseteq S$
- $\delta$  — 转换函数  $S \times \Sigma \cup \{ \epsilon \} \rightarrow 2^S$  ( $2^S \rightarrow S$  的幂集 — S 的子集构成的集合)

北京航空航天大学计算机学院

Compiler

NFA  $M'$  所接受的语言为:

$$L(M') = \{ a \mid \delta(S_0, a) = S' \quad S' \cap Z \neq \emptyset \}$$

例: NFA  $M' = (\{1, 2, 3, 4\}, \{a, b, c\} \cup \{\epsilon\}, \delta, 1, \{4\})$

状态 \ 符号	$\epsilon$	a	b	c
1	{4}	{2, 3}	$\emptyset$	$\emptyset$
2	$\emptyset$	{2}	{4}	$\emptyset$
3	$\emptyset$	$\emptyset$	$\emptyset$	{3, 4}
4	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$

北京航空航天大学计算机学院

Compiler

符号 \ 状态	$\epsilon$	a	b	c
1	{4}	{2, 3}	$\emptyset$	$\emptyset$
2	$\emptyset$	{2}	{4}	$\emptyset$
3	$\emptyset$	$\emptyset$	$\emptyset$	{3, 4}
4	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$

上例题相应的状态图为:

$M'$  所接受的语言 (用正则表达式)  $R = aa^*b|ac^*c| \epsilon$

北京航空航天大学计算机学院

Compiler

复习:

1. 正则表达式与有穷自动机, 给出了两者的定义。  
用3型文法所定义的语言都可以用正则表达式描述。  
用正则表达式描述单词是为了自动生成词法分析程序。  
有一个正则表达式则对应一个正则集合。  
若  $V$  是正则集合, iff  $V = L(M)$   
即一个正则表达式对应一个 DFA  $M$
2. NFA 的定义,  $\delta$  非单值函数, 且有  $\epsilon$  弧, 表现为非确定性  
如:  $\delta(s, a) = \{s_1, s_2\}$   
 $\delta(s, a) = \{s_1, s_3\}$

北京航空航天大学计算机学院

Compiler

### 11.2.3 NFA 的确定化

正如我们所学到的, 用计算机程序实现 DFA 是很容易的。但在多数计算机硬件并不能正确猜测测路径的情况下, NFA 的实现就有些困难了。

已证明: 不确定的有穷自动机与确定的有穷自动机从功能上来说是等价的, 也就是说能够从:

$$NFA \ M' \xrightarrow{\text{构造成一个}} DFA \ M$$

使得  $L(M) = L(M')$

北京航空航天大学计算机学院

Compiler

为了使得 NFA 确定化, 首先给出两个定义:

**定义1、集合  $I$  的  $\epsilon$ -闭包:**

令  $I$  是一个状态集的子集, 定义  $\epsilon\text{-closure}(I)$  为:

- 1) 若  $s \in I$ , 则  $s \in \epsilon\text{-closure}(I)$ ;
- 2) 若  $s \in I$ , 则从  $s$  出发经过任意条  $\epsilon$  弧能够到达的任何状态都属于  $\epsilon\text{-closure}(I)$ 。

状态集  $\epsilon\text{-closure}(I)$  称为  $I$  的  $\epsilon$ -闭包。

可以通过一例子来说明状态子集的  $\epsilon$ -闭包的构造方法

北京航空航天大学计算机学院

Compiler

例:

如图所示的状态图:

令  $I = \{1\}$ ,  
求  $\epsilon\text{-closure}(I) = ?$

根据定义:  
 $\epsilon\text{-closure}(I) = \{1, 3\}$

北京航空航天大学计算机学院

**定义2:** 令  $I$  是 NFA  $M'$  的状态集的一个子集,  $a \in \Sigma$   
 定义:  $I_a = \epsilon\text{-closure}(J)$   
 其中  $J = \bigcup_{s \in I} \delta(s, a)$

—  $J$  是从状态子集  $I$  中的每个状态出发, 经过标记为  $a$  的弧而达到的状态集合。

—  $I_a$  是状态子集, 其元素为  $J$  中的状态, 加上从  $J$  中每一个状态出发通过  $\epsilon$  弧到达的状态。

同样可以通过一例子来说明上述定义, 仍采用前面给定的状态图为例

北京航空航天大学计算机学院

**例:** 令  $I = \{1\}$   
 $I_a = \epsilon\text{-closure}(J)$   
 $= \epsilon\text{-closure}(\delta(1, a))$   
 $= \epsilon\text{-closure}(\{2, 4\})$   
 $= \{2, 4, 6\}$

根据定义1, 2, 可以将上述的  $M'$  确定化 (即可构造出状态转换矩阵)

北京航空航天大学计算机学院

**例:** 有 NFA  $M'$

$I = \epsilon\text{-closure}(\{1\}) = \{1, 4\}$   
 $I_a = \epsilon\text{-closure}(\delta(1, a) \cup \delta(4, a))$   
 $= \epsilon\text{-closure}(\{2, 3\} \cup \phi)$   
 $= \epsilon\text{-closure}(\{2, 3\})$   
 $= \{2, 3\}$   
 $I_b = \epsilon\text{-closure}(\delta(1, b) \cup \delta(4, b))$   
 $= \epsilon\text{-closure}(\phi)$   
 $= \phi$   
 $I_c = \epsilon\text{-closure}(\delta(1, c) \cup \delta(4, c))$   
 $= \phi$   
 $I = \{2, 3\}, I_a = \{2\}, I_b = \{4\}, I_c = \{3, 4\} \dots$

北京航空航天大学计算机学院

$I$	$I_a$	$I_b$	$I_c$
$\{1, 4\}$	$\{2, 3\}$	$\phi$	$\phi$
$\{2, 3\}$	$\{2\}$	$\{4\}$	$\{3, 4\}$
$\{2\}$	$\{2\}$	$\{4\}$	$\phi$
$\{4\}$	$\phi$	$\phi$	$\phi$
$\{3, 4\}$	$\phi$	$\phi$	$\{3, 4\}$

北京航空航天大学计算机学院

将求得的状态转换矩阵重新编号

DFA  $M$  状态转换矩阵:

符号	a	b	c
0	1	—	—
1	2	3	4
2	2	3	—
3	—	—	—
4	—	—	4

北京航空航天大学计算机学院

DFA  $M$  的状态图:

注意: 原初始状态的  $\epsilon\text{-closure}$  为 DFA  $M$  的初态  
 包含原终止状态 4 的状态子集为 DFA  $M$  的终态。

北京航空航天大学计算机学院

**复习:**

1. 正则表达式与有穷自动机, 给出了两者的定义。  
用3型文法所定义的语言都可以用正则表达式描述,  
用正则表达式描述单词是为了自动生成词法分析程序。  
有一个正则表达式则对应一个正则集合。
2. NFA  $M'$  的定义、确定化  $\rightarrow$  对任何一个NFA  $M'$ , 都可以构造出一个DFA  $M$ , 使得  $L(M) = L(M')$

构造出来的DFA  $M$  唯一吗?

北京航空航天大学计算机学院

**11.2.4 DFA的简化(最小化)**

“对于任一个DFA, 存在一个唯一的状态最少的等价的 DFA”

一个有穷自动机是化简的  $\Leftrightarrow$  它没有多余状态并且它的状态中没有两个是互相等价的。

一个有穷自动机可以通过消除多余状态和合并等价状态而转换成一个最小的与之等价的有穷自动机

北京航空航天大学计算机学院

**定义:**

(1) 有穷自动机的多余状态: 从该自动机的开始状态出发, 任何输入串也不能到达那个状态

例:

	0	1
$S_0$	$S_1$	$S_5$
$S_1$	$S_2$	$S_7$
$S_2$	$S_2$	$S_5$
$S_3$	$S_5$	$S_7$
$S_4$	$S_5$	$S_6$
$S_5$	$S_3$	$S_1$
$S_6$	$S_6$	$S_0$
$S_7$	$S_0$	$S_1$
$S_8$	$S_3$	$S_6$

画状态图可以看出  $S_4, S_6, S_8$  为不可达状态应该消除

	0	1
$S_0$	$S_1$	$S_5$
$S_1$	$S_2$	$S_7$
$S_2$	$S_2$	$S_5$
$S_3$	$S_5$	$S_7$
$S_5$	$S_3$	$S_1$
$S_7$	$S_0$	$S_1$

北京航空航天大学计算机学院

(2) 等价状态  $\Leftrightarrow$  状态  $s$  和  $t$  的等价条件是:

- 1) 一致性条件: 状态  $s$  和  $t$  必须同时为可接受状态或不接受状态。
- 2) 蔓延性条件: 对于所有输入符号  $c$ , 状态  $s$  和  $t$  必须转换到等价的状态里。

对于所有输入符号  $c$ ,  $L(s) = L(t)$ , 即状态  $s$ 、 $t$  对于  $c$  具有相同的后继, 则称  $s$ 、 $t$  是等价的。  
(任何有后继的状态和任何无后继的状态一定不等价)

有穷自动机的状态  $s$  和  $t$  不等价, 称这两个状态是可区分的。

北京航空航天大学计算机学院

**“分割法”:** 把一个DFA(不含多余状态)的状态分割成一些不相关的子集, 使得任何不同的两个子集状态都是可区分的, 而同一个子集中的任何状态都是等价的。

例: 最小化

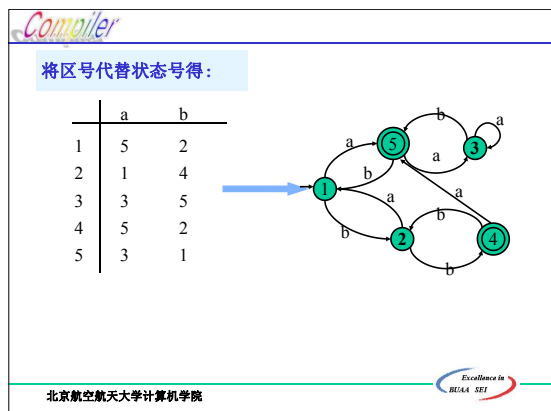
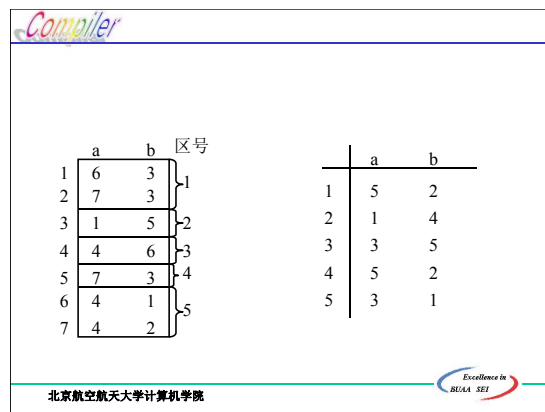
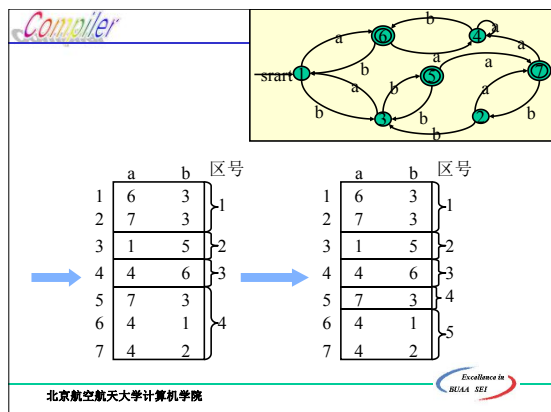
北京航空航天大学计算机学院

解: (一) 区分终态与非终态

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

北京航空航天大学计算机学院





Compiler

**复习：**

1. 正则表达式与有穷自动机，给出了两者的定义。  
用3型文法所定义的语言都可以用正则表达式描述，  
用正则表达式描述单词是为了自动生成词法分析程序。  
有一个正则表达式则对应一个正则集合。
2. NFA  $M'$  的定义、确定化  $\rightarrow$  对任何一个NFA  $M'$ ，都可以构造出一个DFA  $M$ ，使得  $L(M) = L(M')$

**下面：**  
**正则表达式与DFA的等价性**  
我们证明了对任何一个正则表达式，都可以构造出等价的NFA。

北京航空航天大学计算机学院

Compiler

### 11.2.5 正则表达式与DFA的等价性

**定理：**在  $\Sigma$  上的一个子集  $V (V \subseteq \Sigma^*)$  是正则集合，当且仅当存在一个DFA  $M$ ，使得  $V = L(M)$

$V$  是正则集合，  
 $R$  是与其相对应的正则表达式  $\Leftrightarrow$  DFA  $M$   
 $V = L(R) \quad L(M) = L(R)$

所以 正则表达式  $R \Rightarrow$  NFA  $M' \Rightarrow$  DFA  $M$   
 $L(R) = L(M') = L(M)$

**证明：**根据定义。

北京航空航天大学计算机学院

Compiler

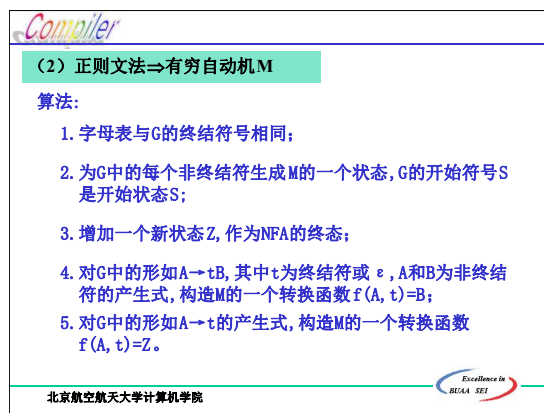
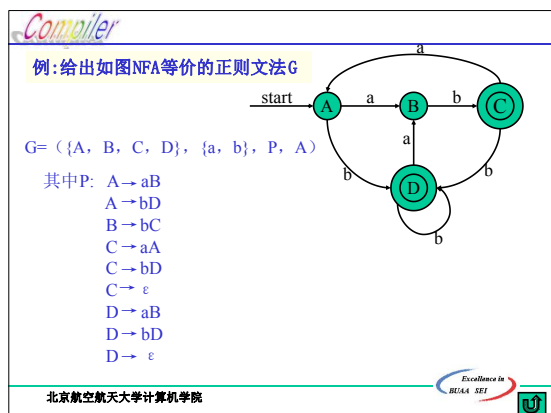
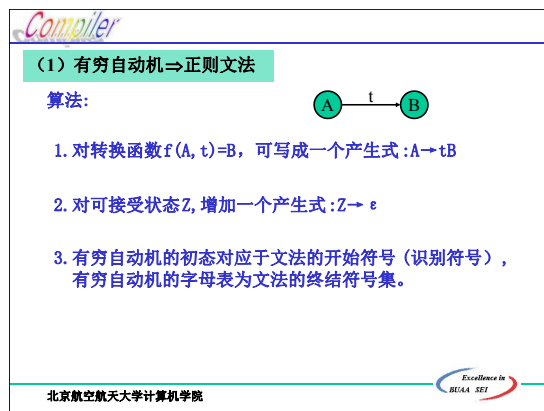
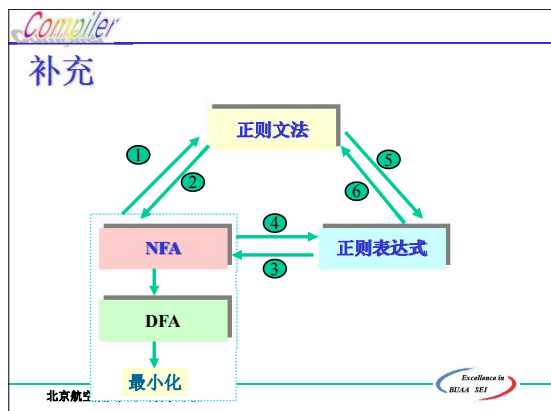
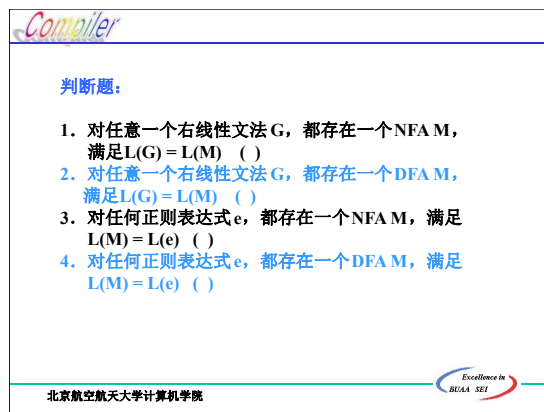
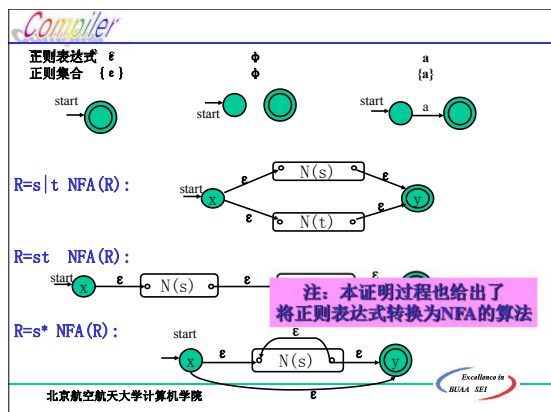
### 复习

**正则表达式和正则集合的递归定义**

有字母表  $\Sigma$ ，定义在  $\Sigma$  上的正则表达式和正则集合递归定义如下：

1.  $\epsilon$  和  $\phi$  都是  $\Sigma$  上的正则表达式，其正则集合分别为： $\{\epsilon\}$  和  $\phi$ ；
2. 任何  $a \in \Sigma$ ， $a$  是  $\Sigma$  上的正则表达式，其正则集合为： $\{a\}$ ；
3. 假定  $U$  和  $V$  是  $\Sigma$  上的正则表达式，其正则集合分别记为  $L(U)$  和  $L(V)$ ，那么  $U|V$ ， $U \cdot V$  和  $U^*$  也都是  $\Sigma$  上的正则表达式，其正则集合分别为  $L(U) \cup L(V)$ 、 $L(U) \cdot L(V)$  和  $L(U)^*$ ；
4. 任何  $\Sigma$  上的正则表达式和正则集合均由 1、2 和 3 产生。

北京航空航天大学计算机学院



Compiler

例: 求与文法  $G[S]$  等价的 NFA

$G[S]:$

- $S \rightarrow aA \mid bB \mid \epsilon$
- $A \rightarrow aB \mid bA$
- $B \rightarrow aS \mid bA \mid \epsilon$

求得:

北京航空航天大学计算机学院

Compiler

左线性正规文法和右线性正规文法的等价

左线性正规文法例

$Z \rightarrow U0 \mid V1$

$U \rightarrow Z1 \mid 1$

$V \rightarrow Z0 \mid 0$

$L(G[Z]) = \{ B^n \mid n > 0 \}$

其中  $B = \{01, 10\}$

$R = (01|10)(01|10)^*$

北京航空航天大学计算机学院

Compiler

右线性正规文法例

$S \rightarrow 0U \mid 1V$

$U \rightarrow 1S \mid 1$

$V \rightarrow 0S \mid 0$

$L(G[S]) = \{ B^n \mid n > 0 \}$

其中  $B = \{01, 10\}$

$R = (01|10)(01|10)^*$

北京航空航天大学计算机学院

Compiler

(3) 正则式  $\Rightarrow$  有穷自动机

语法制导方法

1.(a) 对于正则式  $\phi$ , 所构造 NFA:

(b) 对于正则式  $\epsilon$ , 所构造 NFA:

(c) 对于正则式  $a, a \in \Sigma$ , 则 NFA:

北京航空航天大学计算机学院

Compiler

2. 若  $s, t$  为  $\Sigma$  上的正则式, 相应的 NFA 分别为  $N(s)$  和  $N(t)$ :

(a) 对于正则式  $R = s|t$ , NFA (R)

(b) 对正则式  $R = st$ , NFA (R)

或:

北京航空航天大学计算机学院

Compiler

(c) 对于正则式  $R = s^*$ , NFA (R)

(d) 对  $R = (s)$ , 与  $R = s$  的 NFA 一样.

北京航空航天大学计算机学院

Compiler

例: 为  $R=(a|b)^*abb$  构造NFA  $N$ , 使得  $L(N)=L(R)$

从左到右分解  $R$ , 令  $r_1=a$ , 第一个  $a$ , 则有  $\Rightarrow$

令  $r_2=b$ , 则有  $\Rightarrow$

令  $r_3=r_1|r_2$ , 则有  $\Rightarrow$

北京航空航天大学计算机学院

Compiler

$R=(a|b)^*abb$

令  $r_4=r_3^*$  则有:

令  $r_5=a$ ,  
令  $r_6=b$ ,  
令  $r_7=b$ ,  
令  $r_8=r_5r_6$ ,  
令  $r_9=r_8r_7$  则有  $\Rightarrow$

北京航空航天大学计算机学院

Compiler

$R=(a|b)^*abb$

令  $r_{10}=r_4r_9$  则最终得到NFA  $N$ :

分解  $R$  的方法有很多种, 下面给出另一种分解方式和所构成的 NFA

北京航空航天大学计算机学院

Compiler

$R=(a|b)^*abb$

$\Rightarrow$  (a)

$\Rightarrow$  (b)

$\Rightarrow$  (c)

$\Rightarrow$  (d)

北京航空航天大学计算机学院

Compiler

(4) 有穷自动机  $\Rightarrow$  正则式  $R$

算法:

- 1) 在  $M$  上加两个结点  $x, y$ , 从  $x$  结点用  $\epsilon$  弧到  $M$  的所有初态, 从  $M$  的所有终态用  $\epsilon$  到  $y$  结点形成与  $M$  等价的  $M'$ .  $M'$  只有一个初态  $x$  和一个终态  $y$ .
- 2) 逐步消去  $M'$  中的所有结点, 直至剩下  $x$  和  $y$  结点, 在消结过程中, 逐步用正则式来记弧, 其消结规则如下:

1. 对于 代之为
2. 对于 代之为
3. 对于 代之为

北京航空航天大学计算机学院

Compiler

例:  $M$ :

解: (1) 加  $x, y$

北京航空航天大学计算机学院

Compiler

(2) 消除M中的所有结点

解: (1) 加xy

北京航空航天大学计算机学院

Compiler

(5) 正则文法 $\Rightarrow$ 正则式

利用以下转换规则,直至只剩下一个开始符号定义的产生式,并且产生式的右部不含非终结符。

规则	文法产生式	正则式
规则1	$A \rightarrow xB, B \rightarrow y$	$A = xy$
规则2	$A \rightarrow xA   y$	$A = x^*y$
规则3	$A \rightarrow x, A \rightarrow y$	$A = x   y$

北京航空航天大学计算机学院

Compiler

规则	文法产生式	正则式
规则1	$A \rightarrow xB, B \rightarrow y$	$A = xy$
规则2	$A \rightarrow xA   y$	$A = x^*y$
规则3	$A \rightarrow x, A \rightarrow y$	$A = x   y$

例:有文法G[s]

$S \rightarrow aA|a,$   
 $A \rightarrow aA|dA|a|d$

于是:  $S = aA|a$   
 $A = (aA|dA)|(a|d) \Rightarrow A = (a|d)A|(a|d)$

由规则二:  $A = (a|d)^*(a|d)$   
 代入:  $S = a(a|d)^*(a|d)|a$   
 于是:  $S = a((a|d)^*(a|d) | \epsilon)$

北京航空航天大学计算机学院

Compiler

(6) 正则式 $\Rightarrow$ 正则文法

算法:

- 1) 对任何正则式r,选择一个非终结符S作为识别符号,并产生产生式  $S \rightarrow r$
- 2) 若x,y是正则式,对形为 $A \rightarrow xy$ 的产生式,重写为  $A \rightarrow xB, B \rightarrow y$ ,其中B为新的非终结符,  $B \in V$ 。  
 同样: 对于  $A \rightarrow x^*y \Rightarrow A \rightarrow xA, A \rightarrow y$   
 $A \rightarrow x|y \Rightarrow A \rightarrow x, A \rightarrow y$

例:将 $R = a(a|d)^*$ 转换成相应的正则文法

解: 1)  $S \rightarrow a(a|d)^*$   
 2)  $S \rightarrow aA, A \rightarrow (a|d)^*$   
 3)  $S \rightarrow aA, A \rightarrow (a|d)A, A \rightarrow \epsilon$   
 4)  $S \rightarrow aA, A \rightarrow aA|dA, A \rightarrow \epsilon$

北京航空航天大学计算机学院

Compiler

### 11.3 词法分析程序的自动生成器—LEX (LEXICAL)

**LEX的原理:**

正则表达式与DFA的等价性  
 根据给定的正则表达式自动生成相应的词法分析程序。

**LEX的功能:**

LEX源程序  $\rightarrow$  LEX  $\rightarrow$  词法分析程序L

S.P.(字符串)  $\rightarrow$  L  $\rightarrow$  S.P.单词(字符串)

北京航空航天大学计算机学院

Compiler

#### 11.3.1 LEX源程序

一个LEX源程序主要由三个部分组成:

1. 辅助定义式
2. 识别规则
3. 用户子程序

各部分之间用%%隔开

北京航空航天大学计算机学院

**辅助定义式**是如下形式的 LEX 语句:

$$\begin{aligned} D_1 &\longrightarrow R_1 \\ D_2 &\longrightarrow R_2 \\ &\vdots \\ D_n &\longrightarrow R_n \end{aligned}$$

其中:

$R_1, R_2, \dots, R_n$  为正则表达式。  
 $D_1, D_2, \dots, D_n$  为正则表达式名字, 称简名。

限定: 在  $R_i$  中只能出现字母表  $\Sigma$  中的字符, 以及前面已定义的正则表达式名字, 我们用这种辅助定义式 (相当于规则) 来定义程序语言的单词符号。

北京航空航天大学计算机学院

例: 标识符:  $\text{letter} \rightarrow A|B|\dots|Z$   
 $\text{digit} \rightarrow 0|1|\dots|9$   
 $\text{iden} \rightarrow \text{letter}(\text{letter}|\text{digit})^*$

带符号整数:  $\text{integer} \rightarrow \text{digit}(\text{digit})^*$   
 $\text{sign} \rightarrow +|-|\epsilon$   
 $\text{sign\_integer} \rightarrow \text{sign integer}$

北京航空航天大学计算机学院

**识别规则:** 是一串如下形式的 LEX 语句:

$$\begin{aligned} P_1 &\{A_1\} \\ P_2 &\{A_2\} \\ &\vdots \\ P_m &\{A_m\} \end{aligned}$$

$P_i$ : 定义在  $\Sigma \cup \{D_1, D_2, \dots, D_n\}$  上的正则表达式, 也称词形。  
 $\{A_i\}$ :  $A_i$  为语句序列, 它指出, 在识别出词形为  $P_i$  的单词以后, 词法分析器所应作的动作。  
 其基本动作是返回单词的类别编码和单词值。

北京航空航天大学计算机学院

下面是识别某语言单词符号的 LEX 源程序:

例: LEX 源程序

RETURN 是 LEX 过程, 该过程将单词传给语法分析程序  
 $\text{RETURN}(C, \text{LEXVAL})$   
 其中  $C$  为单词类别编码  
 LEXVAL:  
 标识符: TOKEN (字符数组)  
 整数: DTB (数值转换函数, 将 TOKEN 中的数字串转换二进制值)  
 其他单词: 无定义

AUXILIARY DEF  
 $\text{letter} \rightarrow A|B|\dots|Z$   
 $\text{digit} \rightarrow 0|1|\dots|9$   
 %%  
 RECOGNITION RULES  
 1.BEGIN {RETURN(1,-)}  
 2.END {RETURN(2,-)}  
 3.FOR {RETURN(3,-)}

北京航空航天大学计算机学院

4.DO	{RETURN(4,-)}
5.IF	{RETURN(5,-)}
6.THEN	{RETURN(6,-)}
7.ELSE	{RETURN(7,-)}
8.letter(letter digit)*	{RETURN(8,TOKEN)}
9.digit(digit)*	{RETURN(9,DTB)}
10.:	{RETURN(10,-)}
11.+	{RETURN(11,-)}
12."*"	{RETURN(12,-)}

北京航空航天大学计算机学院

13.,	{RETURN(13,-)}
14."("	{RETURN(14,-)}
15.")"	{RETURN(15,-)}
16.:="	{RETURN(16,-)}
17.=	{RETURN(17,-)}

北京航空航天大学计算机学院

**11.3.2 LEX的实现**

LEX的功能是根据LEX源程序构造一个词法分析程序，该词法分析器实质上是一个有穷自动机。

LEX生成的词法分析程序由两部分组成：

词法分析程序	状态转换矩阵(DFA)
	控制执行程序

∴ LEX的功能是根据LEX源程序生成状态转换矩阵和控制程序

北京航空航天大学计算机学院

**LEX的工作过程：**

- 扫描每条识别规则  $P_i$ ，构造
- 将各条规则的有穷自动机  $M_i$
- 确定化  $NFA \Rightarrow DFA$
- 生成该DFA的状态转换矩阵和控制执行程序

北京航空航天大学计算机学院

如: begin, :=

**LEX二义性问题的两条原则：**

- 最长匹配原则**  
在识别单词过程中，有一字符串  $\overbrace{x \ x \ x}^{P_i}$  根据最长匹配原则，应识别为这是一个符合  $P_i$  规则的单词，而不是  $P_j$  和  $P_i$  规则的单词。
- 优先匹配原则**  
如有一字符串，有两条规则可以同时匹配时，那么用规则序列中位于前面的规则相匹配，所以排列在最前面的规则优先权最高。

北京航空航天大学计算机学院

例：字符串  $\overbrace{b \ e \ g \ i \ n}^{P_i}$

根据原则，应该识别为关键字 begin，所以在写LEX源程序时应注意规则的排列顺序。

此外，**优先匹配原则**是在符合**最长匹配**的前提下执行的。

可以通过一个例子来说明这些问题：

北京航空航天大学计算机学院

例：LEX源程序

```

a { }
abb { }
a*bb* { }

```

**一.读LEX源程序，分别生成NFA，用状态图表示为：**

**二.合并成一个NFA：**

北京航空航天大学计算机学院

**三.确定化 给出状**

状态	a	b	到达终态所识别的单词
初态 {0,1,3,7}	{2,4,7}	{8}	
终态 {2,4,7}	{7}	{5,8}	a
终态 {8}	φ	{8}	a*bb*
终态 {7}	{7}	{8}	
终态 {5,8}	φ	{6,8}	a*bb
终态 {6,8}	φ	{8}	abb

在此DFA中 初态为{0,1,3,7}

终态为{2,4,7},{8},{5,8},{6,8}

北京航空航天大学计算机学院

Compiler

词法分析程序的分析过程

读入字符 | 进入状态

开始

a

b

a

{0,1,3,7}

{2,4,7}

{5,8}

无后继状态(退掉输入字符a)

令输入字符串为aba...

(1) 吃进字符ab

(2) 按反序检查状态子集

检查前一次状态是否含有原NFA的终止状态

即检查{5,8},含有终态8, 因此断定所识别的单词ab是属于 $a^*bb^*$ 中的一个。

若状态子集中无NFA的终态, 则要从识别的单词再退掉一个字符(b), 然后再检查上一个状态子集。

若一旦吃进的字符都退完, 则识别失败, 调用出错程序, 一般是跳过一个字符然后重新分析。(应打印出错信息)

北京航空航天大学计算机学院

Excellence in BEAA SEE

Compiler

三点说明:

1) 以上是LEX的构造原理, 虽然是原理性的, 但据此就不准将LEX构造出来。

2) 所构造出来的LEX是一个通用的工具, 用它可以生成各种语言的词法分析程序, 只需要根据不同的语言书写不同的LEX源文件就可以了。

3) LEX不但能自动生成词法分析器, 而且也可以产生多种模式识别器及文本编辑程序等。

北京航空航天大学计算机学院

Excellence in BEAA SEE

Compiler

第十一章作业:

P254-255 1, 2, 4, 5;

北京航空航天大学计算机学院

Excellence in BEAA SEE

16