

编译系统

第三章

词法分析



哈尔滨工业大学 陈鄞

## 第3讲(词法分析)要点

- ▶词法分析的任务: 识别单词
  - ▶单词的描述
  - > 单词的识别

### 1、单词的描述

- ▶RG (正则文法): G
  - $\succ G \rightarrow L(G)$



- ▶RE (正则表达式): r 更直观、更紧凑
  - $> r \rightarrow L(r)$
- ► RD (正则定义): d 给RE命名
  - ▶形式: d →r

### 2、单词的识别

▶理论基础: 有穷自动机 (FA)

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

**▶**S: 状态集合

 $\Sigma$ : 输入符号集合

 $\triangleright \delta$ : 转换函数,  $S \times \Sigma \rightarrow 2^S$ 

 $\succ s_0$ : 初始状态,  $s_0 \in S$ 

 $\triangleright F$ : 终止状态集合,  $F \subseteq S$ 

①FA的表示 
$$M = (S, \Sigma, \delta, s_{\theta}, F)$$

- > 转换图
  - > 节点

状态

>边(有向) 输入符号

户转换表

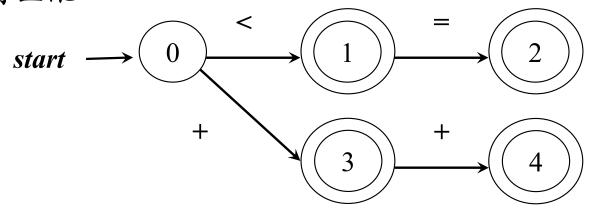
S	$a_1$	$a_2$	•••	$a_n$
$s_{I}$				
$s_2$				
• • •				
S <sub>m</sub>				

## ②FA定义(接收/识别)的语言

- $> M \rightarrow L(M)$ 
  - ▶给定输入串x,如果存在一个对应于串x的从初始状态到某个 终止状态的转换序列,则称串x被该FA接收
  - $\triangleright$ 由一个有穷自动机M接收的所有串构成的集合称为是该FA定义(或接收)的语言,记为L(M)

## ②FA定义(接收/识别)的语言

- $> M \rightarrow L(M)$
- > 最长子串匹配原则
  - ▶ 当输入串的多个前缀与一个或多个模式匹配时,总是选择最长的前缀进行匹配



▶ 在到达某个终态之后,只要输入带上还有符号,DFA就继续前进,以 便寻找尽可能长的匹配

③FA的分类 
$$M = (S, \Sigma, \delta, s_{\theta}, F)$$

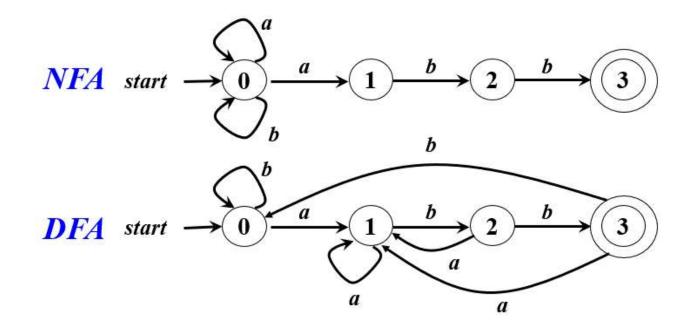
▶DFA (确定的有穷自动机) 更适用

 $\triangleright \delta: S \times \Sigma \rightarrow S$ 



>NFA(非确定的有穷自动机)更直观

 $\triangleright \delta$ :  $S \times \Sigma \rightarrow 2^S$ 



③FA的分类  $M = (S, \Sigma, \delta, s_{\theta}, F)$ 

►DFA (确定的有穷自动机)

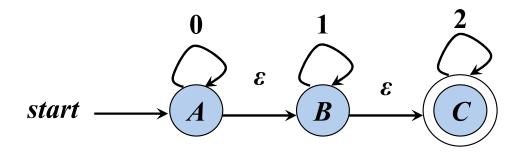
$$\triangleright \delta: S \times \Sigma \rightarrow S$$



>NFA(非确定的有穷自动机)

 $\triangleright \delta$ :  $S \times \Sigma \rightarrow 2^S$  不带 " $\varepsilon$ -边"

 $\triangleright$  δ:  $S \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^S$  # "ε-边"



③FA的分类 
$$M = (S, \Sigma, \delta, s_0, F)$$

start

►DFA (确定的有穷自动机)

$$\triangleright \delta: S \times \Sigma \rightarrow S$$

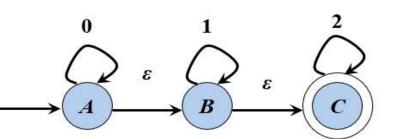


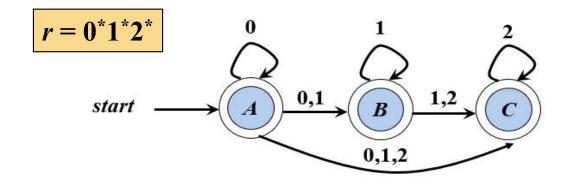
>NFA(非确定的有穷自动机)

$$\triangleright \delta$$
:  $S \times \Sigma \rightarrow 2^S$ 

不带 
$$^{\prime\prime}\varepsilon$$
-边"

$$\triangleright$$
  $\delta$ :  $S \times \Sigma \to 2^S$  不帯 " $\varepsilon$ -边"   
  $\triangleright$   $\delta$ :  $S \times (\Sigma \cup \{\varepsilon\}) \to 2^S$  帯 " $\varepsilon$ -边" 更直观





# ④DFA的实现 $M = (S, \Sigma, \delta, s_0, F)$

▶函数体由一个循环构成

```
s = s_0;

c = nextChar ();

while (c! = eof) {

s = move (s, c);

c = nextChar ();

}

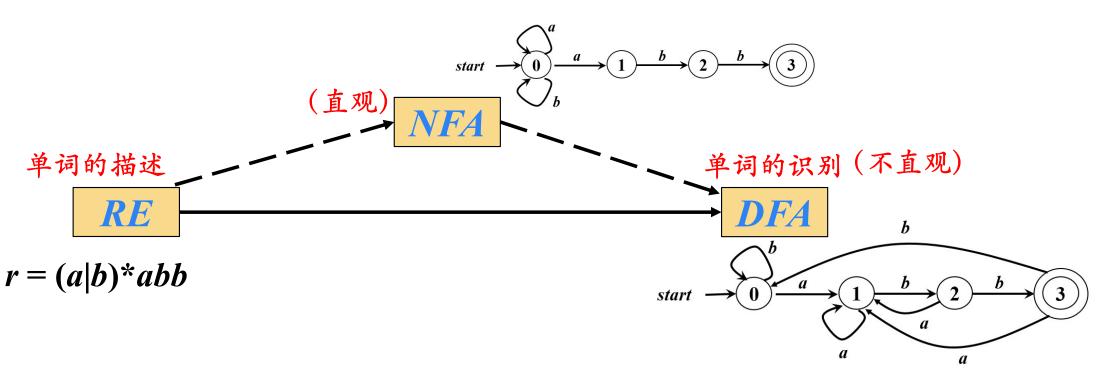
if (s \triangle F +) return "yes";

else return "no";
```

$S$ $\Sigma$	$a_1$	$a_1$	•••	$a_n$
$s_1$				
$s_2$				
• • •				
$S_{m}$				

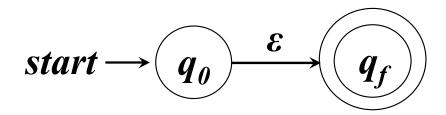
# ⑤FA与RE的等价性

- $\triangleright RE \Leftrightarrow FA$ 
  - $r \Leftrightarrow M$
- ▶词法分析器 (scanner) 的实现思路



## ⑥从RE 到NFA

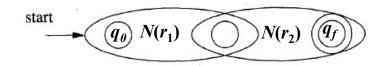
➤ E对应的NFA



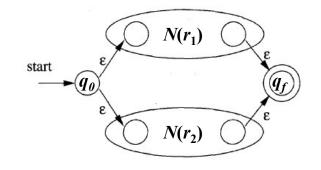
 $\triangleright$ 字母表 $\Sigma$ 中符号 $\alpha$ 对应的NFA

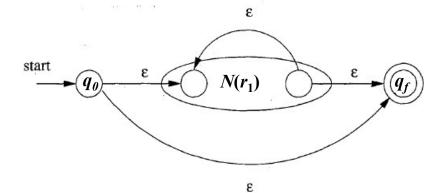
$$start \rightarrow \boxed{q_0} \xrightarrow{a} \boxed{q_f}$$

 $r = r_1 r_2$ 对应的NFA

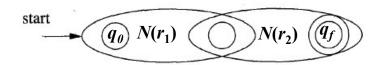


 $r = r_1 | r_2$ 对应的NFA

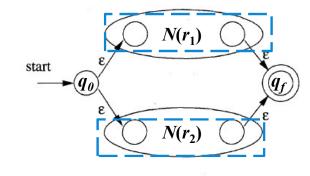


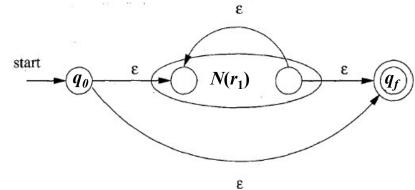


 $r = r_1 r_2$ 对应的NFA

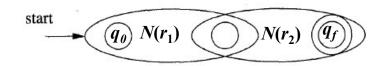


 $r = r_1 | r_2$ 对应的NFA





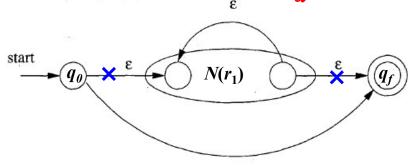
 $r = r_1 r_2$ 对应的NFA



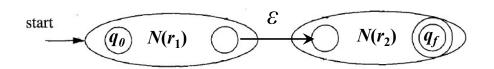
 $N(r_1)$ 

 $r = r_1 | r_2$ 对应的NFA

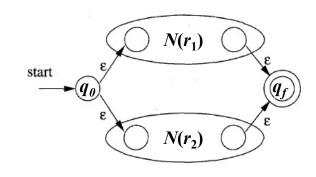
 $N(r_2)$ 

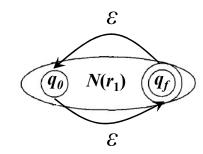


 $r = r_1 r_2$ 对应的NFA



 $r = r_1 | r_2$ 对应的NFA





## ⑦从NFA到DFA:子集构造法

