



中山大學

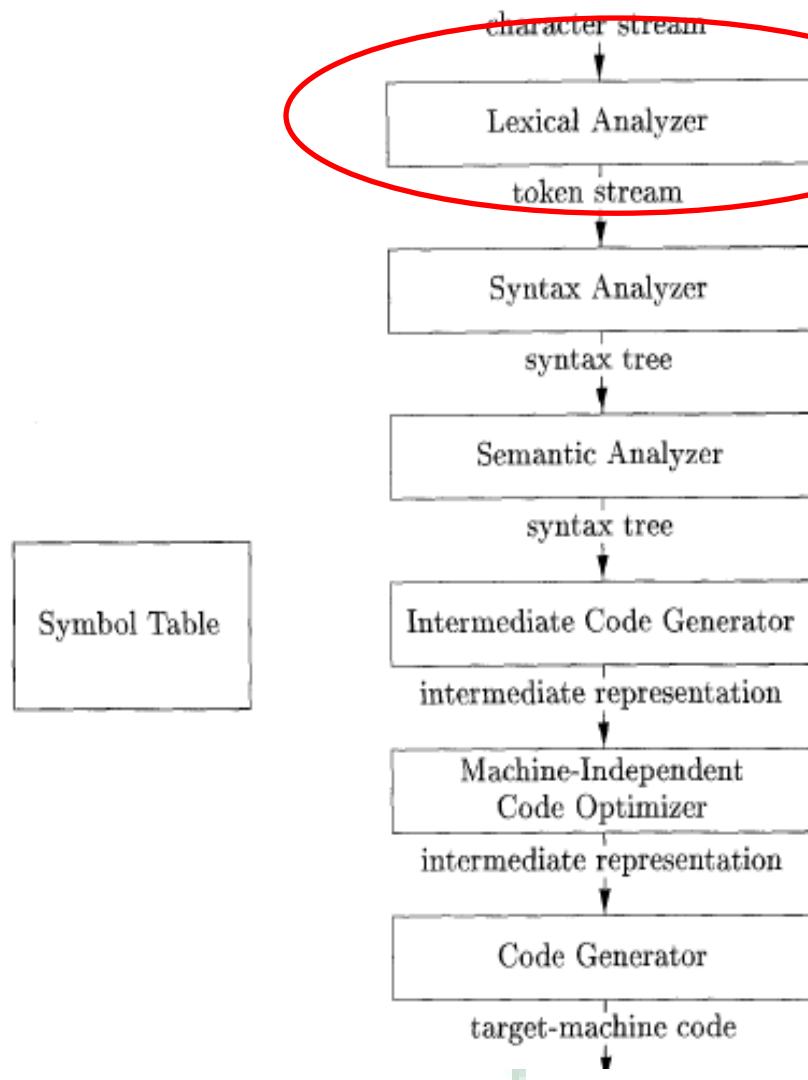
Principles of Compiler Construction

Lecture 2 Lexical Analysis (I)



中山大學

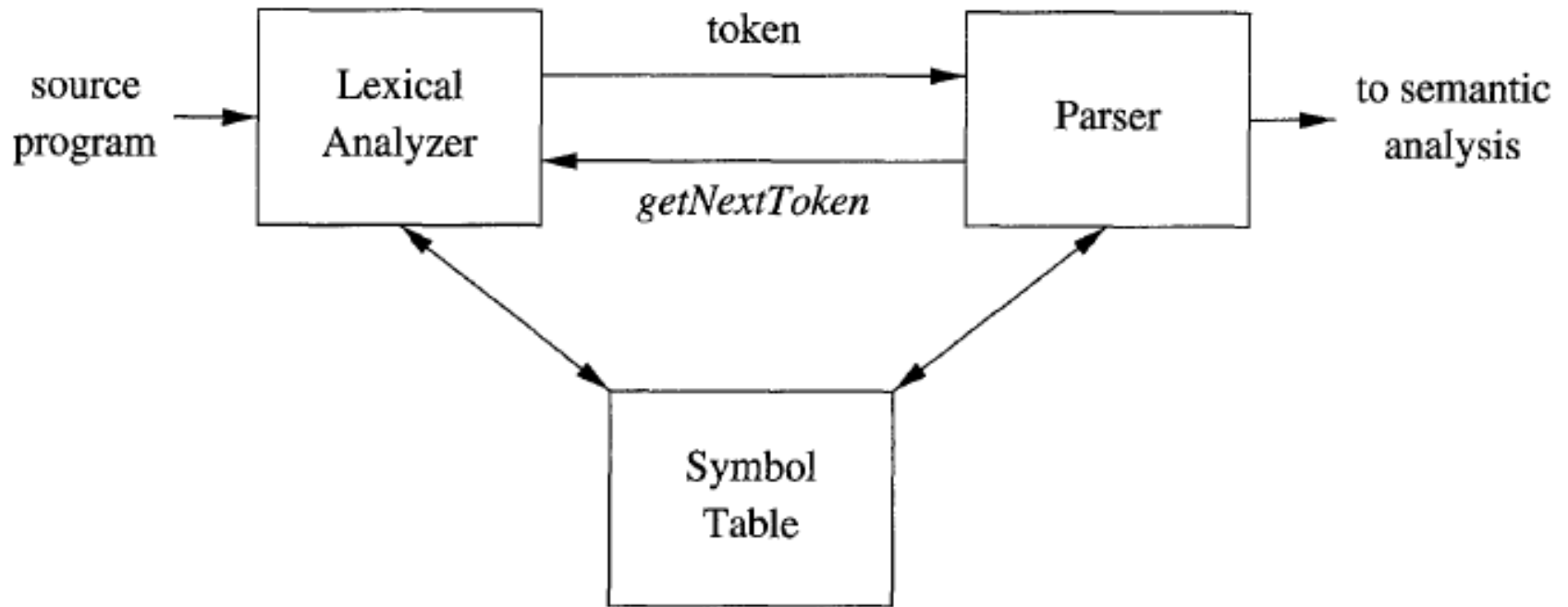
回顾：编译器的结构





中山大學

Lexical Analyzer vs. Syntax Analyzer (Parser)



中山大學

词法分析的任务

源程序在经词法分析之前是字符的序列

例如：

```
if (month == March)
```

```
    nday = 31;
```

是字符序列：\tif (month == March)\n\t\tnday = 31;



中山大學

词法分析的任务

词法分析器将其转化为token的序列

- token是指带有附加信息的字符串

例如<“month”, id>

单纯的字符串我们称为lexeme

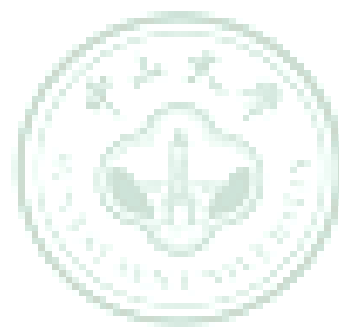
- 在自然语言中：名词，动词，形容词……
- 程序语言中：标识符，关键字，常量，运算符，分界符
- 类别属性是包括在token中的基本信息
- 语法分析器每次从词法分析器获得一个token，对不同类别token的处理截然不同



中山大學

关键问题

- 如何识别出一个token?
 - 如何识别出一个lexeme?
 - 如何知道lexeme的类别属性?



中山大學



中山大學

类别定义

- 关键字或保留字：预先指定的单词，如if, for, else等
- 标识符：
 - 字母开头的字母或数字序列
 - 非保留字
- 数字：
 - 正整数：首字符非0的数字序列
 - 负数：' - ' 和正整数的连接
 -



中山大學



中山大學

类别定义

- 自然语言定义方式的缺陷：
 - 繁琐
 - 不精确
 - 难以被计算机处理
- 怎么办？
 - 借助形式化的语言！



中山大學



中山大學

3.3.1 *Strings and Languages*

字母表: An alphabet is any finite set of symbols.

语言: 字母表 Σ 上的语言, 是由 Σ 中字符组成的字符串的集合

例如:

$$\Sigma = \{0, 1\}$$

则 $\{001, 100100\}$, $\{\}$, $\{1, 11, 111, 1111, \dots\}$ 都是定义在字母表 Σ 的语言

中山大學



语言 (*Language*)

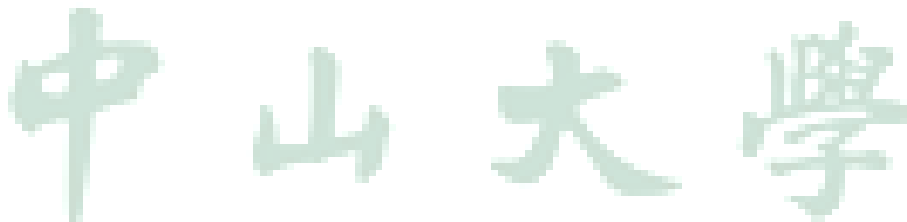
字母表：字符的集合

语言：字母表 Σ 上的语言，是由 Σ 中字符组成的字符串的集合

例如：

$$\Sigma = \{0, 1\}$$

则 $\{001, 100100\}$, $\{\}$, $\{1, 11, 111, 1111, \dots\}$ 都是定义在字母表 Σ 的语言



程序设计语言

程序：告诉计算机“做什么”和“怎么做”

程序设计：编写程序的过程

程序设计语言：编写程序使用的语言

语句：定义在一个字符集合上的按一定的语法规则连接的有意义的字符串

语言：由语句组成的全体

作用：交换信息的工具

自然语言与形式语言

- 自然语言:

人与人交换信息的工具。不精确，有二义性

汉语、英语

手提包

- 形式语言:

数学语言，精确语言，无二义性。

计算机语言为人与计算机交换信息的工具

计算机语言

文法、自动机、表达式



中山大學

语言的运算

语言是一种集合，所以集合的运算也适用于语言

| OPERATION | DEFINITION AND NOTATION |
|---|---|
| <i>Union of L and M</i> | $L \cup M = \{s \mid s \text{ is in } L \text{ or } s \text{ is in } M\}$ |
| <i>Concatenation of L and M</i> | $LM = \{st \mid s \text{ is in } L \text{ and } t \text{ is in } M\}$ |
| <i>Kleene closure of L</i> | $L^* = \bigcup_{i=0}^{\infty} L^i$ |
| <i>Positive closure of L</i> | $L^+ = \bigcup_{i=1}^{\infty} L^i$ |



中山大學



中山大學

语言运算的例子

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

1. $L \cup D$ is the set of letters and digits — strictly speaking the language with 62 strings of length one, each of which strings is either one letter or one digit.
2. LD is the set of 520 strings of length two, each consisting of one letter followed by one digit.
3. L^4 is the set of all 4-letter strings.
4. L^* is the set of all strings of letters, including ϵ , the empty string.
5. $L(L \cup D)^*$ is the set of all strings of letters and digits beginning with a letter.
6. D^+ is the set of all strings of one or more digits.



正则语言 (*Regular Language*)

3.3.3 Regular Expressions

正则语言的优点:

- 容易理解
- 能高效实现
- 具有坚实的理论基础



中山大學

正则表达式及其描述的语言

递归定义:

- **Basis:**

ϵ is a regular expr; $L(\epsilon) = \{\epsilon\}$.

$a \in \Sigma$ is a regular expr; $L(a) = \{a\}$.

- **Induction: if r and s are regular exprs,**

$r \mid s$ is a regular expr; $L(r \mid s) = L(r) \cup L(s)$.

rs is a regular expr; $L(rs) = L(r) L(s)$.

r^* is a regular expr; $L(r^*) = (L(r))^*$.

(r) is a regular expr; $L((r)) = L(r)$.



中山大學

一些规定

为了去除不必要的符号，我们规定：

- 一元运算符*具有最高的优先级
- 连接运算符具有次高优先级
- |运算符的优先级最低
- 上述运算符都是左结合的 (left associative)



中山大學



中山大學

例子

令字母表为 $\{a, b\}$

1. The regular expression $\mathbf{a|b}$ denotes the language $\{a, b\}$.
2. $\mathbf{(a|b)(a|b)}$ denotes $\{aa, ab, ba, bb\}$, the language of all strings of length two over the alphabet Σ . Another regular expression for the same language is $\mathbf{aa|ab|ba|bb}$.
3. $\mathbf{a^*}$ denotes the language consisting of all strings of zero or more a 's, that is, $\{\epsilon, a, aa, aaa, \dots\}$.
4. $\mathbf{(a|b)^*}$ denotes the set of all strings consisting of zero or more instances of a or b , that is, all strings of a 's and b 's: $\{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$. Another regular expression for the same language is $\mathbf{(a^*b^*)^*}$.
5. $\mathbf{a|a^*b}$ denotes the language $\{a, b, ab, aab, aaab, \dots\}$, that is, the string a and all strings consisting of zero or more a 's and ending in b .



中山大學

正则语言

- A language that can be defined by a regular expression is called a *regular set*.
- If two regular expressions r and s denote the same regular set, we say they are *equivalent* and write $r = s$.
 - 例如: $(a|b) = (b|a)$



中山大學



中山大學

一些正则语言的等价规则

| LAW | DESCRIPTION |
|----------------------------------|--|
| $r s = s r$ | $ $ is commutative |
| $r (s t) = (r s) t$ | $ $ is associative |
| $r(st) = (rs)t$ | Concatenation is associative |
| $r(s t) = rs rt; (s t)r = sr tr$ | Concatenation distributes over $ $ |
| $\epsilon r = r\epsilon = r$ | ϵ is the identity for concatenation |
| $r^* = (r \epsilon)^*$ | ϵ is guaranteed in a closure |
| $r^{**} = r^*$ | $*$ is idempotent |

中山大學



中山大學

例子：

C语言的标识符：

$(a \mid b \mid \dots \mid z \mid A \mid B \mid \dots \mid Z \mid _)(a \mid b \mid \dots \mid z \mid A \mid B \mid \dots \mid Z \mid _)$
 $\mid (0 \mid 1 \mid \dots \mid 9)^*$

还是太繁琐了，有没更简洁的表示？



中山大學



中山大學

例子：

C语言的标识符：

$(a \mid b \mid \dots \mid z \mid A \mid B \mid \dots \mid Z \mid _)$
 $\mid (0 \mid 1 \mid \dots \mid 9)^*$

还是太繁琐了，有没更简洁的表示？

letter \rightarrow **A** | **B** | ... | **Z** | **a** | **b** | ... | **z**

digit \rightarrow **0** | **1** | ... | **9**

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

中山大學

正则定义 (*Regular Definition*)

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

...

$$d_n \rightarrow r_n$$

where:

1. Each d_i is a new symbol, not in C and not the same as any other of the d'_s , and
2. Each r_i is a regular expression over the alphabet $C \cup \{d_1, d_2, \dots, d_{i-1}\}$.



中山大學

例子

例1: C语言的标识符

$letter_ \rightarrow A \mid B \mid \cdots \mid Z \mid a \mid b \mid \cdots \mid z \mid _$
 $digit \rightarrow 0 \mid 1 \mid \cdots \mid 9$
 $id \rightarrow letter_ (letter_ \mid digit)^*$



中山大學



中山大學

例子

例2：无符号浮点数

digit $\rightarrow 0 \mid 1 \mid \cdots \mid 9$

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

optionalFraction $\rightarrow . \text{digits} \mid \varepsilon$

optionalExponent $\rightarrow (E (+ \mid - \mid \varepsilon) \text{digits}) \mid \varepsilon$

number $\rightarrow \text{digits optionalFraction optionalExponent}$

中山大學



中山大學

正则表达式的扩展

- $r^+ = rr^*$
- $r? = r | \epsilon$
- $[abc] = a | b | c,$
- $[a-z] = a | b | \dots | z$



中山大學



中山大學

例子

C语言标识符:

letter_ → [A-Za-z_]

digit → [0-9]

id → *letter_* (*letter* | *digit*)*

无符号浮点数:

digit → [0-9]

digits → *digit*⁺

number → *digits* (. *digits*)? (E [+-]? *digits*)?



中山大學



中山大學

练习

下列正则表达式描述什么语言?

- $a(a|b)^*a$
- $(a|b)^*a(a|b)(a|b)$
- $a^*ba^*ba^*ba^*$
- $((\epsilon|a)b^*)^*$



中山大學



中山大學

练习

用正则表达式描述下列语言：

- 所有由按词典递增序排列的小写字母组成的字符串

例如：add, low都符合要求，而zzg则不符合

$a*b*c*d*e*f*g*h*i*....$

- 不以ab开头的含有字母a和b的字符串.



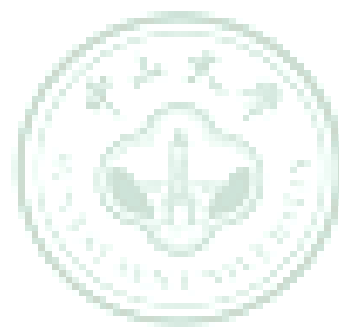
中山大學



中山大學

正则表达式的其它应用

- 文本搜索和模糊匹配
- 合法性检查
- 文本的自动更正和编辑
- 信息提取
-



中山大學



中山大學

Chomsky hierarchy

| Grammar | Languages | Automaton | Production rules (constraints) |
|---------|------------------------|---|--|
| Type-0 | Recursively enumerable | Turing machine | $\alpha \rightarrow \beta$ (no restrictions) |
| Type-1 | Context-sensitive | Linear-bounded non-deterministic Turing machine | $\alpha A \beta \rightarrow \alpha \gamma \beta$ |
| Type-2 | Context-free | Non-deterministic pushdown automaton | $A \rightarrow \gamma$ |
| Type-3 | Regular | Finite state automaton | $A \rightarrow a$ and $A \rightarrow aB$ |



中山大學



中山大學

A Big QUESTION

如何让计算机识别用正则
表达式描述的语言?



中山大學



中山大學

非确定有限自动机 (*NFA*)

A *nondeterministic finite automaton* (NFA) consists of:

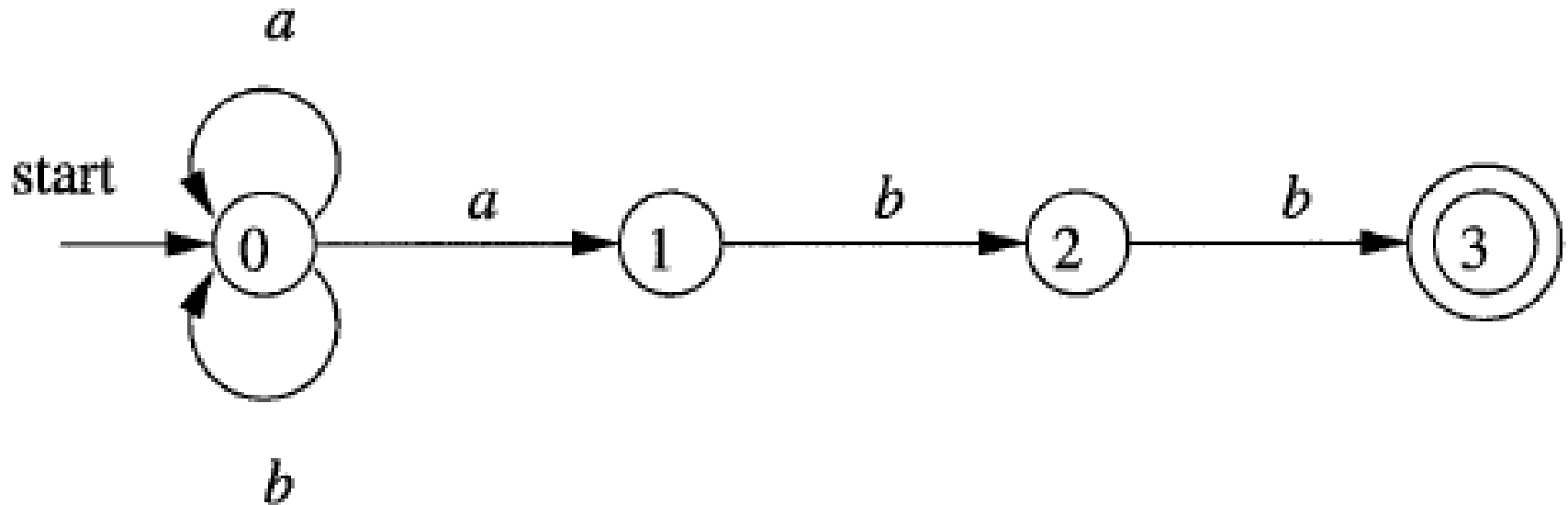
1. A finite set of states S .
2. A set of input symbols Σ , the *input alphabet*. We assume that ϵ , which stands for the empty string, is never a member of Σ .
3. A *transition function* that gives, for each state, and for each symbol in $\Sigma \cup \{\epsilon\}$ a set of *next states*.
4. A state s_0 from S that is distinguished as the *start state* (or *initial state*).
5. A set of states F , a subset of S , that is distinguished as the *accepting states* (or *final states*).

The *language defined* (or *accepted*) by an NFA is the set of strings labeling **some** path from the start to an accepting state.



中山大學

例子：NFA识别语言



识别语言 $L((a|b)^*abb)$.

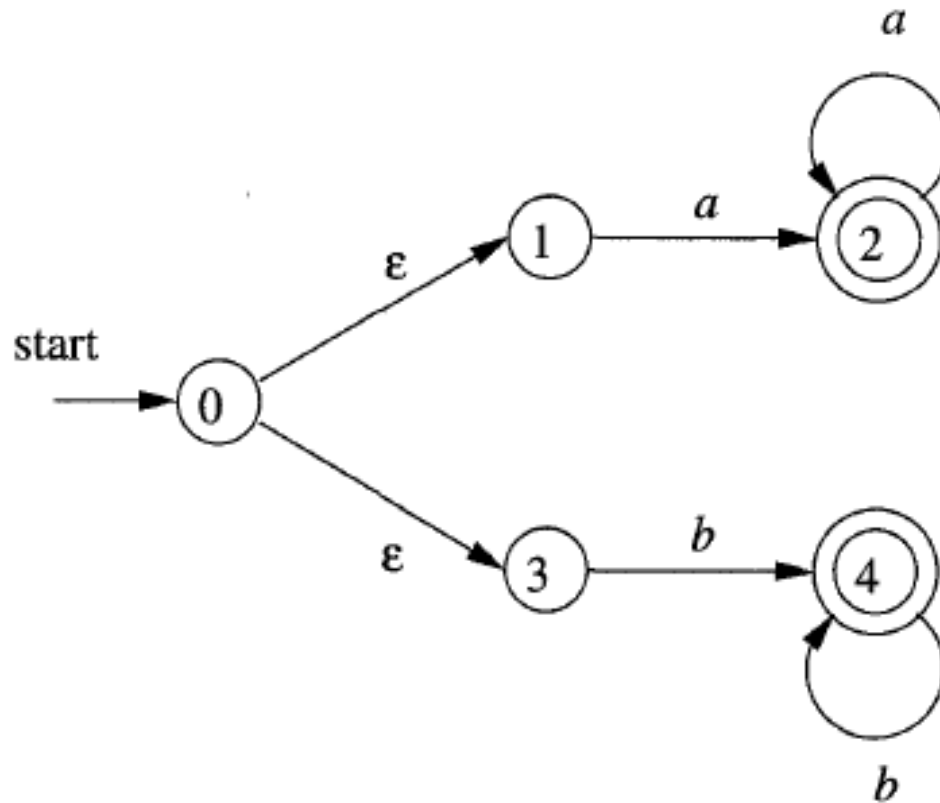


中山大學



中山大學

例子：带 ϵ 的NFA



识别语言 $L(aa^* | bb^*)$.

中山大學



确定有限自动机 (DFA)

A deterministic finite automaton (DFA) is a special case of an NFA where:

1. There are no moves on input ϵ , and
2. For each state s and input symbol a , there is exactly one edge out of s labeled a .



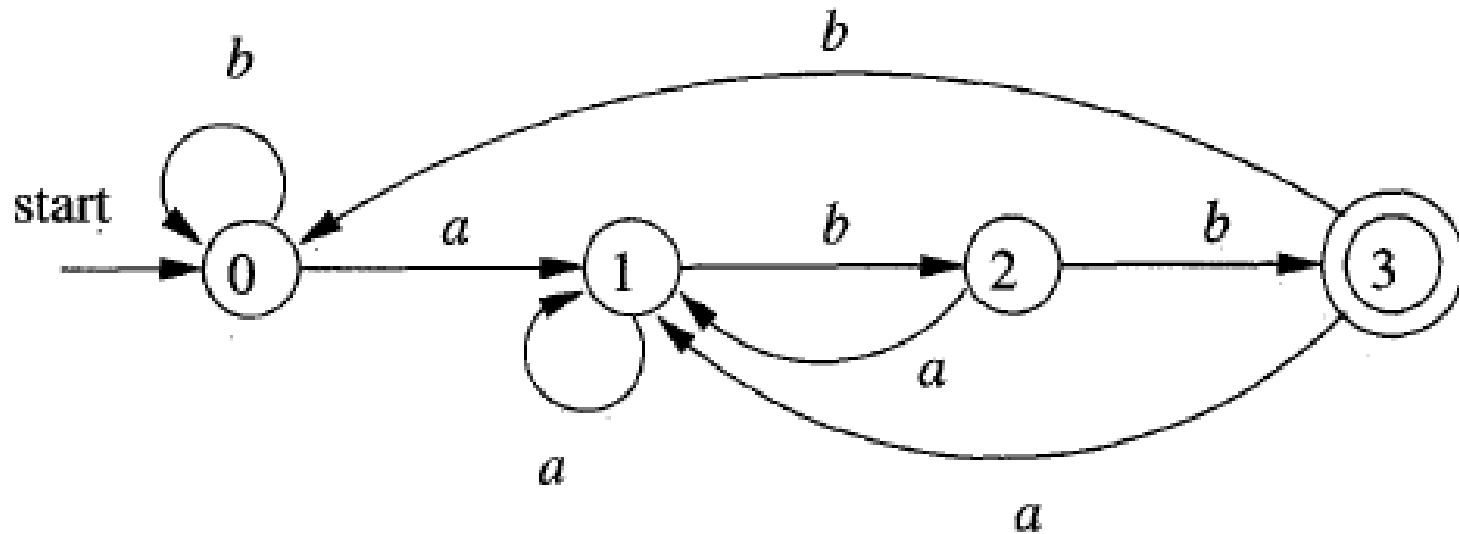
浙江大学



中山大學

例子

识别 $(a|b)^*abb$ 的DFA



中山大學



中山大學

基于 DFA 的识别算法

```
 $s = s_0;$   
 $c = nextChar();$   
while (  $c \neq eof$  ) {  
     $s = move(s, c);$   
     $c = nextChar();$   
}  
if (  $s$  is in  $F$  ) return "yes";  
else return "no";
```

中山大學



中山大學

See you next time!



中山大學