



第二章 词法分析

词法分析的基本功能

正则表达式

词法分析的基本功能

词法分析程序 是编译程序的一部分，
是整个编译过程的第一步工作。

词法分析器 读取源程序的字符序列，逐个
拼出单词并构造相应的内部表示。同时检查
源程序中的词法错误。它的核心作用即为将
字符序列转化为计算机内部表示。

抽取单词序列的例子



抽取单词
序列的例子

```
if (position > 10) rate = 3.14 * initial;
```



```
<$if,->,<$open,->,<$id,position>,<$gt, ->,  
<$num, 10>,<$close, ->,<$id, rate>,<$eq,->,  
<$num, 3.14>,<$mult,->,<$id, initial>,<$semi,  
->
```

词法分析的基本功能

抽取单词
序列的例子

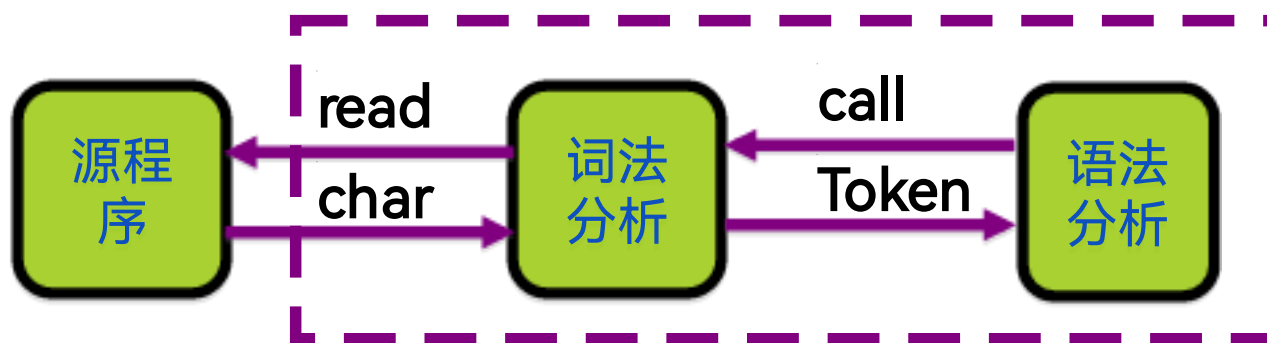
词法分析器
的接口

词法分析器的接口



词法分析器有两类

一类是仅作为语法分析的子程序:



另一类是作为编译器的独立一遍处理器:



单词及单词类型的划分



单词

是指语言中具有独立含义的最小的语义单位。

```
if (position > 10) rate = 3.14 * initial;
```

例如 3.14*initial就可以划分成:

3.14, *,initial这三个单词。

但是3.14不可以继续划分成3,,14

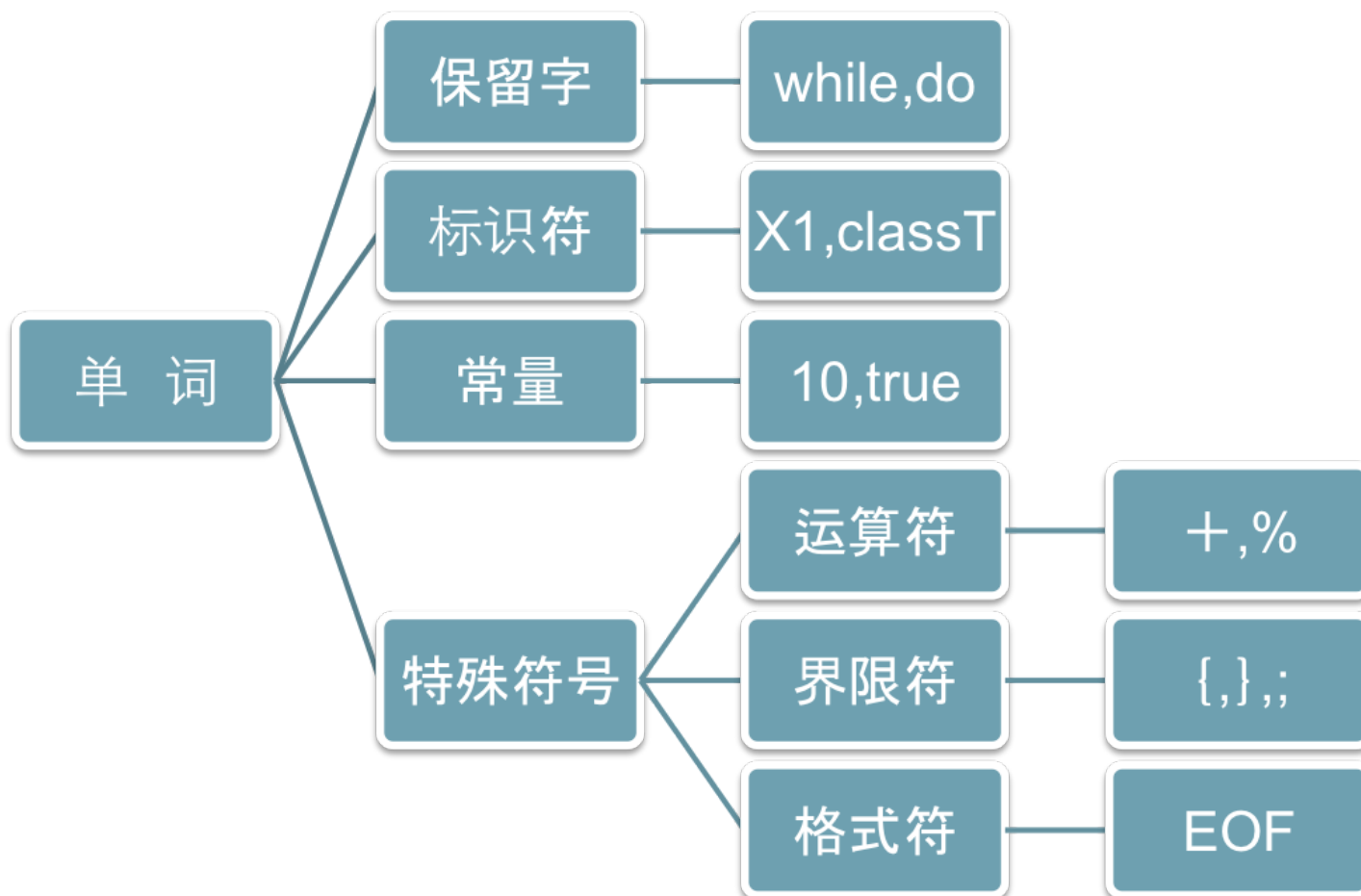
词法分析的基本功能

抽取单词
序列的例子

词法分析器
的接口

单词及单词 类型的划分

单词及单词类型的划分



词法分析的基本功能

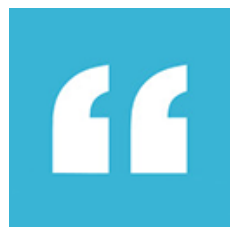
抽取单词
序列的例子

词法分析器
的接口

单词及单词
类型的划分

思考如何实现
词法分析

思考如何实现词法分析



你怎么做？

怎么做合理？



 把问题分析清楚

 采用何种描述方式

 设计算法

词法分析的基本功能

抽取单词
序列的例子

词法分析器
的接口

单词及单词
类型的划分

思考如何实现
词法分析

单词的描述工具

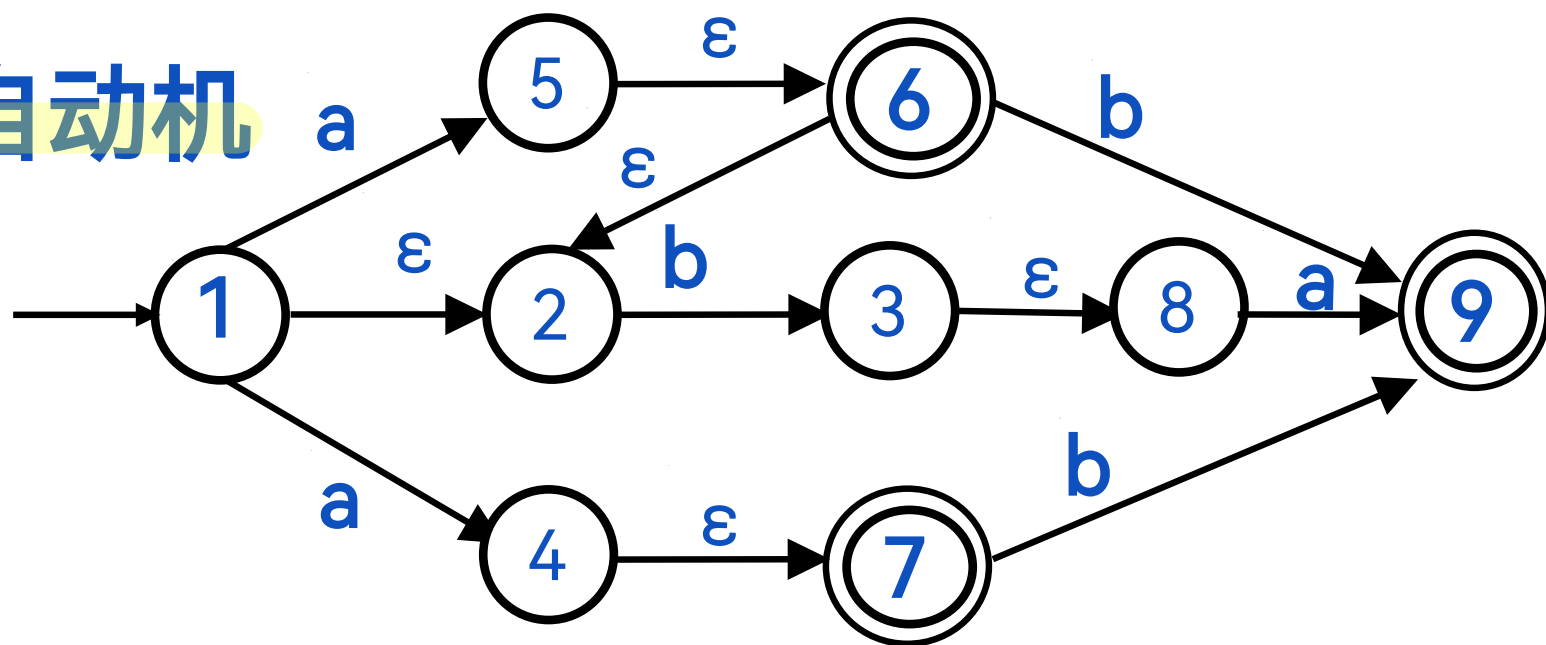
单词的描述工具



正则表达式

$((y|z)^*x(y|z)^*x)^*(y|z)^*$

自动机



基本概念

NEW

字母表 Σ , 元素的非空有穷集合。

符号串 由字母表中的符号组成的任何有穷序列。或者如下定义:

1. **空符号串** (用 ε 表示) 是 Σ 上的符号串
2. 若 α 是 Σ 上的符号串, x 是 Σ 的元素, 则 $x\alpha$ 是 Σ 上的符号串
3. β 是 Σ 上的符号串, 当且仅当它可以由 1 和 2 导出.

基本概念



符号串的连接

设 α 和 β 均是字母表 Σ 上的符号串，它们的连接是把 β 的所有符号顺序接在 α 的符号之后所得到的符号串。

例如：

有： $\alpha=abc$ $\beta=def$

则： $\alpha\beta=abcdef$

特殊情况 对于空串 ε 来说

$$\varepsilon\alpha=\alpha\varepsilon$$

基本概念



符号串的方幂

基本概念

设 a 是字母表 Σ 上的符号串，把 a 自身连接 n 次得到的符号串 a^n ，称作符号串 a 的 n 次幂，记作 $a^n = a^n$ 。

$$a^0 = \varepsilon$$

$$a^1 = a$$

$$a^2 = aa$$

$$a^3 = a^2a = aa^2 = aaa$$

...

$$a^n = a^{n-1}a = aa^{n-1} = aa \dots a \quad (n \text{ 个 } a)$$

基本概念



从这里开始集合.

符号串集合的乘积

设A、B 是两个符号串集合，AB表示A与B的乘积，则定义 $AB = \{xy | (x \in A) \wedge (y \in B)\}$

例如：

✓ $A = \{ab, cd\}, B = \{ef, gh\}$

$AB = \{abef, abgh, cdef, cdgh\}$

特殊：

$A\emptyset = \emptyset A = \emptyset$

基本概念



符号串集合的方幂

设 A 是符号串集合，则称 A^i 是符号串集合 A 的方幂，其中 i 是非负整数。

$$A^0 = \{\epsilon\}, A^1 = A, A^2 = AA, \dots, A^n = AA \dots A$$

例如：

$A = \{a, b\}$ A^3 就是由 a, b 组成的任意长度为 3 的串集。如 aab 、 aba 等等。

$\{a, b\} \{a, b\} \{a, b\}$ 分别从每个里取一个元素

基本概念



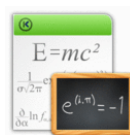
符号串集合的正闭包

设A是符号串集合，则称 A^+ 是符号串集合A的正闭包 $A^+ = A^1 \cup A^2 \cup A^3 \dots \cup A^n \dots$

符号串集合的星闭包

设A是符号串集合，则称 A^* 是符号串集合A的星闭包 $A^* = A^0 \cup A^1 \cup A^2 \cup A^3 \dots \cup A^n \dots = A^0 \cup A^+$

正则表达式



定义

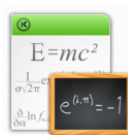
Σ 上的所有正则表达式的集合

设 Σ 为字母表，RE 是定义在 Σ 上的正则表达式集，则有：

证：

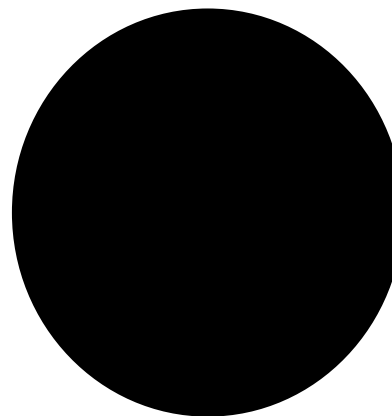
1. $\emptyset, \varepsilon \in \text{RE}$;
2. 对于任意符号 $a \in \Sigma$ ，则 $a \in \text{RE}$;
3. 若 $r, s \in \text{RE}$ ，则 $r|s \in \text{RE}$ ，
 $r \cdot s \in \text{RE}$ ， $r^* \in \text{RE}$;

正则表达式



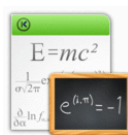
什么是解释，什么是语义？

1+1=?



定义的三点就是正则表达式的形式，它有什么含义需要我们赋予它相应的解释

正则表达式

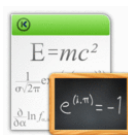


正则表达式对应的语义解释被称作**正则集**
正则集也是正则表达式所表示的**语言**。

在词法分析中，正则表达式是针对单词进行描述的。为此，我们要建立一种由正则表达式到字符串集合的**映射关系**，使得正则表达式的语义解释被描述成字符串的形式。

正则表达式 \rightarrow 字符串集

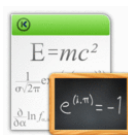
正则表达式



若设 e 、 e_1 、 e_2 为 Σ 上的正则表达式，则 e 所对应的正则集 $L(e)$ 取值如下：

- 当 $e=\emptyset$ 时， $L(e)=\emptyset$;
- 当 $e=\varepsilon$ 时， $L(e)=\{\varepsilon\}$;
- 对于 Σ 中一个字符 a ，若 $e=a$ ，则 $L(e)=\{a\}$;
- 当 $e=e_1 \cdot e_2$ 时， $L(e)=L(e_1)L(e_2)$;
- 当 $e=e_1|e_2$ 时， $L(e)=L(e_1) \cup L(e_2)$;
- $L(e^*)=L(e)^*$;
- $L(e^+)=L(e)^+$.

正则表达式



若用 RE 表示 Σ 上的正则表达式, $L(RE)$ 表示 RE 的正则集, 且 A 、 B 都表示正则表达式, a 表示字母表中的任意符号。有:

$$1) \emptyset \in RE \mid (\emptyset) = \{\} \quad 2) \varepsilon \in RE \mid (\varepsilon) = \{\varepsilon\}$$

$$3) a \in RE \mid L(a) = \{a\} \quad 4) (A) \in RE \mid L((A)) = L(A)$$

$$5) A|B \in RE \mid L(A|B) = L(A) \cup L(B)$$

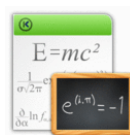
$$6) A \cdot B \in RE \mid L(A \cdot B) = L(A)L(B)$$

$$7) A^* \in RE \mid L(A^*) = L(A)^*;$$

$$8) A^+ \in RE \mid L(A^+) = L(A)^+;$$

.....如 $A?$, $[\text{chi}..\text{chk}]$, $[\text{abc}]$

正则表达式



正则表达式的性质

$$+ = * > . > |$$

$$+ \cdot |$$

运算优先级

$$A | B = B | A$$

| 的可交换性

$$A | (B | C) = (A | B) | C$$

| 的可结合性

$$A (B | C) = (A B) | A C$$

连接的可结合性

$$A (B | C) = A B | A C$$

连接的可分配性

$$(A | B) C = A C | B C$$

连接的可分配性

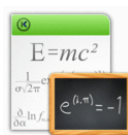
$$A^{**} = A^*$$

幂的等价性

$$A \varepsilon = \varepsilon A = A$$

同一律

正则表达式



用正则表达式描述词法

$L = A|B|\dots|a|b|\dots|z;$

$D = 0|1|\dots|9;$ $D1 = 1|2|\dots|9;$

标识符: $L(L|D)^*$ 正·负·无符号

常数: 整数: $(+|-|\epsilon)(D1D^*)|0$

实数: $(+|-|\epsilon)(D1D^*|0).D^*$

特殊符号: 用枚举的方式来表示

保留字: `while|if|for|...`

运算符: `+|-|*|...`

分界符: `{ }|;|...`

控制符: `\t|\\0|...`

正则表达式的应用



手机中常用的号码地区识别软件

软件的安全监测方法

程序分析技术

正则表达式的局限性

缺乏**对称性字符串**的表达能力

例1: $A = \{a^n b a^n \mid n > 0\}$

例2:

$n \in \mathbb{N};$
 $(AE) \in AE;$
 $AE + AE \in AE$

正则表达式

词法分析
正则表达式
正则表达式
的应用
正则表达式
的局限

例子

例子

 $\Sigma = \{a, b\}$ ab^* $abbb\dots$

Σ 上所有以a为首后跟任意多个
(包括0个) b的符号串集

 $a(ab)^*$ $ab \quad aab \quad abbb \dots$

Σ 上所有以a为首的符号串集

 $\Sigma = \{0, 1\}$ 所有 Σ 上定义的串的正则表达式 $(1|0)^*$

二进制数

能被二整除的二进制数

 $1(1|0)^*0$ $1(1|0)^*0|0$

习题

设字母表 $\Sigma=\{x, y, z\}$, 求:

1、包含偶数个 x 的所有符号串。

$(y/z)^*$

2、不包含连续两个 y 的所有符号串集合。

$((x/z)^* y (x/z)^*)^*$



限定符 (Quantifier)

a* a出现0次或多次 ∞
a+ a出现1次或多次 $1+$
a? a出现0次或1次 有无?
a{6} a出现6次
a{2,6} a出现2-6次
a{2,} a出现两次以上

或运算符 (OR Operator)

(a|b) 匹配a或者b
(ab)|(cd) 匹配ab或者cd

字符类 (Character Classes)

[abc] 匹配a或者b或者c
[a-c] 同上
[a-zA-F0-9] 匹配小写+大写英文字符以及数字
[^0-9] 匹配非数字字符

元字符 (Meta-characters)

\d 匹配数字字符
\D 匹配非数字字符
\w 匹配单词字符(英文、数字、下划线)
\W 匹配非单词字符
\s 匹配空白符(包含换行符、Tab)
\S 匹配非空白字符
. 匹配任意字符(换行符除外)
\bword\b \b标注字符的边界 (全字匹配)
^ 匹配行首
\$ 匹配行尾

贪婪/懒惰匹配 (Greedy / Lazy Match)

<.+> 默认贪婪匹配 “任意字符”
<.+?> 懒惰匹配 “任意字符”

断言 (Lookbehind Assertions)