



中山大學
SUN YAT-SEN UNIVERSITY



国家超级计算广州中心
NATIONAL SUPERCOMPUTER CENTER IN GUANGZHOU

DCS290

Compilation Principle 编译原理

第二章：语言和文法基础

郑馥丹

zhengfd5@mail.sysu.edu.cn

CONTENTS

目录

01

语言和文法的
直观概念

02

符号和
符号串

03

文法和语言的
形式定义

04

文法的
类型

05

上下文无关文
法及其语法树

06

句型的
分析

07

有关文法实用
中的一些说明

1. 程序设计语言

- 程序设计语言包括：语法和语义
 - 语法(syntax)：是一组**规则**，用它可以形成和产生一个合适的程序
 - 语义(semantics)：定义程序的**意义**
 - ✓ 静态语义：程序在语义上要遵守的规则
 - **数组下标越界**
 - **声明和使用的函数没有定义**
 - **零作除数**
 -
 - ✓ 动态语义：表明程序要做什么

2. 文法[Grammar]的直观概念

- 如何来描述一种语言?
 - 如果语言是有穷的（只含有有穷多个句子）：可以将句子逐一系列出来表示
 - 如果语言是无穷的，要找出语言的有穷表示
 - **文法[Grammar]**:
 - ✓ 是语言**语法**的描述工具，实现**用有穷的规则把语言的无穷句子集描述出来**
 - ✓ 严格定义句子的结构，是判断句子结构合法与否的依据
 - 例：“**我是大学生**”是汉语的一个句子

汉语句子的部分构成规则可表示为：

- $\langle \text{句子} \rangle ::= \langle \text{主语} \rangle \langle \text{谓语} \rangle$
- $\langle \text{主语} \rangle ::= \langle \text{代词} \rangle \mid \langle \text{名词} \rangle$
- $\langle \text{代词} \rangle ::= \text{我} \mid \text{你} \mid \text{他}$
- $\langle \text{名词} \rangle ::= \text{王明} \mid \text{大学生} \mid \text{工人} \mid \text{英语}$
- $\langle \text{谓语} \rangle ::= \langle \text{动词} \rangle \langle \text{直接宾语} \rangle$
- $\langle \text{动词} \rangle ::= \text{是} \mid \text{学习}$
- $\langle \text{直接宾语} \rangle ::= \langle \text{代词} \rangle \mid \langle \text{名词} \rangle$

2. 文法[Grammar]的直观概念

• 由规则推导句子

– 方法: 用一条规则, 用**::=右端的符号串代替::=的左端**

✓ **<句子> ::= <主语> <谓语>**

– 表示: 用 “ \Rightarrow ” 表示推导, 其含义是, 使用一条规则, 代替掉 \Rightarrow 左边的某个符号, 产生 \Rightarrow 右端的符号串

– 例如: 句子“我是大学生”的推导过程如下:

<句子>

\Rightarrow **<主语> <谓语>**

\Rightarrow **<代词> <谓语>**

\Rightarrow **我 <谓语>**

\Rightarrow **我 <动词> <直接宾语>**

\Rightarrow **我是 <直接宾语>**

\Rightarrow **我是 <名词>**

\Rightarrow **我是大学生**

汉语句子的部分构成规则可表示为:

- **<句子> ::= <主语> <谓语>**
- **<主语> ::= <代词> | <名词>**
- **<代词> ::= 我 | 你 | 他**
- **<名词> ::= 王明 | 大学生 | 工人 | 英语**
- **<谓语> ::= <动词> <直接宾语>**
- **<动词> ::= 是 | 学习**
- **<直接宾语> ::= <代词> | <名词>**

CONTENTS

目录

01

语言和文法的直观概念

02

符号和符号串

03

文法和语言的形式定义

04

文法的类型

05

上下文无关文法及其语法树

06

句型的分析

07

有关文法实用中的一些说明

1. 字母表[Alphabet] (符号集[a set of symbols])

- 定义：字母表是元素的非空有穷集合

例： $\Sigma=\{0,1\}$ $A=\{a,b,c\}$

- 元素也称为符号，字母表也称符号集
- 不同的语言有不同的字母表
- 程序语言的字母表由**字母、数字和若干专用符号**组成

2. 符号串[String]

- 定义：符号串是由字母表中的符号组成的任何有穷序列

例：0,00,10,011是字母表 $\Sigma=\{0,1\}$ 上的符号串

a,ab,aaca是 $A=\{a,b,c\}$ 上的符号串

- 在符号串中，符号是有顺序的，顺序不同，代表不同的符号串 例：ab和ba不同
- 不含任何符号的符号串称为空串，用 ϵ 表示

注意： $\{\epsilon\}$ 并不等于空集合 $\{\}$ —— Φ

- 符号串长度：是符号串中含有符号的个数

例： $|abc|=3$ $|\epsilon|=0$

- 符号串的头尾，固有头和固有尾：如果 $z=xy$ 是一符号串，则 x 是 z 的头， y 是 z 的尾。如果 x 是非空的，则 y 是固有尾；如果 y 非空，则 x 是固有头。

例： $z=abc$ ， z 的头： ϵ, a, ab, abc ，除 abc 外，其他都是固有头； z 的尾： ϵ, c, bc, abc ，除 abc 外，其他的都是固有尾。

头、尾 \leftrightarrow 前缀、后缀

固有头、固有尾 \leftrightarrow 真前缀、真后缀

3. 符号串的运算

- **符号串的连接**: 设 x 、 y 是符号串, 它们的连接是把 y 的符号写在 x 的符号之后得到的符号串 xy

例如 $x="ST"$, $y="abu"$, 则 $xy="STabu"$

显然 $\epsilon x = x\epsilon = x$

- **符号串的方幂**: 把符号串 a 自身连接 n 次得到的符号串 $a^n = aa\dots aa$

例1 $a^0 = \epsilon$, $a^1 = a$, $a^2 = aa$

例2 $x=AB$, $x^0 = \epsilon$, $x^1 = AB$, $x^2 = ABAB\dots$

4. 符号串集合

- 定义：若集合A中所有元素都是某字母表 Σ 上的符号串，则称A为字母表 Σ 上的符号串集合。
- 符号串集合的乘积：符号串集合A和B的乘积定义为：
 - $AB = \{xy \mid x \in A \text{ 且 } y \in B\}$ ，即AB是由A中的串x和B中的串y连接而成的所有串xy组成的集合。

例：若集合 $A = \{ab, cde\}$ $B = \{0, 1\}$

则 $AB = \{ab0, ab1, cde0, cde1\}$

显然 $\{\varepsilon\}A = A\{\varepsilon\} = A$

4. 符号串集合

- 符号串集合的方幂：设A是符号串的集合，则称 A^i 为符号串集A的方幂，其中i是非负整数。具体定义如下：

$$- A^0 = \{\varepsilon\}, A^1 = A, A^2 = AA, A^k = AA\dots A(k\text{个})$$

例：

$$\text{若集合 } A = \{a, b\} \quad \text{则 } A^0 = \{\varepsilon\} \quad A^1 = A = \{a, b\}$$

$$A^2 = AA = \{a, b\}\{a, b\} = \{aa, ab, ba, bb\}$$

$$A^3 = AAA = \{a, b\}\{a, b\}\{a, b\} = \{aaa, aab, aba, abb, baa, bab, bba, bbb\}$$

5. 集合的闭包[Closure]

- 闭包[Kleene Closure]

- 集合 Σ 的闭包 Σ^* 定义如下: $\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \cup \dots$

- 例: 设有字母表 $\Sigma = \{0, 1\}$**

- 则 $\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \dots$**

- $= \{\epsilon, 0, 1, 00, 01, 10, 11, 000, \dots\}$**

- 即 Σ^* 表示 Σ 上所有有穷长的串的集合。**

- 正闭包[Positive Closure]

- $\Sigma^+ = \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \cup \dots$ 称为 Σ 的正闭包。

- Σ^+ 表示 Σ 上的除 ϵ 外的所有有穷长串的集合

- $\Sigma^* = \Sigma^0 \cup \Sigma^+$**

- $\Sigma^+ = \Sigma \Sigma^* = \Sigma^* \Sigma$**

6. 语言[Language]

- 例如: $\Sigma=\{a,b\}$ $\Sigma^*=\{\epsilon,a,b,aa,ab,ba,bb,aaa,aab,\dots\}$
 1. 集合 $\{ab,aabb,aaabbb,\dots,a^n b^n,\dots\}$ 或 $\{w \mid w \in \Sigma^* \text{ 且 } w=a^n b^n, n \geq 1\}$ 为字母表 Σ 上的一个语言。
 2. 集合 $\{a,aa,aaa,\dots\}$ 或 $\{w \mid w \in \Sigma^* \text{ 且 } w=a^n, n \geq 1\}$ 为字母表 Σ 上的一个语言。
 3. $\{\epsilon\}$ 是一个语言。
 4. Φ 即 $\{\}$ 是一个语言。
 5.

Σ^* 上任意字符串的集合
均是该字母表 Σ 上的语言

CONTENTS

目录

01

语言和文法的
直观概念

02

符号和
符号串

03

文法和语言的
形式定义

04

文法的
类型

05

上下文无关文
法及其语法树

06

句型的
分析

07

有关文法实用
中的一些说明

1. 文法的定义

• 文法定义:

– 文法 **G** 定义为 **四元组** (V_N, V_T, P, S)

- ✓ V_N (Nonterminal): 非终结符集
- ✓ V_T (Terminal): 终结符集
- ✓ P (Production): 产生式 (规则) 集合
- ✓ S (Start): 开始符号或识别符号

- V_N, V_T 和 P 是非空有穷集
- $V = V_N \cup V_T$, V 称为文法 G 的字母表
- $V_N \cap V_T = \Phi$
- S 是一个非终结符, 且至少要在一条产生式的左部出现

• 产生式 (规则) :

– 产生式是一个有序对 (α, β) , 通常写作 $\alpha \rightarrow \beta$ (或 $\alpha ::= \beta$) (读作: α 定义为 β)

$\langle \text{句子} \rangle ::= \langle \text{主语} \rangle \langle \text{谓语} \rangle$

P 中产生式形如: $\alpha \rightarrow \beta$, 其中: $\alpha \in V^+$ 且至少含一个非终结符, $\beta \in V^*$

1. 文法的定义

• 例1：文法 $G=(V_N, V_T, P, S)$ ，其中：

– $V_N = \{\text{句子, 主语, 代词, 名词, 谓语, 动词, 直接宾语}\}$

– $V_T = \{\text{我, 你, 他, 王明, 大学生, 工人, 英语, 是, 学习}\}$

– $P = \{$

<句子> \rightarrow <主语><谓语>

<主语> \rightarrow <代词> | <名词>

<代词> \rightarrow 我 | 你 | 他

<名词> \rightarrow 王明 | 大学生 | 工人 | 英语

<谓语> \rightarrow <动词><直接宾语>

<动词> \rightarrow 是 | 学习

<直接宾语> \rightarrow <代词> | <名词>

}

– $S = \text{句子}$

➤ <句子> ::= <主语><谓语>

➤ <主语> ::= <代词> | <名词>

➤ <代词> ::= 我 | 你 | 他

➤ <名词> ::= 王明 | 大学生 | 工人 | 英语

➤ <谓语> ::= <动词><直接宾语>

➤ <动词> ::= 是 | 学习

➤ <直接宾语> ::= <代词> | <名词>

1. 文法的定义

- 例2：文法 $G=(V_N, V_T, P, S)$ ，其中：
 - $V_N = \{S\}$, $V_T = \{0, 1\}$,
 - $P = \{S \rightarrow 0S1, S \rightarrow 01\}$ ，开始符为 S

- 例3：文法 $G=(V_N, V_T, P, S)$ ，其中：
 - $V_N = \{\text{标识符}, \text{字母}, \text{数字}\}$,
 - $V_T = \{a, b, c, \dots, x, y, z, 0, 1, \dots, 9\}$,
 - $P = \{$
 - $\langle \text{标识符} \rangle \rightarrow \langle \text{字母} \rangle$, $\langle \text{标识符} \rangle \rightarrow \langle \text{标识符} \rangle \langle \text{字母} \rangle$
 - $\langle \text{标识符} \rangle \rightarrow \langle \text{标识符} \rangle \langle \text{数字} \rangle$,
 - $\langle \text{字母} \rangle \rightarrow a, \dots, \langle \text{字母} \rangle \rightarrow z$,
 - $\langle \text{数字} \rangle \rightarrow 0, \dots, \langle \text{数字} \rangle \rightarrow 9$ } ,
 - $S = \langle \text{标识符} \rangle$

随堂练习(1)

- 为只包含数字、加号和减号的表达式，例如 $9-2+5$, $3-1$, 7 等构造一个文法

随堂练习(1)

- 为只包含数字、加号和减号的表达式，例如 $9-2+5, 3-1, 7$ 等构造一个文法

- **G[式子]:**

- ✓ $\langle \text{式子} \rangle \rightarrow \langle \text{数字} \rangle | \langle \text{式子} \rangle \langle \text{运算符} \rangle \langle \text{式子} \rangle$
- ✓ $\langle \text{数字} \rangle \rightarrow 0 | 1 | 2 | 3 | 4 \dots | 9$
- ✓ $\langle \text{数字} \rangle \rightarrow \langle \text{数字} \rangle \langle \text{数字} \rangle$
- ✓ $\langle \text{运算符} \rangle \rightarrow + | -$

2. 文法的简化表示法

- 简化：通常不用将文法的四元组表示出来，**只写出产生式**
- 约定：
 - 默认第一条产生式的左部的符号是开始符号，或用 **$G[S]$** 表示 **S 是开始符号**；
 - 用**大写字母**（或用尖括号括起来）表示**非终结符**；
 - 用**小写字母**表示**终结符**；
 - 左部相同的产生式 **$A \rightarrow \alpha$** ， **$A \rightarrow \beta$** 可以记为； **$A \rightarrow \alpha | \beta$** ，其中“|”表示“或”
- 例如：

文法 $G[S]$:

$S \rightarrow A | SA | SD$

$A \rightarrow a | b | \dots | z$

$D \rightarrow 0 | 1 | \dots | 9$

$V_N = \{S, A, D\}$
 $V_T = \{a, b, \dots, z, 0, 1, \dots, 9\}$
 $P = \{$
 $S \rightarrow A | SA | SD$
 $A \rightarrow a | b | \dots | z$
 $D \rightarrow 0 | 1 | \dots | 9 \}$
 S ——开始符号

3. 推导[Derivation]与归约[Reduction]

(1) 直接推导和直接归约

- $\alpha \rightarrow \beta$ 是文法 G 的产生式, 若有 v, w 满足: $v = \gamma \alpha \delta, w = \gamma \beta \delta$, 其中 $\gamma, \delta \in V^*$, 则称:
 - ✓ v **直接推导** 到 w ,
 - ✓ 或称: w **直接归约** 到 v ,
 - ✓ 记作: $v \Rightarrow w$
- 直接推导: 用产生式的右部替换掉产生式的左部
- 直接归约: 用产生式的左部替换掉产生式的右部

例 文法 G : $S \rightarrow 0S1, S \rightarrow 01$ 有直接推导:

0S1	\Rightarrow	00S11	($S \rightarrow 0S1$)
0S1	\Rightarrow	0011	($S \rightarrow 01$)
00S11	\Rightarrow	000S111	($S \rightarrow 0S1$)
000S111	\Rightarrow	00001111	($S \rightarrow 01$)
S	\Rightarrow	0S1	($S \rightarrow 0S1$)

3. 推导[Derivation]与归约[Reduction]

(2) 推导和归约

– 若存在 $v=w_0 \Rightarrow w_1 \Rightarrow \dots \Rightarrow w_n=w, (n>0)$ (即 v 经过**多步**推到到 w) , 则称:

✓ v 推导出 w ,

✓ 或称: w 归约到 v ,

✓ 记为: $v \Rightarrow^+ w$

– 若有 $v \Rightarrow^+ w$, 或 $v=w$, 则记作 $v \Rightarrow^* w$

$$a \Rightarrow 0\beta, \beta \Rightarrow 2\gamma$$

$$a \Rightarrow 0\beta \Rightarrow 02\gamma$$

$$a \Rightarrow^+ 02\gamma$$

$$a \Rightarrow^* 02\gamma$$

例 文法 G : $S \rightarrow 0S1, S \rightarrow 01$

$$\underline{S} \Rightarrow \underline{0S1} \Rightarrow 00\underline{S11} \Rightarrow 000\underline{S111} \Rightarrow 0000\underline{1111}$$

$$S \Rightarrow^+ 00001111$$

$$S \Rightarrow^* 00001111$$

4. 句型、句子、语言的定义

(2) 语言的定义

– 语言：文法 $G[S]$ 的一切**句子的集合**称为语言，记做 $L(G)$

例 文法 G ：S \rightarrow 0S1, S \rightarrow 01

$S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 0^3S1^3 \Rightarrow \dots \Rightarrow 0^{n-1}S1^{n-1} \Rightarrow 0^n1^n$

$L(G) = \{0^n1^n | n \geq 1\}$

4. 句型、句子、语言的定义

(2) 语言的定义

$$G = (V_N, V_T, P, S), V_N = \{S, B, E\}, V_T = \{a, b, e\}$$

(1) $S \rightarrow aSBE$ 使用(1) $n-1$ 次, 得到 $S \xRightarrow{*} a^{n-1}S(BE)^{n-1}$

(2) $S \rightarrow aBE$ 使用(2) 1次, 得到 $S \xRightarrow{*} a^n(BE)^n$

(3) $EB \rightarrow BE$ 使用(3) 若干次, 得到 $S \xRightarrow{*} a^n B^n E^n$

(4) $aB \rightarrow ab$ 使用(4) 1次, 得到 $S \xRightarrow{*} a^n b B^{n-1} E^n$

(5) $bB \rightarrow bb$ 使用(5) $n-1$ 次, 得到 $S \xRightarrow{*} a^n b^n E^n$

(6) $bE \rightarrow be$ 使用(6) 1次, 得到 $S \xRightarrow{*} a^n b^n e E^{n-1}$

(7) $eE \rightarrow ee$ 使用(7) $n-1$ 次, 得到 $S \xRightarrow{*} a^n b^n e^n$

所以G产生的语言 $L(G) = \{a^n b^n e^n | n \geq 1\}$

5. 文法的等价

- 若 $L(G1)=L(G2)$ ，即，若两个文法所定义的语言是一样的，则称文法 $G1$ 和 $G2$ 是等价的。

例如：文法 $G1[A]: A \rightarrow 0R \quad A \rightarrow 01 \quad R \rightarrow A1$

$G2[S]: S \rightarrow 0S1 \quad S \rightarrow 01$

所定义的语言都是 0^n1^n

因此，两文法等价。