



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. 0型文法

- 通过对产生式施加不同的限制，Chomsky将文法分为**四种类型**：
 - **0型**
 - **1型 (上下文有关文法)**
 - **2型 (上下文无关文法)**
 - **3型 (正规文法)**
- **0型文法(短语文法)**：对任一产生式 $\alpha \rightarrow \beta$ ，都有 $\alpha \in (V_N \cup V_T)^+$ 且至少含有一个非终结符， $\beta \in (V_N \cup V_T)^*$ **(与原文法定义一致)**
- **任何文法都是0型文法。**

2. 1型文法(上下文有关文法[Context-Sensitive Grammar])

• **1型文法 (上下文有关文法)** : 0型文法的特例

- 设文法 $G=(V_N, V_T, P, S)$, 对 P 中的任一产生式 $\alpha \rightarrow \beta$, 都有 $|\beta| \geq |\alpha|$, 仅 $S \rightarrow \varepsilon$ 除外

例 文法 $G[S]$:

$$S \rightarrow aSBE \quad S \rightarrow aBE \quad EB \rightarrow BE$$

$$aB \rightarrow ab \quad bB \rightarrow bb \quad bE \rightarrow be \quad eE \rightarrow ee$$

- 1型文法产生式的一般形式是 $\alpha A \gamma \rightarrow \alpha \beta \gamma$, $\alpha, \gamma \in V^*$, $A \in V_N$, $\beta \in V^+$, 它表示当 A 的上文为 α 且下文为 γ 时可把 A 替换成 β , 因此称1型文法为**上下文有关文法**。

3. 2型文法(上下文无关文法[Context-Free Grammar])

- **2型文法 (上下文无关文法[Context-Free Grammar, CFG])** : 1型文法的特例

- 对任一产生式 $\alpha \rightarrow \beta$, 都有 $\alpha \in V_N$, $\beta \in (V_N \cup V_T)^*$

例 文法G[S]: $S \rightarrow AB$ $A \rightarrow BS|0$ $B \rightarrow SA|1$

- 2型文法产生式的一般形式是: $A \rightarrow \beta$, 它表示不管A的上下文如何即可把A替换成 β , 因此被称为上下文无关文法。
- **通常程序设计语言的文法——2型文法**

4. 3型文法 (正规文法[Regular Grammar])

- **3型文法(正规文法)**: 2型文法的特例

- 任一产生式 $\alpha \rightarrow \beta$ 的形式都为 $A \rightarrow aB$ 或 $A \rightarrow a$, 其中 $A, B \in V_N, a \in V_T$

例如 文法 $G[S]$: $S \rightarrow 0A | 1B | 0$

$A \rightarrow 0A | 1B | 0S$

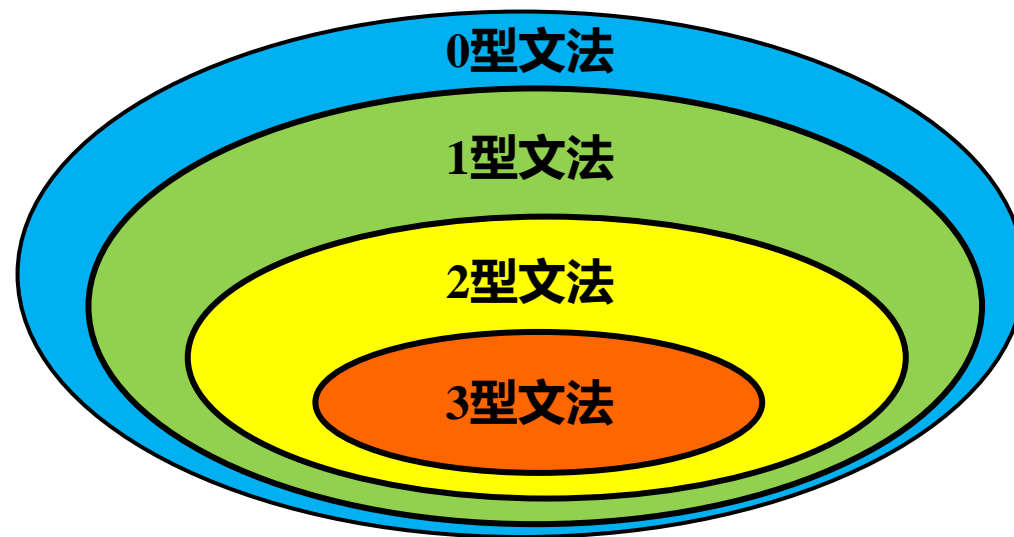
$B \rightarrow 1B | 1 | 0$

- 在程序设计语言中, 3型文法通常用来描述单词的结构

小结

文法类别	产生式形式	说明
0型文法 (短语文法)	$\alpha \rightarrow \beta$, $\alpha \in V^+$, 且至少含一个非终结符, $\beta \in V^*$	对产生式基本无限制
1型文法 (上下文有关文法)	$\alpha \rightarrow \beta$, $ \beta \geq \alpha \geq 1$ 或 $\alpha A \gamma \rightarrow \alpha \beta \gamma$, $\alpha, \gamma \in V^*$, $A \in V_N$, $\beta \in V^+$	将 A 替换成 β 时, 必须考虑 A 的上下文 α, γ
2型文法 (上下文无关文法)	$A \rightarrow \beta$, $A \in V_N$, $\beta \in V^*$	无需考虑 A 在上下文中的出现情况
3型文法 (正规文法)	$A \rightarrow aB$ 或 $A \rightarrow a$, $A, B \in V_N$, $a \in V_T$	产生式全部是规定的形式

四种文法之间的逐级“包含”关系



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. 上下文无关文法 (2型)

- 上下文无关文法有足够的描述能力描述现今程序设计语言的语法结构
- 因此, 我们只关心上下文无关文法形成的语言中的句子的分析

例1: 算术表达式: $E \rightarrow i | E + E | E * E | (E)$

例2: $\langle \text{赋值语句} \rangle \rightarrow i := E$

$\langle \text{条件语句} \rangle \rightarrow \text{if } \langle \text{条件} \rangle \text{ then } \langle \text{语句} \rangle$

$| \text{if } \langle \text{条件} \rangle \text{ then } \langle \text{语句} \rangle \text{ else } \langle \text{语句} \rangle$

2. 规范推导和规范句型

- 如果在推导的任何一步 $\alpha \rightarrow \beta$, 其中 α 、 β 是句型, 都是对 α 中的最左非终结符进行替换, 则称这种推导为**最左推导**; 同理, 如果是总对最右的非终结符进行替换, 则称为**最右推导**
- **最右推导**被称为**规范推导**
- 由规范推导所得的句型称为**规范句型 (右句型)**

例 文法G: $E \rightarrow E+T \mid T$

$T \rightarrow T \times F \mid F$

$F \rightarrow (E) \mid i$

句子 $i+i \times i$ 的推导过程如下:

最左推导: $\underline{E} \Rightarrow \underline{E}+T \Rightarrow \underline{T}+T \Rightarrow \underline{F}+T \Rightarrow i+\underline{T} \Rightarrow i+\underline{T} \times F \Rightarrow i+\underline{F} \times F \Rightarrow i+i \times \underline{F} \Rightarrow i+i \times i$

最右推导: $\underline{E} \Rightarrow \underline{E}+\underline{T} \Rightarrow E+T \times \underline{F} \Rightarrow E+\underline{T} \times i \Rightarrow E+\underline{F} \times i \Rightarrow \underline{E}+i \times i \Rightarrow \underline{T}+i \times i \Rightarrow \underline{F}+i \times i \Rightarrow \underline{i}+i \times i$

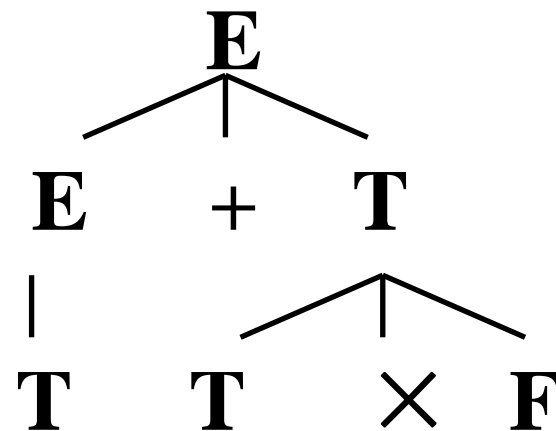
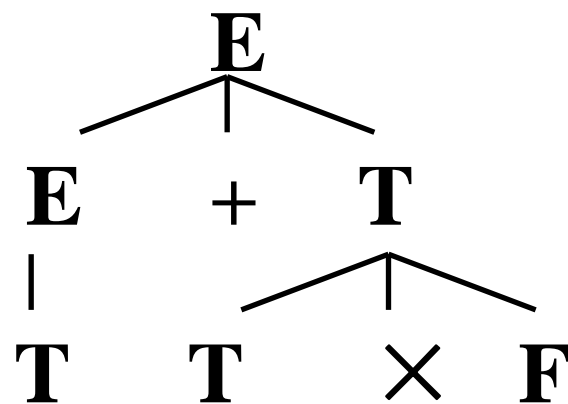
3. 语法树(推导树)[Parse Tree]

- 作用：直观地描述上下文无关文法的**句型推导过程**。给定文 $G=(V_N, V_T, P, S)$ ，对于 G 的任何句型都能构造与之关联的语法树

文法 G : $E \rightarrow E+T \mid T$
 $T \rightarrow T \times F \mid F$
 $F \rightarrow (E) \mid i$

$\underline{E} \Rightarrow \underline{E+T} \Rightarrow \underline{E+T} \times F \Rightarrow T+T \times F$

$\underline{E} \Rightarrow \underline{E+T} \Rightarrow T+\underline{T} \Rightarrow T+T \times F$



从左到右读出叶子结点得到的符号串，为文法的句型。也把该语法树称为该句型的语法树。

从语法树中看不出句型中的符号被替代的顺序

3. 语法树(推导树)[Parse Tree]

- 语法树定义:

- 给定文法 $G=(V_N, V_T, P, S)$, 若一棵树满足下列4个条件, 则称此树为 G 的语法树:

1. 每个结点都有一个标记, 此标记是 V 的一个符号

- 2. 根的标记是 S**

3. 若一结点 n 至少有一个它自己除外的子孙, 并且有标记 A , 则肯定 $A \in V_N$

4. 如果结点 n 有标记 A , 其直接子孙结点从左到右的次序是 n_1, n_2, \dots, n_k , 其标记分别为 A_1, A_2, \dots, A_k , 那么 $A \rightarrow A_1 A_2 \dots A_k$ 一定是 P 中的一个产生式

随堂练习(2)

$S \rightarrow SS^* | SS+ | a$

通过此文法如何生成串 $aa+a^*$ ，并为该串构造**语法树**。

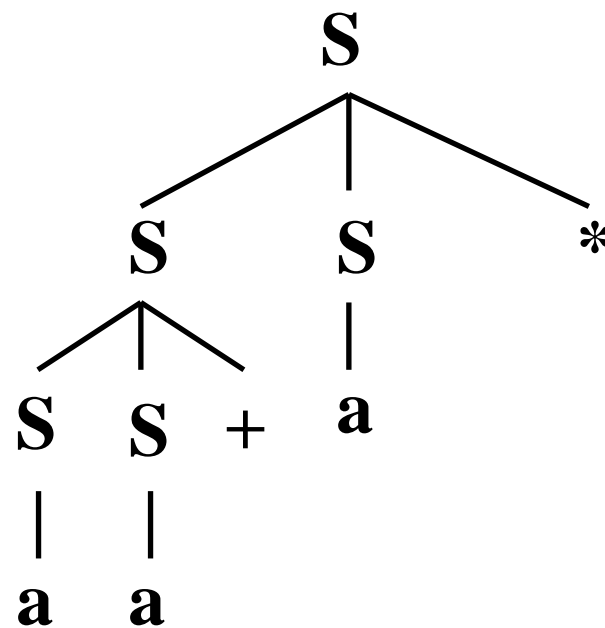
$S \Rightarrow SS^*$

$\Rightarrow SS+S^*$

$\Rightarrow aS+S^*$

$\Rightarrow aa+S^*$

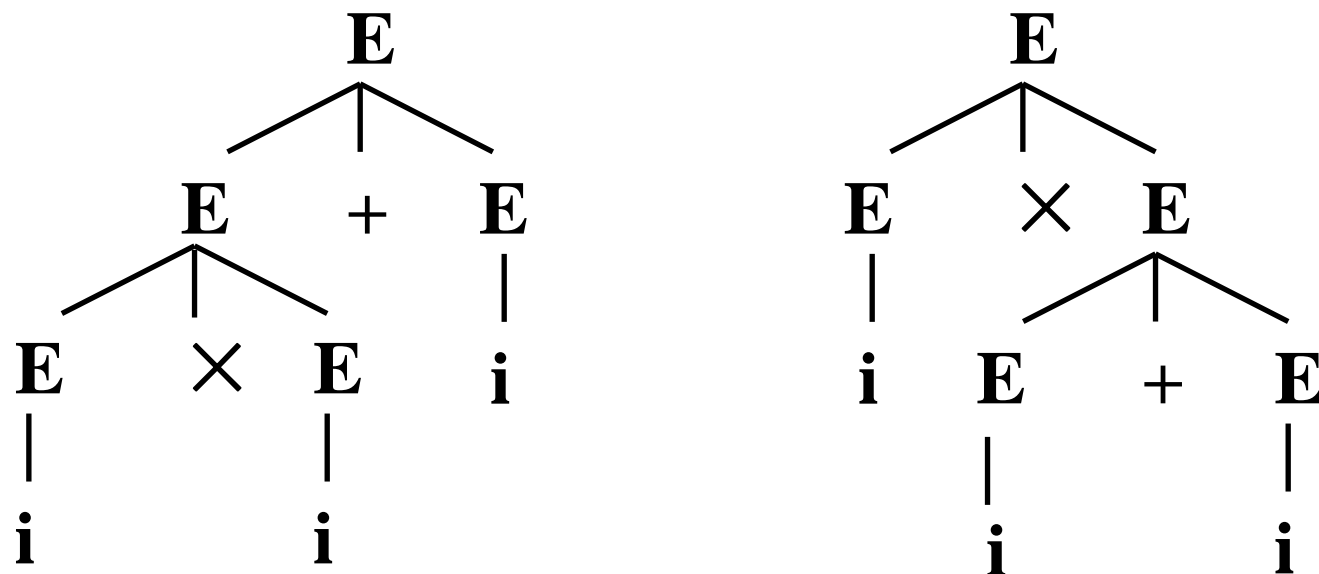
$\Rightarrow aa+a^*$



4. 文法的二义性[Ambiguity]

- 文法G: $E \rightarrow E + E \mid E \times E \mid (E) \mid i$

句子 $i \times i + i$ 对应的语法树



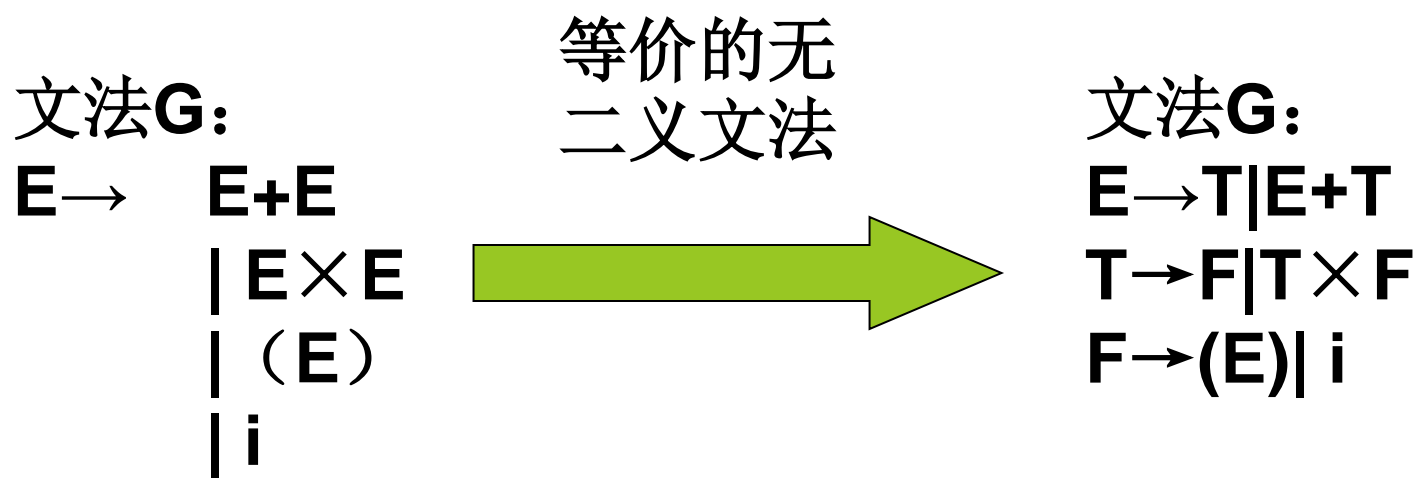
存在两个不同的最左推导:

推导1: $\underline{E} \Rightarrow \underline{E} + E \Rightarrow \underline{E} \times E + E \Rightarrow i \times \underline{E} + E \Rightarrow i \times i + \underline{E} \Rightarrow i \times i + i$

推导2: $\underline{E} \Rightarrow \underline{E} \times E \Rightarrow i \times \underline{E} \Rightarrow i \times \underline{E} + E \Rightarrow i \times i + \underline{E} \Rightarrow i \times i + i$

4. 文法的二义性[Ambiguity]

- 定义：如果一个文法存在某个句子对应**两棵不同的语法树**，则说这个文法是**二义的**
- 二义性文法存在某个句子，它有**两个不同的最左（右）推导**



任何一个二义性的文法，都可以转换成一个等价的无二义性文法。

随堂练习(3)

- 考虑文法 $S \rightarrow aSbS \mid bSaS \mid \varepsilon$, 试说明此文法是二义性的。
- 提示: 可以从对于句子 $abab$ 有两个不同的最左推导来说明。
- 解:
 - $S \Rightarrow aSbS \Rightarrow a\varepsilon bS \Rightarrow a\varepsilon baSbS \Rightarrow a\varepsilon ba\varepsilon bS \Rightarrow a\varepsilon ba\varepsilon b\varepsilon \Rightarrow abab$
 - $S \Rightarrow aSbS \Rightarrow abSaSbS \Rightarrow ab\varepsilon aSbS \Rightarrow ab\varepsilon a\varepsilon bS \Rightarrow ab\varepsilon a\varepsilon b\varepsilon \Rightarrow abab$
 - 存在两个不同的最左推导。
 - 因此文法是二义的。

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. 句型分析

- 对于程序设计语言来说，句型分析就是一个识别输入符号串是否为语法上正确的程序的过程。
- **任务**：句型分析就是**识别一个符号串是否为某文法的句型**，是某个推导的构造过程。
- 由文法开始符号 S 推导出的符号串 α （即 $S \Rightarrow \alpha$ ），称为文法 $G[S]$ 的句型。
- **若能根据该文法构造出该符号串的语法树，则该符号串就是该文法的句型。**
- 句型分析算法采用**从左到右的分析算法**，即总是从左到右地识别输入符号串
- 句型的分析算法分类
 - **自上而下分析法 (Top-Down parsing)**
 - **自下而上分析法 (Bottom-Up parsing)**

2. 自上而下的分析方法[Top-Down parsing]

• 定义:

– 从文法的开始符号出发, 反复使用文法的产生式, 寻找与输入符号串匹配的推导。

• 语法树的构造:

– 将文法的开始符号作为语法树的根, 向下逐步建立语法树, 使语法树的末端结点符号串正好是输入符号串。

推导过程: $S \Rightarrow cAd \Rightarrow cabd$

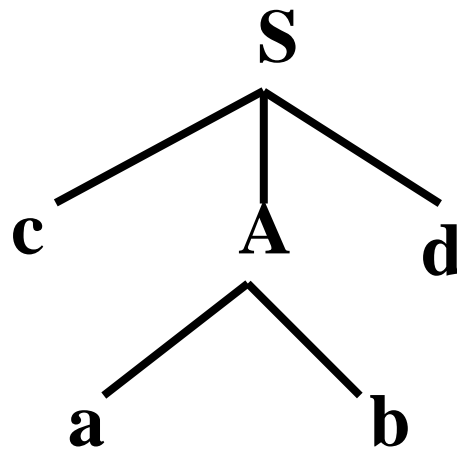
例: 文法G:

$S \rightarrow cAd$

$A \rightarrow ab$

$A \rightarrow a$

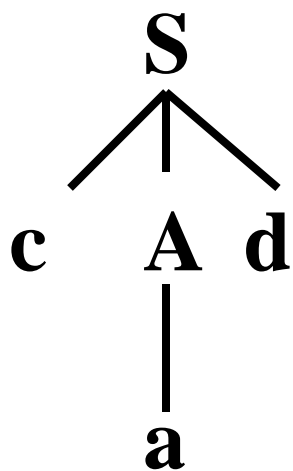
识别输入串 $w=cabd$ 是否为该文法的句子



2. 自上而下的分析方法[Top-Down parsing]

- 自上而下方法的主要问题

例：对输入串cabd自上而下构造语法树的另一过程



不成功!

不成功的原因是

选错产生式 $A \rightarrow a$

$S \rightarrow cAd$

$A \rightarrow ab$

$A \rightarrow a$

自上而下分析的**主要问题是选择产生式**:

假定要被代换的最左非终结符号是 B ，且有 n 条规则： $B \rightarrow A_1 | A_2 | \dots | A_n$ ，那么如何确定用哪个右部去替代 B ？

3. 自下而上的分析方法 [Bottom-Up parsing]

- 定义:

- 从输入符号串开始, 逐步进行归约, 直至归约到文法的开始符号。

- 语法树的构造:

- 从输入符号串开始, 以它作为语法树的末端结点符号串, 自底向上的构造语法树。

归约过程: cabd |- cAd |- S

用 “|-” 表示归约, 下划线部分为被归约符号

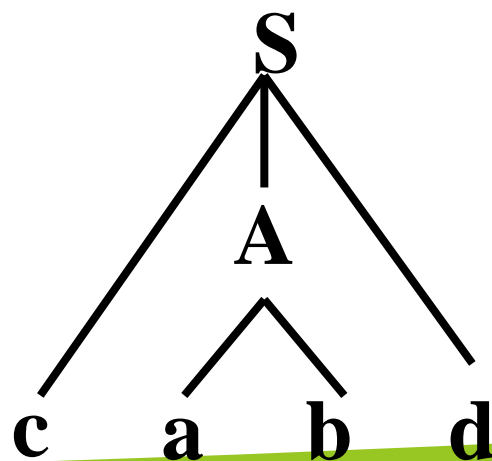
例: 文法G:

$S \rightarrow cAd$

$A \rightarrow ab$

$A \rightarrow a$

识别输入串 $w = cabd$ 是否为该文法的句子



3. 自下而上的分析方法 [Bottom-Up parsing]

- 自下而上分析的主要问题

例：对输入串cabd的两种归约过程

(1) **cabd**|-**cAd**|-S 归约到开始符

(2) **cabd**|-cAbd 不能归约到开始符

S → cAd

A → ab

A → a

- 在自下而上的分析方法中，每一步都是从当前串中选择一个子串加以归约，该子串暂称“**可归约串**”。
- 自下而上分析的**主要问题**：**如何确定“可归约串”——“句柄”**。

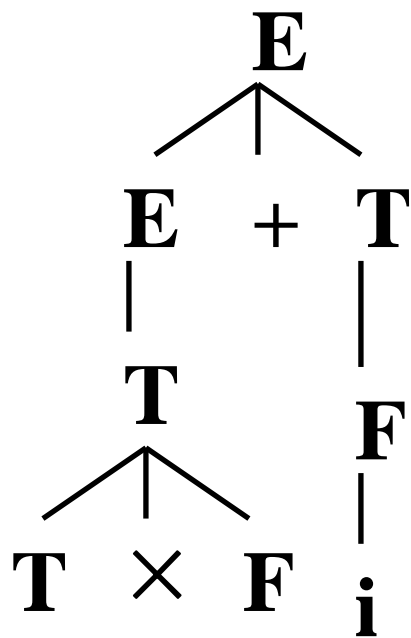
4. 短语[phrase], 直接短语和句柄[handle] (重要!)

- 设 $\alpha\beta\delta$ 是文法 $G[S]$ 中的一个句型, 如果有 $S \overset{*}{\Rightarrow} \alpha A \delta$ 且 $A \Rightarrow \beta$, 则称 β 是句型 $\alpha\beta\delta$ 相对于非终结符 A 的**短语[phrase]**。
- 特别的如有 $A \Rightarrow \beta$, 则称 β 是句型 $\alpha\beta\delta$ 相对于规则 $A \rightarrow \beta$ 的**直接短语 (简单短语)**。
- 一个句型的最左直接短语称为该句型的**句柄[Handle]**。
- **句柄就是“可归约串”**。
- 注意: 在短语的定义中包括了三个条件, 都必须满足:
 - (1) $\alpha\beta\delta$ 是文法的一个句型;
 - (2) $S \overset{*}{\Rightarrow} \alpha A \delta$;
 - (3) $A \Rightarrow \beta$ 。
 - ✓ (1) (2) 说明 $\alpha\beta\delta$ 、 $\alpha A \delta$ 都必须是句型;
 - ✓ (2) (3) 说明, 将 $\alpha\beta\delta$ 中的 β 归约 A 后, 得到的 $\alpha A \delta$ 一定要是句型。假如符号串 β , 将其归约成 A 后得到的符号串不能由开始符号推出, 则 β 不是短语。

4. 短语[phrase], 直接短语和句柄[handle] (重要!)

例: 文法 $G[E]: E \rightarrow E+T|T, T \rightarrow T \times F|F, F \rightarrow (E) | i$ 的一个句型是 $T \times F+i$, 相应的语法树见右图:

1. 因为 $E \overset{*}{\Rightarrow} T+i$ 且 $T \Rightarrow T \times F$, 所以 $T \times F$ 是句型相对于 T 的短语, 且是相对于 $T \rightarrow T \times F$ 的直接短语;
2. 因为 $E \overset{*}{\Rightarrow} T \times F+i$ 且 $F \Rightarrow i$, 所以 i 是句型相对于 F 的短语, 且是相对于 $F \rightarrow i$ 的直接短语;
3. 因为 $E \overset{*}{\Rightarrow} E$ 且 $E \overset{+}{\Rightarrow} T \times F+i$, 所以 $T \times F+i$ 是句型相对于 E 的短语;
4. $T \times F$ 是最左直接短语, 即句柄。



注意: 虽然 $F+i$ 是句型 $T \times F+i$ 的一部分, 但**不是短语**, 因为尽管有 $E \overset{+}{\Rightarrow} F+i$, 但是**不存在从文法开始符 $E \overset{*}{\Rightarrow} T \times E$ 的推导**

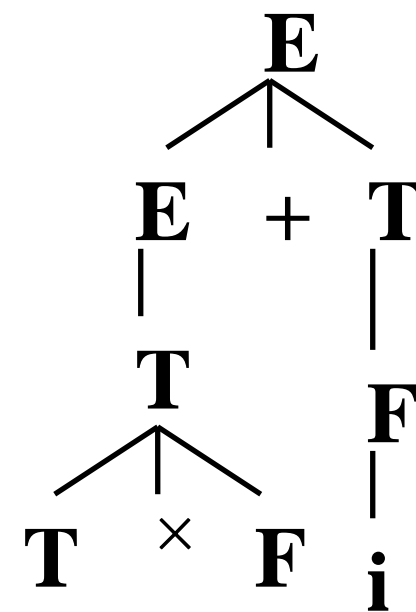
若有 $S \overset{*}{\Rightarrow} \alpha A \delta$ 且 $A \overset{+}{\Rightarrow} \beta$, 则称 β 是句型 $\alpha \beta \delta$ 相对于非终结符 A 的短语。

5. 短语与语法树

- 从句型的语法树上很容易找出句型的短语——语法树中每棵子树的末端结点构成相对于子树根的短语

对于句型 $T \times F + i$ 来说:

- ✓ 五棵子树对应五个**短语** (其中2个重复): $T \times F$, i , $T \times F + i$
- ✓ 两层子树的末端结点构成直接短语, 两棵两层子树对应两个**直接短语**: $T \times F$, i
- ✓ 位于最左边的两层子树的末端结点构成**句柄**: $T \times F$



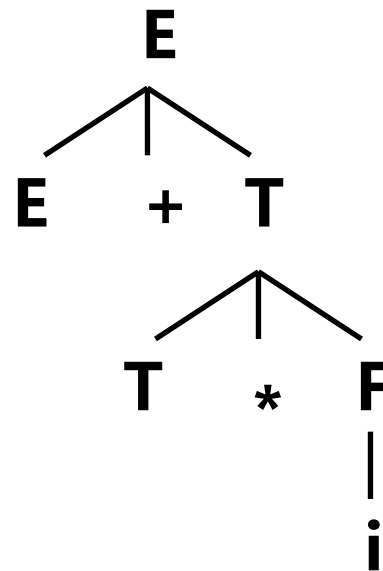
5. 短语与语法树

文法G[E]:

- $E \rightarrow T \mid E + T \mid E - T$
- $T \rightarrow F \mid T * F \mid T / F$
- $F \rightarrow (E) \mid i$ 的一个句型是 $E + T * i$

3棵子树，对应3个短语：

- $E + T * i$
- $T * i$
- i (直接短语、句柄)



5. 短语与语法树

文法 $G[E]: E \rightarrow E+T \mid T, T \rightarrow T \times F \mid F, F \rightarrow (E) \mid i$ 的一个句型是 $F \times i + i$

➤ 7棵子树, 对应有7个短语:

$F \times i + i$ ($E \xrightarrow{*} \underline{E}, E \xrightarrow{*} F \times i + i$)

$F \times i$ ($E \xrightarrow{*} \underline{E} + T, E \xrightarrow{*} F \times i$)

$F \times i$ ($E \xrightarrow{*} \underline{T} + T, T \xrightarrow{*} F \times i$)

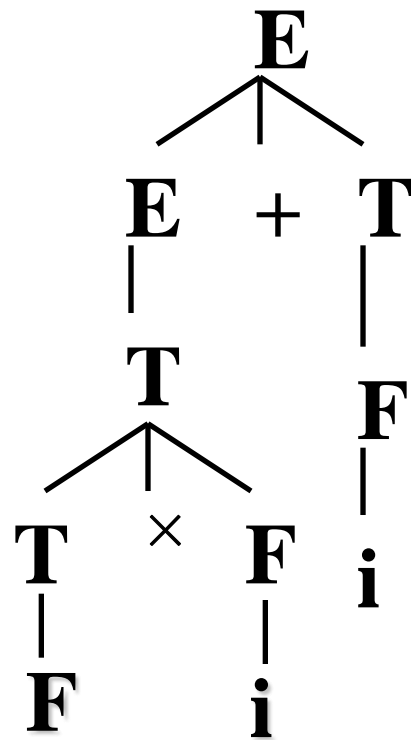
句柄

F ($E \xrightarrow{*} \underline{T} \times i + i, T \xrightarrow{*} F$) (直接短语)

i ($E \xrightarrow{*} F \times \underline{F} + i, F \xrightarrow{*} i$) (直接短语)

i ($E \xrightarrow{*} F \times i + \underline{T}, T \xrightarrow{*} i$)

i ($E \xrightarrow{*} F \times i + \underline{F}, F \xrightarrow{*} i$) (直接短语)



5. 短语与语法树

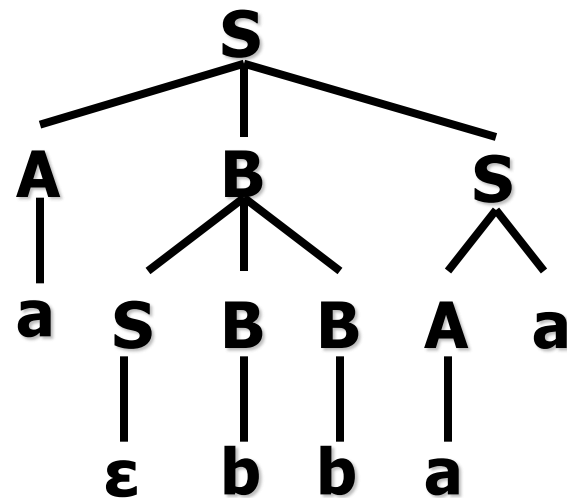
- 小结：找短语、直接短语、句柄的方法：
 - 找短语、直接短语、句柄时，可以通过**语法树**来进行。
 - 句型或句子才有所谓的短语等，所以要**先确定句型或句子**，然后再找它的短语、直接短语、句柄。
 - **要求某一句型的句柄等，语法树只需要分析到该句型。**

随堂练习(4)

求句子abbaa的短语、直接短语、句柄。

8棵子树，对应8个短语：

- abbaa
- a (直接短语)
- bb → (句柄)
- aa
- ϵ (直接短语)
- b (直接短语)
- b (直接短语)
- a (直接短语)



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. 有关文法的实用限制

• 有害规则和多余规则

- 有害规则： $U \rightarrow U$ ，无用且引起文法的二义性
- 多余规则：所有句子推导都用不到的规则，表现在：
 - ✓ 不可到达：文法中某些非终结符不在任何规则的右部出现，所以任何句子的推导中都无法用到它。
 - ✓ 不可终止：不可推导出终结符号串的非终结符。

例：文法 $G[S]$ ：

(1) $S \rightarrow Be$ ~~(2) $B \rightarrow Ce$~~ (3) $B \rightarrow Af$

(4) $A \rightarrow Ae$ (5) $A \rightarrow e$ ~~(6) $C \rightarrow Cf$~~

~~(7) $D \rightarrow f$~~

多余规则为：

- 不可到达 (7)
- 不可终止 (6)
- (2) 也是不可终止的

1. 有关文法的实用限制

- 为避免有害规则和多余规则:

- 非终结符A必须在某句型中出现。即有 $S \Rightarrow^* \alpha A \beta$, 其中 α, β 均属于 $(V_N \cup V_T)^*$
- 必须能从A推出终结符号串t来。即 $A \Rightarrow^* t$, 其中 $t \in V_T^*$

2. 上下文无关文法中的 ϵ 规则

- 上下文无关文法中允许使用 $A \rightarrow \epsilon$ 产生式, $A \rightarrow \epsilon$ 称为 ϵ 规则

第二章课后作业

- 1. 令文法 $G[E]$ 为:

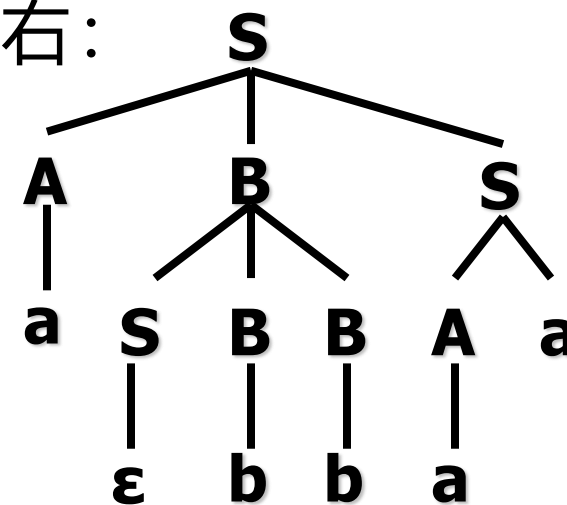
$$E \rightarrow T \mid E + T \mid E - T$$

$$T \rightarrow F \mid T * F \mid T / F$$

$$F \rightarrow (E) \mid i$$

证明 $E + T * F$ 是它的一个规范句型，指出这个句型的所有短语、直接短语和句柄。

- 2. 一个上下文无关文法生成句子 $abbaa$ 的唯一语法树如右:



- (1) 给出该句子相应的最左推导和最右推导。
- (2) 该文法的产生式集合 P 可能有哪些元素?
- (3) 找出该句子的所有短语、简单短语和句柄。

第二章课后作业

- 提交要求：

- 文件命名：学号-姓名-第二章作业；
- 文件格式：.pdf文件；
- 手写版、电子版均可；若为手写版，则拍照后转成pdf提交，但**须注意将照片旋转为正常角度，且去除照片中的多余信息**；电子版如word等转成pdf提交；
- 提交到超算习堂（第二章作业）处；
- 提交ddl：**3月11日晚上12:00**；
- **重要提示：不得抄袭！**

课前/课后 复习/提问

1. 根据右侧的文法推导出句子“他是王明”

〈句子〉

⇒ 〈主语〉 〈谓语〉

⇒ 〈代词〉 〈谓语〉

⇒ 他 〈谓语〉

⇒ 他 〈动词〉 〈直接宾语〉

⇒ 他是 〈直接宾语〉

⇒ 他是 〈名词〉

⇒ 他是 王明

- <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>

2. 进行如下字符串运算

(1) $x = \text{"abc"}$, $y = \text{"defg"}$, 则 $xy = ?$

$xy = \text{"abcdefg"}$

(2) $x = \text{"abc"}$, 则 $x^2 = ?$

$x^2 = \text{"abcabc"}$

(3) $A = \{a, b, c\}$ $B = \{d, e\}$, 则 $AB = ?$

$AB = \{ad, ae, bd, be, cd, ce\}$

(4) $B = \{d, e\}$, 则 $B^2 = ?$

$B^2 = BB = \{d, e\} \{d, e\} = \{dd, de, ed, ee\}$

3. 文法G: $S \rightarrow 0S1, S \rightarrow 01$

由文法G推导出句子0000011111

S \Rightarrow 0S1

\Rightarrow 00S11

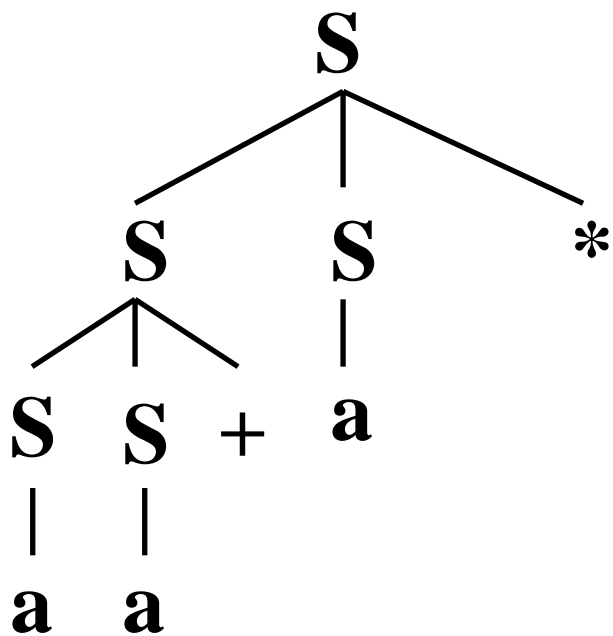
\Rightarrow 000S111

\Rightarrow 0000S1111

\Rightarrow 0000011111

课前/课后 复习/提问

- 4. 找出句型 $aa+a^*$ 的短语、直接短语和句柄



5棵子树, 对应5个短语:

- ✓ $aa+a^*$
- ✓ $aa+$
- ✓ a (直接短语)
- ✓ a (直接短语) (句柄)
- ✓ a (直接短语)