

## 第4章 语法分析

符号流 -> 语法树

### ■ 语法分析器的类型：

从左到右扫描字符

- 通用型
- 自顶向下：通常处理LL文法
- 自底向上：通常处理LR文法

L：寻找串的一个最左推导

R：反向地构造出一个最右推导序列

### 上下文无关文法

组成

- 终结符号
  - 组成串的基本符号，与“词法单元名字”同义
- 非终结符号
  - 语法变量，表示特定串的集合
  - 给出了语言的层次结构，这种层次结构是语法分析和翻译的关键
- 一个开始符号
  - 某个特定的非终结符号，其表示的串集合是这个文法生成的语言
- 一组产生式
  - 描述将终结符号和非终结符号组合成串的方法
  - 产生式左部（头）是一个非终结符号
  - 符号“ $\rightarrow$ ”
  - 一个由零个或多个终结符号与非终结符号组成的产生式右部（体）

句型/句子/语言：

- 句型（sentential form）：
  - 如果 $S \xRightarrow{*} \alpha$ ，那么 $\alpha$ 就是文法的一个句型
  - 可能既包含非终结符，又包含终结符号；可以是空串
- 句子（sentence）
  - 文法的句子就是不包含非终结符号的句型
- 语言
  - 文法G的语言就是G的句子的集合，记为 $L(G)$
  - $\omega$ 在 $L(G)$ 中当且仅当 $\omega$ 是G的句子，即 $S \xRightarrow{*} \omega$

最左推导，最右推导，最左句型，最右句型，语法分析树

## 推导与语法分析树的例子

$$E \Rightarrow - E \Rightarrow - ( E ) \Rightarrow - ( E + E ) \Rightarrow - ( \text{id} + E ) \Rightarrow - ( \text{id} + \text{id} )$$

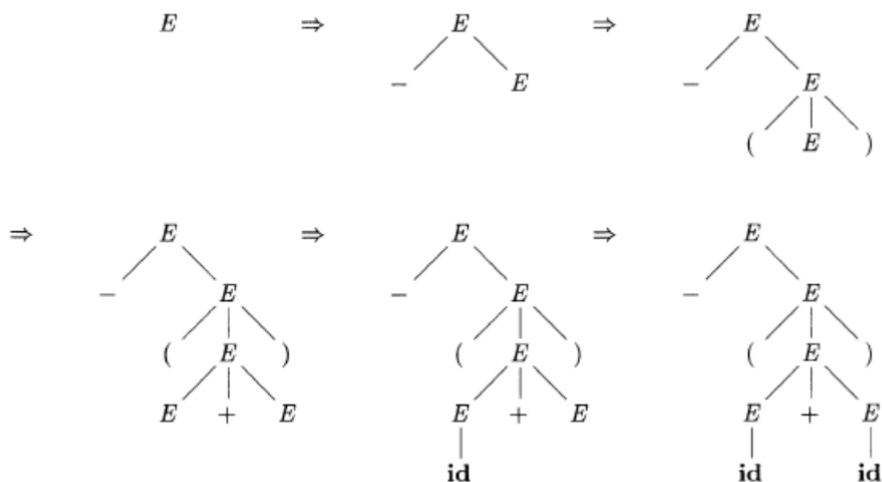


图 4-4 推导(4.8)的语法分析树序列

## 验证文法生成的语言

- 首先证明 $L(G) \subseteq L$ : G生成的每个串都在L中
- 然后证明 $L \subseteq L(G)$ : L的每个串都能由G生成
- 一般使用数学归纳法

不考

## 词法分析和语法分析比较

阶段	输入	输出	描述体系
词法分析	源程序符号串	词法单元序列	正则表达式
句法分析	词法单元序列	语法树	上下文无关文法

- 文法比正则表达式描述能力更强
- 正则表达式描述词法单元比较简洁
- 基于正则表达式构造的词法分析器效率更高
- 正则表达式适合描述词法结构，文法适合描述嵌套结构

上下文无关文法和正则表达式：

上下文无关文法比正则表达式的能力更强，即：

- 一些用文法描述的语言不能用正则文法描述
- 所有的正则语言都可以使用文法描述

不考为NFA构造等价文法

## 消除二义性

不考：二义性的消除方法没有规律可循，通常并不是通过改变文法来消除二义性

## 提取左公因子

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma$$

替换为

$$A \rightarrow \alpha A' \mid \gamma$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

## 消除左递归 (重点)

### 左递归的定义

- 左递归的定义
  - 文法中一个非终结符号A使得对某个串 $\alpha$ ，存在一个推导  $A \Rightarrow A\alpha$ 。则称这个文法是左递归的
- 立即左递归
  - 如果存在  $A \rightarrow A\alpha$ ，则称为立即左递归
- 为什么要消除左递归？
  - 自顶向下的语法分析技术不能处理左递归的文法
- 立即左递归的消除实例：

$$\boxed{A \rightarrow A\alpha \mid \beta} \longrightarrow \begin{cases} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \epsilon \end{cases}$$

### 立即左递归的消除

- 假设非终结符号A存在立即左递归的情形，假设以A为左部的规则有：

$$A \rightarrow \boxed{A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m} \boxed{\beta_1 \mid \beta_2 \mid \dots \mid \beta_n}$$

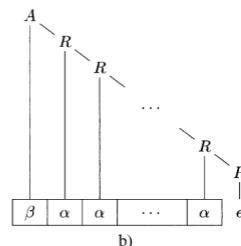
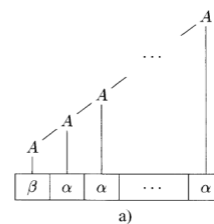
首先对A产生式**分组**(所有  $\alpha_i$  不等于  $\epsilon$ ， $\beta_i$  不以A开头)

可以**替换**为

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$$

- 由A生成的串总是以某个  $\beta_i$  开头，然后跟上零个或者多个  $\alpha_i$  的重复



### 消除多步左递归

- 输入：没有环  $A_i \Rightarrow^+ A_i$  和  $A \rightarrow \epsilon$
- 输出：一个等价的无左递归的文法
- 算法原理：
  - 给非终结符号**排序**， $A_1, A_2, \dots, A_n$
  - 如果只有  $A_i \Rightarrow A_j$  ( $i < j$ )，则不会有左递归
  - 如果发现  $A_i \Rightarrow A_j$  ( $i > j$ )，**代入**  $A_j$  的当前产生式，若替换后有  $A_i$  的直接左递归，再**消除**

伪代码：

- 输入：没有环和 $\epsilon$ 产生式的文法G
- 输出：等价的无左递归的文法
- 步骤：
  - 将文法的非终结符号任意排序为 $A_1, A_2, \dots, A_n$
  - for  $i=1$  to  $n$  do {
    - for  $j = 1$  to  $i-1$  do
      - {
      - 将形如 $A_i \rightarrow A_j \gamma$ 的产生式替换为 $A_i \rightarrow \delta_1 \gamma | \delta_2 \gamma | \dots | \delta_k \gamma$ ,
      - 其中 $A_j \rightarrow \delta_1 | \delta_2 | \dots | \delta_k$ 是以 $A_j$ 为左部的所有产生式;
      - }
    - 消除 $A_i$ 的立即左递归;

## 示例

- $S \rightarrow Aa | b, \quad A \rightarrow Ac | Sd | \epsilon$
- 排序:  $S \ A$
- 替换:
  - $i=1$ , 没有处理
  - $i=2$ , 替换  $A \rightarrow Sd$  中的  $S$ , 得到  $A \rightarrow Ac | Aad | bd | \epsilon$
- 消除立即左递归

$$S \rightarrow Aa | b$$

$$A \rightarrow b d A' | A'$$

$$A' \rightarrow c A' | a d A' | \epsilon$$

## 自顶向下语法分析

在推导的每一步，对非终结符号A，应用哪个产生式，以可能产生于输入串相匹配的终结符号串

### 预测分析技术

通过在输入中向前看固定多个符号来选择正确的产生式

通常情况下，我们只需要向前看一个符号

递归下降分析一般只适合于每个子表达式的第一个终结符能够为产生式选择提供足够信息的那些文法

给出两个与文法相关的两个函数 □ FIRST □ FOLLOW

基于上述两个函数，可以根据下一个输入符号来 选择应用哪个产生式

### First函数计算

First定义:

- 定义：可从 $\alpha$ 推导得到的串的首符号的集合，其中 $\alpha$ 是任意的文法符号串。如果 $\alpha \xRightarrow{*} \epsilon$ ，那么 $\epsilon$ 也在 $FIRST(\alpha)$ 中
- FIRST函数的意义
  - 如果两个A产生式  $A \rightarrow \alpha | \beta$ ，其中 $First(\alpha)$ 和 $First(\beta)$ 是不相交的集合。下一个输入符号是a，若 $a \in First(\alpha)$ ，则选择 $A \rightarrow \alpha$ ，若 $a \in First(\beta)$ ，则选择 $A \rightarrow \beta$
  - 用于预测产生式的选择

First计算:

- 对于文法符号 $X$ 的 $FIRST(X)$ , 通过不断应用下列规则, 直到没有新的终结符号或者 $\epsilon$ 可以加入到任何 $FIRST$ 集合为止
  - 如果 $X$ 是终结符号, 那么 $FIRST(X)=\{X\}$
  - 如果 $X$ 是非终结符号, 且有规则 $X \rightarrow a...$ , 那么将 $a$ 添加到 $FIRST(X)$ 中; 如果 $X \rightarrow \epsilon$ , 那么 $\epsilon$ 也在 $FIRST(X)$ 中
  - 对于规则 $X \rightarrow Y_1Y_2...Y_n$ , 把 $FIRST(Y_1)$ 中的非 $\epsilon$ 符号添加到 $FIRST(X)$ 中。如果 $\epsilon$ 在 $FIRST(Y_1)$ 中, 把 $FIRST(Y_2)$ 中的非 $\epsilon$ 符号添加到 $FIRST(X)$ 中...; 如果 $\epsilon$ 在 $FIRST(Y_n)$ 中, 把 $\epsilon$ 添加到 $FIRST(X)$ 中

示例:

- |                                       |                                     |
|---------------------------------------|-------------------------------------|
| $E \rightarrow T E'$                  | ■ $First(F)=First(T)=First(E)$      |
| $E' \rightarrow + T E' \mid \epsilon$ | $=\{ (, id \}$                      |
| $T \rightarrow F T'$                  | ■ $First(E')=\{ +, \epsilon \}$     |
| $T' \rightarrow * F T' \mid \epsilon$ | ■ $First(T')=\{ *, \epsilon \}$     |
| $F \rightarrow ( E ) \mid id$         | ■ $First(TE')=First(T)=\{ (, id \}$ |
|                                       | ■ $First(+TE')=\{ + \}$             |
|                                       | ■ .....                             |

## Follow函数计算

Follow定义:

- 对于非终结符号 $A$ ,  $FOLLOW(A)$ 定义为可能在某些句型中紧跟在 $A$ 右边的终结符号的集合
  - 例如 $S \xRightarrow{*} \alpha A a \beta$ , 终结符号 $a \in Follow(A)$
- 如果 $A$ 是某些句型的最右符号, 那么 $\$ \in Follow(A)$ 。\$是特殊的输入串“结束标记”
- **FOLLOW函数的意义:**
  - 如果 $A \rightarrow \alpha$ , 当 $\alpha \rightarrow \epsilon$ 或 $\alpha \xRightarrow{*} \epsilon$ 时,  $FOLLOW(A)$ 可以帮助我们做出选择恰当的产生式
  - 例如: 如果 $A \rightarrow \alpha$ ,  $b$ 属于 $FOLLOW(A)$ , 如果 $\alpha \xRightarrow{*} \epsilon$ , 则若当前输入符号是 $b$ , 可以选择 $A \rightarrow \alpha$ , 因为 $A$ 最终到达了 $\epsilon$ , 而且后面跟着 $b$

Follow计算:

- 计算各个非终结符号 $A$ 的 $FOLLOW(A)$ 集合, 不断应用下列规则, 直到没有新的终结符号可以被加入到任意 $FOLLOW$ 集合中
  - 将 $\$$ 放入 $FOLLOW(S)$ ,  $S$ 是开始符号, 而 $\$$ 是输入串的结束标记
  - 如果存在产生式 $A \rightarrow \alpha B \beta$ , 那么 $First(\beta)$ 中除 $\epsilon$ 之外的所有符号都在 $Follow(B)$ 中
  - 如果存在一个产生式 $A \rightarrow \alpha B$ , 或存在产生式 $A \rightarrow \alpha B \beta$ 且 $First(\beta)$ 包含 $\epsilon$ , 那么 $Follow(A)$ 中的所有符号都在 $Follow(B)$ 中

示例:

□ 文法

$$\begin{aligned}
 E &\rightarrow T E' & E' &\rightarrow + T E' \mid \epsilon & T &\rightarrow F T' \\
 T' &\rightarrow * F T' \mid \epsilon & F &\rightarrow ( E ) \mid i
 \end{aligned}$$

□ FOLLOW(E)={, , \$}       $\leftarrow \text{FIRST}( ) \cup \{ \$ \}$   
 □ FOLLOW(E')={, , \$}       $\leftarrow \text{FOLLOW}(E)$   
 □ FOLLOW(T)={+, ), \$}       $\leftarrow \text{FIRST}(E') \cup \text{FOLLOW}(E')$   
 □ FOLLOW(T')={+, ), \$}       $\leftarrow \text{FOLLOW}(T)$   
 □ FOLLOW(F)={\*, +, ), \$}       $\leftarrow \text{FIRST}(T') \cup \text{FOLLOW}(T)$

- 将\$放入FOLLOW(S)，S是开始符号，而\$是输入串的结束标记
- 如果存在产生式 $A \rightarrow \alpha B \beta$ ，那么First( $\beta$ )中除 $\epsilon$ 之外的所有符号都在Follow(B)中
- 如果存在一个产生式 $A \rightarrow \alpha B$ ，或存在产生式 $A \rightarrow \alpha B \beta$ 且First( $\beta$ )包含 $\epsilon$ ，那么Follow(A)中的所有符号都在Follow(B)中

## 预测分析表

- 输入：文法G
- 输出：预测分析表M
- 方法：对于文法G的每个产生式 $A \rightarrow \alpha$ ，进行如下处理
  - 对于First( $\alpha$ )中的每个终结符号a，将 $A \rightarrow \alpha$ 加入到M[A,a]
  - 如果 $\epsilon$ 在First( $\alpha$ )中，那么对于Follow(A)中的每个终结符号b，将 $A \rightarrow \alpha$ 加入到M[A,b]中
  - 如果 $\epsilon$ 在First( $\alpha$ )中，且\$在Follow(A)中，将 $A \rightarrow \alpha$ 加入到M[A,\$]中
- 完成上述操作后，若M[A,a]中没有产生式，填为Error

## 分析表驱动预测分析

已匹配	栈	输入	动作
	E\$	id + id * id\$	
	TE'\$	id + id * id\$	输出 $E \rightarrow TE'$
	FT'E'\$	id + id * id\$	输出 $T \rightarrow FT'$
	id T'E'\$	id + id * id\$	输出 $F \rightarrow id$
id	T'E'\$	+ id * id\$	匹配 id
id	E'\$	+ id * id\$	输出 $T' \rightarrow \epsilon$
id	+ TE'\$	+ id * id\$	输出 $E' \rightarrow + TE'$
id +	TE'\$	id * id\$	匹配 +

## 自底向上语法分析

从叶子（输入串中的终结符号，将位于分析树的底端）开始，向上到达根结点

在每一步的归约中，一个与某产生式体相匹配的特定子串被替换为该产生式头部的非终结符号，一次归约实质上是一个推导的反向操作。

## 句柄



## 句柄

重要



- 最右句型 $\gamma$ 的一个句柄
  - 满足下述条件的产生式 $A \rightarrow \beta$ 及串 $\beta$ 在 $\gamma$ 中出现的位置
  - 条件：将这个位置上的 $\beta$ 替换为 $A$ 之后得到的串是 $\gamma$ 的某个最右推导序列中出现在位于 $\gamma$ 之前的最右句型
- 通俗地说
  - 句柄是最右推导的反向过程中被规约的那些部分
- 句柄的作用
  - 对句柄的规约，代表了相应的最右推导中的一个反向步骤

### 移进归约语法分析框架

- 移入：将下一个输入符号移动到栈顶
- 归约：将句柄归约为相应的非终结符号
  - 句柄总是在栈顶
  - 具体操作时弹出句柄，压入被归约到的非终结符号
- 接受：宣布分析过程成功完成
- 报错：发现语法错误，调用错误恢复子程序

例子：

栈	输入	动作
\$	$id_1 * id_2 \$$	移入
$\$ id_1$	$* id_2 \$$	按照 $F \rightarrow id$ 归约
$\$ F$	$* id_2 \$$	按照 $T \rightarrow F$ 归约
$\$ T$	$* id_2 \$$	移入
$\$ T *$	$id_2 \$$	移入
$\$ T * id_2$	$\$$	按照 $F \rightarrow id$ 归约
$\$ T * F$	$\$$	按照 $T \rightarrow T * F$ 归约
$\$ T$	$\$$	按照 $E \rightarrow T$ 归约
$\$ E$	$\$$	accept

### LR语法分析技术



## ■ LR(k)的语法分析概念

- L表示最左扫描，R表示反向构造出最右推导
- k表示最多向前看k个符号

## ■ 当k的数量增大时，相应的语法分析器的规模急剧增大

- K=2时，程序设计语言的语法分析器的规模通常非常庞大
- 当k=0、1时已经可以解决很多语法分析问题，因此具有实践意义
- 因此，我们只考虑k≤1的情况

### LR(0)项



#### LR(0) 项



- 文法的一个产生式加上在其产生式体中某处的一个点
  - $A \rightarrow .XYZ$ ,  $A \rightarrow X.YZ$ ,  $A \rightarrow XY.Z$ ,  $A \rightarrow XYZ.$
  - 注意:  $A \rightarrow \epsilon$ 只对应一个项 $A \rightarrow .$
- 直观含义
  - 项 $A \rightarrow \alpha.\beta$ 表示已经扫描/归约到了 $\alpha$ ，并期望接下来的输入中经过扫描/归约得到 $\beta$ ，然后把 $\alpha\beta$ 归约到A
  - 如果 $\beta$ 为空，表示我们可以把 $\alpha$ 归约为A
- 项也可以用一对整数表示
  - (i,j)表示第i条规则，点位于右部第j个位置

### 增广文法

## ■ 增广文法

- G的增广文法G'是在G中增加新开始符号S'，并加入产生式 $S' \rightarrow S$ 而得到的
- G'和G接受相同的语言，且按照 $S' \rightarrow S$ 进行归约实际上就表示已经将输入符号串归约成为开始符号
- 引入的目的是告诉语法分析器何时宣布接受输入符号串，即用 $S' \rightarrow S$ 进行归约时，表明分析结束。

即用 $S' \rightarrow S$ 进行归约时，表明分析结束



- 如果I是文法G的一个项集，那么CLOSURE(I)就是根据下列规则从I构造得到的项集
  - 将I中的各个项加入到CLOSURE(I)中
  - 如果 $A \rightarrow \alpha.B\beta$ 在CLOSURE(I)中，那么对B的任意产生式 $B \rightarrow \gamma$ ，将 $B \rightarrow \gamma$ 加到CLOSURE(I)中
  - 不断重复第二步，直到收敛
- 第二步的意义
  - 项 $A \rightarrow \alpha.B\beta$ 表示期望在接下来的输入中归约到B
  - 显然，要归约到B，首先要扫描归约到B的某个产生式的右部
  - 因此对每个产生式 $B \rightarrow \gamma$ ，加入 $B \rightarrow \gamma$ 
    - 表示它期望能够扫描归约到 $\gamma$

GOTO:

- I是一个项集，X是一个文法符号，GOTO(I,X)定义为I中所有形如的项 $[A \rightarrow \alpha \cdot X\beta]$ 所对应的项 $[A \rightarrow \alpha X \cdot \beta]$ 的集合的闭包
  - 根据项的历史-期望的含义，GOTO(I,X)表示读取输入中的X或者归约到一个X之后的情况
  - GOTO(I,X)定义了LR(0)自动机中状态I在X之上的转换
- 例如：
  - $I = \{[E' \rightarrow E.], [E \rightarrow E.+T]\}$
  - GOTO(I,+)计算如下
    - I中只有一个项的点后面跟着+，对应的项为 $[E \rightarrow E+.T]$
    - $CLOSURE(\{[E \rightarrow E+.T]\}) = \{[E \rightarrow E+.T], [T \rightarrow T*F], [T \rightarrow F.], [F \rightarrow .(E)], [F \rightarrow .id]\}$

LR(0)自动机

- 构造方法
  - 规范LR(0)项集族中的项集可以作为LR(0)自动机的状态
  - $GOTO(I,X)=J$ ，则从I到J有一个标号为X的转换
  - 初始状态为 $CLOSURE(\{S' \rightarrow .S\})$ 对应的项集
  - 接受状态：包含形如 $A \rightarrow \alpha.$ 的项集对应的状态

作用:

- 假设文法符号串 $\gamma$ 使LR(0)自动机从开始状态运行到状态（项集） $j$ 
  - 如果 $j$ 中有一个形如 $A \rightarrow \alpha.$ 的项，那么
    - 在 $\gamma$ 之后添加一些终结符号可以得到一个最右句型
    - $\alpha$ 是 $\gamma$ 的后缀，且 $A \rightarrow \alpha$ 是这个句型的句柄
    - 表示可能找到了当前最右句型的句柄
  - 如果 $j$ 中存在一个项 $B \rightarrow \alpha.X\beta$ ，那么
    - 在 $\gamma$ 之后添加 $X\beta$ ，然后再添加一个终结符号串可得到一个最右句型
    - 在这个句型中 $B \rightarrow \alpha X\beta$ 是句柄
    - 此时表示还没有找到句柄，需要移入

#### LR语法分析表

- 两个部分：动作ACTION，转换GOTO
- ACTION有两个参数：状态 $i$ 、终结符号 $a$ 
  - 移入 $j$ ： $j$ 是一个状态。把 $j$ 压入栈
  - 归约 $A \rightarrow \beta$ ：把栈顶的 $\beta$ 归约为 $A$
  - 接受：接受输入、完成分析
  - 报错：在输入中发现语法错误
- 状态集上的GOTO函数
  - 如果 $\text{GOTO}[i, A] = j$ ，那么 $\text{GOTO}[i, A] = j$

例子：



#### LR分析

- 输入： $\text{id} * \text{id} + \text{id}$

	栈	符号	输入	动作
(1)	0		$\text{id} * \text{id} + \text{id} \$$	移入
(2)	0 5	$\text{id}$	$* \text{id} + \text{id} \$$	根据 $F \rightarrow \text{id}$ 归约
(3)	0 3	$F$	$* \text{id} + \text{id} \$$	根据 $T \rightarrow F$ 归约
(4)	0 2	$T$	$* \text{id} + \text{id} \$$	移入
(5)	0 2 7	$T *$	$\text{id} + \text{id} \$$	移入
(6)	0 2 7 5	$T * \text{id}$	$+ \text{id} \$$	根据 $F \rightarrow \text{id}$ 归约
(7)	0 2 7 10	$T * F$	$+ \text{id} \$$	根据 $T \rightarrow T * F$ 归约
(8)	0 2	$T$	$+ \text{id} \$$	根据 $E \rightarrow T$ 归约
(9)	0 1	$E$	$+ \text{id} \$$	移入
(10)	0 1 6	$E +$	$\text{id} \$$	移入
(11)	0 1 6 5	$E + \text{id}$	$\$$	根据 $F \rightarrow \text{id}$ 归约
(12)	0 1 6 3	$E + F$	$\$$	根据 $T \rightarrow F$ 归约
(13)	0 1 6 9	$E + T$	$\$$	根据 $E \rightarrow E + T$ 归约
(14)	0 1	$E$	$\$$	接受

状态	ACTION						GOTO		
	id	+	*	(	)	\$	E	T	F
0	s5			s4			1	2	3
1		s6				acc			
2	r2	s7		r2	r2				
3	r4	r4		r4	r4				
4	s5			s4			8	2	3
5	r6	r6		r6	r6				
6	s5			s4			9	3	
7	s5			s4				10	
8		s6			s11				
9	r1	s7		r1	r1				
10	r3	r3		r3	r3				
11	r5	r5		r5	r5				