

# Análisis sintáctico

Pedro O. Pérez M., PhD

Diseño de compiladores  
Tecnológico de Monterrey

*pperezm@tec.mx*

07-2021

## 4.1 - 4.3

4.1 Introducción

4.2 Gramáticas libres de contexto

4.3 Escribiendo una gramática

## 4.4 - 4.6

4.4 Analisis Top-Down

4.5 Análisis Bottom-Up

4.6 Analizadores LR

## 4.7 - 4.9

4.7 Analizadores LR más poderosos

4.9 Generadores de analizadores sintácticos

## 4.1 Introducción

La sintaxis de las construcciones de un lenguaje de programa deben ser especificados a través de una gramática libre de contexto o en la notación BNF (Backus-Naur Form). Pero, ¿cuáles son sus beneficios?

- ▶ Una gramática da una especificación precisa y fácil de entender de un lenguaje de programación.
- ▶ Para ciertas gramáticas, nosotros podemos construir automáticamente un analizador sintáctico eficiente; con el beneficio adicional de que el proceso de construcción nos puede revelar ambigüedades sintácticas.
- ▶ La estructura de una gramática es útil para una correcta traducción de código fuente a código objeto y para la detección de errores.
- ▶ La gramática permite que un lenguaje se desarrolle o evolucione iterativamente.

## 4.1.1 El rol del analizador sintáctico

Existen tres tipos generales de analizadores sintácticos:

- ▶ Universal (Algoritmos Cocke-Younger-Kasami y Earley).
- ▶ Top-down (LL).
- ▶ Bottom-up (LR).

En cualquier caso, la entrada del analizador sintáctico es revisada de izquierda a derecha, un símbolo a la vez.

## 4.1.2 Gramáticas

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow ( E ) \mid \mathbf{id} \end{aligned}$$

Ejemplo de gramática LR

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

Ejemplo de gramática LL

## 4.1.3 Manejo de errores sintácticos

Los errores de programación comúnmente pueden ocurrir en los siguientes niveles:

- ▶ Errores léxicos.
- ▶ Errores sintácticos.
- ▶ Errores semánticos.
- ▶ Errores lógicos.

El manejador de errores de un analizador sintáctico tiene metas sencillas de describir pero difíciles de realizar:

- ▶ Reportar la presencia de errores de forma clara y acertada.
- ▶ Recuperarse de los errores rápidamente y así poder detectar los siguientes errores.
- ▶ Agregar la mínima sobrecarga al procesamiento de programas correctos.

## 4.1.4 Estrategias de recuperación de errores

Existen 4 estrategias principales para la recuperación de errores:

- ▶ Modo pánico: Cuando se descubre un error, el analizador sintáctico descarta todos los tokens de entrada hasta encontrar un *token de sincronización*.
- ▶ Nivel de frase: Reemplazar un prefijo de la entrada restante por alguna cadena que permita al analizador continuar.
- ▶ Producciones erróneas.
- ▶ Corrección global.



## 4.2.1 Definición formar de una gramática libre de contexto

Un gramática libre de contexto consiste de:

- ▶ **Terminales.** Son los símbolos básicos de cuales están formados las cadenas (strings).
- ▶ **No terminales.** Son las variables sintácticas que hacen referencia a un conjunto de cadenas.
- ▶ **Símbolo inicial.**
- ▶ **Producciones.** Una producción consiste en;
  - ▶ Un no terminal llamado *cabeza* o *lado izquierdo*.
  - ▶ El símbolo  $\rightarrow$ .
  - ▶ Un *cuerpo* o *lado derecho* que consiste en cero o más terminales y no terminales.

$$\begin{aligned} \text{expression} &\rightarrow \text{expression} + \text{term} \\ \text{expression} &\rightarrow \text{expression} - \text{term} \\ \text{expression} &\rightarrow \text{term} \\ \text{term} &\rightarrow \text{term} * \text{factor} \\ \text{term} &\rightarrow \text{term} / \text{factor} \\ \text{term} &\rightarrow \text{factor} \\ \text{factor} &\rightarrow ( \text{expression} ) \\ \text{factor} &\rightarrow \text{id} \end{aligned}$$

Figure 4.2: Grammar for simple arithmetic expressions

## 4.2.2 Notación

1. These symbols are terminals:

- (a) Lowercase letters early in the alphabet, such as *a*, *b*, *c*.
- (b) Operator symbols such as  $+$ ,  $*$ , and so on.
- (c) Punctuation symbols such as parentheses, comma, and so on.
- (d) The digits  $0, 1, \dots, 9$ .
- (e) Boldface strings such as **id** or **if**, each of which represents a single terminal symbol.

2. These symbols are nonterminals:

- (a) Uppercase letters early in the alphabet, such as  $A$ ,  $B$ ,  $C$ .
- (b) The letter  $S$ , which, when it appears, is usually the start symbol.
- (c) Lowercase, italic names such as *expr* or *stmt*.
- (d) When discussing programming constructs, uppercase letters may be used to represent nonterminals for the constructs. For example, nonterminals for expressions, terms, and factors are often represented by  $E$ ,  $T$ , and  $F$ , respectively.

3. Uppercase letters late in the alphabet, such as  $X$ ,  $Y$ ,  $Z$ , represent *grammar symbols*; that is, either nonterminals or terminals.
4. Lowercase letters late in the alphabet, chiefly  $u, v, \dots, z$ , represent (possibly empty) strings of terminals.
5. Lowercase Greek letters,  $\alpha, \beta, \gamma$  for example, represent (possibly empty) strings of grammar symbols. Thus, a generic production can be written as  $A \rightarrow \alpha$ , where  $A$  is the head and  $\alpha$  the body.
6. A set of productions  $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k$  with a common head  $A$  (call them *A-productions*), may be written  $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$ . Call  $\alpha_1, \alpha_2, \dots, \alpha_k$  the *alternatives* for  $A$ .
7. Unless stated otherwise, the head of the first production is the start symbol.

## 4.2.3 Derivaciones

To understand how parsers work, we shall consider derivations in which the nonterminal to be replaced at each step is chosen as follows:

1. In *leftmost* derivations, the leftmost nonterminal in each sentential is always chosen. If  $\alpha \Rightarrow \beta$  is a step in which the leftmost nonterminal in  $\alpha$  is replaced, we write  $\alpha \underset{lm}{\Rightarrow} \beta$ .
2. In *rightmost* derivations, the rightmost nonterminal is always chosen; we write  $\alpha \underset{rm}{\Rightarrow} \beta$  in this case.

$$E \rightarrow E + E \mid E * E \mid - E \mid ( E ) \mid \mathbf{id} \\ -(\mathbf{id} + \mathbf{id})$$

## 4.2.4 Árboles de análisis y derivaciones

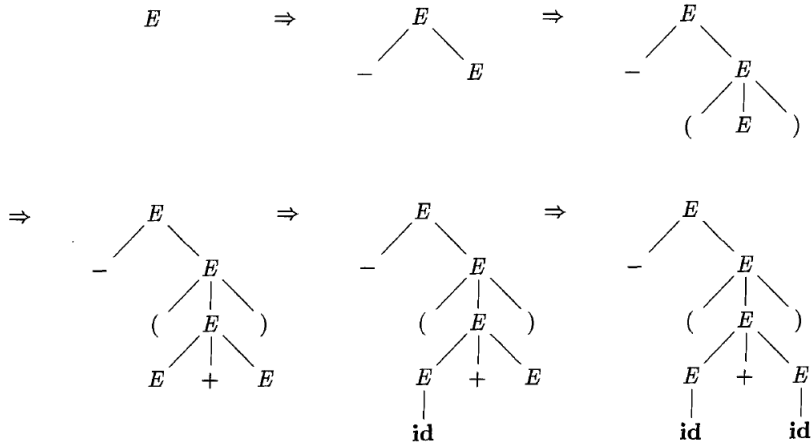


Figure 4.4: Sequence of parse trees for derivation (4.8)



## 4.2.5 Ambigüedad

Si una gramática produce más de un árbol de análisis para una misma entrada, entonces es una gramática ambigua.

## 4.2.7 Gramáticas libres de contexto vs. expresiones regulares.

Cada construcción que pueda ser descrita por una expresión regular puede ser descrita con una gramática libre de contexto, pero no lo contrario.

Alternativamente, cada lenguaje regular es un lenguaje libre de contexto, pero no lo contrario.

- ▶ Para cada estado  $i$  de un AFN, crea un no terminal  $A_i$ .
- ▶ Si el estado  $i$  tiene una transición al estado  $j$  con la entrada  $a$ , agregar la producción  $A_i \rightarrow aA_j$ . Si el estado  $i$  va al estado  $j$  con una transición  $\epsilon$ , agrega la producción  $A_i \rightarrow A_j$ .
- ▶ Si  $i$  es un estado de aceptación, agrega  $A_i \rightarrow \epsilon$ .
- ▶ Si  $i$  es el estado inicial, haz  $A_i$  el símbolo inicial de la gramática.

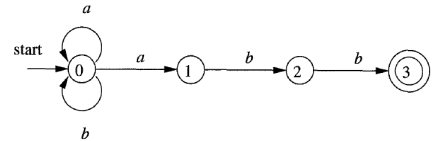


Figure 3.24: A nondeterministic finite automaton

## Ejercicio 4.1

**Exercise 4.2.1:** Consider the context-free grammar:

$$S \rightarrow S S + \mid S S * \mid a$$

and the string  $aa + a*$ .

- a) Give a leftmost derivation for the string.
- b) Give a rightmost derivation for the string.
- c) Give a parse tree for the string.
- ! d) Is the grammar ambiguous or unambiguous? Justify your answer.
- ! e) Describe the language generated by this grammar.

## 4.3.1 Análisis léxico vs sintáctico

Si todo lo que puede ser descrito por una expresión regular también puede ser descrito por una gramática, ¿porqué usamos expresiones regulares para definir el análisis léxico de un lenguaje?

- ▶ Separando la estructura sintáctica de un lenguaje en un componente léxico y otro no léxico provee de una conveniente manera de modularizar el “Front End” en dos componentes manejables.
- ▶ Las reglas léxicas de un lenguaje son frecuentemente muy simples, no se requiere una notación tan poderosa como las gramáticas.
- ▶ Las expresiones regulares generalmente proveen de una notación más concisa y simple para los tokens.
- ▶ Analizadores léxico eficientes pueden ser construidos automáticamente a partir de expresiones regulares que de gramáticas.

## 4.3.2 Eliminando ambigüedad

$stmt \rightarrow$  **if** *expr* **then** *stmt*  
          | **if** *expr* **then** *stmt* **else** *stmt*  
          | **other**

**if**  $E_1$  **then**  $S_1$  **else if**  $E_2$  **then**  $S_2$  **else**  $S_3$

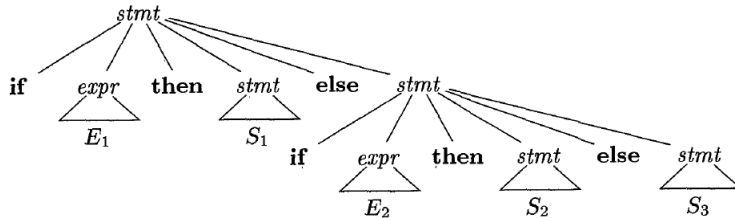


Figure 4.8: Parse tree for a conditional statement

if  $E_1$  then if  $E_2$  then  $S_1$  else  $S_2$

(4.15)

has the two parse trees shown in Fig. 4.9.

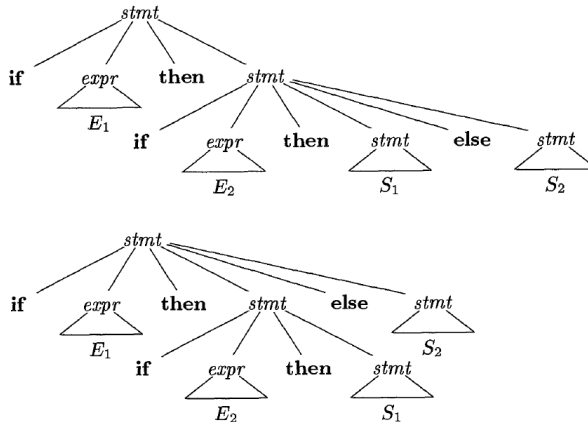


Figure 4.9: Two parse trees for an ambiguous sentence

$$\begin{array}{lll} stmt & \rightarrow & matched\_stmt \\ & | & open\_stmt \\ matched\_stmt & \rightarrow & \text{if } expr \text{ then } matched\_stmt \text{ else } matched\_stmt \\ & | & \text{other} \\ open\_stmt & \rightarrow & \text{if } expr \text{ then } stmt \\ & | & \text{if } expr \text{ then } matched\_stmt \text{ else } open\_stmt \end{array}$$

Figure 4.10: Unambiguous grammar for if-then-else statements



### 4.3.3 Eliminando la recursión por la izquierda

$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' \mid \epsilon \end{aligned}$$

$$\begin{aligned} E &\rightarrow E + T \mid E - T \mid T \\ T &\rightarrow T * F \mid T / F \mid F \\ F &\rightarrow ( E ) \mid \mathbf{id} \end{aligned}$$

**Algorithm 4.19:** Eliminating left recursion.

**INPUT:** Grammar  $G$  with no cycles or  $\epsilon$ -productions.

**OUTPUT:** An equivalent grammar with no left recursion.

**METHOD:** Apply the algorithm in Fig. 4.11 to  $G$ . Note that the resulting non-left-recursive grammar may have  $\epsilon$ -productions.  $\square$

$$\begin{array}{l} S \rightarrow A a \mid b \\ A \rightarrow A c \mid S d \mid \epsilon \end{array}$$

- 1) arrange the nonterminals in some order  $A_1, A_2, \dots, A_n$ .
- 2) **for** ( each  $i$  from 1 to  $n$  ) {
- 3)     **for** ( each  $j$  from 1 to  $i - 1$  ) {
- 4)         replace each production of the form  $A_i \rightarrow A_j \gamma$  by the  
              productions  $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$ , where  
               $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  are all current  $A_j$ -productions
- 5)     }
- 6)     eliminate the immediate left recursion among the  $A_i$ -productions
- 7) }

Figure 4.11: Algorithm to eliminate left recursion from a grammar

## 4.3.4 Factorización izquierda

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \cdots \mid \alpha\beta_n \mid \gamma;$$

↓

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

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

$$\begin{aligned} S &\rightarrow i E t S \mid i E t S e S \mid a \\ E &\rightarrow b \end{aligned}$$

## Ejercicio 4.2

**Exercise 4.3.1:** The following is a grammar for regular expressions over symbols  $a$  and  $b$  only, using  $+$  in place of  $|$  for union, to avoid conflict with the use of vertical bar as a metasymbol in grammars:

$$S \rightarrow S S + \mid S S * \mid a$$

- a) Left factor this grammar.
- b) Does left factoring make the grammar suitable for top-down parsing?
- c) In addition to left factoring, eliminate left recursion from the original grammar.
- d) Is the resulting grammar suitable for top-down parsing?

**Exercise 4.3.1:** The following is a grammar for regular expressions over symbols  $a$  and  $b$  only, using  $+$  in place of  $|$  for union, to avoid conflict with the use of vertical bar as a metasyMBOL in grammars:

a)  $S \rightarrow 0 S 1 \mid 0 1$  with string 000111.

- a) Left factor this grammar.
- b) Does left factoring make the grammar suitable for top-down parsing?
- c) In addition to left factoring, eliminate left recursion from the original grammar.
- d) Is the resulting grammar suitable for top-down parsing?

**Exercise 4.3.1:** The following is a grammar for regular expressions over symbols  $a$  and  $b$  only, using  $+$  in place of  $|$  for union, to avoid conflict with the use of vertical bar as a metasyMBOL in grammars:

! c)  $S \rightarrow S ( S ) S \mid \epsilon$  with string  $((()())$ .

- a) Left factor this grammar.
- b) Does left factoring make the grammar suitable for top-down parsing?
- c) In addition to left factoring, eliminate left recursion from the original grammar.
- d) Is the resulting grammar suitable for top-down parsing?

## 4.4.1 Análisis de descenso recursivo

Revisar la sección 4.4.1 del libro de  
[AHO]



## 4.4.2 FIRST y FOLLOW

- ▶  $\text{FIRST}(\alpha)$ , donde  $\alpha$  es cualquier cadena de símbolos de la gramática, es el conjunto de terminales que empieza una secuencia de cadenas que derivan de  $\alpha$ .
- ▶  $\text{FOLLOW}(A)$ , para un no terminal  $A$ , es el conjunto de terminales  $a$  que pueden aparecer inmediatamente a la derecha de  $A$  en alguna sentencia.

$$\begin{array}{lcl} E & \rightarrow & T E' \\ E' & \rightarrow & + T E' \mid \epsilon \\ T & \rightarrow & F T' \\ T' & \rightarrow & * F T' \mid \epsilon \\ F & \rightarrow & ( E ) \mid \text{id} \end{array}$$

1. If  $X$  is a terminal, then  $\text{FIRST}(X) = \{X\}$ .
2. If  $X$  is a nonterminal and  $X \rightarrow Y_1 Y_2 \cdots Y_k$  is a production for some  $k \geq 1$ , then place  $a$  in  $\text{FIRST}(X)$  if for some  $i$ ,  $a$  is in  $\text{FIRST}(Y_i)$ , and  $\epsilon$  is in all of  $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$ ; that is,  $Y_1 \cdots Y_{i-1} \xRightarrow{*} \epsilon$ . If  $\epsilon$  is in  $\text{FIRST}(Y_j)$  for all  $j = 1, 2, \dots, k$ , then add  $\epsilon$  to  $\text{FIRST}(X)$ . For example, everything in  $\text{FIRST}(Y_1)$  is surely in  $\text{FIRST}(X)$ . If  $Y_1$  does not derive  $\epsilon$ , then we add nothing more to  $\text{FIRST}(X)$ , but if  $Y_1 \xRightarrow{*} \epsilon$ , then we add  $\text{FIRST}(Y_2)$ , and so on.
3. If  $X \rightarrow \epsilon$  is a production, then add  $\epsilon$  to  $\text{FIRST}(X)$ .

$$\begin{array}{ll} E & \rightarrow T E' \\ E' & \rightarrow + T E' \mid \epsilon \\ T & \rightarrow F T' \\ T' & \rightarrow * F T' \mid \epsilon \\ F & \rightarrow ( E ) \mid \text{id} \end{array}$$

1. Place \$ in FOLLOW( $S$ ), where  $S$  is the start symbol, and \$ is the input right endmarker.
2. If there is a production  $A \rightarrow \alpha B \beta$ , then everything in FIRST( $\beta$ ) except  $\epsilon$  is in FOLLOW( $B$ ).
3. If there is a production  $A \rightarrow \alpha B$ , or a production  $A \rightarrow \alpha B \beta$ , where FIRST( $\beta$ ) contains  $\epsilon$ , then everything in FOLLOW( $A$ ) is in FOLLOW( $B$ ).

## 4.4.23 Gramáticas LL(1)

- ▶ Los analizadores predictivos, o analizadores de descenso recursivo sin “backtracking” pueden ser contruidos para una clase de gramáticas conocidas como LL(1).
- ▶ Una gramática  $G$  es LL(1), sí y solo si siendo  $A \rightarrow \alpha | \beta$  son distintas producciones de  $G$ , se dan las siguientes condiciones:
  - ▶ Para ningún terminal  $a$ , ambas cadenas empiezan con el mismo símbolo.
  - ▶ A lo mas un no terminal,  $\alpha$  o  $\beta$ , pueden derivar en  $\epsilon$ .
  - ▶ Si  $\beta \Rightarrow \epsilon$ , entonces  $\alpha$  no debe derivar ninguna cadena que empiece con una terminal en FOLLOW( $A$ ). De igual forma,  $\alpha \Rightarrow \epsilon$ , entonces  $\beta$  no puede derivar ninguna cadena que empiece con un termina en FOLLOW( $A$ ).

## Ejercicio 4.3

Para la siguiente gramática, genera:

- ▶ Los conjuntos de FIRST y FOLLOW.
- ▶ ¿Es una gramática LL(1)?

$$S \rightarrow A a$$

$$A \rightarrow B D$$

$$B \rightarrow b \mid \epsilon$$

$$D \rightarrow d \mid \epsilon$$

**Algorithm 4.31:** Construction of a predictive parsing table.

**INPUT:** Grammar  $G$ .

**OUTPUT:** Parsing table  $M$ .

**METHOD:** For each production  $A \rightarrow \alpha$  of the grammar, do the following:

1. For each terminal  $a$  in  $\text{FIRST}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, a]$ .
2. If  $\epsilon$  is in  $\text{FIRST}(\alpha)$ , then for each terminal  $b$  in  $\text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$ . If  $\epsilon$  is in  $\text{FIRST}(\alpha)$  and  $\$$  is in  $\text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, \$]$  as well.

## Ejercicio 4.4

Para la siguiente gramática:

- Genera los conjuntos de FIRST y FOLLOW.

$$\begin{aligned} S &\rightarrow iEtSS' \mid a \\ S' &\rightarrow eS \mid \epsilon \\ E &\rightarrow b \end{aligned}$$

NON - TERMINAL	INPUT SYMBOL					
	$a$	$b$	$e$	$i$	$t$	$\$$
$S$	$S \rightarrow a$			$S \rightarrow iEtSS'$		
$S'$			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
$E$		$E \rightarrow b$				

Figure 4.18: Parsing table  $M$  for Example 4.33



## 4.4.4 Análisis Predictivo No Recursivo

**METHOD:** Initially, the parser is in a configuration with  $w\$$  in the input buffer and the start symbol  $S$  of  $G$  on top of the stack, above  $\$$ . The program in Fig. 4.20 uses the predictive parsing table  $M$  to produce a predictive parse for the input.  $\square$

```
set  $ip$  to point to the first symbol of  $w$ ;  
set  $X$  to the top stack symbol;  
while (  $X \neq \$$  ) { /* stack is not empty */  
    if (  $X$  is  $a$  ) pop the stack and advance  $ip$ ;  
    else if (  $X$  is a terminal )  $error()$ ;  
    else if (  $M[X, a]$  is an error entry )  $error()$ ;  
    else if (  $M[X, a] = X \rightarrow Y_1 Y_2 \cdots Y_k$  ) {  
        output the production  $X \rightarrow Y_1 Y_2 \cdots Y_k$ ;  
        pop the stack;  
        push  $Y_k, Y_{k-1}, \dots, Y_1$  onto the stack, with  $Y_1$  on top;  
    }  
    set  $X$  to the top stack symbol;  
}
```

Figure 4.20: Predictive parsing algorithm

NON - TERMINAL	INPUT SYMBOL					
	id	+	*	(	)	\$
$E$	$E \rightarrow TE'$			$E \rightarrow TE'$		
$E'$		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
$T$	$T \rightarrow FT'$			$T \rightarrow FT'$		
$T'$		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
$F$	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

Figure 4.17: Parsing table  $M$  for Example 4.32

id + id \* id

## Ejercicio 4.5

Para la siguiente gramática, genera:

- ▶ La tabla de análisis predictivo. Emplea el algoritmo 4.31.
- ▶ ¿Se reconoce la entrada **bda**? Emplea el algoritmo 4.34.
- ▶ ¿Se reconoce la entrada **dba**? Emplea el algoritmo 4.34.

$$S \rightarrow A a$$

$$A \rightarrow B D$$

$$B \rightarrow b \mid \epsilon$$

$$D \rightarrow d \mid \epsilon$$

## 4.4.5 Recuperación de errores en análisis predictivo

Modo pánico:

- ▶ Agrega todos los símbolos del  $FOLLOW(A)$  en el conjunto de sincronización de  $A$ .
- ▶ Utiliza la estructura jerárquica de las construcciones.
- ▶ Agrega todos los símbolo del  $FIRST(A)$  a los conjuntos de sincronización del no terminal  $A$ .
- ▶ Si un no terminal genera un  $\epsilon$ , entonces la producción que deriva  $\epsilon$  es tomada por omisión.
- ▶ Si hay un terminal en el tope de la pila, se hace pop de ese terminal.

NON - TERMINAL	INPUT SYMBOL					
	id	+	*	(	)	\$
$E$	$E \rightarrow TE'$			$E \rightarrow TE'$	synch	synch
$E'$		$E \rightarrow +TE'$			$E \rightarrow \epsilon$	$E \rightarrow \epsilon$
$T$	$T \rightarrow FT'$	synch		$T \rightarrow FT'$	synch	synch
$T'$		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
$F$	$F \rightarrow \text{id}$	synch	synch	$F \rightarrow (E)$	synch	synch

Figure 4.22: Synchronizing tokens added to the parsing table of Fig. 4.17

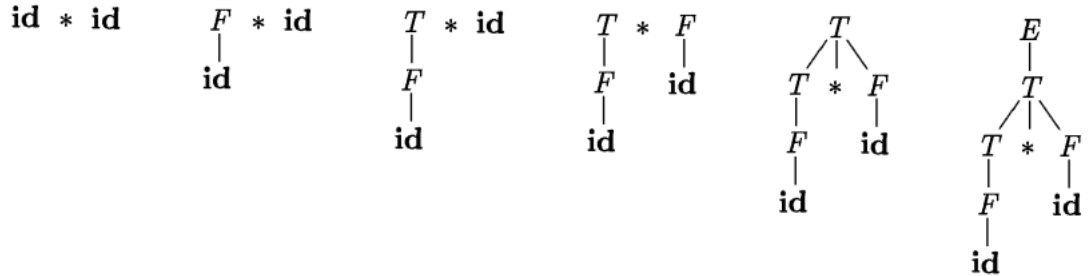


Figure 4.25: A bottom-up parse for  $\text{id} * \text{id}$

## 4.5.1 Reducciones - 4.5.2 Manejo de las reducciones

RIGHT SENTENTIAL FORM	HANDLE	REDUCING PRODUCTION
$\mathbf{id_1 * id_2}$	$\mathbf{id_1}$	$F \rightarrow \mathbf{id}$
$F * \mathbf{id_2}$	$F$	$T \rightarrow F$
$T * \mathbf{id_2}$	$\mathbf{id_2}$	$F \rightarrow \mathbf{id}$
$T * F$	$T * F$	$E \rightarrow T * F$

Figure 4.26: Handles during a parse of  $\mathbf{id_1 * id_2}$

### 4.5.3 Analizador Shift-Reduce

STACK	INPUT	ACTION
\$	<b>id<sub>1</sub></b> * <b>id<sub>2</sub></b> \$	shift
\$ <b>id<sub>1</sub></b>	* <b>id<sub>2</sub></b> \$	reduce by $F \rightarrow \text{id}$
\$ $F$	* <b>id<sub>2</sub></b> \$	reduce by $T \rightarrow F$
\$ $T$	* <b>id<sub>2</sub></b> \$	shift
\$ $T$ *	<b>id<sub>2</sub></b> \$	shift
\$ $T$ * <b>id<sub>2</sub></b>	\$	reduce by $F \rightarrow \text{id}$
\$ $T$ * $F$	\$	reduce by $T \rightarrow T * F$
\$ $T$	\$	reduce by $E \rightarrow T$
\$ $E$	\$	accept

Figure 4.28: Configurations of a shift-reduce parser on input **id<sub>1</sub>\*id<sub>2</sub>**



Mientras que las operaciones principales son shift y reduce, en realidad hay 4 posibles acciones que debe realizar un analizador shift-reduce:

- ▶ Shift. Mueve el siguiente símbolo de la entrada al tope de la pila.
- ▶ Reduce. El lado derecho de la cadena a ser reducida debe estar en el tope de la pila. Localizar el lado izquierdo de la cadena dentro del stack y decidir con qué no terminar se reemplaza la cadena.
- ▶ Aceptar. Anuncia una terminación exitosa del análisis.
- ▶ Error. Descubre un error de sintaxis y llama a una rutina de recuperación de error.

## 4.5.4 Conflictos durante el análisis shift-reduce

Existen gramáticas libres de contexto para las cuales este tipo de analizadores no puede ser usado. Cada analizador shift-reduce puede alcanzar una configuración en la que el analizador, conociendo el contenido de la pila y el siguiente símbolo de entrada no puede decidir si desea reducir o cambiar (conflicto shift-reduce), o no puede decidir cual reducir aplicar (conflicto reduce-shift).

# Introducción a los analizadores LR

- ▶ El tipo más prevalente de analizadores bottom-up son llamados LR( $k$ ).
- ▶ Entre los analizadores de este tipo, revisaremos los Simple LR (SLR), los canonical-LR y los LALR. Estos últimos son los más empleados en las herramientas de generadores de analizadores.

## 4.6.1 ¿Por qué analizadores LR?

- ▶ Los analizadores LR pueden ser contruidos para reconocer virtualmente todos los lenguajes de programación para los cuales una gramática libre de contexto puede ser escrita.
- ▶ El método de análisis LR es el método de análisis shift-reduce sin retroceso más general conocido.
- ▶ Un analizador LR puede detectar un error sintáctico tan pronto como sea posible hacerlo en un escaneo de izquierda a derecha de la entrada.
- ▶ Las gramáticas que puede ser analizadas usando métodos LR son es un superconjunto propio de la clase de gramáticas que pueden ser analizadas con un método predictivo o LL.

Pero... construir un analizador LR para una gramática típica de un lenguaje de programación requiere de mucho trabajo para ser construido a mano.  
Afortunadamente, existen muchas herramientas para eso.

## 4.6.2 Item y el autómata LR

Un analizador LR realiza las decisiones shift-reduce manteniendo estados para determinar en dónde estamos en un análisis. Un *item* $LR(0)$  de una gramática  $G$  es una producción de  $G$  con un punto en alguna posición del cuerpo de la producción. Por ejemplo, la producción  $A \rightarrow XYZ$  tiene cuatro items:

$$A \rightarrow \cdot XYZ$$

$$A \rightarrow X \cdot YZ$$

$$A \rightarrow XY \cdot Z$$

$$A \rightarrow XYZ \cdot$$

Una colección de conjuntos de items  $LR(0)$ , llamados la colección *canónica*  $LR(0)$ , provee la base para la construcción de un autómata finito determinista que es usado para realizar decisiones de análisis. Este autómata es llamado *autómata  $LR(0)$* .

# CLOSURE(I)

```
SetOfItems CLOSURE(I) {  
    J = I;  
    repeat  
        for ( each item  $A \rightarrow \alpha \cdot B \beta$  in J )  
            for ( each production  $B \rightarrow \gamma$  of G )  
                if (  $B \rightarrow \cdot \gamma$  is not in J )  
                    add  $B \rightarrow \cdot \gamma$  to J;  
    until no more items are added to J on one round;  
    return J;  
}
```

$$\begin{array}{lcl} E & \rightarrow & E + T \mid T \\ T & \rightarrow & T * F \mid F \\ E & \rightarrow & (E) \mid \text{id} \end{array}$$

Figure 4.32: Computation of CLOSURE



# GOTO( $I, X$ )

```
void items( $G'$ ) {  
     $C = \text{CLOSURE}(\{[S' \rightarrow \cdot S]\})$ ;  
    repeat  
        for ( each set of items  $I$  in  $C$  )  
            for ( each grammar symbol  $X$  )  
                if ( GOTO( $I, X$ ) is not empty and not in  $C$  )  
                    add GOTO( $I, X$ ) to  $C$ ;  
    until no new sets of items are added to  $C$  on a round;  
}
```

Figure 4.33: Computation of the canonical collection of sets of LR(0) items

# Autómata LR(0)

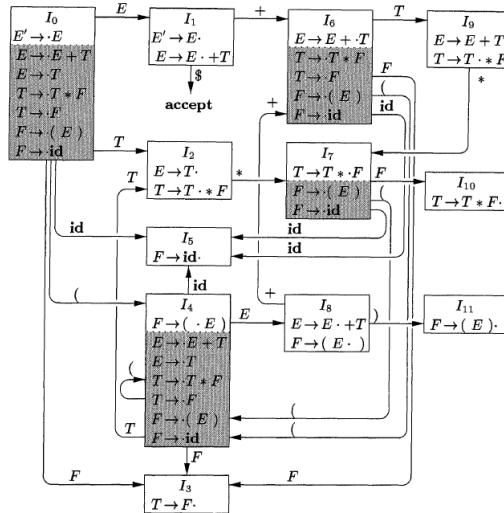


Figure 4.31: LR(0) automaton for the expression grammar (4.1)

## Uso del autómata LR(0)

Si la cadena  $S$  de símbolos gramaticales hace que el autómata LR(0) de un estado inicial  $i$  a otro estado  $j$ , entonces cambiamos el siguiente símbolo  $a$  de la entrada. De lo contrario, elegimos reducir, el item en el estado  $j$  nos dirá que producción usar.

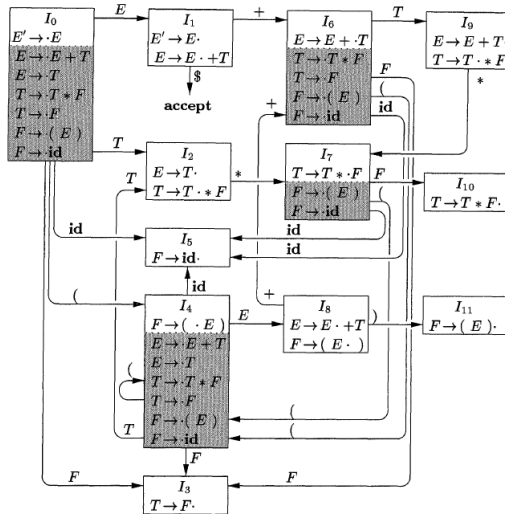


Figure 4.31: LR(0) automaton for the expression grammar (4.1)

id \* id

## 4.6.3 Algoritmo de análisis LR

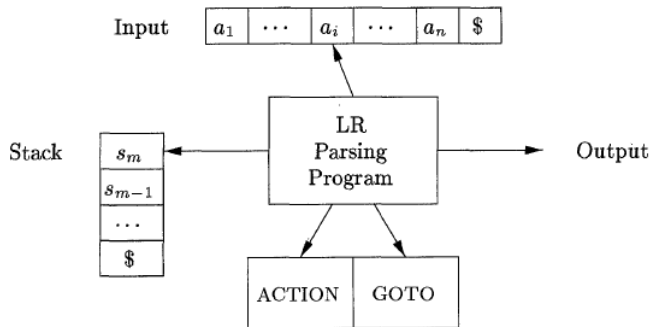


Figure 4.35: Model of an LR parser

**Algorithm 4.46:** Constructing an SLR-parsing table.

**INPUT:** An augmented grammar  $G'$ .

**OUTPUT:** The SLR-parsing table functions ACTION and GOTO for  $G'$ .

**METHOD:**

1. Construct  $C = \{I_0, I_1, \dots, I_n\}$ , the collection of sets of LR(0) items for  $G'$ .
2. State  $i$  is constructed from  $I_i$ . The parsing actions for state  $i$  are determined as follows:
  - (a) If  $[A \rightarrow \alpha \cdot a \beta]$  is in  $I_i$  and  $\text{GOTO}(I_i, a) = I_j$ , then set  $\text{ACTION}[i, a]$  to “shift  $j$ .” Here  $a$  must be a terminal.
  - (b) If  $[A \rightarrow \alpha \cdot]$  is in  $I_i$ , then set  $\text{ACTION}[i, a]$  to “reduce  $A \rightarrow \alpha$ ” for all  $a$  in  $\text{FOLLOW}(A)$ ; here  $A$  may not be  $S'$ .
  - (c) If  $[S' \rightarrow \cdot S]$  is in  $I_i$ , then set  $\text{ACTION}[i, \$]$  to “accept.”

If any conflicting actions result from the above rules, we say the grammar is not SLR(1). The algorithm fails to produce a parser in this case.

3. The goto transitions for state  $i$  are constructed for all nonterminals  $A$  using the rule: If  $\text{GOTO}(I_i, A) = I_j$ , then  $\text{GOTO}[i, A] = j$ .
4. All entries not defined by rules (2) and (3) are made “error.”
5. The initial state of the parser is the one constructed from the set of items containing  $[S' \rightarrow \cdot S]$ .

Si durante la construcción de la tabla existe algún conflicto con las reglas anteriores, entonces la gramática no es SLR.

STATE	ACTION						GOTO		
	id	+	*	(	)	\$	<i>E</i>	<i>T</i>	<i>F</i>
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			

Figure 4.37: Parsing table for expression grammar

Generar el conjunto de elementos SLR, la función GOTO y la tabla de análisis de:

$$\begin{array}{lcl} S & \rightarrow & L = R \mid R \\ L & \rightarrow & *R \mid \mathbf{id} \\ R & \rightarrow & L \end{array}$$



## Ejercicio 4.6

Para la siguiente gramática:

$$\begin{aligned} S &\rightarrow SS* \\ S &\rightarrow SS+ \\ S &\rightarrow a \end{aligned}$$

- ▶ Construye el conjunto de elementos SLR.
- ▶ Calcula las funciones *GOTO*.
- ▶ ¿Es una gramática SLR?

**Algorithm 4.44:** LR-parsing algorithm.

**INPUT:** An input string  $w$  and an LR-parsing table with functions ACTION and GOTO for a grammar  $G$ .

**OUTPUT:** If  $w$  is in  $L(G)$ , the reduction steps of a bottom-up parse for  $w$ ; otherwise, an error indication.

**METHOD:** Initially, the parser has  $s_0$  on its stack, where  $s_0$  is the initial state, and  $w\$$  in the input buffer. The parser then executes the program in Fig. 4.36.

□

```
let  $a$  be the first symbol of  $w\$$ ;  
while(1) { /* repeat forever */  
    let  $s$  be the state on top of the stack;  
    if ( ACTION[ $s, a$ ] = shift  $t$  ) {  
        push  $t$  onto the stack;  
        let  $a$  be the next input symbol;  
    } else if ( ACTION[ $s, a$ ] = reduce  $A \rightarrow \beta$  ) {  
        pop  $|\beta|$  symbols off the stack;  
        let state  $t$  now be on top of the stack;  
        push GOTO[ $t, A$ ] onto the stack;  
        output the production  $A \rightarrow \beta$ ;  
    } else if ( ACTION[ $s, a$ ] = accept ) break; /* parsing is done */  
    else call error-recovery routine;  
}
```

Figure 4.36: LR-parsing program

STATE	ACTION						GOTO		
	id	+	*	(	)	\$	<i>E</i>	<i>T</i>	<i>F</i>
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			

Figure 4.37: Parsing table for expression grammar

$$\text{id} * \text{id} + \text{id}$$

## Ejercicio 4.7

Usando la tabla generada en el ejercicio anterior, muestra las acciones de análisis para la entrada  $aa * a +$ .

Cómo se mencionó antes, existen dos métodos adicionales de análisis LR:

- ▶ “canonical-LR” (LR), el cual hace uso de los símbolos “lookahead”. Este método usa un conjunto más grande de *items* llamado *items LR(1)*.
- ▶ “lookahead-LR” (LALR), el cual está basado en el conjunto de items LR(0), pero tiene menos estados que los de un analizador basado en items LR(1). Puede manejar más gramáticas que un SLR.

## 4.7.1 Elementos LR(1) canónicos

En el método SLR, cuando el estado  $i$  realizar una reducción por la producción  $A \rightarrow \alpha$ , si el conjunto de items  $I_i$  contiene  $[A \rightarrow \alpha.]$  y  $a$  está en el  $FOLLOW(A)$ . En algunas situaciones, sin embargo, cuando un estado  $i$  aparece en el tope de la pila, el prefijo viable  $\beta\alpha$  está en la pila es tal que  $\beta\alpha$  no puede ser seguido por  $a$  en cualquier producción del lado derecho.

¿Recuerdas que sucedió con esta gramática?

$$\begin{array}{lcl} S & \rightarrow & L = R \mid R \\ L & \rightarrow & *R \mid \mathbf{id} \\ R & \rightarrow & L \end{array}$$

- ▶ Es posible agregar más información al estado que nos permita determinar cuál de las reducciones aplicar a  $A \rightarrow \alpha$ . Esto se logra dividiendo los estados cuando sea necesario, haciendo que cada estado indique exactamente cuáles símbolos de entrada pueden seguir a  $\alpha$  para el cuál hay una reducción de  $A$ .
- ▶ La información extra es incorporada dentro del estado redefiniendo los elementos para que tengan un segundo componente  $[A \rightarrow \alpha.\beta, a]$ , donde  $A \rightarrow \alpha\beta$  es una producción y  $a$  es un terminal o el marcado \$. Llamamos a esto, un elemento  $LR(1)$ .

## 4.7.2 Construcción del conjunto de elementos LR(1)

```
SetOfItems CLOSURE( $I$ ) {  
    repeat  
        for ( each item  $[A \rightarrow \alpha \cdot B \beta, a]$  in  $I$  )  
            for ( each production  $B \rightarrow \gamma$  in  $G'$  )  
                for ( each terminal  $b$  in FIRST( $\beta a$ ) )  
                    add  $[B \rightarrow \cdot \gamma, b]$  to set  $I$ ;  
    until no more items are added to  $I$ ;  
    return  $I$ ;  
}  
  
SetOfItems GOTO( $I, X$ ) {  
    initialize  $J$  to be the empty set;  
    for ( each item  $[A \rightarrow \alpha \cdot X \beta, a]$  in  $I$  )  
        add item  $[A \rightarrow \alpha X \cdot \beta, a]$  to set  $J$ ;  
    return CLOSURE( $J$ );  
}
```



```
void items( $G'$ ) {  
    initialize  $C$  to CLOSURE( $\{[S' \rightarrow \cdot S, \$]\}$ );  
    repeat  
        for ( each set of items  $I$  in  $C$  )  
            for ( each grammar symbol  $X$  )  
                if ( GOTO( $I, X$ ) is not empty and not in  $C$  )  
                    add GOTO( $I, X$ ) to  $C$ ;  
    until no new sets of items are added to  $C$ ;  
}
```

Construir el conjunto de elementos LR(1) y la función GOTO para la siguiente gramática:

$$E \rightarrow E + T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow F$$

$$F \rightarrow (E) + T$$

$$F \rightarrow id$$

## 4.7.3 Construcción de tablas de análisis LR(1)

**Algorithm 4.56:** Construction of canonical-LR parsing tables.

**INPUT:** An augmented grammar  $G'$ .

**OUTPUT:** The canonical-LR parsing table functions ACTION and GOTO for  $G'$ .

**METHOD:**

1. Construct  $C' = \{I_0, I_1, \dots, I_n\}$ , the collection of sets of LR(1) items for  $G'$ .
2. State  $i$  of the parser is constructed from  $I_i$ . The parsing action for state  $i$  is determined as follows.
  - (a) If  $[A \rightarrow \alpha \cdot a\beta, b]$  is in  $I_i$  and  $\text{GOTO}(I_i, a) = I_j$ , then set  $\text{ACTION}[i, a]$  to "shift  $j$ ." Here  $a$  must be a terminal.
  - (b) If  $[A \rightarrow \alpha \cdot, a]$  is in  $I_i$ ,  $A \neq S'$ , then set  $\text{ACTION}[i, a]$  to "reduce  $A \rightarrow \alpha$ ."
  - (c) If  $[S' \rightarrow S \cdot, \$]$  is in  $I_i$ , then set  $\text{ACTION}[i, \$]$  to "accept."

If any conflicting actions result from the above rules, we say the grammar is not LR(1). The algorithm fails to produce a parser in this case.

3. The goto transitions for state  $i$  are constructed for all nonterminals  $A$  using the rule: If  $\text{GOTO}(I_i, A) = I_j$ , then  $\text{GOTO}[i, A] = j$ .
4. All entries not defined by rules (2) and (3) are made "error."
5. The initial state of the parser is the one constructed from the set of items containing  $[S' \rightarrow \cdot S, \$]$ .

## Ejercicio 4.8

Construir el conjunto de elementos LR(1) y la función GOTO para la siguiente gramática:

$$E \rightarrow EE+$$

$$E \rightarrow EE*$$

$$E \rightarrow a$$

Construir la tabla de análisis.

## 4.7.4 Construcción de tablas de análisis LALR

**Algorithm 4.59:** An easy, but space-consuming LALR table construction.

**INPUT:** An augmented grammar  $G'$ .

**OUTPUT:** The LALR parsing-table functions ACTION and GOTO for  $G'$ .

**METHOD:**

1. Construct  $C = \{I_0, I_1, \dots, I_n\}$ , the collection of sets of LR(1) items.
2. For each core present among the set of LR(1) items, find all sets having that core, and replace these sets by their union.
3. Let  $C' = \{J_0, J_1, \dots, J_m\}$  be the resulting sets of LR(1) items. The parsing actions for state  $i$  are constructed from  $J_i$  in the same manner as in Algorithm 4.56. If there is a parsing action conflict, the algorithm fails to produce a parser, and the grammar is said not to be LALR(1).
4. The GOTO table is constructed as follows. If  $J$  is the union of one or more sets of LR(1) items, that is,  $J = I_1 \cap I_2 \cap \dots \cap I_k$ , then the cores of  $\text{GOTO}(I_1, X)$ ,  $\text{GOTO}(I_2, X)$ ,  $\dots$ ,  $\text{GOTO}(I_k, X)$  are the same, since  $I_1, I_2, \dots, I_k$  all have the same core. Let  $K$  be the union of all sets of items having the same core as  $\text{GOTO}(I_1, X)$ . Then  $\text{GOTO}(J, X) = K$ .

# Construcción eficiente de tablas LALR

Existen varias modificaciones que podemos realizar al algoritmo anterior y así evitar la construcción completa de un conjunto LR(1):

- ▶ Primero, podemos representar cualquier conjunto LR(0) o LR(1) por su elemento inicial ( $[S' \rightarrow S]$  or  $[S' \rightarrow S, \$]$ )
- ▶ Construimos el conjunto de kernels LALR(1) a partir del LR(0) usando un algoritmo de propagación y generación espontánea de lookaheads.
- ▶ Ya que tenemos los kernels LALR(1), podemos generar una tabla de análisis empleado el algoritmo que se presenta a continuación.

**METHOD:** The algorithm is given in Fig. 4.45.  $\square$

```
for ( each item  $A \rightarrow \alpha \cdot \beta$  in  $K$  ) {  
     $J := \text{CLOSURE}(\{[A \rightarrow \alpha \cdot \beta, \#]\})$ ;  
    if (  $[B \rightarrow \gamma \cdot X \delta, a]$  is in  $J$ , and  $a$  is not  $\#$  )  
        conclude that lookahead  $a$  is generated spontaneously for item  
             $B \rightarrow \gamma X \cdot \delta$  in  $\text{GOTO}(I, X)$ ;  
    if (  $[B \rightarrow \gamma \cdot X \delta, \#]$  is in  $J$  )  
        conclude that lookaheads propagate from  $A \rightarrow \alpha \cdot \beta$  in  $I$  to  
             $B \rightarrow \gamma X \cdot \delta$  in  $\text{GOTO}(I, X)$ ;  
}
```

```
SetOfItems CLOSURE( $I$ ) {  
    repeat  
        for ( each item  $[A \rightarrow \alpha \cdot B \beta, a]$  in  $I$  )  
            for ( each production  $B \rightarrow \gamma$  in  $G'$  )  
                for ( each terminal  $b$  in FIRST( $\beta a$ ) )  
                    add  $[B \rightarrow \cdot \gamma, b]$  to set  $I$ ;  
    until no more items are added to  $I$ ;  
    return  $I$ ;  
}
```



$$\begin{array}{lcl} S' & \rightarrow & S \\ S & \rightarrow & L = R \mid R \\ L & \rightarrow & *R \mid \text{id} \\ R & \rightarrow & L \end{array}$$
$$I_0: S' \rightarrow \cdot S$$
$$I_1: S' \rightarrow S \cdot$$
$$I_2: S \rightarrow L \cdot = R \\ R \rightarrow L \cdot$$
$$I_3: S \rightarrow R \cdot$$
$$I_4: L \rightarrow * \cdot R$$
$$I_5: L \rightarrow \text{id} \cdot$$
$$I_6: S \rightarrow L = \cdot R$$
$$I_7: L \rightarrow * R \cdot$$
$$I_8: R \rightarrow L \cdot$$
$$I_9: S \rightarrow L = R \cdot$$

**METHOD:**

1. Construct the kernels of the sets of LR(0) items for  $G$ . If space is not at a premium, the simplest way is to construct the LR(0) sets of items, as in Section 4.6.2, and then remove the nonkernel items. If space is severely constrained, we may wish instead to store only the kernel items for each set, and compute GOTO for a set of items  $I$  by first computing the closure of  $I$ .
2. Apply Algorithm 4.62 to the kernel of each set of LR(0) items and grammar symbol  $X$  to determine which lookaheads are spontaneously generated for kernel items in  $\text{GOTO}(I, X)$ , and from which items in  $I$  lookaheads are propagated to kernel items in  $\text{GOTO}(I, X)$ .
3. Initialize a table that gives, for each kernel item in each set of items, the associated lookaheads. Initially, each item has associated with it only those lookaheads that we determined in step (2) were generated spontaneously.
4. Make repeated passes over the kernel items in all sets. When we visit an item  $i$ , we look up the kernel items to which  $i$  propagates its lookaheads, using information tabulated in step (2). The current set of lookaheads for  $i$  is added to those already associated with each of the items to which  $i$  propagates its lookaheads. We continue making passes over the kernel items until no more new lookaheads are propagated.

FROM	TO
$I_0: S' \rightarrow \cdot S$	$I_1: S' \rightarrow S \cdot$ $I_2: S \rightarrow L \cdot = R$ $I_2: R \rightarrow L \cdot$ $I_3: S \rightarrow R \cdot$ $I_4: L \rightarrow * \cdot R$ $I_5: L \rightarrow \text{id} \cdot$
$I_2: S \rightarrow L \cdot = R$	$I_6: S \rightarrow L = \cdot R$
$I_4: L \rightarrow * \cdot R$	$I_4: L \rightarrow * \cdot R$ $I_5: L \rightarrow \text{id} \cdot$ $I_7: L \rightarrow * R \cdot$ $I_8: R \rightarrow L \cdot$
$I_6: S \rightarrow L = \cdot R$	$I_4: L \rightarrow * \cdot R$ $I_5: L \rightarrow \text{id} \cdot$ $I_8: R \rightarrow L \cdot$ $I_9: S \rightarrow L = R \cdot$

Figure 4.46: Propagation of lookaheads

SET	ITEM	LOOKAHEADS			
		INIT	PASS 1	PASS 2	PASS 3
$I_0$ :	$S' \rightarrow \cdot S$	\$	\$	\$	\$
$I_1$ :	$S' \rightarrow S \cdot$		\$	\$	\$
$I_2$ :	$S \rightarrow L \cdot = R$		\$	\$	\$
	$R \rightarrow L \cdot$		\$	\$	\$
$I_3$ :	$S \rightarrow R \cdot$		\$	\$	\$
$I_4$ :	$L \rightarrow * \cdot R$	=	=/\$	=/\$	=/\$
$I_5$ :	$L \rightarrow \text{id} \cdot$	=	=/\$	=/\$	=/\$
$I_6$ :	$S \rightarrow L = \cdot R$			\$	\$
$I_7$ :	$L \rightarrow * R \cdot$		=	=/\$	=/\$
$I_8$ :	$R \rightarrow L \cdot$		=	=/\$	=/\$
$I_9$ :	$S \rightarrow L = R \cdot$				\$

Figura 4.47. Construcción de los ítems LR(0) para el ejemplo 4.4.1

## Ejercicio 4.9

Empleando el algoritmo anterior, construye el conjunto de elementos LALR para la siguiente gramática:

$$S \rightarrow SS+$$

$$S \rightarrow SS*$$

$$S \rightarrow a$$

## Ejercicio 4.10

Responde las siguientes preguntas:

- ▶ ¿Cuándo una gramática es LALR(1) pero no SLR(1)?
- ▶ ¿Cuándo una gramática es LR(1) pero no LALR(1)?

## 4.9 Generadores de analizadores sintácticos

Revisar la sección 4.9 del libro de [AHO]