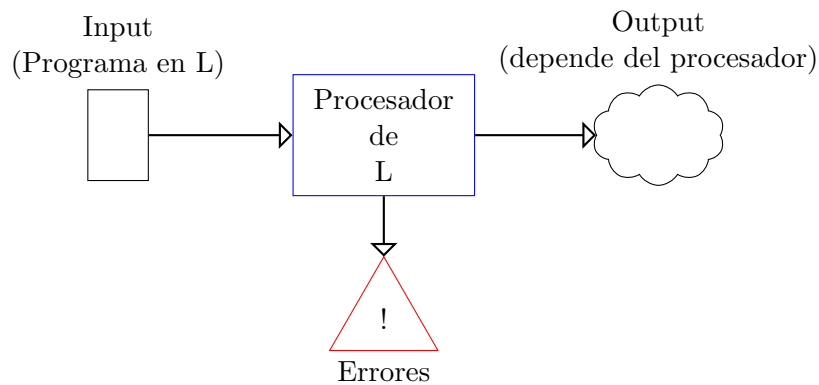


Procesadores de Lenguajes

David Antuña Rodríguez

Contenidos

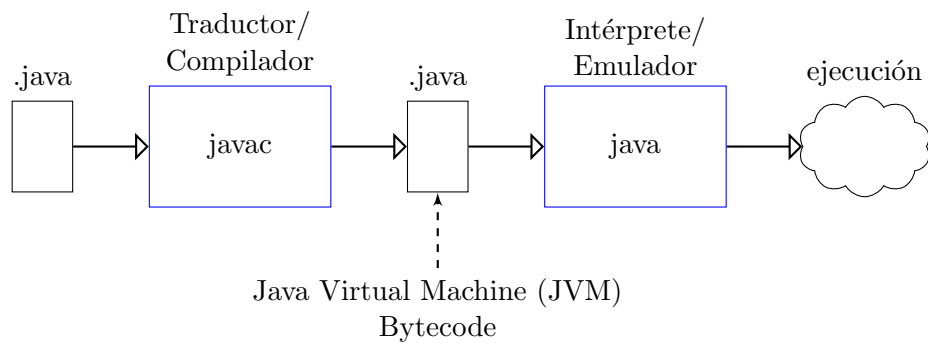
1 Introducción



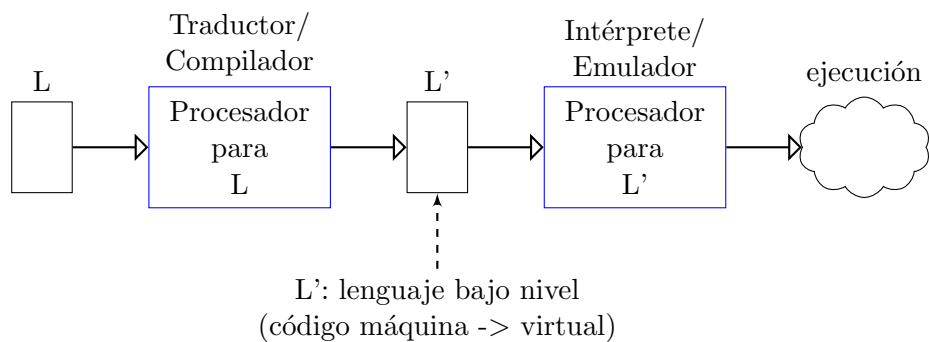
Compilador: Genera una representación en otro L (Alto nivel -> Bajo nivel).

Intérprete: Ejecuta una representación.

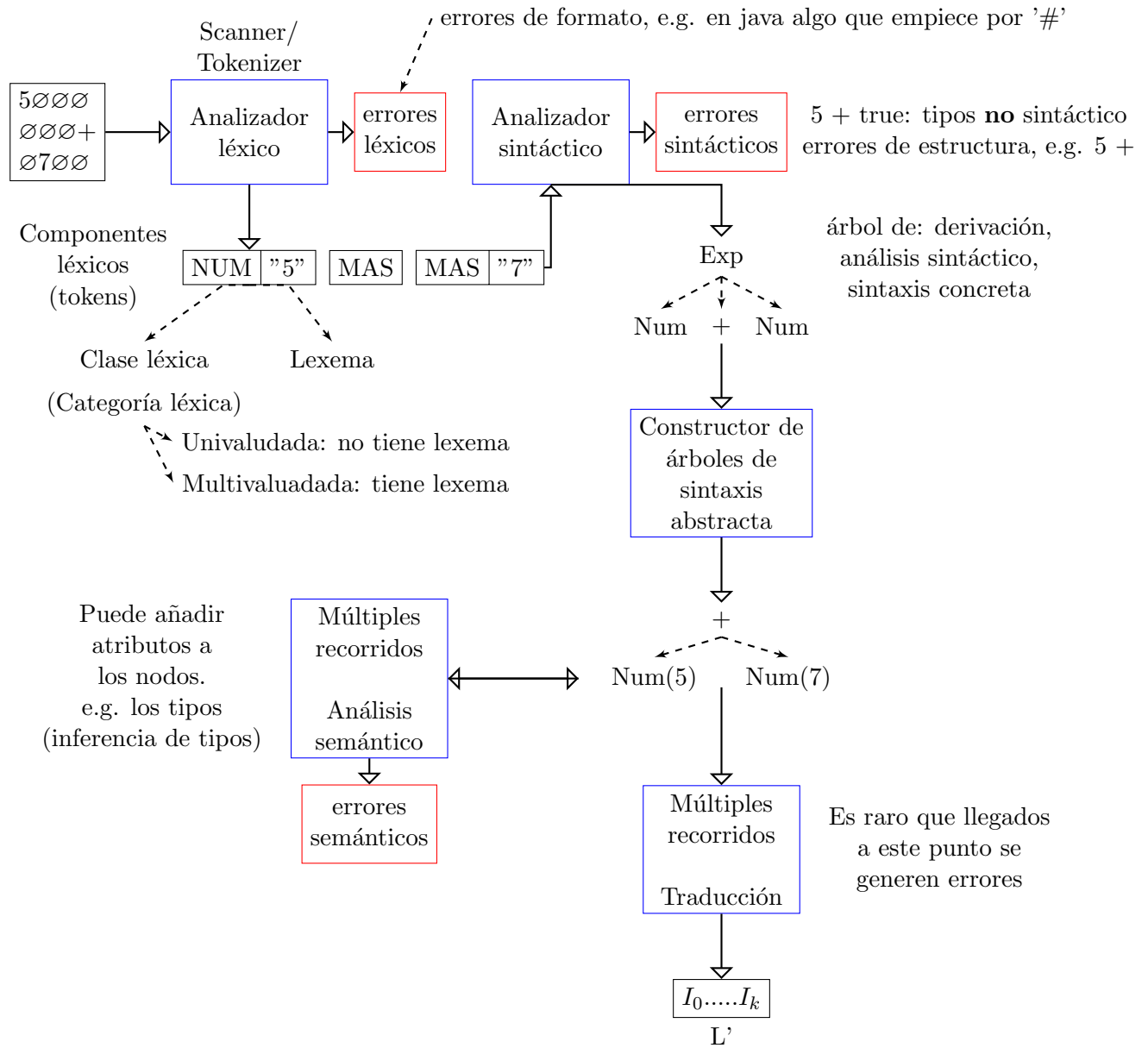
Java



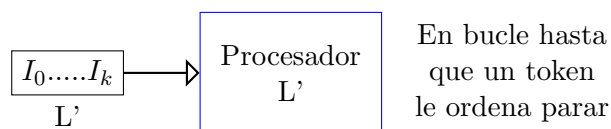
Abstracción



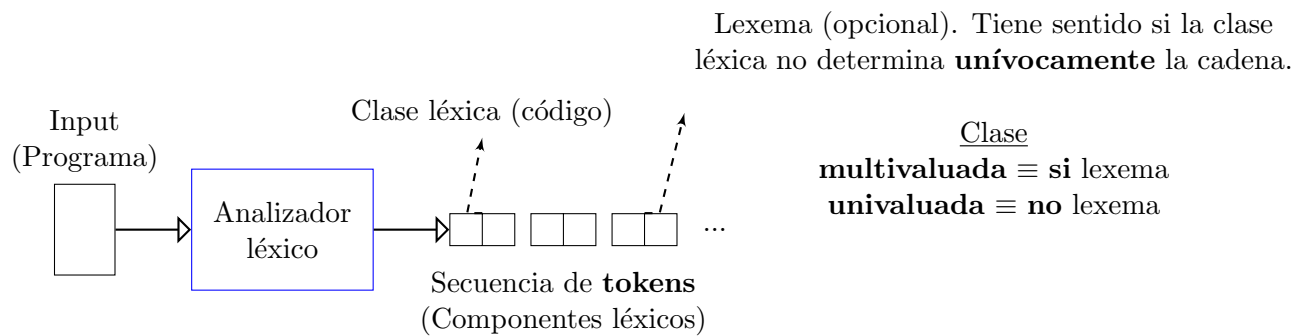
1.1 Traductor



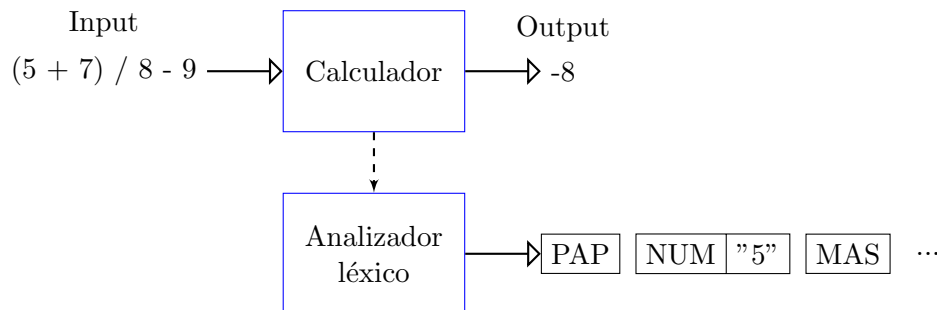
1.2 Intérprete



2 Análisis léxico



2.1 e.g. prototípico



Una expresión aritmética es...

- ...un **número entero**.
- ...una **expresión** seguida de un **operador** seguido de otra **expresión**.
- ... (seguido de una **expresión** seguido de).

Los operadores son: +, -, * y /.

2.2 Método de desarrollo

2.2.1 Identificación de las clases léxicas

Normalmente:

- Los **signos** (puntos, operadores...) conviene que sean clases **univaluadas**.
- **Variables** (NUM, LETRA...) son **multivaluadas**.
- **Palabras reservadas** (bool, int...) son **univaluadas**.

Clases léxicas de ??.

- PAP (paréntesis apertura)
- PCIERRE (paréntesis cierre)
- NUM
- MAS
- MENOS
- POR
- DIV

e.g.

```
class C {  
    int x;  
    ....
```

- class \equiv univaluada
- C \equiv id
- { \equiv univaluada
- int \equiv univaluada
- x \equiv id

2.2.2 Describir las clases léxicas

Cada clase léxica es un lenguaje formal de los posibles lexemas.

Suposición: son lenguajes regulares, por tanto:

- Se pueden describir mediante ER_s (expresiones regulares).
- Se pueden reconocer utilizando AFD_s (autómatas finitos deterministas).

e.g. NENT

- Descripción informal.
Empieza con un signo (+ o -) opcional.
A continuación aparecen uno o más dígitos.
e.g. +5, +007, 000, -0, -08, -280.
- Descripción formal.
 $(\backslash+ | \backslash-)?[0-9]^+$

Notación

- ER sobre el alfabeto Σ (e.g. UNICODE, ASCII...).
- Una "letra" de Σ , e.g. a. Denota $\{a\}$
- Si E_0 y E_1 son ER_s , también lo son:
 - * $(E_0 \mid E_1)$. Denota la unión, $L(E_0) \cup L(E_1)$.
 - * $(E_0 \bullet E_1)$. Denota la concatenación, $\{W_0W_1 \mid W_0 \in L(E_0), W_1 \in L(E_1)\}$.
 - * (E_0^*) . Denota $\{\epsilon, W_0, W_0W_1, W_0W_1W_2\ldots \mid W_0, W_1, W_2, \ldots \in L(E_0)\}$.
- ϵ denota el lenguaje formado por la cadena vacía.
- Convenios.
 - $*$ tiene mayor prioridad que \mid y \bullet .
 - \bullet tiene mayor prioridad que \mid .
 - \bullet se puede omitir.
 - Los $()$ se usan para cambiar prioridades.
 - Conjuntos de caracteres, e.g. $[0-9, a]$ es el conjunto formado por los dígitos y la a.
 - * $[e_0, \ldots, e_n]$ donde e_i puede ser...
 - ...una letra.
 - ...a-b, conjunto de caracteres comprendidos entre a y b.
 - * Conjuntos complementados $[\wedge e_0, \ldots, e_n]$.
e.g. $[\wedge 0-9, a-z]$ es el conjunto de caracteres que no son dígitos ni letras minúsculas.
 - * $E^+ \equiv EE^*$. Aparece una o más veces pero mínimo una.
 - * $E? \equiv (E \mid \epsilon)$. Aparece o no, es opcional.
 - * \backslash . La forma de escape.

e.g. Identificadores

- Pueden contener letras, dígitos y $_{-}$.
- Empiezan por letra o $_{-}$.
- Van seguidos de una secuencia de 0 o más caracteres válidos.

$[a-z, A-Z, _] [a-z, A-Z, _, 0-9]^*$

e.g. NENT, como en el ejemplo anterior pero sin 0_s a la izquierda.

$[\backslash +, \backslash -]^?([1-9] [0-9]^* \mid 0)$

Para mayor claridad vamos a utilizar DR_s (definiciones regulares) en lugar de ER_s .

- $(*) \equiv$ clase léxica, marca cual es la definición principal.
- $[I] \equiv$ ignorables, se utiliza para definir que caracteres no se han de tener en cuenta.
- Las palabras subrayadas corresponden a definiciones auxiliares.

e.g. Identificadores

$(*) \text{ IDEN} \equiv \underline{\text{Letra}} (\underline{\text{Letra}} \mid \underline{\text{Dig}})^*$
 $\text{Letra} \equiv [\text{a-z}, \text{A-Z}, _]$
 $\text{Dig} \equiv [0-9]$

e.g. Números enteros

$(*) \text{ LENT} \equiv \underline{\text{Signo}}? (0 \mid \underline{\text{DPos}} \underline{\text{Dig}}^*)$
 $\text{Signo} \equiv [\backslash+, \backslash-]$
 $\text{DPos} \equiv [1-9]$
 $\text{Dig} \equiv [0-9]$

e.g. Literales reales

- Empiezan por un entero.
- Continúan por...
 - ...una parte decimal.
 - ...una parte exponencial.
 - ...una parte decimal seguida de una exponencial.
- Parte decimal. '.' seguido de una secuencia de uno o más dígitos, sin ceros superfluos a la derecha.
- Parte exponencial. E o e seguida de un entero.

e.g. +5.7E-28

$(*) \text{ LREAL} \equiv \underline{\text{LENT}} (\underline{\text{PDEC}} \mid \underline{\text{PEXP}} \mid \underline{\text{PDEC}} \underline{\text{PEXP}})$
 $\text{PDEC} \equiv \backslash. (\underline{\text{Dig}}^* \underline{\text{DPos}} \mid 0)$
 $\text{PEXP} \equiv (\text{E} \mid \text{e}) \underline{\text{LENT}}$
 $\text{LENT} \equiv \underline{\text{Signo}}? (0 \mid \underline{\text{DPos}} \underline{\text{Dig}}^*)$
 $\text{Signo} \equiv [\backslash+, \backslash-]$
 $\text{DPos} \equiv [1-9]$
 $\text{Dig} \equiv [0-9]$

Ignorables

[I] SEP $\equiv [\text{'}, \backslash t, \backslash n, \backslash r, \backslash b]$, posibles separadores.

[I] COM $\equiv \# [\wedge \backslash n]^* \backslash n$, e.g. `# Esto es un comentario.\n`

e.g DR de ??

(*) PAP $\equiv \backslash ($

(*) PCIERRE $\equiv \backslash)$

(*) NUM $\equiv \text{Signo? } (0 \mid \underline{\text{DPos}} \underline{\text{Dig}}^*)$

Signo $\equiv [\backslash +, \backslash -]$

DPos $\equiv [1-9]$

Dig $\equiv [0-9]$

(*) MAS $\equiv \backslash +$

(*) MENOS $\equiv \backslash -$

(*) POR $\equiv \backslash *$

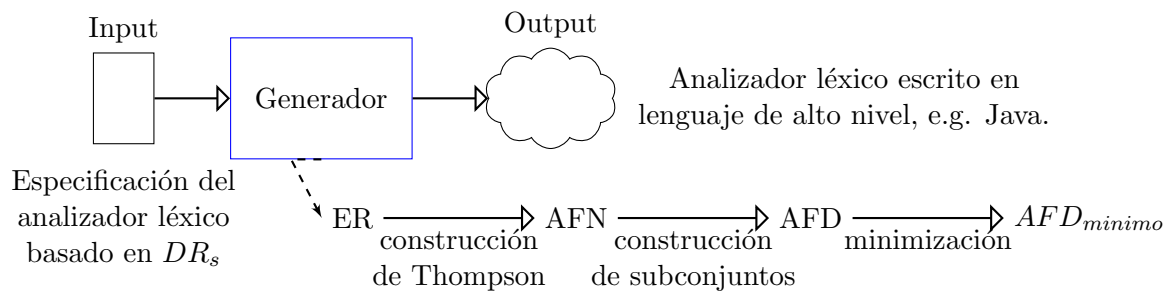
(*) DIV $\equiv \backslash /$

[I] SEP $\equiv [\text{'}, \backslash t, \backslash n, \backslash r, \backslash b]$

[I] COM $\equiv \# [\wedge \backslash n]^* \backslash n$

2.3 Implementación

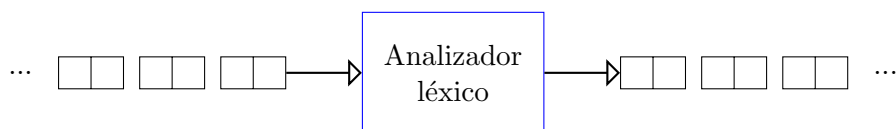
2.3.1 Generadores de analizadores léxicos



2.3.2 Implementación manual

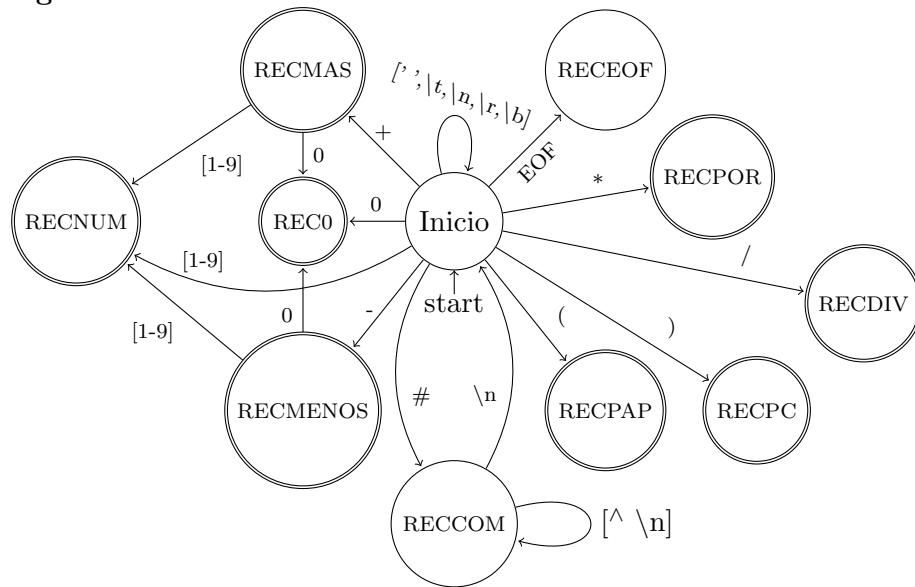
Se implementa utilizando AFD_s donde las transiciones a estados de error no se dibujan, son implícitas.

El AFD no reconoce la entrada completamente sino que utiliza una **arquitectura pull**, solicita los token uno a uno y los va procesando.



Siempre hay que tener una clase léxica que represente el EOF.

e.g. AFD de ??



Inicialización: Conseguir el primer caracter en sigCar.

```

SigToken
Estado ← EstadoInicial
lexema ← ""
loop {
    switch(Estado) {
        case  $S_0$ 
        ...
        case  $S_n$ 
    }
}

case  $S_i$ 
si sigCar  $\in T_0$ 
    Estado ← SigEstado
    lexema ← lexema + sigCar
    (lexema solo cambia en estados no ignorables)
    Actualizar sigCar
si no si sigCar  $\in T_1$ 
    Estado ← SigEstado
    lexema ← lexema + sigCar
    Actualizar sigCar
    ...
si no
    //Si es un estado final
    return token
    //Si no es un estado final
    error

```

Dos variables extra, fila y columna, para poder dar información extra en el mensaje de error.

Codificación parcial (pseudocódigo)

```
Estado ← Inicio
lexema ← ""
loop {
  switch (Estado) {
    Inicio: if sigCar = ( then transita(RECPAP)
           else if ...
           ...
           else if sigCar = # textbf then transitaIgnorando(RECCOM)
           else if sigCar ∈ [1-9] then transita(RECNUM)
           else error()

    RECNUM: if sigCar ∈ [1-9] then transita(RECNUM)
           else return token(NUM, lexema)

    ...

    RECCOM: if sigCar ≠ \n then transitaIgnorando(RECCOM)
           else transitaIgnorando(Inicio)
  }
}

transita(S) {
  Estado ← S
  lexema ← lexema + sigCar
  Actualiza sigCar
}

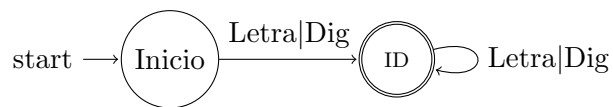
transitaIgnorando(S) {
  Estado ← S
  Actualiza sigCar
}
```

e.g.

Especificación léxica

EVALUA	(*) ID \equiv <u>Letra</u> (<u>Letra</u> <u>Dig</u>)*
5 + x	Letra \equiv [a-z,A-Z,-]
DONDE	Dig \equiv [0-9]
x = 27	(*) EVALUA \equiv [E,e][V,v][A,a][L,l][U,u][A,a]
	(*) DONDE \equiv [D,d][O,o][N,n][D,d][E,e]

En el return que reconoce los id se comprueba si el lexema es una palabra reservada, de serlo se devuelve la palabra reservada en lugar del id.

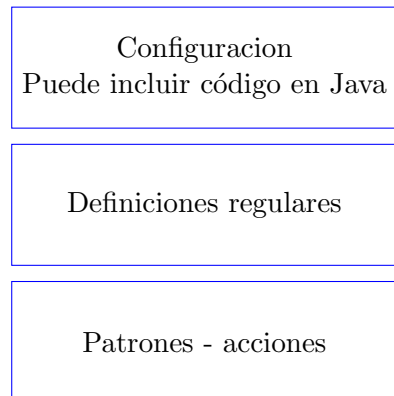


Hay un ejemplo de codificación en el campus, tiene una errata en el diagrama de transiciones.

2.3.3 Implementación mediante herramientas

La herramienta que vamos a utilizar es **JLex**, hay un ejemplo en el CV.

Formato de JLex:

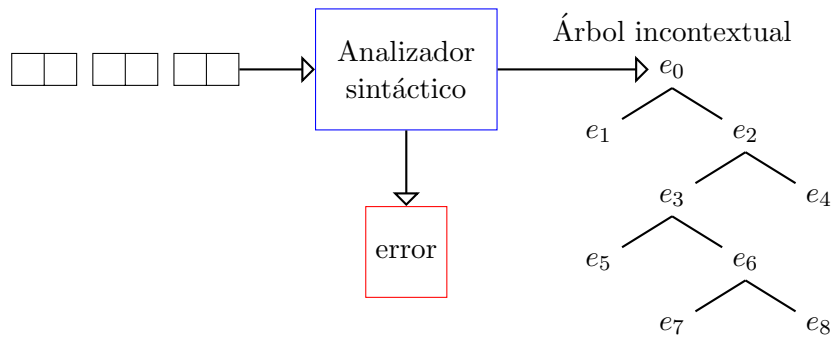


Para que una DR pueda referirse a otras tienen que haberse definido antes, lo que en las DR_s subrayamos tiene que ir entre {}, hacen sustituciones literales (no ponen paréntesis).

Las palabras reservadas se ponen sin mas, el orden de reconocimiento equivale al de aparición en el fichero, en la parte de los return.

Compilar: java -cp jlex.jar JLex.Main input(DR_s)

3 Análisis sintáctico



Hipótesis: Las cadenas de clases léxicas generadas por el analizador léxico forman un lenguaje incontextual.

3.1 Recordatorio

3.1.1 Gramática

Sea una GI (Gramática Incontextual), $GI(N,T,P,S)$:

- $N \equiv$ alfabeto de no terminales \Rightarrow clases sintácticas.
- $T \equiv$ alfabeto de terminales \Rightarrow clases léxicas.
- $P \equiv$ conjunto de reglas de la forma $A \rightarrow \alpha$, donde $A \in N$ y $\alpha \in (NUT)^*$.
- $S \in N$ y es el símbolo inicial.

Sea $G \equiv (N,T,P,S) \rightarrow G$ denota un lenguaje $L(G)$

- Relación de derivación \Rightarrow_G ($0 \Rightarrow$ si G se sobreentiende) $\Rightarrow \subseteq (NUT)^* \times (NUT)^*$
 $\alpha A \beta \wedge A \rightarrow \gamma \in P$ entonces $\alpha A \beta \Rightarrow \alpha \gamma \beta$
- Se considera \Rightarrow_* aplicar cero o más veces \Rightarrow .

$$L(G) = \{w \in T^* \mid S \Rightarrow_* w\}$$

e.g. Número binario

$$N = \{N,B\}$$

$$T = \{0,1\}$$

$$P = N \rightarrow B$$

$$N \rightarrow NB$$

$$B \rightarrow 0$$

$$B \rightarrow 1$$

$$S = N$$

Esto sería equivalente a dar tan solo las reglas(P), donde:

- El símbolo azul de la primera regla es el símbolo inicial(S).
- El conjunto de símbolos azules son los no terminales(N).
- El conjunto de símbolos rojos son los terminales(T).

La derivación se denomina *mas a la izquierda* si siempre se reescribe el no terminal que está mas a la izquierda, equivalente para *mas a la derecha*.

Derivación

$N \Rightarrow NB \Rightarrow NBB \Rightarrow N1B \Rightarrow B1B \Rightarrow B10 \Rightarrow 010$

Derivación mas a la izquierda

$N \Rightarrow NB \Rightarrow NBB \Rightarrow BBB \Rightarrow 0BB \Rightarrow 01B \Rightarrow 010$

Derivación mas a la derecha

$N \Rightarrow NB \Rightarrow N0 \Rightarrow NB0 \Rightarrow N10 \Rightarrow B10 \Rightarrow 010$

3.1.2 Árbol de análisis sintáctico

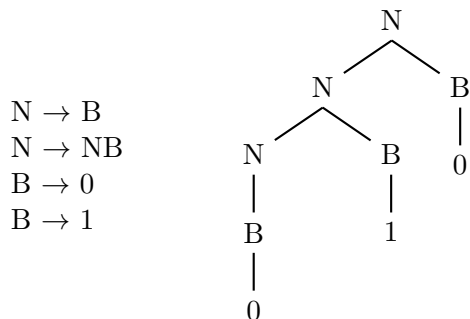
La idea es obtener una representación única para cada sentencia.

Árboles

- La raíz está etiquetada con el símbolo inicial.
- Están ordenados, hay un orden en los hijos de los nodos (primer hijo, segundo hijo...).
- Los nodos internos están etiquetados por no terminales.
- Los nodos hoja están etiquetados por terminales o por ϵ .
- Si un nodo está etiquetado por A y sus hijos por α entonces $A \rightarrow \alpha \in P$.

$w \in L(G) \Leftrightarrow \begin{array}{c} s \\ \triangle \\ w \end{array} \text{ (es un árbol de análisis sintáctico)}$

e.g.



Cada nodo es equivalente a un array en el que se referencian sus hijos en orden.

3.2 Especificación sintáctica

3.2.1 Determinar las clases sintácticas

¿Cómo? A partir de la especificación informal.

De arriba a abajo:

- Identificar las clases complejas.
- Descomponerlas en clases más simples hasta llegar a las clases léxicas.

Descomposición → Las clases resultantes deben estar al mismo nivel, el más alto posible.

e.g. L que describe libros

Titulo [...]
Autores [Autor]...[Autor]
Capitulo [...][...]
...
Capitulo [...][...]

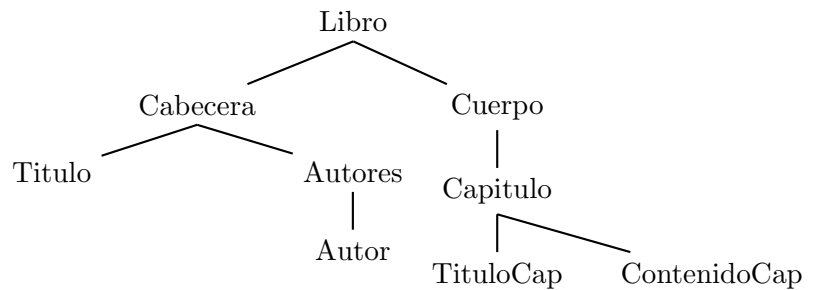


Figure 3.2.1.1 : Árbol informal.

Clases léxicas
Autores
Capitulo
CAP
CCIE
Texto

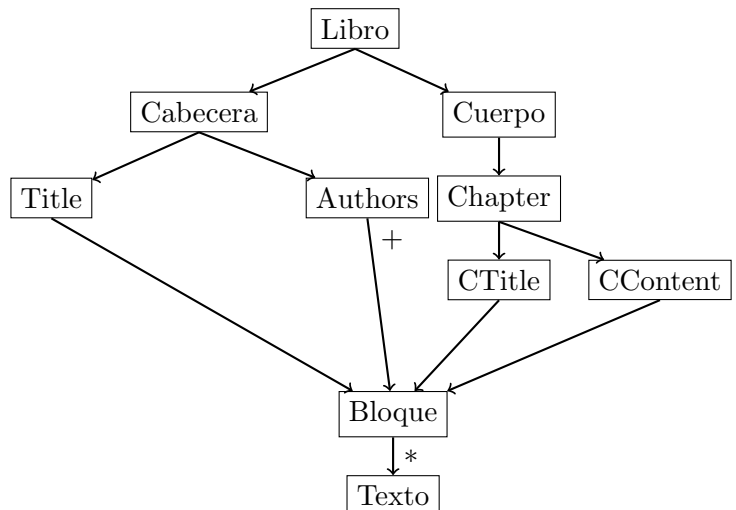


Figure 3.2.1.2 : Árbol formal.

3.2.2 Patrones para secuencias

- Secuencia de 1 o más I_s , separados por \square .

e.g. $I \mid I \square I \mid I \square I \square I \dots$

$LI \rightarrow I$

$LI \rightarrow LI \square I$

e.g.

$NUM \rightarrow B$

$NUM \rightarrow NUM , B$

- Secuencia de 0 o más I_s , separados por \square .

$S \rightarrow \varepsilon$

$S \rightarrow LI$

$LI \rightarrow I$

$LI \rightarrow LI \square I$

e.g. Declaraciones

$Decs \rightarrow \varepsilon$

$Decs \rightarrow LDEC$

$LDEC \rightarrow Dec$

$LDEC \rightarrow LDEC ; Dec$

Como caso particular de los patrones anteriores tenemos aquellos en que \square es ε , no están separados por nada.

- Secuencia de 1 o más I_s .

$LI \rightarrow I$

$LI \rightarrow LI I$

- Secuencia de 0 o más I_s .

$S \rightarrow \varepsilon$

$S \rightarrow LI$

$LI \rightarrow I$

$LI \rightarrow LI I$

También es posible usar un terminador en lugar de un separador para cada ítem.

e.g. $I \square I \square I \square I \square \dots$ donde $I \square$ es un ítem.

- Secuencia de 1 o más I_s terminados en \square .

$LI \rightarrow BI$

$LI \rightarrow LI BI$

$BI \rightarrow I \square$

- Secuencia de 0 o más I_s terminados en \square .

$S \rightarrow \varepsilon$

$S \rightarrow LI$

$LI \rightarrow BI$
 $LI \rightarrow LI BI$
 $BI \rightarrow I \square$

e.g. Libro

Libro \rightarrow Cabecera Cuerpo
 Cabecera \rightarrow Title Authors
 Title \rightarrow Titulo Bloque
 Authors \rightarrow Autores LBloques
 LBloques \rightarrow Bloque
 LBloques \rightarrow LBloques Bloque
 Cuerpo \rightarrow LChapter
 LChapter \rightarrow Chapter
 LChapter \rightarrow LChapter Chapter

Chapter \rightarrow CTitle CContent
 CTitle \rightarrow Capitulo Bloque
 CContent \rightarrow Bloque
 Bloque \rightarrow [CTexto]
 CTexto $\rightarrow \varepsilon$
 CTexto \rightarrow LTexto
 LTexto \rightarrow Texto
 LTexto \rightarrow LTexto Texto

3.2.3 Patrones para operadores

La prioridad indica el orden en que los operadores se evalúan, si dos tienen la misma prioridad se consulta la asociatividad, a izquierdas o a derechas.

e.g.

$\underline{\underline{5 + 6 - 7 * 8 + 9}}$ * es el más prioritario.
 + y - son igual de prioritarios pero asocian a izquierdas.

Cada operador tiene un nivel de prioridad, pueden existir múltiples operadores con un mismo nivel.

- Operadores binarios infijos.
 Pueden asociar a izquierdas
e.g. $\underline{\underline{5 + 6 + 7}}$
 O a derechas.
e.g. $\underline{\underline{5 + 6 + 7}}$
 O no asociar, no está permitido encadenar operadores.
e.g. $5 + 6 + 7$
- Operadores unarios prefijos.
e.g. -5
 Pueden ser...
 - ...**asociativos**, pueden encadenarse varios.
 - ...**no asociativos**, no pueden encadenarse.
 Asocian siempre a derechas.

- Operadores unarios postfijos.

e.g. $5+$, $x++$

Pueden ser...

- ...**asociativos**, pueden encadenarse varios.
- ...**no asociativos**, no pueden encadenarse.

Asocian siempre a izquierdas.

e.g. Lenguaje

- Numeros y variables.

- $+$, $-$, $*$, $/$, $-$ unario.

Operador	Aridad	Tipo	Prioridad	Asociatividad
$+, -$	2		0	izquierda
$*, /$	2		1	izquierda
$-$	1	prefijo	2	si

$$\underline{\underline{5 + 6 - \underline{\underline{-7 * 8 + 9}}}}$$

Si existen varios operadores con la misma prioridad pero que asocian distinto porque existen múltiples interpretaciones.

Operador	Aridad	Tipo	Prioridad	Asociatividad
$+$	2		0	izquierda
$-$	2		0	derecha
$*, /$	2		1	izquierda
$-$	1	prefijo	2	si

$$\underline{\underline{5 + 6 + 7}}$$

$$\underline{\underline{5 + 6 + 7}}$$

$\text{Exp} \rightarrow \text{Variable}$

$\text{Exp} \rightarrow \text{Numero}$

$\text{Exp} \rightarrow \text{OpUn Exp}$

$\text{Exp} \rightarrow \text{Exp OpBin Exp}$

$\text{OpUn} \rightarrow -$

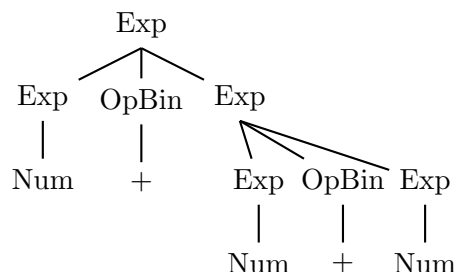
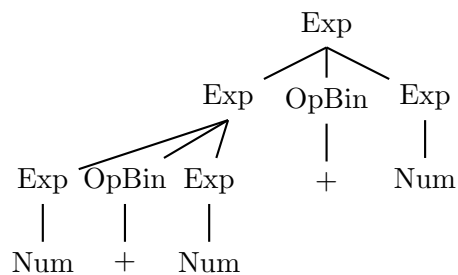
$\text{OpBin} \rightarrow +$

$\text{OpBin} \rightarrow -$

$\text{OpBin} \rightarrow *$

$\text{OpBin} \rightarrow /$

$\text{Num} + \text{Num} + \text{Num}$



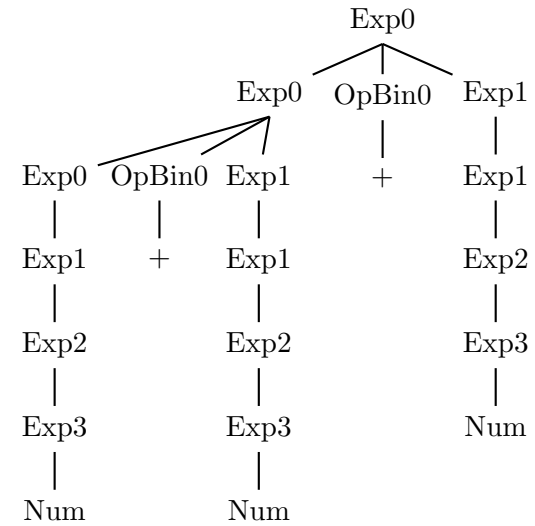
Como hemos encontrado dos posibles árboles de derivación la gramática es ambigua y no vale.

Cada nivel de prioridad tendrá su propio terminal.

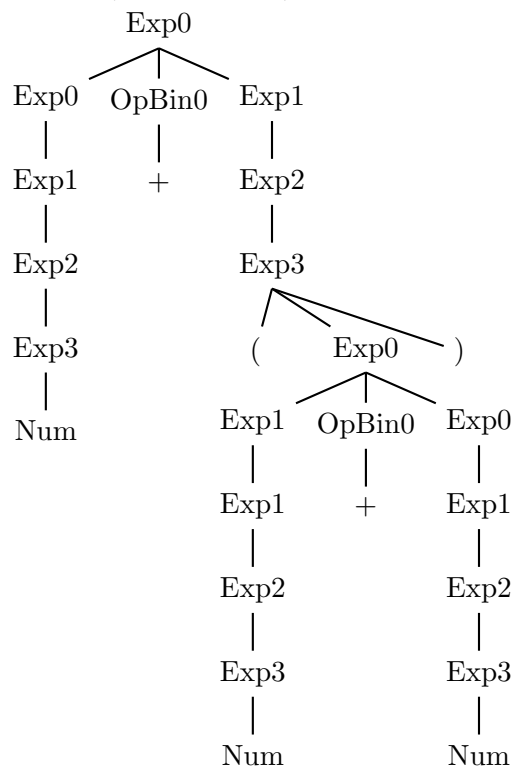
Si crece a izquierdas reduzco el nivel de prioridad de la derecha, si es a derechas baja el de la izquierda y si no asocia ambos aumentan.

$\text{Exp0} \rightarrow \text{Exp0 OpBin0 Exp1}$
 $\text{Exp0} \rightarrow \text{Exp1}$
 $\text{Exp1} \rightarrow \text{Exp1 OpBin1 Exp2}$
 $\text{Exp1} \rightarrow \text{Exp2}$
 $\text{Exp2} \rightarrow \text{Exp2 OpBin2 Exp3}$
 $\text{Exp2} \rightarrow \text{Exp3}$
 $\text{Exp3} \rightarrow (\text{Exp0})$
 $\text{Exp3} \rightarrow \text{Num}$
 $\text{Exp3} \rightarrow \text{Var}$
 $\text{OpBin0} \rightarrow +$
 $\text{OpBin0} \rightarrow -$
 $\text{OpBin1} \rightarrow *$
 $\text{OpBin1} \rightarrow /$
 $\text{OpUn2} \rightarrow -$

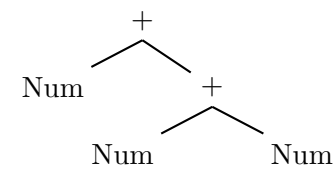
$\text{Num} + \text{Num} + \text{Num}$



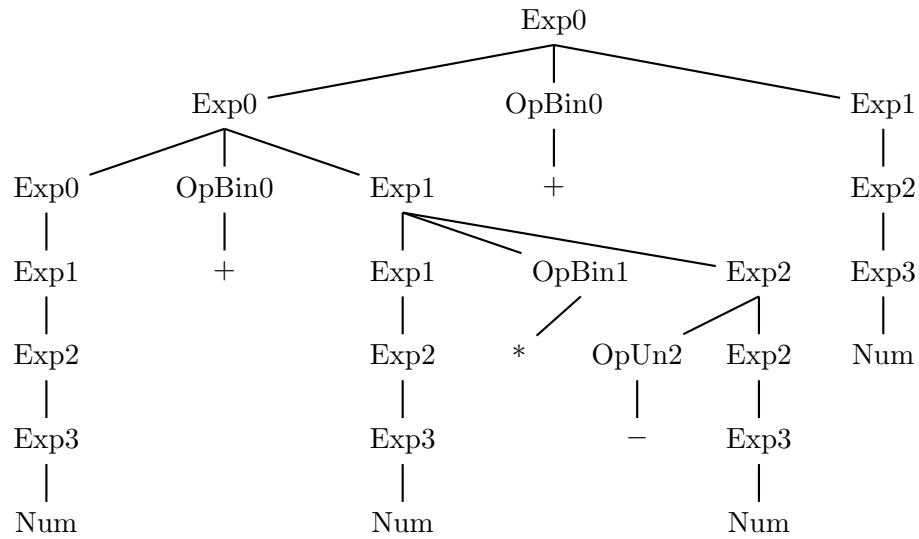
$\text{Num} + (\text{Num} + \text{Num})$



El árbol de la izquierda sería equivalente al siguiente.



$Num - Num * -Num + Num$



e.g. Lenguaje

Expresiones básicas: números e identificadores, se pueden usar ().

Operador	Prioridad	Aridad	Tipo	Asociatividad
==, !=	0	2		no
+	1	2		izq.
*	1	1	prefijo	no
=	2	2		der.
++	2	1	postfijo	si
~	3	1	prefijo	si
<<	3	1	postfijo	no

$Exp0 \rightarrow Exp1 \ Op0 \ Exp1$

$Exp0 \rightarrow Exp1$

$Exp1 \rightarrow Exp1 \ + \ Exp2$

$Exp1 \rightarrow * \ Exp2$

$Exp1 \rightarrow Exp2$

$Exp2 \rightarrow Exp3 \ = \ Exp2$

$Exp2 \rightarrow Exp2 \ ++$

$Exp2 \rightarrow Exp3$

$Exp3 \rightarrow \sim \ Exp3$

$Exp3 \rightarrow Exp4 \ <<$

$Exp3 \rightarrow Exp4$

$Exp4 \rightarrow Num$

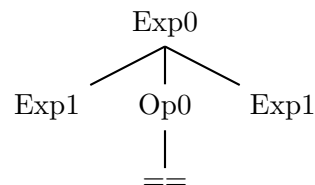
$Exp4 \rightarrow Id$

$Exp4 \rightarrow (\ Exp0 \)$

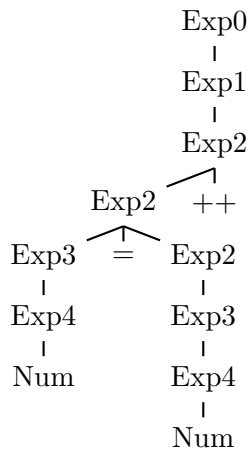
$Op0 \rightarrow ==$

$Op0 \rightarrow !=$

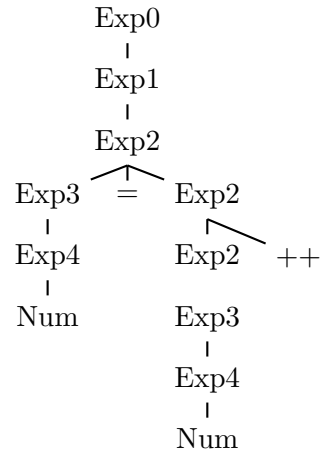
$Num == Num == Num$, está prohibido.



5 = 6 ++

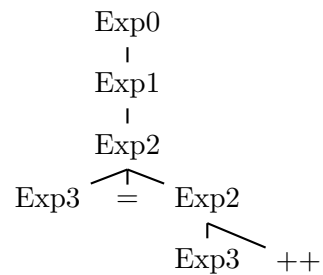


5 = 6 ++



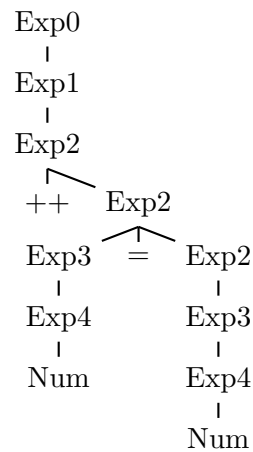
Supongamos que el operador ++ no asocia.
 Ahora, $\text{Exp2} \rightarrow \text{Exp3} ++$

5 = 6 ++



Supongamos que el operador ++ es prefijo
 y asocia.
 Ahora, $\text{Exp2} \rightarrow ++ \text{Exp2}$

5 = 6 ++



- Los operadores prefijos que asocian son incompatibles con los operadores que asocian a izquierdas.

- Los operadores postfijos que asocian son incompatibles con los operadores que asocian a derechas.

e.g. Dada g obtenemos su tabla de operadores.

$\text{Exp0} \rightarrow \text{Exp0} \rightarrow \text{Exp1}$

$\text{Exp0} \rightarrow \text{Exp1} \sim \text{Exp1}$

$\text{Exp0} \rightarrow \# \text{Exp1}$

$\text{Exp0} \rightarrow \text{Exp1}$

$\text{Exp1} \rightarrow \text{Exp1} @ \text{Exp2}$

$\text{Exp1} \rightarrow \text{Exp2} \wedge \text{Exp2}$

$\text{Exp1} \rightarrow \text{Exp2}$

$\text{Exp2} \rightarrow \text{Exp2} *$

$\text{Exp2} \rightarrow \text{Exp3}$

$\text{Exp3} \rightarrow \text{Num}$

$\text{Exp3} \rightarrow (\text{Exp0})$

Operador	Prioridad	Aridad	Tipo	Asociatividad
\rightarrow	0	2		izq.
\sim	0	2		no
$\#$	0	1	prefijo	no
$@$	1	2		izq.
\wedge	1	2		no
$*$	2	1	postfijo	si

Supongamos que en un mismo nivel de prioridad existen:

- $+$, binario que asocia a izquierdas.
- $-$, unario prefijo que no asocia.
- $*$, unario postfijo que asocia.

Hay problemas con estos operadores por conflictos de crecimiento, si un operador crece a izquierdas y el otro a derechas los árboles que generan son incompatibles.

No	A izquierdas	A derechas	Si	Estado
$+, -, *$				Ok
$+, -$			$*$	Ok
$+, *$			$-$	Ok
$-, *$	$+$			Ok
$-, *$		$+$		Ok
$-$	$+$		$*$	Ok
$*$		$+$	$-$	Ok
	$+$		$*, -$	X
$-$		$+$	$*$	X
$*$	$+$		$-$	X
		$+$	$*, -$	X

e.g. Lenguaje, especificación completa.

Un programa es una secuencia de una o más instrucciones separadas por $;$.

Se dispone de los siguientes tipos de instrucciones:

- Asignación.
Una variable seguida de \leftarrow seguida de una expresión.
- Si.
IF expresión Instrucciones **FI**
- Si-Sino.
IF expresión
 I_0
ELSE
 I_1
FI
- Bucle.
WHILE expresion **DO**
 I
OD
- Bloque.
{ Secuencia de I_s separadas por ; }

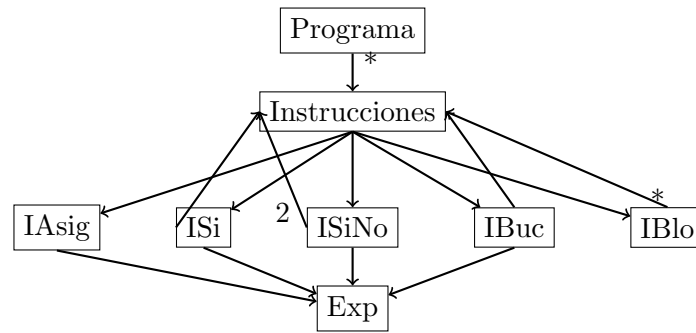
Las expresiones básicas las componen números (literales) y variables.
Se dispone de los siguientes operadores:

- Comparación. $<$, $>$, $==$, $!=$, $<=$, $>=$
Binarios, menor nivel de prioridad y no asocian.
- Aritméticos. $+$, $-$, $*$, $/$
Binarios, más prioridad que los de comparación y asocian a izquierdas.
- Lógicos. **and**, **or**
Binarios, más prioridad que los aritméticos y asocian a derechas.
- Unarios. **!**, **-**
Prefijos, más prioritarios que los lógicos y asocian.

Se pueden utilizar paréntesis, (...).

Un posible programa sería:

```
x ← 25
WHILE x <= 50 DO
    x ← x + 1
OD
```



Operador	Prioridad	Aridad	Tipo	Asociatividad
<, >, ==,				
!=, <=, >=	0	2		no
+, -, *, /	1	2		izq.
and, or	2	2		der.
!, -	3	1	prefijo	si

Programa \rightarrow LI

LI \rightarrow I

LI \rightarrow LI ; I

I \rightarrow IAsig

I \rightarrow ISi

I \rightarrow ISino

I \rightarrow IBuc

I \rightarrow IBlo

IAsig \rightarrow Var \leftarrow Exp0

ISi \rightarrow IF Exp0 I FI

ISino \rightarrow IF Exp0 I ELSE I FI

IBuc \rightarrow WHILE Exp0 DO I OD

IBlo \rightarrow { S }

S \rightarrow ϵ

S \rightarrow LI

Exp0 \rightarrow Exp1 Op0 Exp1

Exp0 \rightarrow Exp1

Exp1 \rightarrow Exp1 Op1 Exp2

Exp1 \rightarrow Exp2

Exp2 \rightarrow Exp3 Op2 Exp2

Exp2 \rightarrow Exp3

Exp3 \rightarrow Op3 Exp3

Exp3 \rightarrow Exp4

Exp4 \rightarrow (Exp0)

Exp4 \rightarrow Num

Exp4 \rightarrow Var

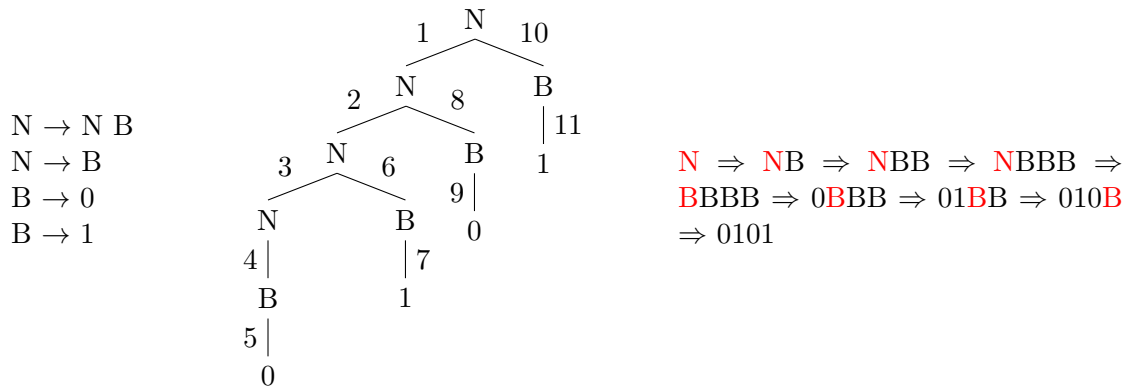
OpX se corresponde con los operadores del nivel X.

3.3 Implementación de analizadores descendentes

Vamos a ver analizadores descendentes (top-down parsers), concretamente implementaremos analizadores sintácticos predictivos recursivos descendentes.

Los árboles expanden de izquierda a derecha equivalente a una derivación más a la izquierda, los algoritmos que veremos se basan en este tipo de derivación.

e.g. 0101



Pila de análisis (\leftarrow)	Entrada	Movimiento
N	0101	Expandir $N \rightarrow N B$
NB	0101	Expandir $N \rightarrow N B$
NBB	0101	Expandir $N \rightarrow N B$
NBBB	0101	Expandir $N \rightarrow B$
BBBB	0101	Expandir $B \rightarrow 0$
0BBB	0101	Desplazar 0
BBB	0 01	Expandir $B \rightarrow 1$
1BB	0 01	Desplazar 1
BB	01 01	Expandir $B \rightarrow 0$
0B	01 01	Desplazar 0
B	010 1	Expandir $B \rightarrow 1$
1	010 1	Desplazar 1
	0101	Aceptar

Gramática	Los no terminales se interpretan como funciones que reconocen partes de la entrada.
$N' \rightarrow N \vdash$	Las decisiones se basan en un token, el algoritmo siempre lleva un token adelantado.
$N \rightarrow B R$	Los terminales consumen el token si coincide y rechazan la entrada si no.
$R \rightarrow B R$	
$R \rightarrow \varepsilon$	
$B \rightarrow 0$	Para asegurarse de que la cadena pertenece y no solo un prefijo siempre se añade una producción adicional, N' , que llama al símbolo inicial y luego reconoce el fin de fichero.
$B \rightarrow 1$	

Pseudocódigo

```

Np() {
  N();
  rec(EOF);
}

N() {
  B();
  R();
}

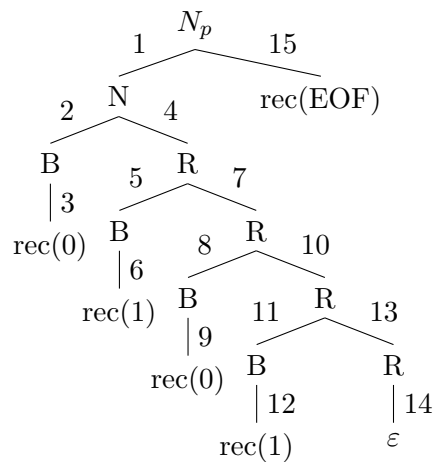
B() {
  if sigToken == 0 then
    rec(0);
  else
    rec(1);
  fi
}

R() {
  if sigToken ∈ {0,1} then
    B();
    R();
  fi
}

rec(T) {
  if sigToken == T then
    sigToken ← sigToken();
  else
    error
  fi
}

```

e.g. 0101-

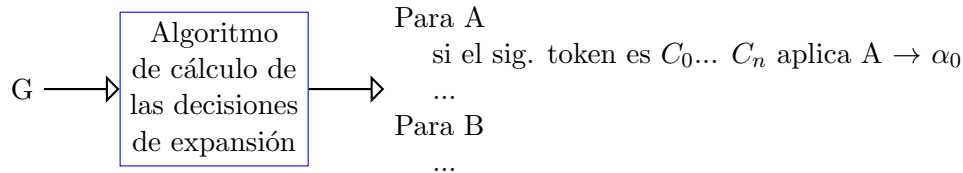


¿Funciona este método para cualquier gramática?

Como se puede ver con un único token no es posible decidir que producción usar.

$$\begin{array}{ll} A \rightarrow 0 X & L \rightarrow L i \\ A \rightarrow 0 Y & L \rightarrow i \end{array}$$

En general la recursión a izquierdas es problemática para esta metodología. La solución es transformar estas gramáticas con métodos que veremos más adelante.



Con este método podemos implementar gramáticas LL(1).

Si tenemos $A \rightarrow \alpha$ tenemos que preguntarnos por qué token pueden empezar las cadenas de α y comprobar sigToken.

Vamos a definir una función que emplearemos más adelante:

PRIM(α), primeros de α .

Son todos los terminales que pueden comenzar sentencias derivables desde α . Además, $\varepsilon \in \text{PRIM}(\alpha)$ si $\alpha \Rightarrow^* \varepsilon$.

e.g.

$$\begin{array}{ll} N \rightarrow B R & \\ R \rightarrow B R & \\ R \rightarrow \varepsilon & \text{PRIM}(N) = \{0, 1\} \\ B \rightarrow 0 & \text{PRIM}(B R) = \{0, 1\} \\ B \rightarrow 1 & \end{array}$$

Luego $A \rightarrow \alpha$ se podrá aplicar durante la expansión si el siguiente símbolo $\in \text{PRIM}(\alpha) - \{\varepsilon\}$. Al conjunto de estos símbolos se le denomina directores, $\text{Dir}(A \rightarrow \alpha)$.

$$\text{PRIM}(\alpha) - \{\varepsilon\} \subseteq \text{Dir}(A \rightarrow \alpha)$$

Introducimos un operador, concatenación en modulo ε denotada por el símbolo \oplus .

$$\begin{array}{ll} \Gamma_0 \oplus \Gamma_1 = (\Gamma_0 - \{\varepsilon\}) \cup \Gamma_1 & \text{si } \varepsilon \in \Gamma_0 \\ & (\Gamma_0 - \{\varepsilon\}) \cup \emptyset \quad \text{c.o.c.} \end{array}$$

e.g.

$$\begin{array}{l} \{a, b\} \oplus \{c, d\} = \{a, b\} \\ \{a, b, \varepsilon\} \oplus \{c, d\} = \{a, b, c, d\} \\ \{a, b, \varepsilon\} \oplus \{c, d, \varepsilon\} = \{a, b, c, d, \varepsilon\} \end{array}$$

Sea **a** un terminal y **A** un no terminal.

$$\text{PRIM}(a) = \{a\}$$

$$\begin{aligned}\text{PRIM}(\varepsilon) &= \{\varepsilon\} \\ \text{PRIM}(A) &= \cup \text{PRIM}(\alpha_i) \quad \text{donde } \{\alpha \mid A \rightarrow \alpha\}\end{aligned}$$

En general:

$$\text{PRIM}(X_0 \dots X_n) = \text{PRIM}(X_0) \oplus \text{PRIM}(X_1) \oplus \dots \oplus \text{PRIM}(X_n)$$

e.g.

$$\text{PRIM}(B \ R) = \text{PRIM}(B) \oplus \text{PRIM}(R)$$

e.g.

$$\begin{aligned}N &\rightarrow B \ R & \text{PRIM}(N) &= \text{PRIM}(B) \oplus \text{PRIM}(R) \\ R &\rightarrow B \ R & \text{PRIM}(R) &= \text{PRIM}(B) \oplus \text{PRIM}(R) \cup \text{PRIM}(\varepsilon) \\ R &\rightarrow \varepsilon & \text{PRIM}(B) &= \text{PRIM}(0) \cup \text{PRIM}(1) \\ B &\rightarrow 0 & \text{Por tanto:} & \\ B &\rightarrow 1 & \text{PRIM}(B) &= \{0, 1\} \\ & & \text{PRIM}(R) &= \{0, 1\} \cup \{\varepsilon\} \\ & & \text{PRIM}(N) &= \{0, 1\}\end{aligned}$$

Los primeros no resuelven $R \rightarrow \varepsilon$ porque el cuerpo de la regla se puede anular, necesitamos más terminales.

Introducimos el operador siguientes de A, $\text{SIG}(A)$.

El operador se define como el conjunto de posibles terminales que siguen en las sentencias del lenguaje los fragmentos derivados de A.

$\vdash \in \text{SIG}(S)$, donde S es el símbolo inicial.

Dada la producción $A \rightarrow \alpha B \beta$ sabemos que:

- $\text{PRIM}(\beta) - \{\varepsilon\} \subseteq \text{SIG}(B)$
- Si $\varepsilon \in \text{PRIM}(\beta) \Rightarrow \text{SIG}(A) \subseteq \text{SIG}(B)$

e.g.

$$\begin{aligned}\text{Obj} &\rightarrow \text{Desig} \ \text{RParams} & \text{SIG}(\text{Desig}) &= \text{PRIM}(\text{RParams}) - \{\varepsilon\} \cup \text{SIG}(\text{Obj}) \\ \text{RParams} &\rightarrow (\ \text{LParams} \) & \text{SIG}(\text{Obj}) &= \{;\} \\ \text{RParams} &\rightarrow \varepsilon & \text{PRIM}(\text{RParams}) &= \{(\ , \ \varepsilon\} \\ & & \text{SIG}(\text{Desig}) &= \{(\ , ;\}\end{aligned}$$

e.g.

$$\begin{array}{ll}
N \rightarrow B R & \text{SIG}(N) = \{\vdash\} \\
R \rightarrow B R & \text{SIG}(R) = \text{PRIM}(\varepsilon) \oplus \text{SIG}(N) \cup \text{PRIM}(\varepsilon) \oplus \text{SIG}(R) \\
R \rightarrow \varepsilon & = \{\varepsilon\} \oplus \text{SIG}(N) \cup \{\varepsilon\} \oplus \text{SIG}(R) \\
B \rightarrow 0 & = \text{SIG}(N) \cup \text{SIG}(R) \\
B \rightarrow 1 & = \{\vdash\} \cup \text{SIG}(R) \\
& = \{\vdash\} \\
& \text{SIG}(B) = \text{PRIM}(R) \oplus \text{SIG}(N) \cup \text{PRIM}(R) \oplus \text{SIG}(R) \\
& = \{0, 1, \varepsilon\} \oplus \text{SIG}(N) \cup \{0, 1, \varepsilon\} \oplus \text{SIG}(R) \\
& = \{0, 1, \vdash\} \cup \{0, 1, \vdash\}
\end{array}$$

Revisemos ahora los el conjunto de directores, teniamos que:

$$\text{DIR}(A \rightarrow \alpha) = \text{PRIM}(\alpha) - \{\varepsilon\}$$

Pero ahora al anularse podemos definir qué ocurre, la formula quedaria:

$$\text{DIR}(A \rightarrow \alpha) = \text{PRIM}(\alpha) \oplus \text{SIG}(A)$$

e.g.

$$\begin{array}{ll}
N \rightarrow B R & \text{DIR}(N \rightarrow B R) = \text{PRIM}(B R) \oplus \text{SIG}(N) \\
R \rightarrow B R & = \text{PRIM}(B) \oplus \text{PRIM}(R) \oplus \text{SIG}(N) \\
R \rightarrow \varepsilon & = \{0, 1\} \\
B \rightarrow 0 & \text{DIR}(R \rightarrow B R) = \{0, 1\} \\
B \rightarrow 1 & \text{DIR}(R \rightarrow \varepsilon) = \text{PRIM}(\varepsilon) \oplus \text{SIG}(R) \\
& = \{\varepsilon\} \oplus \text{SIG}(R) \\
& = \{\vdash\} \\
& \text{DIR}(B \rightarrow 0) = \{0\} \\
& \text{DIR}(B \rightarrow 1) = \{1\}
\end{array}$$

Sea G una gramática LL1 de la forma

$$A \rightarrow \alpha$$

$$A \rightarrow \beta$$

siempre se cumple que $\text{DIR}(A \rightarrow \alpha) \cap \text{DIR}(A \rightarrow \beta) = \emptyset$

Revisemos ahora el pseudocódigo de R.

```

R() {
  if sigToken ∈ {0,1} then
    B();
    R();
  elif sigToken ∈ {⊢} then
    //TODO
  else
    error({0, 1, ⊢})
  fi
}

error(Γ) {
  print Token inesperado sigToken, se esperaba Γ.
}

```

e.g.

$$\begin{array}{ll}
L \rightarrow [E_s] & \text{PRIM}(L) = \text{PRIM}([\] \oplus \text{PRIM}(E_s) \oplus \text{PRIM}([\] \\
E_s \rightarrow LI_s & \quad = \{\} \\
E_s \rightarrow \varepsilon & \text{PRIM}(E_s) = \text{PRIM}(LI_s) \cup \text{PRIM}(\varepsilon) \\
LI_s \rightarrow I RLI_s & \quad = \{a, [, \varepsilon\} \\
RLI_s \rightarrow , I RLI_s & \text{PRIM}(LI_s) = \text{PRIM}(I) \oplus \text{PRIM}(RLI_s) \\
RLI_s \rightarrow \varepsilon & \quad = \{a, [\} \\
I \rightarrow a & \text{PRIM}(RLI_s) = \text{PRIM}(,) \oplus \text{PRIM}(I) \oplus \text{PRIM}(RLI_s) \cup \text{PRIM}(\varepsilon) \\
I \rightarrow L & \quad = \{, , \varepsilon\} \\
& \text{PRIM}(I) = \text{PRIM}(a) \cup \text{PRIM}(L) \\
& \quad = \{a, [\} \\
[a, a, [], [a, a]] &
\end{array}$$

$$\begin{aligned}
\text{SIG}(L) &= \{\vdash\} \cup \text{PRIM}(\varepsilon) \oplus \text{SIG}(I) \\
&= \{\vdash, , [\} \\
\text{SIG}(E_s) &= \text{PRIM}([\] \oplus \text{SIG}(L) \\
&= \{\} \\
\text{SIG}(LI_s) &= \text{PRIM}(\varepsilon) \oplus \text{SIG}(E_s) \\
&= \{\} \\
\text{SIG}(RLI_s) &= \text{PRIM}(\varepsilon) \oplus \text{SIG}(LI_s) \cup \text{PRIM}(\varepsilon) \oplus \text{SIG}(RLI_s) \\
&= \{\} \\
\text{SIG}(I) &= \text{PRIM}(RLI_s) \oplus \text{SIG}(LI_s) \cup \text{PRIM}(RLI_s) \oplus \text{SIG}(RLI_s) \\
&= \{, , []\}
\end{aligned}$$

$$\begin{aligned}
\text{DIR}(L \rightarrow [E_s]) &= \text{PRIM}([E_s]) \oplus \text{SIG}(L) \\
&= \text{PRIM}([\] \oplus \text{PRIM}(E_s) \oplus \text{PRIM}([\] \oplus \text{SIG}(L) \\
&= \{\} \\
\text{DIR}(E_s \rightarrow LI_s) &= \{a, [\} \\
\text{DIR}(E_s \rightarrow \varepsilon) &= \{\} \\
\text{DIR}(LI_s \rightarrow I RLI_s) &= \{a, [\} \\
\text{DIR}(RLI_s \rightarrow , I RLI_s) &= \{, \} \\
\text{DIR}(RLI_s \rightarrow \varepsilon) &= \{\} \\
\text{DIR}(I \rightarrow a) &= \{a\} \\
\text{DIR}(I \rightarrow L) &= \{\}
\end{aligned}$$

3.3.1 Resolución por aproximaciones sucesivas

La resolución por sustitución nos ha servido hasta ahora pero es un problema ante gramáticas con múltiples ecuaciones del tipo:

$$\begin{aligned}
\vec{x} &= \Phi [\vec{x}], \vec{x} \text{ in } U \text{ con } U \text{ finito.} \\
\vec{x} &\text{ es un punto fijo de } \Phi
\end{aligned}$$

Planteamos ahora como método la resolución por aproximaciones sucesivas.

podemos crear G' , la gramática transformada, sacando el factor común

$$A \rightarrow \alpha RA$$

$$RA \rightarrow \beta_0$$

...

$$RA \rightarrow \beta_n$$

e.g.

$$\begin{aligned} E_0 &\rightarrow \textcolor{red}{E}_1 + E_0 \\ E_0 &\rightarrow \textcolor{red}{E}_1 \end{aligned}$$

$$\begin{aligned} E_0 &\rightarrow E_1 R_0 \\ R_0 &\rightarrow + E_0 \\ R_0 &\rightarrow \varepsilon \end{aligned}$$

Veamos ahora como tratar la recursión. Sea G una gramática de la forma

$$A \rightarrow \textcolor{red}{A} \alpha_0$$

...

$$A \rightarrow \textcolor{red}{A} \alpha_n$$

$$A \rightarrow \beta_0$$

...

$$A \rightarrow \beta_n$$

podemos crear G' cambiando el tipo de recursión, de izquierdas a derechas.

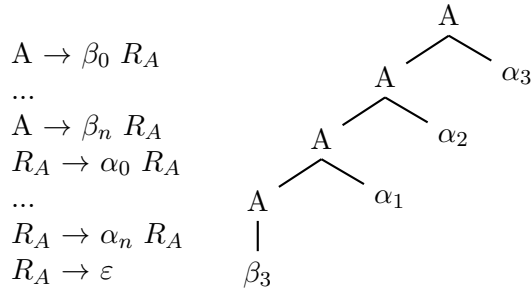


Figure 3.3.2.1 : Árbol de G .

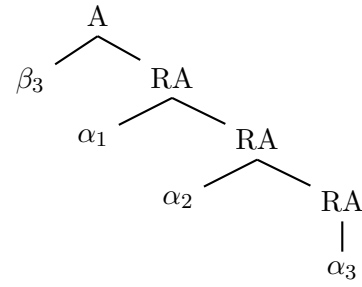


Figure 3.3.2.2 : Árbol de G' .

e.g.

$$\begin{aligned} E_0 &\rightarrow E_0 + E_1 \\ E_0 &\rightarrow \textcolor{red}{E}_1 \\ E_1 &\rightarrow E_2 * E_1 \\ E_1 &\rightarrow E_2 \\ E_2 &\rightarrow \text{num} \\ E_2 &\rightarrow (E_0) \end{aligned}$$

Recursión a izquierdas
 Factor común

$$\begin{aligned} E_0 &\rightarrow E_1 R_0 \\ R_0 &\rightarrow + E_1 R_0 \\ R_0 &\rightarrow \varepsilon \\ E_1 &\rightarrow E_2 R_1 \\ R_1 &\rightarrow * E_1 \\ R_1 &\rightarrow \varepsilon \\ E_2 &\rightarrow \text{num} \\ E_2 &\rightarrow (E_0) \end{aligned}$$

Por tanto la especificación se haria utilizando G pero implementamos G'.

Calculamos los primeros de G'

$$\begin{aligned} \text{PRIM}(E_0) &= \text{PRIM}(E_1) \oplus \text{PRIM}(R_0) \\ \text{PRIM}(R_0) &= \text{PRIM}(+) \oplus \text{PRIM}(E_1) \oplus \text{PRIM}(R_0) \cup \text{PRIM}(\varepsilon) \\ \text{PRIM}(E_1) &= \text{PRIM}(E_2) \oplus \text{PRIM}(R_1) \\ \text{PRIM}(R_1) &= \text{PRIM}(*) \oplus \text{PRIM}(E_1) \cup \text{PRIM}(\varepsilon) \\ \text{PRIM}(E_2) &= \text{PRIM}(\text{num}) \cup \text{PRIM}(() \oplus \text{PRIM}(E_0) \oplus \text{PRIM}()) \end{aligned}$$

	0	1	2	3	4
PRIM(E_0)	\emptyset	\emptyset	\emptyset	num (num (
PRIM(R_0)	\emptyset	$+\varepsilon$	$+\varepsilon$	$+\varepsilon$	$+\varepsilon$
PRIM(E_1)	\emptyset	\emptyset	num (num (num (
PRIM(R_1)	\emptyset	$*\varepsilon$	$*\varepsilon$	$*\varepsilon$	$*\varepsilon$
PRIM(E_2)	\emptyset	num (num (num (num (

Ahora los siguientes de G'

$$\begin{aligned} \text{SIG}(E_0) &= \{\vdash\} \cup \text{PRIM}()) \oplus \text{SIG}(E_2) \\ \text{SIG}(R_0) &= \text{PRIM}(\varepsilon) \oplus \text{SIG}(E_0) \cup \text{PRIM}(\varepsilon) \oplus \text{SIG}(R_0) \\ \text{SIG}(E_1) &= \text{PRIM}(R_0) \oplus \text{SIG}(E_0) \cup \text{PRIM}(R_0) \oplus \text{SIG}(R_0) \\ \text{SIG}(R_1) &= \text{PRIM}(\varepsilon) \oplus \text{SIG}(E_1) \\ \text{SIG}(E_2) &= \text{PRIM}(R_1) \oplus \text{SIG}(E_1) \end{aligned}$$

	0	1	2	3	4
SIG(E_0)	\emptyset	\vdash	\vdash	\vdash	\vdash
SIG(R_0)	\emptyset	\emptyset	\vdash	\vdash	\vdash
SIG(E_1)	\emptyset	$+$	$+\vdash$	$+\vdash$	$+\vdash$
SIG(R_1)	\emptyset	\emptyset	$+$	$+\vdash$	$+\vdash$
SIG(E_2)	\emptyset	$*$	$*+$	$*+\vdash$	$*+\vdash$

Por último empleamos los siguientes y los primeros para calcular los directores.

$$\begin{aligned} \text{DIR}(E_0 \rightarrow E_1 R_0) &= \{\text{num}, \{\} \\ \text{DIR}(R_0 \rightarrow + E_1 R_0) &= \{+\} \\ \text{DIR}(R_0 \rightarrow \varepsilon) &= \{\vdash, \{\} \\ \text{DIR}(E_1 \rightarrow E_2 R_1) &= \{\text{num}, \{\} \\ \text{DIR}(R_1 \rightarrow * E_1) &= \{*\} \\ \text{DIR}(R_1 \rightarrow \varepsilon) &= \{+, \vdash, \} \\ \text{DIR}(E_2 \rightarrow \text{num}) &= \{\text{num}\} \\ \text{DIR}(E_2 \rightarrow (E_0)) &= \{(\} \end{aligned}$$

Los directores para las gramáticas con múltiples producciones son disjuntos, por lo que no vamos a tener problemas de decisión.

3.3.3 Tabla de análisis sintáctico descendente

Sea G una gramática de la forma

$N \rightarrow B \text{ RN}$	$\text{DIR}(N \rightarrow B \text{ RN}) = \{0, 1\}$
$\text{RN} \rightarrow B \text{ RN}$	$\text{DIR}(\text{RN} \rightarrow B \text{ RN}) = \{0, 1\}$
$\text{RN} \rightarrow \varepsilon$	$\text{DIR}(\text{RN} \rightarrow \varepsilon) = \{\vdash\}$
$B \rightarrow 0$	$\text{DIR}(B \rightarrow 0) = \{0\}$
$B \rightarrow 1$	$\text{DIR}(B \rightarrow 1) = \{1\}$

La tabla de análisis de G sería la siguiente.

	0	1	\vdash
N	$N \rightarrow B \text{ RN}$	$N \rightarrow B \text{ RN}$	
RN	$\text{RN} \rightarrow B \text{ RN}$	$\text{RN} \rightarrow B \text{ RN}$	$\text{RN} \rightarrow \varepsilon$
B	$B \rightarrow 0$	$B \rightarrow 1$	

Cada celda indica qué producción aplicar al detectar cierto símbolo, si una celda está vacía daría error.

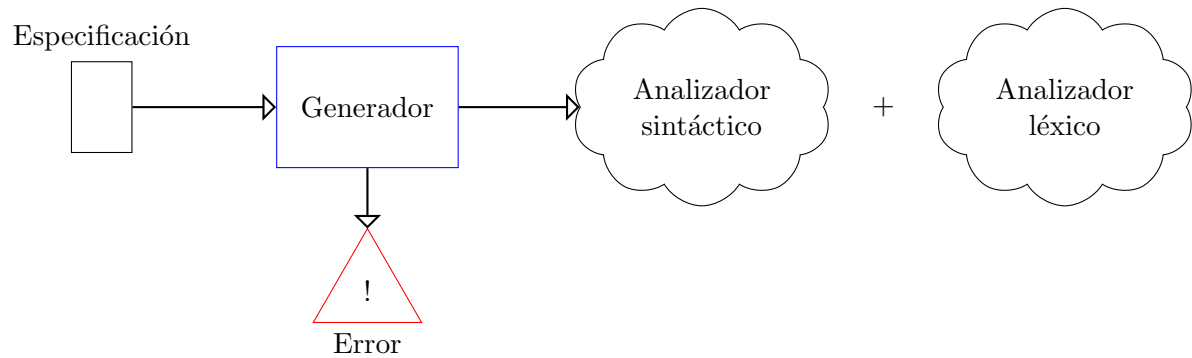
La gramática es $\text{LL}(1)$ si en cada celda hay a lo sumo una producción.

Se puede implementar un programa genérico que utilice la tabla para implementar el analizador, se llama analizador sintáctico descendente predictivo no-recursivo, o dirigido por tabla.

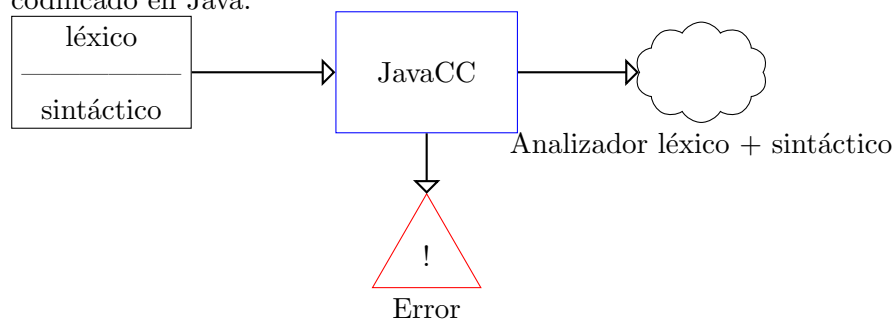
e.g. Queremos reconocer $010\vdash$, utilizando la tabla de análisis de G .

Pila (\leftarrow)	Entrada	Acción
N	010 \vdash	$N[0] \equiv N \rightarrow B \text{ RN}$
B RN	010 \vdash	$B[0] \equiv B \rightarrow 0$
0 RN	010 \vdash	Consume 0
RN	10 \vdash	$\text{RN}[1] \equiv \text{RN} \rightarrow B \text{ RN}$
B RN	10 \vdash	$B[1] \equiv B \rightarrow 1$
1 RN	10 \vdash	Consume 1
RN	0 \vdash	$\text{RN}[0] \equiv \text{RN} \rightarrow B \text{ RN}$
B RN	0 \vdash	$B[0] \equiv B \rightarrow 0$
0 RN	0 \vdash	Consume 0
RN	\vdash	$\text{RN}[\vdash] \equiv \text{RN} \rightarrow \varepsilon$
	\vdash	Pila vacía \rightarrow aceptar

3.3.4 Implementación mediante herramientas



En nuestro caso el generador será JavaCC que dará lugar al analizador sintáctico codificado en Java.



Soporta gramáticas en formato EBNF (Extended BNF) aunque no lo vamos a usar. BNF (Backus Normal Form) es el formato de reglas gramaticales que hemos empleado hasta ahora.

e.g. BNF

$$E_0 \rightarrow E_1 ((\backslash + | \backslash -) RE_0)^*$$

Para compilar usar: `java -cp javacc.jar javacc input.jj`

Genera varios ficheros, *TokenManager.java es el analizador.

Esta herramienta soporta gramáticas LL(k), pueden tener k símbolos adelantados. En el caso de nuestras gramáticas la k es 1.

LL(*) son gramáticas que pueden utilizar cadenas arbitrariamente largas, la herramienta usada para esta tipología es ANTLR.

3.4 Análisis ascendente

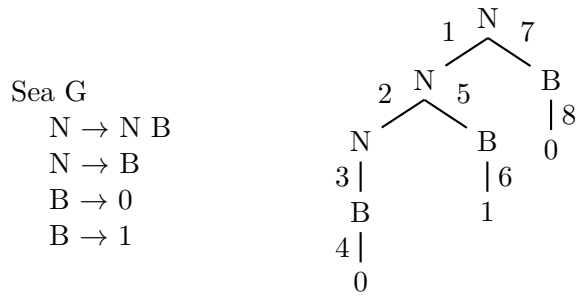


Figure 3.4.1 : Descendente

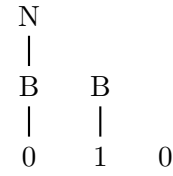
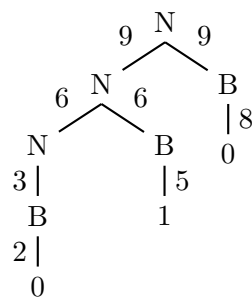


Figure 3.4.2 : Ascendente

El modelo que vamos a seguir es el analizador por desplazamiento-reducción.

N	Pila (\leftarrow)	Entrada	Acción
0		010	
1	0	10	desplazamiento
2	B	10	reducción $B \rightarrow 0$
3	N	10	reducción $N \rightarrow B$
4	N 1	0	desplazamiento
5	N B	0	reducción $B \rightarrow 1$
6	N	10	reducción $N \rightarrow N B$
7	N 0		desplazamiento
8	N B		reducción $B \rightarrow 0$
9	N		reducción $N \rightarrow N B$

La entrada se ha reducido al símbolo inicial por tanto aceptamos la cadena. El siguiente árbol muestra el orden de expansión, los numeros indican el N de la tabla que lo ha generado.



Veamos ahora una derivación más a la derecha de la sentencia.

$N \Rightarrow NB \Rightarrow N0 \Rightarrow NB0 \Rightarrow N10 \Rightarrow B10 \Rightarrow 010$

La secuencia, leída de derecha a izquierda, se corresponde con la construcción ascendente del árbol.

Los posibles contenidos de la pila se llaman prefijo viable.

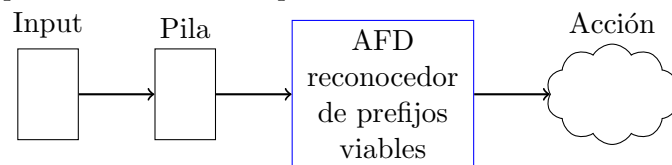
- Asidero: sufijo del prefijo viable que se sustituye en una reducción.

e.g. $n + n + n$

Los asideros serán los marcados de color rojo.

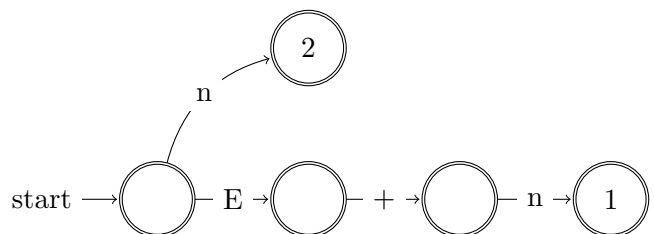
	PILA
	n
	E
$E \rightarrow E + n$	E +
$E \rightarrow n$	E + n
	E
	E +
	E + n
	E

El conjunto de prefijos viables para una gramática G es un conjunto regular, es decir, puede ser reconocido por una AFD.



e.g. AFD

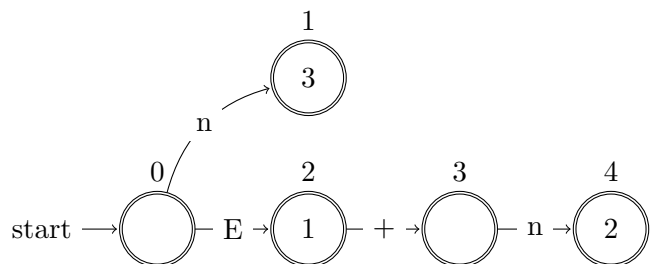
- (1) $E \rightarrow E + n$
- (2) $E \rightarrow n$



El número en un estado indica cuál de las reglas reducirá, los estados sin número representan la acción de desplazar.

e.g. AFD

- (1) $E' \rightarrow E \vdash$
- (2) $E \rightarrow E + n$
- (3) $E \rightarrow n$



	\vdash	n	$+$	E
0		d 1		2
1	$r E \rightarrow E + n$	$r E \rightarrow E + n$	$r E \rightarrow E + n$	
2	a		d 3	
3		d 4		
4	$r E \rightarrow E + n$	$r E \rightarrow E + n$	$r E \rightarrow E + n$	

Solo los terminales dictan acciones, por eso la columna de la E solo tiene transiciones. Las casillas vacías en la tabla indican los errores sintácticos.

Pila	Entrada	Estado AFD	Acción
$n + n + n$	\vdash	0	desplazar
n	$+ n + n$	1	reducir $E \rightarrow n$
E	$+ n + n$	2	desplazar
$E +$	$n + n$	3	desplazar
$E + n$	$+ n$	4	reducir $E \rightarrow E + n$
E	$+ n$	2	desplazar
$E +$	n	3	desplazar
$E + n$	\vdash	4	reducir $E \rightarrow E + n$
E	\vdash	2	aceptar

Estos autómatas pueden optimizarse, ahora la pila tiene siempre en la cima el estado en el que está. Para que el cambio se aprecie mejor los estados se marcaran en rojo.

Pila	Entrada	Estado AFD	Acción
0	$n + n + n$	0	desplazar
0 n 1	$+ n + n$	1	reducir $E \rightarrow n$
0 E 2	$+ n + n$	2	desplazar
0 E 2 + 3	$n + n$	3	desplazar
0 E 2 + 3 n 4	$+ n$	4	reducir $E \rightarrow E + n$
0 E 2	$+ n$	2	desplazar
0 E 2 + 3	n	3	desplazar
0 E 2 + 3 n 4	\vdash	4	reducir $E \rightarrow E + n$
0 E 2	\vdash	2	aceptar

Para que siga funcionando las reducciones deberan desapilar varios elementos.

$$0 \text{ n } 1 \rightarrow 0 \rightarrow 0 \text{ E } 2$$

3.4.1 Métodos de construcción de AFDs

Concretamente queremos implementar AFDs que reconozcan prefijos viables.

En las tablas puede haber más de una acción por celda, se denomina conflicto, esto provocará que el análisis no sea determinista y el método no funcionará.

Los métodos son los siguientes, de menos a más genéricos.

- LR(0) \rightarrow gramáticas LR(0)
- SLR(1) \rightarrow gramáticas SLR(1)
- LALR(1) \rightarrow gramáticas LALR(1)
- LR(1) \rightarrow gramáticas LR(1)

De todos ellos el más empleado es LALR(1).

Existen dos tipos de conflictos.

- Desplazamiento-reducción.
Si en una celda nos encontramos como acciones posibles:
d X
r A $\rightarrow \alpha$
- Reducción-reducción.
Si en una celda nos encontramos como acciones posibles:
r A $\rightarrow \alpha$
r B $\rightarrow \beta$

Gramática LR(0), es posible determinar que acción realizar únicamente a partir del prefijo viable.

A \rightarrow a
A \rightarrow b
Prefijos viables $\equiv \{\varepsilon, A, a, b\}$

Gramática LR(1), emplea el prefijo viable y un símbolo anticipado de la entrada.

A \rightarrow a
A \rightarrow a b
Es LR(1) pero no LR(0).

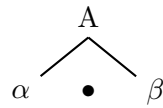
3.4.1.1 Método LR(0)

A $\rightarrow \alpha \bullet \beta$, donde α y β se denominarán elementos LR(0).

En el contexto de una estructura A, como la descrita previamente.

- Hemos "visto" α

- Es posible ver β y, potencialmente, reducir por $A \rightarrow \alpha \beta$

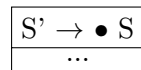


Los estados del AFD LR(0) son conjuntos de elementos LR(0), cada conjunto será válido para cierto prefijo viable.

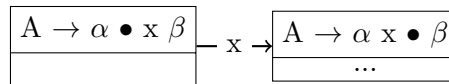
Para contruirlo.

- Partimos de una serie de elementos iniciales, denominados elementos nucleares.
 $A_0 \rightarrow \alpha_0 \beta_0 \dots A_n \rightarrow \alpha_n \beta_n$
- Siguiendo la regla de cierre habrá otros elementos, denominados elementos no nucleares.
 $A \rightarrow \alpha \bullet B \beta \Rightarrow B \rightarrow \bullet \gamma$

Estado inicial

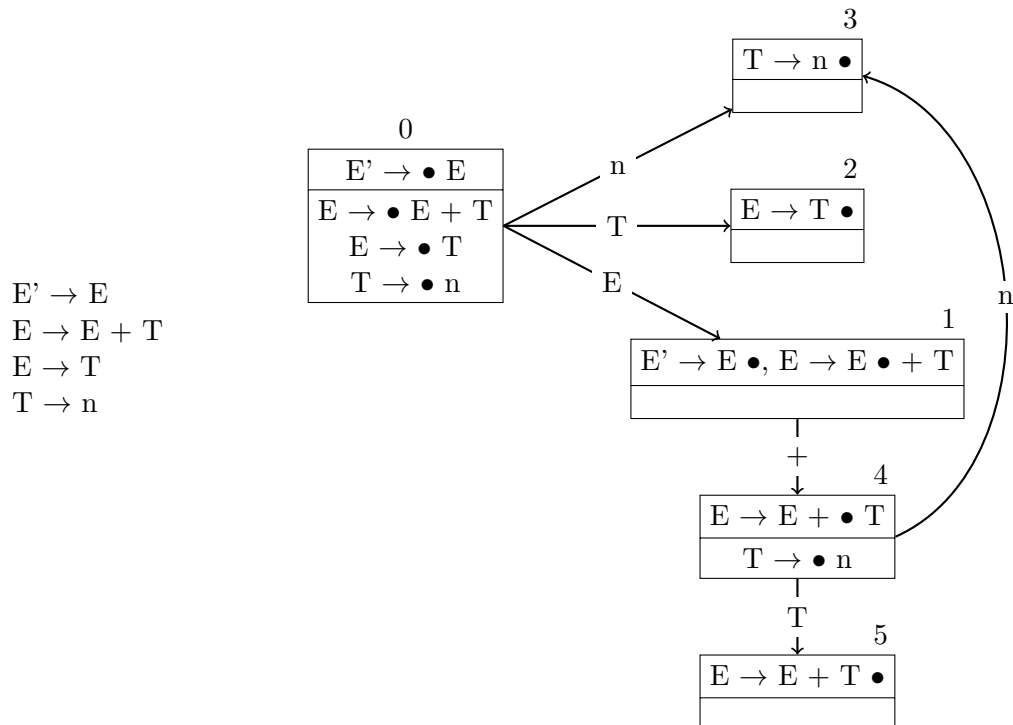


Estado i



... serán las posibles producciones para continuar el análisis

e.g.

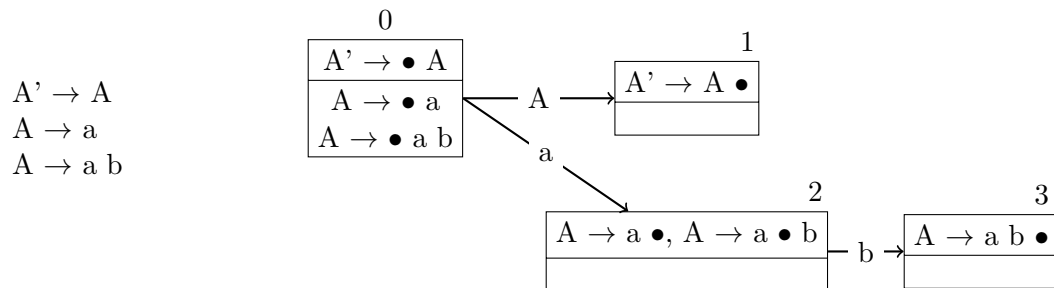


Del automata anterior podemos obtener la siguiente tabla.

	n	⊢	+	E	T
0	d 3			1	2
1		d 4	a		
2	r E → T	r E → T	r E → T		
3	r T → n	r T → n	r T → n		
4	d 3				5
5	r E → E + T	r E → E + T	r E → E + T		

Una gramática LR(0) es aquella para la cual el método LR(0) genera tablas sin conflictos. Entendemos conflicto como una celda no definida de manera unívoca.

e.g. Conflicto desplazamiento-reducción

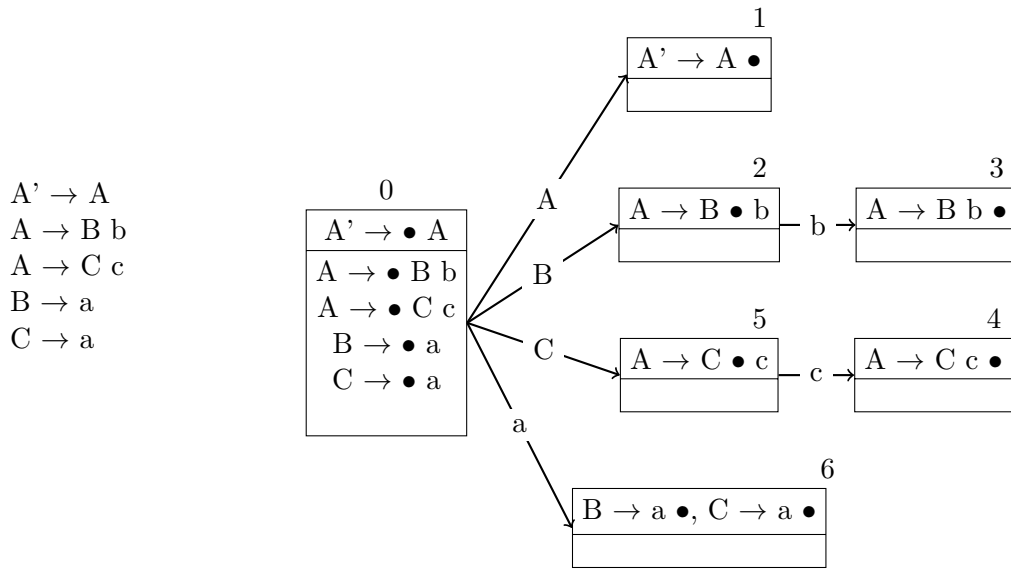


El problema se da en el estado 2 porque tiene dos posibles acciones, conflicto de desplazamiento-reducción.

	a	b	⊢	A
0	d 2			1
1		a		
2	r A → a	r A → a		
3	r A → a b	r A → a b		

La gramática no es LR(0).

e.g. Conflicto reducción-reducción



	a	b	c	⊢	A	B	C
0	d 6				1	2	5
1		a					
2		d 3					
3	r A → B b	r A → B b	r A → B b	r A → B b			
4	r A → C c	r A → C c	r A → C c	r A → C c			
5			d 4				
6	r B → A r C → a	r B → A r C → a	r B → A r C → a	r B → A r C → a			

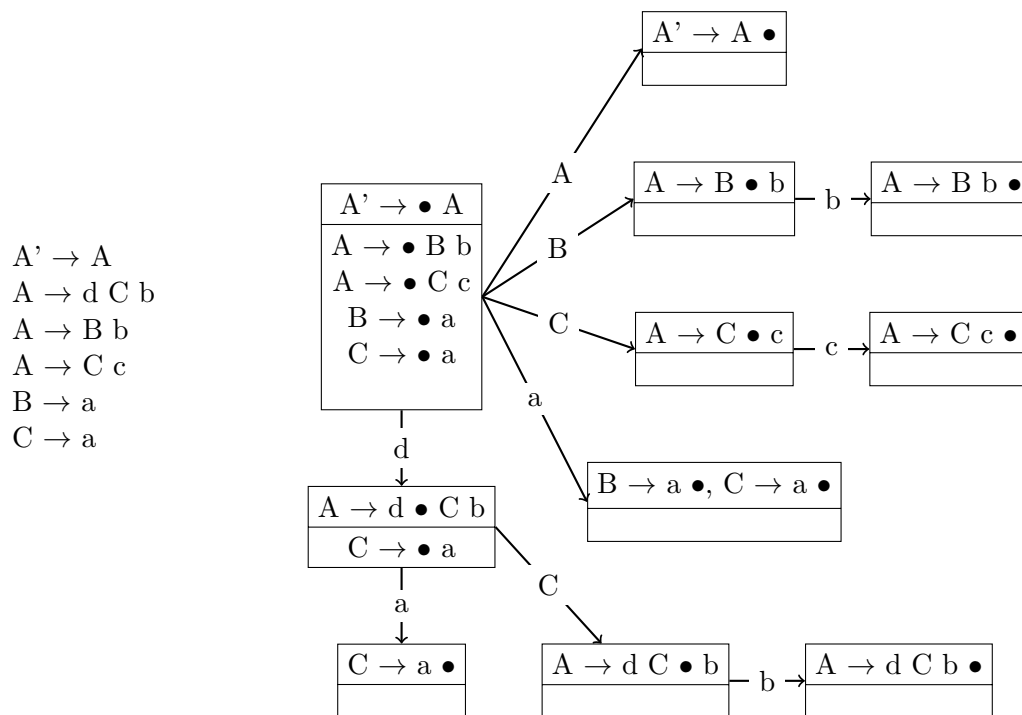
3.4.1.2 Método SLR(1)

Consiste en utilizar los siguientes en prefijos viables para no realizar reducciones inútiles. Al aplicarlo la tabla anterior quedaria de la siguiente manera.

	a	b	c	⊢	A	B	C
0	d 6				1	2	5
1		a					
2		d 3					
3				r A → B b			
4				r A → C c			
5			d 4				
6	r B → A	r C → a					

No hay conflictos por lo que la gramática es SLR(1).

e.g. La siguiente gramática supone un problema incluso para SLR(1), no se suelen encontrar casos de este tipo.



Como hay un conflicto de reducción-reducción la gramática no es SLR(1).

3.4.1.3 Método LR(1)

En este método también se indica el símbolo por el cual se aplica la reducción.

$$[A \rightarrow \alpha \bullet \beta, a]$$

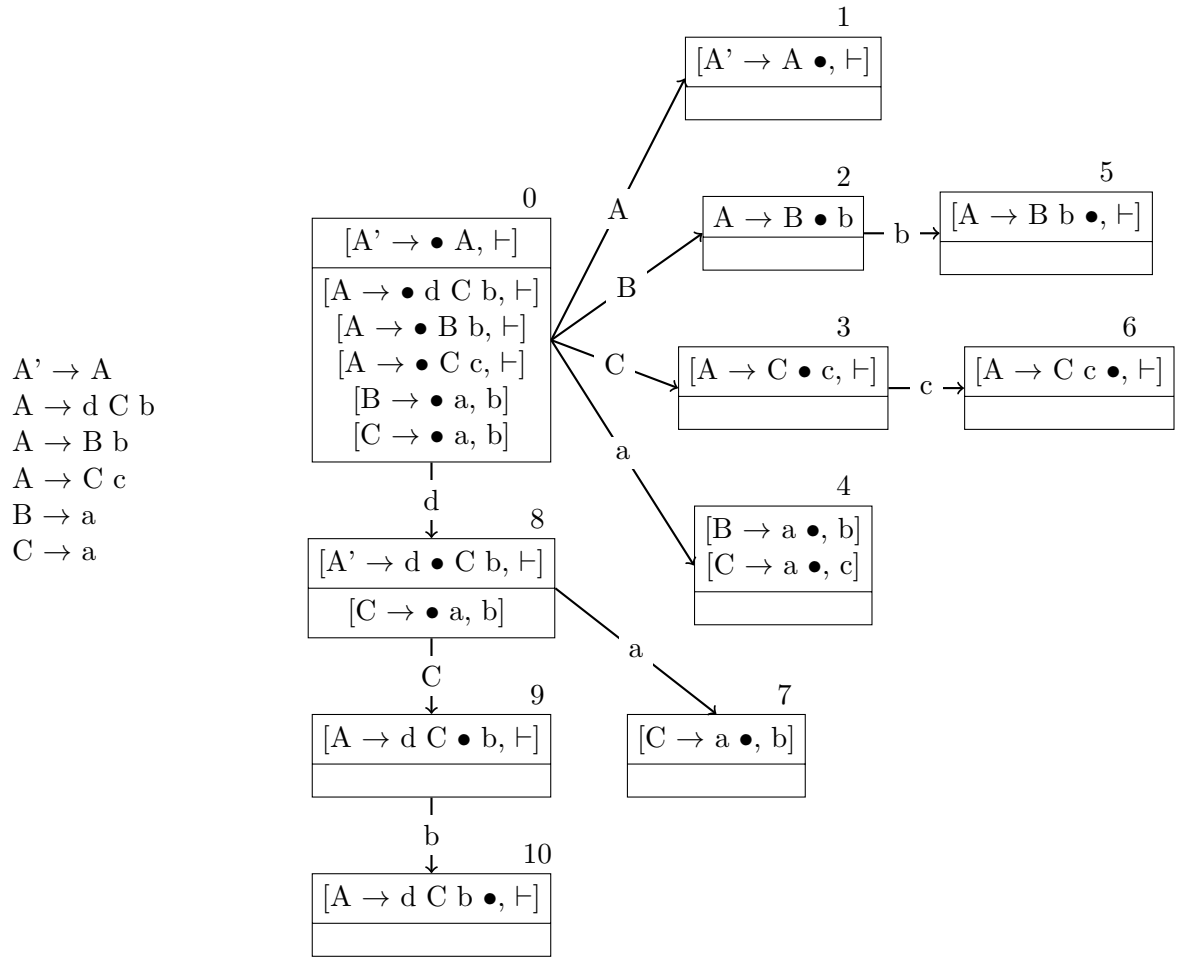
$[S' \rightarrow \bullet S, \vdash]$	$[A \rightarrow \alpha \bullet B \beta, a]$
	$[B \rightarrow \bullet \gamma, b]$

Se expandirán producciones en un nodo para toda producción de la forma $B \rightarrow \gamma$ y para todo $b \in \text{PRIM}(B) \oplus \{a\}$.

e.g.

...	
$[A_0 \rightarrow \alpha_0 \bullet x \beta_0, b_0]$	$\text{---} x \rightarrow$
$[A_n \rightarrow \alpha_n \bullet x \beta_n, b_n]$	
	$[A_0 \rightarrow \alpha_0 x \bullet \beta_0, b_0]$ $[A_n \rightarrow \alpha_n x \bullet \beta_n, b_n]$...

e.g.



	a	b	c	d	⊢	A	B	C
0	d 4			d 8		1	2	3
1					a			
2		d 5						
3			d 6					
4		r B → a	r C → a					
5					r A → B b			
6					r A → C c			
7		r C → a						
8	d 7							9
9		d 10						
10					r A → d C B			

Como la tabla no tiene conflictos podemos concluir que esta gramática es LR(1).

En el ejemplo no ha ocurrido pero el número de estados podría verse incrementado

$[C \rightarrow a \bullet, b]$	$[C \rightarrow a \bullet, \vdash]$

Como puede verse ambos estados coincidirían si utilizáramos el método LR(0).

3.4.1.4 Método LALR(1)

Para mejorar el método anterior podemos fusionar los estados.

$[C \rightarrow a \bullet, b]$
$[C \rightarrow a \bullet, \vdash]$

Como lo único que se ha realizado ha sido dicha fusión podemos asegurar que toda gramática LALR(1) es necesariamente LR(1).

Sin embargo, al realizar dicha fusión pueden aparecer conflictos por lo que no toda gramática LR(1) tiene que ser LALR(1).

- Conflicto desplazamiento-reducción.

$[A \rightarrow \alpha \bullet a \beta, -]$	$[A \rightarrow \alpha \bullet a \beta, -]$
$[B \rightarrow \gamma \bullet, -]$	$[A \rightarrow \alpha \bullet a \beta, -]$
	$[A \rightarrow \alpha \bullet a \beta, -]$
	$[B \rightarrow \gamma \bullet, a]$
	$[B \rightarrow \gamma \bullet, -]$

No pueden aparecer conflictos NUEVOS de este tipo, pero si los había en LR(1) seguirán estando.

- Conflicto reducción-reducción.

Pueden aparecer conflictos nuevos de este tipo.

$[A \rightarrow \alpha \bullet \beta, a]$

$[A \rightarrow \alpha \bullet \beta, b]$

$[A \rightarrow \alpha \bullet \beta, c]$

Lo podemos escribir de la siguiente forma.

$[A \rightarrow \alpha \bullet \beta, \{a, b, c\}]$

$[A \rightarrow \alpha \bullet, \psi_0]$
$[B \rightarrow \beta \bullet, \psi_1]$

$$\begin{aligned}\psi_0 \cap \psi_1 &= \emptyset \\ \psi'_0 \cap \psi'_1 &= \emptyset \\ (\psi_0 \cup \psi'_0) \cap (\psi_1 \cup \psi'_1) &\neq \emptyset\end{aligned}$$

$[A \rightarrow \alpha \bullet, \psi'_0]$
$[B \rightarrow \beta \bullet, \psi'_1]$

Por tanto al fusionar pueden aparecer conflictos de este tipo.

e.g.

$B \rightarrow A b$
 $C \rightarrow B b$
 $C \rightarrow A c$
 $C \rightarrow B b$
 $A \rightarrow a$
 $B \rightarrow a$

$[A \rightarrow a \bullet, b]$
$[B \rightarrow a \bullet, c]$

$[A \rightarrow a \bullet, c]$
$[B \rightarrow a \bullet, b]$

Podemos fusionar ambos estados dando lugar al siguiente.

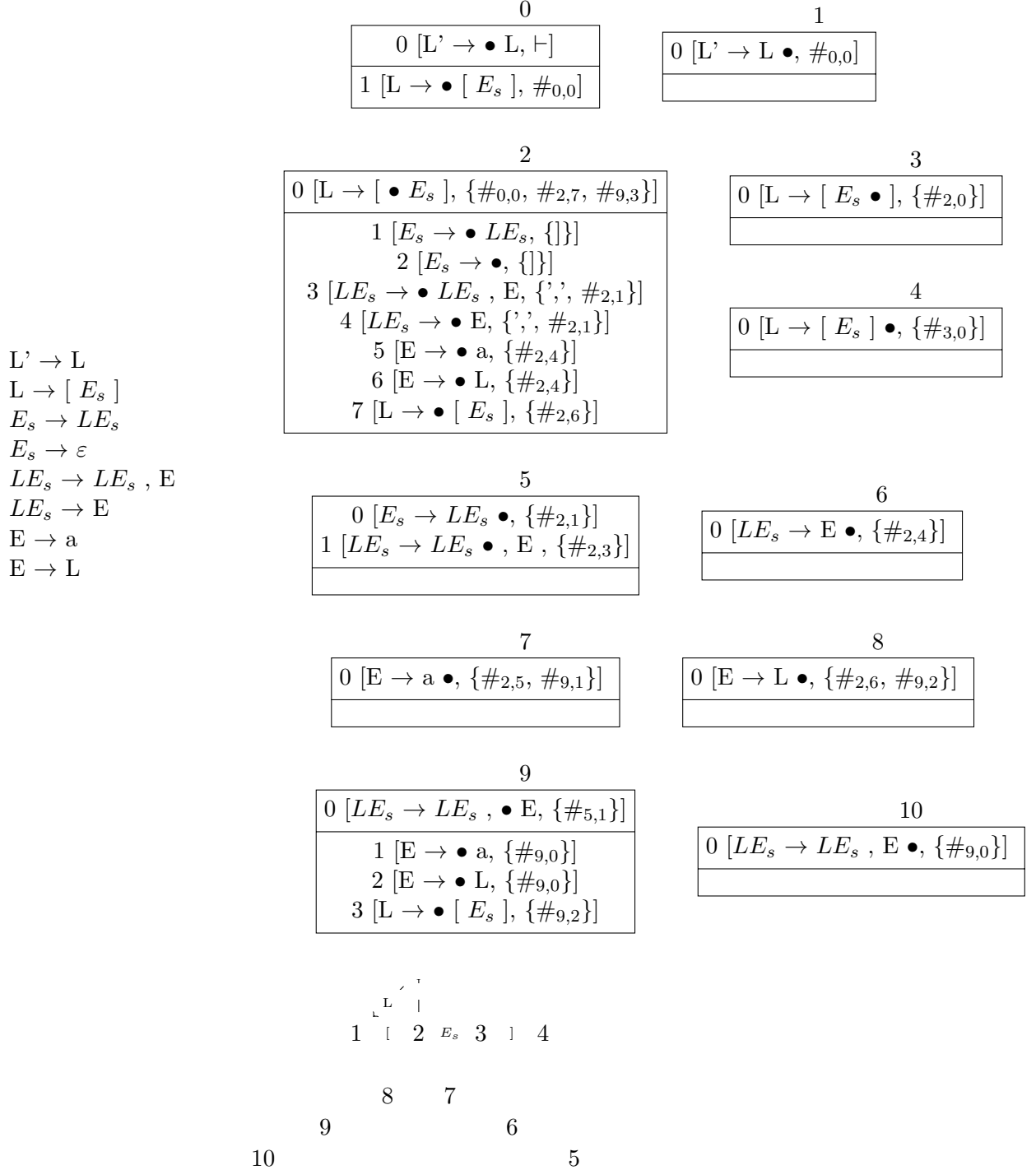
$[A \rightarrow a \bullet, \{b, c\}]$
$[B \rightarrow a \bullet, \{b, c\}]$

Al fusionarlos se han generado conflictos, probando así que no toda gramática LR(1) es LALR(1).

TODO: Pasar automata ejemplo grande.

Construcción directa del autómata LALR(1)

$\#_{n,m} \equiv$ preanálisis del elemento m del estado n .



4 Semántica estática

5 Traducción