

Teoría de Lenguajes

Analizadores Sintácticos Descendentes

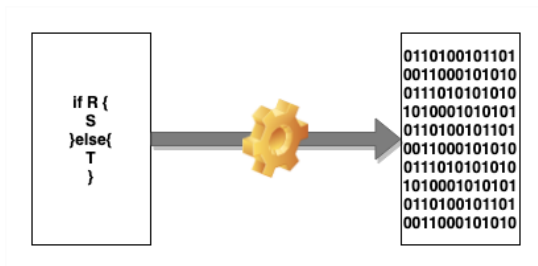
Natalia Pesaresi

DC-UBA

1er. Cuatrimestre 2020

Compiladores

- Tienen como objetivo convertir un código fuente de programa a otro lenguaje: Algunos a código máquina ejecutable, o otros intermedio interpretable (e.g., Java o .NET)
- Tendrá que poder detectar errores en el código, y poder marcarlos de forma de facilitar su corrección



Compiladores

Proceso de un Compilador

- ➊ **Análisis Léxico:** Se encarga de procesar las cadenas y generar los lexemas correspondientes. Por cada lexema se generará la lista de tokens (par token-valor).
`int x = 0;`
- ➋ **Análisis Sintáctico:** A partir de la salida del analizador léxico se generará el árbol sintáctico correspondiente a la código recibido
`if(){...}else{...}`
- ➌ **Análisis Semántico:** Se genera información para las etapas subsiguientes, y en particular se encargará del chequeo de tipos
- ➍ **Generación de Código Intermedio**
- ➎ **Optimización**
- ➏ **Generación de Código ejecutable**

El problema de interés

- Dado un lenguaje L definido por la gramática

$$G = \langle \{E, T, F\}, \{+, *, \textit{num}, (,)\}, P, E \rangle$$

$$E \rightarrow E + T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow F$$

$$F \rightarrow \textit{num}$$

$$F \rightarrow (E)$$

- y dada la cadena

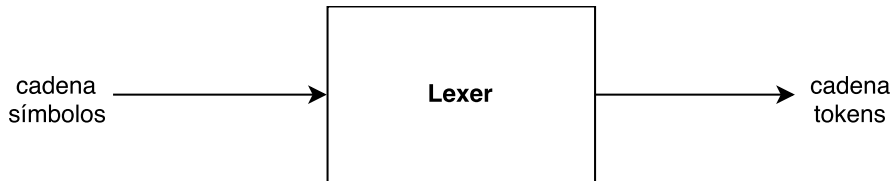
$$\alpha : (14 + 6) * 2$$

- ¿ $\alpha \in L$?

¿Cómo sabe la computadora qué es un **num**?

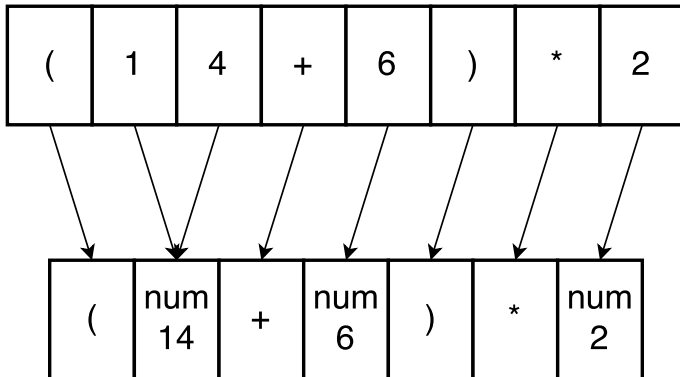
- ¿Reemplazamos al terminal *num* por los terminales 14, 6 y 2?
- ¿Reemplazamos la producción $F \rightarrow num$ por las producciones $F \rightarrow 14$, $F \rightarrow 6$ y $F \rightarrow 2$?
- ¿Podemos hacer esto para todos los números?
- ¡Son infinitos!
- Ahí es donde entran en juego los **analizadores léxicos** (lexers)

Analizador Léxico (Lexer)

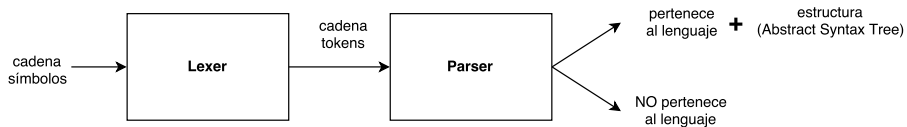


- Usa expresiones regulares
- Agrupa símbolos de la entrada en tokens
- Un token tiene un tipo y un valor

Ejemplo

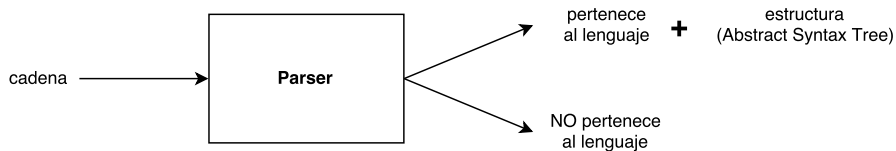
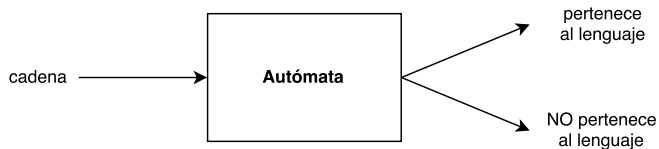


Lexer + Parser

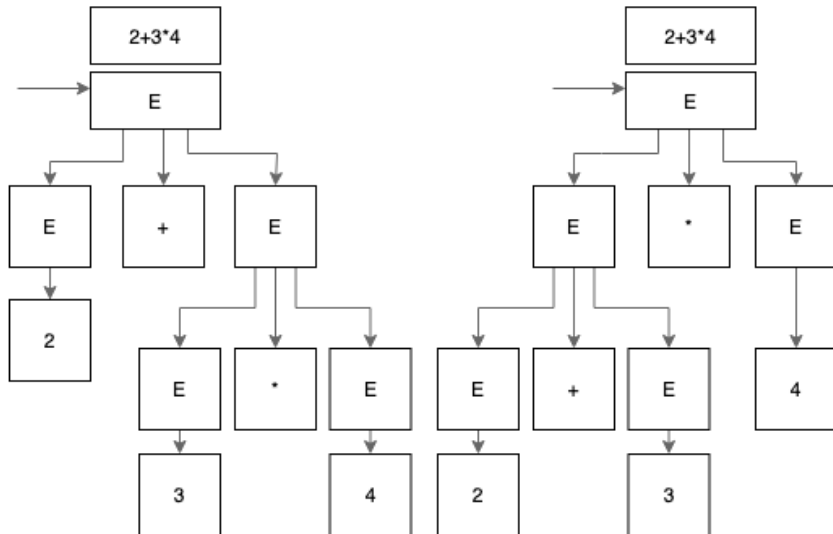


- Cadena de símbolos: $(14+6)*2$
- Cadena de tokens: `(num+num)*num`
- Árbol sintáctico abstracto

Autómata vs. Parser



AST



Tipos de Parsers(Analizadores Sintácticos)

- Descendentes (top-down):
 - ▶ Construyen el árbol sintáctico desde la raíz hasta las hojas (la secuencia de tokens).
 - ▶ Ejemplos: DFS, parser predictivos recursivos, LL(1), LL(k).
- Ascendentes (bottom-up)
 - ▶ Construyen el árbol sintáctico desde las hojas (la secuencia de tokens) hasta la raíz (el símbolo distinguido).
 - ▶ Ejemplos: parsers shift-reduce, LR(0), SLR, LALR, LR(K)
- Otros
 - ▶ Generalized LR (GLR), Earley, CYK

Parsers descendentes(top-down)

- Dada la siguiente gramática:

$$S \rightarrow A$$

$$S \rightarrow B$$

$$A \rightarrow a$$

$$B \rightarrow b$$

- ¿Cómo genero la cadena a?

Parsers predictivos descendentes

- Parten del símbolo distinguido. S en nuestro ejemplo.
- Y la secuencia que vamos a iterar es:
 - ▶ Identificar, el **token corriente**, **tc**. Cual es el token que queremos reconocer en este momento.
 - ▶ Y que producción lo reconoce. Para elegir la producción, calculamos el conjunto de **símbolos directrices**.

Cadena: a.

$$S \rightarrow A$$

$$S \rightarrow B$$

$$A \rightarrow a$$

$$B \rightarrow b$$

¿Cómo elegir una producción?

Primeros, Siguiente, Símbolos Directrices

Vamos a calcular dos funciones y luego SD.

Primeros: de una cadena, es un conjunto que contiene los símbolos terminales con los que 'comienza'.

$$\text{Primeros}(\alpha) : (V_N \cup V_T)^* \rightarrow \mathbf{P}(V_T)$$

$$\text{Primeros}(\alpha) = \{t \in V_T \mid \alpha \xRightarrow{*} t\beta\}$$

Ej. Gramática: *(las minúsculas son tokens)*

$$S \rightarrow A|B|C$$

$$A \rightarrow a|Cm$$

$$B \rightarrow abBb|A$$

$$C \rightarrow \lambda|c$$

$$P(a) = \{a\} \quad P(c) = \{c\} \quad P(abBb) = \{a\} \quad P(\text{lambda}) = \emptyset$$

$P(Cm) = \{c, m\}$, m se agrega porque C es anulable.

$$P(A) = P(a) \cup P(Cm) = \{a, c, m\}$$

¿Cómo elegir una regla?

Primeros, Siguiente, Símbolos Directrices

Siguientes: de un **símbolo no terminal** es el conjunto de tokens, que aparecen 'después' de este símbolo entre las producciones.

$$\text{Siguientes}(N) : V_N \rightarrow \mathbf{P}(V_T)$$

$$\text{Siguientes}(N) = \{t \in V_T \mid S \xRightarrow{*} \dots Nt\dots\}$$

Ej. Gramática:

$$S \rightarrow A|B$$

$$A \rightarrow a|Cm$$

$$B \rightarrow abBb|A|C$$

$$C \rightarrow \lambda|c$$

$$S(C) = \{m\} \text{ y } S(B) = \{b\}$$

¿Cómo elegir una regla?

Primeros, Siguiente, Símbolos Directrices

Función que nos da el conjunto de de Símbolos Directrices (SD):

$$SD(N \rightarrow \beta) = \left\{ \begin{array}{ll} \text{Primeros}(\beta) & \text{Si } \lambda \notin \beta \\ \text{Primeros}(\beta) \cup \text{Siguietes}(N) & \text{Si } \lambda \in \beta \end{array} \right\}$$

Ej. Gramática:

$$S \rightarrow A|B$$

$$A \rightarrow a|Cm$$

$$B \rightarrow abBb|A|C$$

$$C \rightarrow \lambda|c$$

$$SD(B \rightarrow abBb) = P(abBb) = \{a\}$$

$$SD(B \rightarrow C) = P(C) \cup S(B) = \{c\} \cup \{b\} = \{c, b\}$$

Implementaciones de Parsers predictivos descendentes

Tenemos 2 maneras de implementar un parser descendente.

- 1 Programable: uso de pila explícito
- 2 Seguimiento: otro con uso de pila implícito

Programable (1)

Algoritmo general:

- se crea un procedimiento por cada no terminal (ej. S, A, B, C).
Main es el correspondiente al símbolo distinguido.
- dentro de cada **procedimiento** se hacen 3 cosas:
 - ▶ se invocan los procedimientos de los no terminales incluidos
 - ▶ se busca los terminales esperados (`match`)
 - ▶ o no se hace nada (si es un λ)
- si hay más de una producción para un mismo no terminal, se agregan condicionales para ejecutar sólo el procedimiento que contenga al token actual (tc) en sus SD's.
 - ▶ se agrega un `else final` con `error()`, para cubrir el caso en que el tc no está en los SD de ninguna de las reglas posibles

Ej. $B \rightarrow abBb$

Seguimiento (2)

- Se crea una tabla: cadena || pila || producción
- se escribe la cadena a reconocer seguida de \$ en la primer columna.
- Se escribe el símbolo inicial(S) y \$, con el inicial como tope de pila.
- Iteración:
 - ▶ Si el tope de pila es un es un símbolo terminal y coincide con el tc, acción match y despilar.
 - ▶ Si el tope de pila es un es un símbolo terminal y **no** coincide con el tc, ERROR().
 - ▶ Si el tope de pila es un no temrinal, elegir la producción que tiene al tc entre sus símbolos directrices.
 - ▶ Si no, ERROR().

Algunos problemas y cómo resolverlos

- Ambigüedad por SDs no disjuntos
 - ▶ factorización a la izquierda (*left-factorization*)
Ej. $B \rightarrow abBb|am$
- Recursión a izquierda
 - ▶ eliminación de la recursión inmediata
Ej. $B \rightarrow Ba|s|t$
 - ▶ eliminación de la recursión no inmediata
Ej. $B \rightarrow Ca|d$ y $C \rightarrow Bf$

Factorización a la izquierda

Si tenemos producciones de la forma

$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ donde $\alpha \neq \lambda$, las reemplazamos por

$$\begin{aligned} A &\rightarrow \alpha A' \mid \delta_1 \mid \delta_2 \mid \dots \mid \delta_k \\ A' &\rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n \end{aligned}$$

Ej. $B \rightarrow abBb \mid am \Rightarrow B \rightarrow aB' \quad \text{y} \quad B' \rightarrow bBb \mid m$

Eliminación de la recursión inmediata

Si tenemos producciones de la forma $A \rightarrow A\alpha \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_k$ las reemplazamos por

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_k A' \\ A' &\rightarrow \alpha A' \mid \lambda \end{aligned}$$

Ej. $B \rightarrow Ba \mid s \mid t \Rightarrow B \rightarrow sB' \mid tB' \quad y \quad B' \rightarrow a \mid \lambda$

Eliminación de la recursión no inmediata

- 1 Numerar los no terminales A_1, \dots, A_n
- 2 Para $i \rightarrow 1$ a n
 - 1 Para $j \rightarrow 1$ a $i - 1$
 - 1 Reemplazar $A_i \rightarrow A_j \gamma$ por $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$
(donde $A_j \rightarrow \delta_1 \mid \dots \mid \delta_k \in P$)
 - 2 Fin Para
 - 3 Eliminar la recursión inmediata de A_i
- 3 Fin Para

Ej. $B \rightarrow Ca|d \quad y \quad C \rightarrow Bf$
 $B \rightarrow Bfa|d$
 $B \rightarrow dB' \quad y \quad B' \rightarrow fa|\lambda$

Ejemplo

Dada la siguiente gramática para la declaración de tipos en un lenguaje de programación:

$G_1 = \langle \{T, S\}, \{\text{array}, [,], \text{of}, \dots, \text{int}, \text{char}, \text{num}, \uparrow\}, T, P \rangle$, con P :

$$T \longrightarrow S \mid \uparrow S \mid \text{array of } T \mid \text{array } [S] \text{ of } T$$

$$S \longrightarrow \text{int} \mid \text{char} \mid \text{num}.. \text{num}$$

- Hacer un parser descendente recursivo para esta gramática. Si no es posible, corregirla previamente
- Realizar el seguimiento de como se reconoce una cadena.

Ejemplo

$$\begin{aligned}
 T &\longrightarrow S \mid \uparrow S \mid \text{array of } T \mid \text{array } [S] \text{ of } T \\
 S &\longrightarrow \text{int} \mid \text{char} \mid \text{num}..\text{num}
 \end{aligned}$$

Queremos saber,

- ¿La cadena $\text{array } [\text{num}..\text{num}] \text{ of } \uparrow \text{char} \in L(G_1)$?
- ¿Es posible hacer un parser descendente recursivo para esta gramática?

No porque hay solapamiento entre los SDs de dos reglas del no terminal T

Ejemplo: Gramática corregida

Aplicando factorización a la izquierda, para crear G'_1 :

$$T \rightarrow \text{array of } T \mid \text{array } [S] \text{ of } T \Rightarrow$$

$$T \rightarrow \text{array} T'$$

$$T' \rightarrow \text{of } T \mid [S] \text{ of } T$$

entonces,

$$G'_1 = \langle \{T, T', S\}, \{\text{array}, [,], \text{of}, \dots, \text{int}, \text{char}, \text{num}, \uparrow\}, T, P \rangle, \text{ con } P:$$

$$T \longrightarrow S \mid \uparrow S \mid \text{array } T'$$

$$T' \longrightarrow \text{of } T \mid [S] \text{ of } T$$

$$S \longrightarrow \text{int} \mid \text{char} \mid \text{num} \dots \text{num}$$

Extendemos la gramática con,

$$S_0 \rightarrow T\$$$

Ejemplo: Gramática corregida y SDs

Calculemos los Símbolos directrices de todas las producciones.

$$G'_1 = \langle \{S_0, T, T', S\}, \{\text{array}, [,], \text{of}, \dots, \text{int}, \text{char}, \text{num}, \uparrow\}, S_0, P \rangle$$

Producción		SD
S_0	$\longrightarrow T$	$\{\text{int}, \text{char}, \text{num}, \uparrow, \text{array}\}$
T	$\longrightarrow S$	$\{\text{int}, \text{char}, \text{num}\}$
T	$\longrightarrow \uparrow S$	$\{\uparrow\}$
T	$\longrightarrow \text{array } T'$	$\{\text{array}\}$
T'	$\longrightarrow \text{of } T$	$\{\text{of}\}$
T'	$\longrightarrow [S] \text{ of } T$	$\{[\]\}$
S	$\longrightarrow \text{int}$	$\{\text{int}\}$
S	$\longrightarrow \text{char}$	$\{\text{char}\}$
S	$\longrightarrow \text{num} \dots \text{num}$	$\{\text{num}\}$

Ejemplo: Parser descendente recursivo predictivo

```
Proc T()
```

```
  if tc in { int, char, num }
    S();
  else if tc in { ^ }
    match('^');
    S();
  else if tc in { array }
    match('array');
    T_p();
  else
    error();
```

```
End
```

```
Proc T_p()
```

```
  if tc in { of }
    match('of');
    T();
  else if tc in { [ ] }
    match('['); S(); match(']');
    match('of');
    T();
  else
    error();
```

```
End
```

```
Proc S()
```

```
  if tc in { int }
    match('int');
  else if tc in { char }
    match('char');
  else if tc in { num }
    match('num');
    match('..');
    match('num');
  else
    error();
```

```
End
```

```
Proc Main();
```

```
  T();
  match('$');
  accept();
```

```
End
```

Seguimiento: ok

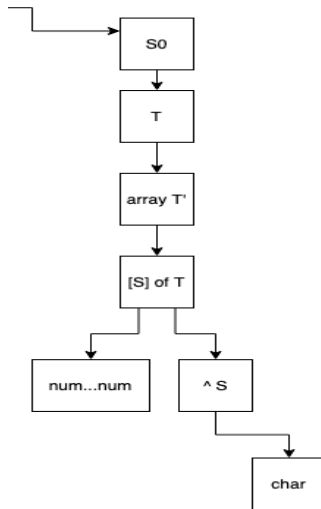
Cadena $array[num..num]of \uparrow char\$$:

Cadena	Pila	Producción / Acción
$array[num..num]of \uparrow char\$$	$T\$$	$T \rightarrow array \quad T'$
$array[num..num]of \uparrow char\$$	$array \quad T'\$$	$match\{array\}$
$[num..num]of \uparrow char\$$	$T'\$$	$T' \rightarrow [S] \text{ of } T$
$[num..num]of \uparrow char\$$	$[S] \text{ of } T\$$	$match\{\{\}$
$num..num]of \uparrow char\$$	$S] \text{ of } T\$$	$S \rightarrow num..num$
$num..num]of \uparrow char\$$	$num..num] \text{ of } T\$$	$match\{num\}$
$..num]of \uparrow char\$$	$..num] \text{ of } T\$$	$match\{..\}$
$num]of \uparrow char\$$	$num] \text{ of } T\$$	$match\{num\}$
$]of \uparrow char\$$	$] \text{ of } T\$$	$match\{\}\}$
$of \uparrow char\$$	$of \quad T\$$	$match\{of\}$
$\uparrow char\$$	$T\$$	

Seguimiento: ok

Cadena	Pila	Producción / Acción
$\uparrow \text{char}\$$	$T\$$	$T \rightarrow \uparrow S$
$\uparrow \text{char}\$$	$\uparrow S\$$	$\text{match}\{\uparrow\}$
$\text{char}\$$	$S\$$	$S \rightarrow \text{char}$
$\text{char}\$$	$\text{char}\$$	$\text{match}\{\text{char}\}$
$\$$	$\$$	$\text{match}\{\$\}$

AST:



Seguimiento: error

Cadena $array[num..num]of$:

Cadena	Pila	Producción / Acción
$array[num..num]of\$$	$T\$$	$T \rightarrow array \ T'$
$array[num..num]of\$$	$array \ T'\$$	$match\{array\}$
$[num..num]of\$$	$T'\$$	$T' \rightarrow [S] \ of \ T$
$[num..num]of\$$	$[S] \ of \ T\$$	$match\{[\]$
$num..num]of\$$	$S] \ of \ T\$$	$S \rightarrow num..num$
$num..num]of\$$	$num..num] \ of \ T\$$	$match\{num\}$
$..num]of\$$	$..num] \ of \ T\$$	$match\{..\}$
$num]of\$$	$num] \ of \ T\$$	$match\{num\}$
$]of\$$	$] \ of \ T\$$	$match\{]\}$
$of\$$	$of \ T\$$	$match\{of\}$
$\$$	$T\$$	ERROR()

Seguimiento: ustedes!

Realizar el seguimiento de las siguientes cadenas:

- *of num...num*
- *array [char] of array of int*

Teoría de Lenguajes

Analizadores Sintácticos Descendentes

Natalia Pesaresi

DC-UBA

1er. Cuatrimestre 2020