

# Teoría de Lenguajes

Intro compiladores + Analizadores Sintácticos Descendentes

Christian G. Cossio Mercado

DC-UBA

1er Cuatrimestre 2017

# Compiladores

- Tienen como objetivo convertir un código fuente de programa a otro lenguaje
- En particular importarán aquellos que realizan la conversión a código máquina ejecutable o en código intermedio interpretable (e.g., Java o .NET)
- Adicionalmente, 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)
- ➋ **Análisis Sintáctico:** A partir de la salida del analizador léxico se generará el árbol sintáctico correspondiente al código recibido
- 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**

## 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

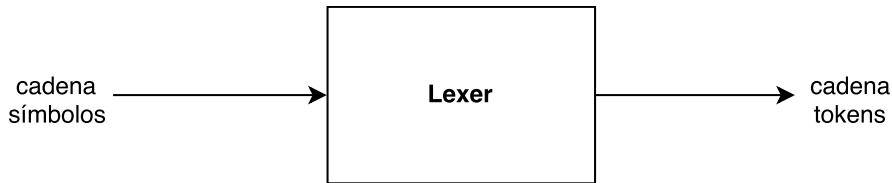
$$\alpha : (14 + 6) * 2$$

- ¿ $\alpha \in L$ ?

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

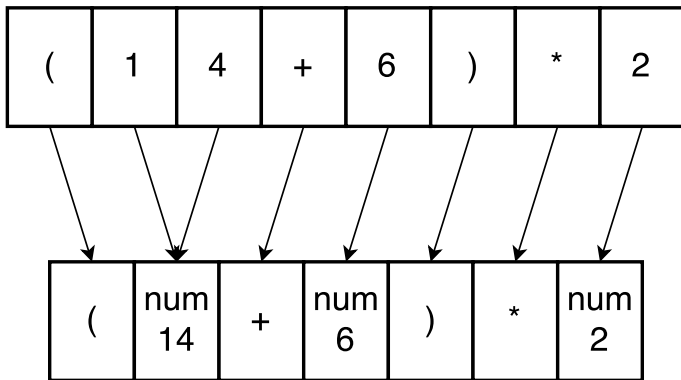
- ¿Reemplazamos al terminal *num* por los terminales 14, 6 y 2?
- ¿Reemplazamos la producción  $F \rightarrow num$  por las producciones  $F \rightarrow 14$ ,  $F \rightarrow 6$  y  $F \rightarrow 2$ ?
- ¿Podemos hacer esto para todos los números?
- ¡Son infinitos!
- Ahí es donde entran en juego los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

# Analizador Léxico (Lexer)

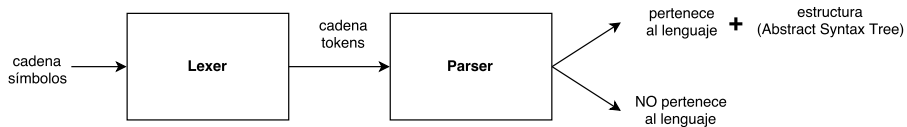


- Usa expresiones regulares
- Agrupa símbolos de la entrada en tokens
- Un token tiene un tipo y un valor

# Ejemplo



# Lexer + Parser





# Analizadores Sintácticos (parsers)

- ¿Para qué sirve un analizador Sintáctico?

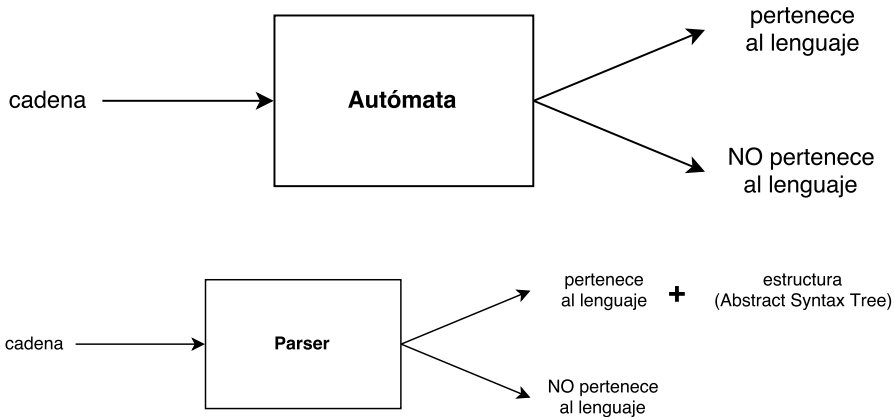
# Analizadores Sintácticos (parsers)

- ¿Para qué sirve un analizador Sintáctico?
- En pocas palabras: para verificar si el código sigue la gramática del lenguaje y obtener el árbol sintáctico correspondiente

# Analizadores Sintácticos (parsers)

- ¿Para qué sirve un analizador Sintáctico?
- En pocas palabras: para verificar si el código sigue la gramática del lenguaje y obtener el árbol sintáctico correspondiente
- Entonces, ¿para qué querría un analizador sintáctico, si ya tengo a los autómatas?

# Autómata vs. Parser



# Tipos de Analizadores Sintácticos

- Descendentes (top-down):

# Tipos de Analizadores Sintácticos

- Descendentes (top-down):
  - ▶ Construyen el árbol sintáctico desde la raíz hasta las hojas (la secuencia de tokens).

# Tipos de 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).

# Tipos de 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)



# Tipos de 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).

# Tipos de 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)

# Tipos de 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

# Tipos de 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

# Analizadores top-down

- Dada la siguiente gramática:

$$\begin{aligned}S &\rightarrow A \mid B \\A &\rightarrow a \\B &\rightarrow b\end{aligned}$$

- ¿Cómo genero la cadena a?

# Analizadores Recursivos descendentes

- Opción A: hacer DFS o BFS

# Analizadores Recursivos descendentes

- Opción A: hacer DFS o BFS
  - ▶ ¿Qué problemas tiene?

# Analizadores Recursivos descendentes

- Opción A: hacer DFS o BFS
  - ▶ ¿Qué problemas tiene?
- Opción B: Usar un analizador predictivo



# Analizadores Recursivos descendentes

- Opción A: hacer DFS o BFS
  - ▶ ¿Qué problemas tiene?
- Opción B: Usar un analizador predictivo
  - ▶ Ya tengo un criterio de selección de un regla

# Analizadores Recursivos descendentes

- Opción A: hacer DFS o BFS
  - ▶ ¿Qué problemas tiene?
- Opción B: Usar un analizador predictivo
  - ▶ Ya tengo un criterio de selección de un regla
  - ▶ Sigue teniendo algunas debilidades/restricciones

# Parsers predictivos descendentes

- Parsers descendentes (parten del símbolo distinguido)
- Utilizan la información del próximo token (i.e., token corriente, o `tc`) y los SD's de cada regla
- Con eso,
  - ▶ se crea un procedimiento por cada no terminal, siendo el procedimiento principal (`Main`) el correspondiente al símbolo distinguido
  - ▶ dentro de cada procedimiento, se invocan los procedimientos de los no terminales incluidos, se busca los terminales esperados (`match`), o no se hace nada (si es un  $\lambda$ ),
  - ▶ si hay más de regla para un mismo no terminal, se agregan condicionales para ejecutar sólo la regla 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

# ¿Cómo elegir una regla?

## Primeros, Siguiente, Símbolos Directrices

- Vamos a definir dos funciones: Primeros y Siguientes

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

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

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

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

- Con eso creamos la función de Símbolos Directrices (SD):

$$SD(A \rightarrow \beta) = \begin{cases} \text{Primeros}(\beta) & \text{si } \beta \text{ no anulable } (\beta \not\xRightarrow{*} \lambda) \\ \text{Primeros}(\beta) \cup \text{Siguientes}(A) & \text{si } \beta \text{ anulable } (\beta \Rightarrow^{*} \lambda) \end{cases}$$

# Algunos problemas y cómo resolverlos

- Ambigüedad por SDs no disjuntos

# Algunos problemas y cómo resolverlos

- Ambigüedad por SDs no disjuntos
  - ▶ factorización a la izquierda (*left-factorization*)

# Algunos problemas y cómo resolverlos

- Ambigüedad por SDs no disjuntos
  - ▶ factorización a la izquierda (*left-factorization*)
- Recursión a izquierda

# Algunos problemas y cómo resolverlos

- Ambigüedad por SDs no disjuntos
  - ▶ factorización a la izquierda (*left-factorization*)
- Recursión a izquierda
  - ▶ eliminación de la recursión inmediata



# Algunos problemas y cómo resolverlos

- Ambigüedad por SDs no disjuntos
  - ▶ factorización a la izquierda (*left-factorization*)
- Recursión a izquierda
  - ▶ eliminación de la recursión inmediata
  - ▶ eliminación de la recursión no inmediata

# 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

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

# 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

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

# Eliminación de la recursión no inmediata

- ➊ Numerar los no terminales  $A_1, \dots, A_n$
- ➋ Para  $i \rightarrow 1$  a  $n$ 
  - ➊ Para  $j \rightarrow 1$  a  $i - 1$ 
    - ➊ 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  $A_j \rightarrow \delta_1 \mid \dots \mid \delta_k \in P$ )
    - ➋ Fin Para
    - ➌ Eliminar la recursión inmediata de  $A_i$
- ➌ Fin Para

## Ejemplo

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

$G_1 = \langle \{S_0, T, S\}, \{\text{array}, [, ], \text{of}, \dots, \text{int}, \text{char}, \text{num}, \uparrow\}, S_0, P \rangle$ , con  $P$ :

$$S_0 \longrightarrow T$$

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

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

- Hacer un parser descendente recursivo para esta gramática. Si no es posible, corregirla previamente
- La cadena `array [num..num] of ↑ char`  $\in L(G_1)$ ?
- ¿Cuál es la secuencia de derivaciones? ¿Cuál es el árbol sintáctico?

## Ejemplo: Gramática corregida

- No es posible hacer un parser descendente recursivo predictivo para la gramática, porque hay solapamiento entre los SDs de dos reglas del no terminal  $T$
- Se lo resuelve aplicando factorización a la izquierda, para crear  $G'_1$

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

$$S_0 \longrightarrow T$$

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

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

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

## Ejemplo: Gramática corregida

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

$$S_0 \longrightarrow T$$

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

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

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

# Ejemplo: Gramática corregida y SDs

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

Regla		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

```

End

```

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

```

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

```

End

```

Proc Main();
  T();
  match('$');
  accept();

```

End

# Teoría de Lenguajes

Intro compiladores + Analizadores Sintácticos Descendentes

Christian G. Cossio Mercado

DC-UBA

1er Cuatrimestre 2017