

Conceptos de Lenguajes de Programación

Alberto Pardo Marcos Viera

Instituto de Computación, Facultad de Ingeniería
Universidad de la República, Uruguay

- En lenguajes de computadores, tanto la **sintaxis** como la **semántica** se deben especificar **sin ambigüedades**
 - Notaciones formales
- Distinguir entre sintaxis y semántica es de utilidad
- Sintaxis
 - Cómo especificar reglas estructurales del lenguaje: expresiones regulares, gramáticas libres de contexto
 - Cómo el compilador identifica la estructura: scanners, parsers

Los **tokens** son los componentes básicos de un programa, los strings más cortos con significado individual.

- Ej. Palabras reservadas, identificadores, símbolos, constantes, etc.

Los tokens se especifican con **expresiones regulares**:

- Un carácter
- El string vacío (ϵ)
- Concatenación de dos expresiones regulares
- Alternativa entre dos expresiones regulares ($|$)
- Estrella de Kleene ($*$)

Ejemplo de Expresión Regular

$number \rightarrow integer \mid real$

$integer \rightarrow digit \, digit^*$

$real \rightarrow integer \, exponent \mid decimal \, (exponent \mid \epsilon)$

$decimal \rightarrow digit^* \, (\, . \, digit \mid digit \, .) \, digit^*$

$exponent \rightarrow (e \mid E) \, (+ \mid - \mid \epsilon) \, integer$

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

Gramáticas Libres de Contexto

La notación de las **gramáticas libres de contexto** (CFG) suele llamarse Backus-Naur Form (BNF)

- Un conjunto de terminales (T)
- Un conjunto de no-terminales (N)
- Un no-terminal de inicio (S)
- Un conjunto de producciones, que pueden ser recursivas

EBNF si agregamos estrella y suma de Kleene.

Ejemplo de CFG:

$$\text{expr} \longrightarrow \text{id} \mid \text{number} \mid - \text{expr} \mid (\text{expr}) \\ \mid \text{expr op expr}$$
$$\text{op} \longrightarrow + \mid - \mid * \mid /$$

1. $expr \longrightarrow term \mid expr \textit{ add_op } term$
2. $term \longrightarrow factor \mid term \textit{ mult_op } factor$
3. $factor \longrightarrow id \mid \textit{ number } \mid - factor \mid (expr)$
4. $add_op \longrightarrow + \mid -$
5. $mult_op \longrightarrow * \mid /$

Análisis Sintáctico: Scanner

El **scanner** es usualmente el encargado de:

- reconocer los tokens
- quitar comentarios
- tratar los pragmas
- guardar el texto de los identificadores, números, strings
- guardar localizaciones en el fuente (archivo, línea, columna), para mensajes de error

Tres formas de construirlos:

- ad-hoc
 - Generalmente más rápidos y compactos
- Construir un DFA de forma semi-automática
 - Usualmente como casos anidados
- Construir un DFA dirigido por una tabla
 - Generadores como lex (flex) y scangen

Análisis Sintáctico: Parser

El **parser** es usualmente el encargado de:

- llamar al scanner para obtener los tokens del programa
- ensamblar los tokens en un árbol de sintaxis
- pasar el árbol a las siguientes fases del compilador

Mientras una CFG es una **generadora** de un lenguaje CF, un parser es un **reconocedor** del lenguaje.

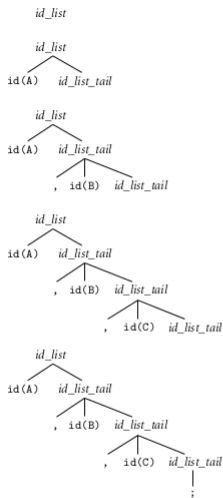
Para cualquier CFG se puede crear un parser $O(n^3)$

Existen **clases de gramáticas** para las que podemos construir parsers lineales. Las más importantes:

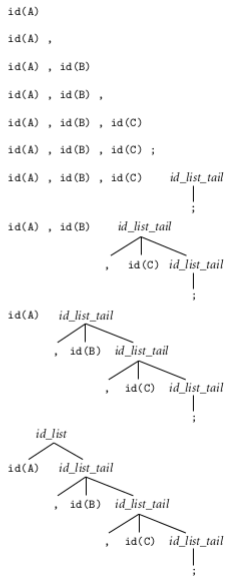
- **LL**: left-to-right scanning, left-most derivation
- **LR**: left-to-right scanning, right-most derivation

- Parsers LL se suelen llamar **top-down** o **predictivos**.
Construyen el árbol desde la raíz, prediciendo en cada paso la siguiente producción a usar.
- Parser LR se suelen llamar **bottom-up**. Construyen el árbol desde las hojas, reconociendo cuando una colección de sub-árboles hijos se pueden unir con un padre.

Parsers LL y LR



<code>id_list</code>	\rightarrow	<code>id id_list_tail</code>
<code>id_list_tail</code>	\rightarrow	<code>, id id_list_tail</code>
<code>id_list_tail</code>	\rightarrow	<code>;</code>



Gramática LL(1)

$$\textit{program} \rightarrow \textit{stmt_list} \ \ \$\$$$
$$\textit{stmt_list} \rightarrow \textit{stmt} \ \textit{stmt_list} \mid \epsilon$$
$$\textit{stmt} \rightarrow \textit{id} \ := \ \textit{expr} \mid \textit{read} \ \textit{id} \mid \textit{write} \ \textit{expr}$$
$$\textit{expr} \rightarrow \textit{term} \ \textit{term_tail}$$
$$\textit{term_tail} \rightarrow \textit{add_op} \ \textit{term} \ \textit{term_tail} \mid \epsilon$$
$$\textit{term} \rightarrow \textit{factor} \ \textit{factor_tail}$$
$$\textit{factor_tail} \rightarrow \textit{mult_op} \ \textit{factor} \ \textit{factor_tail} \mid \epsilon$$
$$\textit{factor} \rightarrow (\ \textit{expr} \) \mid \textit{id} \mid \textit{number}$$
$$\textit{add_op} \rightarrow + \mid -$$
$$\textit{mult_op} \rightarrow * \mid /$$

- Puedo construir un parser **recursivo descendente** (o generarlo con ANTLR)
- O generar un **LL parse table** con un driver program.

- Casi siempre se basan en tablas.
- Mantiene un stack con las raíces de los subárboles parcialmente completados.
- **shifts** tokens al stack
- **reduces** símbolos en tope del stack