



Tecnológico de Monterrey

Entrega 1

Curso: TC3002B.503

Estudiante:

Marco Ottavio Podesta Vezzali - A00833604

Fecha de entrega: 29/04/2025

ANTLR (ANother Tool for Language Recognition) (Target Python)

<https://www.antlr.org/>

ANTLR es una herramienta generadora de analizadores que permite definir gramáticas en un formato declarativo basado en EBNF. Aunque está desarrollada en Java, permite generar código en múltiples lenguajes, incluido Python, facilitando la integración en proyectos Python sin necesidad de reescribir la funcionalidad del analizador.

Formato de entrada

- Las gramáticas se definen en archivos con extensión `.g4`.
- La sintaxis de estas gramáticas es declarativa, basándose en EBNF para especificar reglas léxicas y sintácticas.
- Las reglas para el análisis léxico se definen mediante patrones que utilizan expresiones regulares, mientras que las reglas sintácticas describen la estructura del lenguaje a analizar.
- Las reglas sintácticas se escriben debajo, con nombres en minúsculas, definiendo cómo los tokens se combinan para formar estructuras más complejas.
- Se pueden incluir **acciones semánticas** directamente en la gramática mediante bloques de código en el lenguaje de destino (Python en este caso).
-

Tipo de ejecución

- Para la generación de código, ANTLR toma el archivo `.g4` y genera automáticamente el código fuente necesario (clases de lexer, parser, visitor, listener).
- Una vez generado, el analizador puede integrarse dentro de un proyecto Python para procesar entradas y generar árboles de sintaxis (parse trees) o recorrer el árbol utilizando visitors o listeners.

Técnicas de análisis

- Análisis léxico: Basado en expresiones regulares integradas en la gramática.
- Análisis sintáctico: Utiliza LL(*) ("LL-star"), una mejora de los parsers LL tradicionales, permitiendo predicción basada en la evaluación de múltiples tokens hacia adelante (lookahead) y backtracking opcional.

Tipo de licenciamiento

- ANTLR se distribuye bajo una licencia de código abierto similar a BSD o una licencia permisiva comparable, lo que permite su uso en proyectos tanto comerciales como de código abierto.

Herramientas teóricas empleadas

- Análisis léxico: Se apoya en expresiones regulares definidas directamente en la gramática para identificar tokens.
- Análisis sintáctico: Utiliza métodos de análisis LL(*) que permiten manejar un amplio rango de gramáticas y resolver ambigüedades mediante técnicas predictivas y, cuando es necesario, estrategias de backtracking.

Tipo de interfaz

- Las definiciones se realizan en archivos de texto (.g4) que especifican tanto los tokens (reglas léxicas) como la estructura (reglas sintácticas) de forma declarativa.
- No cuenta con una interfaz gráfica; la configuración se realiza mediante la edición del archivo de gramática, asistido por IDEs que suelen ofrecer realce sintáctico y validación.
- La salida generada es un conjunto de módulos Python que contienen clases y métodos para iniciar el análisis y manejar los árboles de sintaxis resultantes.

Facilidad para añadir código propio

- Dentro del archivo .g4 es posible incrustar bloques de código en Python en secciones designadas (por ejemplo, en las secciones de acción), los cuales se ejecutan al reconocer una regla determinada.
- Se puede separar la definición de la gramática de la lógica de procesamiento adicional, permitiendo incluir funciones o métodos personalizados en el código Python generado o en módulos auxiliares que interactúan con el analizador.
- La generación del código de ANTLR permite modificar o ampliar el comportamiento del analizador sin alterar el proceso de compilación, facilitando la incorporación de validaciones o transformaciones adicionales.

Lark <https://github.com/lark-parser/lark>

Lark es una librería para la generación de analizadores léxicos y sintácticos, completamente implementada en Python.

Formato de entrada

- Las gramáticas se definen utilizando una sintaxis basada en EBNF.
- Se pueden especificar las reglas tanto en cadenas de texto dentro del código como en archivos externos.
- Las reglas léxicas se establecen mediante expresiones regulares integradas en la definición de la gramática, y las reglas sintácticas describen la estructura jerárquica del lenguaje a analizar.

Tipo de ejecución

- LALR(1): Para gramáticas deterministas donde se busca mayor rendimiento.
- Earley: Para gramáticas ambiguas o recursivas, permitiendo manejar estructuras complejas sin la rigidez del análisis determinista.
- Al compilar la gramática, Lark genera de forma inmediata el árbol de derivación (parse tree) durante la ejecución en Python, sin etapas de compilación externa.

Tipo de licenciamiento

- Lark se distribuye bajo licencia MIT, lo que lo hace completamente gratuito para uso comercial y de código abierto, ofreciendo gran flexibilidad en la integración y distribución de proyectos.

Herramientas teóricas empleadas

- Análisis léxico: Se basa en expresiones regulares para reconocer tokens, integradas de forma nativa en la gramática.
- Análisis sintáctico: Implementa técnicas LALR(1) para gramáticas deterministas y el algoritmo Earley para aquellas ambiguas o con recursión profunda, proporcionando una solución adaptable al tipo de lenguaje y complejidad de la gramática.

Tipo de interfaz

- La gramática se define mediante cadenas de texto o archivos externos y se integra mediante una API Python sencilla.

- No posee una interfaz gráfica propia; la definición y configuración se hace directamente en el código, aprovechando la sencillez y flexibilidad de Python.
- La salida principal es la generación de árboles de análisis (parse trees), los cuales pueden ser recorridos, transformados y manipulados mediante callbacks o transformadores definidos por el usuario.

Facilidad para añadir código propio

Permite asociar funciones o transformadores a reglas específicas, mediante la integración de acciones semánticas.

Extensibilidad

- Es sencillo combinar la gramática con código Python adicional. Se pueden definir clases y funciones que procesen o modifiquen el árbol de análisis, integrándose de forma natural en el flujo del programa.
- La configuración y extensión del analizador no requieren pasos de compilación externos, permitiendo iteraciones rápidas durante el desarrollo.

ANTLR es una herramienta madura y poderosa para la generación de analizadores léxicos y sintácticos, con un enfoque declarativo mediante archivos .g4 basados en EBNF. Aunque su núcleo está en Java (por lo que se necesita del JDK). Soporta múltiples lenguajes de destino, incluyendo Python. Destaca por su capacidad de integrar acciones semánticas en la gramática, su análisis sintáctico avanzado basado en LL(*) que permite lookahead múltiple y backtracking. Su enfoque requiere un proceso de compilación previo para generar las clases de lexer, parser, y visitors o listeners en Python.

Lark es una librería nativa de Python, más ligera y directa para construir analizadores, sin necesidad de procesos de generación de código externos. Soporta tanto análisis determinista (LALR(1)) como análisis más flexible y robusto (Earley), adaptándose bien a diferentes tipos de gramáticas, incluidas las ambiguas.

Resumen comparativo:

- ANTLR ofrece más potencia y opciones de optimización para proyectos grandes y complejos, pero hay una curva de configuración más pronunciada.
- Lark es nativo en python, por lo que resulta más fácil de implementar. Además, siendo este un proyecto relativamente sencillo, se puede implementar con Lark.

I. Diseñar las Expresiones Regulares

Token	Descripción	Regex (estándar)
KW_PROGRAM	program	\bprogram\b
KW_VAR	var	\bvar\b
KW_MAIN	main	\bmain\b
KW_VOID	Void	\bvoid\b
KW_INT	Palabra reservada int	\bint\b
KW_FLOAT	Palabra reservada float	\bfloat\b
ID	Nombre (letra seguida de letras/dígitos/_)	[a-zA-Z][a-zA-Z0-9_]*
KW_WHILE	while	\bwhile\b
KW_DO	do	\bdo\b
KW_END	end	\bend\b
KW_IF	if	\bif\b
KW_ELSE	else	\belse\b
KW_PRINT	print	\bprint\b
CTE_INT	Constante entera	[0-9]+
CTE_FLOAT	Constante flotante	[0-9]+\.[0-9]+
CTE_STRING	Cadena entre comillas	"([^\"] \\.)*"
PLUS	Operador suma	\+
MINUS	Operador resta	-
MULT	Operador multiplicación	*
DIV	Operador división	/
LT	Menor que	<
GT	Mayor que	>

NEQ	Diferente (!=)	!=
ASSIGN	Asignación (=)	=
SEMICOLON	Punto y coma	;
COMMA	Coma	,
LPAREN	Paréntesis izquierdo	\(
RPAREN	Paréntesis derecho	\)
LBACE	Llave izquierda	{
RBRACE	Llave derecha	}
COLON	Dos puntos	:

II. Gramática libre de contexto

<PROGRAMA> ::= KW_PROGRAM ID SEMICOLON <VARS> <FUNCS> KW_MAIN <Body> KW_END
<VARS> ::= KW_VAR <ID_LIST> COLON <TYPE> SEMICOLON <VARS> ε
<ID_LIST> ::= ID <ID_LIST_TAIL>
<ID_LIST_TAIL> ::= COMMA ID <ID_LIST_TAIL> ε
<TYPE> ::= KW_INT KW_FLOAT
<CTE> ::= CTE_INT CTE_FLOAT
<EXPRESIÓN> ::= <EXP>

<EXP> GT <EXP> <EXP> LT <EXP> <EXP> NEQ <EXP>
<EXP> ::= <TÉRMINO> <EXP'>
<EXP'> ::= PLUS <TÉRMINO> <EXP'> MINUS <TÉRMINO> <EXP'> ε
<TÉRMINO> ::= <FACTOR> <TÉRMINO'>
<TÉRMINO'> ::= MULT <FACTOR> <TÉRMINO'> DIV <FACTOR> <TÉRMINO'> ε
<FACTOR> ::= LPAREN <EXPRESIÓN> RPAREN PLUS <FACTOR> MINUS <FACTOR> ID <CTE>
<ASSIGN> ::= ID ASSIGN <EXPRESIÓN> SEMICOLON
<Print> ::= KW_PRINT LPAREN <PrintList> RPAREN SEMICOLON
<PrintList> ::= <PrintItem> <PrintListTail>
<PrintListTail> ::= COMMA <PrintItem> <PrintListTail> ε
<PrintItem> ::= <EXPRESIÓN> CTE_STRING
<F_Call> ::= ID LPAREN <ArgList> RPAREN SEMICOLON
<ArgList> ::= <EXPRESIÓN> <ArgListTail>
<ArgListTail> ::= COMMA <EXPRESIÓN> <ArgListTail> ε

$\langle \text{STATEMENT} \rangle ::= \langle \text{ASSIGN} \rangle$ $\quad \langle \text{CONDITION} \rangle$ $\quad \langle \text{CYCLE} \rangle$ $\quad \langle \text{F_Call} \rangle$ $\quad \langle \text{Print} \rangle$
$\langle \text{BODY} \rangle ::= \text{LBRACE} \langle \text{STATEMENT_LIST} \rangle \text{RBRACE}$
$\langle \text{STATEMENT_LIST} \rangle ::= \langle \text{STATEMENT} \rangle \langle \text{STATEMENT_LIST} \rangle$ $\quad \epsilon$
$\langle \text{CONDITION} \rangle ::= \text{KW_IF} \text{LPAREN} \langle \text{EXPRESIÓN} \rangle \text{RPAREN} \langle \text{BODY} \rangle \langle \text{ELSE_OPT} \rangle$ SEMICOLON
$\langle \text{ELSE_OPT} \rangle ::= \text{KW_ELSE} \langle \text{BODY} \rangle$ $\quad \epsilon$
$\langle \text{CYCLE} \rangle ::= \text{KW_WHILE} \text{LPAREN} \langle \text{EXPRESIÓN} \rangle \text{RPAREN} \text{KW_DO} \langle \text{BODY} \rangle$ SEMICOLON
$\langle \text{FUNCS} \rangle ::= \langle \text{FUNC} \rangle \langle \text{FUNCS} \rangle$ $\quad \epsilon$
$\langle \text{FUNC} \rangle ::= \text{KW_VOID} \text{ID} \text{LPAREN} \langle \text{PARAM_LIST_OPT} \rangle \text{RPAREN} \langle \text{VARS} \rangle \langle \text{Body} \rangle$ SEMICOLON
$\langle \text{PARAM_LIST_OPT} \rangle ::= \langle \text{PARAM_LIST} \rangle$ $\quad \epsilon$
$\langle \text{PARAM_LIST} \rangle ::= \text{ID} \text{COLON} \langle \text{TYPE} \rangle \langle \text{PARAM_LIST_TAIL} \rangle$
$\langle \text{PARAM_LIST_TAIL} \rangle ::= \text{COMMA} \text{ID} \text{COLON} \langle \text{TYPE} \rangle \langle \text{PARAM_LIST_TAIL} \rangle$ $\quad \epsilon$

Test Plan para Compiler

1. Pruebas de Programas Completos

Caso	Descripción	Código de prueba	Resultado esperado
1.1	Programa mínimo válido	program a; main {} end	AST con programa vacío.
1.2	Programa con variables y cuerpo vacío	program b; var x: int; main {} end	AST donde vars no esté vacío.

2. Pruebas de Declaraciones de Variables (vars)

Caso	Descripción	Código de prueba	Resultado esperado
2.1	Variables tipo int	var x, y: int;	Lista de variables con tipo int.
2.2	Variables tipo float	var a: float;	Lista de variables con tipo float.
2.3	Sin variables	<i>(no se escribe var)</i>	vars vacío en AST.

3. Pruebas de Statements

Tipo de statement	Código ejemplo	Resultado esperado
Asignación	x = 5;	Nodo assign con ID y valor 5.

Llamada a función	<code>foo(1, 2);</code>	Nodo call con argumentos.
Condicional if-else	<code>if (x != y) {} else {};</code>	Nodo if con neq y rama else.
Ciclo while	<code>while (x > 0) do {};</code>	Nodo while con gt.

4. Pruebas de Expresiones (expr, term, factor)

Expresión	Código	Resultado esperado
Aritmética básica	<code>x + y;</code>	Nodo add.
Aritmética anidada	<code>(x - 2) * 3;</code>	Nodos sub, luego mult.