

# Informe: Construcción de un Generador de Lexer

Nombre del Autor

7 de julio de 2025

## Resumen

El proyecto consiste en el diseño e implementación de un compilador desarrollado en Rust, que utiliza LLVM para la generación de código intermedio y clang para la compilación final. El compilador abarca todas las fases fundamentales: análisis léxico (lexer), sintáctico (parser), semántico y generación de código. El flujo de trabajo parte de la lectura del código fuente, construcción del Árbol de Sintaxis Abstracta (AST), validación semántica y generación de código intermedio (LLVM-IR), hasta la obtención del ejecutable final. El proceso está automatizado mediante un Makefile, facilitando la compilación, ejecución y limpieza de artefactos generados.

**Nota:** Para una comprensión rigurosa, se recomienda precisar las definiciones de cada fase del compilador y los conceptos de autómatas y gramáticas, ya que son fundamentales para el diseño de analizadores léxicos y sintácticos.

## 1. Gramática Formal del Lenguaje

La gramática del lenguaje que se va a compilar se define formalmente mediante una gramática libre de contexto (GLC), especificando los símbolos terminales, no terminales, el símbolo inicial y las reglas de producción. Es importante precisar que una GLC se expresa como una cuádrupla  $(V, \Sigma, R, S)$ , donde  $V$  es el conjunto de símbolos no terminales,  $\Sigma$  el conjunto de símbolos terminales,  $R$  el conjunto de reglas de producción y  $S$  el símbolo inicial.

### Símbolos terminales

- IDENTIFIER, NUMBER, STRING
- FUNCTION, TYPE, LET, IN, IF, ELSE, ELIF, WHILE, FOR, NEW, PRINT, INHERITS, TRUE, FALSE
- ASSIGN\_OP, DESTRUCTIVE\_ASSIGN\_OP, ARROW\_OP, LOGICAL\_OR\_OP, LOGICAL\_AND\_OP
- EQUALITY\_OP, COMPARISON\_OP, TERM\_OP, FACTOR\_OP, POW\_OP, UNARY\_OP, DOT\_OP
- LPAREN, RPAREN, LBRACE, RBRACE, SEMICOLON, COMMA, COLON

*Precisión:* Los símbolos terminales corresponden a los tokens reconocidos por el lexer. Es recomendable definir formalmente las expresiones regulares asociadas a cada token para evitar ambigüedades.

## Símbolos no terminales

- Program, Statement, FunctionDef, FunctionFullDef, FunctionArrowDef
- TypeDef, Inheritance, TypeBody, MemberDef, Expr, DestructiveAssignExpr
- LogicalOrExpr, LogicalAndExpr, EqualityExpr, ComparisonExpr, TermExpr
- FactorExpr, ExponentExpr, UnaryExpr, CompositeExpr, PrimaryExpr
- Literal, FunctionCall, ArgList, LetIn, IfElse, WhileLoop, ForLoop
- CodeBlock, ExprList, PrintExpr, AssignmentList, Assignment, Params, Signature

*Precisión:* Los símbolos no terminales definen la estructura sintáctica del lenguaje. Es recomendable identificar el símbolo inicial explícitamente (en este caso, **Program**).

## Definiciones léxicas

IDENTIFIER [A-Za-z] [A-Za-z0-9\_]\*

NUMBER [0-9]+({0-9}+)?

STRING ‘‘([^\\"\\]\.)\*’’

BOOLEAN true|false

## Operadores

ASSIGN\_OP =

DESTRUCTIVE\_ASSIGN\_OP :=

ARROW\_OP =>

LOGICAL\_OR\_OP |

LOGICAL\_AND\_OP &

EQUALITY\_OP == !=

COMPARISON\_OP < <= > >=

TERM\_OP + - @

FACTOR\_OP \* / %

POW\_OP ^

UNARY\_OP ! -

DOT\_OP .

## **Delimitadores**

LPAREN (

RPAREN )

LBRACE {

RBRACE }

SEMICOLON ;

COMMA ,

COLON :

## **Palabras reservadas**

FUNCTION function

TYPE type

LET let

IN in

IF if

ELSE else

ELIF elif

WHILE while

FOR for

NEW new

PRINT print

INHERITS inherits

## Reglas de producción

Program  $\rightarrow$  (Statement SEMICOLON)\* Statement? SEMICOLON  
Statement  $\rightarrow$  FunctionDef | TypeDef | Expr  
FunctionDef  $\rightarrow$  FunctionFullDef | FunctionArrowDef  
FunctionFullDef  $\rightarrow$  FUNCTION IDENTIFIER Params COLON Signature CodeBlock  
FunctionArrowDef  $\rightarrow$  FUNCTION IDENTIFIER Params COLON Signature ARROW Expr  
Params  $\rightarrow$  LPAREN (IDENTIFIER COLON Signature (COMMA IDENTIFIER COLON Signature))\* RPAREN  
TypeDef  $\rightarrow$  TYPE IDENTIFIER Params? Inheritance? LBRACE TypeBody RBRACE  
Inheritance  $\rightarrow$  INHERITS IDENTIFIER ArgList?  
TypeBody  $\rightarrow$  (MemberDef SEMICOLON)\* MemberDef?  
MemberDef  $\rightarrow$  IDENTIFIER ASSIGN Expr | FunctionDef  
Expr  $\rightarrow$  DestructiveAssignExpr | LogicalOrExpr  
DestructiveAssignExpr  $\rightarrow$  PrimaryExpr DESTRUCTIVE\_ASSIGN\_OP Expr  
LogicalOrExpr  $\rightarrow$  LogicalOrExpr LOGICAL\_OR\_OP LogicalAndExpr | LogicalAndExpr  
LogicalAndExpr  $\rightarrow$  LogicalAndExpr LOGICAL\_AND\_OP EqualityExpr | EqualityExpr  
EqualityExpr  $\rightarrow$  EqualityExpr EQUALITY\_OP ComparisonExpr | ComparisonExpr  
ComparisonExpr  $\rightarrow$  ComparisonExpr COMPARISON\_OP TermExpr | TermExpr  
TermExpr  $\rightarrow$  TermExpr TERM\_OP FactorExpr | FactorExpr  
FactorExpr  $\rightarrow$  FactorExpr FACTOR\_OP ExponentExpr | ExponentExpr  
ExponentExpr  $\rightarrow$  UnaryExpr POW\_OP ExponentExpr | UnaryExpr  
UnaryExpr  $\rightarrow$  UNARY\_OP UnaryExpr | CompositeExpr  
CompositeExpr  $\rightarrow$  LetIn | IfElse | WhileLoop | ForLoop | PrimaryExpr  
PrimaryExpr  $\rightarrow$  FunctionCall | NEW IDENTIFIER ArgList | PrimaryExpr DOT\_OP IDENTIFIER  
PrimaryExpr DOT\_OP FunctionCall | Literal | IDENTIFIER | LPAREN Expr RPAREN  
Literal  $\rightarrow$  NUMBER | STRING | BOOLEAN  
FunctionCall  $\rightarrow$  IDENTIFIER ArgList  
ArgList  $\rightarrow$  LPAREN (Expr (COMMA Expr))\* RPAREN  
LetIn  $\rightarrow$  LET AssignmentList IN CompositeExpr  
IfElse  $\rightarrow$  IF LPAREN Expr RPAREN CodeBlock (ELIF LPAREN Expr RPAREN CodeBlock)\* ELSE LPAREN Expr RPAREN CodeBlock  
WhileLoop  $\rightarrow$  WHILE LPAREN Expr RPAREN CompositeExpr  
ForLoop  $\rightarrow$  FOR LPAREN IDENTIFIER IN RANGE LPAREN Expr COMMA Expr RPAREN CodeBlock  
CodeBlock  $\rightarrow$  LBRACE ExprList RBRACE  
ExprList  $\rightarrow$  (Expr SEMICOLON)\* Expr?  
PrintExpr  $\rightarrow$  PRINT LPAREN Expr RPAREN  
AssignmentList  $\rightarrow$  Assignment (COMMA Assignment)\*  
Assignment  $\rightarrow$  IDENTIFIER ASSIGN Expr  
Signature  $\rightarrow$  IDENTIFIER

## 2. Manejo de Posiciones y Construcción del AST con LALRPOP

En el proceso de análisis sintáctico, cada token reconocido es acompañado por una estructura denominada `Span`, que almacena las posiciones de inicio y fin del token dentro del código fuente (medidas en offsets de byte). Esta información es fundamental para:

- **Rastreo preciso de ubicaciones:** Permite identificar la posición exacta de cada elemento del programa, facilitando la depuración y el análisis posterior.
- **Generación de mensajes de error detallados:** El `Span` posibilita que los errores sintácticos reporten la línea, columna y contexto específico donde ocurre el problema, mejorando la experiencia del usuario.
- **Conservación de metadatos en el AST:** El Árbol de Sintaxis Abstracta (AST) generado por el compilador incluye estos metadatos de ubicación en todos sus nodos, lo que resulta útil para etapas posteriores como el análisis semántico o la generación de código.

La generación de `Span` es gestionada automáticamente por LALRPOP en el archivo `parser.lalrpop`, utilizando las anotaciones `@L` y `@R`, que corresponden a los offsets de los tokens en el texto fuente.

## Teoría del Parsing LR(1) y LALRPOP

LALRPOP genera por defecto un parser de tipo LR(1), empleando una variante llamada *lane table*, que es más compacta en memoria. Los analizadores LR(1) realizan un análisis ascendente (bottom-up) con un símbolo de *lookahead*, lo que garantiza:

- **Reconocimiento eficiente:** El análisis es de complejidad  $O(n)$ , sin necesidad de backtracking.
- **Cobertura completa:** Permite parsear cualquier gramática LR(1), que es una clase amplia y expresiva de gramáticas libres de contexto.

*Precisión:* Es recomendable explicar la diferencia entre los analizadores LR(0), SLR(1), LALR(1) y LR(1), y por qué LALRPOP utiliza LR(1) con optimizaciones de memoria.

Internamente, LALRPOP construye un autómata LR(0), calcula los conjuntos LR(1) y genera las tablas de parsing (shift/reduce). El código Rust resultante implementa el parser como una máquina de estados con un stack, siguiendo el paradigma clásico de los analizadores LR.

## 3. Recomendaciones para la Precisión Conceptual