



Primer proyecto – Compiladores e intérpretes  
Analizador léxico

Profesor:

Allan Rodríguez Dávila

Estudiantes:

Jairo González Hidalgo - 2024178535

Rafael Odio Mendoza - 2024083499

26 de diciembre – Verano 2025

## Contenido

|  |    |
|--|----|
| Repositorio de github.....   | 3  |
| Manual de usuario .....  | 3  |
| Compilación.....   | 3  |
| Ejecución .....  | 3  |
| Uso .....  | 3  |
| Pruebas de funcionalidad .....   | 4  |
| Descripción del problema .....   | 7  |
| Diseño del programa .....  | 7  |
| Decisiones de diseño .....   | 7  |
| Algoritmos usados.....   | 8  |
| Resumen y análisis del proceso que hace jflex para generar el analizador<br>léxico ..... | 9  |
| Librerías usadas .....   | 9  |
| Análisis de resultados.....  | 10 |
| Objetivos alcanzados .....   | 10 |
| Bitácora.....  | 11 |
| Bibliografía.....  | 13 |

# Repositorio de github

Enlace: <https://github.com/JairoGH16/proyecto-compiladores-e-interpretes>

## Manual de usuario

### Compilación

Para que el programa funcione, es necesario tener instalado el JDK de Java (se recomienda la versión 11 o superior). Si se usa VS Code, hay que instalar el entorno de Java, se recomienda el "Extension Pack for Java" de Microsoft. También pueden usar otros IDEs como IntelliJ IDEA o Eclipse sin problema.

El proyecto usa JFlex y CUP para generar el analizador léxico. Estos se encuentran como archivos .jar en la carpeta lib/ y no necesitan instalación adicional.

El proyecto fue diseñado para que no hubiera que estar usando comandos para compilar el Lexer y el Parser cada vez que se cambiara algo. En el menú principal, la Opción 1 se encarga de borrar los archivos .java viejos en la carpeta parserlexer, ejecutar el jar de cup para generar la clase de símbolos (sym.java) y ejecutar el jar de JFlex para generar el motor del scanner (Lexer.java).

Para crearlos manualmente sin usar main, primero hay que estar dentro de la carpeta programa: `\proyecto-compiladores-e-interpretes> cd programa`

Posteriormente se ejecutan

Lexer.java: `java -jar lib/jflex-full-1.9.1.jar parserlexer/Lexer.jflex`

sym.java (que también compila Parser.java): `java -jar lib/java-cup-11b.jar -destdir parserlexer/ -parser Parser parserlexer/Parser.cup`

Los archivos generados quedan en la carpeta parserlexer/

### Ejecución

Ejecutar el archivo Main.java, desplegará un menú.

### Uso

La opción 0 cierra el programa.

Para crear los archivos Lexer.java y sym.java se usa la opción 1.

Para ejecutar el análisis se usa la opción 2, se escogen los archivos, esta muestra una lista de todos los archivos .txt en la carpeta archivos\_prueba/, puede seleccionar un archivo específico o analizar todos mediante la opción 0,

Se pueden ver más pruebas del uso en el siguiente apartado.

## Pruebas de funcionalidad

Se pueden apreciar completamente las pruebas en el video:

<https://www.youtube.com/watch?v=kATyjgQVy2M>

## 1. Generación de archivos .java con la opción 1

```
> lib
  ↳ parserlexer
    ↳ Lexer.jflex
    ↳ Parser.cup
  J Main.java
  info.txt M

99         } else {
100             System.out.println("Opción fuera de rango");
101         }
102     }
103
104     // Analizar el archivo txt y crear el reporte
Windsurf Refactor | Explain | X
105     private static void ejecutarLexer(File archivoFuente) {
106         String nombreSinExt = archivoFuente.getName().replace(".txt", "");
107         String rutaReporte = RUTA_SALIDA + "reporte_" + nombreSinExt + ".txt";
108
109         // Lista para ir guardando los errores para ponerlos al final del txt
110         List<String> listaErrores = new ArrayList<>();
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS GITLENS QUERY RESULTS SQL HISTORY

- PS C:\Users\Jairo\Desktop\Compiladores e interpretes\proyecto\proyecto-compilador...
  - PS C:\Users\Jairo\Desktop\Compiladores e interpretes\proyecto\proyecto-compilador...
    - Error: no se ha encontrado o cargado la clase principal Main
  - PS C:\Users\Jairo\Desktop\Compiladores e interpretes\proyecto\proyecto-compilador...
    - Error: no se ha encontrado o cargado la clase principal Main.java
  - PS C:\Users\Jairo\Desktop\Compiladores e interpretes\proyecto\proyecto-compilador...
    - Error: no se ha encontrado o cargado la clase principal Main
  - PS C:\Users\Jairo\Desktop\Compiladores e interpretes\proyecto\proyecto-compilador...
    - Error: no se ha encontrado o cargado la clase principal Main

----- Menu de análisis léxico -----

1. Generar Lexer.java, Parser.java y sym.java
2. Realizar Análisis Léxico
0. Salir

Seleccione una opción: 1

The screenshot shows an IDE with a project named 'parserlexer'. The file explorer on the left lists 'Lexer.java', 'Lexer.jflex', 'Parser.cup', 'Parser.java', 'sym.java', and 'Main.java'. The main editor displays Java code for 'ejecutarLexer' and 'analizar' methods. The terminal window at the bottom shows the following output:

```

PS C:\Users\Jairo\Desktop\Compiladores e interpretes\proyecto\proyecto-compiladores-e-interpretes; & 'C:\Program Files\Java\jre1.8.0_471\bin\java'
2. Realizar Análisis Léxico
0. Salir
Seleccione una opción: 1
Borrando archivos antiguos...
Generando sym y Parser con CUP...
Generando Lexer con JFlex...
Archivos generados

----- Menu de análisis léxico -----
1. Generar Lexer.java, Parser.java y sym.java
2. Realizar Análisis Léxico
0. Salir
Seleccione una opción:

```

Los errores que marca Parser.java es normal, aún no se desarrolla esa parte porque no es del análisis léxico, no afecta.

## 2. Análisis de un archivo sin errores

The screenshot shows a terminal window with the output of a lexical analysis report for 'prueba1.txt'. The report is titled 'REPORTE DE TOKENS: prueba1.txt' and contains a table of tokens. The table has four columns: IDENTIFICADOR (TIPO), LEXEMA, LINEA, and COLUMNA. The tokens are listed as follows:

| IDENTIFICADOR (TIPO) | LEXEMA           | LINEA | COLUMNA |
|----------------------|------------------|-------|---------|
| 2 (WORLD)            | world            | 1     | 1       |
| 4 (INT)              | int              | 1     | 7       |
| 51 (ASSIGN)          | =                | 1     | 11      |
| 3 (LOCAL)            | local            | 1     | 13      |
| 5 (FLOAT)            | float            | 1     | 19      |
| 59 (ID)              | x                | 1     | 25      |
| 51 (ASSIGN)          | =                | 1     | 27      |
| 56 (FLOAT_LITERAL)   | 3.14             | 1     | 29      |
| 24 (ENDL)            | endl             | 1     | 34      |
| 11 (GIFT)            | gift             | 3     | 1       |
| 10 (COAL)            | coal             | 3     | 6       |
| 9 (NAVIDAD)          | navidad          | 3     | 11      |
| 29 (OPEN_PAREN)      | {                | 3     | 19      |
| 30 (CLOSE_PAREN)     | }                | 3     | 21      |
| 27 (OPEN_BLOCK)      | {                | 3     | 23      |
| 22 (SHOW)            | show             | 4     | 5       |
| 57 (STRING_LITERAL)  | "Hola mundo 123" | 4     | 25      |
| 50 (NOT)             | Σ                | 4     | 30      |
| 24 (ENDL)            | endl             | 4     | 32      |
| 28 (CLOSE_BLOCK)     | !                | 5     | 1       |
| 2 (WORLD)            | world            | 5     | 3       |

The report concludes with the message: 'No se encontraron errores léxicos.' and 'FIN DEL ANALISIS'.

### 3. Análisis de un archivo con errores léxicos

```
reporte_prueba2.txt
programa > archivos_salida > reporte_prueba2.txt
You, 4 hours ago | 2 authors (You and one other)
You, 8 hours ago * feat: manejo de errores léxicos

1  REPORTE DE TOKENS: prueba2.txt
2  =====
3  IDENTIFICADOR (TIPO) | LEXEMA | LINEA | COLUMNA
4  -----
5  2 (WORLD) | world | 1 | 1
6  4 (INT) | int | 1 | 7
7  51 (ASSIGN) | = | 1 | 11
8  3 (LOCAL) | local | 1 | 13
9  5 (FLOAT) | float | 1 | 19
10 59 (ID) | x | 1 | 25
11 51 (ASSIGN) | = | 1 | 27
12 56 (FLOAT_LITERAL) | 3.14 | 1 | 29
13 24 (ENDL) | endl | 1 | 34
14 11 (GIFT) | gift | 3 | 1
15 10 (COAL) | coal | 3 | 6
16 60 (ERROR) | # | 3 | 11
17 22 (SHOW) | show | 4 | 5
18 57 (STRING_LITERAL) | "Hola mundo 123" | 4 | 25
19 50 (NOT) | ! | 4 | 30
20 24 (ENDL) | endl | 4 | 32
21 60 (ERROR) | # | 5 | 1
22 28 (CLOSE_BLOCK) | ! | 6 | 1
23 2 (WORLD) | world | 6 | 3
24 -----
25 ERRORES ENCONTRADOS:
26 - Caracter ilegal <#> en línea 3
27 - Caracter ilegal <#> en línea 5
28 -----
29 FIN DEL ANALISIS
30
```

### 4. Análisis usando todos los tokens del lenguaje

```
reporte_prueba_completa.txt
programa > archivos_salida > reporte_prueba_completa.txt
You, 4 hours ago | 2 authors (You and one other)
You, 8 hours ago * feat: manejo de errores léxicos

1  REPORTE DE TOKENS: prueba_completa.txt
2  =====
3  IDENTIFICADOR (TIPO) | LEXEMA | LINEA | COLUMNA
4  -----
5  2 (WORLD) | world | 1 | 1
6  4 (INT) | int | 1 | 7
7  51 (ASSIGN) | = | 1 | 11
8  3 (LOCAL) | local | 1 | 13
9  5 (FLOAT) | float | 1 | 19
10 59 (ID) | x | 1 | 25
11 51 (ASSIGN) | = | 1 | 27
12 56 (FLOAT_LITERAL) | 3.14 | 1 | 29
13 24 (ENDL) | endl | 1 | 34
14 11 (GIFT) | gift | 3 | 1
15 10 (COAL) | coal | 3 | 6
16 60 (ERROR) | # | 3 | 11
17 22 (SHOW) | show | 4 | 5
18 57 (STRING_LITERAL) | "Hola mundo 123" | 4 | 25
19 50 (NOT) | ! | 4 | 30
20 24 (ENDL) | endl | 4 | 32
21 60 (ERROR) | # | 5 | 1
22 28 (CLOSE_BLOCK) | ! | 6 | 1
23 2 (WORLD) | world | 6 | 3
24 -----
25 ERRORES ENCONTRADOS:
26 - Caracter ilegal <#> en línea 3
27 - Caracter ilegal <#> en línea 5
28 -----
29 FIN DEL ANALISIS
30
```

```
reporte_prueba_completa.txt
programa > archivos_salida > reporte_prueba_completa.txt
You, 4 hours ago | 2 authors (You and one other)
You, 8 hours ago * feat: manejo de errores léxicos

1  REPORTE DE TOKENS: prueba_completa.txt
2  =====
3  IDENTIFICADOR (TIPO) | LEXEMA | LINEA | COLUMNA
4  -----
5  2 (WORLD) | world | 1 | 1
6  4 (INT) | int | 1 | 7
7  51 (ASSIGN) | = | 1 | 11
8  3 (LOCAL) | local | 1 | 13
9  5 (FLOAT) | float | 1 | 19
10 59 (ID) | x | 1 | 25
11 51 (ASSIGN) | = | 1 | 27
12 56 (FLOAT_LITERAL) | 3.14 | 1 | 29
13 24 (ENDL) | endl | 1 | 34
14 11 (GIFT) | gift | 3 | 1
15 10 (COAL) | coal | 3 | 6
16 60 (ERROR) | # | 3 | 11
17 22 (SHOW) | show | 4 | 5
18 57 (STRING_LITERAL) | "Hola mundo 123" | 4 | 25
19 50 (NOT) | ! | 4 | 30
20 24 (ENDL) | endl | 4 | 32
21 60 (ERROR) | # | 5 | 1
22 28 (CLOSE_BLOCK) | ! | 6 | 1
23 2 (WORLD) | world | 6 | 3
24 -----
25 ERRORES ENCONTRADOS:
26 - Caracter ilegal <#> en línea 3
27 - Caracter ilegal <#> en línea 5
28 -----
29 FIN DEL ANALISIS
30
```

## 5. Análisis con todos los tokens pero incluyendo errores

|   |           |       |         |
|---|-----------|-------|---------|
| 1 REPORTE DE TOKENS: prueba_completa_errores.txt                    |           |       |         |
| 2 =====   |           |       |         |
| 3 IDENTIFICADOR (TIPO)  | LEXEMA    | LINEA | COLUMNA |
| 4 =====   |           |       |         |
| 5 2 (WORLD)   | world     | 2     | 1       |
| 6 3 (LOCAL)   | local     | 2     | 7       |
| 7 4 (INT)   | int       | 2     | 13      |
| 8 5 (FLOAT)   | float     | 2     | 17      |
| 9 6 (BOOL)  | bool      | 2     | 23      |
| 10 60 (ERROR)   | #         | 2     | 28      |
| 11 9 (NAVIDAD)  | navidad   | 3     | 1       |
| 12 10 (COAL)  | coal      | 3     | 9       |
| 13 11 (GIFT)  | gift      | 3     | 14      |
| 14 12 (DECIDE)  | decide    | 3     | 19      |
| 15 13 (OF)  | of        | 3     | 26      |
| 16 14 (ELSE)  | else      | 3     | 29      |
| 17 15 (END)   | end       | 3     | 34      |
| 18 16 (LOOP)  | loop      | 3     | 38      |
| 19 17 (EXIT)  | exit      | 3     | 43      |
| 20 18 (WHEN)  | when      | 3     | 48      |
| 21 19 (FOR)   | for       | 4     | 1       |
| 22 20 (RETURN)  | return    | 5     | 1       |
| 23 21 (BREAK)   | break     | 5     | 8       |
| 24 22 (SHOW)  | show      | 5     | 14      |
| 25 23 (GET)   | get       | 5     | 19      |
| 26 24 (ENDL)  | endl      | 5     | 23      |
| 27 60 (ERROR)   | o         | 5     | 28      |
| 28 26 (FALSE)   | false     | 6     | 1       |
| 29 60 (ERROR)   | (         | 7     | 1       |
| 30 42 (LTEQ)  | <=        | 8     | 1       |
| 31 43 (GTEQ)  | >=        | 8     | 4       |
| 32 44 (EQEQ)  | ==        | 8     | 7       |
| 33 45 (NEQ)   | !=        | 8     | 10      |
| 34 46 (LT)  | <         | 8     | 13      |
| 35 47 (GT)  | >         | 8     | 15      |
| 36 48 (AND)   | @         | 8     | 17      |
| 37 49 (OR)  | ~         | 8     | 19      |
| 38 50 (NOT)   | Σ         | 8     | 21      |
| 39 51 (ASSIGN)  | =         | 8     | 23      |
| 40 52 (ARROW)   | ->        | 8     | 25      |
| 41 53 (COMMA)   | ,         | 8     | 28      |
| 42 54 (SEMICOLON)   | ;         | 8     | 30      |
| 43 55 (INT_LITERAL)   | 12345     | 9     | 1       |
| 44 56 (FLOAT_LITERAL)   | 3.14      | 9     | 7       |
| 45 57 (STRING_LITERAL)  | "Hola"    | 9     | 17      |
| 46 58 (CHAR_LITERAL)  | A         | 9     | 19      |
| 47 59 (ID)  | perro_123 | 9     | 23      |
| 48 29 (OPEN_PAREN)  | {         | 9     | 33      |
| 49 29 (OPEN_PAREN)  | {         | 9     | 34      |
| 50 60 (ERROR)   |           |       |         |
| 51 -----  |           |       |         |
| 52 ERRORES ENCONTRADOS:   |           |       |         |
| 53 - Caracter ilegal <#> en línea 2                                 |           |       |         |
| 54 - Caracter ilegal <@> en línea 5                                 |           |       |         |
| 55 - Caracter ilegal <{> en línea 7                                 |           |       |         |
| 56 - Caracter ilegal <Comentario multilinea no cerrado> en línea 11 |           |       |         |
| 57 -----  |           |       |         |
| 58 FIN DEL ANALISIS   |           |       |         |
| 59 -----  |           |       |         |

## Descripción del problema

El objetivo de este primer proyecto es desarrollar la fase de análisis léxico para un lenguaje de programación imperativo diseñado para configuración de chips en sistemas empotrados.

El lenguaje funcionará con ciertos símbolos no convencionales, como (¡ y !) para delimitar bloques, (¿ y ?) como paréntesis, también Σ, @ y ~ como operadores NOT, AND y OR, (ε y ϑ) como comentarios multilinea, también deben haber operadores de división entera, potencia y diferente.

El analizador léxico debe ser capaz de leer archivos fuente escritos en este lenguaje nuevo que se está desarrollando, debe identificar y clasificar cada elemento del código en tokens, reportar errores léxicos sin que el análisis se detenga y generar reportes detallados con información de los tokens y errores.

## Diseño del programa

### Decisiones de diseño

El programa sigue una arquitectura modular separada en tres componentes principales, los cuales son Main.java, Lexer.jflex y Parser.cup.

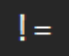
Se utiliza UTF-8 para aceptar los símbolos especiales “ε, ϑ, Σ, ¡, ¿”, al inicio se presentaban problemas porque no detectaba esos caracteres y los tomaba como

errores por ser elementos desconocidos para el lenguaje, por lo que la solución fue forzar UTF-8 al ejecutar JFlex.

Para manejar errores se implementó un estado llamado "PANIC\_MODE", al encontrar un carácter ilegal, se entra en este modo, donde ignora el resto la línea actual y pasa a la siguiente, en la cual vuelve al estado normal, esto evita que el análisis se caiga por errores.

En lugar de cargar Lexer.java directamente, se usa reflexión, que permite cargarlo en tiempo de ejecución, invocar métodos que no tengan referencias directas en tiempo de compilación, extraer los nombres de los tokens de sym.java y regenerar el Lexer.java sin recompilar todo el proyecto desde cero.

Para los reportes se usó un formato de tablas, que facilita encontrar tokens específicos y leer su información.

Para el token de "diferente", que en la mayoría de lenguajes de programación se representa su función con , se decidió usar  $\Sigma=$ , esto debido a que el  $\Sigma$  es el NOT en este lenguaje, así que se decidió que lo mejor era cambiar el NOT del diferente por un  $\Sigma$  también.

## Algoritmos usados

1. Construcción de StringBuffer para cadenas: cuando el lexer encuentra una comilla doble, inicializa StringBuffer vacío, agrega la comilla inicial, cambia el estado a STRING, a partir de ahí si encuentra un \ lo toma como secuencia de escape, si encuentra un " agrega la comilla final y retorna el token, cualquier otro elemento se agrega al buffer y al terminar se retorna el STRING\_LITERAL que posee todo el contenido.
2. Recuperación de errores: cuando el lexer se encuentra un carácter ilegal, imprime el error, entra en el modo pánico (PANIC\_MODE) que ya fue mencionado anteriormente, ignora todo hasta el salto de línea (\n) y al llegar a este vuelve a estado YYINITIAL y continua el análisis normalmente.
3. Búsqueda de nombres de tokens: usa reflexión, carga el sym.java, obtiene todos los campos, compara su valor con el id buscado y retorna su nombre al encontrarlo, si no lo encontrara retornaría un "UNKNOWN"
4. Análisis de un archivo txt: abre el archivo con UTF-8, se escribe un encabezado para la tabla del reporte, se entra en un bucle hasta que no se llegue al EOF (fin del archivo) y mientras que está en ese bucle obtiene el siguiente token, si es un error lo agrega a la lista de errores, extrae el lexema, el tipo de token y su línea y columna, lo escribe en una fila de la tabla y repite con el siguiente token... Al final se escriben los errores en la sección de errores si los hubiera y finalmente se cierran los archivos txt.



5. Menú: consiste en switch y en if y else básicos que leen un número y siguen el recorrido indicado.

## Resumen y análisis del proceso que hace jflex para generar el analizador léxico

Lo primero que JFlex hace es una lectura de las especificaciones que le damos en el archivo `lexer.jflex` y lo divide en código de usuario, opciones y macros y reglas léxicas, dichas reglas funcionan de una forma patrón y acción: "patron" { acción },

Ejemplo: "world" { return symbol(sym.WORLD); }

Significa que cuando se encuentra el patrón indicado, realiza la acción que está entre las llaves { }.

Luego de esto se convierten las expresiones regulares, JFlex usa el algoritmo de Thompson para construir el Autómata Finito No Determinista (AFN), por ejemplo para el patrón {identifier} = [a-zA-Z\_][a-zA-Z0-9\_]\* se crea un estado inicial, se crean las transiciones y ciclos entre estados intermedio y se marca cuál será el estado final como aceptador.

A partir de esto se convierte ese AFN en un AFD (Autómata Finito Determinista), primero se tiene el conjunto de estados alcanzables del AFN, se calcula a qué conjunto de estados se puede llegar y si es nuevo se agrega como estado del AFD y esto se va a repetir hasta que no haya conjuntos nuevos. Luego este estado se minimiza con un algoritmo llamado minimización de Hopcroft que básicamente lo que hace es agrupar estados equivalentes que se comportan de manera idéntica separando aceptadores de no aceptadores hasta reducir el grafo todo lo que sea posible. Así quedan muchos menos estados.

Al final a partir de esto se genera el código `Lexer.java`, que es el que permite hacer los análisis léxicos, tiene una tabla de transiciones y otra de acciones y lo que hace es ir leyendo los caracteres, consultar la tabla de transiciones, ver si es estado aceptador para ejecutar una acción y si hay transiciones inválidas reporta errores. También hay un buffer que maneja la lectura.

## Librerías usadas

### 1. Librerías Externas (Archivos .jar)

**JFlex 1.9.1** (`jflex-full-1.9.1.jar`): generar el analizador léxico a partir de `Lexer.jflex`.

**Java CUP 11b - Binario** (java-cup-11b.jar): procesa Parser.cup para generar sym.java con los IDs de tokens, también genera Parser.java (no usado en esta fase).

**Java CUP 11b - Runtime** (java-cup-11b-runtime.jar): librería necesaria en tiempo de ejecución.

## 2. Paquetes del JDK (Biblioteca Estándar de Java)

**java.io:** gestión de rutas y archivos, lectura de archivos y generación de reportes con UTF-8.

**java.util:** para el menú interactivo por consola y poder usar listas.

**java.lang.reflect:** para cargar dinámicamente las clases generadas y extraer nombres de tokens de sym.java

## Análisis de resultados

Se cumplieron todos los objetivos y requerimientos del proyecto.

### Objetivos alcanzados

- **Construcción de un scanner funcional:** se logró desarrollar un analizador léxico completo utilizando JFlex, capaz de identificar todos los terminales especificados, incluyendo palabras reservadas, operadores, delimitadores y literales.
- **Gestión de caracteres especiales:** se implementó exitosamente el soporte para caracteres exigidos por el lenguaje, como los delimitadores de bloque (¡, !), paréntesis (¿, ?), operadores lógicos ( $\Sigma$ , @, ~) y símbolos de comentarios (€, ə), usando UTF-8.
- **Recuperación de errores en modo pánico:** se logró no detener el análisis ante errores. El sistema identifica el error, lo reporta y sigue el análisis saltando a la siguiente línea del archivo fuente.
- **Generación de los reportes:** el programa produce archivos de salida estructurados que vinculan cada lexema con su identificador numérico, línea y columna, además de mostrar los errores si los hay.

# Bitácora

|  |                                  |                  |
|--|----------------------------------|------------------|
| Commits on Dec 26, 2025  |                                  |                  |
| Docs: documentación pdf parte 1  | Rodiohs committed 17 minutes ago | cf7c26d          |
| docs: Video de pruebas   | JairoGH16 committed 2 hours ago  | b9f868a          |
| fix: incluir comillas en lexemas de strings y reportar errores de cierres de comentarios | JairoGH16 committed 5 hours ago  | 0629bff          |
| Feat: reporte de errores detallado   | Rodiohs committed 6 hours ago    | 9a5c5a6          |
| feat: manejo de errores léxicos en modo pánico y token ERROR                             | JairoGH16 committed 9 hours ago  | 2bf0f52          |
| Feat: main y menú  | Rodiohs committed 10 hours ago   | 69cf1ec          |
| Commits on Dec 25, 2025  |                                  |                  |
| feat: implementacion de especificaciones léxicas y terminales en lexer y parser          | JairoGH16 committed yesterday    | 9187074          |
| vincular librerías de CUP/JFlex  | JairoGH16 committed yesterday    | fe38501          |
| Feat: Inicializar estructura código y JFlex/CUP  | Rodiohs committed yesterday      | 83f34d7          |
| Commits on Dec 24, 2025  |                                  |                  |
| estructura de carpetas   | JairoGH16 committed 2 days ago   | e3ec1c8          |
| Commits on Dec 20, 2025  |                                  |                  |
| Initial commit   | JairoGH16 authored last week     | Verified 35944f0 |

| Fecha                   | Actividad / Commit   | Encargado |
|-------------------------|--|-----------|
| 20 de diciembre de 2025 | Initial commit: Creación del repositorio inicial.  | Jairo     |
| 24 de diciembre de 2025 | Estructura de carpetas: se estructuran las carpetas a cómo se van a usar en el proyecto.   | Jairo     |
| 25 de diciembre de 2025 | Feat: Inicializar estructura código y JFlex/CUP: se incluyen las bases para Lexer.jflex y Parser.cup, se agregan las librerías externas (JFlex y CUP) en el directorio lib | Rafael    |

|                         |   |        |
|-------------------------|---|--------|
| 25 de diciembre de 2025 | <b>Vincular librerías de CUP/JFlex:</b> se conectan al proyecto las librerías externas .jar.  | Jairo  |
| 25 de diciembre de 2025 | <b>Feat: implementación de especificaciones léxicas y terminales en lexer y parser:</b> se definieron todos los tokens y operadores, se configuraron los terminales en Parser.cup para la generación de sym.java, se generaron las clases Lexer.java y sym.java, se creó un archivo de prueba inicial (fuente.txt) para usar después.           | Jairo  |
| 26 de diciembre de 2025 | <b>Feat: Main y menú:</b> Se implementó un main para generar los archivos lexer, sym y parser y hacer el análisis y reportes, se hizo un menú básico con soporte para excepciones básicas como error en el análisis o carpeta de archivos vacía.  | Rafael |
| 26 de diciembre de 2025 | <b>Feat: manejo de errores léxicos en modo pánico y token ERROR:</b> Se agregó el terminal ERROR en Parser.cup para identificar lexemas inválidos, se implementó la recuperación en modo pánico en Lexer.jflex para ignorar el resto de la línea y saltar a la que sigue.   | Jairo  |
| 26 de diciembre de 2025 | <b>Feat: Reporte de errores detallado:</b> Se agregó al main agregar los errores a una lista conforme van saliendo y añadirlos al final de cada reporte.  | Rafael |
| 26 de diciembre de 2025 | <b>fix: incluir comillas en lexemas de strings y reportar errores de cierres de comentarios:</b> Se corrigió el estado STRING para incluir las comillas como parte del lexema, se implementó el estado MULTICOMMENT para validar cierres para comentarios multilinea, reportando un error específico si se alcanza el EOF sin cerrar el bloque. | Jairo  |

|  |   |        |
|--|---|--------|
| 26 de diciembre de 2025                    | <b>Docs: Video de pruebas:</b> video donde se prueba el analizador.   | Rafael |
| 26 de diciembre de 2025                    | <b>Docs: Documentación Parte 1:</b> se adjunta la primera mitad de la documentación.  | Rafael |
| 26 de diciembre de 2025<br>(commit actual) | <b>Docs: Documentación Parte 2:</b> se adjunta la segunda mitad de la documentación externa y se completa la documentación interna. | Jairo  |

## Bibliografía

**JFlex Manual** (Version 1.9.1): <https://jflex.de/manual.html>

**CUP Manual.** Hudson, Scott E: <http://www2.cs.tum.edu/projects/cup/>

**Oracle Java Documentation.** Oracle Corporation:  
<https://docs.oracle.com/en/java/>