

## 1. Expresiones Regulares (Léxico)

**Identificadores (id):** Comienzan con una letra, seguida de letras o números.

LETRA = [a-zA-Z]

DIGITO = [0-9]

ID = {LETRA}({LETRA}|{DIGITO})\*

**Constantes Enteras (cte\_int):** Uno o más dígitos.

CTE\_INT = {DIGITO}+

**Constantes Flotantes (cte\_float):** Dígitos, un punto decimal, y más dígitos.

CTE\_FLOAT = {DIGITO}+\.{DIGITO}+

**Constantes String (cte\_string):** Caracteres entre comillas dobles.

CTE\_STRING = \"^[^']\*\"

**Palabras Reservadas (Keywords):** Son literales fijos.

program  
main  
end  
var  
int  
float  
void  
if  
else  
while  
do  
print

**Operadores:**

+ (Suma)  
- (Resta)  
\* (Multiplicación)  
/ (División)  
> (Mayor que)  
< (Menor que)  
!= (Diferente de)  
= (Asignación)

## **Separadores / Puntuación:**

; (Punto y coma)  
: (Dos puntos)  
, (Coma)  
{ (Llave izquierda)  
} (Llave derecha)  
( (Paréntesis izquierdo)  
) (Paréntesis derecho)

## **2. Lista de Tokens**

Estos son los nombres de los tokens que el analizador léxico (scanner) generaría:

TOKEN\_PROGRAM (para program)  
TOKEN\_MAIN (para main)  
TOKEN\_END (para end)  
TOKEN\_VAR (para var)  
TOKEN\_INT (para int)  
TOKEN\_FLOAT (para float)  
TOKEN\_VOID (para void)  
TOKEN\_IF (para if)  
TOKEN\_ELSE (para else)  
TOKEN\_WHILE (para while)  
TOKEN\_DO (para do)  
TOKEN\_PRINT (para print)  
TOKEN\_ID (para identificadores, ej. miVariable)  
TOKEN\_CTE\_INT (para constantes enteras, ej. 123)  
TOKEN\_CTE\_FLOAT (para constantes flotantes, ej. 3.1416)  
TOKEN\_CTE\_STRING (para constantes string, ej. "Hola")  
TOKEN\_PLUS (para +)  
TOKEN\_MINUS (para -)  
TOKEN\_MULT (para \*)  
TOKEN\_DIV (para /)  
TOKEN\_GT (para >)  
TOKEN\_LT (para <)  
TOKEN\_NE (para !=)  
TOKEN\_ASSIGN (para =)  
TOKEN\_SEMICOLON (para ;)  
TOKEN\_COLON (para :)  
TOKEN\_COMMA (para ,)  
TOKEN\_LBRACE (para {)  
TOKEN\_RBRACE (para })  
TOKEN\_LPAREN (para ()  
TOKEN\_RPAREN (para ))

### 3. Gramática Libre de contexto

$\langle \text{programa} \rangle ::= \text{TOKEN\_PROGRAM} \text{ TOKEN\_ID} \text{ TOKEN\_SEMICOLON} [\langle \text{vars} \rangle] [\langle \text{funcs} \rangle] \text{ TOKEN\_MAIN} \langle \text{cuerpo} \rangle \text{ TOKEN\_END}$

$\langle \text{vars} \rangle ::= \text{TOKEN\_VAR} \langle \text{lista\_vars\_decl} \rangle \mid \epsilon$

$\langle \text{lista\_vars\_decl} \rangle ::= \langle \text{declaracion} \rangle [\langle \text{lista\_vars\_decl} \rangle]$

$\langle \text{declaracion} \rangle ::= \text{TOKEN\_ID} \text{ TOKEN\_COLON} \langle \text{tipo} \rangle \text{ TOKEN\_SEMICOLON}$

$\langle \text{tipo} \rangle ::= \text{TOKEN\_INT} \mid \text{TOKEN\_FLOAT}$

$\langle \text{funcs} \rangle ::= \langle \text{funcion} \rangle [\langle \text{funcs} \rangle] \mid \epsilon$

$\langle \text{funcion} \rangle ::= \text{TOKEN\_VOID} \text{ TOKEN\_ID} \text{ TOKEN\_LPAREN} [\langle \text{parametros} \rangle] \text{ TOKEN\_RPAREN} \text{ TOKEN\_LBRACE} [\langle \text{vars} \rangle] \langle \text{cuerpo} \rangle \text{ TOKEN\_RBRACE}$

$\langle \text{parametros} \rangle ::= \text{TOKEN\_ID} \text{ TOKEN\_COLON} \langle \text{tipo} \rangle [\langle \text{mas\_parametros} \rangle] \mid \epsilon$

$\langle \text{mas\_parametros} \rangle ::= \text{TOKEN\_COMMA} \text{ TOKEN\_ID} \text{ TOKEN\_COLON} \langle \text{tipo} \rangle [\langle \text{mas\_parametros} \rangle] \mid \epsilon$

$\langle \text{cuerpo} \rangle ::= \text{TOKEN\_LBRACE} [\langle \text{lista\_estatutos} \rangle] \text{ TOKEN\_RBRACE}$

$\langle \text{lista\_estatutos} \rangle ::= \langle \text{estatuto} \rangle [\langle \text{lista\_estatutos} \rangle] \mid \epsilon$

$\langle \text{estatuto} \rangle ::= \langle \text{asignacion} \rangle \text{ TOKEN\_SEMICOLON} \\ \mid \langle \text{condicion} \rangle \text{ TOKEN\_SEMICOLON} \\ \mid \langle \text{ciclo} \rangle \text{ TOKEN\_SEMICOLON} \\ \mid \langle \text{llamada\_func} \rangle \text{ TOKEN\_SEMICOLON} \\ \mid \langle \text{impresion} \rangle \text{ TOKEN\_SEMICOLON}$

$\langle \text{asignacion} \rangle ::= \text{TOKEN\_ID} \text{ TOKEN\_ASSIGN} \langle \text{expresion} \rangle$

$\langle \text{condicion} \rangle ::= \text{TOKEN\_IF} \text{ TOKEN\_LPAREN} \langle \text{expresion} \rangle \text{ TOKEN\_RPAREN} \\ \langle \text{cuerpo} \rangle [\langle \text{parte\_else} \rangle]$

$\langle \text{parte\_else} \rangle ::= \text{TOKEN\_ELSE} \langle \text{cuerpo} \rangle \mid \epsilon$

$\langle \text{ciclo} \rangle ::= \text{TOKEN\_WHILE} \text{ TOKEN\_LPAREN} \langle \text{expresion} \rangle \text{ TOKEN\_RPAREN} \\ \text{TOKEN\_DO} \langle \text{cuerpo} \rangle$

$\langle \text{impresion} \rangle ::= \text{TOKEN\_PRINT} \text{ TOKEN\_LPAREN} \langle \text{lista\_impresion} \rangle \\ \text{TOKEN\_RPAREN}$

$\langle \text{lista\_impresion} \rangle ::= \langle \text{arg\_impresion} \rangle [\langle \text{mas\_args\_impresion} \rangle]$

$\langle \text{arg\_impresion} \rangle ::= \langle \text{expresion} \rangle \mid \text{TOKEN\_CTE\_STRING}$

$\langle \text{mas\_args\_impresion} \rangle ::= \text{TOKEN\_COMMA} \langle \text{arg\_impresion} \rangle$   
 $[\langle \text{mas\_args\_impresion} \rangle] \mid \epsilon$

$\langle \text{llamada\_func} \rangle ::= \text{TOKEN\_ID} \text{TOKEN\_LPAREN} [\langle \text{lista\_argumentos} \rangle]$   
 $\text{TOKEN\_RPAREN}$

$\langle \text{lista\_argumentos} \rangle ::= \langle \text{expresion} \rangle [\langle \text{mas\_argumentos} \rangle] \mid \epsilon$

$\langle \text{mas\_argumentos} \rangle ::= \text{TOKEN\_COMMA} \langle \text{expresion} \rangle [\langle \text{mas\_argumentos} \rangle] \mid \epsilon$

$\langle \text{expresion} \rangle ::= \langle \text{exp} \rangle [ \langle \text{op\_relacional} \rangle \langle \text{exp} \rangle ]$

$\langle \text{op\_relacional} \rangle ::= \text{TOKEN\_GT} \mid \text{TOKEN\_LT} \mid \text{TOKEN\_NE}$

$\langle \text{exp} \rangle ::= \langle \text{termino} \rangle [ \langle \text{op\_suma\_resta} \rangle \langle \text{exp} \rangle ]$

$\langle \text{op\_suma\_resta} \rangle ::= \text{TOKEN\_PLUS} \mid \text{TOKEN\_MINUS}$

$\langle \text{termino} \rangle ::= \langle \text{factor} \rangle [ \langle \text{op\_mult\_div} \rangle \langle \text{termino} \rangle ]$

$\langle \text{op\_mult\_div} \rangle ::= \text{TOKEN\_MULT} \mid \text{TOKEN\_DIV}$

$\langle \text{factor} \rangle ::= \text{TOKEN\_LPAREN} \langle \text{expresion} \rangle \text{TOKEN\_RPAREN}$   
 $\mid \text{TOKEN\_ID}$   
 $\mid \langle \text{cte} \rangle$

$\langle \text{cte} \rangle ::= \text{TOKEN\_CTE\_INT} \mid \text{TOKEN\_CTE\_FLOAT}$

## Comparación de herramientas

Característica	Flex (Scanner)	Bison (Parser)	ANTLR (Scanner + Parser)	JavaCC (Scanner + Parser)	Otros (Ej: PEG, Combinators)
Tipo	Generador de Scanner	Generador de Parser	Generador Scanner+Parser	Generador Scanner+Parser	Bibliotecas/Generadores
Algoritmo Parsing	N/A (Exp. Regulares)	LALR(1), GLR	LL(*), Adaptativo	LL(k)	Recursive Descent, Packrat
Notación Gramática	N/A	BNF (similar a Yacc)	EBNF	EBNF	PEG, Código Host
Lenguaje Salida	C, C++	C, C++, Java, D	Java, C#, Python, JS, C++, Go, Swift, PHP	Java, C++	Varia (Java, Scala, Haskell, C++, Python...)
Integración Lex/Parse	Separado (requiere Parser)	Separado (requiere Scanner)	Integrado	Integrado	Generalmente Integrado
Generación AST	No directa	Requiere acciones semánticas	Soporte integrado (listeners/visitors)	Soporte integrado	Generalmente integrado
Manejo Errores	Básico	Moderado (requiere esfuerzo)	Avanzado	Moderado	Variable, a menudo manual
Soporte Unicode	Limitado/Indirecto	N/A	Sí	Sí	Variable
Licencia	BSD	GPLv3 con excepción	BSD	BSD	Variable (MIT, Apache, etc.)
Madurez/Actividad	Muy Maduro, Estable	Muy Maduro, Estable	Maduro, Muy Activo	Maduro, Menos Activo	Variable
Curva Aprendizaje	Moderada (RegEx)	Moderada/Alta (LR)	Moderada (LL)	Moderada (LL)	Variable
IDE/Herramientas	No específicas	No específicas	Sí (ANTLRWorks, plugins)	No específicas	Variable

**Para el desarrollo del compilador del lenguaje BabyDuck, se seleccionaron las herramientas Flex y Bison.**

**Madurez y Estabilidad:** Flex y Bison son herramientas extremadamente maduras, probadas en innumerables proyectos durante décadas.

**Estándar de Facto en C/C++:** Son las herramientas tradicionales y estándar para generar scanners y parsers en C y C++. Se integran de forma natural en flujos de trabajo basados en C/C++, que es un lenguaje común para la implementación de compiladores.

**Rendimiento:** Los analizadores generados en C por Flex y Bison son conocidos por su alta velocidad de ejecución, lo cual es importante para herramientas de compilación.

**Disponibilidad:** Forman parte del conjunto de herramientas estándar en la mayoría de los sistemas operativos tipo Unix (Linux, macOS) y están fácilmente disponibles para Windows. No suelen requerir dependencias complejas.

**Recursos Educativos:** Existe una vasta cantidad de documentación, tutoriales, libros y ejemplos disponibles para Flex y Bison, lo que facilita el aprendizaje y la resolución de problemas. Son herramientas frecuentemente utilizadas en cursos de teoría de compiladores.

**Adecuación al Problema:** Para un lenguaje relativamente simple como BabyDuck, las capacidades de Flex para el léxico definido con expresiones regulares y Bison para la gramática, son perfectamente adecuadas. No se requieren las características más avanzadas de herramientas como ANTLR.

## 4. Reglas Léxicas (Scanner - Flex)

Las expresiones regulares se definieron en un archivo `scanner.l` utilizando la sintaxis de Flex:

**Sección de Definiciones:** Se usaron alias para patrones comunes (LETRA [a-zA-Z], DIGITO [0-9]).

**Sección de Reglas** (`%% . . . %%`): Se especificaron pares patrón { acción }.

Los patrones corresponden a las expresiones regulares definidas (ej. {ID}, {CTE\_INT}, "program", "+", etc.).

Para tokens con valor asociado (IDs, constantes), se utiliza la variable `yylval` (definida en Bison) para almacenar el valor.

Se incluye una regla para manejar caracteres no reconocidos.

## 5. Reglas Gramaticales (Parser - Bison)

La gramática libre de contexto se definió en un archivo `parser.y` usando la sintaxis de Bison:

**Sección de Declaraciones** (`%{ . . . %}` y `%token, %union, etc.`):

Se declaró la unión `yy1val` para manejar los tipos de datos asociados a los tokens.

Se declararon todos los nombres de los tokens (`%token TOKEN_PROGRAM TOKEN_ID ...`), asociando aquellos con valores a un campo de la unión (`%token <sval> TOKEN_ID`).

Se definió el símbolo inicial de la gramática (`%start programa`).

### **Sección de Reglas Gramaticales** (`%% ... %%`):

Las reglas se escribieron en una notación similar a BNF:

```
simbolo_no_terminal : cuerpo_regla_1 | cuerpo_regla_2 ...  
;
```

El cuerpo de la regla consiste en secuencias de terminales y no terminales.