



Instituto Tecnológico y de Estudios Superiores de Monterrey

TC3002B Desarrollo de Aplicaciones avanzadas de Ciencias Computacionales Gpo 502

Patito – Entrega #1

José Eduardo de Valle Lara

A01734957

Profesora:

Elda Guadalupe Quiroga González

Contenido

Etapa 0.....	3
Tokens y Expresiones Regulares.....	3
Definición de la gramática libre de contexto.....	5
Etapa 1.....	9
Desarrollo del Scanner y Parser.....	9
Desarrollo.....	9
Código del lenguaje.....	9
Workflow.....	11
Test-Plan.....	12
Pruebas temporales.....	12
Hallazgos.....	14

Etapa 0

Tokens y Expresiones Regulares

Se agruparon los tokens con sus respectivas expresiones regulares en una misma sección para simplificar el documento. Las expresiones regulares de las palabras reservadas y algunos operadores son el mismo token ya que son patrones que no cambian, son fijos.

Token	Expresión Regular	Explicación
ID	[a-zA-Z][a-zA-Z0-9]*	Se usa para nombrar variables y funciones, debe comenzar con una letra y puede ser seguido por cualquier combinación de letras y números
Constantes: valores fijos que no cambian durante la ejecución del programa		
Enteros (cte_int)	[0-9]+	Debe haber al menos un número
Flotantes (cte_float)	[0-9]+\.[0-9]+	Empiezan con al menos uno o muchos números seguido de un punto "." y de al menos un número después del punto
Strings (cte.string)	\"[^"]*\"	Empiezan con una comilla doble " seguida de cero o más repeticiones de cualquier caracter que no sea una comilla doble y finaliza con una comilla doble
Operadores y símbolos: se usan para realizar operaciones, comparaciones y agrupaciones		
+	\+	Se usó un backslash \ para indicar la expresión regular de + ya que en expresiones regulares se refiere a uno o más elementos
-	-	
*	*	Se usó un backslash \ para indicar la expresión regular de * ya que en expresiones regulares se refiere a cero o más elementos
/	/	
>	>	

<	<	
!=	!=	
=	=	
;	;	
:	:	
,	,	
(\(Se usó un backslash \ para indicar la expresión regular de (porque se usa para agrupaciones en expresiones regulares
)	\)	Se usó un backslash \ para indicar la expresión regular de) porque se usa para agrupaciones en expresiones regulares
{	\{	Se usó un backslash \ para indicar la expresión regular de { porque se usa para cuantificadores en expresiones regulares
}	\}	Se usó un backslash \ para indicar la expresión regular de } porque se usa para cuantificadores en expresiones regulares
Palabras reservadas: identificadores con un significado especial que no pueden usarse como nombres de variables o funciones		
program	program	
main	main	
end	end	
var	var	
int	int	
float	float	
void	void	
if	if	
else	else	
while	while	

do	do	
print	print	

Definición de la gramática libre de contexto

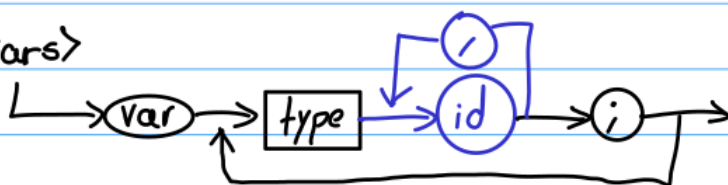
$\langle \text{Programa} \rangle$

$\langle \text{Programa} \rangle \rightarrow \text{program id ; } \langle \text{Opt_vars} \rangle \langle \text{Opt_funcs} \rangle \text{ main } \langle \text{Body} \rangle \text{ end}$

$\langle \text{Opt_vars} \rangle \rightarrow \langle \text{Vars} \rangle \mid \epsilon$

$\langle \text{Opt_funcs} \rangle \rightarrow \langle \text{FUNCS} \rangle \mid \epsilon$

$\langle \text{Vars} \rangle$



$\langle \text{Vars} \rangle \rightarrow \text{var } \langle \text{DEC_VAR} \rangle$

$\langle \text{DEC_VAR} \rangle \rightarrow \langle \text{TYPE} \rangle \langle \text{LIST_id} \rangle ; \langle \text{DV}' \rangle$

↑
puede tener 1
o más ids

↑
puede tener un
ciclo o no

$\langle \text{DV}' \rangle \rightarrow \langle \text{DEC_VAR} \rangle \mid \epsilon$

$\langle \text{FUNCS} \rangle$

$\langle \text{FUNCS} \rangle \rightarrow \langle \text{FUNC_DEF} \rangle$

$\langle \text{FUNC_DEF} \rangle \rightarrow \langle \text{RETURN_TYPE} \rangle \text{id} (\langle \text{PARAMS} \rangle) \langle \text{OPTIONAL_LOCAL_VARS} \rangle$
 $\langle \text{Body} \rangle ;$

$\langle \text{RETURN_TYPE} \rangle \rightarrow \text{void} \mid \langle \text{TYPE} \rangle$

$\langle \text{PARAMS} \rangle \rightarrow \langle \text{PARAM_LIST} \rangle \mid \epsilon$

$\langle \text{PARAM_LIST} \rangle \rightarrow \text{id} : \langle \text{TYPE} \rangle \langle \text{MORE_PARAMS} \rangle$

$\langle \text{MORE_PARAMS} \rangle \rightarrow , \langle \text{PARAM_LIST} \rangle \mid \epsilon$

$\langle \text{OPTIONAL_LOCAL_VARS} \rangle \rightarrow \langle \text{VARS} \rangle \mid \epsilon$

$\langle \text{Body} \rangle$

$\langle \text{Body} \rangle \rightarrow \{ \langle \text{STATEMENT_LIST} \rangle \}$

$\langle \text{STATEMENT_LIST} \rangle \rightarrow \langle \text{STATEMENT} \rangle \langle \text{STATEMENT_LIST} \rangle \mid \epsilon$

$\langle \text{TYPE} \rangle$

$\langle \text{TYPE} \rangle \rightarrow \text{int} \mid \text{float}$

$\langle \text{STATEMENT} \rangle$

Esta regla presenta una ambigüedad en $\langle \text{ASSIGN} \rangle$ y $\langle \text{F_call} \rangle$ ya que las 2 invocan un id, el primero es para asignar variables y el segundo es para funciones, se resolvió con el factor común izquierdo

$\langle \text{STATEMENT} \rangle \rightarrow \text{id} \langle \text{ID_STATEMENT_CONTINUATION} \rangle \mid \langle \text{CONDITION} \rangle$
 $\mid \langle \text{CYCLE} \rangle \mid \langle \text{Print} \rangle$

$\langle \text{ID_STATEMENT_CONTINUATION} \rangle \rightarrow = \langle \text{EXPRESION} \rangle ; \mid (\langle \text{EXPRESION} \rangle) ;$

<ASSIGN>

<ASSIGN> → id = <EXPRESION>;

<CONDITION>

<CONDITION> → if (<EXPRESION>) <Body> <OPTIONAL_ELSE>;

<OPTIONAL_ELSE>; → else <Body> | ε

<CYCLE>

<CYCLE> → while (<EXPRESION>) do <Body>

<F.Call >

<F.Call> → id (<EXPRESION>);

<Print>

<Print> → print (<PRINT_LIST>);

<PRINT_LIST> → <PRINT_ITEM> <PRINT_MORE_ITEMS>

<PRINT_ITEM> → <EXPRESION> | cte.string

<PRINT_MORE_ITEMS> → , <PRINT_ITEM> <PRINT_MORE_ITEMS> | ε

<EXPRESSION>

<EXPRESSION> → <EXP> <OPTIONAL_COMPARISON>

<OPTIONAL_COMPARISON> → <RELATIONAL_OPERATOR> <EXP> | ε

<RELATIONAL_OPERATOR> → > | < | !=

$\langle \text{EXP} \rangle$

$\langle \text{EXP} \rangle \rightarrow \langle \text{TERMINO} \rangle \langle \text{EXP_PRIME} \rangle$

$\langle \text{EXP_PRIME} \rangle \rightarrow + \langle \text{TERMINO} \rangle \langle \text{EXP_PRIME} \rangle \mid - \langle \text{TERMINO} \rangle \langle \text{EXP_PRIME} \rangle \mid \varepsilon$

$\langle \text{TERMINO} \rangle$

$\langle \text{TERMINO} \rangle \rightarrow \langle \text{FACTOR} \rangle \langle \text{PRIME_TERM} \rangle$

$\langle \text{PRIME_TERM} \rangle \rightarrow * \langle \text{FACTOR} \rangle \langle \text{PRIME_TERM} \rangle \mid / \langle \text{FACTOR} \rangle \langle \text{PRIME_TERM} \rangle \mid \varepsilon$

$\langle \text{FACTOR} \rangle$

$\langle \text{FACTOR} \rangle \rightarrow (\langle \text{EXPRESION} \rangle) \mid \text{id} \mid \langle \text{CTE} \rangle$

$\langle \text{CTE} \rangle$

$\langle \text{CTE} \rangle \rightarrow \text{cte_int} \mid \text{cte_float}$

Etapa 1

Desarrollo del Scanner y Parser

Desarrollo

Se usó Python con ANTLR para el desarrollo del compilador de BabyDuck. Los tokens y las reglas de gramática se declaran en el archivo principal con extensión **.g4**. Por convención, el nombre de los tokens tiene que ser en mayúsculas mientras que los nombres de las reglas deben empezar en minúsculas, ANTLR sigue una convención opuesta a lo que se enseñó en clases.

En clases se mostró que la forma de representar un camino vacío es con un ϵ y así se crearon las gramáticas libres de contexto a mano, pero fueron transformadas ligeramente al momento de ingresarlas a ANTLR, todas las reglas que representaban un elemento opcional o una repetición se plasmaron con los símbolos de ANTLR para simplificar la lectura de las reglas, se ignoraron las reglas ...OPCIONAL y ...PRIMA de la primera versión de la gramática libre de contexto por los símbolos $*$, $+$ y $?$ de las expresiones regulares que son soportados por ANTLR.

De momento los tokens y gramáticas se encuentran en el mismo archivo ya que ANTLR solo lee un archivo principal. No se tiene contemplado separar el contenido en distintos archivos debido a que se implementó toda la lógica del lexer y parser en un archivo sencillo de leer. Todas las etapas que siguen son con respecto al análisis semántico y tabla de símbolos y la síntesis (backend), eso ya está afuera del alcance de ANTLR y requiere programación manual en Python

Código del lenguaje

Definición de símbolos terminales

```
53 // Palabras reservadas
54 PROGRAMA: 'programa';
55 INICIO: 'inicio';
56 FIN: 'fin';
57 VARS: 'vars';
58 ENTERO: 'entero';
59 FLOTANTE: 'flotante';
60 SI: 'si';
61 SINO: 'sino';
62 ESCRIBE: 'escribe';
63 MIENTRAS: 'mientras';
64 HAZ: 'haz';
65 NULA: 'nula';
```

```

67 // Operadores y puntuaciones
68 PARENTESIS_IZQUIERDO: '(';
69 PARENTESIS_DERECHO: ')';
70 CORCHETE_IZQUIERDO: '[';
71 CORCHETE_DERECHO: ']';
72 LLAVE_IZQUIERDA: '{';
73 LLAVE_DERECHA: '}';
74 PUNTO_Y_COMA: ',';
75 DOS_PUNTOS: ':';
76 COMA: ',';
77 ASIGNACION: '=';
78 MAS: '+';
79 MENOS: '-';
80 MULTIPLICACION: '*';
81 DIVISION: '/';
82 MAYOR_QUE: '>';
83 MENOR_QUE: '<';
84 DIFERENTE_DE: '!=';
85 IGUAL_QUE: '==';

```

Definición de expresiones regulares

```

87 // IDs y constantes
88 ID: [a-zA-Z_] [a-zA-Z_0-9]*;
89 CTE_ENT: [0-9]+;
90 CTE_FLOT: [0-9]+ '.' [0-9]+;
91 LETRERO: '"' .*? '"';
92
93 // Ignorar espacios en blanco y saltos de línea
94 ESPACIOS: [ \t\r\n]+ -> skip;

```

Definición de reglas

```

1 grammar BabyDuck;
2
3 programa: PROGRAMA ID PUNTO_Y_COMA vars? funcs* INICIO cuerpo FIN ;
4
5 vars: VARS declarar_variables+ ;
6 declarar_variables: declarar_ids DOS_PUNTOS tipo PUNTO_Y_COMA ;
7 declarar_ids: ID (COMA ID)* ;
8
9 funcs: (NULA | tipo) ID PARENTESIS_IZQUIERDO parametros? PARENTESIS_DERECHO LLAVE_IZQUIERDA vars? cuerpo LLAVE_DERECHA PUNTO_Y_COMA ;
10 parametros: ID DOS_PUNTOS tipo (COMA ID DOS_PUNTOS tipo)* ;
11
12 cuerpo: LLAVE_IZQUIERDA estatuto* LLAVE_DERECHA ;
13
14 tipo: ENTERO | FLOTANTE ;
15
16 estatuto: ID continuacion_de_estatuto_id | condicion | ciclo | imprime | CORCHETE_IZQUIERDO estatuto* CORCHETE_DERECHO ;
17 continuacion_de_estatuto_id: ASIGNACION expresion PUNTO_Y_COMA | PARENTESIS_IZQUIERDO (expresion (COMA expresion)*)? PARENTESIS_DERECHO PUNTO_Y_COMA ;
18
19 // La regla <ASIGNA> aparece en todas las reglas como una ambigüedad, por eso su lógica se implementó
20 // manualmente en cada caso. Como consecuencia ninguna otra regla la llama y deja de ser útil para la
21 // gramática. Se eliminó pero se conservó como comentario para mostrar su lugar en el diagrama de sintaxis
22 // asigna: ID ASIGNACION expresion PUNTO_Y_COMA ;

```

```

24 condicion: SI PARENTESIS_IZQUIERDO expresion PARENTESIS_DERECHO cuerpo (SINO cuerpo)? PUNTO_Y_COMA ;
25
26 ciclo: MIENTRAS PARENTESIS_IZQUIERDO expresion PARENTESIS_DERECHO HAZ cuerpo PUNTO_Y_COMA ;
27
28 // La regla <LLAMADA> aparece en todas las reglas como una ambigüedad, por eso su lógica se implementó
29 // manualmente en cada caso. Como consecuencia ninguna otra regla la llama y deja de ser útil para la
30 // gramática. Se eliminó pero se conservó como comentario para mostrar su lugar en el diagrama de sintaxis
31 // llamada: ID PARENTESIS_IZQUIERDO (expresion (COMA expresion)*)? PARENTESIS_DERECHO ;
32
33 imprime: ESCRIBE PARENTESIS_IZQUIERDO imprimir_elementos PARENTESIS_DERECHO PUNTO_Y_COMA ;
34 imprimir_elementos: ((expresion | LETRERO) (COMA (expresion | LETRERO))*);
35
36 expresion: exp ((MAYOR_QUE | MENOR_QUE | DIFERENTE_DE | IGUAL_QUE) exp)? ;
37
38 exp: termino ((MAS | MENOS) termino)* ;
39
40 termino: factor ((MULTIPLICACION | DIVISION) factor)* ;
41
42 factor: PARENTESIS_IZQUIERDO expresion PARENTESIS_DERECHO
43 | (MAS | MENOS)? dato_o_llamada // posible ambigüedad en id, no sabe si elegir entre `id` o `id(...)`
44 ;
45 dato_o_llamada:
46 // permite elegir entre `ID` o `ID(...)`
47 ID ( PARENTESIS_IZQUIERDO (expresion (COMA expresion)*)? PARENTESIS_DERECHO )?
48 | cte
49 ;
50
51 cte: CTE_ENT | CTE_FLOT ;
52

```

Workflow

Se generó con ANTLR los archivos del lexer y parser después de haber definido los elementos principales del lenguaje, se usaron los siguientes comandos:

- Generación de lexer y parser
 - **antlr4 -Dlanguage=Python3 BabyDuck.g4** → generación de BabyDuckLexer.py, BabyDuckParser.py y otros archivos
- Correr el proyecto
 - **antlr4 BabyDuck.g4** → generar archivos de java a partir del contenido del archivo principal de ANTLR
 - **javac *.java** → compilar archivos de java a bytecode
 - **grun BabyDuck programa -tokens < HelloWorld.BabyDuck** → checar la generación de tokens de un archivo del lenguaje
 - **grun BabyDuck programa -gui < HelloWorld.BabyDuck** → visualizar árbol de sintaxis a partir de un archivo del lenguaje

Test-Plan

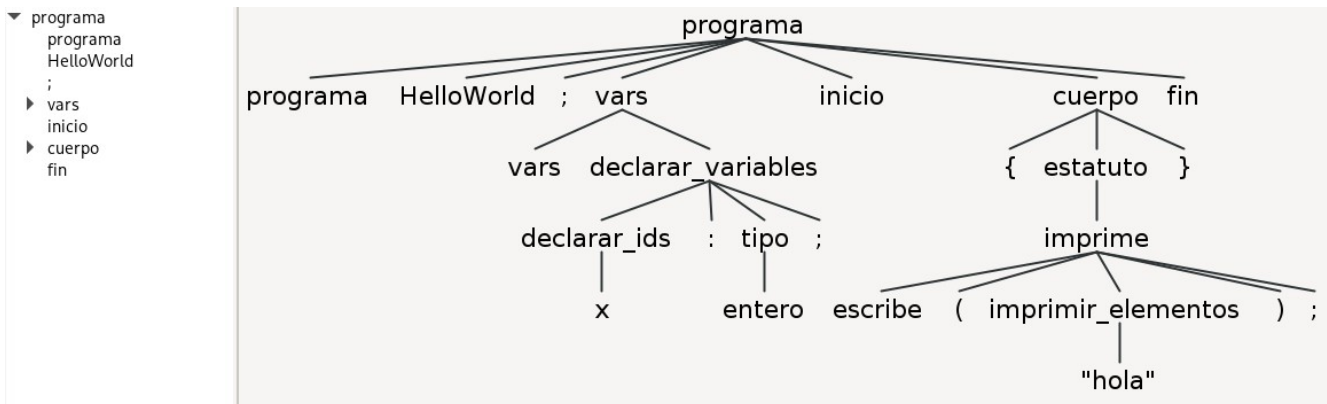
Pruebas temporales

Estas se consideran pruebas temporales debido a que son pruebas ejecutadas manualmente con el fin de demostrar la funcionalidad más básica del proyecto. Se planea usar una herramienta de automatización de pruebas como pytest para entregas a futuro que requieran una integración entre el análisis (frontend) y la síntesis (backend), no el simple funcionamiento de ambas piezas separadas.

Se realizaron 3 pruebas para esta entrega, las primeras 2 pruebas son correctas de acuerdo a la gramática pero la tercera prueba es incorrecta porque se omitió un token a propósito para probar los resultados. Los archivos de prueba se pueden encontrar en el repositorio de [GitHub](#), aquí solo se muestran los resultados al visualizar la generación de tokens y del árbol de sintaxis.

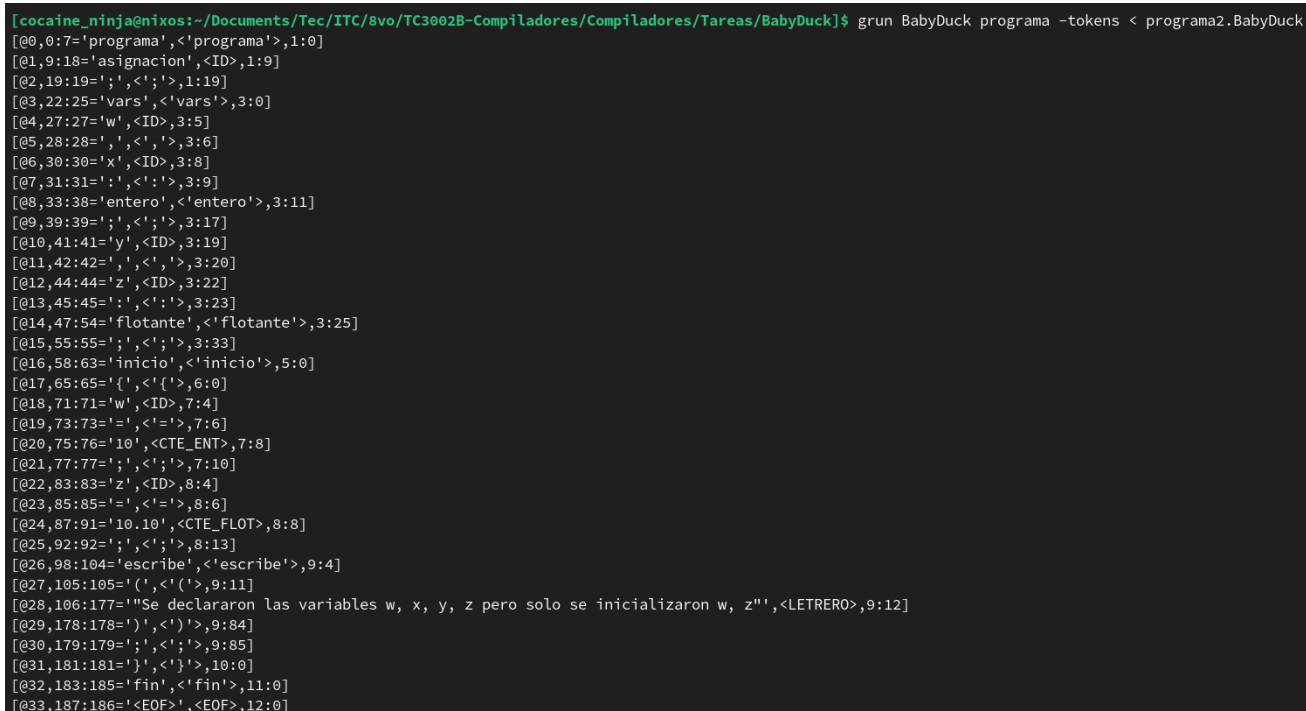
Prueba 1

Archivo: HelloWorld.BabyDuck



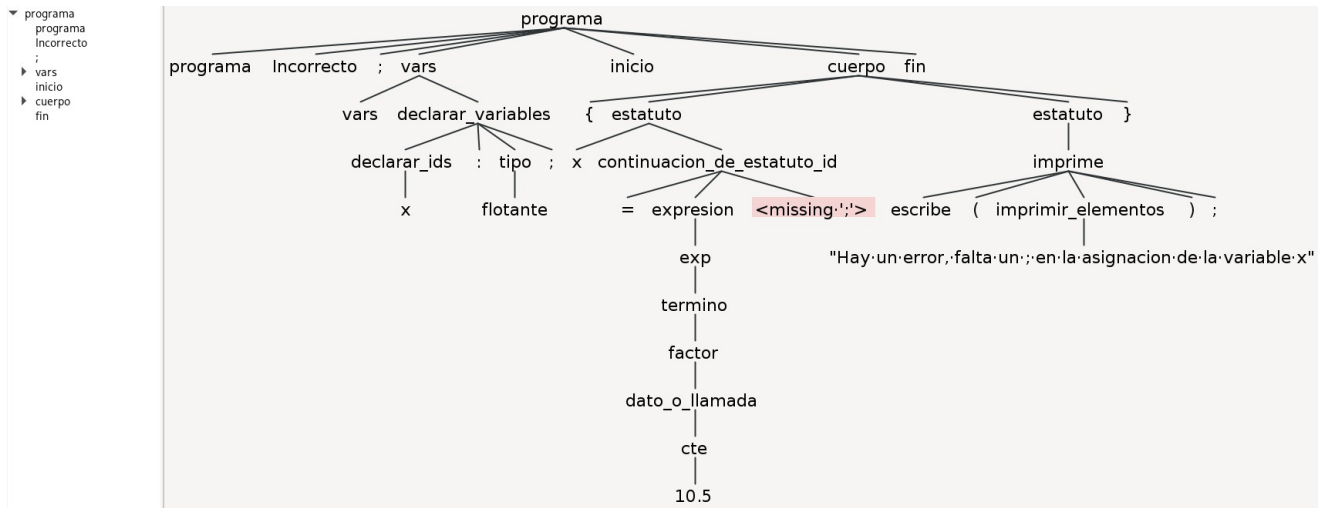
```
[cocaine_ninja@nixos:~/Documents/Tec/ITC/8vo/TC3002B-Compiladores/Compiladores/Tareas/BabyDuck]$ grun BabyDuck programa -tokens < HelloWorld.BabyDuck
[00,0:7='programa',<'programa'>,1:0]
[01,9:18='HelloWorld',<ID>,1:9]
[02,19:19=';',<'>',1:19]
[03,22:25='vars',<'vars'>,3:0]
[04,27:27='x',<ID>,3:5]
[05,28:28=':',<'>',3:6]
[06,30:35='entero',<'entero'>,3:8]
[07,36:36=';',<'>',3:14]
[08,39:44='inicio',<'inicio'>,5:0]
[09,46:46='{',<'{'>,6:0]
[10,52:58='escribe',<'escribe'>,7:4]
[11,59:59='(',<'('>,7:11]
[12,60:65='"hola"',<LETRERO>,7:12]
[13,66:66=')',<')'>,7:18]
[14,67:67=';',<'>',7:19]
[15,69:69='}',<'>',8:0]
[16,71:73='fin',<'fin'>,9:0]
[17,75:74='<EOF>',<EOF>,10:0]
```

Archivo: programa2.BabyDuck



Prueba 3

Archivo: incorrect.BabyDuck



```
[cocaine_ninja@nixos:~/Documents/Tec/ITC/8vo/TC3002B-Compiladores/Compiladores/Tareas/BabyDuck]$ grun BabyDuck programa -tokens < incorrect.BabyDuck
[@0,0:7='programa',<'programa'>,1:0]
[@1,9:18='Incorrecto',<ID>,1:9]
[@2,19:19=';',<'>',1:19]
[@3,21:24='vars',<'vars'>,2:0]
[@4,26:26='x',<ID>,2:5]
[@5,27:27=':',<':'>,2:6]
[@6,29:36='flotante',<'flotante'>,2:8]
[@7,37:37=';',<'>',2:16]
[@8,40:45='inicio',<'inicio'>,4:0]
[@9,47:47='{',<'{'>,5:0]
[@10,51:51='x',<ID>,6:2]
[@11,53:53='=',<'='>,6:4]
[@12,55:58='10.5',<CTE_FLOAT>,6:6]
[@13,62:68='escribe',<'escribe'>,7:2]
[@14,69:69='(',<'('>,7:9]
[@15,70:129='Hay un error, falta un ; en la asignacion de la variable x"',<LETRERO>,7:10]
[@16,130:130=')',<')'>,7:70]
[@17,131:131=';',<'>',7:71]
[@18,133:133='}',<'>',8:0]
[@19,135:137='fin',<'fin'>,9:0]
[@20,139:138='<EOF>',<EOF>,10:0]
line 7:2 missing ';' at 'escribe'
```

Hallazgos

Al inicio elegí Go como el lenguaje para implementar el proyecto BabyDuck con el módulo **participle**, pero decidí cambiar a Python con **ANTLR** debido a una mala estimación de la complejidad de la síntesis (backend). Si bien el esfuerzo de programación manual es menor con **participle** que con **ANTLR**, reflexioné y me di cuenta que me tomaría más tiempo implementar el proyecto porque estaría atrapado más tiempo resolviendo bugs del lenguaje que en entender dónde se encuentra el cuello de botella del compilador, es un lenguaje que apenas empecé a aprender recientemente y no tengo un dominio de su ecosistema como administrador de paquetes, de versiones del lenguaje, uso de herramientas, etc. Elegí **ANTLR** porque es una herramienta muy sólida con bastante documentación, soporta múltiples lenguajes como C#, C++, JavaScript, Java y Python.

Mientras estaba implementando las reglas de gramática en ANTLR encontré ciertos tokens que no había implementado y unos cambios ligeros en el lenguaje que generan ciertas gramáticas, eso me hizo reescribir el lenguaje que generan mis reglas gramaticales a mano antes de plasmarlas en ANTLR. Durante este proceso identifiqué una segunda ambigüedad que no se cubrió en clases:

- <ESTATUTO> : se identificó la primera ambigüedad en clases donde se usa un id en las reglas internas <ASIGNA> y <LLAMADA>, la primera es para asignación de datos a variables y la segunda para invocar la ejecución de funciones, ambas usan un id pero con distintos casos. Se resolvió con el Factor Común Izquierdo indicando explícitamente que después de recibir un id pueden venir distintos casos.
- <FACTOR> : se agregó un camino adicional al diagrama de sintaxis ya existente, dicho camino consiste únicamente en la regla <LLAMADA> que hace uso del token id, la cuestión es que otro camino de la regla actual requiere uso específico de un id. Para solucionar esto también se utilizó el Factor Común Izquierdo porque el token que sigue de un id es diferente en ambos casos.

Esto provocó que el lenguaje de las reglas <ASIGNA> y <LLAMADA> tuviera que implementarse en el Factor Común Izquierdo, y como consecuencia ninguna otra regla llamaba a estas 2 reglas, quedaron huérfanas. Se decidió eliminarlas de la gramática porque podrían confundir al lector, pero se mantuvieron en comentarios del código indicando las razones de su eliminación y que inicialmente eran reglas del diagrama de sintaxis.