

Syntactic Analyzer for the C Minus Language

Victor Emmanuel Guerra Aguado - A01568075

Summary	1
Context & Notation	1
Analysis	7
Requirements	7
The Definition of the final Grammar	8
Making the grammar fit its description	9
Elimination of Left Recursion	11
Elimination of Left Factors	13
Simplification of the Grammar	16
An unexpected visitor	22
FIRST, FOLLOW, & FIRST+ sets & Parsing Table	24
Design	39
Implementation	45
Testing	49
References	56

Summary

The following document will report upon the development of a syntactic analyzer for the proposed C Minus programming language in the Python programming language, including in the complete phases of **Analysis**, **Design**, **Implementation**, and **Testing**. Each phase will contain its appropriate documentation and implementation.

Context & Notation

Context Free Grammars (CFGs)

A grammar is a set of rules and examples that deals with the syntax and word structure of a language via a set of Productions. In the context of compiler design and programming languages, it provides the rules of structure for a given programming language.

A Grammar is said to be Context-Free if all its productions are of the form $\alpha \rightarrow \beta$ where alpha consists of a single variable (non-terminal symbol) and beta is a string of language symbols (terminal symbols or tokens) and non-terminals. [1]

The notation used to specify grammars is known as the Backus Normal Form. A grammar is defined by a 4-tuple set (V, T, P, S) where V is the set of non-terminal symbols, T is the set of terminal symbols, P is the set of productions, and S is the start symbol. The set of terminal symbols is the list of Tokens produced by the lexical analyzer and are represented by:

- Lower case letters: a, b, c, etc.
- Operator Symbols: +, -, =, etc.
- Punctuation Symbols: ,, ., (, {, etc.
- Digits: 0 ... 9
- Keywords: int, main, if, etc.

The set of non-terminals, sometimes called syntactic variables, represents a set of strings; the set of strings defines the language generated by the grammar. Non-terminal symbols are represented by:

- Upper case letters: A, B, C, etc.
- Letter S is the start symbol.
- Italic lower case string: *exp*, *statement*, etc.

One of the non-terminal symbols is designated as the start symbol, and the set of strings it specifies is the language generated by the grammar.

The set of productions of a grammar specifies the manner in which terminals and non-terminals can be combined to form strings. Each productions consists of [1]:

- A non-terminal known as the head or left-side of the production
- The arrow symbol :“ \rightarrow ”
- A body or right-side, which consists of zero or more terminals and non-terminals

Issues with CFG(Ambiguity, Left Recursion, Left Factorization, Simplification of Grammars)

There are some aspects of a CFG that must be dealt with in order to be properly used for a programming language and the corresponding syntax analysis phase of its compiler:

- Ambiguity
Parse Trees uniquely express the structure of the syntax as leftmost or rightmost derivations. Unfortunately it's possible for a grammar to permit a string to have more

than one parse tree, if it does the grammar is called ambiguous. This represents a serious problem for a parser since the grammar does not specify precisely the syntactic structure of a program.

There is NO algorithm that converts an ambiguous grammar to an equivalent unambiguous one in a straightforward manner. The basic method for dealing with ambiguity is to change the grammar into a form that forces the construction of the correct parse tree, thus removing the ambiguity.

- Left Recursion

A grammar is left recursive if it has a non-terminal A such that there is a derivation $A \Rightarrow A\alpha$ for some string α . If this kind of grammar is used in Top-Down parsers, the parser may go into an infinite loop. Top-Down parsers cannot handle left recursive grammars. Therefore, all left recursive grammars must be converted into an equivalent non-left recursive grammar. Left Recursion may appear in a single step of the derivation (immediate left recursion), or it may appear in more than one step of the derivation.

A grammar G is left-recursive if there is a rule of the form: $A \rightarrow A\alpha \mid \beta$ where β does not start with A , i.e., the leftmost terminal on the right hand side is the same as the non-terminal on the left hand side. To eliminate immediate left recursion rewrite the grammar as:

$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' \mid \varepsilon \end{aligned}$$

This is an equivalent grammar G' , which is free of left recursion. In General, for all productions of the form:

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

where $\beta_1, \beta_2, \dots, \beta_n$ do not start with A , Eliminate immediate left recursion by:

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A' \\ A' &\rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \varepsilon \end{aligned}$$

A grammar can be free of immediate left recursion, but it still can be left recursive. Eliminate complete left recursion:

1) Arrange non-terminals in some order: $X_1 \dots X_n$

2) Apply the following procedure:

for $i=1$ to n do {
 for $j=1$ to $i-1$ do {
 replace each production $X_i \rightarrow X_j \beta$ by

$$X_i \rightarrow \alpha_1 \beta \mid \alpha_2 \beta \mid \dots \mid \alpha_k \beta$$

where

$$X_j \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$$

}

eliminate immediate left recursion among X_i .

}

- Left Factoring

Sometimes, the grammar may have common a prefix in many productions like $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n$ where α is a common prefix. While processing α it can not be decided whether to expand A by $\alpha\beta_1$ or by $\alpha\beta_n$. So this needs Backtracking. Eliminating Left Factoring is a grammar transformation technique that is useful to produce grammars that are suitable for predictive or top-down parsing.

To Eliminate left factoring:

For each non-terminal A with two or more alternatives (production rules) with a common non-empty prefix:

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_m$$

where $\gamma_1, \gamma_2, \dots, \gamma_n$ do not begin with α convert it into:

$$A \rightarrow \alpha A' \mid \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_m$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

- Simplification of grammars

All Grammars are not always optimized. A grammar may contain extra or unnecessary symbols, which will increase the length of the grammar. Simplification of the grammar generally includes the following steps:

- 1) Elimination of useless symbols.
- 2) Elimination of ϵ -productions.
- 3) Elimination of unit productions.

1) A symbol is useless if it cannot derive a terminal or if it is not reachable from the start symbol. To eliminate such a symbol: eliminate all productions in which it appears on either side.

2) If a CFL contains the word ϵ , then the CFG must have an ϵ -production. However, if a CFG has an ϵ -production, then the CFL does not necessarily contain the word ϵ .

In a given CFG, a non-terminal X is nullable if:

1. There is a production $X \rightarrow \epsilon$, or
2. There is a derivation that starts at X and leads to ϵ , i.e., $X \Rightarrow \epsilon$

To eliminate ϵ -Productions:

1. Construct the set of all Nullable Variables V_n
 2. For each production $A \rightarrow B$, if B is a Nullable Variable, replace it by ϵ and add it with all possible combinations on the RHS.
 3. Take out the production $B \rightarrow \epsilon$.
- 3) A Unit Production is a production of the form $A \rightarrow B$. To eliminate Unit Productions:

1. For each pair of non-terminals A and B such that:
 - a. There is a production $A \rightarrow B$ and
 - b. The productions of B are:
 $B \rightarrow s_1 \mid s_2 \mid s_3 \mid \dots \mid s_n$, where at least one $s_i \in (T + V)^*$, i.e., s_i is a string of terminals and non-terminals.
2. Create the new productions:
 - a. $A \rightarrow s_1 \mid s_2 \mid s_3 \mid \dots \mid s_n$
3. Remove the unit production $A \rightarrow B$.
4. Repeat for all pairs A and B .

Parsing Algorithms & Top-Down Parsers

Parsing algorithms can be classified into three types: Universal Parsers, which perform the syntax analysis with any grammar, however, this kind of parser is inefficient and therefore it is not used in commercial compilers; Top-down Parsers (TDP), which begin from the start symbol of a grammar, and applies production rules for each non-terminal until it gets a string formed only by terminals; and Bottom-Up Parsers (BUP), in general more powerful than Top-Down parsers, but extremely difficult to write and develop by hand, and are therefore only used in most of the automatic parser generators; bottom-up parsers have none of the restrictions of a top-down parser previously discussed with the exception of an unambiguous grammar; Bottom-Up parsing constructs a parse tree from an input string beginning at the leaves and working towards the root, or start symbol of the grammar.[2]

Top-Down Parsers can be subdivided into Brute Force Parsers and Predictive Parsers. The Brute Force Technique uses full Backtracking, meaning, whenever a non-terminal has more than one alternative, it always expands the first option the first time it encounters it, if the parser does not successfully finish the analysis, the parser backtracks all the way to the first execution of the first alternative, and now tries with the second alternative, the parser repeats this until there is a match with the input string, or all combinations are verified.

Predictive Parsers are capable of deciding the correct alternative in order to expand a non-terminal; the prediction relies on information about what first symbols can be generated from a production or rule: if the first symbol of a production is a non-terminal, then the non-terminal has to be expanded until a set of terminal symbols are obtained; the grammar upon which a predictive parser is designed has three restrictions:

1. It has to be unambiguous.
2. It has to be free of left recursion.
3. It has to be free of left factors.

Predictive Parsers can be divided into two types:

1. Recursive Descent Parser
2. Non-recursive Descent (or LL(1)) Parser

Recursive Descent Parsers are the simplest type of predictive parsers. A Predictive Recursive Descent parser is constructed by writing recursive procedures for each non-terminal; it is a Top-Down process in which the parser attempts to verify whether the syntax of the input string is correct as it is read from left to right, and it always expands a non-terminal by its right hand side of the rule until it gets a string of terminals. The parser reads tokens from the scanner and matches them with terminals from the grammar that describes the syntax of the input, token by token from left to right.

A Non-Recursive Predictive parser can be built by explicitly maintaining a stack instead of implicitly using recursive calls, and a parsing table to dictate the parsing decisions. The LL(1) parser consists of the following elements:

1. The input buffer that contains the string to be parsed, i.e. The scanners' output
2. The stack that contains a sequence of grammar symbols
3. The parsing table $M[X, a]$ that dictates the parsing decisions
4. The parsing program that takes the symbol X at the top of the stack together with "a" the current input symbol from the scanner. If X is a non-terminal then an X -production is chosen according to $M[X, a]$. If X is a token then it checks for a match between X and a .

The Parsing Algorithm Used

For the purpose of the C Minus language compiler, the LL(1) parsing algorithm was selected. In order to be able to implement the algorithm, the following steps will be necessary:

1. Elimination of Ambiguity.
2. Elimination of Left Recursion.
3. Elimination of Left Factors.
4. Elimination of Useless Symbols.
5. Elimination of ϵ -productions.
6. Elimination of Unit Productions.
7. Calculation of the FIRST Set.
8. Calculation of the FOLLOW Set.
9. Calculation of the FIRST+ Set.
10. Construction of the Parsing Table.
11. Implementation of the LL(1) Stack.
12. Implementation of the LL(1) parsing program.

Programming Language

The chosen programming language for the compiler, lexical analyzer, and syntactic analyzer is Python. It was chosen for its simplicity and straightforwardness regarding the management of data structures, memory, and data types; its flexibility would make it possible to take some load from the implementation phase that could be used during other phases.

Analysis

Requirements

1. The Parser must correctly recognize the C minus grammar.
2. The Parser must recognize and report the following Syntactical Errors:
 - a. The program does not contain at least one declaration.
 - b. void type used for a variable.
 - c. Expected a (
 - d. Expected a {
 - e. Expected a [
 - f. Expected a ;
 - g. Expected a)
 - h. Expected a }

- i. Expected a]
 - j. Expected an identifier
 - k. Expected a number constant
 - l. Expected an expression
3. The Parser must recognize and report the following Semantic Errors:
 - a. The programs' last declaration is not main.
 - b. A variable has been used before declaration.
 - c. A function has been used before declaration.
4. The Parser must modify the existing Symbol Tables to add the following columns:
 - a. isVar (is a variable)
 - b. isFun (is a function)
 - c. dataType (data type that it is associated to)
 - d. noArgs (if it is a function, how many arguments does it receive)
 - e. isGlobal (if it is a global declaration)
 - f. isLocal (if it is a local declaration)
5. The Error Reports must provide the number of the line in which the error was encountered.

The Definition of the final Grammar

The Parser will be made to follow the Top-Down LL(1) [2] algorithm, which requires said grammar to comply with several requirements before its implementation into code such as no Left Recursion, no Left Factors, and is unambiguous, while the simplification of the grammar is not necessary for a functional Top-Down parser, some simplification will be done in order to make the process of calculating the FIRST, FOLLOW, and FIRST+ Sets less arduous.

Firstly, the provided grammar for the C minus language contains several inconsistencies with the provided semantics of the language; it must be reworked to make it fit its

description:

1. $program \rightarrow declaration_list$
2. $declaration_list \rightarrow declaration_list \ declaration \mid declaration$
3. $declaration \rightarrow var_declaration \mid fun_declaration$
4. $var_declaration \rightarrow type_specifier \ ID ; \mid type_specifier \ ID [\ NUM] ;$
5. $type_specifier \rightarrow \mathbf{int} \mid \mathbf{void}$
6. $fun_declaration \rightarrow type_specifier \ ID (\ params) \ compound_stmt$
7. $params \rightarrow param_list \mid \mathbf{void}$
8. $param_list \rightarrow param_list , \ param \mid param$
9. $param \rightarrow type_specifier \ ID \mid type_specifier \ ID [\]$
10. $compound_stmt \rightarrow \{ \ local_declarations \ statement_list \}$
11. $local_declarations \rightarrow local_declarations \ var_declaration \mid \epsilon$
12. $statement_list \rightarrow statement_list \ statement \mid \epsilon$
13. $statement \rightarrow assignment_stmt \mid call_stmt \mid compound_stmt \mid selection_stmt$
 $\mid iteration_stmt \mid return_stmt \mid input_stmt \mid output_stmt$
14. $assignment_stmt \rightarrow var = expression ;$
15. $call_stmt \rightarrow call ;$
16. $selection_stmt \rightarrow \mathbf{if} (expression) \ statement$
 $\mid \mathbf{if} (expression) \ statement \ \mathbf{else} \ statement$
17. $iteration_stmt \rightarrow \mathbf{while} (expression) \ statement$
18. $return_stmt \rightarrow \mathbf{return} ; \mid \mathbf{return} \ expression ;$
19. $input_stmt \rightarrow \mathbf{input} \ var ;$
20. $output_stmt \rightarrow \mathbf{output} \ expression ;$
21. $var \rightarrow \mathbf{ID} \mid \mathbf{ID} [\ arithmetic_expression \]$
22. $expression \rightarrow arithmetic_expression \ relop \ arithmetic_expression$
 $\mid arithmetic_expression$
23. $relop \rightarrow <= \mid < \mid > \mid >= \mid == \mid !=$
24. $arithmetic_expression \rightarrow arithmetic_expression \ addop \ term \mid term$
25. $addop \rightarrow + \mid -$
26. $term \rightarrow term \ mulop \ factor \mid factor$
27. $mulop \rightarrow * \mid /$
28. $factor \rightarrow (\ arithmetic_expression \) \mid var \mid call \mid \mathbf{NUM}$
29. $call \rightarrow \mathbf{ID} (\ args)$
30. $args \rightarrow args_list \mid \epsilon$
31. $args_list \rightarrow args_list , \ arithmetic_expression \mid arithmetic_expression$

Figure 1. Original Provided Grammar

Making the grammar fit its description

- 1) *The grammar cannot identify whether an identifier for a function or variable is being used before its declaration.* This problem cannot be solved by a context free grammar, and will be solved later in the coding step, by making sure the first time an identifier is found, it is during its declaration.

- 2) *The grammar does not specify that there is always a final function declaration with the identifier 'main'. This problem can be partially solved by changing the grammar to have at least one declaration; later, during the coding of the parser, code can be added to make sure that said declaration uses the identifier 'main' and is, in fact, a function.*
- 3) *The grammar allows for a variable to be assigned a **void** type. This issue can be solved by creating separate type specifier variables for variables and functions.*
- 4) *The grammar allows for parameters of type **void** to be assigned to a function. This can be solved by using only the type specifier variable for variables in the productions of the non-terminal *param*.*

Any other remaining specification cannot be verified during the syntax analysis, and correspond to further steps in the compilation process. The resulting grammar after the previous modifications is as follows, the corrections made in the previous step highlighted in green:

Produccion				
No. P	LHS	→	RHS	RHS
1	program	→	declaration_list void ID (void) compound_stmt	void ID (void) compound_stmt
2	declaration_list	→	declaration_list declaration	epsilon
3	declaration	→	var_declaration	fun_declaration
4	var_declaration	→	var_type_specifier ID ;	var_type_specifier ID [NUM] ;
5	func_type_specifier	→	int	void
6	var_type_specifier	→	int	
7	fun_declaration	→	func_type_specifier ID (params) compound_stmt	
8	params	→	param_list	void
9	param_list	→	param_list , param	param
10	param	→	var_type_specifier ID	var_type_specifier ID []
11	compound_stmt	→	{ local_declarations statement_list }	
12	local_declarations	→	local_declarations var_declaration	epsilon
13	statement_list	→	statement_list statement	epsilon
14	statement	→	assignment_stmt selection_stmt input_stmt	call_stmt iteration_stmt output_stmt
15	assignment_stmt	→	var = expression ;	
16	call_stmt	→	call ;	
17	selection_stmt	→	if (expression) statement	if (expression) statement else statement
18	iteration_stmt	→	while (expression) statement	
19	return_stmt	→	return ;	return expression ;
20	input_stmt	→	input var ;	
21	output_stmt	→	output expression ;	
22	var	→	ID	ID [arithmetic_expression]
23	expression	→	arithmetic_expression relop arithmetic_expression	arithmetic_expression
24	relop	→	<=	<
		→	>=	>
		→		!=
25	arithmetic_expression	→	arithmetic_expression addop term	term
26	addop	→	+	-
27	term	→	term mulop factor	factor
28	mulop	→	*	/
29	factor	→	(arithmetic_expression) NUM	var call
30	call	→	ID (args)	
31	args	→	args_list	epsilon
32	args_list	→	args_list , arithmetic_expression	arithmetic_expression

Figure 2. The grammar after corrections for semantic compliance.

Elimination of Left Recursion

Now that the grammar appropriately represents the language described by the semantics, it must be refactored to, firstly, eliminate the left-recursion present, highlighted in red in Figure 2:

- 1) In the productions of non-terminal **No. 2** 'declaration_list':

$$declaration_list \rightarrow declaration_list declaration$$

So it is replaced by the productions:

$$declaration_list \rightarrow declaration declaration_list'$$

$$declaration_list' \rightarrow declaration declaration_list' \mid \epsilon$$

- 2) In the productions of non-terminal **No. 9** 'param_list':

$$param_list \rightarrow param_list, param \mid param$$

So it is replaced by the productions:

$$param_list \rightarrow param param_list'$$

$$param_list \rightarrow , param param_list' | \varepsilon$$

- 3) In the productions of non-terminal **No. 12** '*local_declarations*':

$$local_declarations \rightarrow local_declarations var_declaration | \varepsilon$$

So they are replaced by the productions:

$$local_declarations \rightarrow local_declarations'$$

$$local_declarations' \rightarrow var_declarations local_declarations' | \varepsilon$$

- 4) In the productions of non-terminal **No. 13** '*statement_list*':

$$statement_list \rightarrow statement_list statement | \varepsilon$$

So they are replaced by the productions:

$$statement_list \rightarrow statement_list'$$

$$statement_list' \rightarrow statement statement_list' | \varepsilon$$

- 5) In the productions of non-terminal **No. 25** '*arithmetic_expression*':

$$arithmetic_expression \rightarrow arithmetic_expression addop term | term$$

So they are replaced by the productions:

$$arithmetic_expression \rightarrow term arithmetic_expression'$$

$$arithmetic_expression' \rightarrow addop term arithmetic_expression' | \varepsilon$$

- 6) In the productions of non-terminal **No. 27** '*term*':

$$term \rightarrow term mulop factor | factor$$

So they are replaced by the productions:

$$term \rightarrow factor term'$$

$$term' \rightarrow mulop factor term | \varepsilon$$

- 7) In the productions of non-terminal **No. 32** '*args_list*':

$$args_list \rightarrow args_list , arithmetic_expression | arithmetic_expression$$

So they are replaced by the productions:

$$args_list \rightarrow arithmetic_expression args_list'$$

$$args_list' \rightarrow , arithmetic_expression args_list' | \varepsilon$$

With the previous modifications the grammar results as follows:

Produccion				
No. P	LHS	→	RHS	RHS
1	program	→	declaration_list void ID (void) compound_stmt	void ID (void) compound_stmt
2	declaration_list	→	declaration declaration_list	epsilon
4	declaration	→	var_declaration	fun_declaration
5	var_declaration	→	var_type_specifier ID ;	var_type_specifier ID [NUM] ;
6	func_type_specifier	→	int	void
7	var_type_specifier	→	int	
8	fun_declaration	→	func_type_specifier ID (params) compound_stmt	
9	params	→	param_list	void
10	param_list	→	param param_list'	
11	param_list'	→	, param param_list'	epsilon
12	param	→	var_type_specifier ID	var_type_specifier ID []
13	compound_stmt	→	{ local_declarations statement_list }	
14	local_declarations	→	local_declarations'	
15	local_declarations'	→	var_declaration local_declarations'	epsilon
16	statement_list	→	statement_list'	
17	statement_list'	→	statement statement_list'	epsilon
18	statement	→	assignment_stmt	call_stmt
			selection_stmt	iteration_stmt
			input_stmt	output_stmt
19	assignment_stmt	→	var = expression ;	
20	call_stmt	→	call ;	
21	selection_stmt	→	if (expression) statement	if (expression) statement else statement
22	iteration_stmt	→	while (expression) statement	
23	return_stmt	→	return ;	return expression ;
24	input_stmt	→	input var ;	
25	output_stmt	→	output expression ;	
26	var	→	ID	ID [arithmetic_expression]
27	expression	→	arithmetic_expression relop arithmetic_expression	arithmetic_expression
28	relop	→	<=	<
			>=	==
29	arithmetic_expression	→	term arithmetic_expression'	
30	arithmetic_expression'	→	addop term arithmetic_expression'	epsilon
31	addop	→	+	-
32	term	→	factor term'	
33	term'	→	mulop factor term'	epsilon
34	mulop	→	*	/
35	factor	→	(arithmetic_expression)	var
			NUM	call
36	call	→	ID (args)	
37	args	→	args_list	epsilon
38	args_list	→	arithmetic_expression args_list'	
39	args_list'	→	, arithmetic_expression args_list'	epsilon

Figure 3. The grammar after rework for elimination of left recursion.

Elimination of Left Factors

Next, the grammar must be reworked in order to eliminate the left-factorization. Left factorization issues detected are highlighted in red in **Figure 3**:

- 1) In the productions of non-terminal **No. 5** 'var_declaration':

$$\text{var_declaration} \rightarrow \text{var_type_specifier } \mathbf{ID} ; \mid \text{var_type_specifier } \mathbf{ID} [\text{NUM}] ;$$

So it is replaced by the productions:

$$\text{var_declaration} \rightarrow \text{var_type_specifier } \mathbf{ID} \text{ var_declaration'}$$

$$var_declaration' \rightarrow ; | [NUM] ;$$

- 2) In the productions of non-terminal **No. 12** '*param*':

$$param \rightarrow var_type_specifier \textbf{ID} | var_type_specifier \textbf{ID} []$$

So it is replaced by the productions:

$$param \rightarrow var_type_specifier \textbf{ID} param'$$

$$param_list \rightarrow [] | \epsilon$$

- 3) In the productions of non-terminal **No. 21** '*selection_stmt*':

$$selection_stmt \rightarrow \textbf{if} (expression) statement | \textbf{if} (expression) statement \textbf{else} statement$$

So they are replaced by the productions:

$$selection_stmt \rightarrow \textbf{if} (expression) statement selection_stmt'$$

$$selection_stmt' \rightarrow \textbf{else} statement | \epsilon$$

- 4) In the productions of non-terminal **No. 23** '*return_stmt*':

$$return_stmt \rightarrow \textbf{return} ; | \textbf{return} expression ;$$

So they are replaced by the productions:

$$return_stmt \rightarrow \textbf{return} return_stmt'$$

$$return_stmt' \rightarrow ; | expression ;$$

- 5) In the productions of non-terminal **No. 26** '*var*':

$$var \rightarrow \textbf{ID} | \textbf{ID} [arithmetic_expression]$$

So they are replaced by the productions:

$$var \rightarrow \textbf{ID} var'$$

$$var' \rightarrow [arithmetic_expression] | \epsilon$$

- 6) In the productions of non-terminal **No. 27** '*expression*':

$$expression \rightarrow arithmetic_expression \textit{relop} arithmetic_expression | arithmetic_expression$$

So they are replaced by the productions:

$$expression \rightarrow arithmetic_expression expression'$$

$$expression' \rightarrow \textit{relop} arithmetic_expression | \epsilon$$

With all previous corrections, the grammar results as follows:

Produccion				
No.	LHS	→	RHS	RHS
1	program	→	declaration_list void ID (void) compound_stmt	void ID (void) compound_stmt
2	declaration_list	→	declaration declaration_list	epsilon
3	declaration	→	var_declaration	fun_declaration
4	var_declaration	→	var_type_specifier ID var_declaration'	
5	var_declaration'	→	;	[NUM] ;
6	func_type_specifier	→	int	void
7	var_type_specifier	→	int	
8	fun_declaration	→	func_type_specifier ID (params) compound_stmt	
9	params	→	param_list	void
10	param_list	→	param param_list'	
11	param_list'	→	, param param_list'	epsilon
12	param	→	var_type_specifier ID param'	
13	param'	→	[]	epsilon
14	compound_stmt	→	{ local_declarations statement_list }	
15	local_declarations	→	local_declarations'	
16	local_declarations'	→	var_declaration local_declarations'	epsilon
17	statement_list	→	statement_list'	
18	statement_list'	→	statement statement_list'	epsilon
19	statement	→	assignment_stmt selection_stmt input_stmt	call_stmt iteration_stmt output_stmt
20	assignment_stmt	→	var = expression ;	
21	call_stmt	→	call ;	
22	selection_stmt	→	if (expression) statement selection_stmt'	
23	selection_stmt'	→	else statement	epsilon
24	iteration_stmt	→	while (expression) statement	
25	return_stmt	→	return return_stmt'	
26	return_stmt'	→	;	expression ;
27	input_stmt	→	input var ;	
28	output_stmt	→	output expression ;	
29	var	→	ID var'	
30	var'	→	[arithmetic_expression]	epsilon
31	expression	→	arithmetic_expression expression'	
32	expression'	→	relop arithmetic_expression	epsilon
33	relop	→	<= >=	< == !=
34	arithmetic_expression	→	term arithmetic_expression'	
35	arithmetic_expression'	→	addop term arithmetic_expression'	epsilon
36	addop	→	+	-
37	term	→	factor term'	
38	term'	→	mulop factor term'	epsilon
39	mulop	→	*	/
40	factor	→	(arithmetic_expression) NUM	var call
41	call	→	ID (args)	
42	args	→	args_list	epsilon
43	args_list	→	arithmetic_expression args_list'	
44	args_list'	→	, arithmetic_expression args_list'	epsilon

Figure 4. The grammar after elimination of left factors.

Simplification of the Grammar

Now the grammar is appropriate for the LL(1) algorithm, however, some additional simplifications have been made in order to better optimize the grammar, such as removal of useless symbols, removal of unit productions, and elimination of as many ϵ -productions as possible. In **Figure 4**, all ϵ -Productions are highlighted in red, however, most of these cannot be eliminated, since by doing this a left factor will be re-introduced in such a way that, by re-eliminating the left factor will re-introduce an ϵ -production. All such ϵ -productions are highlighted in blue in **Figure 5**:

Produccion				
No. F	LHS	→	RHS	RHS
1	program	→	declaration_list void ID (void) compound_stmt	void ID (void) compound_stmt
2	declaration_list	→	declaration declaration_list	epsilon
4	declaration	→	var_declaration	fun_declaration
5	var_declaration	→	var_type_specifier ID var_declaration'	
6	var_declaration'	→	;	[NUM] ;
7	func_type_specifier	→	int	void
8	var_type_specifier	→	int	
9	fun_declaration	→	func_type_specifier ID (params) compound_stmt	
10	params	→	param_list	void
11	param_list	→	param param_list'	
12	param_list'	→	, param param_list'	epsilon
13	param	→	var_type_specifier ID param'	
14	param'	→	[]	epsilon
15	compound_stmt	→	{ local_declarations statement_list }	
16	local_declarations	→	local_declarations'	
17	local_declarations'	→	var_declaration local_declarations'	epsilon
18	statement_list	→	statement_list'	
19	statement_list'	→	statement statement_list'	epsilon
20	statement	→	assignment_stmt selection_stmt input_stmt	call_stmt iteration_stmt output_stmt
21	assignment_stmt	→	var = expression ;	
22	call_stmt	→	call ;	
23	selection_stmt	→	if (expression) statement selection_stmt'	
24	selection_stmt'	→	else statement	epsilon
25	iteration_stmt	→	while (expression) statement	
26	return_stmt	→	return return_stmt'	
27	return_stmt'	→	;	expression ;
28	input_stmt	→	input var ;	
29	output_stmt	→	output expression ;	
30	var	→	ID var'	
31	var'	→	[arithmetic_expression]	epsilon
32	expression	→	arithmetic_expression expression'	
33	expression'	→	relop arithmetic_expression	epsilon
34	relop	→	<= >=	< > == !=
35	arithmetic_expression	→	term arithmetic_expression'	
36	arithmetic_expression'	→	addop term arithmetic_expression'	epsilon
37	addop	→	+	-
38	term	→	factor term'	
39	term'	→	mulop factor term'	epsilon
40	mulop	→	*	/
41	factor	→	(arithmetic_expression) NUM	var call
42	call	→	ID (args)	ID ()
43	args	→	args_list	
44	args_list	→	arithmetic_expression args_list'	
45	args_list'	→	, arithmetic_expression args_list'	epsilon

Figure 5. All ϵ -productions in the grammar.

However there is one ϵ -production that can be eliminated, seen as productions of non-terminal No. 42 '*args*' in **Figure 4**, and solved in **Figure 5** with the following method:

- 1) $call \rightarrow \mathbf{ID} (args)$ can also be $call \rightarrow \mathbf{ID} ()$, so by adding this production to '*call*' the production $args \rightarrow \epsilon$ can be eliminated.

This does re-introduce a left factor in non-terminal, but eliminating this does not produce a new ϵ -production, thus avoiding an infinite loop of elimination.

- 2) In the productions of non-terminal **No. 42** '*call*':

$$call \rightarrow \mathbf{ID} (args) \mid \mathbf{ID} ()$$

So they are replaced by the productions:

$$call \rightarrow \mathbf{ID} (call'$$

$$call' \rightarrow) \mid args)$$

These changes result in the following grammar:

Produccion				
No. Pr	LHS	→	RHS	RHS
1	program	→	declaration_list void ID (void) compound_stmt	void ID (void) compound_stmt
2	declaration_list	→	declaration declaration_list	epsilon
3	declaration	→	var_declaration	fun_declaration
4	var_declaration	→	var_type_specifier ID var_declaration'	
5	var_declaration'	→	;	[NUM] ;
6	func_type_specifier	→	int	void
7	var_type_specifier	→	int	
8	fun_declaration	→	func_type_specifier ID (params) compound_stmt	
9	params	→	param_list	void
11	param_list	→	param param_list'	
12	param_list'	→	, param param_list'	epsilon
13	param	→	var_type_specifier ID param'	
14	param'	→	[]	epsilon
15	compound_stmt	→	{ local_declarations statement_list }	
16	local_declarations	→	local_declarations'	
17	local_declarations'	→	var_declaration local_declarations'	epsilon
18	statement_list	→	statement_list'	
19	statement_list'	→	statement statement_list'	epsilon
20	statement	→	assignment_stmt selection_stmt input_stmt	call_stmt iteration_stmt output_stmt
21	assignment_stmt	→	var = expression ;	
22	call_stmt	→	call ;	
23	selection_stmt	→	if (expression) statement selection_stmt'	
24	selection_stmt'	→	else statement	epsilon
25	iteration_stmt	→	while (expression) statement	
26	return_stmt	→	return return_stmt'	
27	return_stmt'	→	;	expression ;
28	input_stmt	→	input var ;	
29	output_stmt	→	output expression ;	
30	var	→	ID var'	
31	var'	→	[arithmetic_expression]	epsilon
32	expression	→	arithmetic_expression expression'	
33	expression'	→	relop arithmetic_expression	epsilon
34	relop	→	<= >=	< == > !=
35	arithmetic_expression	→	term arithmetic_expression'	
36	arithmetic_expression'	→	addop term arithmetic_expression'	epsilon
37	addop	→	+	-
38	term	→	factor term'	
39	term'	→	mulop factor term'	epsilon
40	mulop	→	*	/
41	factor	→	(arithmetic_expression) NUM	var call
42	call	→	ID (call'	
43	call'	→	args))
44	args	→	args_list	
45	args_list	→	arithmetic_expression args_list'	
46	args_list'	→	, arithmetic_expression args_list'	epsilon

Figure 6. The grammar after the correction of the ϵ -production at non-terminal 42.

Next, the elimination of unit productions. In **Figure 6**, all unit productions are highlighted in red; most of them will be eliminated, however, some will remain, since the variable is also used in other productions and their elimination might make the grammar more difficult to understand:

- 1) The production '*declaration* \rightarrow *fun_declaration*' is replaced with the production '*declaration* \rightarrow *func_type_specifier* **ID** (*params*) *compound_stmt*'. And, in this case the non-terminal '*fun_declaration*' becomes a useless symbol and is eliminated from the grammar.
- 2) The production '*params* \rightarrow *param_list*' is replaced with the production '*params* \rightarrow *param param_list*'. And, in this case the non-terminal '*param_list*' becomes a useless symbol and is eliminated from the grammar.
- 3) The production '*local_declarations* \rightarrow *local_declarations*' is replaced with the production '*local_declarations* \rightarrow *var_declaration local_declaration*' AND '*local_declarations* \rightarrow ϵ '. To avoid redundant productions the calls to '*local_declarations*' are changed into '*local_declarations*', the symbol '*local_declarations*' becomes useless, and is eliminated.
- 4) The production '*statement_list* \rightarrow *statement_list*' is replaced with the production '*statement_list* \rightarrow *statement statement_list*' AND '*statement_list* \rightarrow ϵ '. To avoid redundant productions the calls to '*statement_list*' are changed into '*statement_list*', the symbol '*statement_list*' becomes useless, and is eliminated.
- 5) The production '*statement* \rightarrow *assignment_stmt*' is replaced with the production '*statement* \rightarrow *var* = *expression* ;'. And, in this case, the non-terminal '*assignment_stmt*' becomes a useless symbol and is eliminated from the grammar.
- 6) The production '*statement* \rightarrow *call_stmt*' is replaced with the production '*statement* \rightarrow *call* ;'. And, in this case, the non-terminal '*call_stmt*' becomes a useless symbol and is eliminated from the grammar.
- 7) The production '*statement* \rightarrow *iteration_stmt*' is replaced with the production '*statement* \rightarrow **while** (*expression*) *statement*'. And, in this case, the non-terminal '*iteration_stmt*' becomes a useless symbol and is eliminated from the grammar.
- 8) The production '*statement* \rightarrow *return_stmt*' is replaced with the production '*statement* \rightarrow **return** *return_stmt*'. And, in this case, the non-terminal '*return_stmt*' becomes a useless symbol and is eliminated from the grammar.

- 9) The production ' $statement \rightarrow input_stmt$ ' is replaced with the production ' $statement \rightarrow \textbf{input } var ;$ '. And, in this case, the non-terminal ' $input_stmt$ ' becomes a useless symbol and is eliminated from the grammar.
- 10) The production ' $statement \rightarrow output_stmt$ ' is replaced with the production ' $statement \rightarrow \textbf{output } expression ;$ '. And, in this case, the non-terminal ' $output_stmt$ ' becomes a useless symbol and is eliminated from the grammar.
- 11) The production ' $args \rightarrow args_list$ ' is replaced with the production ' $args \rightarrow arithmetic_expression args_list$ '. Now the symbol ' $args_list$ ' is a useless symbol and is eliminated from the grammar, afterwards, for the sake of readability, the symbol ' $args_list$ ' is changed into ' $args_list$ ' in all its appearances.

With these changes, the resulting grammar is the following:

Produccion				
No. Pr	LHS	→	RHS	RHS
1	program	→	declaration_list void ID (void) compound_stmt	void ID (void) compound_stmt
2	declaration_list	→	declaration declaration_list	epsilon
3	declaration	→	var_declaration	func_type_specifier ID (params) compound_stmt
4	var_declaration	→	int ID var_declaration'	
5	var_declaration'	→	;	[NUM] ;
6	func_type_specifier	→	int	void
7	params	→	param param_list	void
8	param_list	→	, param param_list	epsilon
9	param	→	int ID param'	
10	param'	→	[]	epsilon
11	compound_stmt	→	{ local_declarations statement_list }	
12	local_declarations	→	var_declaration local_declarations	epsilon
13	statement_list	→	statement statement_list	epsilon
14	statement	→	var = expression ; selection_stmt input var ;	call ; while (expression) statement output expression ;
15	selection_stmt	→	if (expression) statement selection_stmt'	
16	selection_stmt'	→	else statement	epsilon
17	return_stmt'	→	;	expression ;
18	var	→	ID var'	
19	var'	→	[arithmetic_expression]	epsilon
20	expression	→	arithmetic_expression expression'	
21	expression'	→	relop arithmetic_expression	epsilon
22	relop	→	<=	<
		→	>=	==
		→		!=
23	arithmetic_expression	→	term arithmetic_expression'	
24	arithmetic_expression'	→	addop term arithmetic_expression'	epsilon
25	addop	→	+	-
26	term	→	factor term'	
27	term'	→	mulop factor term'	epsilon
28	mulop	→	*	/
29	factor	→	(arithmetic_expression)	var
		→	NUM	call
30	call	→	ID (call'	
31	call'	→	args))
32	args	→	arithmetic_expression args_list	
33	args_list	→	, arithmetic_expression args_list	epsilon

Figure 7. The grammar after eliminating Unit Productions.

An unexpected visitor

The next step is to eliminate deep left factors. Why just now, as opposed to doing so along with the normal left factors? Deep left factor is not covered in class material, so students had to find out this problem by themselves. Accordingly, this step is covered after several attempts of constructing the Parsing Table and encountering multiple productions in a single cell. The method is not guaranteed to not affect the grammar, since it was an educated guess:

- 1) The production '*program* \rightarrow *declaration_list* **void** **ID** (**void**) *compound_stmt* | **void** **ID** (**void**) *compound_stmt*' is impossible to fix without modifying the grammar outside of its definition, so it is replaced with '*program* \rightarrow *declaration_list*' and '*declaration_list* \rightarrow *declaration* *declaration_list*' and '*declaration_list* \rightarrow *declaration* *declaration_list*' and '*declaration_list* \rightarrow ϵ '. This way there is still at least one declaration in the program, and the ID and parameters will be checked via code.
- 2) The productions '*declaration* \rightarrow *var_declaration* | *func_type_specifier* **ID** (*params*) *compound_stmt*' are replaced by the productions '*declaration* \rightarrow **int** **ID** *declaration*' | **void** **ID** (*params*) *compound_stmt*' and '*declaration* \rightarrow *var_declaration*' | (*params*) *compound_stmt*'. This makes the symbol '*var_declaration*' a useless symbol, and is eliminated from the grammar.
- 3) The productions '*statement* \rightarrow *var* = *expression* ; | *statement* \rightarrow *call* ;' are replaced with a single production: '*statement* \rightarrow **ID** *var_or_call_stmt*' and the following productions were added: '*var_or_call_stmt* \rightarrow *var* = *expression* ; | (*call* ;'.
- 4) The productions '*factor* \rightarrow *var* | *call*' were replaced with a single production: '*factor* \rightarrow **ID** *var_or_call*' and the following productions were added: '*var_or_call* \rightarrow *var*' | (*call*'.

Figure 8 depicts the grammar after the previous changes:

Produccion				
No. Pr	LHS	→	RHS	RHS
1	<i>program</i>	→	<i>declaration_list</i>	
2	<i>declaration_list</i>	→	<i>declaration declaration_list'</i>	
3	<i>declaration_list'</i>	→	<i>declaration declaration_list'</i>	epsilon
4	<i>declaration</i>	→	int <i>ID</i> <i>declaration'</i>	void <i>ID</i> (<i>params</i>) <i>compount_stmt</i>
5	<i>declaration'</i>	→	<i>var_declaration'</i>	(<i>params</i>) <i>compount_stmt</i>
6	<i>var_declaration'</i>	→	;	[NUM] ;
7	<i>params</i>	→	<i>param param_list</i>	void
8	<i>param_list</i>	→	, <i>param param_list</i>	epsilon
9	<i>param</i>	→	int <i>ID</i> <i>param'</i>	
10	<i>param'</i>	→	[]	epsilon
11	<i>compound_stmt</i>	→	{ <i>local_declarations statement_list</i> }	
12	<i>local_declarations</i>	→	int <i>ID</i> <i>var_declaration'</i> <i>local_declarations</i>	epsilon
13	<i>statement_list</i>	→	<i>statement statement_list</i>	epsilon
14	<i>statement</i>	→	<i>ID</i> <i>var_or_call_stmt</i>	<i>compound_stmt</i>
		→	while (<i>expression</i>) <i>statement</i>	<i>return</i> <i>return_stmt'</i>
		→	<i>output</i> <i>expression</i> ;	<i>input</i> <i>ID</i> <i>var'</i> ;
15	<i>var_or_call_stmt</i>	→	<i>var</i> = <i>expression</i> ;	(<i>call'</i> ;
16	<i>var_or_call</i>	→	<i>var'</i>	(<i>call'</i>
17	<i>selection_stmt</i>	→	if (<i>expression</i>) <i>statement</i> <i>selection_stmt'</i>	
18	<i>selection_stmt'</i>	→	else <i>statement</i>	epsilon
19	<i>return_stmt'</i>	→	;	<i>expression</i> ;
20	<i>var'</i>	→	[<i>arithmetic_expression</i>]	epsilon
21	<i>expression</i>	→	<i>arithmetic_expression expression'</i>	
22	<i>expression'</i>	→	<i>relop arithmetic_expression</i>	epsilon
23	<i>relop</i>	→	<=	<
		→	>=	=
				!=
24	<i>arithmetic_expression</i>	→	<i>term arithmetic_expression'</i>	
25	<i>arithmetic_expression'</i>	→	<i>addop</i> <i>term arithmetic_expression'</i>	epsilon
26	<i>addop</i>	→	+	-
27	<i>term</i>	→	<i>factor term'</i>	
28	<i>term'</i>	→	<i>mulop factor term'</i>	epsilon
29	<i>mulop</i>	→	*	/
30	<i>factor</i>	→	(<i>arithmetic_expression</i>)	<i>ID</i> <i>var_or_call</i>
				NUM
31	<i>call'</i>	→	<i>args</i>))
32	<i>args</i>	→	<i>arithmetic_expression args_list</i>	
33	<i>args_list</i>	→	, <i>arithmetic_expression args_list</i>	epsilon

Figure 8. The grammar after the elimination of deep left factors.

Finally one last simplification is made to the grammar:

- 1) The non-terminal symbol '*param*' will be eliminated by replacing all its calls with its contents, since '*param*' only appears in the left-hand-side of a production once, it is easy to expand its calls. The productions: '*params* → *param param_list*' and '*param_list* → , *param param_list*' will be replaced by the productions: '*params* → **int** *ID* *param'* *param_list*' and '*param_list* → , **int** *ID* *param'* *param_list*' respectively, and the production '*param* → **int** *ID* *param*' since the non-terminal '*param*' has become a useless symbol.

The grammar results as follows:

No. Pro	Produccion			
	LHS	→	RHS	RHS
1	program	→	declaration_list	
2	declaration_list	→	declaration declaration_list'	
3	declaration_list'	→	declaration declaration_list'	epsilon
4	declaration	→	int ID declaration'	void ID (params) compound_stmt
5	declaration'	→	var_declaration'	(params) compound_stmt
6	var_declaration'	→	;	[NUM] ;
7	params	→	int ID param' param_list	void
8	param_list	→	, int ID param' param_list	epsilon
9	param'	→	[]	epsilon
10	compound_stmt	→	{ local_declarations statement_list }	
11	local_declarations	→	int ID var_declaration' local_declarations	epsilon
12	statement_list	→	statement statement_list	epsilon
13	statement	→	ID var_or_call_stmt while (expression) statement output expression ;	compound_stmt return return_stmt' input ID var' ;
14	var_or_call_stmt	→	var' = expression ;	(call' ;
15	var_or_call	→	var'	(call'
16	selection_stmt	→	if (expression) statement selection_stmt'	
17	selection_stmt'	→	else statement	epsilon
18	return_stmt'	→	;	expression ;
19	var'	→	[arithmetic_expression]	epsilon
20	expression	→	arithmetic_expression expression'	
21	expression'	→	relop arithmetic_expression	epsilon
22	relop	→	<=	<
		→	>=	==
				>
				!=
23	arithmetic_expression	→	term arithmetic_expression'	
24	arithmetic_expression'	→	addop term arithmetic_expression'	epsilon
25	addop	→	+	-
26	term	→	factor term'	
27	term'	→	mulop factor term'	epsilon
28	mulop	→	*	/
29	factor	→	(arithmetic_expression)	ID var_or_call
30	call'	→	args)	NUM
31	args	→	arithmetic_expression args_list	
32	args_list	→	, arithmetic_expression args_list	epsilon

Figure 9. The final grammar.

This grammar still has one issue which will become clearer during the creation of the parsing table, and can be resolved during that step.

FIRST, FOLLOW, & FIRST+ sets & Parsing Table

Once the grammar is ready, the FIRST and FOLLOW sets must be calculated, in order to calculate the FIRST+ set, which is the one that will allow the LL(1) algorithm to make the parsing decisions and accept or reject the string of C minus code provided via the Parsing Table.

The FIRST (FIR) set is calculated for the terminals, non-terminals, and ϵ , according to the following rules:

- 1) If $X \rightarrow aY$, then $\text{First}(X) = \text{First}(X) \cup \{a\}$.
- 2) If $X \rightarrow \epsilon$, then $\text{First}(X) = \text{First}(X) \cup \{\epsilon\}$.
- 3) If $X \rightarrow Y_1Y_2...Y_n$, then :
 - a. $\text{First}(X) = \text{First}(X) \cup \text{First}(Y_1) - \{\epsilon\}$
 - b. If $Y_1 \Rightarrow \epsilon$, then
 $\text{First}(X) = \text{First}(X) \cup \text{First}(Y_2) - \{\epsilon\}$
 - c. If $Y_1 \Rightarrow \epsilon \wedge Y_2 \Rightarrow \epsilon \wedge ... \wedge Y_i \Rightarrow \epsilon \wedge i < n$, then
 $\text{First}(X) = \text{First}(X) \cup \text{First}(Y_{i+1}) - \{\epsilon\}$
 - d. If $Y_1 \Rightarrow \epsilon \wedge Y_2 \Rightarrow \epsilon \wedge ... \wedge Y_n \Rightarrow \epsilon$, then
 $\text{First}(X) = \text{First}(X) \cup \{\epsilon\}$

Figure 10. Rules of calculating the FIRST Set [2].

So for every terminal, non-terminal, and ϵ , the FIRST set is calculated:

TERMINALS & e	
else	{ else }
if	{ if }
input	{ input }
int	{ int }
output	{ output }
return	{ return }
void	{ void }
while	{ while }
+	{ + }
-	{ - }
*	{ * }
/	{ / }
>	{ > }
>=	{ >= }
<	{ < }
<=	{ <= }
==	{ == }
!=	{ != }
=	{ = }
;	{ ; }
,	{ , }
({ (}
)	{) }
{	{ { }
}	{ } }
[{ [}
]	{] }
IDENTIFIER	{ IDENTIFIER }
NUM_CONSTANT	{ NUM_CONSTANT }
epsilon	{ epsilon }

Figure 11. The FIRST set of the terminals and ϵ

For the non-terminals:

- 1) $FIR(program) = FIR(declaration_list) = \{\mathbf{int}, \mathbf{void}\}$
- 2) $FIR(declaration_list) = FIR(declaration) = \{\mathbf{int}, \mathbf{void}\}$
- 3) $FIR(declaration_list') = FIR(declaration) \cup FIR(\epsilon) = \{\epsilon, \mathbf{int}, \mathbf{void}\}$
- 4) $FIR(declaration) = FIR(int) \cup FIR(void) = \{\mathbf{int}, \mathbf{void}\}$
- 5) $FIR(declaration') = FIR(var_declaration') \cup FIR() = \{ ;, [, (\}$

- 6) $\text{FIR}(\text{var_declaration}') = \text{FIR}(;) \cup \text{FIR}([\] = \{ ;, [\] \}$
- 7) $\text{FIR}(\text{params}) = \text{FIR}(\text{int}) \cup \text{FIR}(\text{void}) = \{ \text{int}, \text{void} \}$
- 8) $\text{FIR}(\text{param_list}) = \text{FIR}(,) \cup \text{FIR}(\epsilon) = \{ \epsilon, , \}$
- 9) $\text{FIR}(\text{param}') = \text{FIR}([\]) \cup \text{FIR}(\epsilon) = \{ \epsilon, [\] \}$
- 10) $\text{FIR}(\text{compound_stmt}) = \text{FIR}(\{ \}) = \{ \{ \}$
- 11) $\text{FIR}(\text{local_declarations}) = \text{FIR}(\text{int}) \cup \text{FIR}(\epsilon) = \{ \epsilon, \text{int} \}$
- 12) $\text{FIR}(\text{statement_list}) = \text{FIR}(\text{statement}) \cup \text{FIR}(\epsilon) = \{ \epsilon, \text{ID}, \{, \text{if}, \text{while}, \text{return}, \text{input}, \text{output} \}$
- 13) $\text{FIR}(\text{statement}) = \text{FIR}(\text{ID}) \cup \text{FIR}(\text{compound_stmt}) \cup \text{FIR}(\text{selection_stmt}) \cup \text{FIR}(\text{while})$
 $\cup \text{FIR}(\text{return}) \cup \text{FIR}(\text{input}) \cup \text{FIR}(\text{output}) = \{ \text{ID}, \{, \text{if}, \text{while}, \text{return}, \text{input}, \text{output} \}$
- 14) $\text{FIR}(\text{var_or_call_stmt}) = \text{FIR}(\text{var}') \cup \text{FIR}([\] - \{ \epsilon \} \cup \text{FIR}(=) = \{ [, (, = \}$
- 15) $\text{FIR}(\text{var_or_call}) = \text{FIR}(\text{var}') \cup \text{FIR}([\] = \{ \epsilon, [, \{ \}$
- 16) $\text{FIR}(\text{selection_stmt}) = \text{FIR}(\text{if}) = \{ \text{if} \}$
- 17) $\text{FIR}(\text{selection_stmt}') = \text{FIR}(\text{else}) \cup \text{FIR}(\epsilon) = \{ \epsilon, \text{else} \}$
- 18) $\text{FIR}(\text{return_stmt}') = \text{FIR}(;) \cup \text{FIR}(\text{expression}) = \{ ;, (, \text{ID}, \text{NUM} \}$
- 19) $\text{FIR}(\text{var}') = \text{FIR}([\]) \cup \text{FIR}(\epsilon) = \{ \epsilon, [\] \}$
- 20) $\text{FIR}(\text{expression}) = \text{FIR}(\text{arithmetic_expression}) = \{ (, \text{ID}, \text{NUM} \}$
- 21) $\text{FIR}(\text{expression}') = \text{FIR}(\text{relop}) \cup \text{FIR}(\epsilon) = \{ \epsilon, <=, <, >, >=, ==, != \}$
- 22) $\text{FIR}(\text{relop}) = \text{FIR}(<=) \cup \text{FIR}(<) \cup \text{FIR}(>) \cup \text{FIR}(>=) \cup \text{FIR}(==) \cup \text{FIR}(!=) = \{ <=, <, >, >=, ==, != \}$
- 23) $\text{FIR}(\text{arithmetic_expression}) = \text{FIR}(\text{term}) = \{ (, \text{ID}, \text{NUM} \}$
- 24) $\text{FIR}(\text{arithmetic_expression}') = \text{FIR}(\epsilon) \cup \text{FIR}(\text{addop}) = \{ \epsilon, +, - \}$
- 25) $\text{FIR}(\text{addop}) = \text{FIR}(+) \cup \text{FIR}(-) = \{ +, - \}$
- 26) $\text{FIR}(\text{term}) = \text{FIR}(\text{factor}) = \{ (, \text{ID}, \text{NUM} \}$
- 27) $\text{FIR}(\text{term}') = \text{FIR}(\text{mulop}) \cup \text{FIR}(\epsilon) = \{ \epsilon, *, / \}$
- 28) $\text{FIR}(\text{mulop}) = \text{FIR}(*) \cup \text{FIR}(/) = \{ *, / \}$
- 29) $\text{FIR}(\text{factor}) = \text{FIR}([\]) \cup \text{FIR}(\text{ID}) \cup \text{FIR}(\text{NUM}) = \{ (, \text{ID}, \text{NUM} \}$
- 30) $\text{FIR}(\text{call}') = \text{FIR}(\text{args}) \cup \text{FIR}([\]) = \{ (,), \text{ID}, \text{NUM} \}$
- 31) $\text{FIR}(\text{args}) = \text{FIR}(\text{arithmetic_expression}) = \{ (, \text{ID}, \text{NUM} \}$
- 32) $\text{FIR}(\text{args_list}) = \text{FIR}(,) \cup \text{FIR}(\epsilon) = \{ \epsilon, , \}$

Calculating the FIRST set:	
NON-TERMINALS	
<i>program</i>	{ int, void }
<i>declaration_list</i>	{ int, void }
<i>declaration_list'</i>	{ int, void, epsilon }
<i>declaration</i>	{ int, void }
<i>declaration'</i>	{ :, [, (}
<i>var_declaration'</i>	{ :, [}
<i>params</i>	{ int, void }
<i>param_list</i>	{ ,, epsilon }
<i>param'</i>	{ [, epsilon }
<i>compound_stmt</i>	{ { }
<i>local_declarations</i>	{ int, epsilon }
<i>statement_list</i>	{ ID, {, if, while, return, input, output, epsilon }
<i>statement</i>	{ ID, {, if, while, return, input, output }
<i>var_or_call_stmt</i>	{ (, [, = }
<i>var_or_call</i>	{ (, [, epsilon }
<i>selection_stmt</i>	{ if }
<i>selection_stmt'</i>	{ else, epsilon }
<i>return_stmt'</i>	{ :, (, ID, NUM }
<i>var'</i>	{ [, epsilon }
<i>expression</i>	{ (, ID, NUM }
<i>expression'</i>	{ <, <=, >, >=, ==, !=, epsilon }
<i>relop</i>	{ <, <=, >, >=, ==, != }
<i>arithmetic_expression</i>	{ (, ID, NUM }
<i>arithmetic_expression</i>	{ +, -, epsilon }
<i>addop</i>	{ +, - }
<i>term</i>	{ (, ID, NUM }
<i>term'</i>	{ *, /, epsilon }
<i>mulop</i>	{ *, / }
<i>factor</i>	{ (, ID, NUM }
<i>call'</i>	{ (,), ID, NUM }
<i>args</i>	{ (, ID, NUM }
<i>args_list</i>	{ ,, epsilon }

Figure 12. The FIRST sets of the non-terminals of the grammar.

Next, the FOLLOW (FOL) set is calculated for the non-terminals, according to the rules in **Figure 13**:

- 1) For the start symbol, **Follow(S) = Follow(S) \cup { \$ }.**
- 2) If there is a production **$X \rightarrow \alpha Y \beta$** , then
Follow(Y) = First(β) - { ϵ }.
- 3) If there is a production **$X \rightarrow \alpha Y \beta \theta$** , where **$\epsilon \in \text{First}(\beta)$** then,
Follow(Y) = First(β) - { ϵ } \cup First(θ) - { ϵ }
- 4) If there is a production **$X \rightarrow \alpha Y$** or **$X \rightarrow \alpha Y \beta$** and **$\epsilon \in \text{First}(\beta)$** then,
Follow(Y) = Follow(Y) \cup Follow(X).

Figure 13. Rules for calculating the FOLLOW set of a set of non-terminals.

For every non-terminal, the FOLLOW set is calculated:

- 1) $\text{FOL}(\text{program}) = \text{FOL}(\text{program}) \cup \{ \$ \} = \{ \$ \}$
- 2) $\text{FOL}(\text{declaration_list}) = \text{FOL}(\text{program}) = \{ \$ \}$
- 3) $\text{FOL}(\text{declaration_list}') = \text{FOL}(\text{declaration_list}) \cup \text{FOL}(\text{declaration_list}') = \{ \$ \}$
- 4) $\text{FOL}(\text{declaration}) = \text{FIR}(\text{declaration_list}') - \{ \epsilon \} \cup \text{FOL}(\text{declaration_list}) \cup \text{FOL}(\text{declaration_list}') = \{ \text{int, void, } \$ \}$
- 5) $\text{FOL}(\text{declaration}') = \text{FOL}(\text{declaration}) = \{ \text{int, void, } \$ \}$
- 6) $\text{FOL}(\text{var_declaration}') = \text{FOL}(\text{declaration}') = \{ \text{int, void, } \$ \}$
- 7) $\text{FOL}(\text{params}) = \text{FIR}() = \{) \}$
- 8) $\text{FOL}(\text{param_list}) = \text{FOL}(\text{params}) \cup \text{FOL}(\text{param_list}) = \{) \}$
- 9) $\text{FOL}(\text{param}') = \text{FIR}(\text{param_list}) - \{ \epsilon \} \cup \text{FOL}(\text{params}) \cup \text{FOL}(\text{params_list}) = \{ ,,) \}$
- 10) $\text{FOL}(\text{compound_stmt}) = \text{FOL}(\text{declaration}) \cup \text{FOL}(\text{declaration}') \cup \text{FOL}(\text{statement}) = \{ \text{int, void, } \$, \text{ID, } \{, \}, \text{if, else, while, return, input, output} \}$
- 11) $\text{FOL}(\text{local_declarations}) = \text{FIR}(\text{statement_list}) - \{ \epsilon \} \cup \text{FIR}() \cup \text{FOL}(\text{local_declarations}) = \{ \text{ID, } \{, \}, \text{if, while, return, input, output} \}$
- 12) $\text{FOL}(\text{statement_list}) = \text{FIR}() \cup \text{FOL}(\text{statement_list}) = \{ \}$
- 13) $\text{FOL}(\text{statement}) = \text{FIR}(\text{statement_list}) - \{ \epsilon \} \cup \text{FOL}(\text{statement_list}) \cup \text{FOL}(\text{statement}) \cup \text{FIR}(\text{selection_stmt}') \cup \text{FOL}(\text{selection_stmt}') = \{ \text{ID, } \{, \text{if, else, while, return, input, output} \}$
- 14) $\text{FOL}(\text{var_or_call_stmt}) = \text{FOL}(\text{statement}) = \{ \text{ID, } \{, \text{if, else, while, return, input, output, } \}$
- 15) $\text{FOL}(\text{var_or_call}) = \text{FOL}(\text{factor}) = \{ *, /, +, -,], <, <=, >, >=, ==, !=,), ;, ,, \}$

- 16) $\text{FOL}(\text{selection_stmt}) = \text{FOL}(\text{statement}) = \{\text{ID}, \{, \text{if}, \text{else}, \text{while}, \text{return}, \text{input}, \text{output}, \}\}$
- 17) $\text{FOL}(\text{selection_stmt}') = \text{FOL}(\text{selection_stmt}) = \{\text{ID}, \{, \text{if}, \text{else}, \text{while}, \text{return}, \text{input}, \text{output}, \}\}$
- 18) $\text{FOL}(\text{return_stmt}') = \text{FOL}(\text{statement}) = \{\text{ID}, \{, \}, \text{if}, \text{else}, \text{while}, \text{return}, \text{input}, \text{output}\}$
- 19) $\text{FOL}(\text{var}') = \text{FIR}(;) \cup \text{FIR}(=) \cup \text{FOL}(\text{var_or_call}) = \{;, =, *, /, +, -,], <, <=, >, >=, ==, !=,), , \}$
- 20) $\text{FOL}(\text{expression}) = \text{FIR}()) \cup \text{FIR}(;) = \{), , \}$
- 21) $\text{FOL}(\text{expression}') = \text{FOL}(\text{expression}) = \{), , \}$
- 22) $\text{FOL}(\text{relop}) = \text{FIR}(\text{arithmetic_expression}) = \{ (, \text{ID}, \text{NUM} \}$
- 23) $\text{FOL}(\text{arithmetic_expression}) = \text{FIR}([]) \cup \text{FIR}(\text{expression}') \cup \text{FOL}(\text{expression}) \cup \text{FOL}(\text{expression}') \cup \text{FIR}(\text{args_list}) \cup \text{FOL}(\text{args}) \cup \text{FOL}(\text{args_list}) = \{], <=, <, >, >=, ==, !=,), , , \}$
- 24) $\text{FOL}(\text{arithmetic_expression}') = \text{FOL}(\text{arithmetic_expression}) = \{], <=, <, >, >=, ==, !=,), , , \}$
- 25) $\text{FOL}(\text{addop}) = \text{FIR}(\text{term}) = \{ (, \text{ID}, \text{NUM} \}$
- 26) $\text{FOL}(\text{term}) = \text{FIR}(\text{arithmetic_expression}') \cup \text{FOL}(\text{arithmetic_expression}') = \{ +, -,], <=, <, >, >=, ==, !=,), , , \}$
- 27) $\text{FOL}(\text{term}') = \text{FOL}(\text{term}) \cup \text{FOL}(\text{term}') = \{ +, -,], <=, <, >, >=, ==, !=,), , , \}$
- 28) $\text{FOL}(\text{mulop}) = \text{FIR}(\text{factor}) = \{ (, \text{ID}, \text{NUM} \}$
- 29) $\text{FOL}(\text{factor}) = \text{FIR}(\text{term}') - \{ \epsilon \} \cup \text{FOL}(\text{term}) \cup \text{FOL}(\text{term}') = \{ *, /, +, -,], <, <=, >, >=, ==, !=,), , , \}$
- 30) $\text{FOL}(\text{call}') = \text{FIR}(;) \cup \text{FOL}(\text{var_or_call}) = \{ *, /, +, -,], <, <=, >, >=, ==, !=,), , , \}$
- 31) $\text{FOL}(\text{args}) = \text{FIR}()) = \{) \}$
- 32) $\text{FOL}(\text{args_list}) = \text{FOL}(\text{args}) \cup \text{FOL}(\text{args_list}) = \{) \}$

Calculating the FOLLOW set:	
NON-TERMINALS	
<i>program</i>	{ \$ }
<i>declaration_list</i>	{ \$ }
<i>declaration_list'</i>	{ \$ }
<i>declaration</i>	{ int, void, \$ }
<i>declaration'</i>	{ int, void, \$ }
<i>var_declaration'</i>	{ int, void, \$ }
<i>params</i>	{) }
<i>param_list</i>	{) }
<i>param'</i>	{ ,,) }
<i>compound_stmt</i>	{ \$, int, void, ID, {, }, if, while, return, input, output, else }
<i>local_declarations</i>	{ ID, {, }, if, while, return, input, output }
<i>statement_list</i>	{ }
<i>statement</i>	{ ID, {, }, if, while, return, input, output, else }
<i>var_or_call_stmt</i>	{ ID, {, }, if, while, return, input, output, else }
<i>var_or_call</i>	{ *, /, +, -,], <, <=, >, >=, ==, !=,), ::, , }
<i>selection_stmt</i>	{ ID, {, }, if, while, return, input, output, else }
<i>selection_stmt'</i>	{ ID, {, }, if, while, return, input, output, else }
<i>return_stmt'</i>	{ ID, {, }, if, while, return, input, output, else }
<i>var'</i>	{ =, *, /, +, -,], <, <=, >, >=, ==, !=,), ::, , }
<i>expression</i>	{), ; }
<i>expression'</i>	{), ; }
<i>relop</i>	{ (, ID, NUM }
<i>arithmetic_expression</i>	{], <, <=, >, >=, ==, !=,), ::, , }
<i>arithmetic_expression'</i>	{], <, <=, >, >=, ==, !=,), ::, , }
<i>addop</i>	{ (, ID, NUM }
<i>term</i>	{ +, -,], <, <=, >, >=, ==, !=,), ::, , }
<i>term'</i>	{ +, -,], <, <=, >, >=, ==, !=,), ::, , }
<i>mulop</i>	{ (, ID, NUM }
<i>factor</i>	{ *, /, +, -,], <, <=, >, >=, ==, !=,), ::, , }
<i>call'</i>	{ *, /, +, -,], <, <=, >, >=, ==, !=,), ::, , }
<i>args</i>	{) }
<i>args_list</i>	{) }

Figure 14. The FOLLOW sets for the non-terminals of the grammar.

Finally, the FIRST+ (FIP) set is calculated for the non-terminals according to the rules in Figure 15:

For each production $X \rightarrow \beta$, its **augmented First Set** $\text{First}^+(X \rightarrow \beta)$ is defined as follows:

- 1) If $\epsilon \notin \text{First}(\beta)$, then $\text{First}^+(X \rightarrow \beta) = \text{First}(\beta)$.
- 2) If $\epsilon \in \text{First}(\beta)$, then

$$\text{First}^+(X \rightarrow \beta) = \text{First}(\beta) \cup \text{Follow}(X).$$

Figure 15. Rules for calculating the FIRST+ set of a grammar [2].

For every production, the set is calculated:

- 1) $\text{FIP}(1) = \text{FIR}(\text{declaration_list}) = \{\text{int}, \text{void}\}$
- 2) $\text{FIP}(2) = \text{FIR}(\text{declaration}) = \{\text{int}, \text{void}\}$
- 3) $\text{FIP}(3) = \text{FIR}(\text{declaration}) = \{\text{int}, \text{void}\}$
- 4) $\text{FIP}(4) = \text{FIR}(\epsilon) \cup \text{FOL}(\text{declaration_list}') = \{\epsilon, \$\}$
- 5) $\text{FIP}(5) = \text{FIR}(\text{int}) = \{\text{int}\}$
- 6) $\text{FIP}(6) = \text{FIR}(\text{void}) = \{\text{void}\}$
- 7) $\text{FIP}(7) = \text{FIR}(\text{var_declaration}') = \{;, []\}$
- 8) $\text{FIP}(8) = \text{FIR}() = \{ (\}$
- 9) $\text{FIP}(9) = \text{FIR}(;) = \{ ; \}$

- 10) $FIP(10) = FIR([]) = \{ \} \}$
- 11) $FIP(11) = FIR(int) = \{int\}$
- 12) $FIP(12) = FIR(void) = \{void\}$
- 13) $FIP(13) = FIR(,) = \{ , \}$
- 14) $FIP(14) = FIR(\epsilon) \cup FOL(param_list) = \{\epsilon, \}$
- 15) $FIP(15) = FIR([]) = \{ [\}$
- 16) $FIP(16) = FIR(\epsilon) \cup FOL(param') = \{\epsilon, ,, \}$
- 17) $FIP(17) = FIR(\{ \}) = \{ \{ \}$
- 18) $FIP(18) = FIR(int) = \{int\}$
- 19) $FIP(19) = FIR(\epsilon) \cup FOL(local_declarations) = \{\epsilon, ID, \{, \}, if, while, return, input, output\}$
- 20) $FIP(20) = FIR(statement) = \{ID, \{, if, while, return, input, output\}$
- 21) $FIP(21) = FIR(\epsilon) \cup FOL(statement_list) = \{\epsilon, \}$
- 22) $FIP(22) = FIR(ID) = \{ ID \}$
- 23) $FIP(23) = FIR(compound_stmt) = \{ \{ \}$
- 24) $FIP(24) = FIR(selection_stmt) = \{if\}$

- 25) FIP(25) = FIR(while) = { while }
- 26) FIP(26) = FIR(return) = { return }
- 27) FIP(27) = FIR(input) = { input }
- 28) FIP(28) = FIR(output) = { output }
- 29) FIP(29) = FIR(*var* ') u FIR(=) = { ϵ , [, =}
- 30) FIP(30) = FIR() = { (}
- 31) FIP(31) = FIR(*var* ') u FOL(*var_or_call*) = { ϵ , [, *, /, +, -,], <, <=, >, >=, ==, !=,), ;, ,}
- 32) FIP(32) = FIR() = { (}
- 33) FIP(33) = FIR(if) = { if }
- 34) FIP(34) = FIR(else) = { else }
- 35) FIP(35) = FIR(ϵ) u FOL(*selection_stmt* ') = { ϵ , ID, {, }, if, else, while, return, input, output}
- 36) FIP(36) = FIR(;) = { ; }
- 37) FIP(37) = FIR(*expression*) = { (, ID, NUM }
- 38) FIP(38) = FIR[] = { [}
- 39) FIP(39) = FIR(ϵ) u FOL(*var* ') = { ϵ , ;, =, *, /, +, -,], <, <=, >, >=, ==, !=,), ;}
- 40) FIP(40) = FIR(*arithmetic_expression*) = {(, ID, NUM}
- 41) FIP(41) = FIR(*relop*) = {<, <=, >, >=, ==, !=}
- 42) FIP(42) = FIR(ϵ) u FOL(*expression* ') = { ϵ ,), ;}
- 43) FIP(43) = FIR(<=) = {<= }
- 44) FIP(44) = FIR(<) = {< }
- 45) FIP(45) = FIR(>) = {> }
- 46) FIP(46) = FIR(>=) = {>= }
- 47) FIP(47) = FIR(==) = {== }
- 48) FIP(48) = FIR(!=) = {!= }
- 49) FIP(49) = FIR(*term*) = {(, ID, NUM}
- 50) FIP(50) = FIR(*addop*) = {+, -}
- 51) FIP(51) = FIR(ϵ) u FOL(*arithmetic_expression* ') = { ϵ ,], <, <=, >, >=, ==, !=,), ;, ;}
- 52) FIP(52) = FIR(+) = {+ }
- 53) FIP(53) = FIR(-) = {- }
- 54) FIP(54) = FIR(*factor*) = {(, ID, NUM}
- 55) FIP(55) = FIR(*mulop*) = {*, /}
- 56) FIP(56) = FIR(ϵ) u FOL(*term* ') = { ϵ , +, -,], <, <=, >, >=, ==, !=,), ;, ;}

$$57) \text{FIP}(57) = \text{FIR}(*) = \{ * \}$$

$$58) \text{FIP}(58) = \text{FIR}(/) = \{ / \}$$

$$59) \text{FIP}(59) = \text{FIR}() = \{ (\}$$

$$60) \text{FIP}(60) = \text{FIR}(\text{ID}) = \{ \text{ID} \}$$

$$61) \text{FIP}(61) = \text{FIR}(\text{NUM}) = \{ \text{NUM} \}$$

$$62) \text{FIP}(62) = \text{FIR}(\text{args}) = \{ (, \text{ID}, \text{NUM} \}$$

$$63) \text{FIP}(63) = \text{FIR}()) = \{) \}$$

$$64) \text{FIP}(64) = \text{FIR}(\text{arithmetic_expression}) = \{ (, \text{ID}, \text{NUM} \}$$

$$65) \text{FIP}(65) = \text{FIR}(,) = \{ , \}$$

$$66) \text{FIP}(66) = \text{FIR}(\epsilon) \cup \text{FOL}(\text{args_list}) = \{ \epsilon,) \}$$

Calculating the FIRST+ set:		
PRODUCTIONS		
1	$program \rightarrow declaration_list$	{ int, void }
2	$declaration_list \rightarrow declaration\ declaration_list'$	{ int, void }
3	$declaration_list' \rightarrow declaration\ declaration_list'$	{ int, void }
4	$declaration_list' \rightarrow \epsilon$	{ epsilon , \$ }
5	$declaration \rightarrow int\ ID\ declaration'$	{ int }
6	$declaration \rightarrow void\ ID\ (params)\ compound_stmt$	{ void }
7	$declaration' \rightarrow var_declaration'$	{ :, [}
8	$declaration' \rightarrow (params)\ compound_stmt$	{ (}
9	$var_declaration' \rightarrow ;$	{ ; }
10	$var_declaration' \rightarrow [NUM] ;$	{ [}
11	$params \rightarrow int\ ID\ param'\ param_list$	{ int }
12	$params \rightarrow void$	{ void }
13	$param_list \rightarrow , int\ ID\ param'\ param_list$	{ , }
14	$param_list \rightarrow \epsilon$	{ epsilon ,) }
15	$param' \rightarrow []$	{ [}
16	$param' \rightarrow \epsilon$	{ epsilon , ,,) }
17	$compound_stmt \rightarrow \{ local_declarations\ statement_list \}$	{ { }
18	$local_declarations \rightarrow int\ ID\ var_declaration'\ local_declarations$	{ int }
19	$local_declarations \rightarrow \epsilon$	{ epsilon , ID, {, }, if, while, return, input, output }
20	$statement_list \rightarrow statement\ statement_list$	{ ID, {, if, while, return, input, output }
21	$statement_list \rightarrow \epsilon$	{ epsilon , } }
22	$statement \rightarrow ID\ var_or_call_stmt$	{ ID }
23	$statement \rightarrow compound_stmt$	{ { }
24	$statement \rightarrow selection_stmt$	{ if }
25	$statement \rightarrow while\ (expression)\ statement$	{ while }
26	$statement \rightarrow return\ return_stmt'$	{ return }
27	$statement \rightarrow input\ ID\ var' ;$	{ input }
28	$statement \rightarrow output\ expression ;$	{ output }
29	$var_or_call_stmt \rightarrow var' = expression ;$	{ epsilon , [, = }
30	$var_or_call_stmt \rightarrow (call' ;$	{ (}
31	$var_or_call \rightarrow var'$	{ epsilon , [, *, /, +, -, <, <=, >, >=, ==, !=,), :, , }
32	$var_or_call \rightarrow (call'$	{ (}
33	$selection_stmt \rightarrow if\ (expression)\ statement\ selection_stmt'$	{ if }

34	selection_stmt' → else statement	{ else }
35	selection_stmt' → epsilon	{ epsilon , ID, {, }, if, else, while, return, input, output }
36	return_stmt' → ;	{ ; }
37	return_stmt' → expression ;	{ (, ID, NUM }
38	var' → [arithmetic_expression]	{ [}
39	var' → epsilon	{ epsilon , ;, =, (, ID, NUM }
40	expression → arithmetic_expression expression'	{ (, ID, NUM }
41	expression' → relop arithmetic_expression	{ <, <=, >, >=, ==, != }
42	expression' → epsilon	{ epsilon ,), ; }
43	relop → <=	{ <= }
44	relop → <	{ < }
45	relop → >	{ > }
46	relop → >=	{ >= }
47	relop → ==	{ == }
48	relop → !=	{ != }
49	arithmetic_expression → term arithmetic_expression'	{ (, ID, NUM }
50	arithmetic_expression' → addop term arithmetic_expression'	{ +, - }
51	arithmetic_expression' → epsilon	{ epsilon ,], <, <=, >, >=, ==, !=,), ;, , }
52	addop → +	{ + }
53	addop → -	{ - }
54	term → factor term'	{ (, ID, NUM }
55	term' → mulop factor term'	{ *, / }
56	term' → epsilon	{ epsilon , +, -,], <, <=, >, >=, ==, !=,), ;, , }
57	mulop → *	{ * }
58	mulop → /	{ / }
59	factor → (arithmetic_expression)	{ (}
60	factor → ID var_or_call	{ ID }
61	factor → NUM	{ NUM }
62	call' → args)	{ (, ID, NUM }
63	call' →)	{) }
64	args → arithmetic_expression args_list	{ (, ID, NUM }
65	args_list → , arithmetic_expression args_list	{ , }
66	args_list → epsilon	{ epsilon ,) }

Figure 16. The FIRST+ Sets of the grammar.

Next, the Parsing Table is calculated according to the algorithm:

```

3) Fill the parsing table as follows:
for each non-terminal  $X$  do:
  for each terminal  $a$  do: // initialize table.
    Table[ $X, a$ ] = error; // empty cells will be errors.
  end
  for each production  $p: X \rightarrow \beta$  do:
    for each terminal  $w \in \text{First}^+(X \rightarrow \beta) - \epsilon$  do:
      Table[ $X, w$ ] =  $p$ ;
    end
    if  $\$ \in \text{First}^+(X \rightarrow \beta)$  then
      Table[ $X, \$$ ] =  $p$ ;
    end
  end
end

```

63

Figure 17. Algorithm for filling the Parsing Table.

The Parsing Table is filled using the FIRST+ sets, and results as follows:

	else	if	input	int	output	return	void	while	+	-	*	/	>	>=	<	<=	==	!=	=	;	,	()	{	}	[]	ID	NUM	\$
program				1			1																							
declaration_list				2			2																							
declaration_list'				3			3																							
declaration				5			6																						4	
declaration'																					7		8				7			
var_declaration'																				9						10				
params				11			12																							
param_list																						13		14						
param'																					16		16				15			
compound_stmt																								17						
local_declarations		19	19	18	19	19		19																19	19			19		
statement_list		20	20		20	20		20																20	21			20		
statement		24	27		28	26		25																23				22		
var_or_call_stmt																				29			30				29			
var_or_call									31	31	31	31	31	31	31	31	31	31			31	31	32	31			31	31		
selection_stmt		33																												
selection_stmt'	34, 35	35	35		35	35		35																				35		
return_stmt'																												37	37	
var'									39	39	39	39	39	39	39	39	39	39	39	39	39	39		39		38	39			
expression																							40				40	40		
expression'																							42							
relop																														
arithmetic_expression																												49	49	
arithmetic_expression'									50	50				51	51	51	51	51	51		51	51					51			
addop									52	53																				
term																												54	54	
term'																														
mulop									56	56	55	55	56	56	56	56	56	56	56		56	56		56			56			
factor											57	58																		
call'																												60	61	
args																												62	62	
args_list																												64	64	

Figure 18. The parsing table for the resulting C Minus grammar.

Note that one cell, corresponding to $M[\text{selection_stmt}', \text{else}]$ has two results, this is because of the ‘dangling else’ ambiguity. In this case the correct derivation for the grammar is the one where $M[\text{selection_stmt}', \text{else}] = 34$ so the cell is changed accordingly. The following unambiguous Parsing Table will be translated into code:

	else	if	input	int	output	return	void	while	+	-	*	/	>	>=	<	<=	==	!=	=	;	,	()	{	}	[]	ID	NUM	\$
program				1				1																						
declaration_list				2				2																						
declaration_list'				3				3																					4	
declaration				5				6																						
declaration'																														
var_declaration'																				7		8				7				
params																				9						10				
param_list				11				12																						
param'																						13		14						
compound_stmt																					16					15				
local_declarations		19	19	18	19	19		19																17		19		19		
statement_list		20	20		20	20		20																20	21		20			
statement		24	27		28	26		25																23				22		
var_or_call_stmt																				29			30			29				
var_or_call'										31	31	31	31	31	31	31	31	31	31		31	31	32	31		31	31			
selection_stmt			33																											
selection_stmt'	34	35	35		35	35		35																	35	35		35		
return_stmt'																												37	37	
var'										39	39	39	39	39	39	39	39	39	39	39	39	39	39	39		38	39			
expression																												40	40	
expression'													41	41	41	41	41	41				42		42						
relop												45	46	44	43	47	48													
arithmetic_expression																								49				49	49	
arithmetic_expression'																												51		
addop																														
term																														
term'																								54				54	54	
mulop																														
factor																														
call'																														
args																														
args_list																														
				</																										

Figure 19. The unambiguous Parsing Table for the resulting C Minus grammar.

Design

During the design phase, diagrams made during the development of the Lexical Analyser were updated for compatibility with the design made for this phase of compilation, including: Class Diagram, Sequence Diagram, and the pseudo code provided by the professor.

A diagram of the architecture for an LL(1) Parser provided by the professor, shown in **Figure 20**, was used to create classes for expanding the original class diagram:

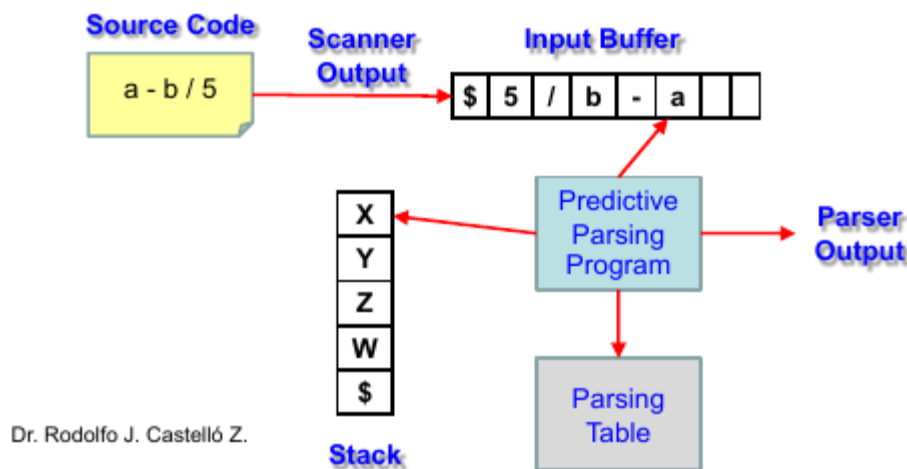


Figure 20. The basic architecture of an LL(1) Parser.

The diagram in **Figure 21** shall illustrate the interactions between the classes described previously; the classes for steps of compilation other than lexical analysis have been omitted as they are not relevant at this stage of compilation.

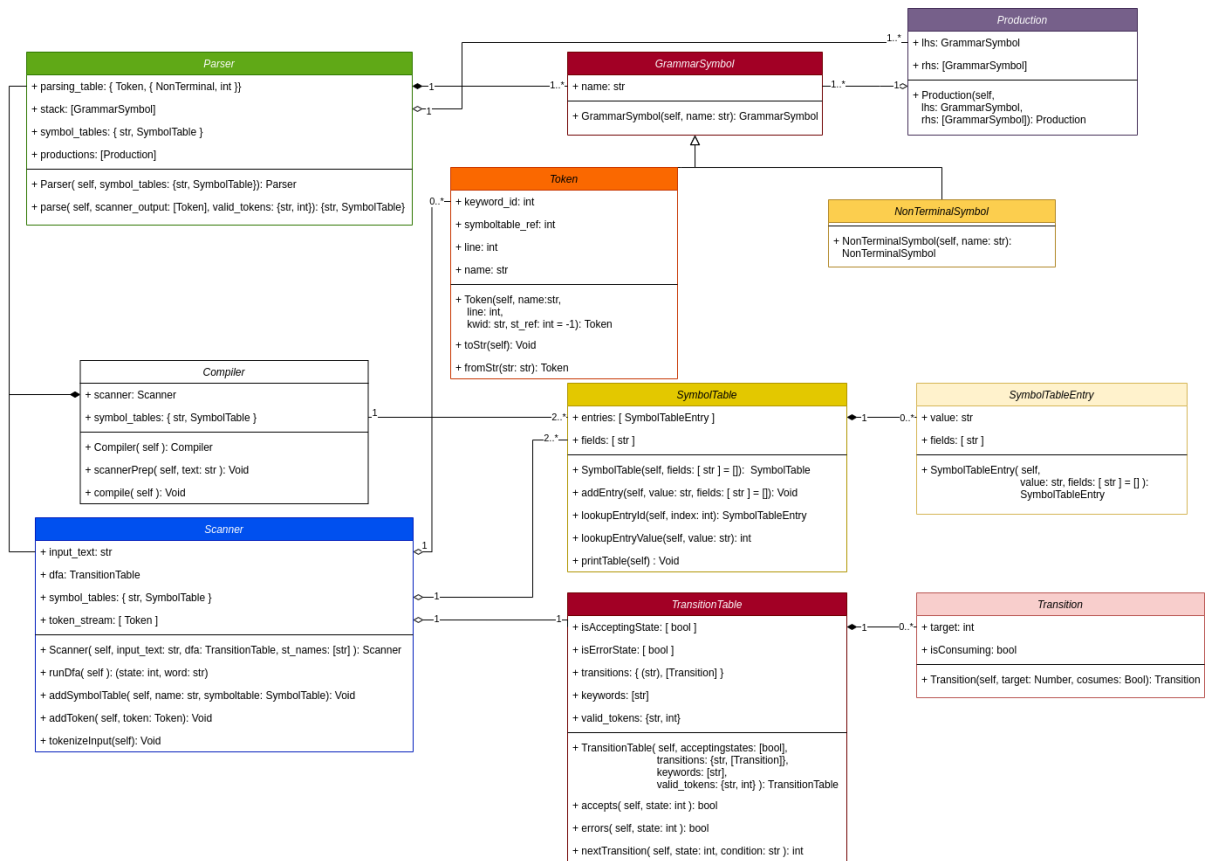


Figure 21. The entire Class Diagram for the compiler program, including updates.

As considered, the *Compiler* class is composed of a *Scanner* object, and now of a *Parser* object, too; a *Production* Class, *GrammarSymbol* class, and a *NonTerminalSymbol* class. Now, the *Token* class inherits from the *GrammarSymbol* class, along with *NonTerminalSymbol*, so that the stack of the parser can be treated as a homogenous list of *GrammarSymbols*. The list of productions from which the Parsing Table will take is a list of *Production* Objects whose right-hand-side is a list of *GrammarSymbols*: *Token* (terminal) objects and *NonTerminalSymbol* (non-terminal) objects that are pushed onto the stack.

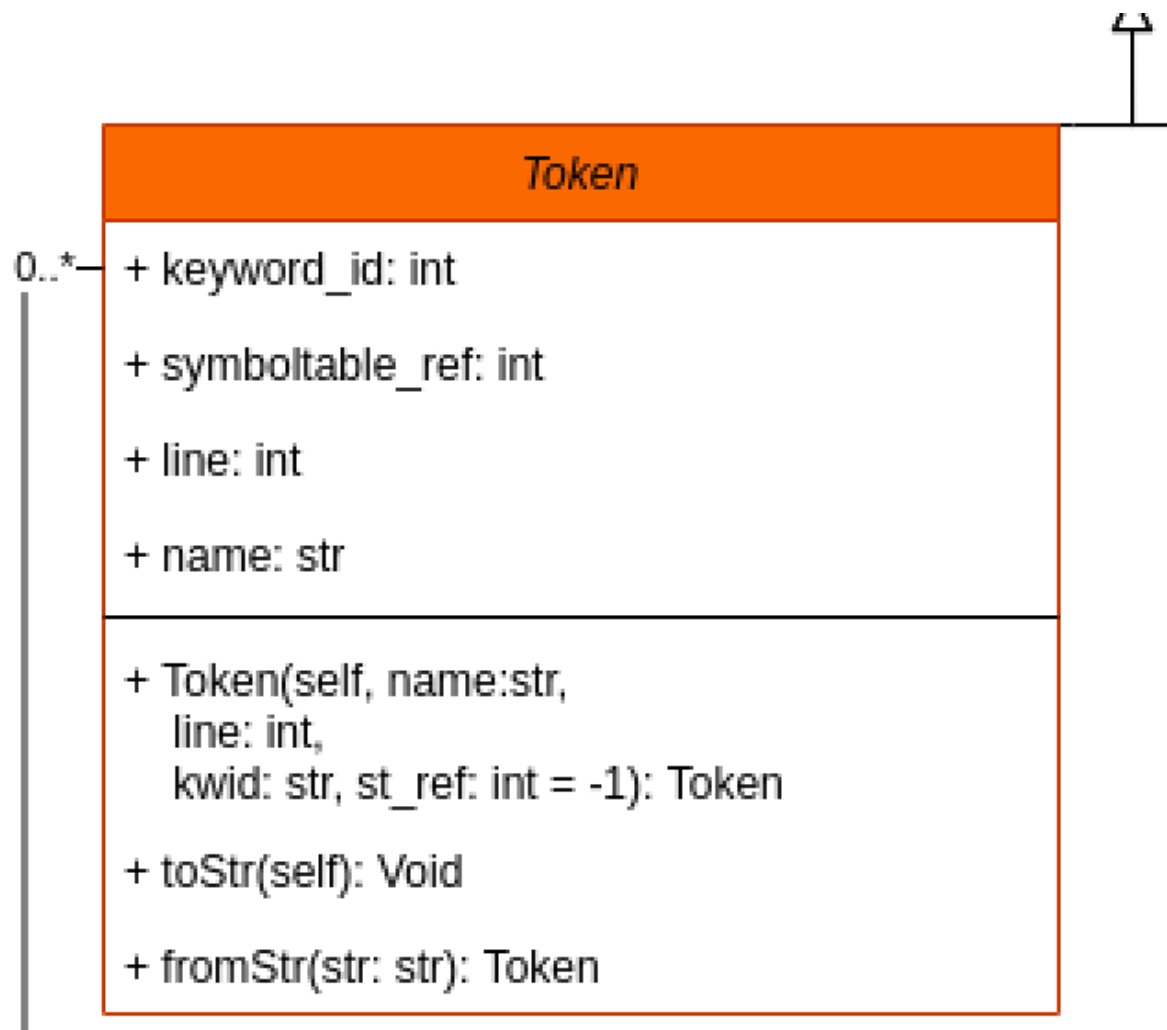


Figure 22. The updated Token class.

The Token class was updated to receive an additional two values in its constructor: the *name* field, and the *line* field. The *name* field serves as a shorthand when comparing the output from the scanner and the elements in the stack, as well as creating a point of comparison between other GrammarSymbol objects that are not necessarily Token objects. The *line* field is a new value passed from the scanner phase that allows the syntax analyser to know the line of the token in case an error is detected.

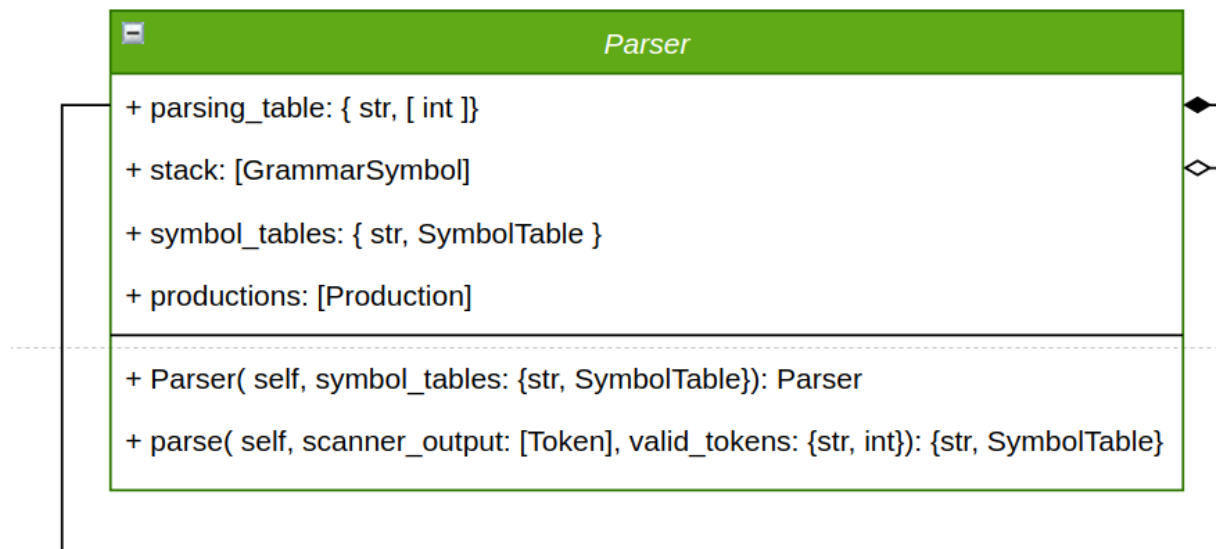


Figure 23. The Parser Class.

The Parser class contains the elements of a parser as described by the LL(1) parser algorithm: a stack, an input of Tokens, the parsing program in the form of the *parse* method, and the parsing table. The parsing table is a dictionary that receives the name of a NonTerminalSymbol in the form of a string, and with an index that corresponds to the Token ID - 1, returns a number corresponding to the number of a production. The Stack is a list which, as the parsing begins, is appended the grammar symbols '\$' and 'program' (the start symbol of the grammar). The symbol tables contain the information about the ID Tokens that have been found, and, during parsing, are updated with new information about them. The Production list is a list of productions, whose order corresponds to the output of the parsing table minus one (- 1), so that the parser knows which GrammarSymbols to add to the stack.

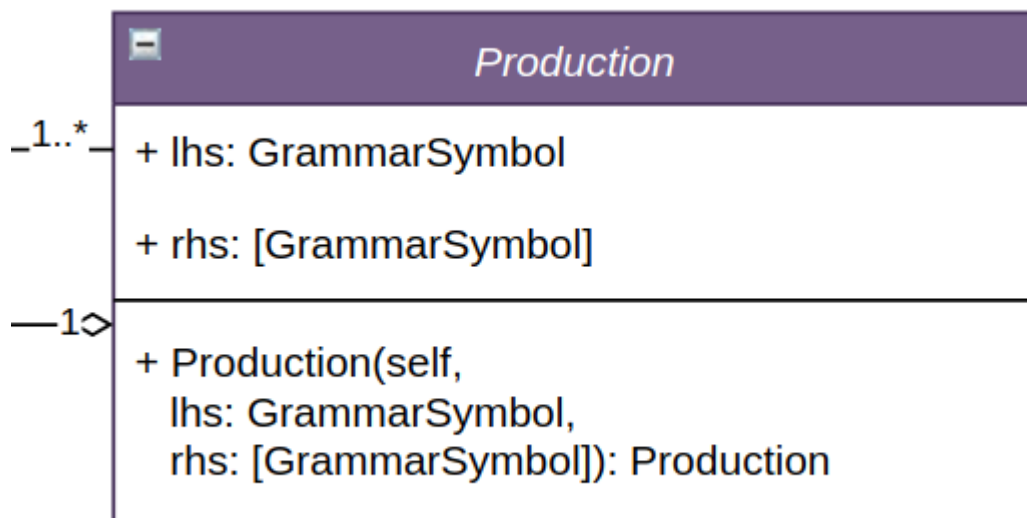


Figure 24. The Production class.

The production class serves as an abstraction for the productions of a grammar, in this case the list of productions from the resulting C Minus grammar are passed into code via the Production class.

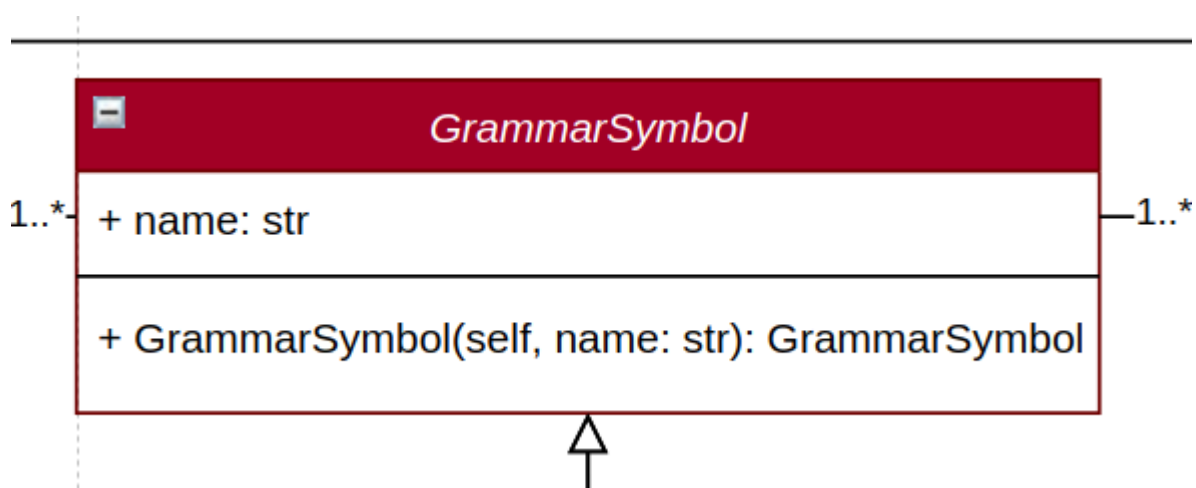


Figure 25. The GrammarSymbol class.

The GrammarSymbol class serves as an abstraction for all terminals, non-terminals, and pseudo-tokens that are used in the definition of a grammar. Pseudo tokens such as ‘\$’ and ‘ε’ do not have a token id, since they are not Tokens, and are not non-terminal symbols, but can be represented as simple grammar symbols with their own unique name.

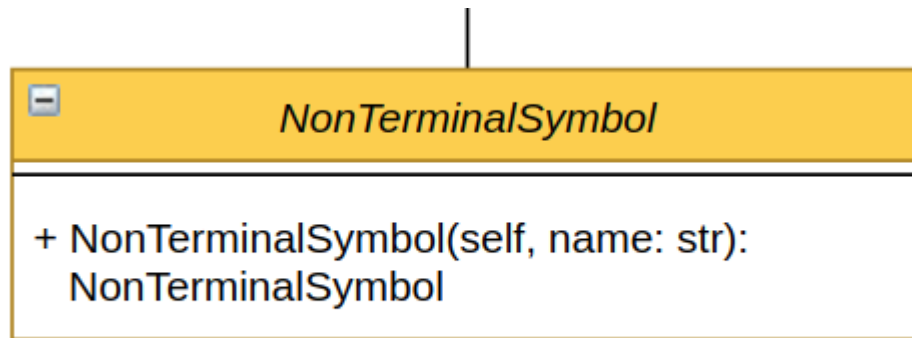


Figure 26. The *NonTerminalSymbol* class.

The *NonTerminalSymbol* class is a child class to the *GrammarSymbol* class, and its only purpose is to serve as an abstraction for the non-terminal symbols in the grammar. The *name* field is inherited from the *GrammarSymbol* class and is the only information needed from the non-terminal symbols themselves.

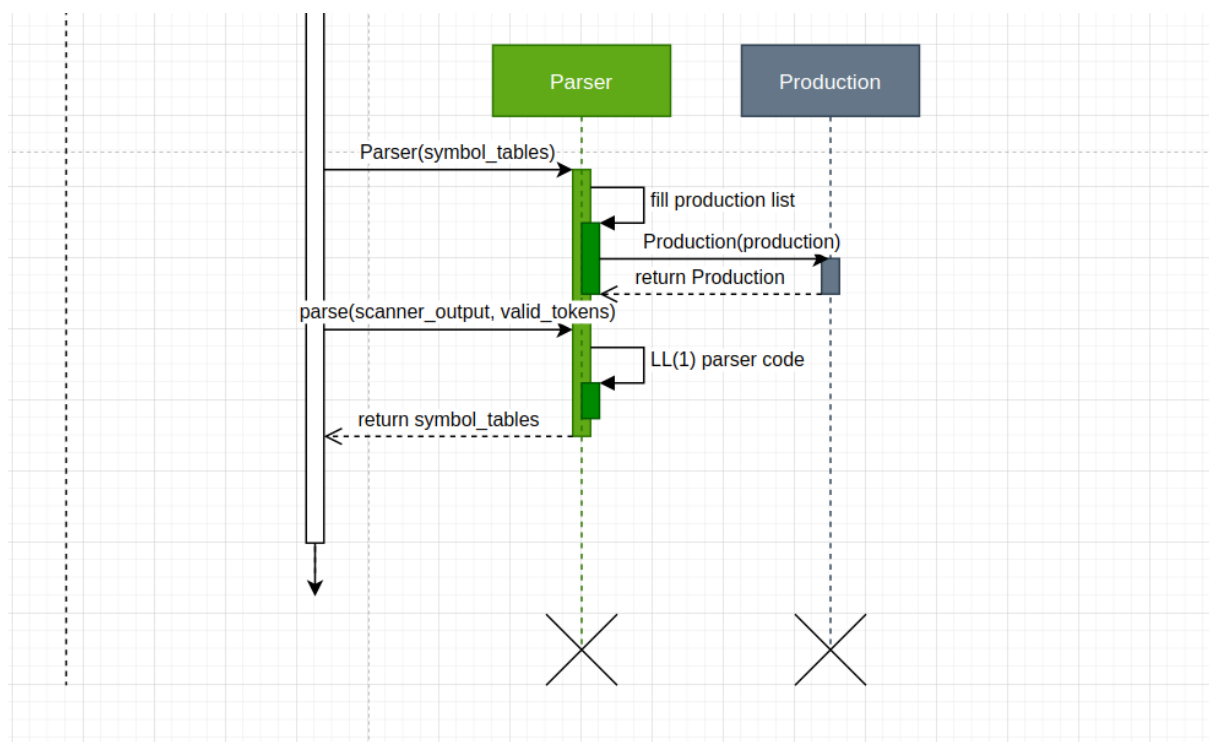


Figure 27. The update to the compilers' sequence diagram.

The section in **Figure 27** is added to the sequence diagram for the compiler, and depicts the behavior of classes involved.

Finally, the following is the pseudocode that was used for the implementation of the LL(1) parsing algorithm for the Syntactic Analyzer for the C Minus language:

```

1) Get first token from the scanner:
   token = nextToken().

2) Push $ into stack.           // EOF into Stack.
   Push S into stack.           // Start non-terminal
                                   // symbol into Stack.

3) while ( TopStack ≠ $ ){      // stack is not empty.
   if ( TopStack = token ) {    // There is a match.
       pop();                   // Take the symbol from top of stack.
       token = nextToken();     // Get following token.
   }
   else if ( TopStack ∈ T ) error ();           // TopStack ≠ token.
   else if ( Table[TopStack, token] = error ) error (); // TopStack ∈ V.
   else if ( Table[TopStack, token] = X → Y1 Y2 ... Yn ){
       pop();                   // Take out the non-terminal from top of stack.
       push(Yn Yn-1 Yn-2 ... Y1); // Push RHS of X inverse order so that.
                                   // Y1 is at top of stack.
                                   // DO NOT push IF X → ε, only pop.
   }
} // end while.
if ( TopStack = $ ∧ token = $ )
    Syntax_Analysis_OK
else
    error ();

```

69

Figure 28. The pseudocode for the LL(1) parsing algorithm. [2]

Implementation

The following images will illustrate how the updates to the compiler were implemented:

```

1  """
2  GrammarSymbol Class
3  author: Victor Emmanuel Guerra Aguado
4
5  Description:
6  |   The GrammarSymbol class provides an abstraction of the grammar symbols present in
7  |   a Context Free Grammar.
8  Attributes:
9  |   + name: str;
10 |   |   The name of the GrammarSymbol object
11 """
12 class GrammarSymbol:
13     def __init__(self, name: str):
14         self.name = name
15
16     def __eq__(self, o: object):
17         if isinstance(o, GrammarSymbol):
18             return self.name == o.name
19         return False
20
21     """
22     GrammarSymbol toStr method
23
24     Description:
25     |   Provides a way to transform the GrammarSymbol Object into a string of characters
26     |   that can be printed onto the screen for visualization or written into a
27     |   file for storage.
28
29     toStr(self) -> str
30     """
31     def toStr(self) -> str:
32         return f"({self.name})"
33
34 class NonTerminalSymbol(GrammarSymbol):
35     def __init__(self, name: str):
36         self.name = name

```

Figure 29. The definition of classes: Grammar Symbol and NonTerminalSymbol.

```

1  from grammarSymbol import GrammarSymbol
2
3  """
4  Production Class
5  author: Victor Emmanuel Guerra Aguado
6
7  Description:
8  |   The Production Class provides a structure for the entries made into the list
9  |   of Productions of the language grammar.
10
11  Attributes:
12  |   + lhs: GrammarSymbol;
13  |   |   The symbol at the left-hand-side of the production.
14  |   + rhs: [GrammarSymbol];
15  |   |   The list of symbols that appear at the right-hand-side of the production.
16  """
17  class Production:
18  |   def __init__(self, lhs: GrammarSymbol, rhs: [GrammarSymbol]):
19  |       self.lhs = lhs
20  |       self.rhs = rhs
21

```

Figure 30. The definition of Production class.

```

1  from grammarSymbol import GrammarSymbol
2  """
3  Token Class
4  author: Victor Emmanuel Guerra Aguado
5
6  Description:
7      The Token Class provides a value object for the tokens created in the
8      Lexical Analysis phase of a compiler.
9
10 Attributes:
11     + name: str;
12     | The name of the token.
13     + line: int;
14     | The line in which this token was detected during the lexical analysis
15     | phase.
16     + keyword_id: int;
17     | The number that tells what type of valid word for the language the token
18     | represents.
19     + symboltable_ref: int;
20     | The ID number that indicates an entry to a Symbol Table, if not an
21     | IDENTIFIER or NUMERICAL_CONST, it defaults to -1.
22
23 """
24 class Token(GrammarSymbol):
25     """
26     Token Class Constructor
27
28     __init__(self, kwid: int, st_ref: int) -> Token
29     """
30     def __init__(self, name: str, line: int, kwid: int, st_ref: int = None):
31         self.name = name
32         self.line = line
33         self.keyword_id = kwid
34         self.symboltable_ref = st_ref
35

```

Figure 31. The update to the Token class.


```

49      """
50      Compiler compile method
51
52      Description:
53          Executes all of the main components' methods to produce an output as
54          indicated by the architecture diagram. Said output is usually passed
55          as an argument as input to the next method.
56
57      compile(self) -> None
58      """
59      def compile(self):
60
61          # Lexical Analysis Phase
62          self.scanner.tokenizeInput()
63          token_stream = self.scanner.token_stream
64          self.symbol_tables = self.scanner.symbol_tables
65
66          token_stream.append(GrammarSymbol("$"))
67          self.parser = Parser(self.symbol_tables)
68          self.parser.parse(token_stream, self.scanner.dfa.valid_tokens)

```

Figure 32. The update to the compile method of class Compiler.

The rest of the code is in file: Parser.zip.

Testing

Test cases were developed around the requirements established during the analysis phase of development, namely:

1. Syntactic Errors that should be detected.
2. Semantic Errors that should be detected.
3. Additions to the Symbol Tables.

CMCTC-08

Description/Purpose:

To verify that the parser adequately parses a correct structure.

Test Script:

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-08 test input source code.
3. Provide the expected results to the test function.
4. Run the compiler.
5. Compare the result with the expected result.

Test Input:

```
“int x[10];
```

```
int miniloc(int a[], int low, int high){
    int i; int x; int k;

    k = low;
    x = a[low];
    i = low + 1;
    while(i < high){
        if(a[i] < x){
            x = a[i];
            k = i;
        }
        i = i + 1;
    }
    return k;
} /* END of miniloc() */
```

```
void sort(int a[], int low, int high){
    int i; int k;

    i = low;
    while(i < high - 1){
        int t;
        /* minloc / I */
        k = miniloc(a,i,high);
        t = a[k];
        a[k] = a[i];
        a[i] = t;
        i = i + 1;
    }
} /* END of sort() */
```

```

void main(void){
    int i;
    i = 0;

    while(i<10){
        input x[i];
        i = i + 1;
    }

    sort(x,0,10);
    i = 0;
    while (i<10){
        output x[i];
        i = i + 1;
    }
} /* END of main() 戰勝利 */

```

Expected Result:

“IDENTIFIER

| ID | VALUE | isVar|isFunc|dataType|noArgs|isGlobal|isLocal

0 | x | True | False | int | 0 | True | False

1 | miniloc | False | True | int | 3 | True | False

2 | a | True | False | int | 0 | False | True

3 | low | True | False | int | 0 | False | True

4 | high | True | False | int | 0 | False | True

5 | i | True | False | int | 0 | False | True

6 | k | True | False | int | 0 | False | True

7 | sort | False | True | void | 3 | True | False

8 | t | True | False | int | 0 | False | True

9 | main | False | True | void | 0 | True | False”

Actual Result:

IDENTIFIER

| ID | VALUE | isVar|isFunc|dataType|noArgs|isGlobal|isLocal

0 | x | True | False | int | 0 | True | False

1 | miniloc | False | True | int | 3 | True | False

2 | a | True | False | int | 0 | False | True
 3 | low | True | False | int | 0 | False | True
 4 | high | True | False | int | 0 | False | True
 5 | i | True | False | int | 0 | False | True
 6 | k | True | False | int | 0 | False | True
 7 | sort | False | True | void | 3 | True | False
 8 | t | True | False | int | 0 | False | True
 9 | main | False | True | void | 0 | True | False

Evaluation: PASS

CMCTC-09

Description/Purpose:

To verify that the parser adequately recognizes and raises an error for: missing main declaration.

Test Script:

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-09 test input source code.
3. Provide the expected results to the test function.
4. Run the compiler.
5. Compare the result with the expected result.

Test Input:

```
“void sort(void) {
    int i;
}”
```

Expected Result:

No main function was declared.

Actual Result:

Evaluation: PASS

CMCTC-10

Description/Purpose:

To verify that the parser adequately recognizes and raises an error for: main function is not the last declaration.

Test Script:

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-10 test input source code.
3. Provide the expected results to the test function.
4. Run the compiler.
5. Compare the result with the expected result.

Test Input:

```
“void main(void) {
int i;
}
void sort(void) {
    int i;
}”
```

Expected Result:

main function is not last declaration.

Actual Result:

main function is not last declaration.

Evaluation: PASS

CMCTC-11

Description/Purpose:

To verify that the parser adequately recognizes and raises an error for: main function receives more than: 0 arguments.

Test Script:

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-11 test input source code.

3. Provide the expected results to the test function.
4. Run the compiler.
5. Compare the result with the expected result.

Test Input:

```
“void main(int i, int a[]) {
int j;
}
”
```

Expected Result:

main function receives more than: 0 arguments.

Actual Result:

main function receives more than: 0 arguments.

Evaluation: **PASS**

CMCTC-12

Description/Purpose:

To verify that the parser adequately recognizes and raises an error for: a variable has been assigned a void type.

Test Script:

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-12 test input source code.
3. Provide the expected results to the test function.
4. Run the compiler.
5. Compare the result with the expected result.

Test Input:

```
“void main(void) {
void i;
}
”
```

Expected Result:

main function receives more than: 0 arguments.

Actual Result:

main function receives more than: 0 arguments.

Evaluation: **PASS**

CMCTC-13

Description/Purpose:

To verify that the parser adequately recognizes and raises an error for: a variable or function has been used before declaration.

Test Script:

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-13 test input source code.
3. Provide the expected results to the test function.
4. Run the compiler.
5. Compare the result with the expected result.

Test Input:

```
“void main(void) {
sort();
}
”
```

Expected Result:

Use of Identifier: sort before declaration. In line: 2

Actual Result:

Use of Identifier: sort before declaration. In line: 2

Evaluation: **PASS**

CMCTC-14

Description/Purpose:

To verify that the parser adequately recognizes and raises an error for: unexpected characters and tokens.

Test Script:

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-14 test input source code.
3. Provide the expected results to the test function.
4. Run the compiler.
5. Compare the result with the expected result.

Test Input:

```
“void main void){
    int i;
    i = 0;

    }
”
```

Expected Result:

Expected a: (, received: void on line 1

Actual Result:

Use of Identifier: sort before declaration. In line: 2

Evaluation: PASS

References

1. R. Castelló, Class Lecture, Topic: “Chapter 3 - Syntax Analysis - Part 1.” TC3048, School of Engineering and Science, ITESM, Chihuahua, Chih, April, 2022.
2. R. Castelló, Class Lecture, Topic: “Chapter 4 - Syntax Analysis - Part 2.” TC3048, School of Engineering and Science, ITESM, Chihuahua, Chih, April, 2022.
3. Thomas Hamilton, “How to Write Test Cases: Sample Template with Examples”, *Guru99*, 2 April 2022 [journal on-line]; available from <https://www.guru99.com/test-case.html>; Internet; accessed 18 April 2022.