



ECOLE
POLYTECHNIQUE
DE BRUXELLES

2023-2024

INFOF403 - Introduction to Language Theory and Compilation

PascalMaisPresque Parser

Ferreira Brenno

Mutkowski Philippe

Gilles GEERAERTS

2023



Table des matières

1	Parser	1
1.1	Introduction	1
1.2	Grammar Transformation	1
1.2.1	Unproductive and Unreachable Rules	1
1.2.2	Ambiguity, Priority and Associativity	1
1.2.3	Left-Recursion	3
1.2.4	Note	3
1.3	Action Table	3
1.3.1	Python Script	3
1.3.2	Terminals, Variables and Grammar	3
1.3.3	First Set Computation Algorithm	4
1.3.4	Follow Set	5
1.3.5	Action Table	6
1.4	Parser	7
1.4.1	Recursive Descent LL(1) Parser	7
1.4.2	Changes in LexicalUnit.java	10
1.4.3	Main Class	10
1.5	Tests	11
1.5.1	All Lexical Units test	11
1.5.2	Euclid	11
1.5.3	ExprArith	11
1.5.4	Fibonacci	11
1.5.5	Comments	11

.1	PascalMaisPresque Grammar	12
----	-------------------------------------	----

1.1 Introduction

This report will talk about the implementation of a parser in java for an imaginary language called PascalMaisPresque. The goal of the second part of this project was changing the original grammar into something usable, writing the corresponding action table and finally writing the parser to checks if the code respects the language. The lexer was already build in the first part of this project, and we choose to keep our own instead of the one provided on the U.V.. A minor flaw was discovered. It's talked about in the parser section.

1.2 Grammar Transformation

First, the grammar had to be changed :

1. To remove unproductive and unreachable variables
2. To avoid ambiguity (and respect priority and associativity of operators)
3. To remove left-recursion

1.2.1 Unproductive and Unreachable Rules

At the original grammar, there was no unproductive and unreachable rules except for the $\langle \text{For} \rangle$ variable at rule 9 which is unproductive. Indeed, in the grammar there is no other rule $\langle \text{For} \rangle \Rightarrow_G^* w$.

1.2.2 Ambiguity, Priority and Associativity

In order to solve ambiguity in the grammar and apply priority and associativity principles, we have to take a look at rule 14-22, $\langle \text{ExprArith} \rangle$, and rule 25-29, $\langle \text{Cond} \rangle$.

<ExprArith> : Rules 14-22

First of all, we consider that an arithmetic expression can produce either an $\langle \text{Expr} \rangle \pm \langle \text{Prod} \rangle$, or just $\langle \text{Prod} \rangle$. This allows to set priority of $*$, $/$ over $+$, $-$. Then, we can produce the multiplication and division, or resolve the variables to terminals by producing an atom. After this steps, the priority is set, the last step was to delete left recursion. This was made by adding the "prime" variables.

14	$\langle \text{ExprArith} \rangle$	$\rightarrow [\text{VarName}] := \langle \text{ExprArith} \rangle$
15		$\rightarrow [\text{Number}]$
16		$\rightarrow (\langle \text{ExprArith} \rangle)$
17		$\rightarrow - \langle \text{ExprArith} \rangle$
18		$\rightarrow \langle \text{ExprArith} \rangle \langle \text{Op} \rangle \langle \text{ExprArith} \rangle$
19	$\langle \text{Op} \rangle$	$\rightarrow +$
20		$\rightarrow -$
21		$\rightarrow *$
22		$\rightarrow /$

FIGURE 1.1 – Version original des règles 14-22.

$\langle \text{ExprArith} \rangle$	$\rightarrow \langle \text{Exp} \rangle + \langle \text{Prod} \rangle$
	$\rightarrow \langle \text{Exp} \rangle - \langle \text{Prod} \rangle$
	$\rightarrow \langle \text{Prod} \rangle$
$\langle \text{Prod} \rangle$	$\rightarrow \langle \text{Prod} \rangle * \langle \text{Atom} \rangle$
	$\rightarrow \langle \text{Prod} \rangle / \langle \text{Atom} \rangle$
	$\rightarrow \langle \text{Atom} \rangle$
$\langle \text{Atom} \rangle$	$\rightarrow [\text{VarName}]$
	$\rightarrow [\text{Number}]$
	$\rightarrow (\langle \text{ExprArith} \rangle)$
	$\rightarrow - \text{Atom}$

FIGURE 1.2 – Ambiguity lifted, priority and associativity set.

14	$\langle \text{ExprArith} \rangle$	$\rightarrow \langle \text{Prod} \rangle \langle \text{ExprArith}' \rangle$
15	$\langle \text{ExprArith}' \rangle$	$\rightarrow + \langle \text{Prod} \rangle \langle \text{ExprArith}' \rangle$
16		$\rightarrow - \langle \text{Prod} \rangle \langle \text{ExprArith}' \rangle$
17		$\rightarrow \epsilon$
18	$\langle \text{Prod} \rangle$	$\rightarrow \langle \text{Atom} \rangle \langle \text{Prod}' \rangle$
19	$\langle \text{Prod}' \rangle$	$\rightarrow * \langle \text{Atom} \rangle \langle \text{Prod}' \rangle$
20		$\rightarrow / \langle \text{Atom} \rangle \langle \text{Prod}' \rangle$
21		$\rightarrow \epsilon$
22	$\langle \text{Atom} \rangle$	$\rightarrow [\text{VarName}]$
23		$\rightarrow [\text{Number}]$
24		$\rightarrow (\langle \text{ExprArith} \rangle)$
25		$\rightarrow - \langle \text{Atom} \rangle$

FIGURE 1.3 – Left-recursion lifted.

One can also notice the definition of the unary $-$. The rule 25 on the last image can turn any terminal or expression between parenthesis into a negative one. It is also put at the lowest level of recursion, very close to a terminal, to ensure that it has the maximum priority, over all other operators.

<Cond> : Rules 25-29

For $\langle \text{Cond} \rangle$ variable, the task was almost the same. We have to ensure that 'and' has priority over 'or'. This was made by producing first an 'or' then 'and'. The 'or' becomes the main operation, thus giving priority to 'and'.

25	$\langle \text{Cond} \rangle$	$\rightarrow \langle \text{Cond} \rangle \text{ and } \langle \text{Cond} \rangle$
26		$\rightarrow \langle \text{Cond} \rangle \text{ or } \langle \text{Cond} \rangle$
27		$\rightarrow \{ \langle \text{Cond} \rangle \}$
28		$\rightarrow \langle \text{SimpleCond} \rangle$
29	$\langle \text{SimpleCond} \rangle$	$\rightarrow \langle \text{ExprArith} \rangle \langle \text{Comp} \rangle \langle \text{ExprArith} \rangle$

FIGURE 1.4 – Original version of rules 25-29.

29	$\langle \text{Cond} \rangle$	$\rightarrow \langle \text{And} \rangle \langle \text{Cond}' \rangle$
30	$\langle \text{Cond}' \rangle$	$\rightarrow \text{or} \langle \text{And} \rangle \langle \text{Cond}' \rangle$
31		$\rightarrow \epsilon$
32	$\langle \text{And} \rangle$	$\rightarrow \langle \text{CondAtom} \rangle \langle \text{And}' \rangle$
33	$\langle \text{And}' \rangle$	$\rightarrow \text{and} \langle \text{CondAtom} \rangle \langle \text{And}' \rangle$
34		$\rightarrow \epsilon$
35	$\langle \text{CondAtom} \rangle$	$\rightarrow \{ \langle \text{Cond} \rangle \}$
36		$\rightarrow \langle \text{ExprArith} \rangle \langle \text{Comp} \rangle \langle \text{ExprArith} \rangle$

FIGURE 1.5 – Ambiguity lifted, priority and associativity set.

1.2.3 Left-Recursion

In addition to rules `<ExprArith>` and `<Cond>`, we had to modify the rules `<InstList>`, `<If>` because they contain a left-recursion to their definitions. This is done by factoring the production rules, adding a tail to each expression.

4	<code><InstList></code>	\rightarrow	<code><Instruction><InstTail></code>
5	<code><InstTail></code>	\rightarrow	<code>...<InstList></code>
6		\rightarrow	ϵ

FIGURE 1.6 – Modification de la règle `InstList`.

26	<code><If></code>	\rightarrow	<code>if<Cond>then<Instruction>else<ElseTail></code>
27	<code><ElseTail></code>	\rightarrow	<code><Instruction></code>
28		\rightarrow	ϵ

FIGURE 1.7 – Modification de la règle `If`

The whole modified grammar can be found at the end of the report.

1.2.4 Note

The production rules of `<Instruction>` could directly produces the correct terminals, deleting a few rules (like 39, 40, 41). The reason of this choice is because it made the code cleaner and easier to understand.

1.3 Action Table

1.3.1 Python Script

At first, a draft of an action table has been written by hand and constructed by checking the first and follows on the fly. The calculation of First and Follow sets turned out to be very time-consuming and not reliable if hand-made. With the idea to speed up the process, we wrote a python script to compute the First and Follow sets and then, based in those information, construct the action table. By doing so, we could allocate more time in double checking the grammar and its rules. With the hand-made action table already pretty complete, we made sure to check discrepancies and made our script reliable.

1.3.2 Terminals, Variables and Grammar

The terminals and variables were encoded as Enum elements in order to ease the identification process.

The grammar was hard-coded inside the script, in the form of a dictionary. The keys were the left-handside variables and the values were the list of rules.

The algorithm also includes detailed debugging logs to aid in understanding its execution. The logging provides insights into the decisions made at each step of the computation. We invite you to check the python script and its output files located in the folder more/playground in order to see its code.

```

55 class Terminals(Enum):
56     VARNAME = "[VarName]",
57     NUMBER = "[Number]",
58     BEGIN = "begin",
59     END = "end",
60     DOTS = "...",
61     ASSIGN = ":",
62     LPAREN = "(",
63     RPAREN = ")",
64     MINUS = "-",
65     PLUS = "+",
66     TIMES = "*",
67     DIVIDE = "/",
68     IF = "if",
69     THEN = "then",
70     ELSE = "else",
71     AND = "and",
72     OR = "or",
73     LBRACK = "{",
74     RBRACK = "}",
75     EQUAL = "=",
76     SMALLER = "<",
77     WHILE = "while",
78     DO = "do",
79     PRINT = "print",
80     READ = "read",
81     EPSILON = "eps",
82     EOF = "$"

```

FIGURE 1.8 – Terminal Enum Class.

```

29 class Variables(Enum):
30     PROGRAM = "<Program>",
31     CODE = "<Code>",
32     INSTLIST = "<InstList>",
33     INSTTAIL = "<InstTail>",
34     INSTRUCTION = "<Instruction>",
35     ASSIGN = "<Assign>",
36     EXPRARITH = "<ExprArith>",
37     EXPRARITHPRIME = "<ExprArith'>",
38     PROD = "<Prod>",
39     PRODPRIME = "<Prod'>",
40     ATOM = "<Atom>",
41     IF = "<If>",
42     ELSETAIL = "<ElseTail>",
43     COND = "<Cond>",
44     CONDPRIME = "<Cond'>",
45     AND = "<And>",
46     ANDPRIME = "<And'>",
47     CONDATOM = "<CondAtom>",
48     COMP = "<Comp>",
49     WHILE = "<While>",
50     PRINT = "<Print>",
51     READ = "<Read>",

```

FIGURE 1.9 – Variable Enum Class.

```

86 grammar = {
87     Variables.PROGRAM: [
88         [Terminals.BEGIN, Variables.CODE, Terminals.END]
89     ],
90     Variables.CODE: [
91         [Variables.INSTLIST],
92         [Terminals.EPSILON]
93     ],
94     Variables.INSTLIST: [
95         [Variables.INSTRUCTION, Variables.INSTTAIL]
96     ],
97     Variables.INSTTAIL: [
98         [Terminals.DOTS, Variables.INSTLIST],
99         [Terminals.EPSILON]
100     ],
101     Variables.INSTRUCTION: [
102         [Variables.ASSIGN],
103         [Variables.IF],
104         [Variables.WHILE],
105         [Variables.PRINT],
106         [Variables.READ],
107         [Terminals.BECL, Variables.INSTLIST, Terminals.END]
108     ],
109     Variables.ASSIGN: [
110         [Terminals.VARNAME, Terminals.ASSIGN, Variables.EXPRARITH]
111     ],
112     Variables.EXPRARITH: [
113         [Variables.PROD, Variables.EXPRARITHPRIME]
114     ],
115     Variables.EXPRARITHPRIME: [
116         [Terminals.PLUS, Variables.PROD, Variables.EXPRARITHPRIME],
117         [Terminals.MINUS, Variables.PROD, Variables.EXPRARITHPRIME],
118         [Terminals.EPSILON]
119     ],

```

FIGURE 1.10 – Grammar.

1.3.3 First Set Computation Algorithm

The algorithm implemented in `compute_first(variable)` of the script is derived from the pseudo-code presented in the lecture notes. It begins by checking whether the first set for a given variable has already been computed. This step is crucial due to the recursive nature of the function. If the first set for the variable is not yet computed, the algorithm initializes and applies the following rules :

1. Terminal Handling :

- If the variable is a terminal, it is included in the first set. It is noteworthy that a terminal can have a first set containing itself. While this might seem non-sensical, it plays a crucial role in the recursive functionality of the function.

2. Production Rule Iteration :

- The algorithm iterates through each production rule of the variable, performing the following steps :

- (a) If the first symbol in the production rule is a terminal, it is added to the first set.
- (b) If the first symbol is a variable, the first set is updated by taking the union with the first set of the respective variable.
- (c) Whenever ϵ is present in the first set of the variable, the first set of the subsequent symbol needs to be computed.

The computed first set is then returned.

Variable	FIRST
PROGRAM	{Terminals BEGIN: 'begin';}
CODE	{Terminals READ: 'read'; <Terminals EPSILON: 'eps'; <Terminals PRINT: 'print'; <Terminals WHILE: 'while'; <Terminals VARNAME: {[VarName]}; <Terminals BEGIN: 'begin'; <Terminals IF: 'if';}
INSTLIST	{Terminals READ: 'read'; <Terminals PRINT: 'print'; <Terminals WHILE: 'while'; <Terminals VARNAME: {[VarName]}; <Terminals BEGIN: 'begin'; <Terminals IF: 'if';}
INSTTAIL	{Terminals READ: 'read'; <Terminals DOTS: '.';}
INSTRUCTION	{Terminals READ: 'read'; <Terminals PRINT: 'print'; <Terminals WHILE: 'while'; <Terminals VARNAME: {[VarName]}; <Terminals BEGIN: 'begin'; <Terminals IF: 'if';}
ASSIGN	{Terminals VARNAME: {[VarName]};}
EXPRARITH	{Terminals NUMBER: {[Number]}; <Terminals MINUS: '-'; <Terminals LPAREN: '('; <Terminals VARNAME: {[VarName]};}
EXPRARITHPRIME	{Terminals PLUS: '+'; <Terminals EPSILON: 'eps'; <Terminals MINUS: '-';}
PROD	{Terminals NUMBER: {[Number]}; <Terminals MINUS: '-'; <Terminals LPAREN: '('; <Terminals VARNAME: {[VarName]};}
PRODPRIME	{Terminals DIVIDE: '/'; <Terminals EPSILON: 'eps'; <Terminals TIMES: '*';}
ATOM	{Terminals NUMBER: {[Number]}; <Terminals MINUS: '-'; <Terminals LPAREN: '('; <Terminals VARNAME: {[VarName]};}
IF	{Terminals IF: 'if';}
ELSETAIL	{Terminals READ: 'read'; <Terminals EPSILON: 'eps'; <Terminals PRINT: 'print'; <Terminals WHILE: 'while'; <Terminals VARNAME: {[VarName]}; <Terminals BEGIN: 'begin'; <Terminals IF: 'if';}
COND	{Terminals LBRACK: '['; <Terminals NUMBER: {[Number]}; <Terminals VARNAME: {[VarName]}; <Terminals LPAREN: '('; <Terminals MINUS: '-';}
CONDPRIME	{Terminals OR: 'or'; <Terminals EPSILON: 'eps';}
AND	{Terminals LBRACK: '['; <Terminals NUMBER: {[Number]}; <Terminals VARNAME: {[VarName]}; <Terminals LPAREN: '('; <Terminals MINUS: '-';}
ANDPRIME	{Terminals EPSILON: 'eps'; <Terminals AND: 'and';}
CONDATOM	{Terminals LBRACK: '['; <Terminals NUMBER: {[Number]}; <Terminals VARNAME: {[VarName]}; <Terminals LPAREN: '('; <Terminals MINUS: '-';}
COMP	{Terminals SMALLER: '<'; <Terminals EQUAL: '=';}
WHILE	{Terminals WHILE: 'while';}
PRINT	{Terminals PRINT: 'print';}
READ	{Terminals READ: 'read';}

FIGURE 1.11 – First Set Result.

1.3.4 Follow Set

Algorithm

The `compute_follow(variable)` function calculates the follow set of a given variable in the grammar while preventing infinite recursion by tracking computed follow sets. The algorithm follows these steps :

1. Base Case Check :

- Verify if the follow set for the variable is already computed. If so, return the computed follow set.

2. Initialization :

- Initialize an empty set for the follow set of the current variable.

3. Start Symbol Check :

- If the variable is the start symbol (`Variables.PROGRAM`), add `Terminals.EOF` to its follow set.

4. Iterate Over Grammar Productions :

- For each variable in the grammar, iterate through its production rules.

5. Production Rule Analysis :

— For each production rule containing the current variable :

(a) **Case 1 : Variable at the End :**

— If the variable is at the end of the production, and it's not the same as the current variable, update its follow set with the computed follow set of the variable where the production rule is found.

(b) **Case 2 : Variable not at the End :**

— If the variable is not at the end of the production :

— If the next symbol is a terminal, add it to the follow set.

— If the next symbol is a variable :

i. If epsilon is in the first set of the next symbol :

— Update the follow set by excluding epsilon from the first set of the next symbol.

— Consider applying the rule with epsilon in the next iteration.

— Update the follow set with the computed follow set of the next symbol.

— If the next symbol is the last in the production, update the follow set with the computed follow set of the variable where the production rule is found.

ii. If epsilon is not in the first set, update the follow set with the first set of the next symbol.

6. **Logging :**

— Log debugging information for each step of the computation.

7. **Return Result :**

— Return the computed follow set for the variable.

This algorithm ensures that the follow sets are computed accurately and efficiently for the given grammar.

1.3.5 Action Table

Algorithm

— **Initialization :**

— Initialize an empty action table, associating each variable and terminal with a None value.

Variable	FOLLOW
1 PROGRAM	<Terminals EOF: '\$'>
2 CODE	<Terminals END: 'end'>
4 INSTLIST	<Terminals END: 'end'>
5 INSTAIL	<Terminals END: 'end'>
6 INSTRUCTION	<Terminals ELSE: 'else'>, <Terminals DOTS: '.'>, <Terminals END: 'end'>
7 ASSIGN	<Terminals ELSE: 'else'>, <Terminals DOTS: '.'>, <Terminals END: 'end'>
8 EXPRARITH	<Terminals OR: 'or'>, <Terminals SMALLER: '<'>, <Terminals RPAREN: ')'>, <Terminals ELSE: 'else'>, <Terminals END: 'end'>, <Terminals DO: 'do'>, <Terminals EQUAL: '='>, <Terminals RBRACK: '['>, <Terminals AND: 'and'>, <Terminals DOTS: '.'>
9 EXPRARITHPRIME	<Terminals OR: 'or'>, <Terminals SMALLER: '<'>, <Terminals RPAREN: ')'>, <Terminals ELSE: 'else'>, <Terminals END: 'end'>, <Terminals DO: 'do'>, <Terminals EQUAL: '='>, <Terminals RBRACK: '['>, <Terminals AND: 'and'>, <Terminals DOTS: '.'>
10 PROD	<Terminals OR: 'or'>, <Terminals SMALLER: '<'>, <Terminals PLUS: '+'>, <Terminals RPAREN: ')'>, <Terminals ELSE: 'else'>, <Terminals MINUS: '-'>, <Terminals DO: 'do'>, <Terminals DOTS: '.'>, <Terminals THEN: 'then'>, <Terminals END: 'end'>
11 PRODPRIIME	<Terminals OR: 'or'>, <Terminals SMALLER: '<'>, <Terminals PLUS: '+'>, <Terminals RPAREN: ')'>, <Terminals ELSE: 'else'>, <Terminals END: 'end'>, <Terminals MINUS: '-'>, <Terminals DO: 'do'>, <Terminals EQUAL: '='>, <Terminals RBRACK: '['>
12 ATOM	<Terminals OR: 'or'>, <Terminals SMALLER: '<'>, <Terminals PLUS: '+'>, <Terminals RPAREN: ')'>, <Terminals ELSE: 'else'>, <Terminals MINUS: '-'>, <Terminals DO: 'do'>, <Terminals DOTS: '.'>, <Terminals TIMES: '*'>, <Terminals THEN: 'then'>
13 IF	<Terminals ELSE: 'else'>, <Terminals DOTS: '.'>, <Terminals END: 'end'>
14 ELSETAIL	<Terminals ELSE: 'else'>, <Terminals DOTS: '.'>, <Terminals END: 'end'>
15 COND	<Terminals DO: 'do'>, <Terminals RBRACK: '['>, <Terminals THEN: 'then'>
16 CONDPRIIME	<Terminals DO: 'do'>, <Terminals RBRACK: '['>, <Terminals THEN: 'then'>
17 AND	<Terminals DO: 'do'>, <Terminals RBRACK: '['>, <Terminals OR: 'or'>, <Terminals THEN: 'then'>
18 ANDPRIME	<Terminals DO: 'do'>, <Terminals RBRACK: '['>, <Terminals OR: 'or'>, <Terminals THEN: 'then'>
19 CONDATOM	<Terminals OR: 'or'>, <Terminals DO: 'do'>, <Terminals RBRACK: '['>, <Terminals AND: 'and'>, <Terminals THEN: 'then'>
20 COMP	<Terminals NUMBER: '[Number]'>, <Terminals MINUS: '-'>, <Terminals LPAREN: '('>, <Terminals VARNAME: '[VarName]'>
21 WHILE	<Terminals ELSE: 'else'>, <Terminals DOTS: '.'>, <Terminals END: 'end'>
22 PRINT	<Terminals ELSE: 'else'>, <Terminals DOTS: '.'>, <Terminals END: 'end'>
23 READ	<Terminals ELSE: 'else'>, <Terminals DOTS: '.'>, <Terminals END: 'end'>

FIGURE 1.12 – Follow Set Result.

— Action Table Filling :

— Iterate over each variable in the grammar.

— For each production of the variable :

1. Compute the first set of the first symbol in the production.

2. For each terminal in the first set :

— If the terminal is not epsilon, assign the absolute number of the rule to the corresponding cell in the action table.

3. If epsilon is in the first set :

— Compute the follow set of the current variable.

— For each terminal in the follow set, assign the absolute number of the rule to the corresponding cell in the action table.

— Production Number Update :

— Increment the production number for each processed production.

— Return Result :

— Return the computed action table.

This algorithm constructs an action table by determining the rule numbers for each variable and terminal combination. It efficiently handles both terminals in the first set and terminals in the follow set, ensuring a comprehensive and accurate action table for the grammar.

1.4 Parser

1.4.1 Recursive Descent LL(1) Parser

The design of the parser was largely inspired by the lecture note's example of a recursive descent LL(1) parser in python, shown in 1.14.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB
1		VARNAME	NUMBER	BEGIN	END	DOTS	ASSIGN	LPAREN	RPAREN	MINUS	PLUS	TIMES	DIVIDE	IF	THEN	ELSE	AND	OR	LBRACK	RBRACK	EQUAL	SMALLER	WHILE	DO	PRINT	READ	EPSILON	EOF
2	PROGRAM			1																								
3	CODE	2			2	3							2									2		2	2			
4	INSTLIST	4			4								4									4		4	4			
5	INSTAIL					6	5																					
6	INSTRUCTION	7			12								8										9		10	11		
7	ASSIGN	13																										
8	EXPANRITH	14	14					14				14																
9	EXPANRITHPRIME					17	17					17	15				17	17	17	17			17	17	17		17	
10	PROD	18	18									18																
11	PRODPRIME					21	21					21	21	19	20		21	21	21	21			21	21	21		21	
12	ATOM	22	23									24			25													
13	IF															26												
14	ELSETAIL	27			27	28	28									27		28							27		27	27
15	COND	29	29						29			29								29								
16	CONDPRIME																31			30		31				31		
17	AND	32	32						32			32								32								
18	ANDPRIME																											
19	CONDATOM	36	36						36			36					34		33	34			34				34	
20	COMP																			35								
21	WHILE																					37	38					
22	PRINT																											
23	READ																									40		
																											41	

FIGURE 1.13 – Action Table Result.

```

1 def S():
2     n = get_next_character() # This is a look-ahead
3
4     # If the look-ahead is in First(A b Follow(S))
5     if n == 'a1' or n == 'a2' or ... or n == 'an' :
6         read_next_character() # Discard the look-ahead
7         r = A()
8         if (r): return False
9         n = read_next_character()
10        if (n == 'b'): return True
11        else: return False
12
13    # If the look-ahead is in First(B c Follow(S))
14    if n == 'b1' or n == 'b2' or ... or n == 'bn' :
15        read_next_character() # Discard the look-ahead
16        r = B()
17        if (r): return False
18        n = read_next_character()
19        if (n == 'c'): return True
20
21    else: return False
22
23    # Otherwise, the input cannot be valid.
24    return False

```

FIGURE 1.14 – Recursive Descent LL(1) Parser in Python as described in the lecture notes.

Parser attributes The most important attributes of the Parser class are the lexer, the currentToken, usedRules array and the root ParseTree.

- The lexer attribute holds an instance of the LexicalAnalyzer, responsible for scanning the input file and segmenting it into tokens. During the second phase of the project, we opted to utilize the lexer developed in part 1. However, this decision led us to identify a minor flaw in its implementation. Specifically, the token DOTS ('...') was absent from the original lexer. Upon detecting this flaw, we addressed it by adding the missing token into the lexer's code.
- The currentToken holds a symbol given by the lexer. This attribute is equivalent to a look-ahead from the theory of PDAs. The choice of which rule to apply is based in this symbol.
- The usedRules arrays contains a sequence of number rules used in order to accept the word. Whenever a rule is applied, its number is append to the end of the array.
- The root ParseTree is just the <Program> ParseTree, but it was names root because it's the beginning of the whole tree.

Parser Methods Each production rule of the grammar was coded as a method of the parser class. Whenever one have to apply rule 13, `<Assign>`, for example, one have to call the method `assign()`. Other important methods are `match()`, `nextToken()` and `syntaxError()` ensures the running of the parser.

- The `program()` method serves as the initial step in the parser. It is called whenever a file needs parsing, as it corresponds to the first rule and represents the start symbol. This method creates the `ParseTree` root and attempts to match the first token, which is `Terminals.BEG`. If the token is successfully matched, the node for `Terminals.BEG` is inserted into the root's children list.

Given that the next token is a `Variable.CODE`, a node is created and inserted into the root children list but the subtree of this node is generated by calling the `code(ParseTree parentTree)` function and passing the code tree as an argument. Consequently, `code(ParseTree parentTree)` is responsible for creating its children nodes and inserting them into their respective children lists. This pattern is repeated whenever the current token is a variable.

Upon the completion of the `code()` execution, it returns the complete parse tree. The `program()` method can then attempt to match the next token, which is `Terminal.END`. If successful, the end node is inserted into the root's children list. After processing `Terminal.END`, the parser should generate the desired output, namely a sequence of applied rules for the derivation. This sequence is precisely maintained by the `usedRules` array.

To conclude the parser's execution, the `program()` method returns the root `ParseTree`. This root can be obtained by the Main Class to subsequently write it to a file.

All production rules were coded in this very same way. First, the current token is matched with some `LexicalUnit` according to the production rule and the node is inserted into the tree. When the token is a variable, the corresponding tree node is created and the method of this variables is responsible to generate the subtree below.

- The `match(Symbol expected)` is the method responsible of checking if the syntax is correct, which means checking if the current token is the token needed to continue the parsing. If it is successful, it reads the token from the input file stream thanks to the method `nextToken()`. If it fails, it calls the `syntaxError(Symbol token, LexicalUnit[] expected)` method.
- The `nextToken()` uses the lexer to obtain the next token of the file. When this function is called, the lexer give the next token and the last token is consumed. It fails if the lexer could not scan the file.

- The `syntaxError(Symbol token, LexicalUnit[] expected)` raises an error to the standard output alerting that the current token is not valid. It shows the line and column of the token and what were the expected tokens that the parser attempt to identify. It gives also the last state of the parser which is the last rule or method executed, but this is more useful for debugging purposes.

Parse Tree By having each method create its children nodes and append them to their respective children lists, then returning their tree alongside the recursive calls of methods, we can construct the parse tree for any input file.

1.4.2 Changes in LexicalUnit.java

With the idea of detecting if the current token is within the grammar, the `LexicalUnit.java` file was updated to include the variables names too. It allowed to compare the type of the token with the expected token using `currentType.getType() == LexicalUnit.AND`, for example. There was also a need to include the new variables generated by the grammar transformation like `<Prod>`, `<ExprArithPrime>`, ...

1.4.3 Main Class

The main class has two commands possible :

- `java -jar ./dist/part2.jar <inputFilePath>`
 - The first one is to only pass as argument the input file path. Thus, the main class calls the method `parse(String filepath)` that initialize an object `Parser` with the file path and calls the method `program()` of it, which is the first rule. When the parser execution is done, the function `parse(String filepath)` returns the parse tree generated by the run. In this case, nothing is done with it
- `java -jar ./dist/part2.jar -wt <LatexFilePath>.tex <inputFilePath>`
 - If the option `-wt <latexFilePath>.tex` is passed in the arguments, the execution is exactly the same, except that the main class write down a file containing the LaTeX of the parse tree.
- In both case, the standard output will contain the sequence of the rules used for this derivation if it is successful. If not, a syntax error will raise and stop the execution. The main execution could be also stopped whenever some other major errors, not related to the parser algorithm, occurs, like wrong command or wrong filepath.

1.5 Tests

1.5.1 All Lexical Units test

This test file has all lexical units written with the correct grammar. This can also be used to check that the parser respects the priority of operations.

1.5.2 Euclid

This was the files provided with the project statement. It was the main file used to debug.

1.5.3 ExprArith

This highlights the priority principle by looking at tree level of operators (* is lower than + then, * has more priority over +).

It highlights also the change of priority in one of the expression by the simple addition of a parenthesis Lastly, it highlights the left associativity when two operators of same priority are together ($2*3/1$, where the * should be lower than the / because it is on the left)

1.5.4 Fibonacci

This is another code that should be runnable.

1.5.5 Comments

Here it highlights that all comments are completely ignored because there are not even a token

.1 PascalMaisPresque Grammar

1	<Program>	→	begin <Code> end
2	<Code>	→	<InstList>
3		→	ϵ
4	<InstList>	→	<Instruction><InstTail>
5	<InstTail>	→	...<InstList>
6		→	ϵ
7	<Instruction>	→	<Assign>
8		→	<If>
9		→	<While>
10		→	<Print>
11		→	<Read>
12		→	begin<InstList>end
13	<Assign>	→	[VarName] :=<ExprArith>
14	<ExprArith>	→	<Prod> <ExprArith>'
15	<ExprArith'>	→	+<Prod><ExprArith'>
16		→	-<Prod><ExprArith'>
17		→	ϵ
18	<Prod>	→	<Atom><Prod'>
19	<Prod'>	→	*<Atom><Prod'>
20		→	/<Atom><Prod'>
21		→	ϵ
22	<Atom>	→	[VarName]
23		→	[Number]
24		→	(<ExprArith>)
25		→	-<Atom>
26	<If>	→	if<Cond>then<Instruction>else<ElseTail>
27	<ElseTail>	→	<Instruction>
28		→	ϵ
29	<Cond>	→	<And><Cond'>
30	<Cond'>	→	or<And><Cond'>
31		→	ϵ
32	<And>	→	<CondAtom><And'>
33	<And'>	→	and<CondAtom><And'>
34		→	ϵ
35	<CondAtom>	→	{<Cond>}
36		→	<ExprArith> <Comp> <ExprArith>
37	<Comp>	→	=
38		→	<
39	<While>	→	while<Cond>do<Instruction>
40	<Print>	→	print([varName])
41	<Read>	→	read([varName])