

MINISTRY OF EDUCATION AND RESEARCH OF REPUBLIC OF MOLDOVA
TECHNICAL UNIVERSITY OF MOLDOVA
FACULTY OF COMPUTERS, INFORMATICS AND MICROELECTRONICS
DEPARTMENT OF SOFTWARE ENGINEERING AND AUTOMATICS

Laboratory Work 6:

Parser & Abstract Syntax Tree.

Course: Formal Languages & Finite Automata

Student: Gusev Roman

Group: FAF-222

Professor: Cojuhari Irina

University Assistant: Crețu Dumitru

Chișinău, 2024

Content

Theory	3
Objectives	3
Implementation	3
Conclusions	20
Bibliography	26

Theory

- **Definitions:**

- **Parser** - The parser is that phase of the compiler which takes a token string as input and with the help of existing grammar, converts it into the corresponding Intermediate Representation(IR). The parser is also known as Syntax Analyzer [1];
- **Abstract Syntax Tree** - Abstract Syntax Tree is a kind of tree representation of the abstract syntactic structure of source code written in a programming language. Each node of the tree denotes a construct occurring in the source code [2].

Objectives

- Get familiar with parsing, what it is and how it can be programmed;
- Get familiar with the concept of AST;
- In addition to what has been done in the 3rd lab work do the following:
 - In case you didn't have a type that denotes the possible types of tokens you need to:
 - * Have a type TokenType (like an enum) that can be used in the lexical analysis to categorize the tokens;
 - * Please use regular expressions to identify the type of the token;
 - Implement the necessary data structures for an AST that could be used for the text you have processed in the 3rd lab work;
 - Implement a simple parser program that could extract the syntactic information from the input text.

Implementation Description

- For the start, I decided to implement the Lexer that was developed in Laboratory Work 3 - Lexer
- For the current Laboratory Work, I decided to implement Parser for the same PBL DSL that me and the team I work with decided to develop. During the work on this Laboratory Work, I decided to implement an example Parser that I found in some guides on YouTube [3]. Specifically, I followed the structure of the presented in the video Lexer, but for the extension I used only my knowledge.
- The Lexer I developed pass through the entire Input symbol by symbol until the End of File, and transform each combination of the Input into Lexemes. The mathematical expressions are not required to be separated by Whitespace and may be written in both, for example "2+2" and "2 + 2", forms, and will be tokenized correctly. For the "words" type Lexemes, they are required to be separated by whitespace. This is done in order to greedily tokenize them and eliminate cases in inputs such as "Formula1", that may be tokenized like:

```
1 FORMULA_TOKEN : Formula
2 NUMBER: 1
```

and tokenize that as a single word, that might give theoretically output as:

```
1 IDENTIFIER: Formula1
```

This was done in order to have a defined structure for the Tokens, that will lead to a whitespace-sensitive Language.

- First thing modified was the Token class themselves. I have several Tokens that I decided to hold in Enum classes, and that have been used in the Lexer to classify them and, therefore, by Parser. I added several new Tokens:

```
1 from enum import Enum
2
3 class OperatorType(Enum):
4     MULTIPLY = "*"
5     DIVISION = "/"
6     ADDITION = "+"
7     SUBTRACTION = "-"
8     POWER = "**"
9
10 class ComparatorType(Enum):
11     EQUALS = "=="
12     NOT_EQUALS = "!="
13     LESS_THAN = "<"
14     GREATER_THAN = ">"
15     LESS_THAN_OR_EQUALS = "<="
16     GREATER_THAN_OR_EQUALS = ">="
17
18 class FileType(Enum):
19     CSV = "csv"
20     TEXT = "txt"
21     JSON = "json"
22     EXCEL = "excel"
```

```

23     CONSOLE = "console"
24
25 class ImageType(Enum):
26     JPG = "jpg"
27     PNG = "png"
28
29 class PlotType(Enum):
30     GRAPH = "graph"
31     BAR = "bar"
32     PIE = "pie"
33     PLOT = "plot"
34     HIST = "hist"
35
36 class VariableType(Enum):
37     FORMULA = "Formula"
38     DATA = "Data"
39     DATASET = "dataset"
40     NAME = "name"
41
42 class ExportToType(Enum):
43     EXPORT_TO_IMAGE = "ExportToImage"
44     EXPORT_TO_FILE = "ExportToFile"
45
46 class VisualizationType(Enum):
47     VISUALIZE_DATA = "VisualData"
48     VISUALIZE_FORMULA = "VisualFormula"

```

- At the same time, I modified the Token class, that will describe what type of Lexeme I have in input. I added 3 methods for representation in the Abstract Syntax Tree, for checking if the Token is of a specific kind and if the Token has a specific string value in it.

```

1 class Token:
2     ...
3     def ast_repr(self, _):

```

```

4         return f"{self.kind}          {self.string}"
5
6     def has(self, *strings):
7         return self.string in strings
8
9     def of(self, *kinds):
10        return self.kind in kinds
11
12    ...

```

- Next Class developed for the Parser implementation was "Node", that will hold the structure of the Abstract Syntax Tree. This class is an Abstract class, that will be inherited by other classes that will represent the Nodes of the AST. The class has several methods that will be used in the subclasses, such as:

- **line**: that will return the line of the Node,
- **nodes**: that will return the Nodes of the current Node,
- **ast_repr**: that will represent the Node in the AST,
- **mark**: that will mark the Nodes in the AST,
- **construct**: that will construct the Node from the Parser.

```

1 from abc import ABC, abstractmethod
2 class Node(ABC):
3     @property
4     def line(self):
5         return self.nodes()[0].line
6
7     @abstractmethod
8     def nodes(self):
9         pass
10
11     def __repr__(self):
12         return self.ast_repr()
13
14     def ast_repr(self, prefix="\t"):
15         ast_string = "("

```

```

16     ast_string += type(self).__name__ + "("
17     string = type(self).__name__
18     nodes = self.nodes()
19
20     for i, node in enumerate(nodes):
21         at_last = (i == len(nodes) - 1)
22         symbol = "          " if at_last else "          "
23         prefix_symbol = "" if at_last else "    "
24
25         node_string = node.ast_repr(f"{prefix}{prefix_symbol}\t"
26                                     ")")
27         string += f"\n{prefix}{symbol} {node_string}"
28
29     ast_string += ")"
30     return string
31
32     def mark(self):
33         for node in self.nodes():
34             node.mark()
35
36     @classmethod
37     @abstractmethod
38     def construct(cls, parser):
39         pass

```

- Next class developed was the class for the Primary Nodes, that will hold the Tokens that are not Binary, i.e., that are represented by a single derivation. This class will hold the Token itself and will have the methods to return the Nodes of the current Node and to represent the Node in the AST.
- Next class - BinaryNode, that will hold the Tokens that are Binary, i.e., that are represented by 3 derivations. This class will hold the Left Node, Operator Node and Right Node. Also, it will have the methods to return the Nodes. Besides that it will have a method to construct the Binary Node from the Parser.

```

1 from abc import ABC, abstractmethod

```

```

2  ...
3  class PrimaryNode(Node, ABC):
4      def __init__(self, token):
5          self.token = token
6
7      def nodes(self):
8          return [self.token]
9
10     def ast_repr(self, prefix="\t"):
11         return f"{self.token.kind}          {self.token.string}"
12  ...
13  class BinaryNode(Node, ABC):
14      def __init__(self, left, op, right):
15          self.left = left
16          self.op = op
17          self.right = right
18
19      def nodes(self):
20          return [self.left, self.op, self.right]
21
22      @classmethod
23      def construct_binary(cls, parser, make, part, ops):
24          node = part.construct(parser)
25
26          while parser.next().has(*ops):
27              op = parser.take()
28              right = part.construct(parser)
29              node = make(node, op, right)
30
31          return node

```

- After that, I developed the class for Parser that will hold the tokens from the Lexer and will construct the AST from them. This class holds the Tokens and the Index of the current Token. It has several methods that will be used in the Parser:

- **expression**: that will return the Expression Node,
- **formula_content** : *that will return the Additive Expression Node*, **next** : *that will return the next Token*,
- **take**: that will take the next Token,
- **expecting_has** : *that will check if the next Token has a specific string value*,
- Also, this class has 2 properties that will return the Expression and Formula Content Nodes, that will be used further in expressions to avoid the circular imports in classes.

```

1  ...
2  class Parser:
3      @property
4      def expression(self):
5          return Expression
6
7      @property
8      def formula_content(self):
9          return Additive_Expression
10
11     def __init__(self):
12         self.tokens = None
13         self.index = 0
14
15     def next(self):
16         return self.tokens[self.index]
17
18     def take(self):
19         token = self.next()
20         self.index += 1
21         return token
22
23     def expecting_has(self, *strings):
24         if self.next().has(*strings):
25             return self.take()
26
27         raise ParserError(self.next(), f"Expected one of: {strings}")

```

```

    ")
28
29     def expecting_of(self, *kinds):
30         if self.next().of(*kinds):
31             return self.take()
32
33         raise ParserError(self.next(), f"Expected one of the type:
           {kinds}")
34
35     def build_ast(self, tokens):
36         self.tokens = tokens
37         node = Expression.construct(self)
38         if self.next().of("EOF"):
39             return node
40
41         raise ParserError(self.next(), f"Unexpected token: {self.
           next()}")

```

- Then, I started to define the classes for each expression in the Language that I developed. I started with the Program_Expression class, that will hold the Commands of the Program. This class will have the methods to return the Nodes of the current Node and to represent the Node in the AST. This class will add new Command Expressions to the list of Commands until the End of File.

```

1     ...
2     class Program_Expression(Node):
3         def __init__(self, commands):
4             self.commands = commands
5
6         def nodes(self):
7             return self.commands
8
9         @classmethod
10        def construct(cls, parser):
11            commands = []

```

```

12
13     while True:
14         if parser.next().of("EOF"):
15             break
16         command = Command_Expression.construct(parser)
17         commands.append(command)
18
19     return Program_Expression(commands)
20     ...

```

- Next Expression developed was the Command_Expression class, that will hold the Commands of the Program. This class has the methods to return the Nodes of the current Node and to represent the Node in the AST. Also, this class may derive into 3 different classes - Variable_Expression, Comment_Expression and Command_Expression.

```

1     ...
2     class Command_Expression(Node):
3         def __init__(self, command, right):
4             self.command = command
5             self.right = right
6
7         def nodes(self):
8             return [self.command, self.right]
9
10        @classmethod
11        def construct(cls, parser):
12            if parser.next().of("COMMENT_LINE", "COMMENT_BLOCK"):
13                return Comment_Expression(parser.take())
14
15            command = Variable_Expression.construct(parser)
16
17            right = parser.expecting_has(";")
18
19            return Command_Expression(command, right)

```

- Here I provide 3 classes that represent Commenting methods - Block and Line comments. They are taking the token that are of Type COMMENT_LINE and COMMENT_BLOCK and construct their derivations in the AST.

```
1  ...
2  class Comment_Expression(Node):
3      def __init__(self, comment):
4          self.comment = comment
5
6      def nodes(self):
7          return [self.comment]
8
9      @classmethod
10     def construct(cls, parser):
11         if parser.next().of("COMMENT_LINE"):
12             return Comment_Line_Expression(parser)
13         else:
14             return Comment_Block_Expression(parser)
15
16 class Comment_Block_Expression(Node):
17     def __init__(self, comment_line):
18         self.comment_line = comment_line
19
20     def nodes(self):
21         return [self.comment_line]
22
23     @classmethod
24     def construct(cls, parser):
25         return Comment_Block_Expression(parser.expecting_of("
26 COMMENT_BLOCK"))
27
28 class Comment_Line_Expression(Node):
29     def __init__(self, comment_line):
30         self.comment_line = comment_line
```

```

30
31     def nodes(self):
32         return [self.comment_line]
33
34     @classmethod
35     def construct(cls, parser):
36         return Comment_Line_Expression(parser.expecting_of("
COMMENT_LINE"))

```

- The next class developed was the Variable_Expression class, that will hold the Variables of the Program. This class derives into 2 different ways - Data declaration and Formula declaration. This class has the methods to return the Nodes of the current Node and to represent the Node in the AST.
- If the keyword is "Formula", then the Expression will be Additive_Expression, otherwise, it will be Read_From_Path_Expression.

```

1  ...
2  class Variable_Expression(Node):
3      def __init__(self, keyword, identifier, op, expression):
4          self.keyword = keyword
5          self.identifier = identifier
6          self.op = op
7          self.expression = expression
8
9      def nodes(self):
10         return [self.keyword, self.identifier, self.op, self.
expression]
11
12     @classmethod
13     def construct(cls, parser):
14         # Expect 'Formula' or 'Data' keyword
15         keyword = parser.expecting_has("Formula", "Data")
16
17         # Expect Identifier (ID)
18         identifier = Identifier.construct(parser)

```

```

19
20     # Expect '=' symbol
21     op = parser.expecting_has("=")
22
23     # Parse Expression
24     if keyword.has("Formula"):
25         expression = Additive_Expression.construct(parser)
26     else:
27         expression = Read_From_Path_Expression.construct(parser
28             )
29
30     return Variable_Expression(keyword, identifier, op,
31         expression)

```

- The next class is Read From Path expression that will hold the Path of the Data that will be read from. It is designed in the same manner as previous classes, that will have the methods to return the Nodes of the current Node and keeping track of the Tokens and their types. Also, this class includes another class - Path_Expression, that will hold the specific path of the Data that will be read from.

```

1 ...
2 class Read_From_Path_Expression(Node):
3     def __init__(self, keyword, left, path, right):
4         self.keyword = keyword
5         self.left = left
6         self.path = path
7         self.right = right
8
9     def nodes(self):
10         return [self.keyword, self.left, self.path, self.right]
11
12     @classmethod
13     def construct(cls, parser):
14         keyword = parser.expecting_of("READ_FROM")
15         left = parser.expecting_of("LPAREN")

```

```

16     path = Path_Expression.construct(parser)
17     right = parser.expecting_of("RPAREN")
18
19     return Read_From_Path_Expression(keyword, left, path, right
20         )
21 ...
22 class Path_Expression(Node):
23     def __init__(self, expression):
24         self.expression = expression
25
26     def nodes(self):
27         return [self.expression]
28
29     @classmethod
30     def construct(cls, parser):
31         expression = parser.expecting_of("PATH")
32         return Path_Expression(expression)

```

- For the Formula declaration, I implemented parse methods for the Additive Expression, that will hold the mathematical addition or Subtraction of the Variables. This class has the methods to return the Nodes of the current Node and to represent the Node in the AST. Then, if the Operator is "+" or "-", then the right Node will be the Additive Expression, otherwise, it will be the Multiplicative Expression. Then, if the Operator is "*" or "/", then the right Node will be the Multiplicative Expression, otherwise, it will be the Exponential Expression. The same structure is persisted across all the classes that represent the mathematical expressions, such as Unary Expressions, Primary Expressions that involves parenthesis and so on.

```

1 ...
2 class Additive_Expression(BinaryNode):
3     @classmethod
4     def construct(cls, parser):
5         return cls.construct_binary(parser, cls,
6             Multiplicative_Expression, ["+", "-"])

```

```

7
8 class Multiplicative_Expression(BinaryNode):
9     @classmethod
10     def construct(cls, parser):
11         return cls.construct_binary(parser, cls,
12                                     Exponential_Expression, ["*", "/"])
13     ...
14 class Exponential_Expression(BinaryNode):
15     @classmethod
16     def construct(cls, parser):
17         node = Unary_Expression.construct(parser)
18
19         if not parser.next().has("**"):
20             return node
21
22         op = parser.take()
23         right = Exponential_Expression.construct(parser)
24         return Exponential_Expression(node, op, right)
25     ...
26 class Unary_Expression(Node):
27     def __init__(self, op, expression):
28         self.op = op
29         self.expression = expression
30
31     def nodes(self):
32         return [self.op, self.expression]
33
34     @classmethod
35     def construct(cls, parser):
36         if parser.next().has("+", "-"):
37             op = parser.take()
38             expression = Unary_Expression.construct(parser)
39             return Unary_Expression(op, expression)

```


39

40

```
return Primary_Expression.construct(parser)
```

- For the Primary Expressions class, I implemented the classes for the Number, Identifier and IdentifierOrNumber, that will be required to be in the Primary Expressions. The Number class will hold the Number Token, the Identifier class will hold the Identifier Token, and last class will be a general if the input is not met the previous 2 classes. This class has the methods to return the Nodes of the current Node.

1

```
...
```

2

```
class Primary_Expression(Node):
```

3

```
    def __init__(self, left, expression, right):
```

4

```
        self.left = left
```

5

```
        self.expression = expression
```

6

```
        self.right = right
```

7

8

```
    def nodes(self):
```

9

```
        return [self.left, self.expression, self.right]
```

10

11

```
    @classmethod
```

12

```
    def construct(cls, parser):
```

13

```
        if not parser.next().has("("):
```

14

```
            if parser.next().of("ID"):
```

15

```
                return Identifier.construct(parser)
```

16

```
            elif parser.next().of("NUMBER"):
```

17

```
                return Number.construct(parser)
```

18

```
            else:
```

19

```
                return IdentifierOrNumber.construct(parser)
```

20

21

```
        left = parser.take()
```

22

```
        expression = parser.formula_content.construct(parser)
```

23

```
        right = parser.expecting_has("(" if left.has("(") else left
            .string)
```

24

```

25         return Primary_Expression(left, expression, right)
26     ...
27 class Number(PrimaryNode):
28     @classmethod
29     def construct(cls, parser):
30         return Number(parser.expecting_of("NUMBER"))
31
32     ...
33 class Identifier(PrimaryNode):
34     @classmethod
35     def construct(cls, parser):
36         return Identifier(parser.expecting_of("ID"))

```

- In the Main class, I decided to modify how User choose between 3 methods of Input and will display the AST of the input:

```

1 parser = Parser()
2 tree = parser.build_ast(tokens)
3 print("TOKENS:")
4 for token in tokens:
5     print(token)
6
7 print("AST:")
8 print(tree.ast_repr())

```

- Also, by the use of Tkinter libray, I decided to leave the User to choose manually a .txt File from the Explorer in order to tokenize its insides as it was in the previous Lexer Laboratory Work:

```

1 ...
2 def select_file():
3     root = tk.Tk()
4     root.withdraw()
5     root.attributes("-topmost", True)
6
7     file_path = filedialog.askopenfilename(initialdir="./Laboratory

```

```
        -Work-6-Parser-AST/ExamplePrograms",  
8                                     title="Select a file")  
9  
10     return file_path  
11     . . .
```

Conclusions / Results

I present here one output for the Laboratory Work nr.6.

First part of the console output is the general information about the laboratory work, student and group:

```
1 Laboratory Work 6 - Parser and AST
2 Student: Gusev Roman
3 Group: FAF-222
```

After that, user is asked to choose the method of input they want:

```
1 MY LEXER AND PARSER:
2 Choose the type of the input:
3 F - File
4 C - Console all lines together
5 L - Console line-by-line
```

After that is the prompt for user:

```
1 Your choice: ____
```

I will provide example of inputs for each type of Input (Each Line from the console starts with ">" symbol to show where to start):

- Manual line-by-line:

Enter lines of code. Type 'exit' to finish.

```
> Data data = ReadFrom("/path1/file1");
```

TOKENS:

VARIABLE_TYPE: 'Data'

ID: 'data'

ASSIGN_OPERATOR: '='

READ_FROM: 'ReadFrom'

LPAREN: '('

PATH: '"/path1/file1"'

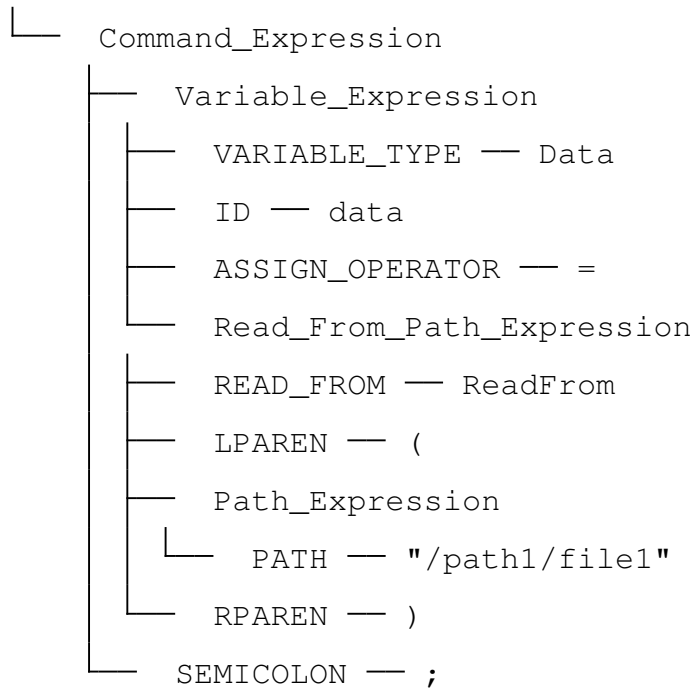
RPAREN: ')'

SEMICOLON: ';'

EOF: 'EOF'

AST:

Program_Expression



- **Manual all line at once:**

Enter lines of code. Type 'exit' to finish.

> Formula x = (-5+6);

> Data y = ReadFrom("/path1");

> exit

TOKENS:

VARIABLE_TYPE: 'Formula'

ID: 'x'

ASSIGN_OPERATOR: '='

LPAREN: '('

OPERATOR: '-'

NUMBER: '5'

OPERATOR: '+'

NUMBER: '6'

RPAREN: ')'

SEMICOLON: ';'

VARIABLE_TYPE: 'Data'

ID: 'y'

ASSIGN_OPERATOR: '='

READ_FROM: 'ReadFrom'

LPAREN: '('

PATH: '"/path1"'

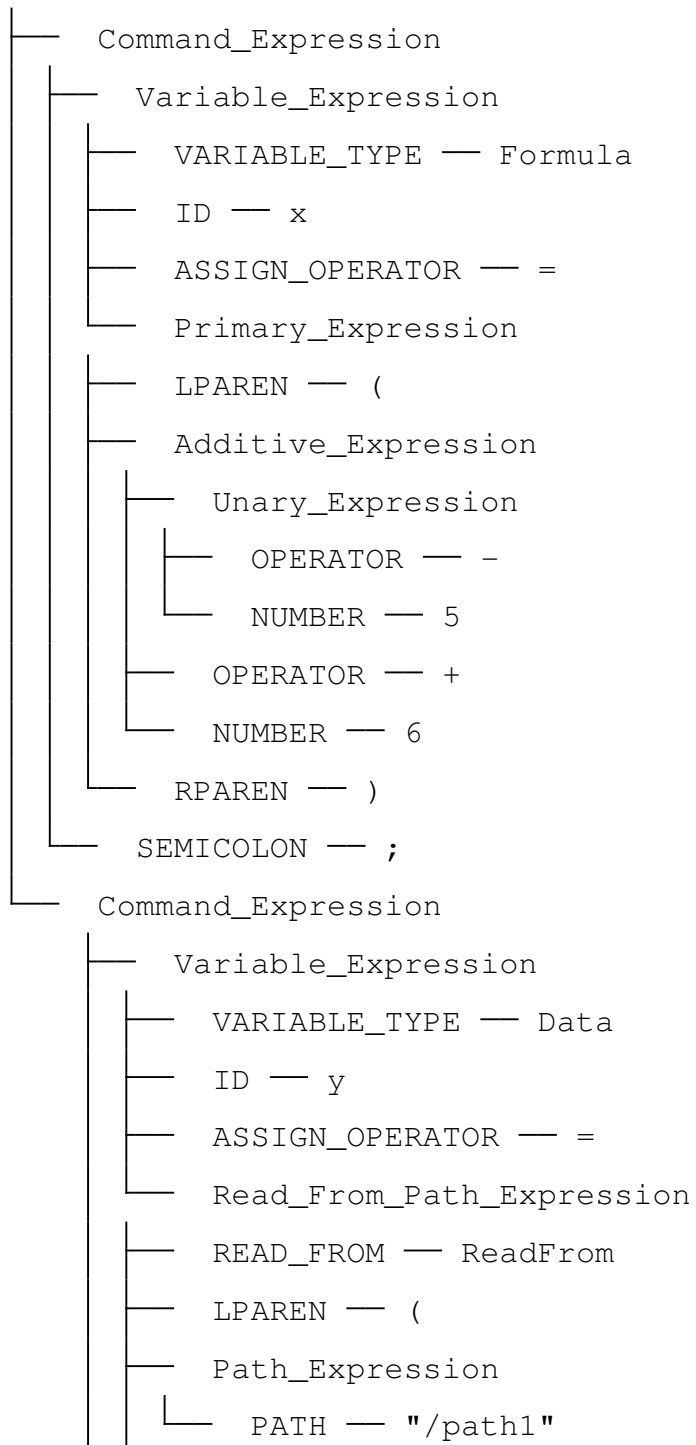
RPAREN: ')'

SEMICOLON: ';'

EOF: 'EOF'

AST:

Program_Expression



```

    └── RPAREN — )
    └── SEMICOLON — ;

```

- File input:

- Here are the contents of the file:

```

1 Formula x = (-5+6);
2 # This is Comment Line :)
3 Data y = ReadFrom("/path1");
4 Data x = ReadFrom("/path1/path/folder");
5 Formula x = x**2+1+2*x*(x+1);

```

- Console Output:

...

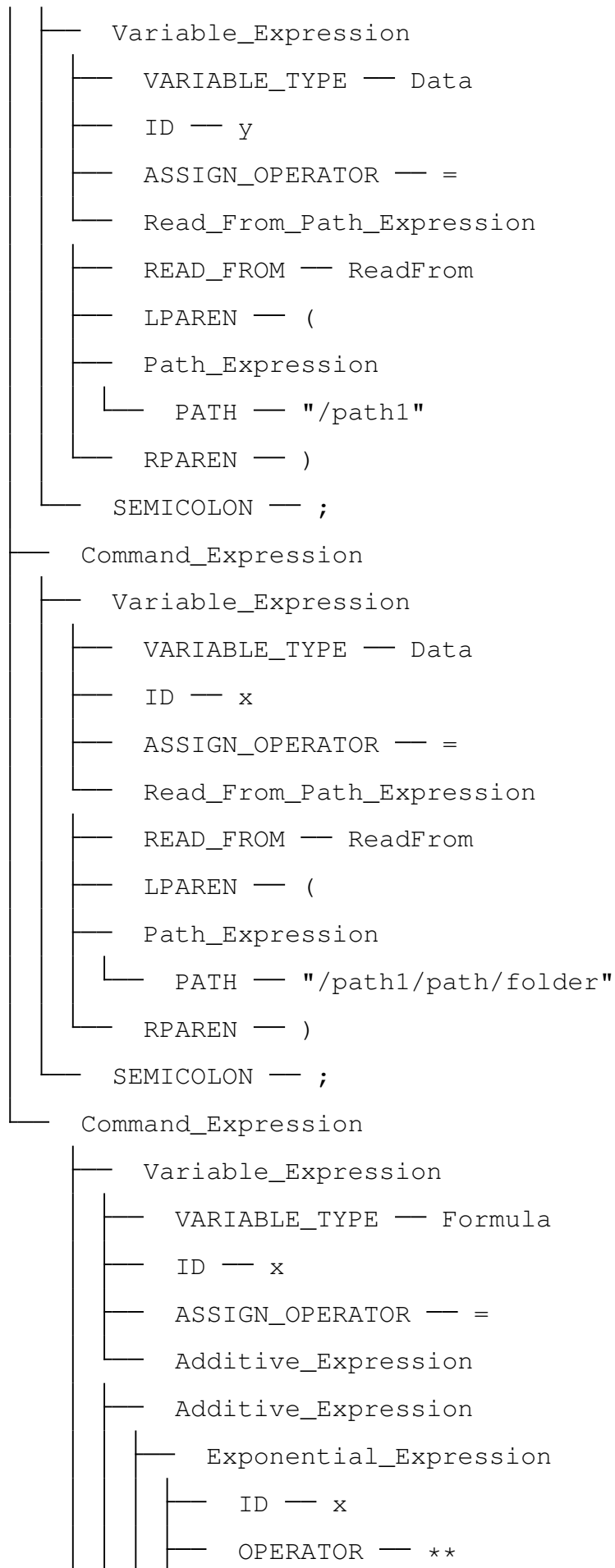
AST:

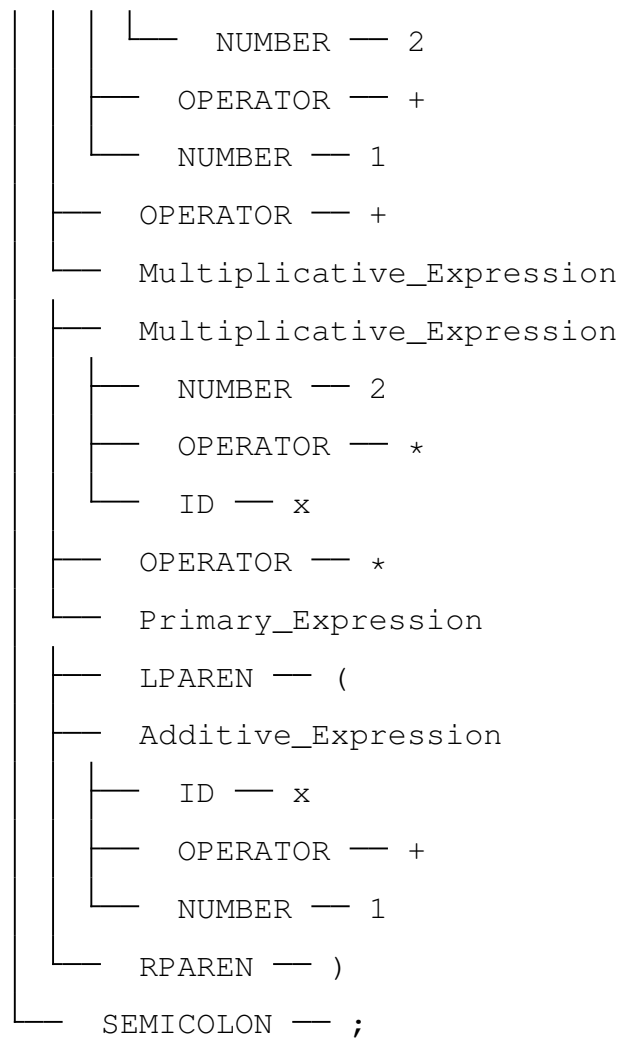
Program_Expression

```

└── Command_Expression
    └── Variable_Expression
        ├── VARIABLE_TYPE — Formula
        ├── ID — x
        ├── ASSIGN_OPERATOR — =
        └── Primary_Expression
            ├── LPAREN — (
            ├── Additive_Expression
            │   ├── Unary_Expression
            │   │   ├── OPERATOR — -
            │   │   └── NUMBER — 5
            │   ├── OPERATOR — +
            │   └── NUMBER — 6
            └── RPAREN — )
    └── SEMICOLON — ;
└── Comment_Expression
    └── COMMENT_LINE — # This is Comment Line :)
└── Command_Expression

```





As a conclusion to this Laboratory Work nr.6, I can say that I accomplished the given task, specifically:

- Implement your own Lexer.
- Build the Abstract Syntax Tree.

Also, I managed to understand better the concept of Parsers, Abstract Syntax Trees, Tokens and Lexemes and their main usages. At the same time, I understood how any Compiler "reads" the code and separates each Lexeme it finds when it analyzes the written code.

Bibliography

- [1] “Types of Parsers in Compiler Design,” GeeksforGeeks, Jul. 26, 2019. Available: <https://www.geeksforgeeks.org/types-of-parsers-in-compiler-design/>. [Accessed: Apr. 26, 2024].
- [2] “Abstract Syntax Tree (AST) in Java,” GeeksforGeeks, Aug. 11, 2021. Available: <https://www.geeksforgeeks.org/abstract-syntax-tree-ast-in-java/>. [Accessed: Apr. 26, 2024].
- [3] Introspective Thinker, How to Create Your Own Programming Language - Episode 3: The Parser, (Jan. 28, 2023). Available: https://www.youtube.com/watch?v=yEBUod_e8Yk. [Accessed: Apr. 27, 2024]