

Paul Kummer

Dr. Hanku Lee

Compilers CSIS 455-01

10/03/2021

Lab 04: Lexing & Parsing

Purpose of Lab

The purpose of this lab is to show how to implement a Lexer and Parser. The Lexer will read a file and convert the characters from the file into tokens that will have a specified meaning. In the process of Lexing, whitespace and newlines will be removed. The line numbers will be tracked, which will be used in future debugging processes. Each token that is created will have a value saved to that token and is used in Parsing. A token can be a word, number, or just a token.

After Lexing, the parser will create an abstract syntax tree (AST). This AST will determine what should be done when evaluating the tokens. The lab demonstrates this by using the Visitor Pattern to visit each node of the tree as it is being created.

Example Code

Main

Main.java

This runs the main line of logic of the compiler. First, the Lexer will be called to run on the input.txt file. Then the lexer object is passed into the Parser, which will begin the parsing.

```
C:\Users\pakum\Desktop\compilers\lab04\code> Main.java > ...
1  import java.io.*;
2  import java.util.*;
3
4  public class Main
5  {
6      Run | Debug
      public static void main(String[] args) throws IOException
7      {
8          Lexer lexer = new Lexer();
9          Parser parser = new Parser(lexer);
10     }
11 }
12
```

Lexer

Lexer.java

The Lexer is responsible for turning characters into meaningful objects. First it will open a file and begin reading each character of the file. If the character is a number, more characters are read until there isn't a number. Then a number token will be created with an integer value. If the character being read is a letter, more letters or numbers will be read and concatenated until a space or operator is encountered. Then a word token will be created with a lexeme value.

While the Lexer is creating new tokens, each token is saved into a hashtable. This table is useful for preventing duplicates of tokens. Also, there can be reserved words too, which will be indicated with a tag value of "true". Otherwise, the words tag value will be "false", indicating it is not a reserved word.

```
C:\Users\pauum\Desktop\compilers\lab04> code > Lexer.java > % Lexer > scan()
1 import java.io.*;
2 import java.util.*;
3
4 public class Lexer
5 {
6     public int line = 1;
7     private char peek = ' ';
8
9     private FileInputStream in;
10    private BufferedInputStream bin;
11
12    private Hashtable<String, Word> words = new Hashtable<String, Word>();
13
14    public Lexer()
15    {
16        reserve(new Word("true", Tag.TRUE));
17        reserve(new Word("false", Tag.FALSE));
18        setupIOStream();
19    }
20
21    void reserve(Word w)
22    {
23        words.put(w.lexeme, w);
24    }
25
26    void setupIOStream()
27    {
28        try
29        {
30            in = new FileInputStream("input.txt");
31            bin = new BufferedInputStream(in);
32        }
33        catch(IOException e)
34        {
35            System.out.println("IOException: ");
36        }
37    }
38
39    void readch() throws IOException
40    {
41        peek = (char) bin.read();
42    }
43
44    public Token scan() throws IOException
45    {
46        for(;; readch())
47        {
48            if(peek == ' ' || peek == '\t')
49            {
50                continue;
51            }
52
53            else if(peek == '\n')
54            {
55                line++;
56            }
57
58            else
59            {
60                break;
61            }
62        }
63
64        if(Character.isDigit(peek))
65        {
66            int v = 0;
67
68            do
69            {
70                v = 10 * v + Character.digit(peek, 10);
71                readch();
72            } while(Character.isDigit(peek));
73
74            return new Num(v);
75        }
76
77        if(Character.isLetter(peek))
78        {
79            StringBuffer b = new StringBuffer();
80
81            do
82            {
83                b.append(peek);
84                readch();
85            } while(Character.isLetterOrDigit(peek));
86
87            String s = b.toString();
88            Word w = words.get(s);
89
90            if (w != null)
91            {
92                return w;
93            }
94
95            w = new Word(s, Tag.ID);
96            words.put(s, w);
97
98            return w;
99        }
100
101        Token t = new Token(peek);
102        peek = ' ';
103
104        return t;
105    }
106 }
107
108 }
```

Token.java

Tokens are used to be the smallest meaningful unit in the compilation process. Each token has a value assigned to it that is held within a variable called tag. This tag value is important for Parsing.

```
C: > Users > pakum > Desktop > compilers > lab04 > code > Token.java > Token
1  public class Token
2
3      public final int tag;
4
5      public Token (int t)
6      {
7          tag = t;
8      }
9
10     public String toString()
11     {
12         return "" + (char)tag;
13     }
14 }
```

Tag.java

The Tag class is used to hold special identifying values that are used in the Parsing phase.

```
C:\Users > pakum > Desktop > compilers > lab04 > code > Tag.java > Tag > NUM
1  // File Tag.java
2
3  public class Tag
4  {
5      public final static int
6          NUM = 256,
7          ID = 257,
8          TRUE = 258,
9          FALSE = 259,
10         EOF = 65535;
11 }
```

Num.java

Num is a subclass of Token, and it stores a tag value of indicating it is a number. Also, a Num has an additional value, which is an integer associated with the Num token.

```
C:\Users > pakum > Desktop > compilers > lab04 > code > Num.java > Num
1 public class Num extends Token
2
3     public final int value;
4
5     public Num(int v)
6     {
7         super(Tag.NUM);
8         value = v;
9     }
10
11    public String toString()
12    {
13        return "" + value;
14    }
15 }
```

Word.java

Word is a subclass of Token, and it stores a tag value indicating it is a word. Also, a word has an additional value, which is a lexeme stored as a string associated with the Word token.

```
C:\Users > pakum > Desktop > compilers > lab04 > code > Word.java > Word
1 public class Word extends Token
2
3     public String lexeme = "";
4
5     public static final Word True = new Word("true", Tag.TRUE);
6     public static final Word False = new Word("false", Tag.FALSE);
7
8     public Word (String s, int tag)
9     {
10         super(tag);
11         lexeme = s;
12     }
13
14     public String toString()
15     {
16         return lexeme;
17     }
18 }
```

Parser

Parser.java

The parser is responsible for evaluating what the tokens should do based on their tags. This will result in an abstract syntax tree (AST) being created.

First, the parser will begin to read the tokens from the Lexer, using the Lexer's "scan" method, and put them into a queue. Whenever a newline is encountered, the tokens that have been read so far will be considered an expression and are sent to a new compilation unit, where they will be evaluated. This process will continue until the end of file (EOF) tag is encountered.

```
C:\Users > pakum > Desktop > compilers > lab04 > code > Parser.java > Parser > Parser(Lexer)
1  import java.io.*;
2  import java.util.*;
3
4  class Parser
5  {
6      public Parser(Lexer lex) throws IOException
7      {
8          Queue<Token> tokens = new LinkedList<Token>();
9          Token tok = lex.scan();
10         int lineNum = lex.line;
11
12
13         while(tok.tag != Tag.EOF)
14         {
15             if(tok != null)
16             {
17                 tokens.add(tok);
18             }
19
20             tok = lex.scan();
21
22             //This will preserve original newlines
23             while(lex.line > lineNum)
24             {
25                 Node root = new CompilationUnit(tokens);
26                 root.accept(new ASTVisitor());
27                 System.out.print("\n\n\t[ New Expression ]\n");
28
29                 //lineNum++; //This will add as many lines as there were in the document
30                 lineNum = lex.line; //This will add only one line
31
32                 tokens.clear();
33             }
34         }
35
36         Node root = new CompilationUnit(tokens);
37         root.accept(new ASTVisitor());
38     }
39 }
```


Nodes

The Nodes are used in the visitor pattern and allow each node to have its own stored values and private methods, like a print method. Some nodes have a unique identifier, which is somewhat redundant because the token already has the value stored. More work should be done to utilize the tokens into nodes.

Node.java

Node is the base class used to build an AST and is utilized by other classes as a parent class.

```
C:\Users > pakum > Desktop > compilers > lab04 > code > Node.java > ...
1  import java.util.*;
2
3  public class Node
4  {
5      Token tok = null;
6
7      public Node()
8      {
9      }
10
11     public Node(Token t)
12     {
13         tok = t;
14     }
15
16     public void accept(ASTVisitor v)
17     {
18         System.out.println("Node: " + tok.toString());
19     }
20 }
```

LiteralNode.java

A literal node cannot be broken down any further and has a string representation.

```
C: > Users > pakum > Desktop > compilers > lab04 > code > LiteralNode.java > LiteralNode
1  //will not call accept because the visitor method is not used
2
3  public class LiteralNode extends Node
4  {
5      public String literal;
6
7      public LiteralNode()
8      {
9
10     }
11
12     public LiteralNode (String literal)
13     {
14         this.literal = literal;
15     }
16
17     @Override
18     public void accept(ASTVisitor v)
19     {
20         printNode();
21     }
22
23     void printNode()
24     {
25         System.out.println("Literal: " + literal);
26     }
27 }
```

AssignmentNode.java

An assignment node is used to assign a value to a node on the left-hand-side of the equals operator. The right-hand-side can be any expression represented as a node.

```
C:\Users > pakum > Desktop > compilers > lab04 > code > AssignmentNode.java > AssignmentNode
1 public class AssignmentNode extends Node
2
3     String ident = "=";
4
5     Node left;
6     Node right;
7
8     public AssignmentNode()
9     {
10
11     }
12
13     public AssignmentNode (Node l, Node r)
14     {
15         this.left = l;
16         this.right = r;
17     }
18
19     @Override
20     public void accept(ASTVisitor v)
21     {
22         v.visit(this);
23     }
24
25     void printNode()
26     {
27         System.out.println("Operand: " + ident);
28     }
29 }
```

CompilationUnit.java

The compilation unit is the starting point of building an expression. It stores a queue of tokens that will become nodes and be evaluated into an AST.

```
C:\Users > pakum > Desktop > compilers > lab04 > code > CompilationUnit.java > ...
1  import java.util.*;
2
3  public class CompilationUnit extends Node
4  {
5      //public AssignmentNode assign;
6      Queue<Token> toks;
7
8      public CompilationUnit(Queue<Token> t)
9      {
10         this.toks = t;
11     }
12
13     @Override
14     public void accept(ASTVisitor v)
15     {
16         v.visit(this);
17     }
18 }
```

AdditionNode.java

The addition node holds two nodes, a left-hand-side, and a right-hand-side, that should be combined to form a new value.

```
C:\> Users > pakum > Desktop > compilers > lab04 > code > AdditionNode.java > AdditionNode

1  public class AdditionNode extends Node
2
3      String ident = "+";
4
5      Node left;
6      Node right;
7
8      public AdditionNode()
9      {
10
11      }
12
13     public AdditionNode (Node l, Node r)
14     {
15         this.left = l;
16         this.right = r;
17     }
18
19
20     @Override
21     public void accept(ASTVisitor v)
22     {
23         v.visit(this);
24     }
25
26     void printNode()
27     {
28         System.out.println("Operand: " + ident);
29     }
30 }
```

SubtractionNode.java

The subtraction node holds two nodes, a left-hand-side, and a right-hand-side, that should be subtracted to form a new value.

```
C: > Users > pakum > Desktop > compilers > lab04 > code > SubtractionNode.java > SubtractionNode
1  public class SubtractionNode extends Node
2
3      String ident = "-";
4
5      Node left;
6      Node right;
7
8      public SubtractionNode()
9      {
10
11      }
12
13      public SubtractionNode (Node l, Node r)
14      {
15          this.left = l;
16          this.right = r;
17      }
18
19
20      @Override
21      public void accept(ASTVisitor v)
22      {
23          v.visit(this);
24      }
25
26      void printNode()
27      {
28          System.out.println("Operand: " + ident);
29      }
30
```

MultiplicationNode.java

The multiplication node holds two nodes, a left-hand-side, and a right-hand-side, that should be multiplied to form a new value.

```
> Users > pakum > Desktop > compilers > lab04 > code > MultiplicationNode.java > MultiplicationNode > MultiplicationNode(Node, Node)
1 public class MultiplicationNode extends Node
2
3     String ident = "*";
4
5     Node left;
6     Node right;
7
8     public MultiplicationNode()
9     {
10
11     }
12
13     public MultiplicationNode (Node l, Node r)
14     {
15         this.left = l;
16         this.right = r;
17     }
18
19
20     @Override
21     public void accept(ASTVisitor v)
22     {
23         v.visit(this);
24     }
25
26     void printNode()
27     {
28         System.out.println("Operand: " + ident);
29     }
30
```

DivisionNode.java

The division node holds two nodes, a left-hand-side, and a right-hand-side, that should be divided to form a new value.

```
C:\Users > pakum > Desktop > compilers > lab04 > code > DivisionNode.java > DivisionNode > DivisionNode(Node, Node)
1  public class DivisionNode extends Node
2
3      String ident = "/";
4
5      Node left;
6      Node right;
7
8      public DivisionNode()
9      {
10
11      }
12
13      public DivisionNode (Node l, Node r)
14      {
15          this.left = l;
16          this.right = r;
17      }
18
19
20      @Override
21      public void accept(ASTVisitor v)
22      {
23          v.visit(this);
24      }
25
26      void printNode()
27      {
28          System.out.println("Operand: " + ident);
29      }
30
```


ModuloNode.java

The modulo node holds two nodes, a left-hand-side, and a right-hand-side, that should be divided and have the remainder returned to form a new value.

```
C:\Users > pakum > Desktop > compilers > lab04 > code > ModuloNode.java > ModuloNode > ModuloNode(Node, Node)
1  public class ModuloNode extends Node
2
3      String ident = "%";
4
5      Node left;
6      Node right;
7
8      public ModuloNode()
9      {
10
11      }
12
13      public ModuloNode (Node l, Node r)
14      {
15          this.left = l;
16          this.right = r;
17      }
18
19
20      @Override
21      public void accept(ASTVisitor v)
22      {
23          v.visit(this);
24      }
25
26      void printNode()
27      {
28          System.out.println("Operand: " + ident);
29      }
30
```

Visitors

ASTVisitor.java

The visitor is where the parsing really happens. When each node is visited, it will call `accept` on a node, which begins the processing of that node. To begin the parsing, the base node of compilation unit is visited first. This node will check if the next tokens in the queue are part of an operation or a literal. Based on what the first and second tokens are, switch statements will create a new node of whatever type is appropriate. For example, if the first node is a literal and the second node is '+', then an assignment node will be created and visited. However, if the first token is a literal and the second is null, then the literal node will be created and visited. If parentheses are encountered, a new compilation unit is created with a new queue of only the tokens within the parentheses. This compilation unit can be used as a node in any other operation.

When nodes other than the compilation unit are visited, the node's private values, which are nodes, will be accepted. Also, as this occurs, the level of the AST will be shown with dots and the values of the node will be displayed. That way a visualization of the AST can be seen.

ASTVisitor.java cont.

```

1  import java.util.*;
2
3  public class ASTVisitor
4  {
5      int level = 0;
6
7      public void visit(CompilationUnit n)
8      {
9          System.out.print("Compilation Unit:\n");
10
11          level++;
12          Node left;
13          Node right;
14          Node operation;
15
16          Queue<Token> newToks = new LinkedList<Token>();
17          Token first = n.toks.poll();
18          Token op = n.toks.peek();
19
20          if(op != null)
21          {
22              right = new CompilationUnit(n.toks);
23
24              switch(first.tag)
25              {
26                  case Tag.NUM: //the token is a number
27                      left = new LiteralNode(first.toString());
28                      break;
29                  case Tag.ID: //the token is a word
30                      left = new LiteralNode(first.toString());
31                      break;
32                  case (int) '(':
33                      while((char)n.toks.peek().tag != ')' && n.toks.peek() != null && (char)first.tag != ')')
34                      {
35                          newToks.add(n.toks.poll());
36                      }
37                      n.toks.poll(); //remove the '(' token
38                      left = new CompilationUnit(newToks);
39                      op = n.toks.peek();
40                      break;
41                  default:
42                      left = new LiteralNode(first.toString());
43                      break;
44              }
45
46              switch((char)op.tag)
47              {
48                  case '=':
49                      operation = new AssignmentNode(left, right);
50                      break;
51                  case '+':
52                      operation = new AdditionNode(left, right);
53                      break;
54                  case '-':
55                      operation = new SubtractionNode(left, right);
56                      break;
57                  case '*':
58                      operation = new MultiplicationNode(left, right);
59                      break;
60                  case '/':
61                      operation = new DivisionNode(left, right);
62                      break;
63                  case '%':
64                      operation = new ModuloNode(left, right);
65                      break;
66                  default:
67                      operation = left;
68                      break;
69              }
70              n.toks.poll(); //remove the operation symbol
71              operation.accept(this);
72          }
73
74          else
75          {
76              switch(first.tag)
77              {
78                  case Tag.NUM: //the token is a number
79                      left = new LiteralNode(first.toString());
80                      dots();
81                      left.accept(this);
82                      break;
83                  case Tag.ID: //the token is a word
84                      dots();
85                      left = new LiteralNode(first.toString());
86                      left.accept(this);
87                      break;
88                  default:
89                      if(first.tag != Tag.EOF)
90                      {
91                          dots();
92                          left = new LiteralNode(first.toString());
93                          left.accept(this);
94                      }
95                      break;
96              }
97          }
98          level--;
99      }

```

ASTVisitor.java cont.

```
G:\Users\paku > Desktop > compilers > lab04 > code > ASTVisitor.java > ASTVisitor > visit(CompilationUnit)
180
181 public void visit(AssignmentNode n)
182 {
183     dots();
184     System.out.print("Assignment:\n");
185
186     level++;
187     dots();
188     n.left.accept(this);
189     level--;
190
191     level++;
192     dots();
193     n.printNode();
194     level--;
195
196     level++;
197     dots();
198     n.right.accept(this);
199     level--;
200 }
201
202 public void visit (AdditionNode n)
203 {
204     dots();
205     System.out.print("Addition:\n");
206
207     level++;
208     dots();
209     n.left.accept(this);
210     level--;
211
212     level++;
213     dots();
214     n.printNode();
215     level--;
216
217     level++;
218     dots();
219     n.right.accept(this);
220     level--;
221 }
222
223 public void visit (SubtractionNode n)
224 {
225     dots();
226     System.out.print("Subtraction:\n");
227
228     level++;
229     dots();
230     n.left.accept(this);
231     level--;
232
233     level++;
234     dots();
235     n.printNode();
236     level--;
237
238     level++;
239     dots();
240     n.right.accept(this);
241     level--;
242 }
243
244 public void visit (MultiplicationNode n)
245 {
246     dots();
247     System.out.print("Multiplication:\n");
248
249     level++;
250     dots();
251     n.left.accept(this);
252     level--;
253
254     level++;
255     dots();
256     n.printNode();
257     level--;
258
259     level++;
260     dots();
261     n.right.accept(this);
262     level--;
263 }
264
265 public void visit (DivisionNode n)
266 {
267     dots();
268     System.out.print("Division:\n");
269
270     level++;
271     dots();
272     n.left.accept(this);
273     level--;
274
275     level++;
276     dots();
277     n.printNode();
278     level--;
279
280     level++;
281     dots();
282     n.right.accept(this);
283     level--;
284 }
285 }
```

ASTVisitor.java cont.

```
207     public void visit (ModuloNode n)
208     {
209         dots();
210         System.out.print("Modulo:\n");
211
212         level++;
213         dots();
214         n.left.accept(this);
215         level--;
216
217         level++;
218         dots();
219         n.printNode();
220         level--;
221
222         level++;
223         dots();
224         n.right.accept(this);
225         level--;
226     }
227
228     public void visit(LiteralNode n)
229     {
230         level++;
231         dots();
232         n.accept(this);
233         level--;
234     }
235
236     public void visit(Node n)
237     {
238         level++;
239         dots();
240         n.accept(this);
241         level--;
242     }
243
244     private void dots()
245     {
246         System.out.print(new String(new char[level*4]).replace('\0', '.'));
247     }
248 }
```

Code Execution

As the code executes, the Lexer begins to read in “input.txt”, which contains some expression written out on each line. Then the lexer is passed into the parser, which starts creating nodes and making an AST. This AST is displayed as output. Below is the expressions in the file and what the result is of the program’s evaluation.

input.txt

c=b+1+input

c=l+mmmm

input = b + result

x = b * (c - 2 * x) + 33

program execution

```
cx3645kg@smaug:~/Documents/CSIS455_compilers/lab04/code$ rm *.class
cx3645kg@smaug:~/Documents/CSIS455_compilers/lab04/code$ javac *.java
cx3645kg@smaug:~/Documents/CSIS455_compilers/lab04/code$ java Main
Compilation Unit:
...Assignment:
.....Literal: c
.....Operand: =
.....Compilation Unit:
.....Addition:
.....Literal: b
.....Operand: +
.....Compilation Unit:
.....Addition:
.....Literal: 1
.....Operand: +
.....Compilation Unit:
.....Literal: input

[ New Expression ]
Compilation Unit:
...Assignment:
.....Literal: c
.....Operand: =
.....Compilation Unit:
.....Addition:
.....Literal: 1
.....Operand: +
.....Compilation Unit:
.....Literal: mmmm

[ New Expression ]
Compilation Unit:
...Assignment:
.....Literal: input
.....Operand: =
.....Compilation Unit:
.....Addition:
.....Literal: b
.....Operand: +
.....Compilation Unit:
.....Literal: result

[ New Expression ]
Compilation Unit:
...Assignment:
.....Literal: x
.....Operand: =
.....Compilation Unit:
.....Multiplication:
.....Literal: b
.....Operand: *
.....Compilation Unit:
.....Addition:
.....Compilation Unit:
.....Subtraction:
.....Literal: c
.....Operand: -
.....Compilation Unit:
.....Multiplication:
.....Literal: 2
.....Operand: *
.....Compilation Unit:
.....Literal: x
.....Operand: +
.....Compilation Unit:
.....Literal: 33
cx3645kg@smaug:~/Documents/CSIS455_compilers/lab04/code$
```

Conclusion

After this lab, I have a much better understanding of how the visitor pattern works. Also, I learned how to create a parser in a different way than I have done before. I have made a parser in C++ that uses recursive descent, which this parser utilizes a similar concept by using a compilation unit as a node that is visited within a different compilation unit. Initially, I found the code hard to follow because it is spread out across many different files, but many of the files are used in similar ways, so the processing only occurs in one place.

I found this lab to be difficult to do and hard to figure out how to make the previous assignments work with this lab. I think the connection between this lab and the other assignments was weak. I was able to see some similarities, but I wasn't sure how to implement it with the parser. I was able to make it work, but I don't know if I did it the best way. Overall, this lab challenged me, and I enjoyed figuring out how to make it work. I had to go back and watch other lecture videos, consult the Compilers textbook, and review the visitor pattern.