**Yan Peng**
**921759056**
**September 26, 2021**
# Assignment 2 Documentation

# GitHub Repository

Link for assignment 2:

https://github.com/sfsu-csc-413-fall-2021/assignment-2---lexer-yuqiao1205

## Project Introduction and Overview

The aim of this project is to understand the function of lexical analysis in the compilation process. The purpose of the lexical analyzer in this assignment is to read the input characters of the source files and convert them into a sequence of tokens, which is the first stage of a compiler.

The Lexer project required me to extend the Lexer class by adding additional tokens and reserved words to perform complete lexical analysis of files with the suffix ".x", meanwhile, to distinguish lexemes starting with a digit token and decide if it is an Integer, NumberLit (any series of numbers followed by a decimal point) or DateLit (one or two number representing the month, followed by "~".) Also, print out the number of lines and token type and include error reporting if invalid tokens are found.

# Scope of Work

| Task | Completed |
|---|:---:|
| Created the getDigitToken method of the Lexer class to distinguish the kinds of tokens starting with a digit (eg, Integer, NumberLit and DateLit) | ✓ |
| Refactored nextToken to delegate discrimination of tokens starting with a digit to a method ( getDigitToken) | ✓ |
| Refactored IdToken's try /catch block to a method (getIdToken) | ✓ |
| Test the implementation with expressions that test all possible cases. The following expressions were used: | ✓ |
| all .x files under folder of sample_files | ✓ |
| program { int i int j<br>    i = (i * j + 7.1 + 7 / i)-2<br>    date = 11~22~2021<br>    date1 = 09~09~89<br>    number = 3.12<br><br>} | ✓ |
| program { int i int j<br>  i = i + j + 7<br>  j = write(i)<br>} | ✓ |
| Updated newNumberToken method of the Lexer class | ✓ |
|     Token newNumberToken( String number, Tokens kind) | ✓ |
| Created readInteger method of the Lexer class | ✓ |
|     String readInteger() | ✓ |
| Created (isNumberLit) method of the Lexer class to check valid floating number, return boolean | ✓ |
| Created (isDateLit) method of the Lexer class to check valid date type, return boolean | ✓ |
| Added additional token kinds to "tokens" file and regenerated Tokens and TokenTypes classes | ✓ |
| Updated Token class to include the line number | ✓ |
|     Constructors(int and Symbol parameters) | ✓ |
|     int getLineNumber() | ✓ |
| Removed initial debug text that shows the file information | ✓ |
| Updated the Lexer to allow input via a filename provided as a command line argument | ✓ |
| Implement the main method in Lexer class to print out the required format | ✓ |

# Execution and Development Environment

I used the IntelliJ IDE on my Mac to finish this assignment and tested in both the IDE and shell.

# Compilation Result

**Using the instructions provided in the assignment one specification:**

```
> javac lexer/Lexer.java
> java lexer.Lexer
```



```
> java lexer.Lexer sample_files/simple.x
```



No error messages or warnings were displayed, and the application ran as expected

## Assumptions

I assumed some tokens may be malformed and implemented some error handling for these invalid tokens.

## Implementation

### Lexer

The Lexer class contains main () method of our application. We first read the input file name and pass it on to the SourceReader for reading. The input file is read character by character and Token is generated when a valid Symbol or Identifier is found for initial analysis. The Lexer reads file name from the command line and outputs usage if no file name is found.

### Token

The Token class is updated to include line numbers at which the token was found. This line number is provided by the Lexer while creating the Token object.

### TokenSetup

The TokenSetup class is responsible for reading tokens from the tokens file and automatically generating TokenType and Tokens classes.

### TokenType

TokenType class is a HashMap, mapping from the type Tokens to Symbol.

### Tokens

The Tokens class is the enumeration of all the tokens.

### SourceReader

SourceReader class is responsible for reading the input file. The file is scanned by a BufferedReader. The position of character and the line number is in the SourceReader.

### Symbol

Symbol class is responsible for generating a symbol from the string of the new token and kind of type Tokens. This class also stores the symbol name.
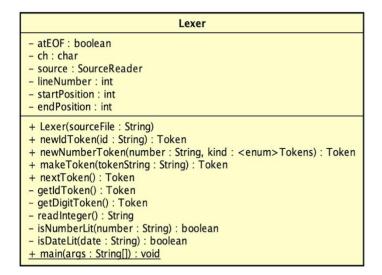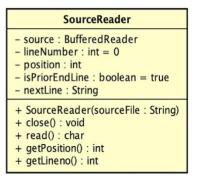
## Code Organization

I put processing of all tokens starting with a digit into the method getDigitToken in Lexer class to localize the functionality and avoid repetition. I also noticed that the logic to check if the input is integer was duplicated in a lengthy try/catch block, I refactored this part into a method readInteger() to handle the case of the Integer(also part of NumberLit/Date token kinds).  Also, I refactored the idToken try/catch block to a method getIdToken(), to check if the token is id token by checking every character is JavaIdentifierStart.

I also noticed some unnecessary parameters such like (int startPosition, int endPosition) when in defining the function such like newIdToken(), makeToken() as their parameters, so I removed those parameters because they are already available as fields in the Lexer class, accessible to the methods above. Though I understand the original code with these parameters don't have to get the object into a certain state, with certain values in the fields to test, and I can see how the method in original code might have been written in a functional style, but I prefer to have a single/source path for information in this case.

# Class Diagram

The following class diagrams show the details of all the classes in this project, including the inheritance hierarchy
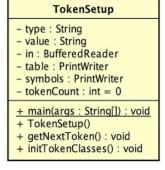
| Lexer |
| --- |
| – atEOF : boolean<br>– ch : char<br>– source : SourceReader<br>– lineNumber : int<br>– startPosition : int<br>– endPosition : int |
| + Lexer(sourceFile : String)<br>+ newIdToken(id : String) : Token<br>+ newNumberToken(number : String, kind : <enum>Tokens) : Token<br>+ makeToken(tokenString : String) : Token<br>+ nextToken() : Token<br>– getIdToken() : Token<br>– getDigitToken() : Token<br>– readInteger() : String<br>– isNumberLit(number : String) : boolean<br>– isDateLit(date : String) : boolean<br>+ main(args : String[]) : void |

| SourceReader |
| --- |
| – source : BufferedReader<br>– lineNumber : int = 0<br>– position : int<br>– isPriorEndLine : boolean = true<br>– nextLine : String |
| + SourceReader(sourceFile : String)<br>+ close() : void<br>+ read() : char<br>+ getPosition() : int<br>+ getLineno() : int |

| TokenType |
| --- |
| |
| + TokenType() |

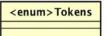| Token |
| --- |
| – leftPosition : int<br>– rightPosition : int<br>– symbol : Symbol<br>– lineNumber : int |
| + Token(leftPosition : int, rightPosition : int, lineNumber : int, symbol : Symbol)<br>+ toString() : String<br>+ getLeftPosition() : int<br>+ getRightPosition() : int<br>+ getLineNumber() : int<br>+ getKind() : <enum>Tokens |

| TokenSetup |
| --- |
| – type : String<br>– value : String<br>– in : BufferedReader<br>– table : PrintWriter<br>– symbols : PrintWriter<br>– tokenCount : int = 0 |
| + main(args : String[]) : void<br>+ TokenSetup()<br>+ getNextToken() : void<br>+ initTokenClasses() : void |

| Symbol |
| --- |
| – name : String<br>– kind : <enum>Tokens |
| – Symbol(name : String, kind : <enum>Tokens)<br>+ toString() : String<br>+ getKind() : <enum>Tokens<br>+ symbol(newTokenString : String, kind : <enum>Tokens) : Symbol |

| <enum>Tokens |
| --- |
| |
| |

# Results and Conclusion

This project helped me to understand the concept of Lexer, which I had not encountered before. The Lexer's job is the process of converting a sequence of characters such as the source code of

a computer program into a sequence of tokens. I successfully implemented all the required features defined in the project specification.  If there is invalid token identified, the Lexer will stop processing the next line and report the error with its position with line numbers.

## Challenges

It was challenging for me to implement the different kinds of digit tokens. My first idea was to set the default value to Integer and use it to check integer part of the number case and date case. Later I introduced the helper methods: isNumberLit() to validate the number cases and isDateLit() to validate Date cases.

## Future Work

I think in future this module will be used by the Parser and Compiler modules in the project. Some interesting future work could be to extend the lexer with other token types, e.g. ++.