

Compiler Project #1, 2022

The goal of the first term-project is to implement a lexical analyzer (a.k.a., scanner) as we have learned. More specifically, you will implement the lexical analyzer for a simplified Java programming language with the following lexical specifications.

<Lexical Specification>

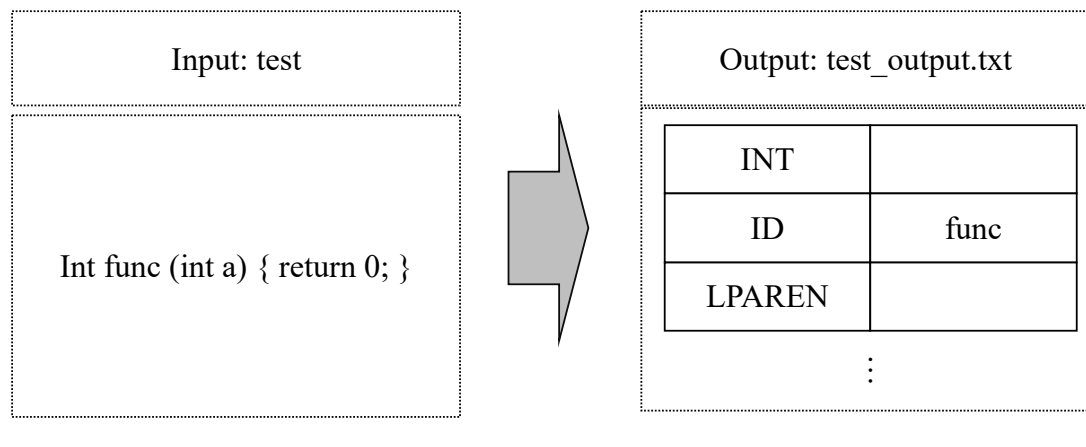
- **Variable type**
 - **int** for a signed integer
 - **char** for a single character
 - **boolean** for a Boolean string
 - **string** for a literal string
- **Signed integer**
 - A single zero digit (e.g., 0)
 - A non-empty sequence of digits, starting from a non-zero digit
 - e.g., 1, 22, 123, 56, ... any non-zero positive integers
 - e.g., 001 is NOT allowed
 - A non-empty sequence of digits, starting from a minus sign symbol and a non-zero digit
 - e.g., -1, -22, -123, -56, ... any non-zero negative integers
- **Single character**
 - A single digit, English letter, block, or any symbol, starting from and terminating with a single symbol
 - e.g., 'a', '1', ' ', '&', ...
- **Boolean string:** true and false
- **Literal string**
 - Any combination of digits, English letters, and blanks, starting from and terminating with "
 - e.g., "Hello World", "My student ID is 20200000"
- **An identifier of variables and functions**
 - A non-empty sequence of English letters, digits, and underscore symbols, starting from an English letter or a underscore symbol (not digits)
 - e.g., i, j, k, abc, ab_123, func1, func_, __func_bar__
- **Keywords for special statements**
 - **if** for if statement
 - **else** for else statement
 - **while** for while-loop statement
 - **class** for class statement
 - **return** for return statement
- **Arithmetic operators:** +, -, *, and /
- **Assignment operator:** =
- **Comparison operators:** <, >, ==, !=, <=, and >=
- **A terminating symbol of statements:** ;
- **A pair of symbols for defining area/scope of variables and functions:** { and }
- **A pair of symbols for indicating a function/statement:** (and)
- **A pair of symbols for using an array:** [and]
- **A symbol for separating input arguments in functions:** ,
- **Whitespaces:** a non-empty sequence of \t, \n, and blanks

According to the above specification, you will (1) define tokens (e.g., token names) for a simplified Java programming language, (2) make regular expressions which describe the patterns of the tokens, (3) construct a NFA (Nondeterministic Finite Automata) for the regular expressions, (4) translate the NFA to DFA (Deterministic Finite Automata), especially in the form of a table, and (5) implement a program which does a lexical analysis (recognizing tokens).

Important Note

1. You MUST build regular expressions, NFAs, and DFAs by writing or drawing by hand. (DO NOT use any kind of computer programs like Lex for this procedure.)
2. You can use C, C++, JAVA, or Python to implement your lexical analyzer.
3. Your lexical analyzer MUST run on Linux or Unix-like OS without any error. (DO NOT use Windows-only APIs and please test your program on Linux machines before submission)
4. Your lexical analyzer should work as follows:

- **On a command line, your analyzer must run with the following command**
lexical_analyser <input_file_name>
- **Input:** A program written in a simplified Java programming language
(You don't need to think about the syntax grammar of the program yet)
- **Output:** <input_file_name>_output.txt
 - o (If an input program has no error) A symbol table which stores the information of all tokens including their names and optional values
 - This output will be used as an input of your next term-project (syntax analyzer)
 - o (Otherwise) An error report which explains why and where the error occurred (e.g., line number)



5. Do not include “WHITESPACE” tokens in the output token list
6. There will be some issues with the symbol ‘-’. Please consider the syntax of the given code to address this.

Term-Project Schedule and Submission

- Deadline: 5/3, 23:59 (through the e-class system)
 - For a delayed submission, you will lose $0.1 \times$ “your original project score” per each delayed day
- Submission:
 - Filename
 - <your student ID>_<your name>.zip or .tar.gz
 - e.g., 2020xxxx_Jinsung_Kim.tar.gz
 - The compressed file should include
 - The source code of your lexical analyzer with detailed comments
 - The executable binary file of your lexical analyzer
 - Document (the most important thing!)
 - It must include (1) the definition of tokens and their regular expressions, (2) the DFA transition graph or table for recognizing the regular expressions, (3) all about how your lexical analyzer works for recognizing tokens (for example, overall procedures, implementation details like algorithms and data structures, working examples, and so on)
 - Test input files and outputs which you used for this project
 - The test input files are not given. Please make some test files, by yourself, which can examine all the token patterns.
- If there exist any error in the given lexical specification, please send an e-mail to me (kimjsung@cau.ac.kr)