

Kennesaw State University

College of Computing and Software Engineering

Department of Computer Science

CS 4308-W02: Concepts of Programming Languages

Deliverable 1 - Scanner

Jesse Schultheis



Katlin Scott



William Stigall



Hailey Walker



September 19, 2023

Initial Problem Statement

The problem is to create a Python-based interpreter for a subset of the SCL language, involving the development of a scanner to parse source code into tokens (including keywords, identifiers, operators, constants, and special characters), specifying the grammar using BNF/EBNF, and providing both console output and a JSON file as output when processing a designated SCL source file via a command-line input.

Summary/Purpose

The purpose of this deliverable is to implement a Scanner file that reads the source code of a given SCL file and then creates a list of tokens in another file.

Detailed Description/Work Done

Python code defines the keywords ("DISPLAY", "IF", "THEN", "ENDIF", "FUNCTION", "IS", "ENDFUN", "PARAMETERS", "INTEGER", "FLOAT", "CHAR", "NOT") and token types ("IDENTIFIER", "UNSIGNICON", "SIGNICON", "PLUS", "MINUS", "STAR", "DIVOP", "EQUOP", "RELOP", "LB", "RB", "LP", "RP", "COMMA", "OF"), takes in six SCL files and reads them. It then tokenizes all of the files and puts each token with its correct token type in the output.json file.

Input Files

arduino_ex1.scl

arrayex1b.scl

bitops1.scl → Used in the screenshots in this report

datablistp.scl

linkedg.scl

welcome.scl

Output Files

output.json

tokens.json → Used in the screenshots in this report

Grammar

$\langle \text{program} \rangle ::= \langle \text{statements} \rangle$

$\langle \text{statements} \rangle ::= \langle \text{statement} \rangle \mid \langle \text{statement} \rangle \langle \text{statements} \rangle$

$\langle \text{statement} \rangle ::= \langle \text{var_declaration} \rangle \mid \langle \text{expression} \rangle \mid \langle \text{print_statement} \rangle \mid \langle \text{if_statement} \rangle \mid$

$\langle \text{function_declaration} \rangle \mid \langle \text{function_call} \rangle$

$\langle \text{var_declaration} \rangle ::= \text{IDENTIFIER EQUOP} \langle \text{expression} \rangle \mid \text{IDENTIFIER} \langle \text{array_def} \rangle$

$\langle \text{array_def} \rangle ::= \text{LB} \langle \text{expression} \rangle \text{RB}$

$\langle \text{expression} \rangle ::= \langle \text{term} \rangle \mid \langle \text{term} \rangle \text{ PLUS } \langle \text{term} \rangle \mid \langle \text{term} \rangle \text{ MINUS } \langle \text{term} \rangle \mid \langle \text{term} \rangle$
 $\text{STAR } \langle \text{term} \rangle \mid \langle \text{term} \rangle \text{ DIVOP } \langle \text{term} \rangle$

$\langle \text{term} \rangle ::= \text{IDENTIFIER} \mid \text{UNSIGNICON} \mid \text{SIGNICON} \mid \langle \text{expression} \rangle \mid \langle \text{function_call} \rangle$

$\langle \text{print_statement} \rangle ::= \text{DISPLAY } \langle \text{expression_list} \rangle$

$\langle \text{expression_list} \rangle ::= \langle \text{expression} \rangle \mid \langle \text{expression} \rangle \text{ COMMA } \langle \text{expression_list} \rangle$

$\langle \text{if_statement} \rangle ::= \text{IF } \langle \text{condition} \rangle \text{ THEN } \langle \text{statements} \rangle \text{ ENDIF}$

$\langle \text{condition} \rangle ::= \langle \text{expression} \rangle \text{ REOP } \langle \text{expression} \rangle \mid \text{NOT } \langle \text{condition} \rangle$

$\langle \text{function_declaration} \rangle ::= \text{FUNCTION IDENTIFIER parameters IS } \langle \text{statements} \rangle$
 ENDFUN IDENTIFIER

$\langle \text{parameters} \rangle ::= \text{PARAMETERS param_list} \mid \epsilon$

$\langle \text{param_list} \rangle ::= \text{IDENTIFIER OF } \langle \text{data_type} \rangle \mid \text{IDENTIFIER OF } \langle \text{data_type} \rangle$
 $\text{COMMA } \langle \text{param_list} \rangle$

$\langle \text{data_type} \rangle ::= \text{INTEGER} \mid \text{FLOAT} \mid \text{CHAR}$

$\langle \text{function_call} \rangle ::= \text{IDENTIFIER parguments}$

$\langle \text{parguments} \rangle ::= \text{LP } \langle \text{expression_list} \rangle \text{ RP}$

Limitations of the above specification and design of the system

For our scanner, we only used a subset of the grammar rather than the entire grammar because we were given that option. So, because of that, there will be limitations on that end.

Discussion of how the solution can be improved and extended

One way the solution could be improved would be to include all the grammar instead of just the subset. However, since we were given the option, we only did a subset to save us some time on that part.

Source Code Screenshots

scl_scanner.py Source Code:

```

1  import re
2  import json
3
4  #Define Keywords and token types
5  KEYWORDS = ["DISPLAY", "IF", "THEN", "ENDIF", "FUNCTION", "IS", "ENDFUN", "PARAMETERS", "INTEGER", "FLOAT", "CHAR", "NOT"]
6  TOKEN_TYPES = [
7      ("IDENTIFIER", r"\b[A-Za-z_][A-Za-z0-9_]*\b"),          #Alphanumeric Identifier
8      ("UNSIGNICON", r"\b\d+\b"),                             #Unsigned integers
9      ("SIGNICON", r"[-+]?[0-9]+\b"),                          #Signed integers
10     ("PLUS", r"\+"),                                          #Addition Operator
11     ("MINUS", r"\-"),                                         #Subtraction Operator
12     ("STAR", r"\*"),                                          #Multiplication Operator
13     ("DIVOP", r"[/]\b"),                                     #Division Operator
14     ("EQUOP", r"="),                                         #Assignment Operator
15     ("RELOP", r"(==|!=|<=|>=|<|>)"),                        #Relational Operators
16     ("LB", r"\["),                                           #Left Bracket
17     ("RB", r"\]"),                                           #Right Bracket
18     ("LP", r"\("),                                           #Left Parenthesis
19     ("RP", r"\)"),                                           #Right Parenthesis
20     ("COMMA", r","),                                         #Comma
21     ("OF", r"OF")                                            #Of Keyword
22 ]
23
24 #Function to tokenize the source code
25 def tokenize(source_code):
26     tokens = [] #list to store tokens
27     position = 0 #Initialize Position in source code
28     #loop through the source code for tokenization
29     while position < len(source_code):
30         match = None
31         #loop through each token type and its corresponding regex pattern
32         for token_type, pattern in TOKEN_TYPES:
33             regex = re.compile(pattern)
34             match = regex.match(source_code, position)
35             #if match append to tokens list
36             if match:
37                 token_value = match.group(0)
38                 if token_type == "IDENTIFIER" and token_value in KEYWORDS:
39                     tokens.append((token_value, token_value))
40                 else:
41                     tokens.append((token_type, token_value))
42                 position = match.end()
43                 break
44             #if no match is found increment position
45             if not match:
46                 position += 1
47
48     return tokens #returns the list of tokens

```

```

49
50 #Main function to execute the program
51 def main():
52     import sys
53     if len(sys.argv) != 2: #check for the correct number of command line arguments
54         print("Usage: python scl_scanner.py <filename>")
55         return
56     #read source code from file
57     with open(sys.argv[1], 'r') as f:
58         source_code = f.read()
59     #tokenize source code
60     tokens = tokenize(source_code)
61     #print each token
62     for token in tokens:
63         print(token)
64     #Save tokens to JSON
65     with open("tokens.json", "w") as outfile:
66         json.dump(tokens, outfile)
67
68 if __name__ == "__main__":
69     main()
70

```

Source Code Execution Screenshot

tokens.json File Source Code (is also printed in the Terminal):

```
1 [{"IDENTIFIER", "import"}, {"IDENTIFIER", "sysplm"}, {"IDENTIFIER", "h"}, {"IDENTIFIER", "description"}, {"IDENTIFIER", "Program"}, {"IDENTIFIER", "bitops1"}, {"IDENTIFIER", "scl"}, {"IDENTIFIER", "Test"}, {"IDENTIFIER", "the"}, {"IDENTIFIER", "Sysplm"}, {"IDENTIFIER", "with"}, {"IDENTIFIER", "several"}, {"IDENTIFIER", "bitwise"}, {"IDENTIFIER", "operations"}, {"IDENTIFIER", "j"}, {"IDENTIFIER", "Garrido"}, {"IDENTIFIER", "March"}, {"UNSIGNICON", "2018"}, {"STAR", ""}, {"IDENTIFIER", "symbol"}, {"IDENTIFIER", "MM"}, {"DIVOP", " /"}, {"UNSIGNICON", "0110"}, {"UNSIGNICON", "0110"}, {"IDENTIFIER", "symbol"}, {"IDENTIFIER", "MN"}, {"DIVOP", " /"}, {"UNSIGNICON", "1011"}, {"UNSIGNICON", "0011"}, {"IDENTIFIER", "global"}, {"IDENTIFIER", "declarations"}, {"IDENTIFIER", "variables"}, {"IDENTIFIER", "define"}, {"IDENTIFIER", "a"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "type"}, {"IDENTIFIER", "unsigned"}, {"IDENTIFIER", "integer"}, {"IDENTIFIER", "define"}, {"IDENTIFIER", "b"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "type"}, {"IDENTIFIER", "unsigned"}, {"IDENTIFIER", "short"}, {"IDENTIFIER", "define"}, {"IDENTIFIER", "c"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "type"}, {"IDENTIFIER", "unsigned"}, {"IDENTIFIER", "long"}, {"IDENTIFIER", "implementations"}, {"IDENTIFIER", "description"}, {"IDENTIFIER", "this"}, {"IDENTIFIER", "is"}, {"IDENTIFIER", "the"}, {"IDENTIFIER", "main"}, {"IDENTIFIER", "function"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "the"}, {"IDENTIFIER", "application"}, {"STAR", ""}, {"IDENTIFIER", "function"}, {"IDENTIFIER", "main"}, {"IDENTIFIER", "is"}, {"IDENTIFIER", "variables"}, {"IDENTIFIER", "define"}, {"IDENTIFIER", "vara"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "type"}, {"IDENTIFIER", "byte"}, {"IDENTIFIER", "define"}, {"IDENTIFIER", "varb"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "type"}, {"IDENTIFIER", "byte"}, {"IDENTIFIER", "define"}, {"IDENTIFIER", "varc1"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "type"}, {"IDENTIFIER", "byte"}, {"IDENTIFIER", "define"}, {"IDENTIFIER", "varc2"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "type"}, {"IDENTIFIER", "byte"}, {"IDENTIFIER", "define"}, {"IDENTIFIER", "varc3"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "type"}, {"IDENTIFIER", "byte"}, {"IDENTIFIER", "define"}, {"IDENTIFIER", "varc4"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "type"}, {"IDENTIFIER", "byte"}, {"IDENTIFIER", "define"}, {"IDENTIFIER", "d1"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "type"}, {"IDENTIFIER", "byte"}, {"IDENTIFIER", "define"}, {"IDENTIFIER", "d2"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "type"}, {"IDENTIFIER", "byte"}, {"IDENTIFIER", "begin"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "vara"}, {"EQUOP", "="}, {"IDENTIFIER", "MN"}, {"DIVOP", " /"}, {"UNSIGNICON", "0110"}, {"UNSIGNICON", "0110"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "varb"}, {"EQUOP", "="}, {"IDENTIFIER", "MN"}, {"DIVOP", " /"}, {"UNSIGNICON", "1011"}, {"UNSIGNICON", "0011"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "varc1"}, {"EQUOP", "="}, {"IDENTIFIER", "vara"}, {"IDENTIFIER", "band"}, {"IDENTIFIER", "varb"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "varc2"}, {"EQUOP", "="}, {"IDENTIFIER", "vara"}, {"IDENTIFIER", "bor"}, {"IDENTIFIER", "varb"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "varc3"}, {"EQUOP", "="}, {"IDENTIFIER", "vara"}, {"IDENTIFIER", "bxor"}, {"IDENTIFIER", "varb"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "varc4"}, {"EQUOP", "="}, {"IDENTIFIER", "negate"}, {"IDENTIFIER", "vara"}, {"IDENTIFIER", "display"}, {"IDENTIFIER", "varc1"}, {"COMMA", ","}, {"COMMA", ","}, {"IDENTIFIER", "varc2"}, {"COMMA", ","}, {"COMMA", ","}, {"IDENTIFIER", "varc3"}, {"COMMA", ","}, {"COMMA", ","}, {"IDENTIFIER", "varc4"}, {"COMMA", ","}, {"COMMA", ","}, {"IDENTIFIER", "Using"}, {"IDENTIFIER", "a"}, {"IDENTIFIER", "mask"}, {"IDENTIFIER", "to"}, {"IDENTIFIER", "select"}, {"IDENTIFIER", "or"}, {"IDENTIFIER", "alter"}, {"IDENTIFIER", "bits"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "varc1"}, {"EQUOP", "="}, {"IDENTIFIER", "vara"}, {"IDENTIFIER", "band"}, {"DIVOP", " /"}, {"IDENTIFIER", "clear"}, {"IDENTIFIER", "bit"}, {"UNSIGNICON", "0"}, {"LP", "("}, {"IDENTIFIER", "lowest"}, {"IDENTIFIER", "bit"}, {"RP", ")"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "varc2"}, {"EQUOP", "="}, {"IDENTIFIER", "vara"}, {"IDENTIFIER", "band"}, {"DIVOP", " /"}, {"IDENTIFIER", "clean"}, {"IDENTIFIER", "all"}, {"IDENTIFIER", "except"}, {"IDENTIFIER", "bit"}, {"UNSIGNICON", "0"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "varc3"}, {"EQUOP", "="}, {"IDENTIFIER", "vara"}, {"IDENTIFIER", "bor"}, {"DIVOP", " /"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "bit"}, {"UNSIGNICON", "0"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "varc4"}, {"EQUOP", "="}, {"IDENTIFIER", "vara"}, {"IDENTIFIER", "bxor"}, {"DIVOP", " /"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "bit"}, {"UNSIGNICON", "0"}, {"IDENTIFIER", "display"}, {"IDENTIFIER", "varc1"}, {"COMMA", ","}, {"COMMA", ","}, {"IDENTIFIER", "varc2"}, {"COMMA", ","}, {"COMMA", ","}, {"IDENTIFIER", "varc3"}, {"COMMA", ","}, {"COMMA", ","}, {"IDENTIFIER", "varc4"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "d1"}, {"EQUOP", "="}, {"LP", "("}, {"IDENTIFIER", "vara"}, {"IDENTIFIER", "lshift"}, {"UNSIGNICON", "3"}, {"RP", ")"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "d2"}, {"EQUOP", "="}, {"LP", "("}, {"IDENTIFIER", "varb"}, {"IDENTIFIER", "rshift"}, {"UNSIGNICON", "2"}, {"RP", ")"}, {"IDENTIFIER", "display"}, {"IDENTIFIER", "d1"}, {"COMMA", ","}, {"IDENTIFIER", "d1"}, {"COMMA", ","}, {"IDENTIFIER", "d2"}, {"COMMA", ","}, {"IDENTIFIER", "divop", " /"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "a"}, {"EQUOP", "="}, {"IDENTIFIER", "d1"}, {"DIVOP", " /"}, {"IDENTIFIER", "move"}, {"IDENTIFIER", "d1"}, {"IDENTIFIER", "to"}, {"LP", "("}, {"IDENTIFIER", "low"}, {"IDENTIFIER", "byte"}, {"RP", ")"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "display"}, {"IDENTIFIER", "a"}, {"COMMA", ","}, {"IDENTIFIER", "a"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "a"}, {"EQUOP", "="}, {"IDENTIFIER", "a"}, {"IDENTIFIER", "lshift"}, {"UNSIGNICON", "0"}, {"DIVOP", " /"}, {"IDENTIFIER", "shift"}, {"IDENTIFIER", "to"}, {"IDENTIFIER", "high"}, {"IDENTIFIER", "byte"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "a"}, {"IDENTIFIER", "display"}, {"IDENTIFIER", "a"}, {"COMMA", ","}, {"IDENTIFIER", "a"}, {"IDENTIFIER", "set"}, {"IDENTIFIER", "a"}, {"EQUOP", "="}, {"IDENTIFIER", "a"}, {"IDENTIFIER", "bor"}, {"IDENTIFIER", "d2"}, {"DIVOP", " /"}, {"IDENTIFIER", "copy"}, {"IDENTIFIER", "d2"}, {"IDENTIFIER", "to"}, {"IDENTIFIER", "low"}, {"IDENTIFIER", "byte"}, {"IDENTIFIER", "of"}, {"IDENTIFIER", "a"}, {"IDENTIFIER", "display"}, {"IDENTIFIER", "a"}, {"COMMA", ","}, {"IDENTIFIER", "a"}, {"IDENTIFIER", "exit"}, {"IDENTIFIER", "endfun"}, {"IDENTIFIER", "main"}]
```

Comments and Conclusion

This program successfully tokenizes each SCL file and uses the type definitions to correctly categorize each token. They were all stored in the tokens.json, and the categorized tokens were stored in outputs.json. We learned about the basic inner workings of a scanner, and we met every requirement of Deliverable 1. In future deliverables, we will continue our education on how scanners work and develop a fully functional one.

References

Sebesta, R. W. (2012). *Concepts of Programming Languages* (10th ed.). Pearson.