



SINGAPORE UNIVERSITY OF
TECHNOLOGY AND DESIGN

50.051 Programming Language Concepts

Project Final Report

Python-to-C

Nicholas Gandhi (1005295)

Cheh Kang Ming (1005174)

Goh Yu Fan (1005054)

Link to codebase:

<https://github.com/fish-r/Python-To-C>

1.Problem Statement

In the realm of programming languages, interoperability and code portability have become essential aspects of software development. Bridging the gap between Python and C, this project introduces a novel tool designed to compile Python source code into equivalent C code. By leveraging the capabilities of the C programming language, known for its efficiency and performance, this tool aims to provide a pathway for integrating Python code seamlessly into C-based applications.

2.Methodology

The task for the problem statement is as such:

Given an input Python(input.py) file, we should generate the output file(s), output.c and output.h that reflects the execution of the Python code but in C.

Constraints

To achieve this, we have several underlying constraints. These constraints must be observed in order for the programme to be working accurately.

Syntax and Semantics of input Python file

The input python file must be syntactically and semantically correct. This means that the Python file should be able to execute successfully without errors. Additionally, there should be no redeclaration of existing variables.

Syntax and Semantics of output C file

The output C file should be syntactically correct, but not expected to be semantically correct.

Python Keywords Coverage

There will be several keywords that will not be implemented. They are:

1. Global
2. Yield
3. And
4. As
5. Assert
6. Class
7. Del
8. Except
9. Exec
10. Finally
11. From
12. Import
13. Is
14. Lambda
15. Not
16. Or

17. Pass
18. Raise
19. Try
20. With

There will be several operators that will not be implemented. They are:

1. `**`
2. `//`
3. `@`
4. `:=`

There will be several delimiters that will not be implemented. They are:

1. `{`
2. `}`
3. `//`
4. `@=`
5. `**=`

Python Data Types

The input Python file will only have these data types:

1. String
2. Integer
3. Float
4. Lists of a single data type (No type mixing and no nested list)

Python Function Return Data Types

The input Python file can only return literals (excluding lists) but not identifiers (e.g. return “asd”)

Python Programming Patterns

The input Python file will not implement any Object Oriented Programming. The file will also not contain list comprehensions.

3.Implementation

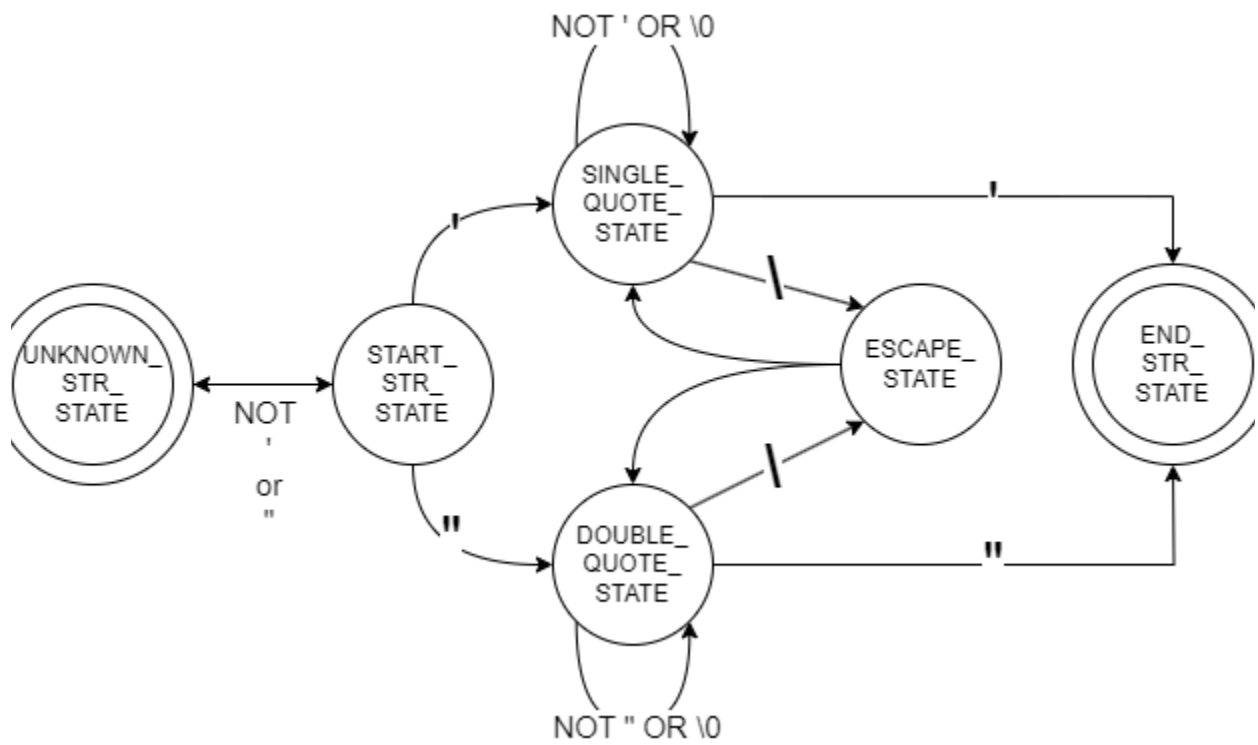
Lexer

The lexer reads the input python code and breaks it down into small, meaningful chunks called tokens.

The lexer has to have the ability to differentiate between literals, python keywords, identifiers, operators, and comments from the given python source code.

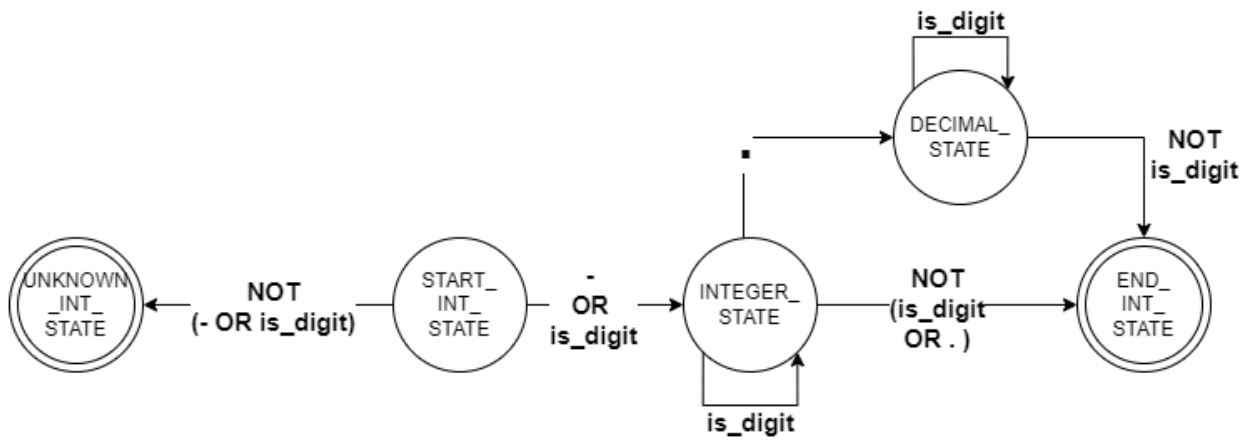
Literal FSM

For the checking of literals, we implemented Finite State Machines (FSMs) to detect whether it is a string literal, integer/float literal and list literals (list_int , list_float, list_string). The FSM diagrams corresponding to the PythonTokenType function are illustrated in the figures below.

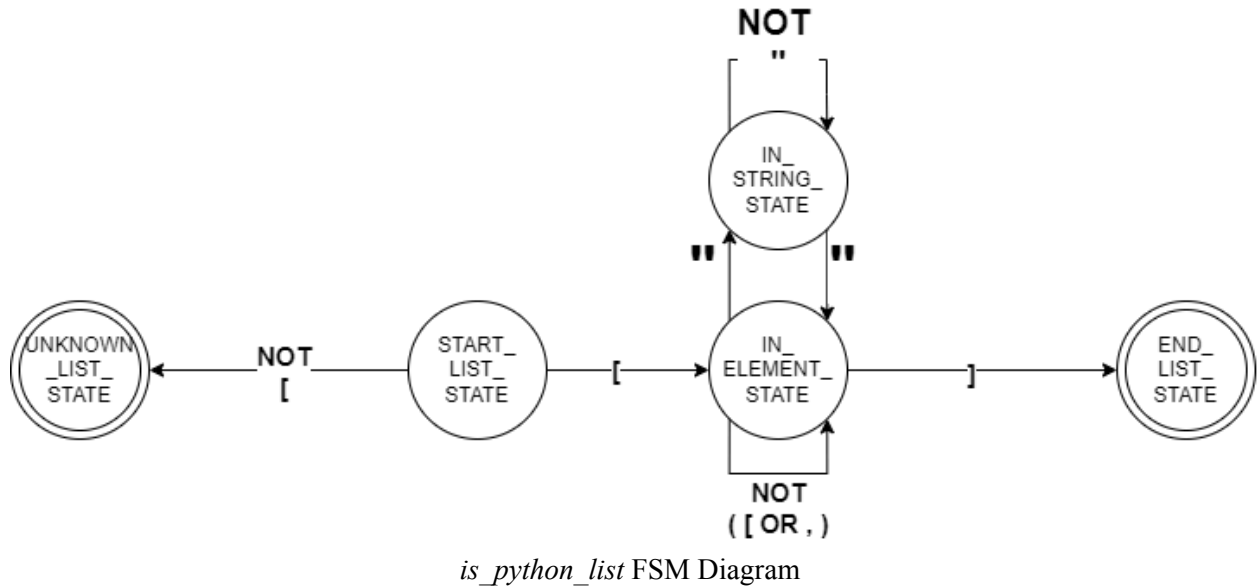


is_python_string FSM Diagram

```
int is_digit(char c) { return c >= '0' && c <= '9'; }
```



is_python_numeric FSM Diagram



Parser

The parser then takes in a stream of python tokens generated by the lexer and builds an Abstract Syntax Tree (AST) based on the type of token. An AST is used in this step to abstract out the syntax of the python code (parenthesis, colons, commas, semi-colons etc), ensuring a smoother code conversion process.

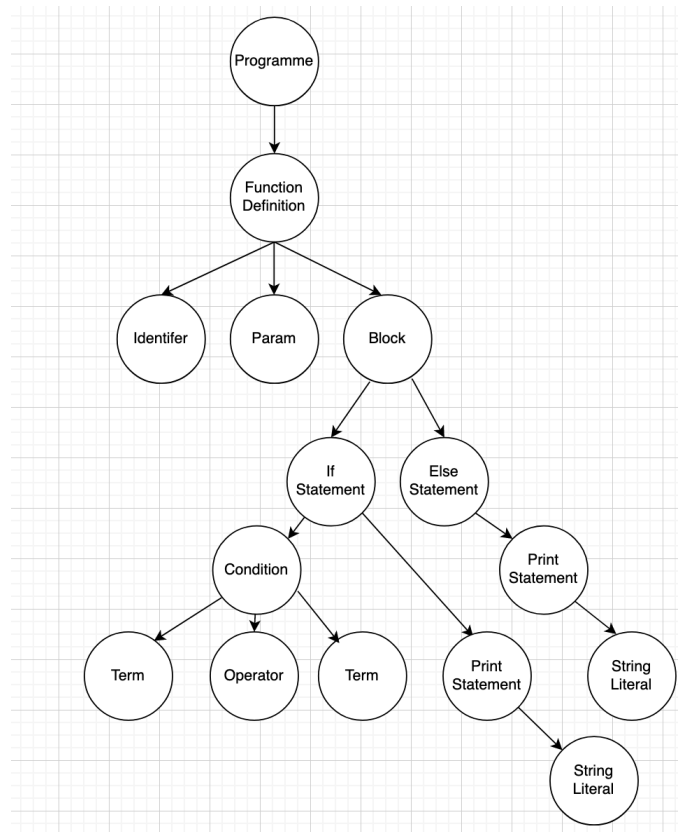
The parser creates different tree nodes of type struct `TreeNode` depending on the type of python token. The type of node then determines how the writer should process and write the corresponding C code to the output file.

The `TreeNode` Struct contains the following variables:

```

6  typedef struct TreeNode
7  {
8      char *label;
9      struct TreeNode *parent;
10     struct TreeNode **children;
11     size_t num_children;
12     Token *token;
13 } TreeNode;
  
```

The children of a Node consist of the essential components required to . For example, the children of a 'IfStatement' node would consist of the 'Condition' node, representing the condition of the if statement.



Sample structure of the AST

C-Writer

The main function of the C Writer (contained in `c_writer.c` file) is to convert the AST into C code. The AST is traversed recursively, per the pseudocode below. There are thus three main components, which are traverse tree, setting state and processing a node.

Traverse Tree

To write the Python code in C, we traversed the AST recursively and performed operations in each traversal, such as setting state and processing the node. This ensures that the C code will respect the scope and order of the input Python code.

```
traverse_tree(TreeNode * root, State *state, NodeQueue * queue, ...):
```

```
    // save nodes in global scope to process later
```

```
    If node in global scope:
```

```
        Enqueue_node( )
```

```
    Return
```

```

// update state based on current node
set_state(root, state, ...)

// process the node and write to file
process_node(root, state, ...)

// traverse and process children recursively
For child in root.children:
    traverse_tree(child, state, queue, ...)

// on recursion exit, close scopes
If node is block:
    Write_to_file( " } ")

// on root node of entire tree, write the main function
if node is root:
    write_main( root, state, queue, ...)

```

Setting State

Firstly, due to the reuse of various `TreeNode` labels (Identifiers, Parameters, EOL etc..), a finite state machine (FSM) is required to keep track of what kind of `TreeNode` is currently being processed, and thus write the respective syntax in C.

For example we would want to write a function definition in C. From the AST, given a parent `TreeNode` with label “FunctionDefinition” with respective child nodes with labels “Identifier” and “Parameter”, the way to write the code in C is to declare the return type of the function and the subsequent return type of the parameters.

However, given a similar case where we want to write a function call, represented by a parent node with label “FunctionCall” and child nodes “Identifier” and “Parameter”, it would be erroneous to declare the return type of the function call (and its parameters) in C.

Therefore, the use of an FSM is required to keep track of what kind of Node is currently being processed. The input for each state is the current state and the current node. In the subsequent `process_node` function, the states will determine what action will be taken (what C code will be written).

```

6  typedef enum {
7      STATE_INIT,
8      WRITE_GLOBAL_VARS,
9      WRITE_FN_DEF,
10     WRITE_IF_STMT,
11     WRITE_CONDITION,
12     WRITE_PRINT_STMT,
13     WRITE_ELSE_STMT,
14     WRITE_RETURN,
15     WRITE_EXPRESSION,
16     WRITE_FOR_LOOP,
17     WRITE_WHILE_LOOP,
18     WRITE_RANGE,
19     WRITE_ITERABLE,
20     WRITE_FN_CALL,
21     WRITE_COMMENT,
22     STATE_END
23 } State;

```

Various States of the c_writer FSM

Process Node

In the `process_node` function, C code will be written depending on what the current state and current node is.

4.Challenges

Through the implementation of this project, we encountered several challenges, some of which we overcame and some that can be improved on. Below are the challenges and the approach we took in the attempt to overcome them.

Closing blocks in C

Python uses indentation to determine which lines of codes belong to the same block, which is a different approach to C that uses curly braces. This then posed a challenge to us on how to close blocks in C as we would not have information on when the block closes when we are traversing the tree.

Identifiers and Return types

In order to print variables in C, we require the type to be declared, which is not required in Python. Hence, we implemented a lookahead to determine the identifier type. This is also done for function definition identifiers, in order to determine the return type.

Printing For Loops

In Python, elements of a List can be iterated over easily, such as in the example: `for i in List`. To convert this into C, we will require not only the identifier *i* but also the Iterable *List*. To achieve this, we initiate temporary nodes which store *i* and *List* to be accessed in the Print function later.

Main Function

Python does not require the use of a main function to execute code existing in the global scope (zero indented code). Thus, this presents a challenge, especially when there may be instances of global code mixed with function declarations in varying orders. In our attempt to solve this, we created a NodeQueue that will enqueue all nodes in the global scope (where parent node is “Program”) and delay processing until the recursion arrives at the root node of the tree at the final traversal.

Here, the nodes in NodeQueue are once again traversed recursively and processed in the same fashion as before to write the code into a main function.

5.Future Improvements

In the future, we aim to enhance our capability to accommodate additional Python data types like dictionaries. Additionally, we aspire to extend our parsing capabilities to encompass more intricate Python functionalities, including list comprehensions, ternary operators, and string formatting.

Appendix

A. Contributions

Goh Yu Fan: *Wrote parts of the lexer, and most of the c_writer.*

Cheh Kang Ming: *Wrote parseTree which creates the AST.*

Nicholas Gandhi: *Wrote the FSMs for parsing literals in lexer and most of lexer and helped in c_writer.*

B. List of Labels Implemented for Abstract Syntax Tree (In no order)

- **FunctionDefinition**
- **FunctionCall**
- **EOL**
- **ForStatement**
- **Identifier**
- **Iterable**
- **Range**
- **Start**
- **Stop**
- **Step**
- **Parameter**
- **PrintStatement**
- **IfStatement**
- **Condition**
- **Term**
- **Operator**
- **ElseStatement**
- **ReturnStatement**
- **Expression**
- **Block**
- **WhileStatement**
- **Char - under Iterable node**
- **String - under Iterable node**
- **Comment**
- **ElifStatement**
- **ListFloat - for both under Iterable and normal list token**
- **ListInt - for both under Iterable and normal list token**
- **ListString - for both under Iterable and normal list token**

C. SAMPLE INPUT AND OUTPUT

```
input > helloworld.py > ...
1 # Integer Literal
2 a = 1
3 # Float Literal
4 b = 2.2
5
6 # String Literal
7 c = "abc"
8 d = 'abc'
9 # Char Literal
10 e = "a"
11 f = 'a'
12 # Empty string
13 g = ""
14 h = ''
15
16 # List Int Literal
17 i = [1, 2, 3]
18 # List Float Literal
19 j = [1.1, 2.5, 3.7]
20 # List String Literal
21 k = ["abc", "def", "ghi"]
22
23
24

C output.c > main()
5 int main() {
6     /* # Integer Literal */
7     int a = 1;
8     /* # Float Literal */
9     float b = 2.2;
10    /* # String Literal */
11    char c[] = "abc";
12    char d[] = "abc";
13    /* # Char Literal */
14    char e = 'a';
15    char f = 'a';
16    /* # Empty string */
17    char g = "";
18    char h = "";
19    /* # List Int Literal */
20    int i[] = {1, 2, 3};
21    /* # List Float Literal */
22    float j[] = {1.1, 2.5, 3.7};
23    /* # List String Literal */
24    char* k[] = {"abc", "def", "ghi"};
25    return 0;
26 }
27
```

Input of different literal types and its respective output

```
input > helloworld.py > ...
1 def print_function(x, y):
2     list = [x, y]
3     for i in list:
4         print(i)
5     for n in range(1, 2, 3):
6         print(n)
7
8
9 for i in range(1, 2):
10     print(i)
11 print_function(1, 2)
12
13 x = 3
14 z = "hel"
15 print(z)

C output.c > ...
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include "output.h"
4
5 void print_function(int x, int y)
6 {
7     int list[] = {x, y};
8     for (int i= 0; i < 1; i++) {
9         printf("%d", list[i]);
10    }
11    for (int n = 1; n < 2; n += 3)
12    {
13        printf("%d", n);
14    }
15 }
16 int main() {
17     for (int i = 1; i < 2; i += 1)
18     {
19         printf("%d", i);
20     }
21     print_function(1, 2);
22     int x = 3;
23     char z[] = "hel";
24     printf(z);
25     return 0;
26 }
27
```

Input of different type of for loop

```
input > helloworld.py > ...
1  def func1():
2      return 1
3  def func2():
4      return "abc"
5  def func3():
6      return 5.5
7  def func4():
8      return 'a'
9
10 func1()
11 func2()
12 func3()
13 func4()

C output.c > ...
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include "output.h"
4
5  int func1()
6  {
7      return 1;
8  }
9  char* func2()
10 {
11     return "abc";
12 }
13 float func3()
14 {
15     return 5.5;
16 }
17 char func4()
18 {
19     return 'a';
20 }
21 int main() {
22     func1();
23     func2();
24     func3();
25     func4();
26     return 0;
27 }
28
```

Different return type