

# Lexical Analyzer for the C Minus Language

Victor Emmanuel Guerra Aguado - A01568075

## Summary

The following document will report upon the development of a lexical analyzer for the proposed C Minus programming language in the Python programming language, including in the complete phases of **Analysis**, **Design**, **Implementation**, and **Testing**. Each phase will contain its appropriate documentation and implementation.

## Context & Notation

### Regular Expressions (RegEx, or REs)

Regular Expressions are a special notation that represent patterns of strings or characters, and facilitate the task of consistent and precise representation of the lexical rules of a language.[1]

A language is composed by a set of strings, and the strings are formed from a set of characters, these characters are the **legal symbols** or **alphabet** of the language and are represented by the Greek letter  $\Sigma$  (sigma).[1]

### Finite State Machines(FSMs)

Also known as Finite Automata (FAs), they are a mathematical way of describing algorithms for behavioral systems. Finite State Machines can be used to describe the process of recognizing patterns in input strings, and can therefore be used for a lexical analyzer. They are strongly related to Regular Expressions; an FSM can always be constructed from a valid Regular Expression.[1]

In a Finite State Machine, the circles represent states which are locations in the process of recognition that record how much of the pattern has been recognized; the arrowed lines represent the transitions that record the change from one state (initial) to another (final), it is labeled with the character or characters that trigger said transition. Accepting states are

indicated by a circle with double lined border, and they signify the recognition of a valid string.[1]

A Deterministic Finite Automaton (DFA) is one where any final state is **uniquely** given by the current initial state and the current input character of its transition.[1]

### **Transition Tables**

A DFA can be expressed as a data structure together with a generic code that will use the information stored in order to take proper action. The transition table is one of such data structures; implemented as a two dimensional array, the transition table is indexed by a state number and an input character, which provide the values for the transition function.[1]

### **Programming Language**

The chosen programming language for the compiler and lexical analyzer is Python. It was chosen for its simplicity and straightforwardness regarding the management of data structures, memory, and data types; its flexibility would make it possible to take some load from the implementation phase that could be used during other phases.

## **Analysis**

### **Requirements**

1. The automata of the language.
  - a. The automaton of the language shall categorize each word read according to the specification of the C Minus Language.
  - b. The system shall recognize comments, as defined by the C Minus language specification and ignore their contents.
2. Tokens and their identification.
  - a. The system shall receive an input file to scan.
  - b. The system shall recognize every given valid word according to the C Minus Language Specification.
  - c. The system shall create a token for every valid word recognized in the input file.

- d. The system shall recognize identifiers and numerical constants as defined by the C Minus language specification and create a token for each found in the input file.
  - e. The system shall recognize the keywords and special symbols of the C Minus language and create a special token for each found in the input file.
  - f. The system shall create tokens from the input file completely, correctly, and in the order recognized, according to the C Minus language specification.
  - g. The system shall create and output an ordered token list for the valid words of the C Minus language found in the input file.
3. Transition Table.
- a. The transition table shall represent completely and correctly the definition of the Deterministic Finite Automaton for the C Minus Language.
  - b. The transition table shall be implemented in code, in order to be able to apply the language DFA in the compiler.
4. Symbol Tables.
- a. The system shall create and use different Symbol Tables for Identifiers and Numeric Constants.
  - b. The system shall create an entry in a Symbol Table for every unique, valid identifier or numeric constant found in the input file.
  - c. The Symbol Tables shall be able to handle the addition of columns when necessary.
  - d. The Symbol Tables shall be able to handle entry Id and Value.
  - e. The system shall output all Symbol Tables created and required.
5. Error Messages generated by the scanner.
- a. The automaton and scanner shall be able to recognize an invalid character and produce an error appropriate to the situation.
  - b. The automaton and scanner shall be able to recognize an open comment symbol with no corresponding closing symbol, and produce an appropriate error.
  - c. The automaton and scanner shall be able to recognize an invalid Identifier or Keyword and produce an appropriate error.
  - d. The automaton and scanner shall be able to recognize an invalid Numerical Constant and produce an appropriate error.

- e. The scanner shall be able to provide the number of the line in code where an error was found.

### The Definition

In order to create the necessary transition table for recognizing valid words of the C Minus language, it is useful to first define the DFA that accepts all of the RegEx, keywords, and symbols shown in the definition of the language. The RegEx and valid words and symbols for the Language are defined as follows:

- a) The Keywords of the language are:

**else**  
**if**  
**int**  
**return**  
**void**  
**while**  
**input**  
**output**

All keywords are reserved words and they are NOT case sensitive [ ... ] .

- b) Special symbols are the following:

<b>!=</b>	logic operator different
<b>=</b>	assignation
<b>;</b>	semicolon
<b>,</b>	coma
<b>(</b>	open parenthesis
<b>)</b>	close parenthesis
<b>[</b>	open square brackets
<b>]</b>	close square brackets
<b>{</b>	open curly brackets
<b>}</b>	close curly brackets
<b>/*</b>	open comment
<b>*/</b>	close comment
<b>+</b>	arithmetic addition operation
<b>-</b>	arithmetic subtraction operation
<b>*</b>	arithmetic multiplication operation
<b>/</b>	arithmetic division operation
<b>&lt;</b>	logic operator less than
<b>&lt;=</b>	logic operator less or equal than
<b>&gt;</b>	logic operator greater than
<b>&gt;=</b>	logic operator greater or equal than
<b>==</b>	logic operator equal

- c) Other tokens are **ID** and **NUM**, their corresponding Regular Expressions definitions are as follows:

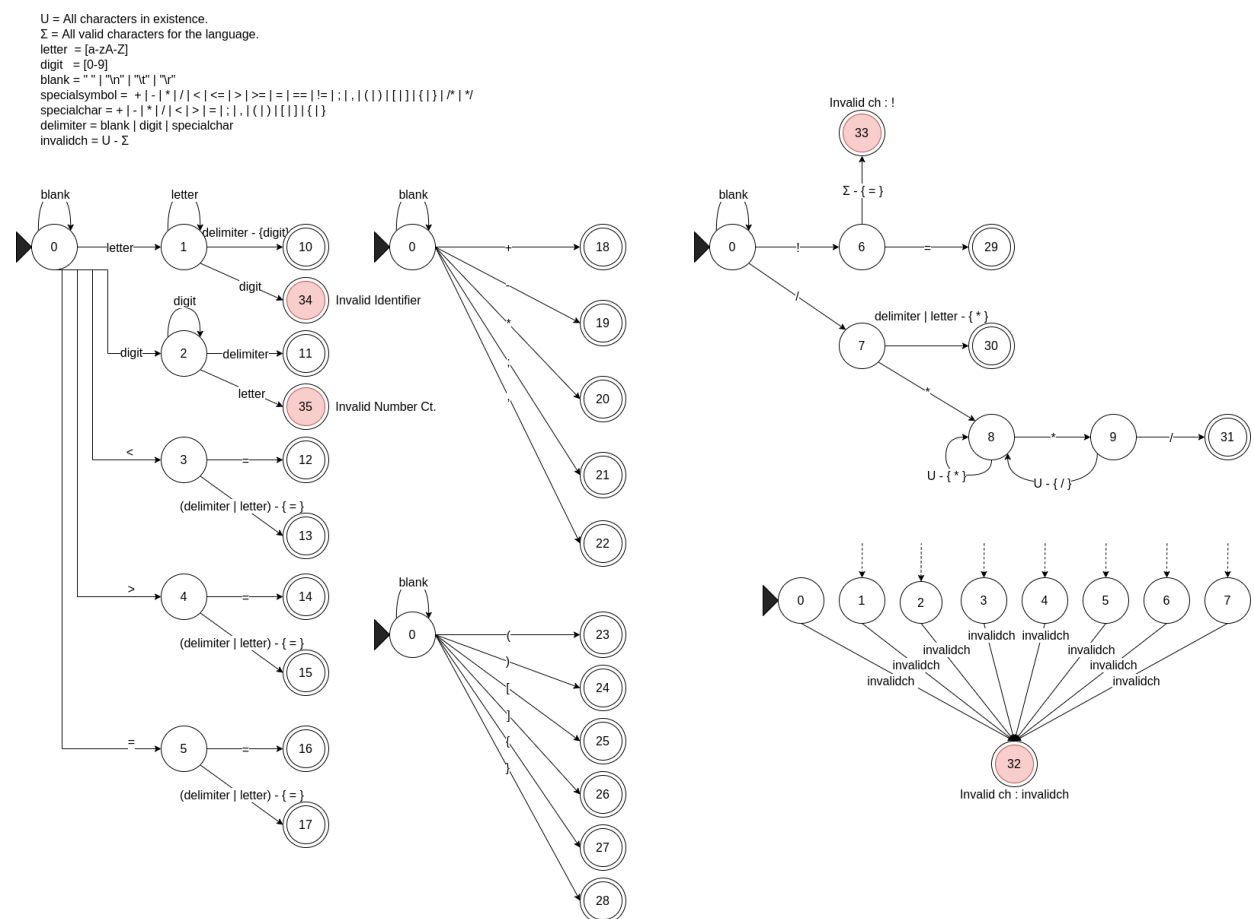
**ID** = *letter*<sup>+</sup>  
**NUM** = *digit*<sup>+</sup>  
*letter* = [a-zA-Z]  
*digit* = [0-9]

Identifiers are letter sensitive, i.e., lower and uppercase letters are distinct.

d) White space consists of *blanks*, *newlines*, and *tabs*. Whitespace is ignored, but it **MUST** be recognized. White space together with **ID**'s, **NUM**'s, and **keywords**, are considered as delimiters.

e) Comments are **C** language style, i.e., they are enclosed by `/* ... */`. Comments can be placed anywhere white space can appear, i.e., comments cannot be placed within tokens. Comments may include more than one line.

Given these specifications the DFA in **Figure 1.1** was designed:



**Figure 1.1 The Whole DFA for the C Minus Language**

The definitions given on the top left of **Figure 1.1** are as follows:

U = All characters in existence.

$\Sigma$  = All valid characters for the C Minus Language.

$letter = [a-zA-Z]$

$digit = [0-9]$

$blank = " " | "\n" | "\t" | "\r"$

$specialchar = + | - | * | / | < | > | = | ; | , | ( | ) | [ | ] | \{ | \}$

$specialsymbol = specialchar | <= | >= | == | != | /* | */$

$delimiter = blank | digit | specialchar$

$invalidch = U - \Sigma$

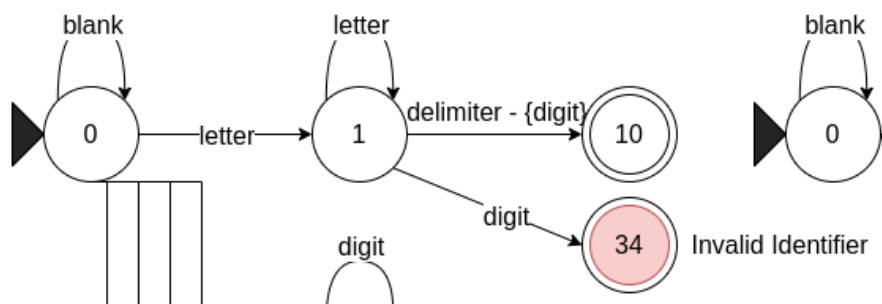
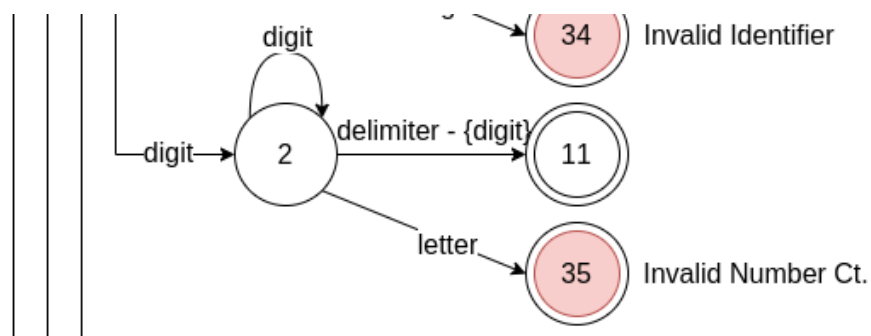


Figure 1.2 The keyword and identifier recognition of the DFA for the C Minus Language.

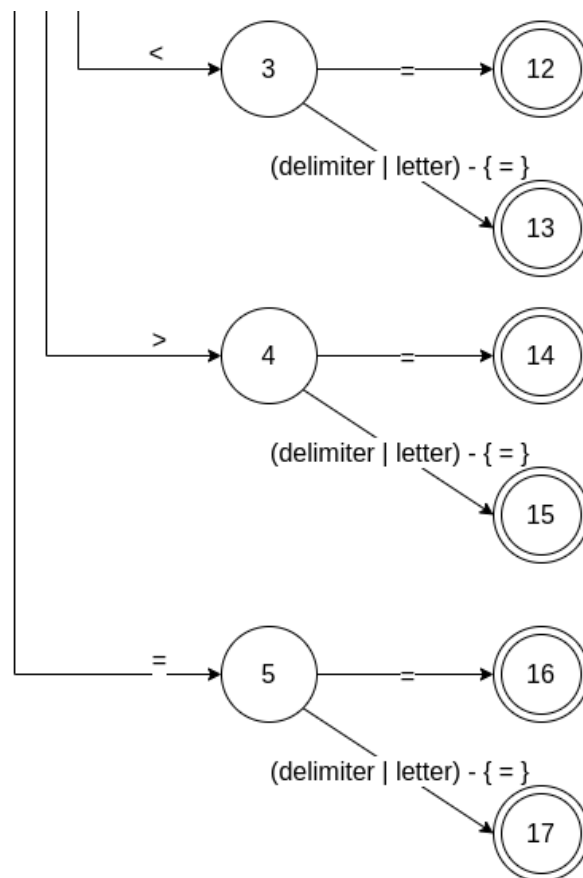
State 0 of the DFA, shown in **Figure 1.2**, is where the lexical recognition begins. It loops back to itself whenever it encounters a BLANK, and starts recognition of all other valid words in the language. It has transitions to states: 1, 2, 3, 4, 5, 6, 7, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27 and 28.

States 1, 10, and 34 of the DFA, shown in **Figure 1.2**, are in charge of recognizing valid *keywords* and *Identifiers*; state 1 loops for every *letter* found, whenever a non-digit delimiter is found the automaton changes its state to 10, an accepting state, and finishes recognition; if, while in state 1, it receives a *digit* the DFA changes to accepting state 34, thus indicating that an invalid *identifier* was found.



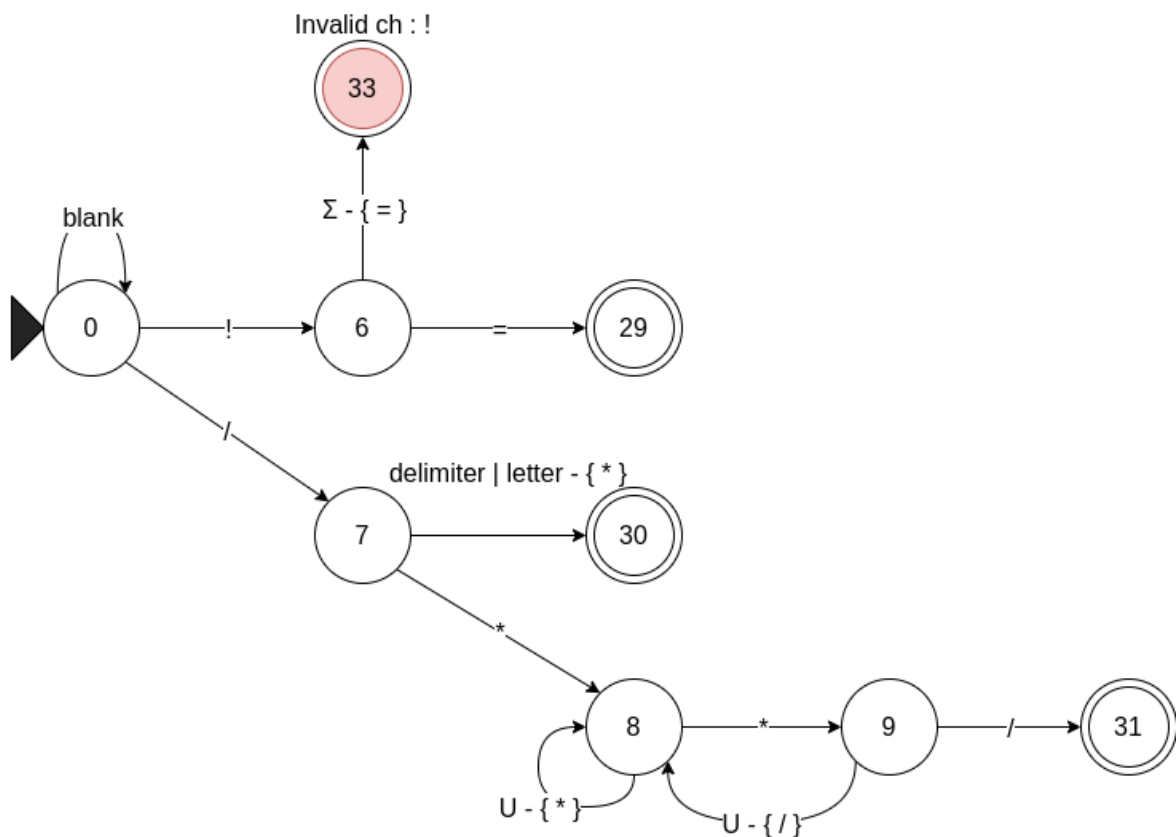
**Figure 1.3 The numerical constant recognition of the DFA for the C Minus Language.**

States 2, 11, and 35 of the DFA, shown in **Figure 1.3**, are concerned with recognizing *Numerical Constants*; state 2 loops for every *digit* found, whenever a non-digit delimiter is found the automaton changes its state to 11, an accepting state, and finishes recognition; if, while in state 2, it receives a *letter* the DFA changes to accepting error state 35, thus indicating that an invalid *identifier* was found.



**Figure 1.4 Recognition of Compound Symbols of the DFA for the C Minus Language.**

States 3, 4, and 5, shown in **Figure 1.4**, are concerned with the recognition of Compound Symbols: <=, >=, and ==; Upon receiving one of the symbols: <, >, or = ( all of which are valid symbols in the C Minus Language ) the automaton looks at the following character, if it is anything but a "=", the automaton moves to accepting state 12, 14, or 16 respectively, and indicates recognition of one of the single-character-symbols shown previously; otherwise, the automaton moves to accepting state 13, 15, or 17, respectively, and indicates the recognition of one of the compound symbols shown previously.

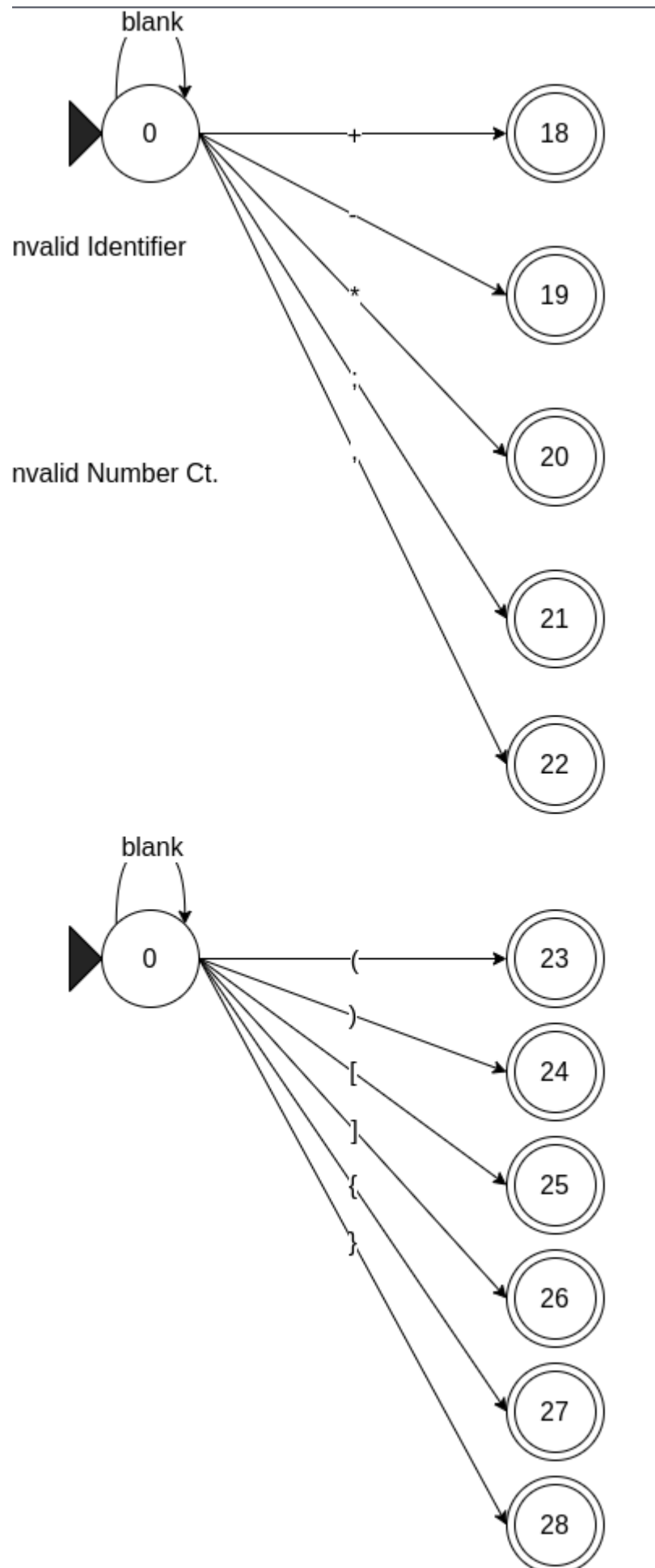


**Figure 1.5 Other compound symbols and comment recognition.**

State 6, shown in **Figure 1.5**, is in charge of recognizing a “!=” symbol, however, since the exclamation point by itself is not a valid word in the C Minus language if it is encountered, an error should be raised.

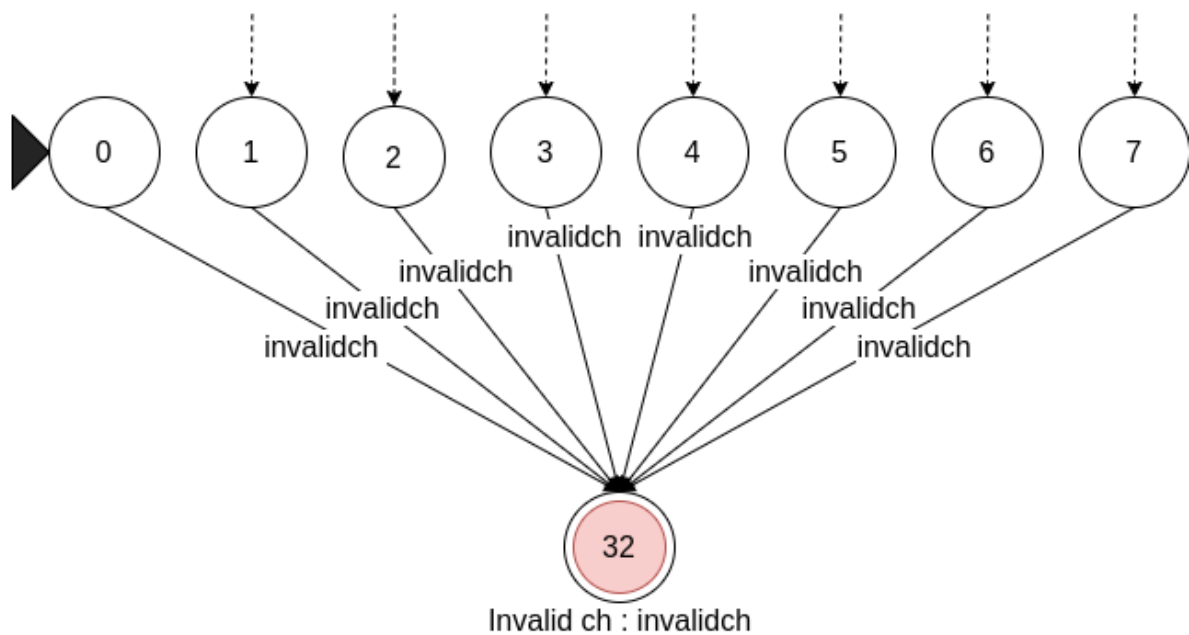
State 7, in **Figure 1.5**, recognizes the “/” symbol, but it also has a very important role with the recognition of comments; if the forward slash (“/”) is followed by anything but an asterisk (“\*”) the forward slash is then recognized as one of the valid words of the language, but if it is followed by an asterisk then a comment has been opened, and states 8 and 9 have the responsibility of looking for the closing tag for said comment. If the closing tag is found, then all collected characters are ignored, and the accepting state of 31 is reached, ending the iteration, otherwise it is an error that shall be notified to the user.





**Figure 1.6 States 18 through to 28.**

States 18 through 28, shown in **Figure 1.6**, provide the recognition of all remaining one-character symbols of the language. The recognition is simple since it only requires one transition, from initial state to accepting state; after reading one of the following characters: “+”, “-”, “\*”, “;”, “:”, “(”, “)”, “[“, “]”, “{“, “}”, the state is accepting, thus, there is little to no opportunity for errors to arise during these recognitions.



**Figure 1.7 The Error State 32, for invalid characters.**

Finally, states 32 33, 34, and 35, 32 pictured in **Figure 1.7** and the rest can be seen in **Figures 1.5, 1.3, and 1.2** respectively, are the error states. State is a common error state to all non-accepting states of the DFA; it recognizes when a given character is not a valid character according to the C Minus language specification, that is any character outside of the 26 letters of the alphabet (uppercase and lowercase), the 10 digits, and any of the following: “+”, “-”, “\*”, “/”, “<”, “>”, “=”, “;”, “:”, “(”, “)”, “[“, “]”, “{“, “}”. Any other character is outside the definition of the language and, therefore, an error in lexical analysis.

## Transition Table

The Following considerations were made, in order to best translate the DFA into a Transition Table:

1. The DFA shall be expressed as a transition table data structure.
2. The ID's and keywords are going to be differentiated within the code, not the DFA.
3. Non-accepting states were numbered first, followed by Accepting States, and Error Accepting States, in order to better organize the Transition Table.

Given the generated DFA and taking the previous considerations. the following Transition Table was generated:

STATE	ACCEPTING	LETTER	DIGIT	BLANK	<	>	=	!	/	+	-	*	;	,	(	)	[	]	{	}	INVALIDCH
0	FALSE	1	2	0	3	4	5	6	7	18	19	20	21	22	23	24	25	26	27	28	32
1	FALSE	1 [34]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	[10]	32
2	FALSE	[35]	2 [11]	[11]	[11]	[11]	[11]	[11]	[11]	[11]	[11]	[11]	[11]	[11]	[11]	[11]	[11]	[11]	[11]	[11]	32
3	FALSE	[13]	[13]	[13]	[13]	[13]	12 [13]	[13]	[13]	[13]	[13]	[13]	[13]	[13]	[13]	[13]	[13]	[13]	[13]	[13]	32
4	FALSE	[15]	[15]	[15]	[15]	[15]	14 [15]	[15]	[15]	[15]	[15]	[15]	[15]	[15]	[15]	[15]	[15]	[15]	[15]	[15]	32
5	FALSE	[17]	[17]	[17]	[17]	[17]	16 [17]	[17]	[17]	[17]	[17]	[17]	[17]	[17]	[17]	[17]	[17]	[17]	[17]	[17]	32
6	FALSE	33	33	33	33	33	29 33	33	33	33	33	33	33	33	33	33	33	33	33	33	32
7	FALSE	[30]	[30]	[30]	[30]	[30]	[30]	[30]	[30]	[30]	[30]	8 [30]	[30]	[30]	[30]	[30]	[30]	[30]	[30]	[30]	32
8	FALSE	8	8	8	8	8	8	8	8	8	8	8	9	8	8	8	8	8	8	8	8
9	FALSE	8	8	8	8	8	8	8	8	31	8	8	8	8	8	8	8	8	8	8	8
10	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
11	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
12	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
13	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
14	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
15	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
16	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
17	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
18	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
19	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
20	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
21	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
22	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
23	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
24	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
25	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
26	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
27	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
28	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
29	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
30	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
31	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
32	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
33	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
34	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
35	TRUE	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	

**Figure 1.8 The Transition Table from the C Minus DFA.**

In **Figure 1.8**, the rows in yellow correspond to the non-accepting states, the rows in green to accepting states, and the rows in orange to accepting error states. The first column determines whether the state is an accepting state, however, since this column does not represent a transition, it shall not be reflected in the transition table itself, but as a different data point.

Accepting states have no specified transition since once the automaton is at one it shall stop execution and return its state and recognized string, therefore there is no need for an accepting state to have a transition.

In addition, the transitions that are displayed surrounded by square braces (“[“, “]”) are the **non-consuming** transitions, that is when triggered, these transitions will not “consume” the character that was read.

### Tokens

The valid words of the language shall be organized into the following structure:

1. else
2. if
3. input
4. int
5. output
6. return
7. void
8. while
9. +
10. -
11. \*
12. /
13. >
14. >=
15. <
16. <=
17. ==
18. !=
19. =
20. ;
21. ,
22. (
23. )
24. {
25. }

- 26. [
- 27. ]
- 28. IDENTIFIER
- 29. NUM\_CONSTANT
- 30. COMMENT

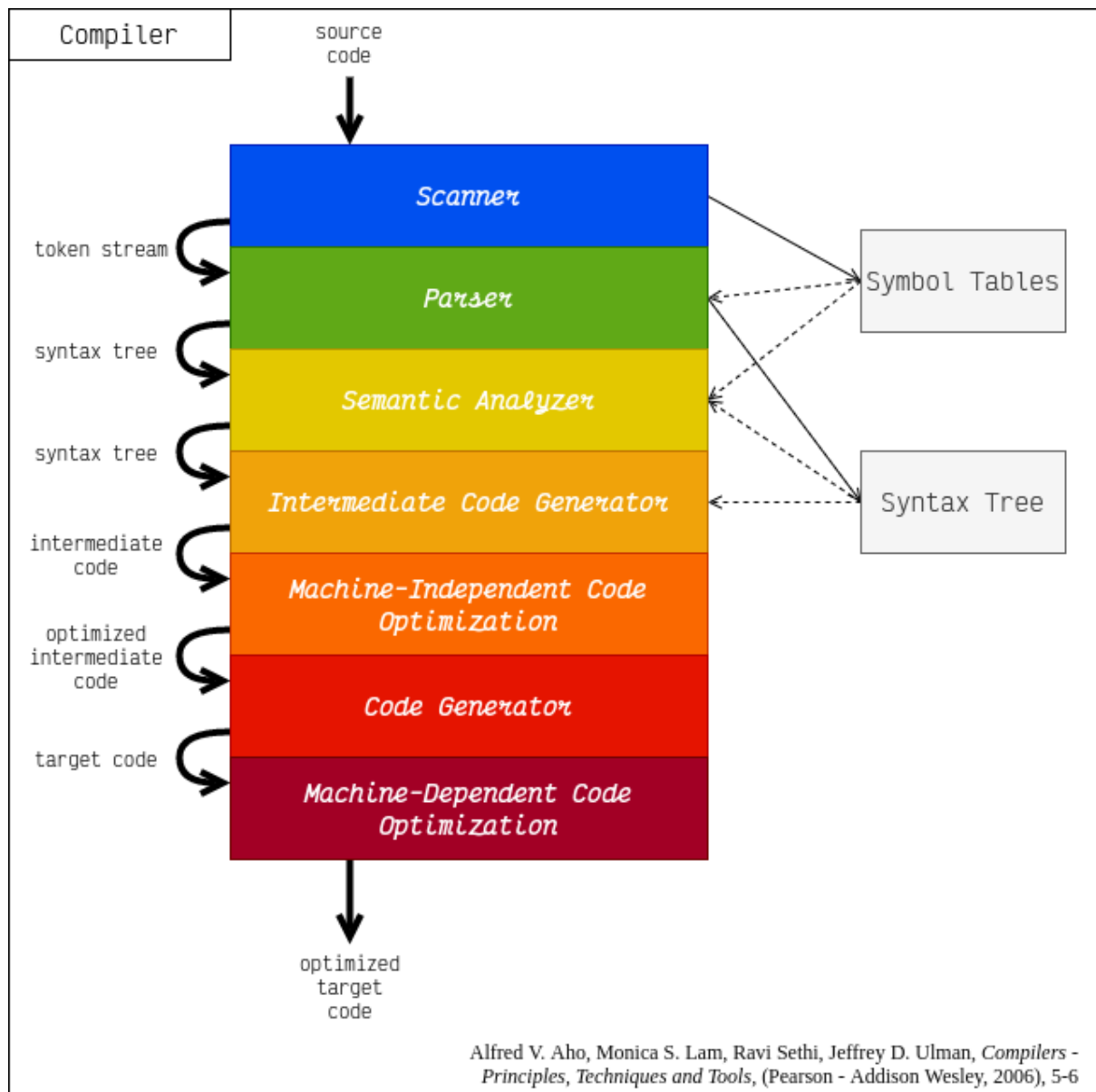
Where the number represents the id that said word will receive, and will identify it during the rest of the compilation process. The order was decided upon alphabetical order for the keywords; as presented in the specification for the special symbols (exceptuating the last four); and as seen in class[1] for Identifiers, Numerical Constants, and Comments.

## Design

During the design phase, multiple diagrams were made, including: Class Diagram, Sequence Diagram, and a Flow Chart, along with a preliminary diagram of the architecture of the entire Compiler. With respect to the Requirements, the following considerations were made:

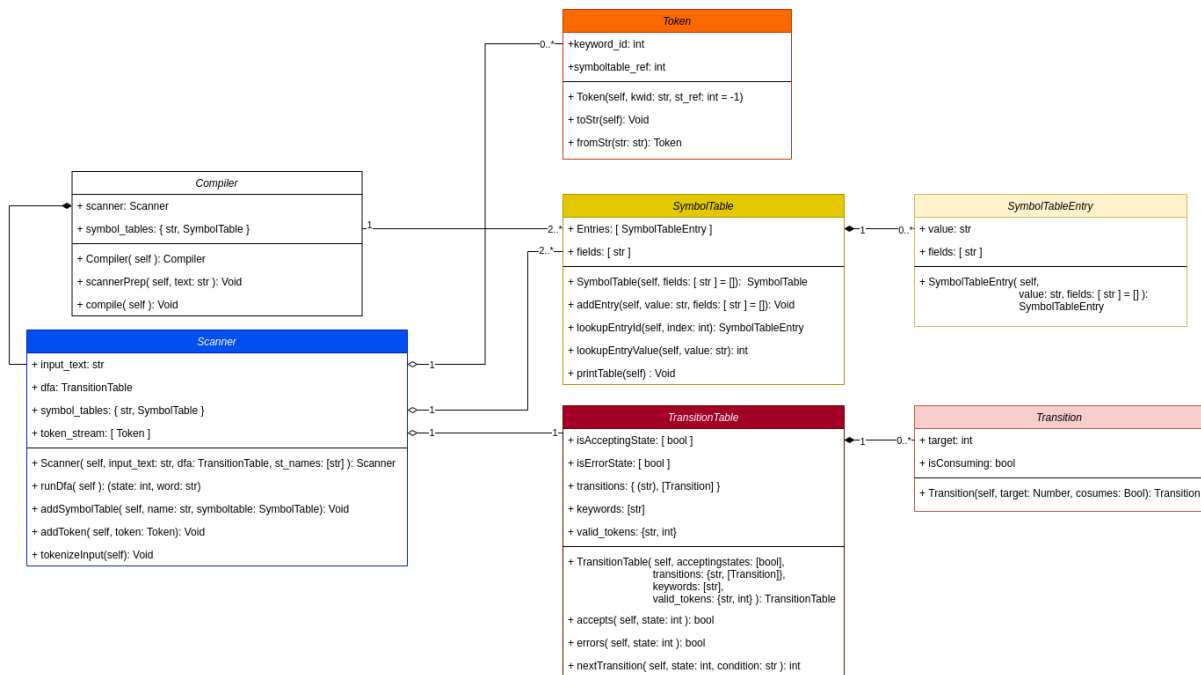
1. The entire scanner for the compiler can be abstracted into a class.
2. The Transition Table can be abstracted into a class, composed of transitions, which themselves are put into arrays with size adequate to index by state number; these arrays will be indexed, with a tuple of characters serving as key, within a dictionary.
3. The Symbol Tables can be abstracted into a class, composed of entries abstracted into classes, with the attributes for their value and any possible fields that might be required.
4. The Symbol Tables will have an attribute for the possible additional fields as an array of the values as required; said array will be expandable as needed.
5. The Tokens that are to be the output of the Scanner can be abstracted into a class, leaving them as a tuple was considered, but the necessary methods for adapting and manipulating these tokens demand a more complicated data structure.
6. The Compiler itself can be abstracted into a class, with its multiple phases abstracted into classes and added as attributes, which can then be called in the order necessary.

A diagram of a preliminary architecture for the entire compiler, seen in **Figure 2.1**, was made with the intention of establishing a structure in which the steps of compilation and their inputs and outputs can be visualized and the shared data that each class utilizes; therefore more easily placed within the *compile* method:



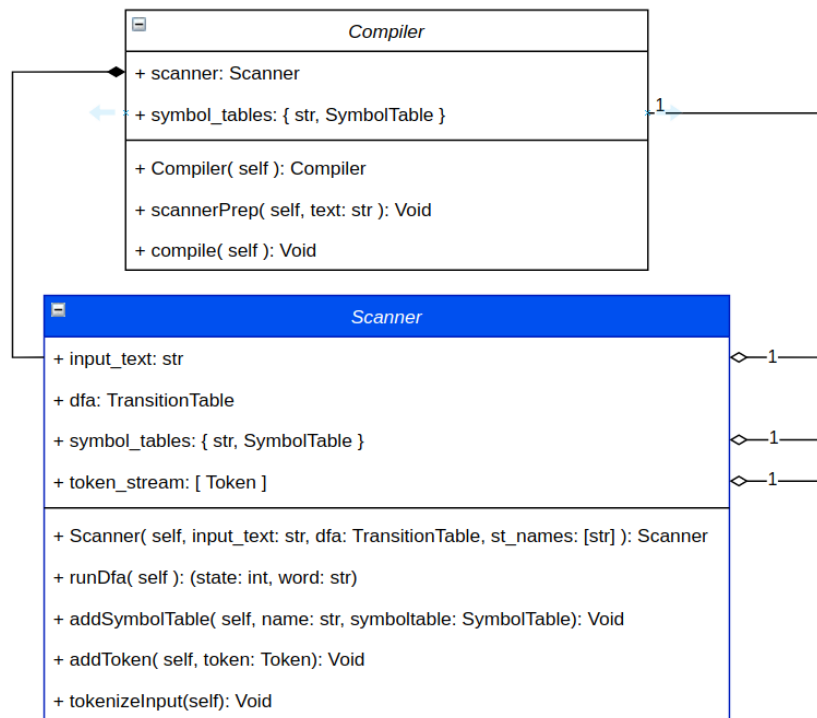
**Figure 2.1** The preliminary architecture diagram of the compiler.

The diagram in **Figure 2.2** shall illustrate the interactions between the classes described previously; the classes for steps of compilation other than lexical analysis have been omitted as they are not relevant at this stage of compilation.



**Figure 2.2 - The entire Class Diagram for the considerations described.**

As considered, the *Compiler* class is composed of a *Scanner* object, as seen in **Figure 2.3**, and will be composed of a *Parser* object, a *SemanticAnalyzer* class, etc. It also has a ‘symbol\_tables’ attribute which incorporates all the resulting *SymbolTable* objects from the Lexical Analysis phase. Its only methods so far are the *scannerPrep* method, which sets up the scanner with the necessary information; and the *compile* method, which starts the compilation process by calling the scanner to produce its output.



**Figure 2.3 - The Compiler Class and the Scanner Class.**

The *Scanner* Class itself, as seen in **Figure 2.3**, is aggregated from the *Token*, *SymbolTable*, and *TransitionTable* classes, and incorporates them as the attributes *token\_stream*, *symbol\_tables*, and *dfa* respectively, where *token\_stream* represents the output of the Scanner as a list of *Token* objects; *symbol\_tables* stores the scanners' *SymbolTable* objects by name via a dictionary structure; and *dfa* stores the C Minus Language specifications as a *TransitionTable* object. Its methods include: *runDfa*, which runs the "DFA" through its transition table exactly once; *addSymbolTable*, which receives a string and a *SymbolTable* object, and saves them within its *symbol\_tables* attribute with the string as a key to access the object in the dictionary; *addToken*, which appends the given *Token* object to the *token\_stream* attribute; and *tokenizeInput*, which executes the whole scanning phase and fills the *token\_stream* and *symbol\_tables* attributes with their respective entries.



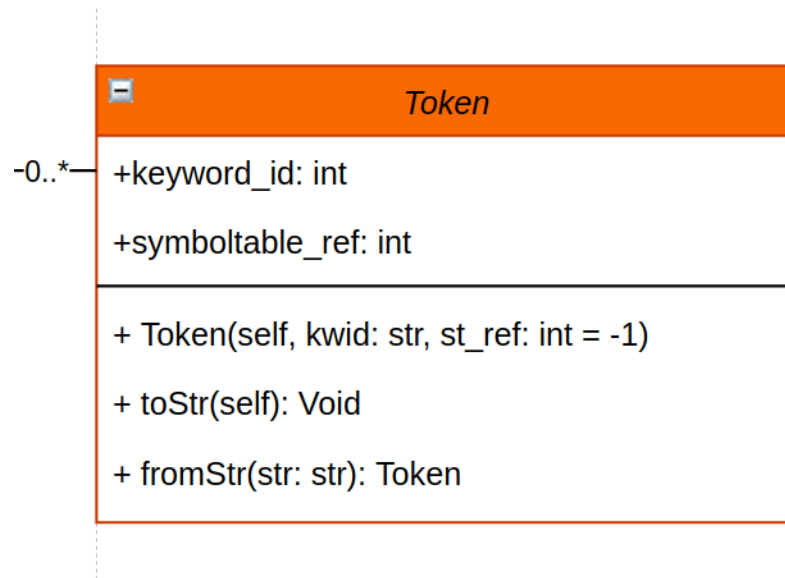


Figure 2.4 - The Token Class

The *Token* Class, seen in **Figure 2.4**, abstracts the components of the tokens generated by a compilers' scanner, and provides some relevant methods for input and output of said tokens.

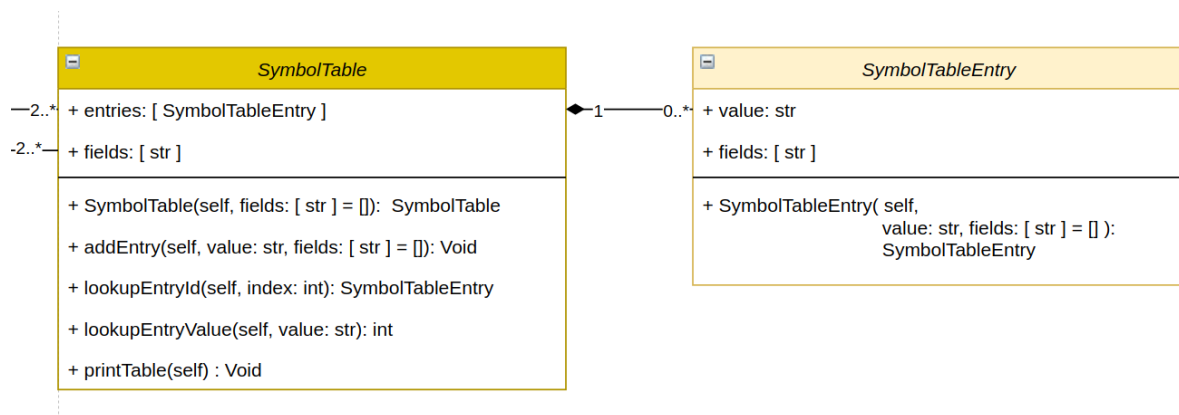
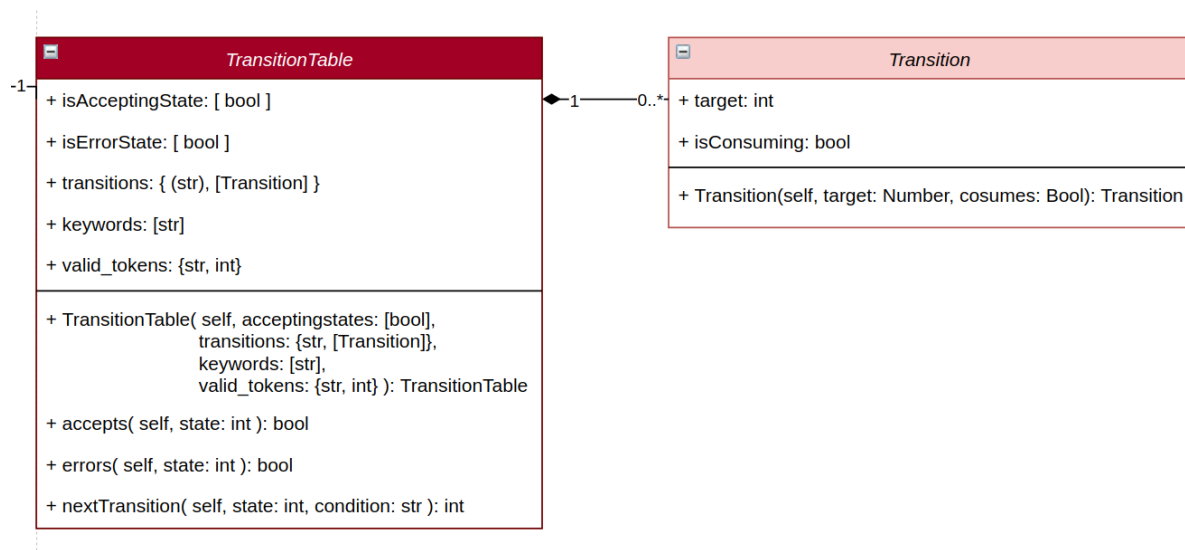


Figure 2.5 - The SymbolTable and SymbolTableEntry Classes.

The *SymbolTable* Class, seen in **Figure 2.5**, acts as an abstraction of the symbol table structure used during the compilation process. It is composed of a list of *SymbolTableEntry* Objects, which represent the unique identifiers or numerical constants found within the source code. Its methods include: *addEntry*, which appends an entry to the *entries*; *lookupEntryId*, which returns an entry given an index; *lookupEntryValue*, which returns an entry given a value (this works since every entry is supposed to have a unique *value*); and *printTable*, which provides visualization of the tables' contents.



**Figure 2.6 - The TransitionTable and Transition Classes.**

The *TransitionTable* Class, as seen in **Figure 2.6**, works as a container for the definition of a language, it abstracts the components of a DFA and language specifications necessary for the lexical analysis. It contains the attributes: *isAcceptingState*, which is an array of boolean values intended to list all states of a DFA, and differentiate those that are accepting states; *isErrorstate* is similar to the *isAcceptingState* attribute, instead of differentiating accepting states, it differentiates the error states; the *keywords* attribute is a list of all the reserved words, or keywords, of a given language, the attribute is necessary due to the adjustments made to the DFA of the C Minus Language, namely, the differentiation between keywords and identifiers being made in code, as opposed to the DFA; *valid\_tokens*, which lists the valid words of the language and assigns them an ID number by which the tokens created by the scanner will refer to; finally the *transitions* attribute is a dictionary with tuples of strings as keys, and *Transition* objects as values.

The *Transition* object acts as a value object for the transitions used in a transition table, and provides no additional functionality.

The sequence diagram in **Figure 2.7** was made to project the path of function calls along the instances of the previously illustrated object classes for the scanning process. It only takes into account the classes: *Compiler*, *Scanner*, *TransitionTable*, *SymbolTable*, and *Token* since the other classes act just as structures for the data in the shown classes:

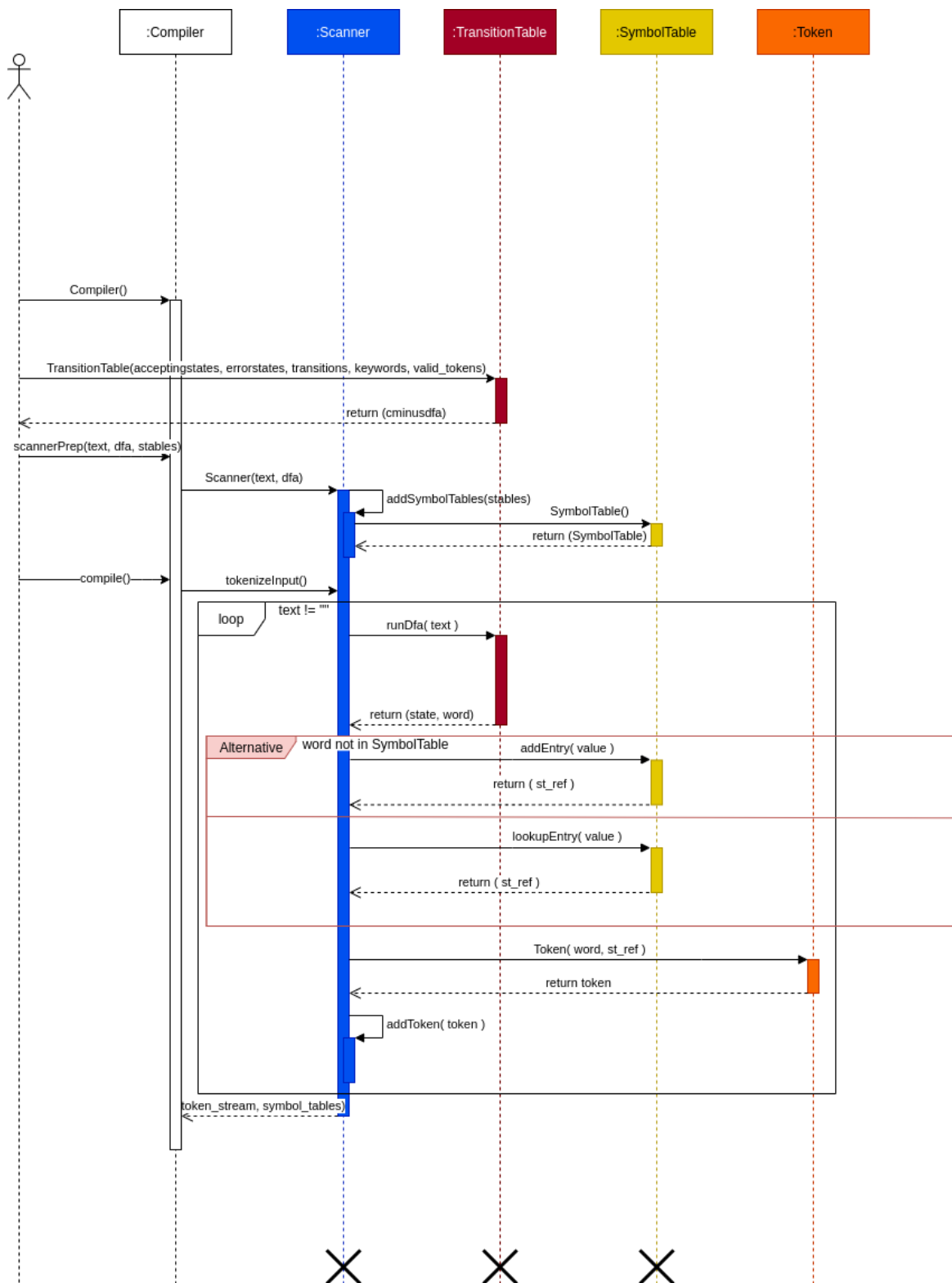
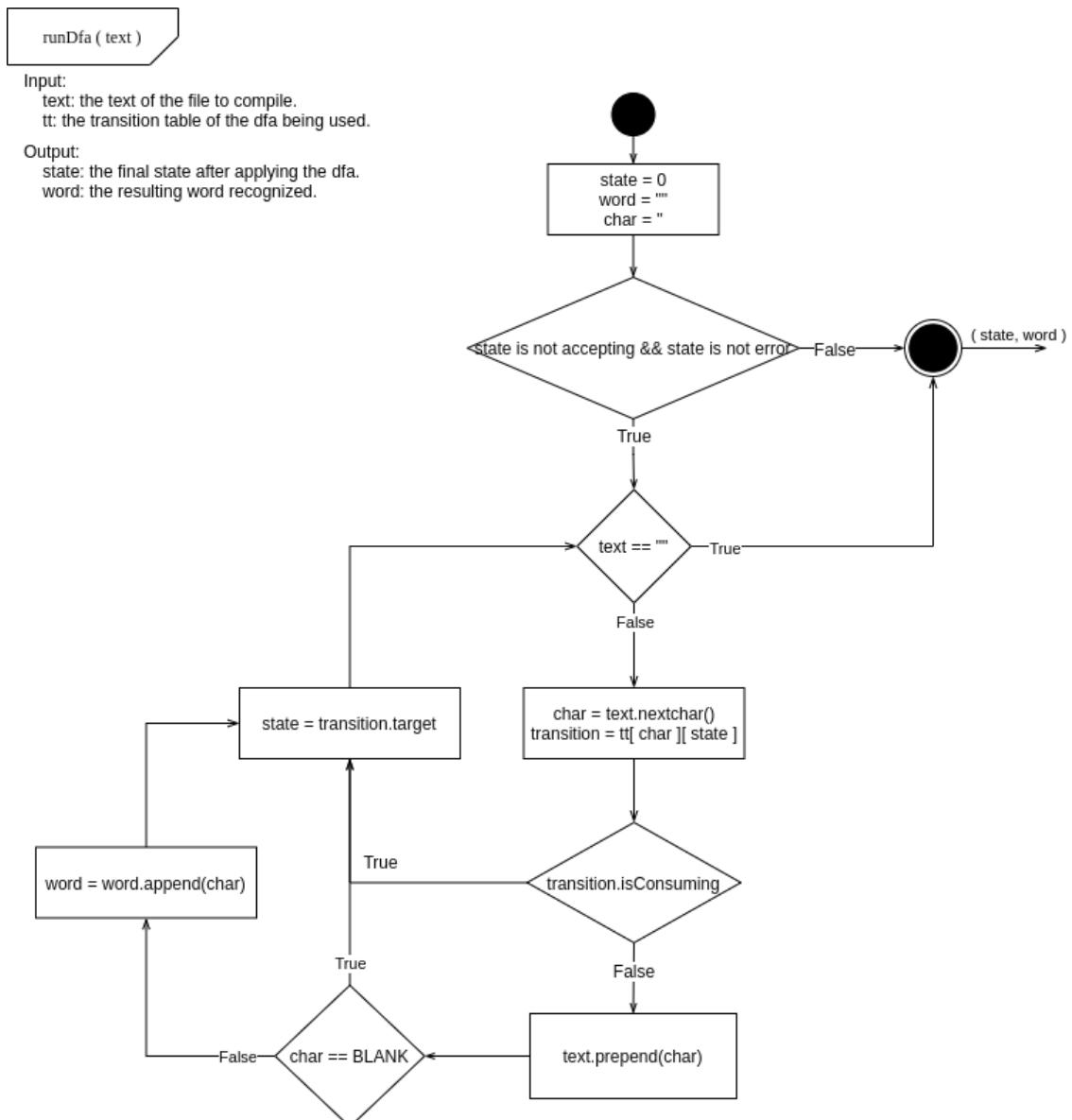


Figure 2.7 - The Sequence Diagram.

The diagram begins with the assumption that a user will call the compiler class from a main program. Within the main program, the user will also provide the necessary information about the definition of the language in the form of a *TransitionTable* object; this information

is then passed to the *Scanner* object along with the source code and the names of the required *SymbolTable* objects for the scanner. After the *Scanner* object is set up, the *compile* method of the *Compiler* object can be called, which calls the *tokenizeInput* method of the scanner which returns the resulting *token\_stream* and *symbol\_tables*. Inside the *tokenizeInput* method, *TransitionTable* objects' method *runDfa* will be run until the *text* field is empty, creating tokens for every recognized word, and creating entries in the symbol tables if necessary. Since at this point the only compilation possible is the lexical analysis of the given source code, the diagram stops after said process.

Although the previous diagram provides a “big picture” of the scanning process, a further breakdown of the *runDfa* method is needed; for this purpose the *runDfa* flowchart, shown in **Figure 2.8**, was made:



**Figure 2.8 The Flowchart for the runDfa method.**

The flowchart illustrates the specifics of what happens inside the *runDfa* method; it provides a pseudocode implementation of any given language specification DFA. The pseudocode equivalent would be the following:

```

1  method runDfa( text, tt ) {
2      state = 0
3      word = ""
4      char = ''
5
6      while state is not accepting AND state is not error {
7
8          if text == "" {
9              break
10         }
11
12         char = text.nextchar()
13         transition = tt [ char ][ state ]
14
15         if transition.isConsuming {
16             state = transition.target
17             continue
18         }
19         else {
20             text.prepend( char )
21             if char == BLANK {
22                 state = transition.target
23                 continue
24             }
25             else {
26                 word = word.append( char )
27                 state = transition.target
28                 continue
29             }
30         }
31     }
32
33     return ( state, word )

```

**Figure 2.9 The pseudocode for the flowchart.**

For the case of Error Recognition, these will be pointed out by the *TransitionTables'* method *runDfa*; after executing, the method will return the state it ended with and based in the states' status, the *Scanner* object will raise an error message that, with the help of a line

counter, includes the line in code in which the error occurred. For the purpose of reporting an unclosed comment this counter will have to take into account the amount of lines that have passed after a comment is opened, in case the comment does not close, so that the reported line number is accurate to where the comment was first opened.

## Implementation

```

1  from scanner import Scanner
2  from symboltable import SymbolTable
3  from transitiontable import TransitionTable
4
5  """
6  Compiler Class
7  author: Victor Emmanuel Guerra Aguado
8
9  Description:
10 |   The Compiler Class serves as an interface to the multiple steps of compilation.
11
12  Attributes:
13 |   + scanner: Scanner; The compilers Scanner, in charge of the lexical analysis.
14 |   + symbol_tables: {str, SymbolTable}; The symbol tables created during the
15 |   lexical analysis.
16  """
17  class Compiler:
18      """
19      Compiler Class Constructor
20
21      __init__(self) -> Compiler
22      """
23      def __init__(self):
24          self.scanner: Scanner = None
25          self.symbol_tables: {str, SymbolTable} = {}
26
27          # The Compiler Class expects further expansion, to account for the
28          # following phases of compilation
29          self.parser = None
30          self.syntax_tree = None
31

```

```

32     """
33     Compiler scannerPrep method
34
35     Description:
36     |     Executes the necessary preparations for the compilers' scanner given a
37     |     text to scan, and a dfa to set the languages' constraints.
38
39     scannerPrep(self, text: str, dfa: TransitionTable, st_names: [str]) -> None
40     """
41     def scannerPrep(self, text: str, dfa: TransitionTable, st_names: [str]):
42     |     self.scanner = Scanner(text, dfa)
43     |     for name in st_names:
44     |         self.scanner.addSymbolTable(name, SymbolTable())
45
46     """
47     Compiler compile method
48
49     Description:
50     |     Executes all of the main components' methods to produce an output as
51     |     indicated by the architecture diagram. Said output is usually passed
52     |     as an argument as input to the next method.
53
54     compile(self) -> None
55     """
56     def compile(self):
57
58     |     # Lexical Analysis Phase
59     |     self.scanner.tokenizeInput()
60     |     token_stream = self.scanner.token_stream
61     |     self.symbol_tables = self.scanner.symbol_tables
62
63     |     # Visualization of both outputs of the Lexical Analysis
64     |     for token in token_stream:
65     |         print(token.toStr())
66
67     |     for table in self.symbol_tables:
68     |         print(table)
69     |         self.symbol_tables[table].printTable()

```

Figure 3.1 The Compiler Class implementation.

```

1  import sys
2  from symboltable import SymbolTable
3  from transitiontable import TransitionTable
4  from ctoken import Token
5  """
6  Scanner Class
7  author: Victor Emmanuel Guerra Aguado
8
9  Description:
10     The Scanner class provides management and operation of the Scanner phase of a
11     Compiler.
12
13  Attributes:
14     + input_text: str;
15     | The text to scan and tokenize; It is consumed along with the scanning and
16     | recognition.
17     + dfa: TransitionTable;
18     | The DFA, or rather its Transition Table implemented as a TransitionTable
19     | Object.
20     + symbol_tables: {str, SymbolTable};
21     | The symbol tables created during the lexical analysis.
22     + token_stream: [Token];
23     | The list of tokens that have been recognized by the scanner.
24     + line: int;
25     | The counter for the lines read during scanning.
26     + comment_displacement: int;
27     | The counter for the lines that occurred inside of a comment, in case said
28     | comment is not closed.
29  """
30  class Scanner:
31      """
32      Scanner Class Constructor
33
34      __init__(self, input_txt: str, dfa: TransitionTable ) -> Scanner
35      """
36      def __init__(self, input_txt: str, dfa: TransitionTable):
37          self.input_text: str = input_txt
38          self.dfa: TransitionTable = dfa
39          self.symbol_tables: {str, SymbolTable} = {}
40          self.token_stream: [Token] = []
41          self.line: int = 1
42          self.comment_displacement = 0

```



```
43     """
44     Scanner addSymbolTable method
45
46     Description:
47     |     Allows for the addition of a new distinct Symbol Table for the tokens
48     |     produced by the Scanner.
49
50     addSymbolTable(self, name: str, symboltable: SymbolTable) -> None
51     """
52     def addSymbolTable(self, name: str, symboltable: SymbolTable):
53     |     self.symbol_tables.setdefault(name, symboltable)
54
55     """
56     Scanner addToken method
57
58     Description:
59     |     Allows for the addition of a new token to the token stream of the scanner.
60
61     addToken(self, token: Token) -> None
62     """
63     def addToken(self, token: Token):
64     |     self.token_stream.append(token)
65
```

```

66     """
67     Scanner runDfa method
68
69     Description:
70         Runs the TransitionTable equivalent of the DFA provided exactly once,
71         and returns the resulting word and end state.
72
73     runDfa(self) -> (state: int, word: str)
74     """
75     def runDfa(self) -> (int, str):
76         self.input_text = self.input_text[::-1]
77         state = 0
78         word = ""
79         current_ch = ""
80         while not self.dfa.accepts(state) and not self.dfa.errors(state):
81             if not self.input_text:
82                 return (state, word)
83             current_ch = self.input_text[-1]
84             #print("current_ch: " + current_ch)
85             trans = self.dfa.nextTransition(state, current_ch)
86
87             if trans.isConsuming:
88                 self.input_text = self.input_text[:-1]
89                 # if current_ch is not BLANK
90                 if current_ch not in ["\n", "\t", " ", "\r"]:
91                     word = word + current_ch
92
93             state = trans.target
94             #print("state: " + str(state))
95
96             if current_ch == '\n':
97                 self.line = self.line + 1
98                 if state == 8 or state == 9:
99                     self.comment_displacement = self.comment_displacement + 1
100
101         self.input_text = self.input_text[::-1]
102         #print("word: " + word)
103         return (state, word)
104

```

```

105 """
106 Scanner tokenizeInput method
107
108 Description:
109     Runs the runDfa Class function until the input_text provided is empty,
110     and adds a token to the token_stream for every valid word recognized.
111
112 tokenizeInput(self) -> None
113 """
114 def tokenizeInput(self):
115     while self.input_text != "":
116         state, word = self.runDfa()
117
118         # In state 10 the word is either a keyword or an identifier
119         if state == 10:
120             # If word is a keyword
121             if word.lower() in self.dfa.valid_tokens.keys():
122                 tok = Token(self.dfa.valid_tokens[word.lower()])
123                 self.addToken(tok)
124
125             # If word is an identifier
126             else:
127                 ## Name of the SymbolTable should also be dynamic, not hardcoded
128
129                 valid_token_id = self.dfa.valid_tokens["IDENTIFIER"]
130                 stid = self.symbol_tables["IDENTIFIER"].lookupEntryValue(word)
131                 # If there is already an entry in the ST with this value
132                 if stid != None:
133                     tok = Token(valid_token_id, stid)
134                     self.addToken(tok)
135                 # If there is NO entry with this value
136                 else:
137                     self.symbol_tables["IDENTIFIER"].addEntry(word)
138                     st_ref = self.symbol_tables["IDENTIFIER"].lookupEntryValue(word)
139                     tok = Token(valid_token_id, st_ref)
140                     self.addToken(tok)
141

```

```

142 # In state 11 the word is a numerical constant
143 elif state == 11:
144     valid_token_id = self.dfa.valid_tokens["NUM_CONSTANT"]
145     stid = self.symbol_tables["NUM_CONSTANT"].lookupEntryValue(word)
146     # If there is already an entry in the ST with this value
147     if stid != None:
148         tok = Token(valid_token_id, stid)
149         self.addToken(tok)
150     # If there is NO entry with this value
151     else:
152         self.symbol_tables["NUM_CONSTANT"].addEntry(word)
153         st_ref = self.symbol_tables["NUM_CONSTANT"].lookupEntryValue(word)
154         tok = Token(valid_token_id, st_ref)
155         self.addToken(tok)
156
157 elif state == 32:
158     sys.exit(f"Invalid Character Error in line: {self.line}")
159 elif state == 33:
160     sys.exit(f"Invalid Character '!' Error in line: {self.line}")
161 elif state == 34:
162     sys.exit(f"Invalid Identifier Error in line: {self.line}")
163 elif state == 35:
164     sys.exit(f"Invalid Number Constant Error in line: {self.line}")
165 elif state == 8 or state == 9:
166     sys.exit(f"Incomplete Comment Error in line: {self.line - self.comment_displacement}")
167 elif state == 31:
168     self.comment_displacement = 0
169     pass
170 elif state in range(12,31):
171     tok = Token(self.dfa.valid_tokens[word.lower()])
172     self.addToken(tok)

```

**Figure 3.2 The Scanner Class implementation.**

```

1  """
2  Token Class
3  author: Victor Emmanuel Guerra Aguado
4
5  Description:
6      The Token Class provides a value object for the tokens created in the
7      Lexical Analysis phase of a compiler.
8
9  Attributes:
10     + keyword_id: int;
11     | The number that tells what type of valid word for the language the token
12     | represents.
13     + symboltable_ref: int;
14     | The ID number that indicates an entry to a Symbol Table, if not an
15     | IDENTIFIER or NUMERICAL_CONST, it defaults to -1.
16  """
17  class Token:
18      """
19      Token Class Constructor
20
21      __init__(self, kwid: int, st_ref: int) -> Token
22      """
23      def __init__(self, kwid: int, st_ref: int = None):
24          self.keyword_id = kwid
25          self.symboltable_ref = st_ref
26
27      """
28      Token toStr method
29
30      Description:
31          Provides a way to transform the Token Object into a string of characters
32          that can be printed onto the screen for visualization or written into a
33          file for storage.
34
35      toStr(self) -> str
36      """
37      def toStr(self) -> str:
38          if self.symboltable_ref is None:
39              return f"({self.keyword_id})"
40          return f"({self.keyword_id}, {self.symboltable_ref})"
41

```

```
42     """
43     Token fromStr method
44
45     Description:
46     | Provides a way to transform Strings of characters created by the
47     | Token.print method into an actual Token Object in order to work with
48     | previously printed Token Objects.
49
50     fromStr(str: str) -> Token
51     """
52     def fromStr(tokenstr: str):
53         temp = tokenstr[1: -1:].split(",")
54         return Token(int(temp[0]), int(temp[0]))
55
```

**Figure 3.3** The Token Class implementation.

```

1  from transition import Transition
2  """
3  TransitionTable Class
4  author: Víctor Emmanuel Guerra Aguado
5
6  Description:
7      The TransitionTable Class represents the implementation of a Deterministic
8      Finite Automaton for a given language into code. It performs everything the
9      DFA would.
10
11  Attributes:
12      + isAcceptingState: [bool];
13      | The list that represents the Accepting States of the Automaton.
14      + isErrorState: [bool];
15      | The list that represents the Error States of the Automaton.
16      + transitions: {str, [Transition]};
17      | The list of transitions of the DFA, the lists in the dictionary act as the
18      | columns of a transition table, The specific transition is then obtained by
19      | using the index as the row of the table.
20      + keywords: [str];
21      | The Keywords of the language, exclusively keywords, no symbols.
22      + valid_tokens: {str, int};
23      | The enumeration of the valid tokens of the language
24  """
25  class TransitionTable:
26      """
27      TransitionTable Class Constructor
28
29      __init__(self,
30              acceptingstates: [bool],
31              errorstates: [bool],
32              transitions: {str, [Transition]},
33              keywords: [String],
34              valid_tokens: {String: Number}) -> TransitionTable
35      """
36      def __init__(self, acceptingstates: [bool], errorstates: [bool],
37                  transitions: dict, keywords: [str], valid_tokens: {str, int}):
38          self.isAcceptingState = acceptingstates
39          self.isErrorState = errorstates
40          self.transitions = transitions
41          self.keywords = keywords
42          self.valid_tokens = valid_tokens
43

```

```

44     """
45     TransitionTable accepts method
46
47     Description:
48     |     Method for finding out if the state provided is an Accepting State
49
50     accepts(self, state: int) -> bool
51     """
52     def accepts(self, state: int) -> bool:
53         return self.isAcceptingState[state]
54
55     """
56     TransitionTable errors method
57
58     Description:
59     |     Method for finding out if the state provided is an Error State
60
61     errors(self, state: int) -> bool
62     """
63     def errors(self, state: int) -> bool:
64         return self.isErrorState[state]
65
66     """
67     TransitionTable nextTransition method
68
69     Description:
70     |     Method for finding the transition that results by the given state and
71     |     condition.
72
73     nextTransition(self, state: int, condition: str) -> Transition
74     """
75     def nextTransition(self, state: int, condition: str) -> Transition:
76         transition = self.transitions["INVALIDCH"][state]
77         for col in self.transitions:
78             if condition in col:
79                 transition = self.transitions[col][state]
80                 break
81
82         return transition

```

**Figure 3.4 The TransitionTable Class implementation**



```

1  """
2  Transition Class
3  author: Victor Emmanuel Guerra Aguado
4
5  Description:
6      The Transition Class provides a structure for the entries made into the
7      Transition Table made as an implementation of the DFA of a language.
8
9  Attributes:
10     + target: int;
11     | The state to which this transition points.
12     + isConsuming: bool;
13     | Whether this specific transition causes the DFA to advance to the next
14     | character.
15 """
16 class Transition:
17     """
18     Transition Class Constructor
19
20     __init__(self, target: int, consumes: bool) -> Transition
21     """
22     def __init__(self, target: int, consumes: bool):
23         self.target = target
24         self.isConsuming = consumes

```

Figure 3.5 The Transition Class implementation.

```

1  import symboltableentry as ste
2  """
3  SymbolTable Class
4  author: Victor Emmanuel Guerra Aguado
5
6  Description:
7      The SymbolTable Class represents one of the Symbol Tables created during the
8      Lexical Analysis phase of a compiler.
9
10 Attributes:
11     + fields: [str];
12     | Names or descriptions of the fields contained within an entry to the
13     | SymbolTable.
14     + entries: [SymbolTableEntry];
15     | The entries of the Symbol Table.
16 """
17 class SymbolTable:
18     """
19     SymbolTable Class Constructor
20
21     __init__(self, fields: [str]) -> SymbolTable
22     """
23     def __init__(self, fields: [str] = []):
24         self.fields = fields
25         self.entries = []
26
27     """
28     SymbolTable addEntry method
29
30     Description:
31         Method for adding an entry to the Symbol Table, given its value and fields.
32
33     addEntry(self, value: str, fields: [str]) -> None
34     """
35     def addEntry(self, value: str, fields: [str] = []):
36         if len(fields) != len(self.fields):
37             print("Entries' fields do not match the Tables' fields.")
38             exit()
39         else:
40             temp = ste.SymbolTableEntry(value, fields)
41             self.entries.append(temp)
42

```

```

43     """
44     SymbolTable lookupEntryId method
45
46     Description:
47     | Method for looking up an entry in the SymbolTable given an index or
48     | entry ID.
49
50     lookupEntryId(self, index: int) -> SymbolTableEntry
51     """
52     def lookupEntryId(self, index: int) -> ste.SymbolTableEntry:
53     |     return self.entries[index]
54
55     """
56     SymbolTable lookupEntryValue method
57
58     Description:
59     | Method for looking up an entry in the SymbolTable given a String value.
60     | Returns the entry index.
61
62     lookupEntryValue(self, value: str) -> int
63     """
64     def lookupEntryValue(self, value: str) -> int:
65     |     for entry in self.entries:
66     |         if entry.value == value:
67     |             return self.entries.index(entry)
68
69     |     return None
70
71     """
72     SymbolTable printTable method
73
74     Description:
75     | Method for printing the contents of the SymbolTable, in case visualization
76     | of the table is needed.
77
78     printTable(self) -> None
79     """
80     def printTable(self):
81     |     print("| ID | VALUE | " + "|".join(self.fields))
82     |     for entry in self.entries:
83     |         print(str(self.lookupEntryValue(entry.value)) + " | " + entry.value
84     |             + "|".join(entry.fields))

```

**Figure 3.6 The SymbolTable Class Implementation.**

```

1  """
2  SymbolTableEntry Class
3  author: Victor Emmanuel Guerra Aguado
4
5  Description:
6      The SymbolTableEntry Class provides a structure for the entries made into a
7      Symbol Table during the Lexical Analysis phase of a Compiler.
8
9  Attributes:
10     + value: str;
11         The value to be stored as an entry of the symbol table, represents a lexeme
12         recognized by the scanner.
13     + fields: [str];
14         Fields where any necessary extra information about the lexeme will be stored.
15 """
16 class SymbolTableEntry:
17     """
18     SymbolTableEntry Class Constructor
19
20     __init__(self, value: str, fields: [str]) -> SymbolTableEntry
21     """
22     def __init__(self, value: str, fields: [str] = []):
23         self.value = value
24         self.fields = fields

```

**Figure 3.7 The SymbolTableEntry Class implementation.**

```

1  import sys
2  from compiler import Compiler
3  from transitiontable import TransitionTable
4  from transition import Transition
5
6  """
7      Main Program File for the C Minus Compiler
8      author: Victor Emmanuel Guerra Aguado
9
10     Description:
11         The main file of the entire compiler program, where the language definition is
12         provided, and the source code text is obtained. So far, the hardcoded language
13         definition is compliant to both that of the C Minus language, and that of the
14         class diagrams designed during the design phase.
15
16     Attributes:
17         None
18 """
19 if __name__ == "__main__":
20     # Obtain a Compiler Object
21     comp = Compiler()
22
23     # Get the file to parse
24     if len(sys.argv) < 2:
25         sys.exit("Input file not provided.")
26
27     filename = sys.argv[1]
28     input_file = open(filename, 'r')
29
30     input_txt = input_file.read()
31
32     input_file.close()
33
34

```

```
36  # To prepare the Scanner, define a transition table of a DFA
37  acceptingstates = [
38      0,0,0,0,0,0,
39      0,0,0,0,0,0,
40      1,1,1,1,1,1,
41      1,1,1,1,1,1,
42      1,1,1,1,1,1,
43      1,1,1,1,1,1,
44      1,1,1,1,1,1,
45      1
46  ]
47
48  errorstates = [
49      0,0,0,0,0,0,
50      0,0,0,0,0,0,
51      0,0,0,0,0,0,
52      0,0,0,0,0,0,
53      0,0,0,0,0,0,
54      0,0,0,0,0,0,
55      0,0,1,1,1,1,
56      1
57  ]
58
```



[illegible]



[illegible]

```

243 keywords = ["else", "if", "input", "int", "output", "return", "void", "while"]
244
245 valid_tokens = {
246     "else": 1,
247     "if": 2,
248     "input": 3,
249     "int": 4,
250     "output": 5,
251     "return": 6,
252     "void": 7,
253     "while": 8,
254     "+": 9,
255     "-": 10,
256     "/": 11,
257     ">": 12,
258     ">=": 13,
259     "<": 14,
260     "<=": 15,
261     "==" : 16,
262     "!=": 17,
263     "=": 18,
264     ";": 19,
265     ",": 20,
266     "(": 21,
267     ")": 22,
268     "{": 23,
269     "}": 24,
270     "[": 25,
271     "]": 26,
272     "IDENTIFIER": 27,
273     "NUM_CONSTANT": 28,
274     "COMMENT": 29,
275 }
276
277 cminusdfa = TransitionTable(acceptingstates, errorstates, transitions, keywords, valid_tokens)
278
279 # We add the DFA and the Input text to the scanner
280 comp.scannerPrep(input_txt, cminusdfa, ["IDENTIFIER", "NUM_CONSTANT"])
281
282 # Start the compilation process
283 comp.compile()

```

Figure 3.8 The main file for the compiler program.

## Verification & Validation

The Verification and Validation phase consists of making sure that the product is being well developed and that the developed product is the correct one. In this section a test model, consisting of several test cases will be presented, along with their purpose, input, expected result, result, and evaluation. The Test Case templates are based on the ones suggested by Thomas Hamilton.[2]

### Test Cases

#### CMCTC-01

**Description/Purpose:**

To verify that the scanner object creates tokens completely and correctly from the given input text and according to the C Minus language definition.

**Test Script:**

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-01 test input source code.
3. Run the compiler.
4. Compare the result with the expected result.

**Test input:**

```
“void main(){
    else if input int output return void while;
    i = 0;
    i + 1
    i - 1
    i * 1
    i / 1
    i > 1 >= 0
    i < 1 <= 0
    i == 1
    i != 1
    [, ]
    /**/
}“
```

**Expected Result:**

(7), (28, 0), (22), (23), (24), (1), (2), (3), (4), (5), (6), (7), (8), (20), (28, 1), (19), (29, 0), (20), (28, 1), (9), (29, 1), (28, 1), (10), (29,1), (28,1), (11), (29,1), (28,1), (12), (29,1), (28,1), (13), (29,1), (14), (29,0), (28,1), (15), (29,1), (16), (29,0), (28,1), (17), (29,1), (28,1), (18), (29,1), (26), (21), (27), (25)

IDENTIFIER

| ID | VALUE |

0 | main

1 | i

NUM\_CONSTANT

| ID | VALUE |

0 | 0

1 | 1

**Actual Result:**

(7)

(28, 0)

(22)

(23)

(24)

(1)

(2)

(3)

(4)

(5)

(6)

(7)

(8)

(20)

(28, 1)

(19)

(29, 0)

(20)

(28, 1)

(9)

(29, 1)

(28, 1)

(10)

(29, 1)

(28, 1)

(11)

(29, 1)

(28, 1)

(12)

(29, 1)

(28, 1)

(13)

(29, 1)

(14)

(29, 0)

(28, 1)

(15)

(29, 1)

(16)

(29, 0)

(28, 1)

(17)

(29, 1)

(28, 1)

(18)

(29, 1)

(26)

(21)

(27)

(25)

IDENTIFIER

| ID | VALUE |

0 | main

1 | i

NUM\_CONSTANT

| ID | VALUE |

0 | 0

1 | 1

**Evaluation: PASS**

```

273     Test Case CMCTC-01
274
275     Description/Purpose:
276         To verify that the scanner object creates tokens completely and correctly
277         from the given input text and according to the C Minus language definition.
278     """
279     def test_token_creation(self):
280         input_txt = """void main(){
281     else if input int output return void while;
282     i = 0;
283     i + 1
284     i - 1
285     i * 1
286     i / 1
287     i > 1 >= 0
288     i < 1 <= 0
289     i == 1
290     i != 1
291     [, ]
292     /**/
293 }"""
294
295     self.compiler.scannerPrep(input_txt, self.dfa, self.symbol_table_names)
296
297     expected_token_stream = [Token(7),Token(28,0),Token(22),Token(23), Token(24),
298     Token(1), Token(2), Token(3), Token(4), Token(5), Token(6), Token(7), Token(8),
299     Token(20), Token(28,1), Token(19), Token(29,0), Token(20), Token(28,1), Token(9),
300     Token(29,1), Token(28,1), Token(10), Token(29,1), Token(28,1), Token(11),
301     Token(29,1), Token(28,1), Token(12), Token(29,1), Token(28,1), Token(13),
302     Token(29,1), Token(14), Token(29,0), Token(28,1), Token(15), Token(29,1),
303     Token(16), Token(29,0), Token(28,1), Token(17), Token(29,1), Token(28,1),
304     Token(18), Token(29,1), Token(26), Token(21), Token(27), Token(25)]
305
306     self.compiler.compile()
307
308     expected_tokens = []
309     for tok in expected_token_stream:
310         expected_tokens.append(tok.toStr())
311
312     actual_tokens = []
313     for tok in self.compiler.scanner.token_stream:
314         actual_tokens.append(tok.toStr())
315
316     self.assertListEqual(expected_tokens, actual_tokens)

```

Figure 4.1 CMCTC-01 Test Case Implementation.

## CMCTC-02

### Description/Purpose:

To verify that the scanner completely and correctly creates Symbol Table entries of any valid given type.

**Test Script:**

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-02 test input source code.
3. Run the compiler.
4. Compare the result with the expected result.

**Test Input:**

```
“void main() {
    one = 0;
    two = 1;
    three = 2;
}”
```

**Expected Result:**

(7), (28,0), (22), (23), (24), (28,1), (19), (29,0), (20), (28,2), (19), (29,1), (20), (28,3), (19),  
(29,2), (20), (25)

IDENTIFIER

ID | VALUE

0 | main

1 | one

2 | two

3 | three

NUM\_CONSTANT

ID | VALUE

0 | 0

1 | 1

2 | 2

**Actual Result:**

(7)

(28, 0)

(22)

(23)

(24)

(1)

(2)

(3)

(4)  
(5)  
(6)  
(7)  
(8)  
(20)  
(28, 1)  
(19)  
(29, 0)  
(20)  
(28, 1)  
(9)  
(29, 1)  
(28, 1)  
(10)  
(29, 1)  
(28, 1)  
(11)  
(29, 1)  
(28, 1)  
(12)  
(29, 1)  
(28, 1)  
(13)  
(29, 1)  
(14)  
(29, 0)  
(28, 1)  
(15)  
(29, 1)  
(16)  
(29, 0)  
(28, 1)  
(17)



(29, 1)

(28, 1)

(18)

(29, 1)

(26)

(21)

(27)

(25)

IDENTIFIER

| ID | VALUE |

0 | main

1 | i

NUM\_CONSTANT

| ID | VALUE |

0 | 0

1 | 1

**Evaluation: PASS**

```

319     """
320     Test Case CMCTC-02
321
322     Description/Purpose:
323     |   To verify that the scanner completely and correctly creates SymbolTable
324     |   entries of any valid given type.
325     """
326     def test_symbol_table_entries(self):
327         input_txt = """void main() {
328         one = 0;
329         two = 1;
330         three = 2;
331     }
332     """
333         self.compiler.scannerPrep(input_txt, self.dfa, self.symbol_table_names)
334
335         expected_token_stream = [Token(7), Token(28,0), Token(22), Token(23), Token(24),
336         Token(28,1), Token(19), Token(29,0), Token(20), Token(28,2), Token(19),
337         Token(29,1), Token(20), Token(28,3), Token(19), Token(29,2), Token(20), Token(25)]
338
339         self.compiler.compile()
340
341         expected_tokens = []
342         for tok in expected_token_stream:
343             expected_tokens.append(tok.toStr())
344
345         actual_tokens = []
346         for tok in self.compiler.scanner.token_stream:
347             actual_tokens.append(tok.toStr())
348
349         self.assertListEqual(expected_tokens, actual_tokens)

```

Figure 4.2 The CMCTC-02 Test Case Implementation.

## CMCTC-03

### Description/Purpose:

To verify that the scanner raises an error when encountering an unclosed comment.

### Test Script:

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-03 test input source code.
3. Run the compiler.
4. Compare the result with the expected result.

### Test Input:

```

“void main() {
    /* Main Function */
    if(true){
        /* Comment

```

```

    }
}
”

```

**Expected Result:**

Incomplete Comment Error in line: 4

**Actual Result:**

Incomplete Comment Error in line: 4

**Evaluation: PASS**

```

352     """
353     Test Case CMCTC-03
354
355     Description/Purpose:
356     | To verify that the scanner raises an error when encountering an unclosed
357     | comment.
358     """
359     def test_unclosed_comment(self):
360         input_txt = """void main() {
361             /* Main Function */
362             if(true){
363                 /* Comment
364                 */
365             }
366         """
367         self.compiler.scannerPrep(input_txt, self.dfa, self.symbol_table_names)
368
369         with self.assertRaises(SystemExit) as res:
370             self.compiler.compile()
371         print(res.exception)
372         self.assertEqual(res.exception.__str__(), "Incomplete Comment Error in line: 4")

```

Figure 4.3 The CMCTC-03 Test Case Implementation

## CMCTC-04

**Description/Purpose:**

To verify that the scanner raises an error when encountering an invalid character for the language.

**Test Script:**

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-04 test input source code.
3. Run the compiler.
4. Compare the result with the expected result.

**Test Input:**

```

“void main() {
    int 編譯器;
    編譯器 = 10;
}
”

```

**Expected Result:**

Invalid Character Error in line: 2

**Actual Result:**

Invalid Character Error in line: 2

**Evaluation: PASS**

```

374     """
375     Test Case CMCTC-04
376
377     Description/Purpose:
378         To verify that the scanner raises an error when encountering an invalid
379         character for the language.
380     """
381     def test_invalid_character(self):
382         input_txt = """void main() {
383             int 編譯器;
384             編譯器 = 10;
385         }
386     """
387     self.compiler.scannerPrep(input_txt, self.dfa, self.symbol_table_names)
388
389     with self.assertRaises(SystemExit) as res:
390         self.compiler.compile()
391     print(res.exception)
392     self.assertEqual(res.exception.__str__(), "Invalid Character Error in line: 2")

```

Figure 4.4 The CMCTC-04 Test Case Implementation.

## CMCTC-05

**Description/Purpose:**

To verify that the scanner accepts any character, even invalid characters for the language, if it is placed inside a comment.

**Test Script:**

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-05 test input source code.
3. Run the compiler.

4. Compare the result with the expected result.

**Test Input:**

```
“void main() {
    /* 冰淇淋 */
}
”
```

**Expected Result:**

[(7), (28,0), (22), (23), (24), (25)]

IDENTIFIER

ID | VALUE

0 | main

NUM\_CONSTANT

ID | VALUE

**Actual Result:**

(7)

(28, 0)

(22)

(23)

(24)

(25)

IDENTIFIER

| ID | VALUE |

0 | main

NUM\_CONSTANT

| ID | VALUE |

**Evaluation: PASS**

```

394     """
395     Test Case CMCTC-05
396
397     Description/Purpose:
398         To verify that the scanner accepts any character, even invalid characters
399         for the language, if it is placed inside a comment.
400     """
401     def test_characters_in_comment(self):
402         input_txt = """void main() {
403             /* 冰淇淋 */
404         }
405     """
406     self.compiler.scannerPrep(input_txt, self.dfa, self.symbol_table_names)
407
408     expected_token_stream = [Token(7), Token(28,0), Token(22), Token(23), Token(24),
409                             Token(25)]
410
411     self.compiler.compile()
412
413     expected_tokens = []
414     for tok in expected_token_stream:
415         expected_tokens.append(tok.toStr())
416
417     actual_tokens = []
418     for tok in self.compiler.scanner.token_stream:
419         actual_tokens.append(tok.toStr())
420
421     self.assertListEqual(expected_tokens, actual_tokens)

```

Figure 4.5 The CMCTC-05 Test Case Implementation.

## CMCTC-06

### Description:

To verify that the scanner can raise an error when recognizing an invalid identifier.

### Test Script:

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-06 test input source code.
3. Run the compiler.
4. Compare the result with the expected result.

### Test Input:

```

“void main() {
    int i2;
    i2 = 0;
}
”

```

**Expected Result:**

Invalid Identifier Error in line: 2

**Actual Result:**

Invalid Identifier Error in line: 2

**Evaluation: PASS**

```

423     """
424     Test Case CMCTC-06
425
426     Description/Purpose:
427     |   To verify that the scanner can raise an error when recognizing an invalid
428     |   identifier.
429     """
430     def test_invalid_identifier(self):
431         input_txt = """void main() {
432             |   int i2;
433             |   i2 = 0;
434         }
435     """
436     self.compiler.scannerPrep(input_txt, self.dfa, self.symbol_table_names)
437
438     with self.assertRaises(SystemExit) as res:
439         self.compiler.compile()
440     print(res.exception)
441     self.assertEqual(res.exception.__str__(), "Invalid Identifier Error in line: 2")

```

Figure 4.6 The CMCTC-06 Test Case Implementation.

## CMCTC-07

**Description:**

To verify that the scanner can raise an error when recognizing an invalid numerical constant.

**Test Script:**

1. Provide the C Minus Language specification to the compiler.
2. Provide the CMCTC-06 test input source code.
3. Run the compiler.
4. Compare the result with the expected result.

**Test Input:**

```

“void main() {
    int i;
    i = 2b;
}
”

```

**Expected Result:**

Invalid Number Constant Error in line: 3

**Actual Result:**

Invalid Number Constant Error in line: 3

**Evaluation: PASS**

```

443     """
444     Test Case CMCTC-07
445
446     Description/Purpose:
447         To verify that the scanner can raise an error when recognizing an invalid
448         numerical constant.
449     """
450     def test_invalid_num_constant(self):
451         input_txt = """void main() {
452             int i;
453             i = 2b;
454         }
455     """
456     self.compiler.scannerPrep(input_txt, self.dfa, self.symbol_table_names)
457
458     with self.assertRaises(SystemExit) as res:
459         self.compiler.compile()
460     print(res.exception)
461     self.assertEqual(res.exception.__str__(), "Invalid Number Constant Error in line: 3")

```

**Figure 4.7 The CMCTC-07 Test Case Implementation.**

**References**

1. R. Castelló, Class Lecture, Topic: “Chapter 2 - Lexical Analysis.” TC3048, School of Engineering and Science, ITESM, Chihuahua, Chih, April, 2022.
2. Thomas Hamilton, “How to Write Test Cases: Sample Template with Examples”, *Guru99*, 2 April 2022 [journal on-line]; available from <https://www.guru99.com/test-case.html>; Internet; accessed 18 April 2022.