

BIRLA INSTITUTE OF TECHNOLOGY & SCIENCE, PILANI
HYDERABAD CAMPUS
SECOND SEMESTER 2020 – 2021

COMPILER CONSTRUCTION (CS F363)
Programming Assignment – 1

Introduction

You will implement a compiler for a miniature programming language, too small a language that you can give a name to it. This language will help getting hands on practice to those concepts you learn in this course. This language must have a set of sequential statements, a conditional construct, loop construct, functions and arrays.

The compiler project will be divided into a two phases and implemented in stages.

Notes on implementation

Implementation of the project **"must" be done in C**. This is to ensure that all data structures and algorithms are hand-coded without the use of high level libraries and implementation must run on Linux.

Assignment Administration

- Project may be worked in **teams of four**. Choose your own team but you will not be allowed to change your team-mate later.
- **Register your team details along with the language you wish to implement in the following Google spreadsheet**

https://docs.google.com/spreadsheets/d/1-o_3ht4dsgxYfJlXQk0iKoyyrGpn8XcTc-nw8YYxziU/edit?usp=sharing

- The project will be evaluated twice once before mid sem and post mid sem. Each stage will be evaluated through a viva-voce.
- **Marking will be based not only on the implementation but also on your understanding of the implementation and the ability to explain your code and answer questions on your part of the work.**
- Each stage has specific deadline. Any submission beyond the one-day extension will carry an automatic depreciation by **2 marks - per day of delay**.

Fair Practice

- Teams are permitted to discuss the project with each other but not allowed to see nor use each other's solutions.
- Plagiarism in any form is unacceptable. Project submissions will be rigorously scrutinized for plagiarism and the team members will be questioned to verify the ownership of the

solution.

Phase-1 (Requirement for Scanner, Scanner Deliverables, Test Suite, How-to-approach Phase-1)
(**Deadline : 15-Feb-2021** , mode of submission: CMS) **(Marks 10)**

Lexer / Scanner

In this assignment you write a DFA-based Lexical Analyser that recognizes some of the basic lexemes. Design, write, and thoroughly test the Language constructs. Write a driver program(**parser**) that calls your Scanner repeatedly, returning each token found by the Scanner in the input stream.

- Requirements Specification:
 - Input: Program File (example shown at the end)
 - Output: Return tokens either in the form of some number or as TK-identifier (example shown at the end)
 - Side Effects: White spaces removed
 - Exceptions: Invalid tokens

Our language reserves all the key words that can appear in the language.

Scanner Deliverables

- **C Code** for scanner.
- Test cases and output files for test cases

Scanner sample Test Suite

Formulate your own test set / programs from the token list given as examples. The test set to be used for evaluation

Token List

- **Keywords:** int, float, boolean, string, while, until, if else, true, false
- **Operators:** +, -, *, /, %, :=, ==, >, <, >=, <=, !=, &&, ||, !, ?, :
- **Delimiters:** {, }, (,), [,], ;, ,
- **Identifiers:** must start with a letter (upper or lower case), and may contain zero or more additional characters as long as they are letters, digits, or underscores
- **Integer Literals:** may begin with an optional plus or minus followed by a sequence of one or more digits, provided that the first digit can only be zero for the number zero (which should not have a plus or minus before it).
- **Floating Point Literals:** may begin with an optional plus or minus followed by a sequence of one or more digits with the same provision above for integers, followed by a decimal point and one or more digits after the decimal point.

- **String literals:** start and end with a double quote followed by zero or more characters that may not be newlines, carriage returns, double quotes, or backslashes. The only exceptions are reserved escape sequences which are limited to the following: \t, \n, \r, \", and \\.

Phase - 1 How-to

1. Read the language Specification: the overview, the grammar (natural form) and the tokens to gain an over-all understanding.
2. Apply your understanding to write Regular Expression and convert them into DFA.
3. Implement the DFA.
4. Test your Scanner with the given test cases.
5. Write your own test cases and document your code
6. **To make your DFA as small as possible and manage the code you may store all reserved words in an string array and first check if the input is matching with any string in the array you can declared it as keyword otherwise you search in your DFA.**
7. **For a smooth transition into the next phase first write the Context Free Grammar / BNF for the mini language and from the grammar identify the terminals which can be grouped into tokens. Draw DFA for each tokens individually and then combine followed by implementing it.**

Your code will be evaluated for 10 marks on the following

1. Identifying the tokens - 3
2. Eliminate comments and while spaces. - 2
3. Generate errors - 1
4. Keep track of line numbers – 1
5. Viva – 3

If I run this on the file

/* A program to compute factorials */

```
int fact( int n) {
if (n <= 1)
return 1;
else
return n*fact(n-1);
}
void main(void) {
int x;
x = 1;
while (x <= 10) {
write(x);
write(fact(x));
writeln();
}
```

```
x = x + 1;  
}  
}
```

I get the following token stream:

Token 100, string int, line number 3

Token 132, string fact, line number 3

Token 128, string (, line number 3

Token 100, string int, line number 3

Token 132, string n, line number 3

Token 129, string), line number 3

Token 114, string {, line number 3

Token 103, string if, line number 4

Token 128, string (, line number 4

Token 132, string n, line number 4

Token 123, string <=, line number 4

Token 130, string 1, line number 4

Token 129, string), line number 4

Token 106, string return, line number 5

Token 130, string 1, line number 5

Token 110, string :, line number 5

Token 104, string else, line number 6

Token 106, string return, line number 7

Token 132, string n, line number 7

Token 118, string *, line number 7

Token 132, string fact, line number 7

Token 128, string (, line number 7

Token 132, string n, line number 7

Token 117, string -, line number 7

Token 130, string 1, line number 7

and so forth. Your token kinds, being constants, are likely to be different from mine, but your string values and line numbers should be the same.