

# Implementation of a compiler for an imperative language IMP

Remy Detobel & Denis Hoornaert

October 17, 2017

## 1 Introduction

The project aim is to implement a compiler for a 'simple' imperative language named *IMP*. Like any imperative programming language, *IMP* is structured of mainstream features such as *keywords* (*if*, *while*, ... statements), the use of *variables*, the use *numbers* and the use of *comments*. The form of these features follows some defined rules :

- a *variable* is a sequence of alphanumeric characters that must start by a letter.
- a *number* is a sequence of one or more digits.
- a *comment* must start by the combination (\*) and ends by the reversed combination (\*)).

The compilation scheme is generally divided in three main phases : analysis, synthesis and optimization. The phases are themselves composed of different steps. For instance, the analysis phase is composed of *lexical analysing* step (or *scanning*), a *syntax analysing* step (or *parsing*) and a *semantic analysing* step. In this assignment, the focus is set on the *analysis phase*.

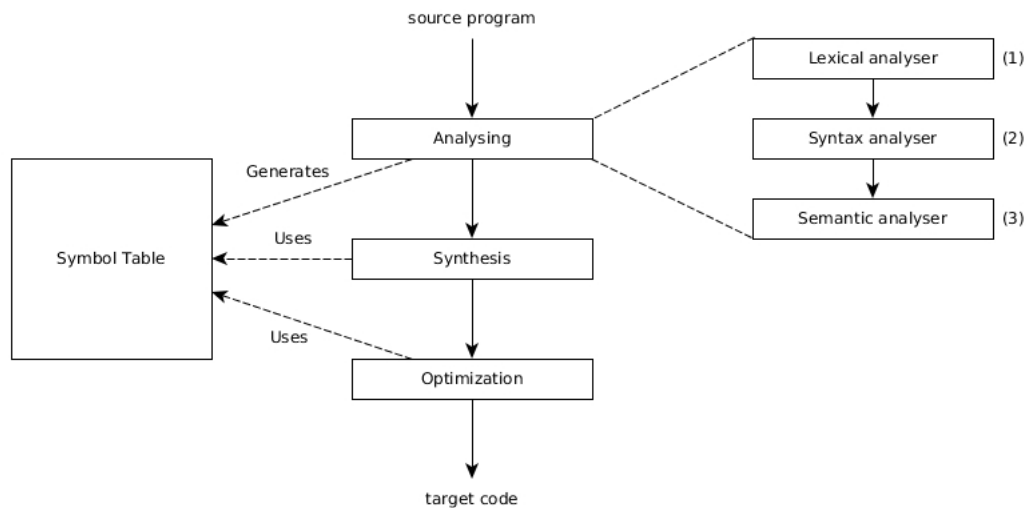


Figure 1 - Compilation phases.

## 2 Implementation of the lexical analyser

In the so called "Dragon book"<sup>1</sup> the *lexical analyser* is defined as follow :

«The *lexical analyser* reads the stream of characters making up the source program and groups the character into a meaningful sequence called *lexemes*.»

<sup>1</sup>V. Aho, A., 2007. *Compilers : Principles, techniques, & Tools*. 2nd ed. New York : Pearson.

A *lexeme* can be defined as a tuple which contains both a *token name* and the associated value. The sequence of *lexemes* generated by the *lexical analyser* will be used by the following step. In addition, the *lexical analyser* will generate a very useful tool, that will be used by all the other steps (as shown in fig 1.), called a *symbol table*. The role of the *symbol table* is to store every variable encountered while scanning the source code and the line where it appears for the first time.

## 2.1 The use of a lexical analyser generator

In order to ease the process of recognizing the lexemes defined in the given `LexicalUnits.java` many *lexical analysers* have been developed. Among them, the most well known generator is the flex family and all its derived versions. In the present project, `jflex` is used as it has been decided to implement the project with the `java` programming language. Using a *lexical analyser generator* eases the analysis of any source because it enables the programmers to describe every *regular expression* by using the *Regex* writing convention and then to generate a `.java` file that will match them all. This generated `.java` file can then be used as any other `java` class.

## 2.2 Lexical analyser structure

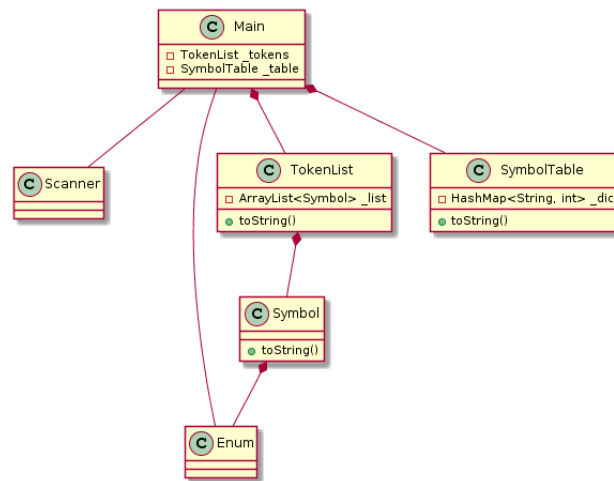


Figure 2 - Class model.

## 2.3 Regular expressions

## 2.4 Dealing with nested comments

Nowadays, most of common programming languages can handle *nested comments*. However, state of the art regular languages cannot boast the same feature. For example, let us try to analyse the following nested comment :

```

1      (* I am a (*nested*) comment *)
2

```

A simple regular language will, once an opening statement (here : `(*)`) encountered, overlook the following characters until it encounters a closing statement (here : `*)`). Unfortunately, if it does so on the example above, it will overlook the second opening statement (columns 9 & 10) and will stop when it comes across the first closing statement (columns 17 & 18) having for consequence that the third part of the *nested comment* will remain. To overcome this problem, the analyser must know how many opening statement it has came across and how many closing statement it should expects in order to know whether it is still in a comment.

The most obvious and smartest way to implement it is to use a counter (i.e. a memory) that will be incremented for every opening encountered and decreased for every closing statement encountered. However, from a theoretical point of view, by using a memory the language cannot be considered as

regular any more. In the present project, it is not a problem and `jflex` allows us to implement such a language.

## 2.5 Tests

```
1 begin
2   s := [45, 68, 23] ;
3   bi := 0 ;
4   len := 3 ;
5   bs := len ;
6   m := (bi+bs)/2 ;
7   while bi < bs and x != s[m] do
8     m := (bi+bs)/2 ;
9     if s[m] < x then
10      bi := m+1 ;
11    else
12      bs := m
13    endif
14  done
15  if len <= m or s[m] != x then
16    m := -1 ;
17  endif
18  print(m)
19 end
```

```
1 begin
2   s := [45, 68, 23];
3   n := 3 ;
4   for i from 0 to n do
5     save := s[i] ;
6     j := i-1 ;
7     while j >= 0 and s[j] > save do
8       s[j+1] := s[j] ;
9       j := j-1 ;
10    done
11    s[j+1] := save ;
12  done
13 end
```

## 3 How to set up the project