

# Implementation of a compiler for an imperative language IMP

Remy Detobel & Denis Hoornaert

November 24, 2017

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Implementation of the lexical analyser</b>	<b>2</b>
2.1	Use of a lexical analyser generator . . . . .	2
2.2	Regular expressions . . . . .	3
2.3	Hypothesis on regular expressions . . . . .	3
2.4	Dealing with nested comments . . . . .	5
2.5	Tests and results . . . . .	5
<b>3</b>	<b>Implementation of the syntax analyser</b>	<b>5</b>
3.1	Use and implementation of a parser generator . . . . .	6
3.2	Parser phases . . . . .	7
3.2.1	Transforming the grammar to LL(1) grammar . . . . .	7
3.2.2	Action table . . . . .	8
3.2.3	Syntax checking . . . . .	8
3.3	Resulting parser generator architecture . . . . .	8
3.4	Tests and results . . . . .	8
3.4.1	Ambiguity . . . . .	9
3.4.2	Useless symbols removal . . . . .	9
3.4.3	Left-recursion removal . . . . .	10
3.4.4	Factorisation . . . . .	10
3.4.5	Resulting action table . . . . .	11
3.4.6	Syntax checking result . . . . .	11
<b>4</b>	<b>How to set up the project</b>	<b>11</b>
4.1	Compilation . . . . .	11
4.2	Execution . . . . .	11
4.3	Test . . . . .	11
4.4	Javadoc . . . . .	12
<b>5</b>	<b>Annexe</b>	<b>12</b>
5.1	Grammar factorisation . . . . .	12
5.2	Removing left-recursion . . . . .	12
5.3	Removing useless variables . . . . .	13
5.4	Ambiguous grammar . . . . .	13
5.5	Basic IMP grammar . . . . .	14
5.6	IMP action table . . . . .	14

# 1 Introduction

The aim of project is to implement a compiler for a 'simple' imperative language named *IMP*. Like any imperative programming language, *IMP* is composed of mainstream features such as *keywords* (*if*, *while*, ... statements), *variables*, *numbers* and *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 optimisation. 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 as shown in fig.1. 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 characters into a meaningful sequence called *lexemes*.»

A *lexeme* can be defined as a tuple which contains both a *token name* and the associated value (not always mandatory). 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 used during all the other steps (as shown in fig.1 .) and 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 Use of a lexical analyser generator

In order to ease the process of recognizing the lexemes defined in the given `LexicalUnits.java` file many *lexical analysers* have been developed. Among them, the most well known generator is the flex program and all its derived versions. In the present project, `jflex` is used as it has been decided to

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

implement the project using the `java` programming language. Using a *lexical analyser generator* eases the analysis of any input 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 recognise all of them. This generated `.java` file can then be used as any other `java` class.



Figure 2: Model class (Pseudo-UML). The TokenList class is the sequence of lexemes and the Scanner class is the file generate by jflex

## 2.2 Regular expressions

Based on the content of `LexicalUnits.java`, we can easily divide the set of lexical units into two distinct groups : the *keyword* group and the variable/constant group.

The implementation of the *keyword* group using regular expressions is pretty straightforward as simply writing the *keyword* is sufficient. For instance, the regular expression of the *keyword* `if` is simply `if`. On the other hand, the implementation of the variable/constant group requires slightly more work. This small group is composed of two elements the variables and the numbers.

The structure of *variables* given in the assignment statements is "a sequence of alphanumeric characters that must start by a letter". Thus, the equivalent regular expression is :

`[a-zA-Z][a-zA-Z0-9]*`

The structure of *numbers* given in the assignment statements is "a sequence of one or more digits". thus, the equivalent regular expression is :

`[0-9]+`

## 2.3 Hypothesis on regular expressions

The only hypothesis that has been made throughout the realisation of the project concerns the behaviour of the *lexical analyser* when a character not specified in either the structure of a *number*, a *variable* or a *keyword* is encountered. Typically, amongst this set there are the following characters :

`}, {, _, |, &, [, ], (, )`

In order words, the question is : What does the *lexical analyser* do for the following line :

```
1 index_of_loop := list[x]
```

Four ideas have been considered :

- Considering these characters as a new lexical unit identified by the following regular expression (not exhaustive) :

```
SpecialChar = ["}", "{", "_", "|", "&", "[", "]", "(", ")"]
```

The example above would be transposed by the *lexical analyser* into the following token list :

token: index	lexical unit: VARNAME
token: _	lexical unit: SPECIALCHAR
token: of	lexical unit: VARNAME
token: _	lexical unit: SPECIALCHAR
token: loop	lexical unit: VARNAME
token: :=	lexical unit: EQUAL
token: [	lexical unit: SPECIALCHAR
token: x	lexical unit: VARNAME
token: ]	lexical unit: SPECIALCHAR

Unfortunately, the assignment statement disallows us to modify the `LexicalUnits.java` file. Consequently, this possibility is not relevant.

- Considering these characters as normal characters. This would mean that they could be part of a *variable* name. Thus, the regular expression for identifying *variables* has to be modify to look like :

```
SpecialChar = ["}", "{", "_", "|", "&", "[", "]", "(", ")"]
[a-zA-Z|SpecialChar][a-zA-z0-9|SpecialChar]*
```

As a consequence, the token list generated by the *lexical analyser* will behave as follow :

token: index_of_loop	lexical unit: VARNAME
token: :=	lexical unit: EQUAL
token: [x]	lexical unit: VARNAME

Even though, using characters like `_` in variable name is common in many programming languages, the other characters are generally not used for this purpose. Given this fact and the fact that the assignment statement does not explicitly mention such a possibility, this idea has been overlooked. Moreover, this implies that variable such as `{!^$'#` would be considered as valid even though having such *variables* is not handy.

- Not considering them. In this possibility, we just overlook them like if they were equivalent to a space character. Implementing this idea is quick but does not really make sense because theses characters would then cause many problems in the following steps.
- Throwing an error. This idea consists simply on printing a warning when an unexpected character is encountered and on stopping the program as resuming does not make any sense. Moreover, this behaviour is pretty common in many programming languages. This solution is the one who fits the best the assignment statements. Therefore, this solution has been preferred over the two others.

## 2.4 Dealing with nested comments

The management of comments using regular language is quite simple. Once an opening statement (here : `'(*)'`) has been encountered, it overlooks the following characters until it encounters a closing statement (here : `'*)'`).

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

Unfortunately, applying the same mechanism on a nested comment will result in a ill-formed outcome. Indeed, in the case of the example above, the analyser 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 statements it came across and how many closing statements it should expect to encounter 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 statement 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 and results

In this section, the results of the implementation are analysed and tested throught three *IMP* source codes : one given in the assignment statements and the two others inspired by algorithms from the Syllabus of Thierry Massart<sup>2</sup>. The aim of testing the *lexical analyser* on these three tests is to ensure that a maximum of the *keywords* and the *variables* are recognized because the set of keywords in the `Euclid.imp` (fig.3) file does not cover every possibilities. This is why these two codes have been chosen. As explained above in the hypothesis subsection 2.3, the program in fig.5 simply stops its execution as it encounters an undefined character at line 2 (`'['`). Unfortunately, it is difficult to cover all the different statements as finding interesting samples of code that do not use `list(s)` is hard.

```
1 begin
2   read(a) ;
3   read(b) ;
4   while b <> 0 do
5     c := b ;
6     while a >= b do
7       a := a-b
8     done ;
9     b := a ;
10    a := c
11  done ;
12  print(a)
13 end
```

Figure 3: *IMP* code to compute the gcd of two numbers

## 3 Implementation of the syntax analyser

The *syntax analysis* (or *parsing*) is the step of the *analysis* phases that aims to verify the structure *syntax* of the source code and then reporting the potential errors using both the list of tokens generated previously by the *scanner* and a grammar (see definition later) given by the language designer. The outcome of the parser is a *syntax tree* that will be used by the following phase (as shown in the fig.X). We distinguish two types of *parsers* : the *top-down* and the *bottom-up*.

---

<sup>2</sup>Thierry Massart, 2014. *Programmation*. Release 3.3.3 .

```

1  (*
2    Algorithmme to calcul Fibonacci
3  *)
4  begin
5    read(n) ;
6    a := 0 ;
7    b := 0 ;
8    for n from 0 to n do (* loop from 0 to n *)
9      tmp := b ;
10     b := a ;
11     a := tmp
12   done ;
13   print(b)
14 end

```

Figure 4: Implementation of the Fibonacci "*algorithm*" using *IMP*

```

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

```

Figure 5: Implementation of a sorting algorithm using *IMP*

As previously mentioned, *parsers* uses a specific structure called *grammar* which, similarly to spoken languages, aims to describe the allowed structures that a language can display. The *grammar* is composed of a set of variable to which are associated one or more rule(s). Typically, a programming language *grammar* is written following a fixed convention (see fig.8a or any other grammar). Generally, in the case of a programming language, *context-free grammar* are used because of its user-friendly aspect and the fact that it allows the use of an iterative development. However, *context-free* grammars are not adapted to implementation. Therefore, *parser generators* usually transform them into *LL(K)* grammars. In the assignment statement, it is asked to transform the *IMP* grammar into a *LL(1)*.

A *LL(1)* grammar is a class of grammars that can be defined by a predictive top-down parser. Such grammars define context-free languages which can be recognised by push-down automaton (i.e. automaton that use a stack as memory). We define a predictive parser as a type of parser that has an access to the input. This access is characterised by the possibility to know what is the following character to be read (this is very different from reading on the input) in order to enable the parser to take deterministic decisions throughout its execution.

### 3.1 Use and implementation of a parser generator

Nowadays, when designing a programming language, one might find convenient to use a program dedicated to the determinisation of the grammar and the syntax tree generation. Such programs are well known and have been developed a long time ago (yacc in the 1970's).

However, the assignment statement does not allow the use of such a program. Henceforth, it has been decided to implement a dedicated *parser generator* so that the given grammar will be modified without suffering from human mistakes. All the steps that a *parser generator* must achieved are shown

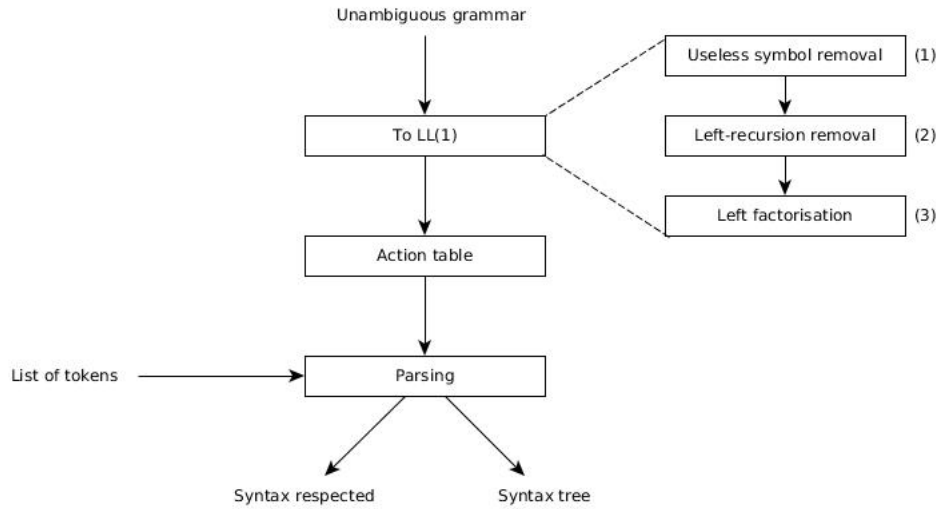


Figure 6: Phases of a parser generator

in the fig.6 and explained the section 3.2.1.

## 3.2 Parser phases

The present section is divided in three sections each explaining the utility of the different mechanisms or sub-mechanisms used by a *parser generator*. The first subsection concerns the transformation of a *context-free* grammar into an *LL(1)* grammar, the second subsection focuses on the creation of an *action table* and the last subsection explains how to check the syntax of an input.

A reader aware of the mechanisms used can easily skip the present section and can directly read the results (section 3.4).

### 3.2.1 Transforming the grammar to LL(1) grammar

In order to transform a given grammar that can be either deterministic or non-deterministic into a LL(1) — which is deterministic —, one must apply four transformations on this grammar. These transformations aim to make the grammars deterministic and thus implementable.

**Ambiguity removal :** Consists of ensuring, by the introduction of rule layers, that for a same input only one interpretation/derivation is possible. It also allows to force and to fix the priority and the associativity of some terminals in the grammar. (Further explanations see annex 5.4)

**Useless variable removal :** Consists of removing every variable that does not appear in any other rule of the grammar (as they will never be called/used) and removing variable that does not contain a rule capable of stopping the recursion of the other rules of the same variable. (Further explanations see annex 5.3)

**Left-recursion removal :** Consists of modifying every rules where the recursion occurs at the first element so that a grammar can still accept non-finite languages but becomes deterministic because the parser will eventually always encounter a terminal. (Further explanations see annex 5.2)

**Left-factorisation :** Consists of finding every rules of a variable that have a common prefix (i.e. same sequence of tokens) and of modifying it in consequence so that the parser can take deterministic decision. For example : which rule to choose between  $S \rightarrow \alpha A$  and  $S \rightarrow \alpha Z$  when the look ahead is  $\alpha$ . (Further explanations see annex 5.1)

### 3.2.2 Action table

An action table is a structure that emphasises the relation between a given variable and a terminal. Actually, it helps answering the following question : Which rule of a variable has to be applied if the next symbol to be read is  $x$  ?

To construct this structure, one has to introduce the function  $first(\alpha)$  that returns all the symbols that can be reached in one step (i.e. using only the first element of a rule). If the first element of the rule is a terminal,  $first(\alpha)$  adds this terminal to the returned set (the set of reachable variables). Otherwise, if the element is a variable, it calls  $first(\alpha)$  on the first element and then merges the returned set with the current one. Unfortunately, there is a special case : the epsilon-rule<sup>3</sup>. In fact, because  $\epsilon$  is neither a variable or a terminal and is — by definition — not expected on the input, one must, when encountering an epsilon, apply the function  $first(\alpha)$  on any element following an appearance of the variable that owns the epsilon-rule. The research of those elements is done by a function called  $follow(\alpha)$ .

### 3.2.3 Syntax checking

Once one has an action table, checking the syntax of the input is possible through the use of a simple stack that will only contain variables and terminals. The mechanism used to check the syntax of the input works as follows. Firstly, one starts by pushing on the stack the initial variable of the grammar. Then, if the element on the *tos* (top of the stack) is a variable, one pops it and identifies, using the action table, which rule to apply given the variable and the look ahead. Once identified, we push the rule on the stack (with the left-most element on the *tos*). Otherwise, if the element on the *tos* is a terminal, one pops it and compares it to the input. If it matches, one resumes, otherwise, it means that the syntax has not been respected.

In addition to that, one can add new accepting or rejecting conditions based on the "state" of the stack and/or the input. For instance in the case of IMP, the parser will rejects if there is no character to read on the input and that the stack is not empty or is there are characters left on the input and that the stack is empty. The only accepting configuration is when the stack is empty and there is no character left on the input.

For example, let's consider the final IMP grammar and the following input **begin a := b end**, the sequence of used rules and the stack utilisation will be as in the fig.7.

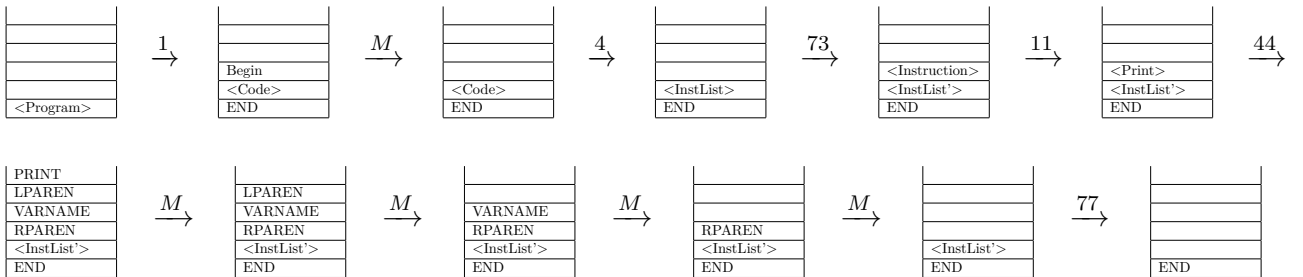


Figure 7: Stack utilisation for a simple program.

## 3.3 Resulting parser generator architecture

### 3.4 Tests and results

As previously mentioned, being able to transform a non deterministic — but yet rather simple to write — grammar into a implementable deterministic grammar, is something one wants to do when asked to implement a parser based on this grammar. In our case, all the methods previously presented in the section 3.2.1 have been applied to the IMP grammar (which is available in the annex fig.12) through the use of a package containing an implemented version of these methods.

<sup>3</sup>a variable rule which is only composed of an epsilon ( $\epsilon$ )



The present section is composed of six subsections. The four first subsections describe and analyse the steps of  $LL(1)$  transformation whereas the two last subsections focus respectively on the action table generated by the implementation and the results of the syntax checker.

For a matter of readability, it has been decided to only display the modified part of the grammar. To see the full version of the modified IMP grammar one must used the process described in the subsection X.X.

### 3.4.1 Ambiguity

Following the assignment statement, identifying the rules where ambiguity occurs is rather simple. The two ambiguous part of the grammar involved both arithmetic expressions and conditions. For both the former and the latter, atomic rules have been identified and categorised as such. Then, two layers have been introduced in order to force the derivation.

In the case of arithmetic expression, the operator  $-$  has been given the highest priority. Thus it has been introduced directly in the set of atomic rules (9 to 12). After that, the  $*$  and  $/$  operators were given the highest priority, this is why they are part of the first layer (5 to 8). Finally, the second and last layer (1 to 4) is composed of the rules involving the remaining operators ( $+$  and  $-$ ).

Notice that, even though it has not been explicitly specified in the assignment statement, any arithmetic expression surrounded by parenthesis has been considered as having a priority as high as the  $-$  operator.

```

1 <ExprArith>      -> <ExprArith> <OpAdd> <ExprArithMul>
2 <ExprArith>      -> <ExprArithMul>
3 <OpAdd>          -> +
4 <OpAdd>          -> -
5 <ExprArithMul>   -> <ExprArithMul> <OpMul> <ExprArithAtom>
6 <ExprArithMul>   -> <ExprArithAtom>
7 <OpMul>          -> *
8 <OpMul>          -> /
9 <ExprArithAtom>  -> VarName
10 <ExprArithAtom> -> [Number]
11 <ExprArithAtom> -> ( <ExprArith> )
12 <ExprArithAtom> -> - <ExprArithAtom>

```

The same modifications have been applied on conditions than on arithmetic expressions.

```

1 <Cond>           -> <Cond> or <CondAnd>
2 <Cond>           -> <CondAnd>
3 <CondAnd>        -> <CondAnd> and <CondAtom>
4 <CondAnd>        -> <CondAtom>
5 <CondAtom>       -> not <SimpleCond>
6 <CondAtom>       -> <SimpleCond>
7 <SimpleCond>     -> <ExprArith> <Comp> <ExprArith>
8 <Comp>           -> =
9 <Comp>           -> >=
10 <Comp>          -> >
11 <Comp>          -> <=
12 <Comp>          -> <
13 <Comp>          -> <>

```

### 3.4.2 Useless symbols removal

The unambiguous IMP grammar does not contain any useless symbols. One can be easily convinced of the reachability by drawing a graph where every variable is represented by a node and every edge represents the fact that a variable appears at least once in one the rules of the other. One can also be convinced of the grammar productiveness by observing that for each recursive rule, there is another that stops the recursivity. These intuitions were proved right by the implementation.

### 3.4.3 Left-recursion removal

When one wants to remove the left-recursion, one knows, following the definition given above, that one has to look for rules where the left-hand side is also the first element of the right-hand side.

Doing so on the unambiguous and useless symbols free IMP grammar returns once again the arithmetic expressions and the conditions. But this is not surprising given the trick used to make the grammar unambiguous. This explains the reasons behind the order of the steps.

In both cases, there are the introduction of a suffixed U variable and a suffixed V variable. These variables respectively used to remove indirect left-recursion and to transform each direct left-recursion into a right-recursion. Notice that the introduced right-recursion are productive as the algorithm stipulates that the suffixed V variable must own an epsilon-rule<sup>4</sup>.

```
1 <ExprArith>      -> <ExprArithU> <ExprArithV>
2 <ExprArithU>     -> <ExprArithMul>
3 <ExprArithV>     -> <OpAdd> <ExprArithMul> <ExprArithV>
4 <ExprArithV>     -> eps
5 <ExprArithMul>   -> <ExprArithMulU> <ExprArithMulV>
6 <ExprArithMulU>  -> <ExprArithAtom>
7 <ExprArithMulV>  -> <OpMul> <ExprArithAtom> <ExprArithMulV>
8 <ExprArithMulV>  -> eps
```

```
1 <Cond>           -> <CondU> <CondV>
2 <CondU>          -> <CondAnd>
3 <CondV>          -> or <CondAnd> <CondV>
4 <CondV>          -> eps
5 <CondAnd>        -> <CondAndU> <CondAndV>
6 <CondAndU>       -> <CondAtom>
7 <CondAndV>       -> and <CondAtom> <CondAndV>
8 <CondAndV>       -> eps
```

### 3.4.4 Factorisation

Looking in a grammar for rules to factorise consists in identifying variables that have two or more rules that share a common prefix (i.e. a same sequence of tokens). These rules are then modified following the mechanism explained in the subsection 5.1.

In the case of the unambiguous, useless symbols free and left-recursion free IMP grammar, only three variables need to see their rules factorised : **InstList**, **If** and **For** (respectively line 4, 22 and 37 in fig.12).

```
1 <InstList>       -> <Instruction> <InstList'>
2 <InstList'>      -> ; <InstList>
3 <InstList'>      -> eps
```

The factorisation of **InstList** is an instance of the particular case mentioned in subsection 5.1. In fact, one of the rule to factorised does not really diverge as it is equal to the suffix. Thus an epsilon-rule<sup>4</sup> is introduced (rule 3).

```
1 <If>             -> if <Cond> then <Code> <If'>
2 <If'>           -> else <Code> endif
3 <If'>           -> endif
```

```
1 <For>            -> for VarName from <ExprArith> <For'>
2 <For'>           -> by <ExprArith> to <ExprArith> do <Code> done
3 <For'>           -> to <ExprArith> do <Code> done
```

The factorisation of **If** and **For** are quite mainstream as they present a real divergence. Thus after the prefix, a new variable is introduced and associated to the remaining of the common prefixed rules.

---

<sup>4</sup>a variable rule which is only composed of an epsilon ( $\epsilon$ )

### 3.4.5 Resulting action table

### 3.4.6 Syntax checking result

## 4 How to set up the project

In order to simplify the compilation and the support of external libraries, it has been decided to use a well known *java* project manager named *Maven*. Its configuration file (`pom.xml`) defines the `main` file, defines the source folder, manages the *JFlex* library and the package that must be compiled with this library.

### 4.1 Compilation

Compiling the project with *Maven* is easy as the user only needs to execute : `mvn clean compile`. However, at the first execution, the user needs to execute `mvn install` so that *Maven* can install the required library.

If the user does not want to use *Maven*, he can execute different commands from the root project :

```
java -jar jflex-1.6.1.jar -d src/be/ac/ulb/infof403/ src/be/ac/ulb/infof403/lex/Scanner.flex
```

Where `jflex-1.6.1.jar` is the path to the `.jar` executable library, `-d` is the output folder path specifier and the last parameter is the path to the `.flex` file.

Then, the user can compile the java source codes and can create the corresponding `.class` files. The bash command to compile all the *java* files located in the `src/` folder is the following :

```
javac -d target $(find ./src/* | grep .java)
```

This command generates the corresponding `.class` files and put them in the `target/` folder. You must create the "target" folder if it does not currently exist. Finally, the *jar* file can be generated by using the command :

```
jar cvfe dist/INFO-F403-IMP.jar be/ac/ulb/infof403/Main -C target/ .
```

Where `INFO-F403-IMP.jar` is the name of the generated *jar* file and *target* is the folder where are located the `.class` files.

### 4.2 Execution

To execute the resulting jar file, the user only has to type :

```
java -jar dist/INFO-F403-IMP.jar <sourceFile>
```

Where `<sourceFile>` is the path to the IMP file. If the source file is not specified then the program will use the file `test/Euclid.imp`.

### 4.3 Test

The program has a system which automatically compares each output file (`.out`) to the result of the execution of the corresponding `.imp` file. The execution of the test can be specified by adding the parameter `-test` at the program execution instruction, like this :

```
java -jar dist/INFO-F403-IMP.jar <sourceFile> -test <testFile>
```

Where `<testFile>` is the name of the output file. If not specified, the program will automatically load a file test based on the *source file* name. It will only change the file extension from `.imp` to `.out`.

## 4.4 Javadoc

The javadoc is located in the `doc/` folder. To generate the javadoc with Maven you must execute `mvn javadoc:javadoc`. If you do not want to use Maven, you could execute the following command :

```
javadoc -d doc/javadoc/ -keywords -sourcepath src -subpackages be
```

Where `doc/javadoc/` is the output folder. The option `-keywords` enable HTML in the javadoc.

## 5 Annexe

### 5.1 Grammar factorisation

This mechanism is applied every time a given variable has two (or more) rules that have a common prefix. The aim is to reduce the number of repetitions. To achieve this, each variable that has two (or more) rules with a common prefix sees these rules replaced by concatenation of the prefix and a new variable. This new variable has for rules the remaining of the factorized rules (i.e. the rules that have a common prefix without this prefix).

		$S \rightarrow relationship$	(5)
$S \rightarrow friendship$	(1)	$\rightarrow friendS'$	(6)
$\rightarrow friend$	(2)	$S' \rightarrow ship$	(7)
$\rightarrow relationship$	(3)	$\rightarrow ly$	(8)
$\rightarrow friendly$	(4)	$\rightarrow \epsilon$	(9)
(a) Unmodified grammar		(b) Factorisation outcome	

Figure 8: Caption place holder

For instance, in the fig.8a, the rule (3) has no prefix while the other rules whereas (1), (2) and (3) have a common prefix : '*friend*'. Thus, following the mechanism explained above, we replace these three rules by a new one (rule (6)) composed of the prefix ('*friend*') and the new variable ( $S'$ ). The variable  $S'$  is then associated with the remaining of each former rule with a common prefix of  $S$ . Notice that the rule (2) is a particular case as it matches exactly the prefix. To overcome this issue, the created rule is formed of  $\epsilon$  (rule (6)).

Such a technique is used to ensure that the parser will be deterministic. In our case, we want to implement a parser with a look ahead of one. Therefore, if a variable has two (or more) rules like  $S \rightarrow fA$  and  $S \rightarrow fZ$ , the parser won't be able to decide which one to apply.

### 5.2 Removing left-recursion

Even though recursion is a main feature of grammars as it allows them to recognise non-finite language, it also introduces non-determinism when the recursion occurs at the very first element of the right-hand side. To make a grammar (and thus the parser) deterministic but keep the recursivity, one must execute some manipulations.

First, one wants to transform every indirect left-recursion into direct left-recursion. Achieving that is quite simple as one only has to take a rule and replace every variable located at the very beginning of the left-hand side and replace it by all of its own rules. For instance, in fig.9a, the grammar is indirectly recursive because  $S$  calls  $S'$  which, when applying rule (11), calls  $S$ . The outcome of this transformation (see fig.9b) recognises the same language but is now directly left-recursive.

Secondly, one wants to transform every left-recursion by a right-recursion for determinism purpose (similar to factorisation). One can achieve it by introducing two new variables. The first variable will

$S \rightarrow S'b$	(10)		$V \rightarrow aV'$	(15)
$S' \rightarrow Sa$	(11)	$S \rightarrow Sab$	$V' \rightarrow abV'$	(16)
$\rightarrow \epsilon$	(12)	$\rightarrow a$	$\rightarrow \epsilon$	(17)
(a) Unmodified grammar	(b) Indirect recursion removal	(c) Transformation to right-recursive		

Figure 9: Caption place holder

be associated to a set of rules each composed of the concatenation of a non-recursive rule and the newly created second variable. This second variable will be associated with a set of rules composed of every recursive rules where the first element (the recursive variable) has been removed concatenated with this exact second variable. Doing so transforms every left-recursion in a right-recursion. However, this right-recursion will never stops. This is why a rule composed of  $\epsilon$  is associated to the second variable.

### 5.3 Removing useless variables

When speaking of *useless* variables, we distinguish two types of variables :

**The *unproductive* ones :** An unproductive variable is a variable that never leads to any formation of a word. Typically, such a variable does not have any non-recursive rule. Thus, forming a word using this variable leads to an infinite recursion.

**The *unreachable* ones :** An unreachable variable is a variable that is not called by any other rule of the grammar in which it belongs.

So far, the best way to find both unproductive and unreachable variables is to look respectively for productive and reachable variables and remove them from the grammar afterwards. However, eventually, we are only interested in productive and reachable variables. Thus, once the former and the latter are found, we consider them as the final grammar.

Determining the set of productive symbols consists of first considering every terminal as productive. Then, for each variable, we look at each rule and add the variable to the set if and only if every symbols appearing in the rule are already in the set. The resulting set is the set of every reachable symbols of the given grammar.

Retrieving the set of reachable symbols from a given grammar can be achieved by using a similar method to the one explained above. In fact, one must consider first a set containing only the initial variable of the grammar which is — without lost of generality — always considered as reachable. Then, for each variable of the grammar, one must check whether the variable is in the set of reachable symbols. If yes, one can then add all the symbols appearing in the rules of this variable.

### 5.4 Ambiguous grammar

Ambiguity occurs when, for a given word/input, multiple interpretations (or trees) can be derived due to an *ambiguity* in the rules the parser has to choose making it non-deterministic and thus in proper for any implementation. Unfortunately, there does not exist any algorithm resolving this issue as the given grammar gives little information. Henceforth, extra information only known by the language designer must be integrated. The most common example of ambiguity is the arithmetic priority (Reminder : the multiplication has an higher priority than the addition).

For example, applying the grammar of fig.10a on the word  $id + id * id$  will result in two different interpretation as shown in Fig.11a and Fig.11b.

To address this issue, the language designer must 'force' the derivation (and hence the priority) by introducing new variables that could be seen as extra layers. For instance, on fig.10a, the grammar is composed of two *atomic* terminals : *Cst* and *Id*. These terminals will be encapsulated in a new

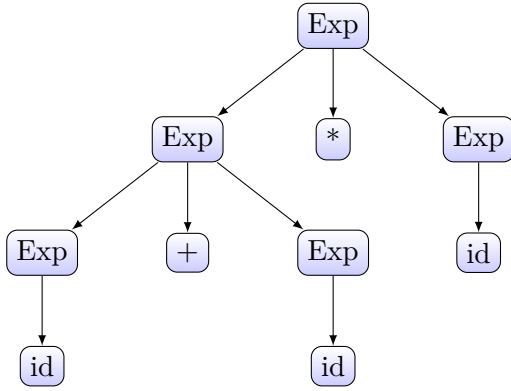
$$\begin{aligned}
Exp &\rightarrow Exp + Exp & (18) \\
&\rightarrow Exp * Exp & (19) \\
&\rightarrow Cst & (20) \\
&\rightarrow Id & (21)
\end{aligned}$$

(a) Ambiguous grammar

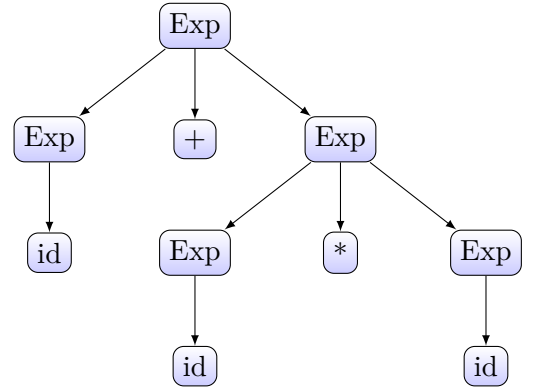
$$\begin{aligned}
Exp &\rightarrow Exp + Prod & (22) \\
&\rightarrow Prod & (23) \\
Prod &\rightarrow Prod * Atom & (24) \\
&\rightarrow Atom & (25) \\
Atom &\rightarrow Cst & (26) \\
&\rightarrow Id & (27)
\end{aligned}$$

(b) Unambiguous grammar

Figure 10: Caption place holder



(a) First possible derivation of the input.



(b) Second possible derivation of the input.

Figure 11: Derivations of the input  $id + id * id$  using the grammar in fig.10a

variable called *Atom*. In addition, we decide — based on the arithmetic priority — that multiplication has an higher priority than addition. Therefore, as for atomic elements, we introduce a new variable called *Prod* that has for rules a single atomic value (25) and the product of a multiplication and a atomic element (24). Finally, the same mechanism is once again applied to addition. Resulting in the rules (22) and (23).

As previously mentioned, there does not exist an algorithm that resolves grammar ambiguity. However, there exists many ambiguity detection algorithms with have their own proprieties as mentioned in this article<sup>a</sup>. The reader is invited to read this document for further information.

<sup>a</sup>H.J.S. Basten, August 17, 2007. *Ambiguity Detection Methods for Context-Free Grammars*. Master's Thesis, Universiteit Van Amsterdam.

## 5.5 Basic IMP grammar

## 5.6 IMP action table

	(	<=	begin	if	do	-	read	while	and	then	:=	VarName	)	else	end	done	not	by	from	>=	<>	to	print	for	or	+	[Number]	;	=	/	*	<	endif	>		
If				67																																
CondAndU	61				61							61					61										61									
InstList				73			73	73				73											73	73												
Assign												13																								
ExprArith	50				50							50															50									
Program			2																																	
ExprArithMulV		53			53	53			53	53				53	53	53		53	53	53	53	53		53	53	53	53	53	53	53	53	53	53	53	53	
CondV				58						58																57										
OpMul																																				
ExprArithMulU	51				51							51																51								
ExprArithMul	55				55							55																55								
For'																			81			83														
If'														69																					71	
CondAtom	33				33							33					32											33								
ExprArithU	46				46							46																46								
InstList'														77	77	77														75				77		
Instruction				8			12	9				7											11	10												
CondU	56					56						56					56											56								
OpAdd					17																						16									
ExprArithV		48		48	47				48	48				48	48	48			48	48	48	48		48	47											
CondAndV				63					62	63															63											
Read							45																													
While								41																												
SimpleCond	34				34							34																34								
ExprArithAtom	24				25							22																23								
Comp		38																																		
Cond	60				60							60					60											60								37
For																							79													
Code				4			4	4				4		3	3	3							4	4											3	
CondAnd	65				65							65					65																			
Print																							44													

Table 1: IMP action table.

```

1 <Program>      -> begin <Code> end
2 <Code>         -> eps
3               -> <InstList>
4 <InstList>     -> <Instruction>
5               -> <Instruction> ; <InstList>
6 <Instruction>  -> <Assign>
7               -> <If>
8               -> <While>
9               -> <For>
10              -> <Print>
11              -> <Read>
12 <Assign>       -> [VarName] := <ExprArith>
13 <ExprArith>    -> [VarName]
14               -> [Number]
15               -> ( <ExprArith> )
16               -> - <ExprArith>
17               -> <ExprArith> <Op> <ExprArith>
18 <Op>           -> +
19               -> -
20               -> *
21               -> /
22 <If>           -> if <Cond> then <Code> endif
23               -> if <Cond> then <Code> else <Code> endif
24 <Cond>         -> <Cond> <BinOp> <Cond>
25               -> not <SimpleCond>
26               -> <SimpleCond>
27 <SimpleCond>   -> <ExprArith> <Comp> <ExprArith>
28 <BinOp>        -> and
29               -> or
30 <Comp>         -> =
31               -> >=
32               -> >
33               -> <=
34               -> <
35               -> <>
36 <While>        -> while <Cond> do <Code> done
37 <For>          -> for [VarName] from <ExprArith> by <ExprArith> to <ExprArith> do <
    Code> done
38               -> for [VarName] from <ExprArith> to <ExprArith> do <Code> done
39 <Print>        -> print ( [VarName] )
40 <Read>         -> read ( [VarName] )
41

```

Figure 12: The basic IMP grammar as given in the assignment statement.