

# ALGOL-0

## Introduction to language theory and compiling

### Project – Part 3

Gilles GEERAERTS

Raphaël BERTHON

Sarah WINTER

November 11<sup>th</sup>, 2019

## Statement

For this third and last part of the project, we ask you to augment the recursive-descent LL(1) parser you have written during the first and second parts to let it *generate code* that corresponds to the semantics of the ALGOL-0 program that is being compiled (those semantics are informally provided in Appendix A). The output code must be LLVM intermediary language (LLVM IR), which you studied during the practicals. This code must be accepted by the `llvm-as` tool, so that it can be converted to machine code. It is **not allowed** to first compile to another language (*e.g.* C) and then use the language compiler to get LLVM IR code.

When generating the code for arithmetic and boolean expressions, pay extra attention to the associativity and priority of the operators.

## Bonus

For this last part of the project, you can enrich ALGOL-0 with several features:

- syntactic sugar/simple extensions, *e.g.* allowing parentheses in boolean expressions
- variable scoping, *i.e.* handle programs such as:

```
BEGIN
  x := 42;                co this is a global variable
  FOR x FROM 0 BY 1 TO 5 DO co this is a local variable
    PRINT(x)
  ENDWHILE;
  PRINT(x);
END
```

which should then output 0 1 2 3 4 5 42

- functions (you can allow `BEGIN/END` to appear multiple times, or introduce additional keywords)
- additional types: boolean<sup>1</sup>, real numbers (`double`) and strings, or even arrays, lists, ...
- recursive functions

---

<sup>1</sup>This type is already implicitly present in ALGOL-0 with `<SimpleCond>`, but you could allow assigning boolean to variables. You could then implement a basic type-checker.

They are sorted by increasing difficulty, but you can choose any, or all. If you would like to add a feature that is not in the list, tell us first. You can also provide compiling optimizations (*dead code elimination, inlining, ...*), but this is less rewarding for you<sup>2</sup> since such optimisations are most likely provided by LLVM.

You have entire freedom of implementation for these bonuses; in particular, you can enrich the keywords and syntax, as long as ALGOL-0 is a subset of your language (any ALGOL-0 program must compile correctly). You are however required to explain what you did and how you did it in your report (we are not supposed to guess how your program works, nor what bonus you implemented), and you should provide test files to demonstrate your additional features. If you have any questions, please send us an email.

This bonus can get you *up to one point*, which will be added to your overall project grade (*i.e.* it can compensate for points you lost on previous parts). However, if you obtain more than 8/8, it will *not* be passed on to your exam grade. Also, note that this is only a *bonus*: it is better to provide a working ALGOL-0 compiler than a buggy Mega-ALGOL-0 one.

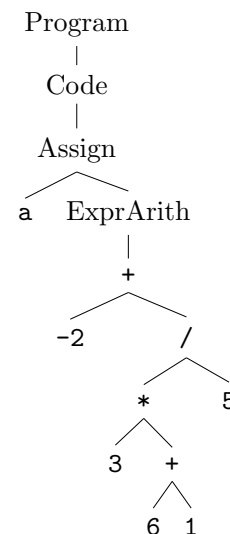
## Guidelines

Those are only guidelines, and you have entire freedom of implementation for the code generator (except of course using compiling tools such as CUP).

To generate code from your parser, you can first modify it so that it builds an an Abstract Syntax Tree, either by doing it directly during the parsing or by post-processing the Parse Tree. Abstract Syntax Trees are meant to abstract away some non-terminals to obtain the very structure of the program. For instance, consider the program in Fig. 1a:

```
BEGIN
  a := -2 + 3 * (6 + 1) / 5
ENDPROG
```

(a) A very simple program



(b) An abstract representation of the parse tree of the program in Fig. 1a.

Figure 1: A simple program and its corresponding AST.

Its parse tree can then be converted to the AST of Figure 1b. Note that here, since  $*$  and  $/$  have the same priority, the parsing is done left-associatively, *i.e.* the arithmetic expression should be parenthesized as  $-2 + (3 * (6 + 1)) / 5 = 2$  and not  $-2 + 3 * ((6 + 1) / 5) = 1$  (this matters because “/” is the *euclidean* division).

---

<sup>2</sup>In terms of interest, not in terms of grading: this will give you as many bonus points as the previous ones.

You should draw inspiration from this example to design an abstract representation of <If>, <For>, <While>, and the like. Note in particular that some non-terminals were omitted, since they have no semantic meaning; they are only dummy non-terminals used to simulate lists, *e.g.* <InstList>, enforce priorities of operators, or remove left-recursion.

That being done, you can walk the AST and generate code. For example, when walking a node



Tree1 Tree2, you should first evaluate the expression represented by Tree1, then evaluate the expression represented by Tree2, then compute the result of the addition and store it in an unnamed variable (whose number should be deduced from the ones used in Tree1 and Tree2). It should give something of the form:

```

[sequence of instructions to evaluate Tree1]
[%n is assigned the result of Tree1]
[sequence of instructions to evaluate Tree2, unnamed variables start at n + 1]
[%m is assigned the result of Tree2]
%p3 = %n + %m
  
```

Of course, you can also directly generate the code from the parse tree, or obtain the AST directly from the parser, or even do everything at once, but this might be more complex: this decomposition separates the difficulties.

## Requirements

You must hand in:

- A PDF report containing the necessary justifications, choices and hypotheses;
- The source code of your compiler;
- The ALGOL-0 example files you used to test your compiler;
- All required files to evaluate your work.

You must structure your files in five folders:

- **doc** contains the JAVADOC and the PDF report;
- **test** contains all your example files;
- **dist** contains an executable JAR **that must be called part3.jar**;
- **src** contains your source files;
- **more** contains all other files.

Your implementation must contain:

1. your scanner if you were able to use it for parsing, or the one we provided;
2. your parser, except if it does not work properly, in which case you may use ours<sup>4</sup>;

---

<sup>3</sup>At this point, you can deduce the value of  $p$ .

<sup>4</sup>It will be available on Université Virtuelle in the beginning of the last week of November

3. an executable public class `Main` that reads the file given as argument and writes on the standard output stream the LLVM intermediary code. You can provide, as **options**<sup>5</sup>, the possibility to output it to a `.ll` file, and to interpret the code (for instance, `java -jar part3.jar inputFile.alg -o outputFile.ll` would output the LLVM IR code to `outputFile.ll`, and `java -jar part3.jar -exec inputFile.alg` would interpret the code).

The command for running your executable must be: `java -jar part3.jar inputFile`

You need to compress your folder (in the *zip* format—no *rar* or other format), **following the same naming convention as for Part 2: Part3.Surname1(\_Surname2)?.zip** where `Surname1` and, if you are in a group, `Surname2` are the last names of the student(s) of the group (in alphabetical order). You need to submit it on the Université Virtuelle before **December, 22<sup>nd</sup>**.

## A Informal semantics of the ALGOL-0 language

We only provide an informal description of the semantics, since a formal one would needlessly complicate the matter for such a simple language. There is nothing surprising here, since those semantics are similar to the ones of everyday languages. The value to which a nonterminal `<NT>` evaluates will be denoted by  $\llbracket \text{NT} \rrbracket$ , *e.g.*  $\llbracket \text{ExprArith} \rrbracket$ .

- The code represented by `<Program>` should be the result of the processing of `<Code>` (in other words, the `BEGIN` and `END` markers are just markers).
- `<Code>` is a list of instructions `<Instruction>`, which should be executed sequentially.
- `<Assign>`: `[VarName] := <ExprArith>` means the program should store  $\llbracket \text{ExprArith} \rrbracket$  in the variable `VarName` (which should be stored in a memory location, not simply in an LLVM variable).
- `<If>`: `IF <Cond> THEN <Code> ENDIF` means that if the condition computed by `<Cond>` (*i.e.*  $\llbracket \text{Cond} \rrbracket$ ) is true, then `<Code>` should be executed, otherwise the program should go to the next instruction.
- `<If>`: `IF <Cond> THEN <Code1> ELSE <Code2> ENDIF` means that if  $\llbracket \text{Cond} \rrbracket$  is true, then `<Code1>` should be executed, otherwise `<Code2>` should be executed instead.
- `<While>`: `WHILE ( <Cond> ) DO <Code> ENDWHILE` means that the program should test `<Cond>`, then execute `<Code>` if  $\llbracket \text{Cond} \rrbracket$  is true and then repeat, otherwise it should do nothing<sup>6</sup>.
- `<For>`: `FOR [Varname] FROM <ExprArith1> BY <ExprArith2> TO <ExprArith3> DO <Code> ENDWHILE` means that the program should initialize the variable named `VarName` to  $\llbracket \text{ExprArith1} \rrbracket$ . Then, check if the value of `Varname` is less or equal than  $\llbracket \text{ExprArith2} \rrbracket$ , if so, execute `<Code>` and increment the value `VarName` by  $\llbracket \text{ExprArith2} \rrbracket$ , and so on until the value of `VarName` exceeds  $\llbracket \text{ExprArith3} \rrbracket$ .
- `<Print>`: `PRINT([Varname])` should print the value of `[Varname]` to `stdout`.

<sup>5</sup>It is important that, by default, your program outputs the LLVM IR code on `stdout` (even if you use a temporary file to store the LLVM IR code, in which case you just have to print it on screen).

<sup>6</sup>Giving formal semantics to `WHILE` is actually very hard. Here is a hint of how it could be done: `<While>` can be unrolled as `IF <Cond> THEN <Code> WHILE <Cond> DO <Code> ENDWHILE ENDIF`.

- <Read>: `read([Varname])` should read an integer from `stdin`, and store it in the corresponding `[Varname]`.
- <ExprArith>: those are the semantics of usual arithmetic expressions written in infix notation, with the conventional precedence of operators (given in Part 2).
- <Cond>: those are the semantics of usual boolean expressions with only operators *and*, *or* and *not*, with the conventional precedence of operators (given in Part 2).