

DPCC: DParo's Own C-Alike Compiler

Davide Paro

December 2020

Assignment Description

This project is the implementation of an assignment for a course on **Compilers** at the department of Computer Engineering Master Degree Padova (ITA).

The assignment consists in implementing a toy compiler (mostly the frontend side) for a toy language. We are free to design the syntax of this toy language however we like.

The assignment specs out the how the compiler should be composed. We can in fact distinguish these macro components:

- **Input Stage** deals with the input byte stream that composes the source of the program.
- **Lexer/Scanner** has the purpose of grouping characters (lexical analysis) together to compose compounded structures (called tokens). For the project assignment we can use **Flex** to aid in the code generation for the scanner.
- **Parser** for performing the syntax analysis. It is what defines the look & fell (grammar) of the language. For the project assignment we can use **Bison** to aid in the implementation of a good parser.
- **Intermediate Code Generator**. The ultimate purpose of a compiler is to produce something useful. In this project assignment we are not requested to implement a proper backend. Thus, we need to emit a 3AC representation of our input program. More in this later.

In particular the final Intermediate representation that we need to emit is based on Three Address Code (3AC), that is, each statement can only have 1 operand at the left hand side of the assignment, and 2 operands at the right hand side of the assignment, and an operator driving the operation that should be performed.

You can view more about 3AC at the following wikipedia link.

In practice the emitted 3AC code is on itself a partially valid C program, it's only missing variable declarations at the top for the temporary variables.

For the specification of the Intermediate Code that is generated please refer to appendix A

So the project requires to produce this kind of 3AC / C hybrid. Control flow is allowed to be implemented through the usage of C labels and simple if conditional followed by a goto statement. Inside the if conditional there can only be a single element composing the expression.

The assignment requires the following features from the programming language that we should develop:

- Variables declaration, initialization and assignment
- Handling of variable scopes. Variable names can be reused if out of scope. Variable shadowing may or may not warn/fail/pass depending on the design choices.
- Only 2 types of variables: integers, booleans
- Assignment statements, print statements, if statements, and at least 1 loop statement at will (while, for, ...)
- Handling of simple mathematical expressions that we can encounter in common programming languages, addition, subtraction, multiplication, division, modulo, etc ...
- **Function definition, function calls, and custom user definable types are not required**

A quick peak at the language and at the compiler

This section is concerned with showing the reader, what the final language, compiler proposed in this implementation will look like. All the internal details are instead deferred to later sections.

It is best to show a simple example

Structure of the language

The proposed language is very similar to a C and JavaScript Hybrid. Most of the syntax is derived from C-alike language, and some keyword concepts are instead taken from Javascript. That is newlines and spaces mostly do not matter, they simply introduce token boundaries.

```
1 // C style Multi-line comments like /* ... */ pair are not supported
2
3 // Print statement with immediate C-style strings. C-style strings are not
4 // a valid type and can only be passed to print
5 print("Hello world\n");
6
7 // Variable declaration and initialization
8 let a = 10;          // Integer Type deduced
9 let f = 10.0;        // Float type deduced
10 let b = false;      // Boolean type deduced
11
12 // Explicit types
13 let i: int = 0xffff & ~0xb00111;
14 let f: float = 10.0 + 20.0 ** 2;
15
16 // Variables can be printed
17 print(i);           // Print integer
18 print(f);           // Print float
19
20 // Scoping and restricting variable declarations to the current scope
21 {
22     // Simple single dimension arrays declaration
23     let buf_i: int[100];        // Known size integer array
24     let buf_f: float[100];      // Known size float array
25
26     // Integer array with deduced size from the RHS initializer list
27     let buf = [ 10, 20, 30, 40, 50 ];
28
29     // Arrays can be printed
30     print(buf);
31 }
32
33 let buf: int[100];
34
35 // Control flow
36 for (let i = 0; i < 100; i++) {
37     buf[i] = i ** 2;
38
39     if (buf[i] == 10) {
40         print("buf[i] is 10!!!\n");
41     }
42     else if (buf[i] == 20) {
43         print("buf[i] is 20!!!\n");
44     }
45     else {
46         print("None of above\n");
47     }
48 }
```

So as the reader can see from above there's not much fancy stuff about the syntax, it is mostly C-alike. The cool thing though about this implementation is that we can mostly declare variables with a simple 'let' keyword like in Javascript, but unlike Javascript the compiler has a way to deduce types whenever possible and still provides strong type semantics to each identifier like in C (or any other strongly typed language).

For now only 5 types are implemented: **bool**, **int**, **float**, **string**, **bool[]**, **int[]**, **float[]**. Only single dimensions array are for now supported. So arrays do not generalize to any dimensions.

Most of these types have full on semantics, that is the compiler can deduce a type of an expression starting from the type of its operand. In some cases it can reject the code maybe because an operator is used with invalid typed operands.

String for now is kind of a quirk type, meaning that it doesn't have a full type tracking inside the compiler. The compiler still knows what a string is, but strings cannot be assigned to any variables and they can instead be used only as a parameter to the print statement.

Calling the compiler

The proposed implementation

In the following sections it is described how the proposed

- **Input Stage.** The compiler only supports file. The input stream is implemented in the following way. The file is loaded and entirely copied into a memory buffer. This is more than adequate for what it's necessary to do for the project assignment.
- Un **Lexer/Scanner** per la tokenizzazione del sorgente da caratteri a tipi di dati strutturati. La scelta ricade su **Flex** per la gestione dell'analisi lessicale
- Un **Parser** per implementare la sintassi del linguaggio e gestire l'analisi sintattica. La scelta ricade su **Bison** per la gestione di questa componente
- Un semplice **generator di codice intermedio**. Il progetto prevede di generare un ibrido Assembly/C/3AC come semplice esempio di gestione di generazione

Implementation Details

Custom log is implemented to override default bison behaviour

Common types used withing the compiler

Most notable functions used in the compiler

How the AST is traversed

How type checking and type deduction is implemented

A brief introduction at the code implementing the compiler

- The order of the childs in each node **is very important** and it drives the correctness of the final output. In this way we can use simple printf to generate the code each time we find something while traversing the tree

Typescript used for code generation of some part of the compiler

Encountered problems

Testing framework

Problem analysis

...

Program design

...

Evaluation of the program

...

Process description

...

Conclusions

...

Appendix A: Structure of the Intermediate Code

Appendix: program text

```
1 let a: int[] = 10;  
2 let s = "Hello world";  
3 {  
4  
5 }
```

```
1 int main() {  
2  
3 }  
4 char **argv;
```

Appendix B: Example Program Iterative Merge Sort

```
1  let array = [  
2      15, 59, 61, 75, 12, 71, 5, 35, 44,  
3      6, 98, 17, 81, 56, 53, 31, 20, 11,  
4      45, 80, 8, 34, 71, 83, 64, 28, 3,  
5      88, 50, 48, 80, 5  
6  ];  
7  
8  
9  for (let curr_size = 1; curr_size < len; curr_size = 2 * curr_size) {  
10     for (let left_start = 0; left_start < len - 1; left_start = left_start + 2 * curr_size) {  
11         let mid = len - 1;  
12  
13         if ((left_start + curr_size - 1) < len - 1) {  
14             mid = left_start + curr_size - 1;  
15         }  
16  
17         let right_end = len - 1;  
18  
19         if ((left_start + 2 * curr_size - 1) < len - 1) {  
20             right_end = left_start + 2 * curr_size - 1;  
21         }  
22  
23         {  
24             let l = left_start;  
25             let m = mid;  
26             let r = right_end;  
27             let n1 = m - l + 1;  
28             let n2 = r - m;  
29  
30             let L: int[1024];  
31             let R: int[1024];  
32  
33             for (let i = 0; i < n1; i++) {  
34                 L[i] = array[l + i];  
35             }  
36  
37             for (let i = 0; i < n2; i++) {  
38                 R[i] = array[m + 1 + i];  
39             }  
40  
41  
42             let i = 0;  
43             let j = 0;  
44             let k = l;  
45  
46             while (i < n1 && j < n2) {  
47                 if (L[i] <= R[j]) {  
48                     array[k++] = L[i++];  
49                 } else {  
50                     array[k++] = R[j++];  
51                 }  
52             }  
53  
54             while (i < n1) {  
55                 array[k++] = L[i++];  
56             }  
57             while (j < n2) {  
58                 array[k++] = R[j++];  
59             }  
60         }  
61     }  
62 }  
63  
64 print("Sorted array\n");  
65 print(array);
```