

Output C Code

Computer Language Processing '16 Final Report

Cédric Viaccoz Bastian Nanchen

EPFL

cedric.viaccoz@epfl.ch

1. Introduction

The first part of the lab was to implement a compiler for Toy Object-Oriented Language (Tool), which is a small object-oriented programming language.

This has led us to conceive a Lexer, a $LL(1)$ grammar, a Parser, a Type Checking system and to generate binary code for the Java Virtual Machine.

The aim of this project is to replace code generation for the Java Virtual Machine with a pretty printer, which outputs C code.

The difficulty is to make object-oriented code into C code which does not have built in data types to support classes, as it is an imperative programming language.

2. Examples

The example below presents how a class and the inheritance could be depicted in the C programming language.

```
1 class Animal{
2     var isAPet: Bool;
3     var speed: Int;
4     def sleep(): String = {
5         return "zzz...";
6     }
7 }
8
9 class Dog extends Animal{
10     def bark(): String = {
11         return "Woof!";
12     }
13
14     //override
15     def sleep(): String = {
16         return "zzz... wouaf... zzz...";
17     }
18 }
```

This Tool code will be pretty-printed in C like this:

```
1 struct Animal{
2     int isAPet;
3     int speed;
4     char * (*sleep)(void *);
5 };
6
7 struct Dog{
8     int isAPet;
9     int speed;
10    char * (*sleep)(void *);
11    char * (*bark)(void *);
12 };
13
14 char * Animal_sleep (void * this) {
15     // some code that will be depicted in the
16     // implementation part
17 }
18
19 char * Dog_sleep (void * this) {
20     // some code that will be depicted in the
21     // implementation part
22 }
23
24 char * Dog_bark (void * this) {
25     // some code that will be depicted in the
26     // implementation part
27 }
```

Classes are represented in C by a structure. The variables of the class are members of the structure. They are declared in the same fashion.

The variables of the class “Animal”, are also located in the “Dog” structure because the class “Dog” inherits from the class “Animal”.

Contrariwise to Tool, the methods of the class are declared outside of the structure. In C, a structure can't have a function declaration as a member, thus the functions should be declared outside, and the structure holds a pointer to the function.

The function is declared with the name of the class as a prefix followed by an underscore. Thus the belong-

ing to the “Dog” class is suggested, and this allows overriding (a same function with two implementations according to the class).

Therefore this is how the C code mimics the behavior of an object-oriented class.

There exists other dilemmas like method chaining, dynamic dispatch,..., which will be discussed in the implementation section.

3. Implementation

3.1 Theoretical Background

The only background needed for this extension is knowing the AST of Tool and the C language. However since we are trying to emulate an object-oriented language into one that is not, the book “Object-oriented programming with ANSCI-C”[?] has been useful for some concepts and it will be referenced throughout this report.

3.2 Implementation Details

Two files were created to implement this project.

3.2.1 Data Type

The first one is “CDataType.scala”. It contains a class that depicts C structures.

```
1 class StructDef(val name: String,
2   val membersList: ListBuffer[StructMember])
```

A trait “StructMember” is extended by two classes to implement the two kinds of structure members, which are a variable or a pointer to a function (as explained in section 2). The pointer to a function holds two fields, one for an internal representation of a C function pointer (the class “FunctionPtr”), and one for the AST version of the method declaration that should be linked to the pointer.

```
1 class StructVar(val name: String,
2   val tpe: CType) extends StructMember {
3   // some code
4 }
5
6 class StructFunctionPtr(val ptr: FunctionPtr,
7   var mtDcl: MethodDecl) extends StructMember {
8   // some code
9 }
10
11 class FunctionPtr(val name: String,
12   val retType: CType, val args: List[CType]) {
13   // some code
```

14 }

In this file is also defined the CType abstract class. This allows us to represent a C version of the Tool types. We only have 4 object that extends this abstract class : CInt, CIntArray, CStruct and CString. There is no equivalent of the Bool type since boolean values are treated as ints in C. All classes of Tool are represented by the CStruct type, since structs are used to emulate class in our printed C code.

3.2.2 Code Generation

The second file is “COutputGeneration.scala”. Its skeleton and behavior are based on the “CodeGeneration.scala” completed to output binary code for the Java Virtual Machine.

The program uses the same architecture of methods like: “cGenMethod”, “cGenStats”, “cGenExpr”. Each method returns a StringBuilder.

The final output are two files, one .h and one .c, which contains all the concatenated StringBuilders.

Overview The main method of the object

COutputGeneration is `def run(ctx: Context)(prog: Program)`, which takes as a parameter the AST of the code.

The method starts by creating a StringBuilder for the preprocessor directives. It includes the headers `stdio.h`, `string.h` and `stdlib.h`. It contains `#define INT_MAX_LENGTH 12`, which is useful in order to have correct Int to String character concatenation (12 being the number of digit of -2^{32} , the longest number an int can hold.). It also contains the definition of default constructor macros for every class of the program. We will get back to this point later in this section.

```
1 #define nAnimal 0
2 #define nDog 1
```

Listing 1. In the case of the example used in section 2, this would be the produced macros

Finally, a StringBuilder contains the inclusion of the file header, which is directly created.

Afterward a method takes care of generating the structures corresponding to the classes of Tool:

`def genStructDef(ct: ClassDecl): StructDef`. The method takes as only parameter a `ClassDecl` and returns the corresponding `StructDef`, which was pre-

sented in section 3.2.1 Data Type. Then the method `toStringRepr` of `StructDef` is called to return the structures as Strings to be printed.

Thereafter the method manages the program methods.

```
1 def genMethods(ct: ClassDecl): StringBuilder =
2   (for(mt <- ct.methods)yield(cGenMethod(ct,
   mt))).foldLeft(new
   StringBuilder())((a,b) => a append b)
```

The method `cGenMethod` traverses the AST node `MethodDecl` to translate the statements and the expressions in C, using the methods `cGenStats` and `cGenExpr` like in the Code Generation lab. The method returns a `StringBuilder` representing a function, that corresponds to a class method of Tool.

Then the main method of is translated into C.

```
1 def genMainMethod(main: MainObject):
   StringBuilder = {
2   val mainMethod = new StringBuilder("int
   main(void){\n")
3   main.stats.foldLeft(mainMethod)((sB, stmt) =>
   sB append(cGenStat(stmt)(1, None)))
4   return mainMethod.append("\treturn 0;\n")
5   }
```

The `genMainMethod` method receives as a parameter the AST node `MainObject`. It is converted as the main function in C.

```
1 int main(void){
2   return 0;
3 }
```

`MainObject` contains a list of `StatTree`. Each `StatTree` is evaluated by the method `cGenStat` and the returned `StringBuilder(s)` is the body of the C main function. In the run method, there is an object `defaultConstructor`, that represents the construction of the default constructor. At the end of the run method, it will produce a C helper function to be printed with every program. The function is called `new`. It takes as only parameter an integer. It is one of the integers defined for each Tool class, that we have seen before, in Listing 1. Subsequently a method `addStructConstructor` takes care of adding a `case` element to the new function for a correct initialization and allocation of a Tool object in C.

```
1 void * new(int type){
2   void * object;
3   switch(type){
```

```
4   case nAnimal:
5     object = malloc(sizeof(struct Animal));
6     ((struct Animal *) object)->sleep =
       Animal_sleep;
7     break;
8   case nDog:
9     object = malloc(sizeof(struct Dog));
10    ((struct Dog *) object)->bark = Dog_bark;
11    ((struct Dog *) object)->sleep = Dog_sleep;
12  default:
13    return NULL;
14  }
15  return object;
16 }
```

Listing 2. In the case of the example used in the section 2.

Since this new function will be called each time a new was written in Tool, class methods need to have it visible. Because new references all the functions defined, it can't be put before them in our outputted C program, so its prototype is included in the `.h` file. And it is the sole purpose of the `.h`.

In the end, the functions

`void helper_reverse_plus(char str[], int len)`, `char* itoa(int num)` and `int * arrayAlloc(int size)` are constructed. Their usefulness will be explained in the subsections to come.

Finally all these snippets of C code are concatenated and written in a `.c` file.

Dynamic dispatch Dynamic dispatch is accomplished during the default construction of a new Instance of a class. In Tool, there is only on type of constructor, the default one. One idea we borrowed from our reference book[?] (chapter 11, section 4, paragraph “Plugging the Memory Leaks”) was to have a default method called `new` be outputted which would take care of creating a new Instance of a struct and returning it as a generic pointer `void *`. This function was already presented in the overview. `void *` is the type used in C to represent classes, getting class members or functions is then accomplished by casting this generic pointer to the corresponding struct.

According to the overriding of the methods or not, this method also takes care of setting the value of the function pointers of the struct to the corresponding function. To be sure that calling `sleep` on an instance

of `Animal` will get the correct pointer, for inherited structs, we make sure they hold the same fields in the same place than their parent struct. If we look at the struct definitions of example 2, we clearly see that the function pointer `sleep` is held in the same place in both structs. Therefore, if we consider this Tool code snippet:

```

1 program Sleep{
2   println(new Farm().getPet(0).sleep());
3   //prints 'zzz...'
4   println(new Farm().getPet(1).sleep());
5   //prints 'zzz... wouaf... zzz...'
6 }
7
8 class Farm{
9   def getPet(sel: Int): Animal = {
10     if(sel == 0){
11       return new Animal();
12     }else{
13       return new Dog();
14     }
15   }
16 }
```

Running it will result in the `sleep` message of `Animal` being printed, followed by the `sleep` message of `Dog`. This is an example of dynamic dispatch. Using casting, the same result can be obtained in C :

```

1 int main(void){
2   void * farm0 = new(nFarm)
3   void * animal0 = ((struct Farm
4     *)farm0)->getPet(0);
5   printf("%s\n", ((struct Animal *)
6     animal0)->sleep())
7   //prints 'zzz...'
8   void * farm1 = new(nFarm)
9   void * animal1 = ((struct Farm
10     *)farm1)->getPet(1);
11   printf("%s\n", ((struct Animal *)
12     animal1)->sleep())
13   //prints 'zzz... wouaf... zzz...'
14 }
```

Casting to `Animal` is made in both cases because `getPet` returns an instance of `Animal` so the compiler treats the value returned to be of type “`Animal`”. Even if a `Dog` is returned it will use the `sleep` function of the `dog` since the pointer to `sleep` in `Dog` struct is in the same place as the one of `Animal` struct.

Length of Arrays Another problem encountered was retrieving the length of an array. In Tool, getting it is made in the same way as in Java, by accessing the field

length of the array. However this cannot be translated easily in C if the array is not statically allocated. The way we treat arrays of ints in our extension is by defining them as pointers to a value of type `int`. Then creating a new array of `n` elements is made by using `calloc` with `n` and `sizeof(int)` as arguments. The value returned points now to a region of the memory that can hold `n` ints. We will call this pointer “`x`”. And since our array is just a pointer, trying the trick of `sizeof(x)/sizeof(int)` to get its length will either be equal to 1 or 2 (depending on the word size of the system) because `sizeof(x)` in this case returns the size of the pointer, not the size of the memory this pointer was allocated to.

To solve this problem, we defined a helper function called `arrayAlloc[?]` printed with every program.

```

1 int * arrayAlloc(int size){
2   int * smrtArray = calloc(size + 1, sizeof(int));
3   smrtArray[0] = size;
4   return (smrtArray + 1);
5 }
```

What this function does is pretty simple but ingenious. The array is allocated for one more element than the required size, and the length of this array is stored at the first place of the allocated array. We then return the value of the pointer incremented by one. Thus, when we need to get the length of the array “`x`”, we can get it from the value stored 1 position before the pointer “`x`”

```

1 int * x = arrayAlloc(42);
2 int xLength = *(x - 1) //xLength is equal to 42.
```

Concatenation The AST node “Plus” needs some attention. The Tool language allows to apply the “+” operator to `Int` and `String` operands.

The situation, where the two operands are `Ints`, is straightforward to emulate in C, since it’s also just an addition.

The event where the two operands are `Strings` is a little bit more tricky. It is necessary to allocate memory before concatenating the two `Strings` of characters.

```

1 strcpy(malloc(strlen(lhsString) +
2   strlen(rhsString) + 1), lhsString, rhsString)
```

Finally when the two operands are different, we were obliged to create two C functions:

```
void helper_reverse_plus(char str[], int len)
```

and `char* itoa(int num)`. These functions transform an integer into a string of characters and return it. These two functions are written in every .c file generated by the compiler.

Method chaining This problem was one that we had not seen coming.

In Tool, there is the possibility to write multiple calls to a method on one line. It works because each method returns an object. We can't directly translate it in C with the model we chose to emulate object oriented programming. Our "methods" need to have as a first argument a pointer to the struct calling, so to make a method call on a struct we would need to write : `void * a = new(nA); ((struct A *)a)->foo(a)`. Our variable a needs to be referenced two times, first to get the good function pointer and then as the first argument if `foo` needs to access fields defined in A. Therefore chained method calls need intermediate variables to work correctly. To tackle this problem, we created a simple object: "tmpVarGen".

```

1 object tmpVarGen{
2   private var counter = 0
3   private var lastSuffix = ""
4   def getFreshVar(suffix: Option[String]): String
5     = {
6       counter += 1
7       val sffx = suffix match{
8         case Some(s) => s
9         case None => ""
10      }
11      lastSuffix = sffx
12      return "tmp"+sffx+counter
13    }
14   def getLastVar: String =
15     "tmp"+lastSuffix+counter
16 }

```

It is used to have a unique variable name. It has two methods. "getLastVar" returns the last created variable name. "getFreshVar" returns a new variable name with a suffix given as an argument to make the temporary variable hold the information of which expression it represents.

Now at each evaluation of an AST node, which is an expression, it is important first to evaluate the expressions and second to assign it to intermediate variables using the object "tmpVarGen". Next the expression is written using the intermediate variables.

```

1 case Equals(lhs: ExprTree, rhs: ExprTree) =>
2   val lhsString = cGenExpr(lhs)

```

```

3   val lhsLastVar = tmpVarGen.getLastVar
4   val rhsString = cGenExpr(rhs)
5   val rhsLastVar = tmpVarGen.getLastVar
6   val andExprResultVar = genTabulation(indentLvl)+
7     CInt.toString()+ " "+tmpVarGen.getFreshVar+" =
8     "+
9     lhsLastVar+" == "+rhsLastVar+";\n"
10  return
11    lhsString.append(rhsString).append(andExprResultVar)

```

Here is an example with the evaluation of an expression AST node: Equals. The left-hand side and right-hand side expressions are first evaluated using the "cGenExpr()" method and are stored in an intermediate variable using the "tmpVarGen" object. Afterward the equal expression as known in C: "lhs == rhs" is written using the intermediate variable. At the end we append all elements together. Thus one Tool line to write an equality corresponds to three lines in C programming language.

The same procedure has been used for all other expressions and also for statements because they contain expression(s).

4. Possible Extensions

The aim of the project, to output C, has been reached. The C programs outputted can be compiled on gcc and running them will output the same result than if the original Tool was run on the JVM.

However enhancements could be done to produce better programs. Firstly, the problem we encountered about method chaining forced us to print intermediate variables for each expression. Even though it doesn't change the functionality of the program, it makes the output more difficult to be read by humans. A new model for object-oriented in C could be conceived to make method chaining more straight-forward. The struct methods might not have to hold a pointer to the struct itself as first argument, and instead accessing a registry of all classes/structs produced to get the reference to their instance and access the fields.

Another enhancement would be to include a garbage collector[?] to our program. In our implementation, instances of classes and arrays are allocated but never freed, thus creating memory leaks. A simple garbage collector mechanic could be printed with the program to deallocate the pointers when reaching the end of the scope.