

Handout 7 (Compilation)

The purpose of a compiler is to transform a program a human can read and write into code the machine can run as fast as possible. The fastest code would be machine code the CPU can run directly, but it is often good enough for improving the speed of a program to target a virtual machine. This produces not the fastest possible code, but code that is often pretty fast. This way of producing code has the advantage that the virtual machine takes care of things a compiler would normally need to take care of (like explicit memory management).

As a first example in this module we will implement a compiler for the very simple While-language. It will generate code for the Java Virtual Machine (JVM). Unfortunately the Java ecosystem does not come with an assembler which would be handy for our compiler-endeavour (unlike Microsoft's Common Language Infrastructure for the .Net platform which has an assembler out-of-the-box). As a substitute we use in this module the 3rd-party programs Jasmin and Krakatau

- <http://jasmin.sourceforge.net>
- <https://github.com/Storyyeller/Krakatau>

The first is a Java program and the second a program written in Python. Each of them allow us to generate *assembly* files that are still readable by humans, as opposed to class-files which are pretty much just (horrible) zeros and ones. Jasmin (respectively Krakatau) will then take an assembly file as input and generate the corresponding class file for us.

Good about the JVM is that it is a stack-based virtual machine, a fact which will make it easy to generate code for arithmetic expressions. For example when compiling the expression $1 + 2$ we need to generate the following three instructions

```
ldc 1
ldc 2
iadd
```

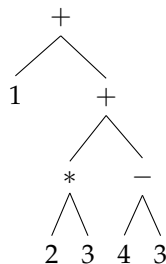
The first instruction loads the constant 1 onto the stack, the next one loads 2, the third instruction adds both numbers together replacing the top two elements of the stack with the result 3. For simplicity, we will throughout consider only integer numbers. Therefore we can use the JVM instructions `iadd`, `isub`, `imul`, `idiv` and so on. The `i` stands for integer instructions in the JVM (alternatives are `d` for doubles, `l` for longs and `f` for floats).

Recall our grammar for arithmetic expressions (E is the starting symbol):

$$E ::= T + E \mid T - E \mid T$$
$$T ::= F * T \mid F \setminus T \mid F$$

$$F ::= (E) \mid Id \mid Num$$

where *Id* stands for variables and *Num* for numbers. For the moment let us omit variables from arithmetic expressions. Our parser will take this grammar and given an input produce abstract syntax trees. For example we will obtain for the expression $1 + ((2 * 3) + (4 - 3))$ the following tree.



To generate JVM code for this expression, we need to traverse this tree in post-order fashion and emit code for each node—this traversal in post-order fashion will produce code for a stack-machine (what the JVM is). Doing so for the tree above generates the instructions

```

ldc 1
ldc 2
ldc 3
imul
ldc 4
ldc 3
isub
iadd
iadd
  
```

If we “run” these instructions, the result 8 will be on top of the stack (I leave this to you to verify; the meaning of each instruction should be clear). The result being on the top of the stack will be an important convention we always observe in our compiler. Note, that a different bracketing of the expression, for example $(1 + (2 * 3)) + (4 - 3)$, produces a different abstract syntax tree and thus also a different list of instructions. Generating code in this post-order-traversal fashion is rather easy to implement: it can be done with the following recursive *compile*-function, which takes the abstract syntax tree as argument:

```

compile(n)       $\stackrel{\text{def}}{=}$    ldc n
compile(a1 + a2)  $\stackrel{\text{def}}{=}$    compile(a1) @ compile(a2) @ iadd
compile(a1 - a2)  $\stackrel{\text{def}}{=}$    compile(a1) @ compile(a2) @ isub
compile(a1 * a2)  $\stackrel{\text{def}}{=}$    compile(a1) @ compile(a2) @ imul
compile(a1 \ a2)  $\stackrel{\text{def}}{=}$    compile(a1) @ compile(a2) @ idiv
  
```

However, our arithmetic expressions can also contain variables. We will represent them as *local variables* in the JVM. Essentially, local variables are an array or pointers to memory cells, containing in our case only integers. Looking up a variable can be done with the instruction

```
iload index
```

which places the content of the local variable *index* onto the stack. Storing the top of the stack into a local variable can be done by the instruction

```
istore index
```

Note that this also pops off the top of the stack. One problem we have to overcome, however, is that local variables are addressed, not by identifiers, but by numbers (starting from 0). Therefore our compiler needs to maintain a kind of environment where variables are associated to numbers. This association needs to be unique: if we muddle up the numbers, then we essentially confuse variables and the consequence will usually be an erroneous result. Our extended *compile*-function for arithmetic expressions will therefore take two arguments: the abstract syntax tree and an environment, *E*, that maps identifiers to index-numbers.

$$\begin{aligned} \text{compile}(n, E) & \stackrel{\text{def}}{=} \text{ldc } n \\ \text{compile}(a_1 + a_2, E) & \stackrel{\text{def}}{=} \text{compile}(a_1, E) @ \text{compile}(a_2, E) @ \text{iadd} \\ \text{compile}(a_1 - a_2, E) & \stackrel{\text{def}}{=} \text{compile}(a_1, E) @ \text{compile}(a_2, E) @ \text{isub} \\ \text{compile}(a_1 * a_2, E) & \stackrel{\text{def}}{=} \text{compile}(a_1, E) @ \text{compile}(a_2, E) @ \text{imul} \\ \text{compile}(a_1 \setminus a_2, E) & \stackrel{\text{def}}{=} \text{compile}(a_1, E) @ \text{compile}(a_2, E) @ \text{idiv} \\ \text{compile}(x, E) & \stackrel{\text{def}}{=} \text{iload } E(x) \end{aligned}$$

In the last line we generate the code for variables where $E(x)$ stands for looking up the environment to which index the variable x maps to. This is similar to an interpreter, which also needs an environment: the difference is that the interpreter maintains a mapping from variables to current values (what is the currently the value of a variable), while compilers need a mapping from variables to memory locations (where can I find the current value for the variable in memory).

There is a similar *compile*-function for boolean expressions, but it includes a “trick” to do with *if*- and *while*-statements. To explain the issue let us first describe the compilation of statements of the *While*-language. The clause for *skip* is trivial, since we do not have to generate any instruction

$$\text{compile}(\text{skip}, E) \stackrel{\text{def}}{=} ([], E)$$

whereby $[]$ is the empty list of instructions. Note that the *compile*-function for statements returns a pair, a list of instructions (in this case the empty list) and an environment for variables. The reason for the environment is that assignments in the *While*-language might change the environment—clearly if a variable is

used for the first time, we need to allocate a new index and if it has been used before, then we need to be able to retrieve the associated index. This is reflected in the clause for compiling assignments, say $x := a$:

$$\text{compile}(x := a, E) \stackrel{\text{def}}{=} (\text{compile}(a, E) @ \text{istore } \text{index}, E')$$

We first generate code for the right-hand side of the assignment and then add an `istore`-instruction at the end. By convention the result of the arithmetic expression a will be on top of the stack. After the `istore` instruction, the result will be stored in the index corresponding to the variable x . If the variable x has been used before in the program, we just need to look up what the index is and return the environment unchanged (that is in this case $E' = E$). However, if this is the first encounter of the variable x in the program, then we have to augment the environment and assign x with the largest index in E plus one (that is $E' = E(x \mapsto \text{largest_index} + 1)$). To sum up, for the assignment $x := x + 1$ we generate the following code

```

iadd
istore nx

```

where n_x is the index (or pointer to the memory) for the variable x . The code for looking-up the index for the variable is as follow:

$$\text{index} = E.\text{getOrElse}(x, |E|)$$

In case the environment E contains an index for x , we return it. Otherwise we “create” a new index by returning the size $|E|$ of the environment (that will be an index that is guaranteed to be not used yet). In all this we take advantage of the JVM which provides us with a potentially limitless supply of places where we can store values of variables.

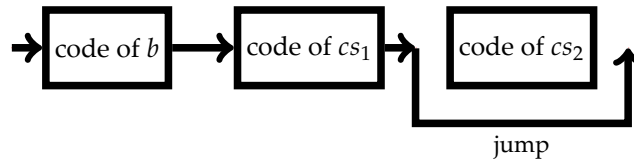
A bit more complicated is the generation of code for `if`-statements, say

```

if b then cs1 else cs2

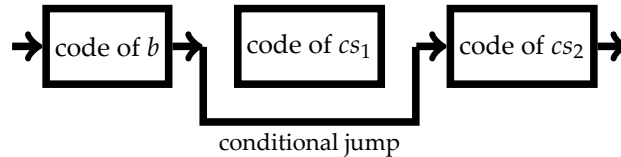
```

where b is a boolean expression and where both $cs_{1/2}$ are the statements for each of the `if`-branches. Lets assume we already generated code for b and $cs_{1/2}$. Then in the true-case the control-flow of the program needs to behave as



where we start with running the code for b ; since we are in the true case we continue with running the code for cs_1 . After this however, we must not run

the code for cs_2 , but always jump after the last instruction of cs_2 (the code for the `else`-branch). Note that this jump is unconditional, meaning we always have to jump to the end of cs_2 . The corresponding instruction of the JVM is `goto`. In case b turns out to be false we need the control-flow



where we now need a conditional jump (if the if-condition is false) from the end of the code for the boolean to the beginning of the instructions cs_2 . Once we are finished with running cs_2 we can continue with whatever code comes after the if-statement.

The `goto` and the conditional jumps need addresses to where the jump should go. Since we are generating assembly code for the JVM, we do not actually have to give (numeric) addresses, but can just attach (symbolic) labels to our code. These labels specify a target for a jump. Therefore the labels need to be unique, as otherwise it would be ambiguous where a jump should go to. A label, say L , is attached to code like

```
L :
  instr1
  instr2
  ⋮
```

where a label is indicated by a colon. The task of the assembler (in our case Jasmin or Krakatau) is to resolve the labels to actual addresses, for example jump 10 instructions forward, or 20 instructions backwards.

Recall the “trick” with compiling boolean expressions: the `compile`-function for boolean expressions takes three arguments: an abstract syntax tree, an environment for variable indices and also the label, lab , to where an conditional jump needs to go. The clause for the expression $a_1 = a_2$, for example, is as follows:

$$\text{compile}(a_1 = a_2, E, lab) \stackrel{\text{def}}{=} \text{compile}(a_1, E) @ \text{compile}(a_2, E) @ \text{if_icmpne } lab$$

where we are first generating code for the subexpressions a_1 and a_2 . This will mean after running the corresponding code there will be two integers on top of the stack. If they are equal, we do not have to do anything (except for popping them off from the stack) and just continue with the next instructions (see control-flow of ifs above). However if they are *not* equal, then we need to (conditionally) jump to the label lab . This can be done with the instruction

```
if_icmpne lab
```

Other jump instructions for boolean operators are

\neq	\Rightarrow	<code>if_icmpeq</code>
$<$	\Rightarrow	<code>if_icmpge</code>
\leq	\Rightarrow	<code>if_icmpgt</code>

and so on. I leave it to you to extend the *compile*-function for the other boolean expressions. Note that we need to jump whenever the boolean is *not* true, which means we have to “negate” the jump condition—equals becomes not-equal, less becomes greater-or-equal. If you do not like this design (it can be the source of some nasty, hard-to-detect errors), you can also change the layout of the code and first give the code for the else-branch and then for the if-branch. However in the case of while-loops this “upside-down-inside-out” way of generating code still seems the most convenient.

We are now ready to give the compile function for if-statements—remember this function returns for statements a pair consisting of the code and an environment:

$$\text{compile}(\text{if } b \text{ then } cs_1 \text{ else } cs_2, E) \stackrel{\text{def}}{=} \begin{array}{l} L_{\text{ifelse}} \text{ (fresh label)} \\ L_{\text{ifend}} \text{ (fresh label)} \\ (is_1, E') = \text{compile}(cs_1, E) \\ (is_2, E'') = \text{compile}(cs_2, E') \\ (\text{compile}(b, E, L_{\text{ifelse}}) \\ @ is_1 \\ @ \text{goto } L_{\text{ifend}} \\ @ L_{\text{ifelse}} : \\ @ is_2 \\ @ L_{\text{ifend}} :, E'') \end{array}$$

In the first two lines we generate two fresh labels for the jump addresses (just before the else-branch and just after). In the next two lines we generate the instructions for the two branches, is_1 and is_2 . The final code will be first the code for b (including the label just-before-the-else-branch), then the `goto` for after the else-branch, the label L_{ifelse} , followed by the instructions for the else-branch, followed by the after-the-else-branch label. Consider for example the if-statement:

```
if 1 = 1 then x := 2 else y := 3
```

The generated code is as follows:

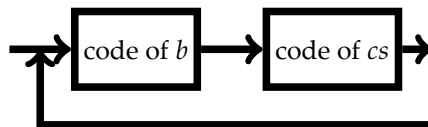
```

1   ldc 1
2   ldc 1
3   if_icmpne L_ifelse
4   ldc 2
5   istore 0
6   goto L_ifend
7 L_ifelse:
8   ldc 3
9   istore 1
10  L_ifend:

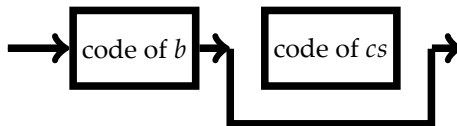
```

The first three lines correspond to the the boolean expression $1 = 1$. The jump for when this boolean expression is false is in Line 3. Lines 4-6 corresponds to the if-branch; the else-branch is in Lines 8 and 9. Note carefully how the environment E is threaded through the recursive calls of *compile*. The function receives an environment E , but it might extend it when compiling the if-branch, yielding E' . This happens for example in the if-statement above whenever the variable x has not been used before. Similarly with the environment E'' for the second call to *compile*. E'' is also the environment that needs to be returned as part of the answer.

The compilation of the while-loops, say `while b do cs` , is very similar. In case the condition is true and we need to do another iteration, and the control-flow needs to be as follows



Whereas if the condition is *not* true, we need to jump out of the loop, which gives the following control flow.



Again we can use the *compile*-function for boolean expressions to insert the appropriate jump to the end of the loop (label L_{wend} below).

$$\text{compile}(\text{while } b \text{ do } cs, E) \stackrel{\text{def}}{=} \\
L_{wbegin} \text{ (fresh label)} \\
L_{wend} \text{ (fresh label)} \\
(is, E') = \text{compile}(cs_1, E) \\
(L_{wbegin} : \\
@ \text{compile}(b, E, L_{wend}) \\
@ is \\
@ \text{goto } L_{wbegin} \\
@ L_{wend} :, E')$$

I let you go through how this clause works. As an example you can consider the while-loop

```
while x <= 10 do x := x + 1
```

yielding the following code

```
1 L_wbegin:
2   iload 0
3   ldc 10
4   if_icmpgt L_wend
5   iload 0
6   ldc 1
7   iadd
8   istore 0
9   goto L_wbegin
10 L_wend:
```

I leave it to you to read the code and follow its controlflow.

Next we need to consider the statement `write x`, which can be used to print out the content of a variable. For this we need to use a Java library function. In order to avoid having to generate a lot of code for each `write`-command, we use a separate helper-method and just call this method with an argument (which needs to be placed onto the stack). The code of the helper-method is as follows.

```
1 .method public static write(I)V
2   .limit locals 1
3   .limit stack 2
4   getstatic java/lang/System/out Ljava/io/PrintStream;
5   iload 0
6   invokevirtual java/io/PrintStream/println(I)V
7   return
8 .end method
```

The first line marks the beginning of the method, called `write`. It takes a single integer argument indicated by the (I) and returns no result, indicated by the V. Since the method has only one argument, we only need a single local variable (Line 2) and a stack with two cells will be sufficient (Line 3). Line 4 in-


```

1 .class public XXX.XXX
2 .super java/lang/Object
3
4 .method public <init>()V
5     aload_0
6     invokevirtual java/lang/Object/<init>()V
7     return
8 .end method
9
10 .method public static main([Ljava/lang/String;)V
11     .limit locals 200
12     .limit stack 200
13
14     ...here comes the compiled code...
15
16     return
17 .end method

```

Figure 1: Boilerplate code needed for running generated code.

structs the JVM to get the value of the field out of the class `java/lang/System`. It expects the value to be of type `java/io/PrintStream`. A reference to this value will be placed on the stack. Line 5 copies the integer we want to print out onto the stack. In the next line we call the method `println` (from the class `java/io/PrintStream`). We want to print out an integer and do not expect anything back (that is why the type annotation is `(I)V`). The return-instruction in the next line changes the control-flow back to the place from where `write` was called. This method needs to be part of a header that is included in any code we generate. The helper-method `write` can be invoked with the two instructions

```

iload E(x)
invokestatic XXX/XXX/write(I)V

```

where we first place the variable to be printed on top of the stack and then call `write`. The `XXX` need to be replaced by an appropriate class name (this will be explained shortly).

By generating code for a `While`-program, we end up with a list of (JVM assembly) instructions. Unfortunately, there is a bit more boilerplate code needed before these instructions can be run. The complete code is shown in Figure 1. This boilerplate code is very specific to the JVM. If we target any other virtual machine or a machine language, then we would need to change this code. Lines 4 to 8 in Figure 1 contain a method for object creation in the JVM; this method is called *before* the `main`-method in Lines 10 to 17. Interesting are the Lines 11 and 12 where we hardwire that the stack of our programs will never be larger than 200 and that the maximum number of variables is also 200. This seem to be

conservative default values that allow is to run some simple While-programs. In a real compiler, we would of course need to work harder and find out appropriate values for the stack and local variables.

To sum up, in Figure 2 is the complete code generated for the slightly non-sensical program

```
x := 1 + 2;
write x
```

I let you read the code and make sure the code behaves as expected. Having this code at our disposal, we need the assembler to translate the generated code into JVM bytecode (a class file). This bytecode is then understood by the JVM and can be run by just invoking the java-program.

Arrays

Maybe a useful addition to the While-language would be arrays. This would let us generate more interesting While-programs by translating BF*** programs into equivalent While-code. So in this section lets have a look at how we can support the following three constructions

```
new arr [15000]
x := 3 + arr [3 + y]
arr [42 * n] := ...
```

The first construct is for creating new arrays, in this instance the name of the array is arr and it can hold 15000 integers. The second is for referencing an array cell inside an arithmetic expression—we need to be able to look up the contents of an array at an index determined by an arithmetic expression. Similarly in the line below, we need to be able to update the content of an array at an calculated index.

For creating a new array we can generate the following three JVM instructions:

```
ldc number
newarray int
astore loc_var
```

First we need to put the dimension of the array onto the stack. The next instruction creates the array. With the last we can store the array as a local variable (like the “simple” variables from the previous section). The use of a local variable for each array allows us to have multiple arrays in a While-program. For looking up an element in an array we can use the following JVM code

```
aload loc_var
index_aexp
iaload
```

The first instruction loads the “pointer” to the array onto the stack. Then we have some instructions corresponding to the index where we want to look up

```

.class public test.test
.super java/lang/Object

.method public <init>()V
    aload_0
    invokevirtual java/lang/Object/<init>()V
    return
.end method

.method public static write(I)V
    .limit locals 1
    .limit stack 2
    getstatic java/lang/System/out Ljava/io/PrintStream;
    iload 0
    invokevirtual java/io/PrintStream/println(I)V
    return
.end method

.method public static main([Ljava/lang/String;)V
    .limit locals 200
    .limit stack 200
    ldc 1
    ldc 2
    iadd
    istore 0
    iload 0
    invokestatic test/test/write(I)V
    return
.end method

```

Figure 2: Generated code for a test program. This code can be processed by an Java assembler producing a class-file, which can be run by the java-program.

the array. The idea is that these instructions will leave a concrete number on the stack, which will be the index into the array we need. Finally we need to tell the JVM to load the corresponding element onto the stack. Updating an array at an index with a value is as follows.

```

aload loc_var
index_aexp
value_aexp
istore

```

Again the first instruction loads the “pointer” to the array onto the stack. Then we have some instructions corresponding to the index where we want to update the array. After that come the instructions for with what value we want to update the array. The last line contains the instruction for updating the array.

Next we need to modify our grammar rules for our While-language: it seems best to extend the rule for factors in arithmetic expressions with a rule for looking up an array.

$$\begin{aligned}
 E &::= T + E \mid T - E \mid T \\
 T &::= F * T \mid F \setminus T \mid F \\
 F &::= (E) \mid \underbrace{Id[E]}_{new} \mid Id \mid Num
 \end{aligned}$$

There is no problem with left-recursion as the E is “protected” by an identifier and the brackets. There are two new rules for statements, one for creating an array and one for array assignment:

$$\begin{aligned}
 Stmt &::= \dots \\
 &\mid \text{new } Id[Num] \\
 &\mid Id[E] := E
 \end{aligned}$$

With this in place we can turn back to the idea of creating While programs by translating BF programs. This is a relatively easy task because BF only has eight instructions (we will actually only have seven because we can omit the read-in instruction from BF). But also translating BF-loops is going to be easy since they straightforwardly can be represented by While-loops. The Scala code for the translation is as follows:

```

1  def instr(c: Char) : String = c match {
2    case '>' => "ptr := ptr + 1;"
3    case '<' => "ptr := ptr - 1;"
4    case '+' => "field[ptr] := field[ptr] + 1;"
5    case '-' => "field[ptr] := field[ptr] - 1;"
6    case '.' => "x := field[ptr]; write x;"
7    case '[' => "while (field[ptr] != 0) do {"
8    case ']' => "skip};"
9    case _ => ""
10 }

```

The idea behind the translation is that BF-programs operate on an array, called `field`. The BP-memory pointer into this array is represented as the variable `ptr`. The BF-instructions `>` and `<` increase, respectively decrease, `ptr`. The instructions `+` and `-` update a cell in `field`. In Line 6 we need to first assign a `field`-cell to an auxiliary variable since we have not changed our write functions in order to cope with writing out any array-content directly. Lines 7 and 8 are for translating BF-loops. Line 8 is interesting in the sense that we need to generate a `skip` instruction just before finishing with the closing `"}`". The reason is that we are rather pedantic about semicolons in our While-grammar: the last command cannot have a semicolon—adding a `skip` works around this snag. Putting all this together is we can generate While-programs with more than 400 instructions and then run the compiled JVM code for such programs.