

Handout 10 (Static Analysis)

If we want to improve the safety and security of our programs, we need a more principled approach to programming. Testing is good, but as Edsger Dijkstra famously wrote:

“Program testing can be a very effective way to show the presence of bugs, but it is hopelessly inadequate for showing their absence.”

While such a more principled approach has been the subject of intense study for a long, long time, only in the past few years some impressive results have been achieved. One is the complete formalisation and (mathematical) verification of a microkernel operating system called seL4.

<http://sel4.systems>

In 2011 this work was included in the MIT Technology Review in the annual list of the world’s ten most important emerging technologies.¹ While this work is impressive, its technical details are too enormous for an explanation here. Therefore let us look at something much simpler, namely finding out properties about programs using *static analysis*.

Static analysis is a technique that checks properties of a program without actually running the program. This should raise alarm bells with you—because almost all interesting properties about programs are equivalent to the halting problem, which we know is undecidable. For example estimating the memory consumption of programs is in general undecidable, just like the halting problem. Static analysis circumvents this undecidability-problem by essentially allowing answers *yes* and *no*, but also *don’t know*. With this “trick” even the halting problem becomes decidable...for example we could always say *don’t know*. Of course this would be silly. The point is that we should be striving for a method that answers as often as possible either *yes* or *no*—just in cases when it is too difficult we fall back on the *don’t-know*-answer. This might sound all like abstract nonsense. Therefore let us look at a concrete example.

A Simple, Idealised Programming Language

Our starting point is a small, idealised programming language. It is idealised because we cut several corners in comparison with real programming languages. The language we will study contains, amongst other things, variables holding integers. Using static analysis, we want to find out what the sign of these integers (positive or negative) will be when the program runs. This sign-analysis seems like a very simple problem. But even such simple problems, if approached naively, are in general undecidable, just like Turing’s halting problem. I let you think why?

© Christian Urban, King’s College London, 2014

¹<http://www2.technologyreview.com/tr10/?year=2011>

Is sign-analysis of variables an interesting problem? Well, yes—if a compiler can find out that for example a variable will never be negative and this variable is used as an index for an array, then the compiler does not need to generate code for an underflow-check. Remember some languages are immune to buffer-overflow attacks, but they need to add underflow and overflow checks everywhere. According to John Regehr, an expert in the field of compilers, overflow checks can cause 5-10% slowdown, in some languages even 100% for tight loops.² If the compiler can omit the underflow check, for example, then this can potentially drastically speed up the generated code.

What do programs in our simple programming language look like? The following grammar gives a first specification:

$\begin{aligned} \textit{Stmt} ::= & \textit{label} : \\ & \textit{var} := \textit{Exp} \\ & \textit{jmp? Exp label} \\ & \textit{goto label} \\ \textit{Prog} ::= & \textit{Stmt} \dots \textit{Stmt} \end{aligned}$	$\begin{aligned} \textit{Exp} ::= & \textit{Exp} + \textit{Exp} \\ & \textit{Exp} * \textit{Exp} \\ & \textit{Exp} = \textit{Exp} \\ & \textit{num} \\ & \textit{var} \end{aligned}$
---	--

I assume you are familiar with such grammars.³ There are three main syntactic categories: *statements* and *expressions* as well as *programs*, which are sequences of statements. Statements are either labels, variable assignments, conditional jumps (*jmp?*) and unconditional jumps (*goto*). Labels are just strings, which can be used as the target of a jump. We assume that in every program the labels are unique—if there is a clash, then we do not know where to jump to. The conditional jumps and variable assignments involve (arithmetic) expressions. Expressions are either numbers, variables or compound expressions built up from +, * and = (for simplicity reasons we do not consider any other operations). We assume we have negative and positive numbers, ...-2, -1, 0, 1, 2... An example program that calculates the factorial of 5 is in our programming language as follows:

```

a := 1
n := 5
top:
  jmp? n = 0 done
  a := a * n
  n := n + -1
  goto top
done:
```

As can be seen each line of the program contains a statement. In the first two lines we assign values to the variables *a* and *n*. In line 4 we test whether *n* is zero, in which case we jump to the end of the program marked with the label *done*.

²<http://blog.regehr.org/archives/1154>

³http://en.wikipedia.org/wiki/Backus-Naur_Form

```

n := 6
m1 := 0
m2 := 1
loop:
  jmp? n = 0 done
  tmp := m2
  m2 := m1 + m2
  m1 := tmp
  n := n + -1
  goto top
done:

```

Figure 1: A mystery program in our idealised programming language. Try to find out what it calculates!

If n is not zero, we multiply the content of a by n , decrease n by one and jump back to the beginning of the loop, marked with the label `top`. Another program in our language is shown in Figure ?? . I let you think what it calculates.

Even if our language is rather small, it is still Turing complete—meaning quite powerful. However, discussing this fact in more detail would lead us too far astray. Clearly, our programming is rather low-level and not very comfortable for writing programs. It is inspired by real machine code, which is the code that is executed by a CPU. So a more interesting question is what is missing in comparison with real machine code? Well, not much...in principle. Real machine code, of course, contains many more arithmetic instructions (not just addition and multiplication) and many more conditional jumps. We could add these to our language if we wanted, but complexity is really beside the point here. Furthermore, real machine code has many instructions for manipulating memory. We do not have this at all. This is actually a more serious simplification because we assume numbers to be arbitrary small or large, which is not the case with real machine code. In real machine code, basic number formats have a range and might over-flow or under-flow this range. Also the number of variables in our programs is potentially unlimited, while memory in an actual computer, of course, is always limited. To sum up, our language might look ridiculously simple, but it is not too far removed from practically relevant issues.

An Interpreter

Designing a language is like playing god: you can say what names for variables you allow; what programs should look like; most importantly you can decide what each part of the program should mean and do. While our language is quite simple and the meaning of statements, for example, is rather straightforward, there are still places where we need to make real choices. For example

consider the conditional jumps, say the one in the factorial program:

```
    jmp? n = 0 done
```

How should they work? We could introduce Booleans (`true` and `false`) and then jump only when the condition is `true`. However, since we have numbers in our language anyway, why not just encoding `true` as one, and `false` as zero? In this way we can dispense with the additional concept of Booleans.

I hope the above discussion makes it already clear we need to be a bit more careful with our programs. Below we shall describe an interpreter for our programming language, which specifies exactly how programs are supposed to be run...at least we will specify this for all *good* programs. By good programs I mean where all variables are initialised, for example. Our interpreter will just crash if it cannot find out the value for a variable when it is not initialised. Also, we will assume that labels in good programs are unique, otherwise our programs will calculate “garbage”.

First we will pre-process our programs. This will simplify the definition of our interpreter later on. By pre-processing our programs we will transform programs into *snippets*. A snippet is a label and all the code that comes after the label. This essentially means a snippet is a *map* from labels to code.⁴

Given that programs are sequences (or lists) of statements, we can easily calculate the snippets by just traversing this sequence and recursively generating the map. Suppose a program is of the general form

$$stmt_1 \, stmt_2 \, \dots \, stmt_n$$

The idea is to go through this sequence of statements one by one and check whether they are a label. If yes, we add the label and the remaining statements to our map. If no, we just continue with the next statement. To come up with a recursive definition for generating snippets, let us write \square for the program that does not contain any statement. Consider the following definition:

$$\begin{aligned} snippets(\square) &\stackrel{\text{def}}{=} \emptyset \\ snippets(stmt \, rest) &\stackrel{\text{def}}{=} \begin{cases} snippets(rest)[label := rest] & \text{if } stmt = label: \\ snippets(rest) & \text{otherwise} \end{cases} \end{aligned}$$

In the first clause we just return the empty map for the program that does not contain any statement. In the second clause, we have to distinguish the case where the first statement is a label or not. As said before, if not, then we just “throw away” the label and recursively calculate the snippets for the rest of the program (the otherwise clause). If yes, then we do the same, but also update the map so that *label* now points to the rest of the statements (the if clause). This looks all relatively straightforward, but there is one small problem we need to overcome: our two programs shown so far have no label as *entry point*—that

⁴Be sure you know what maps are. In a programming context they are often represented as association list where some data is associated with a key.

is where the execution is supposed to start. We usually assume that the first statement will be run first. To make this the default, it is convenient if we add to all our programs a default label, say "" (the empty string). With this we can define our pre-processing of programs as follows

$$\text{preproc}(\text{prog}) \stackrel{\text{def}}{=} \text{snippets}(\text{"": prog})$$

Let us see how this pans out in practice. If we pre-process the factorial program shown earlier, we obtain the following map:

<code>""</code> \mapsto	<pre style="margin: 0;">a := 1 n := 5 top: jmp? n = 0 done a := a * n n := n + -1 goto top done:</pre>	<code>top</code> \mapsto	<pre style="margin: 0;">jmp? n = 0 done a := a * n n := n + -1 goto top done:</pre>	<code>done</code> \mapsto	<code>[]</code>
---------------------------	--	----------------------------	---	-----------------------------	-----------------

I highlighted the *keys* in this map. Since there are three labels in the factorial program (remember we added ""), there are three keys. When running the factorial program and encountering a jump, then we only have to consult this snippets-map in order to find out what the next statements should be.

We should now be in the position to define how a program should be run. In the context of interpreters, this “running” of programs is often called *evaluation*. Let us start with the definition of how expressions are evaluated. A first attempt might be the following recursive function:

$$\begin{aligned}
 \text{eval_exp}(n) &\stackrel{\text{def}}{=} n && \text{if } n \text{ is a number like } \dots -2, -1, 0, 1, 2\dots \\
 \text{eval_exp}(e_1 + e_2) &\stackrel{\text{def}}{=} \text{eval_exp}(e_1) + \text{eval_exp}(e_2) \\
 \text{eval_exp}(e_1 * e_2) &\stackrel{\text{def}}{=} \text{eval_exp}(e_1) * \text{eval_exp}(e_2) \\
 \text{eval_exp}(e_1 = e_2) &\stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } \text{eval_exp}(e_1) = \text{eval_exp}(e_2) \\ 0 & \text{otherwise} \end{cases}
 \end{aligned}$$

While this should look all rather intuitive, still be very careful. There is a subtlety which can be easily overlooked: The function *eval_exp* takes an expression of our programming language as input and returns a number as output. Therefore whenever we encounter a number in our program, we just return this number—this is defined in the first clause above. Whenever we encounter an addition, well then we first evaluate the left-hand side e_1 of the addition (this will give a number), then evaluate the right-hand side e_2 (this gives another number), and finally add both numbers together. Here is the subtlety: on the left-hand side of the $\stackrel{\text{def}}{=}$ we have a + (in the teletype font) which is the symbol for addition in our programming language. On the right-hand side we have + which stands for the arithmetic operation from “mathematics” of adding two

numbers. These are rather different concepts—one is a symbol (which we made up), and the other a mathematical operation. When we will have a look at an actual implementation of our interpreter, the mathematical operation will be the function for addition from the programming language in which we implement our interpreter. While the + is just a symbol that is used in our programming language. Clearly we have to use a symbol that is a good mnemonic for addition, otherwise we will confuse the programmers working with our language. Therefore we use +. A similar choice is made for times in the third clause and equality in the fourth clause.

Remember I wrote at the beginning of this section about being god when designing a programming language. You can see this here: we need to give meaning to symbols. At the moment however, we are a poor fallible god. Look again at the grammar of our programming language and our definition. Clearly, an expression can contain variables. So far we have ignored them. What should our interpreter do with variables? They might change during the evaluation of a program. For example the variable `n` in the factorial program counts down from 5 down to 0. How can we improve our definition above to give also an answer whenever our interpreter encounters a variable in an expression? The solution is to add an *environment*, written *env*, as an additional input argument to our *eval_exp* function.

$$\begin{aligned}
 \text{eval_exp}(n, \text{env}) &\stackrel{\text{def}}{=} n && \text{if } n \text{ is a number like } \dots -2, -1, 0, 1, 2\dots \\
 \text{eval_exp}(e_1 + e_2, \text{env}) &\stackrel{\text{def}}{=} \text{eval_exp}(e_1, \text{env}) + \text{eval_exp}(e_2, \text{env}) \\
 \text{eval_exp}(e_1 * e_2, \text{env}) &\stackrel{\text{def}}{=} \text{eval_exp}(e_1, \text{env}) * \text{eval_exp}(e_2, \text{env}) \\
 \text{eval_exp}(e_1 = e_2, \text{env}) &\stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } \text{eval_exp}(e_1, \text{env}) = \text{eval_exp}(e_2, \text{env}) \\ 0 & \text{otherwise} \end{cases} \\
 \text{eval_exp}(x, \text{env}) &\stackrel{\text{def}}{=} \text{env}(x)
 \end{aligned}$$

This environment *env* also acts like a map: it associates variables with their current values. For example after evaluating the first two lines in our factorial program, such an environment might look as follows

$$\boxed{\mathbf{a}} \mapsto 1 \quad \boxed{\mathbf{n}} \mapsto 5$$

Again I highlighted the keys. In the clause for variables, we can therefore consult this environment and return whatever value is currently stored for this variable. This is written *env(x)*—meaning we query this map with *x* we obtain the corresponding number. You might ask what happens if an environment does not contain any value for, say, the variable *x*? Well, then our interpreter just “crashes”, or more precisely will raise an exception. In this case we have a “bad” program that tried to use a variable before it was initialised. The programmer should not have done this. In a real programming language, we would of course try a bit harder and for example give an error at compile time, or design our language in such a way that this can never happen. With the second version of *eval_exp* we completed our definition for evaluating expressions.

Next comes the evaluation function for statements. We define this function in such a way that we recursively evaluate a whole sequence of statements. Assume a program p (you want to evaluate) and its pre-processed snippets sn . Then we can define:

$$\begin{aligned}
 eval_stmts([], env) & \stackrel{\text{def}}{=} env \\
 eval_stmts(label: rest, env) & \stackrel{\text{def}}{=} eval_stmts(rest, env) \\
 eval_stmts(x := e rest, env) & \stackrel{\text{def}}{=} eval_stmts(rest, env[x := eval_exp(e, env)]) \\
 eval_stmts(goto lbl rest, env) & \stackrel{\text{def}}{=} eval_stmts(sn(lbl), env) \\
 eval_stmts(jmp? e lbl rest, env) & \stackrel{\text{def}}{=} \begin{cases} eval_stmts(sn(lbl), env) & \\ \quad \text{if } eval_exp(e, env) = 1 & \\ eval_stmts(rest, env) & \text{otherwise} \end{cases}
 \end{aligned}$$

The first clause is for the empty program, or when we arrived at the end of the program. In this case we just return the environment. The second clause is for when the next statement is a label. That means the program is of the form $label: rest$ where the label is some string and $rest$ stands for all following statements. This case is easy, because our evaluation function just discards the label and evaluates the rest of the statements (we already extracted all important information about labels when we pre-processed our programs and generated the snippets). The third clause is for variable assignments. Again we just evaluate the rest for the statements, but with a modified environment—since the variable assignment is supposed to introduce a new variable or change the current value of a variable. For this modification of the environment we first evaluate the expression e using our evaluation function for expressions. This gives us a number. Then we assign this number to the variable x in the environment. This modified environment will be used to evaluate the rest of the program. The fourth clause is for the unconditional jump to a label, called lbl . That means we have to look up in our snippets map sn what are the next statements for this label. Therefore we will continue with evaluating, not with the rest of the program, but with the statements stored in the snippets-map under the label lbl . The fifth clause for conditional jumps is similar, but to decide whether to make the jump we first need to evaluate the expression e in order to find out whether it is 1. If yes, we jump, otherwise we just continue with evaluating the rest of the program.

Our interpreter works in two stages: First we pre-process our program generating the snippets map sn , say. Second we call the evaluation function with the default entry point and the empty environment:

$$eval_stmts(sn(""), \emptyset)$$

It is interesting to note that our interpreter when it comes to the end of the program returns an environment. Our programming language does not contain any constructs for input and output. Therefore this environment is the only effect we can observe when running the program (apart from that our interpreter

Figure 2: The entire code of the interpreter for our idealised programming language.

might need some time before finishing the evaluation of the program and the CPU getting hot). Evaluating the factorial program with our interpreter we receive as “answer”-environment

$$\boxed{a} \mapsto 120 \quad \boxed{n} \mapsto 0$$

While the discussion above should have illustrated the ideas, in order to do some serious calculations, we clearly need to implement the interpreter.

Scala Code for the Interpreter

Functional programming languages are very convenient for implementations of interpreters. A good choice for a functional programming language is Scala, a programming language that combines functional and object-oriented programming-styles. It has received in the last five years or so quite a bit of attention. One reason for this attention is that, like the Java programming language, Scala compiles to the Java Virtual Machine (JVM) and therefore Scala programs can run under MacOSX, Linux and Windows.⁵ Unlike Java, however, Scala often allows programmers to write very concise and elegant code. Some therefore say Scala is the much better Java. A number of companies, The Guardian, Twitter, Coursera, FourSquare, LinkedIn to name a few, either use Scala exclusively in production code, or at least to some substantial degree. If you want to try out Scala yourself, the Scala compiler can be downloaded from

<http://www.scala-lang.org>

Let us have a look at the Scala code shown in Figure ?? . It shows the entire code for the interpreter, though the implementation is admittedly no frills.

Static Analysis

Finally we can come back to our original problem, namely finding out what the signs of variables are

⁵There are also experimental backends for Android and JavaScript.