

CS 322 Spring 2014 Homework 3: Generating X86-64 Assembly Code

Due (via DTL) on Friday, June 6, 2014 at noon

This assignment asks you to make extensions and improvements to a code generator that converts 3-address IR (similar to that used in Homework 2) into X86-64 assembly code. The assignment assumes familiarity with the contents of Lab 7, where graph-coloring register allocation was discussed.

To complete the assignment, you will need a Java 1.6+ compiler (as usual) and access to an X86-64 compatible processor. With appropriate minor tweaks, the generated code should work on either Linux (e.g. the `linuxlab.cs.pdx.edu` machines) or MacOS. If you want to use Windows, you are on your own.

Your homework solution should be submitting using the D2L dropbox in the form of a single file called `sol3.zip`, which contains your revised implementations of `IR.java` and `Assignment.java`, together with any other auxiliary source files you created, and a `README` file. This zip file can be created using a command like:

```
zip sol3.zip IR.java Assignment.java README
```

More details on what should be in these files are given below.

This assignment will be scored out of 30 points, following the distribution indicated below. Your score on this assignment contributes 15% towards your final grade. A non-zero score is necessary to avoid receiving an automatic F for the assignment and hence for the class.

Provided Materials

You should download the “Homework 3 Materials” (`hw3.zip`) from the D2L website. Upon unzipping, you should see a `hw3` directory with the following items:

<code>X86Gen.java</code>	main program for generator
<code>IR.java</code>	IR representation extended with methods for X86-64 code generation
<code>Assignment.java</code>	simplistic register assignment
<code>Liveness.java</code>	computation of live-out sets, ranges and intervals (similar to Lab 7)
<code>X86.java</code>	utility support library for generating X86-64 code
<code>irParser.jj</code>	parsing for IR
<code>lib.c</code>	runtime library for generated X86-64 code
<code>IRInterp.jar</code>	executable version of IR interpreter
<code>Makefile</code>	for the generator
<code>tst/</code>	directory with sample test programs and reference outputs
<code>run</code>	sample script for running tests

The IR used here is very similar to that of Homework 2, but there may be a few minor differences.

A working X86 code generator can be built using the `Makefile`. On MacOS, first edit function `globalize` at the top of file `X86.java` in the obvious way.

The code generator reads an `.ir` file specified on the command line, uses the parser specified in `irParser.jj` to convert it to internal form, and then generates assembly code on standard output, which is normally captured into a `.s` file. The generator should be run with the `-ea` flag to enable Java assertion checking, e.g. as follows:

```
java -ea X86Gen f.ir 1> f.s
```

The generated `.s` file can then be assembled and linked via `gcc` with the “standard library” `lib.c` to make an X86-64 executable, as follows on Linux:

```
gcc -o foo foo.s lib.c
```

or on MacOS:

```
gcc -Wl,-no_pie -o foo foo.s lib.c
```

When executed, programs should behave as they do when evaluated using the IR interpreter (and, for those derived from `miniJava` sources, as they do when compiled using `javac`).

The `tst` directory contains examples of IR generated from `miniJava` source files, which are also included for your reference. The output to be expected by running each of these tests is in the corresponding `.out.ref` file. The assembly language outputs from a version of the generator with task 1 completed (see below) are in the corresponding `.s.ref` files. (These outputs are for linux; the MacOS versions are in a subdirectory `tst/mac`, so if you are developing on a Mac you'll want to copy these up one level to replace the linux versions.)

You will certainly wish to develop additional IR files to use as tests, which is most easily done by modifying one of the existing files. You can test the behavior of your IR code by running the interpreter; for example, to run file `foo.ir`, type

```
java -jar IRInterp.jar foo.ir
```

The `run` script can be used to automate some simple regression testing: read the comment at the top to see what it does.

Your Tasks

Make the following two improvements to the existing generator implementation in `IR.java` and `Assignment.java`. It is recommended that you attempt these improvements in the order listed here. Combine all your changes into single new versions of the files `IR.java` and `Assignment.java`. You may introduce additional Java source files if you find it useful. Also, you must write a brief `README` file indicating which improvements you attempted, and giving a summary of what you did for improvement 2 (below).

1. [20 pts.] The code generator assigns a unique X86 register to each `IR.Reg`. It does this once and for all for the entire function; this approach doesn't always produce great code, but it is simple and adequate. However, the register allocator in `Assignment.java` is much *too* simplistic: once an X86 register is used for one IR register, it is never re-used for another one, even if the two IR registers have disjoint live ranges. Thus, the generator quickly runs out of registers when given functions of even modest size (e.g. `test18`, `test19`, `test34`, and `test38`).

Your task is to improve this situation by implementing a graph-coloring register allocator instead, using the heuristic approach described in Lab 7. Use the existing `Assignment.java` file as the basis for your implementation. The segment of code to replace is clearly marked. Your new code should continue to use the existing `findAssignment` to choose registers. In more detail:

1. Build a register interference graph.
2. Initialize a stack which will contain pairs of the form (node, neighbors).
3. While the graph is non-empty:
 - (a) Find a minimum-degree node n . (If there is more than one such node, pick any one).
 - (b) Push n and a list of its neighbors onto the stack.
 - (c) Delete n and its incident edges from the graph.
4. While the stack is non-empty:
 - (a) Pop a node and its list of neighbors from the stack.
 - (b) Calculate the set of available registers (based on the registers already assigned to neighbors).
 - (c) Call `findAssignment` to pick a register for this node. If this call fails, raise an assertion and die.

To check the behavior of your assignment algorithm, you will wish to write new IR test programs and, most likely, to add printing facilities to debug the graph construction and traversal process. Note that the register assignment is included as a comment at the top of the generated `.s` files. The files `tst/*.s.ref` contain generated code that

was produced using our graph-coloring assignment algorithm. Your assignment output should be the same, modulo possible (consistent) renaming of registers.

2. [10 pts.] The current generator produces distinctly sub-optimal code for IR code derived from miniJava array operations. For example, consider the code generated for the first evaluation of `A[j]` in the `selectionSort` method in `test38` (after part 1, above, has been completed). It requires six X86 instructions, including an expensive (multicycle) `imulq`, something like this (your register names may differ):

```
# 9.    t2 = j * 4
        movq %rdx,%rsi
        imulq $4,%rsi
# 10.   t3 = A + t2
        movq %rsi,%r10
        movq %rdi,%rsi
        addq %r10,%rsi
# 11.   t4 = [t3]:I
        movslq (%rsi),%r8
```

This operation can certainly be done with fewer instructions; your task is to alter the generator so that it produces shorter and/or cheaper code sequences for array dereferences like this one. In fact, a one instruction sequence suffices:

```
movslq (%rdi,%rdx,4),%r8
```

To obtain full credit, you should achieve this optimization. But this is quite difficult! *Any* improvement in the sequence will do to obtain some credit; the more you improve it, the better your score. Note that your generator needs to produce improved code for *arbitrary* array dereferences, not just this particular one! But there's no point in spending effort optimizing other kinds of code: just focus on array dereferences.

Obviously, this task is rather open-ended, and is meant to encourage your creativity. Here are some hints and ideas to get you started:

- There are several cheaper ways to encode multiplication by powers of 2.
- The existing generator doesn't notice that most binary operations are commutative, and thus misses an easy way to shorten the sequence generated at line 10.
- One obvious problem is that the IR code generator itself produced three instructions (and two temporaries) from the original miniJava source, but the current X86 generator examines only one IR instruction at a time. You can alter the generator to examine multiple IR instructions at once to look for particular patterns. (This is awkward to do in an object-oriented style: be prepared to make heavy use of `instanceof`.)
- Alternatively, although you cannot change the form of the IR fed into your generator, you can alter the IR internally by *pre-processing* it into a variant form, perhaps having a large instruction set, before doing X86 instruction generation.
- In principle, another approach would be to *post-process* the generated X86 code, e.g. using peep-hole optimization to improve the instruction sequences. But since the generator currently just prints out X86 assembler directly rather than generating an internal form of it, this will be a lot of work, and is not recommended.
- For general inspiration, try writing simple C programs and compiling them using `gcc -S` to see what the generated `.s` assembly file looks like.

Whatever you do, your generator must continue to produce valid output on *all* IR inputs! Be careful: some changes that are innocent on these particular test IR files will cause other IR files to break. In particular, be sure to make use of

the `liveOutSets` information to make sure that you do not remove instructions that set registers whose values may be needed in the future.

To avoid getting bogged down in large implementation efforts that might not get finished (or yield improvements), it is strongly recommended that you start by making simple changes, and proceed to fancier ideas only as your time allows. Don't forget to record what you did in the `README` file.

Other Information

You may wish to refer back to the Lab 6 handouts on the X86-64 instruction architecture and using `gdb` for debugging assembly code. Remember to add a `.file` directive to the top of any `.s` file that you want to debug under `gdb`.

The existing code generator uses a convention whereby all values in registers are valid 64-bit quantities. (This is unlike the approach used in the Lab 6, where we tracked how much of each register was valid at each point, and selected the correct instruction forms accordingly.) This has the advantage that we can just use quad (64-bit) instruction forms throughout—except when loading or storing from memory, or for certain instructions that don't exist in quad form.

The existing code generator completely avoids the possibility of storing arguments or locals on the stack:

- Generation fails on functions with more than six arguments.
- Generation fails if there are insufficient registers for the register assignment algorithm we're using (i.e., we don't attempt to spill).

You should not worry about fixing these limitations (although of course they would be unacceptable in a real compiler). But of course one of the purposes of task 1 is to reduce the demand for registers.

The generated code obeys the X86-64 calling and register usage conventions. It does not use a frame pointer, since there is (essentially) no frame! Any callee-save registers used by a function need to be pushed onto the stack at the start of the function and popped back off at the end. Symmetrically, caller-save registers are pushed just before the call and popped back off just after the call.

Two registers (`%r10` and `%r11`) are reserved as local temporaries for use *inside* the translation of a single IR instruction.

Each IR file should contain a function called `_main`, representing the “main program” where the generated program should start running. The file `lib.c` contains an ordinary C `main` function that simply invokes `_main`. It also contains implementations of all the “built-in” functions used by `miniJava` programs.

MacOS configures `gcc` to associate the assembly-level name `_foo` with the C-level name `foo`. The function `globalize` at the top of `X86.java` compensates for this.