

CS 40 Assignment: Code Improvement Through Profiling

Note for Spring 2023:

Because of an operating system upgrade on CS machines, we are using `qcachegrind`. If there are problems running it, then you can `ssh` into `red-giant` or, from outside, into `blue-giant.cs.tufts.edu`, which are still running the older OS version, and run `kcachegrind`. If the newer version has any problems, please report them on Piazza, and we will send the information to EECS IT.

There is no design document for this assignment.

Contents

1 Purpose and overview	1
2 What we expect of you	1
2.1 Your starting point	2
2.2 Tracking changes as you make them	2
2.3 Laboratory notes	3
2.4 Analysis of assembly code	5
2.5 Performance of the final stages	5
2.6 What to submit	6
3 Methods of improving performance	7
4 Partial solution to the adventure game	8

1 Purpose and overview

The purpose of this assignment is to learn to use profiling tools to apply your knowledge of machine structure and assembly-language programming.

2 What we expect of you

Use code-tuning techniques to improve the performance of your choice of your um emulator. You are strongly encouraged to work with the same partner (if any) you had for HW6; that way, you don't have to learn about someone else's um implementation. If, however, you are working with a different partner, you may use either of your um emulators as a starting point. You will tune your Universal Machine running the large "sandmark" benchmark. If you and your partner do not have a working Universal Machine between you, it will be acceptable to borrow a Universal Machine from another student, but *only* if you have already submitted what you have for the Universal Machine assignment.

The key parts of the assignment are to *identify bottlenecks using qcachegrind* and to *improve the code by increments*. You will therefore want to *do most of your profiling on small inputs, as long as they are sufficient to give usefully measurable execution times and to exercise the code paths of interest*.

Your grade will be based on three outcomes:

- Your *laboratory notes* about the *initial state* of your program and *each stage of improvement*, including *differences from the previous stage*. This is important: improving your code is not the only goal; we also want to see that you have a computer scientist's approach to documenting the details that would allow someone else to understand the state of your system when you started, and why each of your changes had the effect that it did.
- Your analysis of the assembly code of the most expensive procedure in your final program.
- The *performance of your final-stage program*, measured as follows: Your Universal Machine running a large benchmark, not identical to the sandmark but closely related to it. That is: for your testing, you will have the same version of sandmark you had while developing your um. We will benchmark your program running a test that is similar but not quite identical.

2.1 Your starting point

Please begin with your code in the state it was after the Universal Machine assignment. If your code did not work, you may fix it, or you may start with another student's Universal Machine. If you have not yet completed the UM, you may not look at another student's UM until you have submitted.

Please take baseline measurements of your code *as submitted*.

If you fixed your UM before starting to tune it, also take baseline measurements using that version, and use those as the point of comparison for the rest of this assignment.

2.2 Tracking changes as you make them

During this assignment, you may run into a dead end that requires you to go back to a previous version of your code. We *recommend*, but do not require, that you use *git* to *keep track of each stage of improvement* in your code.

Unfortunately, *git* is a sophisticated professional's tool, with many, many options you won't need. Your first source for all things *git* should be a fellow student who has learned *git* in one of Ming Chow's courses, or failing that, Ben Lynn's online tutorial *Git Magic*. As a guide, you should be learning to *init* a repository, something you do once, and then to add files that have significant changes, and to label groups of additions by using *commit*. *git* will record every version of every file you have committed. One task you might need extra help with is if you want to go back to an older version of your files. There are several ways to do it. Look on the Web or ask for help.

Remember, per the course policy, you must *not* store your in a publicly accessible place. If you want to share a *git* repository with your partner, then you could either

- use our department GitLab instance at <https://gitlab.cs.tufts.edu>. Every student will have an account automatically created the first time you log in. Use your EECS username and password as you do logging into the homework servers. You can then share the repository with your partner (after they have logged in to create their account.)
- use a *private* repository on GitHub, which you may share with your partner. It must be private!

If you do not choose to learn *git* then you should definitely find some other means of keeping the original versions of your code and Makefiles, and every important intermediate version as well. One option is to make a new directory for each working version. Put into it a README that will remind you what is in the directory, what previous version it was based on, how well it worked, etc.

2.3 Laboratory notes

Begin by *choosing a data set*. For the Universal Machine, you will use the small midmark benchmark, the large sandmark benchmark, and a partial solution to the adventure game.

For *each stage* and for *each input*, please

- Report the user-mode time required to execute the program on the input, as measured with the `time` command (for information, try `man 1 time` and see the examples below). Be aware that *the C shell has a built-in time command*, and it stupidly writes to standard output instead of standard error. If you are using the C shell, you will need to use `/usr/bin/time`.

Each input is a Universal Machine binary, and “executing the program” means running `um` on that binary, supplying a suitable test input on standard input.

- For small inputs, report the total number of instruction fetches, which you can measure by running the program under `valgrind --tool=callgrind --dump-instr=yes`.
- Report *two different relative time values*:
 - The user-mode time of this stage divided by the user-mode time of the initial stage (time relative to start)
 - The user-mode time of this stage divided by the user-mode time of the *previous* stage (time relative to previous stage)

Lower relative times are better.

- In clear, correct English, say what was the bottleneck from the previous stage and how you improved the code.

You can see some sample reports (using made-up data, for an early version of the assignment that focused on `ppmtrans`) in Table 1.

We are asking you to record *user-mode* time, which is (roughly) the time the CPU actually spent running your code. It does not include time spent running background tasks or other programs that may be active on the same machine, and it does not include time spent waiting for disk (e.g. for page faults) or for network traffic.

When you change the code, it is critical that *each set of changes be small and isolated*. Otherwise you will not know what changes are responsible for what improvements.

1. Your *starting point* should be your code as submitted, compiled and linked with your original `Make-file`.
2. Your *first change* should be to compile with `-O1` and to link with `-lci40-01`, which must come *before* other libraries.
3. Your *second change* should be to compile with `-O2` and to link with `-lci40-02`.
4. After that you can start profiling with `callgrind` and `qcache` (tools covered in the profiling lab session) and improving your code based on the results.

Keep in mind that `-O2` is *usually* but not always better than `-O1` even though the GCC documentation suggests that it should be. The compiler tries to do more elaborate optimizations with `-O2`, but rarely those actually make things worse.

<i>Benchmark</i>	<i>Time</i>	<i>Instructions</i>	<i>Rel to start</i>	<i>Rel to prev</i>	<i>Improvement</i>
big	30s	—	1.000	1.000	No improvement (starting point)
small	1s	75.02×10^6	1.000	1.000	
big	28s	—	0.933	0.933	Compiled with optimization turned on and linked against -lcii-01
small	900ms	69.21×10^6	0.920	0.920	
big	28s	—	0.933	1.000	Compiled with optimization turned on and linked against -lcii-02
small	900ms	69.18×10^6	0.920	1.000	
big	25s	—	0.833	0.893	Removed Array_width() call from for loop and placed result in local variable instead
small	833ms	62.01×10^6	0.833	0.926	
big	22s	—	0.733	0.880	Removed array->blocks expression from loop and placed result in local variable
small	800ms	56.16×10^6	0.800	0.960	
big	60s	—	2.000	2.727	Used explicit for loop instead of blocked-array mapping function. Time improved for small image but got worse for big image—undid change.
small	650ms	49.20×10^6	0.650	0.813	
big	18s	—	0.600	0.300	Changed representation of blocks so that access to elements within the blocked mapping function uses unsafe pointer arithmetic without bounds checking
small	600ms	44.89×10^6	0.600	0.923	

Table 1: Sample report for blocked image rotation (made-up data)

2.4 Analysis of assembly code

Once you have improved the code as much as you can, use `valgrind` and `qcache-grind` to find the part of your program that takes the most time. This may be a function or a loop or a `switch` statement (or something else). The answer to this question isn't "main": we want you to identify a program hot spot. (You can find it by clicking on the Self tab in `qcache-grind`.) For this final phase you may want to use the `--dump-instr=yes` option so you can see the assembly code in `qcache-grind`. The advantage of `qcache-grind` over `objdump -d` is that it will tell you how many times each instruction was executed.

Once you've found the place where the most time is spent, examine the assembly code and either *identify specific ways (changes to the assembly code itself) in which it could be improved* or *argue that there is not an obvious way to improve the assembly code*.

Things to look for include:

- Do you see opportunities to keep data in registers, where currently there are unnecessary memory accesses?
- Do you see unnecessary computation in loops?

Be alert for a *horrible* idiom in the Intel assembly language: the instruction

```
mov %esi, %esi
```

looks redundant, but it's not. This idiom stands for an instruction that zeroes out the most significant 32 bits of the 64-bit register `%rsi`. We on the CS 40 staff think this is a bad bit of language design, but do watch out for it.

For this assignment, *there is no need to modify assembly code*.

2.5 Performance of the final stages

All profiling and measurements should be done on the machines in JCC 235 or 240. As of Spring 2019, all the machines in 235 and 240 are the same and those in Halligan 120 are different. The machines in Halligan 120 are currently usable. Since the machines in different rooms may be different, you must not mix reports of measurements run on one type of machine with those of another. You must do all your work on a single type of machine or you will go mad!

If you have trouble finding seats, you might make some progress by using another machine to build and test experimental versions, *being sure to keep all the intermediate versions of the source and Makefiles*. When the machine you need frees up, go back and measure and record all your intermediate results on that machine. Usually but not always, what makes an improvement on one machine will help on another. *That said, there is much less chance for confusion if you just choose a room and stick to it!*

To find out what kind of machine you're running on, use the command

```
grep 'model name' /proc/cpuinfo
```

and check whether the output is:

```
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
```

```
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
model name : Intel(R) Core(TM) i7-10700T CPU @ 2.00GHz
```

or:

```
model name : Intel(R) Core(TM) i7-7700T CPU @ 2.90GHz
model name : Intel(R) Core(TM) i7-7700T CPU @ 2.90GHz
model name : Intel(R) Core(TM) i7-7700T CPU @ 2.90GHz
model name : Intel(R) Core(TM) i7-7700T CPU @ 2.90GHz
model name : Intel(R) Core(TM) i7-7700T CPU @ 2.90GHz
model name : Intel(R) Core(TM) i7-7700T CPU @ 2.90GHz
model name : Intel(R) Core(TM) i7-7700T CPU @ 2.90GHz
model name : Intel(R) Core(TM) i7-7700T CPU @ 2.90GHz
```

The first represents the machines in JCC 235 and 240, the second represents the machines in Halligan 120 (as of November 2022).

To find out the RedHat release, use the command

```
cat /etc/redhat-release
```

which will say something like

```
Red Hat Enterprise Linux release 8.6 (Ootpa)
```

And you can check the gcc version using `gcc --version`, and will show something like `gcc (GCC) 11.2.0`

Measure your code with both `gcc -O1` and `gcc -O2`. Neither one is faster for all problems; report the better of the two results. In your final Makefile, use whichever gives the best results.

You will be evaluated both on your improvement relative to the code you start with and on the absolute performance of your final results. Your laboratory notes will record all your improvements and the performance of your final stages.

2.6 What to submit

Please use the `submit40-profile` script to submit the following items.

1. A Makefile, which when run with `make`, compiles all your source code and produces a um binary.
2. All benchmarks you used as test data.

3. A README file which

- Identifies you and your programming partner by name
- Acknowledges help you may have received from or collaborative work you may have undertaken with others
- Explains what routine in the final `um` takes up the most time, and says whether the assembly code could be improved
- Says approximately how many hours you have spent *analyzing the problems posed in the assignment*
- Says approximately how many hours you have spent *solving the problems after your analysis*

Last but not least, please submit your `labnotes.pdf` file via Gradescope.

3 Methods of improving performance

In performance, really big wins usually come from better algorithms which provide superior asymptotic complexity. But for our programs, there is not much algorithmic action; everything should be pretty much linear in the number of UM instructions executed. You may find some non-linearities relating to the patterns of access to segments, and those might be worth addressing. You can often improve things by *changing your data structures*.

Here are some trite but true observations:

- *Memory references are expensive*, especially when data is not in the cache. In fact, compared with memory references, arithmetic with values in registers is practically free. If you give `valgrind` the `--simulate-cache=yes` option, it will count loads and stores and also simulate the cache. I don't see how to get the load/store data without also running an expensive cache simulation.
- In general you will be improving performance of your base implementation, but it's surprisingly easy to gain speed at the expense of excessive memory usage. *As with the original UM assignment, you should make sure you're not letting your memory grow without bound.*
- On AMD64, *calls to leaf procedures* are pretty cheap, but *calls to non-leaf procedures* can be expensive.

What if your program is nothing *but* memory references and procedure calls?! How can you make progress?

- To know what to improve, *you must profile*. Measure, measure, and measure again. Your best friends are `valgrind --tool=callgrind` and the `qcachegrind` visualizer.¹

Nothing is more frustrating than to spend a lot of time improving code that is rarely executed.

- The C compiler can be stupid about memory references. Because of pointer aliasing, if you write to memory, the C compiler may assume that *all* values in memory have changed, and may have to be reloaded.

¹For some programmers in some cases, `gprof` can be a pretty good friend, but it is useful only if you have access to all the source code, including libraries. And `gprof` does not have a good visualization tool like `qcachegrind`. In fact, the damn thing won't even report all the data it has because it uses only two digits after the decimal point. Beastly `gprof` is not *my* friend.

- The C compiler has no idea when multiple calls to a function will return the same value. If you do have an idea, you can help out the C compiler by putting results in local variables.
- The C compiler has to assume that a function call could scribble all over memory. After a function call, values referenced through pointers may have changed. If *you* know the values haven't changed, make sure those value are sitting in local varaibles, so that the compiler knows it too.
- The C compiler is *staggeringly good* at managing local variables and putting them in machine registers. All you have to do is get your values into local variables; the compiler will do the rest. This is a big change from the 1970s!
- If a lot of time is spent in one procedure, like say `UArray_at`, you often have two choices: make each call of the procedure run faster, or change code somewhere else so the procedure is called less often.
- Your book by O'Hallaron and Bryant devotes an entire chapter to code improvement (Chapter 5). There's about 15 pages' worth of really good low-hanging fruit, and then there are a lot of details.
 - The first part, through the end of Section 5.1, gives an excellent explanation of aliasing and will help you understand the pessimism with which the compiler must treat memory references.
 - Sections 5.2 and 5.3 present a basic framework and example. If you like toy benchmark programs and graphs with lines on them, these sections are for you.
 - Sections 5.4 to 5.6, which comprise only ten pages, give more detailed explanations of the most important of the techniques I've sketched above.
 - Section 5.7 tells a complicated story that explains some of the complexities of modern processors. However, the focus on the Intel Core i7 may not be helpful for understanding the behavior of the Core 2 Quad processors on the lab machines, as there are some differences.
 - Sections 5.13 and 5.14 discuss the use of a profiler. Unfortunately the chapter refers to `gprof`, which is a legacy tool that I recommend against using unless you are stuck with a problem for which `valgrind` is just too slow.
 - Sections 5.8 to 5.10 describe program transformations which, for the most part, a good optimizing compiler can do better than you can.
 - Section 5.15 summarizes material in earlier sections. Perhaps you will find the summaries useful for review?
 - Sections 5.11 and 5.12 present material that I consider interesting and important but well beyond the scope of CS 40. This material is more likely to be taught in a 100-level architecture course aimed at juniors, seniors, and beginning graduate students.

4 Partial solution to the adventure game

Figure 1 gives a a partial solution to the adventure game. The order goes down by column. This solution can be made into a benchmark that is intermediate in difficulty between the midmark and the sandmark.

n	take screw
take bolt	take motherboard
take spring	comb motherboard screw
inc spring	take A-1920-IXB
take button	comb A-1920-IXB bolt
take processor	comb A-1920-IXB processor
take pill	comb A-1920-IXB radio
inc pill	take transistor
take radio	comb A-1920-IXB transistor
take cache	comb motherboard A-1920-IXB
comb processor cache	take keypad
take blue transistor	comb keypad motherboard
comb radio transistor	comb keypad button
take antenna	s
inc antenna	

Figure 1: Partial solution to the adventure game