

# Coding for Researchers

May 25, 2016

Northwestern University

Peter Frazier

# How I learned (then later forgot) to write good code

- I earned my undergraduate degree from Caltech in 2000, with a dual major in computer science and physics.
- As a CS student, I wrote code in various languages (C, C++, Java, Lisp, UNIX shell scripting) for:
  - class assignments
  - one semester long course project
  - summer jobs
  - projects for fun

# How I learned (then later forgot) to write good code

- Upon graduation, I joined a small startup company founded by a friend of mine.
- I mostly developed a web client in PHP, with bits of C and Perl along the way.
- I learned that working in a group on a big project is different than working individually on small projects.

# How I learned (then later forgot) to write good code

- Our startup eventually folded, and I joined another startup.
- I wrote code in C for the backend of a relational database.
- My code was reviewed by senior programmers with excellent skills, who taught me to improve my code.
- This is where I really learned how to code.

# How I learned (then later forgot) to write good code

- This startup eventually laid most people off, and I joined a big company.
- I wrote code in C for the filesystem of a proprietary UNIX maintained by a relational database company.
- I learned that writing code at a big company is very different from writing code at a small company.

# How I learned (then later forgot) to write good code

- I got bored, and decided to go get a masters & PhD in Operations Research at Princeton.
- I was a student from 2005 until 2009, and wrote code in Matlab and Java for my dissertation research, along with bits of C and Java.
- I learned that writing code for research as a graduate student is very different than writing code for customers as a software engineer.

# How I learned (then later forgot) to write good code

- In 2009, I graduated and became an assistant professor at Cornell's Operations Research department.
- I stopped writing lots of code myself, and started supervising students who write lots of the code.
- I learned that supervising students who write code is very different than writing code yourself.
- My coding skills slowly declined as I fell out of practice...

# How I learned (then later forgot) to write good code

- In 2015, I went on sabbatic leave, and joined Uber's data science team, first as an individual contributor, and then as the manager for data science on uberPOOL.
- I wrote part of the pricing system for uberPOOL in python, then handed the codebase off to my team.
- I learned that knowing how to code makes you a better data scientist. I also learned some things about how to do useful research.
- My coding skills haven't fully recovered from years of disuse, but I still help to tell you some things you can use...



# Is it worth your time to improve your coding?

- In this talk I will make some suggestions about things you can do to improve your coding.
- Acting on any one of these suggestions would require time & energy.
- This is time & energy that you could instead be using to “be productive”.
- You will have to decide whether it is worth it.

# Improving your coding has these advantages

- 1. You will be able to do research more quickly.
- 2a. Your code will be easier for others to understand and use.
- 2b. This will make others more likely to use your code.
- 2c. This will increase the impact of your research.
- 3. Coding well may help you get a job, and will make you more effective in that job after you get it.

# Part I: Coding for Researchers

- Learn from people who code well
- Split your code into functions
- Create unit tests
- Document your code
- Add error handling
- Use a good text editor or IDE
- Use version control
- Code in Python instead of Matlab
- Save the input & output to expensive experiments
- Save the raw data you use to generate plots
- Use profiling to observe which parts of your program are slow.
- Understand a bit about computer architecture.
- Use system monitoring tools to observe CPU, disk and network to find the bottleneck.

# Part 2: Some things I've learned at Uber about doing useful research

- Impact is a choice
- Ease of implementation matters
- Communicate well
- Theory is not convincing
- Own your research

# Part I: Coding for Researchers

# Learn from people who code well

- If your office mates code well, ask them questions.
- If your office mates code badly, help them.
- If you have the opportunity:
  - Look at well-written code.
  - Encourage a good coder to review your code & make suggestions.

# Split your code into functions

- Each function should do a specific task, transforming some well-defined inputs into some well-defined outputs, in a well-defined way.
- Aim for no more than 200 lines of code in a function.
- If a function is longer, ask whether it can be broken into smaller functions.

# Split your code into functions

- “A program should not require its readers to hold more than a handful of facts in memory at once”
- Breaking code into functions reduces the number of things you need to remember when writing code.
- This reduces bugs, and allows you to write code faster.



# Split your code into functions

- Splitting your code into functions makes it easier for others to understand your code.
- It also makes it easier for you & others to reuse your code.

# Create unit tests

- A unit test is a piece of code that tests a single function or unit of your code.
- They are worth the work because it is easier to fix bugs if you find them early.
- Imagine being close to submitting a paper, or finishing a chapter of your dissertation, and discovering that your main numerical result is invalid due to a bug.
- Find bugs early!

# Here are some example unit tests

- For code that solves a dynamic program, simulate the optimal policy and compare its estimated value to what the solution code claims it should be.
- For code that estimates something from data, generate lots of fake data where you know the something and see if your code can estimate it correctly.

# Here are some more unit tests

- If you wrote some old code that computes something in a special case, and new code that computes it more generally, run both codes on random instances where both should work, and compare the output.
- If you wrote old slow code and new fast code that compute the same thing, run them both on random instances and compare the output.
- Run generic code on random inputs and see if it crashes.

# Turn bugs into unit tests

- If a bug arises outside of a unit test:
  - Write a new unit test that reproduces the condition under which that bug arose, and checks whether that bug is still happening.

# Here's how I use unit tests

- I make my unit tests functions `test1()`, `test2()`, etc. that take no arguments, print out a message “OK” or “FAIL”, and return a single bit.
- I then make another function `test()` that runs all the unit tests, one after the other, and reports whether everything passed or not.
- This makes it easy to test whether everything still works after a software upgrade, or on a new system.

# Document your code

- Describe inputs & outputs at the beginning of your functions.
- Use variable and function names that mean something (without being too long).
- If you have a lot of files, separate them into directories that make sense, and add a README.txt file in each.
- This is useful because:
  - it makes it easier for others to read your code
  - it makes it easier for you to read your code, several months down the road when you've forgotten what you were thinking when you wrote it.
  - it makes you think about the structure of your code.

# Here is some example documented code

```
%[impl,x,y,extras] = SimAdaptiveCKG(N,sample,dim,n0a,n0b,MPerDim,MLE)
%
% Simulates CKG algorithm through N samples, and collects at each sample time
% 1<=n<=N the measurement and implementation decision. Uses the passed
% MLE to estimate hyperparameters, including the noise variance.
%
% N -- number of measurements to take.
% sample(x) -- function that, given a (flattened) measurement location x, returns a
%   noisy sample of the function value at that point.
% dim -- number of dimensions in the search domain.
% n0a -- number of measurements to put in the first stage latin hypercube.
% n0b -- number of measurements to put in the second stage, when we measure
%   again the n0b best points.
% MPerDim -- number of discretizations in each dimension. A scalar.
% MLE -- function handle of the form,
%   [cov0hat, noisevarhat, scalehat, betahat] = MLE(xd, y)
%
% impl -- vector of implementation decisions
% x -- vector of measurement decisions
% y -- vector of observations
% extras.cov0 -- array of estimated cov0
% extras.mu0 -- array of estimated mu0
% extras.scale -- array of estimated scale
% extras.noisevar -- array of estimated noise variances
% extras.logkgfactor -- array of max(logkgfactors)
%
function [impl,x,y,extras] = SimAdaptiveCKG(N,sample,dim,n0a,n0b,MPerDim,MLE)
    M = MPerDim^dim;
    impl = zeros(1,N); % preallocate for speed.
    x = zeros(1,N);
    y = zeros(1,N);
    extras.cov0 = NaN*ones(1,N);
    extras.mu0 = NaN*ones(1,N);
    extras.scale = NaN*ones(dim,N);
    extras.noisevar = NaN*ones(1,N);
    extras.logkgfactor = NaN*ones(1,N);

    toZd = @(z) ZToZd(z,MPerDim,dim);

    % Sample n0a points using latin hypercube sampling.
    n0 = n0a+n0b;
    assert(n0a>0);
    xd(1:n0a,:) = DiscreteLHS(MPerDim,dim,n0a);
```



# Documenting your code is useful

- Documenting your code is useful because:
  - it makes it easier for others to read your code
  - it makes it easier for you to read your code, several months down the road when you've forgotten what you were thinking when you wrote it.
  - it makes you think about the structure of your code.

# Documenting your code is useful some of the time

- Sometimes we write code that we are just going to throw away right afterward.
- If you really will throw it away, it is ok not to document it.
- It is often hard to know, when you write it, whether you will throw it away.
- If you find yourself continuing to use old code that is undocumented, take the time to clean it up & document it.

# Add error handling

- If a function assumes something about its inputs, add code that checks whether these assumptions are met.
- If they are not, print out an error message and return from the function.
- This makes it easier for you to fix it when your code crashes.
- This makes your code catch errors that would have silently caused incorrect output.
- This makes your code more usable by others.

# Use a good text editor or IDE

- Throwbacks like me are familiar with one of two reasonable choices:
  - vi / vim
  - emacs
- I also know good programmers that use IDEs (integrated desktop environments) like eclipse.
- It should be possible to move the cursor anywhere on the screen, or cut & paste from any location to any other location, with 5 keystrokes or less (and no mouse), most of the time.

# If you use vim, also try the Vimium chrome extension

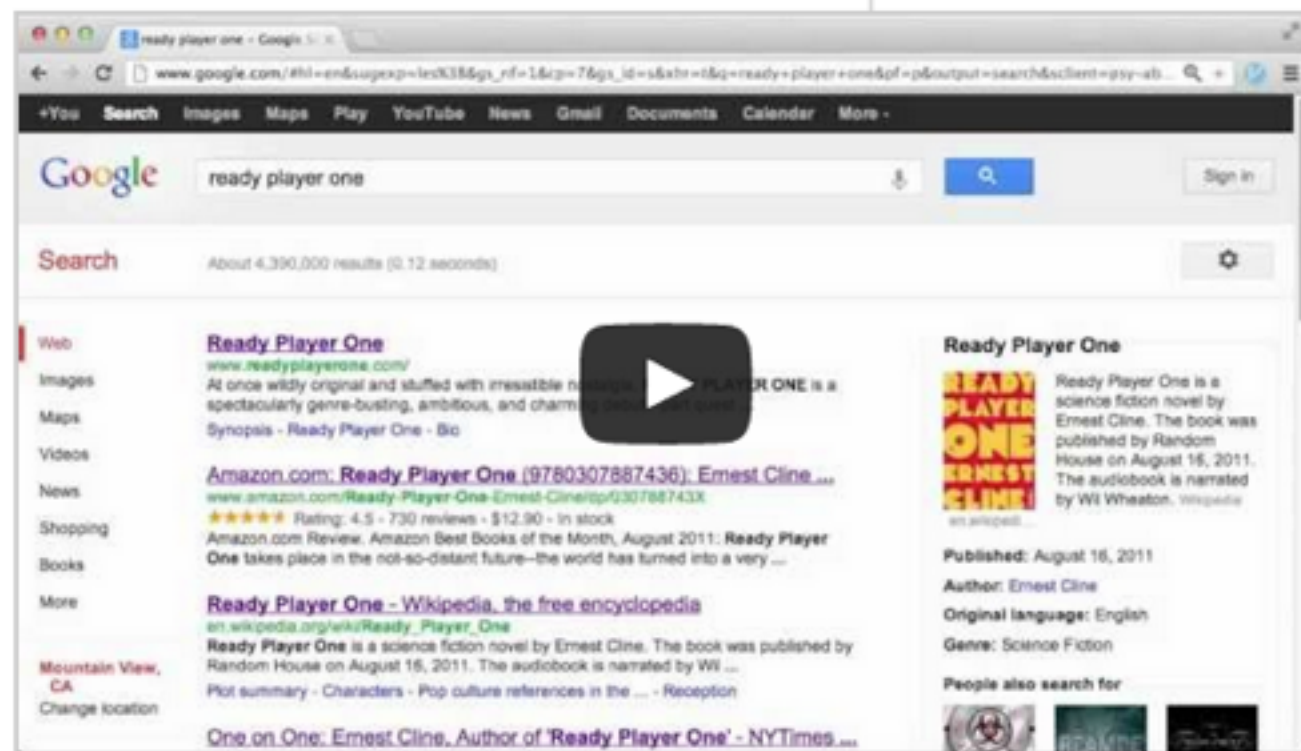


[View source on GitHub](#)

Vimium is a Google Chrome extension which provides keyboard shortcuts for navigation and control in the spirit of the Vim editor. Watch this video to get a feel for how it works:

## Features

- Helps you navigate the web without touching the mouse.
- Uses a clever highlighting method to navigate using links.
- Customizable keyboard shortcuts
- Has an in-page help dialog to remind you of your personalized shortcuts.



## Keyboard shortcuts

### Navigating history

**H** : Go back in history  
**L** : Go forward in history

### Navigating the page

**?** : Show the help dialog  
**j** : Scroll down

**u**, **<c-u>** : Scroll a half page up

**d**, **<c-d>** : Scroll a half page down


**r** : Reload the page

**gs** : View page source

# Use version control (git w/ bitbucket or github)

[US] [https://bitbucket.org/peter\\_frazier/vascular](https://bitbucket.org/peter_frazier/vascular)

Bitbucket Teams Repositories Create  ?


 **vascular**  
peter\_frazier Share

Clone Branch Pull request


Overview Source Commits Branches Pull requests Downloads

Research project in post-operative surveillance for patients undergoing surgery to repair a loss of blood flow. This repository was originally hosted on Cornell sourceforge, and was migrated to bitbucket.

Recent activity

 peter\_frazier pushed 3 commits to peter\_frazier/vascular 22 hours ago

- c02f23e - Adding the additional claudicant statistical fit parameters to the
- fa7ff89 - Modifying the DMHI MLE code to also run on claudicants. Also mod...
- c8b8c6a - Continuing to work on the MDM paper, while camping.

 peter\_frazier pushed 3 commits to peter\_frazier/vascular

SSH

1 Branch 0 Tags 0 Forks 1 Watcher

Owner	peter_frazier
Access level	Private
Type	Git
Last updated	22 hours ago
Created	2013-06-04
Size	17.6 MB (download)

# There is a nice tutorial on the bitbucket website

The screenshot shows a web browser window with the URL <https://confluence.atlassian.com/display/BITBUCKET/Bitbucket+101>. The page is titled "Bitbucket 101" and is part of the "Atlassian Documentation" for "Bitbucket". The left sidebar contains a "Bitbucket Documentation" menu with a "Bitbucket 101" section expanded, listing various topics like "Create a repository", "Import code from an existing project", "Repository privacy, permissions, and more", "Manage an individual account or a team", "Administer a repository", "Use your repository", "Use the SSH protocol with Bitbucket", "Mark up comments, issues, and commit messages", "Use the issue tracker", "Use a wiki", "Convert from other version control systems", "Use the JIRA DVCS Connector Plugin", "Use the Bitbucket REST APIs", "Write brokers (hooks) for Bitbucket", "Frequently Asked Questions", and "Support and other resources".

The main content area is titled "Bitbucket 101" and includes the following text:

Bitbucket / Bitbucket Documentation Home

## Bitbucket 101

If you are new to hosting your code, code management with distributed version control systems (DVCS), or either Git or Mercurial, this Bitbucket 101 tutorial gives you a taste of all of them. In this tutorial, you'll first install both Git and Mercurial for your operating system. You'll do some work using both Git and then Mercurial. Throughout, you'll use the hosted code management system that is Bitbucket. The tutorial consists of the following pages:

- [Set up Git and Mercurial on Windows](#) (or [Mac OSX](#) or [Linux](#))
- [Create an Account and a Git Repo](#)
- [Clone Your Git Repo and Add Source Files](#)
- [Fork a Repo, Compare Code, and Create a Pull Request](#)
- [Add Users, Set Permissions, and Review Account Plans](#)
- [Set up a Wiki and an Issue Tracker](#)
- [Set up SSH for Git](#)
- [Mac Users: SourceTree a Free Git and Mercurial GUI](#)
- [Sourcetree Git GUI on Windows](#)

The tutorial teaches you some simple DVCS workflows using basic Git and Mercurial commands.

On the right side of the page, there is a blue rounded rectangle containing the text:

**60**  
minutes or less

Looks like a lot of work? No worries, it does not take as long as you think.

# Code in Python instead of Matlab

- Python is just as easy to learn as Matlab.
- Matlab is commercial. Python is free.  
If you release Python code, anyone can use it.  
If you release Matlab code, only Matlab users can use it.  
[You can release your Matlab as a binary, or port to octave, but for this to be useful to users, you need to actually do it.]
- If you go on the job market in Silicon Valley,  
interviewers will view knowledge of Python positively.  
They will view knowledge of Matlab negatively.



# Code in Python instead of Matlab

- Matlab doesn't make it easy to split code into functions.
- Matlab doesn't make it easy to split functions across files.
- Matlab doesn't make it easy to group files into modules.
- Matlab doesn't make it easy to write unit tests.
- Python makes all of this easy.
- When you are writing Matlab, and are pressed for time, it is very tempting to write bad code.  
(e.g., 2000 lines of code in one Matlab script, with no unit tests.)
- When you are writing Python, and are pressed for time, you write better code.

# Code in Python instead of Matlab

- Reasonable reasons to write Matlab:
  - Your target users are Matlab users.
  - Your code relies on a legacy Matlab codebase.
  - You know Matlab and nothing else and you have a deadline. (but learn Python after the deadline is done.)
- If you want to write in R or Java or C or C++, then that is ok, but you should still learn Python.

# Save the output from expensive experiments

- Suppose your computational experiments take several hours to run.
- Have your code output & save **everything** you might reasonably be interested in the future, subject to space constraints.
- For example, if you are calculating mean & variance using Monte Carlo, save each replicate.
- Have your code save its inputs & outputs to a file.
  - In Matlab, .mat format is convenient.
  - In many languages (R, C/C++, Matlab), .csv is convenient and portable.
  - In object-oriented languages (C++, Python, Java), you can create classes that are serializable, which allows writing and reading complex objects directly to files.

# Saving output from expensive experiments has advantages

- Separate your post-processing code by having it read from these stored data files, rather than directly from the expensive experiment code.
- Later, if you want to perform some other analysis on the data, or fix a bug in your post-processing code, you don't need to re-run the experiment.
- You can write your post-processing code in a different language than the expensive experiment code.
- For example, write your expensive experiment in C, and your post-processing in R.

# Save the input to expensive experiments

- Several months later, when you come back to a saved output file, you won't remember the inputs or version of the code that you used to generate it.
- Have your experimental code print out & save:
  - all inputs
  - the current date and time
  - a version number for the code.

# If you use git, here is a nice way to get a “version number”

- Run “git rev-parse HEAD”
- This returns the hash (unique identifier) of the most recent commit.
- You can then look up this commit later in git’s version history.

# If you use git, here is a nice way to get a “version number”

I ran this on the command line: `git rev-parse HEAD`  
`cb85f1e8687de718ae41d8f4e9fadb3970d508b3`

The screenshot shows the Bitbucket web interface for a repository named 'vascular' owned by 'peter\_frazier'. The interface includes a navigation bar with 'Teams', 'Repositories', and a 'Create' button. Below the repository name, there are tabs for 'Overview', 'Source', 'Commits', 'Branches', 'Pull requests', and 'Downloads'. The 'Overview' tab is selected, showing a description of the repository and a 'Recent activity' section. In the 'Recent activity' section, a commit by 'peter\_frazier' is highlighted with a red box around the commit hash 'cb85f1e'. A red arrow points from the text 'I can then compare the first 7 characters to the website:' to this commit hash.

Bitbucket Teams Repositories Create

vascular peter\_frazier Share

Clone Branch

Overview Source Commits Branches Pull requests Downloads

Research project in post-operative surveillance for patients undergoing surgery to repair a loss of blood flow. This repository was originally hosted on Cornell sourceforge, and was migrated to bitbucket.

Recent activity

peter\_frazier pushed 1 commit to peter\_frazier/vascular 3 minutes ago

cb85f1e Adding some documentation to execute.m

peter\_frazier pushed 2 commits to peter\_frazier/vascular 2013-12-05

SSH 1 Branch

Owner Access level Type Last updated Created

I can then compare the first 7 characters to the website:

# It might also be a good idea to save the output of “git status”

- In theory, one would always commit all of your changes before running an expensive experiment (keep in mind that you can use a branch if this is just a “throw away” experiment).
- In practice, one might forget to do this.
- Saving the output of “git status” will let you know whether you have any uncommitted changes.



# Here is some example output from “git status”

```
en-or-apf98-a-2% git status
# On branch master
# Changes not staged for commit:
#   (use "git add <file>..." to update what will be committed)
#   (use "git checkout -- <file>..." to discard changes in working directory)
#
#       modified:   ../presentations/2013.06.INFORMS_Healthcare/2013.06.INFORMS_Healthcare_Frazier_Vascular.key
#
# Untracked files:
#   (use "git add <file>..." to include in what will be committed)
#
#       ../code/opt/my_fatality.m
#       ../code/opt/my_fminconop.m
#       ../code/stat/Test_stats.m
#       ../data/Copy of WCMC_LEDBS_safe.xlsx
#       ../data/Copy of WCMC_LEDBS_safe9-27_2.xlsx
#       ../data/ResearchData.xlsx
#       ../data/WCMC_LEDBS_safe.xlsx
#       ../data/WCMC_LEDBS_safe10-18.xlsx
#       ../data/WCMC_LEDBS_safe10-18_plus_cols_AP_AQ.xlsx
#       analysis/.plotOptimalVsEndTime.m.swp
#       stat/.Oct18_erlang_exp_likelihood.m.swp
#       ../mdm_paper/MDM_paper_AJM.docx
no changes added to commit (use "git add" and/or "git commit -a")
```

# Similar things can be done with subversion

- Use “svn info” to get the current revision.
- Use “svn status” to get a list of modified files.

# Save the raw data you use to generate plots

- When you generate a figure that you might want to later include in a paper or presentation, save the raw data you used to generate it.
- When (inevitably) you later decide to change the formatting, you will be able to do so easily.

# Here's how I make it easy to regenerate figures

- When I get a figure I might want to re-use, I write a script that reads in the raw data, creates the figure, and writes the figure to a file.
- I save the data, script, and figure in the repository.
- Later it is easy to find the data & script to regenerate the figure because:
  - The script has the figure filename in it, and so will turn up in searches.
  - I usually name the script with a similar name to the figure.

```

% Creates the following plots, all in the fig/ directory:
%     plotOptimalVsEndTime_improvement_rest.pdf
%     plotOptimalVsEndTime_improvement_tissue.pdf
%     plotOptimalVsEndTime_prob_rest.pdf
%     plotOptimalVsEndTime_prob_tissue.pdf
% The "improvement" plots show the percentage improvement of the optimal over
% status quo, as a function of the time horizon, and letting the optimal
% schedule use the same number of checkups as the status quo over that
% time-span.
%
% The "prob" plots show the the probability of an emergency, again as a
% function of the time horizon, letting the optimal schedule use the same
% number of checkups as the status quo.
%
% There is one "improvement" plot and one "prob" plot for each patient sub-group
% (rest pain and tissue loss).
%
function plotOptimalVsEndTime()

% Values in paper for rest pain
% lambda = 0.0338;
% alpha = 0.450;

% New values for rest pain
lambda_rest = 0.0336;
alpha_rest = 0.242;

% Values in paper for tissue loss
% lambda_tissue = 0.0530;
% alpha_tissue = 0.9022;

% New values for tissue loss
lambda_tissue = 0.0389;
alpha_tissue = 0.6798;

make_plot(lambda_tissue,alpha_tissue,'tissue', 'Tissue Loss');
make_plot(lambda_rest,alpha_rest,'rest','Rest Pain');
end

function make_plot(lambda,alpha,file_extension,plot_title)
endTime = [2:52];
for i=1:length(endTime)
    foptimal(i) = CalculateOptimal(endTime(i),lambda,alpha);
    foriginal(i) = CalculateOriginal(endTime(i),lambda,alpha);
    improvement(i) = (foriginal(i) - foptimal(i)) / foriginal(i);
end

plot(endTime,foriginal,'r--', endTime,foptimal,'k-')
set(gca,'fontsize',18);
legend('optimal', 'status quo','Location','SouthEast');
title(plot_title);
xlabel('Time Horizon (4wk periods)');
ylabel('Prob(undetected failure)');
axis([1 52 0 0.8])
print('-dpdf', sprintf('fig/%s_%s.pdf', 'plotOptimalVsEndTime_prob', file_extension)
);

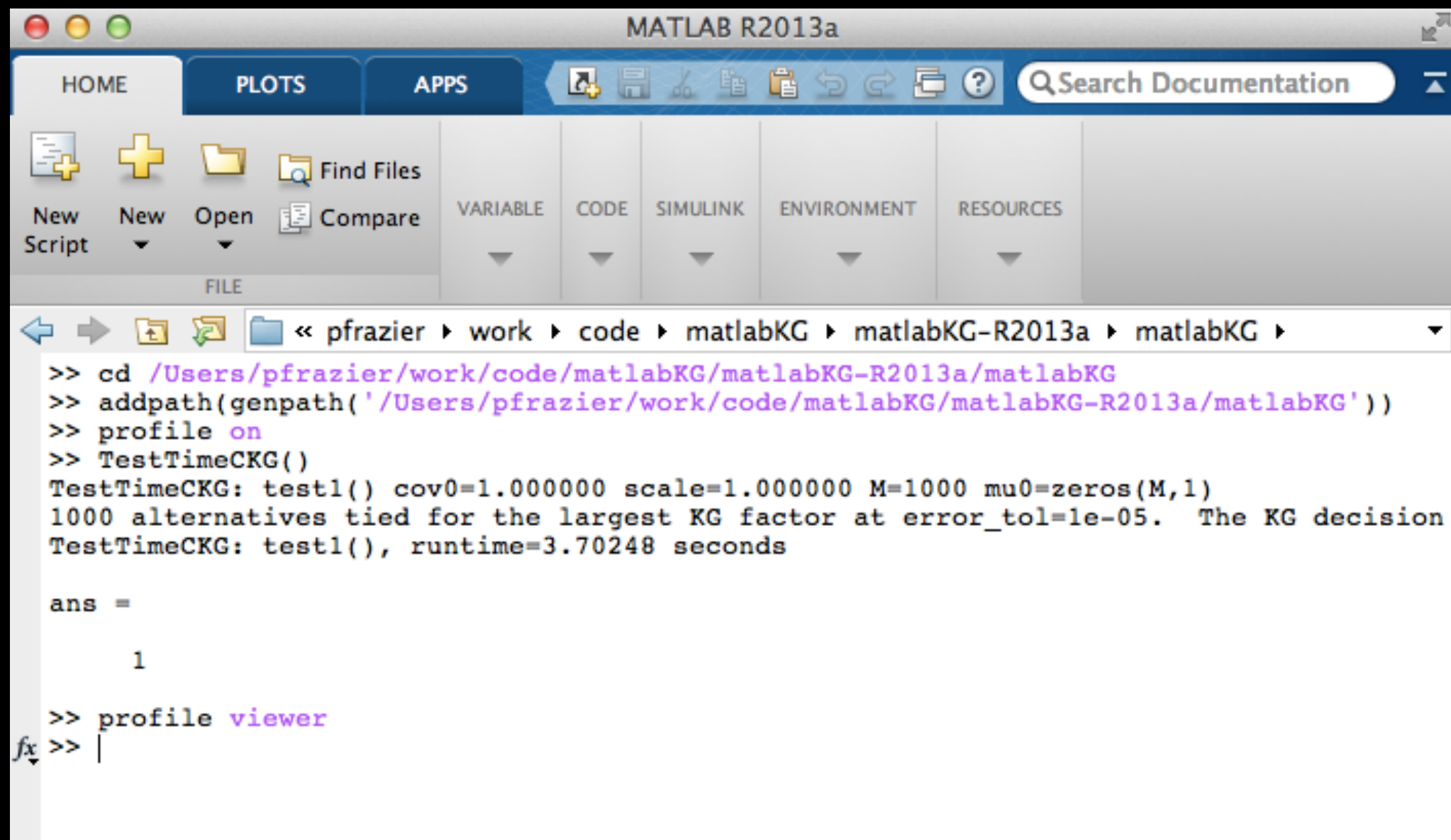
plot(endTime,improvement*100)
set(gca,'fontsize',18);
title(plot_title);
xlabel('Time Horizon (4wk periods)');
% ylabel('Percent improvement (status quo - optimal) / (status quo)');
ylabel('Percent improvement');
axis([1 52 0 max(improvement)*100*1.1])

```

# To make your program faster, first understand why it is slow

- Use **profiling** to observe which parts of your program are slow.
- Understand a bit about **computer architecture**.
- Use **system monitoring tools** to observe CPU, disk and network to find which is the bottleneck.

# Here's how to profile in Matlab



The image shows the MATLAB R2013a interface. The top toolbar includes tabs for HOME, PLOTS, and APPS, along with various icons for file operations and documentation. Below the toolbar, there are buttons for 'New Script', 'New', 'Open', and 'Compare', followed by a 'FILE' tab. The main workspace area displays a command window with the following text:

```
>> cd /Users/pfrazier/work/code/matlabKG/matlabKG-R2013a/matlabKG
>> addpath(genpath('/Users/pfrazier/work/code/matlabKG/matlabKG-R2013a/matlabKG'))
>> profile on
>> TestTimeCKG()
TestTimeCKG: test1() cov0=1.000000 scale=1.000000 M=1000 mu0=zeros(M,1)
1000 alternatives tied for the largest KG factor at error_tol=1e-05. The KG decision
TestTimeCKG: test1(), runtime=3.70248 seconds

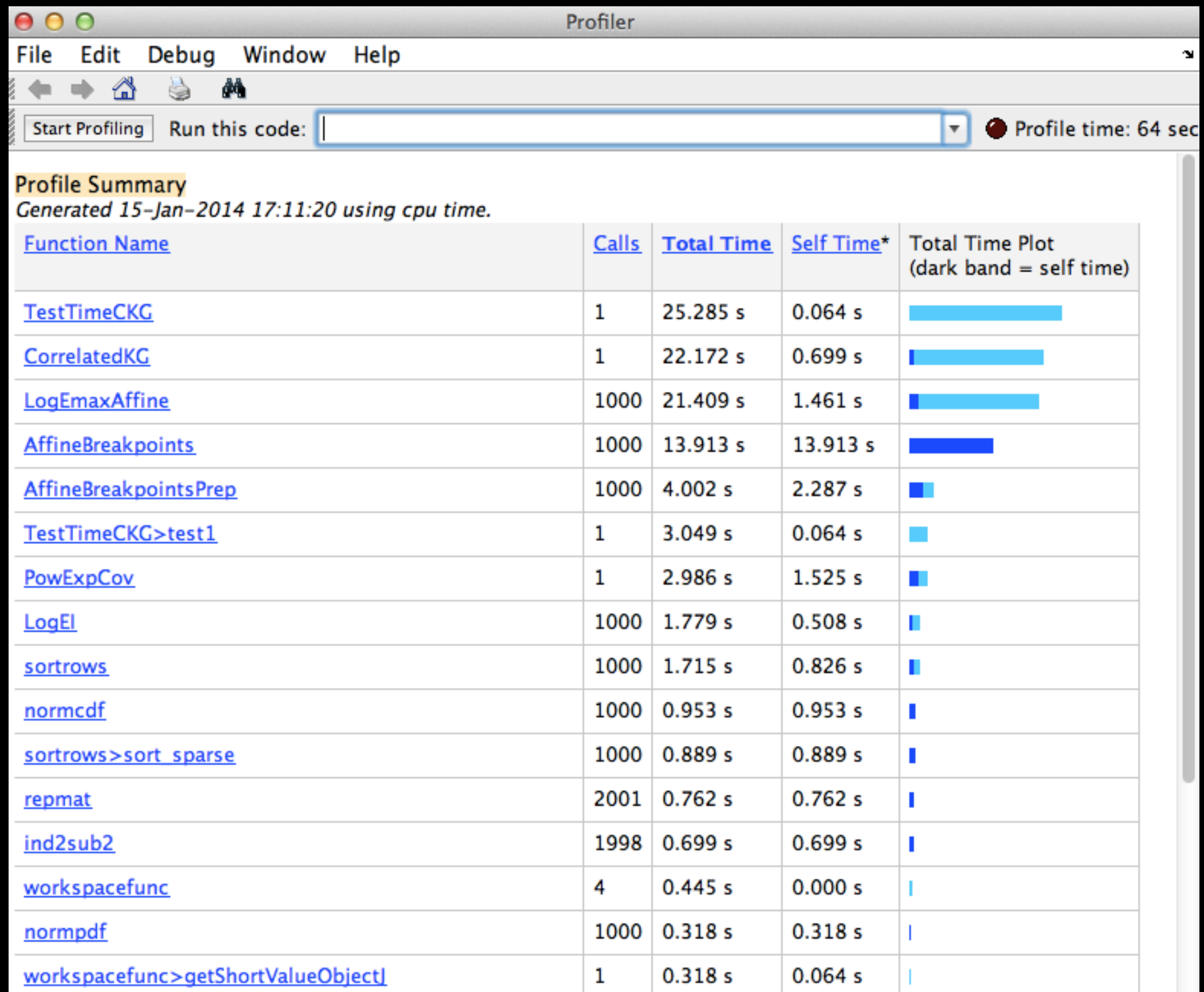
ans =

     1

>> profile viewer
fx >> |
```

- Run “profile on” to start profiling.
- Run “profile off” to stop.
- Run “profile viewer” to see the results.

# Here's how to profile in Matlab





# Here's how to profile in Matlab

The screenshot shows the MATLAB Profiler window. At the top, there's a menu bar (File, Edit, Debug, Window, Help) and a toolbar with navigation icons. Below the toolbar, there's a "Start Profiling" button and a "Run this code:" dropdown menu. On the right, it says "Profile time: 64 s".

The main content area displays the profile results for the function **AffineBreakpoints** (1000 calls, 13.913 sec). It includes a timestamp "Generated 15-Jan-2014 17:14:19 using cpu time." and a link to the function file: [/Users/pfrazier/work/code/matlabKG/matlabKG-R2013a/matlabKG/ckg/AffineBreakpoints.m](#). There's also a link to "Copy to new window for comparing multiple runs".

Below this, there's a "Refresh" button and a set of checkboxes for various options: "Show parent functions", "Show busy lines", "Show child functions", "Show Code Analyzer results", "Show file coverage", and "Show function listing". All checkboxes are checked.

The "Parents (calling functions)" section shows a table with one entry: [LogEmaxAffine](#) (function) with 1000 calls.

The "Lines where the most time was spent" section shows a table with the following data:

Line Number	Code	Calls	Total Time	% Time	Time Plot
<a href="#">44</a>	<code>c(1+j) = (a(j) - a(i+1))/(b(i+...</code>	55118	8.322 s	59.8%	
<a href="#">43</a>	<code>j = A(Alen); % jindex = Alen</code>	55118	0.953 s	6.8%	
<a href="#">58</a>	<code>end</code>	28059	0.699 s	5.0%	
<a href="#">50</a>	<code>Alen = Alen-1; % Remove last e...</code>	27059	0.699 s	5.0%	
<a href="#">53</a>	<code>break % quit while(1) loop</code>	28059	0.635 s	4.6%	
All other lines			2.605 s	18.7%	
Totals			13.913 s	100%	

Below the table, there's a section for "Child (called functions)" which is partially visible.

# If your Matlab program is too slow...

- Profile it first to find out what parts are causing speed issues.
- Replace loops by operations on matrices and vectors.
- Try the matlab compiler.
- Can you use more efficient algorithms?
- Rewrite the slow parts of your code, or all of your code, in C.

# Here's how to profile in C and C++

- use the -pg and -g flags when compiling your source code into object files, and when linking your object files into an executable.
- When you run your program, it will create a file in the current directory called “gmon.out”
- Then run “gprof” with your executable's name as the first argument.

# Here's an example in C

## First, here's my test program

```
[pf98@asimov]~/profiling% cat myprog.c
#include <stdio.h>
#include <math.h>

double my_function1(double x) { return x+2; }
double my_function2(double x) { return log(x); }

int main() {
    int i;
    double x=0;
    for(i=0;i<100000000;i++) {
        x=my_function1(x);
        x=my_function2(x);
        if (i<20)
            printf("%f\n", x);
    }
    printf("%f\n", x);
}
```

# Here's an example in C

compiling with —  
profiling on using  
-g and -pg (and  
also with the  
math library -lm)

run gprof to get  
profiling output

```
[pf98@asimov]~/profiling% gcc -o myprog myprog.c -g -pg -lm
[pf98@asimov]~/profiling% ./myprog
0.693147
0.990710
1.095511
1.129953
1.141018
1.144547
1.145670
1.146027
1.146140
1.146176
1.146188
1.146192
1.146193
1.146193
1.146193
1.146193
1.146193
1.146193
1.146193
1.146193
1.146193
[pf98@asimov]~/profiling% gprof myprog
Flat profile:
```

running my program

my program computes some numbers

# here's the profiling output

```

Each sample counts as 0.01 seconds.
  %   cumulative    self         self         total
time seconds  seconds   calls   ns/call   ns/call   name
44.74      0.31      0.31 100000000      3.13      3.13  my_function2
28.87      0.52      0.20                main
23.09      0.68      0.16 100000000      1.62      1.62  my_function1
  4.33      0.71      0.03                frame_dummy

```

# If your C or C++ program is too slow...

- understand whether your code is CPU or I/O bound
- improve efficiency and use better algorithms in portions of the code identified as bottlenecks by profiling
- don't forget to turn off profiling once you no longer need it --- it slows down the code
- compile with the -O3 flag to tell the compiler to aggressively optimize for speed.
- rewrite your code to run in parallel on more cores, or on a GPU
- understand CPU cache, and design your code to reduce cache misses

# Understand a bit about computer architecture

- Your computer has different ways of storing data, each with its own speed & size
  - Level 1 cache on the CPU: ~64KB, 0.5ns
  - Level 2 and higher caches, on and near the CPU: ~2MB, 5ns
  - RAM: ~4GB, 100ns
  - SSD Hard Disk: ~100GB, 150,000ns=150us
  - Traditional Hard Disk: ~1TB, 10,000,000ns=10ms

# The time to access data affects your program's speed

- If you write your program to get more of its data from the CPU cache, and less from RAM, it will be faster.
- If you write your program to get more of its data from the RAM, and less from disk, it will be faster.



# The time to access data affects your program's speed

- The operating system automatically figures out when & where a particular memory address will be stored (except for SSD vs. traditional hard disks).
- If you are only accessing a few KB of data at a time, this will mostly be served from CPU cache and will be fast.
- If you are jumping around within many MB of data, this will be slower as it will need to be served from RAM.
- If you are jumping around within many GB of data, this will need to be served from disk, even if you don't explicitly store it on disk in your code.

# Use performance monitoring tools to check which resources are being used

- On UNIX and Mac OSX, I use top and sar. vmstat is also useful.
- On Windows, I use performance monitor.

# Learn UNIX

so you can leverage high-performance computing,  
and so you can be more productive with your Mac

- Parallel and high-performance computing environments (e.g., through Amazon Web Services) let you do things you can't do otherwise. These environments usually run UNIX.
- Mac OS X is a flavor of UNIX. Knowing UNIX makes Mac users much more productive (through Terminal).
- Being a command-line wiz saves time, but takes time to learn.
- At Cornell there is a short course in UNIX scripting that students can take.
- Perhaps there is something similar here.

# Don't be afraid to rewrite code

- Write code in a high-level language first, then rewrite in a low-level language later if needed for speed.
- Your code will eventually grow unwieldy. Depending on circumstances, it may be worth taking the time to refactor and redesign.

# Further Reading

- Wilson et al. 2013, “Best Practices for Scientific Computing”,  
<http://arxiv.org/pdf/1210.0530v4.pdf>

# Part 2: Some things I've learned at Uber about doing useful research

- Impact is a choice
- Ease of implementation matters
- Communicate well
- Theory is not convincing
- Own your research

# Ease of Implementation Matters

- At Uber (and everywhere), some algorithmic changes are easy to make. Others are hard to make.
- Easy changes get made quickly.  
Hard changes get made slowly, or not at all.
- The difficulty depends on how well the change fits with the engineering infrastructure, product roadmap, and organizational structure.
- Successful data scientists propose algorithmic changes that are algorithmically good and easy to implement.
- If you work with a company, they will be more likely to implement your ideas if they are easy to implement.

# Communicate Well

- At Uber, data scientists must communicate with other data scientists, software engineers, and product managers.
- Having an impact requires:
  - understanding from others what problems are important.
  - understanding from others what changes will be easy.
  - convincing people that your solutions are good.
  - coordinating about who will do the work.
- Engineers and product managers think differently from data scientists and it is easy to be bad at communicating with them.



# Communication Matters

- These forms of communication matter:
  - Listening.
  - Speaking in conversations, meetings, and presentations.
  - Slide decks.
  - Writing in emails and technical documents.
  - Being empathetic.

# Theory is not convincing

- In academia, we prove guarantees for algorithms to demonstrate that they are good.
- At Uber, theoretical guarantees rarely influence decisions.
- Instead, we test algorithms in simulation and in live experiments.
- Testing a solution in the real world gives a pretty good indication of whether it will work in the real world.

# Theory is not convincing

- Most theoretical guarantees are actually not very useful for making decisions:
  - Asymptotics are only useful when the system studied is asymptotic.
  - Finite-time bounds are only useful if the bound is tight enough to be meaningful.
  - Worst-case bounds for models that only approximate reality don't provide much reassurance about the actual worst-case performance.
  - Making an algorithm worse empirically so that it is easier to prove a performance guarantee about it does not make it likely to be adopted in the real world.
- We do use tight finite-time performance bounds to make decisions about allocating effort to projects.
- Theory is also useful if it helps us understand an algorithm.
- I imagine that for some problems (e.g., Google AdWords) that are easy to model, worst-case performance guarantees may be comforting.

# Theory is not convincing (But theory still matters)

- Here are some reasonable reasons to prove a theorem:
  - Because it is fun.
  - To get a paper published.
  - To graduate.
  - Because it will help you understand the system you are studying, or mathematics itself.
  - Because the theorem you are proving could actually inform a decision that you or someone else will make down the road.

# Thinking about projects' impact before starting them will help you choose projects that matter (Impact is a choice)

- Let me tell you three sad stories:
  - The story of a PhD student who developed a better method for solving an optimization problem, without thinking about whether anyone cared about that optimization problem.
  - The story of a PhD student who developed a queuing method robust to errors in estimating a particular parameter, without thinking about whether that parameter was easier or harder to estimate than the other inputs required by the method.
  - The story of a data scientist solving the problem that best leveraged his skills, instead of the one that was most important to the business.
- And one happy story:
  - The story of another data scientist listing 5 projects and thinking about their business impact before choosing which one to work on.

# Own your research

- Be an owner, not a renter.
- Own your research.
- You are ultimately responsible for the decisions you make.
- If you see a problem, fix it.  
Don't think, "My advisor didn't complain, so it must be ok."
- If you want your research to be useful, then:  
Go learn about the real world,  
how people do things,  
what practitioners know how to do and don't know how to do,  
what is broken.  
Ask, "What do you worry about?"  
Figure out what questions are relevant,  
and try to answer some of them.

Thank you!  
Any Questions?