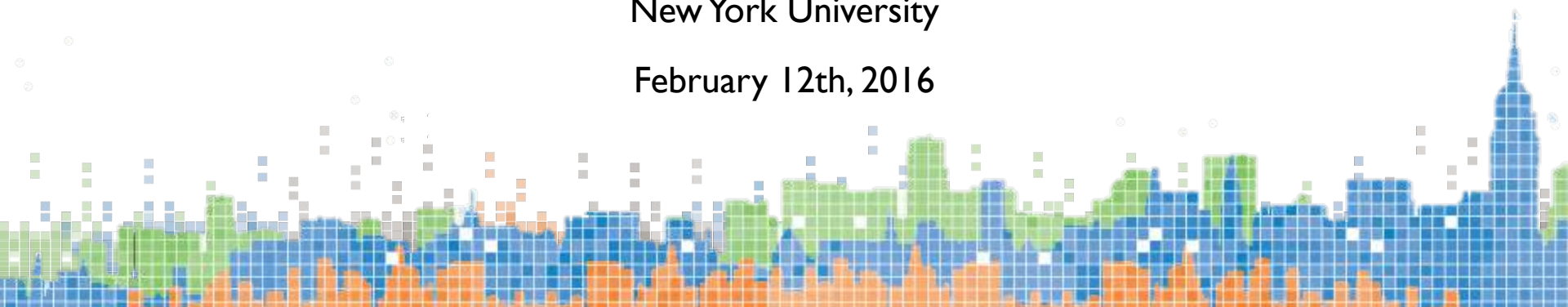# How Not to Lose Your Code, Your Degree, and Your Future Job

Justin Salamon

Center for Urban Science and Progress
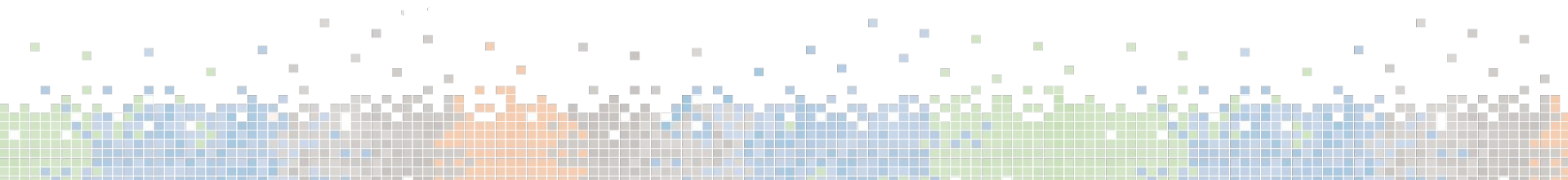
New York University

February 12th, 2016

# How Not To Lose Your Code

(or: How To Write Good Code, Manage It, Share It, and Be Awesome)

# Horror Story 1

Bob: Oh no! Someone broke into our apartment and stole my laptop! All of my code is gone…

Alice: I thought you kept a backup of your code on an external hard drive, no?

Bob: They took that too…

# Horror Story 2

Bob: dammit, I've made some changes to my code and now it doesn't work anymore

Alice: that's ok Bob. Just revert to a version of the code you know worked.

Bob: ehm… version?

# Horror Story 3

Bob: I was collaborating with Rob on a project, and we were emailing each other the latest version of the code. Over spring break we both worked separately, and now we don't know how to merge our work!

Alice: that's ok Bob. Just use Git to compare the two repositories, identify conflicting sections of the code and help you merge them back into one.

Bob: ehm… Git?

# Horror Story 4
## (true!)

Ex-master student: Hey Justin, I was reading your paper from a few years ago , and I was wondering how exactly you normalize the data in Section 5?

Justin: Oh god… let me find my PhD backup drive…

[finds drive]

Justin: Oh god… let me find the code used for that paper on the drive…

[finds code]

Justin: Oh god… let me figure out which version of the code w   d when the paper was written…

[after much suffering, figures out which version of the code and how normalization was applied. Fail averted. Luck

[kicks himself for not using version control at the tin

# Version Control

Software to help keep track of changes made to files:

- Tracks the **history** of your work

- Helps you **collaborate** with others

- Helps keep an **online backup** of your code

# Keeping Track of History

- How do you get back to that **working version** you had yesterday?

- How do you get from "it's not working" to **understanding what went wrong**?

- How would you **repeat** the experiments from that paper (or project) you worked on last year?

# Collaborating

With yourself:

- Need to run same code on laptop and on CUSP server:
  - How do you get the code onto both?
  - How do **verify** that you have the same code on both?

With others:

- Working on a project with classmates:
  - How do you find out when they changes something?
  - How do you merge your changes without making a mess?
  - How can you find out which of you introduced a bug, and when?

# Version Control

A version control system:

- Records you files' **history**

- Shows the **differences** between versions

- Handles **sync** between copies on different computes and **the cloud**

# How is that Different from Dropbox?

Guarantees consistency:

- Code needs to be **exactly as written** across all files
- Files changed together can be **updated together**
- Changes to a file by different people must be **merged correctly**

Keeps records:

- Publish and replicate history reliably
- Interrogate past **changes** and find out what version you are looking at

# Version Control Evolution

- Revision Control System (RCS): 1982
  - Each file versioned independently, manual sync ☹

- Concurrent Versions System (CVS): 1986, 1990
  - Multiple file versioning ☺
  - Can only change latest version, clumsy networking, poor support for binary files ☹

- Subversion (SVN): 2000
  - Similar to CVS but with many improvements ☺
  - Versioning done on server-side: local dev is tricky ☹
  - Single point of failure (the server) ☹

# Git
## [Torvalds, 2005]

- Features:
  - **Distributed** version control system (DVCS)
  - Does not require a centralized server
    - But you can still have one, if you want
  - Easy local incremental development
  - No single point of failure
    - Every developer has local copy with full history of the repository
- Drawbacks:
  - A bit of a learning curve…
- Other DVCs
  - Mercurial (`hg`)
  - Bazaar (`bzr`)

# Using Git: client-server

1. `git clone`      Make local copy of the repo

2.    (edit files)

3. `git commit`      Register your changes locally

4. `git push`      Share changes upstream

5. `git pull`      Get updates from upstream

# Demo

# Advanced Usage: Tags

- Some revisions are special:
  - Initial paper submission
  - Camera ready submission
  - Public software releases
- Tagging links **semantic versions** to **revisions**
- Example:
  - `git tag —a v1.0`
  - `git push origin --tags`

# Advanced Usage: Branches

- What if you want to develop new features, but retain version control on a stable codebase?

- Working in a **branch** of the source tree

- Merge back when you're ready

- Especially useful for collaborations

# Advanced Usage: Branches

- Example: create a new branch
  - `git checkout —b unstable`
  - `(edits, commits, pushes)`
- Switch to master, bug fix, switch back
  - `git checkout master`
  - `(edits, commits, pushes)`
  - `git checkout unstable`
- Merge <span style="color:red">unstable</span> back into master
  - `git checkout master`
  - `git merge unstable`

unstable

# Advanced Usage: Branches



| Graph | Description | Commit |
|---|---|---|
| | implemented a first cut at mirex tempo eval | 84e414e |
| | `⑂ upstream/master` `⑂ origin/master` `⑂ origin/HEAD` Adding web service link to readme | 72bb77a |
| | `⑂ master` `1 behind` Merge branch 'master' of https://github.com/craffel/mir_eval | 15b857c |
| | Formatting issue | c7b156b |
| | Setting version number to 0.1 | 7f88b75 |
| | Post manual merge -- io.load mods | e9ac436 |
| | Merge branch 'master' of https://github.com/craffel/mir_eval | 1446977 |
| | Fixed join bass note bug | ae747d3 |
| | Merge branch 'master' of https://github.com/craffel/mir_eval | 5cedb94 |
| | Manual merge success. | c5fc21f |
| | Added newline catch to load_intervals to keep moving past erroneous gaps in the file (occurs in Billboard dataset?) | ed43f65 |
| | Added additional quality support (1, 5, X); fixed bug in minor6 interval definition. | 189fee9 |
| | `⑂ upstream/gh-pages` `⑂ origin/gh-pages` Adding API service link | 5b34a32 |
| | Upping version number | ec64f7b |
| | New docs for 0.1 | d9793f9 |
| | Merging | abd4a25 |
| | `🏷0.1` Using README as long_description | 6c7e390 |
| | Using rst | 160d0e4 |
| | Switching from md to rst | ca464c6 |
| | Changing to MIT license | f2a1525 |
| | Files for pypi | a5d1391 |
| | Merge branch 'master' of https://github.com/craffel/mir_eval | abedc4c |
| | Explicitly handle infinite SIR and silent sources for #91 | 716d5ad |
| | Resolves #93 | 65cee31 |
| | Not just a collection of scripts... | b2e8c11 |
| | Adding ref info | a1d1fdb |
| | Typo, also making actual code snippet | c58ecf6 |
| | Adding note for anaconda | 97d2f34 |

# Hosting Git Online

# GitHub

- Free hosting for open source projects
  - Free organization accounts for academics
- Social network integration
- **Extra usability tools**
  - User management
  - Pull requests
  - Issue tracking, comments, wiki
  - Release management
  - Webhooks & services!
- **Portfolio for potential employers!**
  - More on this later

# How Not To Lose Your Code : GitHub

# How Not To Lose Your Code : GitHub

# How Not To Lose Your Code : GitHub

# How Not To Lose Your Code : GitHub

# How Not To Lose Your Code :  GitHub

# What About Private Repositories?

# What About Sensitive Code/Data?



## CUSP Data Facility can host Git repo's!

# What if I'm not sure?



## CUSP Data Facility

# Git In Practice
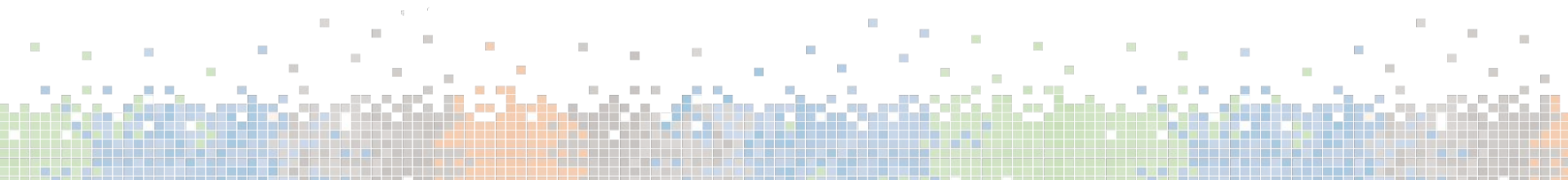
# What if I'm More of a Visual Type?

# GitHub: Example Workflow

- Pull from GitHub
  - Either `develop` or `master` branch, depends…
- Develop locally
  - First on ipython notebook
  - Then on versioned source
  - <span style="color:red">Run unit tests</span>
  - Commit
  - Keep editing, pulling changes from collaborators
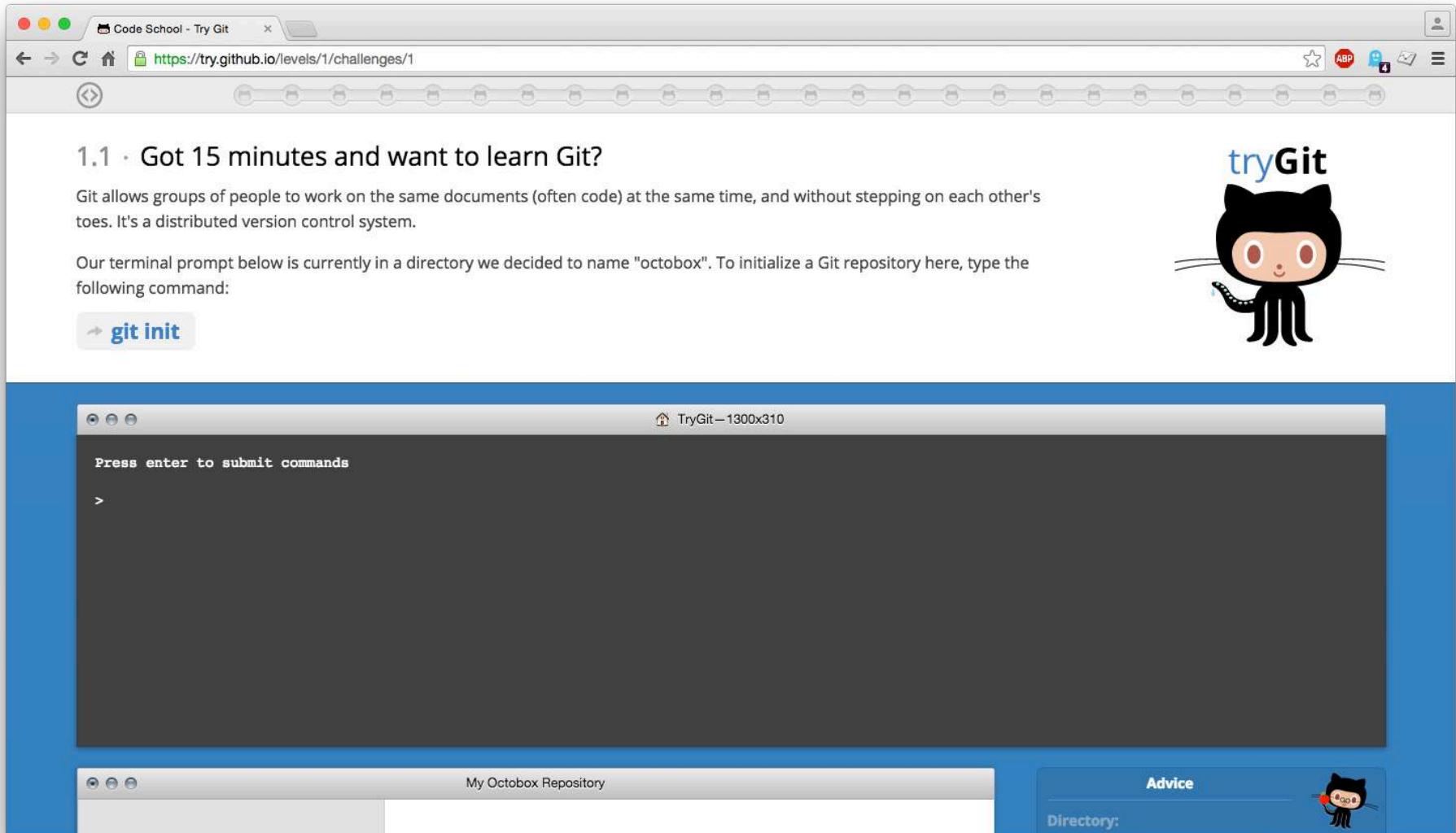- When it's ready:
  - Push back to GitHub

# Typical Version Control Questions

- When should I start using version control for my project?

- Which files should I track in the repository?

- How often should I commit?

- How often should I push changes to a shared repo?

# Gitting Started

https://try.github.io/

# Gitting Started

https://datahub.cusp.nyu.edu/computing.html

# Advanced GitHub

- **Webhooks & services** allow integration of 3$^{rd}$ party "apps" into your GitHub repo:
  - **Travis CI**: continuous integration, runs all unit tests on every pull request, for multiple builds (e.g. python 2.7, 3.4, 3.5)
  - **ReviewNinja**: code review, must get ninja star from someone who's reviewed your code before you can merge PR
  - **Coveralls**: shows percentage of code covered by unit tests, highlights code not covered by any test.
  - Many more…

# Advanced GitHub

# Publishing Code

- Make sure it has a license!
  - BSD/MIT-style is a good choice for research code
  - GPL for complete applications or code with possible commercial value
  - Ensure license is at least described in a README file

- Make it citable, get a DOI
  - Zenodo.org
  - https://guides.github.com/activities/citable-code/

- Code implementing research?
  - Tell users what they should cite if they use it

# 3 Things To Do Tomorrow

1. Get your current research code into a version control repository & push it to a hosting site (can be private)

2. Pull it onto another computer, get it to build and run

3. Open source? Choose a license!
   - Your code includes someone else's code? Make sure the licenses are compatible!

# Unit Testing

Unit testing is awesome.

We don't have time to cover it.

If you're not familiar with it, look it up.

# Unit Testing: What is it?

- A "unit test" is a bit of code that calls one of your functions, gives it some input, and tells you whether it returned the right result

- Write a set of these, and you have a "test suite"

- A "test framework" can help you write them more quickly; there's at least one for every programming language and environment

- Should be set up so you can run all tests in one go

# Unit Testing: What is it for?

An automated way of ensuring:

- That your code's API works

- That the individual parts of your code work correctly

- That you don't break your code when changing it

- Also useful when developing a tricky algorithm (test-driven development)

# Unit Testing Questions

How do I write tests when I don't know what results to expect?

- Break it down into functions whose behaviour you can predict

- Test individual components, not the whole thing

- Testable code is also more readable code (and so more reviewable code, and…)

Unit testing is about trying to ensure that *the code implements the method*—not that *the method is the right one*
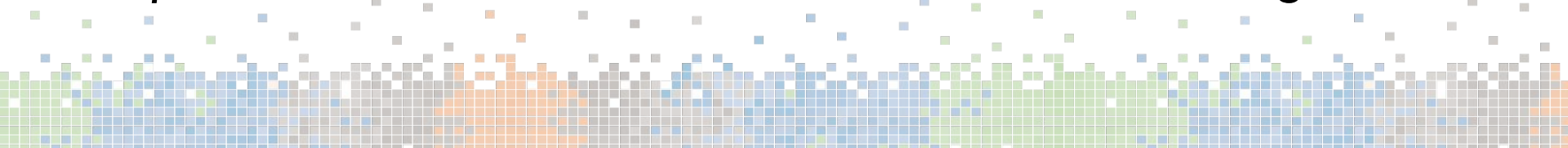
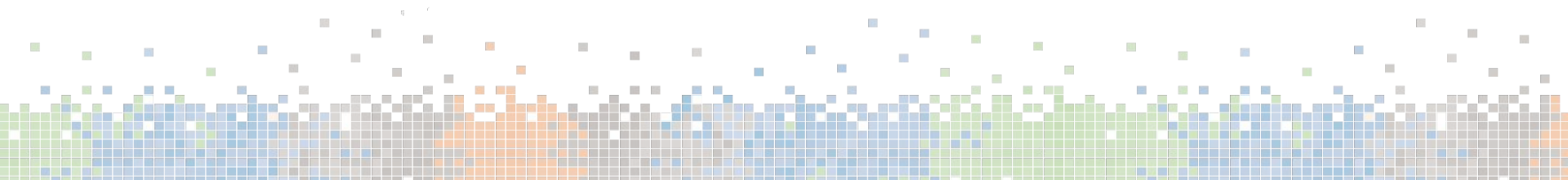# Unit Testing Questions

What sort of test data and test cases should I write?
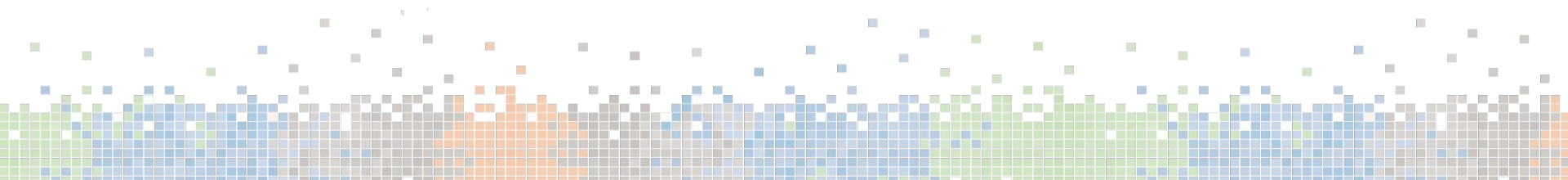
* The simplest possible ones!

But I have big data sets and complex results!

* Don't use real-world data: that's a different kind of test

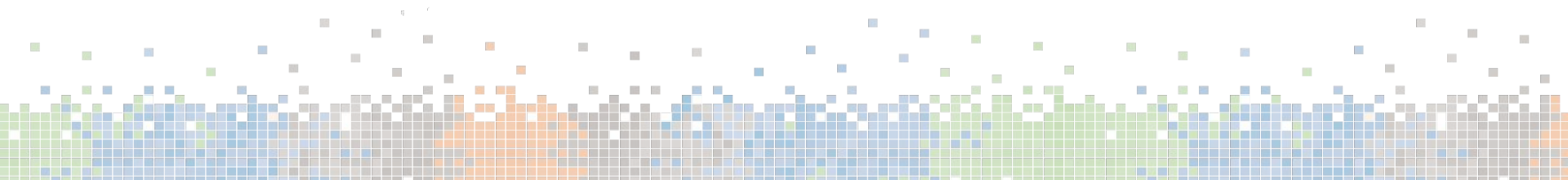* Look for the smallest possible input to test a given behaviour
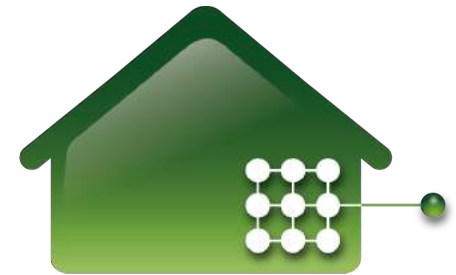
# Unit Testing Questions

Example…

# How Not To Lose Your Degree

(or: How To Backup Your Data, Share It, and Be Awesome)
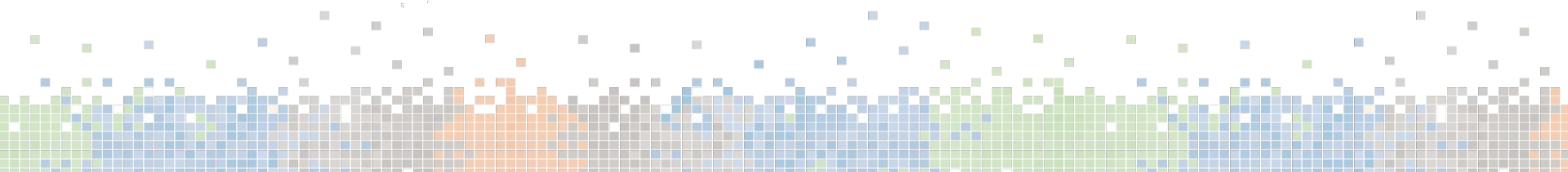
# A Show of Hands

# Horror Story 1

*Hi, a friend of mine just overwrote two months of her PhD thesis with an older version. I know recovery of overwritten data is possible, but wonder if I'd need special hardware to do it. Does anyone know something about this ?*

*Thank You.*

5 October 2005 Linux Forums - http://tinyurl.com/8t7uaop

WORKING COPY IS NOT ENOUGH

KEEP BACKUPS

# Horror Story 2

*A tiny television sits where a big screen used to, and a Christmas tree stands with little underneath it...*

*Even worse than the gifts, the crooks stole a MacBook Pro laptop and a LaCie hard drive.*

*The hard drive had ... her dissertation and nearly seven years of research for her doctoral degree she was set to finish in a few weeks.*

*Osuna had everything backed up on a separate hard drive in a safe, but burglars made off with that too.*

*"All I could think about is that all that time is gone, all that effort, everything is gone," Osuna said.*

*22 December 2010 KRQE - http://tinyurl.com/9a5j56f*

LOCAL COPY IS NOT ENOUGH

BACKUP TO THE CLOUD

# Horror Story 3

*...her car was broken into and her chrome Mac book pro was stolen.*

*She has a back-up for all but the last six months of research, but the most important part of the research had happened recently.*
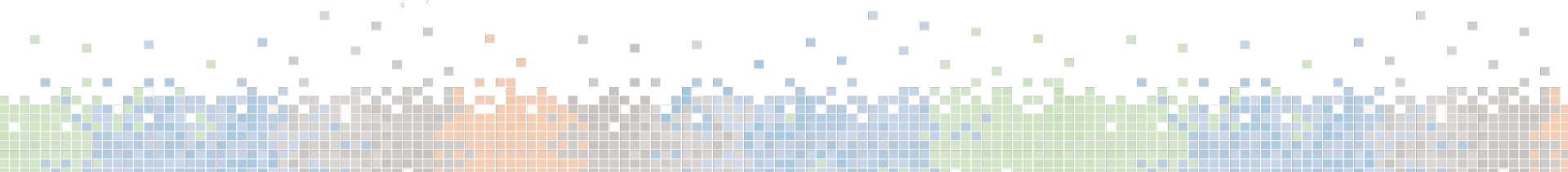
NBC4 January 06 2011 - http://tinyurl.com/92jf2lr

**MANUAL BACKUP IS NOT ENOUGH**

**SCHEDULE YOUR BACKUPS**

# Archiving Your Data

- Project (or course, or degree) is over, now what?
- Archive!
  - Allow follow-on research
  - Allow validation of your results

- Beware!
  - Don't use obscure formats
  - Don't use obscure media
    - BBC Domesday Project [1986] used laserdisc!
  - Don't rely on technology being available
  - Keep original source material

# Archiving Your Data

When someone looks at your data, will they understand:

- Why you created it?

- What the data is useful for?

- What column 27 in table 15 actually means?

- What are the UNITS?! (Hertz? Meters? Miles?)

- How the data was created?

- What the source data was on which this data is based?

**THERE'S NO DATA WITHOUT METADATA!**

# Archiving Your Data



# CUSP Data Facility

# Archiving Your Data

https://datahub.cusp.nyu.edu/services.html

# Archiving Your Data

https://datahub.cusp.nyu.edu/services.html

# How Not To Lose Your Degree: Publishing Data

"Many researchers believe that if scientists set out to reproduce preclinical work published over the past decade, a majority would fail. This, in short, is the reproducibility crisis."



THE CHRONICLE OF HIGHER EDUCATION
March 27, 2015

Log In | Events | Store

Subscribe Today

Home | **News** | Global | Opinion & Ideas | Facts & Figures | Blogs | Advice | Forums | Jobs

Search The Chronicle | Go

## Research

March 16, 2015

### Amid a Sea of False Findings, the NIH Tries Reform

David Banks, Bloomberg

Science needs to get its house in order, says Francis Collins, director of the NIH. "We can't afford to waste resources and produce nonreproducible conclusions."

**Most Popular**

Most Viewed | Most Commented

1. Is 'Design Thinking' the New Liberal Arts?
2. You're Distracted. This Professor Can Help.
3. Video: 'The Athletic Machine Is in Charge of the University'
4. Video: How an Elite Women's College Lost Its Base and Found Its Mission
5. Social Networks for Academics Proliferate, Despite Some Scholars' Doubts

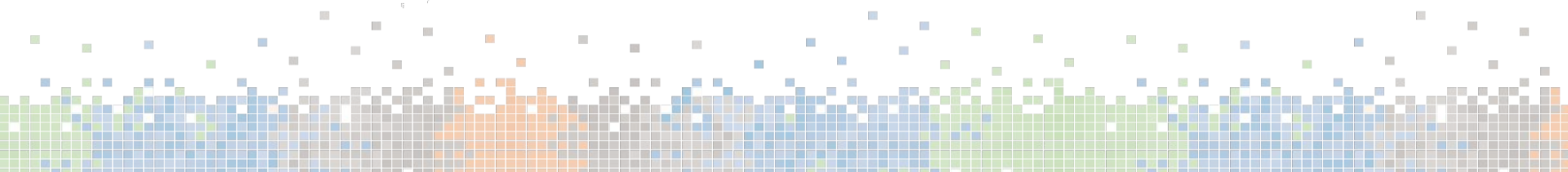Jobs on **Vitae**          Search 7,904 opportunities

**Browse by Position Type**

All Types | Faculty/Research | Administrative | Executive | Jobs Outside Academe
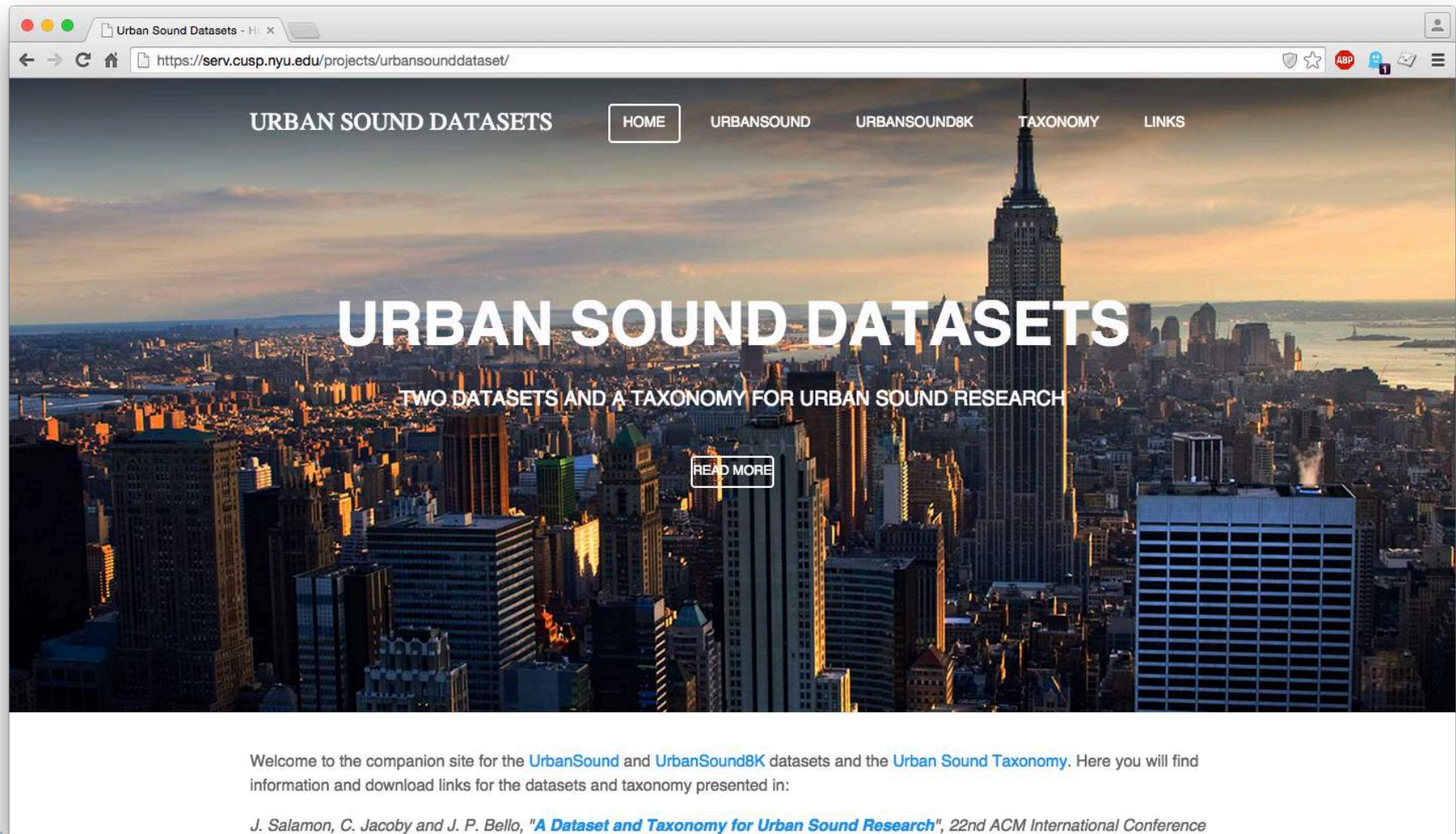
**Search by Keyword**

# Publishing Your Data

Why publish my data?

- Help others to:
  - Validate your research
  - Validate their implementation of your algorithm
  - Advance the state of the art
  - Combine with other datasets

- So what's in it for me?
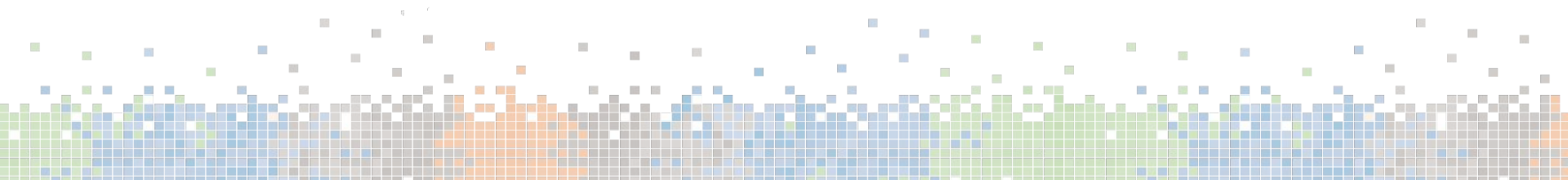  - Citations!
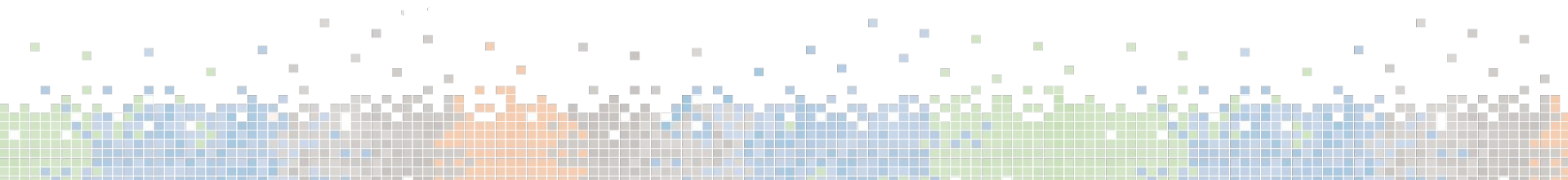  - Impact! (job)

# Publishing Your Data
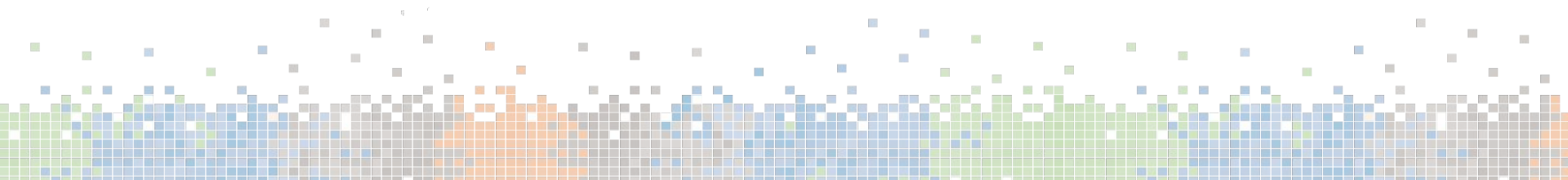
# Publishing Your Data

- Make sure you are allowed to publish!
  - Who owns the data?
    - NYU / Industry partner / Funding body / Creative Commons
  - Adhere to CUSP policies
    - Privacy: are you allowed to publish the data?
    - Publication: are you expected to publish the data?
    - Repositories: where should you publish the data?
    - Licenses: who should be allowed to access the data?

- Make sure it has a license!
  - Example: Creative Commons
    - CC0? CC Attribution Non-Commercial? Something else?

# How Not To Lose Your Future Job

(or: How To Maximize The Impact Of Your Work and Be Awesome)

# Get Found

# Git Found?

- Employers/recruiters increasingly searching/asking for GitHub profiles

- But wait… there's more…

# Git Found?

# Git Found?
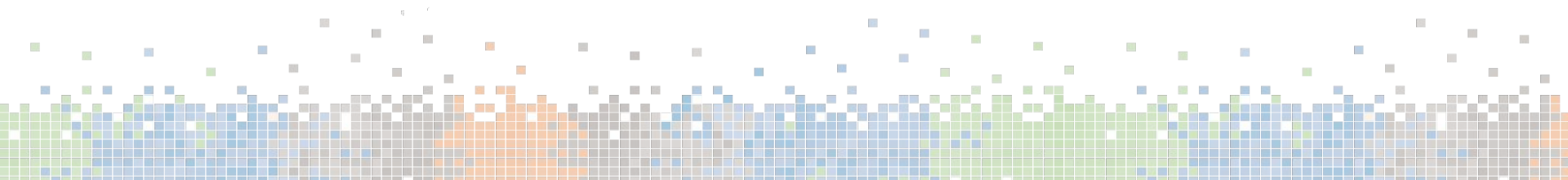
So is GitHub my new CV?

# Git Found?

So is

# Git Found?

- Even a stellar GitHub account is not a replacement for:
  - Your resume
  - Interview skills
  - Presentation skills


- But… it can help you get the interview in the first place!
- Using Git (not specifically GitHub) is a basic skill every programmer should have nowadays
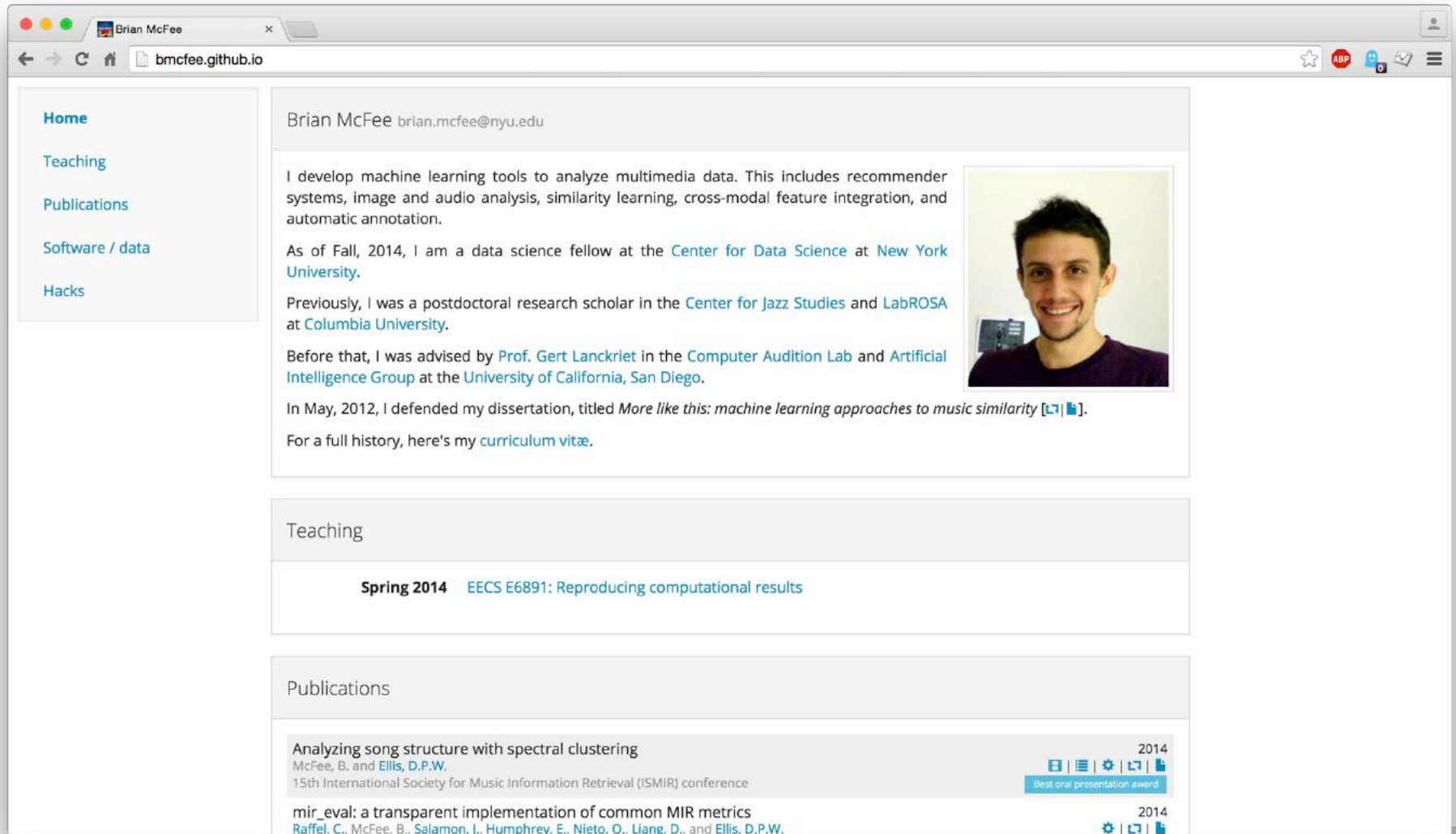
# Take Control of Your Online Presence

- LinkedIn, GitHub, etc.
  - Limited to displaying a specific type of information

- A personal website is where you can really showcase your best!
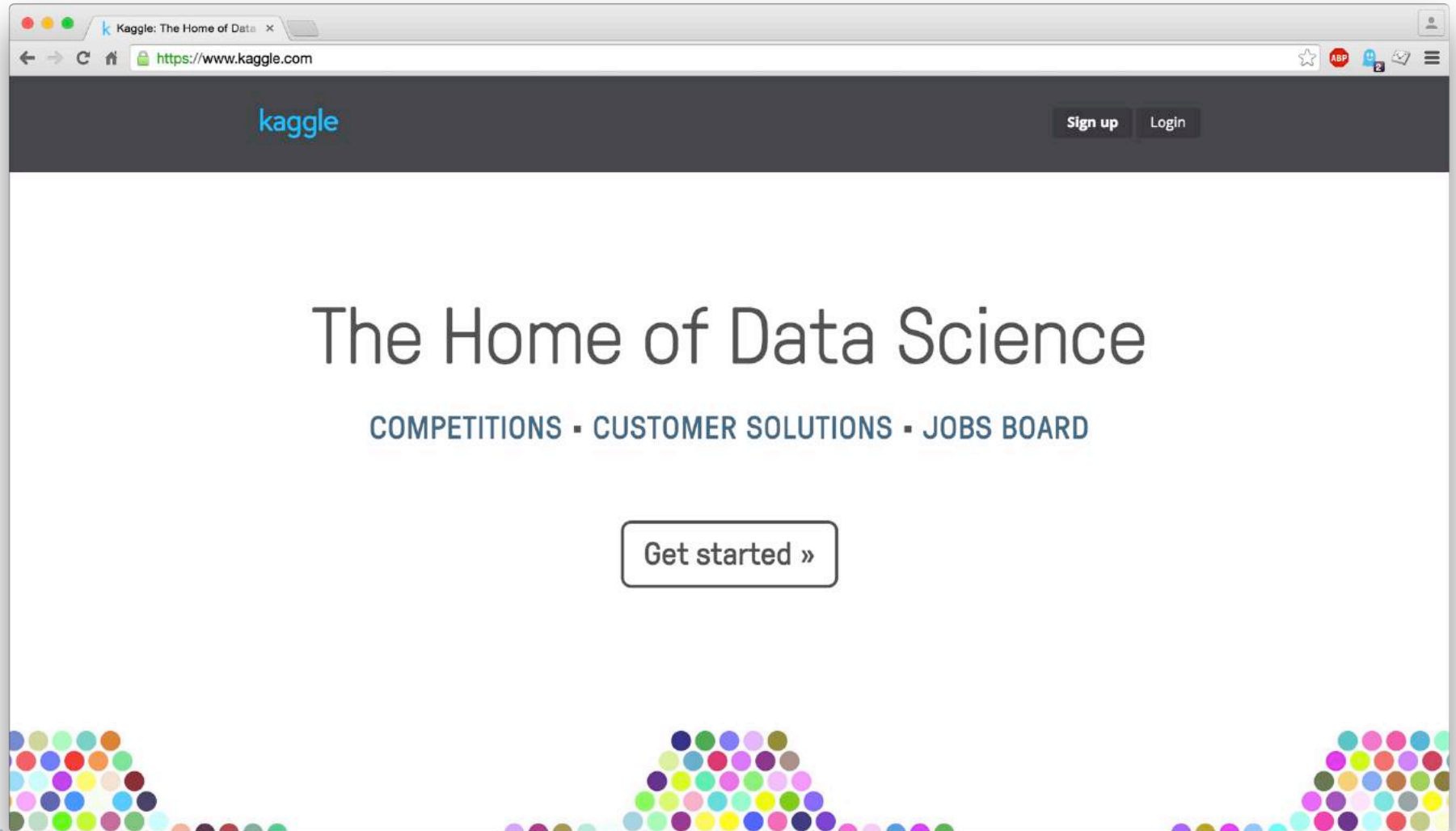
# Yes

# No

# Personal Website

- Showcase your best

- If an employer searches for you (online), have control over the first thing they see

- An interesting blog-post (or ipython notebook!) can be a great source of traffic to your site!
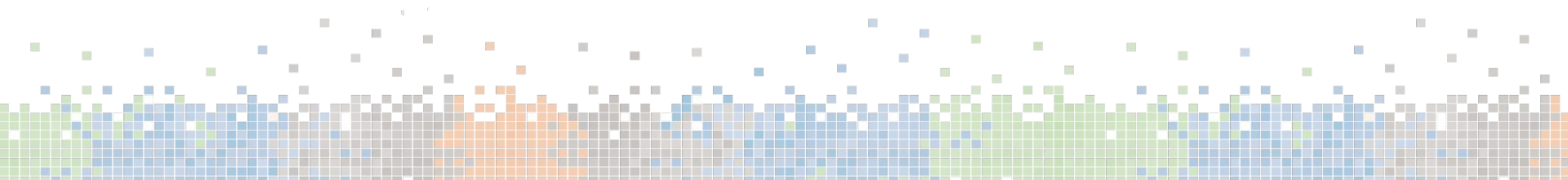
# Kaggle

# Kaggle

- Excellent experience…
- But will it land me a job?
    - If you win…

# Credits

Parts of this presentation (version control, GitHub, unit testing, data publishing/archiving and some horror stories) taken from:

- Chris Cannam et. al (Queen Mary, University of London): ISMIR 2012 tutorial on Reusable software and reproducibility in music informatics research
  - http://soundsoftware.ac.uk/videos
- Dan Ellis & Brian McFee (Columbia): Version Control & Github:
  - http://www.ee.columbia.edu/~dpwe/e6891/outline.html

# RECAP

- Code
  - Use version control (probably Git)
  - Write unit tests
  - Use online hosting (e.g. GitHub / BitBucket)
- Data
  - Back it up!
  - Publish it / archive it
- Online presence
  - Great way to get attention from employers
    - (but not a replacement for the classic skill-set)
  - Be in control of your online presence!

Q&A Time