

Software Archaeology

17-313 Spring 2025

Foundations of Software Engineering

<https://cmu-313.github.io>

Michael Hilton, Austin Henley, and Nadia Nahar

Administrivia

- Part (b) is due Thursday, Jan 23rd, 11:59pm.
- If you haven't: **PLEASE FILL OUT TEAMWORK SURVEY!**
- Get started early, ask for help, and check the #technical-support channel; chances are your questions have been asked by others!

Smoking Section

- Last full row

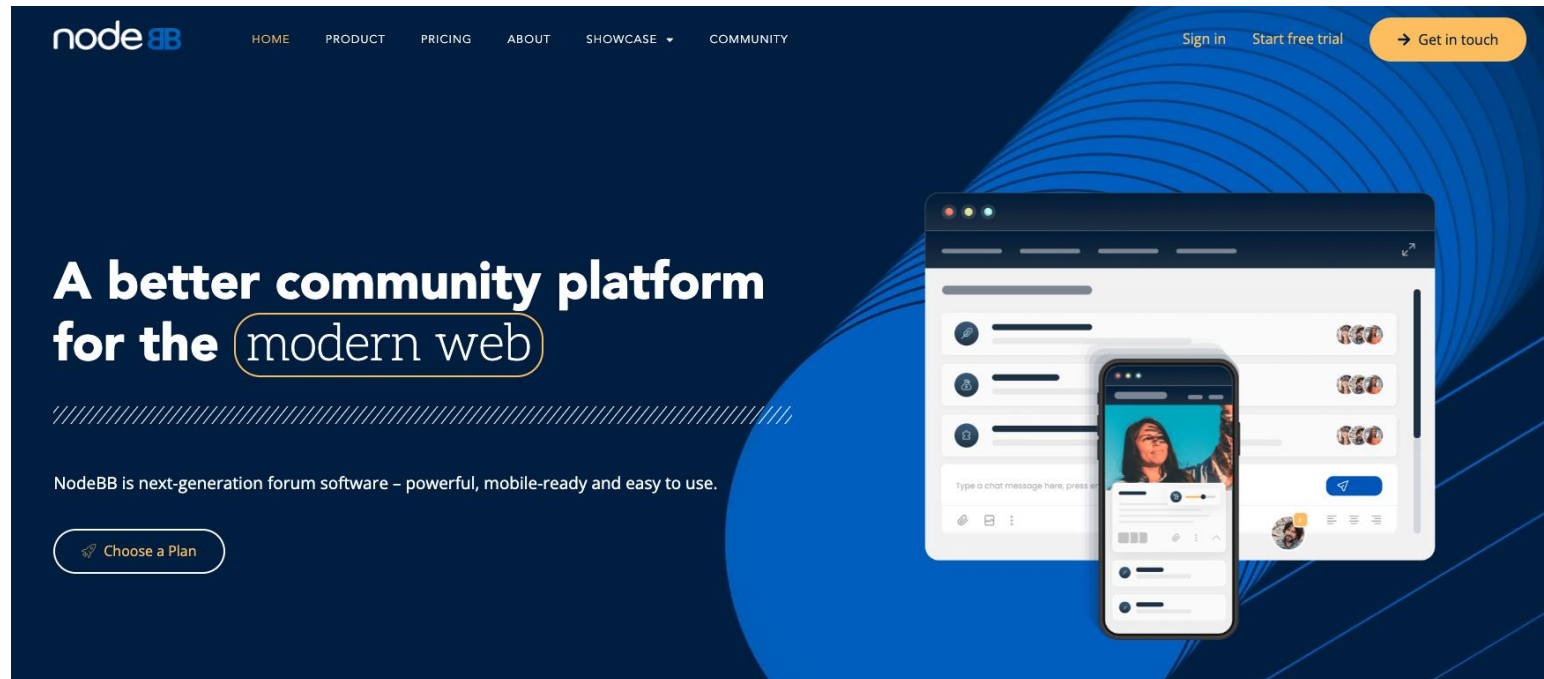


Learning Goals

- Understand and scope the task of taking on and understanding a new and complex piece of existing software
- Appreciate the importance of configuring an effective IDE
- Contrast different types of code execution environments including local, remote, application, and libraries
- Enumerate both static and dynamic strategies for understanding and modifying a new codebase

Context: big ole pile of code

- ... do something with it!

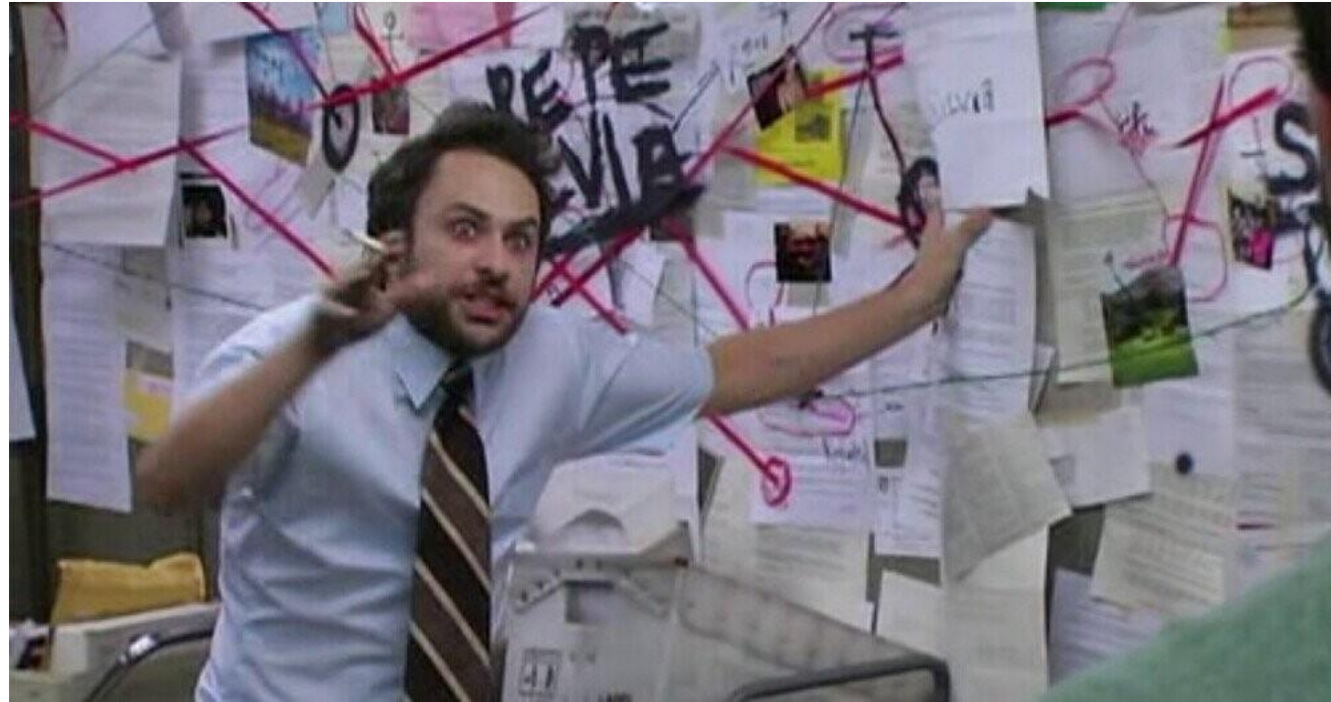


Participation Activity—Part 1

- Take out a piece of paper (or ask for one).
- Write down the **challenges you've faced trying to understand someone else's code.**
- Pair with your neighbor and discuss your answers. Do you agree?
- Share with the class!
- Write your own andrewID on the paper, leave it at the end of class.

**You will never
understand the
entire system!**

Challenge: How do I tackle this codebase?



Participation Activity—Part 2

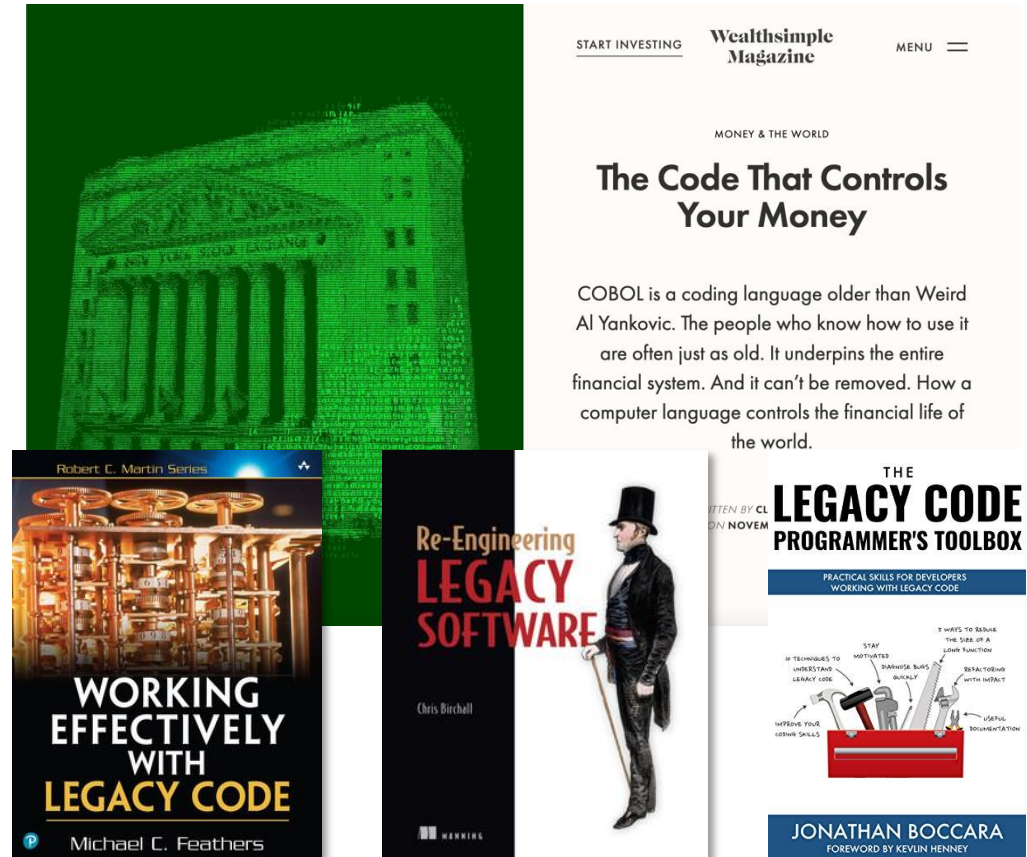
- Write down **strategies to understand a large codebase that is unfamiliar to you.**

Challenge: How do I tackle this codebase?

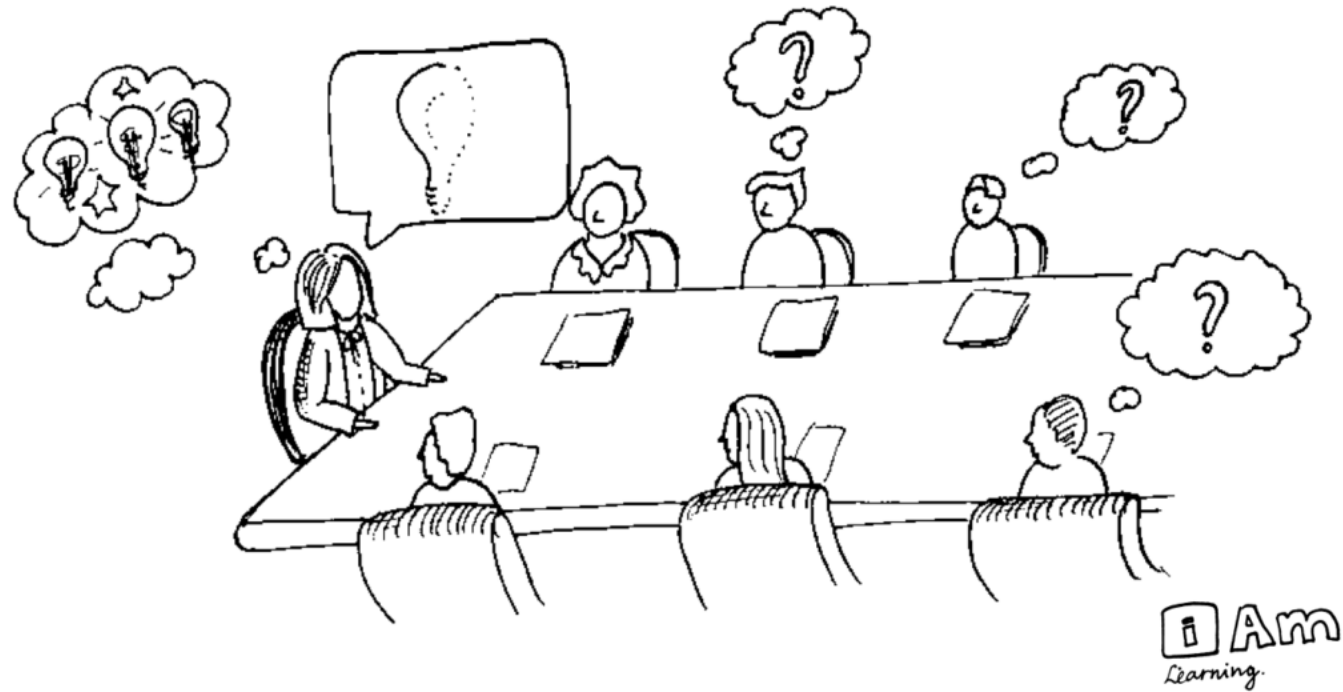
- Leverage your previous experiences (languages, technologies, patterns)
- Consult documentation, whitepapers
- Talk to experts, code owners
- Follow best practices to build a working model of the system

Bad news: There are few helpful resources!

- **Working Effectively with Legacy Code.**
Michael C. Feathers. 2004.
- **Re-Engineering Legacy Software.**
Chris Birchall. 2016.
- **The Legacy Code Programmer's Toolbox.**
Jonathan Boccara. 2019.



Why? Because of Tacit Knowledge



Today: How to tackle codebases

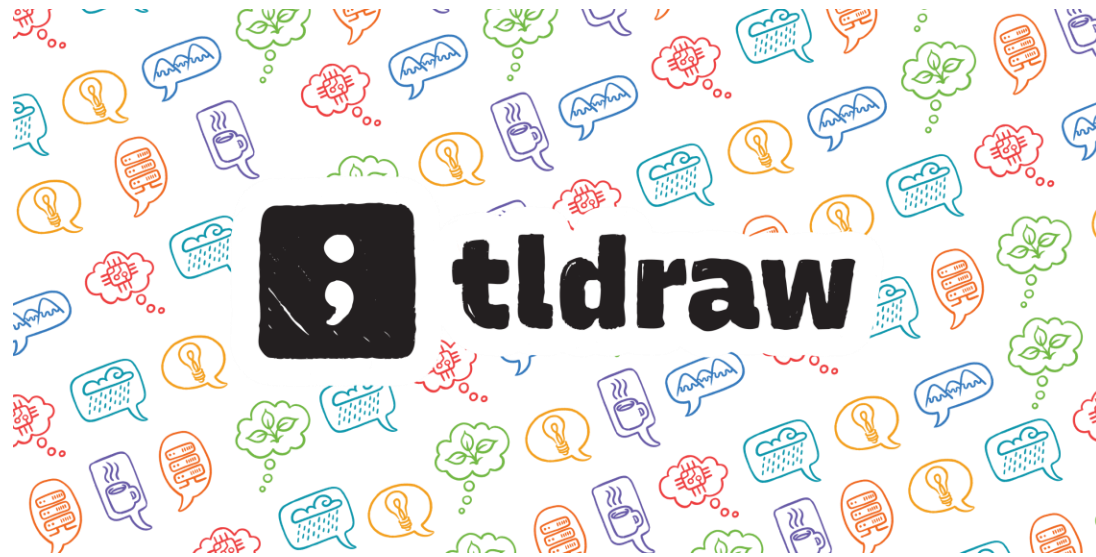
- Goal: develop and test a working model or set of working hypotheses about how (some part of) a system works
- Working model: an understanding of the pieces of the system (components), and the way they interact (connections)
- Observation, probes, and hypothesis testing
 - Helpful tools and techniques!



essentially,
all models are wrong,
but some are useful

George E. P. Box

Live Demonstration: tldraw



<https://github.com/tldraw/tldraw>

Participation Activity—Part 3

- Write down **what you would do if you wanted to modify the Duplicate functionality.**

Steps to Understand a New Codebase

- Look at README.md
- Clone the repo.
- Build the codebase.
- Figure out how to make it run.
- What do you want to mess with?
- Traceability - Attach a debugger
 - View Source
 - Find the logs.
 - Search for constants (strings, colors, weird integers (#DEADBEEF))

My experiences (headaches) at companies

- Documentation was **ALWAYS** out of date—often the core devs didn't know
- Had to ask someone to ask someone to help me get the project building (i.e., sit beside me for hours)
- Better take notes... not unusual to break something and need to do it all again
- Often the authors are no longer there
- So many design decisions are never written down, or they are trapped in old Jira tickets, commit messages, and emails

Program comprehension strategies

Novice

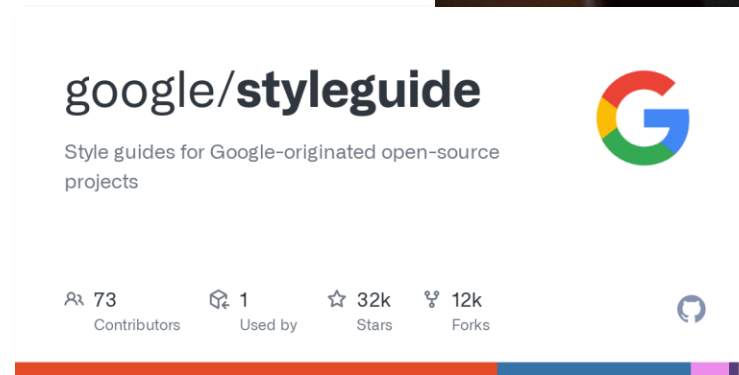
- Reads code line by line
- Revisits same code repeatedly
- Trial and error
- Only tests “happy path”

Expert

- “Top down”
- Recognizes patterns
- Forms hypotheses
- Checks up/downstream consequences

Observation: Software is full of patterns

- File structure
- System architecture
- Code structure
- Names
- ...



```
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
def __init__(self, *args, **kwargs):
    self.file = None
    self.fingerprints = set()
    self.logdupes = True
    self.debug = debug
    self.logger = logging.getLogger(__name__)
    if path:
        self.file = open(os.path.join(path, 'requests.json'),
                        'w')
        self.file.seek(0)
        self.fingerprints.update(e.request for e in self.requests)

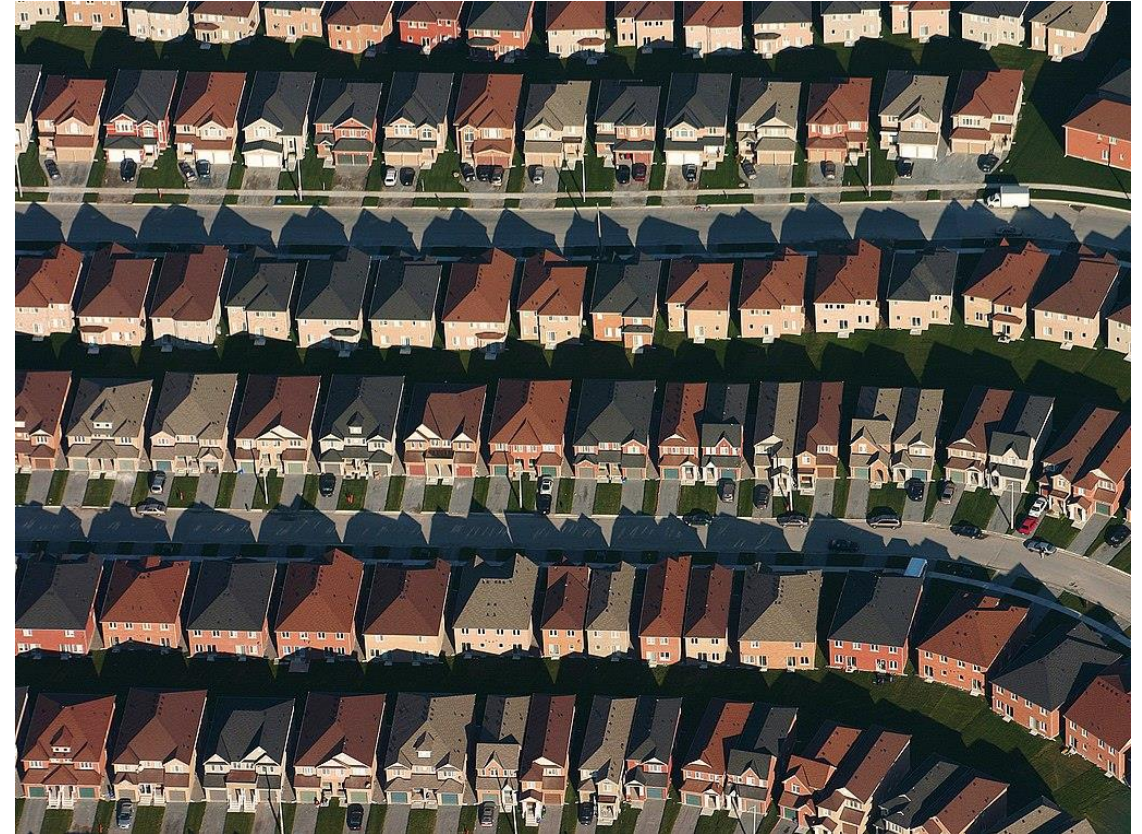
    @classmethod
    def from_settings(cls, settings):
        debug = settings.getbool('SUPERFICIAL_DEBUG')
        return cls(job_dir(settings), debug)

    def request_seen(self, request):
        fp = self.request_fingerprint(request)
        if fp in self.fingerprints:
            return True
        self.fingerprints.add(fp)
        if self.file:
            self.file.write(fp + os.linesep)

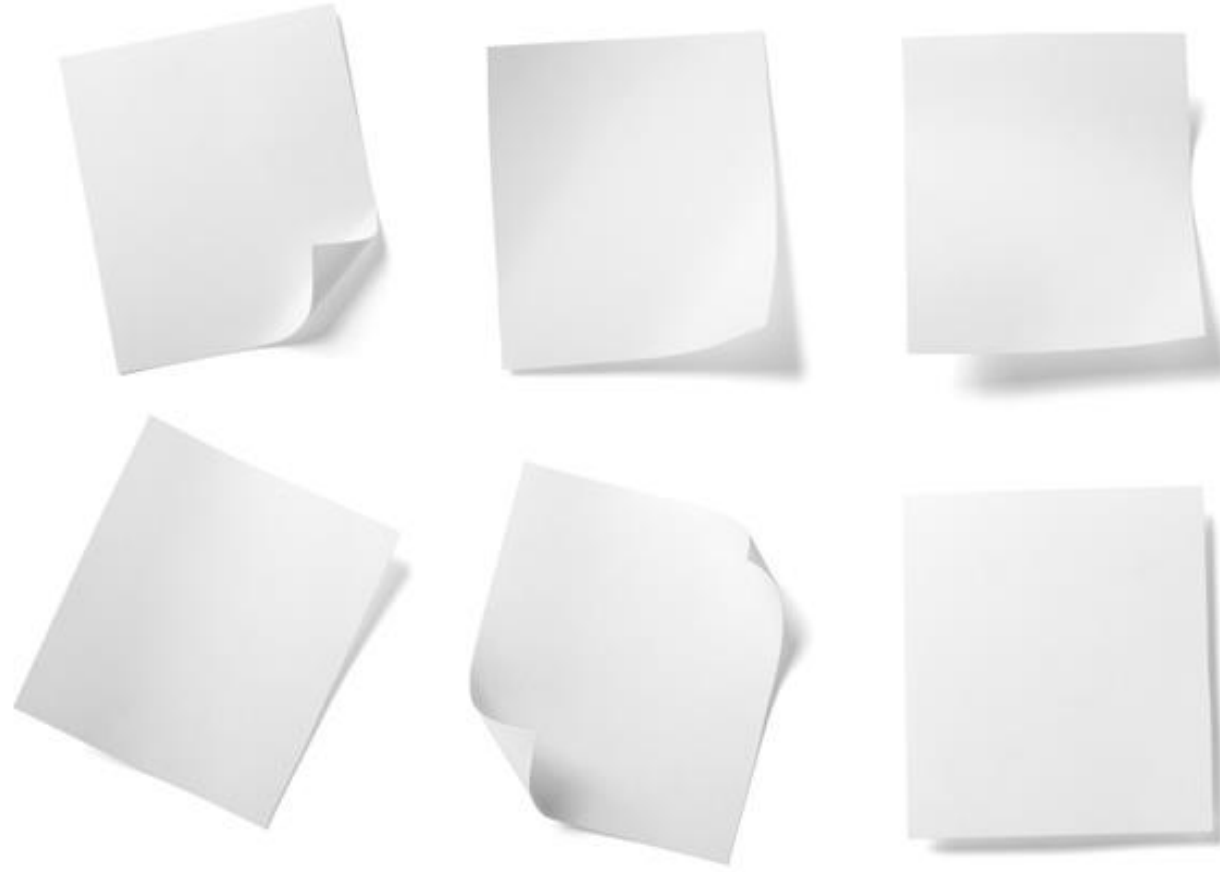
    def request_fingerprint(self, request):
        return request_fingerprint(request)
```

Observation: Software is massively redundant

- There's always something to copy/use as a starting point!



Observation: Code must run to do stuff!



Ask me about the
11,000-line file

Observation: If code runs, it must have a beginning...



Observation: If code runs, it must exist...

```
0x08048416 <+18>: jg     DWORD PTR [ebp+0x8], 0x1
0x08048419 <+21>: mov   eax, DWORD PTR [ebp+0xc]
0x0804841b <+23>: mov   ecx, DWORD PTR [eax]
0x08048420 <+28>: mov   edx, 0x8048520
0x08048425 <+33>: mov   eax, ds:0x8049648
0x08048429 <+37>: mov   DWORD PTR [esp+0x8], ecx
0x0804842d <+41>: mov   DWORD PTR [esp+0x4], edx
0x08048430 <+44>: mov   DWORD PTR [esp], eax
0x08048435 <+49>: call  0x8048338 <fprintf@plt>
0x0804843a <+54>: mov   eax, 0x1
0x0804843c <+56>: jmp   0x8048459 <main+85>
0x0804843f <+59>: mov   eax, DWORD PTR [ebp+0xc]
0x08048442 <+62>: add   eax, 0x4
0x08048444 <+64>: mov   eax, DWORD PTR [eax]
0x08048448 <+68>: mov   DWORD PTR [esp+0x4], eax
0x0804844c <+72>: lea  eax, [esp+0x10]
0x0804844f <+75>: mov   DWORD PTR [esp], eax
0x08048454 <+78>: call  0x8048338 <fprintf@plt>
```

Code must exist. But where?

- Locally installed programs: run cmd, OS launch, I/O events, etc.
 - Binaries (machine code) on your computer
- Local applications in dev: build + run, test, deploy (e.g., docker)
 - Source code in repository (+ dependencies)
- Web apps server-side: Browser sends HTTP request (e.g., GET, POST)
 - Code runs remotely (you can only observe outputs)
- Web apps client-side: Browser runs JavaScript, event handlers
 - Source code is downloaded and run locally (see: browser dev tools!)

Can running code be Probed/Understood/Edited?

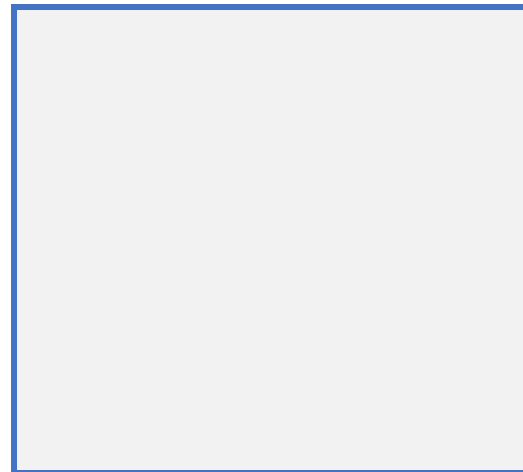
White-box



Source code built locally

(P+U+E)

Grey-box



Binaries running locally

Open source

(P+U)

Closed source

(P)

Black-box



Server-side apps running remotely

Open source

(U)

Closed source

(Talk to NSA)

Creating a model of unfamiliar code



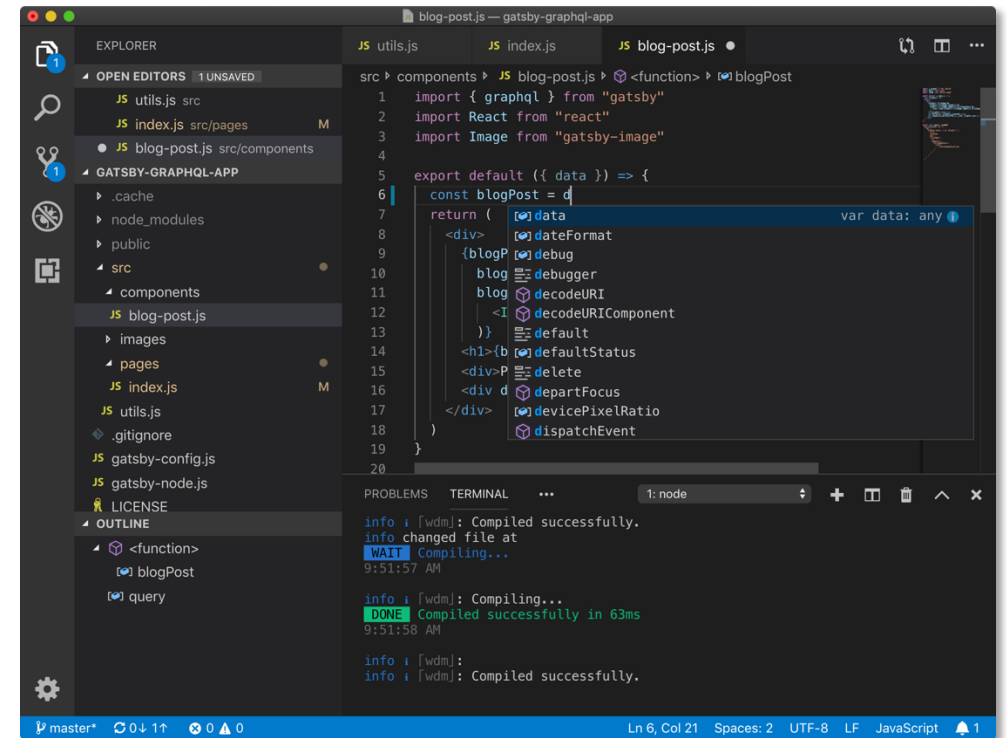
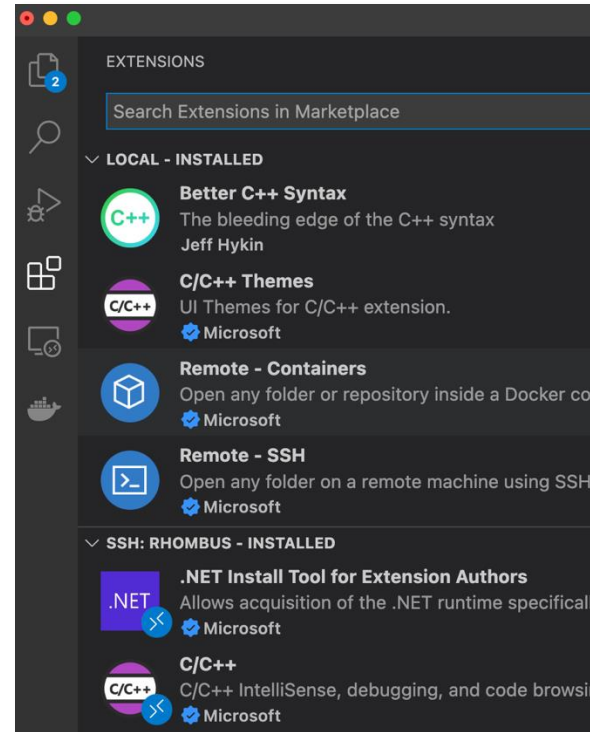
Source code built
locally

Information Gathering

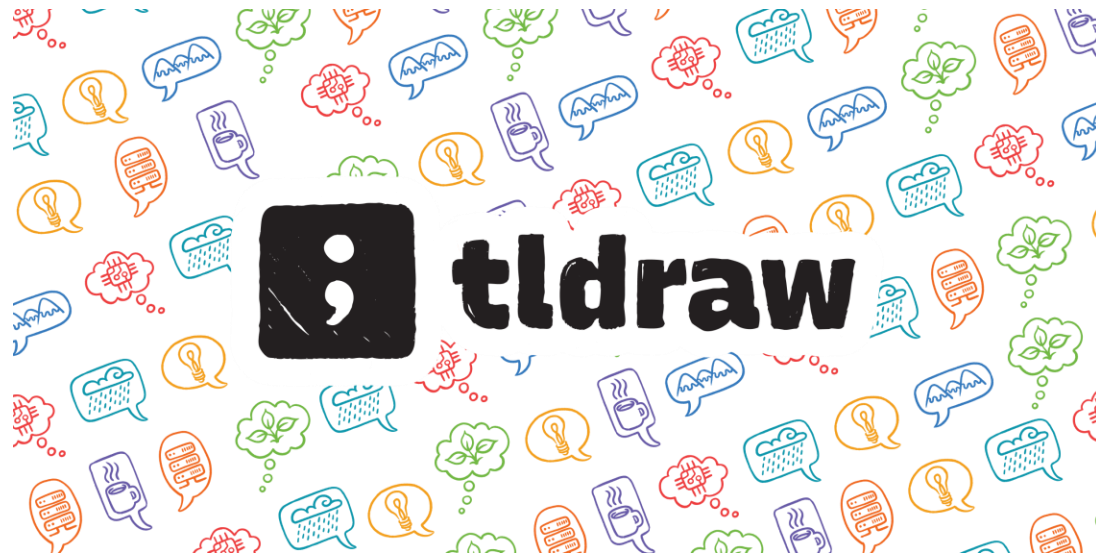
- Basic needs:
 - Code/file search and navigation
 - Code editing (probes)
 - Execution of code, tests
 - Observation of output (observation)
- Many choices here on tools! Depends on circumstance.
 - grep/find/etc. Knowing Unix tools is invaluable
 - A decent IDE
 - Debugger
 - Test frameworks + coverage reports
 - Google (or your favorite web search engine)
 - ChatGPT

Static Information Gathering: Use an IDE!

Real software is too complex to keep in your head



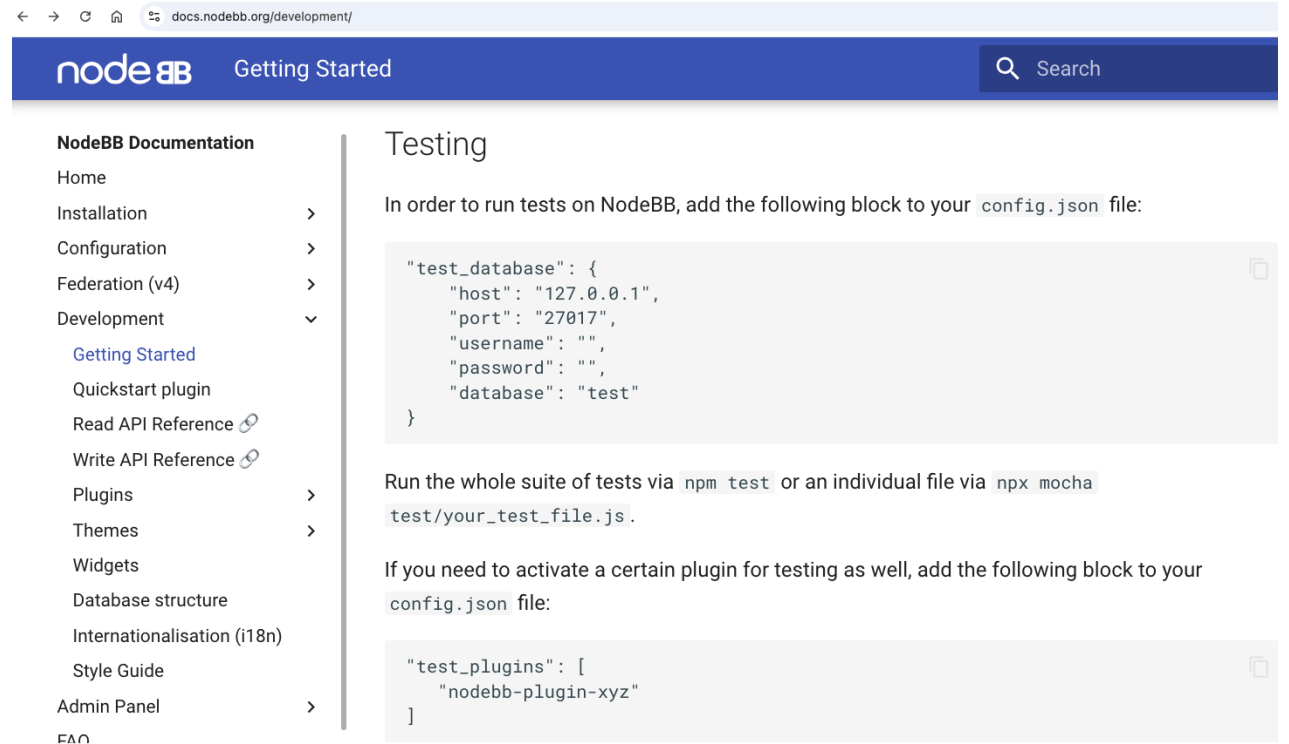
Live Demonstration: tldraw



<https://github.com/tldraw/tldraw>

Consider documentation and tutorials judiciously

- Great for discovering entry points!
- Can teach you about general structure, architecture (more on this later in the semester)
- Often out of date.
- As you gain experience, you will recognize more of these, and you will immediately know something about how the program works
- Also: discussion boards; issue trackers



The screenshot shows a web browser window with the URL `docs.nodebb.org/development/`. The page title is "nodeBB Getting Started" and there is a search bar. The left sidebar contains a "NodeBB Documentation" menu with items: Home, Installation, Configuration, Federation (v4), Development, Getting Started (highlighted), Quickstart plugin, Read API Reference, Write API Reference, Plugins, Themes, Widgets, Database structure, Internationalisation (i18n), Style Guide, and Admin Panel. The main content area is titled "Testing" and contains the following text:

In order to run tests on NodeBB, add the following block to your `config.json` file:

```
"test_database": {  
  "host": "127.0.0.1",  
  "port": "27017",  
  "username": "",  
  "password": "",  
  "database": "test"  
}
```

Run the whole suite of tests via `npm test` or an individual file via `npx mocha test/your_test_file.js`.

If you need to activate a certain plugin for testing as well, add the following block to your `config.json` file:

```
"test_plugins": [  
  "nodebb-plugin-xyz"  
]
```

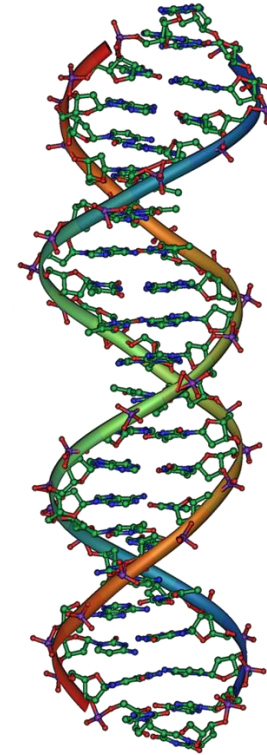
Discussion Boards and Issue Trackers

- Software is written by people.
- How can we talk to them?
- Fortunately, they probably aren't dead.
- So, you can report problems on GitHub.
- Or, ask them questions on StackOverflow.

The screenshot shows the Stack Overflow search results page for the query "java on mac". The page features a navigation bar with "stackoverflow", "About", "Products", "For Teams", a search bar containing "java on mac", and "Log in" and "Sign up" buttons. A left sidebar contains navigation links for "Home", "Questions", "Tags", "Users", "Companies", "Collectives", and "Teams". The main content area displays "Search Results" for "java on mac" with 500 results. Three results are visible, each with a title, vote count, answer count, view count, and tags. The first result is "How to set or change the default Java (JDK) version on mac-OS?" with 1311 votes and 36 answers. The second is "How to install Java 8 on Mac" with 1271 votes and 34 answers. The third is "Where is Java Installed on Mac OS X?" with 861 votes and 20 answers. A right sidebar contains an advertisement for "stackoverflow FOR TEAMS" and a section titled "Hot Network Questions" with several question links.

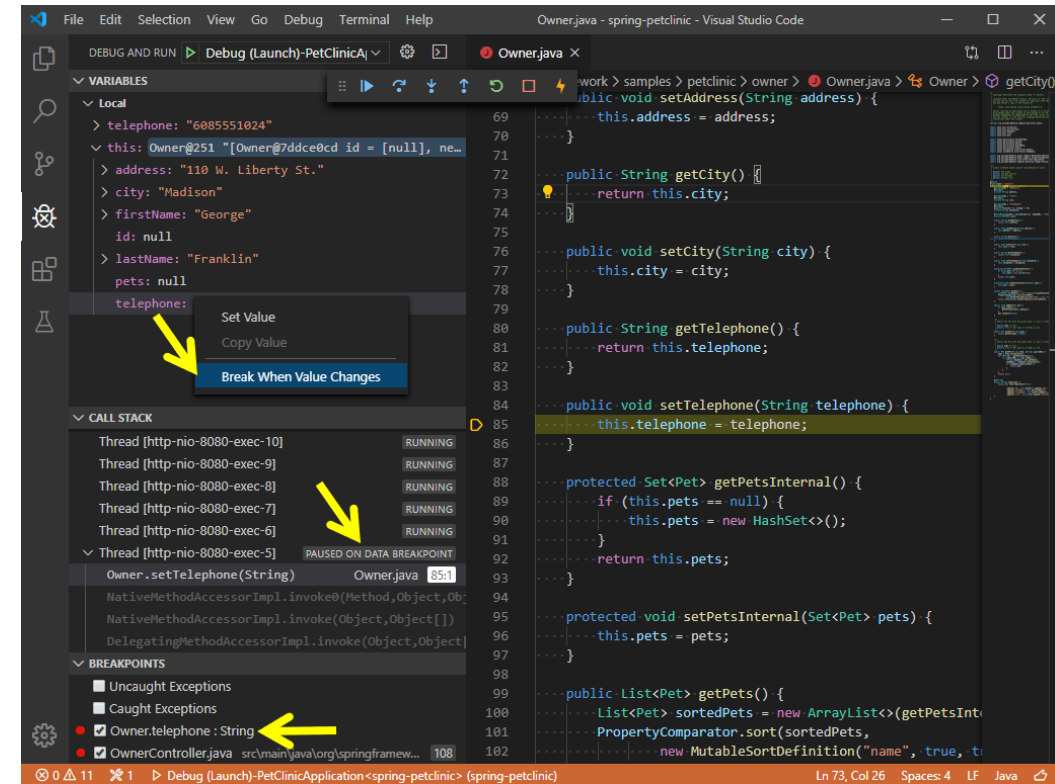
Dynamic Information Gathering Change helps to inform and refine mental models

- Build it.
- Run it.
- Change it.
- Run it again.
- How did the behavior change?



Probes: Observe, control or “lightly” manipulate execution

- print(“this code is running!”)
- Structured logging
- Debuggers
 - Breakpoint, eval, step through / step over
 - (Some tools even support remote debugging)
- Delete debugging
- Chrome Developer Tools

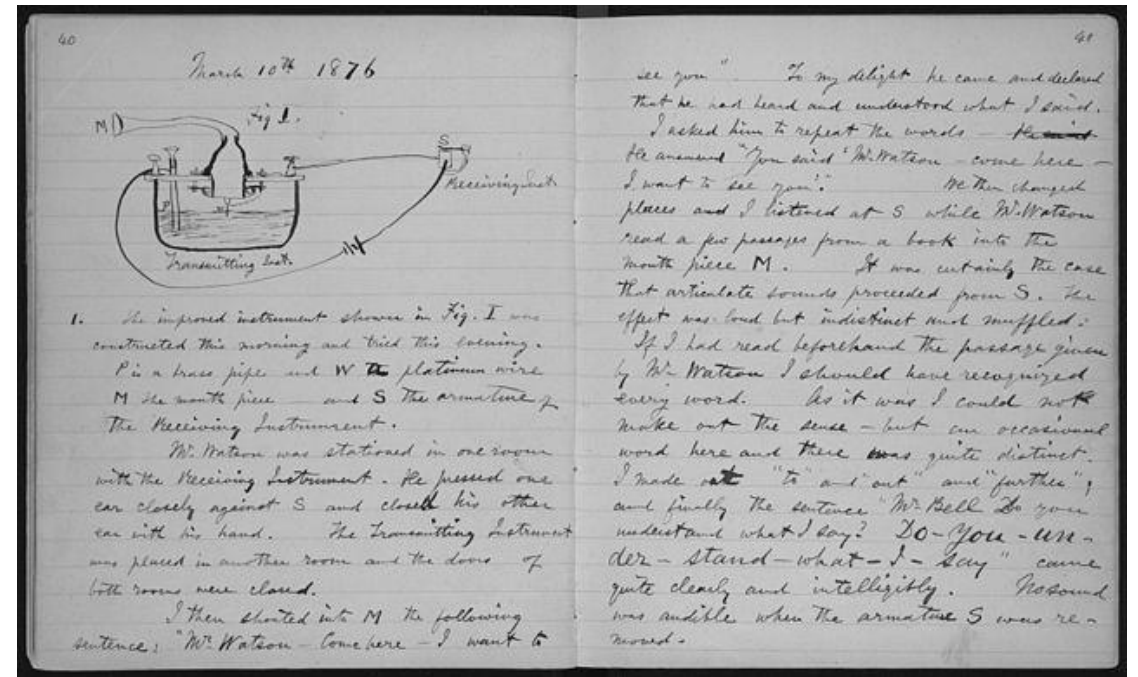


Step 0: Sanity check basic model + hypotheses

- Confirm that you can build and run the code.
 - Ideally both using the tests provided, and by hand.
- Confirm that the code you are running is the code you built
- Confirm that you can make an externally visible change
- How? Where? Starting points:
 - Run an existing test, change it
 - Write a new test
 - Change the code, write or rerun a test that should notice the change
- Ask someone for help

Document and share your findings!

- Update README and docs
 - Or better: use a Developer Wiki
 - Use [Mermaid](#) for diagrams
- Screencast on Twitch
- Collaborate with others
- Include negative results, too!



Next time...

- Metrics and Measurement