

# Software Archaeology and Anthropology

17-313 Fall 2023

Foundations of Software Engineering

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

Andrew Begel and Rohan Padhye

# Administrivia

- Slack
  - Please add a profile picture.
  - Ask questions in #general or #technical-questions. Please use threads.
- Office hours can be found on the course home page: <http://cmu-313.github.io>
- For those of you who requested to swap recitations, stay tuned.
- COVID or other health issues? Please stay home.

# Smoking Section

- Last full row



# Homework

- Homework 1 is released.
  - Part (a) is due Friday Sept 1, 11:59 pm. **That's tomorrow!**
  - Part (b) is due Thursday, Sept 7, 11:59pm.
  - This is an individual assignment; we will compose groups next week.
  - Get started early, ask for help, and check the #technical-questions channel; chances are your questions have been asked by others!

# Team Formation Survey Due Friday

- Team formation survey is posted on Canvas and in Slack #announcements.
- Please fill in by TONIGHT!

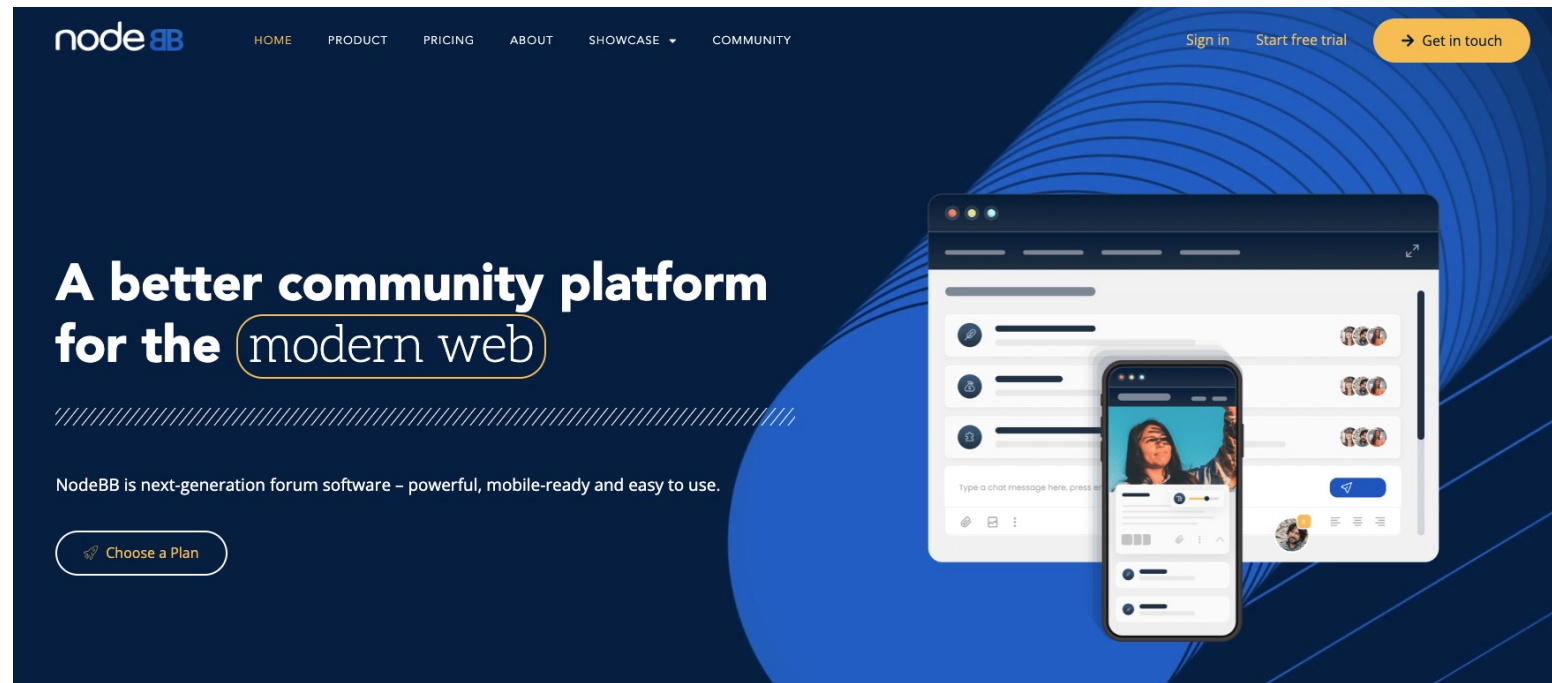


# Learning Goals

- 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

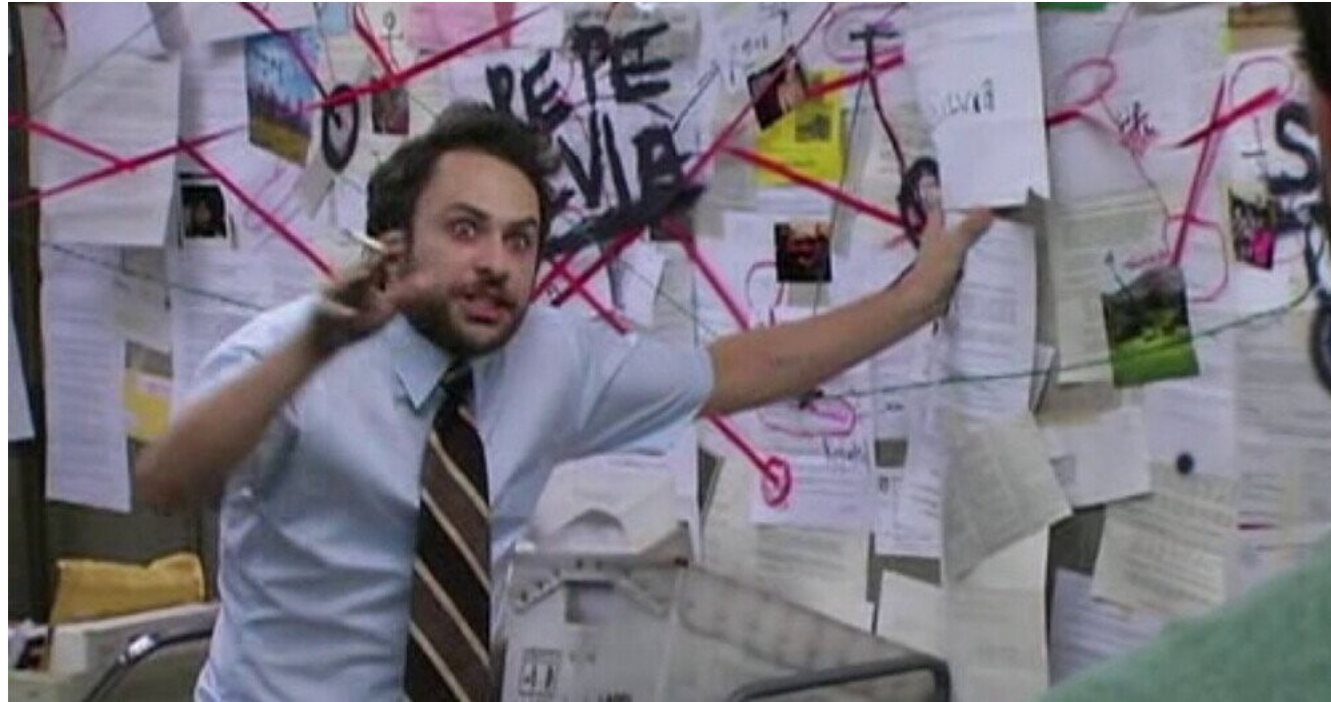
- ... do something with it!



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



# Challenge: How do I tackle this codebase?

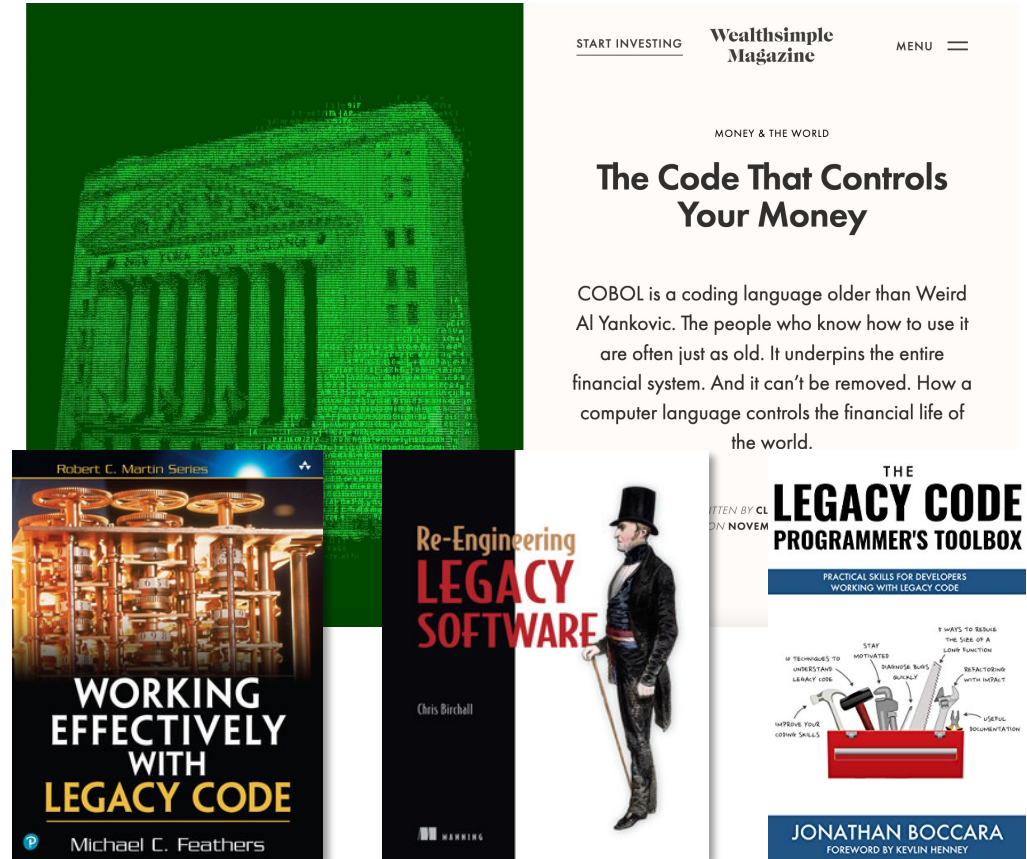


# Challenge: How do I tackle this codebase?

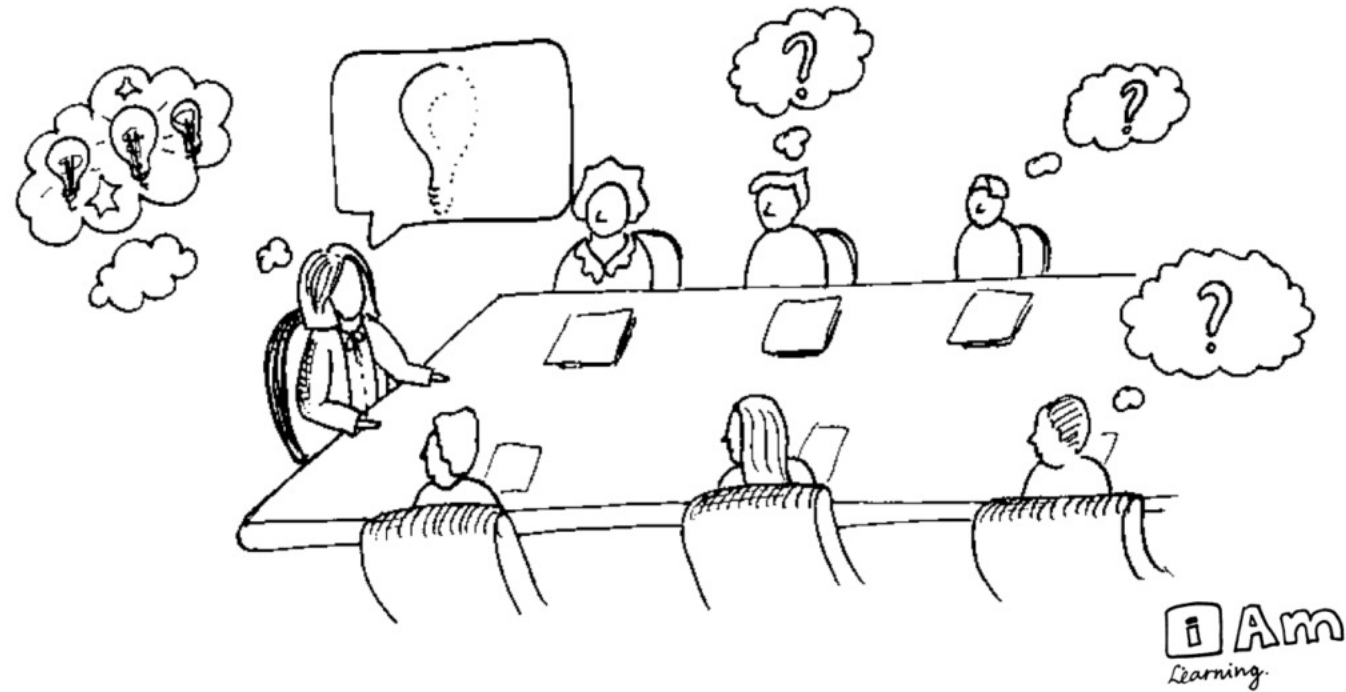
- Leverage your previous experiences (languages, technologies, patterns)
- Consult documentation, whitepapers
- Talk to experts, code owners
- Follow best practices to build a working model of the system

# 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

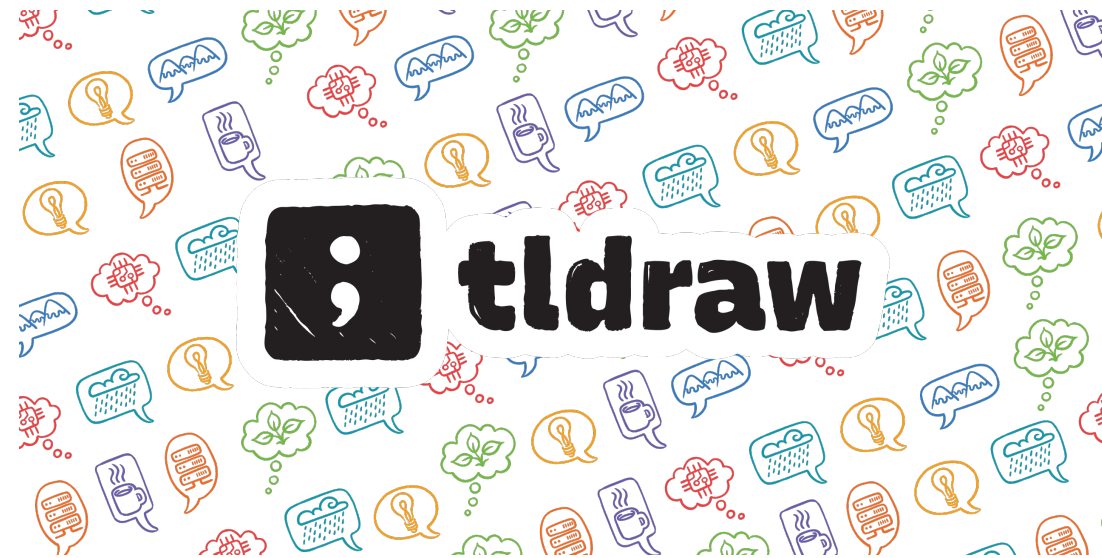
- Goal: develop and test a working model or set of working hypotheses about how (some part of) a system works
- Working model: an understanding of the pieces of the system (components), and the way they interact (connections)
- Focus: Observation, probes, and hypothesis testing
  - Helpful tools and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



essentially,  
all models are wrong,  
but some are useful

George E. P. Box

# Live Demonstration: tldraw



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

# 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?
  - Clone and own
- Traceability - Attach a debugger
  - View Source
  - Find the logs.
  - Search for constants (strings, colors, weird integers (#DEADBEEF))

# Participation Activities

- Pull out your phone.
- Download the Gradescope app.
- Log into Gradescope.
  - Use email login, not SSO.
- Click “+” button to add 17-313.
  - Entry Code: G24487



Download Gradescope  
for iPhone



Download Gradescope  
for Android

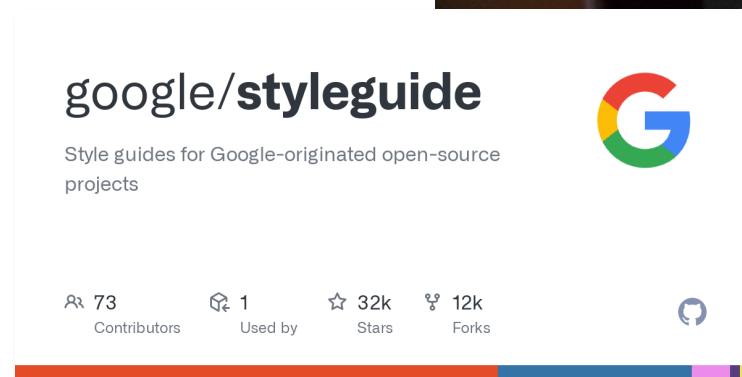


# Participation Activity

- Take out a piece of paper.
- Write down one pro and one con about trying to understand a new codebase by compiling and building it vs. just reading the code.
- Pair with your neighbor and discuss your answers. Do you agree?
- Share with the class!
- Submit it on Gradescope by the end of class.
  - Under Not Submitted (Assignments), click on August 31 Activity.
  - Take a picture of your paper.
  - Assign the picture to Question 1.
  - Submit.

# 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

self.file = None
self.fingerprints = self()
self.logdupes = True
self.debug = debug
self.logger = logging.getLogger(__name__)
if path:
    self.file = open(os.path.join(path, "requests.json"),
                    "a")
    self.file.seek(0)
    self.fingerprints.update({e: str(fp)} for e in self())

@classmethod
def from_settings(cls, settings):
    debug = settings.getbool("SUPPRESS_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!



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>
```

# The Beginning: Entry Points

- Locally installed programs: run cmd, OS launch, I/O events, etc.
- Local applications in dev: build + run, test, deploy (e.g., docker)
- Web apps server-side: Browser sends HTTP request (GET/POST)
- Web apps client-side: Browser runs JavaScript, event handlers

# 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?

Transparent



Source code built locally

(P+U+E)

Translucent



Binaries running locally

Open source

(P+U)

Closed source

(P)

Opaque



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

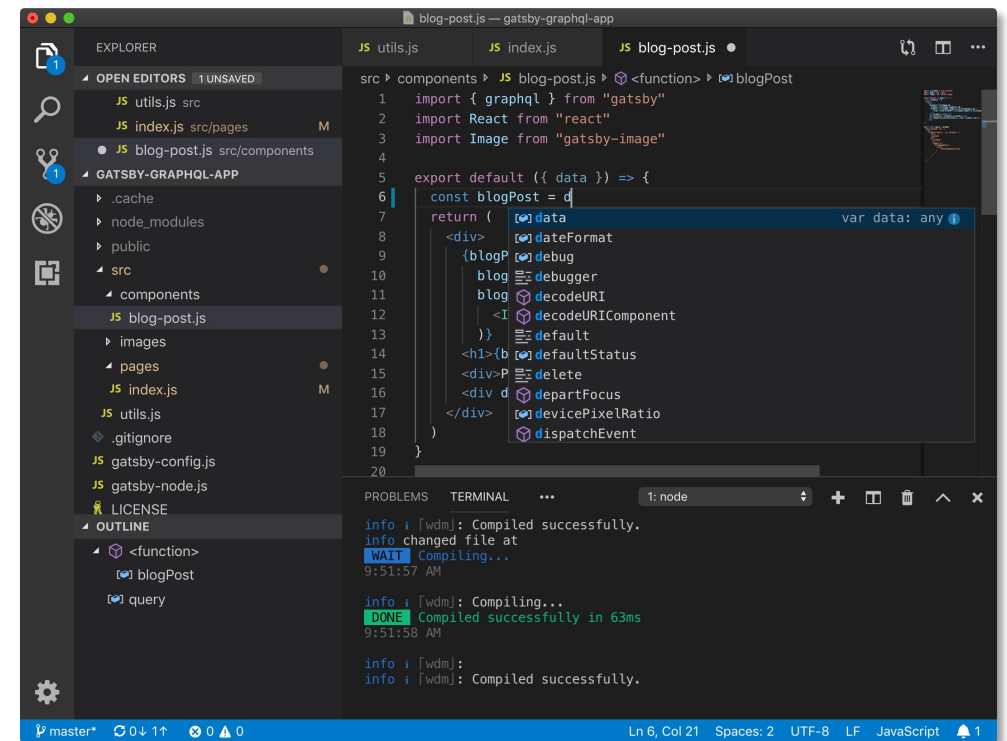
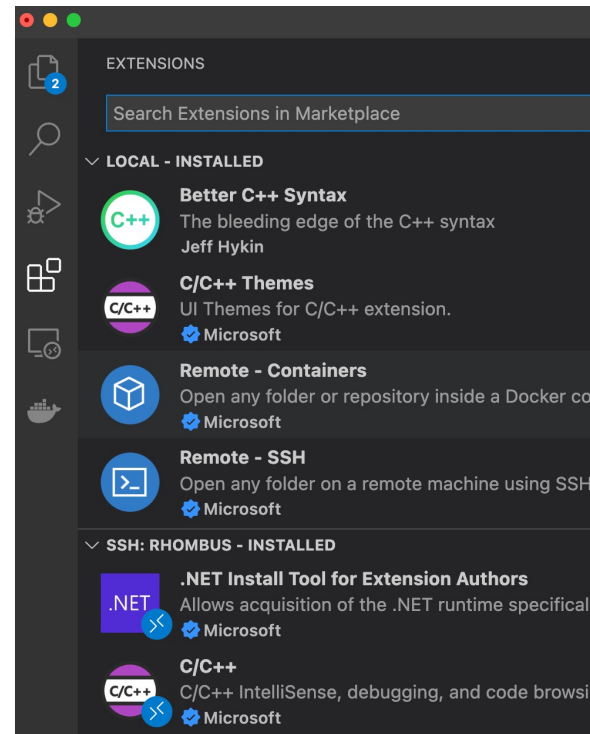
- Basic needs:
  - Code/file search and navigation
  - Code editing (probes)
  - Execution of code, tests
  - Observation of output (observation)
- Many choices here on tools! Depends on circumstance.
  - grep/find/etc. Knowing Unix tools is invaluable
  - A decent IDE
  - Debugger
  - Test frameworks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



At the command line: **grep** and **find**!  
(Google for tutorials)

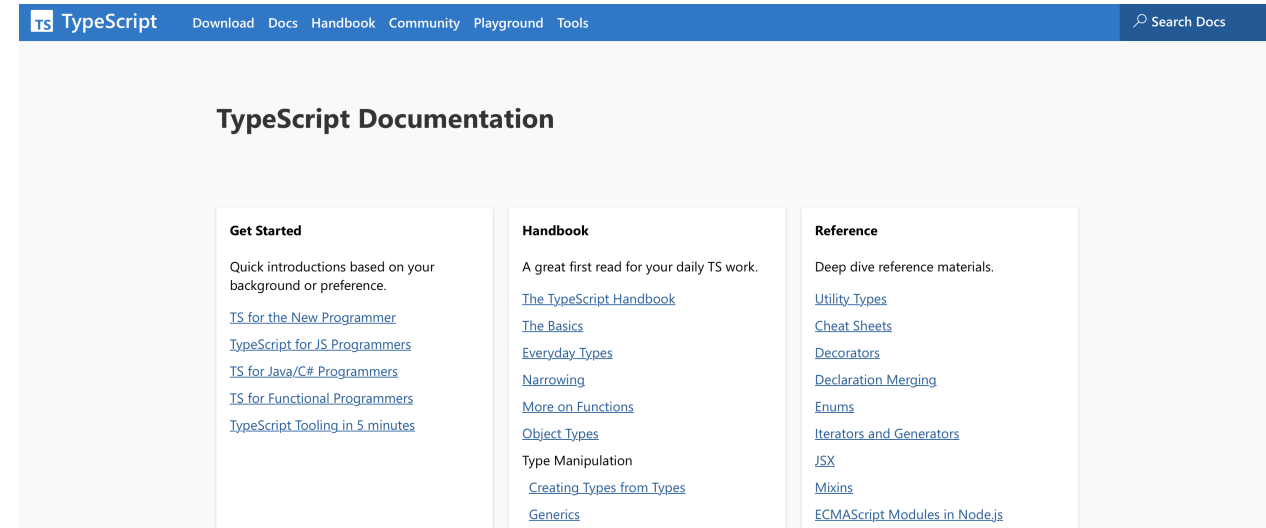
# Static Information Gathering: Use an IDE!

## Real software is too complex to keep in your head



# 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 the TypeScript Documentation website. The header is blue with the TypeScript logo and navigation links: Download, Docs, Handbook, Community, Playground, Tools. A search bar is on the right. The main content is titled "TypeScript Documentation" and is divided into three columns:

- Get Started**: Quick introductions based on your background or preference. Links include: [TS for the New Programmer](#), [TypeScript for JS Programmers](#), [TS for Java/C# Programmers](#), [TS for Functional Programmers](#), and [TypeScript Tooling in 5 minutes](#).
- Handbook**: A great first read for your daily TS work. Links include: [The TypeScript Handbook](#), [The Basics](#), [Everyday Types](#), [Narrowing](#), [More on Functions](#), [Object Types](#), [Type Manipulation](#), [Creating Types from Types](#), and [Generics](#).
- Reference**: Deep dive reference materials. Links include: [Utility Types](#), [Cheat Sheets](#), [Decorators](#), [Declaration Merging](#), [Enums](#), [Iterators and Generators](#), [JSX](#), [Mixins](#), and [ECMAScript Modules in Node.js](#).

# 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 displays the Stack Overflow website interface. At the top, the navigation bar includes the Stack Overflow logo, links for 'About', 'Products', and 'For Teams', a search bar containing 'java on mac', and buttons for 'Log in' and 'Sign up'. The main content area is titled 'Search Results' and shows 'Results for java on mac' with 'Search options not deleted'. It indicates '500 results' and provides sorting options: 'Relevance', 'Newest', and 'More'. The search results list several questions:

- How to set or change the default Java (JDK) version on mac-OS?**: 1311 votes, 36 answers, 1.4m views. Asked by Venkat 13.2k on Feb 23, 2014 at 5:46.
- How to install Java 8 on Mac**: 1271 votes, 34 answers, 1.3m views. Asked by user3763100 13k on Jun 21, 2014 at 15:05.
- Where is Java Installed on Mac OS X?**: 861 votes, 20 answers, 1.0m views. Asked by Thunderforge 19.6k on Apr 5, 2013 at 4:58.
- How do I install Java on Mac OSX allowing version switching?**: 479 votes, 11 answers. Asked by user3763100 13k on Jun 21, 2014 at 15:05.

On the right side, there is an advertisement for 'stackoverflow FOR TEAMS' with the text 'With no silos to hold you back, you could be unstoppable.' and a 'Become Unstoppable' button. Below the ad is a 'Hot Network Questions' section listing various popular questions.

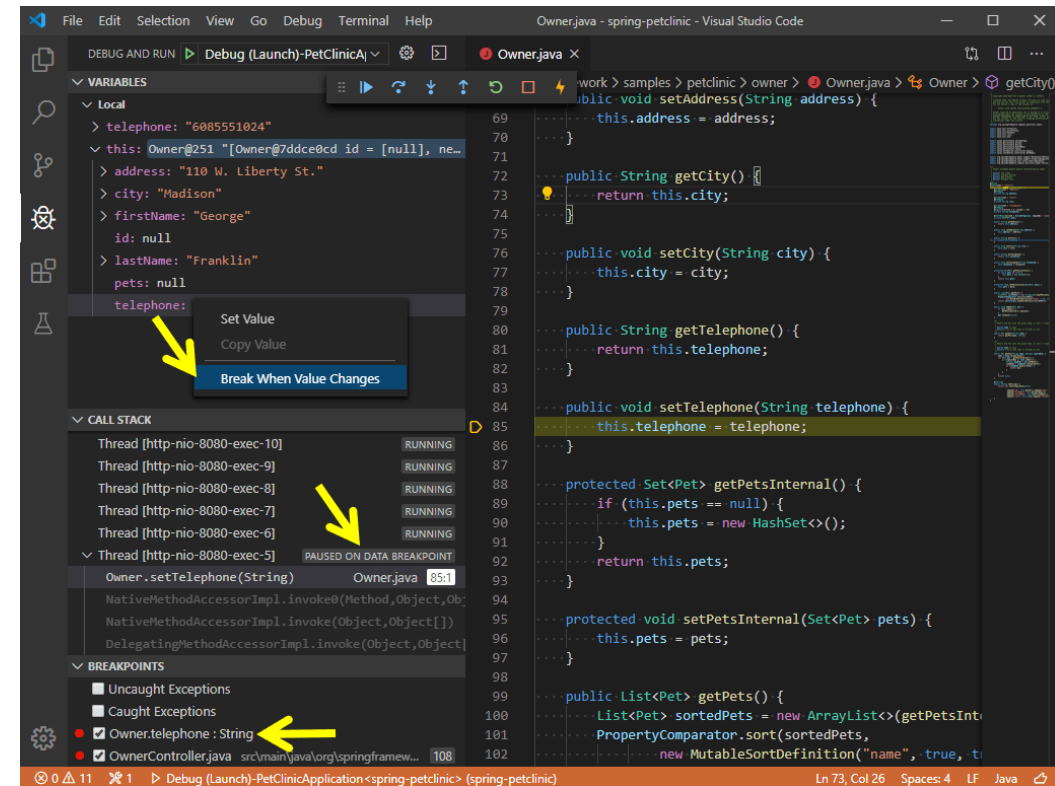
# Dynamic Information Gathering Change helps to inform and refine mental models

- 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



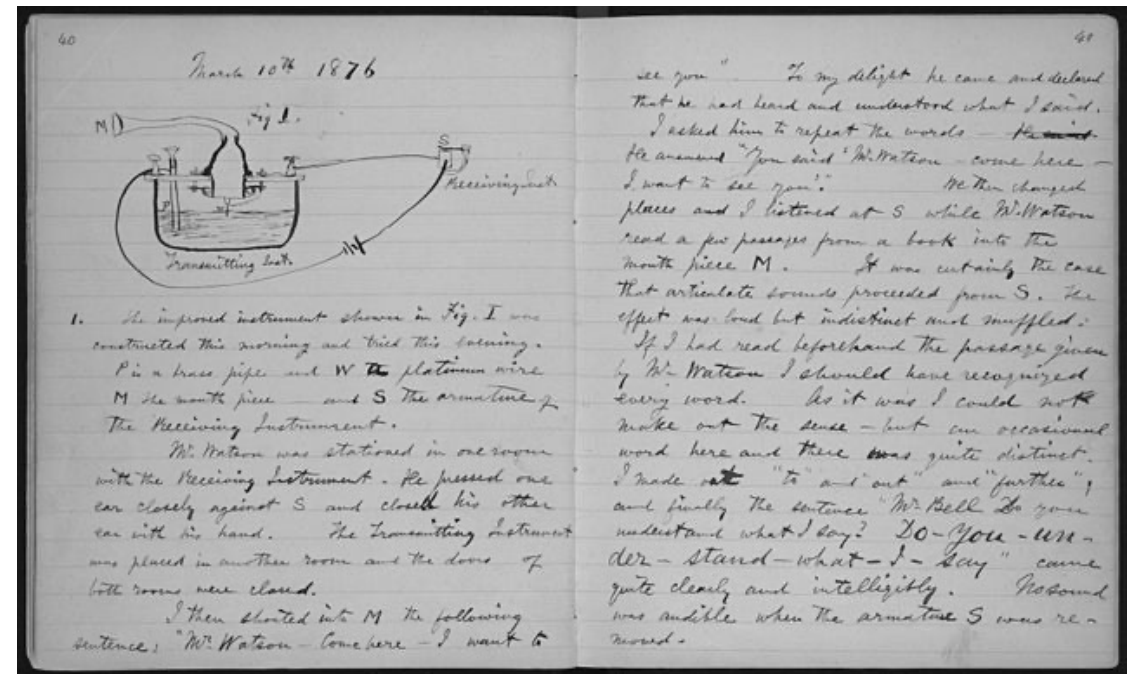


# Step 0: Sanity check basic model + hypotheses

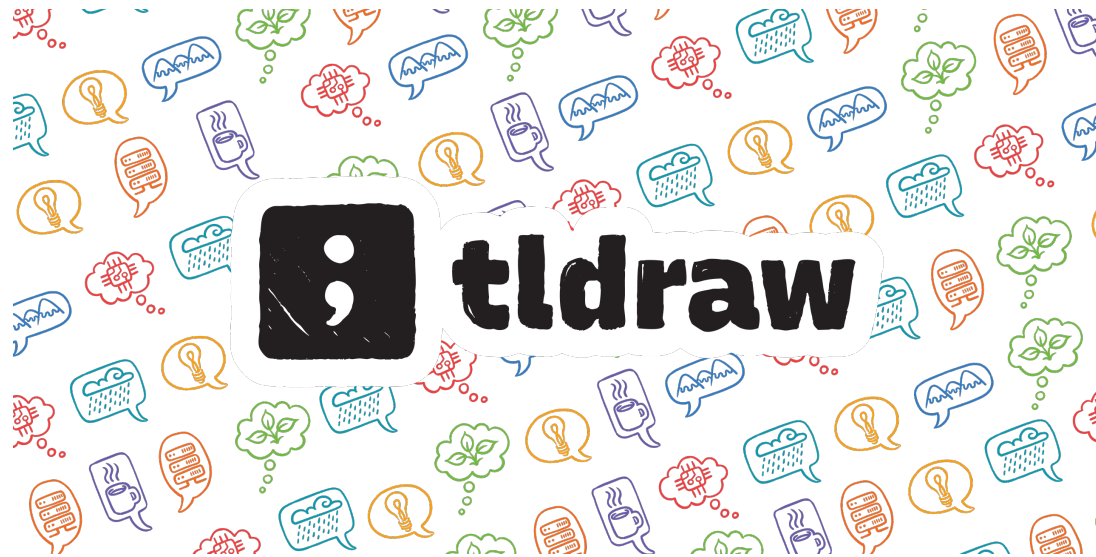
- Confirm that you can build and run the code.
  - Ideally both using the tests provided, and by hand.
- Confirm that the code you are running is the code you built
- Confirm that you can make an externally visible change
- How? Where? Starting points:
  - Run an existing test, change it
  - Write a 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!



Let's try some of these techniques again...



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