# Software Archaeology and Anthropology

17-313 Fall 2025

Foundations of Software Engineering

https://cmu-17313q.github.io

Eduardo Feo Flushing

Carnegie Mellon University

# Administrivia

- Slack

  - Please add a profile picture.

  - Ask questions in `#general` or `#technical-questions`.

    - Please use threads.

    - Use the search tool.

- Office hours can be found on the course home page: http://cmu-17313q.github.io

# Project P1

- **P1A**: Checkpoint due next Sunday (August 31$^{st}$)
  - Only 5% of total P1 points – meant to ensure you start on time
- **P1B**: Due Sunday, September 7$^{th}$
  - Refactor a javascript file to improve its quality
  - It will be posted tonight
  - Start early

# Archaeology vs Anthropology

Artifacts
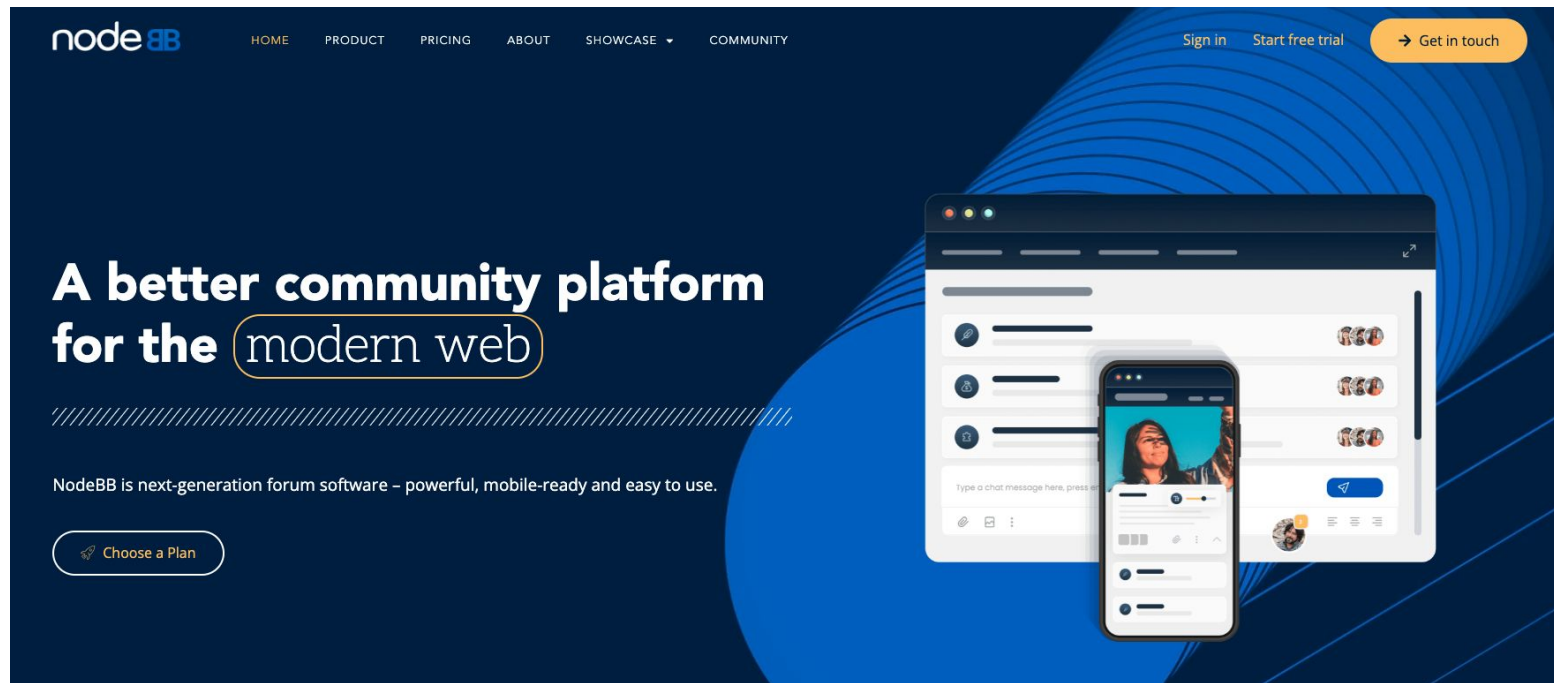
Human societies and culture

# Learning Goals

- Understand and scope the task of working with a new and complex piece of existing software

- Appreciate the importance of configuring an effective IDE

- Contrast different types of code execution environments: local, remote, application, and library-based

- Enumerate both static and dynamic strategies for understanding and modifying a new codebase

S3D

Carnegie
Mellon
University

# Context: big old pile of code

- … do something with it!

# You will never understand the entire system!

# Challenge: How do I tackle this codebase?

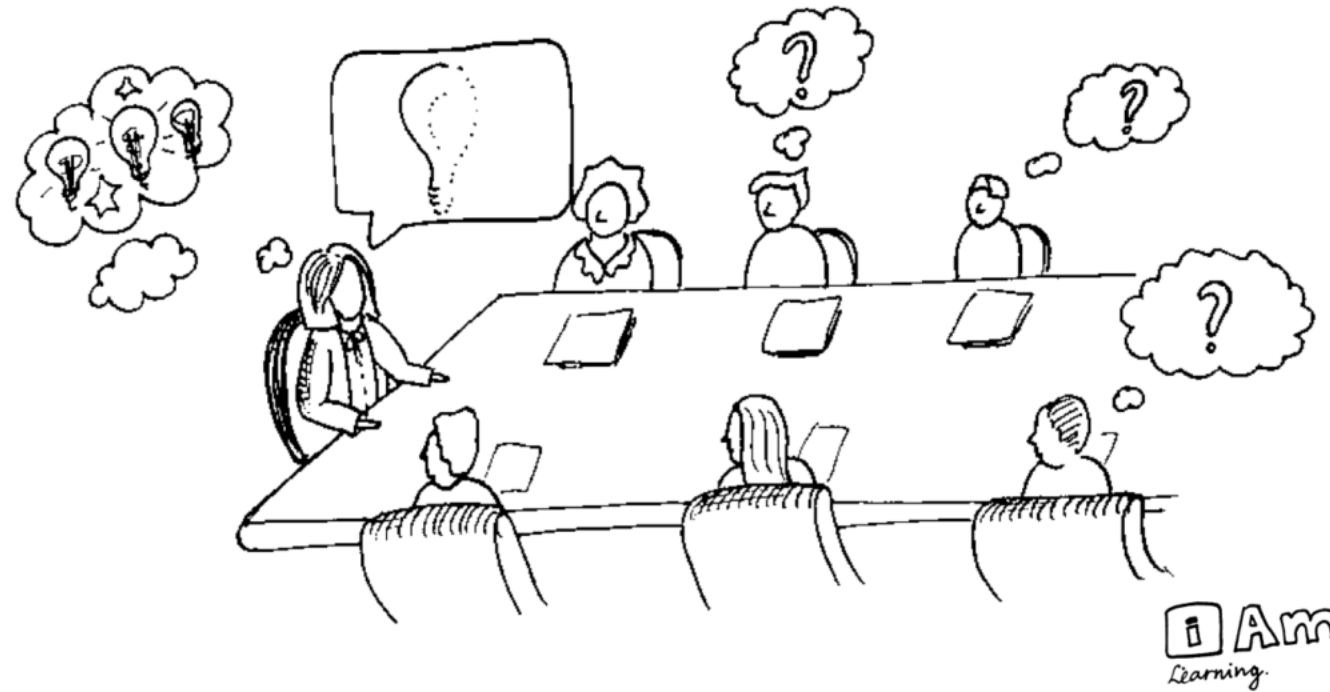# Challenge: How do I tackle this codebase?

- 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

S3D

# 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)

START INVESTING     **Wealthsimple Magazine**     MENU ≡

MONEY & THE WORLD

**The Code That Controls Your Money**

COBOL is a coding language older than Weird Al Yankovic. The people who know how to use it are often just as old. It underpins the entire financial system. And it can't be removed. How a computer language controls the financial life of the world.

# Why? Because of Tacit Knowledge

# Today: How to tackle codebases

- Goal: Develop and test a working model (or set of hypotheses) about how 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](https://github.com/tldraw/tldraw)

Carnegie
Mellon
University

# Academic Honesty

- Standard Collaboration Policy

- In group work, be honest about contribution of group members; do not cover for others

- Unless explicitly prohibited, you may use generative AI (e.g. ChatGPT) to help you write your prose and code. You are responsible for its correctness. Be sure to attribute the content to the service you used. (let us know if you have concerns about teammate's work)

- DO NOT submit participation sheets for people who are not in class. This is considered an <u>academic integrity violation</u>

# 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 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!
- Write your own andrewID on the paper, leave it at the end of class.

S3D

# Observation: Software is full of patterns

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



google/**styleguide**

Style guides for Google-originated open-source projects

73 Contributors   1 Used by   32k Stars   12k Forks

# Observation: Software is massively redundant

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

Carnegie
Mellon
University

# Observation: Code must run to do stuff!

# Observation: If code runs, it must have a beginning…

# The Beginning: Entry Points

- Locally installed programs: run cmd, OS launch, I/O events, etc.

- Web apps server-side: Browser sends HTTP request (GET/POST)

- Web apps client-side: Browser runs JavaScript, event handlers

# Can running code be Probed/Understood/Edited?

| Transparent | Translucent | Opaque |
|:-----------:|:-----------:|:------:|

Source code built locally

Binaries running locally

Server-side apps running remotely

| (P+U+E) | Open source | Closed source | Open source | Closed source |
|:-------:|:-----------:|:-------------:|:-----------:|:-------------:|
|         | (P+U)       | (P)           | (U)         | (Talk to NSA) |

S3D

# **Creating a model of unfamiliar code**

Source code built locally

# Static 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 LLaMA

# Static Information Gathering: Use an IDE!
# Real software is too complex to keep in your head

Carnegie
Mellon
University

# Dependency maps

# 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

Carnegie Mellon University

# Discussion Boards and Issue Trackers

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

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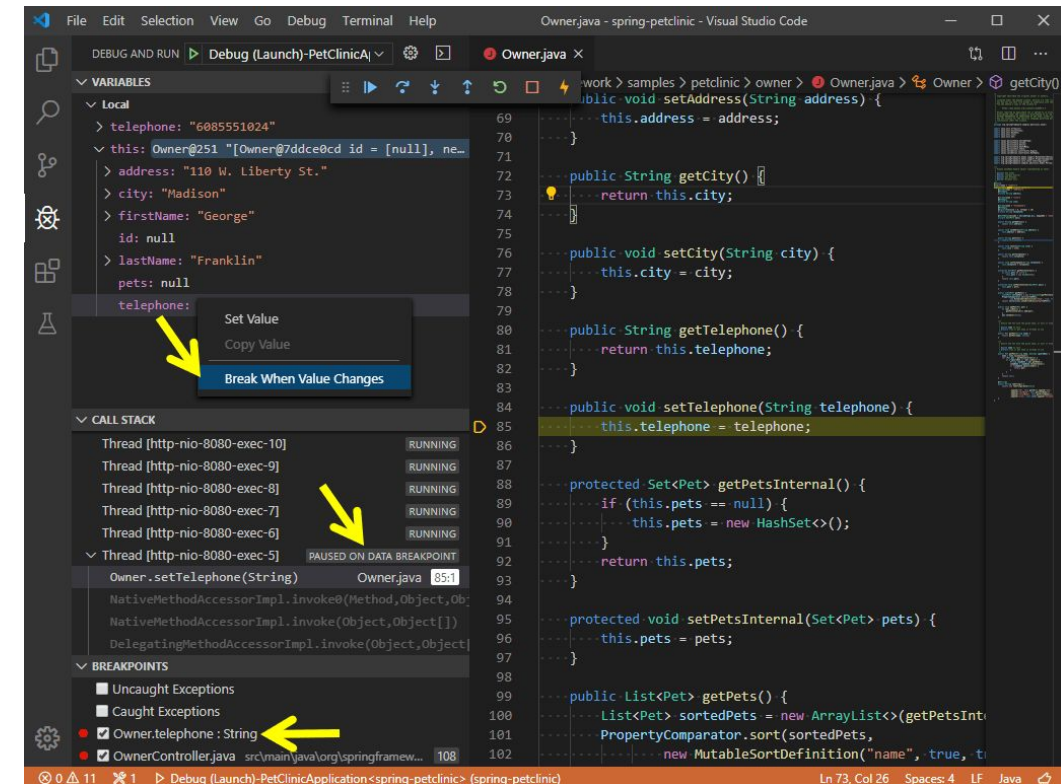- How did the behavior change?

# How to start?

- 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
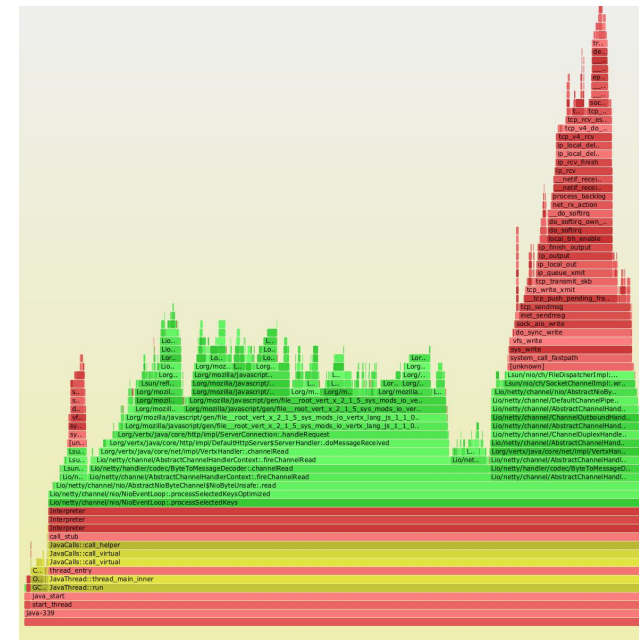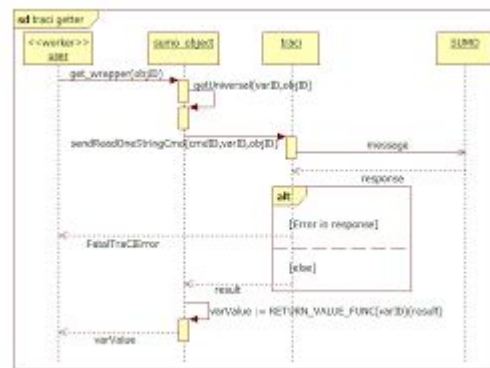- Ask someone for help

# 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

# Runtime code analysis tools

- Collect runtime traces and visualize them
  - Flame graphs
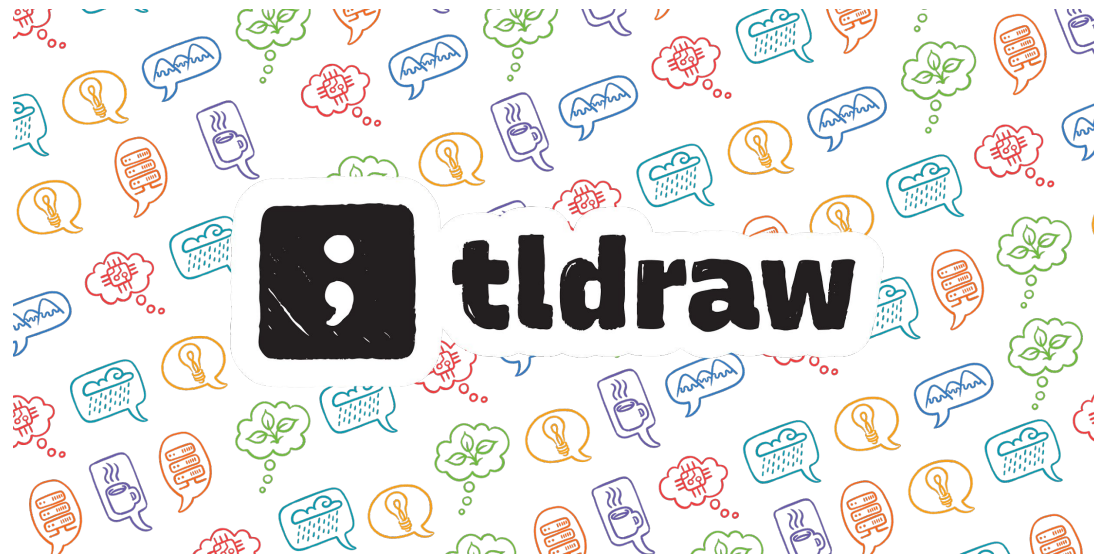  - Sequence diagrams
- Use judiciously

# Tip: Find a particular thing and trace the action backward



E.g.,
Where do categories come from?
How are they stored?
How are they rendered?

S3D

Carnegie
Mellon
University

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



**https://github.com/tldraw/tldraw**

# Remember…

- Reading and understanding code is one of the most important skills you should learn

- It's common to get stuck or feel overwhelmed. **Don't give up!**

- Consider yourself lucky! Things are much easier today

Carnegie
Mellon
University

# 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