# My Claude Code Workflow

## A Framework for AI-Assisted Empirical Research

Scott Cunningham

# The Core Insight

**Typical Use**

"Write me a function that does X"

shift →

**My Use**

"I'm seeing Y in the data. What could explain this?"

**Thinking partner**

The hard part of empirical research isn't writing code—
it's **figuring out what code to write** and
whether **results mean what you think**.

# The Fundamental Problem: Claude Has Amnesia

| Session 1 | forget ➤ | Session 2 | forget ➤ | Session 3 |

Every session starts from **zero context**.

Most people just re-explain everything verbally.
**I build external memory in markdown files.**

# External Memory via Markdown

| README.md | CLAUDE.md | logs/*.md | docs/ |
|:---:|:---:|:---:|:---:|
| What we're working on | Problems & solutions for Claude | Session-by-session progress | Reference materials |

**Institutional memory** that persists even though Claude's doesn't

# Part I

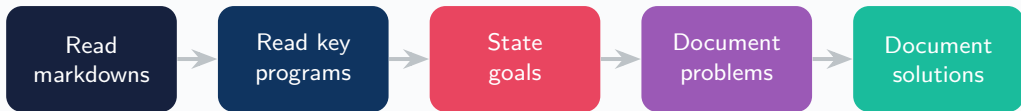The Daily Practices

# The Socratic Method for Alignment

*"Do you see the issue with this specification?"*

*"That's not it. The problem is the standard errors."*
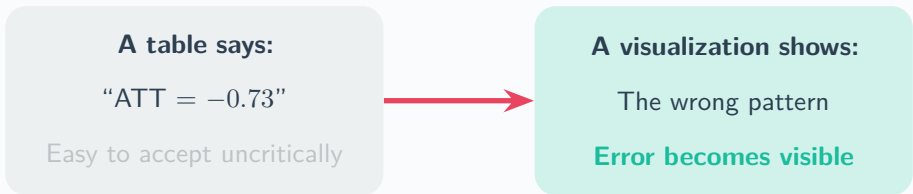
*"Guess at what I'm about to ask you to do."*

**Why?** If Claude guesses wrong,
that reveals a misunderstanding
that needs correcting **before** proceeding.

# Session Startup Routine



Read markdowns → Read key programs → State goals → Document problems → Document solutions

Takes **2 minutes**. Claude starts each session
**informed** rather than ignorant.

# Verification Through Visualization

**A table says:**

"ATT $= -0.73$"

Easy to accept uncritically

**A visualization shows:**

The wrong pattern

**Error becomes visible**

**Trust pictures over numbers.** Always ask for figures.

# Part II

Cross-Software Replication

# Bugs Are Orthogonal Across Languages

R  **bug A**  =  **Stata**  **bug B**  =  **Python**  **bug C**

?

A bug in `dplyr` code is **unlikely** to be the same bug in Stata.

If all three produce the **same answer**, you have **high confidence**.

# Why Cross-Language Validation Works

**Syntax errors are language-specific**

Both must navigate their own syntax

**Default behaviors differ**

NA propagation, factor ordering, etc.

**Implementation paths differ**

`group_by()` vs. `collapse`

Same answer $\Rightarrow$ Same concept $\Rightarrow$ **Correct**

# The Validation Table

| Check | R | Stata | Match? |
|-------|---|-------|--------|
| Sample rows | 5,234 | 5,234 | ✔ |
| Mean age (treatment) | 31.7 | 31.7 | ✔ |
| Mean hourly wage | $5.43 | $5.43 | ✔ |
| N missing earnings | 412 | 412 | ✔ |

If any differ ⇒ **investigate before proceeding**.
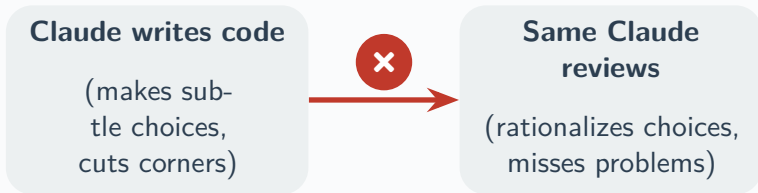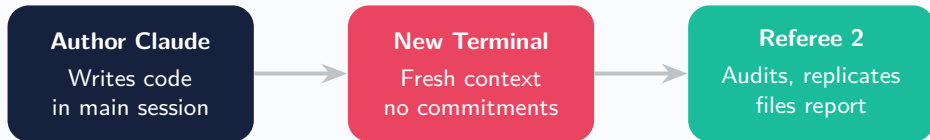
This isn't paranoia. It's **methodology**.

# Part III

Adversarial Review

# You Can't Grade Your Own Homework

**Claude writes code**

(makes subtle choices, cuts corners)

❌ →

**Same Claude reviews**

(rationalizes choices, misses problems)

Asking the same Claude to review its own code is like asking a student to grade their own exam.

# The Solution: Referee 2 Protocol

| Author Claude | New Terminal | Referee 2 |
|---|---|---|
| Writes code in main session | Fresh context no commitments | Audits, replicates files report |

True adversarial review requires **separation**.

Fresh context. No prior commitments. Formal protocol.

# The Five Audits

| Code Audit | Cross-Lang Replication | Directory Audit | Output Audit | Econometrics Audit |
|:---:|:---:|:---:|:---:|:---:|
| Coding errors, missing values, merge diagnostics | R, Stata, Python match to 6 decimal places | Replication package ready? | Tables/figures from code or manual? | Specification coherence, standard errors |

**Formal Referee Report + Replication Scripts**

## What Referee 2 Catches

✔ **Unstated assumptions** — "Did you verify X or just assume it?"

✔ **Alternative explanations** — "Could the pattern come from something else?"

✔ **Documentation gaps** — "Where does it explicitly say this?"

✔ **Logical leaps** — "You concluded A, but evidence only supports B"

✔ **Missing verification** — "Have you actually checked the raw data?"

# The Philosophy

Referee 2 isn't about being negative.

It's about **earning confidence**.

A conclusion that survives rigorous challenge
is stronger than one never questioned.

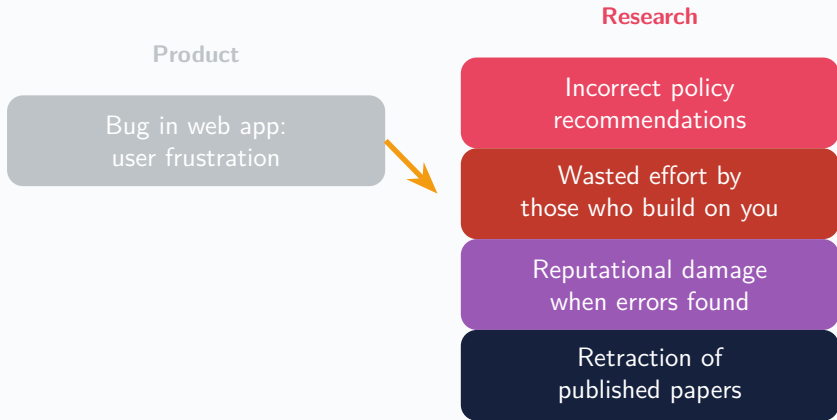Find weak points and either **fix them** or **accept them knowingly**.

# Part IV

Why This Works

# Research Is Not Product Development

| Aspect | Product Dev | Research |
|--------|-------------|----------|
| Goal | Ship working code | Understand correctly |
| Error cost | Bug in production | Wrong conclusion |
| Iteration | Fast ship fix later | Slow careful right |
| Testing | Unit tests CI | Visual inspection |
| Success | Does it run? | Does it mean what we think? |

The difference between "the code runs" and "**the code is correct**."

# The Stakes Are Different

Product

Bug in web app:
user frustration

Research

Incorrect policy
recommendations

Wasted effort by
those who build on you

Reputational damage
when errors found

Retraction of
published papers

## The Workflow in Summary

| Dimension | My Approach |
|---|---|
| Philosophy | Thinking partner, not code generator |
| Memory | External via markdown (Claude has amnesia) |
| Verification | Cross-software: R = Stata = Python |
| Review | Referee 2 protocol in fresh terminal |
| Documentation | First-class output, not afterthought |
| Visualization | Trust pictures over numbers |
| Speed | Correctness over velocity |

## The Key Insight

Claude Code isn't just a code generator—

it's a **thinking partner**.

Ask it questions. Make it explain.
Verify visually. Document everything.

And when stakes are high, spawn **Referee 2**.

# That's How I Use AI for Research

Scott Cunningham

scunning.com · causalinf.substack.com