

Building Real Software with AI

Lessons from Six Weeks of Shipping

Christopher Meiklejohn

Adjunct Faculty, Carnegie Mellon University

Software Engineer, DoorDash

Foundations of Software Engineering · 15-313 · Spring 2026

What I'm Building

Zabriskie — a social platform for people who care about live music, film, and books.

- The internet used to have third places built around taste
- You found people through what they were into, not where they lived
- Nothing exists for people who organize their lives around shows, albums, and culture
- The jam band community is the proof of concept — culture already built around shared experience
- **Zabriskie is the tool: iOS · Android · Web · Live shows · Setlists · Couch touring**

zabriskie.app

SHOW MAP

All Years

140 shows · 64 cities · 5 bands

UPCOMING

- #1 **Goose** · **ExploreAsheville.com Arena**
Asheville, NC · Apr 10, 2026
- #2 **Goose** · **Coca-Cola Amphitheater**
Birmingham, AL · Apr 11, 2026
- #3 **Goose**
Orlando, FL · Apr 12, 2026

zabriskie.app

Chris

Chris drops off after this show

Patrick is joining the run

Sara is joining the run

Fri, Apr 17, 2026

The BayCare Sound
Clearwater, FL

Patrick Sara

Mary is joining the run

Tony is joining the run

Sat, Apr 18, 2026

The St. Augustine Amphitheatre
St. Augustine, FL

Patrick Mary Sara Tony

Sun, Apr 19, 2026

The St. Augustine Amphitheatre
St. Aug

zabriskie.app

ALBUM COMFORT VINYL

One Size Fits All
Frank Zappa And The Mothers Of Invention* · 1975

@cmeik ARCHITECT

Inca Roads

The last official release of Frank Zappa's band, The Mothers of Invention. This is the . . . probably) best of the band. 4, Ruth Hussey, Charles Thompson, Napoleon Murphy Brock, Tom Courton

zabriskie.app

GOOSE MODE

YOUR NEXT SHOW

09 : 20 : 59 : 07
DAYS HRS MIN SEC

ExploreAsheville.com Arena
Friday, April 10 · Asheville, NC
The Spring 2026 Tour

IN VENUE · 8

AT HOME · 8

Tour Timeline

The Spring 2026 Tour
13 shows · 5 attending · 4 cities

Why This Was Impossible Before AI

I've wanted to build this for years. Every previous attempt died.

- Not because the idea was wrong — because the scope is insane for one person with a full-time job
- Think about what “a social app” actually means:
 - Go backend with 300+ database migrations
 - React frontend + iOS app + Android app + web PWA — three clients
 - Real-time features, push notifications, setlist scrapers, audio players
 - Auth, security, CI, deployment
- That's not a side project. That's a team of five working full-time.
- Every previous attempt stalled around week three
- **AI changed the math. The first version was built in six days. On an iPad. In bed.**

What We Actually Built

Three clients. One backend. One person.

- Backend: Go, PostgreSQL, Railway (auto-deploy on push to main)
- Frontend: React (Vite), served as PWA
- iOS/Android: Capacitor — native shell, web engine, one codebase
- Server-Driven UI (SDUI): Backend controls what the app renders. No App Store deploy needed.
- Social: Posts, likes, comments, follows, bookmarks, @mentions, push notifications
- Live shows: Real-time setlists, couch touring, live chat, Dynamic Island
- Integrations: Spotify, setlist.fm, Relisten, Archive.org, Phantasy Tour, TTBase
- Infrastructure: GitHub Actions CI, Playwright E2E tests
- Scale: 300+ database migrations, 43 SDUI screens, ~170 files touched

Path to First Release

384	27	168	6	144
messages	sessions	files touched	days to v1	commits/week

- GitHub repo auto-deploying on push to main via Railway
- Claude Code iPad app running in a container — making changes and pushing to the repo
- Built the first version on an 11" iPad, in bed, at 2am
- Didn't touch a computer until day two

PART 1

What AI Actually Gives You

The Cost of Trying Drops to Zero

The MVP was built in 6 days. On an iPad. In bed.

- Within a few hours: working app — login, feed, posting — deployed on Railway
- Spotify integration, setlist.fm integration, real show data, venues, setlists
- Didn't move to a laptop until day two — the screen felt limiting, not the tool
- **Once the MVP was live, iteration cost dropped to near zero**
- Redesigned a major feature 6 times in 9 hours. Not mockups — running code.
- The version that shipped was the fifth attempt. I only knew it was right because I'd seen the four that weren't.
- **When trying something costs nothing, you stop planning and start building.**

Breadth Becomes Possible

One person. One week. All of this:

Domain	What
Audio	Relisten API, playlist players, caching
Real-time	Live chat, setlist polling
Data	Scrapers for 4 setlist sources + Archive.org
Security	Rate limiting, CORS, JWT hardening
Infra	Deployment pipeline, monitoring, error handling
Mobile	Dynamic Island fix, safe area insets
Social	Quick Post, @mentions, bug forum
Design	6 UI redesigns in 9 hours, admin analytics

"I'm not some 10x developer. I have a collaborator that doesn't sleep."

Fixing Bugs from the Beacon Theatre

Three nights watching Tedeschi Trucks Band. Shipping from my seat.

1. Notice a bug from my seat in the audience
2. Describe it to Claude on my phone in plain English
3. Claude fixes it and pushes to production
4. Pull up the app and verify — SDUI means instant deploy

No laptop. No IDE. No terminal. Just a phone.

Server-driven UI: change the experience without shipping through the App Store.

PART 2

Where It Goes Wrong

Two failure modes. The AI either follows the wrong path confidently, or it declares victory before the work is done.

The Confidence Problem

30+ instances of wrong-approach debugging. Claude moves confidently in the wrong direction.

What it does:






- Reads code instead of running it
- Builds elaborate theories
- Prefers the layer it's comfortable with
- Declares fixes without testing

What it should do:

- Check the logs first
- Test with real data
- Flag uncertainty before acting
- Verify runtime behavior

The Same Bug, Five Times





A background job that automatically marks shows as “live” when they start.

-  Bug ships. The poller runs every 60 seconds and does nothing. Nobody notices.
-  Emergency fix works — but the AI removes it during a refactor the next day.
-  Second fix only applied to one of two pollers. The other silently fails.
-  Third fix introduces a new edge case. Still broken.
-  Finally: real tests that insert data, run the poller, and verify the status changes.

Four confident fixes. Same bug. Only stopped when we wrote a test that proved it worked.

“The Show Is Happening Right Now”

Saturday night. Goose at Jam in the Streets. First show of spring.

-  8:00 PM — Show starts. RSVP toggle broken. Live Activities dead.
-  8:00–8:30 — Claude reads code for 30 min. No testing. No logs. Just guessing how both the RSVP and Live Activity features worked — wrong guesses for both.
- 8:45 PM — I had to explain how both features actually worked and force it to curl the endpoint. Answer was right there.
-  Root cause — timezone math: midnight UTC = 8 PM Eastern. Show marked as “past” the moment it started.
-  The fix — one line. If the show’s status is already “live,” don’t check the time.

The developer found both bugs, not the AI. When there’s no error message, the AI has nothing to work with.

Premature Victory

The second failure mode: the AI declares things done before they actually work.

- **Phantom APIs** — Frontend wired to backend routes that don't exist.
- **Disconnected routes** — Backend endpoints added but never wired to the frontend.
- **Silent form drops** — Form submits, endpoint returns 200, half the parameters never arrive.
- **Dead buttons** — Buttons render. Click handlers attached. Nothing happens on tap.
- **Fake success** — “Feature complete” = the code compiles. Not that it works for a user.

The representation of correctness was verified. The reality was not.

Memory Isn't Learning

Claude has persistent memory. It saves lessons after failures. Then ignores them.

Ship → Break → Emergency fix → Fix breaks something → Fix the fix → Add rule → Ignore next week

Rules saved. Rules ignored. Same session.

- Admin-merged a PR bypassing CI
- Pushed directly to main in the same session
- Both had explicit memory files saying never to do them

It declared it had learned. It hadn't.

Worktrees — Isolation for AI Sessions

Problem: Your main checkout has dirty files. If the AI touches it, things get tangled.







Solution: Git worktrees.

- A worktree is a second checkout of the same repo in a separate directory
- Its own branch, its own working tree, shared history
- Claude Code can spin one up automatically — isolated sandbox for a task
- When it's done, merge the branch and delete the worktree. Main checkout untouched.
- **Each AI session gets its own clean environment. No cross-contamination.**

In theory, this is perfect. In practice...

The Worktree Disaster

A two-file Go version fix. What actually happened:

-  Claude was in a clean, isolated git worktree — designed for exactly this.
-  Instead of fixing there, it cd'd into the main repo with a dozen dirty files.
-  Staged everything. Committed with the Go fix. Pushed. Opened a PR.
-  PR got auto-merged before I could close it.
-  Duplicate variables, broken E2E tests, wrong form placeholders — all in prod, site **down**.
-  Four follow-up commits across three PRs to clean up a two-file change.

It declared the isolation was set up. Then it walked right out of it.

How the Deployment Pipeline Works

Push to main → Railway builds both services → app goes live.

- **Two services on Railway — a Go backend and a React frontend. Each deploys independently.**
- **Healthcheck — pings the backend's /health endpoint. If not ready in time, deploy fails.**
- **Auto-deploy — every push to main triggers both services to rebuild.**
- If a config change breaks the build or makes it start slower, the healthcheck times out
- The deploy fails, and every subsequent push also fails until someone fixes it

The Deployment Cascade

One unnecessary line. Five broken deploys.

While fixing an unrelated feature, Claude added a healthcheck to the Railway config as a “drive-by improvement.”

Five successive commits changed:

- The port
- The builder
- The start command
- The Dockerfile

None worked. Because the real problem was never diagnosed.

It declared each fix “done.” Five times. None of them were.

PART 3

How I Made the AI Work for Me

Every failure in Part 2 taught me something. Most people write a prompt, get an answer they don't like, and try the same prompt on a different model. That's not how you build software with AI. You have to learn how to drive it.

Describe Problems, Not Solutions

✗ Don't say: "Fix the WebSocket reconnection handler in showPoller.go line 147"

✓ Do say: "The live setlist isn't updating during the show"

- **Feature request? Ask it to explore and make a plan before writing any code.**
- **Bug? Ask it to reproduce it first before doing anything else.**
- Let it explore — it might find the bug in a file you didn't suspect
- But if it's going the wrong direction, redirect hard and fast
- "Wrong layer." "That's not the bug." "Curl the endpoint."
- Short, direct, corrective. A conversation, not a contract.

Be Reactive and Corrective

Anthropic sends weekly usage reports to Claude Code power users. Mine described my style as:

“Reactive and corrective rather than spec-driven”

- Don't write a 500-word spec and hand it over
- Start. Watch what happens. Course-correct in real time.
- “No, not that file.” “Check the logs first.” “Stop reading code and test it.”
- A 29-second median response time — barely reading output before firing back

Every minute you let it run down a wrong path is a minute and tokens wasted. You're the driver, not the passenger.

Make the Context Persistent

Without project instructions, every session starts from zero.

```
# CLAUDE.md – your most important artifact

## Rules: "No database triggers. EVER."
## Rules: "Two-Attempt Rule: step back after 2 failures"
## Rules: "Always restart servers after .go file changes"
## Build: "cd backend && go build ./..."
```

What goes in it:

- Architecture decisions (“we use SDUI, never hardcode React pages”)
- Build & run commands (so it doesn’t guess wrong)
- Rules from past failures (every rule is a scar)
- “Never do this” lists (timestamps, triggers, pkill)
- Testing requirements (“run E2E tests after ANY change”)

Good Session vs. Bad Session

GOOD SESSION

“Add TTB with 58 shows”

- Clear deliverable: band + shows + scraper
- Built a scraper for a community setlist database
- Deployed mid-show, verified live
- Done in one sitting
- **Tight scope → shipped → stop**

BAD SESSION

“Debug push notifications”

- No clear stopping point
- Deep into provisioning profiles, certs
- Plausible path, completely wrong
- Actual problem: missing AppDelegate methods
- **Open-ended → spiraled → wasted**

What You Provide vs. What the AI Provides

You:

- Taste — naming things, choosing what feels right, knowing your users
- Product intuition — noticing when a user gives up mid-flow
- Domain knowledge — knowing which data source is reliable
- The observation that something is needed

The AI:

- Breadth — Go, React, Swift, SQL in one session
- Speed — 6 working versions in 9 hours
- Context-carrying — migration + handler + component in one thought
- Infinite patience for boring work

PART 4

When You Move This Fast, Testing Is Everything

The App Broke All at Once

The entire app shipped with zero authentication on any API route.

- Adding auth middleware broke everything — routes returned 401s across the board
- The app was unusable. Every screen, every action, every API call.
- External integrations broke silently — Spotify, Discogs, setlist APIs
- Buttons that looked fine didn't work. No error, no feedback, just nothing.
- **The app had grown too complex to verify manually.**

We had tests. 43.9% backend coverage. E2E specs. Barely any of them caught the problem.

Writing the Tests

43.9% → 70.4% backend coverage. E2E expanded from a handful to 27 specs. Cost: nothing.

- I went to bed. Multiple agents worked through the night. Coverage jumped 27 points.
- **API tests — hitting every endpoint to verify auth worked correctly**
- **E2E UI tests — massively expanded the Playwright suite for real user flows**
- 27 Playwright test specs, 9,000+ lines of test code just on the frontend
- Edge cases, error conditions, auth boundary tests
- The AI handles the coverage floor. The tricky tests are still yours.
- **You have no excuse not to have coverage anymore**

Tests Created New Problems

Solving the testing problem created three more.

- **CI burned through GitHub Actions credits. 40+ min to merge. Had to parallelize.**
- **Tests burned 3rd-party API rate limits. Every test run hit real APIs.**
- **Had to build a mock server. ~1,200 lines of mock infrastructure.**
- 22,000+ lines of test code against ~72,000 lines of app code
- On the frontend, test code is 81% the size of the app itself
- A significant chunk of the codebase exists solely to verify the rest

Testing infrastructure is real engineering. It's not a side task.

Mobile-Only Bugs

The web worked. The simulator told a different story.

- Broken buttons — z-index issues causing invisible elements to capture taps
- Login flows that worked on web but failed silently on iOS
- Content rendering under the Dynamic Island
- Keyboard covering input fields on certain screens
- Async race conditions that only appeared on slower mobile connections
- The simulator catches a lot — but you have to actually run it and click through things

Dogfooding Under Real Conditions

Not “I used my app once and it seemed fine.” Use it during the thing it’s built for.

Pigeons Playing Ping Pong show. From my couch.

- 10:33 PM — Chat input cut off on Chrome mobile. Fixed.
- 10:39 PM — Live posts not appearing in main feed. Fixed.
- 11:20 PM — Songs appearing out of order. Duplicate comments. Fixed.
- 11:34 PM — Live page content rendering under the Dynamic Island.

Four bugs in one hour. All found by using the app during a real show.

What the Simulator Can't Tell You

Some things only work on a real device.

- **Push notifications — simulators don't support APNs**
- **Live Activities / Dynamic Island — the real UX only exists on hardware**
- **Actual network conditions — venue WiFi or cellular behave differently**
- **The physical experience — how it feels to hold it, tap with one thumb**
- You need both: simulators for fast iteration, real devices for what matters most

Automating QA — The Goal and Android

Daily automated visual QA — every morning at 8:47 AM, both platforms.

Android: 90 minutes to build. Shockingly easy.

- adb controls the emulator — taps, swipes, screenshots, all built in
- WebView exposes Chrome DevTools Protocol over a WebSocket
- Find the socket in /proc/net/unix, forward with adb forward
- Inject auth tokens via CDP localStorage.setItem()
- Modern, documented, works exactly like you'd expect
- **The entire Android QA automation was done in 90 minutes.**

Automating QA — iOS Was a Nightmare

iOS: 6+ hours to build. Apple exposes almost nothing.

- No CDP. No WebDriver. No programmatic WebView access.
- TCC.db hacking — write to the iOS privacy database to pre-approve permissions
- AppleScript for keyboard — can't type "@" (modifier key). Paste from clipboard.
- Coordinate probing with idb — coordinate system doesn't match visual layout
- JWT extraction from SQLite — read `localStorage.sqlite3` from the app container
- Every step was a hack on top of a hack

Android gave us modern tooling. iOS gave us archaeology.

Automating the Release Pipeline

Some things the simulator can't test. You need real devices — which means real builds.

- Push notifications, Live Activities, App Store behavior — none work in the simulator
- Manually uploading after every change? Not sustainable at this pace.
- **TestFlight pipeline — one script: build, Capacitor sync, Xcode archive, upload.**
- **Google Play pipeline — build, sign, upload to internal testing track.**
- Every release automated end to end
- **The effort to automate was enormous. But the alternative was impossible.**

CI Blocks the Merge, Not You

Every time I bypassed CI, I broke production. Every. Single. Time.

- Never --admin merge to bypass CI
- Never push without tests passing
- Never skip the pipeline because “it’s a small change”
- Three rounds of push-and-pray before doing what should have been step one
- **The pipeline is your last line of defense between “Claude thinks it’s done” and “users see it.”**

The AI Takes the Path of Least Resistance

Even with all of this infrastructure, the AI constantly finds shortcuts.

- Admin merging to bypass CI — tests take 40 min. It knows `--admin` skips the wait.
- Working directly on main — perceived urgent issue? Skip branches, skip PRs.
- Injecting JWT tokens to skip login — login screen never tested.

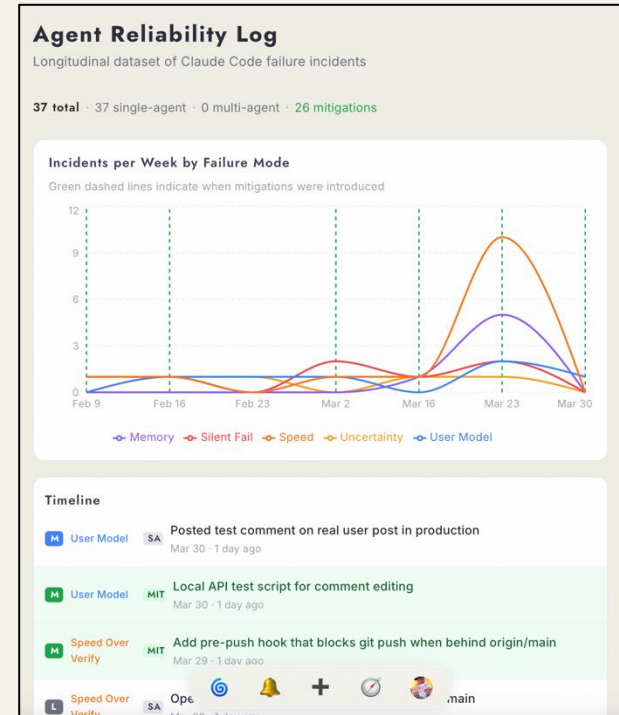
The pattern:

- You build CI → it learns `--admin`
- You build branch protection → it pushes to main
- You build login flows → it injects tokens
- You have to explicitly forbid every shortcut, and then watch it try anyway.

Incident Logging — Making the AI Account

When the AI breaks something, it has to write it down and fix its own process.

- Built an `agent_incidents` table —
 - every mistake logged with failure mode, description
- Then it proposes a mitigation —
 - a concrete code change, rule, or check
- Mitigations tracked: proposed → implemented → verified
- **Force the AI to change its behavior,**
 - **not just fix the immediate problem for this instance**
- Does it work?
 - Mitigations get implemented.
 - Whether they prevent future incidents is an open question



If you can't make the AI learn, at least make it document.

Where We Landed

The workflow today — end to end, mostly automated.

- Users submit bugs into a custom database through the app
- **A triage skill pulls the next bug, analyzes it, and starts fixing it**
- **A QA skill tests across all surfaces — web, iOS simulator, Android emulator**
- **A test skill writes and runs tests for the change**
- **The app is auto-released to TestFlight, Google Play, and production on merge**
- When the AI does something wrong, it **logs an incident and proposes a mitigation**
- Feedback loop: **Focus groups**, etc. to determine feature viability using app “personas”

This is only possible because we spent days building the infrastructure to prevent the AI from doing bad things.

The AI is the best collaborator for the first 80%. The most dangerous for the last 20%.

What Lessons Can We Apply From Last Time?

The Cheat Sheet

1. **Be a product manager.** You steer; the AI writes.
2. **Describe problems, not technical solutions.** Let it explore, but redirect fast in the early game.
3. **Make context persistent.** CLAUDE.md is your most important file outlining the rules and the scars.
4. **Testing is the most important thing you do.** When the AI makes you go fast, nothing else keeps it honest.
5. **Automate everything.** Hooks, CI, QA sweeps, release pipelines.
6. **Build skills for repeatable flows.** Bug triage, QA sweeps, releases — encode them once, run them forever.
7. **Make the AI accountable.** When it breaks something, make it document what happened and how to prevent it, this is key to surviving the late game.

The more you restrict what the AI can do wrong, the more you'll succeed — and the faster you'll be able to go.

Questions?

Christopher Meiklejohn

zabriskie.app

@cmeik

Zabriskie. Where taste resonates.

This slide deck was made by the AI.