

VIBE ARCHITECTURE

When syntax is free,
structure is your only asset.



PETER ISBERG

Dedication

Dedication

To the engineers who spent their weekends fighting the build so the rest of us could ship.

And to the architects who understand that speed without structure is just a faster way to fail.

Table of Contents

Part I — The Paradox of Speed

- **1. The Entropy of the Prompt** — `ch01_entropy_of_the_prompt.md`
- **2. Shift Left, Guide Right** — `ch02_shift_left_guide_right.md`

Part II — Engineering the Guardrails

- **3. Context Diets** — `ch03_context_diets.md`
- **4. Annotation-Driven Design** — `ch04_annotation_driven_design.md`
- **5. Architectural Unit Testing** — `ch05_architectural_unit_testing.md`
- **6. Visualizing the Vibe** — `ch06_visualizing_the_vibe.md`
- **7. Spec-Driven Development: The Spec Is the Guardrail** — `ch07_spec_driven_development.md`
- **8. Spec Kit in Practice: From Prompt to Merged Code** — `ch08_spec_kit_in_practice.md`
- **9. Kiro: When the IDE Is the Spec Engine** — `ch09_kiro.md`
- **10. Zenflow: The Spec as an Executable DAG** — `ch10_zenflow.md`
- **11. Get-Shit-Done: Spec-Driven Plus Context Engineering** — `ch11_get_shit_done.md`
- **12. Skill Forge: Skills Forged From Your Own Codebase** — `ch12_skill_forge.md`

Part III — Handling the Hard Stuff

- **13. The Asynchronous Trap** — `ch13_the_asynchronous_trap.md`
- **14. Schema Drift & API Contracts** — `ch14_schema_drift_and_api_contracts.md`
- **15. The Untrusted Context: Prompt Injection** — `ch15_prompt_injection.md`
- **16. The Attacker Agent: Adversarial Runtime Testing** — `ch16_attacker_agent.md`

- **17. Multi-Agent Orchestration & CI/CD** — `ch17_multi_agent_orchestration_and_cicd.md`
- **18. The Two-AI Gripe: Adversarial Review Pipelines** — `ch18_the_two_ai_gripe.md`
- **19. Prompt Lineage: Git for Intent** — `ch19_prompt_lineage.md`
- **20. Git as a Save-State Engine** — `ch20_git_as_a_save_state_engine.md`

Appendix

- **A. The Lean Context-File Blueprint** — `appendix_a_context_blueprint.md`
- **B. The Quick-Start Checklist** — `appendix_b_quick_start.md`
- **C. Anti-Patterns** — `appendix_c_anti_patterns.md`
- **D. Measuring Vibe Architecture** — `appendix_d_metrics.md`

Back Matter

- Conclusion — `conclusion.md`
 - Epilogue: The Horizon — `epilogue_future_chapters.md`
 - The Zero-Marginal-Cost Baseline — `manifesto.md`
 - About the Author — `about_the_author.md`
 - Glossary — `glossary.md`
 - References — `references.md`
-

About This Book

The core thesis: *Vibe coding gives you speed, but Vibe Architecture gives you scale. When syntax is free, structure is your only asset.*

Lately, "vibe architecture" has been used as a warning — a mocking label for systems where AI lazily strung components together without a plan. This book flips that definition on its head. Vibe Architecture is the intentional, highly disciplined framework that allows vibe coding to scale without collapsing into a mess.

This is a lean, practical handbook for the senior engineer, the tech lead, and the architect asking one question: *How do I let my team use Cursor and Claude Code to move ten times faster, while ensuring our system doesn't turn into a distributed nightmare of unmaintainable, tightly coupled code?*

The answer is a single repeated pattern: human intent in, automated machine enforcement out, zero tax on the flow. Every chapter is one instrument that implements that pattern — architectural unit tests, source-resident annotations, forced-collision concurrency detectors, and Git used as a save-state engine.

The examples lean toward Java and Spring Boot, where the enforcement tooling is most mature. The discipline is language-agnostic.

How to Read This Book

- **Part I — The Paradox of Speed** diagnoses why AI naturally erodes architecture, and reframes the engineer's job around directing structure rather than typing syntax.
- **Part II — Engineering the Guardrails** is the core: context discipline, annotation-driven design, architectural testing, keeping the architecture continuously visible, and the paradigm that unifies them — Spec-Driven Development.
- **Part III — Handling the Hard Stuff** covers the failure modes the guardrails do not fully neutralize — AI-generated concurrency bugs, schema and API contract drift — then scales enforcement into the CI/CD pipeline, puts a second adversarial AI on guard duty so the system reviews itself, and ends with recovering from a destructive prompt.

Companion Reading

- *Vibe Coding: Building Production-Grade Software With GenAI, Chat, Agents, and Beyond* — Gene Kim & Steve Yegge (IT Revolution, 2025)
- *Supercharged Coding with GenAI: From vibe coding to best practices*

Companion Tooling

- **VibeTags** — compile-time AI guardrail annotations for Java:
`https://github.com/PIsberg/vibetags`
 - **async-test** — forced-collision concurrency testing for JUnit 5:
`https://github.com/PIsberg/async-test-lib`
-

Acknowledgments

Acknowledgments

A book like this is never built in isolation. It is the product of hundreds of late-night arguments about architecture, thousands of broken builds, and the collective wisdom of the engineering community trying to figure out how to work alongside our new AI counterparts.

Thank you to the early reviewers who read the messy first drafts and challenged my assumptions, particularly the communities around ArchUnit, Spring Modulith, and modern software architecture. Your rigorous demand for structural integrity in an age of generated code shaped the core of Part II.

To Gene Kim and Steve Yegge, whose conversations around "Vibe Coding" provided the necessary friction and inspiration to articulate what happens when you attempt to scale that vibe to enterprise complexity.

Finally, thank you to the developers who open-source their guardrails and share their failures. The techniques in this book—from contract testing to adversarial AI review—exist because you were willing to document the hard lessons learned in production.

1. The Entropy of the Prompt

"Vibe coding gives you speed, but Vibe Architecture gives you scale. When syntax is free, structure is your only asset."

1.1 The Tipping Point

There is a moment in every AI-assisted project where the magic curdles. For the first few days, vibe coding feels like cheating. You describe a feature, the assistant writes it, the tests go green, and you move on. Velocity is intoxicating. Then the project crosses a threshold — usually somewhere between the tenth and the fiftieth file — and the same workflow that felt like a superpower starts producing something that smells distinctly like spaghetti.

This is not a tooling failure. It is a structural inevitability, and understanding *why* it happens is the foundation of everything else in this book.

An AI coding assistant does not care about your long-term architecture. It cares about making the current prompt work. That is its objective function. When you ask it to "add a discount to the checkout flow," it will find the shortest path from the current state of the code to a state where the discount works. If the shortest path runs directly through three architectural boundaries you spent a week designing, the assistant will run straight through all three without hesitation — and without telling you.

1.2 Why AI Naturally Creates Spaghetti

Entropy is the natural tendency of a system to move toward disorder. In a vibe-coded project, the assistant is an entropy engine pointed directly at your codebase. Every prompt is a small force, and absent a constraint, each force pushes the system toward the lowest-energy state — which is almost never the well-architected one.

[!NOTE] Architect's Note: The Biology of the Context Window The LLM context window degrades in precision remarkably like human working memory. While modern models boast 200,000-token windows, their precision inevitably drops as the context grows. Jamming an entire monolithic repository into the prompt does

not give the AI architectural vision; it just induces artificial cognitive overload. "Context diets" (Chapter 3) are a biological analogy as much as a technical one. Concretely, here is what an unconstrained assistant will happily do as your project grows:

- **Import a database entity directly into your UI layer**, because the entity already has the field the view needs and writing a DTO is "extra work" the prompt did not ask for.
- **Tightly couple two features that should be independent**, because calling the other feature's service directly is fewer lines than introducing an event or an interface.
- **Duplicate a subtly-wrong helper across five packages**, because it generated a fresh one each time rather than discovering the existing abstraction.
- **Reach across a module boundary** to grab a package-private detail, because nothing in the prompt — or the language — told it the boundary was sacred.

A human developer accumulates this kind of debt slowly, and feels the friction with each shortcut. An AI accumulates it at the speed of generation, and feels nothing. The danger of vibe coding in large projects is not that the AI writes *bad* code — line by line, it is often excellent. The danger is that it writes *boundary-violating* code, fluently and tirelessly, until the boundaries no longer exist.

1.3 The Obvious Fix That Backfires: Strict Module Systems

If the problem is boundary violation, the textbook answer in the Java world is the Java Platform Module System (JPMS). Declare your modules, export only what should be public, and the compiler itself will refuse to let the UI layer import a persistence entity. The boundary becomes a law of physics rather than a guideline.

In a hand-written codebase evolving at human speed, this is sound advice. In a vibe-coded project, it backfires.

JPMS turns every new feature into a negotiation with the module path. The assistant adds a class, the build fails because the package is not exported, you (or the assistant, on its next turn) edit `module-info.java`, the build fails again because of a split package,

and the intoxicating velocity you came for has been replaced by a slow grind of `module-info` archaeology. You will find yourself fighting the compiler configuration every single time you want to add a feature. The guardrail works — it stops the AI cold — but it stops *you* just as hard. Teams that try this almost always rip JPMS back out within a sprint, and conclude, wrongly, that architectural enforcement and vibe coding are incompatible.

They are not incompatible. JPMS is simply the wrong instrument: it enforces boundaries at the one moment — every compile, in the build configuration — where friction hurts the vibe most.

1.4 The Thesis: Vibe Architecture

Lately the phrase "vibe architecture" has been used as a warning — a mocking label for systems that AI lazily strung together without a plan. This book flips that definition on its head.

Vibe Architecture is the intentional, disciplined framework that lets vibe coding scale without collapsing into chaos. It is the set of guardrails, signposts, and safety nets that absorb the AI's entropy *without* taxing the human's flow. The rest of this book is a catalogue of those mechanisms. The good ones share a single property: they are written once, by a human, with intent — and thereafter enforced automatically, by a machine, at a moment that does not interrupt the vibe.

You will see this pattern repeatedly:

- An **architectural unit test** (Chapter 5) is written once. It then fails, automatically, in the test run whenever the AI crosses a package boundary. You paste the failure back to the assistant; it corrects itself. You never touched a config file.
- A **source annotation** (Chapter 4) is written once, next to the code it protects. It then regenerates the assistant's guardrail file on every compile. The rule travels with the code.
- A **living architecture diagram** (Chapter 6) is wired into the build once. It then regenerates from the code on every compile, so a dependency edge the AI slipped in is something you see, not something you have to remember.
- A **specification** (Chapter 7) is written once, in version control. The prompt stops being the source of truth; the spec is, and the AI is measured against it.

- A **concurrency detector** (Chapter 13) is configured once. It then forces a race condition the AI couldn't reason about into a reproducible, diagnosable test failure.
- An **adversarial reviewer** (Chapter 18) — a second AI paid to be unhappy — is given the rules once. It then grips on every diff that erodes the structure, so the human is out of the review loop entirely.
- A **micro-commit** (Chapter 20) is made once per working sub-task. It then becomes the save-state you reset to when the next prompt goes sideways.

Each mechanism converts a category of AI entropy into an automated, fast-feedback signal the AI can consume and fix itself. That conversion — human intent in, machine enforcement out, zero flow tax — is the whole discipline.

1.5 Who This Book Is For

This is a practical handbook for the senior engineer, the tech lead, and the architect who are being asked a hard question: *How do I let my team use Cursor and Claude Code to move ten times faster, while ensuring our system doesn't turn into a distributed nightmare of unmaintainable, tightly coupled code?*

The examples lean toward Java and Spring Boot, because that is where the enforcement tooling is most mature and because that ecosystem feels the boundary-erosion problem acutely. But the discipline is language-agnostic. If you are working in TypeScript or Go, the specific annotation processor changes; the principle — single point of human intent, automated machine enforcement, no tax on the flow — does not.

Two companion books are worth reading alongside this one. *Vibe Coding: Building Production-Grade Software With GenAI, Chat, Agents, and Beyond* by Gene Kim and Steve Yegge is the definitive text on leaning into AI-driven prototyping without letting your architecture dissolve. *Supercharged Coding with GenAI: From vibe coding to best practices* covers the practical tactics of pairing fast iteration with proper safety nets. This book is narrower and more opinionated than either: it is specifically about the *architectural guardrails*, and the Java tooling that implements them.

1.6 Where We Go From Here

The book is organized around the lifecycle of an AI-assisted change.

Part I — The Paradox of Speed diagnoses the problem (this chapter) and reframes the engineer's job around it (Chapter 2).

Part II — Engineering the Guardrails is the core. Chapter 3 covers context management so the AI does not lose the plot. Chapter 4 covers annotation-driven design — turning your source code into a set of strict signposts the AI cannot ignore. Chapter 5 covers architectural unit testing — letting your test suite, not your build config, enforce the boundaries. Chapter 6 keeps the architecture *visible* — diagrams that regenerate from the code on every compile, so a human can still hold the mental model when the AI writes twenty files in one prompt. Chapter 7 unifies all of it: Spec-Driven Development, the paradigm where a version-controlled specification — not the chat prompt — is the durable source of truth, and every guardrail in Part II is one clause of it. Chapters 8 through 11 put that paradigm in your hands as four shapes of working tooling: the Spec Kit command loop (Chapter 8), the spec-native Kiro IDE (Chapter 9), Zenflow's executable spec graph (Chapter 10), and the context-engineering framework Get-Shit-Done (Chapter 11). Chapter 12 closes the input side: the agent's own skills, forged from your codebase rather than borrowed from someone else's.

Part III — Handling the Hard Stuff covers the failure modes the guardrails do not fully neutralize, then the climax and the safety net. Chapter 13 is the internal trap: the concurrency bugs the AI writes confidently and wrongly. Chapter 14 is its external mirror: the database migrations and API payload changes the AI makes that break already-deployed consumers. Chapter 15 is when that external boundary turns hostile: prompt injection, where untrusted issue, web, or tool content hijacks the agent's instructions. Chapter 16 is the Attacker Agent, turning the generator's talent against your own contracts. Chapter 17 scales every guardrail up from the local editor into the shared CI/CD pipeline, turning advice into branch-level law. Chapter 18 is the climax — the Two-AI Gripe: a second AI with the opposite objective function stands guard so the system reviews itself and the human is out of the review loop. Chapter 19 version-controls the intent itself — the prompts and system rules — so a breaking change can be traced to where the instruction drifted, not just where the code did. Chapter 20 is the safety net under all of it — Git as a save-state engine that lets you recover when a prompt destroys a working state in two seconds flat.

The next chapter is about the shift in your own role that makes all of this work: moving from typing code to directing it.

2. Shift Left, Guide Right

2.1 The Job Changed While You Were Typing

For two decades the craft of software engineering was, mechanically, the act of typing correct code. Design mattered, review mattered, testing mattered — but the load-bearing skill, the one interviews tested and the one that filled your day, was producing syntactically and semantically correct text in a programming language.

That skill is now largely free. When syntax is free, the scarce, valuable thing is not the code — it is the *structure* the code lives in and the *intent* that shaped it. Your job moved. It shifted left, toward decisions made before any code exists, and it moved up the stack, from typing to guiding.

This chapter is about that move. The mechanics — annotations, architectural tests, concurrency detectors, Git discipline — are the rest of the book. But the mechanics only work if you have internalized the role change first. An architect who still thinks their value is in the keystrokes will fight the AI for the keyboard and lose. An architect who understands their value is in the constraints will let the AI have the keyboard and win.

2.2 Shift Left: Decisions Before Code

"Shift left" traditionally means moving testing and security earlier in the lifecycle. In Vibe Architecture it means something sharper: **the architectural decision must exist, in an enforceable form, before the prompt that would violate it.**

Consider the order-and-billing example. The decision "the `order` package must not depend on the `billing` package" is an architectural constraint. In the old world you might enforce it through code review — a human notices the bad import in a pull request and asks for a change. That is "guide right" only: you catch the violation after it is written, downstream, by inspection.

In a vibe-coded project, review-only enforcement collapses. The AI generates faster than you review, and the violating import is buried in a 400-line diff that also contains the feature you actually wanted. You will approve it. Everyone does.

So the decision has to be shifted left into an artifact that exists *before* the prompt and enforces itself *without* your attention:

```
// An ArchUnit rule, written once, enforced on every test run
classes().that().resideInAPackage("..order..")
    .should().onlyDependOnClassesThat()
    .resideInPackages("..order..", "..common..", "java..");
```

This single rule, committed to the repository, changes the economics. The constraint now precedes every future prompt. When the AI writes `import com.app.billing.BillingService;` inside the `order` package, the test suite goes red — not in review, not in production, but in the same feedback loop the AI is already operating in.

The shift-left move is not "write more tests." It is "encode the architectural decision as the first artifact, so that it is older than every line of feature code that could break it."

2.3 Guide Right: Steering the Generation

Shifting decisions left handles the boundaries you can name in advance. But much of architecture is not a rule you can write down before the work — it is judgment applied during the work. That is where "guide right" comes in: the discipline of steering the AI's generation in real time, through the context and prompts you give it.

Guiding right has three moves, in increasing order of leverage.

Direct the structure, not the syntax. A weak prompt says "write a function to apply discounts." A guiding prompt says "add discount handling. It belongs in the `pricing` package. It must not import anything from `order` or `billing` — depend only on the `PricingInput` value object. Return a `Money`, never a primitive double." You are not describing the code. You are describing the *box the code must fit in*. The AI is extraordinarily good at filling a well-specified box and extraordinarily bad at inventing the box itself.

State the invariant, then ask for the implementation. Instead of "make this thread-safe," say "this class is accessed by the request pool and the scheduler pool simultaneously. It must be safe under concurrent `recordHit` and `snapshot` calls. Use a lock-free strategy backed by atomics; do not introduce a synchronized block." The invariant is the part you, the architect, own. The implementation is the part the AI owns. Naming the invariant explicitly is what keeps the AI from silently trading it away for a simpler-looking solution.

Make the boundary visible in the code itself. The highest-leverage move is to stop repeating yourself. If you find yourself typing "don't change this signature, it's a public contract" into the prompt for the third time, the instruction belongs *in the source*, as an annotation the AI's tooling reads automatically (Chapter 4). Guiding right at scale means converting repeated verbal guidance into permanent, code-resident signposts.

2.4 The New Failure Mode: Beautifully Deceptive Code

There is a specific reason guiding right matters more in the AI era than it did in the pair-programming era, and it is worth stating plainly because it reshapes how you review.

A junior developer writes code that *looks* uncertain when it is wrong. Hesitant naming, a TODO, a comment that says "not sure this is right." The wrongness is legible. An AI assistant writes code that *looks* confident and correct *especially* when it is wrong. It produces clean, idiomatic, well-named, well-commented code that handles the happy path beautifully and contains a subtle race condition, an off-by-one in an edge case, or a silent boundary violation, all wrapped in the same polished prose as the correct parts.

This is the single most important mental adjustment for an AI-native engineer: **fluency is no longer evidence of correctness.** You cannot review AI code the way you reviewed human code, scanning for the parts that look shaky, because nothing looks shaky. The deceptive code and the correct code are stylistically identical.

The consequence is that line-by-line human review, your primary safety net for twenty years, is now your *weakest* one. You cannot out-read the generator. The reliable safety nets are the automated ones — the architectural test that does not care how confident the code looks, the concurrency detector that forces the hidden race into the open, the contract test that fails the instant a frozen signature changes. Human judgment moves from *inspecting the output* to *designing the constraints and the detectors*. That is the entire premise of Part II.

2.5 From Code Monkey to Systems Director

It is worth being honest about what this role change feels like, because it is not purely a gain.

If your professional identity is built on the satisfaction of writing elegant code by hand, vibe architecture will feel like a loss before it feels like a leverage. You will write less code. You will read more diffs. You will spend more time on constraints, tests, and

prompts than on implementations. The dopamine of a clever one-liner is replaced by the slower satisfaction of a system that scaled without collapsing.

The compensation is scope. A systems director who has internalized this discipline can hold a far larger system in their head than a code monkey ever could, because they are no longer tracking implementations — they are tracking constraints, and the constraints are enforcing themselves. You are not directing every keystroke. You are directing the *shape*, and trusting the automated guardrails to keep the generation inside it.

The remaining seventeen chapters are the toolkit for doing exactly that. Part II builds the guardrails: context discipline so the AI does not lose the plot (Chapter 3), source-resident signposts the AI cannot ignore (Chapter 4), architectural tests that enforce the boundaries without taxing the flow (Chapter 5), living diagrams to visualize the architecture (Chapter 6), Spec-Driven Development to unify intent (Chapter 7), and four shapes of spec-driven tooling that put it to work — the Spec Kit command loop (Chapter 8), the spec-native Kiro IDE (Chapter 9), Zenflow's executable spec graph (Chapter 10), the Get-Shit-Done context-engineering framework (Chapter 11), and Skill Forge, which learns the agent's skills from your own codebase (Chapter 12). Part III handles the hard stuff that survives the guardrails: concurrency traps (Chapter 13), schema and contract drift (Chapter 14), prompt injection and untrusted context (Chapter 15), adversarial runtime fuzzing (Chapter 16), multi-agent orchestration (Chapter 17), adversarial AI review pipelines (Chapter 18), prompt lineage that version-controls the intent behind the code (Chapter 19), and the Git discipline that lets you undo a destructive prompt in one command (Chapter 20).