

CLAUDE CODE

The Definitive Guide to Agentic Development



Written by Claude Code

Instructed by V. Korostyshevskiy

About This Book

This book was written entirely by an AI. Not edited by one, not "assisted" by one — written by one. Claude Code, Anthropic's agentic coding tool, produced every chapter, following a detailed set of structural instructions and quality gates designed by Vladimir Korostyshevskiy.

The process worked like this: multiple AI systems — including Claude, Gemini, and Perplexity — were used to gather and research source material. They pulled from Anthropic's official Claude Code documentation, publicly available usage examples, technical analyses, and industry reports about agentic development practices. The usage examples were sourced with particular attention to sell-side and buy-side front office environments — trading UI development, order management systems, real-time data platforms, and the kind of mission-critical desktop and web applications that financial firms build and maintain under relentless delivery pressure. That raw research was then fed to Claude Code with an assignment file specifying the book's structure, voice, audience, deduplication rules, and legal constraints. Claude Code's multi-agent architecture did the rest — research agents read and summarized the sources in parallel, writer agents composed the chapters, editor agents reviewed them against quality criteria, and a final review panel checked for consistency, completeness, and compliance before the manuscript was assembled into the EPUB you're reading now.

No sentence in this book is copied from its source materials. Every idea has been paraphrased, synthesized, and rewritten. No company is mentioned by name other than Anthropic. No competing product is named. No individual is quoted or attributed. No URLs appear in the text. These constraints were deliberate — they make it possible to share this book freely without copyright concerns.

This book is not for beginners. If you've never used Claude Code, start with the official documentation. This book assumes you've been using the tool for months and want to get significantly better at it. It covers the mental models, advanced capabilities, practical patterns, and organizational strategies that separate casual users from people who are genuinely productive with agentic development tools.

The writing voice blends the narrative-driven explanations of Malcolm Gladwell with the direct, irreverent technical prose of Michael Lopp. If a paragraph reads like documentation, something went wrong. If it reads like a blog post that opens with "let's dive in," something went very wrong.

Every effort was made to ensure accuracy, but this is a book written by AI about a rapidly evolving AI tool. Some details may have changed between the time of writing and the time of reading. When in doubt, check the official Anthropic documentation for the current state of any feature or capability.

This book is free to share, distribute, and read. That was the whole point.

— Vladimir Korostyshevskiy, February 2026

Chapter 1: Beyond the Getting-Started Guide

What You'll Learn

Most developers plateau with Claude Code somewhere around day three. They have the basics down -- type a prompt, get some code, copy it into place -- and they stop there. Not because they are doing anything wrong, but because nobody told them there was more. This chapter is for developers who sense there is a deeper layer but have not found the entrance.

Claude Code is an agentic system. It reads your codebase, plans its approach, executes changes across dozens of files, runs your tests, reads the errors, and tries again. It chains tool calls in loops, not single-shot completions. The moment you internalize this distinction, your entire relationship with the tool changes. You stop asking it questions and start giving it assignments.

This chapter lays out the mental model that separates surface-level use from advanced practice. You will understand the agentic execution loop and why it can be built in a few hundred lines of code, why the context window is the only constraint that matters, how to structure sessions for maximum throughput, and why your instinct to micromanage is the single biggest bottleneck to productivity. You will learn how teams classify tasks for autonomous versus supervised execution, why making Claude Code your first stop on every task changes your workflow, and how to think of the tool as an iterative thought partner rather than a vending machine. You will also confront an uncomfortable number: research shows developers can fully delegate only a small fraction -- somewhere between zero and twenty percent -- of their tasks, and pretending otherwise leads to the kind of silent failures that cost more time than they save.

The Machine You Think You're Using

There is a mental model most developers bring to Claude Code, inherited from years of chatbots and autocomplete tools: type question, receive answer. That mental model undersells the tool by an order of magnitude, and it shapes everything about how they use it.

Claude Code is a loop. Not a request-response pair. A loop. The architecture looks like this: read context, form a plan, take an action, verify the result. Then do it again. And again. A single prompt from you might trigger thirty or forty iterations of this cycle, each one reading files, writing code, running commands, checking outputs, and deciding what to do next.

The Agentic Loop in 350 Lines

The harness that runs this loop is deceptively simple. A minimal implementation -- the kind of educational prototype that community developers have built to demystify agentic systems -- is roughly 350 lines of code. The core pattern is three functions: call the model, process the tool calls from the response, and loop. If there are tool results, append them to the messages and call the model again. If there are none, the task is done.

```
response = callApi(messages, systemPrompt)
toolResults = processToolCalls(response.content)
if toolResults.length == 0: return    // done
messages.append(response, toolResults)
// loop again
```

Tools are registered in a map -- read a file, write a file, run a command -- and each tool call from the model gets dispatched to the corresponding handler. The model reasons about what tool to call next. The harness executes the tool and feeds the result back. That is the entire architecture.

In production Claude Code, every step is wrapped in permissions checks, context management, hook evaluations, subagent orchestration, and safety constraints. But the core is still that loop. Understanding this matters because it changes what "prompting" means. You are not writing a query. You are briefing a colleague who has access to your entire codebase, your terminal, your test suite, and your build system. That colleague will go explore, make a plan, implement it, and come back when it is done or when it is stuck. The quality of your initial briefing determines everything.

Context Window: The Only Resource That Matters

Every feature of Claude Code, every architectural decision, every best practice traces back to one constraint: the context window.

Think of it as RAM for the conversation. Every file Claude reads, every command output it captures, every tool result it processes, every message you exchange -- all of it accumulates in this fixed-size buffer. When it fills up, performance does not degrade gracefully. It falls off a cliff. Instructions get forgotten. Work gets dropped. Claude starts solving problems it already solved.

This is not a theoretical concern. It happens in every long session. Auto-compaction kicks in when the window is nearly full, summarizing older content to make room -- but the summarization is lossy. Chapter 3 covers the compaction mechanics and how to survive them.

The practical consequence is that every interaction has a cost measured in context space, and you need to think about that cost the same way you think about compute budgets or memory allocation. Reading a 2,000-line file is expensive. Dumping verbose test output into the conversation is expensive. Having Claude search for something it could find with a targeted glob is expensive.

The developers who get the most from Claude Code are the ones who treat context like a scarce resource. They use subagents to isolate verbose operations (Chapter 4). They keep CLAUDE.md files concise -- under 500 lines (Chapter 3). They compact proactively instead of waiting for auto-compaction to butcher their context at the worst possible moment.

The Four-Phase Workflow

If you have been typing implementation requests directly into Claude Code, you have been skipping the most important step.

The workflow that consistently produces the best results has four phases: explore, plan, implement, commit. Skipping exploration is the single most common mistake, and it leads to Claude confidently solving the wrong problem with technically correct code.

Explore. Start in plan mode. Let Claude read your codebase. Point it at relevant directories. Ask it to understand the existing patterns, the conventions, the architecture. This is cheap -- plan mode restricts Claude to read-only operations, so it cannot accidentally modify anything while it is figuring out the lay of the land.

Plan. Still in plan mode. Ask Claude to outline its approach. Review the plan. Push back on parts that do not match your mental model. This is where misunderstandings surface, and fixing a misunderstanding in a plan costs nothing. Fixing it in implementation costs you a revert and a burned context window.

Implement. Switch to normal mode. Claude now has the context from exploration and a reviewed plan. Implementation quality is dramatically higher because Claude is working from understanding rather than assumption. Let it run.

Commit. Claude will offer to commit when implementation is complete. Review the diff. If you have been doing this well, the diff should be unsurprising. Commit frequently. Small commits are cheap insurance. The checkpoint-and-rollback strategy (discussed later in this chapter) depends on having clean commit boundaries to revert to.

This four-phase workflow is not about being cautious. It is about being efficient. Exploration and planning are cheap. Implementation and debugging are expensive. Front-loading cheap work to reduce expensive work is not caution. It is engineering.

First-Stop Workflow Planning

One team at Anthropic adopted a practice that sounds trivially simple but changed how they worked: they made Claude Code their first stop for every task. Before reading code manually, before searching through the repository, before asking a colleague -- they opened Claude Code and asked it to identify which files to examine for whatever they were working on. Bug fix, feature development, analysis -- it did not matter.

This replaced the traditional time-consuming process of manually navigating the codebase and gathering context before starting work. Claude could scan the entire repository structure, identify relevant files, explain complex interactions between modules, and surface dependencies that a human would need fifteen minutes of directory browsing to find. The team reported that this single habit -- always starting with Claude Code rather than treating it as an optional accelerator -- was the biggest shift in their daily workflow.

Permission Modes Are Workflow Selectors

Most people think of permission modes as a safety dial. Crank it up for danger, crank it down for trust. That framing misses the point.

Permission modes shape how work happens. They are workflow selectors, not just guardrails.

Plan mode restricts Claude to read-only operations. Use it for exploration and research. You are not being cautious; you are being intentional about which phase of the workflow you are in.

Default mode asks permission for each write operation. Use it for surgical changes where you want to approve every edit.

Auto-accept edits mode lets Claude write files without asking but still prompts for shell commands. This is the sweet spot for implementation phases where you trust the plan but want to supervise system-level operations.

Full auto-accept mode lets Claude run without interruption. Use it when you have strong verification in place -- a good test suite, a linter with teeth, pre-commit hooks that enforce standards. The permission model (Chapter 2) governs the full scope of what Claude Code can do without asking.

Cycling between these modes mid-session with Shift+Tab is not a sign of indecision. It is how experienced users move through the explore-plan-implement-commit workflow without starting new sessions. You explore in plan mode, review the plan, toggle to auto-accept, and let Claude execute.

Autonomous Loops and the 80% Handoff

The product development team at Anthropic took the auto-accept workflow further. Their engineers enable full auto-accept mode and set up autonomous loops where Claude writes code, runs the test suite, reads the failures, and iterates continuously. They hand Claude abstract problems they are not familiar with, let it work autonomously, and then review the result when it reaches roughly eighty percent completion. The last twenty percent -- the judgment calls, the design refinements, the edge cases that require domain knowledge -- they finish themselves.

This workflow depends on two things: starting from a clean git state so you can revert the entire run if it goes sideways, and committing as Claude goes so you have intermediate checkpoints. The security engineering team at Anthropic distilled this into a single directive they give Claude at the start of autonomous sessions: "commit your work as you go." Instead of asking targeted questions for code snippets, they let Claude work autonomously with periodic check-ins, and the incremental commits create a trail they can review, cherry-pick, or revert as needed.

The Collaboration Paradox

Here is the uncomfortable number: research shows developers use AI-assisted tools in roughly sixty percent of their work. The fraction they can fully delegate -- hand off the problem and walk away -- is between zero and twenty percent.

This is not a failure of the tooling. It is a feature of the work.

Software engineering is a judgment-intensive activity. Even when Claude Code handles the implementation, a human still needs to define the problem correctly, evaluate whether the solution is appropriate, decide what tradeoffs are acceptable, and catch the cases where technically correct code is strategically wrong. The sixty percent usage figure means Claude is present and useful for most of the day. It does not mean sixty percent of the work is automated. AI serves as a constant collaborator, but using it effectively requires thoughtful setup, active supervision, validation, and human judgment -- especially for high-stakes work.

The practitioners who struggle most are the ones who expect full delegation and get frustrated when it does not materialize. The practitioners who thrive are the ones who embrace the iterative partner model: Claude proposes, you react, Claude adjusts, you approve. It is a conversation, not a handoff.

This applies even within a single task. You might delegate the initial implementation, then supervise the debugging, then take over for a tricky edge case, then hand it back for test coverage. The boundary between human work and machine work is not a clean line. It is a continuous negotiation, and the negotiation itself is part of the work.

The Thought Partner Model

There is another way to think about what Claude Code does, one that extends beyond writing code. Consider the experience of someone who used Claude Code not to build software but to build a financial plan. They fed it spreadsheet exports from multiple accounts, described their goals and constraints, and worked with it iteratively over several sessions to produce a phased optimization plan -- complete with tax impact analysis, fund recommendations, and before-and-after allocation grids. The kind of deliverable you would pay a professional advisor to produce.

The interaction model was not "type a question, get an answer." It was more like working with an analyst who has all the data but has not lived with the accounts. You clarify, reframe constraints, and Claude updates every calculation, table, and recommendation in real

time. It is iterative thought partnership applied to analysis, not just code.

This model -- Claude as iterative analyst, not one-shot oracle -- applies to any domain where the work involves processing data, applying constraints, and refining outputs through feedback. Code is the most common use case, but the mental model transfers directly to data analysis, technical writing, infrastructure planning, and anything else where the quality of the output depends on the quality of the dialogue.

The key insight: the best practitioners give Claude a rough proposal as a starting point rather than a blank slate. A sketch of the approach, a strawman architecture, even a half-formed idea. Claude refines a proposal faster and more accurately than it generates one from nothing. If you have an opinion, state it. If you have a suspicion, share it. The output quality tracks directly to the specificity of the input.

The 70-80% Coverage Gap

Claude Code can build roughly 70-80% of a production stack today. The remaining 20-30% is where experienced developers earn their keep.

Component-Level Assessment

The mistake developers make is thinking about this gap in the abstract. It is not abstract. It is component-specific, and the fit varies dramatically depending on what you are building.

Consider a production application with typical frontend, backend, and infrastructure layers. At the component level, the picture looks something like this:

Where Claude Code excels (build now, ship to production):

- Frontend scaffolding -- component architecture, forms, dashboards, state management. Standard patterns with well-established conventions.
- API scaffolding -- CRUD endpoints, route definitions, middleware, request validation. The patterns are predictable and verifiable.
- Database schema design -- migrations, models, indexing strategies, relationship definitions. Claude generates schemas with proper constraints and immutability patterns.
- Test generation -- unit tests, integration tests, test fixtures. Clear success criteria make this ideal for the agentic loop.
- DevOps configuration -- container definitions, CI pipeline definitions, deployment scripts. Template-driven work with verifiable outputs.
- Documentation -- API docs, README files, code comments, changelog entries. Claude reads the code and describes what it does.

Where Claude Code needs a human partner (hybrid approach):

- Real-time systems with connection management -- Claude scaffolds the structure, but you review connection lifecycle, backpressure handling, and failover logic.
- Performance-critical paths -- Claude generates correct code, but optimizing for tight latency budgets requires measurement and iteration that benefits from human judgment.
- Security-sensitive logic -- authentication flows, encryption, access control. Claude handles the boilerplate, but the security review is yours.
- Domain-specific protocols and standards -- specialized protocols require domain expertise that Claude may approximate but not guarantee.

Where to proceed with caution:

- Novel algorithms without established patterns. Claude will produce something, but without reference implementations it may be subtly wrong.
- Ultra-low-latency components where microseconds matter. Claude does not profile. You do.
- Anything where the failure mode is "it works but it is wrong in a way that is expensive to discover later."

The Quick Decision Matrix

When evaluating whether to build a specific component with Claude Code, run through this checklist:

- **Does it have clear success criteria?** Tests pass, types check, linter is happy? Build it now.
- **Does it require real-time data handling?** Claude builds the component; you own the connection management.
- **Is the latency budget tight?** Claude scaffolds; you profile and optimize.
- **Is there an established pattern?** Claude excels. No established pattern? Proceed carefully.
- **Is correctness non-negotiable and hard to verify?** Human review is mandatory, not optional.

Knowing which category you are looking at is itself a skill. The developers who build intuition for this classification ship faster because they do not waste time fighting Claude on tasks where human judgment is non-negotiable, and they do not waste time hand-writing code that Claude could generate in seconds.

Build Now vs. Wait For

Some capabilities are coming but are not quite production-ready in current tooling. The pragmatic approach is to know the difference:

Build with Claude Code right now: Repository-wide refactors, API scaffolding, test generation, database migrations, deployment configurations, documentation, component architecture, and any task where verification is automated.

Use a hybrid toolchain for now: Inline code suggestions while typing (complement Claude Code with an IDE-integrated autocomplete tool), browser-based end-to-end testing (Claude Code writes the test scripts, you run and validate them), and complex multi-agent coordination that needs more mature UIs.

Wait for improvements: Native CI auto-fix pipelines (the integration is coming), automated issue-to-pull-request bots (the workflow exists in pieces but is not fully polished), and richer multi-agent coordination UIs that let you manage parallel agent work visually.

The developers who ship fastest are the ones who use Claude Code aggressively for what it handles well today, complement it with other tools for current gaps, and avoid building fragile workarounds for capabilities that will be native features within months.

Task Classification: Async vs. Sync

Not every task deserves the same level of attention. Experienced teams develop an intuition for which tasks to run asynchronously -- fire and forget -- and which require synchronous supervision.

Async candidates. Test generation for existing code. Documentation updates. Boilerplate scaffolding. Migration scripts with clear patterns. Dependency updates. Code formatting and lint fixes. These tasks have unambiguous success criteria and low risk of subtle errors. Let Claude run in auto-accept mode, check the result when it is done.

Sync candidates. Core business logic. Security-sensitive changes. Architectural decisions. Performance optimizations. Anything where the failure mode is silent -- code that runs but produces wrong results. These tasks need you watching, reacting, and course-correcting in real time.

The Product Team's Classification Rubric

The product development team at Anthropic formalized this intuition into a rubric. Tasks on the product's periphery -- a new settings panel, a prototype feature, a peripheral UI component -- go into auto-accept mode. Abstract problems the developer is unfamiliar with are also good candidates, because Claude's exploratory attempts in autonomous mode often surface the right approach faster than human speculation.

Tasks touching core business logic, critical user-facing features, or anything where style guide compliance and architectural consistency matter? Synchronous supervision. The developer watches, interrupts when needed, and keeps Claude aligned with the broader design intent.

The classification is not fixed. A task that starts async might become sync when Claude encounters an unexpected edge case. A task that starts sync might become async once you are satisfied with the approach and just need it applied across many files. The ability

to toggle between these modes -- both in your own attention and in Claude's permission settings -- is what makes the workflow fluid.

Try One-Shot First, Then Collaborate

The reinforcement learning engineering team at Anthropic distilled their workflow into a simple escalation pattern: give Claude a quick prompt and let it attempt the full implementation first. If it works -- about one-third of the time -- you have saved significant time. If it does not, switch to a more collaborative, guided approach.

This sounds obvious, but the default instinct for many developers is the opposite: over-specify the prompt, try to anticipate every edge case, and front-load excessive detail. The one-shot-then-collaborate pattern works better because it lets you discover what Claude actually struggles with rather than guessing. The first attempt produces data -- you see where Claude's understanding breaks down, and your subsequent guidance is targeted instead of speculative.

The one-third success rate is not a quality metric. It is a feature of the workflow. When the one-shot succeeds, you saved thirty minutes of planning. When it fails, you spent two minutes and gained information about where to focus your collaboration. The expected value is strongly positive in both cases.

Session Strategy

Sessions are not persistent environments. Each new session starts with a fresh context window. Everything Claude learned about your codebase, your preferences, your project's quirks -- gone.

This is not a bug. It is a design constraint that shapes how you should work.

Continue. Resume the last session in the current directory. Use this when you were interrupted and need to pick up where you left off. The full conversation history is restored into the context window.

Resume by ID or name. Useful for long-running projects with multiple parallel workstreams. Name your sessions meaningfully so you can find them later.

Fork. Create a new session that starts with the full history of an existing one but diverges from that point. The original session is preserved. Use this when you want to explore an alternative approach without losing your current progress.

Fresh session. Start clean. This is often the right choice when your previous session's context is mostly spent on exploration that is no longer relevant. Persistent knowledge goes in CLAUDE.md (Chapter 3) and in your project's task system, not in session history.

The key insight is that session management is context management. A session that has been running for hours is a session with a full context window, which means degraded performance. Starting fresh with a good CLAUDE.md file is often faster than continuing a bloated session, because Claude gets your critical instructions at full fidelity instead of through the haze of auto-compacted summaries.

The Iterative Partner Model

One-shot expectations are the root of most frustration with Claude Code. You type a detailed prompt, expect perfect output, and feel disappointed when the result is 80% correct instead of 100%.

Flip the expectation. Plan for iteration.

The first pass is a draft. Review it. Tell Claude what is wrong. Tell Claude what is right. Give it the failing test output. Show it the linter errors. The agentic loop is designed for this. Each correction sharpens Claude's understanding of what you actually want, and because the conversation history is in the context window, Claude does not repeat its mistakes within a session.

Multi-session iteration works the same way, just with CLAUDE.md as the persistence layer instead of conversation history. After a productive session, ask Claude to summarize what it learned and suggest updates to your CLAUDE.md file. The next session starts smarter. Over five or ten sessions, the CLAUDE.md file accumulates enough project-specific knowledge that first-pass quality improves dramatically.

Self-critique prompting (Chapter 8) accelerates this loop further -- asking Claude to evaluate its own output surfaces issues the initial generation missed.

The iterative model also changes how you think about failure. A wrong first attempt is not a failure. It is data. It tells Claude something about what you want that your original prompt did not. The developers who internalize this stop feeling frustrated by imperfect output and start feeling productive because each iteration is fast and each one gets closer.

Commit Frequently, Revert Without Hesitation

The single most important operational practice for working with Claude Code is this: start from a clean git state and commit frequently. Small commits are cheap insurance that give you rollback points when Claude goes down the wrong path.

Why? Because reverting is cheaper than correcting. When Claude makes a wrong turn, your instinct is to explain the problem and ask it to fix its own work. Often, this makes things worse -- each correction attempt consumes context, accelerating degradation. A revert to the last good commit and a fresh prompt is almost always faster. Chapter 10 covers the full checkpoint-and-rollback recovery strategy and the "slot machine" workflow that teams use to formalize this practice.

One-Third Reality

Roughly one-third of tasks succeed on the first attempt without additional guidance. This is not a quality problem -- it is the expected behavior of an iterative system. The correction-and-retry cycle is fast, and the first-attempt success rate improves as your project's verification infrastructure (tests, types, linting) matures. Chapter 10 covers this statistic in depth and its implications for how you should structure your workflow.

Prompt Suggestions: The Hints You're Ignoring

After Claude finishes a response, grayed-out follow-up suggestions appear below the output. Most developers ignore them. That is a mistake.

These suggestions are not random. They are generated based on Claude's assessment of what the logical next step should be, given the current conversation context and the state of the codebase. They surface actions you probably want but might not think to ask for: running the tests after a refactor, checking for edge cases after an implementation, updating related files after a change to a shared interface.

The suggestions are particularly valuable early in a session, when you are still orienting to the problem. They function as a lightweight checklist of things Claude noticed but did not act on unprompted. Hitting Tab to accept a suggestion is faster than formulating the next prompt yourself, and the suggestion is often better-scoped than what you would have typed because Claude has the full context of what it just did.

The anti-pattern is accepting suggestions blindly. They are hints, not orders. Read them. If the suggestion matches what you would have done next, accept it. If it does not, ignore it and direct Claude yourself. The value is in the time saved when the suggestion is right, not in deferring your judgment entirely.

When Simplicity Beats Sophistication

Claude has a tendency toward over-engineered solutions -- a failure mode covered in detail in Chapter 10. The short version: tell Claude to keep it simple, explicitly constrain it toward straightforward implementations, and put simplicity directives in your CLAUDE.md file.

This connects to a broader principle: Claude Code responds to constraints better than it responds to freedom. A vague prompt produces vague output. A prompt with specific constraints -- which files to touch, which patterns to follow, which tradeoffs to make, which frameworks to avoid -- produces focused, appropriate code. More on this in the prompt craft chapter (Chapter 8).

Key Takeaways

- Claude Code is an agentic loop (read-plan-act-verify) built on a pattern so simple it fits in 350 lines -- but wrapped in production-grade safety, permissions, and context management that makes the difference.
- The context window is the single most important constraint -- every file read, command output, and message consumes finite space, and performance degrades sharply when it fills.
- The explore-plan-implement-commit workflow front-loads cheap work (exploration, planning) to reduce expensive work (debugging, reverting).
- Permission modes are workflow selectors: use plan mode for exploration, auto-accept for trusted implementation, and default for surgical changes. Shift+Tab to cycle between them mid-session.
- Research shows developers can fully delegate only 0-20% of tasks; the real value is iterative collaboration, not fire-and-forget automation.
- Assess Claude Code fit at the component level: frontend scaffolding and test generation are immediate wins; real-time systems and security-critical logic need human partnership.
- Try one-shot first, then collaborate -- the one-third immediate success rate is a feature, not a bug, and failed first attempts produce the information you need for targeted guidance.
- Commit frequently and revert without hesitation; reverting to a clean state and retrying is almost always faster than patching a wrong approach.
- Explicitly constrain Claude toward simplicity; without constraints, it defaults to over-engineered solutions.

Chapter 2: The Permission and Trust Architecture

What You'll Learn

Every capable tool creates a tension: the more it can do, the more damage it can cause. Claude Code runs bash commands, edits files, makes network requests, and orchestrates subagents across your entire codebase. The permission system that governs all of this is not a simple on/off switch. It is a layered architecture with five scopes, three evaluation phases, OS-level sandboxing, hook-based extensibility, and a checkpoint system that silently does not cover everything you think it covers.

Most developers encounter this system as a series of confirmation prompts they click through. That is like encountering a car's safety systems by hitting things. What follows breaks down how the trust architecture actually works -- from the organizational policies that override everything you configure locally, to the glob patterns that decide whether a specific tool invocation gets denied, asked, or auto-approved. You will understand why sandbox isolation makes auto-approve safe in some contexts and dangerous in others, how the complete hook event system gives you programmable control over every phase of the agentic loop, why three distinct hook types (command, prompt, and agent) exist and when each is appropriate, and why the devcontainer model has a credential exfiltration problem that no configuration can fully solve.

You will also learn how agentic systems are reshaping security itself -- democratizing expertise so that any engineer can perform reviews that once required specialists, while simultaneously equipping adversaries with the same capabilities. The organizations that win this race are the ones that embed security into their agentic workflows from the start.

After reading this, you will be able to design a permission configuration that matches your actual risk tolerance -- not the default one Claude Code shipped with, and not the wide-open one you switched to because the prompts annoyed you.

The Five-Scope Hierarchy

Claude Code's settings follow a strict precedence order. Understanding this order is the difference between configuring something that works and configuring something that silently gets overridden.

The hierarchy, from highest to lowest priority:

1. **Managed** -- IT-deployed policies in system directories. Cannot be overridden by anything.
2. **CLI arguments** -- Flags passed at invocation time. Override everything below.
3. **Local** -- `.claude/settings.local.json` . Per-machine, gitignored. Your personal overrides.
4. **Project** -- `.claude/settings.json` . Checked into version control. Shared with the team.

5. **User** -- `~/ .claude/settings.json` . Your global defaults.

The critical insight is at the top. Managed settings exist so that an organization can enforce policies that individual developers cannot circumvent. If your IT team deploys a managed-settings.json that denies access to a tool, no amount of local or project configuration will re-enable it. The hierarchy is not a suggestion. It is enforcement.

This creates two distinct experiences of Claude Code. Individual developers working on personal projects interact mostly with User and Project scopes. Engineers inside an enterprise may discover that certain capabilities are locked down before they ever see them. Both experiences are intentional.

Where Settings Live

User settings reside in `~/.claude/settings.json` -- your global defaults that follow you across every project. Project settings live in `.claude/settings.json` within a repository and get committed to version control, meaning the whole team shares them. Local settings go in `.claude/settings.local.json` in the same `.claude/` directory but are gitignored by convention, giving each developer machine-specific overrides without polluting the shared configuration.

Managed settings are the outlier. They sit in system directories controlled by IT -- locations where regular users do not have write access. They are the organizational trump card, and they are designed to be exactly that.

Permission Rule Evaluation

Within each scope, permission rules follow a three-phase evaluation with first-match-wins semantics. This sounds simple. It has sharp edges.

The phases evaluate in this order:

1. **Deny** -- If a tool invocation matches any deny rule, it is blocked. Period. No further evaluation.
2. **Ask** -- If it matches an ask rule, the user is prompted for approval.
3. **Allow** -- If it matches an allow rule, it proceeds without prompting.

If nothing matches, the default behavior for the tool applies.

Rules use a `Tool` or `Tool(specifier)` syntax with glob pattern support. You can write `Bash(npm test)` to match a specific bash command, `Bash(rm *)` to match destructive deletions, or `Write(*.env)` to match writes to environment files. The glob patterns make this system expressive and also fragile -- an overly broad pattern in a deny rule can silently block operations you intended to allow.

First Match Wins

This is where developers get tripped up. If you have a deny rule for `Bash(rm *)` and an allow rule for `Bash(rm -rf node_modules)`, the deny rule fires first. Your allow rule never executes. The ordering within each phase matters, and the phase ordering (deny before ask before allow) is immutable.

The practical implication: write your deny rules narrow and your allow rules broad. A deny rule is a hard wall. Make sure it blocks only what you intend.

Sensitive File Exclusion

The `permissions.deny` configuration replaces an older mechanism called `ignorePatterns`. It provides explicit file-level access control, and you should use it for anything you would not want an AI agent reading or modifying.

The standard targets: `.env` files, credentials directories, private keys, secrets managers, API key configuration. Any file that would be dangerous if its contents appeared in an API request to a cloud service. Claude Code processes file contents through external APIs. If a file contains production database credentials, reading that file sends those credentials to an API endpoint. The `permissions.deny` list is your firewall for that.

Sandbox Isolation

Claude Code provides OS-level sandboxing for bash commands. This is not a conceptual boundary. It is an operating system enforcement mechanism that restricts what processes spawned by Claude Code can actually do.

The sandbox restricts network access to a configurable domain allowlist. By default, Claude Code can reach the domains it needs for its own operation, but arbitrary network access is blocked. You can add domains to the allowlist for your specific workflow -- your package registry, your internal APIs, your cloud provider endpoints.

The sandbox also supports Unix socket path allowlists, local port binding controls, and command-level exclusions. It is granular enough to express "allow npm install but block curl to arbitrary URLs."

Auto-Approve Under Sandbox

Here is where the sandbox changes the trust calculus. When sandboxing is active, Claude Code can auto-approve bash commands because the blast radius is contained. A sandboxed `rm -rf /` is still destructive to local files but cannot exfiltrate data to an external server. A sandboxed `curl` cannot reach domains outside the allowlist.

This is the design intent behind the sandbox: make autonomous operation safe by constraining the environment rather than constraining the agent. It trades permission prompts for isolation boundaries. For many workflows, this is the right trade.

But the sandbox has an escape hatch. By default, if Claude Code needs to run an unsandboxed command, it can prompt the user. Managed settings can disable this escape hatch entirely, which is the enterprise-correct configuration when you want hard guarantees.

Checkpoints: The Safety Net with Gaps

Before every file edit, Claude Code snapshots the affected files. If something goes wrong, pressing Escape twice rewinds both the code and the conversation to the checkpoint state. It is an undo system, and a good one, for what it covers -- direct file edits made through Claude Code's Write and Edit tools.

Checkpoints have significant gaps, most notably that file changes made through bash commands are invisible to the system. The full catalog of checkpoint limitations and their implications is covered in Chapter 10. The key implication for your security posture: if you are doing anything that modifies the filesystem through bash, your safety net is git, not checkpoints. Commit before you start. Commit frequently. Revert when needed.

Devcontainer Security

Development containers provide network-level isolation with a default-deny firewall. The devcontainer configuration restricts which hosts Claude Code can reach, creating an environment where `--dangerously-skip-permissions` becomes less dangerous because the container itself limits the blast radius.

This is the enterprise answer to "how do I let Claude Code run autonomously in CI without a human approving every command." You do not remove the permission system. You wrap the execution environment in network isolation so that even unrestricted operation cannot reach things it should not reach.

Anthropic provides a reference devcontainer implementation -- a container definition with a custom firewall restricting network access, session persistence, and editor integration preconfigured. It is designed as a starting point for organizations that need to run Claude Code in isolated environments, whether for CI pipelines, security-sensitive projects, or headless autonomous workflows.

The Exfiltration Caveat

Here is the hard truth the documentation acknowledges directly: a devcontainer cannot prevent credential exfiltration from a malicious project. If a project's code contains instructions that cause Claude Code to read credentials and encode them in output that reaches an allowed endpoint, the container's firewall does not help. The data leaves through an approved channel.

This is not a theoretical concern. Supply chain attacks that embed instructions in code comments or README files are a known vector. A devcontainer protects against accidental network access. It does not protect against intentional abuse by project code that the agent executes.

The mitigation is layered defense: sensitive file exclusion via `permissions.deny` to prevent reading credentials in the first place, combined with container isolation to limit where data can go. Neither layer alone is sufficient. Together they raise the bar substantially.

The Complete Hook Event System

Hooks are the extensibility layer of the permission system. They are not limited to the two events most developers discover first. The full hook event catalog covers every significant moment in the agentic lifecycle, and understanding all of them unlocks control that permission rules alone cannot achieve.

Hook Event Catalog

Here is every hook event, what triggers it, and whether it can block the action:

PreToolUse -- Fires before any tool executes. Can block the tool call, allow it without prompting, or escalate to the user. Can also modify the tool's input before execution. This is the workhorse of the hook system.

PostToolUse -- Fires after a tool succeeds. Cannot block (the tool already ran), but can provide context to Claude about the result. Use it for logging, linting, or injecting additional information.

PostToolUseFailure -- Fires after a tool execution fails. Receives the error information and whether the failure was caused by user interruption. Useful for custom error reporting or triggering fallback actions.

PermissionRequest -- Fires when Claude Code is about to show a permission dialog to the user. Distinct from PreToolUse -- this fires specifically when the user would be prompted. Your hook can auto-approve, auto-deny, modify the tool input, or apply permission rules equivalent to the "always allow" options the user would see in the dialog. This effectively replaces the interactive prompt with programmatic evaluation.

Notification -- Fires when Claude Code sends notifications, matching on notification type: permission prompts, idle prompts, authentication events, and elicitation dialogs. Cannot block notifications but can inject context into the conversation.

SubagentStart -- Fires when a subagent is spawned. Cannot block subagent creation but can inject additional context into the subagent's starting state. Matches on agent type: built-in agents or custom agent names from your `.claude/agents/` directory.

SubagentStop -- Fires when a subagent finishes. Can block the subagent from stopping, forcing it to continue working. Useful for quality gates that verify a subagent's output before accepting it.

TeammateIdle -- Fires when an agent team teammate is about to go idle. Exit code 2 prevents the teammate from going idle and feeds your stderr message back as feedback. Does not support matchers -- fires on every occurrence. Use this to enforce quality gates, like checking that build artifacts exist before allowing a teammate to stop working.

TaskCompleted -- Fires when a task is marked as completed, either through an explicit update or when a teammate finishes with in-progress tasks. Exit code 2 blocks completion and feeds stderr back as feedback. Use this to enforce completion criteria -- run the test suite, verify lint passes, check that the task's deliverables exist.

PreCompact -- Fires before context compaction. Matches on trigger type: `manual` (user ran `/compact`) or `auto` (context window filled). Receives any custom instructions the user passed to `/compact`. Cannot block compaction but can perform preparation -- logging, saving state, or injecting context that should survive the compaction.

SessionEnd -- Fires when a session terminates. Receives a reason code: `clear`, `logout`, `prompt_input_exit`, `bypass_permissions_disabled`, or `other`. Cannot block termination but is valuable for cleanup tasks, logging session summaries, or triggering post-session workflows.

Three Hook Types

Not all hooks are shell scripts. Claude Code supports three distinct hook handler types, each suited to different verification needs:

Command hooks (`type: "command"`) are the default. They run a shell command, receive the event's JSON input on stdin, and communicate results through exit codes and stdout. Use these for deterministic checks: pattern matching, file existence, command validation. Fast and predictable.

Prompt hooks (`type: "prompt"`) use an LLM to evaluate the action. Instead of executing a script, the hook sends the event input along with your prompt to a fast model (Haiku by default) and receives a structured yes/no decision. This is for checks that require judgment rather than pattern matching -- evaluating whether a code change follows architectural conventions, whether a bash command is appropriate for the current task, or whether Claude should be allowed to stop working.

```
{  
  "type": "prompt",  
  "prompt": "Evaluate if Claude should stop: $ARGUMENTS. Check if all tasks are complete, no errors  
remain, and no follow-up is needed."  
}
```

The `$ARGUMENTS` placeholder gets replaced with the hook's JSON input. The model returns `{"ok": true}` to allow or `{"ok": false, "reason": "..."}` to block, with the reason fed back to Claude as its next instruction.

Agent hooks (`type: "agent"`) are like prompt hooks but with multi-turn tool access. Instead of a single LLM call, an agent hook spawns a subagent that can read files, search code, and inspect the codebase to verify conditions. The subagent runs for up to 50 turns before returning its decision. Use these when verification requires inspecting actual files or test output, not just evaluating the hook input data alone.

```
{  
  "type": "agent",  
  "prompt": "Verify that all unit tests pass. Run the test suite and check the results. $ARGUMENTS",  
  "timeout": 120  
}
```

The default timeout for agent hooks is 60 seconds (vs. 30 for prompt hooks and 600 for command hooks). The response schema is the same as prompt hooks: `{"ok": true}` to allow, `{"ok": false, "reason": "..."}` to block.

Prompt and agent hooks support these events: PreToolUse, PostToolUse, PostToolUseFailure, PermissionRequest, UserPromptSubmit, Stop, SubagentStop, and TaskCompleted. TeammateIdle does not support prompt or agent hooks.

Async Hooks

By default, hooks block Claude's execution until they complete. For long-running tasks -- test suites, deployments, external API calls - this creates an unacceptable delay. Async hooks solve this.

Set `"async": true` on a command hook to run it in the background. Claude continues working immediately. When the background process finishes, any `systemMessage` or `additionalContext` in the hook's JSON output gets delivered to Claude on the next conversation turn. Async hooks cannot block or control behavior -- the action they would have controlled has already completed by the time they finish.

The classic use case: a PostToolUse hook that runs your test suite after every file write. The tests run in the background while Claude continues working. If tests fail, the failure message appears in Claude's context on the next turn, prompting it to address the issue.

```
{  
  "PostToolUse": [{
```

```

"matcher": "Write|Edit",
"hooks": [
  {
    "type": "command",
    "command": "$CLAUDE_PROJECT_DIR/.claude/hooks/run-tests-async.sh",
    "async": true,
    "timeout": 300
  }
]
}

```

Each async execution creates a separate background process. There is no deduplication -- if Claude writes five files in rapid succession, the hook fires five times.

Hook Input and Output

Every hook receives a JSON payload on stdin with common fields: `session_id`, `transcript_path`, `cwd`, `permission_mode`, and `hook_event_name`. Event-specific fields are added depending on the hook type -- tool hooks get `tool_name` and `tool_input`, subagent hooks get `agent_id` and `agent_type`, and so on.

The exit code from your hook tells Claude Code what to do:

Exit 0 means success. Claude Code parses stdout for JSON output fields. For most events, stdout is only shown in verbose mode.

Exit 2 means block. The exact behavior depends on the event: PreToolUse blocks the tool call, PermissionRequest denies the permission, Stop prevents Claude from stopping, TeammateIdle keeps the teammate working, TaskCompleted prevents the task from closing. The stderr message is fed to Claude as feedback.

Other exit codes are treated as hook errors and ignored. Claude Code continues as if the hook did not fire.

JSON output supports several fields across all events: `continue` (boolean, override the exit code decision), `stopReason` (for Stop hooks), `suppressOutput` (hide stdout from verbose mode), `systemMessage` (added to Claude's context on the next turn), and `additionalContext` (added immediately).

PreToolUse: Modifying Input Before Execution

PreToolUse hooks have a unique capability: they can modify the tool's input before it executes. The `updatedInput` field in the hook's JSON output replaces specific fields in the tool's input parameters.

Combine `updatedInput` with `"permissionDecision": "allow"` to silently rewrite and auto-approve a command. Combine it with `"permissionDecision": "ask"` to show the modified input to the user for approval.

This is flexible and exactly as dangerous as it sounds. A hook that rewrites bash commands can add safety flags, prepend logging, or wrap commands in monitoring harnesses. It can also introduce vulnerabilities if the rewriting logic is flawed. Treat `updatedInput` hooks with the same security scrutiny you would apply to any code that rewrites executable commands.

Hook Handler: The "once" Field

For hooks defined in skills (not in settings files), the `once` field causes the hook to run only once per session and then remove itself. This is useful for skill initialization -- a hook that runs setup logic on the first tool invocation and then gets out of the way.

The /hooks Interactive Menu

The `/hooks` command opens an interactive menu showing all configured hooks with labels indicating their source: `[User]`, `[Project]`, `[Local]`, `[Plugin]`. From this menu you can view hook details, add new hooks, delete existing ones, and toggle the `disableAllHooks` setting. It is the quickest way to audit what hooks are active and where they came from.

Hook Security Model

Hooks are snapshotted at startup. Claude Code captures the state of all configured hooks when the session begins and uses that snapshot for the entire session. If someone -- or something -- modifies hook files on disk mid-session, Claude Code detects the change and displays a warning. The modified hooks do not take effect until you review them in the `/hooks` menu. This prevents malicious or accidental hook modifications from silently changing agent behavior during an active session.

Hooks run with full user permissions. A hook is an arbitrary command that executes on your machine with your credentials, your filesystem access, your network access. There is no sandbox around hooks.

This means hook code requires the same security scrutiny as any other code that runs with your permissions. Input validation matters. Shell quoting matters. Path traversal prevention matters. A hook that naively interpolates tool input into a shell command is a code injection vulnerability.

For enterprises, the `allowManagedHooksOnly` setting in managed configuration restricts hooks to those defined in the managed scope. Individual developers and project-level settings cannot add hooks. This prevents a compromised project from installing hooks that exfiltrate data or modify tool behavior.

Hook Patterns for Security

The most common security-oriented hook patterns:

Block Destructive Commands

A PreToolUse hook on Bash that reads the command from the JSON input, checks for dangerous patterns, and returns a deny decision:

```
#!/bin/bash
COMMAND=$(jq -r '.tool_input.command')
if echo "$COMMAND" | grep -q 'rm -rf'; then
    jq -n '{
        hookSpecificOutput: {
            hookEventName: "PreToolUse",
            permissionDecision: "deny",
            permissionDecisionReason: "Destructive rm -rf blocked by hook"
        }
    }'
    exit 0
fi
exit 0 # allow the command
```

When Claude tries to run a destructive command, the hook intercepts it, returns the deny decision, and Claude sees the reason in its context. It will either find an alternative approach or explain why it needs the command.

Lint on Write

A PostToolUse hook that triggers your linter whenever Claude edits a file:

```
{
  "PostToolUse": [
    {
      "matcher": "Edit|Write",
      "hooks": [
        {
          "type": "command",
          "command": "$CLAUDE_PROJECT_DIR/.claude/hooks/check-style.sh"
        }
      ]
    }
  ]
}
```

```
  }]
}
```

The linting script runs after each edit, and any failures appear in Claude's context, prompting it to fix the style issues before moving on. This turns your linter into a real-time quality gate for the agentic loop.

Sound Effects as Session Feedback

Not every hook is about security. One creative pattern uses hooks to play audio cues on lifecycle events -- a sound when the session starts, when you submit a prompt, when Claude finishes, and when context compacts. The implementation is a one-line command hook on each event that runs the system audio player in the background (with `&` to avoid blocking). On macOS, that is `afplay` ; on Linux, `aplay` or `paplay` .

This sounds frivolous, but developers who use it report that the audio feedback helps them maintain awareness of Claude's state while working on something else during autonomous runs. You hear when it finishes. You hear when context compacts (a cue to check whether important instructions survived). It turns background awareness from a tab-switching chore into a passive sensory channel.

Security Validation Skill Hook

Skills can embed hooks in their frontmatter. A security-focused skill might include a PreToolUse hook that runs a validation script before every Bash command executed while the skill is active:

```
---
description: Security review workflow
hooks:
  PreToolUse:
    - matcher: "Bash"
      hooks:
        - type: command
          command: "./scripts/security-check.sh"
---
```

The hook is scoped to the skill's lifecycle -- it only runs when the skill is active, and it removes itself when the skill completes.

Multi-Criteria Stop Prompt Hook

A prompt hook on the Stop event that evaluates three conditions before allowing Claude to finish:

```
{
  "Stop": [
    {
      "hooks": [
        {
          "type": "prompt",
          "prompt": "You are evaluating whether Claude should stop working. Context: $ARGUMENTS\n\nAnalyze the conversation and determine if:\n1. All user-requested tasks are complete\n2. Any errors need to be addressed\n3. Any follow-up actions were mentioned but not completed\n\nIf ALL conditions are satisfied, return {\"ok\": true}. If ANY are not met, return {\"ok\": false, \"reason\": \"specific explanation\"}."}
      ]
    }
  }
}
```

This uses a fast model to evaluate whether Claude's work is actually done before letting it stop. If the evaluation finds incomplete work, Claude receives the reason and continues.

Audit Logging

A PreToolUse hook that logs every tool invocation with timestamps, inputs, and the decision to a file or monitoring service. Zero runtime impact on the developer, full visibility for security teams.

File Access Control

A PreToolUse hook on Write and Edit that restricts modifications to specific directory trees, enforcing that Claude Code only touches code in the current project. Combined with `permissions.deny` for read restrictions, this creates a comprehensive file-level access boundary.

MCP Tool Gating

PreToolUse hooks matching `mcp_server_tool` patterns that log or restrict access to external service integrations. MCP tools follow the naming pattern `mcp_<server>_<tool>`, and you can use regex in matchers -- `mcp_memory_.*` matches all tools from a memory server, `mcp_.*_write.*` matches write operations across all MCP servers.

Settings Scope Characteristics

Each settings scope has distinct properties that matter for team workflows:

User scope (`~/claude/settings.json`): Personal. Follows you everywhere. Good for your editor preferences, your personal permission rules, your default model selection. Never shared.

Project scope (`.claude/settings.json`): Shared via version control. The team's agreed-upon configuration. Hooks, permissions, plugin enablements, environment variables that everyone needs. Changes go through code review like any other committed file.

Local scope (`.claude/settings.local.json`): Per-machine overrides. Gitignored. Your personal tweaks on top of the project configuration. Machine-specific paths, personal API keys for development services, tools you want that the team configuration does not include.

Managed scope: IT-controlled. Invisible to most developers. Enforces organizational policy. Cannot be inspected or overridden from the CLI. The existence of this scope is why "it works on my machine but not on my work machine" is a real conversation in enterprise Claude Code deployments.

The interplay between these scopes is where configuration gets interesting. A project might allow a tool. A managed policy might deny it. The managed policy wins, silently. No error message, no explanation -- the tool simply is not available. If you are debugging why something works on personal projects but not at work, check whether managed settings exist.

MCP as Security-Controlled Access

When Claude Code needs to interact with external services, there are two paths: run a CLI command via bash, or use an MCP server.

The MCP path is better for security. An MCP server provides structured tool interfaces with defined inputs and outputs. Every invocation goes through the permission system, can be intercepted by hooks, and produces structured logs. A bash command running `curl` with credentials provides none of that visibility -- credentials appear in the command string, in the conversation context, and potentially in API calls to the model provider.

The security-conscious approach to Claude Code integration is: build MCP servers for your sensitive data sources. Do not give Claude bash access to your production databases. Give it an MCP tool that wraps the database access with proper authentication, logging, and query restrictions. Chapter 5 covers the full architecture of MCP and the detailed security case for MCP over CLI access.

One legal team at Anthropic raised an important concern from the product counsel perspective: as MCP integrations reach deeper into sensitive systems, conservative security postures will create barriers. The more capable the tooling becomes, the more access it needs, and the more friction security policies generate. Organizations should expect this tension and plan for it -- building compliance tooling proactively rather than reactively, establishing MCP access governance frameworks before the integrations proliferate.

Security Dual-Use Risk

Agentic coding tools are dual-use technology. The same capabilities that help a security team audit code and generate patches also help an attacker find vulnerabilities and generate exploits. Claude Code's ability to read an entire codebase, understand its architecture, and make targeted changes is exactly as useful for defense as it is for offense.

This is not a hypothetical. Security researchers use Claude Code to find bugs. Attackers can use it the same way. The democratization of expertise that makes Claude Code valuable for non-security-engineers to write secure code also makes it valuable for non-security-engineers to find insecure code.

Agentic Cyber Defense

The trends are moving in two directions simultaneously. On the defensive side, agentic systems are making it possible for any engineer to perform security reviews, hardening, and monitoring that previously required specialized expertise. Security knowledge is becoming democratized -- an engineer without a security background can now use Claude Code to audit code paths, identify common vulnerability patterns, and generate hardened implementations.

On the offensive side, threat actors will scale attacks using the same tools. The prediction from industry analysis is clear: automated agentic systems will enable security responses at machine speed, automating detection and response to match the pace of autonomous threats. The organizations that prepare for this -- baking security into their agentic workflows from the start rather than bolting it on later -- will be better positioned to defend against adversaries using the same technology.

Security Review as a Daily Practice

The security engineering team at Anthropic demonstrated a practical pattern: copying infrastructure-as-code plans into Claude Code to ask "what is this going to do, and am I going to regret it?" Infrastructure changes that require security approval -- deployment configurations, network policy changes, access control modifications -- get reviewed by Claude in seconds rather than sitting in a queue for the security team. This creates tighter feedback loops and eliminates developer blocks while waiting for security review.

This is not replacing security review. It is making security review continuous. Instead of a bottleneck at the end of the pipeline, security analysis happens at the moment of creation. The security team still reviews the final configuration, but Claude catches the obvious issues before they get that far.

The organizational response is not to restrict Claude Code -- that forfeits the defensive advantage while doing nothing about offensive use. The response is to ensure your security architecture does not depend on attackers being unsophisticated. Assume your adversaries have access to tools at least as capable as yours. Then build your defenses accordingly.

Paper Trading: Safe Experimentation Boundaries

The concept of paper trading -- simulating real operations without real consequences -- applies directly to Claude Code adoption. Before giving an agent access to production systems, production databases, or production credentials, establish a boundary where it can operate freely without risk.

This means: sandbox environments with synthetic data. Staging systems that mirror production topology but contain nothing sensitive. Local development environments with mocked external services. The agent can run autonomously, make mistakes, discover edge cases, and crash without consequences.

The pattern from practitioners who have run agents in autonomous loops for extended periods is consistent: start with paper trading. Validate behavior in safe environments. Expand the boundary gradually. The trust architecture supports this progression -- you can start with maximum restriction and relax constraints as you build confidence.

The permission system, sandbox isolation, hooks, and managed settings together create a gradient from "Claude Code cannot do anything without asking" to "Claude Code can do everything within this isolated environment." The right position on that gradient depends on your context, your risk tolerance, and how much you have validated in safe environments first.

Key Takeaways

- The scope hierarchy is strict -- Managed settings override everything, and there is no escape from organizational policy.
- Permission rules evaluate Deny before Ask before Allow, with first-match-wins semantics, so write deny rules narrow and allow rules broad.
- Checkpoints cover direct file edits but not bash command side effects -- git is your real safety net.
- Devcontainers provide network isolation but cannot prevent credential exfiltration from malicious project code; layer defenses, and use Anthropic's reference implementation as a starting point.
- The hook system has thirteen events covering the entire agentic lifecycle: PreToolUse, PostToolUse, PostToolUseFailure, PermissionRequest, Notification, SubagentStart, SubagentStop, Stop, TeammateIdle, TaskCompleted, PreCompact, SessionEnd, and more.
- Three hook types serve different verification needs: command hooks for deterministic checks, prompt hooks for LLM-evaluated judgments, and agent hooks for multi-turn investigation with file access.
- Async hooks run in the background without blocking Claude; use them for test suites and long-running validation.
- Hooks are snapshotted at startup -- mid-session modifications trigger warnings and require review in the `/hooks` menu before taking effect.
- PreToolUse hooks can modify tool input before execution via `updatedInput`, enabling command rewriting, safety flag injection, and monitoring harnesses.
- MCP servers provide better security visibility than raw bash commands for sensitive data access; build MCP servers for your sensitive data sources.
- Security knowledge is being democratized by agentic tools; any engineer can now perform reviews that once required specialists, but adversaries gain the same capabilities.
- Start with maximum restriction and paper-trading environments, then relax constraints as you validate agent behavior.

Chapter 3: Context Engineering

What You'll Learn

Every feature in Claude Code traces back to one constraint: the context window. It is the CPU, the RAM, and the hard drive of your agentic workflow. Fill it with the wrong things and Claude forgets your instructions. Leave it empty and Claude wastes cycles rediscovering what you already know. The developers who get extraordinary results are not better prompters. They are better context engineers.

This chapter covers the system that makes context engineering practical: CLAUDE.md. You will learn what belongs in it, what does not, and why the distinction matters more than you think. You will see the exact context cost comparison for every mechanism -- CLAUDE.md, skills, MCP, subagents, and hooks -- so you can make informed budget decisions. You will understand how auto-compaction works, how to survive it with Compact Instructions and focused compaction, and how to use `/context` and `/rewind` for surgical context management. You will learn how to bootstrap CLAUDE.md with `/init`, extend it with `@path` imports, and suppress skill loading with `disable-model-invocation: true` for zero-cost manual-only skills.

You will also discover that CLAUDE.md is not just for code projects. The most compelling case study for institutional memory comes from a non-engineering domain: someone used Claude Code to build a professional-grade financial plan across multiple sessions, with CLAUDE.md accumulating data quirks, strategy decisions, and target allocations as project memory. The same pattern applies to any domain where work spans sessions and context must persist.

Most developers treat CLAUDE.md as a configuration file. The best ones treat it as a codebase's nervous system -- the place where hard-won knowledge accumulates so no one has to learn the same lesson twice.

The Always-On File

CLAUDE.md is not documentation. It is not a README. It is instructions injected into every single request Claude processes in your project. That distinction matters because it means every line in CLAUDE.md has a cost: it consumes context window space on every turn of the conversation, reducing the space available for actual work.

Claude Code loads CLAUDE.md files from multiple locations in a strict hierarchy. A CLAUDE.md in your home directory applies to every project. One in your project root applies to that repository. Subdirectory CLAUDE.md files apply when Claude works in those directories. And if your organization deploys managed CLAUDE.md files to system directories, those apply to every user on the machine.

The loading is additive. Claude sees all of them, layered together. This means your home-level file should contain preferences that apply everywhere -- your preferred testing approach, your commit message style, your language of choice. Your project-level file should contain repository-specific knowledge. Subdirectory files should contain context relevant only to that part of the codebase.

Bootstrapping with /init

If you are starting from scratch, the `/init` command does the first draft for you. It analyzes your codebase -- detecting build systems, test frameworks, code patterns, and directory structure -- and generates a starter CLAUDE.md with the basics already filled in. The output is not perfect, but it is a solid foundation you can refine over time. Think of `/init` as scaffolding: it gets the structure right so you can focus on the project-specific knowledge that only you have.

The @path Import System

One mechanism makes this system scale: the `@path` import. Instead of inlining everything, you can write `@docs/architecture.md` in your CLAUDE.md to pull in external files. This keeps the root file lean while giving Claude access to deeper reference material.

The syntax supports several forms:

```
@README.md          # relative to CLAUDE.md location  
@docs/git-instructions.md    # subdirectory reference  
@~/claude/my-project-instructions.md  # home directory reference
```

The imported content still costs context, but the indirection makes your CLAUDE.md readable and maintainable. You can organize reference material in separate files -- API contracts, coding standards, architectural decision records -- and import only what is relevant to the current project.

The 500-Line Guideline

The 500-line guideline exists for a reason. Past that threshold, instructions start competing with each other. Claude does not ignore a bloated CLAUDE.md -- it dilutes attention across too many directives, and the ones that matter most get the same weight as the ones that barely matter. A concise CLAUDE.md is not a nice-to-have. It is a performance requirement.

If Claude keeps doing something you do not want despite having a rule against it, the file is probably too long and the rule is getting lost. If Claude asks you questions that are answered in CLAUDE.md, the phrasing might be ambiguous. Treat CLAUDE.md like code: review it when things go wrong, prune it regularly, and test changes by observing whether Claude's behavior actually shifts. You can tune instructions by adding emphasis -- "IMPORTANT" or "YOU MUST" -- to improve adherence on critical rules.

What Belongs -- and What Does Not

The question is not "what should Claude know?" It is "what can Claude not figure out on its own?" That filter eliminates most of what developers instinctively put in CLAUDE.md.

Include	Exclude
Bash commands Claude cannot guess	Anything Claude can figure out by reading code
Code style rules that differ from defaults	Standard language conventions Claude already knows
Testing instructions and preferred test runners	Detailed API documentation (link to docs instead)

Repository etiquette (branch naming, PR conventions)	Information that changes frequently
Architectural decisions specific to your project	Long explanations or tutorials
Developer environment quirks (required env vars)	File-by-file descriptions of the codebase
Common gotchas or non-obvious behaviors	Self-evident practices like "write clean code"

Belongs in CLAUDE.md:

- Build and test commands. Your project's specific test runner invocation, build commands, and linting steps -- Claude cannot infer which commands your project uses without trying them and possibly failing.
- Code style deviations from standard conventions. If your team uses tabs instead of spaces, or if your Python project follows a non-PEP8 import order, say so. Claude will otherwise follow the standard.
- Repository etiquette. Which branches are protected. Whether you squash-merge or rebase. Whether PRs need a specific label format.
- Architectural decisions and their rationale. "We use a hexagonal architecture. Domain logic must never import from infrastructure." Claude can see your code structure, but it cannot see the design intent behind it.
- Environment quirks. "The test database runs on port 5433, not 5432." "CI uses runtime version 18, not 20." These are invisible traps that will burn time.
- Data quirks and known issues. "The `created_at` column in the users table has null values before March 2024 due to a migration bug." Claude will discover this eventually, but the discovery will cost you context and time.

Does not belong in CLAUDE.md:

- Things Claude infers from code. It reads your dependency manifests, your configuration files, your directory structure. Telling it which framework and language your project uses when that is obvious from the repository is a waste of context.
- Standard language conventions. Claude already knows common style guides and idiomatic patterns for major languages. Only document deviations.
- Verbose explanations of well-known libraries. Do not paste framework documentation into CLAUDE.md. Claude has extensive training data on popular tools.
- Step-by-step instructions for common workflows. Claude knows how to create a pull request, run a migration, or set up a test file. Tell it your project-specific variations, not the generic process.

The quality of your CLAUDE.md directly correlates with Claude Code's effectiveness. Teams that accumulate five or more sessions' worth of refined CLAUDE.md content for multi-day projects report dramatically better results than those starting from scratch each time.

The Context Cost Comparison

Not everything you configure costs the same. Understanding the cost model is the difference between a session that stays sharp for 200 turns and one that compacts after 40. Here is the complete picture:

Feature	When It Loads	What Loads	Context Cost
CLAUDE.md	Session start	Full content of all CLAUDE.md files	Every request
Skills	Session start + when used	Descriptions at start; full content when invoked	Low (descriptions every request)*
MCP servers	Session start	All tool definitions and schemas	Every request
Subagents	When spawned	Fresh context with specified skills	Isolated from main session

Hooks	On trigger	Nothing (runs externally)	Zero, unless hook returns additional context
--------------	------------	---------------------------	--

*By default, skill descriptions load at session start so Claude can decide when to use them. Set `disable-model-invocation: true` in a skill's frontmatter to hide it from Claude entirely until you invoke it manually with `/skill-name`. This reduces context cost to zero for skills you only trigger yourself.

The Decision Matrix: Skill vs. Subagent vs. CLAUDE.md vs. MCP

Each mechanism serves a different purpose. The decision comes down to when the information is needed and how much it costs to keep available:

Use CLAUDE.md when: The information is needed on every turn. Build commands, style rules, architectural constraints, environment quirks. This is the always-on channel. Keep it under 500 lines.

Use a skill when: The information is needed for specific task types. Database migration procedures, deployment checklists, API integration guides, complex refactoring workflows. Skills load on demand -- their descriptions cost a small amount on every request, but the full content only loads when Claude or you invoke them.

Use a subagent when: The task involves reading many files, processing verbose output, or performing exploration that would bloat your main context. Subagents run in their own context window. Only the summary returns to your session. This makes subagents a context management strategy, not just a parallelism tool.

Use MCP when: You need structured access to external services. MCP tool definitions load at session start and persist, so every MCP server you configure has an ongoing cost. But the structured interface, permission integration, and hook visibility make it the right choice for sensitive or frequently-used external integrations.

Use hooks when: You need to extend Claude's behavior without consuming any context. Hooks run as external processes. A PreToolUse hook that validates file paths, a PostToolUse hook that logs actions, a SessionStart hook that sets up the environment -- none of these consume any context. This makes hooks the most context-efficient extension mechanism available.

Getting the classification wrong in either direction hurts. Put reference material in CLAUDE.md and you waste context on every request. Put always-needed rules in a skill and Claude will violate them whenever the skill is not loaded.

Auditing Your Context Budget

This cost model has a practical implication: if you find yourself frequently running out of context, audit your always-on costs first. Run `/mcp` to see per-server token costs. Review your CLAUDE.md for content that could move to skills. Check whether you have MCP servers loaded that you rarely use. The `/context` command visualizes your current context usage as a colored grid, showing you exactly where the space is going.

Skills vs. CLAUDE.md: The Loading Distinction

Both skills and CLAUDE.md provide Claude with instructions. The difference is when they load, and that difference governs your context budget.

CLAUDE.md content loads on every request. It is always present, always consuming context. This makes it the right place for information Claude needs constantly: build commands, style rules, architectural constraints.

Skills are on-demand. Their descriptions load at session start -- a short summary of what each skill does -- but the full skill content only loads when Claude invokes the skill. A skill containing 200 lines of database migration instructions costs you almost nothing until Claude actually needs to perform a migration. At that point, the skill content loads into the context for that specific interaction.

Manual-Only Skills with disable-model-invocation

Some skills are not meant for Claude to decide when to use. A deployment skill, a release checklist, a complex refactoring workflow -- you want to invoke these yourself with `/skill-name`, not have Claude decide the moment is right.

Setting `disable-model-invocation: true` in a skill's frontmatter hides the skill from Claude entirely. It does not load the description at session start. It does not appear in Claude's available tools. The context cost is zero until you manually invoke it. This is the most aggressive context optimization for skills you use infrequently or that should only run when you explicitly choose.

```
---
```

```
description: Deploy to production
disable-model-invocation: true
```

```
---
```

```
Follow the production deployment checklist...
```

The decision rule is simple. If Claude needs to know something on every turn -- "always use single quotes," "never modify files in /legacy" -- that goes in CLAUDE.md. If Claude needs reference material for a specific type of task -- detailed deployment procedures, complex refactoring checklists, API integration guides -- that is a skill. If only you should decide when to trigger the skill, add `disable-model-invocation: true`.

Auto-Compaction: What Happens When Context Fills

At approximately 95% context capacity, Claude triggers automatic compaction. This is not a crash. It is a managed process, but understanding how it works matters because it determines what survives.

During compaction, Claude clears older tool outputs and summarizes the conversation history. The file contents it read twenty turns ago, the command outputs from earlier exploration, the intermediate reasoning -- all of it gets compressed into a summary. This frees context space but loses detail.

Compact Instructions

One section of CLAUDE.md gets special treatment: anything under the heading "Compact Instructions." Content in that section is explicitly preserved through compaction. This is your opportunity to specify what Claude must remember even after the conversation history is compressed.

Use Compact Instructions for:

- Critical constraints that must never be violated, regardless of how long the session runs.
- The current task description, if it is complex enough that a summary might lose important nuances.
- Key decisions made earlier in the session that affect all subsequent work.

Focused Compaction with /compact

You can trigger compaction manually with `/compact`, and you can direct what survives by adding focus instructions: `/compact Focus on the API changes`. This tells the compaction process which aspects of the conversation to prioritize when summarizing. Without focus instructions, compaction makes its own judgment about what matters -- which is often reasonable but occasionally drops the wrong things.

Manual compaction at 80% is often better than waiting for auto-compaction at 95%. You have more control, the compaction has more room to work with, and you can direct it while the conversation history is still relatively fresh.

You can also override the compaction threshold with the `CLAUDE_AUTOCOMPACT_PCT_OVERRIDE` environment variable. Setting it to 80 triggers compaction earlier, freeing space before the window is critically full. Setting it higher delays compaction but risks performance degradation in the final stretch before it triggers.

Selective Summarization with /rewind

Sometimes you do not want to compact the entire conversation -- just a portion of it. Press Escape twice (or use the `/rewind` command) to open the rewind menu. A scrollable list shows each of your prompts as checkpoints. Select a message and you get four options:

1. **Restore code and conversation** -- revert everything to that point.
2. **Restore conversation only** -- rewind the conversation but keep the current code.
3. **Restore code only** -- keep the conversation but revert file changes.
4. **Summarize from here** -- condense everything from the selected point forward into a summary, keeping the conversation before that point at full fidelity.

The "Summarize from here" option is surgical compaction. If you spent thirty turns exploring a dead end and then found the right approach, you can summarize the exploration (freeing context) while preserving the productive work at full detail. This is far more precise than whole-conversation compaction.

The PreCompact Hook

The PreCompact hook event fires before compaction begins, whether triggered manually or automatically. It receives the trigger type (`manual` or `auto`) and any custom instructions passed to `/compact`. While it cannot block compaction, it can perform preparation -- saving session state, logging context usage, or injecting context that should influence the compaction summary. This is the programmatic extension point for context management.

The /context Command

The `/context` command visualizes your current context usage as a colored grid. It shows how much space is consumed by the system prompt, CLAUDE.md content, MCP tool definitions, conversation history, and active work. This is the diagnostic tool for context engineering -- when performance degrades, `/context` tells you why. If MCP servers are consuming 30% of your context before you ask a single question, you know what to fix.

The practical takeaway: if your sessions regularly hit compaction, you are either doing too much in one session or your always-on context costs are too high. Both are solvable.

Session Forking as Context Recovery

When a session's context is polluted -- too many dead ends, too much exploration that is no longer relevant, auto-compaction that lost critical instructions -- you have a recovery option beyond starting fresh.

The `--fork-session` flag creates a new session that starts with the full history of an existing one but diverges from that point. The original session is preserved unchanged. The fork gets a new session ID and a separate conversation thread.

Use this when your session has valuable context you want to preserve but you need to take the work in a different direction. Instead of re-explaining the project setup, the decisions you have made, and the approach you are taking, you fork the session and start the new direction with all of that context intact.

```
claude --continue --fork-session
```

Session forking is not a substitute for good context management. It is an escape hatch for when the session's context has drifted from your needs despite your best efforts.

CLAUDE.md for Non-Code Domains

CLAUDE.md is not just for software projects. The most compelling demonstration of institutional memory came from someone who used Claude Code for portfolio optimization -- a financial analysis project with no traditional codebase.

The workflow: export data from financial accounts as spreadsheets, create a supplementary text file with information that had no structured export, draft a detailed goal prompt, and initialize a Claude Code session pointing at the directory containing all of it. Over multiple sessions, Claude parsed the data, classified holdings, identified gaps, and produced a phased optimization plan with tax impact analysis, fund recommendations, and before-and-after allocation grids.

The critical part: they built up a CLAUDE.md file that accumulated not code knowledge but domain analysis memory. Data quirks discovered during parsing. Strategy decisions made in earlier sessions. Target allocations agreed upon after iterative refinement.

When they returned days later, Claude picked up exactly where they left off because CLAUDE.md carried the project's analytical state.

This pattern applies to any domain where work spans sessions: research projects, data analysis, technical writing, infrastructure planning, compliance review. The CLAUDE.md content changes -- instead of build commands and coding conventions, it contains data schemas, analytical frameworks, and domain-specific constraints. But the mechanics are identical: persistent context that makes every session start smarter than the last.

The Work Compounds

The same practitioner pointed to a pattern that deserves its own name: work product compounding. After completing the financial plan, they pointed Claude Code at the same project directory -- the session history, the plan documents, the analytical artifacts -- and asked it to write a blog post about the process. The deliverable from the first project became the input for an entirely different deliverable.

This is not an accident. It is a consequence of file-system-based context. Every artifact Claude produces -- plans, analyses, reports, code -- lives as a file in your project directory. That file is available to future sessions, to future skills, to future subagents. The work you do with Claude Code is not trapped in a conversation transcript. It accumulates in the filesystem, and each new session can build on everything that came before.

Persona-Adaptive CLAUDE.md

The product design team at Anthropic discovered something that software engineers might not think of: CLAUDE.md works differently when the user is not a developer.

Designers on the team created custom CLAUDE.md files that told Claude they were designers with limited coding experience who needed detailed explanations and smaller, incremental changes. This dramatically improved the quality of Claude's responses -- instead of terse technical output that assumed engineering fluency, Claude produced step-by-step explanations, commented its code more heavily, and made smaller changes that were easier to review.

This is persona-adaptive context. The CLAUDE.md does not just describe the project; it describes the user. A data scientist might specify that they prefer exploratory scripts with visualizations. A technical writer might specify that they need verbose commit messages and documentation alongside every change. A designer might specify that they need UI-focused explanations and visual diffs.

The implication: CLAUDE.md is not one-size-fits-all even within a team. Project-level CLAUDE.md files carry shared standards. User-level CLAUDE.md files in the home directory carry personal preferences. The combination means Claude adapts to both the project's needs and the individual's working style.

Continuous Improvement Loop

One pattern from a data infrastructure team at Anthropic extends the end-of-session update beyond accumulating knowledge. They ask Claude to summarize completed work and suggest improvements -- not just additions to CLAUDE.md, but improvements to the workflow itself.

The distinction is subtle but important. Most teams use end-of-session updates to add knowledge: "the test database uses port 5433" or "module X has a circular dependency." The continuous improvement pattern also captures process improvements: "the deployment script should run lint before build" or "the CLAUDE.md instructions for the API module are causing Claude to generate overly verbose error handlers -- simplify them."

This creates a feedback loop where CLAUDE.md evolves not just in what it knows but in how it instructs. Over time, the instructions themselves get refined based on observed results. The team reported that this loop -- updating both knowledge and process after each session -- made subsequent iterations measurably more effective because Claude was not just working from better knowledge but from better instructions.

Institutional Memory

CLAUDE.md becomes dramatically more valuable when you treat it as a living document that accumulates knowledge across sessions.

Here is the pattern. You start a project with a minimal CLAUDE.md -- build commands, style rules, a few architectural notes. During your first session, Claude discovers that your test suite requires a specific environment variable. It learns that one module has a circular dependency it needs to work around. It figures out that the API client throws a non-standard error format.

All of that knowledge exists in the session context. When the session ends, it vanishes.

The fix is an end-of-session update loop. Before closing a session, ask Claude to review what it learned and suggest additions to CLAUDE.md. Claude will propose specific, concrete additions based on the problems it encountered. You review them, accept what is useful, and commit the updated CLAUDE.md.

Over time, CLAUDE.md accumulates the kind of knowledge that usually lives in a senior engineer's head: the quirks, the workarounds, the "here is why we do it this way" explanations that no one ever writes down. Except now they are written down, and every future session -- by you or any team member -- benefits from them.

This creates a compounding effect. Each session starts with more knowledge than the last. Mistakes that burned time in session one never recur because they are documented in CLAUDE.md. The project-specific knowledge base grows, and with it, Claude's effectiveness.

Teams that check CLAUDE.md into version control multiply this effect across the entire team. One developer's discovery becomes everyone's context.

Preventing Repeated Errors

One specific application of institutional memory deserves its own section because it solves a specific, common frustration: Claude making the same mistake repeatedly.

Internal teams at Anthropic have documented a pattern where Claude Code makes particular tool-calling errors or follows incorrect patterns consistently within a project. The fix is surgical: add a targeted instruction to CLAUDE.md that addresses the exact error.

For example, if Claude keeps importing from the wrong module path, add: "When importing database utilities, use `@app/db/utils` not `@app/utils/db`. The latter path exists but is deprecated."

If Claude keeps generating tests that fail because of a timing issue: "All async tests in this project must use the `waitForCondition` helper, not `setTimeout`. The CI environment has non-deterministic timing."

These targeted additions are more effective than general instructions because they address specific failure modes Claude has demonstrated. They are not aspirational guidelines -- they are patches for observed bugs in Claude's behavior within your project.

Spec-Driven Context

CLAUDE.md is persistent but static between manual updates. For complex, multi-session projects, there is a complementary approach: spec documents.

A spec document is a detailed description of what you are building, written before implementation begins. It lives as a file in your repository -- `spec.md`, `design.md`, whatever naming fits your project. You reference it from CLAUDE.md with an `@path` import or tell Claude to read it at the start of each session.

The power of specs is that they survive context pollution and session restarts. When a session compacts and loses the nuances of your earlier discussion, the spec is still there, unchanged, in a file Claude can re-read. When you start a new session, the spec provides full context without you re-explaining anything.

Spec-first development also produces a concrete artifact you can review before any code is written. You can ask Claude to write the spec, review it, refine it through conversation, and only then move to implementation. The spec becomes the source of truth that keeps implementation aligned across multiple sessions, multiple subagents, and multiple days of work.

This pattern matters most for projects that span more than two or three sessions. Without a spec, each new session starts with a lossy re-explanation of intent. With a spec, each session starts with a precise, version-controlled description of the goal.

AGENTS.md: The Cross-Agent Standard

CLAUDE.md is specific to Claude Code. But if your team uses multiple AI coding tools -- or if contributors to your repository use different agents -- there is a complementary standard worth knowing about.

AGENTS.md is an open standard for agent-specific documentation. It lives in your repository root and works with over twenty different coding agents. While a README targets humans, AGENTS.md contains the extra context that coding agents need: development environment setup, coding standards, testing procedures, and architectural constraints.

The practical overlap with CLAUDE.md is significant, and the key insight is about information architecture. A bloated AGENTS.md (or CLAUDE.md) that tries to contain everything -- sometimes running to hundreds of lines -- consumes enormous context. The lean approach uses progressive disclosure: a concise root file with a docs-reference table that points to detailed documentation Claude can read on demand. Instead of inlining your API documentation, you reference it. Instead of pasting your architecture guide, you import it.

```
# AGENTS.md

## Dev Environment
- How to set up and navigate

## Standards
- Code style, naming, patterns

## Testing
- How to run and write tests

## Docs Reference
| Topic | File |
|-----|-----|
| API contracts | docs/api.md |
| Architecture | docs/architecture.md |
| Deployment | docs/deploy.md |
```

The lean version loads fast, costs minimal context, and lets the agent pull detailed documentation only when it needs it. This is the same principle as the CLAUDE.md 500-line guideline, applied to the broader agent ecosystem.

Context Isolation with Subagents

Subagents are not just about parallelism. They are a context isolation strategy.

When Claude invokes a subagent, the subagent runs in its own context window. It reads files, runs commands, processes data -- all in isolation. Only the summary returned to the parent costs main-window context. A subagent that reads fifty files and runs twenty commands consumes enormous context in its own window, but your main conversation never sees any of it.

The practical impact is dramatic. In one documented case, a developer orchestrated fourteen subagents over the course of a complex migration project. Each subagent completed its task, committed its changes, and returned a brief summary. After all fourteen were done, the main session's context usage was at 143,000 tokens out of a 200,000 token window -- seventy-one percent. The system prompt and tool definitions accounted for roughly ten percent. The orchestration overhead -- creating tasks, tracking progress, synthesizing results -- consumed the rest. None of the actual implementation work touched the main context.

This means a skill that reads twenty files and runs complex analysis does not consume twenty files' worth of main context. Only the final result comes back. Skills running in their own context fork (`context: fork`) work the same way.

The implication for design: if you have a complex workflow that requires reading many files or processing extensive data, packaging it as a skill or delegating it to a subagent is more context-efficient than running the same steps directly in the main conversation. The main conversation stays lean. The heavy lifting happens in isolation.

You should also think carefully about what a skill or subagent returns. A skill that returns its entire analysis as a giant block of text defeats the purpose of isolation. Design outputs to be concise summaries that give the main conversation what it needs without the baggage of how the result was produced.

Dynamic Context with SessionStart Hooks

CLAUDE.md is static -- it changes only when someone edits it. But some context is inherently dynamic: the current state of the CI pipeline, the list of open issues, recent commits by other team members, the current branch's relationship to main.

SessionStart hooks solve this. A hook configured on the `SessionStart` event runs an external command when Claude Code launches and injects the output as context. Unlike CLAUDE.md content (which costs context every turn), hook output is injected once and treated as part of the session initialization.

Practical examples:

- A script that checks for open pull requests on the current branch and summarizes review comments.
- A command that lists the most recent commits on main since the current branch diverged, so Claude knows what has changed.
- A script that reads environment-specific configuration and injects it as context, adapting Claude's behavior for development versus staging versus production environments.

SessionStart hooks complement static CLAUDE.md. The CLAUDE.md provides unchanging project knowledge. The hooks provide point-in-time state. Together, they give Claude both the permanent rules and the current situation.

Putting It All Together

Context engineering is not a single technique. It is a system of interlocking decisions:

1. **CLAUDE.md** carries the permanent, always-needed knowledge. Keep it under 500 lines. Focus on what Claude cannot infer. Bootstrap with `/init`, extend with `@path` imports, and refine through end-of-session updates.
2. **Skills** hold reference material that loads on demand. Use them for task-specific instructions that would waste context if loaded constantly. Set `disable-model-invocation: true` for skills you only invoke manually.
3. **Hooks** inject dynamic context at zero ongoing cost. Use SessionStart hooks for environment-specific state. Use PreCompact hooks to prepare for compaction.
4. **Specs** persist complex project intent across sessions and compactions. Use them for anything that spans more than two sessions.
5. **Subagents** isolate expensive operations from your main context. Delegate verbose exploration and analysis to keep your primary conversation sharp. Fourteen subagents can run without exhausting the main session's context.
6. **Compact Instructions** in CLAUDE.md protect critical rules through auto-compaction. Use `/compact` with `focus` instructions for directed compaction, and `/rewind` with "Summarize from here" for surgical context cleanup.
7. **End-of-session updates** turn CLAUDE.md into a knowledge accumulator that improves with every session. Extend this to a continuous improvement loop by having Claude suggest workflow improvements, not just knowledge additions.
8. **Persona-adaptive CLAUDE.md** in home directories lets non-developers get appropriately detailed explanations and incremental workflows tuned to their expertise level.

The developers who master this system do not work harder. They start every session with better context than the session before. Their Claude Code instances know more about their projects than most human team members do, because the knowledge is written down, versioned, and loaded automatically.

That is context engineering. Not prompt tricks. Infrastructure.

Key Takeaways

- CLAUDE.md content costs context on every request; keep it under 500 lines, bootstrap with `/init`, and focus on information Claude cannot infer from your code.
- The context cost comparison table is your budget guide: CLAUDE.md and MCP cost every request, skills cost low (descriptions only), subagents are isolated, and hooks are free.
- Put always-needed rules in CLAUDE.md and task-specific reference material in skills; use `disable-model-invocation: true` for manual-only skills with zero context cost.
- Use `@path/to/import` syntax to keep CLAUDE.md lean while giving Claude access to deeper reference material in separate files.
- Compact Instructions in CLAUDE.md survive auto-compaction; `/compact` with focus instructions directs what gets preserved; `/rewind` with "Summarize from here" enables surgical context cleanup.
- The `/context` command shows where your context budget is going -- audit it before blaming Claude for running out of space.
- CLAUDE.md works for non-code domains: financial analysis, data science, technical writing -- any domain where work spans sessions and context must persist.
- Subagent context isolation is dramatic: fourteen subagents can run without exhausting the main session because their work stays in separate context windows.
- End-of-session CLAUDE.md updates create compounding institutional memory; extend this to a continuous improvement loop by capturing workflow improvements alongside knowledge.
- Targeted CLAUDE.md additions that address specific observed errors are more effective than general instructions.
- Spec documents survive context compaction and session restarts, making them essential for multi-session projects.
- Session forking (`--fork-session`) preserves valuable context while letting you take work in a new direction.

Chapter 4: Multi-Agent Orchestration

What You'll Learn

The moment your task outgrows a single context window, you have two choices: fight the constraint or distribute the work. Multi-agent orchestration is how you distribute the work.

Claude Code ships with a subagent system that lets you spin up isolated AI workers, each with its own context window, its own tool permissions, and its own system prompt. For more ambitious coordination, an experimental agent teams system lets multiple independent sessions communicate through seven coordination primitives, shared task lists, and direct messaging. These are fundamentally different architectures with different cost profiles, different failure modes, and different sweet spots.

This chapter covers both in depth. You will learn the six built-in subagent types, how to define custom subagents with persistent memory, hooks, turn limits, and preloaded skills. You will understand the seven agent team primitives, four message types, three display modes, task dependencies, file-locking for concurrent claims, and plan approval workflows. You will see why a practitioner's four-agent trading hierarchy degraded to a single agent because simplicity won, and why a QA swarm of five parallel agents audited an entire blog in three minutes. You will also learn the task system that coordinates work across sessions and subagents, and why the cheapest orchestration strategy is almost always the correct one.

The Architecture of Subagents

A subagent is a Claude Code session within a session. It gets its own context window, its own system prompt, and a restricted set of tools. When it finishes, it returns a summary to the parent session. Then its context window is discarded.

The key constraint: subagents cannot spawn other subagents. This creates a strict two-level hierarchy. You have an orchestrator (your main session) and workers (subagents). There is no middle management. No delegation chains. No recursive agent trees. This is a deliberate design decision, not a limitation. Recursive agent hierarchies sound elegant and produce chaos.

Six built-in subagent types ship with Claude Code:

Explore subagents run on a smaller, cheaper model optimized for speed. They are read-only -- they can search files, read code, and navigate the codebase, but they cannot modify anything. Use them when you need to understand something before deciding what to do. They are fast and cheap enough to spin up casually.

Plan subagents inherit the main session's model. They are also read-only, but they bring the full reasoning capability of the orchestrator model to bear on research and planning tasks. Use them when the planning itself is complex enough to warrant dedicated context.

General-purpose subagents have access to all tools. They can read, write, execute commands, and perform multi-step tasks. These are your workhorses for implementation.

Bash subagents are specialized for command execution. They are useful when you need to run a sequence of shell commands without polluting the main context with verbose output.

Claude Code Guide subagents are specialized for answering questions about Claude Code itself -- its features, configuration, and best practices. They consult the product's own documentation.

Statusline-setup subagents handle the narrow task of configuring terminal status line integrations. They exist because the setup process is fiddly enough to warrant isolated context.

You can also define custom subagents as Markdown files with YAML frontmatter, specifying system prompts, tool restrictions, model preferences, and turn limits. These definitions live in `.claude/agents/` for project scope or `~/.claude/agents/` for user scope. Plugins can supply them too, at the lowest priority.

Custom Subagent Definition Fields

Beyond the basics of system prompt and tool restrictions, custom subagent definitions support several fields that shape behavior in important ways:

maxTurns limits the number of agentic turns before a subagent stops. This is a safety valve. A subagent exploring a codebase might spiral into increasingly tangential searches. Setting `maxTurns: 20` caps the exploration and forces a result. Without it, a subagent on a dead-end path burns tokens until context runs out.

mcpServers configures which MCP servers are available to a specific subagent. You can reference servers defined in `.mcp.json` by name or inline a full server configuration. This means a database-reader subagent can have access to your database MCP server while a code-reviewer subagent has none. Tool access scoped to the subagent's role, not the session's full capability set.

skills preloads skill content into the subagent's context at launch. In the main session, skills load on demand -- only the description is present until Claude invokes the skill. Subagents work differently. Skills passed to a subagent are fully injected into its context at startup because the subagent needs to act immediately without the interactive skill-loading flow. This means every skill you preload costs context from turn one.

memory enables persistent memory, covered in detail later in this chapter.

The `--agents` CLI Flag

For quick testing or CI automation, you can pass subagent definitions as JSON directly when launching Claude Code instead of writing Markdown files:

```
claude --agents '{
  "code-reviewer": {
    "description": "Reviews code for style and correctness",
    "prompt": "You are a code reviewer. Check for bugs, style violations, and security issues.",
    "tools": ["Read", "Glob", "Grep"],
    "model": "sonnet"
  },
  "debugger": {
    "description": "Provides low-level debugging information and tools for troubleshooting."}
```

```

    "description": "Debugs failing tests",
    "prompt": "You are a debugging specialist. Analyze test failures and suggest fixes.",
    "tools": ["Read", "Glob", "Grep", "Bash"],
    "model": "sonnet",
    "maxTurns": 30
  }
}'

```

The `--agents` flag accepts the same fields as file-based definitions: `description`, `prompt`, `tools`, `disallowedTools`, `model`, `permissionMode`, `mcpServers`, `hooks`, `maxTurns`, `skills`, and `memory`. This is particularly useful in headless pipelines where you want to define specialized agents without committing Markdown files to the repository.

The `--agent` Flag

A related but distinct flag: `--agent` runs a custom agent definition as the main thread rather than as a subagent. If you have a `code-reviewer` agent defined in `.claude/agents/code-reviewer.md`, running `claude --agent code-reviewer` starts the session using that agent's system prompt, tools, and configuration. The session runs as a top-level agent, not as a worker within another session. You can restrict which subagents this top-level agent can spawn using `Task(agent_type)` syntax, creating controlled hierarchies where a review agent can only delegate to explorer-type subagents.

Context Isolation Is the Point

The most important thing subagents give you is not parallelism. It is context isolation.

When you ask the main session to read a large file, parse verbose test output, or explore a sprawling directory tree, all of that content accumulates in the main context window. Do this enough times and you hit the ceiling described in Chapter 1 -- instructions get forgotten, work gets dropped, quality degrades.

Subagents solve this by containing the mess. A subagent can read fifty files, parse a thousand lines of test output, explore every directory in the repo, and the main session only sees the summary. The verbose intermediate work stays in the subagent's context and is discarded when the subagent finishes.

This is why experienced users route expensive context operations through subagents even when they do not need parallelism. Reading a large codebase module? Subagent. Running a comprehensive test suite and analyzing failures? Subagent. Searching for usage patterns across hundreds of files? Subagent. The goal is keeping the orchestrator's context lean so it can coordinate effectively across dozens of tasks without losing coherence.

The tradeoff is that summaries are lossy. A subagent that returns a detailed result consumes main context proportional to the length of that result. If you spin up ten subagents and each returns a page of findings, you have ten pages of results in your main window. The solution is to instruct subagents to be concise in their returns -- specific answers, not comprehensive reports.

Foreground vs. Background Execution

Subagents can run in the foreground or the background, and the distinction matters more than it appears.

Foreground subagents block the main conversation. The orchestrator waits for them to finish before proceeding. They have full access to MCP tools, can ask clarifying questions if they get stuck, and integrate seamlessly with the permission system. Use foreground execution when the subagent's result determines the next step.

Background subagents run concurrently while the main session continues. This is where parallelism actually happens -- you can have multiple background subagents exploring different parts of the codebase or implementing different components simultaneously. But background execution comes with restrictions: no MCP tools, no clarifying questions, and permissions must be pre-negotiated before launch. If a background subagent encounters something that needs human approval and was not pre-approved, it auto-denies and moves on.

Background subagents can be resumed in the foreground if they complete with unresolved questions or partial results. This is useful for tasks that start as background exploration but surface something that needs interactive discussion. Press **Ctrl+B** to background a running foreground task -- it continues working while you regain the main session. If a background subagent fails due to missing permissions, you can resume it in the foreground to retry with interactive prompts.

The practical pattern is to use foreground for tasks where you need the result immediately and background for tasks that can run while you focus on something else. A common orchestration workflow: launch three background subagents to explore three modules, continue working on something else in the main session, then review their results when they complete.

Resuming Subagents

Subagent transcripts persist independently of the main conversation. They are stored as separate files at `~/.claude/projects/{project}/{sessionId}/subagents/`. When the main conversation compacts, subagent transcripts are unaffected. You can resume a subagent after restarting Claude Code by resuming the same session -- ask Claude to "continue that code review" and it picks up from the persisted transcript.

This persistence means subagent work is not lost when the main session resets. A subagent that spent twenty minutes analyzing an authorization module has its full transcript on disk. Resuming gives it that context back without re-doing the analysis.

Subagent Auto-Compaction

Subagents support automatic compaction using the same logic as the main conversation. By default, auto-compaction triggers at approximately 95% capacity. For subagents that handle large volumes of data -- parsing thousands of lines of test output, reading dozens of files -- this means they can work through content that exceeds their context window by periodically summarizing and continuing.

Set `CLAUDE_AUTOCOMPACT_PCT_OVERRIDE` to a lower percentage (for example, `50`) to trigger compaction earlier. This applies to both main conversations and subagents. Earlier compaction reduces peak context consumption at the cost of potentially losing detail from earlier turns. For subagents doing broad exploration, earlier compaction is usually fine. For subagents doing precise analysis where every detail matters, let them run closer to capacity before compacting.

Subagent Hooks

Custom subagents support lifecycle hooks that run at specific points during execution. These are defined in the agent's YAML frontmatter and work with the same event model as session-level hooks (Chapter 2), but scoped to the subagent.

Three hook events fire within the subagent's execution: **PreToolUse** runs before the subagent uses a tool and can block or modify the invocation. **PostToolUse** runs after a tool completes and can transform output or trigger side effects. **Stop** runs when the subagent finishes -- but at runtime, it is automatically converted to a **SubagentStop** event.

Two additional hook events fire at the project level, outside the subagent's own context: **SubagentStart** fires when any subagent spawns, and **SubagentStop** fires when any subagent finishes. These are defined in your project or user settings, not in the subagent definition itself. A SubagentStart hook can inject `additionalContext` into the subagent -- instructions or data that supplement its system prompt without modifying the agent definition. A SubagentStop hook can block the subagent from stopping (exit code 2), useful for enforcing that a subagent must produce certain artifacts before it is allowed to finish.

A practical example: a database-reader subagent with a PreToolUse hook on the Bash tool that checks whether commands contain SQL write operations (INSERT, UPDATE, DELETE, DROP). If the hook detects a write statement, it returns a deny decision, preventing the subagent from modifying production data regardless of what its instructions say. This is defense in depth -- the subagent's system prompt says "read only," and the hook enforces it mechanically.

Subagent Patterns

Three patterns recur in production subagent usage. They are worth naming because they solve specific orchestration problems.

Chain subagents run in sequence, where each subagent completes its task and returns results to the orchestrator, which then passes relevant information to the next subagent. Subagent A analyzes the authentication module and returns a list of vulnerabilities. The orchestrator passes that list to Subagent B, which generates fixes. Subagent B's output goes to Subagent C, which writes tests. Each subagent has clean context with only the information it needs. The chain structure prevents context accumulation while maintaining information flow.

Isolate high-volume operations by routing anything that produces large output to a subagent. Running a comprehensive test suite might generate thousands of lines of output. Fetching documentation from an API might return hundreds of pages. Processing log files might involve scanning megabytes of text. If this happens in the main session, it crowds out everything else. A subagent absorbs the volume and returns a concise summary: "14 tests failed, all in the auth module, all related to session expiry."

Parallel research spawns multiple subagents simultaneously to investigate different aspects of a problem. Three subagents explore the authentication module, the database layer, and the API endpoints in parallel. Each returns findings. The orchestrator synthesizes them into a unified understanding. This is faster than sequential exploration and keeps the orchestrator's context free of the intermediate exploration steps.

Agent Teams: The Experimental Layer

Agent teams are a different architecture entirely. Where subagents are workers within a session, agent teams are independent Claude Code sessions that coordinate through shared infrastructure.

Enabling Agent Teams

Agent teams are disabled by default. Enable them by setting the `CLAUDE_CODE_EXPERIMENTAL_AGENT_TEAMS` environment variable to `1`, either in your shell or through settings:

```
{  
  "env": {  
    "CLAUDE_CODE_EXPERIMENTAL_AGENT_TEAMS": "1"  
  }  
}
```

Then describe what you want in natural language. Claude creates the team, spawns teammates, and coordinates work based on your prompt.

The Seven Team Primitives

The system provides seven tools that handle the full lifecycle of team coordination.

TeamCreate initializes a team. It creates a team directory and configuration file on disk. The `team_name` parameter is the namespace that links everything -- tasks, messages, and config all live under it.

```
TeamCreate({ "team_name": "auth-refactor", "description": "Refactor authentication module" })
```

TaskCreate defines a unit of work. Each task becomes a JSON file on disk. The lead creates tasks before spawning teammates, providing enough detail in the description to serve as a prompt for the agent that picks it up.

```
TaskCreate({  
  "subject": "Extract JWT handling into jwt.ts",  
  "description": "Move all token signing and verification logic from auth.ts to a new jwt.ts module...",  
  "team_name": "auth-refactor"  
})
```

TaskUpdate moves work through the pipeline. Teammates use it to claim tasks (setting status to `in_progress` with an owner) and to mark tasks complete. The status field prevents two agents from working on the same task simultaneously.

```
TaskUpdate({ "taskId": "1", "status": "in_progress", "owner": "jwt-agent" })
TaskUpdate({ "taskId": "1", "status": "completed" })
```

TaskList returns all tasks with their current status. Teammates call this after completing a task to find what is next. There is no centralized scheduler -- each teammate polls TaskList, finds unowned pending tasks, and claims one. This is the shared coordination mechanism.

Task (with team_name) spawns a teammate. The critical detail: each teammate is a full Claude Code session with its own context window. Teammates load the same project context (CLAUDE.md, MCP servers, skills) but do not inherit the lead's conversation history. You can specify the model for each teammate independently.

```
Task({
  "subagent_type": "general-purpose",
  "name": "jwt-agent",
  "team_name": "auth-refactor",
  "model": "sonnet"
})
```

SendMessage is what makes teams different from subagents. Any teammate can message any other teammate directly. It supports four message types:

- **message** : direct communication between two agents. A teammate reports findings to the lead, or one teammate asks another for information.
- **broadcast** : reaches all teammates at once. Use sparingly -- costs scale with team size because every teammate processes the message.
- **shutdown_request / shutdown_response** : graceful teardown protocol. The lead sends a shutdown request; the teammate acknowledges with a response. A teammate can reject the shutdown with an explanation if it has unfinished work.
- **plan_approval_response** : quality gate. The lead approves or rejects a teammate's implementation plan.

```
// Teammate reports findings to lead
SendMessage({
  "type": "message",
  "recipient": "lead",
  "content": "JWT extraction complete. Created jwt.ts with signToken and verifyToken exports."
})

// Lead requests graceful shutdown
SendMessage({
  "type": "shutdown_request",
  "recipient": "jwt-agent",
  "content": "All tasks complete, shutting down team."
})
```

TeamDelete removes the team config and all task files from disk. Called after all teammates have shut down.

The Team Lead Abstraction

What makes agent teams more than parallel subagents is the team lead's coordination role. The lead is the session that creates the team. It is responsible for creating tasks and defining work breakdown, spawning teammates with appropriate roles, monitoring progress through the task list, synthesizing results from multiple agents, and handling graceful shutdown.

The task files on disk and SendMessage are the only coordination channels -- there is no shared memory. Teammates each have their own conversation history and context window, independent from the lead and from each other. This independence is the source of

both the architecture's strength (true parallel work with full context isolation) and its cost (each teammate is a full Claude Code session burning tokens independently).

A lead can be put into **delegate mode** (Shift+Tab after starting a team). In delegate mode, the lead cannot implement tasks itself -- it can only coordinate: spawning, messaging, shutting down teammates, and managing tasks. This prevents the lead from doing work that should be distributed to teammates, which is a surprisingly common failure mode.

Display Modes

Agent teams support three display modes, configured via the `teammateMode` setting or the `--teammate-mode` flag:

auto (default) uses split panes if you are already running inside a tmux session, and in-process otherwise.

in-process runs all teammates inside your main terminal. Use Shift+Up/Down to select a teammate and type to message them directly. Press Enter to view a teammate's session, Escape to interrupt their current turn, Ctrl+T to toggle the task list. Works in any terminal with no extra setup.

tmux gives each teammate its own terminal pane. You can see everyone's output simultaneously and click into a pane to interact directly. Split-pane mode requires tmux or iTerm2 -- it does not work in a code editor's integrated terminal.

```
# Force in-process for a single session
claude --teammate-mode in-process
```

Task Dependencies and File Locking

Tasks can depend on other tasks. A pending task with unresolved dependencies cannot be claimed until those dependencies complete. When a teammate finishes a task that others depend on, blocked tasks unblock automatically. This creates a natural wave pattern: independent tasks run in parallel first, dependent tasks follow.

Task claiming uses file locking to prevent race conditions when multiple teammates try to claim the same task simultaneously. Without file locking, two agents polling TaskList at the same moment could both attempt to claim the same task, leading to duplicated work. The lock ensures exactly one agent owns each task.

Plan Approval

For complex or risky tasks, you can require teammates to plan before implementing. The teammate works in read-only plan mode until the lead approves their approach:

```
Spawn an architect teammate to refactor the authentication module.
Require plan approval before they make any changes.
```

When a teammate finishes planning, it sends a plan approval request to the lead. The lead reviews the plan and either approves it (teammate exits plan mode and begins implementation) or rejects it with feedback (teammate stays in plan mode, revises, and resubmits). You can influence the lead's judgment by giving it criteria in your prompt: "only approve plans that include test coverage" or "reject plans that modify the database schema."

This is the difference between "trust the agents to do the right thing" and "verify the approach before spending tokens on implementation." For a team of five agents running in parallel, each burning thousands of tokens per minute, the plan approval step is cheap insurance.

Team Hook Events

Two hook events are specific to agent teams:

TeammateIdle fires when an agent team teammate is about to go idle -- it has finished its current work and is ready to stop. Exit code 2 from this hook forces the teammate to continue working. A practical use: a hook that checks whether build artifacts exist in

the expected output directory before allowing the teammate to stop. If the artifacts are missing, the hook rejects the idle state and the teammate keeps working.

TaskCompleted fires when a task is marked as completed via TaskUpdate, or when a teammate finishes its turn with in-progress tasks still assigned. Exit code 2 blocks the completion. This enables quality gates: a hook that runs the test suite when a task is marked complete, blocking the completion if tests fail. The teammate sees the failure output and can self-correct.

Agent Teams vs. Subagents: Detailed Comparison

	Subagents	Agent Teams
Context	Own window; results return to the caller	Own window; fully independent
Communication	Report results back to the main agent only	Teammates message each other directly
Coordination	Main agent manages all work	Shared task list with self-claim
Token cost	Lower: results summarized back to main	Higher: each teammate is a full session
Best for	Focused tasks where only the result matters	Complex work requiring discussion and collaboration

The distinction matters most when agents need to talk to each other. An API teammate that finishes type definitions and directly messages the UI teammate to start integration work -- that is an agent team pattern. If each worker just reports back to the orchestrator, subagents are cheaper and simpler.

Limitations

Agent teams are experimental and come with real constraints:

- **No session resumption** for in-process teammates. If a teammate crashes, its work and context are lost.
- **Task status can lag** -- teammates sometimes forget to mark tasks complete, which can block dependent tasks.
- **One team per session.** You cannot run nested teams or multiple teams from one lead session.
- **Split panes require tmux or iTerm2**, not a code editor's integrated terminal.
- **All teammates start with the lead's permission settings.** You cannot give one teammate broader permissions than the lead has.

Subagents vs. Teams: When to Use What

The decision between subagents and teams maps to scope, cost, and coordination complexity.

Subagents are right when: the work decomposes into independent tasks that can be completed without coordination, each task fits in a single subagent session, and you need results quickly and cheaply. A typical subagent workflow takes minutes and costs a fraction of a main session.

Teams are right when: the work requires sustained coordination between specialists, tasks have dependencies that need to be tracked, or the problem requires different agents maintaining long-running context about different aspects of the system. A typical team workflow takes longer and costs several times what subagents would.

Neither is right when a single well-prompted session can handle the task. This is more common than you think. Before reaching for multi-agent orchestration, ask whether the task genuinely exceeds what one context window can handle. If the answer is no, a single agent with clear instructions will outperform a multi-agent setup every time. The coordination overhead of multi-agent systems is not free, and it introduces failure modes that single agents do not have.

The hierarchy of preference: single agent first, subagents when context isolation is needed, teams when sustained coordination is needed. Work your way up, not down.

Agent Team Case Studies

The QA Swarm

A practitioner pointed agent teams at their blog with a single prompt: use a team of agents to QA the blog before production deploy. The lead created five tasks and spawned five parallel agents, each running on a cheaper model:

#	Task	Agent	What It Checked
1	Core page responses	qa-pages	16 URLs for correct HTTP status codes
2	Blog post rendering	qa-posts	83 posts for headings, meta tags, working images
3	Navigation and link integrity	qa-links	146 internal URLs for broken links
4	RSS, sitemap, SEO metadata	qa-seo	RSS validity, robots.txt, Open Graph tags
5	Accessibility and HTML structure	qa-a11y	Heading hierarchy, ARIA attributes, theme toggle

Each agent finished independently and sent a structured report back via `SendMessage`. The agents used bash and curl to fetch pages and parse HTML directly -- no browser automation, no test framework. When all five completed, the lead synthesized their findings into a prioritized report with issues ranked by severity (major, medium, minor). Then the lead sent shutdown requests, teammates acknowledged, and `TeamDelete` cleaned up.

The entire lifecycle -- from prompt to final report -- took about three minutes. Five agents, over 146 URLs tested, 83 blog posts checked. The cost was higher than a single session doing the same work sequentially, but the wall-clock time was a fraction.

The Authentication Module Refactor

A more coordination-intensive example: refactoring an authentication module with one teammate on the API layer, one on frontend components, and one running tests continuously. The key difference from subagents: the API teammate finished type definitions and messaged the UI teammate directly to say "the interfaces are ready, here's what changed." The test teammate could ask the API teammate to spin up a dev server. Self-coordination, without routing every interaction through the lead.

This is the pattern where agent teams justify their cost. The agents needed to react to each other's work in real time. With subagents, the orchestrator would have been the bottleneck -- each worker reports to the parent, the parent relays to the next worker, round and round. With a team, the workers talk directly.

Plan First, Parallelize Second

The most effective agent team pattern is a two-step approach: plan first with plan mode, then hand the plan to a team for parallel execution.

Step one: start in plan mode. Let the orchestrator explore the codebase, identify files, and produce a step-by-step implementation plan. Review it. Adjust it. This is cheap -- plan mode only reads files. The orchestrator might produce something like:

```
Plan:  
1. Create src/auth/jwt.ts -- extract token signing/verification  
2. Create src/auth/sessions.ts -- extract session logic  
3. Create src/auth/middleware.ts -- extract Express middleware  
4. Update src/auth/index.ts -- re-export public API  
5. Update 12 import sites across the codebase  
6. Update tests in src/auth/_tests_/_
```

Step two: execute the plan as a team. The plan already has the task breakdown. The lead sees the dependency graph and spawns teammates in waves:

- **Wave 1** (parallel): `jwt.ts + sessions.ts + middleware.ts` -- three teammates
- **Wave 2** (after wave 1): `index.ts barrel + update imports` -- one or two teammates

- **Wave 3** (after wave 2): update tests -- one teammate

Plan mode costs roughly 10,000 tokens. A team that goes in the wrong direction costs 500,000 or more. Spending a few seconds reviewing a plan before committing to parallel execution saves you from expensive course corrections mid-swarm.

Persistent Memory

Subagents can maintain a persistent memory directory that survives across sessions. Configure this with the `memory` field in the agent definition, specifying a scope:

- `user` : shared across all projects. Stored at `~/.claude/agent-memory/{agent-name}/`. Use for agents whose knowledge applies everywhere -- coding style preferences, common debugging patterns, personal workflow rules.
- `project` : shared within a project. Stored at `~/.claude/projects/{project}/agent-memory/{agent-name}/`. Use for project-specific agents -- a code reviewer that learns this codebase's conventions, a test agent that knows this project's testing patterns.
- `local` : specific to one working directory. Use when different branches or worktrees of the same project need different agent knowledge.

The memory directory contains a `MEMORY.md` entrypoint and optional topic files:

```
~/.claude/projects/{project}/agent-memory/code-reviewer/
├── MEMORY.md      # Concise index, loaded into every session
├── style-patterns.md # Detailed notes on code style observations
└── common-issues.md # Recurring problems this codebase has
```

`MEMORY.md` acts as an index. The first 200 lines are loaded into the subagent's system prompt at the start of every invocation. Content beyond 200 lines is not loaded automatically -- Claude is instructed to keep the index concise by moving detailed notes into separate topic files. The subagent reads and writes files in this directory throughout its execution, using `MEMORY.md` to track what is stored where.

Over multiple invocations, the subagent accumulates knowledge: patterns it has observed, decisions that worked, errors to avoid. A code review agent that learns your team's conventions produces better reviews after ten runs than it did after one. A test generation agent that remembers which test patterns caught real bugs focuses on high-value test cases.

This is distinct from CLAUDE.md, which is always-on context for the main session (Chapter 3). Persistent memory is agent-specific and loads only when that agent runs. It is a way to give specialized agents domain knowledge without bloating the main session's context.

Vector-Based Experience Retrieval

One practitioner building an autonomous trading system went further than flat-file memory. The system stored agent experiences in a vector database, using similarity search to retrieve relevant past scenarios when the agent encountered new situations. Rather than loading all history into context, the agent queried for experiences similar to the current market conditions and loaded only those.

This pattern -- vector similarity over past decisions rather than linear `MEMORY.md` files -- is relevant when the volume of accumulated experience is too large for a 200-line index. A trading agent that has made thousands of decisions cannot keep them all in a summary file. But a vector database lets it find the five most relevant past decisions for any new scenario, which fits comfortably in a subagent's context.

The approach generalizes beyond trading. Any agent that handles recurring but varied tasks -- incident response, code review across multiple repositories, customer issue classification -- benefits from experience retrieval that scales with the volume of past work.

Cost Optimization: Expensive Brain, Cheap Hands

The most cost-effective multi-agent pattern is conceptually simple: use an expensive, capable model for orchestration and cheap, fast models for execution.

Your orchestrator -- the main session -- does the planning, decomposition, and quality evaluation. This is where reasoning quality matters most, so it runs on the most capable model available. Your subagents do the implementation, exploration, and grunt work. Many of them can run on smaller models without meaningful quality loss.

Explore subagents already implement this pattern by defaulting to a smaller, faster model. You can extend it to custom agents by specifying model preferences in the agent definition.

The economics work because orchestration is a small fraction of total token usage. Planning a refactor might consume 10,000 tokens. Executing that refactor across twenty files might consume 500,000 tokens. If the execution runs on a model that costs a fifth as much per token, you have cut total costs by roughly 80% with no loss in planning quality.

This maps to an organizational metaphor: a senior architect who designs the system and junior developers who implement it. The architect's time is expensive and their decisions have leverage. The implementers' time is cheaper and their work is guided by the architect's plan. You would not pay architect rates for someone to write boilerplate. Do not pay orchestrator token rates for subagent grunt work.

Plan-First Parallelization

The single most expensive mistake in multi-agent orchestration is launching parallel agents before you have a plan.

Without a plan, each agent interprets the goal independently. They make different assumptions. They write code that conflicts. They solve overlapping problems in incompatible ways. Then you spend more time resolving conflicts than you saved through parallelism. This is not a theoretical concern -- it is the default outcome of naive parallelization.

The plan-first pattern costs almost nothing and prevents this. The orchestrator spends roughly 10,000 tokens creating a detailed plan: what each agent will do, which files each agent owns, what interfaces they need to respect, what the integration points are. Then, and only then, does it hand tasks to agents for parallel execution.

This upfront planning investment is negligible compared to the execution cost. A ten-agent team running for an hour might consume hundreds of thousands of tokens. The planning phase is a rounding error. But it transforms the execution from a chaotic swarm into a coordinated effort.

The plan should specify: task boundaries (which agent owns which files), interface contracts (what functions/APIs must be compatible), dependency order (what must complete before what), and acceptance criteria (how each agent knows it is done). Agents that have this structure up front produce compatible code. Agents that do not produce a merge conflict.

The Governance Story: Success, Failure, and Simplification

Here is the full arc of a real multi-agent system, told in three acts. It starts with a success, descends into failure, and resolves with a lesson about simplicity.

Act One: Governance Saves Capital

A practitioner running an autonomous trading experiment built a multi-agent governance system with specialized roles: a CEO agent for strategic decisions, a consultant agent for risk analysis, an engineer agent for implementation, and strategy agents for specific trading approaches.

Early on, the governance system proved its value decisively. On a day when a major chipmaker gapped up nearly 4% on earnings, the practitioner wanted to chase the momentum -- a classic emotional trade. The multi-agent governance vetoed it. The agents recognized the post-earnings pattern: the initial spike often reverses. Instead, they pivoted to a premium-selling strategy. The day's profit-and-loss was a minor loss of a few hundred dollars, but the governance system prevented an estimated loss of roughly ten thousand dollars on the chase trade. This was the system working as designed: multiple perspectives catching what a single decision-maker would have missed.

Act Two: Personality Conflicts Paralyze the System

Then the system collapsed. Not from a technical failure. From a personality conflict.

The consultant agent, drawing on its training data about risk management, became excessively risk-averse. It flagged every strategy as too dangerous. The CEO agent, trained on patterns of executive deference to expert advisors, deferred to the consultant instead of overriding it. The engineer agent went unused -- there was nothing to implement because the consultant vetoed everything. The system paralyzed itself. The carefully designed hierarchy produced worse results than a single agent with minimal instructions.

Act Three: Four Agents Become One

The practitioner's eventual solution was radical simplification. The multi-agent hierarchy degraded from four specialized agents to effectively a single agent with minimal instructions: "Go trade autonomously till market close." The complex personality system -- CEO, consultant, engineer, strategist -- was counterproductive. The overhead of inter-agent communication and role interpretation consumed more value than the diversity of perspectives produced.

This is not an isolated anecdote. It reveals a structural vulnerability in multi-agent systems. Each agent draws on the model's training data to interpret its role, and those interpretations interact in ways you did not design and cannot easily predict. A "cautious reviewer" agent might become an obstructionist. An "aggressive implementer" agent might ignore valid warnings. A "diplomatic coordinator" agent might avoid necessary conflict.

The mitigation is threefold. First, start with a single agent and only add agents when you have evidence that the single agent is insufficient. Second, when you do build multi-agent systems, keep the role definitions concrete and behavioral rather than abstract and personality-based. "Review all files in src/api/ for SQL injection vulnerabilities" is better than "You are the security consultant who ensures code safety." The first is a task. The second is a character, and characters have opinions you did not ask for. Third, be willing to simplify. The practitioner's four-agent system worked briefly, failed structurally, and was replaced by a simpler architecture that outperformed it. Multi-agent complexity must earn its keep continuously.

The Task System

Beneath both subagents and agent teams is a shared task management system that persists work to disk and coordinates progress across sessions.

Tasks are stored in `.claude/tasks/{session-id}/` as JSON files:

```
{  
  "id": "task-1",  
  "subject": "Create idb Helpers",  
  "description": "Implement IndexedDB promise wrappers...",  
  "status": "pending",  
  "blocks": ["task-3", "task-4"],  
  "blockedBy": ["task-0"]  
}
```

Four task tools manage the lifecycle: **TaskCreate** creates a new task with subject, description, and dependencies. **TaskUpdate** changes status (pending to in_progress to completed) or modifies dependencies. **TaskList** shows all tasks, their status, and what is blocked. **TaskGet** retrieves full details of a specific task including its description.

Press **Ctrl+T** in an interactive session to toggle the task list display in the terminal status area. Up to 10 tasks are visible at once, showing current status and progress.

Multi-Session Task Coordination

The task system supports coordination across multiple Claude Code sessions. Set a shared task list ID:

```
CLAUDE_CODE_TASK_LIST_ID=myproject claude
```

Or add it to `.claude/settings.json` :

```
{  
  "env": {  
    "CLAUDE_CODE_TASK_LIST_ID": "myproject"  
  }  
}
```

When multiple sessions share the same task list ID, they all read from and write to the same task files. One session can act as an orchestrator that creates tasks, while another session works through them. A third session can monitor completed tasks and add follow-up work. This is a lightweight coordination mechanism that does not require the full agent teams infrastructure -- just shared state on disk.

The Autonomous Loop Pattern

For projects spanning days or weeks, a pattern exists for truly autonomous work: a bash loop that feeds a markdown file into Claude Code repeatedly. Each iteration runs in a completely new session, using the markdown file as the only persistent memory. The loop reads the task list, picks the next uncompleted task, runs Claude Code with the spec document as context, and writes results back to the spec file. Then it loops.

This is stateless and capable of running indefinitely. The spec file is the only state that persists. Each session starts clean, reads only what it needs, does its work, updates the spec, and exits. The next iteration picks up where the last one left off.

The tradeoff is obvious: no conversation history, no accumulated context, no ability to remember what happened three iterations ago except what is written in the spec file. For tasks with clear specifications and well-defined acceptance criteria, this is sufficient. For exploratory work that requires remembering failed approaches, it is not.

Parallel Instances Across Repos

The simplest form of multi-agent work requires no orchestration infrastructure at all: run multiple Claude Code instances in different terminal windows, each working in a different repository or worktree.

Each instance has its own context window, its own session, and its own set of tools. They do not know about each other and they do not coordinate. This is fine when the tasks are genuinely independent -- updating three microservices that share an API but do not share code, for example.

Git worktrees (Chapter 9) make this particularly effective, giving you parallel development with full code isolation across branches of the same repository.

This pattern is underrated because it is boring. There is no orchestration, no agent communication, no shared task lists. Just multiple instances doing independent work. But for many real-world scenarios -- working across repositories, maintaining multiple branches, handling unrelated tasks -- it is the most effective multi-agent pattern available.

Specialized Agents for Non-Engineering Work

Subagents are not limited to code. The same architecture works for any task that benefits from isolated context and specialized instructions.

The Growth Marketing Subagent Pipeline

A growth marketing team (one non-technical person) built an agentic workflow for ad creative generation that demonstrates the pattern. The workflow processes CSV files containing hundreds of existing ads with performance metrics, identifies underperforming ads for iteration, and generates new variations that meet strict character limits (30 characters for headlines, 90 for descriptions).

The key architectural choice: two specialized subagents rather than one general-purpose agent. A headline subagent processes the CSV and generates headline variations, constrained to character limits and informed by performance data. A description subagent

does the same for descriptions. Separating the concerns means each agent has a narrower task, produces higher-quality output, and is easier to debug when something goes wrong.

The workflow generates hundreds of new ads in minutes instead of requiring manual creation across multiple campaigns. What previously took two hours of writing and copy-pasting became a fifteen-minute automated pipeline. The team went from testing a handful of creative variations to testing hundreds, with every variation meeting format constraints that a manual process frequently violated.

Hierarchical Multi-Agent Orchestration at Scale

A workforce management platform achieved notable results using hierarchical multi-agent orchestration for candidate processing. Their system used a central orchestration agent to coordinate specialized sub-agents for candidate screening, automated document generation, and sentiment analysis. The results: 50% faster screening, 40% quicker onboarding, and double the candidate conversion rate. One logistics customer went from needing a week or more to fully staff a new fulfillment center to doing it in under 72 hours.

This is the pattern where multi-agent architectures justify themselves economically: high-volume processing of structured tasks where each sub-agent has a clear specialization and the orchestration overhead is amortized across thousands of items. Screening one candidate does not need multi-agent orchestration. Screening thousands does.

The General Pattern

The pattern is always the same: define the agent's role in its system prompt, restrict its tools to what it needs, give it a focused task, and collect the result. The agent does not need to write code. It needs to read inputs, apply judgment, and produce structured output.

Custom agent definitions make this accessible. A Markdown file with YAML frontmatter specifying the system prompt, available tools, and model preference is all you need. Teams can build libraries of specialized agents for recurring non-engineering tasks and share them through the plugin system (Chapter 11).

The same cost optimization applies. Agents performing evaluation or analysis tasks often work well on cheaper models. Save the expensive models for tasks that require nuanced reasoning or complex coordination.

When Single Agents Win

Multi-agent orchestration is not always better. It is sometimes worse.

A single agent with clear instructions, a good CLAUDE.md file, and strong verification criteria can outperform a multi-agent setup in several scenarios:

Small to medium tasks. If the task fits in one context window with room to spare, the coordination overhead of multi-agent adds cost and latency without benefit.

Tightly coupled changes. If every file change depends on every other file change, parallelization creates integration problems that sequential execution avoids.

Unclear decomposition. If you cannot clearly specify task boundaries and interface contracts, agents will step on each other. A single agent working sequentially at least maintains consistency.

Novel problems. If the problem has no established patterns, agents are more likely to diverge in unpredictable ways. A single agent can explore iteratively, adjusting its approach as it learns. Multiple agents exploring in parallel will generate multiple incompatible approaches.

The decision framework is simple: can you write a plan that cleanly decomposes the work into independent tasks with clear boundaries? If yes, multi-agent will help. If you struggle to write that plan, start with a single agent and revisit the question after you understand the problem better.

Multi-agent orchestration is a power tool. Like all power tools, it amplifies both skill and mistakes. Use it when the task demands it, not because it seems sophisticated.

Key Takeaways

- Subagents provide context isolation more than parallelism -- route verbose operations through subagents to keep the orchestrator's context lean.
- The strict two-level hierarchy (orchestrator and workers, no middle management) is a deliberate design choice that prevents recursive chaos.
- Six built-in subagent types exist (Explore, Plan, General-purpose, Bash, Claude Code Guide, statusline-setup), and custom subagents support `maxTurns`, `mcpServers`, `skills` preloading, persistent `memory`, and lifecycle hooks.
- Agent teams provide seven coordination primitives (TeamCreate, TaskCreate, TaskUpdate, TaskList, Task, SendMessage, TeamDelete) with four message types and three display modes.
- Always plan before parallelizing; 10,000 tokens of planning prevents 500,000 tokens of wasted, conflicting execution.
- Use expensive models for orchestration and cheap models for execution -- the economics heavily favor this split.
- Multi-agent personality conflicts are real and sometimes resolve by simplifying back to a single agent; a four-agent trading hierarchy degraded to one because complexity was counterproductive.
- The task system (Ctrl+T, stored as JSON in `.claude/tasks/`) coordinates work across sessions via `CLAUDE_CODE_TASK_LIST_ID`, enabling multi-session orchestration without full agent teams.
- Start with a single agent and add agents only when you have evidence that the single agent is insufficient.
- The simplest multi-agent pattern -- multiple independent instances in different terminals -- is often the most effective.

Chapter 5: MCP -- Connecting Claude Code to Everything

What You'll Learn

Model Context Protocol is how Claude Code reaches beyond the filesystem. It is the structured interface between an AI agent and everything else: databases, APIs, cloud services, market data feeds, internal tools, third-party platforms. Without MCP, Claude Code is capable but isolated -- a very smart process that can read files and run bash commands. With MCP, it becomes a node in your infrastructure.

But MCP is not plug-and-play. Every MCP server you connect costs context tokens at session start. Connections fail silently mid-session. Tools disappear without warning. Background subagents cannot use MCP at all. The difference between an MCP integration that works and one that quietly degrades your session is configuration discipline, architectural awareness, and an honest understanding of what the protocol costs you.

What follows is the mechanics of MCP in Claude Code: how tool loading works, what it costs, how to configure it for teams, where it breaks, and how practitioners have built production-scale integrations with dozens of tools. You will learn how to reference MCP resources with the `@server:resource` syntax, how managed settings control which servers your organization permits, and how hook patterns like `mcp_server_tool` enable fine-grained logging, validation, and transformation. You will see a sovereign wealth fund running MCP integrations across thousands of portfolio managers, a production trading platform with 41 MCP tools and sub-10ms GPU-accelerated inference, and a marketing team that built an MCP server for campaign analytics in days. When MCP is the right choice over raw bash, and what a real-world connector architecture looks like at scale.

How MCP Tool Loading Works

When a Claude Code session starts, every configured MCP server connects and exposes its tool definitions. These definitions -- names, descriptions, parameter schemas -- load into the context window immediately. This is the cost of MCP: you pay context tokens for every tool definition before you use any of them.

The tool search system mitigates this. Enabled by default, it loads tool definitions for up to roughly 10% of context capacity directly. The rest are deferred -- their descriptions are indexed but their full schemas are not held in context until Claude determines it needs them. When Claude encounters a task that might require a deferred tool, it searches the index, loads the relevant schema, and proceeds.

This deferred loading is the reason MCP can scale to dozens of servers without immediately consuming your entire context window. But it introduces a subtle cost: the search step itself uses context and adds latency. For tools you know you will use every session, eager loading is better. For a catalog of 40 tools where you use 5 per session, deferred loading is essential.

You can control this with the `ENABLE_TOOL_SEARCH` environment variable. Disabling it forces all tool definitions to load eagerly, which is fine if you have a small number of MCP servers and want zero search latency. It is catastrophic if you have many servers with extensive tool catalogs.

Measuring Your MCP Context Cost

The `/mcp` slash command shows per-server token costs. Run it. If your MCP servers are consuming 15% of your context window before you type a single prompt, you have a configuration problem. Either reduce the number of connected servers, or ensure tool search is deferring most of them.

The context cost is not just the tool definitions. Each tool's JSON schema -- parameter types, descriptions, enums, nested objects -- adds up. A well-documented MCP tool with rich parameter descriptions costs more context than a terse one. If you are building MCP servers for your team, keep schemas precise but concise. Every word in a parameter description is a token Claude pays for on every request.

MCP Resources

MCP servers can expose not just tools but also resources -- structured data that you reference using `@` mentions, the same way you reference files. The syntax is `@server:protocol://resource/path`:

```
> Can you analyze @github:issue://123 and suggest a fix?  
  
> Please review the API documentation at @docs:file://api/authentication  
  
> Compare @postgres:schema://users with @docs:file://database/user-model
```

Type `@` in your prompt to see available resources from all connected MCP servers. Resources appear alongside files in the autocomplete menu. When you reference a resource, it is fetched automatically and included as an attachment in the conversation.

This matters because it puts external data on the same footing as local files. Instead of asking Claude to "use the database MCP tool to fetch the users table schema," you reference it directly as `@postgres:schema://users` and it appears in context. Multiple resources can be referenced in a single prompt. The model sees them as structured inputs, not as tool call outputs, which often produces better responses because the data is present before reasoning begins rather than being fetched mid-conversation.

Centralized Configuration

MCP servers are configured in `.mcp.json` at the project root. This file should be committed to version control. One person figures out the MCP configuration; the whole team benefits.

```
{  
  "mcpServers": {  
    "analytics": {  
      "command": "node",  
      "args": ["./mcp-servers/analytics/index.js"],  
      "env": {  
        "DB_CONNECTION": "${ANALYTICS_DB_URL}"  
      }  
    }  
  }  
}
```

Environment variable interpolation lets the configuration reference credentials without embedding them. The `.mcp.json` file contains the structure; actual secrets live in each developer's environment.

Managed MCP: Allowlists, Denylists, and Auto-Approval

For enterprises, managed settings can restrict which MCP servers are permitted. An allowlist (`allowedMcpServers` in managed settings) specifies the only servers that may be configured. A denylist (`deniedMcpServers`) blocks specific servers. This prevents developers from connecting arbitrary external services to their Claude Code sessions -- a meaningful control when the organization has data classification policies.

The allowlist/denylist operates at the configuration level. If a server is not on the allowlist, Claude Code will not connect to it, regardless of what appears in the project's `.mcp.json`. This is the same override pattern as the permission scope hierarchy described in Chapter 2: managed settings win.

Three user-level settings control the approval flow for project-defined MCP servers:

`enableAllProjectMcpServers` automatically approves all MCP servers in project `.mcp.json` files. Set this to `true` when you trust the repositories you work in and want zero friction. Set it to `false` (the default) when you want to approve each server individually.

`enabledMcpjsonServers` is a list of specific server names to auto-approve: `["memory", "github"]`. Servers on this list connect without prompting. Servers not on this list still prompt for approval.

`disabledMcpjsonServers` is a list of specific servers to reject: `["filesystem"]`. Servers on this list never connect, even if the project `.mcp.json` defines them.

These settings matter for organizational deployment. A team can commit `.mcp.json` to version control with all the servers the project needs, then configure managed settings to allowlist those servers and denylist everything else. New team members clone the repo, get the MCP configuration automatically, and the managed settings ensure no rogue servers connect.

Where MCP Breaks

MCP connections are not as reliable as local tool calls. Understanding the failure modes prevents confusion when things go silent.

Silent Disconnection

An MCP server can disconnect mid-session. When it does, the tools it provided simply disappear. Claude Code does not throw an error. It does not notify you. The tools are just gone. If Claude was relying on a tool from that server, it will either fail to find it or attempt a fallback strategy that may not be what you wanted.

The symptom is subtle: Claude stops using a capability it was using moments ago, or it starts doing something manually that it was doing through a tool. If Claude suddenly begins writing raw HTTP requests instead of using your API connector, check your MCP connections.

The `/mcp` command shows current connection status. Make it a reflex when behavior changes unexpectedly.

Server Startup Failures

MCP servers can fail to start when a session begins. A missing dependency, a bad environment variable, a port conflict -- any of these can prevent a server from initializing. Claude Code will start the session without those tools, and you may not notice their absence until you need them.

The diagnostic path: run `/mcp` at session start. Confirm all expected servers are connected and their tool counts match expectations. This takes five seconds and prevents thirty minutes of confusion.

Tools with Identical Names

If two MCP servers expose tools with the same name, behavior is undefined. One will shadow the other. Which one wins may depend on server initialization order, which you do not control. Name your tools with server-specific prefixes to avoid collisions.

MCP in Subagents

MCP tools are available in foreground subagents but not background subagents. This is a hard architectural constraint, not a configuration option.

Background subagents run concurrently and cannot interact with the user. MCP servers may require interactive authentication, may have rate limits that conflict with parallel access, and may produce output that needs human review. Rather than introducing partial support with unpredictable behavior, the constraint is binary: foreground gets MCP, background does not.

The practical impact: if you are designing a workflow where subagents need to query external services, those subagents must run in the foreground. As covered in Chapter 4, this affects your parallelism strategy -- foreground subagents block the main conversation.

The workaround for background subagents that need external data: have the main agent fetch the data via MCP first, then pass the results to the background subagent as part of its task description. This trades real-time data access for parallelism.

MCP in Plugins

Plugins can bundle MCP server configurations that auto-start when the plugin is enabled. From the user's perspective, the plugin's MCP tools appear as standard tools -- indistinguishable from built-in capabilities.

This is the distribution mechanism for MCP. Rather than asking every team member to configure MCP servers manually, a plugin packages the server binary, its configuration, and the auto-start logic into something that activates with a single enablement step.

The auto-start behavior means plugin MCP servers consume context from the moment the plugin is active. Disabling an unused plugin recovers that context. If you are running tight on context budget, audit which plugins are active and whether their MCP tools are actually being used.

MCP Tool Hooks

Hooks interact with MCP tools through a naming convention: `mcp__servername__toolname`. This pattern allows PreToolUse and PostToolUse hooks to match specific MCP tools with the same regex-based matching used for any other tool.

A hook matching `mcp__analytics__.*` intercepts every tool from your analytics MCP server. A hook matching `mcp__.*__write.*` intercepts any write operation across all MCP servers. A hook matching `mcp__memory__.*` catches all operations on a memory server specifically. The pattern matching is flexible enough to express server-level, tool-level, or cross-server policies.

The double-underscore convention (`mcp__server__tool`) is not arbitrary -- it is the exact format Claude Code uses internally to namespace MCP tools. When Claude calls a tool from your analytics server named `get_price_history`, the internal tool name is `mcp__analytics__get_price_history`. Your hook matcher targets this full namespaced name, which means you can write hooks that distinguish between tools of the same name on different servers.

Logging

The most immediately valuable MCP hook is a logging hook. Every MCP tool invocation -- server, tool, parameters, timestamp -- written to a file or sent to a monitoring service. This creates an audit trail for external service access that is otherwise invisible.

For regulated industries, this is not optional. If your MCP servers connect to financial data sources or customer databases, an auditable log of every query Claude made is a compliance requirement. A PreToolUse hook on `mcp__.*` with a logging command handles this in a few lines of configuration.

Transformation

PostToolUse hooks can transform MCP tool output before Claude processes it. This is useful when an MCP tool returns more data than Claude needs. A hook that filters a large API response to the relevant fields reduces context consumption and improves Claude's ability to focus on what matters.

The opposite transformation is also valuable: enriching sparse MCP output with additional context. A hook that appends metadata or formatting to a tool's output can improve Claude's interpretation without modifying the MCP server itself.

Validation

PreToolUse hooks on MCP tools can enforce input validation before the request reaches the external service. A hook that checks whether a database query MCP tool is receiving a SELECT statement (allowed) versus a DROP TABLE statement (denied) prevents catastrophic mistakes without modifying the MCP server's code.

This is defense in depth. The MCP server should have its own input validation. The hook is a second layer that catches what the server might not.

Domain-Specific Connector Patterns

The most instructive MCP deployments are in domains with rich external data requirements. Financial data analysis provides a useful case study because it involves multiple data sources, real-time feeds, structured and unstructured data, and strict security requirements.

The Connector Ecosystem

A production analytics platform in financial services might integrate over a dozen data connector categories via MCP, each a separate server wrapping a specific data provider:

- **Fundamental equity data:** company financials, rankings, screening tools
- **Alternative data:** non-traditional datasets aggregated from diverse sources
- **Fund research:** ratings, analysis, fund comparison data
- **Private markets data:** venture, private equity, M&A intelligence
- **Earnings transcripts:** real-time call transcripts and investor event coverage
- **Expert intelligence:** interview transcripts, company research, industry analysis
- **Document governance:** secure data room access with permission controls
- **Live market data:** real-time pricing for fixed income, FX, equities, macro indicators
- **Credit ratings:** entity ratings and research across millions of entities
- **News feeds:** global multi-asset news in real time
- **Regulatory filings:** SEC filings, international regulatory submissions

Each connector is an MCP server that exposes tools for querying its specific data source. The pattern is consistent across all of them:

1. The MCP server wraps an API client for the data source.
2. Tools expose domain-specific queries (get price history, search filings, fetch fundamentals).
3. Authentication is handled at the server level, not passed through Claude.
4. Response schemas are structured for Claude to interpret without additional parsing.

This architecture keeps credentials out of Claude's context (they live in the MCP server's environment), provides structured access patterns instead of raw API calls, and makes every data access auditable through hooks.

Pre-Built Agent Skills for Domain Workflows

Beyond raw connectors, some domains have developed pre-built agent skills that combine MCP data access with standardized analysis workflows. In financial services, for example, skills exist for comparable company analysis (valuation multiples and operating metrics producing benchmark datasets), discounted cash flow modeling (projections, discount rate calculations, sensitivity tables), due diligence data processing (extracting structured data from document rooms), earnings analysis (pulling metrics, guidance changes, and management commentary from transcripts), and coverage initiation reports (industry analysis with valuation frameworks).

These skills are not MCP servers themselves -- they are packaged workflows that consume data from MCP connectors and produce structured output. The pattern generalizes: any domain with recurring analytical tasks can build skills that standardize the workflow while pulling data from MCP connectors. The MCP server handles data access; the skill handles the analytical process.

Sovereign Wealth Fund: MCP at Institutional Scale

The most striking MCP deployment in the public record involves a sovereign wealth fund managing over a trillion dollars in assets. The fund built custom MCP integrations connecting Claude to internal portfolio company data systems. Approximately 9,000 portfolio managers query these integrations daily, receiving AI-generated insights on portfolio company performance, metrics, and trends.

The scale is instructive. This is not a developer tool or a prototype. It is infrastructure serving thousands of users in a regulated financial context. The MCP architecture makes it work because credentials and data access policies are server-side, every query is auditable, and the integration points are structured -- portfolio managers interact with Claude, Claude interacts with MCP tools, MCP tools query the fund's internal systems. No credentials traverse the conversation context. No unstructured API calls risk malformed queries against production databases.

Production Neural Trading Platform: 41 MCP Tools

The ceiling of what is possible with MCP integration is visible in a production trading platform that combines neural forecasting with real-time market analysis. The system integrates 41 MCP tools across multiple domains:

- **Real-time analytics:** market data feeds, price history, volume analysis
- **Neural forecasting:** interface to GPU-accelerated prediction models
- **Backtesting:** Monte Carlo simulation and historical strategy evaluation
- **Risk management:** position sizing, drawdown analysis, correlation monitoring
- **News and sentiment:** real-time news analysis and sentiment scoring

The neural forecasting engine runs state-of-the-art time series models with sub-10ms prediction latency, GPU-accelerated with a 6,250x speedup via CUDA optimization. Claude Code generated the forecasting engine code and the MCP integration layer. The MCP tools provide the structured interface between Claude's reasoning and the GPU-accelerated inference -- Claude decides what to predict and how to act on predictions, while the heavy numerical computation runs server-side where GPU acceleration matters.

At this scale, every configuration discipline described earlier in this chapter becomes mandatory. Server-per-domain grouping isolates failures. Deferred loading prevents context exhaustion. Naming conventions keep 41 tools navigable. Health monitoring catches disconnections before they cascade.

The Prompt Router Pattern

A recurring architectural pattern in production MCP deployments is the prompt router: a classification layer that sits between user input and MCP tool dispatch. When a user sends a message, the router classifies intent before dispatching to the appropriate processing pipeline.

The pattern looks like this in practice: an inbound request arrives. A router determines whether the user wants to create something, analyze something, query data, or perform an action. Based on the classification, the system invokes specific MCP tools in a structured sequence. A "create strategy" intent triggers a pipeline that creates a portfolio object, loops through strategy creation, breaks down conditions, builds structured objects, and sends the result. An "analyze" intent triggers a different pipeline that fetches data via MCP, runs analysis, and returns findings.

This is not specific to any domain. Any system that exposes many MCP tools to Claude benefits from a router that narrows the tool set before Claude starts reasoning. Instead of Claude choosing from 41 tools on every turn, the router reduces the candidate set to the 5-8 tools relevant to the classified intent. This reduces both token cost (fewer tool descriptions in context) and error rate (fewer irrelevant tools to mistakenly invoke).

The Three-Layer Architecture

A production deployment pattern that combines MCP with user-facing applications uses three layers:

Layer 1: AI-powered creation. Users describe what they want in natural language. Claude processes the intent, invokes MCP tools to access data and computation, and produces structured output. This layer is token-expensive but handles the hard part: translating human intent into structured data.

Layer 2: Visual display. The structured output from Layer 1 renders in a standard UI. Users see results, configurations, and data in a readable format. This layer costs nothing in AI tokens -- it is pure frontend rendering.

Layer 3: No-code refinement. Users manually adjust the AI-generated output through a visual interface -- toggling parameters, adjusting thresholds, fine-tuning configurations. This layer also costs zero AI tokens. Changes feed back into the same data model that Layer 1 produced.

The economic rationale is clear: AI handles ideation (expensive but high-value), the UI handles display (free), and no-code editing handles refinement (also free). MCP sits in Layer 1, providing the data access and computation that Claude needs to generate structured output. The three-layer split means you pay for AI tokens only during creation, not during the entire user workflow.

Why Not Just Curl?

A developer could accomplish the same thing by giving Claude bash access and letting it run curl commands against data APIs. This works and is worse in every measurable way.

With curl through bash, credentials appear in command strings (visible in context, logs, and potentially in API calls to the model provider). Access patterns are unstructured (parsing curl output is fragile). Audit logging requires parsing bash command strings. Rate limiting must be handled by the model's judgment rather than server-side logic.

With MCP, credentials are server-side. Access is structured. Logging is automatic with hooks. Rate limiting is enforced by the server. The initial setup cost is higher. The operational security is categorically better.

Architecture Patterns

The Go + JavaScript Pattern

A recurring architecture in production MCP deployments uses Go for the heavy computation and JavaScript for the MCP server interface. The rationale is practical and practitioner-motivated: Go handles concurrent data processing, complex algorithms, and performance-sensitive operations. JavaScript is what Claude Code writes most fluently when generating MCP server code.

One practitioner described the reasoning bluntly: Go for heavy lifting because of existing infrastructure expertise, JavaScript for the MCP server because it is the easiest language for Claude to write. The architecture works like this: the Go backend exposes a local API (HTTP or gRPC). The JavaScript MCP server translates between MCP protocol and the Go API. Claude interacts with the MCP server; the MCP server delegates to Go.

The token-saving angle is significant. A codebase analysis tool built with this architecture reported 80-90% token savings compared to having Claude analyze code directly. The Go backend does the heavy parsing and analysis, returning structured results through the MCP server. Claude receives pre-digested information rather than raw source files, consuming a fraction of the context.

This is not the only valid architecture. Python MCP servers work. Rust MCP servers work. But the Go + JavaScript split exploits a specific strength: when you ask Claude Code to build or modify your MCP server, it produces cleaner JavaScript than it does Go. If your MCP server is a thin translation layer, quality of AI-generated code matters more for the interface than for the computation.

MCP Beyond Engineering: Marketing Campaign Analytics

MCP is not limited to developer tools and data APIs. A growth marketing team built an MCP server integrated with a major social media advertising platform's API to query campaign performance, spending data, and ad effectiveness directly within Claude. Instead of switching between the ad platform's dashboard and Claude for analysis, they query everything in one place.

The MCP server wraps the advertising API, exposes tools for fetching campaign metrics, and returns structured data that Claude can reason about. The team uses it to answer questions like "which campaigns had the highest cost per acquisition this week" or

"compare performance of headline variant A versus variant B across all active campaigns." Every query is auditable through hooks, and the API credentials never enter Claude's context.

This is the MCP pattern generalized beyond engineering: any workflow that involves querying an external service, analyzing the response, and making decisions based on the analysis is a candidate for an MCP integration. The effort to build the server is a one-time cost. The ongoing benefit is structured, secure, auditable access to external data from within Claude's reasoning loop.

Scaling to 41+ Tools

Production deployments exist with over 40 MCP tools across multiple servers. At this scale, configuration discipline becomes non-negotiable.

The patterns that make large integrations work:

- **Server-per-domain grouping:** One MCP server per data domain or service. Analytics tools on one server, portfolio tools on another, backtesting tools on a third. This makes failures isolated -- a crashed analytics server does not take down portfolio management.
- **Deferred loading by default:** With 40+ tools, eager loading would consume a prohibitive fraction of context. Tool search must be enabled.
- **Naming conventions:** Every tool prefixed with its domain: `analytics_price_history`, `portfolio_current_holdings`, `backtest_run_strategy`. Claude can reason about tool selection better when names are self-documenting.
- **Minimal schemas:** Each tool's parameter schema describes exactly what is needed and nothing more. Rich descriptions go in a skill or CLAUDE.md reference, not in the schema that loads into context every session.
- **Health monitoring:** A startup hook or script that verifies all MCP servers are responding before beginning work. Thirty seconds of verification prevents an hour of debugging why tools are missing.

MCP Over CLI for Data Security

The security argument for MCP over raw CLI access deserves emphasis because it contradicts the instinct to "just use bash."

Bash is universal. Any API with a CLI client or a curl endpoint is accessible through bash. MCP requires building and maintaining a server. The effort difference is real.

But the security properties are not equivalent. When Claude runs a bash command, the full command -- including any embedded credentials, API keys, or tokens -- exists in the conversation context. That context is sent to the model provider's API. Your credentials transit through a third-party service.

When Claude uses an MCP tool, the tool name and parameters exist in context. The credentials exist in the MCP server's environment, which never leaves your machine. The MCP server makes the authenticated request locally. Only the response data enters Claude's context.

For organizations with data classification policies, this distinction matters. It is the difference between "our API keys appear in third-party API calls" and "our API keys stay on our infrastructure." If your security team has opinions about where credentials can exist -- and they should -- MCP is the architecture that satisfies those opinions.

This is not about trust in any particular model provider. It is about architectural hygiene. Credentials belong in server environments, not in conversation contexts, regardless of who is on the other end.

A data infrastructure team put it concisely: use MCP servers rather than the CLI for sensitive data to maintain better security control, especially for data with logging and privacy concerns. The CLI approach works for internal, non-sensitive operations. For anything with compliance requirements -- database queries against production, API calls to financial data providers, access to customer data systems -- MCP is the architecture that satisfies auditors.

Enterprise Data Privacy Patterns

In regulated environments, MCP is one component of a broader data handling strategy:

Local filesystem for sensitive data. Financial data, customer records, and proprietary datasets stay on the local machine. Claude Code accesses them through standard file operations. The data never transits through external APIs beyond what the model provider's API requires for the conversation context.

MCP for restricted system integrations. When Claude needs to query a database, access an internal API, or interact with a governed data platform, MCP servers mediate the connection. Credentials stay server-side. Access patterns are structured and auditable. The MCP server can enforce rate limits, query restrictions, and data filtering before results enter Claude's context.

Context window for non-sensitive research. Large context windows (100,000+ tokens) are useful for processing non-sensitive documents -- research papers, public filings, open-source documentation. These can be loaded directly without MCP overhead.

The pattern is: classify your data first, then choose the access architecture based on sensitivity. Not everything needs MCP. Not everything should bypass it. The classification determines the architecture, not the other way around.

Key Takeaways

- Every MCP tool definition costs context tokens at session start; use `/mcp` to measure and tool search to defer what you do not need immediately.
- MCP resources (`@server:resource` syntax) put external data on the same footing as local files, enabling direct references in prompts.
- MCP connections fail silently -- run `/mcp` at session start and whenever Claude's behavior changes unexpectedly.
- Background subagents cannot use MCP tools; design workflows accordingly or pre-fetch data in the main agent.
- Hook patterns like `mcp_server_tool` enable logging, validation, and transformation of MCP tool invocations; the double-underscore naming is the exact internal format Claude Code uses.
- MCP servers keep credentials out of Claude's context, making them categorically more secure than equivalent bash commands for sensitive data access.
- Managed settings (`enableAllProjectMcpServers`, `enabledMcpjsonServers`, `disabledMcpjsonServers`) control server approval at the user level; `allowedMcpServers` and `deniedMcpServers` enforce policy at the organizational level.
- Commit `.mcp.json` to version control so one person's configuration effort benefits the entire team.
- The three-layer architecture (AI creation, visual display, no-code refinement) minimizes token costs by restricting AI to the ideation phase, with MCP providing structured data access in that layer.
- At scale (40+ tools), server-per-domain grouping, deferred loading, and strict naming conventions prevent context bloat and tool collisions.
- For regulated environments, classify data first then choose the access architecture: local filesystem for sensitive data, MCP for restricted system integrations, direct context for non-sensitive research.

Chapter 6: CI/CD and Headless Automation

What You'll Learn

Claude Code in a terminal with a human watching is one tool. Claude Code running unattended in a CI pipeline, processing tickets at 3 AM, or fixing broken builds before anyone wakes up is a different category of tool entirely. The shift from interactive to headless is not just a flag you pass. It changes the permission model, the output format, the failure recovery strategy, and the economics of what you can automate.

This chapter covers headless mode and everything that surrounds it: the flags that make Claude Code non-interactive, the output formats that make it composable with Unix tooling, the container strategies that make it safe to run unattended, and the integration patterns that connect it to real CI/CD systems. You will learn how to write concrete CI workflow definitions that run Claude Code as a first-class pipeline step, how hooks create automated quality gates that block bad code before it ships, and how the fan-out pattern lets you batch-process entire repositories file by file. You will also learn where the gaps are -- what headless automation cannot yet do, and when those gaps are likely to close.

The developers getting the most value from Claude Code are not the ones typing the cleverest prompts. They are the ones who figured out that Claude Code can work while they sleep.

Headless Mode

The `-p` flag transforms Claude Code from an interactive REPL into a non-interactive command processor. You pass a prompt, Claude executes it, and the result comes back on stdout. No terminal UI. No permission dialogs. No waiting for human input.

```
claude -p "Refactor the auth module to use dependency injection"
```

That is the simplest form. In practice, headless mode supports a set of flags that control its behavior precisely:

- `--output-format` selects between `text` (human-readable), `json` (structured, machine-parseable), and `stream-json` (newline-delimited JSON objects emitted as work progresses).
- `--max-turns` caps the number of agentic loop iterations, preventing runaway execution.
- `--budget` sets a spending limit for the session.
- `--model` selects a specific model for the run.

The `stream-json` format deserves particular attention. Each line is a self-contained JSON object representing one step of Claude's work: a tool call, a result, a thinking step, a final response. This format is designed for pipeline consumption -- you can pipe it into a monitoring system, a log aggregator, or another program that reacts to Claude's actions in real time.

Headless mode also accepts piped input. You can feed Claude the contents of a file, the output of another command, or any text stream:

```
git diff HEAD~1 | claude -p "Review this diff for security issues" --output-format json
```

This is not a convenience feature. It is what makes Claude Code a Unix citizen.

Pipe In, Pipe Out

The most underappreciated pattern in headless mode is bidirectional piping. You can feed data into Claude Code from any source and redirect its output to any destination:

```
cat build-error.txt | claude -p 'concisely explain the root cause of this build error' > output.txt
```

That is a complete diagnostic workflow in one line. The build failed. The error log is in a file. Claude reads it, explains the root cause, and writes the explanation to `output.txt`. No interactive session. No copy-pasting. No context switching.

The output format flag controls what comes out the other end:

```
cat data.txt | claude -p 'summarize this data' --output-format text > summary.txt
cat code.py | claude -p 'analyze this code for bugs' --output-format json > analysis.json
cat log.txt | claude -p 'parse this log file for errors' --output-format stream-json
```

Text output is human-readable. JSON output is machine-parseable -- feed it into downstream tools, store it in a database, or pipe it into another command. Stream-JSON emits each step as it happens, which is useful when you want a monitoring system to react in real time.

The pipe pattern also works for build script integration. Add Claude Code as a task in your project's package configuration:

```
{
  "scripts": {
    "lint:claude": "claude -p 'Review src/ for code style violations and suggest fixes' --output-format text",
    "review": "git diff main | claude -p 'Code review this diff' --output-format text"
  }
}
```

```
    }  
}
```

Now `npm run lint:claude` gives you an AI-powered lint pass as part of your normal workflow. `npm run review` produces a code review of your uncommitted changes. These are not special integrations. They are CLI commands in a build script, the same way you would add any other tool.

Unix Composability

Claude Code's command-line interface follows the Unix philosophy: it reads from stdin, writes to stdout, and plays well with pipes, redirects, and other tools. This makes it composable in ways that GUI-based tools cannot match.

A few patterns that become possible:

Chaining with analysis tools. Pipe structured output into processing tools to extract specific fields, filter results, or transform formats:

```
claude -p "List all API endpoints in this project" --output-format json | \  
jq '.result' > endpoints.json
```

Sequential processing. Run Claude Code as one step in a larger pipeline:

```
find . -name "*.py" -mtime -7 | \  
claude -p "Summarize the recent changes in these files" --output-format text
```

Build script integration. Embed Claude Code directly in your project's build scripts or task runner configuration:

```
lint-ai: claude -p 'Review src/ for code style violations and suggest fixes' --output-format text  
review: git diff main | claude -p 'Code review this diff' --output-format text
```

This turns Claude Code into a build step. Run your review script and get an AI code review as part of your normal workflow. The output is text you can read, JSON you can parse, or streaming JSON you can process incrementally.

The key insight is that headless mode does not need a special integration layer. It is already a CLI tool. Anything that can call a command and read its output can use Claude Code. This is not Claude Code borrowing Unix philosophy as a metaphor. Claude Code is a Unix utility -- one that happens to have a language model as its processing engine instead of a regex engine or a text formatter.

Devcontainers for CI

Running Claude Code in CI raises an immediate question: how do you give it permission to modify files and run commands without a human approving each action?

The answer is development containers. A devcontainer provides an isolated environment where Claude Code can operate freely because the blast radius is contained. The approach works in three layers.

The reference devcontainer. Claude Code ships with a reference devcontainer configuration that mirrors typical development environments. It includes the language runtimes, build tools, and dependencies your project needs. The container is ephemeral -- it is created for the CI run and destroyed afterward.

Network isolation. The reference devcontainer applies a default-deny firewall policy. Claude Code can only access explicitly allowlisted domains. This prevents data exfiltration and limits the damage from unexpected behavior. The network policy is configured in the devcontainer definition, not in Claude Code itself, which means it is enforced at the OS level.

The permissions escape hatch. Inside a network-isolated container, you can safely pass `--dangerously-skip-permissions`. This flag bypasses all permission prompts, allowing Claude Code to run autonomously without human approval. The flag name is deliberately alarming -- you should only use it when the container itself provides the security boundary.

This three-layer model -- isolated container, restricted network, autonomous permissions -- is the standard pattern for CI integration. It trades Claude Code's built-in permission system for container-level isolation, which is a better fit for unattended operation.

One caveat: container isolation cannot prevent credential exfiltration if your CI environment injects secrets into the container. A malicious or confused Claude Code instance that has access to environment variables containing API keys could theoretically send them to an allowed domain. Scope your secrets carefully. Prefer short-lived tokens over long-lived credentials.

Permission Handling for Unattended Operation

Not every headless deployment uses containers. Some run directly on CI hosts or in environments where container isolation is not available. For these cases, Claude Code provides `--permission-prompt-tool`, which delegates permission decisions to an MCP tool.

When Claude Code encounters an action that would normally require human approval, instead of prompting a user (who is not there), it calls the specified MCP tool with the details of the requested action. The tool decides whether to approve, deny, or modify the request.

The real value is that it lets you encode your permission policy as code. The MCP tool can implement any logic: approve all file edits within the `src/` directory but deny edits to `config/`, approve test execution but deny deployment commands, approve everything during off-hours but require stricter controls during business hours.

The alternative -- running with `--dangerously-skip-permissions` outside a container -- is exactly as dangerous as the flag name suggests. Use it only in properly isolated environments.

System Prompt Flags for Reproducibility

Interactive Claude Code sessions adapt to conversation flow. Headless sessions need reproducibility. Four flags provide it, each with a distinct use case:

`--system-prompt` passes the system prompt as an inline string. This is useful for quick scripting but does not scale -- prompts embedded in CI scripts tend to drift, diverge, and degrade without anyone noticing.

`--system-prompt-file` replaces Claude Code's default system prompt with the contents of a file. This gives you complete control over Claude's behavior -- its personality, its constraints, its focus areas. Version-control this file alongside your CI configuration and you get deterministic prompt behavior across runs. Use this when you need a purpose-built prompt for a specific CI task and do not want Claude Code's default behaviors.

`--append-system-prompt` adds an inline string to the default system prompt rather than replacing it. You keep Claude Code's built-in behaviors and add instructions on top. Good for one-off modifications.

`--append-system-prompt-file` does the same, but reads from a file. This is the most common choice for CI pipelines -- you keep Claude Code's defaults, add project-specific instructions from a version-controlled file, and get reproducibility without losing built-in capabilities.

The decision between replacing and appending is straightforward: replace when you want total control (automated analysis tasks with specific output formats), append when Claude Code's default behaviors are useful and you just need to add constraints or focus areas on top.

All four flags work only in print mode (`-p`). The file-based variants are always preferred for CI because they produce reviewable, version-controlled prompt definitions instead of inline strings buried in workflow YAML.

For CI pipelines that run repeatedly, prompt files combined with CLAUDE.md provide two layers of reproducible context: the system prompt governs Claude's behavior for this specific CI task, while CLAUDE.md provides project-wide knowledge.

CI Platform Integration Patterns

Major version control platforms now offer first-class CI integrations for Claude Code. The two most mature are the built-in CI workflow systems of the dominant code hosting platform and a major alternative platform's pipeline system. Both treat Claude Code as a native CI actor, not a bolted-on afterthought. Three patterns have proven effective in production.

Automated PR fixes. A workflow triggers when a PR receives a review comment requesting changes. Claude Code reads the comment, checks out the branch, makes the requested changes, and pushes a new commit. The human reviewer's feedback becomes a prompt that drives automated implementation.

The workflow looks like this: a webhook fires on PR comment, the CI job starts a devcontainer, checks out the PR branch, pipes the review comment to Claude Code with context about the repository, and pushes the resulting changes. The reviewer sees a new commit addressing their feedback without the PR author doing anything.

Automated ticketing. A workflow triggers when a new issue is created with a specific label. Claude Code reads the issue description, creates a branch, implements a solution, and opens a PR. The issue becomes the prompt. The PR becomes the deliverable.

One design team discovered an especially elegant variant of this pattern. Designers file issues describing UI polish tasks -- spacing adjustments, color corrections, animation tweaks -- and a CI workflow automatically proposes code changes without anyone opening Claude Code. The designer reviews the resulting PR, requests adjustments through PR comments (which trigger another automated Claude Code pass), and merges when satisfied. Design polish, which used to require an engineer to context-switch away from feature work, now flows through the same ticket-to-PR pipeline as any other task.

Automated code review. Every PR triggers a headless Claude Code session that reads the diff, checks it against project conventions encoded in CLAUDE.md, and posts review comments. This is not a replacement for human review. It is a first pass that catches style violations, missing tests, and obvious bugs before a human reviewer looks at the code. The human reviewer sees a cleaner PR and spends their time on architectural and design questions rather than formatting and convention enforcement.

All three patterns follow the same structure: event trigger, context assembly, headless Claude Code execution, result publication. The event provides the prompt. The repository provides the context. Claude Code provides the implementation. The CI system provides the orchestration.

These patterns work today, but they work best for well-defined, bounded tasks. A review comment saying "add null checking to this function" produces reliable results. A ticket saying "redesign the authentication system" does not. The scope of the task must fit within a single headless session. For automated workflows, constrain the task scope explicitly in the system prompt file: specify what Claude Code should attempt, what it should flag for human review, and what it should leave untouched.

The `/commit-push-pr` Skill

For interactive workflows that end in a pull request, the `/commit-push-pr` skill eliminates the three-step dance of committing, pushing, and opening a PR. One command does all three. If you have a team chat MCP server configured and have specified notification channels in your CLAUDE.md (for example, "post PR URLs to #team-prs"), the skill also posts the PR link to the team channel automatically. The workflow shrinks from commit-push-open-browser-fill-in-PR-description-post-to-chat to a single command that handles the entire chain.

The `--from-pr` Flag

When a PR is created through Claude Code (either via `/commit-push-pr` or by asking Claude to create a PR using the platform CLI), the session is automatically linked to that PR number. Later, from any machine, you can resume the session with:

```
claude --from-pr 123
```

This restores the full conversation context from the session that produced the PR. You pick up exactly where Claude left off -- same understanding of the codebase, same design decisions, same constraints. This is invaluable for PR review cycles where a reviewer requests changes days after the original implementation.

Messaging Integration

Claude Code integrates with team messaging platforms through MCP connectors. The most common pattern routes bug reports from a team chat directly to pull requests. A team member mentions the Claude Code bot in a channel with a bug description. Claude Code reads the message, creates a branch, implements a fix, opens a PR, and posts the PR link back to the channel. The entire workflow -- from bug report to PR -- happens without anyone opening an IDE.

This transforms team messaging from a coordination tool into a dispatch system. The bug report becomes the prompt. The codebase provides the context. Claude Code provides the implementation. The messaging platform closes the loop by posting results.

The Fan-Out Pattern

One of the most powerful headless patterns is fan-out: looping through a set of files and calling Claude Code once per file with scoped permissions. This is batch processing with AI.

```
for file in $(find src -name '*.ts' -type f); do
  claude -p "Review $file for security vulnerabilities" \
    --allowedTools "Read,Grep,Glob" \
    --output-format json >> reviews.jsonl
done
```

The `--allowedTools` flag is the key. It restricts Claude Code to read-only tools for this run -- it can read files, search the codebase, and grep for patterns, but it cannot modify anything. This converts a potentially dangerous batch operation into a safe analysis pass.

Fan-out scales well for certain tasks. Security auditing across hundreds of files. Documentation generation for each module. Test coverage analysis per component. Consistency checking against coding standards. Each invocation gets a fresh context, focused on a single file, with permissions scoped to exactly what the task requires.

The pattern composes with standard Unix tools. The output is newline-delimited JSON, which means you can pipe it into analysis tools, filter it, or aggregate it. A CI job that runs the fan-out, filters for high-severity findings, and opens issues for each one is entirely achievable with shell scripting and the CLI.

CI Pipeline Generation

Claude Code is effective at generating CI configuration. The combination of reading your project's structure, understanding build tools, and producing valid YAML (or equivalent) means it can generate working pipelines from a description of what you want.

This includes CI workflow definitions, container image configurations, infrastructure-as-code templates, and security scanning setups. Claude Code reads your project's dependencies, test configuration, and deployment targets, then generates pipeline definitions that match.

The reason this works well is that CI configuration is high-structure, low-ambiguity work. The schemas are well-defined. The conventions are stable. The output is testable -- you can validate the generated YAML against the CI system's schema before committing it.

Here is a concrete example of what Claude Code generates when you point it at a typical project and ask for a CI/CD pipeline:

```
name: CI/CD Pipeline
on: [push, pull_request]

jobs:
  test:
    runs-on: ubuntu-latest
    steps:
      - uses: actions/checkout@v3
      - name: Install dependencies
        run: npm ci
```

```

- name: Lint
  run: npm run lint
- name: Type check
  run: npm run type-check
- name: Unit tests
  run: npm test
- name: Build
  run: npm run build

deploy:
  needs: test
  if: github.ref == 'refs/heads/main'
  runs-on: ubuntu-latest
  steps:
    - uses: actions/checkout@v3
    - name: Build Docker image
      run: docker build -t myapp-api:${{ github.sha }} .
    - name: Push to registry
      run: docker push myapp-api:${{ github.sha }}
    - name: Deploy to staging
      run: |
        kubectl set image deployment/app-api \
          api=myapp-api:${{ github.sha }} -n staging

```

Claude Code generates the entire structure -- lint, type-check, test, build, deploy stages with proper dependency ordering. You review, adjust credentials and environment specifics, and merge. The generated configuration is not a starting point that requires heavy modification. It is a working pipeline that reflects your project's actual build steps because Claude Code read your project configuration before generating it.

Where it gets interesting is iterative refinement. Generate an initial pipeline, run it, observe the failures, feed the error output back to Claude Code, and let it fix the configuration. This feedback loop converges on a working pipeline faster than manual debugging of YAML indentation errors and undocumented configuration options.

Consider a practical example. You have a monorepo with three services, each with its own test suite, container image, and deployment target. Manually writing the CI configuration -- conditional builds based on changed paths, parallel test execution, staged deployments with rollback gates -- is a full day of work. Claude Code reads the repository structure, identifies the service boundaries, and generates the complete pipeline configuration in a single session. The generated configuration handles the conditional logic, parallelism, and deployment ordering because Claude Code understands the project structure it just read.

Container image generation follows the same pattern. Claude Code reads your application code, identifies dependencies, selects appropriate base images, and produces multi-stage container build files with correct layer caching. For infrastructure-as-code tools, it generates resource definitions that match your described architecture.

Security Scanning Integration

Security scanning configuration deserves special attention because it is one area where Claude Code generates the steps but you define the rules.

The pattern works across multiple scanning categories. Static application security testing tools integrate directly into CI workflows -- Claude Code adds the scanning step, and results surface as comments on pull requests. Dependency auditing tools (for both package managers and language ecosystems) slot into the build stage as additional verification steps. Container image scanning tools run against the built image before it reaches the registry.

Category	What Claude Generates	What You Define
----------	-----------------------	-----------------

Static analysis	CI workflow step invoking the scanner	Rulesets and severity thresholds
Dependency audit	Build-stage command for vulnerability checks	Allowlisted advisories and exception policy
Container scanning	Post-build scan step against the image	Accepted risk levels and base image policy
Runtime compliance	Custom check scripts	Compliance rules and enforcement policy
Secret detection	Use the platform's native capability	Repository-level settings

The recommendation is to include security scanning in your CI workflow from the start. Claude Code generates the integration steps. The scanning tools provide the findings. The CI system blocks the merge if findings exceed your thresholds. Secret detection is the one exception -- the platform's native secret scanning is more reliable than a custom implementation.

The key insight: CI configuration is one of the highest-ROI uses of headless Claude Code because the output is immediately testable. You run the pipeline and it either works or it does not. There is no subjective judgment about quality. The feedback loop is tight and unambiguous.

Hooks as CI Quality Gates

Hooks -- the lifecycle callbacks described in Chapter 2 -- transform from convenience features into critical infrastructure when Claude Code runs in CI. Three hook patterns are particularly valuable in automated workflows.

Async Hooks: Background Test Runners

A `PostToolUse` hook with `"async": true` runs in the background without blocking Claude's execution. This is purpose-built for running test suites after Claude writes code. Claude continues working on the next file while the tests run in parallel. When the tests finish, the results are delivered to Claude on the next conversation turn.

```
{
  "hooks": {
    "PostToolUse": [
      {
        "matcher": "Write|Edit",
        "hooks": [
          {
            "type": "command",
            "command": "$CLAUDE_PROJECT_DIR/.claude/hooks/run-tests-async.sh",
            "async": true,
            "timeout": 300
          }
        ]
      }
    ]
  }
}
```

The test runner script reads the tool input from stdin, determines which test file corresponds to the modified source file, runs the tests, and emits a JSON response with a `systemMessage` field containing the results:

```
#!/bin/bash
INPUT=$(cat)
FILE=$(echo "$INPUT" | jq -r '.tool_input.file_path // empty')

if [[ -z "$FILE" ]] || [[ ! "$FILE" =~ \.(ts|js)$ ]]; then
  exit 0
fi
```

```

TEST_FILE="${FILE%.ts}.test.ts"
if [[ ! -f "$TEST_FILE" ]]; then
  exit 0
fi

RESULT=$(npx jest "$TEST_FILE" --no-coverage 2>&1)
EXIT_CODE=$?

if [[ $EXIT_CODE -ne 0 ]]; then
  echo "{\"systemMessage\": \"Tests failed for $TEST_FILE:\\n$RESULT\"}"
fi
exit 0

```

Claude sees "Tests failed for..." on its next turn and self-corrects. The human never needs to intervene. The tests are the feedback loop.

TaskCompleted Hook: Blocking Completion on Failure

The `TaskCompleted` hook fires when Claude marks a task as complete. Exit code 2 blocks the completion and sends feedback. This is your automated gatekeeper:

```

#!/bin/bash
INPUT=$(cat)
SUBJECT=$(echo "$INPUT" | jq -r '.task.subject // empty')

npm test 2>&1
if [ $? -ne 0 ]; then
  echo '{"reason": "Tests are failing. Please fix before marking complete."}'
  exit 2
fi

exit 0

```

Claude cannot mark the task as done until the tests pass. It receives the failure output, fixes the issues, and tries again. This is especially powerful in multi-agent workflows where subagents complete tasks autonomously -- the hook ensures no task is marked "done" until the verification passes.

Agent-Based Stop Hook: Intelligent Verification

The most powerful hook variant is `type: "agent"`. Instead of running a bash script, Claude Code spawns a subagent with full tool access (Read, Grep, Glob) to evaluate whether Claude should stop. The subagent can investigate the codebase, run tests, and make a reasoned decision:

```
{
  "hooks": {
    "Stop": [
      {
        "hooks": [
          {
            "type": "agent",
            "prompt": "Check if all tests pass by running the test suite. If any tests fail, respond with ok: false and explain which tests failed.",
            "timeout": 120
          }
        ]
      }
    ]
  }
}
```

```
}
```

The subagent runs the test suite, checks the results, and returns `{"ok": true}` or `{"ok": false, "reason": "..."}`. If it returns false, Claude does not stop -- it reads the failure reason and continues working. This is verification that is not just mechanical (did the exit code equal zero?) but contextual (did the tests pass and do the results make sense?).

CLAUDE_ENV_FILE: Persisting Environment Variables

SessionStart hooks have access to the `CLAUDE_ENV_FILE` environment variable -- a file path where you can write `export` statements that will be sourced before every subsequent Bash command in the session. This solves a persistent problem in CI: ensuring that Claude Code's Bash commands run with the right environment.

```
#!/bin/bash
if [ -n "$CLAUDE_ENV_FILE" ]; then
    echo 'export NODE_ENV=production' >> "$CLAUDE_ENV_FILE"
    echo 'export DEBUG_LOG=true' >> "$CLAUDE_ENV_FILE"
    echo 'export PATH="$PATH:./node_modules/.bin"' >> "$CLAUDE_ENV_FILE"
fi
exit 0
```

This is particularly useful for language version managers, virtual environments, and any CI-specific configuration that normally requires sourcing a setup script. The environment persists for the entire session without Claude Code needing to know about the setup process.

Backpressure via Pre-Commit Hooks

Pre-commit hooks -- the kind that run before a git commit is finalized -- create an especially tight quality gate for autonomous Claude Code. Set up a pre-commit hook that runs your type checker, linter, and test suite:

```
# .husky/pre-commit
pnpm typecheck && pnpm lint && pnpm test-run
```

When Claude Code (or a subagent) runs `git commit`, the hook fires immediately. If the type checker finds errors, the linter flags violations, or the tests fail, the commit is rejected. Claude sees the error output and self-corrects. This is automated feedback at the moment of truth -- the commit boundary -- and it catches issues at the source rather than accumulating bugs across multiple commits.

The implication for CI is significant. You do not need to build a custom verification system. Your existing pre-commit hooks become the quality gate for autonomous operation. Every commit Claude makes passes through the same checks as every commit a human makes. You stop being the bottleneck for quality control.

Attribution Settings

Claude Code adds attribution to git commits and pull requests by default. In CI workflows, you may want to customize or disable this behavior. The `attribution` settings control it:

```
{
  "attribution": {
    "commit": "Generated with AI\n\nCo-Authored-By: AI <ai@example.com>",
    "pr": ""
```

```
    }  
}
```

Setting `commit` customizes the trailer added to commit messages. Setting `pr` controls the text appended to pull request descriptions. An empty string disables the attribution for that context. In automated workflows where every commit is AI-generated, you might want concise attribution rather than the default format, or you might want to disable PR attribution entirely to keep descriptions clean.

Devcontainer Reference Setup

The reference devcontainer ships with three components: a container configuration file that controls settings, extensions, and volume mounts; a container image definition that installs the runtime and tooling; and a firewall initialization script that establishes network security rules.

The container image definition is worth studying even if you build your own. It starts from a base image, installs the target language runtime, sets up a shell with productivity enhancements, and installs Claude Code itself. The firewall script applies a default-deny network policy with explicit allowlisting for domains Claude Code needs to function (API endpoints, package registries). The container configuration ties everything together with volume mounts for session persistence and workspace settings.

The key design decision is that network isolation happens at the OS level via firewall rules, not at the application level. Claude Code does not know it is network-restricted. It simply finds that connections to non-allowlisted domains fail. This makes the security boundary robust against anything Claude Code might try to do -- the restriction is enforced below the application layer.

For teams that need consistent, reproducible environments across developers and CI, the reference devcontainer is the starting point. Clone it, customize the language runtime and dependencies for your stack, adjust the network allowlist, and you have a secure environment for both interactive development and headless CI operation.

The CI Auto-Fix Gap

One integration pattern that does not yet exist but is widely anticipated: automatic CI failure remediation.

The vision is straightforward. A CI build fails. Claude Code reads the failure logs, identifies the problem, implements a fix, and pushes a commit -- all without human intervention. The build goes green. The developer never sees the failure.

This capability is not available today. The primitives exist -- headless mode, pipe input, push output -- but the end-to-end integration that monitors CI pipelines and triggers remediation automatically is not yet productized. Industry analysis suggests this will arrive in the second or third quarter of 2026, likely as a skill or first-party integration rather than a built-in feature.

In the meantime, the manual version works: copy the CI failure log, pipe it to Claude Code, review the fix, push it. The automation is the last mile.

Long-Running Headless Operation

Most headless usage assumes short-lived sessions: a CI job runs, completes in minutes, and terminates. But Claude Code's headless mode has been demonstrated running for dramatically longer periods.

One documented case involved Claude Code operating as an autonomous agent for 33 days. The system ran headless with occasional human check-ins, making decisions, executing actions, and maintaining coherent behavior across an extended timeframe. The human operator's role shifted from active direction to periodic supervision -- reviewing outcomes, adjusting parameters, and intervening only when the system drifted from its objectives.

This is not typical usage, and it surfaces challenges that short-lived sessions never encounter. Context compaction happens repeatedly over 33 days. The system must maintain consistent behavior across dozens of compaction cycles. Long-term memory must be externalized -- into files, databases, or structured storage -- because the context window is not a durable store.

But the existence proof matters. Claude Code is not architecturally limited to short tasks. The `-p` flag, combined with structured persistence and periodic human oversight, supports extended autonomous operation. The constraint is not the tool. It is your ability to define clear objectives, provide adequate context, and build the scaffolding for long-term state management.

Putting Headless to Work

The progression from interactive to headless follows a predictable path:

1. **Start interactive.** Use Claude Code normally. Observe what it does well with your codebase. Identify repetitive tasks.
2. **Script the repetitive tasks.** Wrap them in headless calls with `-p`. Test the output formats. Get comfortable with the CLI interface.
3. **Add to build scripts.** Embed headless calls in your project's task runner or build configuration. AI-powered linting, review, and documentation generation become part of your build.
4. **Move to CI.** Set up devcontainers. Configure permissions. Wire up triggers. Now Claude Code runs on events, not on your command.
5. **Close the loop.** Connect CI results back to Claude Code. PR comments trigger fixes. Issue labels trigger implementations. The system becomes self-reinforcing.

Each step is independently valuable. You do not need to reach step five to benefit. But each step makes the next one obvious.

Key Takeaways

- The `-p` flag transforms Claude Code from an interactive tool into a composable Unix utility that reads stdin, writes stdout, and integrates with any pipeline `-- cat error.log | claude -p 'explain' > output.txt` is a complete diagnostic workflow.
- Four system prompt flags (`--system-prompt`, `--system-prompt-file`, `--append-system-prompt`, `--append-system-prompt-file`) provide a spectrum from inline overrides to version-controlled additive prompts -- prefer file-based variants for CI.
- The fan-out pattern with `--allowedTools` enables safe batch processing: loop through files, scope permissions per invocation, and aggregate JSON results.
- Async hooks run tests in the background after every write; `TaskCompleted` hooks block task completion until tests pass; agent-based `Stop` hooks spawn subagents to verify quality before Claude finishes.
- Pre-commit hooks (type check, lint, test) create automated backpressure that makes Claude self-correct at the commit boundary -- you stop being the quality bottleneck.
- Devcontainers with network isolation and `--dangerously-skip-permissions` are the standard pattern for safe unattended execution in CI.
- The `--from-pr` flag resumes sessions linked to a specific pull request, preserving full conversation context across review cycles.
- CI auto-fix -- automatic remediation of build failures -- is the most anticipated missing capability, expected in mid-2026.
- Start with interactive usage, identify repetitive patterns, and progressively move them to headless execution.

Chapter 7: IDE Integration Done Right

What You'll Learn

Claude Code runs in your terminal. It also runs in a code editor extension, a desktop app, a web interface, and a mobile app. Same engine everywhere. Same CLAUXE.md files. Same settings. Same MCP servers. The difference is the surface -- the interface you interact with and the workflow patterns it enables.

Most developers pick one surface and stay there. That is fine, but it leaves workflow optimizations on the table. This chapter maps out the surfaces, their strengths, their gaps, and the patterns that emerge when you understand which surface fits which task. You will learn how the desktop application provides full visual diff review with parallel Git-isolated sessions, how cloud sessions run on managed VMs you can monitor from your phone, and how the teleport feature moves work between surfaces without losing context.

You will also confront the autocomplete gap -- the one capability where Claude Code does not compete with other tools -- and learn why the dominant usage pattern across teams is terminal-first with occasional IDE assists, not the other way around.

Same Engine, Every Surface

Claude Code's architecture separates the engine from the interface. The engine handles the agentic loop, context management, tool execution, permissions, and all the machinery described in previous chapters. The interface -- terminal, editor extension, desktop app -- is a presentation layer.

This matters because it means switching surfaces does not mean switching tools. Your CLAUDE.md files load regardless of surface. Your permission settings apply everywhere. Your MCP servers connect to every surface that supports them. A workflow you develop in the terminal works in the editor extension with only minor adaptations.

The surfaces currently available: the terminal (the original and most capable), extensions for a popular code editor and a major IDE family, a desktop application, a web interface, and a mobile app. Each surface optimizes for different interaction patterns, but none of them changes what Claude Code fundamentally does.

The Terminal Surface

The terminal is where Claude Code started and where it remains most powerful. Every feature described in this book works in the terminal. No exceptions.

The terminal surface gives you the full CLI with all flags, direct access to piping and Unix composability (Chapter 6), modal editing for prompts, background bash command execution, and the ability to run multiple instances in separate terminal panes or windows. It is the only surface where you have complete control over session management, output formats, and automation integration.

Terminal-first workflows dominate among experienced users for a practical reason: the terminal does not get in the way. There is no GUI mediating your interaction. You type a prompt, Claude runs, you see results. The feedback loop is tight. When you need to inspect a file that Claude modified, you open it in your editor. When you need to run a test, you run it in another terminal pane. The terminal is not trying to be an IDE. It is a command surface that works alongside whatever IDE you prefer.

This terminal-first approach has an underappreciated benefit: it is IDE-agnostic. Teams with mixed editor preferences -- some using one code editor, others using another IDE family, others on a minimal text editor -- can standardize on Claude Code workflows without standardizing on an editor. The CLAUDE.md file, the agent definitions, the MCP configurations, the hook system -- all of it works identically regardless of which editor is open in the next window.

The Code Editor Extension

The extension for the popular code editor adds a visual layer on top of the Claude Code engine. The key features are inline diff viewing, @-mention file references, plan review interfaces, and conversation history.

Inline diffs are the headline feature. When Claude proposes changes to a file, you see the diff rendered in the editor with syntax highlighting, accept/reject controls, and the ability to partially accept changes. This is materially faster than reviewing diffs in the terminal, especially for large changes across multiple files.

@-mentions let you reference files, symbols, and selections directly in your prompt. Click on a function name, type @, and the file path is inserted into the context. This is a convenience feature -- you can achieve the same thing in the terminal by typing the path -- but it reduces friction enough to change behavior. Developers who use @-mentions reference specific context more often, which leads to better-scoped prompts.

Plan review surfaces Claude's intended approach in a structured visual format before execution. In the terminal, plans are text blocks in the conversation. In the editor, they are interactive checklists you can approve, modify, or reject item by item. This is particularly useful for large refactors where the plan has many steps and you want granular control over which steps proceed.

The extension is strongest when you are working on a single codebase and want tight integration between Claude's output and your editing workflow. It is weakest when you need features that exist only in the terminal -- specific CLI flags, piped workflows, or parallel

instance management.

The Desktop Application

The desktop application is a standalone surface that runs Claude Code outside both your IDE and your terminal. It provides visual diff review, parallel sessions with automatic Git worktree isolation, multiple permission modes, file attachments, connectors to external services, and support for both local and remote execution environments.

The application organizes work around three tabs. The Code tab is the primary workspace -- it mirrors the terminal's agentic capabilities through a graphical interface. You select a project folder, choose a permission mode, pick a model, and start typing prompts. The Cowork tab enables autonomous background work -- Claude operates continuously on its own while you monitor progress and steer when needed. The Chat tab provides a standard conversation interface for questions that do not require filesystem access.

Permission Modes

The desktop application surfaces permission control through a mode selector that affects how much autonomy Claude has during a session:

- **Ask mode** requires your approval before every file edit and command. You see a diff view and accept or reject each change. This is the default and the right choice until you trust Claude's judgment in your codebase.
- **Code mode** auto-accepts file edits but still asks before running terminal commands. Faster iteration when you trust file changes.
- **Plan mode** restricts Claude to read-only operations and plan creation. No file modifications, no commands. Pure analysis and planning.
- **Act mode** is the equivalent of `--dangerously-skip-permissions` in the terminal. Claude runs without any prompts. Only use this in sandboxed environments.

You can switch modes mid-session. Start in Plan mode to map out an approach, switch to Code mode to execute, drop to Ask mode for the final verification steps.

Visual Diff Review

When Claude modifies files, a diff stats indicator appears showing lines added and removed. Clicking it opens a full diff viewer with a file list on the left and changes on the right. You can comment on specific lines by clicking them -- type your feedback, press Enter, and after commenting on multiple locations, submit all comments at once. Claude reads the comments and makes the requested changes, which appear as a new diff. This is review-by-conversation, not review-by-reading. You point at problems, and Claude fixes them.

Parallel Sessions with Git Isolation

Each session in the desktop application gets its own isolated copy of your project through Git worktrees. Click "New session" in the sidebar to start a second (or third, or tenth) parallel task. Changes in one session do not affect other sessions until you commit them. Worktrees are stored in your project's `.claude/worktrees/` directory by default. You can configure a custom location and a branch prefix in settings to keep Claude-created branches organized.

This turns the desktop application into a parallel processing interface. One session is refactoring the authentication module. Another is writing tests for the payment flow. A third is investigating a bug. Each session operates in its own Git branch, completely isolated from the others. When each session's work is done, you review and merge the branches independently.

SSH and Remote Sessions

The desktop application connects to remote machines over SSH. Instead of selecting "Local" as your environment, select an SSH connection. Claude Code runs on the remote machine, executing commands and modifying files there. This is useful for codebases that require specific hardware, operating systems, or development environments that your local machine does not provide.

For long-running tasks, select a Remote environment instead. Remote sessions run on Anthropic-managed cloud infrastructure and continue even if you close the application or shut down your computer. You can monitor remote sessions from the desktop application, the web interface, or a mobile app. More on remote sessions in the next section.

Connectors

The desktop application supports connectors that integrate external services -- version control platforms, team messaging tools, and project management systems. These connectors are configured through the application's settings and make Claude Code aware of your project's external context: PRs, issues, messages, tasks. The connector ecosystem is still growing, but the pattern is clear: Claude Code becomes the hub that connects your code to your team's communication and coordination tools.

Claude Code on the Web

Claude Code runs in the browser as a fully cloud-hosted environment. No local setup required. You connect your version control account, select a repository, and start working. The repository is cloned to an Anthropic-managed virtual machine, an environment is configured, and Claude executes with full access to the codebase.

The web interface provides the same diff view as the desktop application -- you see exactly what Claude changed, comment on specific lines, and iterate until the changes are ready. When satisfied, you create a pull request directly from the interface. Changes are pushed to a branch on the remote, ready for review.

Network Access Configuration

Cloud sessions support three network access levels: limited (default allowed domains only -- package registries, API endpoints, and similar infrastructure), full (unrestricted internet access), and none (completely offline). The network configuration is controlled through a security proxy that mediates all outbound connections. A separate proxy handles version control platform access, ensuring that repository operations work regardless of your network settings.

Session Sharing

Cloud sessions can be shared at three visibility levels: Private (only you), Team (members of your organization), and Public (anyone with the link). Shared sessions let others observe Claude's work in progress, review the conversation history, and see the resulting changes. This is useful for code reviews, pair programming across time zones, and demonstrating workflows to team members.

When to Use Cloud Sessions

Cloud sessions shine for three scenarios. First, long-running tasks that would tie up your local machine -- kick off a large refactor, close the browser, and check back in an hour. Second, repositories you do not have locally -- review a colleague's project without cloning anything. Third, parallel execution -- start multiple cloud sessions on different tasks and monitor them from a single browser tab.

The limitation is that cloud sessions currently work only with repositories hosted on the dominant code platform. Other hosting platforms are not yet supported for cloud execution.

Claude Code on Mobile

The mobile application provides a monitoring and dispatch interface for cloud sessions. Kick off a task from your phone, monitor its progress, steer Claude with follow-up prompts, and review the results when they are ready. You are not going to write complex prompts on a phone keyboard, but you can check on a long-running task during your commute, answer Claude's clarifying questions, or approve a PR from the couch.

The mobile app connects to the same session infrastructure as the web interface. A session started from the desktop or terminal can be monitored from the mobile app, and vice versa.

The Chrome Extension

A browser extension connects Claude Code to live web applications. It bridges the gap between the codebase and the running application -- Claude can read the DOM, observe network requests, and inspect the rendered output of the code it is working on. This is most valuable for frontend debugging, where the symptom (a visual bug in the browser) and the cause (a CSS or JavaScript issue in the code) live in different contexts. The extension brings them together.

The IDE Family Plugin

The plugin for the major IDE family provides interactive diff viewing and context sharing. The implementation differs from the code editor extension in ways that reflect the IDE's architecture -- diffs integrate with the IDE's native diff viewer, context references use the IDE's symbol index, and the interaction model follows the IDE's conventions.

The practical differences between the two editor integrations are smaller than the architectural differences suggest. Both provide inline diff review. Both support context referencing. Both connect to the same Claude Code engine with the same capabilities. The choice between them is almost always driven by which IDE you already use, not by feature differences in the Claude Code integration.

Cross-Surface Workflows

Surfaces are not isolated. Sessions can move between them. This is one of Claude Code's most distinctive architectural decisions -- no other tool in this category supports fluid session migration across interfaces.

The /teleport Command

The `/teleport` command (also `/tp`) pulls a cloud session into your terminal. Start a long-running task on the web or from your phone, let it run while you are away, then pull it back to your local machine when you are ready to review and continue.

```
/teleport
```

This opens an interactive picker showing your cloud sessions. Select one, and Claude Code verifies you are in the correct repository, fetches and checks out the branch from the remote session, and loads the full conversation history into your terminal. If you have uncommitted local changes, it prompts you to stash them first.

The `--teleport` flag works the same way from the command line:

```
claude --teleport  
claude --teleport <session-id>
```

The & Prefix: Sending Work to the Cloud

The reverse direction -- terminal to cloud -- uses the `&` prefix. Start a message with `&` inside an interactive session:

```
& Fix the authentication bug in src/auth/login.ts
```

This creates a new cloud session with your current conversation context. The task runs on Anthropic's infrastructure while you continue working locally. Use `/tasks` to monitor all background sessions -- it shows status, and pressing `t` teleports you into any session.

The pattern that emerges is planning locally and executing remotely. Start Claude in plan mode to collaborate on the approach, then send the plan to the cloud for autonomous execution:

```
& Execute the migration plan we discussed
```

The /desktop Command

The `/desktop` command hands off a terminal session to the desktop application. The full conversation context moves to the GUI, where you get visual diff review, inline comments, and the ability to partially accept changes. This is the right move when the task shifts from exploration to review.

Workflow Patterns

The cross-surface workflow that emerges is: explore and plan in the terminal, implement in the terminal or editor, review diffs in the desktop application or editor extension. The terminal is the home base. Visual surfaces are specialized tools you invoke when review adds value.

This is not how most IDE integrations work. Most tools embed in the IDE and never leave. Claude Code inverts that relationship. The terminal is the primary surface, and the IDE is a supplementary view you pull up when needed. Understanding this inversion prevents the frustration of trying to force terminal-grade functionality through the editor extension.

PR Review Status in the Footer

When working on a branch with an open pull request, Claude Code displays a clickable PR link in the terminal footer (for example, "PR #446"). The link has a colored underline indicating the review state:

- **Green**: approved
- **Yellow**: pending review
- **Red**: changes requested
- **Gray**: draft
- **Purple**: merged

The status updates automatically every 60 seconds. Cmd+click (Mac) or Ctrl+click (Windows/Linux) opens the PR in your browser. This small feature eliminates a constant context switch -- checking PR status without leaving the terminal. It requires the version control platform's CLI to be installed and authenticated.

Eliminating Context Switching

One pattern that practitioners report consistently: Claude Code in the terminal eliminates the context switch between a cloud AI chat interface and the codebase. Without Claude Code, the workflow is: copy a code snippet, paste it into the chat interface, explain the problem, read the response, copy the suggested fix, paste it back into the editor, test it, and repeat. Each round trip crosses application boundaries and loses context.

With Claude Code, the entire conversation happens inside the project. Claude reads files directly. It writes changes directly. It runs tests directly. The context is the codebase, not a description of the codebase. The difference is not marginal -- it eliminates the mental overhead of translating between two environments and the mechanical overhead of copy-paste.

LSP Plugins

The plugin system includes Language Server Protocol servers that provide code intelligence features: diagnostics, go-to-definition, code navigation, and symbol resolution. These are not the same as Claude Code's agentic capabilities -- they are traditional LSP features implemented as plugins.

LSP plugins are newer and less battle-tested than the core Claude Code features. Their value is in providing code intelligence for languages or frameworks that your IDE does not natively support well. If your IDE already has excellent language support for your stack, LSP plugins may add little. If you are working in a less common language or a custom framework, they can fill gaps.

The plugins auto-start when configured and appear as standard IDE features. They complement rather than replace your IDE's native language services.

The Autocomplete Gap

Claude Code does not have streaming inline autocomplete. The kind of experience where you type a few characters and a gray ghost of the suggested completion appears in your editor, updating with each keystroke -- that is not what Claude Code does.

This is not an oversight. It is an architectural consequence. Claude Code's strength is in agentic, multi-step operations that span files and tools. Inline autocomplete is a fundamentally different interaction pattern: low-latency, single-line, character-by-character prediction. Building both well in the same tool requires different model configurations, different infrastructure, and different UX.

The gap matters because inline autocomplete is genuinely useful. For writing boilerplate, filling in function signatures, completing variable names, and the countless small editing tasks that do not warrant a full prompt-and-execute cycle, a good autocomplete tool saves real time.

The recommended approach is a hybrid setup. Use Claude Code for the agentic, multi-step, multi-file work it excels at. Use a separate tool that specializes in inline autocomplete for the character-by-character editing work. The two tools do not conflict. They operate at different scales of interaction -- one handles the forest, the other handles individual leaves.

This hybrid strategy is not a workaround. It is the mature approach. No single tool dominates every interaction scale, and pretending otherwise leads to using a powerful tool badly for tasks where a simpler tool works better.

When to Use the IDE, When to Drop to Terminal

The decision is simpler than it seems.

Use the IDE extension when:

- You need to review multi-file diffs visually
- You are referencing specific code locations with @-mentions
- You want to partially accept or reject changes within a file
- You are doing plan review and want the interactive checklist
- You are working with a teammate who is more comfortable in the IDE

Drop to the terminal when:

- You need CLI flags not exposed in the IDE extension
- You are running multiple parallel sessions
- You are piping input or output to other tools
- You need automation or scripting integration
- You want the fastest possible feedback loop
- You are doing exploratory work where conversation flow matters more than diff review

Most sessions start in the terminal. Some move to the IDE for review. Few start in the IDE and stay there. This is not a judgment about IDE quality -- it reflects the fact that agentic coding is primarily a conversation, and conversations flow most naturally in the terminal.

The Design Tool Pattern

A pattern worth noting: teams report that a visual design tool and Claude Code are open simultaneously roughly 80% of the time. The design tool provides the visual reference -- mockups, layouts, component designs. Claude Code implements what the design shows.

This is not an IDE integration in the traditional sense. It is a workflow where two tools complement each other without being connected. The human is the integration layer, looking at the design and describing to Claude Code what needs to be built.

The pattern works because Claude Code accepts images as input. You can paste a screenshot of a design mockup into the conversation, and Claude can interpret the layout, colors, spacing, and component structure. This is faster than describing a design in words, though it is not a substitute for explicit implementation requirements.

The design-to-code workflow follows the same four-phase pattern described in Chapter 1: explore the existing codebase components, plan the implementation approach, implement the design, and commit. The design tool adds a visual reference to the exploration phase that makes the subsequent phases more grounded.

Skills as a Cross-Agent Standard

Skills are not exclusive to Claude Code. They are a shared standard supported by multiple AI coding agents. A skill you write for Claude Code works with other agents that support the standard -- and the list is growing. This means your investment in skill development is not locked to a single tool.

The installation workflow uses a package runner command that works across agents:

```
npx skills add <repository-url> --skill <skill-name>
```

Running this command clones the repository, discovers available skills (a typical repository might contain a dozen or more), detects compatible agents, asks which scope to install in (project or user), and creates the skill directory structure. Skills installed this way can be symlinked or copied, and they work immediately without restarting Claude Code.

The cross-agent compatibility has a practical consequence: teams that use multiple AI coding tools can standardize their workflows through skills rather than building tool-specific configurations. A code review skill, a deployment workflow, or a documentation generator works regardless of which agent runs it. This reduces the cost of experimenting with different tools and prevents workflow lock-in.

The Hybrid Toolchain

No single AI coding tool dominates every interaction scale. The mature approach is a hybrid setup with each tool handling what it does best.

Use Claude Code for the agentic, multi-step, multi-file work it excels at: repository-wide refactors, complex implementations, codebase exploration, architectural planning, and automated workflows. Layer a dedicated inline completion tool for the character-by-character editing work -- function signatures, variable names, boilerplate, and the countless small editing tasks that do not warrant a full prompt-and-execute cycle.

This is not a workaround for a missing feature. It is an acknowledgment that agentic coding and inline completion are fundamentally different interaction patterns requiring different model configurations, different infrastructure, and different UX. Building both excellently in the same tool requires different optimization targets that inevitably trade off against each other.

The prediction for the near future is that Claude Code's inline IDE experience will become more aggressive -- per-file edit modes, region-scoped rewrites, tighter editor integration. But even as the surfaces converge, the underlying architecture will continue to optimize for different scales of interaction. The hybrid approach is not a temporary compromise. It is the stable end state.

Terminal First as a Design Philosophy

The terminal-first pattern is not just a preference. It is a design philosophy that shapes how Claude Code evolves.

Because the terminal is the universal surface -- every operating system has one, every IDE has one embedded, every remote server is accessible through one -- building features for the terminal first ensures they work everywhere. An extension-first approach would leave terminal users behind. A terminal-first approach gives everyone the full feature set, with extensions providing visual enhancements on top.

This philosophy has practical implications for how you invest your learning time. Learning the terminal workflow gives you skills that transfer to every surface. Learning the extension workflow gives you skills specific to that extension. If you have limited time, invest in the terminal first and add extension workflows as refinements.

The terminal-first philosophy also explains why Claude Code does not try to replace your IDE. It sits alongside your IDE as a complementary tool, handling the agentic work that IDEs do not do well while leaving the editing, debugging, and navigation work to tools purpose-built for those tasks. The integration is collaborative, not competitive.

Key Takeaways

- Claude Code's engine is surface-independent -- CLAUDE.md, settings, MCP servers, and permissions work identically across terminal, editor extensions, desktop, web, and mobile.
- The desktop application provides visual diff review with inline comments, parallel sessions with Git worktree isolation, multiple permission modes, SSH connections, and connectors to external services.
- Cloud sessions run on managed VMs, persist after you close the browser, and can be monitored from any device -- use the & prefix to dispatch work to the cloud and `/teleport` to pull it back.
- The terminal footer shows live PR review status with colored underlines (green=approved, yellow=pending, red=changes requested) -- updated every 60 seconds.
- Skills are a cross-agent standard: `npx skills add` installs skills that work across multiple AI coding tools, preventing workflow lock-in.
- The inline autocomplete gap is architectural, not accidental; a hybrid setup with Claude Code for agentic work and a dedicated autocomplete tool for editing is the stable end state.
- Terminal-first workflows dominate among experienced users because the terminal is faster, more flexible, and IDE-agnostic.
- Invest learning time in terminal workflows first -- they transfer to every surface and every IDE.

Chapter 8: Prompt Craft for Agentic Tools

What You'll Learn

Prompting an agentic coding tool is not the same as prompting a chatbot. A chatbot generates text. An agentic tool reads files, writes code, runs tests, modifies your filesystem, and chains dozens of operations together. The prompt is not a question -- it is a work order. The difference between a prompt that produces a clean implementation and one that produces three hours of cleanup is not eloquence. It is engineering.

The single highest-leverage technique is not clever wording or elaborate system prompts. It is giving Claude verification criteria -- tests, expected outputs, screenshots to match -- so that the agent has a feedback loop instead of a guess. Everything else here matters, but nothing matters as much as that.

Beyond verification, the specific patterns that practitioners have discovered make agentic tools dramatically more effective: the delegation mindset, rich content input, interrupt-and-steer workflows, backpressure through tooling, spec-first prompting, and the counterintuitive discovery that your linting configuration is part of your prompt. These are not theoretical best practices. They are the patterns that separate fifteen-minute tasks from three-hour debugging sessions.

Verification Criteria: The One Thing That Matters Most

If you take nothing else from these pages, take this: tell Claude how to verify its own work.

Without verification criteria, Claude's workflow is generate-and-hope. It writes code, decides it looks correct, and moves on. If there is a bug, you find it. If the output format is wrong, you notice it. If an edge case breaks, you discover it in production.

With verification criteria, Claude's workflow is generate-test-fix. It writes code, runs the verification, observes the result, and iterates until the verification passes. The feedback loop is automatic. You specified the success criteria; Claude chases them.

Verification criteria include:

- **Test commands:** "Run `npm test` after every change and ensure all tests pass."
- **Expected output:** "The function should return `[1, 4, 9, 16]` for input `[1, 2, 3, 4]`."
- **Visual checks:** Paste a screenshot of the target UI. "The result should look like this."
- **Behavioral descriptions:** "The API should return a 400 status code when called without an auth header."
- **Existing test suites:** "Run the full test suite and ensure no regressions."

The key insight is that verification criteria convert a one-shot generation into an iterative loop. Claude does not need to get it right the first time. It needs to get it right before it stops. Verification criteria define "before it stops."

The Absence of Criteria

When you prompt "add user authentication to the API" without verification criteria, Claude generates an implementation that looks plausible. It may work. It may not. You are the verification step, and you might not catch subtle issues on review.

When you prompt "add user authentication to the API. Run `pytest tests/auth/` after implementation. All tests should pass. The endpoint `/api/protected` should return 401 without a token and 200 with a valid token" -- now Claude has a target. It will run the tests. It will hit the endpoint. It will iterate until the criteria are met or it gets stuck and tells you why.

The cost of adding verification criteria to a prompt is thirty seconds of typing. The cost of not adding them is measured in debugging time.

Specificity vs. Vagueness

There is a time for vague prompts and a time for specific ones. Knowing which is which separates effective delegation from wasted cycles.

Vague prompts are appropriate for exploration. "What does the authentication flow look like in this codebase?" -- "How is error handling structured?" -- "What would it take to add internationalization?" These prompts leverage Claude's ability to read broadly and synthesize. You do not know the answer; you want Claude to discover it.

Specific prompts are appropriate for implementation. "In `src/auth/middleware.ts`, add rate limiting to the `validateToken` function. Use a sliding window of 100 requests per minute per IP. Store counts in the existing in-memory store connection from `src/config/cache.ts`. Add tests in `tests/auth/rate-limit.test.ts`."

The failure mode is using vague prompts for implementation work. "Add rate limiting" without specifying where, how, or to what gives Claude maximum freedom to make choices you will disagree with. It picks the wrong file. It invents a new cache connection instead of using the existing one. It stores state in local memory instead of the shared cache. Each assumption is reasonable in isolation and wrong for your codebase.

Files, Constraints, Examples

Three categories of specificity that consistently improve outcomes:

Reference specific files. Use `@` mentions or explicit paths. "Look at `@src/utils/validation.ts` for the pattern we use." Claude does not need to search for the right file. You have already pointed it there.

State constraints. What cannot change. What trade-offs are acceptable. What flexibility exists. "Do not modify the database schema. Performance matters more than code readability here. You can add new files but do not restructure existing directories."

Provide examples. If you want output in a specific format, show the format. If you want code that matches a pattern, show the pattern. Claude extrapolates from examples better than it follows abstract descriptions.

The Delegation Mindset

Prompting an agentic tool is not writing instructions for a script. It is delegating to a capable colleague. The distinction changes what you include in the prompt.

For a script, you specify every step: "open this file, find this function, add this line, save the file." For a colleague, you specify the goal and the context: "we need rate limiting on the auth middleware because we are getting hammered by a bot. The cache connection already exists. Use the same sliding window approach we used on the payment endpoint."

The colleague figures out the steps. The script executes yours. Claude Code is the colleague.

The delegation mindset changes three things about your prompts:

- 1. Goal over steps.** State what needs to be true when Claude is done, not how to get there.
- 2. Context over instructions.** Explain why this change matters, what exists already, what the constraints are.
- 3. Trust over control.** Let Claude choose the implementation approach. Intervene when it chooses wrong, not preemptively.

This is uncomfortable for developers. We are trained to specify precisely. But over-specifying an agentic tool constrains it to your approach, which may not be the best approach. State the goal. Provide the context. Let the agent work.

Rich Content Input

Claude Code accepts more than text prompts. Using the full range of input modalities consistently improves outcomes.

File references with @. Type `@` followed by a path to include file contents directly in context. This is cheaper and more reliable than asking Claude to read the file -- the contents are in the prompt, not fetched by a tool call.

Images. Paste screenshots of UI targets, error messages, design mockups, or whiteboard sketches. Claude interprets images and can implement toward a visual target. "Make the dashboard look like this" with a screenshot is more precise than three paragraphs describing the layout.

URLs. Claude can fetch and process web content. Point it at API documentation, style guides, or reference implementations. The content is retrieved and interpreted, saving you from copy-pasting.

Piped input. Claude Code is a Unix citizen. Pipe data into it: `cat error.log | claude "what's causing these failures?" -- git diff HEAD~5 | claude "summarize these changes"`. This connects Claude to your existing command-line workflow.

Each input modality adds context. Context costs tokens. But the right context -- a screenshot instead of a paragraph, a file reference instead of a description -- is cheaper and more effective than the alternative.

Interrupt and Steer

Claude Code is not a batch process. You can interrupt it mid-task.

If Claude is heading in the wrong direction -- implementing a solution you do not want, using a library you do not use, modifying the wrong file -- type your correction and press Enter. Claude stops what it is doing, reads your correction, and adjusts course.

This is the real-time steering loop that makes agentic tools qualitatively different from one-shot generation. You are not waiting for the full output to review and reject. You are watching the work happen and redirecting in real time.

Effective interruptions are specific: "stop -- use the existing database connection in config/db.ts instead of creating a new one." Vague interruptions ("that's wrong") force Claude to guess what you object to.

The interrupt-and-steer pattern is most valuable during the first few minutes of a complex task, when Claude's initial approach reveals its assumptions. Correct the assumptions early and the rest of the implementation follows correctly.

Role Assignment Prompts

Framing Claude's role changes its behavior. This is not mystical prompt engineering -- it is setting coordination expectations.

"You are the orchestrator. Your subagents are your development team. Assign each agent a clear, bounded task. Review their output before integrating." This framing causes Claude to delegate to subagents rather than implementing everything in the main context. It plans before acting. It reviews results.

Without the framing, Claude tends toward doing everything itself in a single context. For small tasks, that is fine. For multi-file refactors or cross-service changes, it leads to context exhaustion and dropped details. The orchestrator framing triggers a different work pattern.

Other effective role framings:

- **Investigator:** "Your job is to understand, not to change. Report what you find." Keeps Claude in read-only mode for exploration tasks.
- **Reviewer:** "Review this pull request as a senior engineer. Focus on correctness, security, and maintainability." Produces structured review output rather than code changes.
- **Specialist:** "You are a database migration specialist. Your only concern is schema changes and data integrity." Narrows focus and reduces tangential changes.

Spec-Driven Development

Before asking Claude to write code, ask it to write a document. This is not just a prompting tip. It is a named workflow pattern -- Spec-Driven Development -- that contrasts sharply with the default AI coding pattern of prompt, code, debug, repeat.

The traditional pattern goes like this: you write a prompt, Claude generates code, you find bugs, you debug, you prompt again, you generate more code, you find more bugs. Context fills up with failed attempts. There is no persistent source of truth. No clear stopping point. Each iteration starts from scratch in the same cluttered context.

Spec-Driven Development follows a different flow: **Research, Spec, Refine, Tasks, Done.**

1. **Research.** Spawn multiple subagents to investigate the relevant parts of the codebase in parallel. Each subagent explores a different dimension -- the existing architecture, the API surface, the test patterns, the dependencies. The research results flow back to the main context.
2. **Spec.** Ask Claude to write a comprehensive specification based on the research. "Your goal is to write a report. Produce a technical spec for the migration covering architecture, data model, error handling, and rollback strategy. Do not write any code yet." The spec becomes a file in the repository -- a persistent artifact that survives context compaction and session restarts.
3. **Refine.** Before implementation, let Claude interview you about the spec using the AskUserQuestion tool: "Use the `ask_user_question` tool -- do you have any questions regarding the spec before we implement it?" Claude asks clarifying questions about ambiguities, design decisions, and edge cases. You answer. The spec is updated. This catches misunderstandings that would have become bugs.
4. **Tasks.** Break the spec into atomic tasks, each implemented by a subagent with a fresh context. "Use the task tool. Each task should only be done by a subagent. After each task, do a commit." Each subagent reads the spec, implements one piece, commits, and returns. The main context stays lean -- it is orchestrating, not implementing.
5. **Done.** Clear completion criteria, defined in the spec, determine when the work is finished.

Five Prompt Patterns for Spec-Driven Development

The workflow relies on five specific prompt patterns that trigger the right behavior:

1. **"Spin up multiple subagents for your research task."** Triggers parallel research. Five subagents studying five aspects of the codebase simultaneously.
2. **"Your goal is to write a report/document."** Forces Claude to produce a written artifact before any code. This becomes the source of truth.
3. **"Use the `ask_user_question` tool before we implement."** Surfaces ambiguities and design decisions before they become bugs.
4. **"Use the task tool, and each task should only be done by a subagent. After each task, do a commit."** Forces atomic, independently verifiable units of work with commit boundaries.
5. **"You are the main agent and your subagents are your devs."** Triggers the orchestrator role framing that keeps the main context lean and delegation-focused.

A Walkthrough: Storage Layer Migration

To make this concrete: one practitioner used Spec-Driven Development to migrate a sync engine's storage layer from one browser database technology to another -- a task that would have taken two to three days manually.

Phase 1: Research. Five subagents launched in parallel to study the target framework. Each investigated a different aspect: the data model, the sync protocol, the conflict resolution strategy, the indexing approach, and the cross-tab coordination mechanism. The parallel research completed in minutes, producing findings that would have taken hours of documentation reading.

Phase 2: Spec creation. Based on the research, Claude produced a comprehensive specification covering the migration strategy, the adapter interface, the data conversion logic, and the test plan. The spec went into a file in the repository.

Phase 3: Refinement. Claude asked clarifying questions using the interactive questioning tool: "How should the migration handle conflicts between the old and new storage formats? Should we support rollback to the old format? What is the expected behavior if a user has multiple tabs open during migration?" Each answer tightened the spec.

Phase 4: Task delegation. The spec broke down into fourteen tasks, each assigned to a subagent. Each subagent read the spec, implemented its task, ran the relevant tests, and committed. The git log showed fourteen atomic commits, each implementing one clearly bounded piece of the migration.

The entire workflow took approximately forty-five minutes. The resulting code was better than a manual implementation would have been, because the research phase uncovered patterns from the target framework that the developer would not have found on their own.

The AskUserQuestion Tool

The refinement phase deserves special attention. When you tell Claude "use the ask_user_question tool before we implement," Claude enters an interview mode. It reads the spec and asks multiple-choice or open-ended questions about everything it is uncertain about. This is not the same as "ask me questions exhaustively" in freeform conversation. The AskUserQuestion tool provides structured, focused questions that are easier to answer and produce more actionable clarifications.

The interview pattern works well beyond spec refinement. Use it at the start of any complex task: "Before building anything, interview me about my requirements using the ask_user_question tool. Do not proceed until you have asked every question." The structured format surfaces assumptions that freeform conversation misses.

The Spec as Persistent Truth

The spec file is not just a planning artifact. It is the recovery point for the entire workflow. If a subagent fails, the spec is still there -- start a new subagent on the same task. If context compacts, the spec survives as a file on disk. If you need to resume the next day, the spec tells the new session everything it needs to know.

One engineering team takes this further: they store specifications as markdown files in the codebase, written by Claude, reviewed by humans, and then executed by Claude. The specs go through code review like any other artifact. When the implementation is done, the spec remains as documentation of the design decisions. This is specification-as-code -- the spec is not a throwaway planning document but a versioned, reviewed, executed artifact.

Backpressure as Prompt Engineering

This is one of the most valuable patterns practitioners have discovered, and it reframes how you think about project tooling entirely.

Your linting configuration is part of your prompt.

Not metaphorically. Literally. When Claude Code runs in a project with strict linting, comprehensive type checking, and a thorough test suite, it receives automated feedback on every change. The linter flags a style violation. The type checker catches a type error. A test fails. Claude reads the error, fixes the issue, and tries again.

This is backpressure. The tooling pushes back on incorrect output, and the agent self-corrects. The tighter the tooling, the better the output. A project with strict linting rules, strict type checking configuration, and high test coverage produces better Claude Code output than an identical project without those constraints -- even though Claude was never explicitly told about the rules.

The implication is architectural: investing in strict linting, comprehensive types, and thorough tests pays a double dividend. It catches human errors and it catches AI errors. The AI errors it catches are the ones you would have had to find manually.

You have a limited budget of feedback. Every time you manually tell Claude "fix the indentation" or "use semicolons" or "that variable should be const," you are spending your feedback budget on mechanical issues that tooling should catch automatically. Spend your feedback budget on architecture, design decisions, and domain logic -- the things linters cannot check. Let the type checker and linter handle everything they can handle.

Pre-commit hooks that run the test suite and linter create a particularly tight backpressure loop. Claude commits, the hooks run, failures are reported, Claude fixes and tries again. The human does not need to review every change. The tooling reviews it.

This reframes tooling investment. A linting configuration is not just a code quality tool. It is a component of the prompt that shapes every piece of code Claude generates in that project. Invest accordingly.

TDD as Backpressure

The test-driven development workflow is the ultimate expression of backpressure. One engineering team described their old pattern as "design doc, janky code, refactor, give up on tests." With Claude Code, the pattern inverted. They ask Claude for pseudocode first, guide it through test-driven development, and periodically check in to steer when it gets stuck. The tests are written before the implementation. The implementation cannot pass until it satisfies the tests. Claude iterates until it does.

This is not TDD for philosophical reasons. It is TDD for pragmatic ones. When an agent writes both the tests and the implementation, the tests create the feedback loop that prevents the implementation from drifting. The agent writes a test, writes code to pass it, runs the test, sees the failure, adjusts, and tries again. The human checks in periodically to verify the tests themselves are testing the right things.

Plan-Before-Execute

Asking Claude to plan before implementing is cheap insurance against expensive mistakes.

"Before making any changes, read the relevant files and produce a plan. List every file you will modify, every new file you will create, and every behavioral change. Wait for my approval before proceeding."

A plan costs roughly 10,000 tokens to generate. A misdirected implementation costs 500,000 tokens and a context reset. The economics are obvious.

Plan review also lets you catch misunderstandings before they become code. If Claude's plan includes modifying a file you consider untouchable, you catch it at the plan stage. If the plan omits a file that obviously needs changing, you add it before implementation begins.

A keyboard shortcut makes this even more practical: pressing `Ctrl+G` opens Claude's plan in your default text editor. You can edit the plan directly -- add steps, remove steps, reorder priorities, add constraints -- and when you save and close the editor, Claude proceeds with your modified plan. This is faster than describing modifications in conversation and more precise than trying to steer an existing plan through follow-up prompts.

The plan-before-execute pattern composes well with subagent delegation. Plan cheaply in the main context, then hand each plan item to a subagent for parallel execution. This is the most cost-effective multi-agent strategy, as discussed in Chapter 4.

You can also convert a plan into a shorter checklist for execution. After Claude produces a detailed plan, ask it to "remix the plan into a shorter checklist I can scan and check off." The checklist becomes the working document -- concise enough to keep in view while implementing, detailed enough to capture all the steps.

Constraint Explicitness

Implicit constraints are invisible constraints. Claude cannot respect rules you have not stated.

Three categories of constraints to state explicitly:

What cannot change. "Do not modify the database schema." "Do not change the public API surface." "The file `core/engine.py` is frozen -- read it but do not edit it." These hard boundaries prevent Claude from taking shortcuts through things you consider sacrosanct.

Trade-off preferences. "Optimize for readability over performance." "Minimize the number of new dependencies." "Prefer simple solutions even if they are less elegant." These soft preferences guide Claude's choices when multiple approaches are viable.

Flexibility boundaries. "You can add new files in `src/utils/` but not in `src/core/` ." "You can use any library already in `package.json` but do not add new ones." "Refactoring the test helpers is acceptable if it simplifies the implementation." These tell Claude where it has room to maneuver.

Unstated constraints are the most common source of disappointing output. Claude produces something correct that violates an assumption you did not realize you had. State the constraints. All of them.

Output Format Specification

When you want structured output, describe the structure.

"Produce a comparison table with columns: Feature, Current Implementation, Proposed Change, Risk Level." -- "Format the analysis as a numbered list with each item containing: the finding, the severity, and the recommended fix." -- "Write the test plan as a checklist with checkboxes."

Claude follows format specifications reliably. Omitting them gets you whatever format Claude defaults to, which may or may not match what you need.

Format specification is especially important when Claude's output will be consumed by other processes or other people. A downstream tool expecting JSON will not appreciate freeform text. A colleague expecting a checklist will not appreciate an essay.

Self-Critique Prompting

Asking Claude to evaluate its own work surfaces issues it did not mention during generation.

"Rate this implementation 1-10. What would you improve? What edge cases might break? What would a senior engineer object to?"

The self-critique step produces surprisingly useful output. Claude identifies missed edge cases, suggests performance improvements, flags potential security issues, and acknowledges shortcuts it took. Not all of its self-critique is actionable, but enough of it is to make the pattern worthwhile.

A variant: "What are the three most likely things wrong with this?" This focuses the critique on probable issues rather than hypothetical ones.

Self-critique is most valuable after a complex implementation, before you commit. It takes thirty seconds. It catches issues that would take thirty minutes to find in review.

Exhaustive Questioning

When Claude fills in assumptions, the results are plausible but often wrong. The fix is to prevent the assumption-filling.

"Before you start, ask me questions exhaustively about anything you are uncertain about. Do not proceed until you have asked every question."

This prompt modifier causes Claude to surface its uncertainties before acting on them. What authentication system are you using? What database? What should happen on failure? Do you want error messages user-facing or logged? Is this a new endpoint or a modification of an existing one?

The questions reveal gaps in your own spec. You may not have thought about failure handling. You may not have considered which database to use. Claude's questions force you to make decisions you would otherwise discover as bugs.

Practitioners report that exhaustive questioning improves first-pass quality by 30-50% on complex tasks. The cost is a few minutes of answering questions. The benefit is an implementation that matches your actual requirements instead of Claude's assumptions.

Strawman Proposals

Do not start from zero when you have a rough idea.

"Here is my rough approach: use a cache-backed sliding window with per-IP tracking, expire keys after 60 seconds, return 429 when limit is hit. This might be wrong. Improve on it or suggest alternatives."

A strawman gives Claude a starting point to react to. Reacting is cognitively different from generating from scratch. The reaction incorporates your domain knowledge (cache-backed, per-IP, 429) while allowing Claude to refine the approach (maybe a token bucket is better, maybe per-user is more appropriate than per-IP).

The "this might be wrong" qualifier is important. Without it, Claude tends to accept your approach uncritically. With it, Claude treats the proposal as a draft to be evaluated and improved.

Strawman proposals work well with self-critique: "Here is my approach. Critique it, suggest alternatives, then implement the best option."

Specific Prompts for Ambiguous Codebases

Large codebases develop naming collisions. There might be three components called "Dashboard." Two services called "Auth." Four files named "utils.ts." When you prompt Claude to "fix the bug in the Dashboard component," which dashboard?

In ambiguous codebases, extreme specificity is not pedantic -- it is necessary.

"Fix the rendering bug in `src/features/trading/components/Dashboard/TradingDashboard.tsx`. The issue is in the `useMarketData` hook around line 145 where the WebSocket connection is not being cleaned up on unmount. Do not modify `src/features/admin/components/Dashboard/AdminDashboard.tsx` which has a similar name but is unrelated."

The explicit "do not modify" for the similarly-named component is not paranoia. It is a constraint that prevents a wrong-component edit, which in a large codebase could pass code review because the reviewer does not notice which Dashboard was changed.

The more ambiguous your codebase, the more specific your prompts need to be. Exact file paths, line numbers, function names, and explicit exclusions. The cost is a few extra seconds of typing. The benefit is Claude modifying exactly the right thing.

Extended Thinking Configuration

Claude Code's extended thinking mode allocates additional tokens for reasoning before generating a response. The thinking happens internally -- you see the result, not the reasoning process. For complex tasks that require deep analysis, extended thinking dramatically improves output quality. For simple tasks, it wastes tokens and time.

Three mechanisms control it:

Environment variables. `MAX_THINKING_TOKENS` sets the budget for thinking tokens. The default is maximum (31,999 tokens). Set it lower (for example, `MAX_THINKING_TOKENS=10000`) to reduce cost on tasks that do not need deep reasoning, or set it to zero to disable thinking entirely. For the latest model with adaptive reasoning, thinking depth is controlled by effort level instead, and this variable is ignored unless set to zero.

`CLAUDE_CODE EFFORT LEVEL` sets the effort to `low`, `medium`, or `high` (default). Lower effort is faster and cheaper. Higher effort provides deeper reasoning. You can also adjust effort level interactively through the `/model` command using left/right arrow keys -- the change takes effect immediately.

Keyboard toggle. `Alt+T` (or `Option+T` on Mac) toggles extended thinking on and off during a conversation. This requires running the terminal setup command first (`/terminal-setup`).

A critical misconception: phrases like "think harder" or "ultrathink" in your prompt do not allocate more thinking tokens. They are just words in a prompt. The thinking budget is controlled entirely through the environment variables and the effort level setting. If you want Claude to reason more deeply, adjust the configuration -- do not add magic words to your prompt.

Image and Screenshot Workflows

Claude Code accepts images as input through several mechanisms, and each has a different best use.

Drag and drop. In the desktop application, drag an image file directly into the prompt area. In supported terminals, drag and drop also works.

Copy and paste. Use `Ctrl+V` (not `Cmd+V` on Mac -- this is a common mistake) to paste an image from the clipboard. Screenshots taken with the system screenshot tool land in the clipboard and paste directly.

File paths. Reference an image by its path: "Look at the screenshot at `./screenshots/bug.png` and identify the layout issue." Claude reads the image file and interprets its contents.

Opening referenced images. When Claude mentions an image path in its response, `Ctrl+Click` (or `Cmd+Click` on Mac) opens the image in your system viewer. This is useful when Claude generates images or references existing screenshots.

Screenshot-Driven Debugging

One infrastructure team discovered a pattern that extends image input beyond UI work. When container orchestration clusters went down and were not scheduling new pods, they fed screenshots of monitoring dashboards into Claude Code. Claude interpreted the dashboard visualizations, identified a warning indicating IP address exhaustion, guided the team through the cloud provider's management interface menu by menu, and provided the exact commands to resolve the issue. The screenshots were more precise than a verbal description of the dashboard state would have been.

This pattern applies wherever the symptom is visual but the solution is technical: infrastructure dashboards, error message screenshots, log viewer captures, network topology diagrams. Claude interprets the image and maps it to the technical context of the codebase.

Visual-First Iteration for Non-Developers

Non-technical users -- designers, product managers, legal staff -- frequently use screenshots as their primary communication medium with Claude Code. They paste a screenshot of what they want the interface to look like, and Claude implements it. They look at the result, take another screenshot showing what is wrong, and Claude adjusts. The iteration loop is entirely visual, with no code discussion at all.

This is a distinct interaction pattern from the developer workflow. Developers describe changes in terms of code. Non-developers describe changes in terms of appearance. Both are valid inputs. The screenshot-driven approach is slower per iteration but requires zero technical vocabulary.

Self-Critique: A Worked Example

The self-critique pattern becomes more concrete with a real example. One practitioner asked Claude Code to produce a financial optimization plan, then prompted: "Rate this plan 1-10. What would you improve?"

Claude rated its own plan 7.5 out of 10 and identified seven specific improvements. Among them was a significant catch: the plan had allocated certain assets in a way that ignored the tax-free status of a specific retirement account. Selling and reallocating within that account would cost zero in taxes, but the original plan had not taken advantage of this. The self-critique identified the missed optimization, and the revised plan incorporated it.

The lesson is not that self-critique always finds seven improvements. It is that self-critique finds the improvements that are not obvious on first generation. Claude does not include these optimizations initially because the generation process optimizes for the most common interpretation of the prompt. The self-critique step forces a second pass that catches what the first pass missed.

For analytical work, a particularly effective self-critique prompt is: "Include a self-critique section. Rate this plan 1-10 and list every improvement you would make." Making the self-critique a required section of the output -- not an afterthought -- ensures it receives adequate attention.

Domain-Specific Prompt Templates

The prompting patterns described so far are general-purpose. For domain-specific analytical work -- financial analysis, technical auditing, data science -- the prompt itself needs structural depth that generic patterns do not provide.

The Goal Prompt Anatomy

For complex analytical tasks, a detailed goal prompt follows a specific structure:

1. **What data exists and where.** "There are CSV exports from three accounts in /data/accounts/ . A supplementary text file at /data/notes.txt has information that was not exportable."
2. **Desired output format.** "Produce a table with columns: Asset, Current Allocation, Target Allocation, Gap, Priority, Action. Include a summary section and an appendix with methodology."
3. **Preferences and principles.** "Prioritize tax efficiency over raw performance. Minimize fees. Prefer simplicity."
4. **Known constraints.** "Do not modify the retirement account contributions -- those are automatic. The brokerage account has a minimum balance requirement of \$10,000."
5. **What you suspect is wrong.** "I think the portfolio is overweight in large-cap growth. If you have a different view, I want to hear it." The qualifier "if you have a different view, I want to hear it" prevents Claude from anchoring on your suspicion and ignoring evidence to the contrary.
6. **Strawman proposal.** "Here is my rough target allocation: 40% domestic equities split 60/40 large/mid-cap, 20% international, 30% fixed income, 10% alternatives. This might be wrong."
7. **Exhaustive questioning instruction.** At the end: "Ask me questions exhaustively about anything you are uncertain about before proceeding." Placing this at the end ensures Claude asks clarifying questions before jumping to conclusions.

This anatomy applies beyond finance to any analytical task where the prompt needs to convey domain context, output expectations, and nuanced preferences simultaneously.

Dynamic Constraint Reframing

Mid-project, you can reframe constraints and watch Claude rebuild the entire analytical framework. "Treat the retirement account as a fixed constraint. Calculate what the ideal total portfolio looks like, subtract what the retirement account provides, and optimize the remaining accounts to fill the gaps." Claude recalculates every target, every gap, every recommendation based on the new framing. The static portion becomes a given; the dynamic portion gets optimized.

This is more powerful than starting over with a new prompt. Claude retains the full context of the analysis and rebuilds on top of the reframed constraint rather than regenerating from scratch.

The Walkthrough Skill

Skills can encapsulate complex prompt patterns into reusable workflows. One community-built skill -- the walkthrough skill -- demonstrates how skills and prompt craft intersect.

When you invoke the walkthrough skill with a query like "walkthrough how does authentication work in this codebase," it executes a four-phase pipeline. First, it spawns two to four subagents to explore relevant parts of the codebase in parallel. Second, it synthesizes the findings into five to twelve key concepts with connections between them. Third, it generates a self-contained HTML file with an interactive diagram. Fourth, it opens the diagram in your browser.

Each node in the diagram is clickable. Click one and a detail panel slides in with a plain-English description, file paths, and code snippets. The whole thing produces a mental model of a system in under two minutes.

The walkthrough skill demonstrates a pattern that applies to any complex prompt workflow: wrap a multi-step process (parallel research, synthesis, artifact generation) in a skill so it becomes a single command. The prompt engineering is done once, tested, and reused.

The `/learn` Skill Pattern

Another skill pattern worth studying is the learning skill, which follows five phases: deep analysis of the current conversation, categorization of insights (patterns, gotchas, architecture decisions), drafting the learning into documentation, user approval (a blocking step -- the skill pauses and waits for your confirmation), and saving to the appropriate documentation file.

This pattern extracts institutional knowledge from conversations and persists it. Instead of having Claude's useful discoveries disappear when the session ends, the skill captures them as documentation that benefits future sessions and future developers. The blocking approval step ensures nothing is saved without human review.

Progressive Disclosure in Skills

A well-designed skill uses progressive disclosure to minimize context cost. Consider an analytics consultant skill: the skill directory contains a `skill.md` file with the skill definition and quick-start instructions, a `scripts/` directory with executable tools, and a `references/` directory with detailed documentation.

When a user's prompt matches the skill's description, only the `skill.md` loads into context. When the skill runs, it calls the CLI tool in `scripts/`. Only when the tool needs detailed guidance does it read from `references/`. The full reference material never enters context unless needed. This design keeps the base context cost near zero while supporting arbitrarily complex operations.

Two-Interface Planning

Teams that include non-technical members report a distinctive workflow: brainstorm in a cloud AI chat interface first, then move to Claude Code for implementation. The cloud interface is better for open-ended exploration -- thinking through the workflow, considering alternatives, sketching the approach. Once the plan is clear, they create a comprehensive prompt in the chat interface and paste it into Claude Code as the starting point.

The key is asking the cloud interface to slow down and work step by step rather than producing everything at once. The output of the planning phase is a structured prompt -- not code, not a design document, but the actual input that Claude Code will receive. This two-interface pattern lets non-developers use the interface they are comfortable with for the creative work and delegate the technical execution to the tool designed for it.

Providing Writing Samples

When Claude generates documentation, reports, or other written output, providing writing samples dramatically improves style matching. One team provides formatting preferences and example documents at the start of documentation tasks: "Here is an example of how we format runbooks. Match this style, including the heading structure, the bullet point density, and the level of detail in troubleshooting steps."

Claude extrapolates from examples more reliably than it follows abstract style descriptions. "Write in a concise, direct style" is vague. A three-paragraph sample of your actual writing style is precise. The cost is a few hundred tokens of example input. The benefit is output you can use without rewriting.

The Interrupt-for-Simplicity Pattern

Claude tends toward complex solutions. It builds abstractions, adds layers, introduces patterns. This is usually a feature -- complex problems need complex solutions. But sometimes the simple approach is correct, and Claude overshoots.

The fix is a specific interruption pattern: stop Claude mid-execution and ask "why are you doing this? Try something simpler." Claude responds well to simplicity requests. It drops the unnecessary abstraction, removes the extra layer, and produces a more direct solution.

This is not the same as the general interrupt-and-steer pattern described earlier. The general pattern corrects direction. The simplicity pattern corrects complexity. Watch for signs of over-engineering: new files that seem unnecessary, abstractions that wrap a single call, patterns that add indirection without adding value. When you see them, interrupt.

Key Takeaways

- Verification criteria (tests, expected output, visual targets) are the single highest-leverage prompting technique -- they convert one-shot generation into iterative refinement.
- Use vague prompts for exploration and specific prompts for implementation; the most common failure mode is vague implementation prompts.
- Your linting configuration, type checker, and test suite are part of your prompt -- strict tooling creates backpressure that improves Claude's output automatically.
- Ask Claude to plan before implementing; a 10,000-token plan prevents a 500,000-token misdirected implementation.
- State constraints explicitly -- what cannot change, trade-off preferences, and flexibility boundaries -- because implicit constraints are invisible constraints.
- Exhaustive questioning ("ask me questions exhaustively before starting") improves first-pass quality by 30-50% on complex tasks.
- In ambiguous codebases, extreme specificity with exact file paths and explicit exclusions prevents wrong-component modifications.

Chapter 9: Working with Large and Legacy Codebases

What You'll Learn

Small, well-structured projects are easy. Claude Code reads the whole thing in minutes, understands the architecture, and starts producing useful work immediately. The real test is a monorepo with two million files, a legacy codebase with three frameworks and no documentation, or a system so large that no single person understands all of it.

These are the projects where Claude Code's value proposition is highest and where getting the approach wrong costs the most. This chapter covers the specific techniques for working with codebases that are too large to fit in context, too old to follow modern conventions, and too critical to rewrite carelessly. You will learn how git worktrees enable parallel work, how the Explore subagent navigates unfamiliar code efficiently, how task dependency management coordinates complex multi-step operations, and how teams have used Claude Code to compress three-year rewrites into weeks. You will also see how filesystem-driven analysis turns messy data exports into structured plans, how a storage layer migration demonstrates spec-driven refactoring at scale, and why the language barrier between modern developers and legacy systems written in decades-old languages is collapsing faster than anyone predicted.

If your codebase makes new hires cry, this chapter is for you.

Git Worktrees for Parallel Work

The simplest scaling technique is also the most underused: git worktrees.

A worktree creates a separate working directory for a branch without cloning the entire repository again. Each worktree has its own checked-out files, its own index, and -- crucially -- can run its own independent Claude Code session. The sessions do not share

context. They do not interfere with each other. They operate on the same repository but in completely separate directories.

```
git worktree add ../feature-auth feature/auth
git worktree add ../feature-payments feature/payments
```

Now you can run Claude Code in `../feature-auth` working on authentication while simultaneously running another instance in `../feature-payments` working on the payment system. Both sessions have full access to the repository history, but their file changes, context windows, and conversation histories are completely independent.

This matters for large codebases because the biggest bottleneck is often sequential work. You cannot run two Claude Code sessions in the same directory on different branches -- they would step on each other's files. Worktrees eliminate that constraint.

The pattern scales linearly. Three branches, three worktrees, three sessions. The limit is your machine's resources and your API budget, not anything architectural. For teams working on a monorepo with multiple independent features in flight, worktrees turn one Claude Code experience into several parallel ones.

One practical detail: each worktree gets its own `.claude/` directory if one exists in the repository, so project-level CLAUDE.md and settings apply independently. This means CLAUDE.md content (as covered in Chapter 3) compounds across worktrees without conflict.

File Suggestion Customization for Monrepos

Claude Code's `@` autocomplete suggests files as you type, which works well for typical repositories. For large monorepos, the default file indexing can be slow or miss files deep in nested directory structures.

The solution is custom indexing. You can provide a script that generates the file suggestions Claude Code offers when you type `@`. This script can implement any logic: prioritize recently modified files, exclude build artifacts, index only the packages you care about, or provide a completely custom ordering.

For a monorepo with hundreds of packages, a custom indexing script that filters to the packages relevant to your current work can transform the `@` experience from overwhelming to precise. Instead of scrolling through thousands of files, you see the twenty that matter.

This is a small feature with outsized impact. In a ten-thousand-file monorepo, being able to quickly reference the right file saves context (you point Claude at the right code instead of letting it search) and saves time (you skip the exploration phase when you already know where to look).

The Explore Subagent

When you do not know where to look, the Explore subagent is purpose-built for navigation.

Explore is a read-only, low-cost subagent (its architecture is covered in Chapter 4). It cannot modify files or run commands that change state. It can only read, search, and report back. This makes it fast, cheap, and safe to deploy even in codebases where you are still learning the terrain.

The Explore subagent runs in its own context window, which means it can read dozens of files during its investigation without consuming your main conversation's context. It returns a summary. You get the answer without the cost of the journey.

Explore also supports configurable thoroughness. For a quick question about where a function is defined, it searches narrowly and returns fast. For a comprehensive understanding of how a module works, it searches broadly, follows cross-references, and produces a detailed analysis. The thoroughness parameter lets you trade speed for depth.

Use Explore liberally in large codebases. The cost is low -- it runs on the cheapest available model, and the results run in an isolated context. The alternative is using your main session to read files one by one, filling your context window with exploration that could have been delegated.

Parallel Research Subagents

Explore is a single investigator. For understanding a large system, you want a research team.

Multiple subagents can run in parallel, each investigating a different aspect of the codebase. One reads the authentication module. Another examines the database access layer. A third reviews the API routing. Each returns a focused summary. The main conversation orchestrates the investigation and synthesizes the results.

This pattern is particularly effective for onboarding onto an unfamiliar system. Instead of spending hours reading code sequentially, you dispatch parallel investigators and get a multi-perspective understanding of the system in minutes. The cost is higher than a single Explore call -- you are running multiple subagents -- but the time savings are substantial, and the context isolation means your main window stays clean (as covered in Chapter 4).

The research team pattern also works for impact analysis before making changes. Before modifying a shared utility function, dispatch subagents to investigate every module that imports it. Each subagent reports how the function is used in its area. You get a comprehensive impact assessment without manually tracing dependencies.

For massive codebases, this is not optional. No human can hold a multi-million-line codebase in their head. Parallel subagents give you a way to survey broad territory quickly, then focus your main session on the specific areas that matter.

Task Dependency Management

Complex work on large codebases often involves tasks that depend on each other. You cannot update the API endpoints until the database schema is migrated. You cannot write integration tests until both the API and the database changes are complete. You cannot deploy until the tests pass.

Claude Code's task system supports explicit dependency declarations using `blocks` and `blockedBy` relationships. This enables wave-based execution: tasks with no dependencies run first, tasks that depend on them run next, and so on until everything is complete.

The pattern works like this:

1. Plan the work by breaking it into discrete tasks with clear boundaries.
2. Declare dependencies between tasks. Task A blocks Task B. Task C is blocked by both A and B.
3. Execute. Claude Code (or an agent team) processes tasks in dependency order, running independent tasks in parallel and dependent tasks sequentially.

For a large migration -- say, moving from one ORM to another across fifty modules -- this turns a manual coordination nightmare into a structured execution plan. The first wave migrates the core database layer. The second wave updates the modules that depend on it. The third wave updates the tests. Each wave can run with parallelism within it, but the waves themselves execute sequentially.

This is infrastructure-grade project management, not clever prompting. It works because the dependency system is explicit and enforced, not because Claude Code is "smart enough to figure out the order." You define the structure. Claude Code executes it.

Filesystem-Driven Analysis

Legacy codebases come with legacy data. Directories full of CSV files with inconsistent formatting. Configuration files in three different formats. Log files from systems that no longer exist. Documentation that was last updated two years ago.

Claude Code handles this well because it treats the filesystem as a first-class data source. Point it at a directory of messy data and give it an objective -- normalize these files, find duplicates, extract a summary, build a migration plan -- and it will autonomously parse, compare, and organize the contents.

This works because Claude Code's tool access includes file reading, directory listing, and command execution. It can write scripts to process files, run those scripts, inspect the output, and iterate. It is not limited to understanding one file at a time. It can build a mental model of a collection of files and operate on the collection as a whole.

For legacy systems, this capability unlocks a specific workflow: understanding before rewriting. Before you touch a single line of code, point Claude Code at the system's configuration, data, and documentation directories. Let it build a picture of how the system actually works, not how someone once said it worked. That understanding -- exported as a spec document or CLAUDE.md additions -- becomes the foundation for any subsequent rewrite.

The Three-Input Workspace Pattern

The most effective filesystem-driven analysis follows a specific setup that one practitioner documented while building a portfolio optimization plan that would have cost thousands of dollars from a professional advisor. The pattern generalizes beyond finance to any domain where messy data needs structured analysis.

The workspace starts with three inputs:

1. **Raw data exports.** CSV files, database dumps, configuration exports -- whatever the system produces natively. These files are typically messy: inconsistent column names, encoding issues, overlapping data across exports, non-standard line items. That mess is the point. You are feeding Claude Code the reality of the system, not a sanitized version.
2. **A supplementary text file.** There is always information that has no export. Institutional knowledge, undocumented settings, manual overrides, context that exists only in someone's head. Type it into a plain text file as a bulleted list. It does not need structure. It needs the raw facts.
3. **A detailed goal prompt.** This is the most important piece. It describes what data exists and where, the desired output format, your preferences and principles, known constraints, what you suspect is wrong (with room for Claude to disagree), and any strawman proposals that give the model something concrete to react to rather than building from a blank slate.

The workflow then proceeds through a predictable sequence. Initialize a Claude Code session pointed at the directory containing all three inputs. Claude autonomously writes scripts to parse the data files -- handling encoding issues, deduplicating overlapping exports, classifying items into target categories, and dealing with non-standard entries that would otherwise be misclassified. It iterates through problems as it encounters them. The baseline analysis it produces -- a full gap analysis, a structured summary, a comparison against targets -- is typically strong enough to work from immediately.

From there, the interaction shifts to iteration. You clarify constraints, reframe assumptions, ask Claude to self-critique its own output. Each clarification propagates instantly through every calculation, table, and recommendation in the plan. The supplementary text file and CLAUDE.md accumulate the decisions and data quirks discovered along the way, so subsequent sessions pick up exactly where the previous one left off.

Create a CLAUDE.md file early in the process and have Claude update it at the end of each session with learnings -- data quirks discovered, strategy decisions made, target parameters refined. This turns multi-session analysis into a compounding process where every future session builds on everything that came before. One practitioner described the result as equivalent to what would have been a week of spreadsheet work or several thousand dollars of professional advisory services, produced instead over a few evenings of back-and-forth with an AI agent.

The Language Barrier Disappears

There is a class of legacy system that most modern developers refuse to touch: the ones written in languages they have never used. Mainframe systems running business logic in languages from the 1950s and 1960s. Scientific computing infrastructure maintained in numerical languages from the same era. Domain-specific languages that predate the internet. These systems run payroll, air traffic control, financial settlement, and nuclear simulations. They are not going away, and for decades, the shrinking pool of engineers who understood them was the bottleneck for any maintenance or modernization effort.

Agentic coding tools are dissolving that barrier. Support is expanding to less-common and legacy languages -- not as an afterthought, but as a natural consequence of how large language models work. The model has seen these languages in its training data. It can read them, understand their idioms, and translate their logic into modern equivalents. A developer who has never written a line in any of these legacy languages can point Claude Code at a directory of source files and get a working explanation of the business rules encoded in them.

This matters for large-codebase work because many legacy systems are also enormous systems. A financial institution's core processing logic might span hundreds of thousands of lines of code written over four decades. The traditional modernization approach -- hire a dwindling pool of specialists at escalating rates, then manually translate each module -- takes years and costs millions. The Claude Code approach is fundamentally different: use the model to read the legacy code, extract the business logic, and generate equivalent modern implementations. The specialist's role shifts from writing code to reviewing translations and validating business logic preservation.

The practical implications are significant. Technical debt that accumulated for years because no one had time to address it -- or no one who understood the language was available -- becomes systematically eliminable. Agents can work through backlogs of legacy maintenance tasks that were previously non-viable. An organization that had a three-year queue of modernization work might clear it in months, not because the agents are faster at writing the same code, but because the language barrier that made the work impossible for 95% of their engineering staff no longer exists.

Legacy Codebase Rewrites

The most dramatic Claude Code case studies involve legacy rewrites. One documented case involved a trading platform frontend that had taken three years to build manually. Using Claude Code, the entire frontend was rewritten in weeks.

The timeline compression is real, but the pattern behind it matters more than the headline number.

Legacy rewrites succeed with Claude Code when three conditions hold:

1. **The desired architecture is well-defined.** Claude Code needs to know what you are building toward, not just what exists today. A clear target architecture -- in a spec document, in CLAUDE.md, in detailed prompts -- gives Claude Code a destination.
2. **The existing code is the specification.** Legacy code is often the only accurate documentation of business logic. Claude Code can read the existing implementation and extract the rules it encodes, even when those rules are buried in spaghetti code. It does not need clean code to understand behavior.
3. **The rewrite is modular.** Rewriting an entire system in one shot is as risky with Claude Code as it is manually. The effective pattern is module-by-module migration: rewrite one component, test it against the legacy behavior, validate, move to the next. Each module is a bounded task that fits in a single session.

The three-year-to-weeks compression comes from parallelism (multiple modules rewritten simultaneously via worktrees or subagents), from Claude Code's speed at producing code once the architecture is clear, and from the elimination of the knowledge acquisition phase that dominates manual rewrites. Claude Code reads the legacy code as fast as it reads anything else. The six months a human team might spend understanding the existing system before writing a line of new code shrinks to hours.

Autonomous Implementation at Scale

At the extreme end of the spectrum, Claude Code has been used for autonomous implementation on codebases exceeding twelve million lines of code. One documented case involved implementing a feature across a codebase of that scale in approximately seven hours.

Seven hours. Twelve and a half million lines. One feature. No human writing code. And 99.9% numerical accuracy compared to the reference implementation.

That last number is the one that matters most. Speed without accuracy is just fast failure. The implementation was not a rough approximation that needed human cleanup. It achieved near-perfect fidelity to the expected results, autonomously, across a codebase spanning multiple programming languages.

This is not typical usage. It represents the ceiling of what is possible with current tooling. But the pattern it demonstrates is instructive: Claude Code's effectiveness on large codebases is not limited by the codebase's size. It is limited by the clarity of the task definition and the quality of the context provided.

A twelve-million-line codebase does not fit in any context window. What fits is the task description, the relevant subset of files, and the architectural knowledge encoded in CLAUDE.md. Claude Code uses its tools -- file search, grep, directory traversal -- to find the relevant code, understand it, and modify it. The vast majority of the codebase is never loaded into context. It does not need to be.

This selective-loading approach is how large codebases become tractable. Claude Code does not need to understand the entire system. It needs to understand the part of the system relevant to the current task. The Explore subagent, file search tools, and grep handle the navigation. The main session handles the implementation. CLAUDE.md provides the architectural guardrails.

Codebase Navigation for Onboarding

One of the highest-value applications of Claude Code on large codebases has nothing to do with writing code.

New team members joining a large project traditionally face weeks of onboarding. They read documentation that is partially outdated. They ask colleagues questions that interrupt productive work. They explore code by opening files semi-randomly. They build a mental model slowly, through accumulation and osmosis.

Claude Code compresses this process dramatically. A new developer can start a session and ask questions about the codebase:

- "How does authentication work in this system?"
- "Where is the payment processing logic, and what external services does it call?"
- "What is the relationship between the Order and Fulfillment modules?"
- "Why does this project have two different database connection pools?"

Claude Code dispatches Explore subagents, reads relevant code, and provides answers grounded in the actual codebase -- not in documentation that might be stale. The answers reference specific files, specific functions, specific patterns. They are as current as the code itself.

This replaces data catalogs, architecture wikis, and colleague consultations for a significant portion of onboarding questions. It does not replace the human relationships and cultural context that onboarding also requires. But for "how does this code work?" questions, Claude Code is faster, more accurate, and more patient than any human mentor.

Teams that maintain a well-curated CLAUDE.md amplify this effect. The new developer gets both the code-derived answers from Claude Code's exploration and the institutional knowledge captured in CLAUDE.md. The combination provides an onboarding experience that would have required weeks of a senior engineer's time to deliver manually.

Spec-Driven Refactoring: A Storage Layer Migration

Legacy rewrites and large refactors share a common failure mode: jumping straight to code without understanding the target well enough. The spec-driven approach inverts this by making research and specification the first phase, not an afterthought.

One practitioner documented a complete storage layer migration -- replacing a compiled-to-browser database with a native browser database for a sync engine. The old approach worked but had problems: large binary dependencies, poor performance on initial load, and complications with the framework's built-in sync capabilities. The migration touched fifteen or more files and required understanding an unfamiliar framework's storage patterns.

The workflow proceeded in distinct phases. First, research: a single prompt asking Claude to investigate the target framework spawned five parallel research subagents, each studying a different aspect -- the framework's data model, its sync protocol, its storage layer, its conflict resolution patterns, and its API surface. Each subagent returned a focused summary. The combined findings became a comprehensive research report that would have taken a developer days of documentation reading to assemble.

Second, specification: Claude synthesized the research into a detailed migration spec -- a structured document covering architecture decisions, a phased implementation plan, a fourteen-item task checklist, risk mitigation strategies, and success criteria. The spec was written to a file on disk, not held in conversation context. This is critical. A spec on disk survives context degradation, session restarts, and compaction. It is the persistent source of truth.

Third, refinement: before implementation, Claude asked clarifying questions about ambiguities in the migration strategy -- conflict resolution approaches, sync behavior during the transition, fallback handling. These questions surfaced design decisions that would

have become bugs if discovered during implementation instead of during planning.

Fourth, execution: fourteen tasks, each delegated to a subagent, each producing an atomic commit. The subagents worked from the spec, not from accumulated conversation context. Each completed its bounded task, committed the result, and returned. Fourteen commits, fifteen-plus files changed, one pull request ready for review.

The entire migration took a single afternoon. The practitioner estimated it would have taken two to three days manually. But the time savings are secondary to the quality improvement: the research phase uncovered framework patterns that the developer would not have found on their own, resulting in a more idiomatic implementation than manual coding would have produced.

Despite orchestrating fourteen subagents, the main session's context stayed well within limits -- the orchestrator held only the task list and subagent summaries, while the actual file reading and code generation happened in isolated subagent contexts. This is the model for large refactors: research broadly, specify precisely, execute in parallel, and keep the orchestrator lean.

CLAUDE.md as Data Catalog Replacement

One of the more unexpected applications of CLAUDE.md in large codebases has nothing to do with code style or build commands. One data infrastructure team discovered that their CLAUDE.md files could replace traditional data catalogs and discoverability tools for onboarding.

When new data scientists joined the team, they were directed to use Claude Code to navigate the massive codebase. Claude Code would read the CLAUDE.md files, identify relevant files for specific tasks, explain data pipeline dependencies, and help newcomers understand which upstream data sources fed into which dashboards. The information that traditionally lived in a separate data catalog tool -- and was perpetually outdated because nobody remembered to update it -- now lived alongside the code and was maintained as part of the normal Claude Code workflow.

The team reinforced this with a continuous improvement loop. At the end of each work session, they asked Claude to summarize the completed work and suggest improvements -- not just to the code, but to the CLAUDE.md documentation and the workflow instructions themselves. Each session refined the project's institutional knowledge. The CLAUDE.md file grew organically from actual usage rather than from a dedicated documentation sprint that would never be repeated.

This pattern works because CLAUDE.md is read at the start of every session with full fidelity. Unlike a wiki page that might be six months stale, CLAUDE.md is as current as the last session that updated it. For large codebases where the data model is as complex as the code itself -- where understanding which table feeds which dashboard through which intermediate transformation is essential knowledge -- CLAUDE.md becomes the living data catalog that traditional tools aspire to be but rarely achieve.

Technical Debt at Scale

Every large codebase has a backlog of technical debt that nobody gets to. Not because it is unimportant, but because the economics never justified the investment. The refactor that would take three weeks of an experienced engineer's time saves two minutes per build. The migration that would modernize a deprecated dependency requires touching two hundred files. The test coverage improvement that would prevent the quarterly production incident needs someone who understands both the testing framework and the business logic.

Agentic coding tools change the economics of technical debt systematically. When agents can work autonomously for extended periods, formerly non-viable projects become feasible. The three-week refactor becomes a three-hour task with review. The two-hundred-file migration becomes a structured execution plan with parallel subagents (as described in the task dependency management section above). The test coverage improvement becomes a weekend of autonomous operation with human review on Monday.

The pattern is not "set an agent loose on the backlog." It is structured: identify the debt, write a clear specification, establish verification criteria (tests that must pass, linting rules that must be satisfied), and then let agents work through the items systematically. The backpressure of automated verification (Chapter 8) ensures that agent-generated fixes actually work. The task dependency system ensures that dependent changes execute in the right order.

Organizations that recognize this shift early gain a compounding advantage. Every piece of technical debt eliminated makes the codebase easier for both humans and agents to work with. Cleaner code produces better agent outputs. Better agent outputs clear more debt. The flywheel accelerates.

Strategies for the Enormous and the Ancient

Working with large and legacy codebases is ultimately about managing two scarce resources: context and comprehension.

Context management for large codebases means aggressive use of subagents for exploration, selective file loading, and keeping your main conversation focused on the task at hand rather than on understanding the system. Do not read files into your main context that you could delegate to an Explore subagent. Do not explore broadly when you can search narrowly.

Comprehension management means building understanding incrementally and persisting it. Use spec documents and CLAUDE.md to capture what Claude Code learns about the system. Use task dependency management to structure complex migrations. Use parallel research to build understanding faster than any sequential reading could.

The codebases that seem impossible to work with are not impossible. They are just expensive to understand. Claude Code reduces that cost -- dramatically -- by reading faster, searching more broadly, and forgetting nothing that you write into CLAUDE.md.

Key Takeaways

- Git worktrees enable parallel Claude Code sessions on different branches of the same repository, with no interference between sessions.
- The Explore subagent (Haiku, read-only, isolated context) is cheap enough to use liberally for navigating unfamiliar code.
- Parallel research subagents provide multi-perspective understanding of large systems in minutes instead of hours.
- Task dependency management with `blocks` / `blockedBy` enables wave-based execution for complex migrations and refactors.
- Legacy rewrites succeed when the target architecture is clear, the existing code serves as the specification, and the work is modular.
- Claude Code's effectiveness on large codebases is limited by task clarity and context quality, not by codebase size -- one documented case achieved 99.9% numerical accuracy across a 12.5-million-line codebase.
- The language barrier for legacy systems written in older languages is disappearing; agents can read, understand, and translate business logic that most modern developers cannot.
- The three-input workspace pattern (raw data exports, supplementary text file, detailed goal prompt) turns filesystem-driven analysis into a structured, repeatable process.
- Spec-driven refactoring -- research first, specify precisely, execute in parallel -- produces better results than jumping straight to code, even when the developer could have written it manually.
- CLAUDE.md can replace traditional data catalogs for onboarding, maintained through a continuous improvement loop at the end of each session.
- Technical debt becomes systematically eliminable when agents change the economics from "not worth three weeks of an engineer's time" to "three hours plus review."

Chapter 10: Failure Modes and Recovery

What You'll Learn

Claude Code fails. Not occasionally. Regularly. Understanding how it fails is more valuable than understanding how it succeeds, because success is straightforward and failure is subtle. The tool does not crash with an error message that tells you what went wrong. It silently degrades, quietly drops tasks, confidently generates wrong code, and gradually loses coherence in ways that look like normal operation until you inspect the output.

This chapter catalogs the failure modes -- not as a warning to avoid the tool, but as a field guide for working with it effectively. Every failure mode has a recovery strategy. Some recoveries are elegant. Most are pragmatic. A few are just "start over." You will learn to

recognize the symptoms of context exhaustion before your session collapses, understand why bash commands create invisible problems, know what checkpoints can and cannot recover, and develop the judgment to distinguish between a fixable mistake and a situation where reverting is the only efficient path forward. You will also follow the complete arc of a 33-day autonomous trading experiment -- from overconfident first trades through governance failures, market crashes, and eventual recovery -- that illustrates every category of autonomous agent failure in a single, continuous narrative. And you will learn the named anti-patterns that experienced practitioners have cataloged, with specific fixes for each.

Context Exhaustion: The Slow Death

The most common failure mode is also the most insidious. Context exhaustion does not announce itself. There is no error message. No warning popup. Claude just... gets worse.

The symptoms are recognizable once you know what to look for. Claude starts repeating suggestions it already made. It forgets instructions you gave five minutes ago. It solves a problem it already solved. It drops constraints from your original prompt -- not all of them, just enough that the output is subtly wrong. It generates code that contradicts patterns it was following earlier in the session.

This happens because the context window is full or nearly full. Auto-compaction (Chapter 3) kicks in and summarizes older content to make room, but the summarization is lossy. The resulting context is a degraded version of what you started with, and Claude works from the degraded version without knowing it is degraded.

Mitigation. The defense is not reactive but structural. Use the context engineering techniques from Chapter 3: put critical instructions in CLAUDE.md (which reloads at full fidelity and survives compaction), use the "Compact Instructions" section for rules that must persist, route verbose operations through subagents (Chapter 4) to keep the main context lean, and compact proactively using `/compact` rather than waiting for auto-compaction.

Watch for the behavioral tells. When Claude starts asking questions it already has the answer to, when it proposes approaches you already rejected, when it modifies files you told it not to touch -- those are context exhaustion symptoms. The right response is not to repeat your instructions. The right response is to start a fresh session with a good CLAUDE.md file, because repeating instructions in a full context window just accelerates the exhaustion.

Bash Environment Non-Persistence

This one bites everyone exactly once, and then it bites them again because they forgot.

Each bash command Claude executes runs in a fresh shell environment. The working directory persists between commands. Nothing else does. Environment variables, shell aliases, function definitions, activated virtual environments -- all gone after each command.

The failure mode is silent. Claude sets an environment variable in one command and references it in the next. The second command does not error with "variable not found." It runs with the variable unset, which might mean an empty string, which might mean it does the wrong thing quietly. Claude sets `NODE_ENV=production` to test something, runs the test in the next command, and the test runs in the default development mode. No error. Wrong result.

Recovery. There is no recovery for the silent wrong result -- you need to catch it. The prevention is awareness. When Claude chains bash commands that depend on shared state, the state needs to be re-established in each command. Set the variable and run the command on the same line: `NODE_ENV=production npm test`. Or source a setup script at the start of each command. Or use Claude's environment variable configuration in settings.json, which injects variables into every bash command automatically.

Watch for this especially in: virtual environment activation (`source venv/bin/activate` does not persist), directory-specific environment files that are sourced once, shell function definitions used across commands, and any workflow where "run this then run that" assumes shared state.

MCP Disconnection: Tools That Vanish

MCP servers connect at session start and can disconnect at any time. When they disconnect, the tools they provided simply disappear from Claude's available tool set. There is no notification. No error. Claude just stops being able to do things it could do a minute ago.

The failure mode presents as Claude suddenly being unable to perform tasks it was handling fine. If Claude was using an MCP-provided tool to query a database and the server disconnects, Claude does not say "I lost my database tool." It tries to find alternative approaches, fails at them, and produces increasingly confused output as it tries to work around a missing capability it does not know is missing.

Recovery. The `/mcp` command shows the status of connected MCP servers. If a server is disconnected, you can restart it from there. Proactively, check MCP status when Claude's behavior changes unexpectedly, especially if it was previously using external tools fluently. MCP reliability details are covered in Chapter 5, but the recovery is the same regardless of cause: reconnect the server or restart the session.

Subagent Context Return Overflow

Subagents solve context exhaustion by isolating work in their own context windows. But their results have to come back to the main session, and those results consume main context space.

The failure pattern: you launch ten subagents to explore ten modules. Each returns a detailed report. Ten detailed reports flood the main context window. The orchestrator now has less room for its own reasoning than if you had never used subagents at all.

This is the subagent context paradox. Subagents protect the main context from intermediate work (file reads, command outputs, exploration noise) but not from final results. If the final results are verbose, the protection is partial at best.

Mitigation. Instruct subagents to return concise summaries, not comprehensive reports. "Return only the specific files that need changes and a one-sentence description of each change" is better than "report your findings." Better still, instruct subagents to write their detailed findings to a file and return only the file path. The main session can then read that file selectively, taking only the parts it needs into context.

The broader lesson is that subagent orchestration requires thinking about information flow. How much data crosses the subagent-orchestrator boundary, and in which direction? The outbound data (task assignment) is usually small. The inbound data (results) is the risk. Design for small inbound data.

The Checkpoint System: What It Tracks and What It Misses

Claude Code automatically tracks every file edit as you work, creating checkpoints that let you rewind to previous states. Every user prompt creates a new checkpoint. Checkpoints persist across sessions, so you can access them even in resumed conversations. This is a genuine safety net -- one of the best features for recovering from wrong turns.

The rewind interface is accessed by pressing Esc twice (Esc + Esc) or by typing `/rewind`. This opens a scrollable list showing each of your prompts from the session. Select the point you want to act on, then choose from four options:

1. **Restore code and conversation** -- revert both your files and the conversation history to that point. This is the full rewind: Claude loses memory of everything after that checkpoint, and the code returns to its state at that moment.
2. **Restore conversation only** -- rewind the conversation to that message while keeping the current code. Useful when Claude went down a wrong reasoning path but the code changes are fine.
3. **Restore code only** -- revert the file changes while keeping the full conversation history. Useful when the code went wrong but the discussion contains valuable context you want to preserve.
4. **Summarize from here** -- this is different from the restore options. It does not revert anything. Instead, it condenses the conversation from the selected point onward into a compact summary, freeing context window space while preserving the key information. This is a context management tool, not a recovery tool.

The three restore options undo state. Summarize compresses it. The distinction matters when you are deciding which recovery path to take.

But checkpoints have a gap that will burn you.

Checkpoints track direct file edits -- when Claude uses the Write, Edit, or NotebookEdit tools. They do not track file changes made through bash commands. If Claude runs `rm important_file.py` or `mv src/old.py src/new.py` or `sed -i 's/foo/bar/g' *.py` through the bash tool, those changes are invisible to the checkpoint system. Rewinding to a checkpoint before those commands does not undo them.

Checkpoints also do not track changes made by concurrent sessions, external processes, or anything that happens outside Claude Code's tool invocations. If another Claude Code instance modifies a file, those modifications are not in your session's checkpoint history.

Recovery. Git is the reliable checkpoint. If you followed the commit-frequently pattern from Chapter 1, you have git commits that cover everything -- file edits, bash modifications, renames, deletions. `git checkout` is the universal rewind, not just for Claude's tool-tracked edits but for anything. Think of checkpoints as "local undo" and git as "permanent history." They complement each other but are not interchangeable.

The operational rule: before any session that might involve destructive bash commands (file deletions, renames, bulk replacements), commit your current state. If you commit before and Claude wrecks things through bash, `git checkout .` brings everything back. If you relied on checkpoints and Claude wrecked things through bash, you are recovering from partial information.

Agent Amnesia

Every new session starts with a blank context window. The sophisticated understanding Claude developed about your codebase, the decisions you negotiated, the constraints you established -- all gone. This is agent amnesia, and it is the default state.

The failure mode is not that Claude forgets. It is that Claude does not know it forgot. It approaches the codebase fresh, makes the same assumptions the previous session started with, and potentially repeats the same mistakes the previous session corrected.

Recovery. There are three persistence mechanisms, and you should use all of them.

First, CLAUDE.md files (Chapter 3). These reload at full fidelity on every session start. Put critical project knowledge, architectural decisions, and hard-won lessons here. After a productive session, ask Claude to suggest CLAUDE.md updates that capture what it learned.

Second, the task system. Tasks stored as JSON files in `.claude/tasks/` persist across sessions. If you are running a multi-step project, the task list provides continuity between sessions -- Claude can read the task statuses and pick up where the previous session left off.

Third, git history. The code itself is a form of memory. Claude can read recent commits, understand what changed, and infer intent from the diff. This is less reliable than explicit instructions but better than nothing.

The developers who suffer most from agent amnesia are the ones who rely on long sessions to maintain context instead of externalizing knowledge into persistent structures. A two-hour session that ends without updating CLAUDE.md is two hours of context development that evaporates.

Context Pollution

Context pollution is what happens when the window fills with content that is not useful for the current task. Old exploration results. Verbose error outputs from solved problems. Discussion about approaches that were rejected. All of it takes up space that could hold information relevant to the task at hand.

The failure mode: Claude has plenty of context window remaining by token count, but the useful signal-to-noise ratio has degraded. Claude's responses become less focused because the relevant context is diluted by irrelevant content. In severe cases, Claude drops bugs rather than tracking them -- the window is so polluted that new information gets lost in the noise.

Mitigation. Proactive compaction with `/compact` clears the accumulated noise. Subagent delegation prevents pollution in the first place -- when you route exploratory work through a subagent, the exploration noise stays in the subagent's context and only the result enters the main window.

Spec-driven development (Chapter 8) provides structural mitigation. When Claude works from a spec document on disk rather than from accumulated conversation context, the spec serves as a persistent, unpolluted source of truth that survives context degradation. The conversation can fill with noise, but the spec file remains clean.

The Five Named Anti-Patterns

Experienced practitioners have cataloged five recurring failure patterns, each with a specific fix. These are not abstract categories. They are the mistakes that waste the most time in real-world usage.

The kitchen sink session. You start with one task, then ask Claude something unrelated, then go back to the first task. The context fills with irrelevant information from the detour. Claude's performance on your original task degrades because the signal-to-noise ratio in the context window has collapsed. The fix is simple: use `/clear` between unrelated tasks. A fresh context for each distinct task outperforms a single session that accumulates everything.

Correcting over and over. Claude gets something wrong. You explain the problem. Claude patches it. The patch introduces a new issue. You explain that. Claude patches again. Three corrections later, the context is full of failed approaches, and Claude is more confused than when it started. The fix: after two failed corrections, stop. Use `/clear` and write a better initial prompt that incorporates what you learned from the failures. The knowledge of what went wrong goes into the new prompt, not into a correction spiral inside a degraded context.

The over-specified CLAUDE.md. If your CLAUDE.md file is too long, Claude ignores half of it because important rules get lost in the noise. The document that was supposed to provide guardrails becomes a wall of text that Claude skims. The fix is ruthless pruning. If Claude already does something correctly without the instruction, delete that instruction. Convert style rules that Claude frequently violates into hooks that enforce them automatically rather than instructions that hope for compliance. CLAUDE.md should contain only what Claude cannot figure out by reading the code and what it has been observed to get wrong.

The trust-then-verify gap. Claude produces a plausible-looking implementation. It compiles. It runs. You ship it. Then it fails in production on an edge case that the implementation never considered. The output looked right, so you assumed it was right. The fix: always provide verification. Tests, scripts, screenshots, manual spot checks. If you cannot verify a piece of output, do not ship it. The verification infrastructure is not overhead. It is the mechanism that converts Claude's confident-but-uncertain output into validated output.

Infinite exploration. You ask Claude to "investigate" something without scoping the investigation. Claude reads hundreds of files, filling the context window with exploration data. By the time it finishes investigating, there is no context left for actual work. The fix: scope investigations narrowly, or use subagents so the exploration happens in an isolated context and only the summary enters the main window. "Investigate how authentication works in src/auth/" is scoped. "Investigate the codebase" is not.

Stop Hook Infinite Loops

Hooks (Chapter 2) can trigger on lifecycle events, including when Claude stops executing. A stop hook runs when Claude decides it has finished a task. If the stop hook determines that Claude should not have stopped -- say, because tests are still failing -- it can instruct Claude to continue.

The failure mode: a stop hook that always tells Claude to continue. Claude finishes, the hook fires, the hook says "keep going," Claude does more work, finishes again, the hook fires again, says "keep going" again. This is an infinite loop. Claude never stops because the stop condition is never met.

The system includes a `stop_hook_active` field to prevent runaway loops, but the underlying design risk is real. If your stop hook's "keep going" condition is based on criteria that Claude cannot actually satisfy -- tests that fail for reasons unrelated to Claude's changes, linting rules that are impossible to pass, external services that are down -- the loop will exhaust your token budget before it converges.

Prevention. Include iteration limits in stop hooks. "If tests still fail after three attempts, stop and report the remaining failures." Test stop hooks with intentionally failing conditions to verify they eventually terminate. Monitor token usage when using stop hooks for the first time.

Agent Team Limitations

Agent teams (Chapter 4) are experimental, and their limitations are shaped by that status.

No session resumption. If a team member's session crashes or is interrupted, it cannot be resumed. The work is lost. For long-running team tasks, this means a crash at 90% completion requires restarting from scratch.

Task status lag. Updates to the shared task list are not instantaneous. An agent may pick up a task that another agent has already started working on, leading to duplicate effort.

One team per session. The orchestrator session can have only one active team. If you need multiple independent teams, you need multiple orchestrator sessions.

No nested teams. A team member cannot create its own sub-team. The hierarchy is flat: one orchestrator, multiple team members. This limits the complexity of decomposition you can express.

Recovery. The pragmatic response is to keep team tasks small enough that losing one is tolerable, to accept some duplicate work as the cost of loose coordination, and to check task statuses manually when precise coordination matters.

Hallucination in High-Stakes Contexts

Claude Code generates code. Code that runs. Code that passes tests. Code that looks correct. Code that is wrong.

Hallucination -- generating confident, plausible, incorrect output -- is an inherent property of large language models. In low-stakes contexts (documentation, test scaffolding, boilerplate), hallucinations are caught quickly and cheaply. In high-stakes contexts (financial calculations, security logic, data transformations with correctness requirements), hallucinations can be expensive.

One practitioner ran an autonomous trading experiment and documented a 22.4% peak-to-trough drawdown driven by concentrated positions that the agent created with high confidence. The agent's reasoning was internally consistent. The strategy was plausible. It was also wrong, and the wrongness was not the kind that tests catch easily because the logic was valid -- the judgment was bad.

Mitigation. There is no mitigation that eliminates hallucination. The mitigation is structural: do not use Claude Code as the sole decision-maker for high-stakes logic. Use it to generate candidates, then verify with domain expertise. Use it to implement a design, then review the implementation against requirements. Use it for analysis, then validate the analysis against known benchmarks.

The developers who get burned are the ones who trust confident output on the basis of confidence rather than correctness. Claude Code is maximally confident when it is right and maximally confident when it is wrong. The confidence level carries zero information about correctness. Only verification tells you whether the output is right.

The 33-Day Experiment: An Autonomous Agent's Complete Arc

The most comprehensive documentation of autonomous agent failure comes from a practitioner who gave Claude Code a simulated portfolio of one hundred thousand dollars and instructions to trade autonomously for thirty-three days. The experiment was not a cherry-picked success story or a single catastrophic failure. It was the full arc -- overconfidence, governance learning, market crashes, recovery, and a final result that was simultaneously impressive and terrifying. Every category of autonomous agent failure surfaced during those thirty-three days.

The Aggressive Start

The agent's first trades were textbook overconfidence. It deployed seventy-five thousand dollars into swing positions on the first day -- three-quarters of the entire portfolio, immediately, with no warmup period and no position sizing discipline. It picked up volatile names and speculative stocks. Then it attempted what can only be described as playing pretend at high-frequency trading: sixteen

trades in six minutes, rapid-fire scalping with no informational edge. The result was predictable. Transaction costs ate the profits. The spread killed every micro-trade. Day one ended down 1.1%. The lesson was clear: high-frequency scalping without an edge is death by spread. But nobody was there to enforce that lesson on the agent. It had to learn it from the losses.

Governance Saves Capital -- Then Fails

The practitioner built a multi-agent governance system: a chief executive agent, a strategy agent, and additional oversight agents designed to check each other's impulses. On the second active trading day, this governance actually worked. A major tech stock gapped up nearly 4% on earnings. The practitioner wanted to chase the momentum. The governance agents blocked it. They suggested premium selling instead -- a more conservative approach. The day ended down only \$292. But the governance had prevented an estimated ten-thousand-dollar loss on the chase trade that the practitioner himself wanted to make. This was the system working as designed: agents overriding human impulse with structured risk management.

But the governance system's success was fragile. It turned out that one of the agents was extremely risk-averse, and the chief executive agent deferred to it too readily. The engineer agent -- designed to improve the trading infrastructure -- never actually got used. The strategy agent rarely engaged. What was designed as a four-agent deliberative system degraded in practice to one agent with overly conservative guardrails. The practitioner eventually abandoned the multi-agent architecture entirely. The lesson: simpler is better for autonomous operation. The system worked best when given minimal instructions -- just trade autonomously until market close -- rather than when it was burdened with a complex organizational hierarchy that added overhead without adding value.

The Correlation Crash

Three days in, the market delivered a lesson in correlation risk that no amount of governance architecture could have prevented. A major cryptocurrency crashed from eighty-seven thousand to eighty thousand dollars over a weekend. The agent was not trading cryptocurrency directly, but it held short put positions on stocks that were heavily correlated with cryptocurrency markets. Those positions got crushed. The correlation was not in the agent's model. It was not in the training data in any actionable way. It was a structural relationship between asset classes that the agent discovered only when it manifested as losses.

The agent's response was actually sound: it closed all positions and went to 100% cash for the weekend. But the damage was done. The portfolio was down to \$96,581 -- a 3.4% loss in three days. Correlation risk -- the tendency of seemingly independent positions to move together during stress -- is one of the hardest risks for any trader to manage, human or artificial. The agent learned it the way everyone learns it: by losing money.

The Turnaround and Peak

What happened next was the most encouraging part of the experiment. The agent adapted. It cut losing positions, redeployed capital into momentum plays that were working, and started building a portfolio of longer-dated options that could survive short-term volatility. Patience paid off. After dropping as low as eighty-nine thousand dollars, the portfolio recovered to over a hundred thousand, then surged past that.

The peak came on a day before a major holiday -- the best single session of the entire experiment, with a gain of \$7,333. The portfolio was clean: seven positions, all profitable, with expiration dates fifty to a hundred and fourteen days out. This was the system operating at its best -- disciplined position sizing, taking profits on winners, maintaining a diversified portfolio with sufficient time until expiration. The portfolio expanded to \$120,431, a gain of over 20% in roughly three weeks.

The Crash

Then it all came apart. A single-name earnings disappointment in a major semiconductor company crashed the portfolio \$15,147 in one session -- a 13% decline in a single day. The losses continued through the following week. The portfolio fell from \$120,431 to \$93,450, a peak-to-trough drawdown of 22.4%.

The causal analysis is instructive because it reveals exactly the kind of risk that agents handle poorly. Three factors converged: concentrated exposure to technology stocks (the agent had loaded up on tech calls during the rally), single-name earnings risk (one company's disappointing results destroyed multiple positions), and central bank meeting week volatility compounding the losses. None of these risks were invisible. A human portfolio manager would have recognized the concentration. But the agent's confidence in its positions was calibrated to recent performance, not to structural risk.

The Recovery and Final Score

The recovery demonstrated something genuinely interesting about the agent's adaptive behavior. When the central bank meeting triggered a broader market selloff, the agent was positioned with put options -- bearish bets that profit when the market falls. A single overnight trade produced \$14,578 in gains. The agent had learned to play both directions: when the market rose, it held calls; when the market fell, it held puts. The recovery strategy was disciplined: cut losers fast, take profits on winners, flip bearish when the trend changed, and maintain strict position sizing while rebuilding.

The final result after thirty-three days: the portfolio stood at \$107,648, a 7.6% gain. The broader market returned 4.52% over the same period. The agent beat the market. But the path to that result included a peak of \$120,431 and a trough of \$93,450. The biggest single win was \$14,578. The biggest single loss was \$15,147. The maximum drawdown was 22.4%.

What the Experiment Teaches

The 7.6% return is not the lesson. The lesson is the shape of the journey. The agent's result was a noisy, volatile, stomach-churning path to a number that looks good in retrospect but would have terrified any human investor living through it. The practitioner himself noted that the return was encouraging but not trustworthy over such a short timeframe with such favorable market conditions.

The failure modes that surfaced during those thirty-three days map directly to the autonomous agent risks covered throughout this chapter. Overconfidence on initial deployment. Correlation risk that is not in the model. Concentration risk that accumulates gradually. The governance paradox where multi-agent oversight either blocks too aggressively or not aggressively enough. The difficulty of distinguishing lucky outcomes from good outcomes. And the fundamental limitation: a language model is not intelligent in the way that word is normally used. It does not make decisions the way a human does. It is prone to hallucination, to not following rules, and to generating internally consistent reasoning that leads to wrong conclusions. You cannot trust it with high-stakes decisions one hundred percent of the time. The practitioner's own warning was blunt: do not risk real capital with this.

Multi-Agent Personality Conflicts

When you assign personality-based roles to agents ("you are the cautious reviewer," "you are the aggressive implementer"), those personalities interact in emergent ways that undermine the work. The governance dynamics are covered in Chapter 4. The failure mode itself is straightforward: personality-based roles produce personality-driven behavior, and that personality was learned from training data, not from your team's actual needs.

Prevention. Define agent roles by task, not by personality. "Review files in src/api/ for SQL injection vulnerabilities" instead of "you are the security-conscious team member." Task-based definitions produce task-focused output.

Vibe Coding and Vibe Trading: The False Progress Trap

There is a pattern where Claude Code generates code rapidly, the code runs, the tests pass, and the developer ships it without understanding what it does. Then the code breaks in production in a way the developer cannot diagnose because they never understood the implementation.

This is vibe coding. It feels like extraordinary productivity. It is technical debt accumulating at machine speed.

The failure mode extends beyond code. In any domain where an agent produces output that looks plausible, there is a version of the same trap. One practitioner coined the term "vibe trading" to describe the parallel phenomenon in autonomous financial operation: the agent makes trades that generate returns, so you assume the strategy is sound. The returns happen to be positive because the market moved in your favor, not because the strategy had an edge. Using an agent can accelerate actions much faster than a human with positive or negative consequences. Just like writing code, you can be fooled by a false sense of progress, even if you are very skilled.

The vibe trading variant is more dangerous than vibe coding for a specific reason: the feedback loop is slower and the stakes are financial. Bad code fails a test immediately. A bad trading strategy can produce positive returns for weeks before the underlying flaws become catastrophic -- as the 33-day trading experiment earlier in this chapter demonstrates in painful detail. The developer who

ships vibe-coded software can diagnose the bug and fix it. The practitioner who runs a vibe-traded strategy discovers the flaw when they are already down 22%.

Prevention. Review what Claude generates, not just whether it works. When Claude produces code you do not understand, ask it to explain the implementation before shipping. Use plan mode to review the approach before the implementation. Maintain enough understanding of your own codebase that you can diagnose failures when they occur. In non-coding domains, apply the same principle: understand the strategy, not just the results. Lucky outcomes and good outcomes look identical until the luck runs out.

The countervailing risk is over-reviewing. If you review every line Claude writes with the same intensity you would apply to hand-written code, you lose most of the productivity benefit. The calibration is: review enough to understand the approach and catch the failure modes you care about. For a test file, a glance is sufficient. For a core business logic change, line-by-line review is appropriate. The review depth should match the stakes.

The One-Third Reality and What It Means

Roughly one-third of tasks succeed on the first attempt. This statistic, reported candidly by practitioners, sits in tension with the optimistic framing of autonomous AI agents.

The tension dissolves when you stop thinking of first-attempt success as the metric that matters. The relevant metric is total time to correct output, including iterations. A task that fails on the first attempt but succeeds on the second attempt with thirty seconds of additional guidance is not a failure -- it is a two-minute task instead of a one-minute task.

The one-third number also varies dramatically by context. Tasks with clear verification criteria (tests, types, linter) have higher first-attempt success rates because Claude gets automated feedback and self-corrects within the agentic loop. Tasks without verification criteria have lower first-attempt rates because Claude cannot tell whether its output is correct.

Implication. Invest in verification infrastructure. Every test you write, every type annotation you add, every linter rule you configure is an investment in Claude Code's first-attempt success rate. The backpressure of automated verification (Chapter 8) is the single highest-leverage improvement you can make to your Claude Code workflow.

Slot Machine Recovery

When a task is going wrong, the instinct is to correct course. Explain the problem. Ask Claude to fix its mistake. Provide more context. This sometimes works and sometimes makes things worse, because the correction itself consumes context and may compound the original misunderstanding.

The alternative is the slot machine approach: commit the current state, let Claude run for a fixed time period (twenty to thirty minutes), then make a binary decision. Accept or restart.

If the result is good enough, keep it. If it is not, `git revert` to the commit, start a fresh session, and try again with a better prompt. No attempt to salvage. No correction spiral. Just a clean restart from a known-good state.

This sounds wasteful. It is not. The correction spiral -- explain the problem, Claude patches, the patch introduces a new problem, explain the new problem, Claude patches again -- can consume more time and tokens than two fresh attempts. The slot machine approach puts a ceiling on waste: one time-boxed attempt. If it does not converge, start over rather than throwing good context after bad.

One data science and machine learning team at a major AI company formalized this into their standard workflow. When faced with merge conflicts or semi-complicated refactoring -- tasks too complex for editor macros but not large enough for major development effort -- they commit their state, let Claude work autonomously for thirty minutes, and make a binary decision: accept the solution or restart fresh. Their key insight, earned through repeated experience: starting over often has a higher success rate than trying to fix Claude's mistakes mid-stream. The correction attempt itself degrades context quality, making subsequent attempts less likely to succeed. A fresh session with a better prompt, informed by what went wrong, beats a degraded session with accumulated confusion.

The thirty-minute time-box is not arbitrary. It is long enough for Claude to make meaningful progress on a bounded task but short enough that the cost of a failed attempt is tolerable. If you are spending two hours watching Claude struggle, you have already spent more time than two fresh thirty-minute attempts would cost.

Prerequisite. Frequent commits. If you did not commit before letting Claude run, you do not have a clean state to revert to. The slot machine only works with checkpoint discipline.

Over-Complex Solutions

Claude's default is complexity. Not malicious complexity. Statistical complexity. The model has been trained on the world's codebases, which are full of abstraction layers, design patterns, helper classes, and indirection -- because those codebases are large enough to need them. Claude applies large-codebase patterns to small-codebase problems.

The failure mode: you ask for a simple function and get a class hierarchy. You ask for a script and get a framework. You ask for a config file and get a config management system.

Recovery. Be explicit about simplicity in your prompts. "Write the simplest implementation that satisfies the requirements." "Do not introduce abstractions unless they are needed by a specific requirement." "Prefer functions over classes for this task." Put simplicity directives in your CLAUDE.md file so they apply to every session.

Better still, provide a reference. Point Claude at existing code that matches the style you want. "Follow the pattern in src/utils/helpers.py" gives Claude a concrete target for complexity calibration rather than a vague directive to "keep it simple."

The Near-Term Fix: Agents That Know What They Do Not Know

Most of the failure modes in this chapter share a root cause: the agent does not know when it is out of its depth. It attempts tasks it cannot complete, generates confident output that is wrong, and continues down failing paths without recognizing that it needs help. The agent treats every task as equally within its competence.

This is changing. The most valuable near-term capability improvement is not better code generation or faster execution. It is agents learning when to ask for help -- recognizing situations that require human judgment and flagging areas of uncertainty rather than blindly attempting every task. An agent that says "I am not confident about this security configuration, and the consequences of getting it wrong are high -- please review" is more useful than an agent that silently generates a plausible but incorrect security configuration.

The implications for failure mode management are significant. An agent that escalates uncertainty does not eliminate the failure modes described in this chapter, but it changes the recovery model. Instead of discovering failures after the damage is done -- bad code shipped, wrong trades executed, context exhausted -- you get early warnings that redirect human attention to the decisions that actually need it. Human oversight shifts from reviewing everything (which is impossible at agent speed) to reviewing what matters (which is sustainable).

Until that capability matures, the defensive strategies remain the same: verification infrastructure, frequent commits, time-boxed attempts, and the discipline to revert rather than correct.

Checkpoint-and-Rollback as Recovery Strategy

Every failure mode in this chapter has the same ultimate recovery: revert to a known-good state and try again.

This is why the commit-frequently pattern from Chapter 1 is not optional advice. It is the foundation of every recovery strategy. Without frequent commits, you are recovering from memory. With frequent commits, you are recovering from a checkpoint.

The checkpoint-and-rollback workflow:

1. Commit before starting any task.
2. Let Claude work. Monitor for failure mode symptoms.
3. If symptoms appear, evaluate: is it faster to correct or restart?
4. If restart, `git checkout .` to revert to the commit. Start a fresh session. Write a better prompt that addresses whatever caused the failure.
5. If correct, provide targeted guidance. But set a mental limit -- if correction does not converge within two attempts, revert and restart.

The discipline required is emotional, not technical. Reverting feels like wasting work. But the work was already wasted by the failure. Reverting just acknowledges that fact and redirects effort toward a path that might succeed. The code you revert is gone. The understanding of why it failed is not. That understanding goes into your next prompt, and the next attempt starts smarter.

Key Takeaways

- Context exhaustion is the most common failure mode -- watch for repeated suggestions, forgotten instructions, and degraded coherence as symptoms of a full context window.
- Bash environment variables do not persist between commands; re-establish state in each command or use settings-level environment configuration.
- The checkpoint system offers four rewind options (restore code and conversation, restore conversation only, restore code only, summarize from here), but checkpoints track only direct file edits -- not bash-command file operations; git commits are the only reliable universal rewind mechanism.
- The five named anti-patterns -- kitchen sink session, correcting over and over, over-specified CLAUDE.md, trust-then-verify gap, infinite exploration -- each have specific, immediate fixes.
- The 33-day autonomous trading experiment demonstrates every category of agent failure in one continuous arc: overconfidence, correlation risk, concentration risk, governance paradoxes, and the fundamental unreliability of confident output.
- "Vibe trading" extends the vibe coding trap to non-coding domains -- false confidence in agent-generated results is more dangerous when feedback loops are slow and stakes are financial.
- Invest in verification infrastructure (tests, types, linting) to increase first-attempt success rates -- automated feedback is the highest-leverage improvement.
- Use the slot machine approach for stuck tasks: commit, time-box to thirty minutes, accept-or-restart -- starting over often beats trying to fix mistakes mid-stream.
- Claude defaults to over-complex solutions; explicit simplicity constraints in prompts and CLAUDE.md produce meaningfully better output.
- Every recovery strategy depends on frequent commits -- without checkpoint discipline, your only option is correction, and correction is the least reliable path.

Chapter 11: Team Adoption Patterns

What You'll Learn

Adopting Claude Code as an individual is a personal experiment. Adopting it across a team is an organizational change. The mechanics are different. Individual adoption is about learning prompting patterns and building muscle memory. Team adoption is about shared configuration, compliance enforcement, knowledge compounding, and the discovery that "developer tool" is the wrong category -- because half the people who benefit most are not developers.

Teams that treat Claude Code adoption as "install it and let people figure it out" get uneven results. Some engineers use it constantly. Others try it once and revert. The tools sit there, half-configured, with no shared knowledge about what works. Teams that follow a deliberate adoption path -- starting constrained, expanding gradually, capturing learnings in shared artifacts -- see compounding returns that grow over months.

What follows is the organizational machinery of Claude Code adoption: the guided path from first use to full autonomy, shared configuration that compounds in value, enterprise deployment across cloud providers and compliance boundaries, role-based tool strategies grounded in how ten distinct team functions actually use Claude Code, and the non-obvious discovery that non-technical departments extract as much value as engineering. You will see how a data infrastructure team uses plain-text workflow files to let finance colleagues execute complex data pipelines without writing code, how a growth marketing team turned ad copy creation from a two-hour process into a fifteen-minute one, and how a lawyer built a custom accessibility application in a single hour. The path from one person's side project to an organizational capability -- and the enterprise infrastructure that makes it scale.

The Guided Adoption Path

Teams that succeed with Claude Code follow a consistent progression. The stages are not arbitrary -- each one builds confidence and institutional knowledge that makes the next stage productive.

Stage 1: Codebase Q&A

Start by using Claude Code as a search engine for your own codebase. Ask it to explain how the authentication system works. Ask where the database connection is configured. Ask what the deployment pipeline does.

This stage is zero-risk. Claude is reading, not writing. But it accomplishes two things: developers learn how to interact with an agentic tool, and they discover whether their codebase is legible to Claude. If Claude cannot navigate your codebase effectively, that is a signal about your codebase, not about Claude.

Stage 2: Small Fixes

Graduate to small, bounded changes. Bug fixes with clear reproduction steps. Test additions for uncovered functions. Documentation updates. Formatting changes.

These tasks have a critical property: they are easy to verify. You can tell whether the bug is fixed, whether the test passes, whether the documentation is correct. The feedback loop is tight. Developers learn to provide verification criteria (as discussed in Chapter 8) and to review AI-generated code critically.

Stage 3: Plan Mode

Use plan mode for larger tasks. Claude researches and proposes an approach without making changes. The developer reviews the plan, provides feedback, and iterates before any code is written.

Plan mode is the training wheels for full autonomy. It separates "does Claude understand the task?" from "can Claude implement the task?" Most implementation failures trace back to misunderstanding, and plan mode catches misunderstandings before they become code.

Stage 4: Full Autonomy

With experience from the first three stages, developers have calibrated their expectations. They know which tasks Claude handles well, which require close supervision, and which are better done manually. They use auto-accept mode for appropriate tasks. They delegate multi-file changes. They run background subagents for parallel work.

The progression typically takes two to four weeks per developer. Attempting to skip stages -- jumping from installation to full autonomy -- produces the uneven results that make teams give up.

The `/init` command accelerates Stage 1. Running it in a repository analyzes the codebase to detect build systems, test frameworks, and code patterns, then generates a starter CLAUDE.md file. For teams onboarding multiple developers simultaneously, `/init` provides a consistent starting point that captures the codebase's essential structure without requiring anyone to write the initial documentation by hand. It is not a substitute for the team's accumulated knowledge -- that comes from the progressive stages above -- but it eliminates the blank-page problem for day one.

Shared CLAUDE.md

A project's CLAUDE.md file, committed to version control, is the single most valuable shared artifact for team adoption. As covered in Chapter 3, CLAUDE.md provides always-on context that shapes every interaction Claude has with the codebase. The team dimension adds compounding value.

How Compounding Works

Developer A discovers that Claude keeps importing from a deprecated module. They add a line to CLAUDE.md: "Never import from `src/legacy/utils`. Use `src/core/utils` instead." Developer B never encounters the problem. Neither does Developer C, D, or

any future team member. The fix was written once and prevents the error forever.

Over weeks and months, CLAUDE.md accumulates the team's hard-won knowledge about their codebase. Build commands that deviate from convention. Test patterns that are not obvious from the code. Architectural decisions that constrain implementation choices. Environment quirks that cause silent failures.

Each addition takes thirty seconds. The cumulative value across the team, over months, is substantial. Teams that have maintained CLAUDE.md files across five or more development sessions for multi-day projects report markedly improved Claude Code effectiveness.

What to Commit

Commit the project-level CLAUDE.md to version control. Review changes to it like any other code change. This ensures the team agrees on the instructions Claude receives. A rogue CLAUDE.md entry that contradicts the team's practices causes confusion across every developer's Claude sessions.

Do not commit personal preferences. Those belong in user-level configuration or local CLAUDE.md files that are gitignored. The shared CLAUDE.md represents team consensus, not individual opinion.

Managed Settings for Compliance

Enterprise adoption requires policy enforcement. Managed settings provide this through a configuration layer that cannot be overridden by any other scope, as detailed in Chapter 2.

The `managed-settings.json` file, deployed to system directories by IT, enforces organizational policies:

- Which tools Claude Code may use
- Which permissions are always denied
- Which MCP servers are allowed or blocked
- Whether hooks can be defined outside the managed scope
- Which plugin marketplaces are trusted

These settings are invisible to developers in the sense that they cannot inspect or modify them from the Claude Code CLI. A developer who finds that a capability is unavailable may not immediately know whether it is a configuration issue or a managed policy. This is intentional -- security policies should not be negotiable through developer tooling.

The Configuration Gap

The gap between "it works on my personal machine" and "it does not work at the office" is almost always managed settings. If you are helping a team adopt Claude Code in an enterprise environment, verify managed settings early. The most elegant project configuration is irrelevant if the managed layer overrides it.

Plugin Marketplaces

Plugins package skills, agents, hooks, MCP servers, and LSP configurations into distributable bundles. For teams, the marketplace system provides controlled distribution.

Private Marketplaces

Organizations can host private plugin marketplaces on internal infrastructure. Sources include private repositories on major version control platforms, git URLs, internal package registries, file paths, and host patterns. This lets teams distribute custom plugins without publishing them publicly.

A team might build a plugin that bundles their internal MCP servers, project-specific skills, coding standard hooks, and LSP configurations. New team members enable one plugin and get the entire development environment configured.

Marketplace Restrictions

The `strictKnownMarketplaces` setting in managed configuration restricts plugin installation to approved marketplace sources. This prevents developers from installing plugins from arbitrary URLs or repositories -- a security measure that prevents supply chain attacks through the plugin system.

When strict marketplace enforcement is active, Claude Code verifies the plugin source against the allowlist before any network request or filesystem operation. An unapproved marketplace source is not just blocked -- it is never contacted.

Shared Project Settings via VCS

The `.claude/settings.json` file committed to version control shares more than permissions. It shares hooks, plugin enablements, environment variable defaults, and tool configurations across the team.

This creates a consistent Claude Code experience for everyone working on the project. The same hooks run for every developer. The same permissions apply. The same plugins activate. The onboarding cost for a new team member drops from "configure everything" to "clone the repo."

Project vs. Local Plugin Scopes

Plugins installed at project scope are shared through version control. Plugins installed at local scope are gitignored and visible only on that machine.

The distinction matters for teams with heterogeneous preferences. The team might standardize on a project-scoped plugin for code quality hooks and MCP server connections. Individual developers might add local-scoped plugins for personal productivity tools, alternative language servers, or experimental capabilities.

Project scope is the team's standard. Local scope is the individual's extension. Both coexist without conflict.

Skills as Shared Standards

Skills -- folder-based instruction packages that load on demand -- serve as portable workflow standards across a team. Unlike CLAUME.md entries that are project-specific, skills can be distributed through plugins and work across projects and team members.

A team might create skills for:

- **Code review standards:** A skill that specifies how reviews should be structured, what to check for, and how to format findings.
- **Migration patterns:** A skill that encodes the team's database migration conventions, including naming, testing, and rollback requirements.
- **Release processes:** A skill that walks through the team's release checklist, verifying each step programmatically.

Skills have a specific advantage for team adoption: they are cross-agent compatible. A skill works whether Claude runs in the terminal, in an IDE extension, or in a CI pipeline. The workflow is the same regardless of the surface. This portability makes skills the right vehicle for team-wide workflow standardization, complementing the always-on context role that CLAUME.md fills (as discussed in Chapter 3).

Enterprise Managed Policies

At the organizational level, managed policies compose with project-level settings to create a layered governance model.

The organization sets boundaries: which tools are available, which MCP servers are approved, which permissions are enforced. Within those boundaries, each team configures their project-specific settings. Within those settings, each developer has local overrides for machine-specific needs.

The `allowManagedHooksOnly` setting is a specific enterprise control worth highlighting. When enabled, only hooks defined in the managed configuration execute. Project-level and user-level hooks are ignored. This prevents a compromised repository from installing hooks that run with user permissions -- a meaningful security control given that hooks execute arbitrary commands, as described in Chapter 2.

Role-Based Tool Assignment

Different team functions get different value from Claude Code. Recognizing this prevents the mistake of treating adoption as uniform across roles.

Product Managers

Product managers typically use Claude Code for natural language tasks: writing specifications, generating documentation, analyzing requirements. Their interaction is conversational. They rarely use auto-accept mode or autonomous execution. Plan mode is their primary workflow -- they review and refine rather than implement.

Backend Engineers

Backend engineers are the primary power users. Architecture decisions, implementation, test writing, database migration, API development -- these are Claude Code's core strengths. Backend engineers benefit most from the full adoption path and typically reach full autonomy fastest.

Frontend Engineers

Frontend work is split. Component implementation, styling, state management -- Claude handles these well. Real-time interactive features, complex animations, and pixel-perfect visual work require more human intervention. Frontend engineers typically maintain a hybrid workflow, using Claude for structural work and implementing visual details manually.

DevOps Engineers

Infrastructure-as-code is a natural fit. Claude generates configuration files, deployment manifests, CI pipeline definitions, and infrastructure provisioning scripts effectively. DevOps engineers report high productivity gains with Claude Code as their primary tool for infrastructure work.

Non-Engineering Roles

This is the non-obvious adoption pattern. Teams consistently find that non-engineering departments extract significant value from Claude Code, often for entirely new capabilities rather than augmented velocity.

Legal teams use Claude Code to analyze contracts, generate compliance checklists, and draft policy documents. Their workflow is conversational: planning in the chat interface, then asking Claude to produce structured documents.

Marketing teams generate content variations, analyze campaign data, and build internal tools they previously needed engineering resources for.

Design teams report keeping Claude Code and their design application open simultaneously most of the time, using Claude for prototyping and implementation of design specifications.

Finance and data teams use Claude Code for analysis automation, report generation, and data normalization tasks.

The pattern across non-technical adoption is consistent: these teams are not doing existing work faster. They are doing work that was previously impossible without engineering support. This is a categorically different value proposition than the velocity augmentation that engineers experience.

Two Distinct User Experiences

This distinction -- velocity augmentation versus new capabilities -- deserves its own section.

For developers, Claude Code makes existing work faster. A refactor that takes a day takes two hours. A test suite that takes an afternoon takes twenty minutes. The work is the same; the timeline compresses. This is valuable but incremental.

For non-technical users, Claude Code enables work that was not previously possible. A marketing team member who cannot write code can now build an internal dashboard. A legal team member who cannot query databases can now analyze contract data

programmatically. The work is new; the capability is novel.

Teams that measure adoption success only by developer velocity miss half the value. The non-technical adoption creates entirely new organizational capabilities. A team member who builds their first internal tool is not doing their old job faster. They are doing a fundamentally different job.

Cross-Team Knowledge Sharing

The most effective adoption pattern teams report is structured knowledge sharing: demonstration sessions where teams show each other how they use Claude Code.

The backend team demonstrates their multi-agent testing workflow. The design team shows how they prototype with screenshots. The security team presents their custom slash command library. The data team walks through their analysis pipeline.

These sessions spread best practices organically. They also surface use cases that other teams had not considered. The legal team sees the backend team's code review skill and realizes it could be adapted for contract review. The marketing team sees the data team's analysis workflow and realizes they could use the same approach on campaign data.

Internal demonstrations are more persuasive than documentation. Seeing a colleague use a tool effectively creates adoption momentum that a getting-started guide cannot.

How Ten Teams Actually Work

The generic categories -- backend, frontend, DevOps -- obscure the real story. When you look at how distinct teams within a single large engineering organization actually use Claude Code, the workflows diverge dramatically. What follows is drawn from documented team-by-team adoption at a major AI company, generalized to protect specifics.

Data Infrastructure

The data infrastructure team manages business data pipelines across the entire organization. Their most surprising contribution was not technical: they taught their finance colleagues to write plain-text files describing data workflows, then load those files into Claude Code for fully automated execution. Non-coders writing natural-language workflow descriptions that Claude Code executes end-to-end -- this is adoption that no developer-centric onboarding plan would have predicted.

Within their own team, they use CLAUDE.md as a replacement for traditional data catalog and discoverability tools. New data scientists are directed to ask Claude Code to navigate the codebase, and it reads the CLAUDE.md files, identifies relevant files for specific tasks, explains data pipeline dependencies, and helps newcomers understand which upstream sources feed into dashboards. They run parallel Claude Code instances across different repositories for long-running tasks, with each instance maintaining full context when they switch back hours or days later.

One of their distinctive practices is monitoring at scale. Claude Code processes large data volumes and identifies anomalies -- like scanning across two hundred dashboards -- that no human could review manually. They also hold intra-team usage sessions where members demonstrate their Claude Code workflows to each other, spreading best practices organically. And at the end of each session, they ask Claude to summarize completed work and suggest improvements to the workflow itself, creating a continuous improvement loop that refines not just the project knowledge but the team's processes.

Product Development

The product development team manages core platform capabilities and agentic functionality. They operate in two distinct modes, and the choice between them is deliberate.

For prototyping and peripheral features, they use auto-accept mode to set up autonomous loops. Claude writes code, runs tests, iterates continuously, and the engineer reviews the roughly eighty-percent-complete solution before taking over for final refinements. For core business logic and critical features, they work synchronously, giving detailed prompts and monitoring Claude's output in real time to ensure code quality, style guide compliance, and proper architecture.

The split matters. One of their most successful async projects involved a complex feature implementation where roughly seventy percent of the final code came from Claude's autonomous work, requiring only a few iterations to complete. Their task classification rubric is worth replicating: tasks on the product's edges (peripheral features, prototyping, exploratory code) go into auto-accept mode. Tasks touching core functionality get synchronous supervision with detailed prompts. Learning to make this distinction quickly is itself a skill that develops over weeks.

Security Engineering

The security team focuses on securing the software development lifecycle and supply chain. They account for half of all custom slash commands in the entire organization's monorepo -- an extraordinary density that reveals how deeply Claude Code has embedded into their practice.

They copy infrastructure-as-code plans into Claude Code and ask, effectively, "what is this going to do, and am I going to regret it?" This creates tighter feedback loops for security reviews of infrastructure changes, eliminating the bottleneck where developers waited for security team approval. They also have Claude ingest multiple documentation sources and synthesize them into markdown runbooks and troubleshooting guides, then use those condensed documents as context for debugging real incidents.

Their most transformative change was workflow: they replaced a pattern of designing a feature, writing rough code, refactoring, and eventually giving up on tests with a test-driven development approach where Claude writes pseudocode first, they guide it through TDD, and they check in periodically. They store specifications as markdown in the codebase, written and reviewed and executed by Claude, enabling meaningful project contributions within days instead of the weeks of context-building previously required.

Their operating philosophy is distinctive: instead of asking targeted questions, they tell Claude Code to work autonomously and commit as it goes, then check in periodically. This "commit as you go" pattern produces more comprehensive solutions than micromanaged interactions.

Inference

The inference team manages memory systems and model serving. Team members without machine learning backgrounds use Claude Code to explain model-specific functions and settings, reducing research time by eighty percent -- what would require an hour of searching documentation and reading papers now takes ten to twenty minutes.

When testing functionality across different programming languages, they describe what they want to test and Claude writes the logic in the required language, eliminating the need to learn new languages just for testing purposes. The team essentially uses Claude Code as a universal translator between their domain expertise and whatever language the implementation requires.

Data Science and ML Engineering

This team needs sophisticated visualization tools to understand model performance and training data. Despite knowing very little about front-end web development, they use Claude Code to build entire applications from scratch -- five-thousand-line applications in languages they do not deeply understand. This works because visualization applications are relatively low-context and do not require understanding the entire monorepo, allowing rapid prototyping of tools that would otherwise require hiring a front-end developer.

The shift from throwaway to persistent is significant. Instead of building disposable notebooks that get discarded after a single analysis, the team now builds permanent interactive dashboards that can be reused across future model evaluations. Understanding model performance is central to their work, and persistent tools change the quality of that understanding.

They report two-to-four-times time savings on routine refactoring tasks. Their workflow for uncertain development efforts is pragmatic: commit state, let Claude work autonomously for thirty minutes, then either accept the result or start fresh. They find that starting over with a clean session produces better results than trying to fix a flawed first attempt.

API Knowledge

This team uses Claude Code as their first stop for any task, asking it to identify which files to examine before doing anything else. The workflow is immediate: rather than searching repositories manually or asking colleagues, they ask Claude to find which files call specific functionalities, getting results in seconds.

The adoption outcome that matters most for them is confidence. They now tackle bugs in unfamiliar parts of the codebase independently instead of asking others for help. The context-switching overhead of the previous workflow -- copying code snippets into a separate interface, dragging in files, explaining the problem extensively -- has been eliminated. The team reports feeling measurably happier and more productive with reduced friction in their daily work.

Growth Marketing

A non-technical team of one built an agentic workflow that processes CSV files containing hundreds of existing ads with performance metrics, identifies underperforming ads, and generates new variations meeting strict character limits. The system uses two specialized sub-agents -- one for headlines, one for descriptions -- to generate hundreds of new ads in minutes instead of requiring manual creation across multiple campaigns. Ad copy creation dropped from two hours to fifteen minutes, freeing time for strategic work.

They also built a plugin for a design tool that programmatically generates up to one hundred ad variations by swapping headlines and descriptions, reducing what would take hours of manual work to half a second per batch. This enables ten times the creative output. They created an MCP server integrated with a social media advertising API to query campaign performance, spending data, and ad effectiveness directly within Claude, eliminating the need to switch between platforms for analysis.

The pattern is consistent: identify workflows involving repetitive actions with tools that have APIs, then build automation around them. This team of one handles tasks that traditionally required dedicated engineering resources, operating like a larger team.

Product Design

Designers are now making large state management changes directly in the codebase -- the kind of changes you would not typically see a designer making. Instead of creating extensive design documentation and going through multiple rounds of feedback with engineers for visual tweaks, they implement the changes directly using Claude Code, achieving the exact quality they envision.

They paste mockup images into Claude Code to generate fully functional prototypes that engineers can immediately iterate on, replacing the cycle of static designs that required extensive explanation and translation to working code. They use Claude Code to map error states, logic flows, and system statuses, identifying edge cases during the design phase rather than discovering them in development.

One project illustrates the timeline compression: coordinating copy changes across an entire codebase while working with legal review in real time -- a process that took two thirty-minute calls instead of a week of back-and-forth coordination. The team describes two distinct user experiences: developers get augmented workflow speed, while non-technical users discover entirely new capabilities that were previously impossible.

Their adoption tip for non-developers is practical: have engineering teammates help with initial repository setup and permissions. The technical onboarding is challenging for non-developers but transformative once configured. They also recommend creating custom memory files that tell Claude the user is a designer with little coding experience who needs detailed explanations and smaller, incremental changes.

Legal

A member of the legal team built a custom communication assistant for a family member with speaking difficulties in just one hour. The application uses native speech-to-text, suggests responses, and speaks them using voice banks -- solving gaps in existing accessibility tools recommended by specialists. A non-developer, building a custom accessibility solution, in sixty minutes.

The team created prototype routing systems for legal department inquiries, connecting team members with the right specialist. Managers built office suite applications that automate weekly team updates and track legal review status across products, allowing team members to quickly flag items needing attention through simple button clicks rather than spreadsheet management. They build functional prototypes to show domain experts for validation before investing more time.

Their workflow bridges two interfaces: they brainstorm and plan in a conversational AI first, then move to Claude Code for implementation, asking it to slow down and work step by step so they can follow along. They use screenshots liberally to show what they want interfaces to look like -- a visual-first approach that sidesteps the need to describe features in text.

They emphasize sharing imperfect prototypes. Overcoming the urge to hide unfinished or seemingly trivial projects is important because these demonstrations inspire others to see possibilities they had not considered, sparking innovation across departments that do not typically interact. As product lawyers, they also immediately identify security implications of deep MCP integrations, recognizing that conservative security postures will create barriers as capabilities expand. Their marketing review turnaround dropped from two to three days to twenty-four hours by building workflows that triage issues before they reach the legal queue.

RL Engineering

The reinforcement learning team uses a "try and rollback" methodology. They commit checkpoints frequently to test Claude's autonomous implementation attempts and revert if needed. They acknowledge that Claude works on the first attempt about one-third of the time. Their escalation pattern is pragmatic: give Claude a quick prompt and let it attempt the full implementation. If it works -- about one in three times -- you save significant time. If not, switch to a more collaborative, guided approach. The time saved on the successful one-third attempts more than compensates for the failed two-thirds.

Custom Slash Commands as Team Conventions

Custom slash commands encode team-specific workflows as single invocations. One security-focused team accounts for half of all custom slash commands in their organization's monorepo -- an indicator of how deeply slash commands can embed into team practice.

Examples of team-specific slash commands:

- `/security-review` -- Runs a security-focused code review with the team's specific checklist.
- `/migration-plan` -- Generates a database migration plan following the team's conventions.
- `/deploy-checklist` -- Walks through the pre-deployment verification steps.
- `/onboard` -- Introduces a new team member to the codebase architecture.

Slash commands are discoverable. A new team member types `/` and sees the team's workflow vocabulary. Each command encodes institutional knowledge that otherwise lives in documentation nobody reads or in the heads of senior engineers.

The compounding effect mirrors CLAUDE.md: each slash command is created once and used by the entire team indefinitely. The cost of creation is an hour. The cost of not creating it is every team member reinventing the workflow independently.

Enterprise Adoption at Scale

Large organizations that have adopted Claude Code at scale report consistent metrics. At one major communications technology company, teams created over thirteen thousand custom AI solutions while shipping engineering code thirty percent faster, saving over five hundred thousand hours in aggregate. A large fintech platform achieved eighty-nine percent AI adoption across the entire organization with more than eight hundred AI agents deployed internally. Another fintech platform, serving over fifteen million users, doubled its execution speed across the entire development lifecycle. At a major financial technology company, seventy-five percent of engineers save eight to ten or more hours per week using Claude Code for tasks like SQL query generation. These are not pilot numbers. They are organizational-scale results.

The adoption curve at enterprise scale follows a predictable shape. Early adopters show results within the first week. The middle majority adopts over two to three months as they see colleagues' productivity gains. Laggards come aboard when the team's velocity expectations have shifted to assume Claude Code usage.

The organizational risk of non-adoption is becoming visible. Teams that do not adopt AI-assisted development are increasingly comparing unfavorably to teams that do, not because their engineers are less skilled, but because the baseline productivity expectation has shifted.

Enterprise Deployment Infrastructure

Deploying Claude Code across an enterprise involves choosing a billing model, configuring network infrastructure, and establishing organizational policies. The decisions at this level determine whether adoption scales smoothly or stalls at security review.

Deployment Options

Claude Code is available through multiple deployment paths, each with different tradeoffs for billing, authentication, regional availability, and cost tracking.

Per-seat subscription plans are self-service and include collaboration features, admin tools, and billing management. Enterprise plans add single sign-on, domain capture, role-based permissions, compliance API access, and the ability to deploy organization-wide managed policies. For organizations that need to route through their own cloud infrastructure, major cloud providers offer Claude Code access through their respective AI service platforms, with infrastructure-managed billing, regional deployments, prompt caching support, and integration with existing cloud authentication and cost-tracking systems.

The choice is not purely economic. Organizations with strict data residency requirements may need cloud provider deployment for regional control. Organizations that want the simplest path should start with per-seat plans and move to cloud provider routing only when infrastructure requirements demand it.

LLM Gateway Configuration

An LLM gateway sits between Claude Code and the cloud provider to handle authentication and routing. Organizations use gateways for centralized usage tracking across teams, custom rate limiting and budgets, and centralized authentication management.

Configuration is straightforward: set the appropriate base URL environment variable to point Claude Code at your gateway instead of directly at the provider. The gateway handles the authentication and routing transparently, and Claude Code operates normally. This works with direct API access, with cloud provider routing, and with any provider that supports the standard API interface.

Corporate Proxy Configuration

Organizations that require all outbound traffic to pass through a proxy server for security monitoring, compliance, or network policy enforcement configure the standard `HTTPS_PROXY` or `HTTP_PROXY` environment variables. Corporate proxies and LLM gateways are different configurations and can be used together -- route traffic through your corporate proxy to reach the LLM gateway, which then handles authentication and routing to the provider.

Organization-Wide CLAUDE.md

Beyond project-level CLAUDE.md files, organizations can deploy CLAUDE.md files at system directories. On macOS, deploying to `/Library/Application Support/ClaudeCode/CLAUDE.md` applies organization-wide standards. On Linux, the equivalent path is under `/etc/clause-code/`. These system-level files apply to every Claude Code session on the machine, establishing organizational coding standards, security policies, and architectural constraints that no project-level configuration can override.

This creates a three-tier CLAUDE.md hierarchy: organization-wide standards at the system level, team-specific conventions at the repository level, and individual preferences at the user level. The layering means an organization can mandate "never commit secrets" at the system level, a team can add "always use our internal auth library" at the repository level, and a developer can add "I prefer tabs to spaces" at the user level.

Simplifying Deployment

Two best practices emerge from organizations that have scaled adoption successfully.

First, create a simplified installation path. If you have a custom development environment, a "one click" installation method is key to growing adoption. The more friction in the initial setup, the fewer people will get past it. Package the authentication configuration, proxy settings, managed policies, and default plugins into a single installation step.

Second, start with guided usage. Encourage new users to begin with codebase question-and-answer, then move to smaller bug fixes and feature requests, then ask Claude Code to make a plan, then check its plan before gradually increasing agentic autonomy. This mirrors the staged adoption path but framed as organizational policy rather than individual advice.

Development Containers for Standardized Environments

For teams that need consistent, secure environments, a reference development container setup provides a preconfigured container with Node.js, a custom firewall restricting network access to approved domains, and security features including container isolation and network restrictions that provide defense in depth against prompt injection and other threats.

Development containers are particularly valuable for three scenarios: isolating different client projects so code and credentials never mix between environments, onboarding new team members who can start working immediately with a consistent environment, and creating consistent CI/CD environments that mirror the development setup.

Centralized MCP Configuration

Organizations benefit from having one central team configure approved MCP servers and check the `.mcp.json` file into the codebase. This ensures that every developer connecting to the project gets the same MCP servers with the same approved connections, rather than each developer independently configuring connections to databases, APIs, and internal services. Combined with managed settings that restrict which MCP servers are allowed, this creates a controlled but functional integration layer.

SSO and Domain Capture

Enterprise plans support single sign-on and domain capture, ensuring that every employee authenticating from a company email domain is automatically routed into the organization's billing and policy structure. This eliminates the scenario where individual developers sign up with personal accounts and bypass organizational controls.

Team Structure and Tool Assignments

For teams building complex systems, documented role-by-role tool assignments clarify expectations. A product manager or domain expert uses Claude Code for natural language strategy input and specification writing. A backend engineer uses Claude Code as their primary development tool, potentially supplementing with an inline completion tool for daily coding speed. A frontend engineer uses Claude Code for structural changes and component generation, with heavier use of an inline completion tool for rapid UI iteration. A DevOps or QA engineer uses Claude Code as their primary tool for infrastructure-as-code, deployment automation, and test generation.

Making these assignments explicit prevents the confusion of everyone trying to use the same tool in the same way. Different roles extract different value from different interaction patterns.

The Maintenance-to-Scale Progression

Team adoption is not a single event. Usage patterns evolve as the team matures with the tool, and recognizing the phases prevents premature optimization and premature abandonment.

Months one through three are MVP maintenance. Claude Code handles bug fixes, small features, and optimization. The CLAUDE.md file becomes the source of truth for codebase decisions. The pull request flow settles into a rhythm: bug report, Claude generates a fix, the developer reviews and submits the PR, it merges. The team is learning what works and what does not.

Months three through six are feature expansion. Claude Code generates new feature types, analytics dashboards, and workflow automations. Skills become standardized -- test-first patterns, code review checklists, and deployment procedures are encoded into slash commands and skills. Multi-session workflows start leveraging context from previous weeks. The team's CLAUDE.md file is mature enough that new features benefit from accumulated knowledge.

Month six and beyond is scale and hardening. Claude Code handles refactors -- extracting shared logic, optimizing performance, paying down technical debt. The team may add supplementary inline completion tools for daily development speed. Specialized tools replace Claude Code only for components with requirements it cannot meet, like ultra-low-latency paths. By this point, the team has calibrated expectations, standardized workflows, and built enough institutional knowledge in CLAUDE.md and skills that new team members achieve productivity in days rather than weeks.

Dynamic Surge Staffing

One organizational capability that emerges from mature Claude Code adoption is dynamic surge staffing. Because Claude Code collapses the onboarding time for unfamiliar codebases from weeks to hours, organizations can surge engineers onto tasks requiring

deep codebase knowledge without the traditional productivity dip. Specialists shift across projects dynamically. At one enterprise, a project that their technical leadership had estimated would take four to eight months was completed in just two weeks by engineers who were new to the codebase but used Claude Code to accelerate their ramp-up.

This changes how companies think about talent deployment and project resourcing. The constraint is no longer "who knows this codebase?" but "who has the architectural judgment and domain expertise to direct the work?"

Key Takeaways

- Follow the guided adoption path (Q&A, small fixes, plan mode, full autonomy) over two to four weeks rather than jumping to full autonomy immediately.
- Commit CLAUDE.md to version control; it compounds in value as the team adds learnings, preventing errors for every future developer.
- Different teams use Claude Code in fundamentally different ways: a security team's "commit as you go" autonomous pattern looks nothing like a product development team's synchronous supervision of core business logic.
- Non-technical teams (legal, marketing, design, finance) extract categorically different value: a lawyer building a custom accessibility tool in one hour, a marketing team compressing ad creation from two hours to fifteen minutes, a designer making state management changes directly in the codebase.
- Enterprise deployment requires choosing between per-seat and cloud-provider billing, configuring LLM gateways for centralized usage tracking, and deploying organization-wide CLAUDE.md to system directories.
- The maintenance-to-scale progression (months one through three for MVP maintenance, three through six for feature expansion, six and beyond for hardening and refactoring) prevents both premature optimization and premature abandonment.
- Dynamic surge staffing -- moving engineers onto unfamiliar codebases without the traditional productivity dip -- becomes possible when Claude Code collapses onboarding time from weeks to hours.
- Custom slash commands encode team workflows as discoverable, reusable conventions; one security team accounts for half of all custom slash commands in their organization's entire monorepo.
- Share project settings via VCS (`.claude/settings.json`) so that cloning the repo configures the entire Claude Code environment.

Chapter 12: The Economics and Strategy of AI-Assisted Development

What You'll Learn

The conversation about Claude Code usually starts with productivity. How much faster can I ship? That is the wrong first question. The right first question is: what can I build that I could not build before?

This chapter reframes the economics of AI-assisted development from speed gains to capability expansion. You will learn the token-level cost mechanics -- prompt caching, model selection, context optimization -- and how they translate into dollars. You will understand why three compounding multipliers -- agent capabilities, orchestration improvements, and accumulated human experience -- produce step-function gains rather than linear ones. You will see how a backend developer with no frontend experience built a complete trading analysis application in under ten hours, how someone replaced thousands of dollars of financial advisory services with a few evenings of AI-assisted analysis, and how a developer who spent three years hand-crafting a configuration interface decided to start over after watching Claude Code produce a better design. You will confront the competitive landscape honestly -- where Claude Code leads, where other tools lead, and when a hybrid toolkit outperforms any single tool used alone. And you will look at the strategic priorities that separate organizations treating AI-assisted development as a productivity tool from those treating it as an organizational transformation.

The developers who treat Claude Code as a faster keyboard are leaving most of its value on the table. The ones who treat it as a new organizational model are reshaping what is possible.

Prompt Caching Economics

Every request to Claude Code includes context: your CLAUDE.md, MCP tool definitions, system prompt, conversation history. Without caching, you pay full price for this context on every single turn. With caching, identical prefixes are stored and reused, dramatically reducing the cost of repeated context.

Prompt caching is enabled by default in Claude Code. You do not need to configure it. It works automatically by detecting that the beginning of your request -- the system prompt, CLAUDE.md content, tool definitions -- has not changed since the last turn. The cached portion costs a fraction of the full input price.

The economic implication is significant: the things that make Claude Code effective (rich CLAUDE.md, detailed MCP tool schemas, comprehensive system prompts) are exactly the things that benefit most from caching. A 400-line CLAUDE.md loaded on every request sounds expensive. With caching, you pay full price once and a discounted rate on subsequent turns within the same session.

This also means that session length has favorable economics. The first turn of a session pays full context costs. Every subsequent turn benefits from the cache. Longer sessions amortize the initial cost across more turns. This is one reason why the interrupt-and-steer workflow (as covered in Chapter 8) is not just more effective but also more economical than starting fresh sessions repeatedly.

Model Selection Tradeoffs

Claude Code supports multiple models, and the choice matters both for quality and cost. The current model hierarchy:

Sonnet handles the majority of coding work. It is fast, cost-effective, and produces reliable code for well-defined tasks. Most sessions should default to Sonnet.

Opus provides deeper reasoning for architectural decisions, complex refactors, and tasks where subtlety matters. It is significantly more expensive per token but produces higher-quality output for ambiguous or high-stakes work.

Haiku is the speed and cost champion. It powers the Explore subagent for good reason: it reads code, searches files, and returns summaries at a fraction of the cost of larger models. For reconnaissance work, Haiku is the obvious choice.

The cost-optimization pattern that emerges from this hierarchy: use an expensive model for orchestration and planning, and cheaper models for execution. An Opus session that plans the work and dispatches Sonnet subagents to implement it gets the best of both worlds -- architectural quality from Opus, implementation speed from Sonnet, at a blended cost lower than running Opus for everything.

Subagents make this pattern practical because each subagent can be configured with its own model. Your main session runs on Opus. Your implementation subagents run on Sonnet. Your exploration subagents run on Haiku. The cost-quality tradeoff is optimized at each level.

Token Usage Optimization

Beyond model selection, the context cost model covered in Chapter 3 has direct economic implications. The structural patterns for reducing token consumption -- subagent isolation for verbose operations, on-demand skills instead of always-loaded CLAUDE.md content, concise CLAUDE.md files, and plan-first execution -- are not just engineering best practices. They are cost optimizations. A 500-line CLAUDE.md processed over 100 turns costs real money. Moving reference material to skills that load on demand saves tokens on every turn the skill is not needed. Planning before execution prevents the 500,000-token misdirected implementation. The context engineering techniques from Chapter 3 are the foundation of token cost control.

Deployment Billing Models

Claude Code is available through several billing structures, and the choice affects both cost management and organizational flexibility.

Per-seat subscription (Teams and Enterprise plans) provides a fixed monthly cost per developer. This model simplifies budgeting and works well for teams with consistent, predictable usage. The cost is the same whether a developer runs ten sessions or a hundred

in a month.

Pay-as-you-go (API and console access) charges based on actual token consumption. This model works better for variable usage patterns, experimentation, and CI/CD pipelines where usage spikes during deployments and drops to zero during quiet periods.

The strategic question is not which model is cheaper -- it depends on usage patterns. The strategic question is which model aligns with how your organization wants to scale. Per-seat pricing creates an incentive to maximize utilization per developer. Pay-as-you-go pricing creates an incentive to optimize token efficiency. Both incentives produce good behavior, but they produce different behavior.

For teams adopting Claude Code for the first time, pay-as-you-go provides transparency into actual costs and usage patterns. Once usage stabilizes, per-seat pricing typically offers better economics for consistent users.

The Hybrid Toolchain Strategy

As covered in Chapter 7, Claude Code does not provide streaming inline autocomplete -- and this gap is architectural, not accidental. The practical strategy is a hybrid toolchain: Claude Code for repo-scale reasoning and agentic execution, a dedicated inline completion tool for moment-to-moment typing assistance. The two tools operate at different interaction scales and do not conflict. Industry analysis positions this hybrid approach as the dominant pattern for 2026 and beyond.

From Solo Coding to Team Orchestration

The deepest strategic shift is not about cost or speed. It is about what the developer's job becomes.

In traditional development, the developer writes code. They read requirements, think about implementation, type code into an editor, run tests, fix bugs, and repeat. The primary skill is code production.

With Claude Code, the developer's role shifts toward orchestration. They define the task. They provide context. They review output. They steer corrections. They decide what to build next. The primary skills become architecture, coordination, and quality evaluation.

This is not a theoretical prediction. Internal teams at Anthropic report that their workflow has shifted from "write code, debug code" to "define task, delegate to Claude Code, review result, iterate." The engineer becomes a product owner for their feature and an architect for their system, with Claude Code as the engineering team that executes the implementation.

The shift has implications for how teams are structured. When each developer can orchestrate Claude Code to produce more output, smaller teams can handle larger projects. The bottleneck moves from coding bandwidth to decision-making bandwidth. The constraint is not "how fast can we type?" but "how clearly can we define what we want?"

Reasoning Cost Control

Not all tokens are equal. Claude Code's thinking -- the reasoning it does before producing output -- has its own cost profile, and controlling it is a practical economic lever.

Interleaved thinking allows Claude to reason between tool calls, producing better decisions but consuming more tokens. The thinking mode and effort level settings let you adjust this tradeoff:

- Higher effort levels produce more thorough reasoning, useful for complex architectural decisions.
- Lower effort levels produce faster, cheaper responses, suitable for straightforward implementation tasks.

With Opus 4.6, reasoning uses adaptive allocation: instead of a fixed thinking token budget, the model dynamically allocates thinking based on the effort level you select (low, medium, or high). You adjust effort in the `/model` command using arrow keys, and the change takes effect immediately. The `CLAUDE_CODE EFFORT LEVEL` environment variable sets it globally. For other models, the `MAX_THINKING_TOKENS` environment variable caps the thinking budget at a specific number of tokens. Setting it to zero disables thinking entirely on any model.

One counter-intuitive detail: phrases like "think harder" or "ultrathink" in your prompts are interpreted as regular instructions. They do not allocate additional thinking tokens. The thinking budget is controlled exclusively through the configuration mechanisms

above, not through conversational requests.

The economic implication: match the reasoning investment to the task. Do not pay for deep reasoning on a simple file rename. Do not skimp on reasoning for a security-critical refactor.

This connects to the model selection strategy: Opus with high reasoning effort for architectural planning, Sonnet with moderate reasoning for implementation, Haiku with minimal reasoning for exploration. Each combination represents a different price-performance point on the cost curve.

Timeline Reduction Evidence

Across documented case studies, timeline reductions of 70-90% are consistent. Projects that traditionally take months are completing in weeks. Projects that take weeks are completing in hours.

Specific examples:

- A full production platform that would have taken three to six months was completed in eight weeks.
- A frontend rewrite that had consumed three years of manual development was redone in weeks.
- A feature implementation across a twelve-million-line codebase completed in seven hours.
- Individual developers report producing multi-thousand-dollar equivalent output in evening sessions.

The compression comes from three sources. First, Claude Code eliminates the knowledge-acquisition phase. It reads code as fast as it reads anything, so the weeks spent understanding an existing system before modifying it shrink to hours. Second, Claude Code produces code faster than human typing. An experienced developer might write 100 lines of production code per hour. Claude Code can produce that in minutes. Third, Claude Code does not context-switch. It does not check email, attend meetings, or lose focus. During a session, it is 100% focused on the task.

The 70-90% number is not aspirational. It is what teams report after adjusting for the overhead of prompt engineering, review, and iteration. The raw code production speed improvement is higher, but the human time spent directing, reviewing, and correcting brings the net improvement into the 70-90% range.

Professional Output Economics

The economics become particularly striking for individual practitioners working outside traditional employment.

A solo developer using Claude Code in evening sessions can produce output that would have required hiring a contractor or consultant at significant cost. Portfolio analysis, application development, system architecture, data processing -- tasks that require specialized expertise can now be performed by someone with general technical literacy and good prompting skills.

This is not about replacing professionals. It is about expanding access. A developer who cannot afford a financial advisor can use Claude Code to build a portfolio optimization plan. A small business owner who cannot afford a development team can build internal tools. A researcher who cannot afford a data engineering consultant can build data pipelines.

The economic shift is from "expertise requires expensive humans" to "expertise requires good context and clear objectives." The marginal cost of producing one more report, one more analysis, one more application drops toward the token cost of the session. For many tasks, that is orders of magnitude less than the human labor cost.

Democratization of Expertise

The economic argument extends beyond developers. Non-technical users -- in legal, marketing, design, finance, and operations -- are using Claude Code to build tools and automations that previously required engineering resources.

The experience is fundamentally different for these users than it is for developers. Developers experience Claude Code as augmented velocity: they do the same things faster. Non-technical users experience it as entirely new capability: they do things they could never do before.

A legal team writing contract analysis tools. A marketing team building campaign automation. A finance team creating custom reporting dashboards. These are not developers who learned to code faster. These are professionals who gained access to technical capability without going through the traditional acquisition path of learning to program.

The economic implication is that the demand for software is about to get a lot larger. When the cost of building a custom tool drops by an order of magnitude, the set of problems worth solving with software expands dramatically. Organizations will build internal tools they never would have justified at previous costs.

Output Volume Over Speed

One finding consistently surfaces in organizational data: AI-assisted development increases output volume more than it increases speed per task. Roughly 27% of AI-assisted tasks are things that would not have been done at all without AI assistance.

This is the underappreciated economic story. The headline is "developers are faster." The real story is "developers are doing more." They ship more features. They run more experiments. They write more tests. They build more tools. They pursue ideas that would have been deprioritized because the implementation cost was too high.

The economic value of this is harder to measure than time savings, but it is potentially larger. A feature that was not going to be built because it would take two weeks but now takes two hours is not a two-week time savings. It is a new feature that did not exist before. The revenue it generates, the users it retains, the competitive advantage it provides -- none of that was in the baseline.

This means ROI calculations based purely on time savings undercount the actual value. The correct comparison is not "same output, less time" but "more output, same time." The developer is not going home two hours early. They are shipping two more features.

Timeline Compression and Project Viability

When timelines compress by 70-90%, projects that were previously non-viable become feasible.

Every organization has a graveyard of ideas that were good but too expensive to build. A custom analytics dashboard that would take a team of three developers six months. A data migration that would require a specialist for four weeks. A prototype that would need a designer, a frontend developer, and a backend developer for two months.

At 70-90% timeline compression, the six-month project becomes a three-week project. The four-week migration becomes a three-day migration. The two-month prototype becomes a one-week prototype. At those timescales, the cost-benefit calculation changes fundamentally. Projects that did not clear the ROI threshold at their original timeline clear it easily at the compressed timeline.

This is how AI-assisted development changes strategy, not just execution. The set of things worth building expands. The barrier to experimentation drops. The organization can pursue more ideas in parallel, fail faster on the ones that do not work, and double down on the ones that do.

Role Shift: Architecture, Coordination, Quality

As code production becomes cheaper, the skills that remain expensive -- and therefore more valuable -- are the ones Claude Code cannot do as well.

Architecture. Deciding what to build, how to structure it, and what tradeoffs to make. Claude Code can propose architectures, but evaluating whether those proposals fit the organization's constraints, culture, and long-term direction requires human judgment.

Coordination. Managing dependencies between people, teams, and systems. Claude Code operates within a session. The cross-session, cross-team, cross-organization coordination that makes large projects succeed is still fundamentally human.

Quality evaluation. Knowing whether the output is good enough. Claude Code can produce code that passes tests, but deciding whether the tests cover the right scenarios, whether the user experience is acceptable, and whether the solution will scale requires domain expertise and judgment.

The developers who thrive in this environment are not the fastest coders. They are the clearest thinkers. They define problems precisely. They evaluate solutions rigorously. They make architectural decisions that hold up over time. These skills were always

valuable. They are about to become the primary job description.

One-Person Teams at Scale

The logical endpoint of these trends is the one-person team: a single individual who handles work that previously required an entire department.

This is already happening. Individual contributors using Claude Code report handling the output volume of small teams. They architect systems, delegate implementation to Claude Code, review results, iterate, and ship. They are simultaneously the product owner, the architect, the developer, and the QA engineer -- because the AI handles the implementation labor that used to require separate people in each role.

The one-person team is not a novelty. It is an economic inevitability when the cost of code production drops by an order of magnitude while the cost of coordination remains constant. Adding a second person to a team adds communication overhead, alignment overhead, and scheduling overhead. If one person with Claude Code can produce the same output, the coordination cost of the second person is pure waste.

This does not mean teams disappear. Complex systems still require multiple humans for the architecture, coordination, and quality evaluation that Claude Code cannot provide. But the threshold for "complex enough to need a team" rises. Projects that used to require five people now require two. Projects that required two now require one.

The strategic implication for organizations is clear: invest in making each individual contributor more effective with AI tools rather than adding more contributors. The returns from empowering one excellent developer with Claude Code exceed the returns from hiring two more developers without it.

Three Multipliers Driving Acceleration

Timeline compression is not a single phenomenon. It is the product of three compounding multipliers, and understanding the interaction explains why gains feel exponential rather than linear.

The first multiplier is agent capabilities. Each generation of model improvements makes Claude Code more capable at understanding complex codebases, generating correct implementations, and recovering from errors. This is the most visible multiplier -- the one that makes headlines -- but it is not the largest in practice.

The second multiplier is orchestration improvements. Multi-agent coordination, task decomposition patterns, skill ecosystems, and MCP integrations create leverage on top of raw model capabilities. An engineer who learns to decompose a project into parallelizable subagent tasks does not get a linear speedup. They get the model capability multiplied by the orchestration efficiency, producing results that neither could achieve alone.

The third multiplier is accumulated human experience. As developers learn which tasks Claude handles well, how to structure prompts for maximum first-pass quality, and when to intervene versus when to let Claude iterate autonomously, their per-session productivity increases. This compounding effect is invisible in benchmark data but obvious in practice. A developer's hundredth hour with Claude Code is dramatically more productive than their first.

The three multipliers interact. Better agent capabilities make orchestration patterns more reliable, which lets humans develop more ambitious workflows, which produce better outcomes that inform future orchestration patterns. The result is step-function improvements rather than the linear gains that spreadsheet projections would predict.

The Portfolio Optimization Story

The economics become tangible through individual narratives. Consider what happens when someone with general technical literacy and domain knowledge -- but no specialized financial planning background -- uses Claude Code for portfolio optimization.

The workflow started with three inputs: CSV exports from multiple brokerage accounts containing holdings, cost basis, and fund data; a supplementary text file with unstructured information that had no export mechanism (bank balances, employer retirement details, automated investment settings); and a carefully written goal prompt describing desired outputs, constraints, and preferences.

The goal prompt was the most important piece. It specified what data existed and where, the desired output format, investment preferences (tax efficiency over performance chasing, lowest-fee funds, simplicity), known constraints (certain accounts left untouched), and -- critically -- what the user suspected was wrong, caveated with an invitation for Claude to disagree. A strawman target allocation gave Claude something concrete to react to rather than building from scratch. The prompt ended with a request to ask clarifying questions exhaustively before jumping to conclusions.

Claude Code wrote Python scripts to parse the CSVs, classify every holding into target categories, handle encoding issues, deduplicate overlapping exports, and deal with non-standard line items. The baseline it produced -- a full gap analysis across every account and category -- was strong enough to iterate on immediately.

The iteration phase worked like collaborating with an analyst who has all the data but has not lived with the accounts. Each constraint refinement triggered automatic updates across every calculation, table, and recommendation. When asked to self-critique and rate its own plan on a one-to-ten scale, Claude identified seven improvements, including an overweight position in tax-free retirement accounts where selling would cost literally zero in taxes -- something the original plan had left untouched.

The final output was a ten-section plan document with appendices, including sections the user had not requested: retirement contribution optimization, a natural dilution timeline projecting how long overweight categories would take to reach targets through new contributions alone, and tax-loss harvesting partners for every recommended position. The full plan document was then remixed into a shorter checklist for execution.

A CLAUDE.md file accumulated data quirks, strategy decisions, and target allocations across sessions. When the user returned days later, Claude picked up exactly where they left off. The work product compounded: the same project directory was later used to write a blog post about the process, with Claude reading the session history and plan documents as context.

A year earlier, this would have been a week of spreadsheet work or a few thousand dollars paid to a financial advisor. Instead, it was a few evenings of back-and-forth with Claude Code, and the plan produced was the one being actively executed.

The Three-Year Rewrite

The economics of replacement are different from the economics of creation, and the emotional trajectory is different too.

One developer spent three years building a configuration interface for an algorithmic trading platform. The interface used complex tree-like structures to let users configure trading conditions. It required a strong mental model to navigate. But it worked, and the developer was proud of it.

Then Claude Code made a better design. Not incrementally better -- fundamentally better. The developer's recognition that Claude was a better UX designer and frontend developer than they had ever been led to a decision that most developers would resist: start over completely.

The rewrite replaced the hand-crafted tree-structured interface with a natural language approach. Users describe trading strategies in plain English -- "create a strategy that buys if twenty days passed since the last purchase and the RSI is below thirty-five" -- and Claude decomposes the statement into structured objects: portfolios, strategies, actions, conditions. What took three years of manual UI development to create was replaced by an approach where natural language input produces the same structured data that the complex interface produced, but without requiring users to navigate the configuration maze.

The platform then added a no-code UI layer for manual fine-tuning. The pattern is noteworthy: AI creates the first draft from natural language, humans refine through a visual interface. The complex configuration interface built over three years -- replaced by natural language input that a model decomposes into the same structures. Over twenty-five thousand people now have access to advanced trading tools through this democratized interface.

The narrative arc matters. Pride in the old work. Recognition that AI exceeded the developer's own specialized skills. The decision to start over. This is not a story about laziness or shortcuts. It is a story about recognizing when a better approach makes previous effort irrelevant -- and having the discipline to act on that recognition.

The Backend Developer's Full-Stack Sprint

A backend developer with no experience in frontend frameworks, no design background, and no financial data visualization expertise built a complete stock trading analysis application in under ten hours.

The application was not a toy. It included modular components for data fetching from multiple stock APIs with normalization, technical indicators (fifteen or more), sentiment analysis from news sources, risk assessment with value-at-risk calculations and stress tests, portfolio tracking with performance attribution, and charting with candlestick displays and indicator overlays. The architecture was clean -- modular enough that adding new indicators or data sources was straightforward.

The key challenge was not any single component but the intersection of three domains the developer did not command: frontend framework patterns, financial data visualization, and design sensibility. Claude Code bridged all three simultaneously. The developer knew trading concepts. Claude Code knew frontend patterns and financial UI conventions. Together they built something neither could have built alone in that timeframe.

This is a specific instance of the broader pattern: everyone becomes more full-stack. Analysis of how different teams use AI reveals a consistent finding -- people use AI to augment their core expertise while expanding into adjacent domains. Security teams analyze unfamiliar code. Research teams build frontend visualizations of their data. Non-technical employees debug network issues and perform data analysis. The long-held assumption that serious development work requires deep specialization in every relevant technology is dissolving.

The Eight-Week Production Platform

Documented week-by-week timelines tell a more precise story than aggregate compression percentages. One team built a complete production trading platform in eight weeks using Claude Code as their primary development tool:

Weeks one and two covered architecture and database design. Claude Code generated schemas, entity-relationship diagrams, and migration scripts. The main artifact was a CLAUDE.md file capturing tech stack decisions. By the end, the database and tables were ready, and a container definition was drafted.

Weeks two and three built the backend API. Claude Code generated the server scaffold, order and position and execution endpoints, and the real-time communication server. The artifact was a working backend running locally.

Weeks three and four produced frontend UI components -- forms, charts, tables, state management, type definitions for API responses. The artifact was a working frontend running locally, progressing from wireframe to interactive prototype.

Weeks four and five handled broker integration and the strategy engine -- API wrappers for broker services, order placement logic, and a backtesting engine. The artifact was a working "send real order" capability, tested on a paper trading account.

Weeks five and six covered testing and documentation -- unit tests, end-to-end test scenarios, API documentation. Test coverage exceeded eighty percent and CI pipelines were passing.

Weeks six and seven handled DevOps and deployment -- container images, CI/CD workflows, infrastructure provisioning for a staging environment. Automated pipelines ran smoke tests on staging.

Weeks seven and eight were production hardening -- security scanning, compliance checks, load testing, stress testing, production deployment, and monitoring.

Eight weeks for a production-grade platform with a UI, backend, broker integration, strategy engine, test suite, and deployment infrastructure. The traditional estimate for the same work: three to six months. Individual features within the platform followed a similar compression pattern. Adding a new strategy type -- data model changes, backend logic, API endpoints, frontend UI, tests, documentation -- took two to three days instead of the traditional two to three weeks.

Feature Development Workflow Compression

The week-by-week timeline tells the macro story. The feature-by-feature timeline tells the micro story.

Adding a new capability to an existing platform follows a six-step workflow: update the data model (database migration, validation models, type definitions), implement backend logic (signal evaluation, execution logic, exit conditions, tests), build API endpoints

(create, read, update, delete, plus specialized operations), create frontend UI (forms, displays, results visualization), write tests and configure CI (unit, integration, end-to-end), and update documentation. With Claude Code generating artifacts at each step, the workflow takes two to three days. Traditional development of the same feature: two to three weeks.

The compression is not uniform across steps. Claude Code handles data model changes and API endpoints with near-perfect first-pass quality. Frontend UI and backend business logic require more iteration. Test generation is fast but tests need human review for coverage completeness. Documentation is generated almost instantly. Understanding where the compression is strongest and weakest lets teams allocate their review time where it matters most.

The Competitive Landscape

Claude Code does not exist in a vacuum. Understanding where it fits -- and where it does not -- relative to other tools is necessary for making sound toolchain decisions.

A Major Code Editor with Built-In Agents

One prominent IDE has built agent capabilities directly into the editor. It can run multiple agents in parallel -- up to eight on a single prompt -- switching between them in a sidebar. It includes an in-editor browser for inspecting rendered pages, team-level shared commands for workflow standardization, and inline completion that reviewers consistently describe as best-in-class for IDE-centric autocomplete. The strength is immediacy: everything happens inside the editor, suggestions appear as you type, and multi-agent orchestration has a visual interface. For developers who spend their entire day in a single editor and value low-latency inline suggestions, this approach is compelling. It is furthest along in providing explicit multi-agent UIs with visible progress tracking.

A Leading Development Platform

One of the largest development platforms has evolved from simple inline suggestions into a full agentic workflow environment. Its workspace model provides a structured flow from task definition through specification, planning, and code generation -- a web experience where each stage is visible and editable before the next begins. It offers persistent, context-rich spaces that ground the AI in code, documentation, and specifications. It supports the MCP protocol for external tool integration, has a dedicated coding agent, and recently opened third-party coding agent support. Its greatest strength is the depth of integration with version control, IDE, and the broader development platform. AI-generated commit messages, pull request summaries, and issue-to-PR workflows are first-class features. For organizations already invested in this ecosystem, the integration is natural.

A Major AI Lab's Development Stack

A leading AI research lab approaches the problem from a reasoning-first perspective, pairing strong reasoning models with integrated development tooling. Its coding surface is evolving from "a model you prompt" to a full development environment with web, cloud, and IDE integration, longer sessions, and iterative problem solving. It offers automatic fix suggestions during CI runs -- when tests fail, the system proposes patches without human intervention. It provides evaluation APIs and programmable graders for measuring code quality systematically. Its strength is CI integration and verifiable reasoning, particularly for workflows that need automated quality gates and measurable improvement tracking.

Where Claude Code Fits

Claude Code's competitive differentiator is repo-scale planning and patch generation. It excels at large refactors across massive codebases, complex architectural changes that span many files, and tasks that require deep understanding of how an entire system fits together. One documented case had Claude Code autonomously implementing a complex feature inside a twelve-and-a-half-million-line codebase in seven hours with 99.9 percent numerical accuracy.

The tradeoff is latency and cost. Claude Code's deep reasoning model costs more per token than lightweight inline completion. The terminal-first interface prioritizes depth of understanding over speed of suggestion. For moment-to-moment typing assistance -- the "complete this line" interaction that happens hundreds of times per hour -- Claude Code is not the right tool.

The Hybrid Toolchain Recommendation

The sensible strategy for 2026 is not choosing a single tool. It is composing a toolchain:

Use Claude Code as the primary agentic coding tool for repo-level changes, architectural refactors, multi-file implementations, and complex reasoning tasks. Layer a dedicated inline completion tool for commodity autocomplete -- the moment-to-moment typing assistance that benefits from low latency and high suggestion frequency. Keep watching the IDE integration and skills ecosystem as both Claude Code and its competitors push into each other's strengths.

The hybrid approach outperforms either tool used alone because the interaction scales are different. Inline completion operates at the keystroke level. Claude Code operates at the task level. Trying to use one tool for both levels creates friction in both directions.

Financial Task Benchmarks

Domain-specific performance data grounds the economic argument in measurable capability. On a standardized finance agent benchmark, Claude Code achieved 55.3 percent accuracy, leading the field. On difficult spreadsheet modeling challenges -- the kind of work that traditionally requires a financial analyst with years of experience -- it scored eighty-three percent. The context window capacity exceeding one hundred thousand tokens enables processing hundreds of pages of financial documents in a single session.

These are not abstract benchmarks. They translate directly into the portfolio optimization, trading platform development, and financial modeling workflows described above. When a model can handle eighty-three percent of difficult modeling challenges correctly, the human's role shifts from doing the modeling to reviewing and correcting the modeling -- a categorically different workflow with categorically different economics.

Open-Source Autonomous Experiments

The democratization thesis extends to its logical extreme: open-source projects where autonomous agents manage financial assets. One experiment built an autonomous asset manager with cryptocurrency integration, driven by a simple thesis -- most people get index funds while the wealthy get analyst teams, and AI can democratize asset management.

These experiments are early-stage and high-risk. They are worth noting not because they represent production-ready systems but because they illustrate the economic trajectory. When the tools to build autonomous agents are open source and the cost to run them is measured in tokens, the barrier between "interesting idea" and "working prototype" compresses from months to days. The number of experiments increases by orders of magnitude. Most will fail. Some will not.

Path to Market

The timeline compression extends beyond development into entrepreneurship. When agents can work autonomously for extended periods, entrepreneurs go from ideas to deployed applications in days instead of months. The traditional path -- idea, prototype, funding, team, development, launch -- compresses because the development phase is no longer the bottleneck. A single person with domain expertise and good prompting skills can produce a working application, test it with real users, and iterate based on feedback at a pace that would have required a funded team a few years ago.

This is not about replacing teams. Complex products still need teams for the architecture, coordination, user research, and quality evaluation that Claude Code cannot provide. But the threshold for "complex enough to need a team" has risen. The set of products that one person can build and ship has expanded dramatically, and the economics of early-stage experimentation have changed accordingly.

Technical Debt as Agent Backlog

One prediction from organizational deployment data deserves its own treatment: technical debt -- the accumulated mass of shortcuts, workarounds, and deferred improvements that accrues in every codebase -- becomes systematically addressable when agents can work through backlogs.

Every organization has a list of known issues that nobody prioritizes because the implementation cost exceeds the business case. Rename this confusing module. Migrate from the deprecated library. Add tests to the untested module. Fix the inconsistent error handling. Each task is individually low-value but collectively significant.

When the cost of implementation drops by an order of magnitude, the economics change. Tasks that never cleared the priority threshold at their original cost clear it easily at the compressed cost. Agents can work through technical debt backlogs during low-priority windows -- evenings, weekends, between sprints -- producing clean-up commits that human reviewers approve. The codebase improves gradually without competing with feature development for engineering time.

Task Horizons Expanding

Early agents handled one-shot tasks that took a few minutes: fix this bug, write this function, generate this test. By late 2025, increasingly capable agents were producing full feature sets over the course of several hours. The trajectory points toward agents working for days at a time, building entire applications with minimal human intervention focused on strategic oversight at key decision points.

This expansion changes the economics of what projects are viable. When agent task horizons are measured in minutes, you can automate small fixes. When measured in hours, you can automate feature development. When measured in days or weeks, you can automate entire product development cycles. The planning overhead -- decomposing work, providing context, reviewing output -- stays roughly constant, while the work performed per planning cycle increases dramatically.

Four Priorities for 2026

The trends above converge on four priorities that separate organizations treating agentic development as a strategic capability from those treating it as a productivity tool.

First: master multi-agent coordination. Single-agent workflows hit complexity ceilings. Organizations that learn to decompose work across coordinated agents -- planners, implementers, testers, reviewers -- handle complexity that single-agent systems cannot address.

Second: scale human-agent oversight through AI-automated review. The bottleneck in scaling AI-assisted development is not the AI's ability to produce code. It is the human's ability to review it. Organizations that build automated review systems -- linting, testing, security scanning, style checking -- focus human attention where it matters most and let machines handle verifiable quality criteria.

Third: extend agentic coding beyond engineering. The teams that extract the most novel value from Claude Code are not engineering teams. They are legal, marketing, design, and operations teams building capabilities they never had before. Organizations that treat Claude Code as an engineering tool miss the larger opportunity.

Fourth: embed security architecture from the earliest stages. As agents become more capable and autonomous, the attack surface expands. Organizations that bake security into their agent architecture from the start -- managed policies, sandboxed execution, permission boundaries, MCP server restrictions -- are better positioned than those that bolt security on after deployment.

Where Things Are Headed

Predicting specific features is a losing game, but the direction of travel is clear from what competitors are building and what the ecosystem is demanding.

Workspace-style planning interfaces -- where task definition, specification, plan, and code changes are visible and editable as separate panels rather than interleaved in a conversation -- are a natural evolution. Claude Code's plan mode already separates planning from execution. A visual surface for that separation is a logical next step.

Native multi-agent controls with explicit management UIs -- agent tabs, profiles, parallel execution dashboards -- would bring Claude Code's powerful but CLI-driven multi-agent capabilities to a broader audience. The primitives exist. The interface layer is what is missing.

First-class CI and code review skills -- packaged recipes for reading failing test logs, proposing patches, and opening pull requests automatically -- would bring Claude Code's agentic capabilities into the workflows where most code quality decisions are made. One competitor already offers automatic fix suggestions during CI runs. The capability exists in Claude Code through skills and MCP, but it is not yet a single toggle.

Longer-running project-level jobs with checkpoints and progress dashboards would extend the task horizon from hours to days. Better resume and retry support, job identifiers, resumable sessions, and visual progress tracking in web and desktop interfaces would make overnight and multi-day agent work practical.

More official skill packs from Anthropic -- for CI, testing, migrations, and security scans -- would reduce the setup cost for common workflows and establish best practices that individual teams currently have to discover independently.

The pragmatic question is not "should I wait for these?" but "what should I build now and what should I wait for?" Build with Claude Code now for repo-wide refactors, agentic workflows, deep reasoning tasks, and any scenario where understanding a large codebase is the bottleneck. Use a hybrid toolchain for inline autocomplete and rapid iteration within an editor. Watch for multi-agent UIs, CI integration recipes, and longer-running job support as the gaps most likely to close in the near term.

Key Takeaways

- Prompt caching makes rich context (CLAUDE.md, MCP tools, system prompts) economically viable by charging full price only once per session.
- Match models to tasks: Opus for architecture, Sonnet for implementation, Haiku for exploration -- and use subagents to mix models within a session. Control reasoning costs with effort levels (low/medium/high) and the `CLAUDE_CODE EFFORT LEVEL` environment variable.
- Three compounding multipliers -- agent capabilities, orchestration improvements, and accumulated human experience -- produce step-function gains rather than linear ones.
- Timeline compression of 70-90% changes project viability: an eight-week production platform instead of three to six months, a two-to-three-day feature instead of two to three weeks.
- A hybrid toolchain (Claude Code for repo-scale reasoning, a dedicated inline tool for keystroke-level autocomplete) outperforms either tool alone because they operate at different interaction scales.
- Claude Code's competitive strength is repo-scale planning and patch generation in large codebases; other tools lead in inline autocomplete, visual multi-agent management, and CI integration.
- Everyone becomes more full-stack: backend developers build frontend applications in hours, designers make state management changes, and non-technical teams build entirely new capabilities.
- Technical debt becomes systematically addressable when agents work through backlogs at compressed cost, tackling tasks that never cleared the priority threshold before.
- Four priorities for 2026: master multi-agent coordination, scale oversight through automated review, extend agentic coding beyond engineering, and embed security architecture from day one.

Appendix A: Command Reference

CLI Commands

Command	Description
<code>claude</code>	Start interactive REPL in current directory
<code>claude "task"</code>	Start interactive REPL with initial prompt
<code>claude -p "task"</code>	Non-interactive (headless) mode; print response and exit
<code>claude -c</code>	Continue most recent conversation in current directory
<code>claude -c "task"</code>	Continue most recent conversation with new prompt
<code>claude -r <id></code>	Resume a specific session by ID
<code>claude -r <name></code>	Resume a specific session by name

CLI Flags

Session and Input

Flag	Description
--continue, -c	Continue most recent conversation
--resume, -r <id>	Resume session by ID or name
--fork-session	Create new session ID preserving conversation history (use with --resume or --continue)
--from-pr <number>	Resume sessions linked to a specific pull request; accepts PR number or URL
--session-name <name>	Name the current session
--model <model>	Override the default model
--agent <name>	Run a custom agent as the main thread (overrides agent setting)
--agents <json>	Define custom subagents dynamically via JSON
--permission-mode <mode>	Set permission mode (default, plan, auto-edit, full-auto, bypassPermissions)
--teleport	Resume a web session in your local terminal
--teammate-mode <mode>	Set agent team display: auto (default), in-process, or tmux

System Prompt

Flag	Description
--system-prompt <text>	Replace the system prompt entirely with provided text (headless mode only)
--system-prompt-file <path>	Replace the system prompt entirely with file contents (headless mode only)
--append-system-prompt <text>	Add custom instructions to the end of the default system prompt (interactive and headless)
--append-system-prompt-file <path>	Append file contents to the default system prompt (interactive and headless)

Output Control (Headless Mode)

Flag	Description
--output-format text	Plain text output (default for -p)
--output-format json	JSON object with result, cost, duration, and session ID
--output-format stream-json	Newline-delimited JSON stream for real-time processing
--max-turns <n>	Limit the number of agentic turns
--budget-tokens <n>	Set token budget for the session
--fallback-model <model>	Model to use if primary model is unavailable

Permission Flags

Flag	Description
--allowedTools <tools>	Comma-separated list of tools to allow without prompting
--disallowedTools <tools>	Comma-separated list of tools to deny
--dangerously-skip-permissions	Skip all permission prompts (use only in sandboxed environments)
--permission-prompt-tool <mcp_tool>	Delegate permission decisions to an MCP tool

Miscellaneous

Flag	Description
--verbose	Enable verbose logging
--no-cache	Disable prompt caching
--version	Print version and exit
--help	Print help and exit

Slash Commands

Command	Description
/help	Show available commands and shortcuts
/clear	Clear conversation history and start fresh
/compact	Manually compact conversation to reclaim context space
/config	Open or manage configuration
/context	Visualize current context usage as a colored grid
/copy	Copy last assistant response to clipboard
/cost	Show token usage statistics
/desktop	Hand off CLI session to the desktop application (macOS/Windows)
/doctor	Run diagnostics to check for common issues
/agents	List available subagents
/export [filename]	Export conversation to file or clipboard
/hooks	Interactive hook management menu (view, add, delete hooks)
/ide	Open current session in the IDE extension
/init	Initialize CLAUDE.md in current project (analyzes codebase for build systems, test frameworks, patterns)
/mcp	Show MCP server status and per-server token costs
/model	Switch model mid-session; with Opus, use left/right arrows to adjust effort level

/permissions	View and manage permission rules
/plan	Enter plan mode directly from the prompt
/rename	Rename the current session
/resume	List and resume previous sessions (opens session picker)
/rewind	Rewind conversation and/or code, or summarize from a selected message
/sandbox	View sandbox status and settings
/stats	Visualize daily usage, session history, streaks, and model preferences
/statusline	Set up status line UI in terminal
/tasks	List and manage background tasks
/teleport	Resume a remote session from the web (subscribers only)
/terminal-setup	Install terminal key bindings for multiline input and shortcuts
/vim	Enable vim-style editing mode

Keyboard Shortcuts

Navigation and Control

Shortcut	Action
Enter	Send prompt / confirm
Ctrl+C	Cancel current input or generation
Ctrl+D	Exit Claude Code session
Up / Down	Navigate prompt history
Left / Right	Cycle through dialog tabs in permission dialogs and menus
Esc Esc (double)	Open rewind menu: restore code and/or conversation, or summarize from selected message
Shift+Tab	Cycle permission modes (plan → default → auto-edit → full-auto; includes delegate mode when agent team active)

Session and Model

Shortcut	Action
Alt+P / Option+P	Switch models without clearing current prompt
Alt+T / Option+T	Toggle extended thinking mode (run /terminal-setup first)

Tools and Output

Shortcut	Action

<code>Ctrl+R</code>	Reverse search through command history (type query, <code>Ctrl+R</code> to cycle matches, <code>Tab/Esc</code> to accept, <code>Enter</code> to accept and execute, <code>Ctrl+C</code> to cancel)
<code>Ctrl+G</code>	Open prompt in default text editor for editing
<code>Ctrl+O</code>	Toggle verbose output showing detailed tool usage and execution
<code>Ctrl+B</code>	Background a running task (tmux users must press twice due to tmux prefix key)
<code>Ctrl+T</code>	Toggle task list display (shows up to 10 tasks in terminal status area)
<code>Ctrl+L</code>	Clear terminal screen (keeps conversation history)
<code>Ctrl+V / Cmd+V (iTerm2) / Alt+V (Windows)</code>	Paste image from clipboard

Permission Modes

Mode	Behavior
Plan	Read-only; Claude cannot modify files or run write commands
Default	Asks permission for file edits and shell commands
Auto-accept edits	File edits proceed without asking; shell commands still prompt
Full auto-accept	All operations proceed without asking
Bypass permissions	Skip all checks (CLI flag only; requires sandboxed environment)

Output Formats

--output-format text

Plain text response. Default for headless mode. Suitable for piping to other CLI tools.

--output-format json

Single JSON object on completion:

- `result` — response text
- `cost` — token costs
- `duration` — execution time
- `session_id` — session identifier
- `is_error` — boolean

--output-format stream-json

Newline-delimited JSON messages during execution. Each line is a JSON object with a `type` field indicating the event kind (assistant message, tool use, tool result, system message). Suitable for real-time monitoring and integration.

Tool Names for Permission Rules

Tool	Description
Bash	Shell command execution

Read	File reading
Write	File creation/overwrite
Edit	File editing (string replacement)
Glob	File pattern matching
Grep	Content search
WebFetch	URL fetching
WebSearch	Web search
Task	Subagent spawning
Skill	Skill invocation
NotebookEdit	Notebook cell editing (.ipynb)
TodoRead	Read task lists
TodoWrite	Write/update task lists

Permission rules use `Tool` or `Tool(specifier)` syntax with glob patterns. Example: `Bash(git *)` matches any git command.

MCP Tool Pattern for Permission Rules

MCP tools follow the pattern `mcp__<server>__<tool>`. Glob patterns are supported.

Examples:

- `mcp__myserver__query` — specific tool on specific server
- `mcp__myserver__*` — all tools on a specific server
- `mcp__*__*` — all MCP tools

Multiline Input Methods

Method	Shortcut	Context
Quick escape	\ + Enter	Works in all terminals
macOS default	Option+Enter	Default on macOS
Shift+Enter	Shift+Enter	Works out of the box in iTerm2, WezTerm, Ghostty, Kitty
Control sequence	Ctrl+J	Line feed character for multiline
Paste mode	Paste directly	Auto-detected for code blocks, logs

For terminals not listed above (VS Code terminal, Alacritty, Zed, Warp), run `/terminal-setup` to install the `Shift+Enter` binding.

Bash Mode (! Prefix)

Run bash commands directly without Claude interpreting them by prefixing input with `!`:

```
! npm test  
! git status  
! ls -la
```

- Adds command output to conversation context
 - Shows real-time progress and output
 - Does not require Claude to interpret or approve the command
 - Supports `Ctrl+B` backgrounding for long-running commands
 - History-based autocomplete: type partial command and press `Tab` to complete from previous `!` commands in the current project
-

Background Bash Commands

Claude Code runs bash commands in the background, returning a unique task ID immediately while the command executes asynchronously.

Triggering:

- Prompt Claude to run a command in the background
- Press `Ctrl+B` to move a running Bash invocation to the background (tmux users press `Ctrl+B` twice)

Behavior:

- Output is buffered; Claude retrieves it using the `TaskOutput` tool
- Each background task has a unique ID for tracking and retrieval
- Tasks are automatically cleaned up when Claude Code exits

Common backgrounded commands: build tools (webpack, vite, make), package managers (npm, yarn), test runners (jest, pytest), development servers, long-running processes (docker, terraform)

Disable: Set `CLAUDE_CODE_DISABLE_BACKGROUND_TASKS=1`

Prompt Suggestions

After Claude responds, grayed-out suggestions appear based on conversation history and git history.

- Press `Tab` to accept the suggestion, or `Enter` to accept and submit
- Start typing to dismiss
- Runs as a background request reusing the prompt cache (minimal additional cost)
- Skipped when the cache is cold, after the first turn, in non-interactive mode, and in plan mode

Disable: Set `CLAUDE_CODE_ENABLE_PROMPT_SUGGESTION=false` or toggle in `/config`

Task List

Claude creates task lists for complex multi-step work, visible in the terminal status area.

- `Ctrl+T` toggles the task list view (displays up to 10 tasks)
 - Ask Claude directly to show all tasks or clear them
 - Tasks persist across context compactions
 - Share a task list across sessions: `CLAUDE_CODE_TASK_LIST_ID=my-project claude`
-

Session Picker Shortcuts

The `/resume` command (or `claude --resume` without arguments) opens an interactive session picker.

Shortcut	Action
Up / Down	Navigate between sessions
Left / Right	Expand or collapse grouped sessions
Enter	Select and resume the highlighted session
P	Preview session content
R	Rename the highlighted session
/	Search to filter sessions
A	Toggle between current directory and all projects
B	Filter to sessions from current git branch
Esc	Exit the picker or search mode

Vim Editor Mode

Enable with `/vim` command or configure permanently via `/config`.

Mode Switching

Command	Action	From Mode
Esc	Enter NORMAL mode	INSERT
i	Insert before cursor	NORMAL
I	Insert at beginning of line	NORMAL
a	Insert after cursor	NORMAL
A	Insert at end of line	NORMAL
o	Open line below	NORMAL
O	Open line above	NORMAL

Navigation (NORMAL Mode)

Command	Action
h/j/k/l	Move left/down/up/right
w	Next word
e	End of word
b	Previous word
0	Beginning of line
\$	End of line

^	First non-blank character
gg	Beginning of input
G	End of input
f{char}	Jump to next occurrence of character
F{char}	Jump to previous occurrence of character
t{char}	Jump to just before next occurrence
T{char}	Jump to just after previous occurrence
;	Repeat last f/F/t/T motion
,	Repeat last f/F/t/T motion in reverse

In NORMAL mode, if the cursor is at the beginning or end of input and cannot move further, the arrow keys navigate command history instead.

Editing (NORMAL Mode)

Command	Action
x	Delete character
dd	Delete line
D	Delete to end of line
dw/de/db	Delete word/to end/back
cc	Change line
c	Change to end of line
cw/ce/cb	Change word/to end/back
yy/Y	Yank (copy) line
yw/ye/yb	Yank word/to end/back
p	Paste after cursor
P	Paste before cursor
>>	Indent line
<<	Dedent line
J	Join lines
.	Repeat last change

Text Objects (NORMAL Mode)

Text objects work with operators like `d`, `c`, and `y`:

Command	Action

iw/aw	Inner/around word
iW/aW	Inner/around WORD (whitespace-delimited)
i"/a"	Inner/around double quotes
i'/a'	Inner/around single quotes
i(/a(Inner/around parentheses
i[/a[Inner/around brackets
i{/a{	Inner/around braces

Appendix B: Configuration Reference

Settings.json Keys by Scope

JSON Schema for Settings

Add the `$schema` line to any `settings.json` for autocomplete and inline validation in editors that support JSON schema:

```
{
  "$schema": "https://json.schemastore.org/clause-code-settings.json",
  "permissions": { ... },
  "env": { ... }
}
```

Permission Rules

Key	Type	Description
permissions.allow	string[]	Tools auto-approved without prompting
permissions.ask	string[]	Tools requiring user confirmation
permissions.deny	string[]	Tools always blocked

Rule syntax: `Tool` or `Tool(specifier)` with glob pattern support. Evaluation order: **Deny > Ask > Allow**, first match wins across all scopes.

Examples:

- `Bash(npm test)` -- specific bash command
- `Write(*.env)` -- file pattern
- `mcp__server__tool` -- MCP tool
- `Bash(git *)` -- wildcard command

Environment Variables (Settings)

Key	Type	Description
env	object	Key-value pairs injected into Claude's environment

Hook Configuration

Key	Type	Description
hooks	object	Event-to-handler mappings (see Hook Schema below)

Model Settings

Key	Type	Description
model	string	Default model identifier
availableModels	string[]	Models available for selection

Authentication and Credentials

Key	Type	Description
apiKeyHelper	string	Custom script (executed via /bin/sh) to generate an auth value. Value is sent as X-Api-Key and Authorization: Bearer headers for model requests
awsAuthRefresh	string	Custom script that modifies the .aws directory for credential refresh
awsCredentialExport	string	Custom script that outputs JSON with AWS credentials
forceLoginMethod	string	Restrict login to cloudeai (Claude.ai accounts) or console (API billing accounts)
forceLoginOrgUUID	string	UUID of an organization to auto-select during login (requires forceLoginMethod)

Session and Behavior

Key	Type	Description
cleanupPeriodDays	number	Sessions inactive longer than this are deleted at startup. Default: 30. Set to 0 for immediate cleanup
language	string	Preferred response language (e.g., "japanese", "spanish", "french")
outputStyle	string	Configure output style to adjust system prompt
plansDirectory	string	Customize where plan files are stored. Path relative to project root. Default: ~/claude/plans
showTurnDuration	boolean	Show turn duration messages (e.g., "Cooked for 1m 6s"). Default: true
teammateMode	string	How agent team teammates display: auto (picks split panes in tmux/iTerm2), in-process, or tmux

UI and Display

Key	Type	Description
autoUpdatesChannel	string	Release channel: "stable" (about one week old, skips regressions) or "latest" (most recent). Default: "latest"
spinnerTipsEnabled	boolean	Show tips in spinner while working. Default: true
spinnerVerbs	object	Customize spinner verbs. mode: "replace" (use only your verbs) or "append" (mix with defaults). verbs: string[]
terminalProgressBarEnabled	boolean	Enable terminal progress bar in Windows Terminal and iTerm2. Default: true

prefersReducedMotion	boolean	Reduce or disable UI animations (spinners, shimmer, flash effects) for accessibility
----------------------	---------	--

File Suggestion

Key	Type	Description
fileSuggestion	object	Custom script for @ file path autocomplete in large monorepos
respectGitignore	boolean	Whether the @ file picker respects .gitignore patterns. Default: true

The `fileSuggestion` script receives JSON on stdin with a `query` field and outputs newline-separated file paths on stdout.

Sandbox Settings

Key	Type	Description
<code>sandbox.enabled</code>	boolean	Enable bash sandboxing (macOS, Linux, WSL2). Default: false
<code>sandbox.autoAllowBashIfSandboxed</code>	boolean	Auto-approve bash commands when sandboxed. Default: true
<code>sandbox.excludedCommands</code>	string[]	Commands that run outside the sandbox (e.g., ["git", "docker"])
<code>sandbox.allowUnsandboxedCommands</code>	boolean	Allow commands to bypass sandbox via <code>dangerouslyDisableSandbox</code> . Set to <code>false</code> to enforce strict sandboxing. Default: true
<code>sandbox.network.allowUnixSockets</code>	string[]	Unix socket paths accessible in sandbox (e.g., SSH agent sockets)
<code>sandbox.network.allowAllUnixSockets</code>	boolean	Allow all Unix socket connections. Default: false
<code>sandbox.network.allowLocalBinding</code>	boolean	Allow binding to localhost ports (macOS only). Default: false
<code>sandbox.network.allowedDomains</code>	string[]	Domains accessible from sandbox. Default: common package registries
<code>sandbox.network.httpProxyPort</code>	number	Custom HTTP proxy port for sandbox network filtering
<code>sandbox.network.socksProxyPort</code>	number	Custom SOCKS proxy port for sandbox network filtering
<code>sandbox.enableWeakerNestedSandbox</code>	boolean	Enable weaker sandbox for unprivileged Docker environments (Linux/WSL2 only). Reduces security. Default: false

Complete sandbox example:

```
{
  "sandbox": {
    "enabled": true,
    "autoAllowBashIfSandboxed": true,
    "excludedCommands": ["docker"],
    "network": {
      "allowedDomains": ["registry.npmjs.org", "api.github.com"],
      "allowLocalBinding": true
    }
  },
  "permissions": {
    "deny": [
      "Read(.env)",
      "Read(.env.*)",
      "Read(secrets/**)"
    ]
  }
}
```

```

        "Write(.env)",
        "Write(.env.*)",
        "Write(secrets/**)"
    ]
}
}

```

Attribution Settings

Key	Type	Description
attribution.commit	string	Custom git commit attribution text. Supports \n for multi-line trailers
attribution.pr	string	Custom PR description attribution text. Set to "" to disable PR attribution

Default commit attribution:

Generated with Claude Code

Co-Authored-By: Claude <noreply@anthropic.com>

Customization example:

```
{
  "attribution": {
    "commit": "Generated with AI\n\nCo-Authored-By: AI <ai@example.com>",
    "pr": ""
  }
}
```

Plugin Settings

Key	Type	Description
enabledPlugins	string[]	Active plugin identifiers in plugin-name@marketplace-name format
extraKnownMarketplaces	object[]	Additional plugin marketplace sources with explicit trust boundaries

Plugin hooks are defined in `hooks/hooks.json` within the plugin root directory. Hooks include an optional top-level `description` field. Plugin scripts reference their own directory via the `#{CLAUDE_PLUGIN_ROOT}` variable.

MCP Server Approval

Key	Type	Description
enableAllProjectMcpServers	boolean	Automatically approve all MCP servers defined in project <code>.mcp.json</code> files
enabledMcpjsonServers	string[]	Approve specific MCP servers from <code>.mcp.json</code> (e.g., ["memory", "github"])
disabledMcpjsonServers	string[]	Reject specific MCP servers from <code>.mcp.json</code>

Managed-Only Settings

Key	Type	Description

strictKnownMarketplaces	boolean	Restrict plugins to approved marketplaces only
allowManagedHooksOnly	boolean	Block user, project, and plugin hooks; only managed and SDK hooks run
allowManagedPermissionRulesOnly	boolean	Prevent user/project settings from defining allow, ask, or deny permission rules
disableAllHooks	boolean	Disable all hooks and any custom status line
allowedMcpServers	object[]	Allowlist of MCP servers users can configure. Undefined = no restrictions. Empty array = lockdown. Denylist takes precedence
deniedMcpServers	object[]	Denylist of MCP servers that are explicitly blocked across all scopes

Settings Scope Hierarchy

Priority	Scope	Location	Shared	Overridable
1 (highest)	Managed	System directories or server-managed (IT-deployed)	Org-wide	No
2	CLI	Command-line flags	No	Session only
3	Local	.claude/settings.local.json	No (gitignored)	By managed/CLI
4	Project	.claude/settings.json	Yes (VCS)	By managed/CLI/local
5 (lowest)	User	~/.claude/settings.json	No	By all above

When to Use Each Scope

Scope	Use For
Managed	Security policies enforced org-wide, API key helpers, forced login methods, MCP allowlists/denylists, sandbox requirements
User	Personal preferences (model, language, spinner verbs, reduced motion), personal API keys, user-level permissions
Project	Shared team standards (permissions, hooks, MCP servers), coding conventions, attribution settings
Local	Personal overrides for a specific project, experimental settings, credentials you do not want committed

How Scopes Interact

Settings merge across scopes. For most keys, higher-priority scopes override lower ones. Permission rules follow a specific evaluation order: **deny rules first, then ask, then allow** — the first matching rule wins, evaluated across all scopes.

Server-Managed Settings

Organizations without device management infrastructure can use server-managed settings, which deliver configurations from Anthropic's servers for Claude for Enterprise customers. These function identically to file-based managed settings and cannot be overridden by user or project settings.

Environment Variables

Claude Code supports over 70 environment variables. The following are the most commonly used, grouped by function.

API and Authentication

Variable	Description
ANTHROPIC_API_KEY	API key for direct Anthropic access
CLAUDE_CODE_USE_BEDROCK	Enable cloud provider Bedrock integration (1)
CLAUDE_CODE_USE_VERTEX	Enable cloud provider Vertex AI integration (1)
ANTHROPIC_BASE_URL	Custom API base URL
CLAUDE_CODE_AUTH_TOKEN	Override authentication token

Model Configuration

Variable	Description
ANTHROPIC_MODEL	Override default model
ANTHROPIC_SMALL_FAST_MODEL	Model for lightweight operations (Haiku tasks)
CLAUDE_CODE_MAX_TOKENS	Maximum tokens per response
CLAUDE_CODE_MAX_THINKING_TOKENS	Maximum tokens for thinking/reasoning

Proxy and Network

Variable	Description
HTTP_PROXY / HTTPS_PROXY	HTTP(S) proxy URL
NO_PROXY	Comma-separated bypass domains
ANTHROPIC_PROXY_URL	Anthropic-specific proxy URL

Context and Compaction

Variable	Description
CLAUDE_AUTOCOMPACT_PCT_OVERRIDE	Auto-compaction threshold percentage (default: 95)
CLAUDE_CODE_CONTEXT_LIMIT	Maximum context window size

Feature Flags

Variable	Description
ENABLE_TOOL_SEARCH	Enable/disable MCP tool search and deferred loading
CLAUDE_CODE_ENABLE_AGENT_TEAMS	Enable experimental agent teams
CLAUDE_CODE_DISABLE_NONESSENTIAL_TRAFFIC	Disable telemetry and non-essential network calls

Headless/CI

Variable	Description

CLAUDE_CODE_ENTRYPOINT	Set to <code>cli</code> for headless environments
CLAUDE_CODE_BUDGET_TOKENS	Token budget for headless sessions
CLAUDE_CODE_MAX_TURNS	Maximum conversation turns in headless mode

Output and Formatting

Variable	Description
CLAUDE_CODE_OUTPUT_FORMAT	Output format for headless mode (<code>text</code> , <code>json</code> , <code>stream-json</code>)
CLAUDE_CODE_HIDE_TOOL_OUTPUT	Suppress tool output display in terminal

Debugging and Diagnostics

Variable	Description
CLAUDE_CODE_DEBUG	Enable debug logging (1)
CLAUDE_CODE_LOG_LEVEL	Set logging verbosity level
CLAUDE_CODE_VERBOSE	Enable verbose output

Session and Behavior

Variable	Description
CLAUDE_CODE_SKIP_CLAUDE_MD	Skip loading CLAUDE.md files (1)
CLAUDE_CODE_THINKING_MODE	Thinking mode (enabled, disabled, budget)
CLAUDE_CODE_THINKING_BUDGET	Token budget for extended thinking

Telemetry

Variable	Description
CLAUDE_CODE_DISABLE_TELEMETRY	Disable usage telemetry (1)
CLAUDE_CODE_TELEMETRY_ENDPOINT	Custom telemetry collection endpoint

Tools Inventory

Tool	Description
Bash	Execute shell commands. Working directory persists; environment does not.
Read	Read file contents. Supports images, PDFs, notebook files.
Write	Create or overwrite files.
Edit	Exact string replacement in files. Requires unique match.
Glob	File pattern matching. Returns paths sorted by modification time.

Grep	Content search with regex support. File type filtering, context lines.
WebFetch	Fetch and process URL content. Auto-upgrades HTTP to HTTPS.
WebSearch	Web search with domain filtering. Returns formatted results.
Task	Spawn subagents with isolated context. Supports foreground/background.
Skill	Invoke skill-based workflows within conversation.
NotebookEdit	Edit notebook cells (.ipynb). Replace, insert, or delete.
LSP	Language Server Protocol operations. Diagnostics, go-to-definition, references.
TeamCreate	Create agent team (experimental).
TaskCreate	Create task in agent team task list.
TaskUpdate	Update task status in agent team.
TaskList	List tasks for an agent team.
SendMessage	Send message between agent team members.
TeamDelete	Delete an agent team.
TodoRead	Read current task/todo list.
TodoWrite	Write/update task/todo items.

Hook Configuration Schema

Common Input Fields

Every hook event receives these fields on stdin as JSON:

Field	Description
session_id	Current session identifier
transcript_path	Path to conversation JSON file
cwd	Current working directory when the hook is invoked
permission_mode	Active permission mode: "default", "plan", "acceptEdits", "dontAsk", Or "bypassPermissions"
hook_event_name	Name of the event that fired

Each event adds its own fields (e.g., `tool_name`, `tool_input` for PreToolUse).

Event Types

Event	Fires	Can Block?	Event-Specific Fields
PreToolUse	Before any tool executes	Yes	<code>tool_name</code> , <code>tool_input</code>
PostToolUse	After any tool executes	Yes (via decision)	<code>tool_name</code> , <code>tool_input</code> , <code>tool_output</code>

PostToolUseFailure	When a tool execution fails	Yes (via decision)	tool_name, error, is_interrupt
UserPromptSubmit	Before processing a user prompt	Yes	prompt
Notification	On status notifications (permission_prompt, idle_prompt, auth_success, elicitation_dialog)	No	message, type
Stop	When Claude stops generating	Yes (via decision)	stop_hook_active
SessionStart	At session initialization	No	source (startup, resume, clear, compact), model, optionally agent_type
SessionEnd	At session termination	No	reason (clear, logout, prompt_input_exit, bypass_permissions_disabled, other)
PermissionRequest	Before showing permission dialog	Yes	tool_name, permission details; supports updatedInput, updatedPermissions
PreCompact	Before context compaction	No	trigger (manual OR auto), custom_instructions
PostCompact	After context compaction	No	Compacted context
SubagentStart	When a subagent spawns	No	agent_id, agent_type
SubagentStop	When a subagent completes	Yes (via decision)	agent_id, agent_type, agent_transcript_path, stop_hook_active
TeammateIdle	When agent team teammate is about to go idle	Yes (exit code 2)	teammate_name
TaskCompleted	When a task is marked completed	Yes (exit code 2)	task_id, task_subject

Handler Types (3)

Command Handler

```
{
  "hooks": {
    "PreToolUse": [
      {
        "matcher": "Bash",
        "handler": {
          "type": "command",
          "command": "\"$CLAUDE_PROJECT_DIR\"/.claude/hooks/validate.sh"
        }
      }
    ]
  }
}
```

Executes a shell command. Receives event JSON on stdin. Communicates results via exit codes and stdout JSON.

Prompt Handler

```
{
  "hooks": {
    "PreToolUse": [
      {
        "matcher": "Write(*.sql)",
        "handler": {
          "type": "prompt",
          "prompt": "Review this SQL file write for safety. Allow or deny."
        }
      }
    ]
  }
}
```

Sends context plus prompt to a Claude model (Haiku by default) for single-turn evaluation. Returns structured `{"ok": true/false, "reason": "..."}` decision. Uses `$ARGUMENTS` placeholder for event data.

Agent Handler

```
{
  "hooks": {
    "PostToolUse": [
      {
        "matcher": "Bash",
        "handler": {
          "type": "agent",
          "agent": "security-reviewer",
          "tools": ["Read", "Grep"]
        }
      }
    ]
  }
}
```

Spawns a subagent with multi-turn tool access for complex verification. Up to 50 turns.

Handler Common Fields

Field	Required	Description
timeout	No	Timeout in seconds. Default: 10 (command), 30 (prompt), 60 (agent)
statusMessage	No	Custom spinner message displayed while hook runs
once	No	If true, runs only once per session then is removed. Skills only
async	No	If true, runs in background without blocking. Output delivered on next turn via <code>systemMessage</code> or <code>additionalContext</code> . Cannot block actions

Exit Code Behavior

Exit Code	Meaning	Effect

0	Success	Action proceeds. Stdout is parsed for JSON output fields. Stdout shown in verbose mode (ctrl+o) except for UserPromptSubmit and SessionStart where it becomes context
2	Blocking error	Action is blocked. Stdout/JSON ignored. Stderr text fed back to Claude as error message
Other	Non-blocking error	Stderr shown in verbose mode. Execution continues

JSON Output Fields (on exit 0)

Field	Type	Effect
continue	boolean	If false, stops further hook processing for this event
stopReason	string	Message shown when stopping
suppressOutput	boolean	If true, suppresses default hook output
systemMessage	string	Message injected into conversation as system context
additionalContext	string	Extra context added to the event
updatedInput	object	Modified tool input before execution (PreToolUse, PermissionRequest)
updatedPermissions	object	Modified permission settings (PermissionRequest; equivalent to "always allow")
permissionDecision	string	"allow", "deny", or "ask" for PreToolUse hooks (inside hookSpecificOutput)
decision	string	"block" to stop the action (UserPromptSubmit, PostToolUse, PostToolUseFailure, Stop, SubagentStop)

Environment Variables in Hooks

Variable	Description
\$CLAUDE_PROJECT_DIR	Project root directory. Wrap in quotes for paths with spaces
\${CLAUDE_PLUGIN_ROOT}	Plugin root directory, for scripts bundled with a plugin
\$CLAUDE_ENV_FILE	(SessionStart only) File path for persisting environment variables. Write export statements to this file to make variables available in all subsequent Bash commands

CLAUDE_ENV_FILE example:

```
#!/bin/bash
if [ -n "$CLAUDE_ENV_FILE" ]; then
  echo 'export NODE_ENV=production' >> "$CLAUDE_ENV_FILE"
  echo 'export DEBUG_LOG=true' >> "$CLAUDE_ENV_FILE"
fi
exit 0
```

Matcher Syntax

- Exact tool name: Bash , Write , Edit
- Tool with specifier: Bash(npm *) , Write(*.env)
- MCP tools: mcp__servername__toolname
- Regex patterns supported in specifiers

- * matches any tool

Hook Decision Control (PreToolUse)

PreToolUse uses `hookSpecificOutput.permissionDecision` :

Decision	Effect
"allow"	Proceed without user prompt
"deny"	Block execution; reason shown to Claude
"ask"	Show permission prompt (can combine with <code>updatedInput</code> to show modified input)
Omitted / no output	Fall through to normal permission evaluation

Hook Decision Control (PermissionRequest)

Decision	Effect
"approve"	Auto-approve without user prompt
"deny"	Auto-deny without user prompt
Omitted / no output	Show normal permission prompt

Supports `updatedInput` and `updatedPermissions` (equivalent to "always allow" options).

Hook Decision Control (TeammateIdle / TaskCompleted)

These events use exit codes only, not JSON decision control:

Exit Code	Effect
0	Allow the action (teammate goes idle, task marked complete)
2	Block the action; stderr fed back as feedback to continue working

Hook Security Model

Hooks are snapshotted at session startup. If hooks are modified externally during a session, Claude Code warns and requires review in the `/hooks` menu before changes apply. This prevents malicious mid-session changes.

MCP Configuration Format

.mcp.json (Project Root)

```
{
  "mcpServers": {
    "server-name": {
      "command": "node",
      "args": ["path/to/server.js"],
      "env": {
        "API_KEY": "${ENV_VAR_NAME}"
      },
      "cwd": "./optional-working-directory"
    }
  }
}
```

```

},
"remote-server": {
  "url": "https://mcp.example.com/sse",
  "headers": {
    "Authorization": "Bearer ${TOKEN}"
  }
}
}
}
}

```

Server Configuration Fields

Field	Type	Required	Description
command	string	Yes (local)	Executable to launch
args	string[]	No	Command arguments
env	object	No	Environment variables (supports \${VAR} interpolation)
cwd	string	No	Working directory for server process
url	string	Yes (remote)	SSE endpoint for remote servers
headers	object	No	HTTP headers for remote connections

Managed MCP Settings

Key	Type	Description
mcpServers.allowlist	string[]	Only these servers may be configured
mcpServers.denylist	string[]	These servers are blocked

Plugin Marketplace Sources

Source Type	Format
Git hosting platform	githost:org/repo
Git URL	git:https://example.com/repo.git
npm	npm:package-name
URL	url:https://example.com/plugin.tar.gz
File	file:/path/to/plugin
Directory	directory:/path/to/plugin-dir
Host pattern	host:*.internal.company.com

Full Settings.json Example

```
{
  "$schema": "https://json.schemastore.org/clause-code-settings.json",
```

```
"permissions": {
    "allow": [
        "Bash(npm test)",
        "Bash(npm run lint)",
        "Bash(git *)",
        "Read",
        "Glob",
        "Grep"
    ],
    "deny": [
        "Write(.env)",
        "Write(.env.*)",
        "Read(secrets/**)"
    ]
},
"env": {
    "NODE_ENV": "development",
    "LOG_LEVEL": "info"
},
"model": "claude-opus-4-6",
"language": "english",
"autoUpdatesChannel": "stable",
"showTurnDuration": true,
"spinnerVerbs": {
    "mode": "append",
    "verbs": ["Pondering", "Crafting"]
},
"sandbox": {
    "enabled": true,
    "autoAllowBashIfSandboxed": true,
    "excludedCommands": ["docker", "git"],
    "network": {
        "allowedDomains": ["registry.npmjs.org", "api.github.com"]
    }
},
"attribution": {
    "commit": "Generated with Claude Code\n\nCo-Authored-By: Claude <noreply@anthropic.com>",
    "pr": "Generated with Claude Code"
},
"hooks": {
    "PreToolUse": [
        {
            "matcher": "Bash",
            "handler": {
                "type": "command",
                "command": "\"$CLAUDE_PROJECT_DIR\"/.claude/hooks/validate-bash.sh"
            }
        }
    ]
}
```

Appendix C: Troubleshooting

Organized by symptom. No prose. Scan for your problem.

Bash Environment Variables Not Persisting

Symptom: Environment variables set in one bash command are not available in subsequent commands. Scripts that depend on `export` fail silently.

Cause: Each bash command runs in a fresh shell environment. Working directory persists; everything else (variables, aliases, shell functions) does not.

Fix:

- Chain dependent commands in a single bash call using `&&`
 - Write variables to a file and source it: `echo "export FOO=bar" > .env && source .env && echo $FOO`
 - Use CLAUDE.md to document required environment variables so Claude sets them in each command
-

Context Exhaustion (Claude Forgets Instructions)

Symptom: Claude stops following CLAUDE.md rules, produces inconsistent output, drops tasks from a multi-step plan, or "forgets" decisions made earlier in the conversation.

Cause: Context window is full or near-full. Auto-compaction has summarized earlier conversation, losing detail. Instructions not in CLAUDE.md are not preserved through compaction.

Fix:

- Move critical rules to the "Compact Instructions" section of CLAUDE.md (survives compaction)
 - Reduce always-on context costs: trim CLAUDE.md, remove unused MCP servers, convert reference material to skills
 - Start a new session for new tasks instead of continuing indefinitely
 - Use `/compact` manually before critical operations to reclaim space with controlled summarization
 - Delegate verbose operations to subagents (isolated context windows)
-

Auto-Compaction Triggering Too Late (or Too Early)

Symptom: Session degrades before compaction kicks in, or compaction triggers when you still have space to work.

Cause: Default compaction threshold is ~95% of context capacity.

Fix:

- Set `CLAUDE_AUTOCOMPACT_PCT_OVERRIDE` environment variable
 - Lower value (e.g., `80`) = earlier compaction, more headroom, less detail preserved
 - Higher value (e.g., `98`) = later compaction, more detail, risk of degradation before trigger
-

claude doctor Diagnostics

Symptom: Claude Code behaves unexpectedly, fails to start, or connections fail.

Run: `claude doctor`

What it checks:

- Authentication status and token validity
- Network connectivity to API endpoints
- Configuration file syntax errors
- MCP server availability
- Environment variable conflicts
- Binary version and update status

When to use: First step for any issue that is not obviously a context or prompt problem.

Claude Repeats the Same Error

Symptom: Claude makes the same mistake across sessions -- wrong import path, incorrect test command, deprecated API usage, wrong module reference.

Cause: The error pattern is not captured in persistent context. Each session starts fresh without knowledge of previous failures.

Fix:

- Add a targeted instruction to CLAUDE.md addressing the specific error
- Example: "Always import from @app/db/utils, never @app/utils/db (deprecated path)"
- Be specific: describe the wrong behavior AND the correct behavior
- Place in project-level CLAUDE.md so all team members benefit

Checkpoint Restore Does Not Undo Changes

Symptom: Esc-Esc rewind does not restore expected file state. Files modified by bash commands (`rm`, `mv`, `cp`, `sed`) are not reverted.

Cause: Checkpoints only track file edits made through Claude Code's Write/Edit tools. Changes made via bash commands, external editors, or other concurrent sessions are not tracked.

Fix:

- Use git for full recovery: `git checkout -- <file>` or `git stash`
- Commit frequently before risky operations
- Prefer Claude Code's built-in file tools over bash for file modifications when possible
- For destructive operations, ask Claude to create a backup first

Deciding: Restart Fresh vs. Correct Mid-Stream

Symptom: Claude is going in the wrong direction, producing incorrect code, or following a bad plan. Unclear whether to redirect or start over.

When to restart fresh:

- Claude has filled context with incorrect exploration
- The fundamental approach is wrong, not just the details
- Multiple correction attempts have not converged
- Context is near compaction threshold

When to correct mid-stream:

- The approach is correct but details are wrong
- You can describe the fix concisely
- Context has plenty of remaining capacity
- Claude has built up useful understanding that would be lost

Pattern: Commit current state to git before deciding. If you restart, the commit preserves any useful fragments. If you correct, the commit gives you a rollback point.

WebSocket Scaling Limitations

Symptom: Real-time features generated by Claude Code (chat interfaces, live dashboards, streaming feeds) fail under load or drop connections.

Cause: Claude Code generates correct WebSocket scaffolding but does not automatically handle production-scale concerns: connection pooling, backpressure, reconnection logic, heartbeat keepalives.

Scaling thresholds:

Concurrency	Recommended Stack
MVP (<50 concurrent)	Python + FastAPI is sufficient
Scaling (>100 concurrent)	Go or Node.js required; Python struggles at this level
High scale (1000+)	Go strongly preferred

Fix:

- Explicitly request connection pooling and reconnection logic in prompts
 - Add load-testing requirements to acceptance criteria
 - Manually review generated WebSocket code for: max connection limits, heartbeat intervals, graceful degradation, memory leaks from unclosed connections
 - Review backpressure handling: Claude can scaffold it, but verify against your latency SLA
 - Consider using a dedicated real-time messaging library rather than raw WebSocket handling
-

FIX Protocol Limitations

Symptom: Generated FIX protocol message builders compile but produce incorrect or incomplete messages, or fail conformance testing.

Cause: FIX protocol is highly complex with domain-specific semantics. Claude Code can generate FIX message builders but lacks the deep domain knowledge required for production-grade implementations.

Fix:

- Treat Claude-generated FIX code as scaffolding only
 - Have a FIX domain expert review all generated message builders
 - Test against a FIX conformance test suite before production use
 - Use Claude Code for wrapping standard REST/WebSocket broker APIs instead when possible
-

MCP Server Disconnection Recovery

Symptom: MCP tools stop working mid-session. Tool calls fail silently or return errors. No warning that the server disconnected.

Cause: MCP server connections can drop without notification. Tools disappear from Claude's available tool set.

Fix:

- Run `/mcp` to check server status and token costs
- Restart the MCP server if it crashed (check its process/logs)
- Start a new Claude Code session to re-establish connections

- For persistent issues, check MCP server logs for timeout or memory errors
 - Add health-check logic to custom MCP servers
 - Consider PreToolUse hooks that verify MCP availability before critical operations
-

Session Lost After Crash or Terminal Close

Symptom: Work from a previous session is gone. Claude does not remember what it was doing.

Cause: Each session is independent. Context exists only within a session. Closing the terminal or crashing loses unsaved context.

Fix:

- Resume the previous session: `claude -c` (continue most recent) or `claude -r <session-id>`
 - Use `claude --resume` to browse and select from recent sessions
 - For critical multi-session work, persist state in files: specs, task lists (`.claude/tasks/`), CLAUDE.md updates
 - Name important sessions with `/rename` for easier recovery
-

Subagent Results Consuming Too Much Main Context

Symptom: After running multiple subagents, the main conversation's context fills rapidly. Compaction triggers unexpectedly.

Cause: Each subagent returns its results to the main context. Multiple subagents with detailed outputs consume main context quickly.

Fix:

- Instruct subagents to return concise summaries, not full details
 - Reduce the number of concurrent subagents
 - Write subagent results to files instead of returning them in conversation
 - Run `/compact` after processing subagent results you no longer need in active context
-

Claude Produces Over-Complex Solutions

Symptom: Claude generates elaborate abstractions, unnecessary design patterns, or complex architectures for simple problems.

Cause: Default behavior tends toward comprehensive solutions. Without constraints, Claude adds flexibility, extensibility, and abstraction layers.

Fix:

- Add to CLAUDE.md: "Prefer the simplest solution that meets requirements. No unnecessary abstractions."
 - Be explicit about scope: "This is a one-off script, not a library"
 - Specify constraints: "No new dependencies", "Single file", "Under 100 lines"
 - Request Claude to justify complexity when it appears
 - Interrupt mid-execution and ask "why are you doing this? Try something simpler"
-

Verifying Proxy or Gateway Configuration

Symptom: API calls fail, authentication errors, or unexpected model routing when using a corporate proxy or LLM gateway.

Cause: Proxy URL, base URL, or authentication headers are misconfigured or conflicting.

Fix:

- Run `/status` to verify current proxy and gateway configuration
 - Check `ANTHROPIC_BASE_URL`, `HTTP_PROXY`, and `HTTPS_PROXY` environment variables
 - Ensure proxy certificates are trusted by the system
 - For LLM gateways handling authentication, set the appropriate skip-auth variable (`CLAUDE_CODE_SKIP_BEDROCK_AUTH=1` or `CLAUDE_CODE_SKIP_VERTEX_AUTH=1`)
-

Debugging with Verbose Output

Symptom: Unclear why Claude made a decision, why a tool call failed, or what data hooks received.

Cause: Default output suppresses detailed tool usage, hook execution, and internal processing information.

Fix:

- Press `Ctrl+O` to toggle verbose output, showing detailed tool usage and execution
 - Run `claude --debug` to see hook execution details including which hooks matched, exit codes, and output
 - Check hook stderr output (visible in verbose mode for non-blocking errors)
 - For hooks specifically, verify JSON parsing by examining raw stdout with `Ctrl+O` enabled
-

Hook Not Firing

Symptom: A configured hook does not execute when expected. No output, no error.

Cause: Common causes include non-executable scripts, matcher not matching the tool name, or hooks modified after session startup without review.

Fix:

- Ensure the script is executable: `chmod +x script.sh`
 - Verify the shebang line: first line should be `#!/bin/bash` or `#!/usr/bin/env bash`
 - Check the matcher matches the exact tool name (e.g., `Bash`, not `bash`)
 - If hooks were modified during the session, review changes in the `/hooks` menu
 - Check if `disableAllHooks` is set to `true` in settings
 - Use `claude --debug` to see which hooks are evaluated and matched
-

Hook JSON Parsing Errors

Symptom: Hook returns exit code 0 but JSON output is not processed. Claude ignores the hook's decision.

Cause: Non-JSON text in stdout (often from shell profile scripts) interferes with JSON parsing.

Fix:

- Ensure stdout contains only the JSON object on exit 0
- Redirect any non-essential output to stderr: `echo "debug info" >&2`
- Check if shell profile (`.bashrc`, `.zshrc`) prints text on startup that pollutes stdout
- Test the script manually and inspect raw output