# 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

---

# 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

```rust
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(),
```

```rust
                                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

```rust
        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**:

```rust
        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`

```rust
    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

```rust
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

```rust
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(anyhow!("Agent execution failed: {}", e)),
        Err(_) => Err(anyhow!("Agent timeout after {:?}",
task.timeout)),
    }
}
```

**Timeout Handling**:

```rust
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`

```rust
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**:

```rust
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**:

```rust
fn compute_verdict(
    agent_outputs: &[AgentOutput],
    synthesis: &SynthesisResult,
) -> Result<ConsensusVerdict> {
    // Count present agents
```

```rust
        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

```rust
    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::<Vec<_>>();

let tasks = join_all(tasks).await;  // Await all at once
```

### 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

```rust
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**:

```rust
// 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`

---

# 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:                                              │
│ ├─ spec.md (2.3 KB)                                     │
│ ├─ plan.md (5.7 KB)                                     │
│ ├─ tasks.md (3.2 KB)                                    │
│ └─ evidence/ (12 files, 450 KB)                         │
├─────────────────────────────────────────────────────────┤
│ Quality Gates:                                          │
│ ├─ 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

- <u>Native Operations</u> - 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`

---

# 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

```
Layer 1: Agent Selection & Routing (tier-based)
    ↓
Layer 2: Agent Orchestration (sequential vs parallel)
    ↓
Layer 3: MCP Consensus Coordination (retry + fallback)
    ↓
Layer 4: Consensus Synthesis (verdict computation)
    ↓
Layer 5: Response Caching (SQLite persistence)
```

**Core Files**:

| File | LOC | Purpose |
| --- | --- | --- |
| `routing.rs` | 221 | Command dispatch & ACE routing |
| `agent_orchestrator.rs` | 2,208 | Sequential/parallel spawning, execution control |
| `consensus_coordinator.rs` | 194 | MCP retry logic, cost summary |
| `consensus.rs` | 1,160 | Artifact collection, verdict synthesis |
| `consensus_db.rs` | 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`

```rust
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`

```rust
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`

```rust
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`

```rust
pub async fn execute_parallel_consensus(
```

```rust
    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

```rust
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 {
```

```rust
            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.)

```rust
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

```rust
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)
```

```rust
        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

```rust
    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
```

```rust
    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

```rust
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)

```rust
pub fn detect_conflict(
    output_a: &AgentOutput,
```

```rust
    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**:

```json
{
  "output_a": "Recommend 3-layer architecture...",
  "output_b": "Suggest monolithic approach...",
  "aspect": "architecture"
}
```

**Output**:

```json
{
  "has_conflict": true,
  "conflict_description": "Agent A recommends layered, Agent B
monolithic",
  "severity": "Critical"
}
```

---

## Conflict Resolution

**Strategy**: User decision required for Critical conflicts

```rust
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 =
```

```rust
    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
  [4] Abort pipeline
```

---

# Layer 5: Response Caching

## SQLite Schema (SPEC-KIT-072)

**Location**: codex-rs/tui/src/chatwidget/spec_kit/consensus_db.rs:50-150

```sql
-- 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

```rust
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 (6.6× 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)

```rust
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 (?1, ?2, ?3, ?4, ?5, ?6, ?7, ?8, ?9)",
            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**

```rust
let agents = select_agents_for_tier(CommandTier::Tier2Multi,
"validate");
    // Returns: [gemini-flash, claude-haiku, gpt5-medium]
```

**Step 2: Parallel Execution**

```rust
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**

```rust
let artifacts = collect_consensus_artifacts("SPEC-KIT-070",
"validate").await?;
    // Tries SQLite (8.7ms) → MCP (15ms) → files (50ms)
```

**Step 4: MCP Synthesis**

```rust
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**

```rust
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`

---

# 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 | - | $ |
| **Tasks** | 1 (Single) | 1 low | $0.10 | - | - | $ |
| **Implement** | 2 (Code) | 2 specialist | $0.11 | - | - | $ |
| **Validate** | 2 (Multi) | 3 cheap | $0.30 | $0.05 | - | $ |
| **Audit** | 3 (Premium) | 3 premium | $0.75 | $0.05 | - | $ |
| **Unlock** | 3 (Premium) | 3 premium | $0.75 | $0.05 | - | $ |
| **Quality Gates** | 0 (Native) | 0 | $0.00 | $0.15-0.20 | $0.15-0.20 | - |
| **TOTAL** | - | - | **$2.31** | **$0.20** | **$0.19** | - |

### 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-pro** | gemini-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-** | OpenAI | $0.0006 | $0.0006 | Code generation, |

**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`

```rust
pub struct CostTracker {
    pub total_cost: f64,
```

```rust
        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 × $0.20 = $0.80
```

**After** (Plan stage):

```
3 agents: 3 × $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 × $2.70 = $27.00/month
Annual: $324/year
```

**Use Cases**: - Personal projects - Small teams - Experimental features

---

### Medium Usage (50 SPECs/month)

```
50 SPECs × $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**: - <u>Agent Orchestration</u> - Multi-agent coordination details - <u>Template System</u> - PRD and doc templates - <u>Workflow Patterns</u> - 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

---

# 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:

```json
{
  "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"],
```

```json
        "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

```json
    {
      // 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

```json
    {
      // 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"
```

```json
        },
        "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

```json
        {
          // 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"
          },
```

```json
    "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**

```json
{
  // 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"},
```

```json
        {"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

```json
  {
    // 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?
    },
```

```json
    "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",
```

```json
      "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`

```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**: - <u>Cost Tracking</u> - Per-stage cost breakdown and analysis - <u>Agent Orchestration</u> - Multi-agent coordination - <u>Workflow Patterns</u> - 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)

---

# 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

```
Decision Flow:
    ↓
Is the task deterministic? (same input → same output)
    ├─ YES → Native operation ($0, <1s)
    └─ NO  → Agent-based ($0.10-0.80, 3-10min)
          ↓
      Does it require reasoning/judgment?
          ├─ YES → Agents (creativity, analysis)
          └─ NO  → Rethink if native is possible
```

### 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

```rust
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)?;
```

```rust
        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

```rust
    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

```rust
    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

```rust
        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)
```

# /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

```rust
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|might|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|com
            ).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

```rust
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

```rust
    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

```rust
    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

```rust
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`

```rust
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

```rust
    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

```rust
    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!("{} 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

```rust
    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![
```

```rust
                        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

```rust
    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),
```

```rust
                    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`

```rust
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)

```rust
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!("-{:.1} points", points),
                    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!("-{:.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!("-{:.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 {
```

```rust
                    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!("-{:.1} points", coverage_deduction),
                    suggestion: "Add acceptance criteria for all
requirements".to_string(),
                });
            }

            score.max(0.0)
        }
```

---

## Consistency Scoring (20 points)

```rust
        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!("-{:.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)?;
```

```rust
        // 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- | Time | Cost |
| --- | --- | --- | --- | --- |

| | | Based | Saved | 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`

# 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

```
/speckit.auto (user command)
      ↓
pipeline_coordinator.rs (advancement loop)
      ↓
├── Guardrail Execution (stage validation)
├── Agent Orchestration (multi-agent consensus)
├── Consensus Coordination (MCP integration)
├── Quality Gate Handler (checkpoint validation)
└── State Persistence (SQLite + evidence files)
```

**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

```rust
#[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)
                                   ↓
                         Quality gate required?
                                   ↓
                         ┌─────────┴─────────┐
                         NO                 YES
                          ↓                  ↓
                     Next stage
QualityGateExecuting
                                             ↓

QualityGateProcessing
                                             ↓

QualityGateValidating
```

```
                                                       ↓
                                          Pass?
QualityGateAwaitingHuman

                                          ↓              ↓
                                       Next stage    (User
decision)
```

## SpecAutoState Struct

**Location**: `codex-rs/tui/src/chatwidget/spec_kit/state.rs:47-110`

```rust
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 | Purp |
|-------|-------|------|--------|------|------|------|
| 0 | **Plan** | 2 (Multi) | 3 | ~$0.35 | 10-12min | Work breakdo |
| 1 | **Tasks** | 1 (Single) | 1 | ~$0.10 | 3-5min | Task decompo |
| 2 | **Implement** | 2 (Code) | 2 | ~$0.11 | 8-12min | Code generati |
| 3 | **Validate** | 2 (Multi) | 3 | ~$0.35 | 10-12min | Test stra |
| 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`

```rust
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

```rust
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

```rust
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`):

```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):

```rust
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**:

```rust
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):

```rust
pub struct ValidateLifecycleState {
    pub attempts: HashMap<String, ValidateAttempt>,
}

pub struct ValidateAttempt {
    pub payload_hash: String,     // SHA-256 of inputs
```

```rust
        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):

```rust
    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()
```

```rust
                .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

```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):

```rust
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 | To |
|-------|-----------|--------|-----------|--------------|----|
| **Plan** | 5s | 10min | 30s | 2min (AfterSpecify) | ~12 |
| **Tasks** | 5s | 3min | 10s | 1min (AfterTasks) | ~4m |
| **Implement** | 10s | 8min | 30s | - | ~9m |
| **Validate** | 5s | 10min | 30s | - | ~10 |
| **Audit** | 5s | 10min | 30s | - | ~10 |
| **Unlock** | 5s | 10min | 30s | - | ~10 |

**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`

---

# 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)?;
```

```rust
        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**:

```rust
    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**:

```rust
    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**:

```rust
    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

```rust
    fn submit_gpt5_validations(
        widget: &mut ChatWidget,
        majority_issues: &[QualityIssue],
        spec_id: &str,
```

```rust
            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**:

```json
        {
          "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**:

```rust
        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**:

```rust
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()?;
```

```rust
            answers.insert(question.id.clone(), answer);
        }

        Ok(answers)
    }
```

**Auto-Apply Changes**:

```rust
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**:

```rust
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()?
```

```rust
            .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**:

```json
{
  "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

```rust
    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!("docs/{}/spec.md",
spec_id));

        let content = std::fs::read_to_string(spec_path)?;

        let mut ambiguities = Vec::new();
```

```rust
        // 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**:

```rust
    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));
```

```rust
        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**:

```rust
    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
```

```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`

---

# 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)

### FR-003
```

**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

```json
{
  "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}",
```

```json
      "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

```json
{
  "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**:

```rust
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):

```toml
# .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**:

```rust
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**: - <u>Workflow Patterns</u> - 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`

---

# 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 Step-by-Step** | ~$2.70 | 1-3 days | High | Learning, review between stages |
| **3. Iterative Development** | Variable | Variable | High | Resume from failures |
| **4. Quality-Focused** | ~$2.70+ | 2-5 days | Highest | Security, compliance, critical |
| | | | | Simple |

| | | | | |
|---|---|---|---|---|
| **5. Cost-Optimized** | ~$0 | 2-8 hours | Medium | 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 ✓