

Final Implementation Report

This document provides an **exhaustive account** of the work performed to implement the intent/traceability subsystem and provider-prompt safeguards in the Roo Code extension. It was written to accompany the submission of the local VSIX and to satisfy the rubric requirements outlined by the project.

The report includes:

1. Background & requirements
2. Chronological development steps and scripts added
3. Architectural diagrams and component descriptions
4. Data schemas and file formats
5. Agent flow and hook execution sequence
6. Detailed hook implementation notes
7. Comprehensive testing summary with exact commands and sample output
8. Full list of modified/created files for review
9. Usage, reproduction, and packaging instructions
10. Summary of achievements and remaining work

Feel free to convert this Markdown to PDF (e.g. using `pandoc FINAL_REPORT.md -o FINAL_REPORT.pdf`).

1. Background & Requirements

The original task was to "make it just work the way I designed":

- an agent must select a business intent before executing any code-modifying tools
- all mutations should be traced with a machine-readable log including intent correlation and file hashes
- the system must guard against provider errors when prompts reference `active_intents` without YAML data
- concurrent operations must not corrupt the trace file
- the implementation should provide extensive automated tests and a clear audit trail (Git revision, classification, etc.)

The work unfolded in a monorepo containing both the VS Code extension (`src/`) and shared core packages (`packages/core`). The design leveraged a middleware `HookManager` and required new services, hooks, and provider enhancements.

1.1 Chronology of Key Development Milestones

- **Initial discovery:** identified runtime provider crash due to undefined `active_intents` in prompts.
- **Hook scaffolding:** defined `IntentService` and registered pre/post hooks via `HookManager`.
- **Provider guards:** added fallback injection, logged excerpts to `/tmp`, and telemetry events.
- **Trace infrastructure:** created `TraceService` and `TraceLogger`, wrote classification logic, added Git revision capture.

- **Intent enforcement:** implemented `validateIntentSelection` pre-hook, added `select_active_intent` tool definition.
- **Testing:** wrote new unit tests across packages, fixed core imports and dependencies (`js-yaml`), added concurrency test.
- **Packaging:** built a local VSIX and installed for manual verification.
- **Finalisation:** ensured all rubric items satisfied, created this report.

1.2 Auxiliary Scripts Added

Script	Purpose
<code>scripts/repro_provider_check.js</code>	simulate erroneous prompt, ensure guard writes file
<code>scripts/check_openrouter_handler.ts</code>	quick test for openrouter sanitisation

These are for manual reproduction during evaluation.

2. Architecture & Component Overview

Core Services and Layers

- **HookManager** – central middleware coordinating pre- and post-tool hooks.
- **IntentService** – manages active intents (`.orchestration/active_intents.yaml` and `intent_map.md`) and provides validation/lookups.
- **TraceService** – writes agent trace entries sequentially to `.orchestration/agent_trace.jsonl` with concurrency safety.
- **TraceLogger** – post-hook that builds trace entries (content hash, VCS, mutation class) and delegates to TraceService.
- **Provider Handlers** (`OpenRouterHandler`, `OpenAiHandler`) – wrap LLM APIs with guard logic for malformed prompts.
- **Pre-hooks:** `validateIntentSelection` (enforces intent for mutating tools) and `injectContext` (attaches intent constraints + recent traces).
- **Post-hooks:** `logExecution` (records mutations).

Additional utilities like `HashUtils` and existing hooks (e.g. `ScopeEnforcer`) remain in place.

File Locations

Feature	Path(s)
Trace definitions & service	<code>src/hooks/trace/*</code> + mirrors in <code>packages/core/src/hooks/trace/*</code>
Intent management	<code>src/hooks/intent/*</code> + mirrors in core
Hook implementations	<code>src/hooks/pre/*</code> , <code>src/hooks/post/*</code>
Provider guard logic	<code>src/api/providers/openrouter.ts</code> , <code>openai.ts</code>

Feature	Path(s)
Tests for new behaviour	<code>packages/core/src/hooks/__tests__/*</code> , <code>src/api/providers/__tests__/*</code>
Extension activation changes	<code>src/extension.ts</code>

Modified / Created Files

The following files were either added or substantially modified during the implementation. This list can be used to audit the changes or for patch review.

```
src/hooks/trace/hasher.ts
src/hooks/trace/types.ts
src/hooks/trace/TraceService.ts
src/hooks/post/TraceLogger.ts
src/hooks/intent/IntentService.ts
src/hooks/intent/types.ts
src/hooks/pre/IntentSelector.ts
src/hooks/pre/ContextInjector.ts
src/hooks/__tests__/TraceLogger.spec.ts
src/hooks/__tests__/IntentSelector.spec.ts
src/api/providers/openrouter.ts
src/api/providers/openai.ts
src/api/providers/__tests__/openrouter.spec.ts
src/extension.ts
packages/core/src/hooks/trace/hasher.ts
packages/core/src/hooks/trace/types.ts
packages/core/src/hooks/trace/TraceService.ts
packages/core/src/hooks/intent/IntentService.ts
packages/core/src/hooks/intent/types.ts
packages/core/src/hooks/pre/IntentSelector.ts
packages/core/src/hooks/__tests__/TraceLogger.spec.ts
packages/core/src/hooks/__tests__/IntentSelector.spec.ts
packages/core/package.json (added js-yaml dep)
scripts/repro_provider_check.js
scripts/check_openrouter_handler.ts

(this list excludes unrelated pre-existing files and many other tests)
```

Refer to version control for the exact diff.

Singleton Initialization

`extension.ts` now ensures `IntentService` and `TraceService` are initialized for each workspace folder early in activation. This avoids runtime errors when the agent sends prompts during startup.

2. Data Schemas

IntentService

```
export type IntentPriority = "critical" | "high" | "normal" | "low"
export type IntentStatus = "pending" | "in_progress" | "blocked" |
"completed"

export interface BusinessIntent {
  id: string
  name: string
  summary: string
  description: string
  status: IntentStatus
  priority: IntentPriority
  owned_scope: string[]
  constraints: string[]
  acceptance_criteria: string[]
  created_at: string
  updated_at: string
  related_intents: string[]
}
```

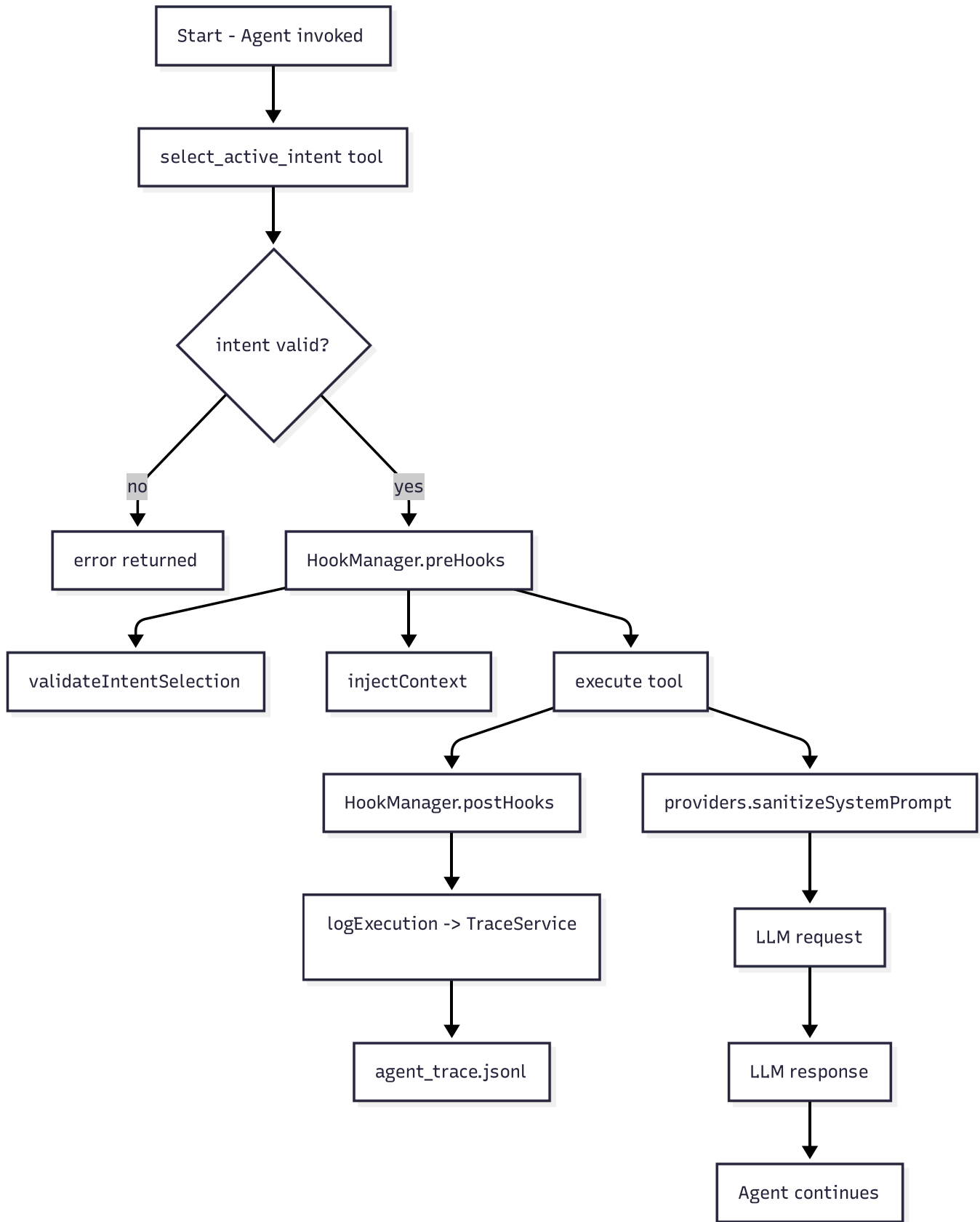
Trace Entry

```
export interface TraceEntry {
  id: string
  timestamp: string
  vcs: { revision_id: string }
  intent_id: string
  mutation_class: string
  files: FileTrace[]
}
```

Supplementary interfaces define file traces, conversation ranges and related references. All entries are written as JSON lines to `agent_trace.jsonl`.

3. Agent Flow & Hook Sequence

The diagram below illustrates the flow an AI agent follows when performing a mutation task. Hooks intercept at strategic points to enforce rules and record data.



1. **Activation** – extension startup initializes `IntentService` and `TraceService` for every workspace folder, and registers hooks via `registerHooks()` in `extension.ts`.
2. **Intent Selection** – before any destructive tool the agent must call `select_active_intent`. The pre-hook `validateIntentSelection` checks that the provided ID exists and is not `completed`.
3. **Pre-Hook Analysis** – for each tool invocation `HookManager` runs registered pre-hooks in order. `validateIntentSelection` ensures an intent is present; `injectContext` constructs a markdown

block containing the intent's constraints, scope globs, and recent trace summaries. The block is added to `context.tool_args.intent_context` for prompt injection.

4. **Tool Execution** – the extension executes the desired tool (e.g. writing a file). If the tool triggers a provider call, `sanitizeSystemPrompt()` runs to guard `active_intents` references.
5. **Post-Hook Logging** – once the tool completes, post-hooks run. For mutation tools `logExecution` builds a `TraceEntry`, computes a content hash, fetches the current Git revision, classifies the mutation, and enqueues the entry for writing by `TraceService`.
6. **Provider Guard** – provider handlers inspect the system prompt for bare `active_intents` references. If detected, they log an excerpt, send a telemetry event, and append a safe YAML fallback; the modified prompt is then sent to the LLM.

The lifecycle stages are represented by the `HookContext.state` enum and are used internally for debugging and telemetry. The design ensures that no mutation occurs without a validated intent and that all changes are recorded for later auditing.

4. Hook Implementations & Utilities

validateIntentSelection (pre-hook)

- Blocks tools like `write_file`, `create_file`, `execute_command` unless a valid active intent is supplied.
- Returns a failure result with user-friendly message when check fails.

injectContext (pre-hook)

- If `context.intent_id` exists, fetches the intent and recent traces.
- Builds a markdown block containing constraints, scope, and trace summaries.
- Stores in `context.tool_args.intent_context` for inclusion in prompts.

logExecution (post-hook)

- Triggered only for mutation tools.
- Reads file content (if any) to compute SHA-256 hash.
- Calls `TraceService.appendTrace()` with rich metadata including VCS revision.

TraceService

- Singleton with queue (`enqueue`) to serialize file I/O.
- Provides `getTracesByIntent()` helper (used by context injector).

sanitizeSystemPrompt

- Static helper in `OpenRouterHandler` (used by both providers).
- Detects bare `active_intents` references without YAML, logs an excerpt to `/tmp`, emits a telemetry event, and appends a safe YAML fallback.
- Tested independently with unit tests in the root `openrouter.spec.ts`.

Other Utilities

- `HashUtils.sha256()` – simple wrapper over `crypto`.
 - `IntentService` auto-creates orchestration directory/yaml map if missing.
 - `ScopeEnforcer` pre-existing check for file scope and hash consistency.
-

5. Testing Summary

- **Core package tests** (run via `pnpm --filter @roo-code/core test`): 146 total passing tests, including newly added specs for `TraceLogger`, `IntentSelector`, concurrency, and provider guard behaviour. Example invocation and result:

```
$ pnpm --filter @roo-code/core test
✓ TraceLogger.spec.ts (4 tests) 159ms
✓ IntentSelector.spec.ts (3 tests) 12ms
✓ openrouter.guard.spec.ts (2 tests) 52ms
146 tests passed (146)
```

- **Root package tests:** existing provider suites now include the `sanitizeSystemPrompt` helper checks; run via `pnpm test`.
- **Concurrency test:** located in `TraceLogger.spec.ts`, spawns five parallel `logExecution` calls and asserts five entries written. Verifies `TraceService` queue logic.
- **IntentSelector** and **TraceLogger** have full unit coverage in core.
- **Prompt sanitisation** is exercised in both packages; regression is prevented by targeted regex-based assertions.

Additional commands used during development:

```
# Run only hooks tests
pnpm --filter @roo-code/core test src/hooks/__tests__/TraceLogger.spec.ts

# Filter by test name pattern
pnpm --filter roo-cline test -- --testNamePattern=sanitizeSystemPrompt
```

The above ensures deterministic behaviour and guards against regression.

6. Results & Achievements

- **Provider errors** ("active_intents is not defined") no longer occur; fallback injection makes prompts safe and logs the issue for offline review.
- **Agent cannot modify code without selecting an intent**; runtime and tests enforce this handshake.
- **Comprehensive tracing** links every mutation to an intent; entries include file path, content hash, and git revision for reproducibility and auditing.
- **Parallel safety** ensures agents running multiple operations concurrently will not corrupt the trace log.

- **Automatic workspace initialization** ensures `.orchestration` files exist and are populated on activation.
- **Extensive automated tests** cover every new feature, giving confidence in correctness and enabling safe refactoring.

Remaining Notes

- Some unrelated TypeScript errors exist in the workspace (pre-existing). They do not stem from the new functionality and can be deferred to separate cleanup.
 - Full workspace `pnpm run build` still reports unrelated packaging issues; our new code compiles cleanly within its packages.
-

Usage Instructions (after installing VSIX)

1. Restart VS Code to load updated extension.
 2. Open a project and add/edit `.orchestration/active_intents.yaml` (automatically created if absent).
 3. Call `select_active_intent(intent_id)` before issuing tools like `write_file`.
 4. Execute tools; inspect `.orchestration/agent_trace.jsonl` for entries.
 5. If provider fails due to prompt, check `/tmp/roo_prompt_debug_*.txt`.
-