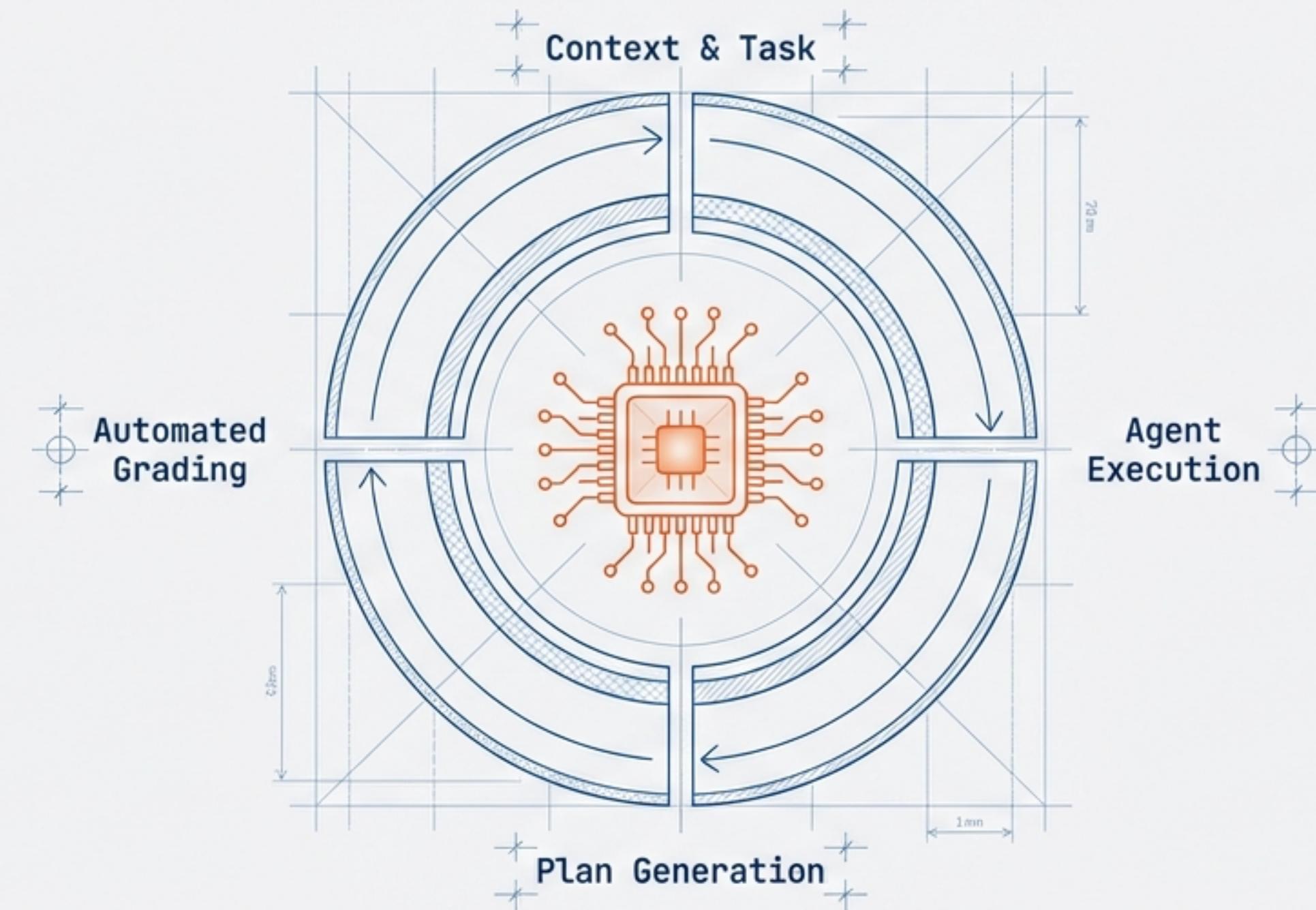


Planning Mode Eval

An Automated Evaluation Pipeline for Claude Code



A systems architecture walkthrough for measuring agentic planning capabilities at scale.

THE MISSION

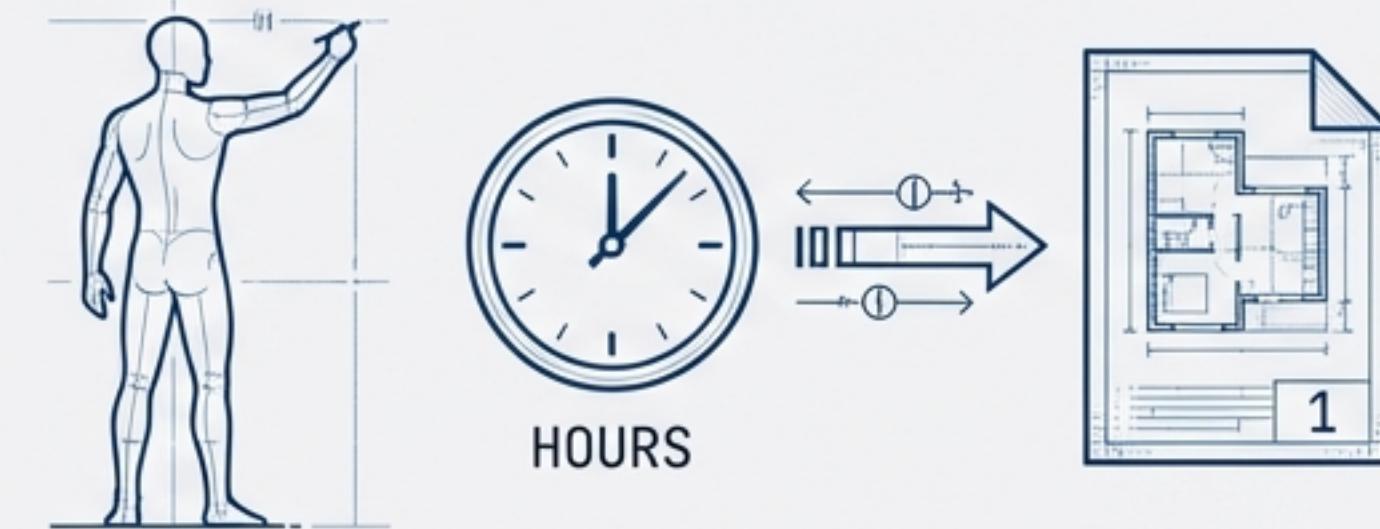
OBJECTIVE

To rigorously test Claude Code's "Plan Mode" capabilities using an end-to-end automated pipeline.

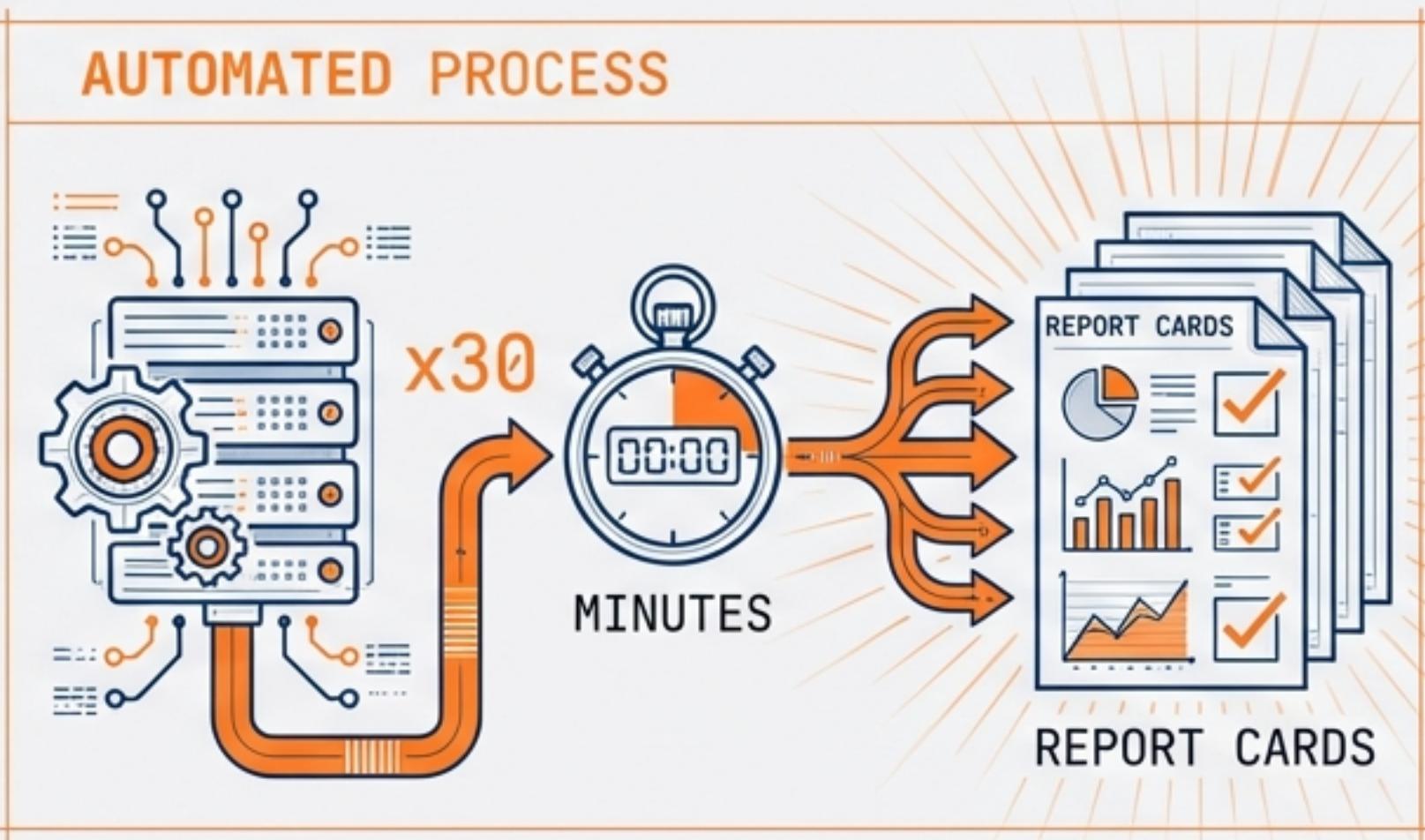
KEY CHALLENGES

- Scale:** Manual evaluation is unscalable. We need to process ~30 tasks per repository rapidly.
- Autonomy:** The system must mimic a human developer's workflow—from understanding the codebase to validating the output—without human hand-holding.

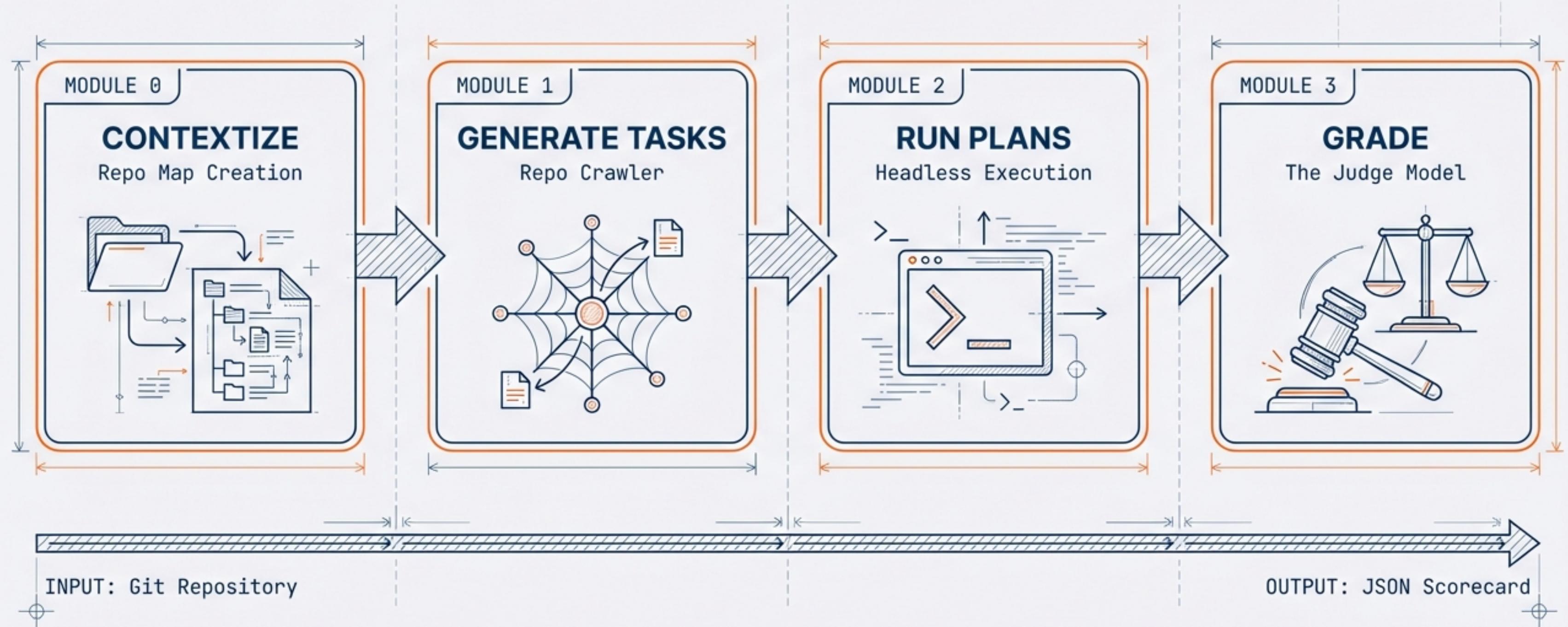
MANUAL PROCESS



AUTOMATED PROCESS



THE EVALUATION PIPELINE ARCHITECTURE



MODULE 0: CONTEXTIZING THE REPOSITORY

CONCEPT

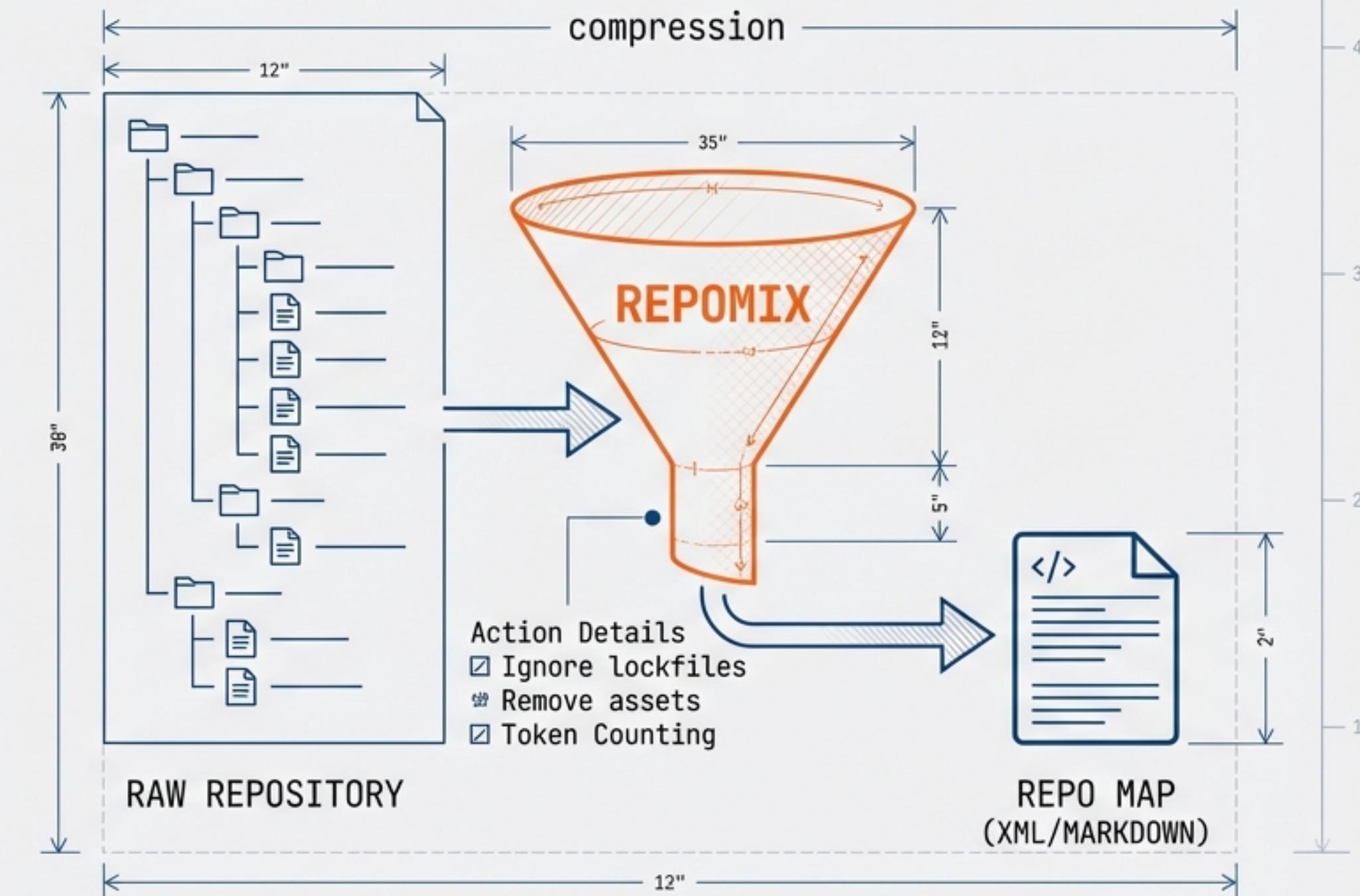
Transforming a directory of files into a single, high-density token stream.

THE TOOL: REPOMIX

Formerly Repopack. A dedicated tool for packing codebases for LLM consumption.

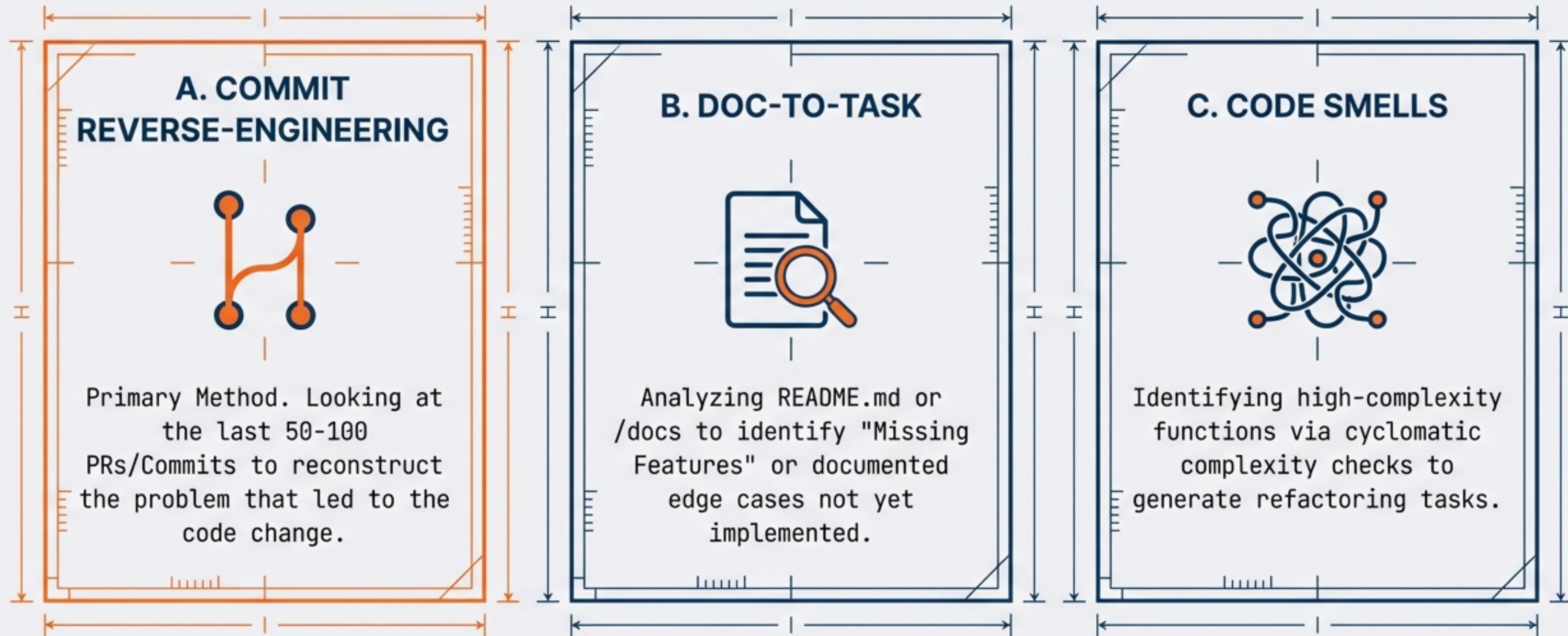
PROCESS

1. Clones the target repository.
2. Filters non-essential data.
3. Generates a compressed map fitting the context window.
4. Serves as ground truth for Agent and Judge.



MODULE 1: TASK GENERATION STRATEGY

We don't write tests; we mine them.

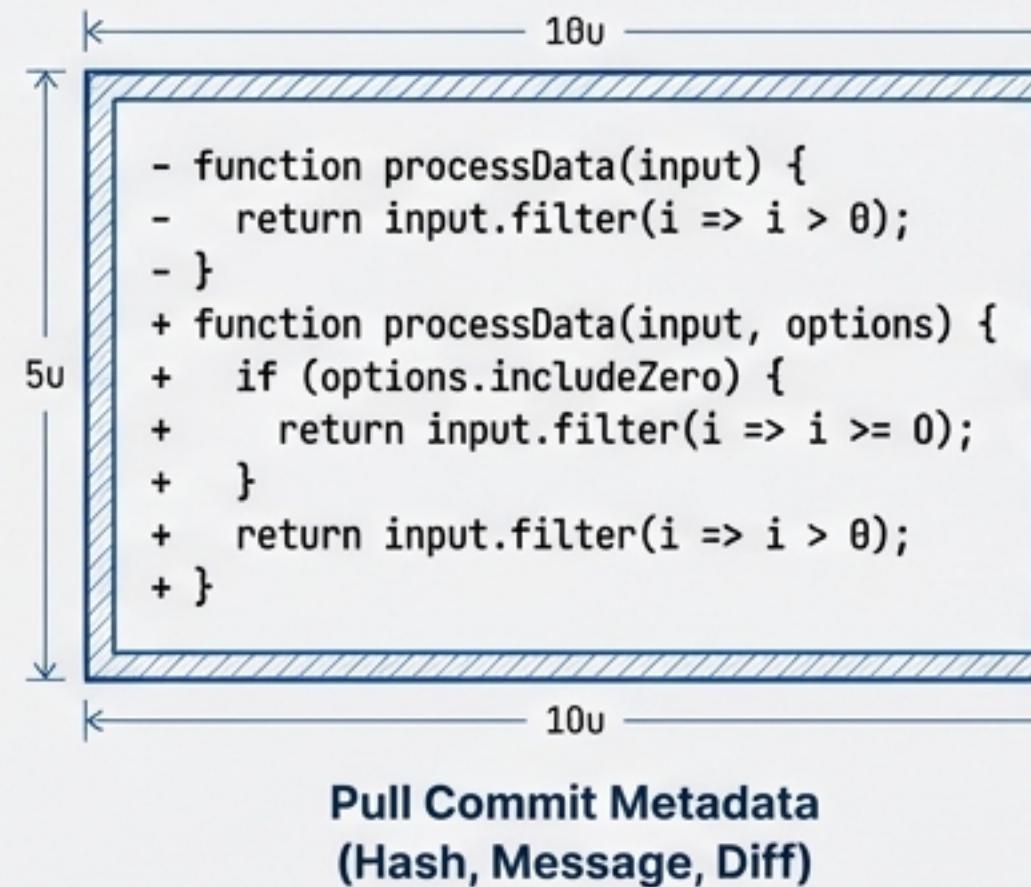


DEEP DIVE: REVERSE-ENGINEERING GIT HISTORY

AUTOMATION PIPELINE SCHEMATIC VOL. 02

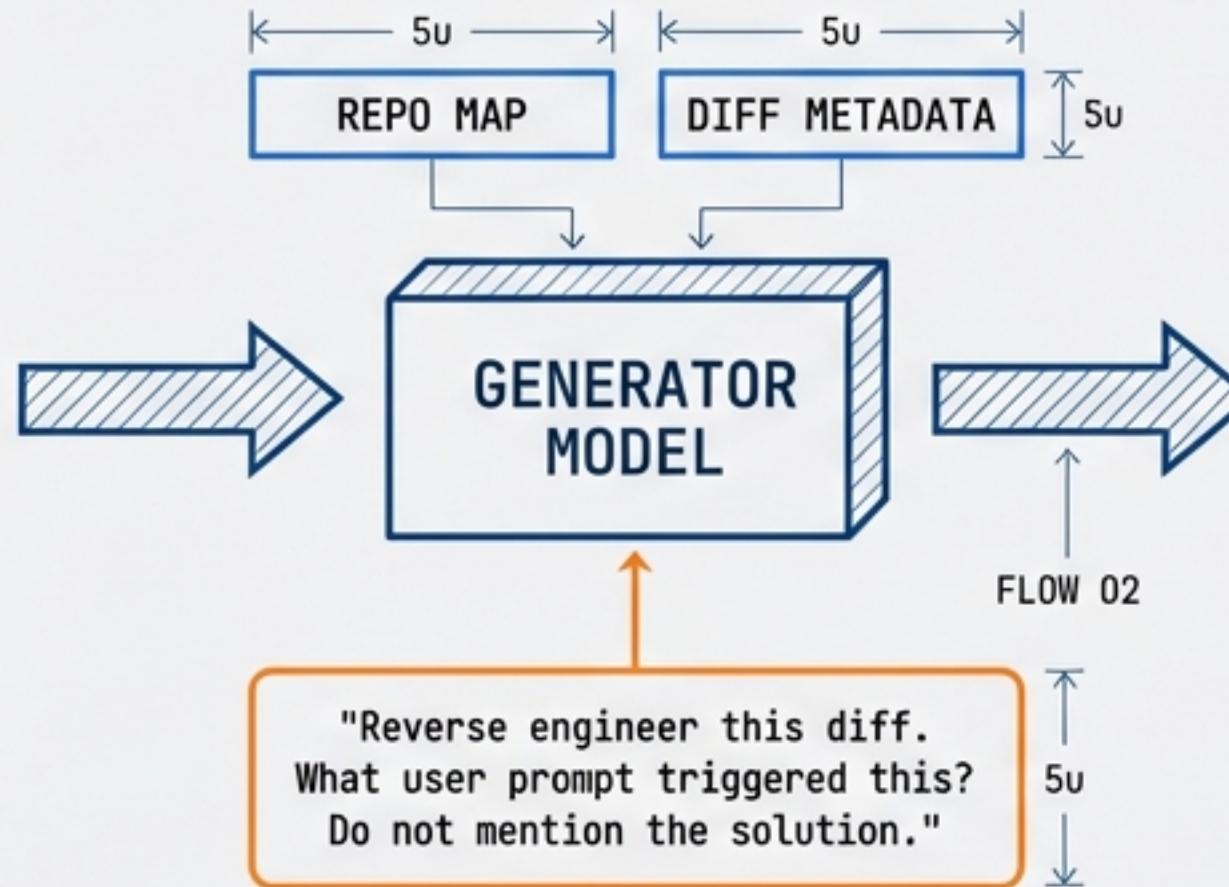
MODULE 1

EXTRACTION



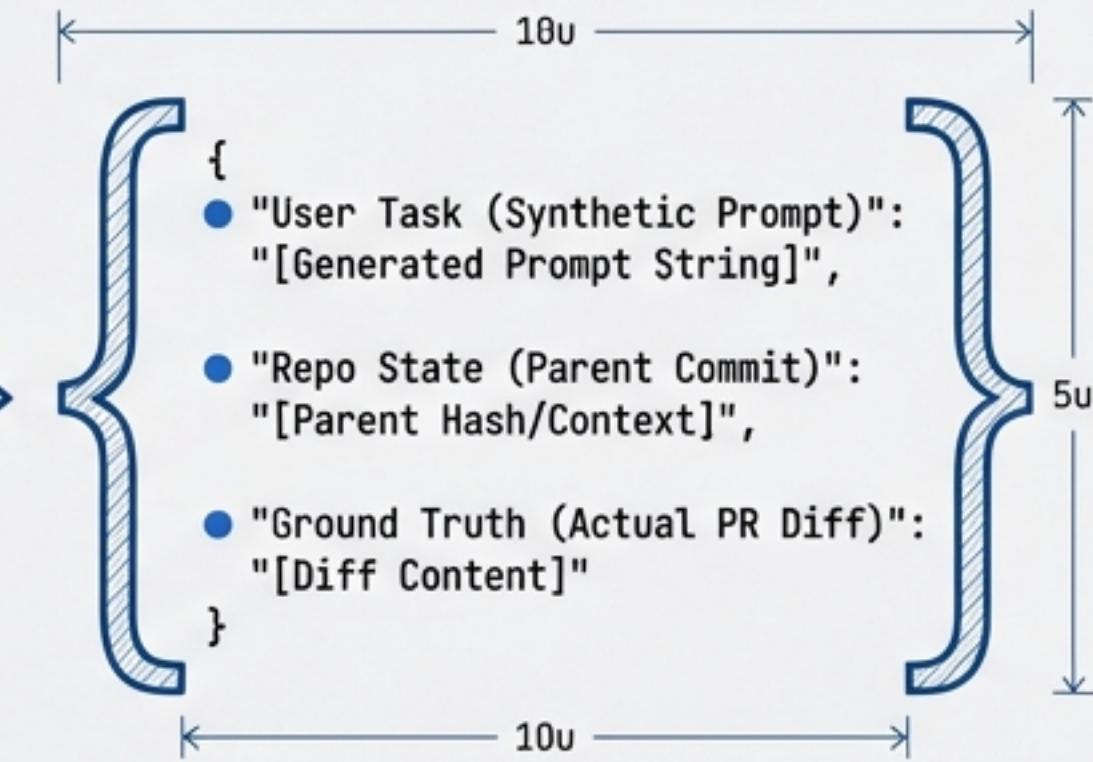
MODULE 2

LLM TRANSFORMATION



MODULE 3

THE OUTPUT OBJECT



Take merge commits, and classify into:

- (1) feature request (50%, 15/30)
- (2) bug fix (30%, 9/30)
- (3) code refactoring (20%, 6/30)
- (4) do-not-use; i.e. not suitable for plan mode

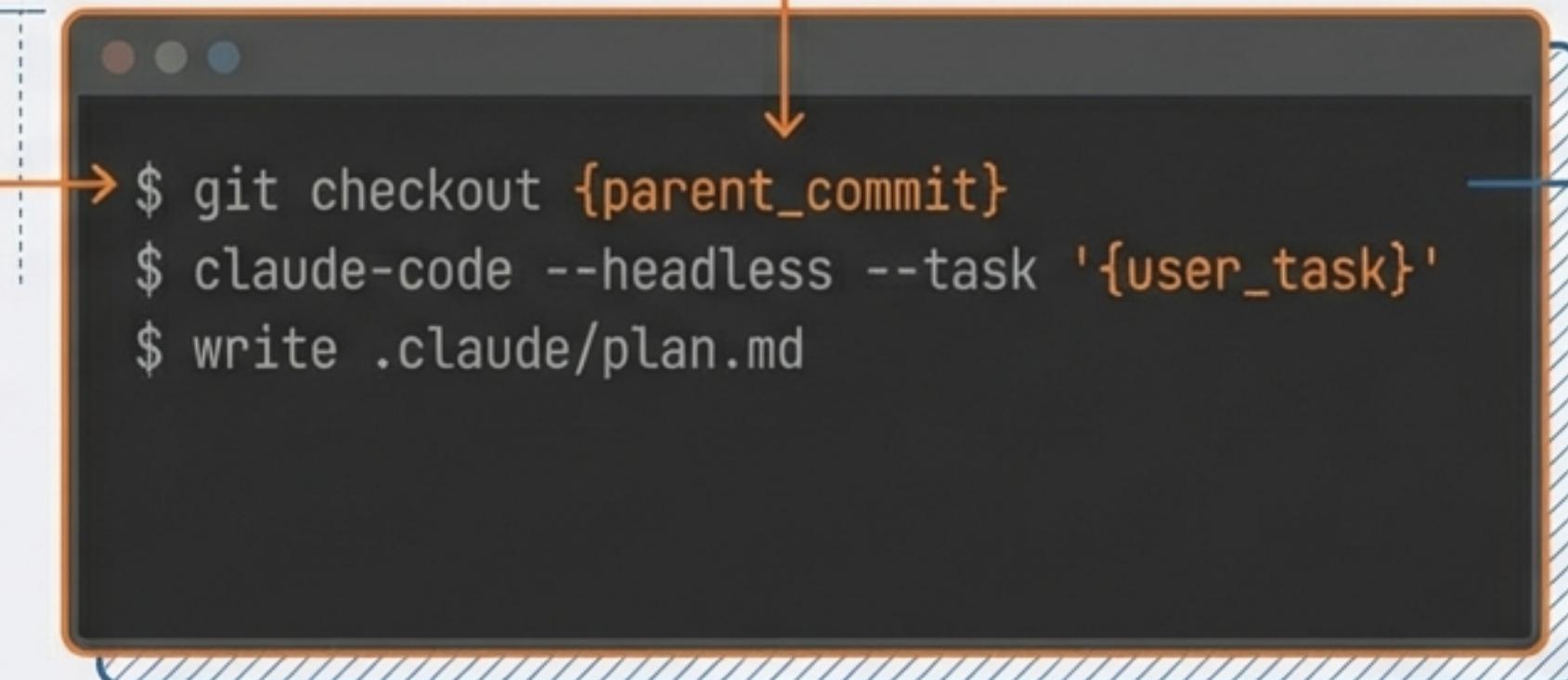


MODULE 2: HEADLESS EXECUTION

Automation Pipeline Schematic Vol. 03

1. State Reset

Clean slate. We revert the repo to the 'broken' or pre-feature state using the parent commit hash.



A terminal window with a dark background and light-colored text. It shows three commands:

```
$ git checkout {parent_commit}  
$ claude-code --headless --task '{user_task}'  
$ write .claude/plan.md
```

2. Execute Agent

Automating the ~~interactive Claude Code CLI~~ via ia wrapper scripts to bypass the terminal UI.

use claude-agent-sdk

3. Capture Plan

Intercepting standard output to save the Markdown plan for the Judge.

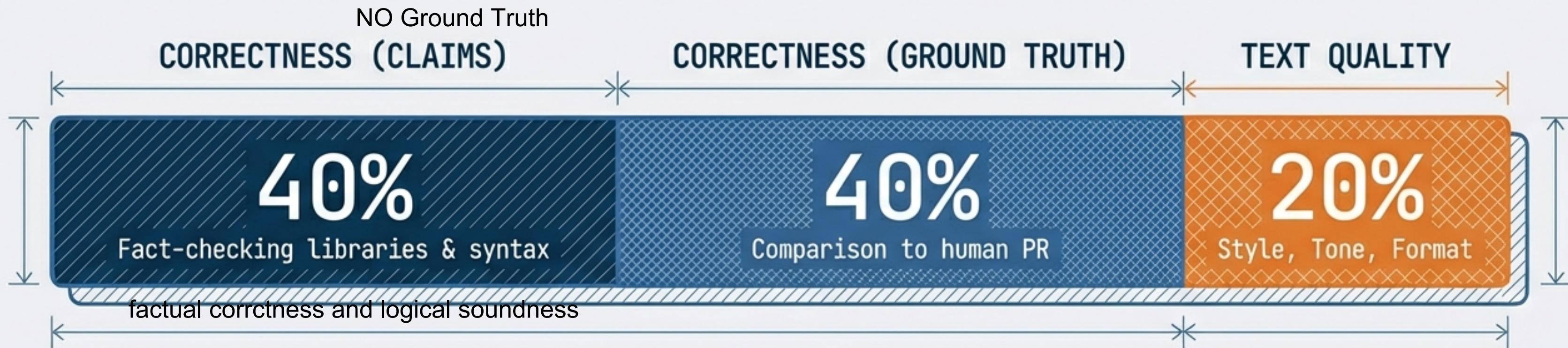
Timing: Asynchronous processing ~6 minutes per task.



MODULE 3: THE JUDGE MODEL

Model-Graded Evaluation using Claude Opus

Opus 4.6 & Sonnet 4.5



JetBrains Mono

WHY OPUS? A stronger model is required to audit the outputs of the smaller, faster agent model.



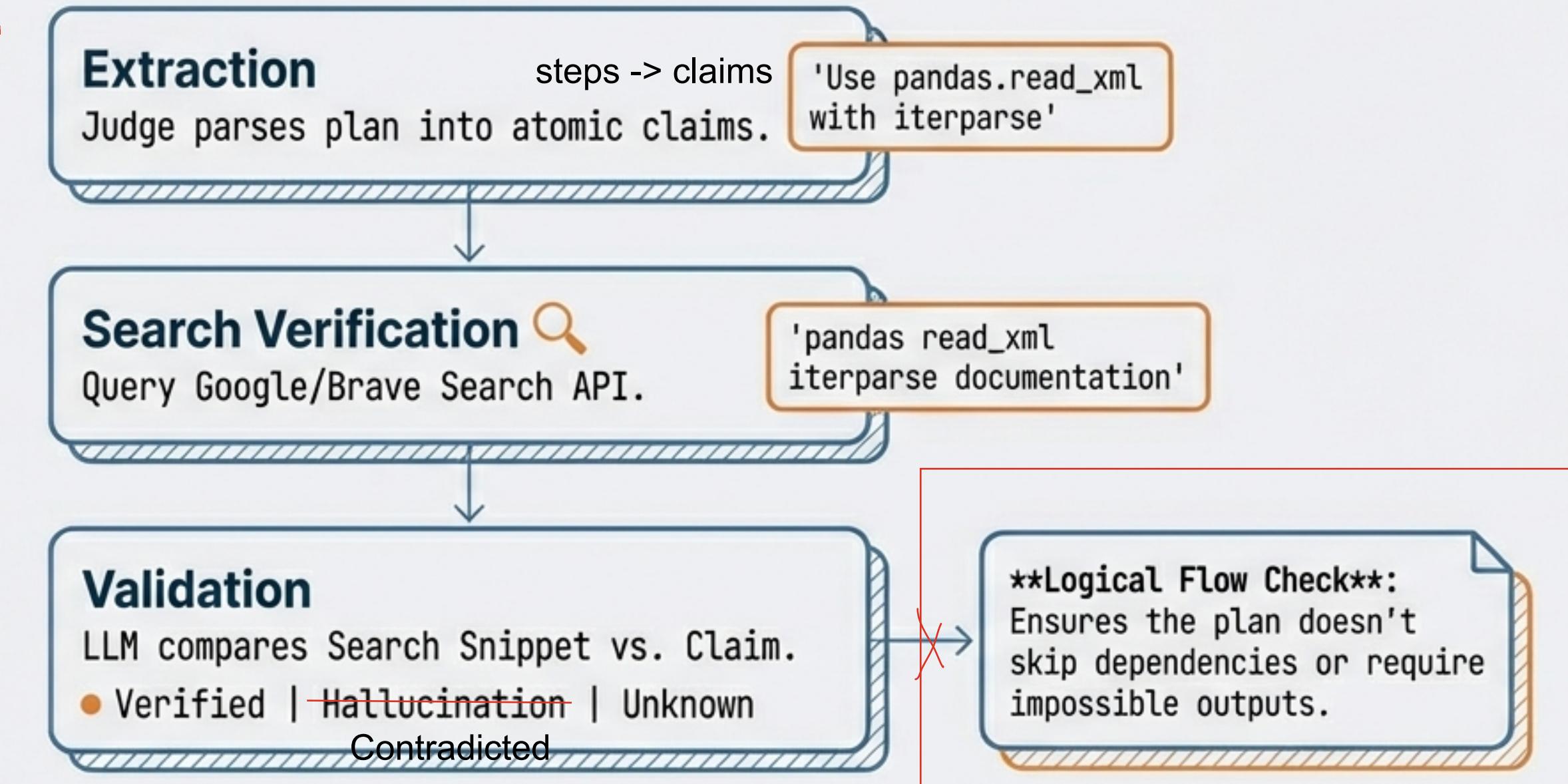
Grading Dimension A: Claim Verification (40%)

J .

Logical soundness:

- (1) Does any step require an output that a previous step fails to produce?
- (2) Is the overall plan logically sound and does it solve the problem?

Opus 4.6

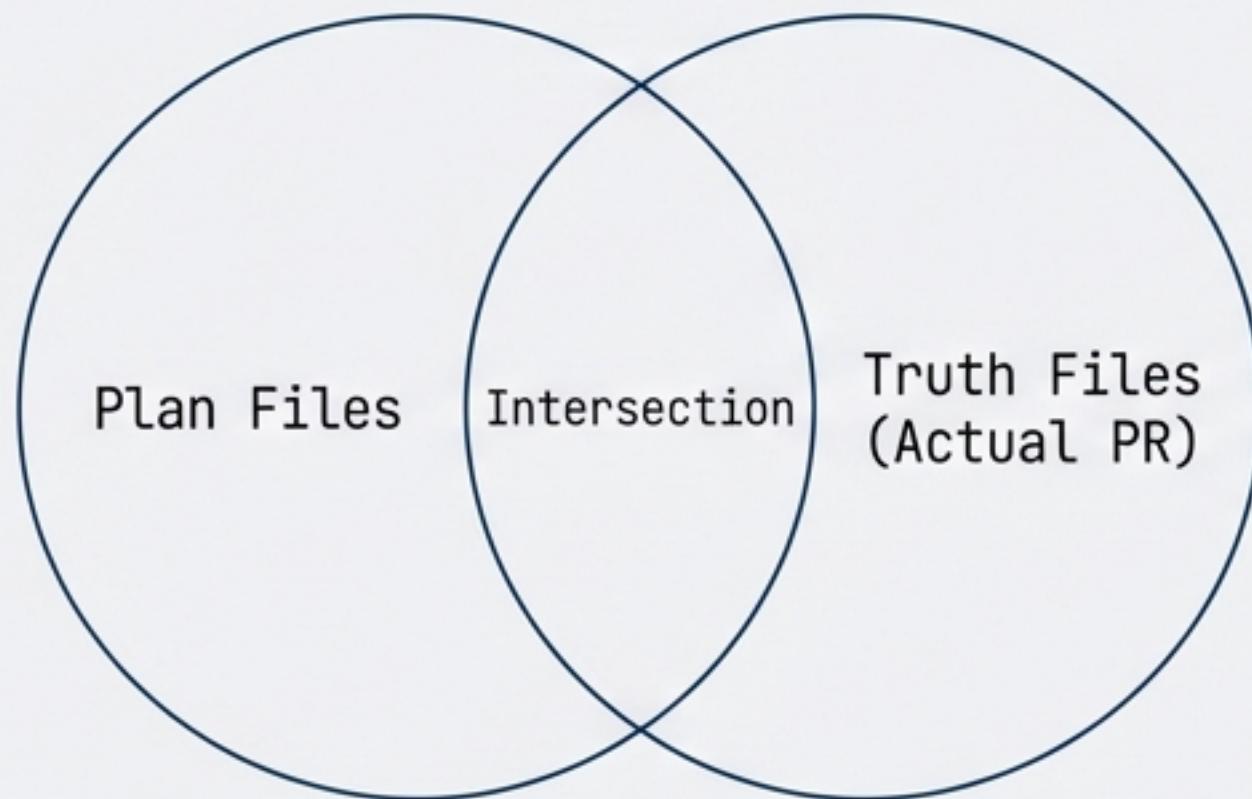


Timing: Processed as part of the Judge Model grading, ~1-2 minutes per task.



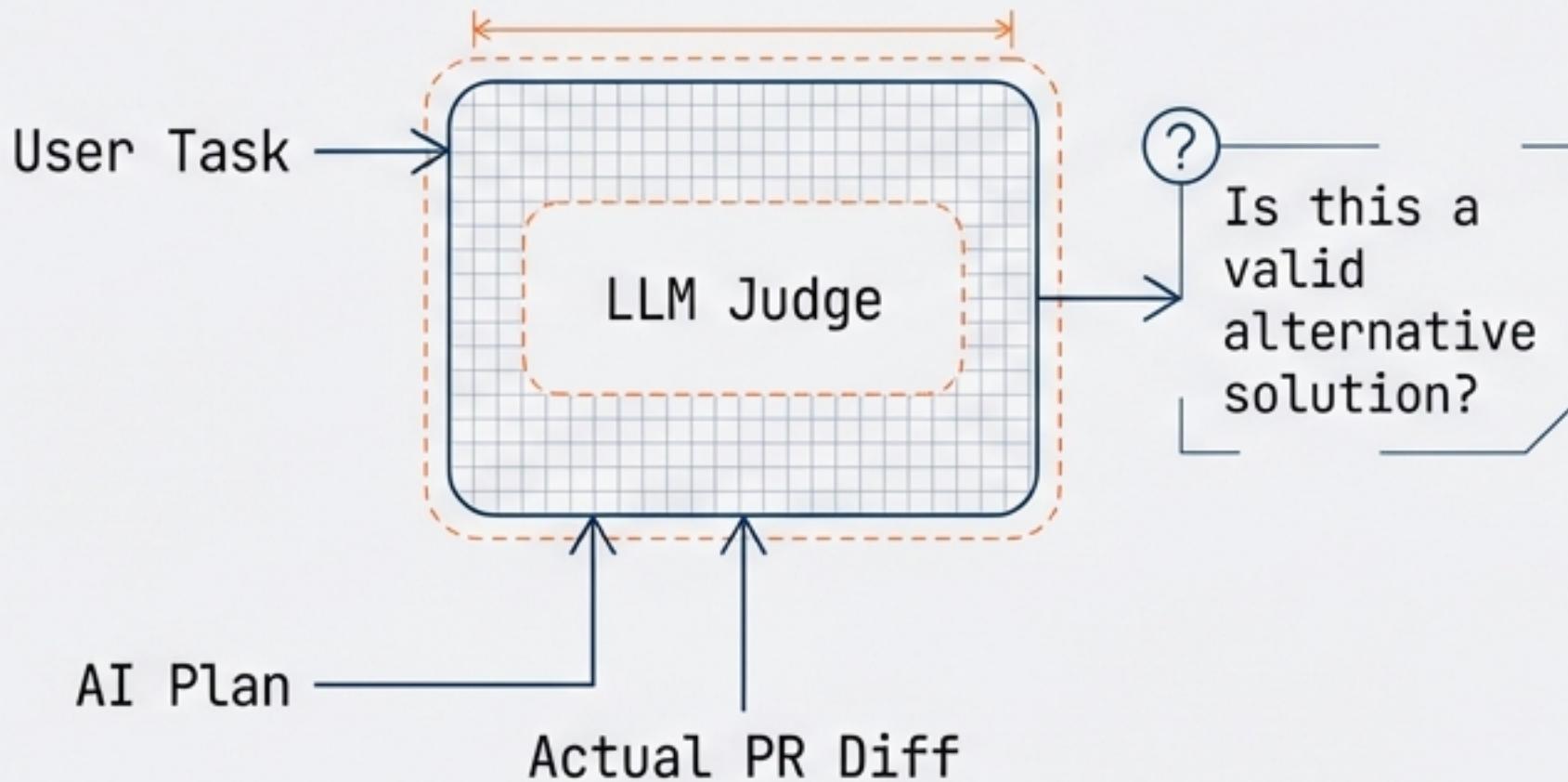
Grading Dimension B: Ground Truth Matching (40%)

Metric 1: File Relevance (Static)



Recall = Intersection / Truth Files
Precision = Intersection / Plan Files

Metric 2: Semantic Equivalency (LLM Judge)



Opus 4.6



Grading Dimension C: Text Quality (20%)

The Style Guide Rubric

Conciseness

High information density. Low fluff.

Failure Mode Example:

"It is important to note that..."

Tone

Professional, action-oriented.

Failure Mode Example:

"We could try to..."

Precision

Specific terminology over vague descriptions.

Failure Mode Example:

"Fix the thing..."

Formatting

Proper markdown usage (bullets, code blocks).

Failure Mode Example:

"Wall of text"



The Final Scorecard

Task ID: 492 - Fix race condition in UserDB

- Claims Verification: **4/5 Verified** (Confirmed `asyncio.lock` usage)
- File Recall: **100%** (Correctly identified `user_db.py`)
- Style Score: **5/5** (Concise, Action-Oriented)

TOTAL SCORE: 92/100

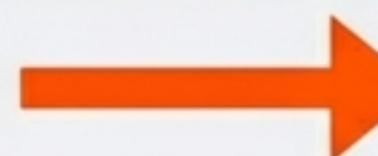
The final output aggregates weighted sub-scores to provide a single signal of confidence for the plan.



Future Roadmap & Optimization

Track 1: Better Context

Current: Static Repo Maps

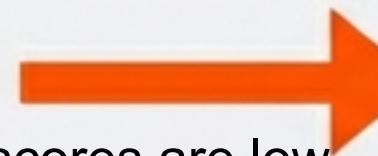


Future: Retrieval Augmented Generation (RAG) for large codebases.

Track 2: Better Tasks

Current: Heuristic Mining

Current GT matching scores are low

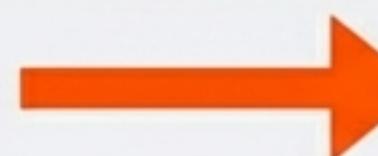


Future: User Research driven tasks & LLM-refined ground truth labeling.

Track 3: Better Grading

some claims are not directly verifiable with web searches

Current: Web Search



Future: Repo-aware verification & Human calibration.



Running the Pipeline

TERMINAL > `src.run_pipeline`

```
# 1. Install & Configure  
pip install -r requirements.txt  
cp .env.example .env # Set ANTHROPIC_API_KEY  
  
# 2. Run Pipeline Steps  
python -m src.run_pipeline contextize      # Build Map  
python -m src.run_pipeline generate-tasks # Mine History  
python -m src.run_pipeline run-plans       # Execute  
python -m src.run_pipeline grade           # Score  
  
# Or run all:  
python -m src.run_pipeline all
```

****Prerequisites:****
Python, Node, Git,
Claude Code CLI.



Technology Stack

Repomix

Context Engine

Claude Code SDK

Agent (Plan Mode)

Claude Opus

Judge Model

Search APIs

Google Custom Search /
Brave

Python

Core Orchestration

Built for modularity and scalability.

