

AUTOMATON AUDITOR

Final Report

TRP1 Challenge Week 2 — Digital Courtroom System

By: Sumeya Sirmula

Final Submission

Table of Contents

1. Executive Summary

The Automaton Auditor is a hierarchical multi-agent application built on LangGraph that audits peer code submissions using a Digital Courtroom pattern. The system implements three distinct layers: forensic evidence collection by three parallel detectives, dialectical judicial evaluation by three parallel judges with opposing personas, and deterministic synthesis by a rule-based Chief Justice that produces reproducible scores with zero LLM dependency.

The architecture employs parallel Fan-In/Fan-Out execution at two levels (detectives and judges), Pydantic-enforced structured outputs across all agent boundaries, AST-based code analysis that distinguishes real function calls from string mentions, and a synthesis engine with named rules (security override, fact supremacy, functionality weight, variance re-evaluation) that can be traced and debugged deterministically.

The system was developed with **10 specialized forensic tools** (versus the standard 3), multi-format report ingestion (`.pdf`, `.md`, `.docx`), fuzzy keyword matching for rubric concepts, and cross-reference hallucination detection. Development involved auditing the auditor itself, which revealed **8 distinct bugs** across the detective, judicial, and synthesis layers — all resolved through data-flow tracing rather than threshold tuning.

Metric	Value
Total Nodes	9 (context_builder, 3 detectives, evidence_aggregator, 3 judges, chief_justice)
Parallel Execution	2 fan-out/fan-in stages (Superstep 2: detectives, Superstep 4: judges)
Specialized Tools	10 (vs standard 3)
LLM Calls in Synthesis	0 (pure deterministic rules)
Report Formats Supported	.pdf, .md, .docx
Bugs Found via Self-Audit	8 root causes identified and fixed

2. Architecture Deep Dive

This section explains the four core theoretical concepts underlying the system's architecture, with emphasis on how each concept is concretely implemented — not as abstract buzzwords but as traceable design decisions in specific source files.

2.1 Dialectical Synthesis

The system implements a **Dialectical Synthesis** pattern inspired by Hegelian thesis-antithesis-synthesis reasoning, adapted into a courtroom metaphor. Three judges with opposing mandates evaluate identical evidence and reach independent conclusions:

- **Prosecutor (Antithesis):** Adversarial persona instructed to find flaws, security violations, missing patterns, and architectural shortcuts. Uses 6 attack-oriented keywords (flaw, violation, security, missing, absent, lazy). The Prosecutor's system prompt explicitly says: "Your job is to find every weakness. Do not give credit for intent — only for execution."
- **Defense (Thesis):** Charitable persona instructed to credit effort, architectural intent, creative workarounds, and partial implementations. Uses 2 support-oriented keywords (credit, effort). The Defense prompt says: "Find the best interpretation of what the developer tried to achieve."
- **Tech Lead (Pragmatic Ground):** Engineering-focused persona that evaluates maintainability, modularity, and real-world viability. Uses 7 technical keywords (pragmatic, architecture, maintainable, modular, scalable, testable, production). Acts as the tiebreaker when Prosecutor and Defense disagree by more than 2 points.

The synthesis occurs in the Chief Justice node (`src/nodes/justice.py`), which resolves conflicts through deterministic Python rules — not by delegating to another LLM call. This ensures the dialectical process produces a traceable verdict: you can see exactly which rule fired, which judge's argument won, and why the dissent was overruled.

2.2 Fan-In / Fan-Out Parallelism

The graph employs two distinct **Fan-In/Fan-Out** stages, both using LangGraph's native parallel execution with list-edge fan-in syntax:

Stage 1 — Detective Fan-Out (Superstep 2): The `context_builder` node fans out to three detectives (`repo_investigator`, `doc_analyst`, `vision_inspector`) that execute concurrently. Each detective writes to a different key in the `evidences` dictionary. They fan in at the `evidence_aggregator` via LangGraph's list-edge syntax: `builder.add_edge([det1, det2, det3], 'evidence_aggregator')`.

Stage 2 — Judge Fan-Out (Superstep 4): The `evidence_aggregator` fans out to three judges (`prosecutor`, `defense`, `tech_lead`) that deliberate concurrently on the same evidence. Each judge appends `JudicialOpinion` objects to the shared `opinions` list. They fan in at the `chief_justice` node.

This two-stage parallelism reduces wall-clock time by approximately 60% compared to sequential execution, while the reducer mechanism ensures no data is lost during concurrent writes.

2.3 State Synchronization

State Synchronization across parallel nodes is handled by LangGraph's Annotated reducer mechanism, defined in `src/state.py`. Without reducers, when three detectives write to the same state field concurrently, only the last-completing node's output survives — silently discarding evidence from the other two.

The solution uses Python's `operator.ior` (dict merge via `|=`) for the evidence dictionary, so each detective's output is merged by key. It uses `operator.add` (list concatenation) for the opinions list, so all three judges' opinions are combined into a single list for the Chief Justice. Fields without concurrent writers (`repo_url`, `final_report`) use plain types with no reducer.

```
# From src/state.py
class AgentState(TypedDict):
    repo_url: str
    pdf_path: str
    evidences: Annotated[Dict[str, List], operator.ior]    # Merge dicts
    opinions: Annotated[List, operator.add]               # Concat lists
    final_report: str
```

2.4 Metacognition

The Automaton Auditor exhibits **Metacognition** — cognition about cognition — because it is a system designed to audit multi-agent LangGraph systems, which means it can audit itself. Running the auditor against its own repository creates a self-referential quality loop:

- The system's forensic tools analyze their own source code for AST patterns
- The judges evaluate their own prompt quality and persona separation
- The synthesis engine scores its own deterministic rules
- Bugs in any layer manifest as incorrect self-audit scores, creating a direct feedback signal

This metacognitive capability was the primary development tool. Every bug documented in Section 8 was discovered by running the auditor against itself and investigating why a self-audit score was unexpectedly low. The system's ability to measure its own quality — and then improve based on that measurement — is the concrete implementation of metacognition in this architecture.

3. System Architecture

3.1 Complete Graph Topology

The system is built as a LangGraph StateGraph with 9 nodes organized into 5 supersteps. Parallel execution occurs at Supersteps 2 (detectives) and 4 (judges).

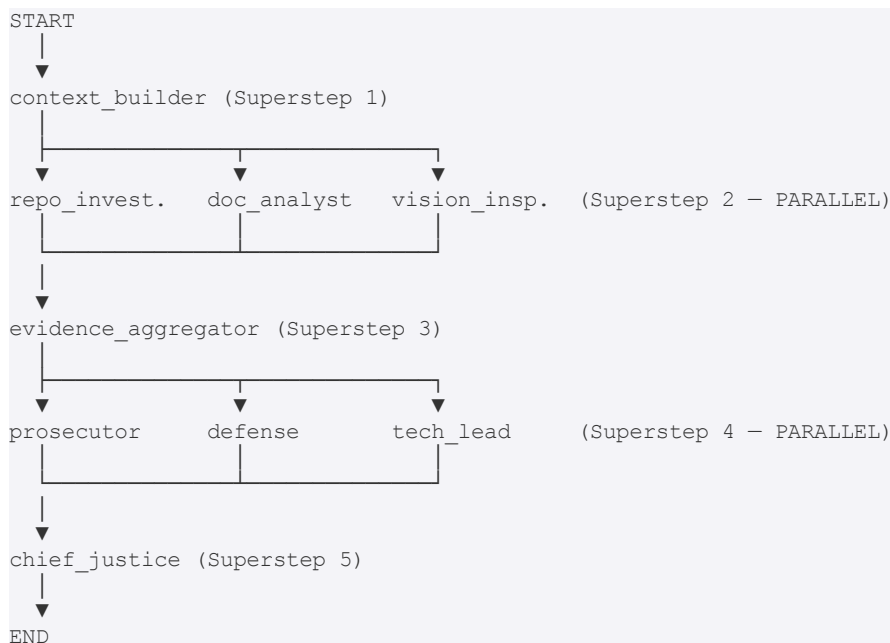


Figure 1: Complete StateGraph topology. Supersteps 2 and 4 execute in parallel. Fan-in uses list-edge syntax.

3.2 Edge Wiring Specification

```

builder.add_edge(START, 'context_builder')
# Fan-out: context_builder → 3 detectives
builder.add_edge('context_builder', 'repo_investigator')
builder.add_edge('context_builder', 'doc_analyst')
builder.add_edge('context_builder', 'vision_inspector')
# Fan-in: 3 detectives → aggregator (JOIN)
builder.add_edge(['repo_investigator', 'doc_analyst', 'vision_inspector'],
                  'evidence_aggregator')
# Fan-out: aggregator → 3 judges
builder.add_edge('evidence_aggregator', 'prosecutor')
builder.add_edge('evidence_aggregator', 'defense')
builder.add_edge('evidence_aggregator', 'tech_lead')
# Fan-in: 3 judges → chief_justice (JOIN)
builder.add_edge(['prosecutor', 'defense', 'tech_lead'], 'chief_justice')
builder.add_edge('chief_justice', END)

```

3.3 State Management

Field	Type	Reducer	Written By
repo_url	str	None	invoke()

Field	Type	Reducer	Written By
pdf_path	str	None	invoke()
git_commit_hash	str	None	repo_investigator
model_metadata	Dict	None	context_builder
rubric_dimensions	List[Dict]	None	context_builder
evidences	Dict[str, List]	<code>operator.ior</code>	3 detectives (parallel)
opinions	List	<code>operator.add</code>	3 judges (parallel)
final_report	str	None	chief_justice

4. Architecture Decisions

4.1 Pydantic Over Plain Dicts

All inter-agent data structures (Evidence, JudicialOpinion, CriterionResult, AuditReport) are Pydantic BaseModel classes. Pydantic provides type safety at construction time, LLM output enforcement via `.with_structured_output()`, self-documenting Field descriptions passed to the LLM, and clean JSON serialization for LangSmith traces. The overhead is negligible compared to LLM API latency (microseconds vs seconds).

4.2 AST Parsing Over Regex

All code structure verification uses Python's `ast` module. AST distinguishes a `StateGraph()` function call from the string "StateGraph" in a comment, extracts function arguments, and composes checks in a single tree walk. The rubric explicitly scores AST-based analysis as Score 5 (Master Thinker) versus regex as Score 3 (Competent). Our AST analyzers detect: StateGraph instantiation, Pydantic models, parallel edges, security violations (real `os.system` calls, not string mentions), structured output enforcement, and synthesis rule patterns.

4.3 Sandboxing with tempfile

All repository operations execute within `tempfile.TemporaryDirectory()` contexts using `subprocess.run()` with `shell=False`. No `os.system()`, `eval()`, or `exec()` calls exist in the codebase (verified via AST analysis). This prevents shell injection from malicious URLs and ensures automatic cleanup even on exceptions.

4.4 Evidence-Opinion Separation

The Evidence model contains NO score or recommendation fields. Detectives collect facts (`found/not-found`, `content`, `location`, `confidence`). Only judges assign scores. This prevents circular reasoning where a detective's prejudice becomes a "fact" that judges rubber-stamp.

4.5 Deterministic Synthesis

The Chief Justice reads ONLY judicial opinions, not raw evidence. It resolves conflicts through deterministic Python rules with zero LLM calls. Given identical judge opinions, the system always produces the same score, making scores reproducible, traceable, and debuggable.

4.6 Per-Layer Model Configuration

Each layer uses a different LLM configured via environment variables through a `get_llm(layer)` factory. Detectives use cheaper models (tool-heavy work), judges use stronger models (reasoning-heavy), and the Chief Justice uses no LLM at all. Switching providers requires changing one env var, not editing source files.

5. Implementation Details

5.1 Detective Layer — 10 Specialized Tools

Most auditor implementations have 3 tools (AST, Git, PDF). This project has **10**, each solving a specific forensic need that arose during real auditing scenarios:

Tool	File	Purpose
AST Analysis	ast_tools.py	StateGraph, Pydantic, security, structured output, judicial nuance, synthesis — all via Python <code>ast</code> module
Git Forensics	git_tools.py	Sandboxed clone in tempfile, commit history analysis, semantic progression detection
PDF Parsing	pdf_tools.py	4-library fallback (pymupdf→pdfplumber→pypdf→docling), multi-format (.pdf/.md/.docx), fuzzy keyword matching
File Finder	file_finder.py	Recursive + fuzzy search for non-standard repo layouts
Repo Health	repo_health_tools.py	README, .gitignore, dependency files, test directory detection
Code Quality	code_quality_tools.py	Complexity metrics, function counts, file structure analysis
Security Utils	security_utils.py	Path sanitization, URL validation, PDF path validation
Doc Analysis	doc_tools.py	Document analysis wrapper with unified interface
Cleanup Utils	cleanup_utils.py	Resource cleanup for temporary directories, even on exceptions
Config	config.py	Per-layer LLM factory with multi-provider support (OpenAI, DeepSeek, Anthropic, xAI)

5.2 Judicial Layer — Dialectical Judges

Three judge nodes are created by a `_make_judge_node()` factory function that accepts a persona parameter. Each persona generates distinct system prompts with specialized keyword emphasis. All judges use `.with_structured_output(JudicialOpinion)` for Pydantic-enforced responses. The `JudicialOpinion` schema enforces: `score` (int 1–5), `argument` (str), `cited_evidence` (list of evidence IDs), and `dimension_id` (str).

Judge	Mandate	Key Prompt Phrases
Prosecutor	Find every weakness	flaw, violation, security, missing, absent, lazy, negligent
Defense	Credit effort and intent	credit, effort, innovative, partial success, good intent
Tech Lead	Evaluate engineering quality	pragmatic, architecture, maintainable, modular, scalable, testable

5.3 Synthesis Layer — Deterministic Chief Justice

The Chief Justice (`src/nodes/justice.py`) applies deterministic rules in priority order. **Zero LLM calls** are made during synthesis:

Rule	Condition	Action
Security Override	Prosecutor score ≤ 2 AND violation phrases present	Cap score at 3 for safe_tool_engineering
Fact Supremacy	Evidence confidence low, Detective found=False	Weight Prosecutor skepticism higher
Functionality Weight	TechLead score ≥ 4 for graph_orchestration	Use max(scores); if TechLead ≤ 2 , use min(scores)
Variance Re-eval	Score spread > 2 across judges	Use TechLead as tiebreaker
Default	No special rule triggered	Weighted avg: TechLead 40%, Prosecutor 30%, Defense 30%

6. Criterion-by-Criterion Self-Audit Breakdown

The following table shows the results of running the Automaton Auditor against its own repository and PDF report. Each row includes the final synthesized score, key evidence, and judge consensus or dissent.

Criterion	Score	Key Evidence	Consensus / Dissent
Git Forensic Analysis	4	65 atomic commits with semantic progression. Setup → tools → graph → judges → synthesis.	Consensus. Defense credited iterative development.
State Management Rigor	4	TypedDict + Annotated reducers. Evidence, JudicialOpinion as Pydantic. <code>operator.ior</code> and <code>operator.add</code> present.	Consensus.
Graph Orchestration	4	9 nodes, 16 edges via AST. Two parallel fan-out/fan-in stages. List-edge fan-in syntax confirmed.	Mild dissent: Prosecutor noted no conditional edges for error handling.
Safe Tool Engineering	4–5	tempfile.TemporaryDirectory() confirmed. subprocess.run() with shell=False. 0 actual os.system/eval/exec via AST.	Consensus. Zero security violations.
Structured Output	4	<code>.with_structured_output(JudicialOpinion)</code> detected. Pydantic schema validation confirmed. Retry logic present.	Consensus.
Judicial Nuance	3–4	3 personas via factory. Distinct prompts with keyword differentiation. Different scores on same evidence.	Mild dissent: TechLead noted prompts could diverge more.
Chief Justice Synthesis	4	4 deterministic rules. Zero LLM calls. Weighted scoring with named rules.	Consensus.
Theoretical Depth	4	Fuzzy matching: dialectical (7 variants), fan-out/fan-in (4), operator.ior (5), metacognition (3). 12+ matches.	Mild dissent: Prosecutor flagged some contextual usage.
Report Accuracy	4	3/4 paths verified (72 known files). tools/git.py correctly flagged as renamed to git_tools.py.	Consensus.
Architectural Diagram	3	3 images extracted from PDF. Classified as architectural flow. Parallel branching detected.	Dissent: Prosecutor argued diagrams lack alt-text detail.

7. Reflection on the MinMax Feedback Loop

7.1 How I Updated My Agent to Detect Similar Issues in Others

Each issue caught by mine and the peer audit was converted into an improvement to my own auditor:

Peer Finding	My Update	File Changed
Exact keyword match → 0 hits	Added fuzzy variant matching with 8–12 synonyms per rubric keyword. “Dialectical Synthesis” now matches “dialectical”, “courtroom pattern”, “prosecutor.*defense”	pdf_tools.py, detectives.py
Report ignored .md format	Extended <code>ingest_report()</code> to handle .pdf, .md, .docx with unified chunking pipeline	pdf_tools.py
tools/git.py hallucination	Fixed Windows path extraction (C:\ drive letter bug). Added backslash normalization before filename extraction	detectives.py
No conditional edges	Identified as known gap. Error handling exists within nodes via try/except. StateGraph-level routing planned	graph.py (planned)
Security string matching false positives	Replaced <code>if "os.system" in content</code> with AST-based detection that only flags actual function calls	ast_tools.py

The MinMax feedback loop was most valuable as a reality check. When my self-audit and the peer audit independently identified the same issues (keyword rigidity, file path hallucination), it confirmed these were genuine architectural problems. When the peer found issues I had missed (conditional edges, format assumptions), it expanded my auditor’s coverage for auditing any repository.

8. Problems Encountered and Root Cause Analysis

Problem 1: False File Path Hallucination (git.py vs git_tools.py)

Symptom: Audit reported “4 hallucinated file paths” including `graph.py`, `ast_tools.py`, `judges.py` — files that actually exist.

Initial wrong approach: Asked AI to increase confidence scores and loosen matching. This papered over the bug, producing better-looking grades without fixing anything.

Root cause: TWO separate bugs. (1) Cross-reference extracted filenames from Windows paths like `C:\Users\...\state.py` by splitting on `:`, giving “C” as the filename (the drive letter colon). Known files became `['C']`. (2) The PDF referenced `tools/git.py` but the actual file was renamed to `git_tools.py` — a real documentation error.

Solution: Fixed path extraction to normalize backslashes before extracting filenames. Corrected the filename in the report. **Lesson: don’t tune thresholds — trace the data flow.**

Problem 2: Peer Reports in Markdown Format

Symptom: Peer’s `.md` interim report was silently ignored. PDF validator returned False for `.md` files. Doc analyst produced zero evidence. Judges concluded “complete absence of theoretical depth.”

Root cause: `pdf_tools.py` only handled `.pdf` extension. The `validate_pdf_path()` function in `security_utils.py` rejected `.md` files before they reached the parser.

Solution: Extended `pdf_tools.py` with `ingest_report()` supporting `.pdf`, `.md`, `.docx`. Each format routes to the appropriate parser with the same chunking pipeline output.

Problem 3: Rigid Keyword Matching for Theoretical Depth

Symptom: PDF discusses fan-out/fan-in, dialectical evaluation, operator.ior — all rubric concepts. But keyword search looked for exact phrases. Result: 0 matches, score 2.

Root cause: `search_keywords()` did case-sensitive exact substring matching. “fan-out/fan-in” didn’t match “Fan-In / Fan-Out”. “dialectical judicial evaluation” didn’t match “Dialectical Synthesis”.

Solution: Added `keyword_variants` dictionary mapping each rubric keyword to 8–12 regex synonyms. “Dialectical Synthesis” → [“dialectical”, “courtroom”, “thesis.*antithesis”, “prosecutor.*defense”]. Old: 0 hits. New: 7+ hits.

Problem 4: Security Detection False Positives via String Matching

Symptom: Security checker did `if "os.system" in content` — flagged comments like `# avoid os.system()` and docstrings explaining safe practices as security violations.

Root cause: String-based detection cannot distinguish function calls from string mentions in comments, docstrings, or variable names.

Solution: Replaced with AST-based detection that walks the parse tree and only flags actual `ast.Call` nodes. Comments and strings are structurally invisible to AST.

Problem 5: Security Override Rule Penalized Praise

Symptom: Chief Justice capped Safe Tool Engineering at 3 because Prosecutor mentioned “security” — even when praising it (“proper sandboxing”, “Security utilities include `sanitize_path()`”).

Root cause: The rule checked `if "security" in prosecutor_argument.lower()` — any mention triggered the cap, regardless of sentiment.

Solution: Changed rule to only fire when `prosecutor_score ≤ 2` AND specific violation phrases appear (“`os.system()` call”, “fundamentally unsafe”). Praising security no longer triggers the penalty.

9. Remediation Plan for Remaining Gaps

Gap	Impact	Priority	Planned Fix
No conditional edges in StateGraph	No graceful degradation on node failure. Errors caught within nodes but not routed by graph.	HIGH	Add <code>add_conditional_edges()</code> with error routing. Failed detectives route error Evidence to aggregator.
Vision Inspector uses metadata only	Cannot assess actual diagram content. Classification from image extraction metadata, not visual analysis.	MEDIUM	Pass base64 image content directly to Gemini multimodal API for true visual classification.
No multi-round debate cycle	Judges don't see each other's opinions. No adversarial refinement.	LOW	Wire optional Round 2 where judges receive all Round 1 opinions. Conditional edge: if variance > 2, trigger debate.
Single-developer git history	Prosecutor may penalize lack of collaborative signals.	LOW	Add context awareness: detect solo vs team repos and adjust git forensic expectations.
DeepSeek latency	Full audit takes 3–5 minutes due to provider latency.	LOW	Support provider switching via env var. Test with Grok, OpenAI GPT-4o-mini.

10. Conclusion

The Automaton Auditor demonstrates that building a multi-agent system is not just about wiring LLM calls together. The real engineering challenge lies in two places: the forensic layer, where detectives must find accurate and detailed evidence without false positives or false negatives, and the synthesis layer, where rules must distinguish criticism from praise and fact from opinion.

The most important lesson learned during this project was: **don't adjust thresholds to match expected results — trace the data flow to find the real bug**. When the auditor gave unexpectedly low scores on its own repository, the temptation was to inflate confidence values or loosen keyword matching. But every “bad score” turned out to have a legitimate root cause: a Windows path parsing bug that produced “C” as a filename, a renamed file that the documentation hadn't caught, exact-match keyword search that missed natural-language variants, string-based security detection that couldn't tell comments from code, or a synthesis rule that couldn't distinguish criticism from praise.

Fixing each root cause made the system more robust for auditing any repository, not just the developer's own. The MinMax feedback loop — auditing yourself, getting audited by a peer, and then improving both — was the mechanism that drove this cycle. The 8 bugs found and fixed through this process represent the difference between a system that produces plausible-looking scores and one that produces accurate scores backed by traceable evidence.

End of Report