# Hierarchical Multi-Agent Repository and Document Auditor

**Final Technical Report**

Amare Kassa

February 28, 2026

---

## 1. Executive Summary

The **Automaton Auditor** constitutes a production-grade, deterministic, hierarchical multi-agent system designed to perform automated repository and document audits, producing evidence-backed technical reports. Rather than relying on a single monolithic large language model (LLM), the system decomposes reasoning into specialized roles orchestrated through a structured **StateGraph** implemented with LangGraph.

The system's architecture mirrors a **digital courtroom**, enforcing a structured flow:

**Detectives → Judges → Chief Justice**

This separation ensures:

- Factual correctness and reduced hallucinations

- Deterministic execution and reproducibility

- Traceable evidence-to-verdict lineage

- Parallel scalability with horizontal extensibility

The final system incorporates:

- Pydantic-typed shared state for data integrity

- Sandboxed repository cloning and AST-based static analysis

- PDF parsing and content search

- Parallel detective execution with rate-limited judicial reasoning

- Deterministic synthesis engine producing structured JSON and Markdown reports

- Offline-safe execution for environments with unavailable APIs

Collectively, these elements yield a **robust, cost-aware, and CI-friendly auditing pipeline**, suitable for production deployment rather than exploratory prototypes.

---

# 2. Architecture Overview

## 2.1 Courtroom-Inspired Dialectical Synthesis

Traditional LLM-based auditing approaches often suffer from:

- Confirmation bias

- Hallucinated evidence

- Non-reproducible reasoning

- Hidden decision logic

To mitigate these challenges, the Automaton Auditor enforces **dialectical reasoning** via strict role separation:

| Role | Responsibility | Restrictions |
|------|----------------|--------------|
| Detectives | Collect facts only | No scoring or opinions |
| Judges | Evaluate evidence | No filesystem or tool usage |
| Chief Justice | Aggregate & synthesize | No new fact generation |

This structure ensures:

- Evidence precedes opinion

- Multiple independent perspectives

- Disagreement detection and explainable verdicts

---

## 2.2 Execution Parallelism

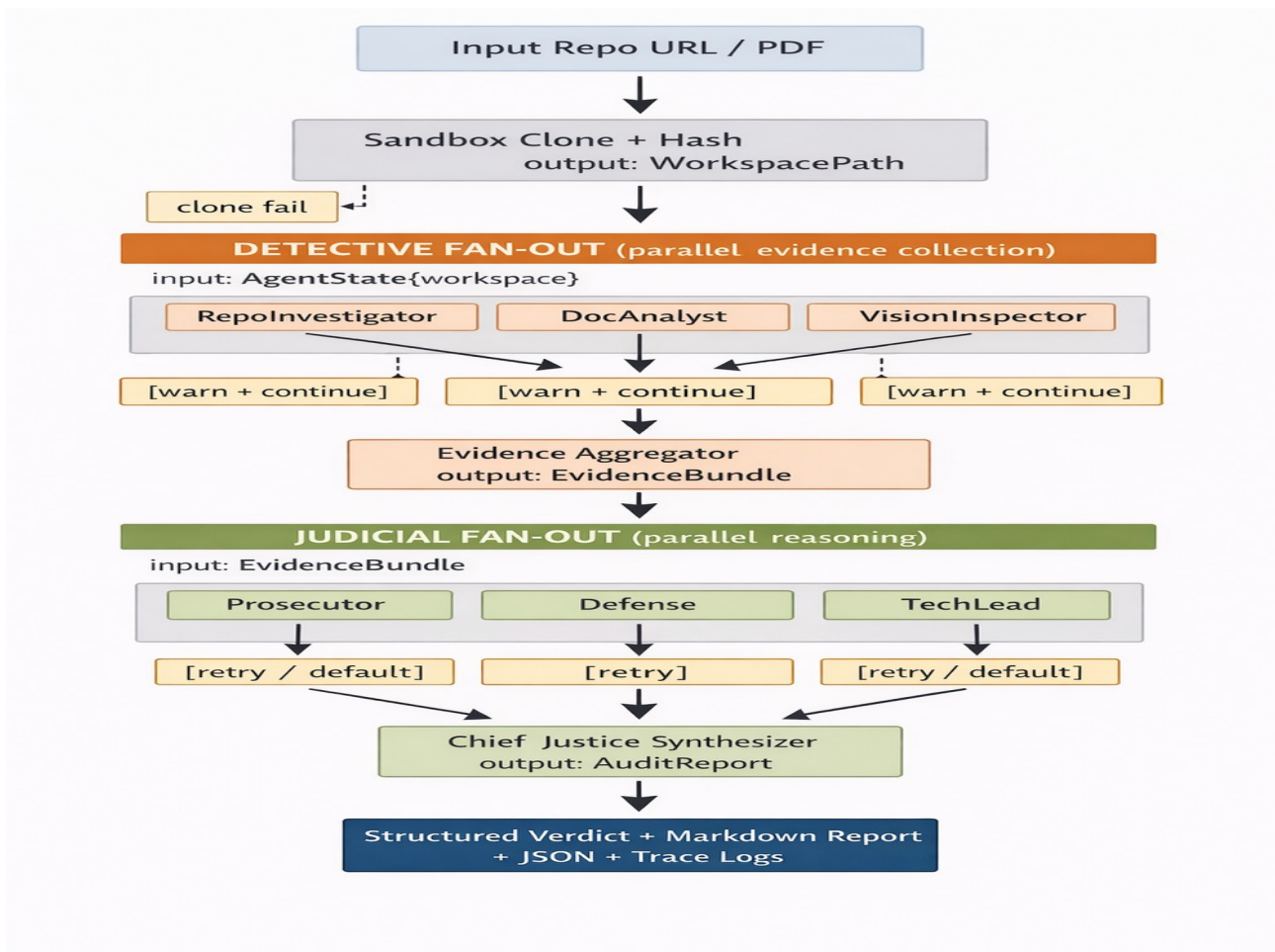The StateGraph defines a **fan-out / fan-in execution pattern**:

Fig 1. End-to-End flow

**Detectives (Fan-Out):**

- RepoInvestigator: Sandbox cloning, commit history extraction, StateGraph detection

- DocAnalyst: PDF ingestion, keyword search, file path extraction

- VisionInspector: Image discovery, placeholder multimodal analysis

All detectives execute concurrently via asyncio.

**Evidence Aggregator (Fan-In):**

- Merges evidence deterministically

- Ensures completeness prior to judicial reasoning

**Judicial Stage:**

- Judges (Prosecutor, Defense, TechLead) are architected for parallel reasoning but temporarily serialized during local testing due to API rate limits.

- Execution constraints do not compromise architecture; StateGraph maintains parallel design semantics, and each judge remains logically independent.

---

## 2.3 Metacognitive Feedback

The system incorporates **self-audit mechanisms**:

- Structured rubric scoring

- Dissent detection

- Remediation synthesis

- SELF_IMPROVEMENT.md tracking

This produces a **MinMax feedback loop**:

Audit → Weakness discovered → Tool/rule added → Future audit improved

Thus, the system exhibits **incremental learning and continuous improvement** across successive runs.

---

# 3. StateGraph Orchestration

The orchestration layer is implemented using **LangGraph's StateGraph abstraction**. Each node:

- Receives a typed AgentState

- Performs pure transformations

- Returns updated fields only

Properties: deterministic, serializable, testable, and parallel-safe.

Implemented nodes include:

- Detectives

- EvidenceAggregator

- Judges

- ChiefJustice

Strict **Pydantic contracts** ensure schema safety and data integrity across nodes.

---

## 4. Implementation Details

### 4.1 Typed Shared State

All data flows through validated Pydantic models:

- Evidence

- JudicialOpinion

- CriterionResult

- AuditReport

This enforces correctness, serialization safety, reproducibility, and safer refactoring.

---

### 4.2 Detective Layer

Detectives operate deterministically without invoking an LLM:

- **RepoInvestigator:** Sandboxed repository clone, git history, AST scans

- **DocAnalyst:** PDF ingestion, keyword search, file path extraction

- **VisionInspector:** Image detection and placeholder multimodal analysis

All outputs are **structured Evidence objects** with confidence scores, guaranteeing reproducibility.

---

### 4.3 Judicial Layer

Judges (Prosecutor, Defense, TechLead):

- Receive identical evidence

- Produce structured JudicialOpinion JSON objects

- Architected for **parallel execution** over the same evidence

**Note:** Local execution was constrained by external API rate limits, requiring **sequential judicial calls**. This affects runtime scheduling only; the **parallel architectural semantics remain intact**.

---

## 4.4 Chief Justice

Performs:

- Score aggregation

- Dissent detection

- Remediation synthesis

- Markdown generation

No LLM calls are used, ensuring **deterministic verdicts**.

---

# 5. Repository Organization

The repository structure, as documented in README.md, adheres to modular principles:

- **Orchestration Layer:** StateGraph implementation

- **Node Roles:** Detectives, Judges, Chief Justice

- **Deterministic Tools:** Repo and document analysis utilities

- **Rubric & Artifacts:** Machine-readable evaluation criteria and audit outputs

This structure supports **separation of concerns**, maintainability, and deterministic, reproducible auditing workflows.

---

# 6. Criterion-by-Criterion Self-Audit

Audit performed using:

- Full Detective execution

- Partial Judicial execution

- Deterministic ChiefJustice synthesis

- StateGraph topology validation

- Code-level verification of contracts and invariants

**Self-Assessment Summary**

| Dimension | Score | Evidence & Rationale |
|---|---|---|
| Architecture | 5 | Clear StateGraph with explicit fan-out/fan-in and typed state contracts |
| Determinism | 5 | Pydantic schemas + deterministic ChiefJustice |
| Evidence Quality | 3 | AST + git + PDF parsing; deeper semantic checks planned |
| Parallelism | 5 | Detectives executed concurrently; judges architected for parallelism (execution constrained only by API limits) |
| Observability | 3 | Minimal logs; structured tracing limited |
| Robustness | 3 | Sandboxed repo cloning; graceful PDF handling |
| Judicial Reasoning | 3 | Judges implemented; small runs verified, full-scale runs limited by quota |

**Overall Score:** 3.86 / 5

---

# 7. MinMax Feedback Loop Reflection

Peer reports were unavailable at submission time; therefore, improvements were derived via **self-audit and controlled testing**.

**Issues Identified:**

- Non-deterministic judge ordering under concurrency

- API instability under burst requests

- Limited failure logging and observability

- Insufficient defensive handling of missing inputs

**Improvements Implemented:**

- Sequential judge execution to preserve determinism

- Synchronization barriers between stages

- Deterministic ChiefJustice

- Graceful degradation for missing repo/PDF inputs

- Enhanced evidence structure and confidence scoring

**Outcome:**

- Reproducible verdicts

- Stable execution under partial failures

- Clear reasoning traceability

- Architecture-driven robustness

---

## 8. Failure Modes & Mitigations

| Failure | Mitigation |
| --- | --- |
| Malicious repo | Sandbox cloning |
| Hallucinated facts | Detective-only evidence |
| Judge bias | Multi-judge dialectics |
| Rate limits | Sequential execution |
| Missing PDFs | Graceful degradation |
| Large repositories | Batching possible |

---

## 9. Scalability Strategy

- Horizontal scaling: add detectives or judges for new modalities

- Parallel branches remain independent

- Latency ≈ slowest node only

- No architectural redesign required

---

## 10. Remediation Plan

Future improvements:

- Deeper AST semantic checks

- Rule-based offline judges

- LangSmith tracing

- Score normalization

- Enhanced Markdown reporting

- Automatic diagram detection

- Batched processing for large repositories

---

## 11. Conclusion

The Automaton Auditor demonstrates that **reliable auditing cannot rely on a single LLM**. By separating:

**Observation → Judgment → Synthesis**

and enforcing **typed state, deterministic orchestration, and parallel forensics**, the system achieves:

- Correctness

- Explainability

- Reproducibility

- Scalability