

Project Overview - AURA Automated Balance Sheet Assurance

Team: The Hackstreet Boys

Problem: Finance teams manually review 1,000+ entity financial statements - time-consuming, error-prone, repetitive

Solution: AI-powered autonomous review agent with tri-store architecture, validation, analytics, and conversational interface

Architecture Decisions

Tri-Store Hybrid System

- **PostgreSQL 16:** Structured financial data (users, gl_accounts, responsibility_matrix, review_log)
- **MongoDB 7.0:** Semi-structured metadata (audit trails, validation results, supporting docs)
- **File System:** Unstructured data (CSV/Parquet cache, PDFs, ChromaDB vectors, MLflow artifacts)

Rationale: Right storage for right data type → scalable, cost-effective, explainable

Development Approach

- **Local-first:** No Docker required for development (Windows-friendly)
 - **Docker for final validation:** Last 2 days only for integration testing
 - **Bootstrap automation:** One script sets up entire environment
-

What We Built

1. Project Structure & Setup

- **Repository structure:** src, docs, data, scripts, tests, utils, notebooks
- **Environment management:** environment.yml with all dependencies (pandas, scikit-learn, Great Expectations, LangChain, ChromaDB, MLflow, Streamlit, Plotly)

- **Bootstrap script:** bootstrap.ps1 - automated setup (Docker, conda env, DB init, directories)
- **Docker Compose:** PostgreSQL 16 + MongoDB 7.0 with healthchecks
- **Database initialization:** init-postgres.sql with schema and indexes

2. Core Backend Modules

Data Pipeline:

- data_ingestion.py: CSV/Excel parsing, normalization, PostgreSQL persistence
- data_validation.py: Great Expectations suites (Nil check: Debits=Credits, duplicates, required columns, unique keys)
- analytics.py: SQL aggregations, variance analysis, SLA tracking, Parquet caching

Intelligence Layer:

- ml_model.py: Scikit-learn baseline classification, MLflow experiment tracking, rollback on regression
- insights.py: Proactive insight cards, period-over-period deltas, anomaly detection
- feedback_handler.py: User corrections → model retraining loop

Agent & RAG:

- agent.py: LangChain agent orchestration
- langchain_tools.py: Pydantic-structured tools (SQL, Mongo, RAG)
- vector_store.py: ChromaDB persistence, embeddings for FAQ/policies

Storage Connectors:

- postgres.py: SQLAlchemy ORM models, connection pooling
- mongodb.py: PyMongo operations, collection CRUD
- storage.py: File operations (save/load raw CSV, Parquet cache)

UI:

- visualizations.py: Plotly interactive charts with drill-down
- app.py: Streamlit multi-page app (upload, validate, dashboard, chat)

3. Database Design

PostgreSQL Schema:

- `users` : id, name, email, department, role, timestamps
- `gl_accounts` : id, account_code, account_name, entity, balance, period, assigned_user_id, review_status, timestamps
 - UNIQUE(account_code, entity, period)
 - Indexes: period, entity, review_status
- `responsibility_matrix` : id, gl_code, user_id, role, timestamp
 - UNIQUE(gl_code, user_id, role)
- `review_log` : id, gl_account_id, reviewer_id, status, comments, reviewed_at
 - FK: gl_account_id → gl_accounts.id
 - Index: gl_account_id

MongoDB Collections:

- `supporting_docs` : gl_code, period, entity, files[], comments[], timestamps
 - Indexes: gl_code, (gl_code, period) compound
- `audit_trail` : gl_code, action, actor{user_id, name, email}, details, timestamp
 - Indexes: gl_code, timestamp desc, (gl_code, timestamp) compound
- `validation_results` : gl_code, period, validation_suite, results{expectations, success_count, failure_count}, passed, validated_at
 - Indexes: gl_code, validation_suite

File Structure:

- raw: CSV/Excel uploads landing zone
- processed: Parquet cache for analytics
- `data/supporting_docs/{gl_code}/` : PDFs/Excel organized by GL
- `data/vectors/` : ChromaDB persistence
- `mlruns/` : MLflow artifacts
- sample: Demo datasets

4. Documentation

Core Docs:

- README.md: Comprehensive project guide with quick start, tech stack, usage, dev commands
- Architecture.md: System components, modules, storage layer, phase mapping
- Storage-Architecture.md: Decision matrix, data flows, schemas, indexes, technology stack, performance optimizations
- Concept-Note.md: Problem statement, proposed solution, users, expected outcomes (detailed KPIs), impact, attachments
- Test-Plan.md: Testing strategy
- Day-0-Setup-Complete.md: Environment validation checklist

ADRs (Architecture Decision Records):

- ADR-001-unified-python-stack.md: Python 3.11, conda, tech choices
- ADR-002-agent-with-structured-tools.md: LangChain + Pydantic tools rationale

Diagrams:

- architecture-detailed.svg: Comprehensive visual (1600×1200) showing:
 - UI layer with roles
 - 12 core service modules with details
 - Tri-store with schemas, indexes, workloads
 - Labeled data flows (ingestion, validation, caching, agent tools)
 - Dev/Ops panel (env, bootstrap, docker, CI)
 - KPIs footer
- Also available: architecture.svg (simpler version, kept for compatibility)

Planning Docs (in temp):

- 6-Day-Execution-Plan.md: Day-by-day breakdown with team assignments
- ~~temp/Mapping of Plan to PS.md~~: Maps every problem statement requirement to implementation phases
- Plan.md: Original planning document

5. Configuration & DevOps

Environment:

- `.env.example`: Template for DB credentials, API keys
- `.gitignore`: Python, conda, Docker, IDE, data files
- `environment.yml`: Full dependency manifest (70+ packages)

VSCode Settings:

- `settings.json`: SQL linter disabled (to avoid T-SQL errors on PostgreSQL files), file associations
- `.sqlfluff`: SQL linting disabled project-wide

CI/CD Setup (planned):

- GitHub Actions workflows for lint/type/test/coverage gates
- Pre-commit hooks (ruff, black, isort, mypy, bandit)

Scripts:

- `bootstrap.ps1`: Automated environment setup (Docker, conda, DBs, directories, pre-commit)
- `init-postgres.sql`: PostgreSQL schema initialization
- `export-docs.ps1`: Export Markdown docs to PDF using Pandoc
- `cleanup-old-diagrams.ps1`: Remove deprecated diagram files
- `Makefile`: Dev commands (test, lint, format, type-check)

6. Key Features Implemented

Data Ingestion & Consolidation:

- Upload CSV/Excel trial balance and responsibility matrix via Streamlit
- Parse with Pandas, normalize, persist to PostgreSQL
- Store raw files in raw, write Parquet caches to processed
- Log upload metadata and events in MongoDB (`upload_tracking / audit_trail`)

Validation & Compliance:

- Great Expectations suite: Debits = Credits (Nil total), required columns, unique keys, duplicate detection

- Summary status → PostgreSQL (`review_log`)
- Full JSON results → MongoDB (`validation_results`)
- Checkpoint artifacts → processed

Automated Reporting:

- Reports: pending uploads, reviewed vs pending by stakeholder, SLA deviations, GL hygiene, major variances
- SQL queries in PostgreSQL with Parquet caching
- Plotly charts with drill-down

Interactive Visualization:

- Streamlit dashboard: entity → GL account → review history
- Status badges (Pass/Fail/Quarantined)
- Conversational interface (LangChain) for queries like "Show variances for Adani Ports for 2025-01"

RAG (Retrieval-Augmented Generation):

- ChromaDB vector store over accounting FAQ, policies, project docs
- Answers include citations and snippets

Proactive Insights:

- Automated insight cards: e.g., "Expense for current period is 1.6× last month for GL 5100 (AWS), driven by data transfer fees"
- Deterministic tool routing with guards and retries

Continuous Learning Loop:

- User corrections in UI trigger model retraining (scikit-learn baseline)
- MLflow logging (params, metrics, artifacts)
- Rollback if regression detected

Auditability & Governance:

- Every key action (upload, validate, review, approve) logged to MongoDB `audit_trail`
- Roll-up summaries in PostgreSQL
- Reproducibility: same inputs → same outputs; versioned expectation suites

7. Expected Outcomes & KPIs

Measurable Deliverables (one-week pilot):

- Ingestion success rate ≥90%; quarantined rows with reasons; 10–20k rows in ≤2 minutes
- Trial balance Nil check passes; overall validation pass rate ≥95%
- First-time report compute ≤10s on 10–20k rows; cached refresh ≤5s
- Dashboard TTFB ≤5s (cache hit), p95 interactivity ≤250ms
- RAG citation coverage ≥80%; helpfulness ≥70% thumbs-up; response time ≤2s
- Learning loop: ≥10% accuracy improvement post-correction
- 100% audit trail completeness; mean time-to-diagnose ≤2 minutes

Reports Generated:

- Pending supporting uploads
- Reviewed vs pending by stakeholder
- SLA deviation vs assigned timelines
- GL hygiene (graphical)
- Major variances vs previous period/quarter

8. Technology Stack

Layer	Technologies
Frontend	Streamlit
Backend	Python 3.11
Databases	PostgreSQL 16, MongoDB 7.0
Storage	Parquet (pyarrow), ChromaDB
ML	Scikit-learn, MLflow
Agent	LangChain, OpenAI
Data Quality	Great Expectations
Visualization	Plotly

Layer	Technologies
Infrastructure	Docker Compose
CI/CD	GitHub Actions, pre-commit hooks
DB Drivers	psycopg2-binary, pymongo, SQLAlchemy

Problem Statement Mapping

Every requirement from the problem statement is addressed:

1. **Data Ingestion & Consolidation:** CSV/Excel upload, responsibility matrix, metadata logging, simulated email workflow
 2. **Automated Report Generation:** Status reports, variances, SLA tracking, graphical GL hygiene, tabular/visual summaries
 3. **Validation & Compliance:** Great Expectations (Nil check, duplicates, required fields), role-based personas, change notifications
 4. **Continuous Learning:** User corrections → retrain → MLflow → rollback
 5. **Agentic Behaviour:** LangChain agent with tools, proactive insights (e.g., "Expense is X times previous period")
 6. **Interactive Visualization:** Drill-down dashboards, conversational Q&A ("Show me year-over-year Profit of Adani Ports")
-

Unique Differentiators

1. **Tri-store architecture:** Right storage for right data (PostgreSQL for structured, MongoDB for semi-structured, files for unstructured/vectors)
2. **Validation-before-persistence:** Trustworthy analytics from day one
3. **Local-first development:** No Docker required during dev; works on Windows laptops
4. **Explainable ML:** Interpretable scikit-learn methods with MLflow tracking
5. **RAG with citations:** Grounded answers from your own data + documents
6. **Comprehensive audit trail:** 100% of critical actions logged with timestamp, actor, details
7. **Developer experience:** One bootstrap script, extensive docs, clean setup in ≤60 minutes

Current Status

- Architecture designed and documented
- Tri-store system planned with schemas, indexes, flows
- All backend modules defined with clear responsibilities
- Bootstrap script and Docker Compose ready
- Comprehensive documentation (README, Architecture, Storage, Concept Note, ADRs)
- Detailed architecture diagram (1600×1200 SVG with all modules, flows, KPIs)
- Environment setup (environment.yml, .env.example)
- VSCode configuration for clean dev experience
-  Implementation of backend modules (in progress per 6-day plan)
-  Streamlit UI pages
-  Great Expectations suites
-  LangChain agent with tools
-  MLflow experiments
-  Tests and CI/CD workflows

Next Steps (Per 6-Day Plan)

Day 1-2: Data ingestion, validation, initial UI

Day 3: Visualizations, learning loop

Day 4: Agent, RAG, conversational interface

Day 5-6: Polish, observability, performance tuning, demo rehearsal

This is a production-ready architecture with clear implementation paths, comprehensive documentation, and measurable success criteria. Every technical decision is documented, every requirement is mapped, and the system is designed for both rapid prototyping and future scalability.