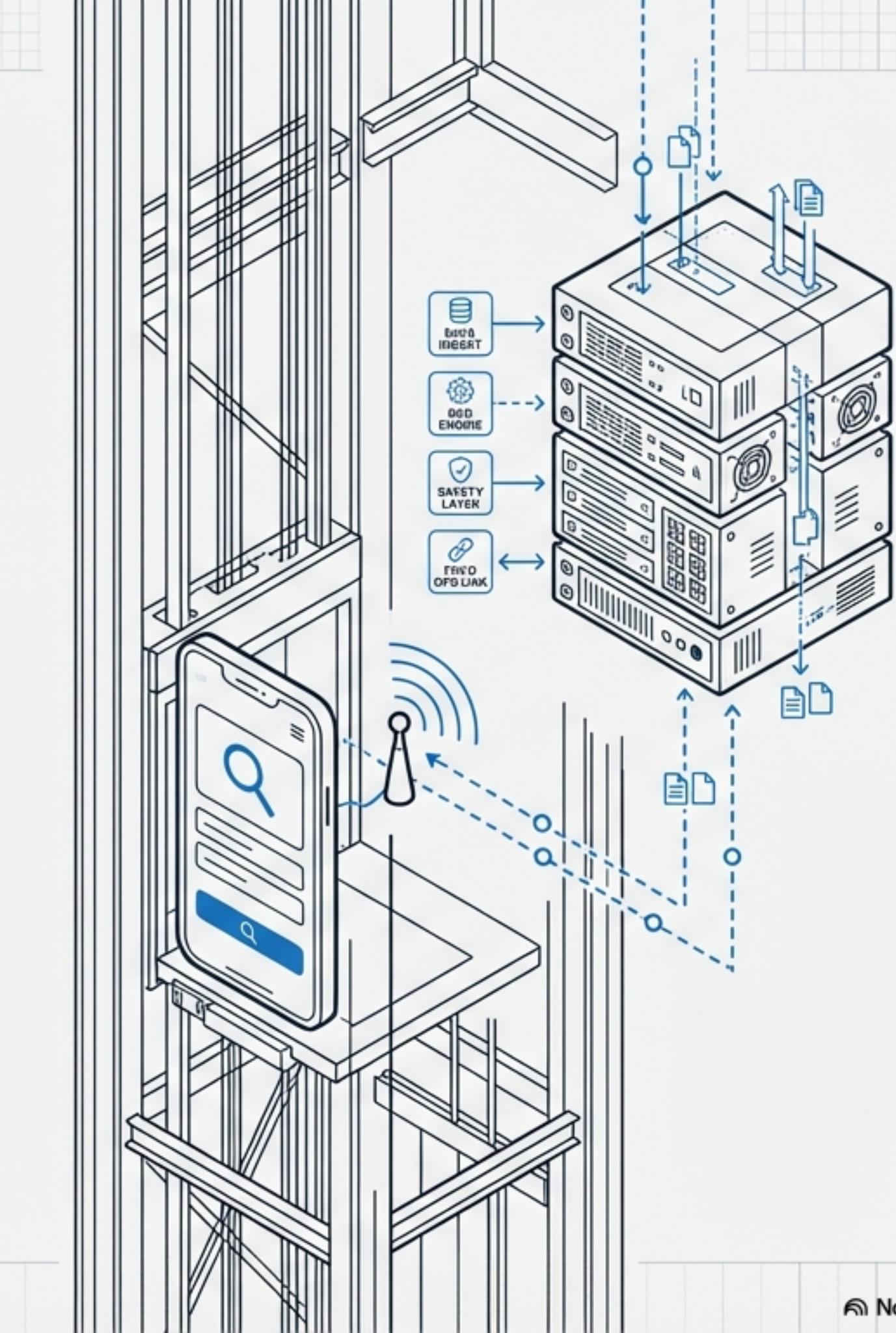


INTERNAL TECHNICAL REVIEW

Enterprise AI Architecture: Technical Implementation Framework

A Safety-First, RAG-Powered Assistant
for Field Operations

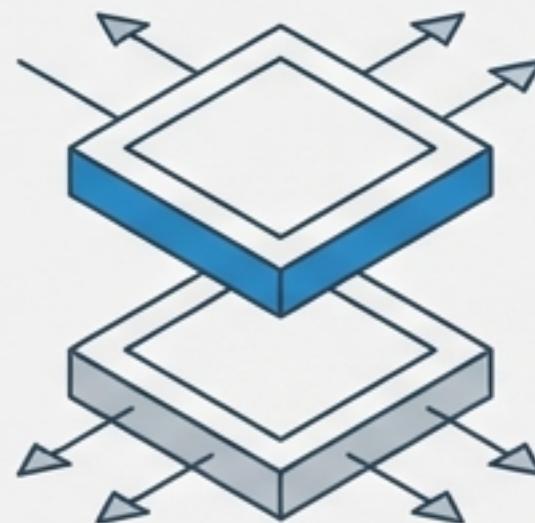
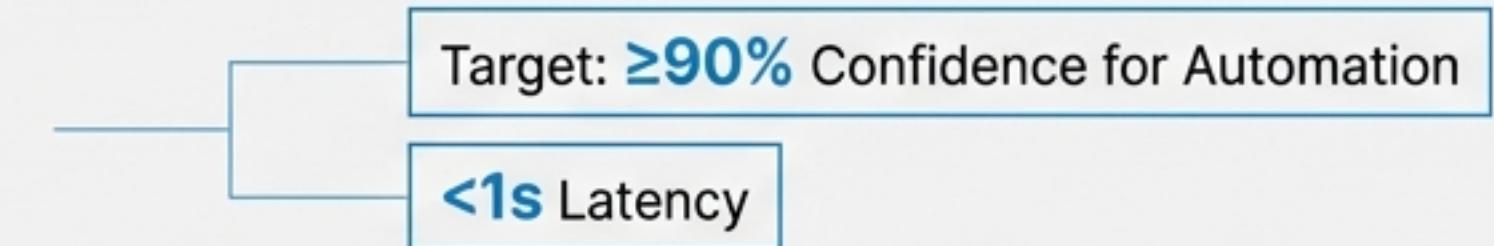
Architecture Validation & Stack Selection



Operational Goals & Strategic Architecture

Core Objective

Reduce mean-time-to-repair while maintaining 100% safety compliance through a 'Human-in-the-Loop' design.



Decoupled Architecture

Frontend/Backend separation allows independent scaling and future mobile migration without refactoring the core logic.



Verifiable Intelligence

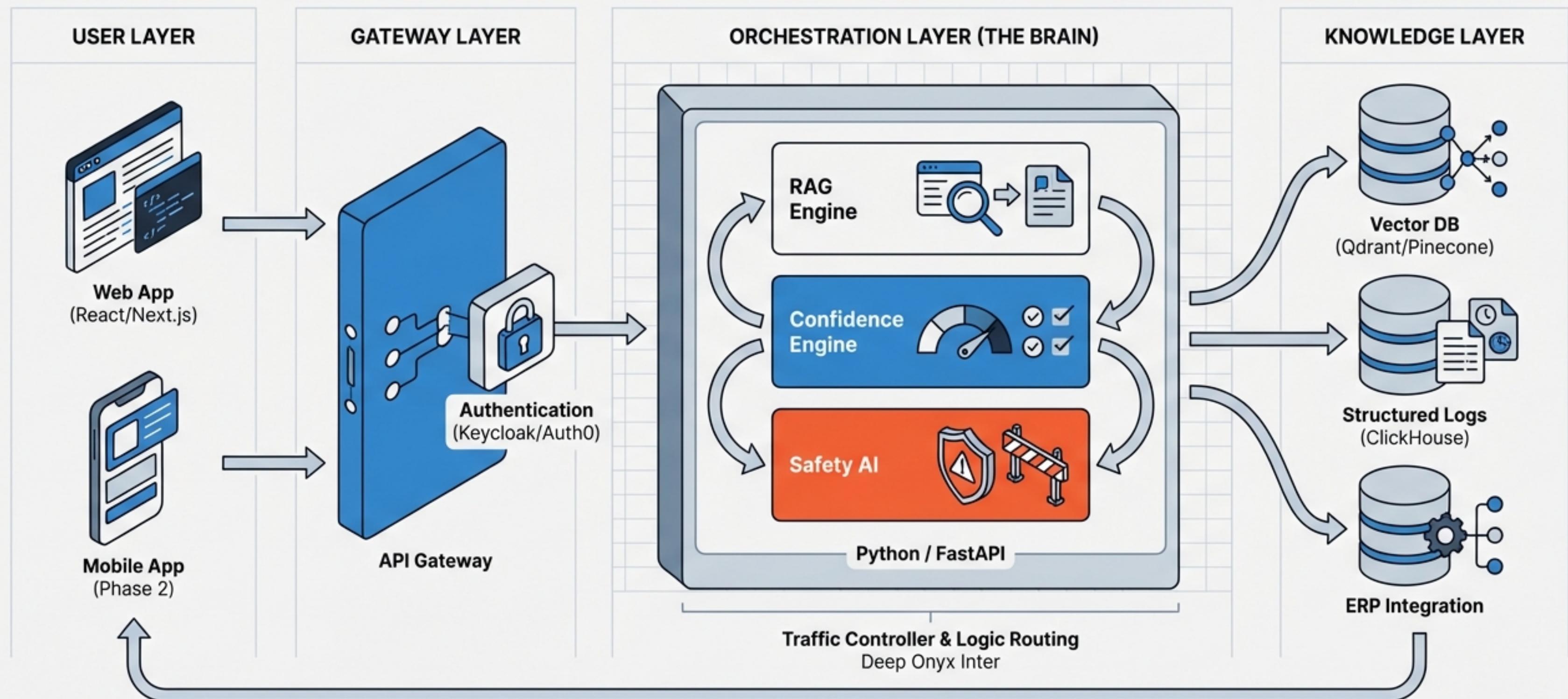
A RAG (Retrieval-Augmented Generation) engine grounded strictly in technical manuals, governed by a multi-vector Confidence Engine.



Non-Negotiable Safety

An always-on 'Sentinel' layer that bypasses AI entirely for emergency keywords (e.g., 'Entrapment', 'Fire') to force human escalation.

High-Level System Architecture



Frontend & User Experience Layer

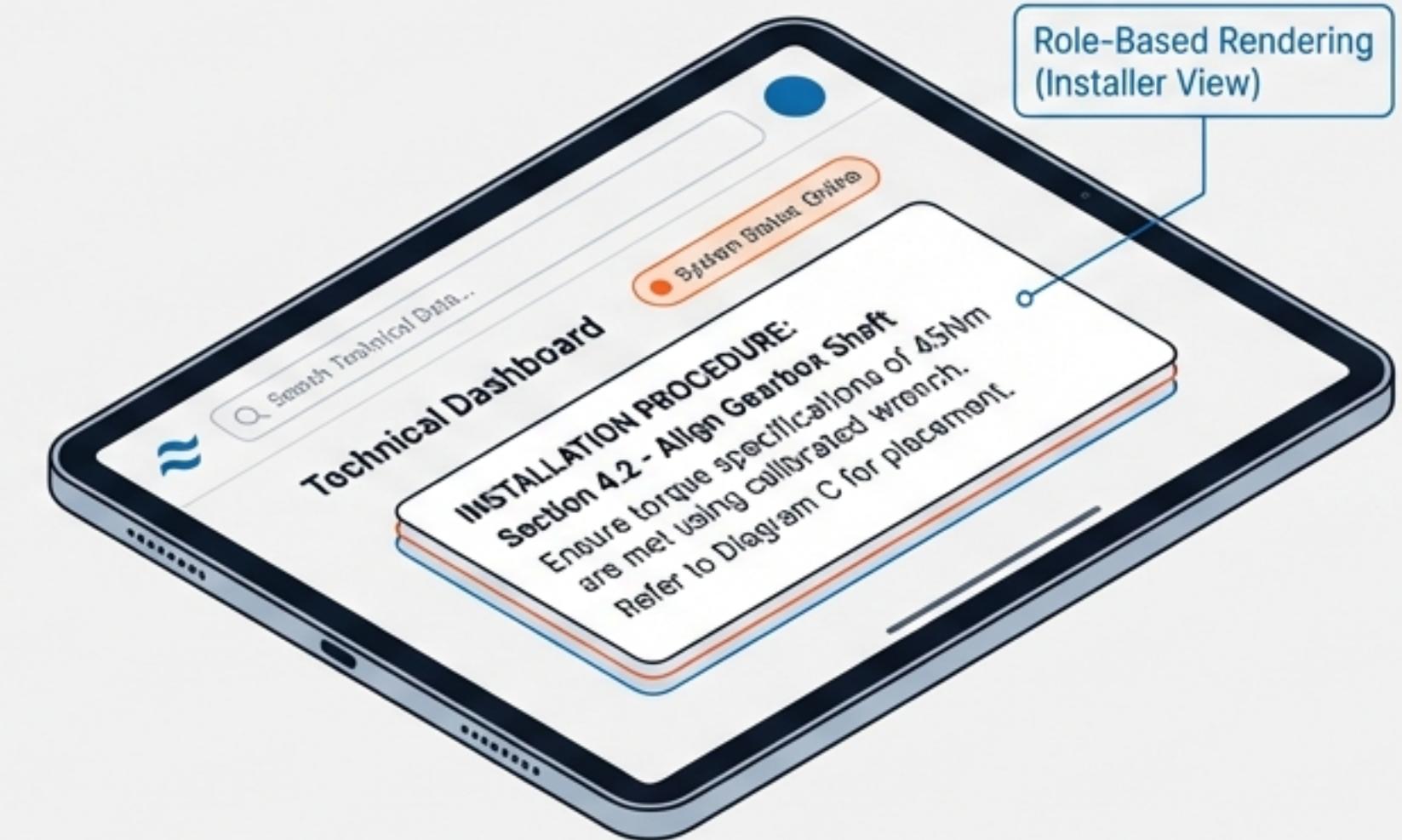
Focus: Performance, SEO, and Field Usability

Primary Stack (Web App)

- Framework: React.js + Next.js (Selected for SSR performance)
- UI Library: Tailwind CSS (Fast, consistent enterprise UI)
- Localization: i18n (react-i18next) for English, German, Greek

Secondary Stack (Mobile App - Phase 2)

- Framework: React Native
- Key Capability: Offline Caching (Critical for shaft work)
- Key Capability: Camera Access (Image Recognition)



Authentication & Access Control Strategy

Technology:

Keycloak
(Self-hosted for sovereignty) or Auth0.

Protocol:

OAuth2 / OpenID Connect.

Compliance:

Full GDPR audit logging.

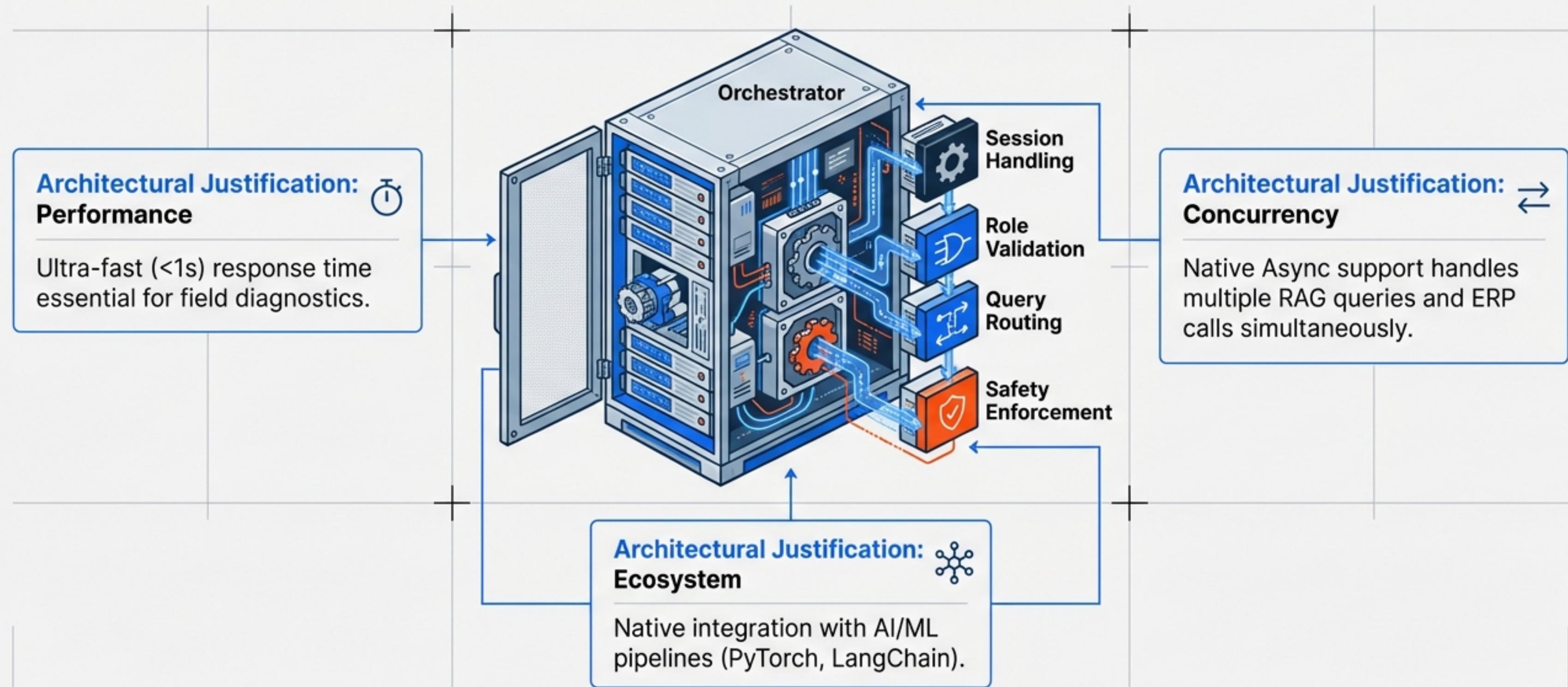
RBAC Matrix

Role Name	System Code	Permissions Scope
Installer	ROLE_FIELD	Read-only access to Manuals & RAG Interface.
Employee	ROLE_INTERNAL	Internal Data Access.
Senior Engineer	ROLE_POWER	Advanced Diagnostics & Override Capabilities.
Admin	ROLE_ADMIN	System Config & User Management.
Management	ROLE_AUDIT	Read-only Access to KPIs & Logs.

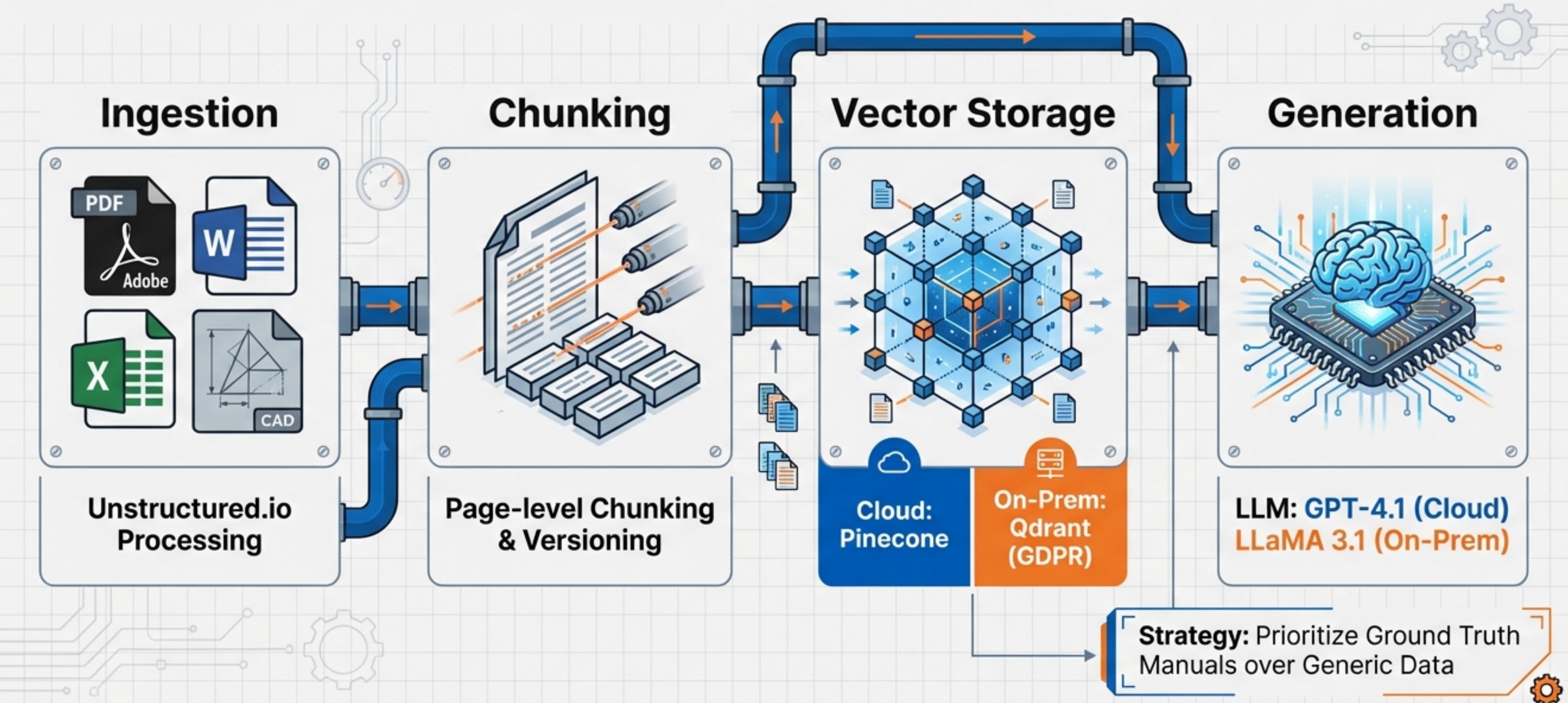


The Orchestration Layer: System Core

Technology: Python + FastAPI



RAG Engine & Knowledge Retrieval Pipeline



The Confidence Scoring Engine

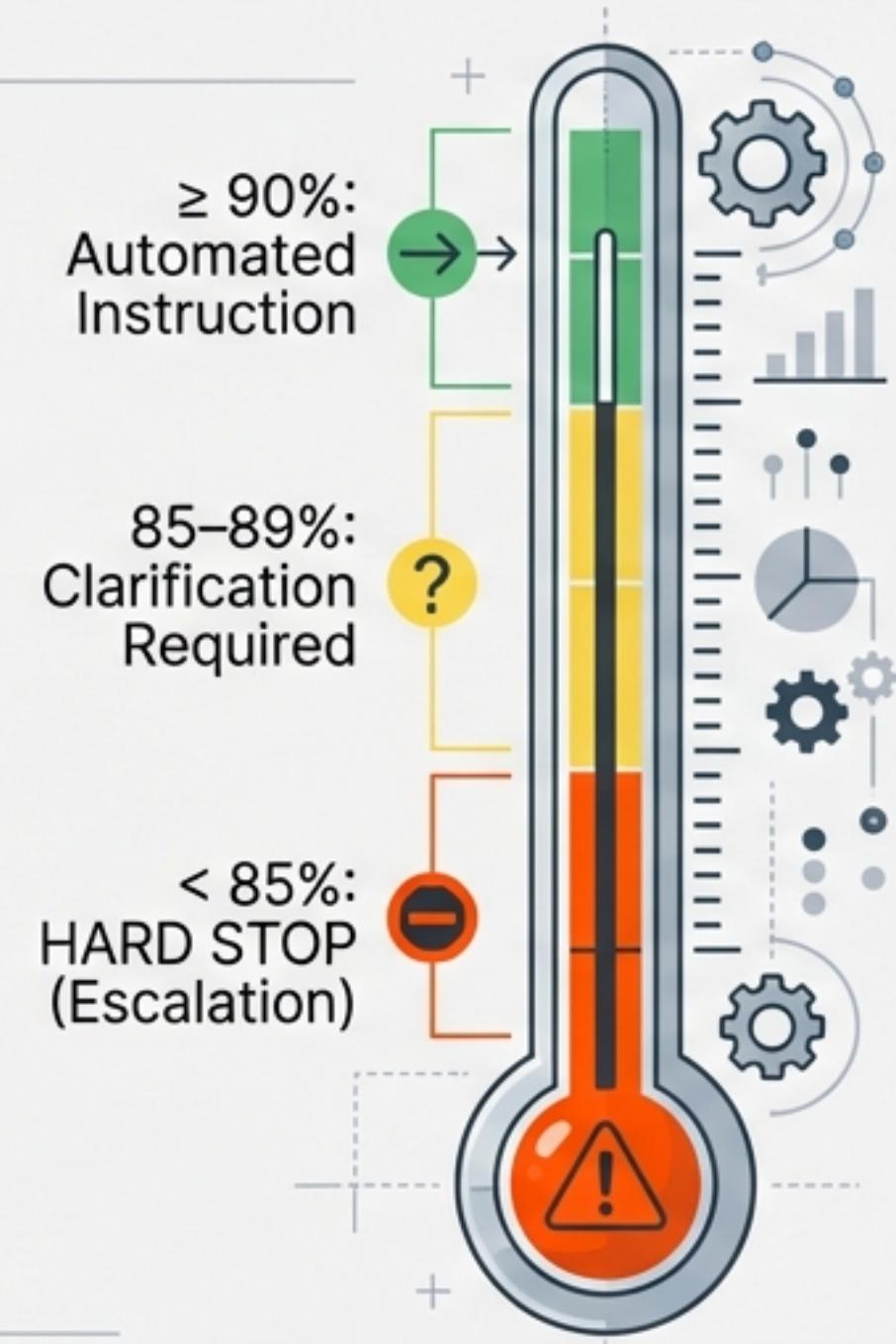
Mathematical Guardrails Against Hallucination

$$\text{Confidence} = (0.4 \times \text{Semantic Similarity}) + (0.4 \times \text{Source Reliability}) + (0.2 \times \text{Model Probability})$$

Semantic Similarity: Cosine score from Vector DB

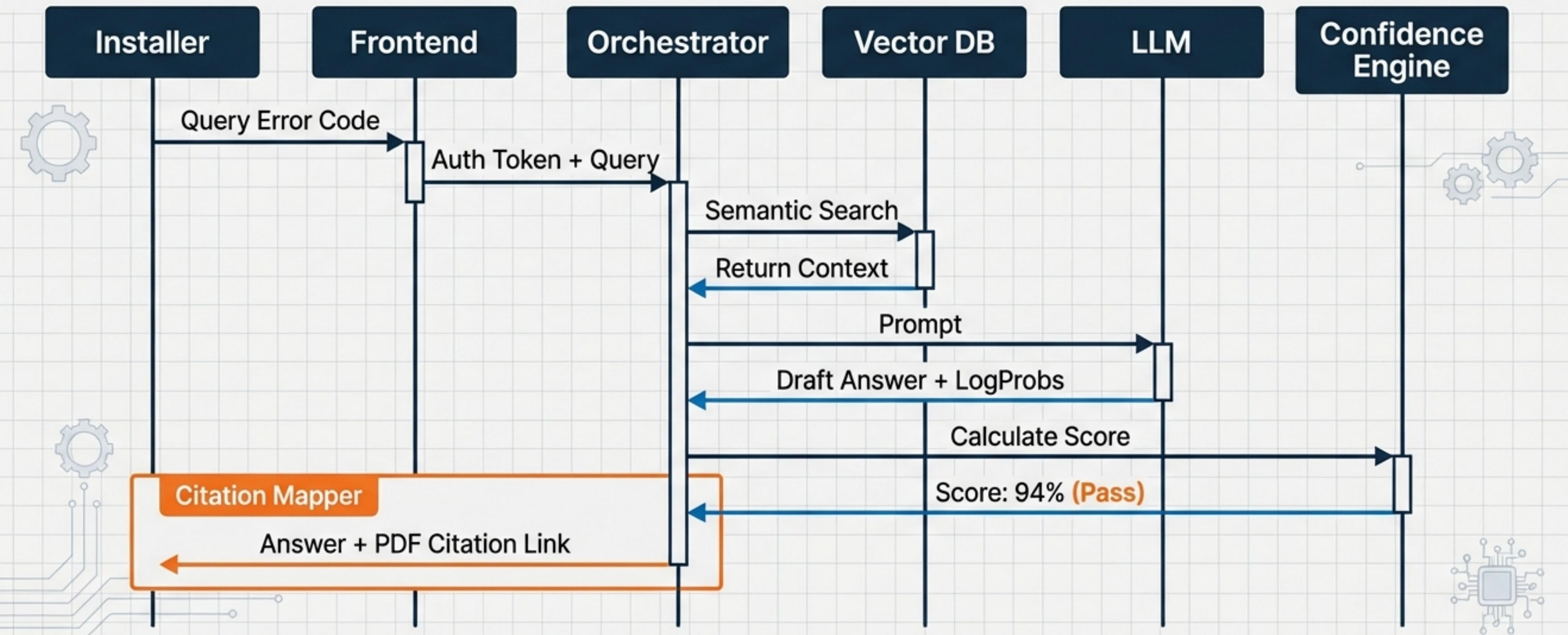
Source Reliability: Metadata weight (Gold Standard PDF > Generic FAQ)

Model Probability: Token entropy / logprobs

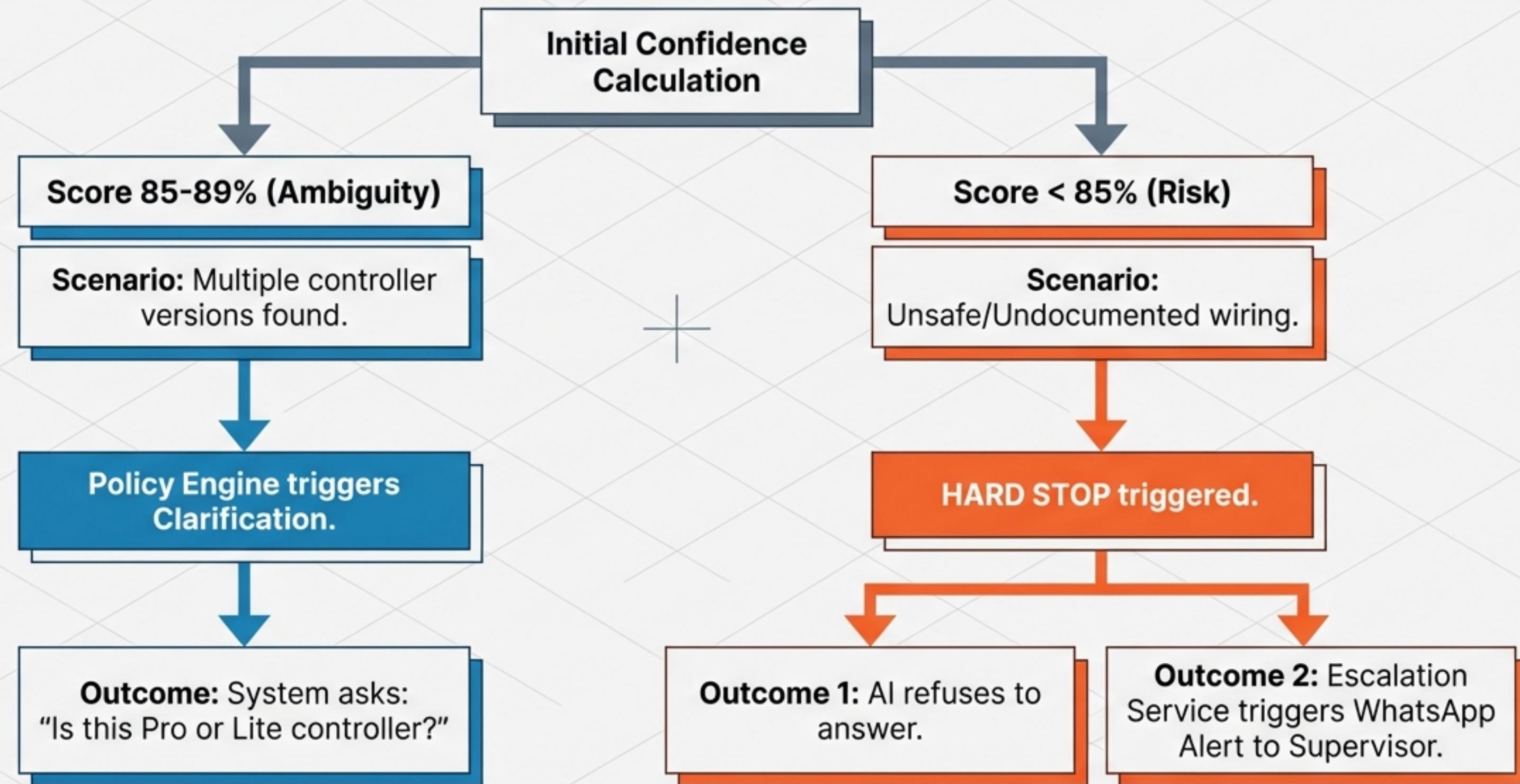


Sequence I: High Confidence Query ('The Happy Path')

Scenario: Installer asks about error code; System matches with $\geq 90\%$ confidence.

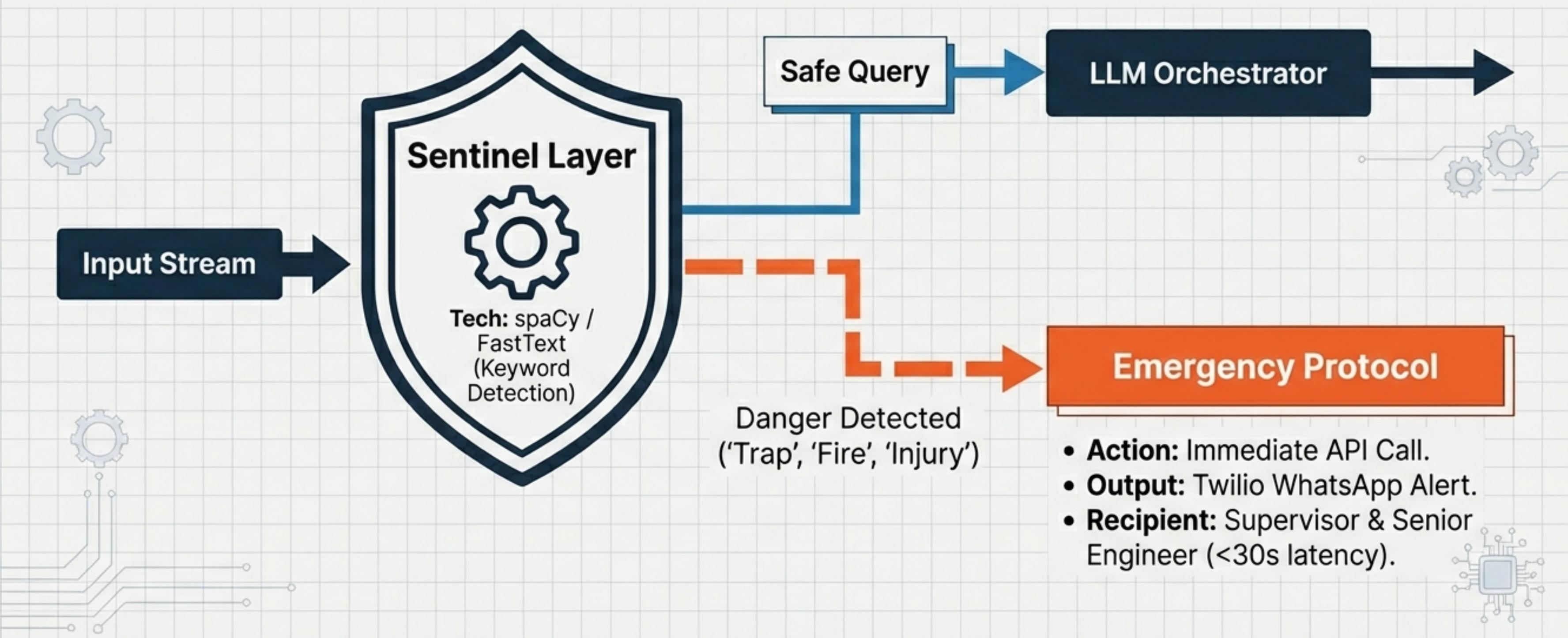


Sequence II: Handling Uncertainty & Hard Stops

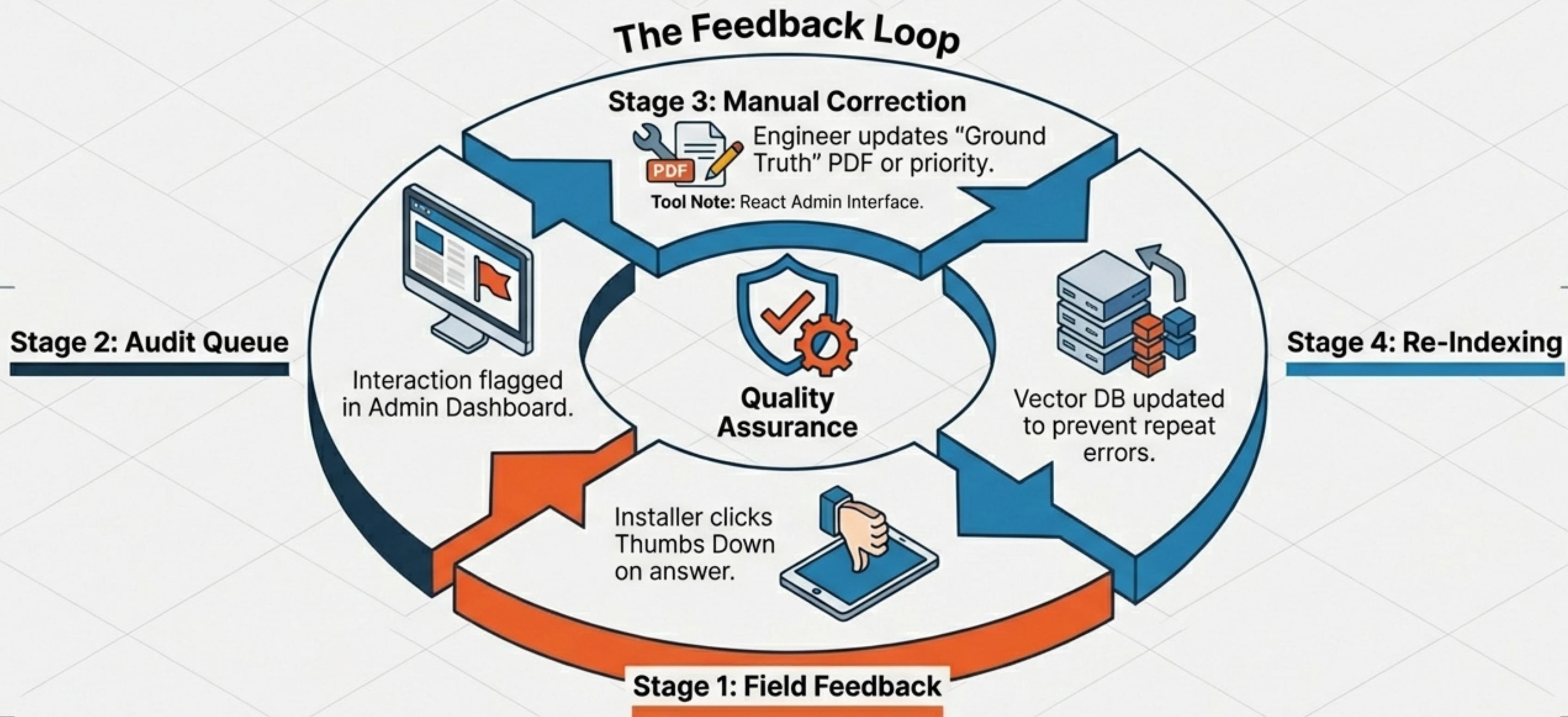


Safety & Emergency AI Layer

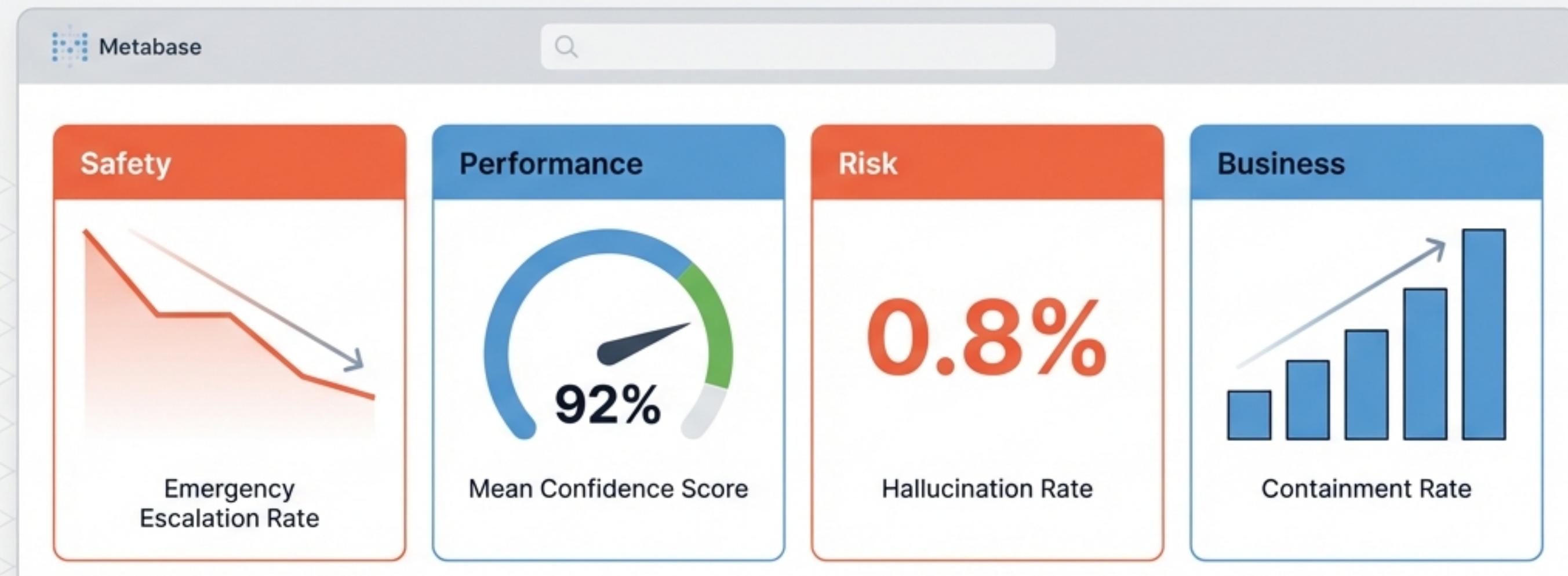
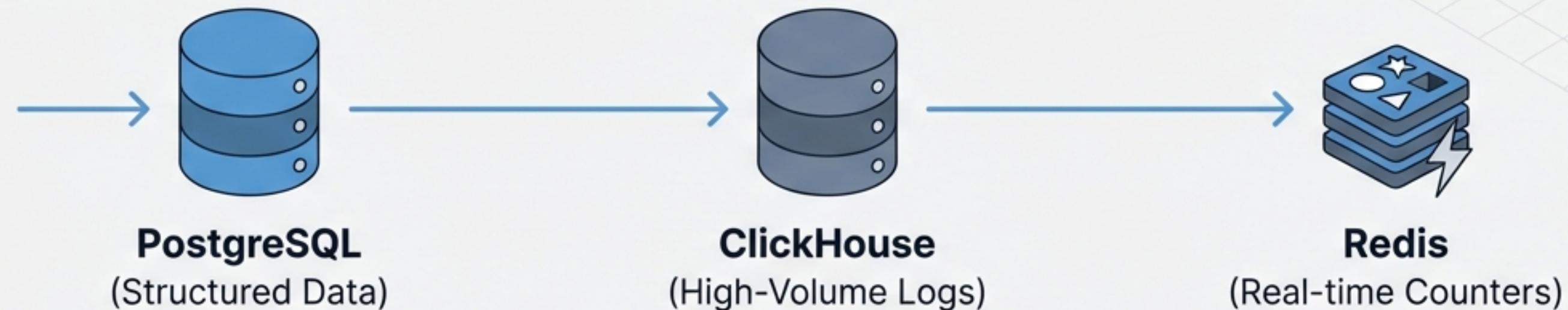
The 'Always-On' Sentinel (Runs Before LLM)



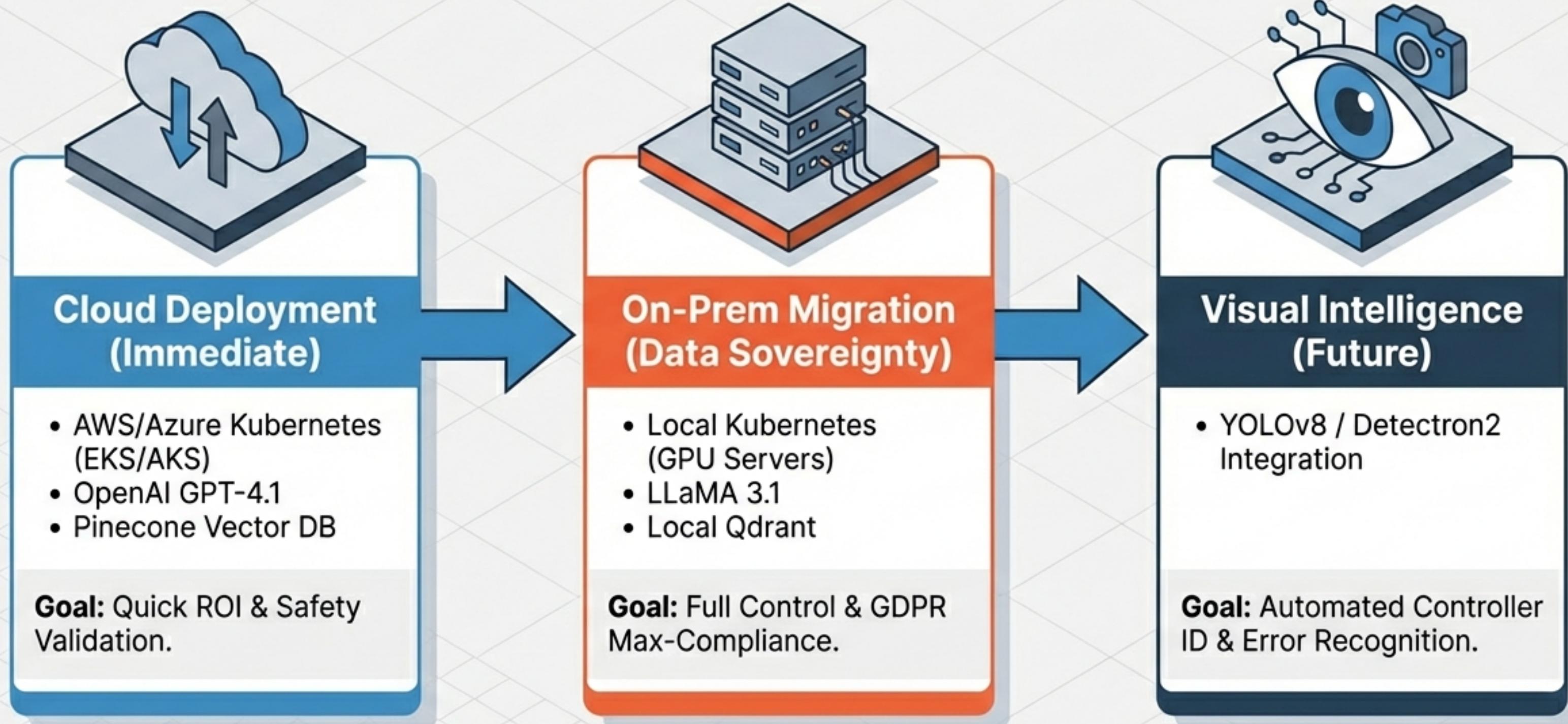
Knowledge Management & Continuous Improvement



Analytics & KPI Engine



Implementation Roadmap: Cloud to On-Prem



Technical Stack Summary & Next Steps

Frontend	React + Next.js (Web), React Native (Mobile)
Backend	Python + FastAPI
Intelligence	GPT-4.1 (Cloud) / LLaMA 3 (On-Prem)
Vector DB	Qdrant / Pinecone
Security	Keycloak / Auth0
Data & Logs	PostgreSQL + Redis + ClickHouse
Safety	spaCy + Twilio WhatsApp
Hosting	AWS / Azure EKS

Recommendation: Proceed immediately with Phase 1 (Cloud) to capture safety data and validate the Confidence Engine before on-prem migration.