

# SaaS Product Intelligence Platform

## Project Summary & Technical Architecture

### 1. Project Overview

The SaaS Product Intelligence Platform is a production-ready ML system designed to answer "Why did [metric] change?" questions using internal company data. It's not a generic chatbot but a sophisticated system that combines retrieval, ranking, constrained LLM reasoning, and continuous learning from user feedback.

### 2. System Architecture

#### Core Pipeline Flow:

Query → Retrieval (Dense + Sparse) → Ranking (ML Model) → LLM Reasoning → Response → Monitoring → Feedback Loop

### 3. Technology Stack

| Component        | Technology                              | Purpose  |
|------------------|---|--|
| API Framework    | FastAPI + Unicorn                       | REST endpoints with async support                    |
| Dense Retrieval  | SentenceTransformers (all-MiniLM-L1242) | Text embeddings, semantic search                     |
| Sparse Retrieval | BM25Okapi (rank-bm25)                   | Keyword-based search                                 |
| Vector DB        | FAISS (Inner Product Search)            | Fast nearest neighbor search                         |
| Ranking Model    | LightGBM with LambdaRank                | Learning-to-rank (pairwise)                          |
| LLM Integration  | Constrained reasoning engine            | Citations, confidence scores, refusal logic          |
| Training         | Feedback-driven retraining              | NDCG label generation, model versioning              |
| Monitoring       | Custom metrics collector                | Latency, recall, NDCG, refusal rate, drift detection |
| Testing          | pytest + pytest-cov                     | Unit tests, integration tests, coverage reporting    |
| CI/CD            | GitHub Actions                          | Automated lint, test, build, deploy                  |
| Containerization | Docker + docker-compose                 | Production deployment                                |
| Data Validation  | Pydantic                                | Schema validation, request/response typing           |
| Logging          | JSON logging (JSONL)                    | Feedback logs, metrics logs, training logs           |

### 4. Core Components

#### Retrieval Layer (Hybrid):

Combines dense semantic search (384-dim embeddings) and sparse keyword search (BM25) to maximize recall. Documents are scored and merged to provide diverse candidate results.

### **Ranking Layer (LambdaRank):**

Uses LightGBM with 6-dimensional features (dense\_score, sparse\_score, doc\_length, term\_overlap, recency\_decay, feedback\_signal). Optimizes for NDCG (ranked relevance metric) rather than binary relevance. Gracefully degrades to weighted combination if LightGBM unavailable.

### **LLM Reasoning (Constrained):**

Synthesizes answers with three critical constraints: (1) Citations required - every claim must cite a source, (2) Confidence scoring - combines doc\_count (40%), rank\_score (40%), answer\_length (20%), (3) Refusal logic - refuses if confidence < 0.5.

### **Feedback Loop:**

All interactions logged to JSONL (query, answer, citations, confidence). User feedback (helpful/not helpful) generates NDCG labels: helpful  $\times$  confidence  $\times$  (not\_refused). Periodic retraining pipeline updates ranking model.

### **Monitoring & Drift Detection:**

Tracks P95 latency, recall, NDCG, refused rate, and mean confidence. Alerts when metrics exceed thresholds: latency > 1000ms, recall < 0.65, refusal > 0.15. Enables production observability.

## 5. API Endpoints

| Endpoint        | Method | Purpose   |
|-----------------|--------|---|
| /query          | POST   | Submit a question, get answer with citations and confidence     |
| /health         | GET    | System health check (status, uptime, drift detection)           |
| /metrics        | GET    | Current performance metrics (latency, recall, NDCG, confidence) |
| /feedback       | POST   | Submit feedback (helpful/not helpful) for a query               |
| /feedback/stats | GET    | Feedback statistics (helpful rate, coverage, refusal rate)      |

## 6. Key Design Decisions

- **Hybrid Retrieval:** Dense + Sparse maximizes recall. Dense captures semantic meaning, sparse captures keyword matches.
- **Learning-to-Rank:** LambdaRank optimizes for ranking order (NDCG), not just relevance classification.
- **Constrained LLM:** Citations + Confidence + Refusal prevents hallucination. System is honest about uncertainty.
- **Feedback-Driven:** Every interaction generates signals for improvement. Continuous learning from users.
- **Monitoring-First:** Drift detection and metrics tracking catch problems in production immediately.
- **Graceful Degradation:** Optional dependencies (LightGBM, TensorFlow) have fallbacks to ensure robustness.

## 7. Performance Profile

**End-to-End Latency:** 200-400ms (typical)

- Dense retrieval: 50-100ms
- Sparse retrieval: 10-20ms
- Ranking: 20-50ms
- LLM synthesis: 100-200ms

**First Query Cold Start:** 2-3 seconds (model loading)

**Throughput:** ~5-10 queries/sec on single machine

**Memory Footprint:** ~4GB (embeddings + FAISS index + models)

## 8. Deployment

**Docker:** Production Dockerfile with multi-stage build, Python 3.11-slim base image

**docker-compose:** Orchestration for local development and testing

**CI/CD:** GitHub Actions pipeline: lint (flake8), test (pytest), build (docker build)

**Scaling:** Stateless API design enables horizontal scaling behind load balancer

## 9. Testing & Quality

**Smoke Tests:** 4/4 passing - validates imports, feedback logging, data validation, configuration

**Unit Tests:** Retrieval, ranking, LLM, pipeline components

**Integration Tests:** End-to-end query flow, feedback loop, retraining pipeline

**Coverage:** pytest-cov for code coverage reporting

**Linting:** flake8 for code quality, black for formatting

Document generated: February 17, 2026