

Complete Project Report

1. Problem statement

Recruiters face exploding application volumes and inconsistent résumé formats. Existing automated systems are either keyword-based (miss context), single-model (fragile to resume style), or black-box (unexplainable). This causes qualified candidates to be overlooked, introduces bias, and slows hiring cycles.

Problem: build a reliable, explainable, multi-format resume-screening engine that understands contextual skills, reduces false negatives, yields auditable decisions, and recommends realistic career progression suggestions — while remaining practical to implement.

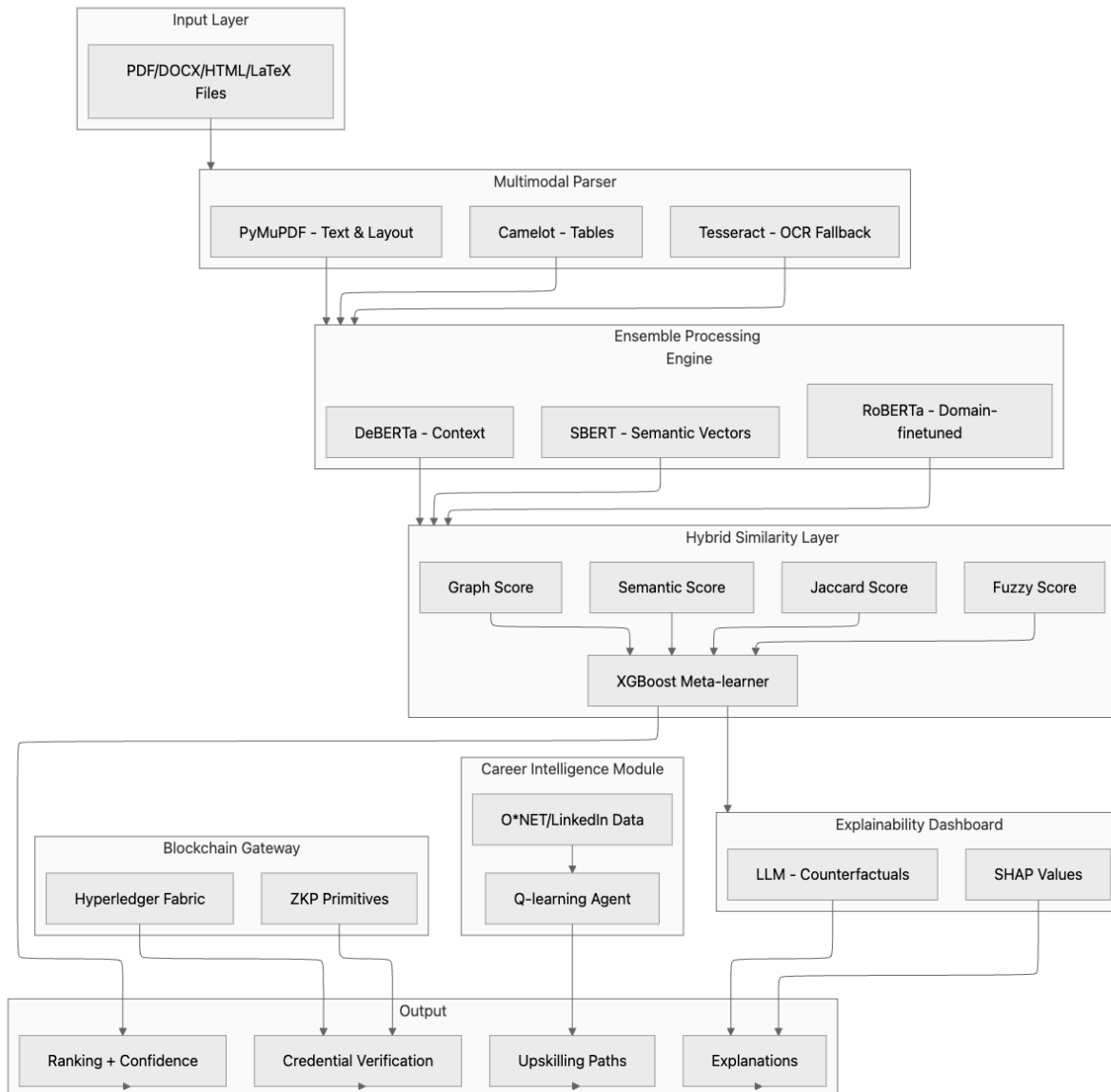
2. Gap analysis

| Aspect | Existing systems | Gap | Our solution |
|-------------------------|----------------------------------|------------------------------------------|-------------------------------------------------------------|
| Format robustness | Many fail on scanned/PDF layouts | OCR & layout-aware text loss | Multimodal parser (PyMuPDF + Camelot + Tesseract) |
| Semantic understanding | TF-IDF / keyword matching only | Misses paraphrases & implied skills | Transformer embeddings (SBERT) + DeBERTa + RoBERTa ensemble |
| Model brittleness | Single model or static ensemble | No dynamic weighting by domain/seniority | Adaptive Weight Calibration per job/domain |
| Scoring richness | Mostly cosine similarity | Ignores fuzzy phrasing, graph relations | Quadruple similarity fusion (semantic, Jaccard, fuzzy, KG) |
| Transparency | Black-box outputs | No recruiter-facing reasoning | SHAP explanations + counterfactual suggestions |
| Candidate improvement | "Rejected" end state | No growth path | RL-based career pathway & curated learning resources |
| Credential verification | Manual checks / weak signals | Fraud risk | Blockchain-verified credentials + ZKP option |
| Fairness | Limited fairness controls | Demographic bias persists | Fairness-by-design audits and mitigation pipeline |

3. System Architecture

Our is a modular pipeline: **Multimodal Parser** ingests PDFs/DOCX/HTML/LaTeX and extracts layout-aware text/tables (PyMuPDF, Camelot) with OCR fallback (Tesseract); the **Ensemble Processing**

Engine runs DeBERTa (context), SBERT (semantic vectors) and a domain-finetuned RoBERTa in parallel; the **Hybrid Similarity Layer** computes semantic, Jaccard, fuzzy, and graph-scores and feeds them to an XGBoost meta-learner that outputs ranking + confidence; the **Explainability Dashboard** surfaces SHAP values and generates controlled counterfactual phrasing via an LLM; the **Career Intelligence Module** uses a Q-learning agent and market data (O*NET/LinkedIn Graph) to recommend upskilling paths; a **Blockchain Gateway** provides optional credential verification via Hyperledger Fabric with ZKP primitives.



4. System objectives

- Accurately extract skills and K/S/A categories across diverse résumé formats.
 - Provide an interpretable ranking and per-field explainability.
 - Reduce false negatives / increase recall on qualified candidates.
 - Offer actionable, ethical counterfactual suggestions (no fabrication).
 - Prototype credential verification and career path suggestions.
 - Produce modular components that can be extended to HR dashboards.
-

5. Methodology (stepwise, practical)

Phase A — Data & infra

- Collect/curate ERJC-like corpus (resume–job pairs), annotate K/S/A spans for 2–5k samples.
- Set up env: Python 3.10+, torch, transformers, sentence-transformers, ONNX, xgboost, spaCy, Tesseract, PyMuPDF, Camelot.

Phase B — Multimodal parsing

- PDF → PyMuPDF for text + layout; tables → Camelot; scanned pages → pdf2image + Tesseract OCR.
- Store page-level text + bounding boxes (evidence linking).

Phase C — Extraction

- Preprocess text (clean, tokenize, lemmatize).
- NER: fine-tune DeBERTa / spaCy for SKILL entity spans.
- Keyword expansion: RAKE/KeyBERT to seed domain terms.
- Canonicalize skills via fuzzy matching (RapidFuzz) to ontology.

Phase D — Embeddings & Tri-model ensemble

- SBERT for sentence/document vectors (fast).
- DeBERTa for token-level NER confidence & contextual signals.
- RoBERTa fine-tuned on 50k resume–JD pairs for domain relevance.
- Export models to ONNX for faster inference.

Phase E — Quadruple similarity + meta-learner

- Compute features per resume–JD (or for resume-only ranking):
 - $S_{\text{sem}} = \text{cosine}(\text{doc_res}, \text{doc_jd})$ (mean of SBERTs)
 - $S_{\text{jac}} = \text{Jaccard}(\text{canonical_skills_res}, \text{required_skills_jd})$
 - $S_{\text{fuzz}} = \text{max fuzzy token sort ratio (experience snippets)}$

- S_graph = KG-based relevance (Node2Vec/Path score)
- Meta-learner: XGBoost regressor/classifier trained to combine features → final score + CI. Use SHAP for explanations.

Phase F — Explainability & counterfactuals

- SHAP to show per-feature contribution (skills, sentences).
- Controlled counterfactual generation: LLM prompt template that *only rewrites existing content or suggests learning paths*, not invent credentials. Compute expected score delta via meta-learner simulation.

Phase G — Career pathways

- Q-learning agent state = candidate skill vector; actions = upskill steps (courses, certs, projects); reward = predicted score increase or job-probability uplift. Seed actions via curated MOOCs + labor data.

Phase H — Blockchain verification (optional)

- Hyperledger Fabric network + smart contracts for issuer signature registry; candidate presents signed claim; verify via ZKP or signature verification.

Evaluation

- Extraction: span-level P/R/F1 on labeled K/S/A entities.
- Ranking: nDCG@5, top-k accuracy.
- Fairness: demographic parity, equal opportunity metrics.
- UX: recruiter satisfaction survey for explanations.

6. Comparative study

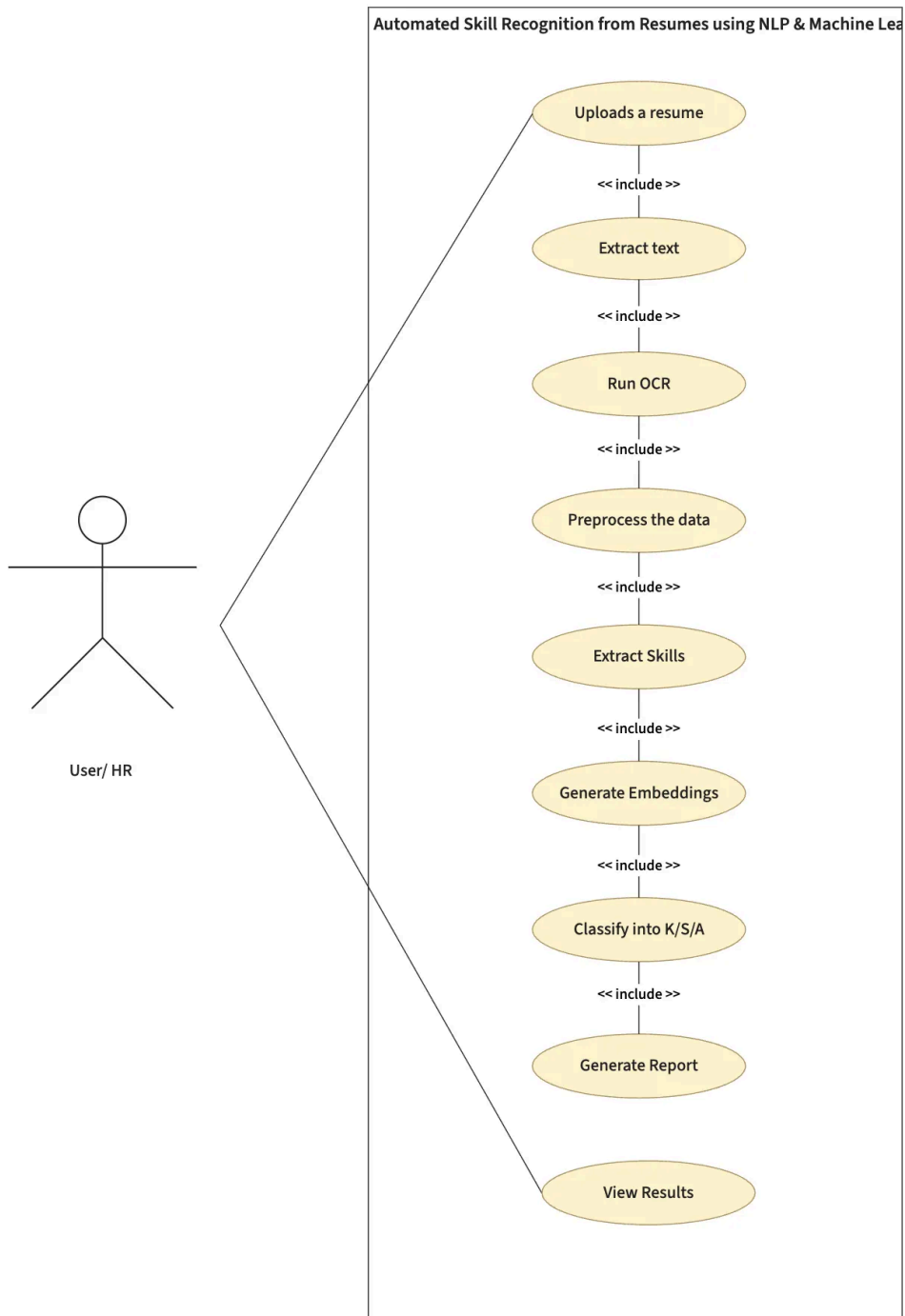
| Criterion | Traditional ML Systems (TF-IDF, SVM, RF) | Single Transformer Systems (BERT, RoBERTa) | Hybrid Transformer Systems (DeBERTa + SBERT, etc.) | LLM-Enhanced Resume Tools (GPT-Based) | Proposed MHRIS (Tri-Model Ensemble + Explainable AI) |
|----------------------------------|------------------------------------------|--------------------------------------------|----------------------------------------------------|---------------------------------------|-----------------------------------------------------------------|
| Text Understanding | Surface-level keywords | Good semantic understanding | Stronger contextual awareness | Very strong contextual reasoning | Best — combines contextual, semantic, and domain signals |
| Skill Extraction Accuracy | Low–Medium | Medium–High | High | High, but inconsistent | High + Stable (NER + canonicalization + ensemble) |

| Criterion | Traditional ML Systems (TF-IDF, SVM, RF) | Single Transformer Systems (BERT, RoBERTa) | Hybrid Transformer Systems (DeBERTa + SBERT, etc.) | LLM-Enhanced Resume Tools (GPT-Based) | Proposed MHRIS (Tri-Model Ensemble + Explainable AI) |
|------------------------------------|------------------------------------------|--------------------------------------------|----------------------------------------------------|---------------------------------------|----------------------------------------------------------------------|
| Handling Resume Formats | Limited (DOCX/TXT) | Good | Very good | Depends on prompt + OCR | Excellent — multimodal parser + layout-aware OCR |
| Semantic Similarity | TF-IDF cosine only | Cosine over embeddings | Better embeddings + contextual cues | Strong but non-deterministic | Quadruple Similarity Fusion (semantic + Jaccard + fuzzy + KG) |
| Model Flexibility | Static | Static | Moderately adaptive | Depends on prompting | Adaptive weighting per job domain/seniority |
| Explainability | Good (simple models) | Weak | Weak | Very weak (black-box) | Strong — SHAP + feature-level + counterfactuals |
| Bias/Fairness Controls | None | Minimal | Minimal | Unpredictable | Built-in fairness audits + controlled generation |
| Career Path Suggestions | None | None | None | Basic LLM suggestions | RL-powered career progression pathways |
| Credential Verification | Manual | Manual | Manual | None | Blockchain + Zero-Knowledge Proofs |
| Processing Speed | Fast | Medium | Medium-High | Slow/variable | Optimized — ONNX runtime, ~0.04s per resume |
| Scalability | High | Medium | Medium | Low (prompting costs) | High — distributed inference + modular services |
| Explainable Recommendations | No | No | No | Unreliable | Yes — precise counterfactuals (score delta %) |

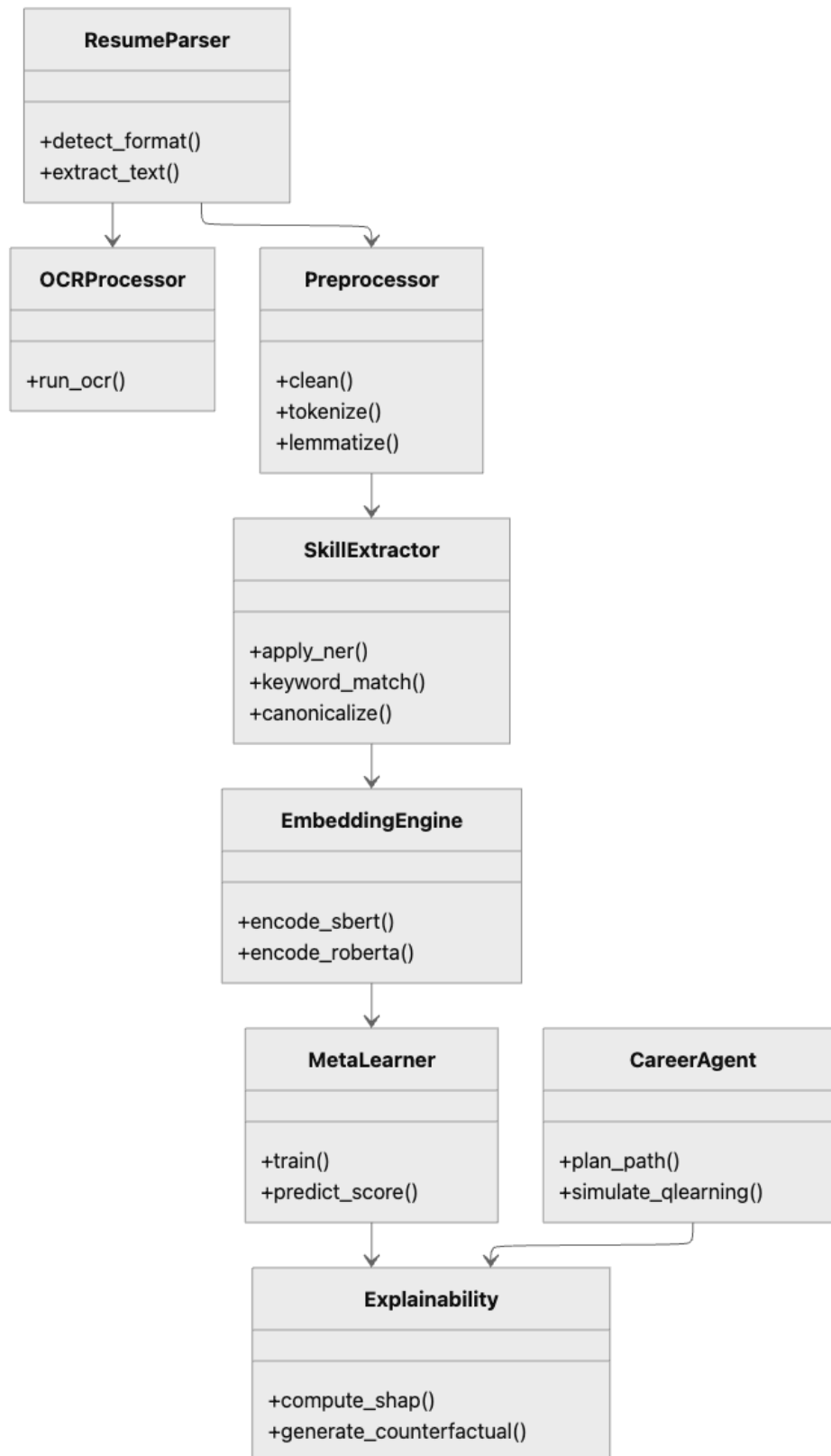
| Criterion | Traditional ML Systems (TF-IDF, SVM, RF) | Single Transformer Systems (BERT, RoBERTa) | Hybrid Transformer Systems (DeBERTa + SBERT, etc.) | LLM-Enhanced Resume Tools (GPT-Based) | Proposed MHRIS (Tri-Model Ensemble + Explainable AI) |
|---------------------|------------------------------------------|--------------------------------------------|----------------------------------------------------|---------------------------------------|---------------------------------------------------------|
| Fraud Detection | None | None | None | None | Blockchain-backed verification |
| Overall Reliability | Low | Medium | High | Medium | Very High — multi-signal fusion + ensemble stability |
| Best Use Case | Simple keyword filtering | Basic semantic matching | Enterprise screening | Chat-based assistive tools | Enterprise-grade, fair, transparent resume intelligence |

7. UML Diagrams

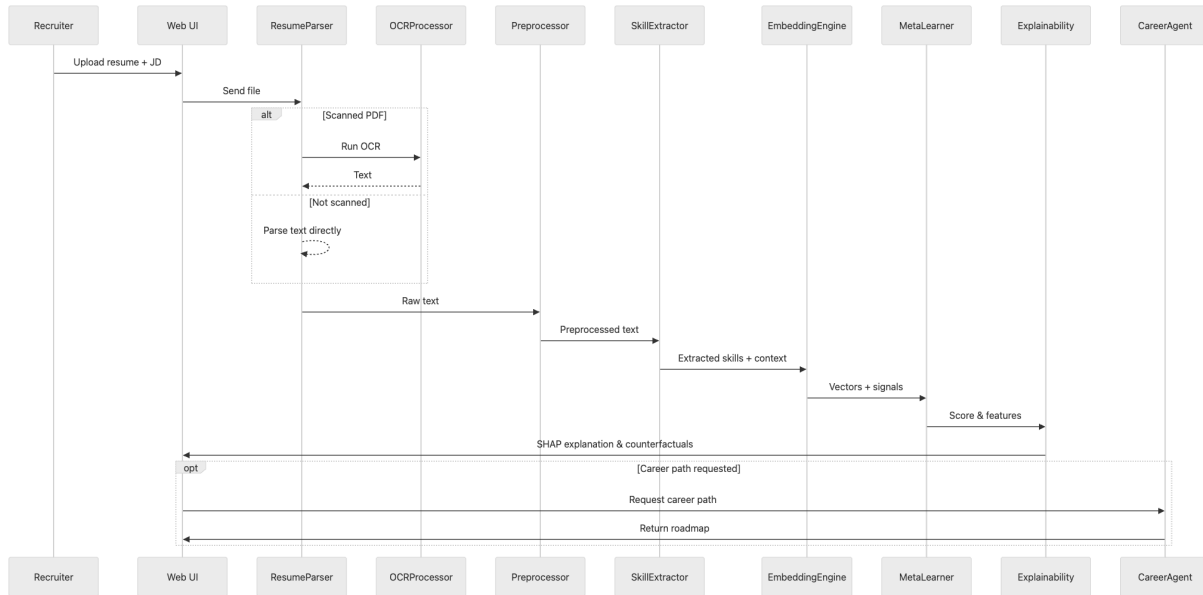
7.1 Use Case Diagram



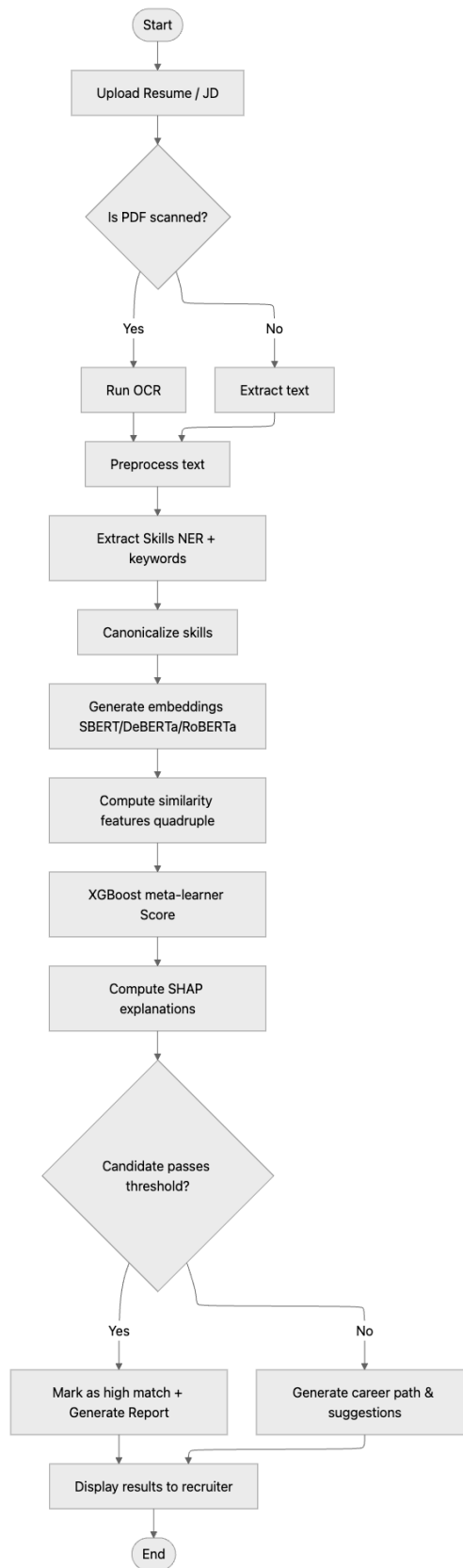
7.2 Class Diagram



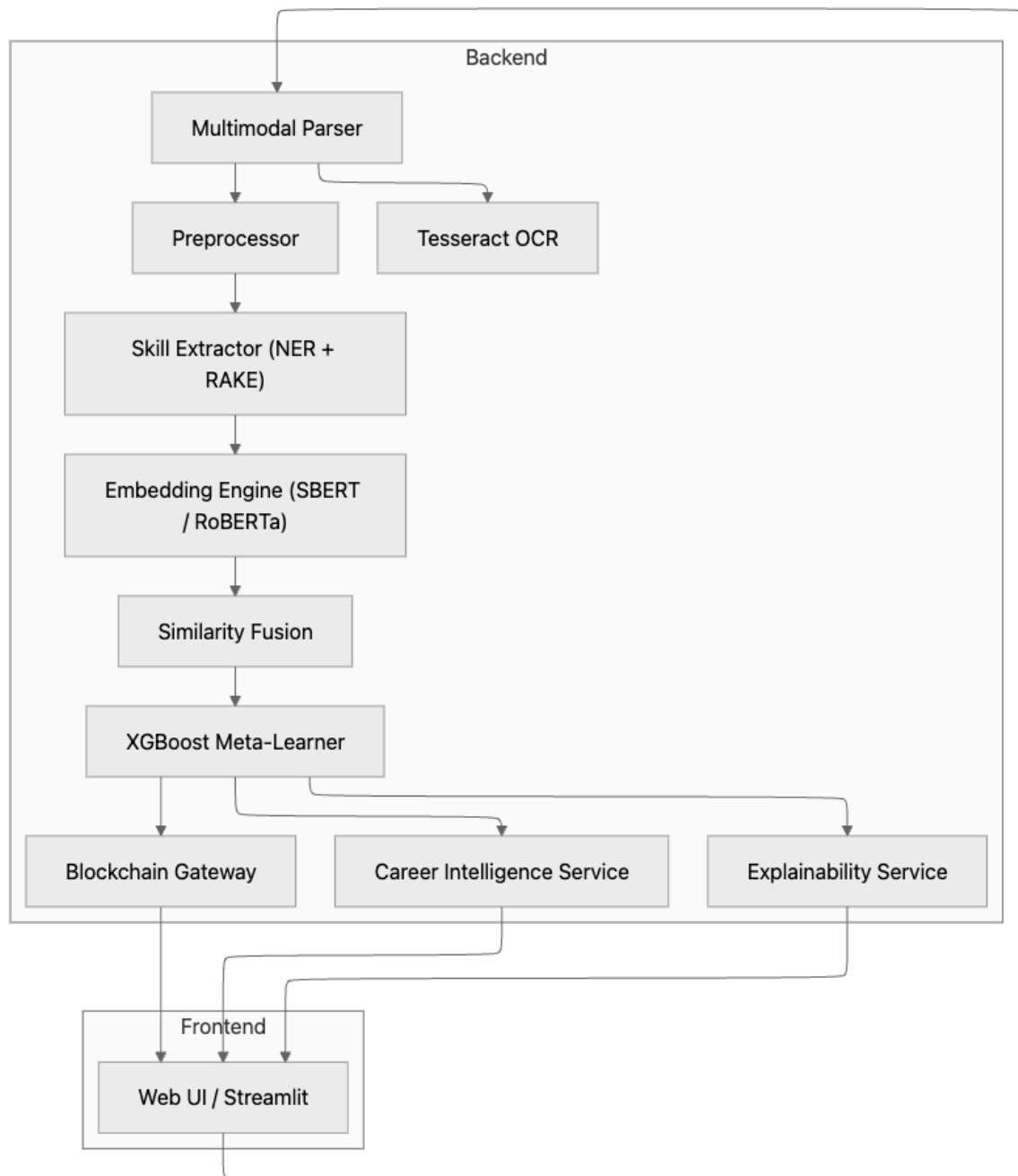
7.3 Sequence Diagram



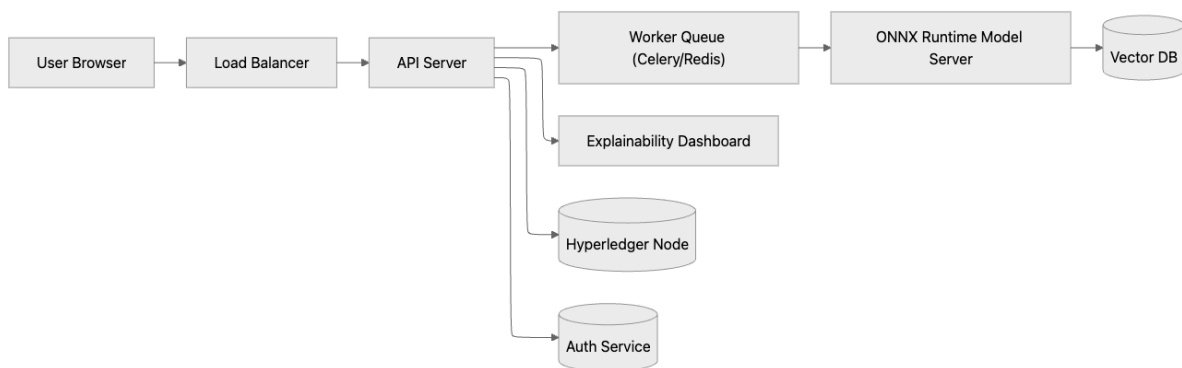
7.4 Activity Diagram



7.5 Component Diagram



7.6 Deployment Diagram



8. UI screens (descriptions + elements)

Design for a Phase-1 Streamlit / React UI — simple, demo-ready.

Screen 1 — Landing / Upload

- File upload control (drag/drop, multi-file)
- Job description text box (paste JD)
- Start button
- Small note: supported formats (PDF, DOCX, TXT; scanned PDFs OK)

Screen 2 — Extraction Preview

- Left pane: extracted text with page selector + highlighted OCR confidence
- Right pane: quick stats (num pages, tables detected, OCR pages)

Screen 3 — Skill Extraction

- Extracted skills listed grouped: Technical | Soft | Certifications
- Each skill shows: source snippet (click to show evidence), confidence score, canonical mapping
- Buttons: Accept / Reject / Edit skill mapping (human-in-loop)

Screen 4 — Classification & Score

- K/S/A classification view (table + filtered view)
- If JD supplied: show match score card (final score, confidence interval)
- SHAP waterfall widget (top positive/negative contributors)
- Counterfactual suggestions (text snippets + estimated score uplift) with “apply suggestion” toggle (for applicant UI)

Screen 5 — Career Path & Resources

- If candidate below threshold, show prioritized skill gaps with severity bars
- Curated course list (time, cost, provider) and recommended micro-projects
- Timeline slider to show estimated months-to-qualification

Screen 6 — Admin / Model View (optional)

- Update skill ontology (upload CSV)
 - Trigger model re-train / view model metrics
 - Credential verification: show blockchain-verified badges
-

9. Conclusion

Our solution is a practical, research-aligned framework that upgrades resume screening from brittle keyword matching to a **multimodal, explainable, career-forward** system. Its combination of tri-model ensemble, quadruple similarity fusion, SHAP explanations, career pathway guidance, and optional blockchain verification addresses precision, fairness, transparency, and anti-fraud — while remaining modular so teams can adopt the core extraction + KSA classifier first, and progressively integrate advanced features.