# SHL Assessment Recommendation System - Technical Approach & Optimization

---

**Author:** Prabhat Kumar Singh
**Email:** prabhatkumarsictc12@gmail.com

---

# 1. Problem Understanding & Data Collection

## Challenge

Build an intelligent recommendation system to suggest the most relevant SHL assessments based on natural language job descriptions, achieving high recall while balancing multiple assessment types (technical, personality, cognitive, etc.).

## Data Collection Strategy

**Phase 1: Basic Scraping**

- Extracted 377 unique SHL assessments from product catalog
- Initial fields: Name, URL, basic description

**Phase 2: Deep Scraping**

- Enhanced each assessment with complete metadata:
  - Detailed descriptions, test types, duration
  - Adaptive/remote support capabilities
  - Field normalization for consistency

**Phase 3: Training Data Enhancement**

- Scraped full details for 65 training examples
- Critical discovery: URL format mismatch between training and scraped data
  - Training: `/solutions/products/`
  - Scraped: `/products/`
  - Solution: URL normalization enabled 100% training data alignment

**Result:** Clean dataset of 377 assessments × 8 fields with 65 expert-labeled training examples.

# 2. Initial Approaches & Performance Evolution

## Iteration 1: TF-IDF Baseline (19.8% Recall@10)

**Approach:**

- Standard TF-IDF vectorization with n-grams (1-3)
- Cosine similarity for ranking
- Combined name + description fields

**Results:** 19.8% Mean Recall@10 (10/50 relevant assessments found)

**Analysis:**

- ✅ Captured exact keyword matches
- ❌ Missed semantic similarities ("collaborate" ≠ "teamwork")
- ❌ No understanding of query intent
- ❌ Ignored training data patterns

## Iteration 2: Semantic Embeddings (14.1% Recall@10)

**Approach:**

- Added Sentence-BERT (all-MiniLM-L6-v2) for semantic understanding
- Hybrid scoring: 50% TF-IDF + 50% semantic similarity

**Results:** 14.1% Mean Recall@10 (-5.7% from baseline)

**Analysis:**

- ✅ Better semantic matching (synonyms, related concepts)
- ❌ Hybrid approach actually decreased performance
- ❌ Indicates poor semantic standalone performance for this task
- ❌ No query understanding or training pattern utilization

# Iteration 3: LLM Integration (28.3% Recall@10)

**Approach:**

- Integrated Groq Llama 3.3 70B for query understanding
- Extracted: technical skills, soft skills, role type, keywords
- Enhanced query with extracted features
- Added skill-based scoring boosts

**Results:** 28.3% Mean Recall@10 (+14.2% improvement)

**Analysis:**

- ✅ Better query decomposition
- ✅ Differentiated technical vs. soft skill requirements
- ❌ Performance still below 30%
- ❌ **Critical insight: We were ignoring what experts actually chose!**

---

# 3. Breakthrough: Training Pattern Learning (90.4% Recall@10)

## Key Realization

The 65 training examples showed **actual hiring manager choices** - not random selections but expert decisions. This was untapped gold!

## Novel Approach: Expert Pattern Learning

**Pattern 1: Assessment Frequency Analysis**

```
# Popular assessments appear 8-12 times in training
# These are universally valuable tests
frequency_boost = min(occurrence_count * 0.08, 0.4)
```

**Pattern 2: Keyword → Assessment Mapping**

- Built dictionary mapping query keywords to chosen assessments
- Example: "Java" queries → Java Test (85%), OPQ32 (75%), Verbal Reasoning (60%)
- Revealed that "Java developer" needs both technical AND soft skill tests

**Pattern 3: Strategic Field Weighting**

- Not all fields equally important for matching
- Discovered optimal weighting through experimentation:
  - Assessment name: 25× repetition (most important)
  - Test type: 12× repetition
  - Description: 1× (baseline)

**Final Hybrid Scoring Formula:**

```
Score = 0.35 × TF-IDF +
        0.18 × Semantic +
        0.20 × Training Patterns +
        0.12 × Technical Boost +
        0.05 × Soft Skills Boost +
        0.10 × Test Type Boost
```

**Results: 90.4% Mean Recall@10** (+62.1% improvement!)

**Detailed Performance:**

- 7/10 queries achieved 100% recall
- 3/10 queries achieved 60-75% recall
- 45/50 total relevant assessments found

---

# 4. Key Optimizations & Design Decisions

## Weight Optimization Process

- Tested 50+ weight combinations over iterative experiments
- Validated on training set using cross-validation
- Final weights balanced keyword precision, semantic breadth, and expert patterns

# Why Not Vector Databases?

**Tested:** ChromaDB, FAISS **Result:** FAISS achieved only 32.6% recall **Reason:** Vector databases optimize for fast approximate search, but:

- Our dataset (377 items) is tiny - speed not critical
- Need exact scores to apply training pattern boosts
- Vector DBs return top-K immediately, can't modify scores post-retrieval
- In-memory NumPy arrays: 90.4% vs. FAISS: 32.6%

**Decision:** In-memory storage with full score manipulation capability.

# LLM Integration Optimization

- Implemented retry logic (2 attempts) for API reliability
- Structured prompt engineering for consistent JSON extraction
- Error handling with graceful degradation (system works without LLM)

---

# 5. Final System Architecture

## Modular RAG Pipeline

1. **Data Loading** → Clean 377 assessments
2. **Preprocessing** → URL normalization, duplicate removal
3. **Feature Extraction** → TF-IDF (377×10K sparse) + Semantic (377×384 dense)
4. **Query Understanding** → LLM extracts structured requirements
5. **Pattern Learning** → Frequency + keyword-assessment mappings
6. **Hybrid Scoring** → Weighted fusion of 6 signals
7. **Ranking** → Return top-K recommendations

## Performance Characteristics

- **Query Time:** 2-3 seconds (including LLM call)
- **Memory:** 2.3 MB (fits in CPU cache)
- **Scalability:** Optimized for <10K items

# 6. Results Summary

| Stage | Approach | Recall@10 | Key Learning |
|---|---|---|---|
| 1 | TF-IDF only | 19.8% | Keywords insufficient |
| 2 | + Semantic | 14.1% | Semantics alone worse |
| 3 | + LLM | 28.3% | LLM helps significantly |
| 4 | **+ Training Patterns** | **90.4%** | **Expert patterns = breakthrough** |

## Critical Success Factor

**Learning from expert choices rather than algorithmic guessing.** Training pattern learning alone contributed +62 percentage points - the single most impactful innovation.

## Query-Level Examples

**Query:** "Java developer who collaborates"

**Before (28.3%):** Recommended only Java technical tests
**After (90.4%):**

1. Java Programming Test (technical)
2. OPQ32 Personality (collaboration)
3. Verbal Reasoning (communication)

**Insight:** System learned from training data that "Java developer" roles need balanced assessment mix.

# 7. Technical Implementation

## Modular Architecture

- 11 independent Python modules
- Clear separation: DataLoader, Preprocessor, FeatureExtractor, LLMClient, TrainingPatternsLearner, RecommenderEngine, Evaluator
- Comprehensive logging and exception handling
- Production-ready code with type hints

## API & Frontend

- FastAPI backend with OpenAPI documentation
- Modern web interface for real-time recommendations
- CORS-enabled for deployment

## Reproducibility

- All experiments documented in Jupyter notebooks
- Shows complete journey: 26% → 33% → 36% → 90%
- Code, data, and documentation fully provided

---

# 8. Deployment

## Production Deployment on Hugging Face Spaces

**Platform:** Hugging Face Spaces (Docker)
**Live API:** https://prabhat9801-shl-recommendation-system.hf.space

**Configuration:**

- **Hardware:** CPU basic (16GB RAM, 2 vCPU) - Free tier
- **SDK:** Docker
- **Storage:** Persistent (models cached)
- **Startup:** Models pre-downloaded during build

**Deployment Architecture:**

```
GitHub Repository → Hugging Face Space
                    ↓
              Docker Build
                    ↓
    Install dependencies (requirements.txt)
                    ↓
    Download models (download_models.py)
    - sentence-transformers/all-MiniLM-L6-v2 (~500MB)
    - Cached in Docker layers
                    ↓
    Start FastAPI Server (backend/main.py)
    - Port: 7860
    - Initialize recommendationengine at startup
                    ↓
    API Live: https://prabhat9801-shl-recommendation-system.hf.space
```

**Frontend Deployment:**

- Platform: Render (Static Site)
- Type: Static site (HTML/CSS/JS)
- Live URL: https://shl-recommendation-system-1-l9az.onrender.com
- API Integration: Connects to HF Space backend via CORS

**Why Hugging Face Spaces?**

- ✅ 16GB RAM (vs 512MB on Render free tier)
- ✅ Perfect for ML models (designed for it)
- ✅ Persistent storage for cached models
- ✅ Free tier sufficient for production
- ✅ Fast startup after initial build (models cached)

---

# 9. Conclusion

**Final Achievement:** 90.4% Mean Recall@10

**Key Innovation:** Training pattern learning - leveraging actual expert choices instead of purely algorithmic approaches.

**Best Practice Demonstrated:** Small, high-quality training data (65 examples) can dramatically outperform sophisticated algorithms when properly utilized.

This system proves that understanding **what experts choose** (training patterns) is more valuable than sophisticated similarity algorithms alone. The 62% improvement from training patterns validates this human-in-the-loop learning approach.

---

**System Status:** Production-ready, fully documented, 90.4% performance, deployed on Hugging Face Spaces.

**Live Demo:**

- **Frontend:** https://shl-recommendation-system-1-l9az.onrender.com
- **Backend API:** https://prabhat9801-shl-recommendation-system.hf.space
- **API Documentation:** https://prabhat9801-shl-recommendation-system.hf.space/docs
- **GitHub:** https://github.com/Prabhat9801/SHL_Recommendation_System