

Project Management RAG System - Complete Setup Summary



Final Working Solution

Your Project Management RAG system is now successfully running! Here's a complete summary of what we accomplished and the startup process.



Issues We Encountered & Fixed

Original Problems:

1. **Missing get_audit_report method** in RAGSystem class
2. **LangChain version conflicts** - _build_model_kwargs import error
3. **Unicode/Transformers errors** - corrupted cache files
4. **Sentence-transformers installation failures**
5. **File naming inconsistency** (.ntml vs .html)

Solutions Applied:

1. **Created compatibility-fixed app.py** with robust error handling
2. **Bypassed problematic embeddings** - used smart keyword search instead
3. **Avoided transformers library entirely** - eliminated Unicode issues
4. **Implemented fallback mechanisms** for different LangChain versions
5. **Added comprehensive file validation** and directory auto-creation



Final Architecture

What's Running:

- **Backend:** Flask server on http://localhost:8081
- **AI Model:** Llama 3.2 3B Instruct (2GB model file)
- **Search:** Smart keyword-based document retrieval
- **Data:** Project management knowledge base (9.8MB text file)
- **Frontend:** HTML interface with real-time status updates

Key Features:

- **No embeddings required** - uses intelligent text matching
- **Works with existing files** - no vector store recreation needed
- **Robust error handling** - graceful degradation
- **Audit logging** - tracks all queries and responses
- **Source attribution** - shows which documents were used



Final File Structure

llm_rag_project/

— app.py	# Fixed Flask backend
— rag-interface.html	# Web interface
— setup_checker.py	# System verification
— quick_fix.py	# Troubleshooting script
— data/	
— pmp_combined.txt	# Training data (9.8MB)
— vector_store/	# FAISS store (not used in final version)
— models/	

```
|   └── llama-3.2-3b-instruct-q4_k_m.gguf # Llama model (2GB)
|   └── logs/                             # Application logs
```

Startup Process

Step 1: Navigate to Project Directory

```
cd llm_rag_project
```

Step 2: Start the Backend Server

```
python3 app.py
```

What happens during startup:






1. **Flask server starts** on port 8081
2. **Initialization thread begins** loading:
 - Project management documents (9.8MB text file)
 - Llama 3.2 3B model (2GB model file)
 - Document chunking and indexing
3. **System status updates** from "Initializing" to "Complete"
4. **Ready for queries** (typically 30-60 seconds)

Step 3: Open Web Interface

Open in browser:

```
file:///path/to/your/llm_rag_project/rag-interface.html
```

Interface Features:

-  **Status indicator** - shows when system is ready
-  **Query input** - ask project management questions
-  **Real-time responses** - powered by Llama 3.2 3B
-  **Source attribution** - shows which documents were used
-  **Response timing** - displays generation time

How It Works

Query Processing Flow:

1. **User submits question** via web interface
2. **Smart search algorithm** finds relevant documents using:
 - Exact phrase matching
 - Keyword intersection
 - Project management term weighting
3. **Context building** - combines top 3 relevant documents
4. **Llama 3.2 3B generates answer** based on context
5. **Response formatting** - clean, structured output
6. **Source attribution** - shows which documents were used

Search Intelligence:

- **Exact phrase matching** (highest priority)
- **Keyword overlap scoring**
- **PM-specific term boosting** (project, risk, scope, etc.)
- **Partial word matching**
- **Document length consideration**

Usage Examples

Effective Queries:

- "What are the key phases of project management?"

- "How do you handle project risks?"
- "What is stakeholder management?"
- "Explain the project scope planning process"
- "What are the cost management techniques?"

System Performance

Typical Response Times:

- **Model loading:** 30-60 seconds (one-time)
- **Query processing:** 2-10 seconds
- **Document search:** <1 second
- **Response generation:** 1-9 seconds

Resource Usage:

- **Memory:** ~4-6GB RAM
- **Model size:** 2GB on disk
- **Data size:** 9.8MB knowledge base

Monitoring & Debugging

API Endpoints:

- GET / - Server status
- GET /api/status - Initialization progress
- POST /api/query - Submit questions
- GET /api/audit - Usage analytics

Log Files:

- **Console output** - real-time status
- **Audit logs** - logs/audit_YYYYMMDD_HHMMSS.json

Success Factors

Why This Version Works:

1. **Simplified architecture** - removed problematic dependencies
2. **Smart fallbacks** - multiple compatibility layers
3. **Robust error handling** - graceful failure modes
4. **Optimized for your setup** - works with existing files
5. **No external dependencies** - just your model and data

The system is now fully operational and ready for project management queries!

