# Knowledge Extraction & Retrieval

| Step | Task | Tools / Models | |
|---|---|---|---|
| 1. Data Ingestion | Text Documents | • PyMuPDF<br>• PDFPlumber<br>• Apache Tika<br>• Azure Form Recognizer<br>• Google Document AI<br>• GPT-4.1 / GPT-4o]<br>• Google Gemini 2.0 Flash / Pro | |
| | Image-Based Documents | • Tesseract<br>• LayoutLM<br>• TrOCR<br>• EasyOCR<br>• PaddleOCR<br>• Donut<br>• Google Vision OCR<br>• Microsoft Azure OCR<br>• Claude 3.5 Sonnet | |
| | Internal System Text | Working on it | |
| 2. Preprocessing | • Cleaning & Normalization | • spaCy<br>• NLTK | |
| 3. Domain Understanding | • Named Entity Recognition (NER)<br>• Key Phrase & Topic Extraction<br>• Relation Extraction<br>• Knowledge Graph Construction | • spaCy custom NER<br>• HuggingFace transformers<br>• Prodigy (annotation tool) | • RoBERTa-NER<br>• BioBERT (medical NER)<br>• FinBERT (finance NER)<br>• med7 (medical entities)<br>• BERT-base multilingual NER |

| 4. Embedding Generation | • Generate Text Embeddings | | • OpenAI text-embedding-3-large / 3-small<br>• Llama 3 Instruct Models<br>• e5-large-v2<br>• all-mpnet-base-v2 |
|---|---|---|---|
| 5. Vector Database | • Store Text Embeddings | • Pinecone<br>• Qdrant<br>• Milvus<br>• Redis Vector Store<br>• Weaviate<br>• ChromaDB | |
| 6. Retrieval System | • Semantic Search<br>• Hybrid Search (Keyword + Semantic)<br>• Contextual Retrieval (RAG) | • LangChain retrievers<br>• LlamaIndex<br>• Qdrant/Pinecone query APIs<br>• Elasticsearch Hybrid Search | ReRankers:<br><br>• Cohere ReRank<br>• bge-reranker-large<br>• cross-encoder/ms-marco-MiniLM |
| 7. Content Generation (LLM Layer) | • Document Summarization<br>• Q&A From Knowledge Base<br>• Drafting New Content<br>• Multi-step Reasoning | | • GPT-4/5<br>• GPT-4o<br>• Llama 3 70B<br>• Claude 3.5 Sonnet |
| 8. Continuous Learning & Feedback Loop | • Reinforcement learning from human feedback<br>• Data augmentations<br>• Active learning<br>• Flow improvements | | • Periodic fine-tuning of BERT/LLM models<br>• Lightweight adapters (LoRA/QLoRA) |