

RAG-GP Project – Phase 2 Instructions (Weeks 1–4)

Objective

Advance from a multi-format retriever prototype to a functional CLI-based RAG system capable of producing LLM-generated, context-aware answers over a comprehensive Alma corpus — including image-aware and XLSX-aware retrieval paths.

Current Status (as of now)

- ✓ Ingestion pipeline supports: **PDF, DOCX, PPTX, TXT, XLSX**
 - ✓ Chunking rules applied per format; FAISS indices built per type
 - ✓ Local and OpenAI embedding modes fully functional
 - ✓ Retrieval pipeline operational with **late fusion**, **per-type recall**, and **query embedding** via `text-embedding-3-large`
 - ✓ CLI command `retrieve` works across mixed-type corpus with debug output
 - ✓ Retrieval metadata preserved (e.g., `doc_id`, `doc_type`, etc.)
-

4-Week Plan (Phase 2: Retrieval→Answer MVP)

Week	Focus Area	Tasks	Output
1	Corpus Expansion	<ul style="list-style-type: none">- Collect 50–100 Alma documents across key modules (Acquisitions, Analytics, Fulfillment, etc.)- Index by format and theme- Store in <code>input/</code> per project	Expanded <code>input/</code> , organized by topic and format
1	Image-Aware Setup	<ul style="list-style-type: none">- In PDF/DOCX chunkers: detect embedded images- Extract them to disk/cache dir- Bind to parent chunk ID/location	Images linked to chunk metadata (<code>image_path</code> , <code>image_page</code>)
2	OCR Extraction for Screenshots	<ul style="list-style-type: none">- Run OCR via <code>tesseractocr</code> or <code>pytesseract</code>- Store text in <code>image_ocr_text</code> field of chunk metadata	OCR data stored alongside chunk content
2	Evaluate Retriever at Scale	<ul style="list-style-type: none">- Run <code>retrieve</code> CLI over new corpus- Save top-K results and log chunk source types	Retrieval report, sample queries with diverse results

Week	Focus Area	Tasks	Output
2	LLM Prompt Template	<ul style="list-style-type: none"> - Draft first prompt format: QA with citations using top-K chunks - Optionally inject OCR text 	<code>prompt_builder.py</code> base implementation
3	Implement <code>ask()</code> CLI Endpoint	<ul style="list-style-type: none"> - Add CLI: <code>ask <project_path> <query></code> - Retrieves top-K chunks, formats prompt, calls OpenAI/GPT 	End-to-end CLI answer generation
3	Image-Inclusive Answering	<ul style="list-style-type: none"> - Concatenate <code>chunk.text</code> + <code>image_ocr_text</code> for embedding + prompt inclusion - Format screenshots descriptively if OCR present 	Answers influenced by visual context
4	Evaluate Answer Quality	<ul style="list-style-type: none"> - Manual analysis of 10–20 queries - Label precision, hallucination, citation quality - Mark XLSX or multi-source failure cases 	Annotated QA sample set, metrics summary
4	Begin XLSX-Aware Chunker	<ul style="list-style-type: none"> - Refactor chunker to process per-table - Capture headers, rows, and normalize content - Optionally summarize table intent 	<code>XLSXChunkerV2</code> , structured chunks with field-level clarity
4	Start Agent Hub Skeleton	<ul style="list-style-type: none"> - Build <code>AgentHub</code> interface with hooks - Add <code>RerankerAgent</code> (e.g., re-sort chunks by alignment score) - Plan <code>SynthesizerAgent</code> (cross-chunk integration) 	Agent interface and minimal reranker stub

Component Ownership and Notes

Component	Owner / Task Lead	Location / Notes
Corpus curation	You	Organize under <code>data/projects/demo_project_full/</code>
OCR integration	Chunker maintainer	Can reuse image utils from PDF chunker, extend to DOCX
Prompt templates	<code>prompt_builder.py</code>	Start with simple QA+source format; plan for citation markup
CLI commands	<code>app/cli.py</code>	<code>ask()</code> should mirror <code>retrieve()</code> in logging/debug output
LLM gateway	<code>scripts/llm/gateway.py</code>	First version can assume OpenAI; later make model pluggable

Component	Owner / Task Lead	Location / Notes
Agent Hub	<code>scripts/agents/hub.py</code> , registry pattern	Each agent gets access to query, chunks, and config

Optional Enhancements (if ahead of schedule)

Task	Description
GPT-4o or Gemini captioning	Generate natural-language captions for key screenshots
Table-to-text converter	Generate per-table descriptions from XLSX sheets
Fusion analysis script	Visualize score distribution across doc types
Streamlit debug viewer	Visualize retrieval + chunks + scores per query

Immediate Next Steps (Week 1 Focus)

1. **Organize and load Alma instructional documents (target: 50–100)**
2. **Enable screenshot extraction + OCR from PDFs and DOCX**
3. **Log and verify image/ocr attachment per chunk**
4. **Prepare 5–10 pilot queries for round-trip `retrieve()` test**