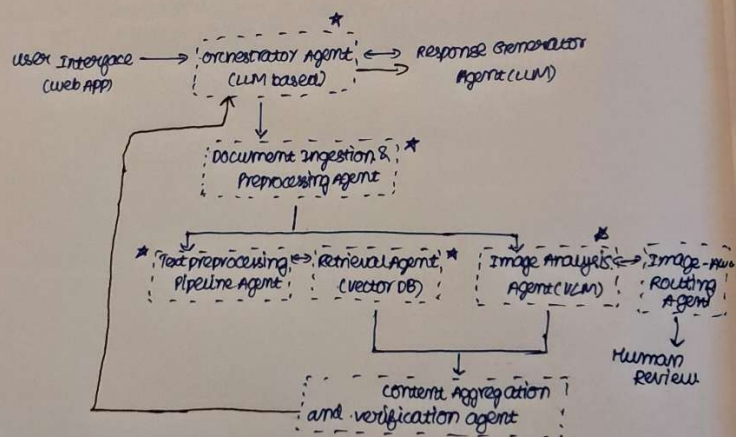


The provided Jupyter Notebooks outline a multi-phase project for building a Retrieval Augmented Generation (RAG) system that can process PDF documents containing both text and images.



The below image states the intended flow of work while creating the multi agent system

**Phase 1 (Imago\_AI\_Assignment\_Phase\_1\_with\_Image\_Handling.ipynb):** This phase focuses on extracting text from PDFs and identifying the presence and location of images.

It defines an EnhancedPyMuPDFLoader to load PDF documents and embed image placeholders (e.g., [IMAGE: img\_1\_1]) directly into the text content. Metadata about these images (ID, page number, position, format, dimensions) is stored.

An ImageAwareTextSplitter is used to chunk the documents while ensuring that image metadata is correctly associated with the chunks containing the corresponding image placeholders.

Text chunks are then embedded using a HuggingFace sentence transformer (intfloat/multilingual-e5-large-instruct) and stored in a Pinecone vector store.

A RetrievalAgent fetches relevant text chunks based on a user query.

An ImageAgent processes these retrieved chunks to extract information about any images mentioned.

A SimpleOrchestrator coordinates the retrieval and image identification steps.

Finally, an EnhancedResponseGenerator uses a Google Generative AI model (Gemini) to generate an answer. The prompt to the LLM includes the retrieved text and a summary of identified images and their context.

```
results = orchestrator.process_query(question)
print(f"Found {len(results['image_info'])} images in the retrieved documents")
for img in results['image_info']:
    print(f"Image ID: {img['image_id']}")
    print(f"Context: {img['document_context']}...")
```

[97]

... Based on the context provided:

Kernels are defined as functions that are called from the CPU code and executed in the GPU. Kernels return void. Different kernels can be launched with different execution instructions. During Found 2 images in the retrieved documents

Image ID: img\_46\_1  
Context: Context not available...

Image ID: img\_46\_2  
Context: Context not available...

**Phase 2 (Imago\_AI\_Assignment\_Phase\_2\_with\_Advanced\_Image\_Handling.ipynb):** This phase enhances the system by performing OCR on the images to extract textual content from them.

Images are first extracted from the PDF and saved.

OCR (using Tesseract) is performed on these images. The extracted OCR text is then associated with the corresponding document page in its metadata (e.g., metadata['image\_data'] = {'image\_p1\_1': 'ocr text from image 1 on page 1'}).

The document pages (now with OCR data in their metadata) are embedded and stored in Pinecone.

The RetrievalAgent retrieves relevant document pages.

A new ResponseGenerator is introduced. This generator:

Retrieves documents based on the query.

For each retrieved document, it calculates the cosine similarity between the page's text content and the OCR text of images found on that page using the same embedding model.

If the similarity is above a certain threshold, the OCR text from those relevant images is appended to the page's main text content before being passed to the LLM (Gemini 1.5 Pro).

This allows the LLM to directly use the textual content of relevant images when formulating an answer.

L2 Access Management)

Displaying: Images/image\_p10\_4.jpeg

| GPU                                 | Kepler GK180      | Maxwell GM200 | Pascal GP100 | Volta GV100              |
|-------------------------------------|-------------------|---------------|--------------|--------------------------|
| Compute Capability                  | 3.5               | 5.2           | 6.0          | 7.0                      |
| Threads / Warp                      | 32                | 32            | 32           | 32                       |
| Max Warps / SM                      | 64                | 64            | 64           | 64                       |
| Max Threads / SM                    | 2048              | 2048          | 2048         | 2048                     |
| Max Thread Blocks / SM              | 16                | 32            | 32           | 32                       |
| Max 32-bit Registers / SM           | 65536             | 65536         | 65536        | 65536                    |
| Max Registers / Block               | 65536             | 32768         | 65536        | 65536                    |
| Max Registers / Thread              | 255               | 255           | 255          | 255 <sup>1</sup>         |
| Max Thread Block Size               | 1024              | 1024          | 1024         | 1024                     |
| FP32 Cores / SM                     | 192               | 128           | 64           | 64                       |
| Ratio of SM Registers to FP32 Cores | 341               | 512           | 1024         | 1024                     |
| Shared Memory Size / SM             | 16 KB/32 KB/48 KB | 96 KB         | 64 KB        | Configurable up to 96 KB |

<sup>1</sup> The per-thread program counter (PC) that forms part of the improved SIMT model typically requires two of the register slots per thread.

Displaying: Images/image\_p5\_1.jpeg

```
graph LR; CPU[Host CPU] --- Bridge[Bridge]; Bridge --- SM[System memory]; Bridge --- HI[Host interface]; HI --- IA[Input assembler]; HI --- VCS[Viewport/clip/setup/raster/zcull]; VCS --- GPU[GPU]
```

The notebook also includes code to display the actual images that were part of the retrieved results.

In summary, the project evolves from identifying image presence (Phase 1) to extracting and utilizing the textual content within images via OCR and relevance filtering (Phase 2) to provide more comprehensive, context-aware answers from PDF documents. Both phases leverage Langchain, Pinecone, HuggingFace embeddings, and Google's Gemini models.

## **Important Research**

During my work on this research project I have found a very useful tool called Llaparse. It is a GenAI-native document parsing platform developed by LlamaIndex, specifically designed to transform complex, unstructured documents into formats optimized for large language model (LLM) applications such as Retrieval-Augmented Generation (RAG) and AI agents.

There are multiple features inherent in it such as using custom language OCR tools and setting document layout along with context so that the parsing can be very much clear.

This assignment taught me the possible background of such a tool and how it can be built.