# Contexts & Machines:
# How Document Parsing Shapes
# RAG results

Alessio Vertemati

Andrea Ponti

ON OFF

oneofftech.xyz

# Retrieval Augmented Generation (RAG)

RAG frameworks enhance large language models by
providing external knowledge to ground their responses.
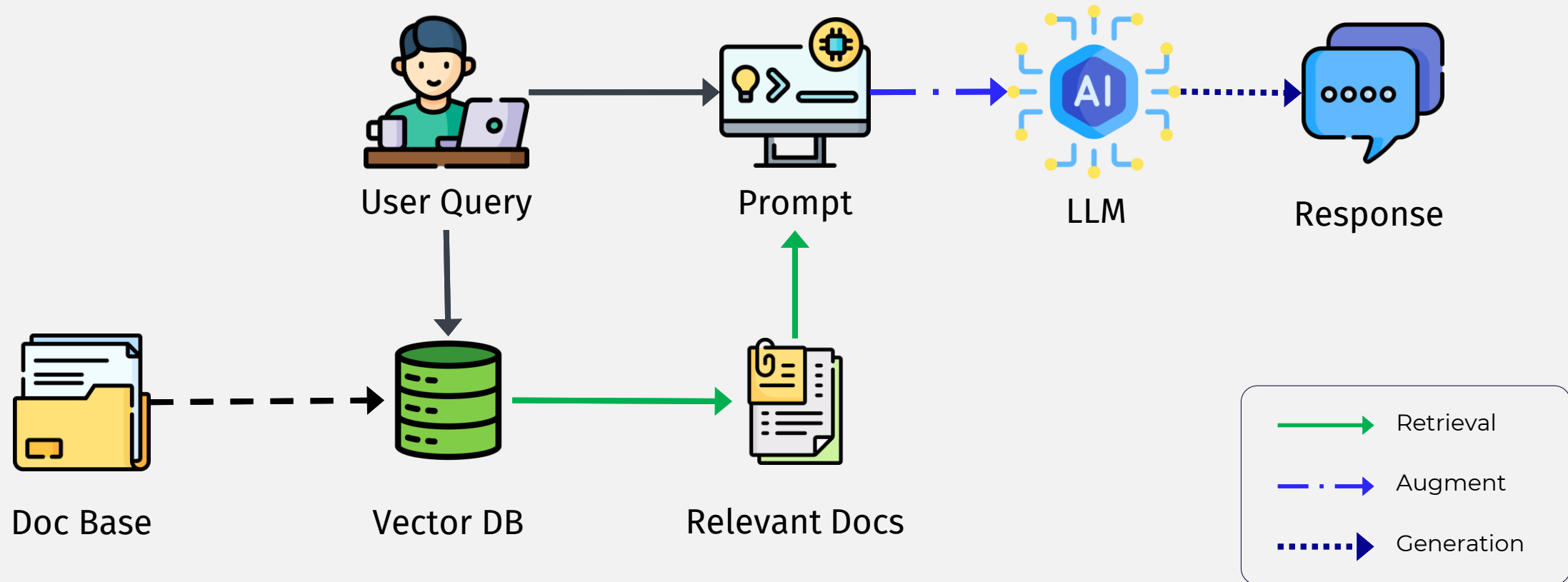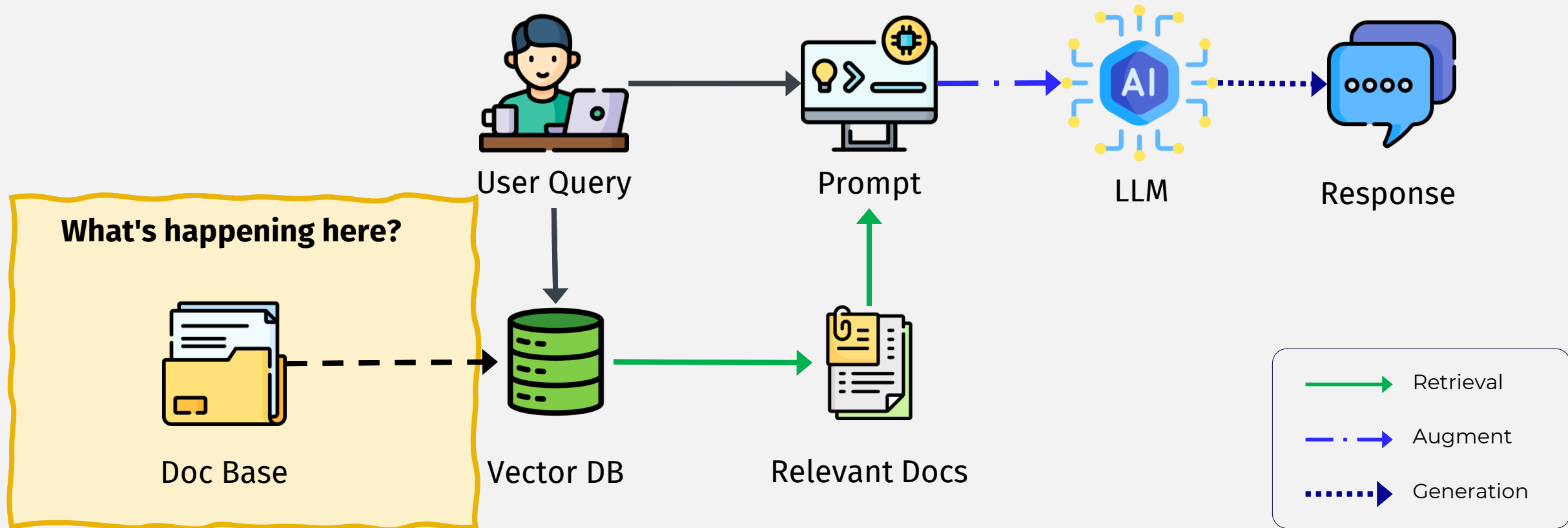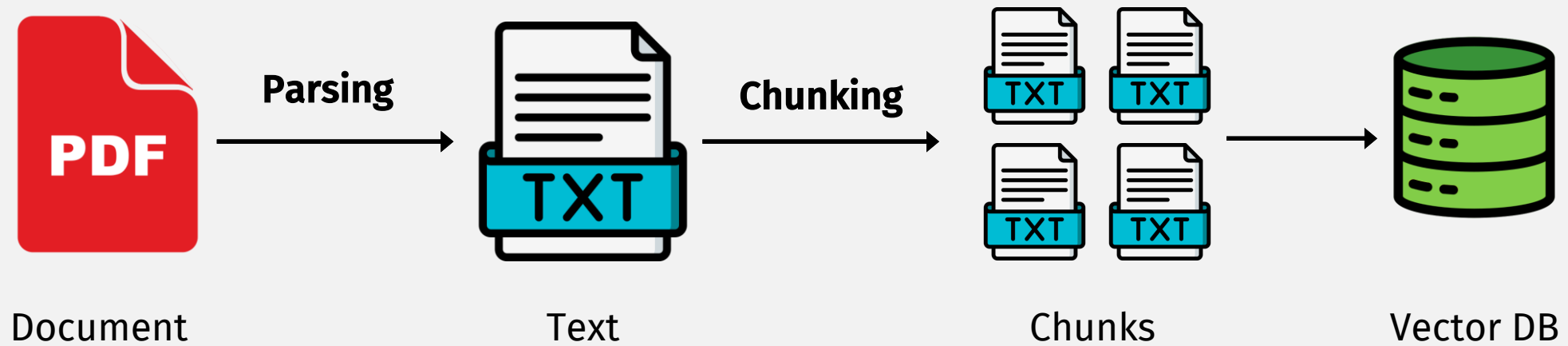
# Retrieval Augmented Generation (RAG)

RAG frameworks enhance large language models by
providing external knowledge to ground their responses.



User Query

Prompt

LLM

Response

**What's happening here?**

Doc Base

Vector DB

Relevant Docs

Retrieval

Augment

Generation

# Document Processing for RAG

Parsing and chunking define what the retriever sees — and if it sees
the wrong thing, the generator fails.



Document      **Parsing** →      Text      **Chunking** →      Chunks      →      Vector DB

# Why is it so hard?

**Common Challenges:** Reading Order, Table Parsing, Headers, Footers, ...



Multi-column documents



Complex Tables



Complex Layout

# Methodology

Legend:
1 - - -> Dataset Building
2 ——> Evaluation

1 Documents → MinerU → ragas → Q&A Dataset

2 Parser (Py, Pdf Act, Unstructured) → Chunker (LangChain, chonkie) → RAG

Evaluation using factual correctness

# Document Processing in RAG

How the document processing strategy impact the performance of RAG

# Our setup

**4 parsers**
- PyMuPDF
- PdfAct
- LlamaParse
- Unstructured

**2 chunkers**
- LangChain – *Rule-Based*
- Chonkie - *Semantic*

**5 documents**
~100 pages each
Projects' report

**255 questions**
single-hop specific
multi-hop specific
multi-hop abstract

**Evaluation**
RAGAS Factual
Correctness

# Chunk Size Effect

*Chunk Size Distribution by parser*



- **Rule-Based chunking** generally creates bigger chunks.

- **Semantic chunking** have higher variance of chunk size.

- ***PdfAct*** (open-source) has a behaviour similar to ***unstructured*** (cloud hosted)

- ***PyMuPDF*** (open-source) has a behaviour similar to ***LlamaParse*** (closed-source)

# Results

Multi-Hop Abstract:
Chunk size invariance

Multi-Hop Specific:
Larger chunks

Single-Hop Specific:
Shorter chunks

Complexity

**Rule-based chunking works great!**



*Factual Correctness by Question Type*

# The Impact of Document Type

How to choose the right document parser?

# Our setup

**parxy**

https://github.com/OneOffTech/parxy

## 4 parsers
- PyMuPDF
- PdfAct
- LlamaParse
- Unstructured

## DocLayNet Dataset
80k pages
6 categories
11 classes

https://huggingface.co/datasets/ds4sd/DocLayNet

Datasets: 😵 ds4sd / **DocLayNet** 🗗   ♡ like 103   Follow 😵 Docling 701

Tasks: ⊡ Object Detection   ⊠ Image Segmentation   Sub-tasks: instance-segmentation   Size: 10K<n<100K   Tags: layout-se

License: 🏛 other

📦 **Dataset card**   ⊞ Data Studio   ›≣ Files and versions   👏 Community 4

👁 **Dataset Preview** ⓘ   </> API   Embed   ⊞ Data Studio

Split (3)
train

# Parser vs Document category

Accuracy calculated using Text Similarity Ratio between the extracted text and the ground truth

# How To Choose the Right Parser?

# Key take-aways

Everything depend on the **document** and **question types**

• Document type (and structure) influence the parser choice

• Question type influence chunking strategy

If your RAG system isn't performing well **look first at what you're retrieving**—and how that content is processed!

# Some links

- https://github.com/data-house/pdfact
- https://unstructured.io/
- https://github.com/oneofftech/awesome-pdf
- https://github.com/opendatalab/MinerU
- https://parxy.eu
- https://docs.cloud.llamaindex.ai/llamaparse/getting_started

ONEOFF

Check our blog oneofftech.xyz/blog

http://www.oneofftech.xyz/

# Icon Credits

- [Scissors icons created by Gulraiz - Flaticon](#)
- [Parsing icons created by Good Ware - Flaticon](#)
- [Ai technology icons created by FACH - Flaticon](#)
- [Embedded icons created by Freepik - Flaticon](#)
- [Database icons created by Creatype - Flaticon](#)
- [Command icons created by Freepik - Flaticon](#)
- [Message icons created by Freepik - Flaticon](#)
- [Document icons created by Freepik - Flaticon](#)
- [Question icons created by Flat-icons-com - Flaticon](#)
- [Screening icons created by Vectors Tank - Flaticon](#)
- [Computer icons created by Freepik - Flaticon](#)
- [Pdf icons created by egorpolyakov - Flaticon](#)
- [Multimedia icons created by surang - Flaticon](#)