# Daniël Spee



https://www.linkedin.com/in/dspee/

https://github.com/Daantie

https://www.luminis.eu/expert/daniel-spee/

daniel.spee@luminis.eu

# Jettro Coenradie

Search and Gen AI Evangelist @ Luminis

https://www.linkedin.com/in/jettro

https://github.com/jettro

https://www.luminis.eu/expert/jettro-coenradie/

Jettro.coenradie@luminis.eu

# Use ChatGPT to answer questions



Who are talking about RAG at JFall this year

At J-Fall 2024, Jettro Coenradie and Daniël Spee will be conducting a hands-on lab titled "Build the best knowledge retriever for your Large Language Model." This session focuses on Retrieval Augmented Generation (RAG) systems, which enhance Large Language Models by integrating external data sources to produce more accurate and contextually relevant responses. The hands-on lab is scheduled for November 7, 2024, from 10:35 am to 12:30 pm. (Jfall)

Additionally, Jettro Coenradie and Daniël Spee are associated with the RAG4j/p project, which aims to teach the power of RAG by providing essential components for building RAG systems in Java and Python. (Rag4J)

🍁 Sources ˅
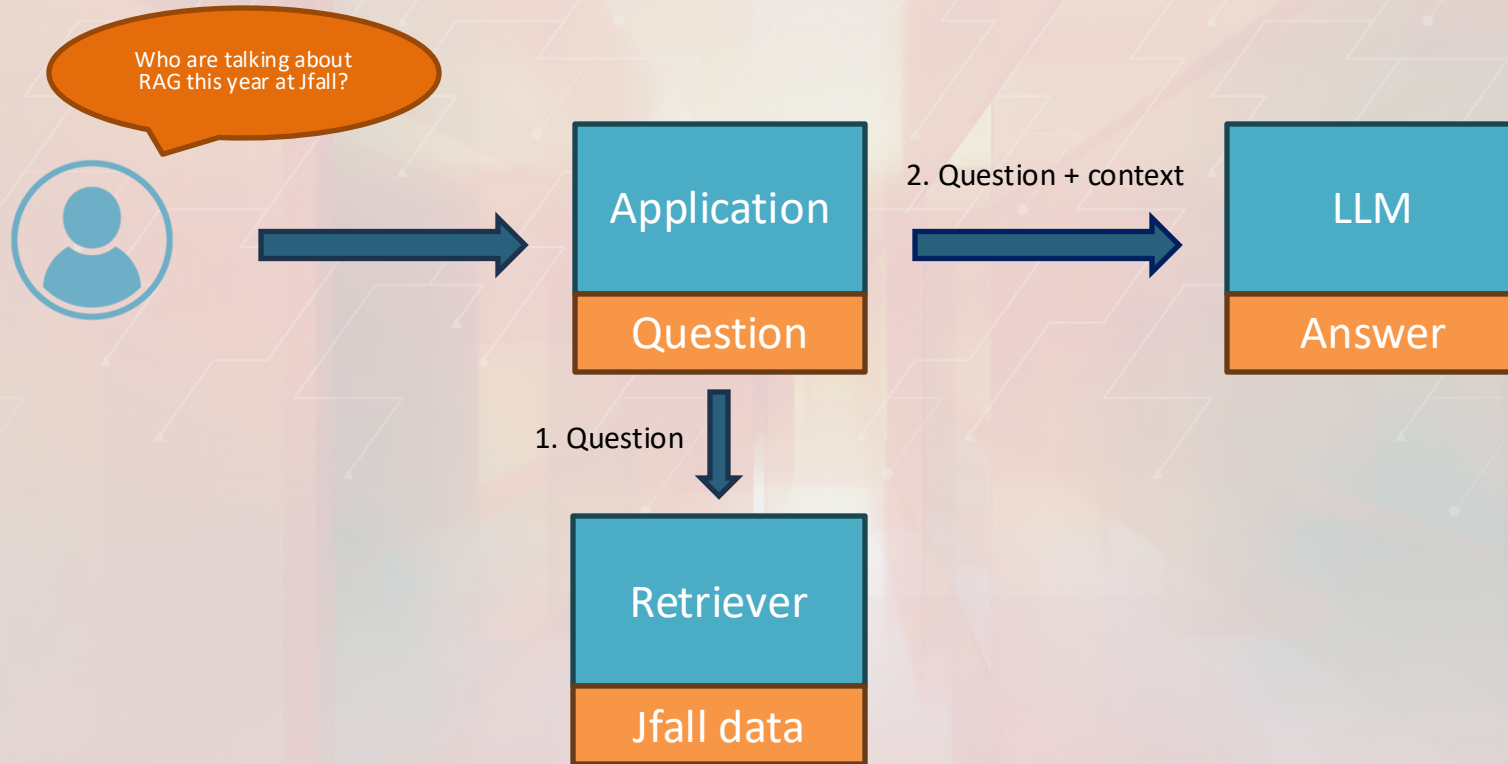
Message ChatGPT

# Use GPT-4o for the same question

# Check the sources for ChatGPT

# Retrieval Augmented Generation (RAG)

# Retrieval Augmented Generation (RAG)

# Retrieval Augmented Generation (RAG)

# Workshop agenda

- Setting up your environment
- Ingestion
- Retrieval strategy
- Retrieval quality
- Overall quality

Ingestion

Why is splitting necessary?

Impact of different splitters

Examples splitters


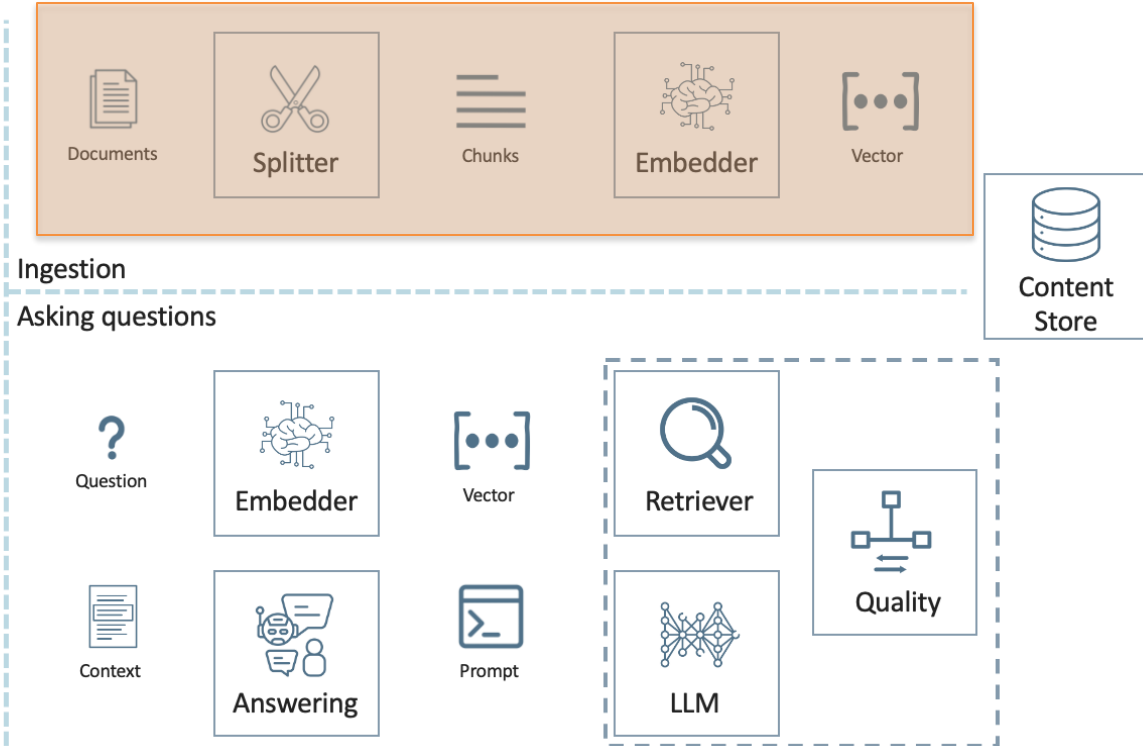DALL E - visualize a splitter that splits text into chunks

# Retrieval Augmented Generation (RAG)

https://rag4j.org
Workshop JFall 2024

https://github.com/RAG4J/rag4j-jfall
https://github.com/RAG4J/rag4p-jfall
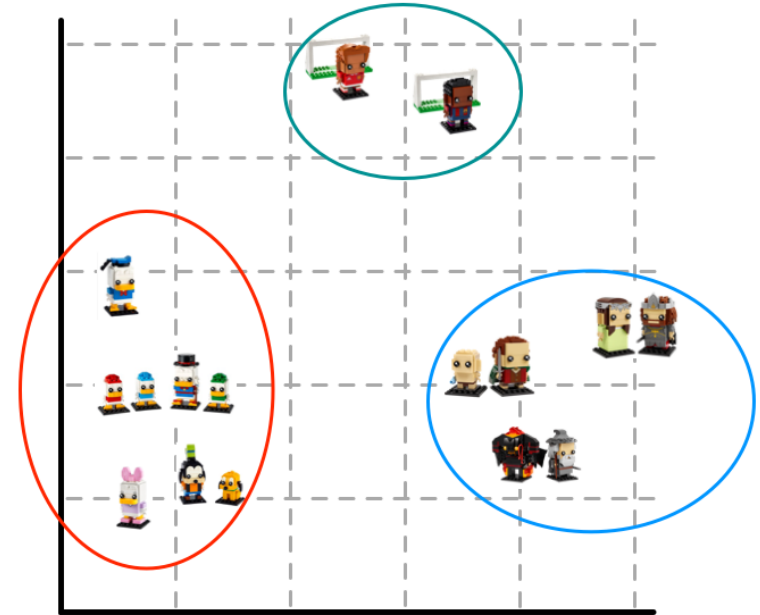
secret_key=
secret-jfallconf-2024-dont-share

```
J: AppStep1Ingestion
P: app_step1_ingestion
```
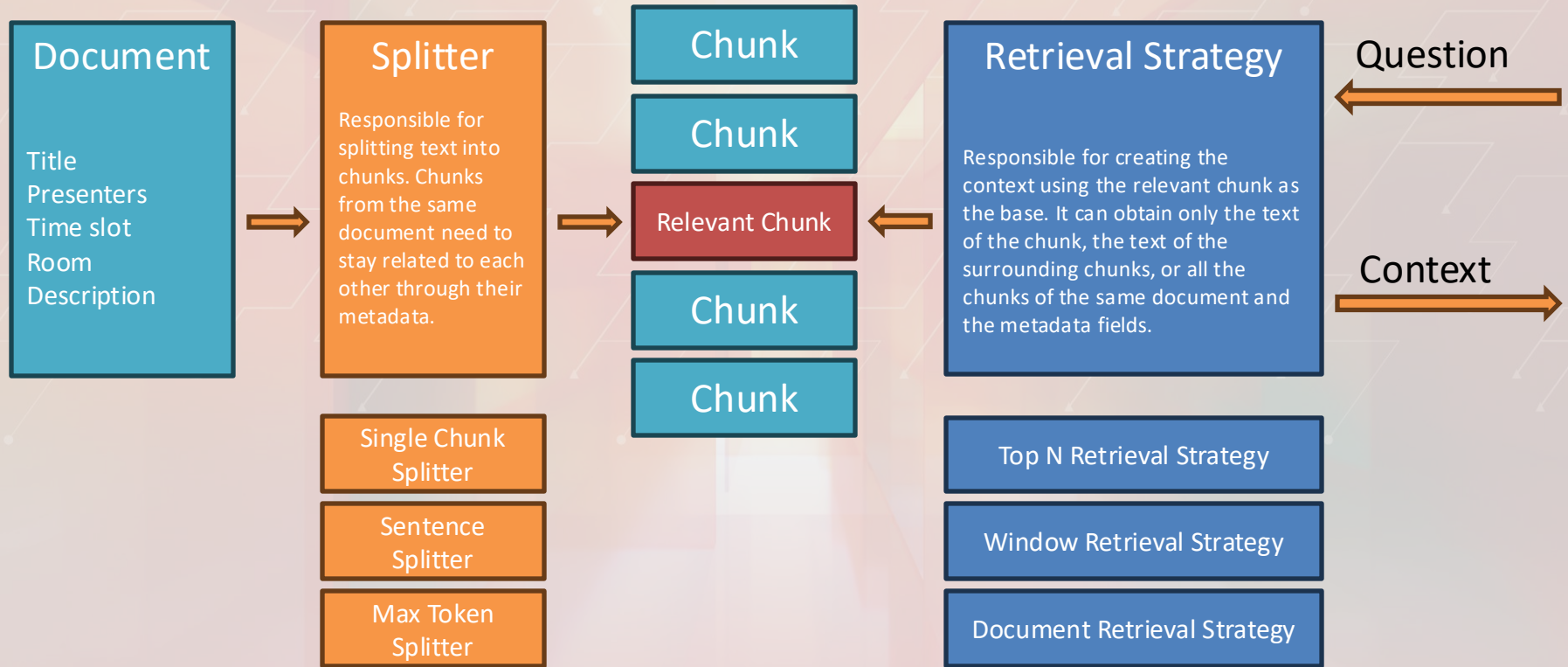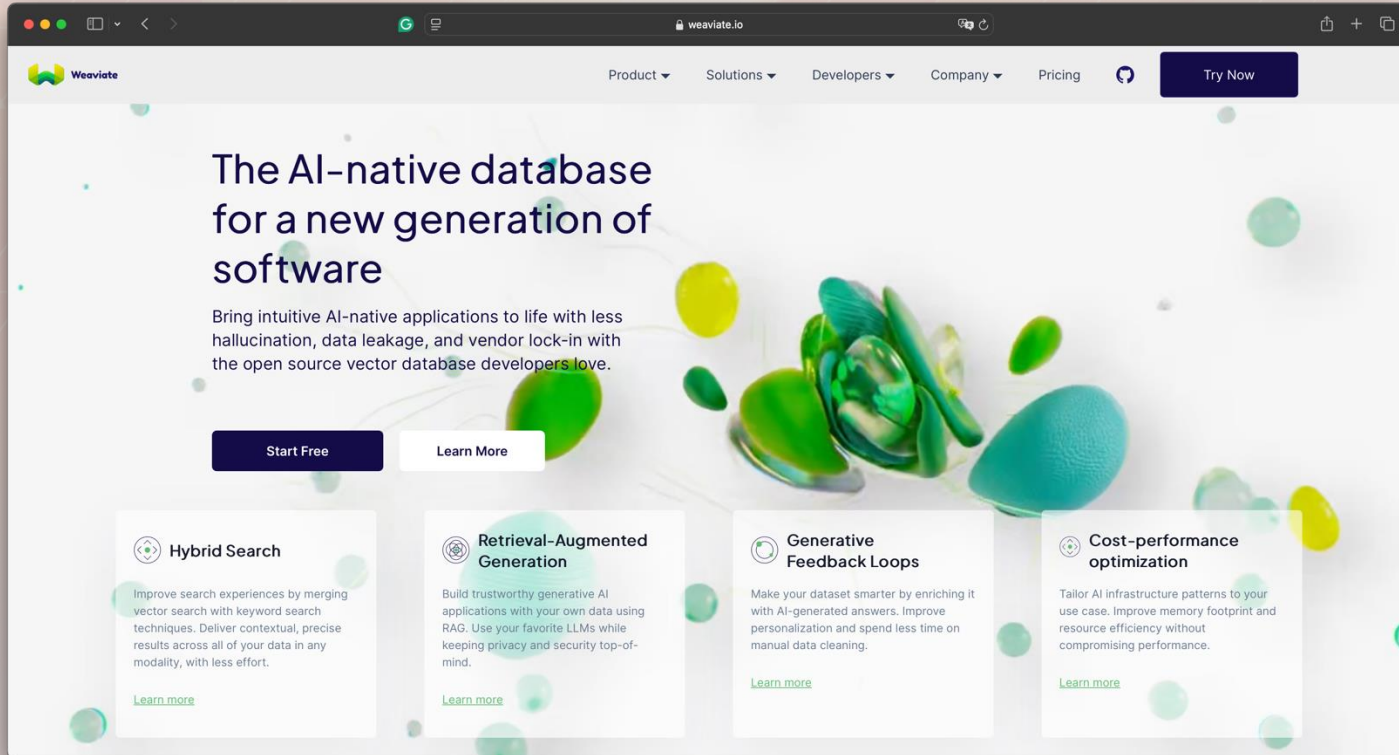
```
Branch: assignment or
assignment_ollama
```

Luminis. / part of yuma

# Retrieving

# Lexical versus Semantic Search

# Retrieving the context

**Document**

Title
Presenters
Time slot
Room
Description

**Splitter**

Responsible for splitting text into chunks. Chunks from the same document need to stay related to each other through their metadata.

Single Chunk Splitter

Sentence Splitter

Max Token Splitter

Chunk

Chunk

Relevant Chunk

Chunk

Chunk

**Retrieval Strategy**

Responsible for creating the context using the relevant chunk as the base. It can obtain only the text of the chunk, the text of the surrounding chunks, or all the chunks of the same document and the metadata fields.

Top N Retrieval Strategy

Window Retrieval Strategy

Document Retrieval Strategy

Question

Context

luminis. / part of yuma

# Weaviate ~ Vector store++

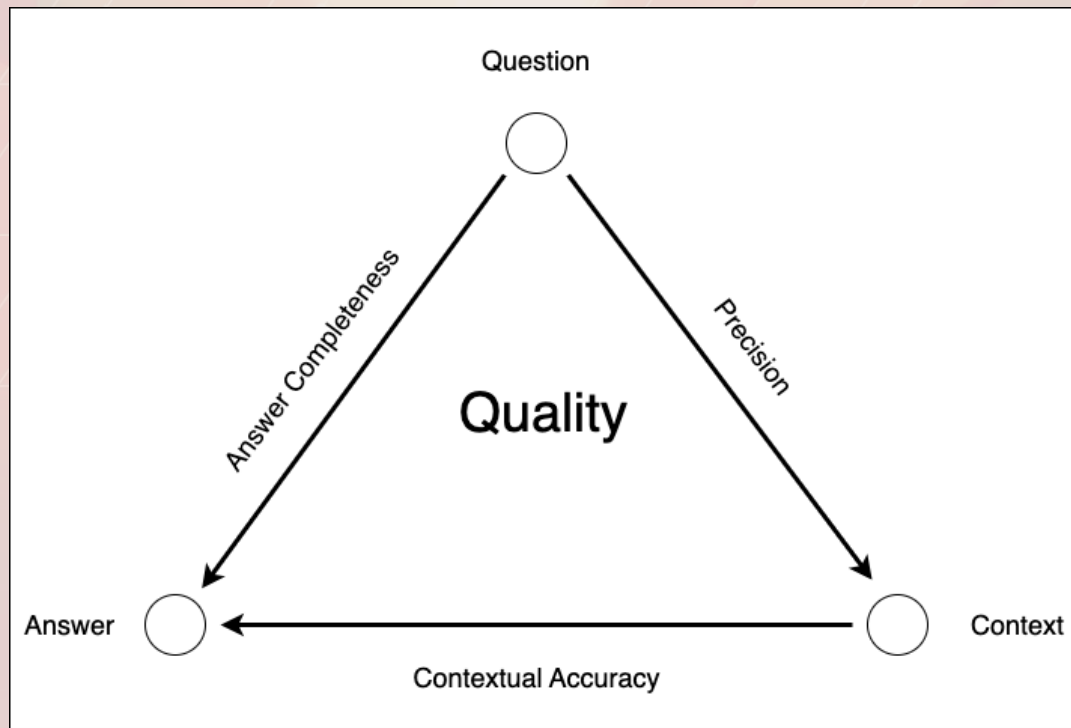# Retrieval Augmented Generation (RAG)

secret_key=
secret-jfallconf-2024-dont-share


```
J: AppStep2Retrieving
P: app_step2_retrieving
```
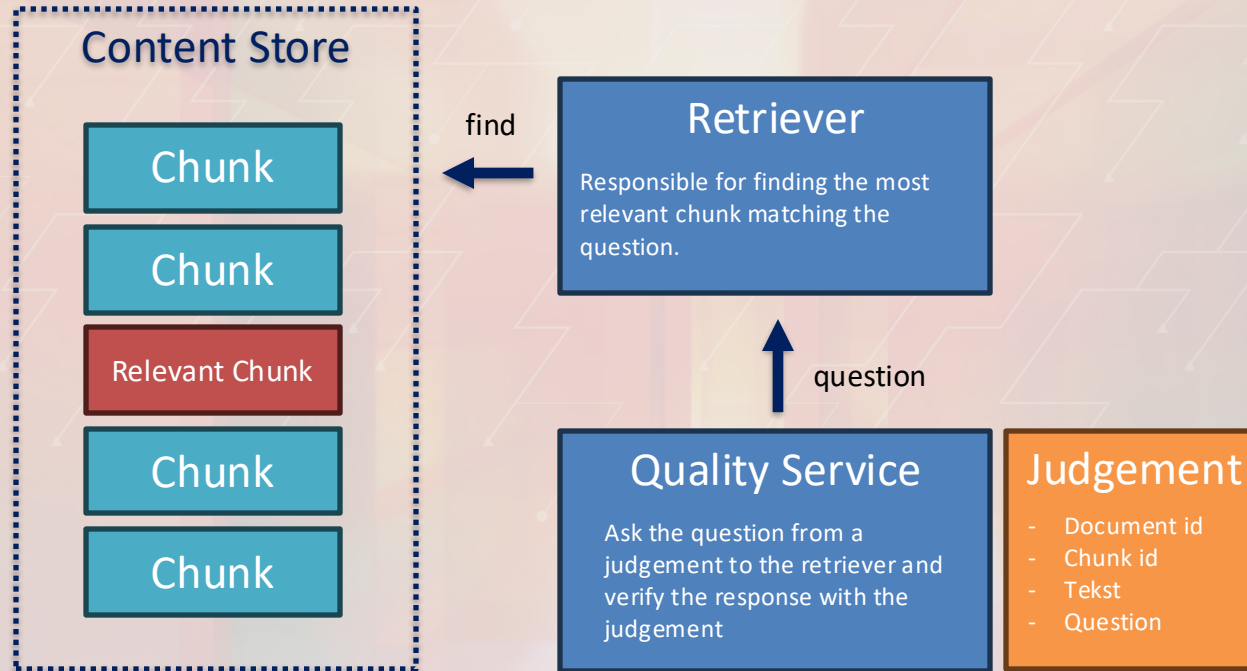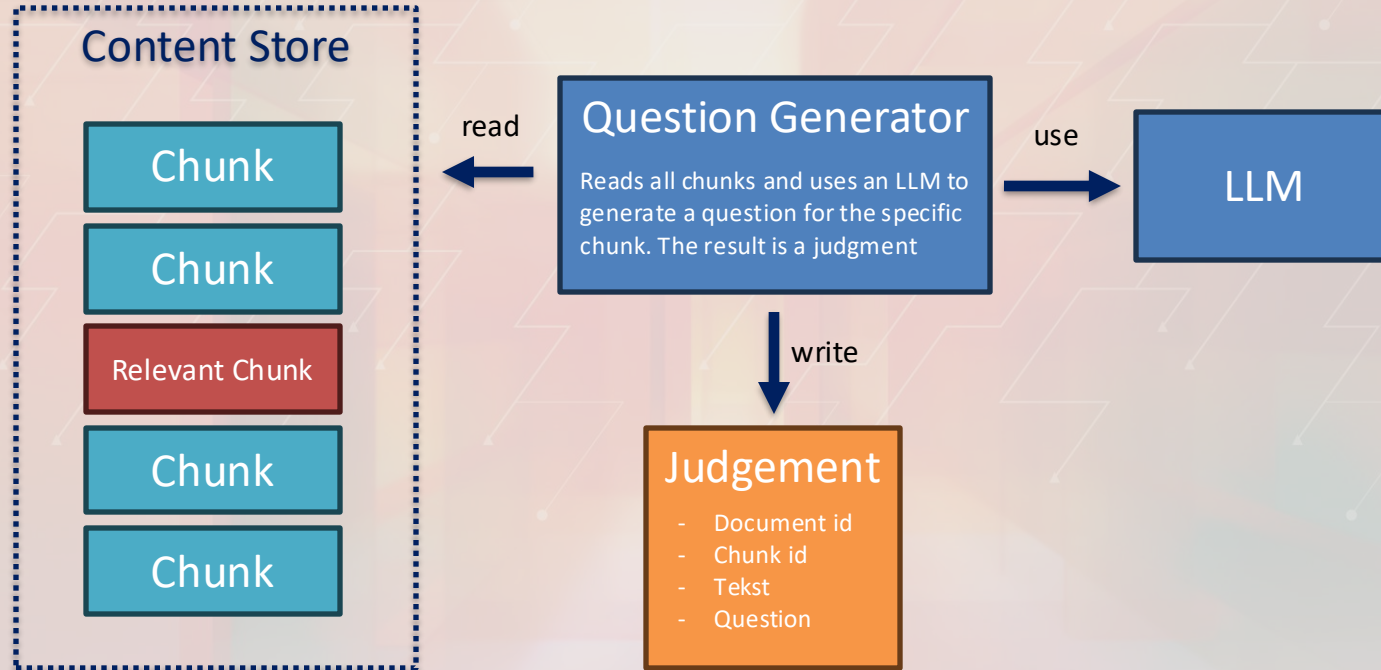
Retrieval Quality

# Quality of RAG

# Retrieval Quality - Precision

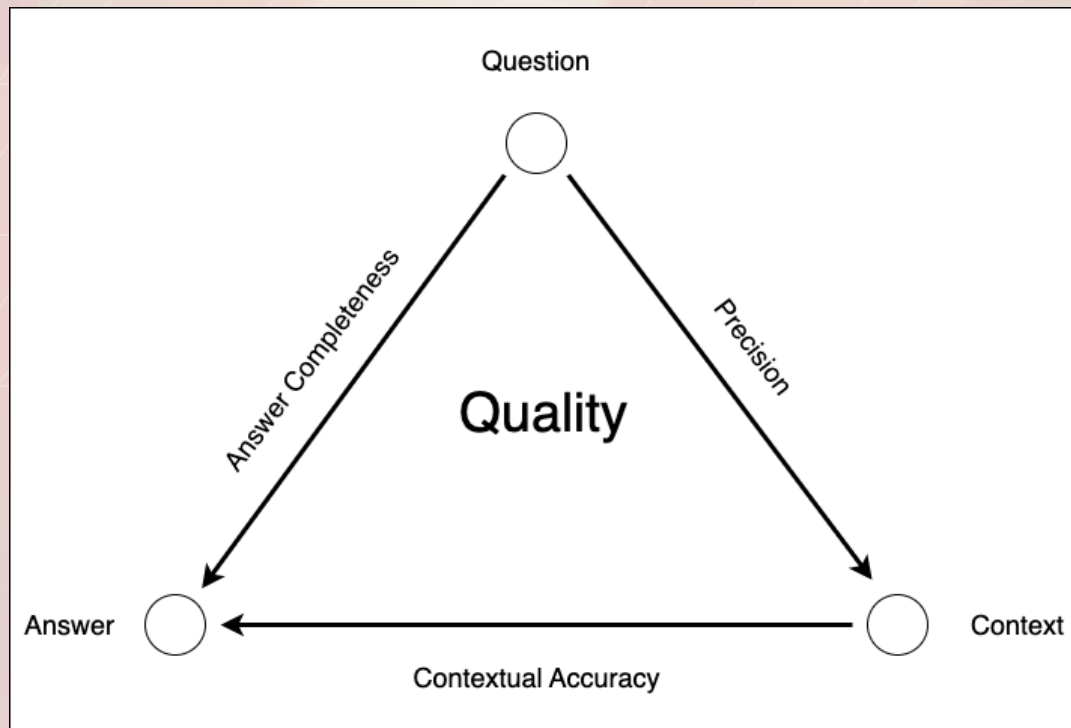# Generate the Judgement list

# Retrieval Augmented Generation (RAG)

secret_key=
secret-jfallconf-2024-dont-share


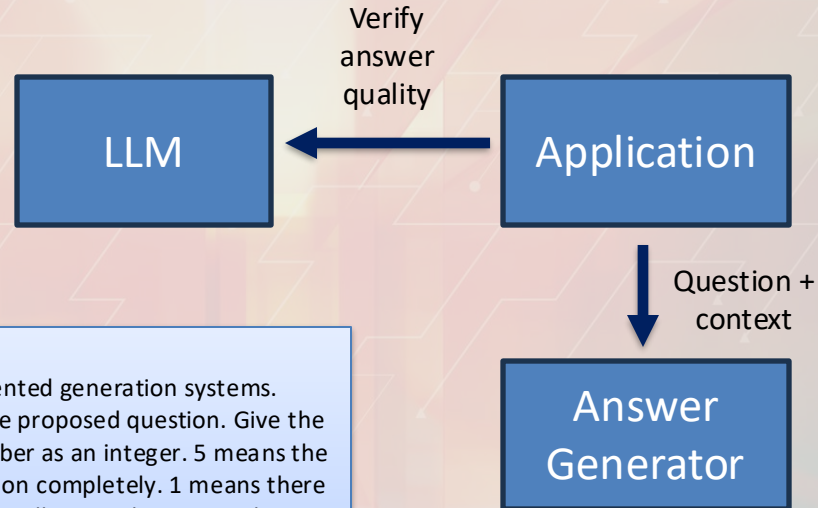J: AppStep3RetrievalQuality
P: app_step3_retrieval_quality

Answer Quality

# Prompt for quality of answer



Verify answer quality

LLM ← Application

Question + context

Answer Generator

You are a quality assistant verifying retrieval augmented generation systems. Your task is to verify a generated answer against the proposed question. Give the answer a score between 1 and 5 and keep the number as an integer. 5 means the answer contains the answer to the proposed question completely. 1 means there is no match between the answer and the question at all. Keep the reason short as in maximum 2 sentences. The question provided after 'question:'. The answer after 'answer:'. Write your answers in json format of:
{{"score": "score", "reason": "reason"}}
An example:
{{"score": 3,  "reason": "The answer is correct but some details are missing."}}

secret_key=
secret-jfallconf-2024-dont-share


```
J: AppStep4AnswerQuality
P: app_step4_answer_quality
```

Text