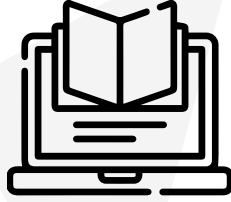


Version: 1.0



RAG _ against _ the

Will you answer my questions ?

Summary

Retrieval Augmented Generation, that's it.
That's the goal of this project .
Made in collaboration with @ldevelle, @pcamaren, @crfernand

#LLM

#AI

42

Intellectual Property Disclaimer

All content presented in this training module, including but not limited to texts, images, graphics, and other materials, is protected by intellectual property rights held by Association 42.

Terms of Use:

- **Personal use:** You are permitted to use the contents of this module solely for personal purpose. Any commercial use, reproduction, distribution, modification, or public display is strictly prohibited without prior written permission from Association 42.
- **Respect for Integrity:** You must not alter, transform, or adapt the content in any way that could harm its integrity.

Protection of Rights:

Any violation of these terms constitutes an infringement of intellectual property rights and may result in legal action. We reserve the right to take all necessary measures to protect our rights, including but not limited to claims for damages.

For any questions regarding the use of the content or to obtain authorization, please contact:
legal@42.fr

Contents

1	Foreword	1
1.1	Context	2
1.1.1	Overview	2
1.1.2	What is RAG ?	2
1.2	Instructions	3
1.2.1	Technical considerations	3
1.2.2	Data Models	4
1.3	Mandatory part	5
1.3.1	Performances	5
1.3.2	Smart Chunking Strategy	6
1.3.3	Retrieving Method	6
1.3.4	CLI Interface	7
1.3.5	Input	8
1.3.6	Output	8
1.3.7	Evaluation	10
1.3.8	Optional part	10

Chapter 1

Foreword

The **birthday paradox** is a classic problem in probability theory that demonstrates how unexpected events can occur more frequently than expected. It shows that even when the probability of an event is very low, it can still happen if there are enough opportunities.

In order to understand it better, let's make a guess :

If we take a classroom of 23 students, what is the probability that at least two of them have the same birthday ?

This problem is a veridical paradox, which regroups a great lists of problems that seem to be false but are in fact, true.

Enough suspense, the probability is : 50%! Who would have guessed ?

$$1 - \left(\frac{364}{365}\right)^{n(n-1)/2}$$

With n being the number of students. The probability is 50% when $n = 23$.

The most surprising is yet to come. If we take a classroom of 70 students, the probability is 99.9% that at least two of them have the same birthday.

"Hey good piece of trivia, but where am I going with this ?" - you may ask.

In cryptography, we now find a type of attack that takes advantage of the birthday paradox to find collisions in a hash function. It has been called : **the birthday attack**.

Now that you know that our instincts can be wrong and that maths leads you to bruteforce, let's move on to the project.

1.1 Context

1.1.1 Overview

A new project is often linked to new techniques and new skills: we've seen **function calling** in *call_me_maybe* and we will continue our exploration into the world of AI in this project. The main topic we're going to approach is **RAG**. But before seeing what **RAG** is about in its substance, let's focus on what it does! To do so, let's take some height.

When you want to work with a model in AI, the first thing to come in mind is to train it. You want the model to be able to use **language, reasoning, structure** and in order to do so you'll feed it a huge amount of data. Once you've done the training, the model "remembers" what you've fed it but it "knows" only the data that you've given to it. If you want to have fresher knowledge, you'll have to retrain it. And it takes a long time.

Training is a *technique*, and **RAG** is another one. Instead of feeding the model data, the RAG will give the model access to an external source of data, and that source is of *your choice*. The techniques can be combined: the model has to be trained on the key concepts such as we've seen before (language, reasoning, structure) to have the basis but for the knowledge, it can combine the trained data and the external connection.

1.1.2 What is RAG ?

Now that we know where we're at, you might ask yourself (if you haven't looked it up yet!) what is **RAG**? To understand it, we'll break it down into its three key concepts:

- **Retrieving**: Since the model is not trained on your specific data, it needs to search the database to *retrieve* the most useful snippets. First, the data has to be prepared. Then the model needs to understand your question. Once that's done, it matches the query with the database to choose the best results and finally pulls out the most relevant pieces of information. This involves **indexing, query encoding, similarity search, ranking**, and **retrieving**.
- **Augmenting**: Once the AI has retrieved the information, it can combine it with what it already "knows." This is called **augmenting**, since the AI is expanding its abilities by adding new information. Starting from the retrieved results, you can clean and filter them to remove irrelevant snippets (to avoid potential **noise**), insert them into the **context window**, and then combine them with the **trained knowledge** (the real augmentation step!).
- **Generating**: Now that you have retrieved the information and augmented it, the AI can finally generate an answer! Whether it's writing text, explaining a concept, or producing code snippets, this is the visible outcome of RAG. To do so, the AI reads the **context window**, understands the task at hand, blends the knowledge, and generates the output. Modern RAG systems often refine while writing, adjusting phrasing on the fly to maintain coherence and match the tone requested in the query.

Now that everything is clear, let's move forward !

1.2 Instructions

1.2.1 Technical considerations

All classes must use `pydantic` for validation and type safety.

- You must use **Python 3.10** for this project.
- You can use the libraries you want. We recommend `dspy`, `fire`, `tqdm`, `langchain`, `bm25s`, `chromadb`, `chonkie` packages.
- The use of `dspy` (or any similar package) is completely forbidden, including `pytorch`, `huggingface`, `transformers`, etc.
- You can use the following models:
 - `ollama_chat/qwen3:0.6b` (default)
 - Feel free to use other models (using the names from the HuggingFace hub) during the beta and let us know!
- We will use `uv` as a project and package manager.
- Your program must be run using the following command (where `src` is the main module):
`uv run python -m src`
- All errors should be handled gracefully. It must never crash unexpectedly, and must always provide a clear error message to the user.
- Your system must provide a CLI interface using Python Fire for easy command-line interaction.
- Progress bars should be implemented for long-running operations using `tqdm`.

1.2.2 Data Models

Your system must implement the following Pydantic models for type-safe data handling. These models ensure data integrity and provide automatic validation throughout the pipeline.

The `MinimalSource` model represents a minimal source of information:

```
class MinimalSource(BaseModel):
    file_path: str
    first_character_index: int
    last_character_index: int
```

The `UnansweredQuestion` model represents an unanswered question:

```
class UnansweredQuestion(BaseModel):
    question_id: str = Field(default_factory=lambda: str(uuid.uuid4()))
    question: str
```

The `AnsweredQuestion` model represents an answered question:

```
class AnsweredQuestion(UnansweredQuestion):
    sources: List[MinimalSource]
    answer: str
```

The `RagDataset` model represents a dataset of RAG questions:

```
class RagDataset(BaseModel):
    rag_questions: List[AnsweredQuestion] | List[UnansweredQuestion]
```

The `MinimalSearchResults` model represents the search results:

```
class MinimalSearchResults(BaseModel):
    search_query: str
    results_final: List[SearchResult]
```

The `MinimalAnswer` model represents an answer:

```
class MinimalAnswer(MinimalSearchResults):
    answer: str
```

The `StudentSearchResults` model represents a search results:

```
class StudentSearchResults(BaseModel):
    search_results: List[MinimalSearchResults]
    k: int
```

The `StudentSearchResultsAndAnswer` model represents a search results and an answer:

```
class StudentSearchResultsAndAnswer(StudentSearchResults):
    search_results: List[MinimalAnswer]
```

Those models **are not exhaustive**. Feel free to add more models and more fields in the existing models (in the search results model for example) if you need to.

1.3 Mandatory part

As you understood from the context, in this project we're going to implement **Retrieval-Augmented Generation**.

The knowledge base you'll be working with is the vLLM project in its version v0.10.1 provided in the resources of the project.

Your system should demonstrate the ability to:

- Build an indexed knowledge base from the repository files (both docs and code).
- Retrieve and rank the most relevant pieces of information.
- Pass them to the LLM within context limitations.
- Generate structured JSON output as described in the output section.
- Implement intelligent chunking strategies for different file types.
- Provide a comprehensive CLI interface for all operations.
- Include evaluation metrics and performance analysis.



Start by measuring your error using the simplest possible approach. Advance to complex methods once your error measurement is improving.

1.3.1 Performances

Your system must respect some minimal performances that are listed as follow :

- Indexing time: 5 minutes maximum
- Retrieval time : 1 minute for 1 question
- Question answering time : 1.8 seconds for 1 question



If you do not meet the expected performances, talk about it on discord

1.3.2 Smart Chunking Strategy

Your system must implement different strategies :

Python Code Chunking:

- Use AST parsing to understand code structure
- Keep functions and classes together as logical units
- Split on scope boundaries (function/class definitions)
- Preserve code integrity while respecting size limits
- Minimal tokenization to preserve code structure

Documentation Chunking:

- Split by headers to maintain semantic structure
- Add overlap between chunks for context preservation
- Handle large sections with sentence-aware splitting
- Full tokenization with stemming and stopword removal



The maximum chunk size is 2000 characters and it has to be configurable.

1.3.3 Retrieving Method

Your system must implement at least one basic retrieving method:

- TF-IDF
- BM25

Performance Target: Your BM25 implementation should achieve at least **75% recall@5** on English questions when evaluated against the provided test set.

1.3.4 CLI Interface

Your system must provide a comprehensive command-line interface using Python Fire.

Required Commands:

- **ingest**: Index the repository

```
uv run python -m src index
```

- **search**: Search the indexed repository

```
uv run python -m src search "OpenAI compatible server" --k 10
```

- **search_dataset**: Search dataset and search results

```
uv run python -m src search_dataset \  
data/datasets/UnansweredQuestions/Dataset_2025-09-21_valid.json
```

- **evaluate**: Evaluate search results

```
uv run python -m src measure_recall_at_k_on_dataset \  
data/output/search_results/Dataset_2025-09-21.json \  
data/datasets/AnsweredQuestions/Dataset_2025-09-21_valid.json
```

- **generate**: Generate answers for dataset

```
uv run python -m src answer_dataset \  
data/output/search_results/Dataset_2025-09-21_valid.json
```

- **answer**: Answer single query with context

```
uv run python -m src answer "How to configure OpenAI server?" --k 10
```

Those commands are not exhaustive. Feel free to add more commands if you need to or to custom them with flags to suit your needs.

1.3.5 Input

Ingestion Options:

- **Full Repository:** Index all files in the vLLM repository
- **Selective Ingestion:** Only process files mentioned in questions.tsv (recommended for testing)

For each query, your system must retrieve relevant chunks of the repository and generate an evidence-based response in the same form as the output.



Linked to the different chunking strategies, you can create different indexes for the different types of files.

1.3.6 Output

Your system must output a comprehensive JSON file containing detailed results and metadata:

- The initial query
- The retrieved chunks of the document which contains:
 - Filename with full path
 - Starting character of the retrieved information
 - Ending character of the retrieved information
 - Corresponding text content
 - File type (python, markdown, etc.)
 - Relevance score from the retrieval algorithm you are using
 - Section information (for documentation files, not needed for code)
- Retrieving strategy being used (TF-IDF, BM25, embedding)
- Comprehensive metrics containing:
 - Recall@5
 - Search time (in milliseconds)
 - Number of chunks processed

Output Format

Your system must provide detailed output as followed:

```
{
  "query": "How do I configure OpenAI compatible server?",
  "retrieved_chunks": [
    {
      "filename": "docs/serving/openai_compatible_server.md",
      "starting_character": 9867,
      "ending_character": 10100,
      "text": "# OpenAI-Compatible Server\n\nvLLM provides an HTTP server...",
      "file_type": "markdown",
      "section": "OpenAI-Compatible Server (level 1)",
      "score": 5.124
    },
    {
      "filename": "vllm/entrypoints/openai/api_server.py",
      "starting_character": 267,
      "ending_character": 400,
      "text": "class OpenAIAPIServer:\n    def __init__(self, args)...",
      "file_type": "python",
      "score": 4.892
    }
  ],
  "strategy": "bm25",
  "metrics": {
    "recall@5": 0.333,
    "search_time_ms": 45,
    "total_chunks": 1504,
  }
}
```

1.3.7 Evaluation

The evaluation of the RAG system is performed using a **recall@k** metric that measures the effectiveness of the retrieval component.

Recall@k Calculation

The recall@k for a given question is calculated by checking how much the retrieved sources overlap with the correct sources.

A source is considered "found" if there is at least 5% overlap between the retrieved source and any correct source.

If there are multiple sources in the question, their retrieval score for that question is

$$\frac{\text{number_found}}{\text{total_sources}}$$

Dataset-Level Evaluation

For the recall on the dataset, it is the average of the recall over all questions.

1.3.8 Optional part

Your system can implement an advanced retrieving strategy:

- Embedding-based retrieval with semantic similarity
- Hybrid approaches combining lexical and semantic search
- Query expansion and refinement techniques

Remember, this is not just about building a search engine – you're creating a comprehensive RAG system that intelligently processes, indexes, and retrieves information to augment language model capabilities. The birthday paradox taught us that our instincts can be wrong; let your implementation prove that systematic approaches and solid engineering can build something truly effective!