

# Python and AI/ML for Weather, Climate and Environmental Applications



Let us enjoy 🚀  
playing 🤖 with  
Python 🐍 and AI/ML!  
 

## Five-Day Schedule Overview

Time	Day 1	Day 2	Day 3	Day 4	Day 5
09:00–10:00	Opening by ECMWF DG, Start: Coding & Science in the Age of AI	Neural Network Architectures	Diffusion and Graph Networks	MLOps Foundations	Model Emulation, AIFS and AICON
10:00–11:00	<b>Lab:</b> Python Startup: Basics	<b>Lab:</b> Feed-forward and Graph NNs	<b>Lab:</b> Graph Learning with PyTorch	<b>Lab:</b> Containers and Reproducibility	<b>Lab:</b> Emulation Case Studies
11:00–12:00	Python, Jupyter and APIs	Large Language Models	Agents and Coding with LLMs	CI/CD for Machine Learning	AI-based Data Assimilation
12:00–12:45	<b>Lab:</b> Work environments, Python everywhere	<b>Lab:</b> Simple Transformer and LLM Use	<b>Lab:</b> Agent Frameworks	<b>Lab:</b> CI/CD Pipelines	<b>Lab:</b> Graph-based Assimilation
12:45–13:30	<b>Lunch Break</b>				
13:30–14:30	Visualising Fields and Observations	Retrieval-Augmented Generation (RAG)	DAWID System and Feature Detection	Anemoi: AI-based Weather Modelling	AI and Physics
14:30–15:30	<b>Lab:</b> GRIB, NetCDF and Obs Visualisation	<b>Lab:</b> RAG Pipeline	<b>Lab:</b> DAWID Exploration	<b>Lab:</b> Anemoi Training Pipeline	<b>Lab:</b> Physics-informed Neural Networks
15:30–16:15	Introduction to AI and Machine Learning	Multimodal Large Language Models	MLflow: Managing Experiments	<b>The AI Transformation</b>	Learning from Observations Only
16:15–17:00	<b>Lab:</b> Torch Tensors and First Neural Net	<b>Lab:</b> Radar, SAT and Multimodal Data	<b>Lab:</b> MLflow Hands-on	<b>Lab:</b> How work style could change	<b>Lab:</b> ORIGEN and Open Discussion
17:00–20:00					Joint Dinner

## Why Retrieval-Augmented Generation (RAG)?

### Limitations of standalone LLMs

- ▶ Fixed knowledge at training time
- ▶ No access to private or local data
- ▶ Risk of hallucinated answers

Who is president of Deutscher Wetterdienst?

The President of the German Weather Service (Deutscher Wetterdienst, DWD) is actually a position that has not existed since 2005. Prior to then, the head of the organization was titled "President" or "Director General". However, in 2005, the role was redefined as a "Chairperson" or "Managing Director".

As of my knowledge cutoff, the current Chairperson (Präsident) of Deutscher Wetterdienst is Prof. Dr. Gerhard Adrian.

Please note that this information might be subject to change, and I'd be happy to update you if new developments occur!

LLMs answer from *learned parameters only*. **This can be completely off reality!**

### Typical real-world needs

- ▶ Large code bases
- ▶ Technical documentation
- ▶ Evolving project knowledge

### Key idea:

Bring the knowledge to the model.

## Core Idea of Retrieval-Augmented Generation

### Standard LLM workflow

Prompt → LLM → Answer

- ▶ No external knowledge
- ▶ No verification

Pure LLMs work well for general reasoning, language understanding, summarization, and creative text generation when the required knowledge is common and static. They perform poorly when accurate, up-to-date, proprietary, or highly technical domain knowledge is required, because they cannot verify facts or access external sources. In such cases, LLMs tend to hallucinate plausible but incorrect answers.

### RAG workflow

Query → Retrieve → LLM → Answer

- ▶ Search relevant documents
- ▶ Inject context into prompt
- ▶ Grounded answers

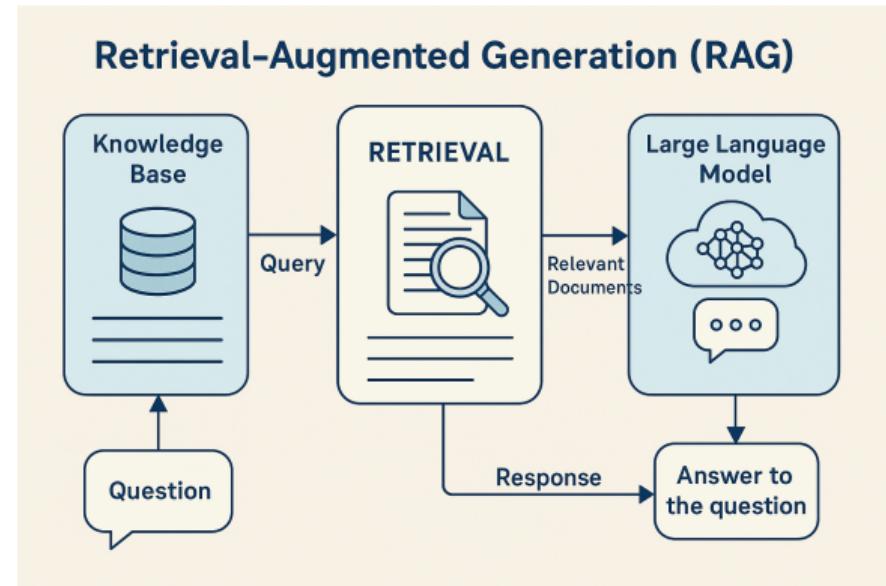
## RAG Architecture Overview

**RAG is more than just text search!**

### Main components

- ▶ Document collection
- ▶ Vector database
- ▶ Retriever
- ▶ Language model

Each component has a *clear role*.



Documents provide context for generation

## Vector Databases and Transformer Embeddings

### Text to vector mapping

Given a tokenized input sequence

$$(w_1, w_2, \dots, w_n)$$

a Transformer maps tokens to embeddings:

$$x_i = E(w_i) \in \mathbb{R}^d$$

After contextualization (self-attention):

$$h_i = \text{Transformer}(x_1, \dots, x_n)$$

A sentence or paragraph embedding is typically:

$$z = \frac{1}{n} \sum_{i=1}^n h_i$$

### Vector database view

Each paragraph is stored as:

$$z_j \in \mathbb{R}^d$$

Similarity search uses a distance measure, e.g.

$$\text{sim}(z_q, z_j) = \frac{z_q \cdot z_j}{\|z_q\| \|z_j\|}$$

or equivalently:

$$\|z_q - z_j\|_2$$

**Key point:** The same embedding space is reused — now *outside* the LLM.

## Elementary Vector Database — Setup

### Toy example

- ▶ Small set of short sentences
- ▶ Different semantic topics
- ▶ Some paraphrases, some unrelated

#### Example sentences

```
1 sentences = [  
2     "The sun is shining brightly today.",  
3     "Heavy rain is falling over the city.",  
4     "Neural networks can learn patterns.",  
5     "Transformers use self-attention."  
6 ]
```

Each sentence will become one vector in a shared space.

We use this to build intuition before scaling to real documents.

$$\text{Sentence} \rightarrow z \in \mathbb{R}^d$$

## Text to Vector Space

### Sentence embedding

A sentence  $s_i$  is mapped to a vector

$$z_i \in \mathbb{R}^d$$

using a pretrained Transformer:

$$z_i = f_\theta(s_i)$$

In our example:

$$d = 384$$

### Embedding matrix

For  $N$  sentences, we obtain:

$$Z = \begin{bmatrix} z_1^\top \\ z_2^\top \\ \vdots \\ z_N^\top \end{bmatrix} \in \mathbb{R}^{N \times d}$$

Each row corresponds to *one sentence*.

## Similarity Search

### Query embedding

A query sentence is mapped to:

$$z_q \in \mathbb{R}^d$$

**Similarity** to stored sentences is measured using cosine similarity:

$$\text{sim}(z_q, z_i) = \frac{z_q \cdot z_i}{\|z_q\| \|z_i\|}$$

See example `1_vector_db_elementary.ipynb` in the `code07/` directory carrying out transforms explicitly. The command `SentenceTransformer("all-MiniLM-L6-v2")` loads a pretrained Transformer-based sentence embedding model that maps each input sentence to a fixed-size vector  $z \in \mathbb{R}^{384}$  capturing its semantic meaning. The model is publicly available at [huggingface.co/sentence-transformers/all-MiniLM-L6-v2](https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2).

### Similarity search (conceptual)

```
1 scores = cosine_similarity(  
2     query_embedding,  
3     embeddings )  
4  
5 top_k = argsort(scores)[-k:]
```

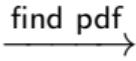
**Result:** Sentences with closest meaning are retrieved — no keywords needed.

## Document Search via Embeddings

### Goal

- ▶ Search *inside* a document
- ▶ Retrieve relevant pages or sections
- ▶ No keyword matching required

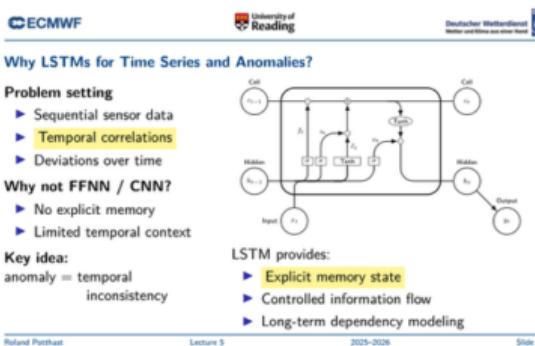
We use the same vector-space idea as in the elementary example.

"temporal pattern"  find pdf

### Input

- ▶ One lecture PDF (Lecture 05)
- ▶ Pages as semantic units
- ▶ Natural-language queries

**Output:** Relevant pages from the PDF



## From PDF Pages to Vectors

### Decomposition

The PDF is split into pages:

$$\{\text{page}_1, \dots, \text{page}_N\}$$

Each page is embedded as:

$$z_i = f_\theta(\text{page}_i) \in \mathbb{R}^d$$

In our example:

$$d = 384$$

### Embedding matrix

All pages form:

$$Z \in \mathbb{R}^{N \times d}$$

Each row corresponds to one PDF page.

**Key point:** Documents are represented as collections of vectors.

#### Extract text from PDF pages

```
1 for i, page in enumerate(reader.pages):
2     text = page.extract_text()
3     if text and len(text.strip()) > 50:
4         pages.append(text.strip())
5         page_ids.append(i + 1)
```

#### Embedding

```
1 page_embeddings = model.encode(
2     pages,
3     convert_to_numpy=True,
4     show_progress_bar=False)
```

## Semantic Search inside a PDF

### Query embedding

A user query is mapped to:

$$z_q \in \mathbb{R}^d$$

Similarity to page embeddings is computed via:

$$\text{sim}(z_q, z_i) = \frac{z_q \cdot z_i}{\|z_q\| \|z_i\|}$$

Top- $k$  most relevant pages are retrieved.

### Semantic page search

```
1 def search_pdf(query, model, pages, embeddings
2                 , page_ids, top_k=3):
3     q_emb = model.encode(
4         query, convert_to_numpy=True)
5     scores = []
6     for i, emb in enumerate(embeddings):
7         score = cosine_similarity(q_emb, emb)
8         scores.append((page_ids[i], score, pages
9                         [i]))
10    scores.sort(key=lambda x: x[1], reverse=
11                  True)
12    return scores[:top_k]
```

## From Retrieval to Evidence

### What we display

- ▶ Retrieved page number
- ▶ Extracted text snippet
- ▶ Original PDF page

This allows direct verification of the search result.

### Why this matters

- ▶ Transparent retrieval
- ▶ No hidden reasoning
- ▶ Trust through inspectable sources

#### Query: What is a graph network

Result 1: Page 9 (similarity = 0.623)

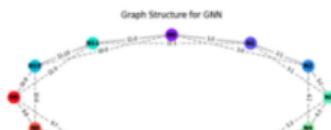
Extracted text (snippet):

Why Graph Neural Networks? When FFNNs are not enough ▶ Data points are not independent ▶ Relations between entities matter ▶ Ordering is not fixed or meaningful  
Typical examples ▶ Physical grids and meshes ▶ Sensor networks ▶ Molecules, social networks  
Graph structure with nodes and edges Core idea ▶ Nodes exchange information ▶ Learning respects graph structure ▶ Inductive bias for relational data  
Roland Potthast Lecture 5 2025–2026 Slide 8

#### Why Graph Neural Networks?

##### When FFNNs are not enough

- ▶ Data points are not independent
- ▶ Relations between entities matter
- ▶ Ordering is not fixed or meaningful



Graph structure with nodes and edges

Next: **Retrieval Augmented Generation:** Retrieve first — generate later.

## Why Do We Need FAISS?

### Above approach

- ▶ Loop over all embeddings
- ▶ Compute similarity one by one
- ▶ Exact but very (!) slow

Cost of one query:

$$\mathcal{O}(N \cdot d)$$

### Problem size

- ▶  $N$ : number of stored vectors (pages, paragraphs, documents)
- ▶  $d$ : embedding dimension (e.g.  $d = 384$ )

Each query compares against *all* vectors.

**Scaling issue:** Large  $N$  makes brute-force search infeasible.

### Typical RAG library sizes

- ▶ **Personal projects:**  $10^2$ – $10^3$  documents →  $N \sim 10^4$  chunks
- ▶ **Team / institutional data:**  $10^4$ – $10^5$  documents →  $N \sim 10^6$ – $10^7$  chunks
- ▶ **Enterprise-scale systems:**  $10^6$ + documents →  $N \gg 10^7$  chunks

Chunking multiplies the number of stored vectors  $N$ .

## FAISS — Vector Search at Scale

- ▶ Facebook AI Similarity Search
- ▶ Optimized nearest-neighbor search
- ▶ Designed for large vector collections

Stores vectors:

$$z_i \in \mathbb{R}^d$$

### Key idea

- ▶ Build an index once
- ▶ Query many times
- ▶ Fast top- $k$  retrieval

Search replaces explicit loops.

### Hierarchical search in FAISS

FAISS accelerates nearest-neighbor search by introducing a *coarse-to-fine hierarchy* in vector space.

#### Step 1: Coarse partitioning

The vector space is partitioned into  $M$  regions using representative centroids:

$$\{c_1, \dots, c_M\}, \quad c_j \in \mathbb{R}^d.$$

Each stored vector  $z_i$  is assigned to its nearest centroid:

$$z_i \mapsto c(z_i).$$

## FAISS in Practice

### Index construction

- ▶ Choose distance metric
- ▶ Add all embeddings
- ▶ Index lives in memory or on disk

Common choice: IndexFlatL2

### Step 2: Query routing

For a query vector  $z_q$ , FAISS finds the closest centroids:

$$\mathcal{C}_q = \arg \underset{L}{\text{top}} \|z_q - c_j\|_2, \quad L \ll M.$$

### Build and query FAISS index

```
1 import faiss
2 d = page_embeddings.shape[1]
3 index = faiss.IndexFlatL2(d)
4
5 index.add(page_embeddings)
6
7 distances, indices = index.search(
8     query_embedding, k)
```

### Step 3: Local search

Exact distances are computed only for vectors stored in the selected regions:

$$z_i \text{ with } c(z_i) \in \mathcal{C}_q.$$

## Streaming LLM Responses

### Standard LLM interaction

prompt → full response

User waits until generation is finished.

### Streaming interaction

prompt →  $(\delta_1, \delta_2, \dots)$

Partial tokens are delivered incrementally.

Recall that

**LLMs generate tokens incrementally!**

### What streaming changes

- ▶ Faster perceived response
- ▶ Progressive rendering
- ▶ Interactive user experience

### What streaming does not change

- ▶ Model architecture
- ▶ Reasoning capability
- ▶ Final content

## OpenAI Streaming

### Streaming API call

- ▶ Same prompt structure
- ▶ `stream=True`
- ▶ Tokens arrive as **deltas**

Each chunk contains new text:

$$\delta_k \subset \text{response}$$

There are several older OpenAI interface versions that language models may still suggest. Be careful: to avoid deprecated APIs, it is often necessary to explicitly provide a current code template.

### OpenAI streaming in Jupyter

```
1 stream = client.chat.completions.create(  
2     model="gpt-4o-mini", messages=[...],  
3     stream=True )  
4  
5 accumulated = ""  
6 handle = display(  
7     Markdown(""), display_id=True)  
8  
9 for chunk in stream:  
10    delta = chunk.choices[0].delta  
11    if delta.content:  
12        accumulated += delta.content  
13        handle.update(  
14            Markdown(accumulated))
```

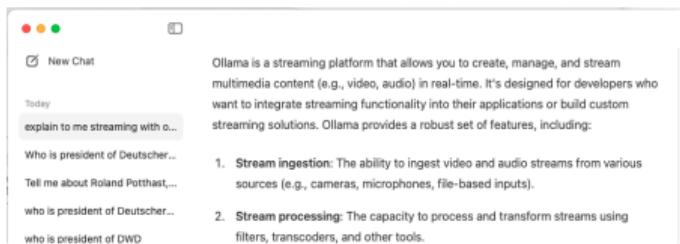
## Streaming with Local LLMs (Ollama)

### Local execution

- ▶ Models run on local hardware
- ▶ No external API calls
- ▶ Full data control

Streaming follows the same principle:

$$(\delta_1, \delta_2, \dots)$$



### Advantages

- ▶ Privacy and compliance
- ▶ Offline usage - train or plane
- ▶ No usage-based cost

### Trade-offs

- ▶ Smaller models
- ▶ Hardware dependent speed

**List of all available models:**  
<https://ollama.ai/library>

## Ollama Streaming in Python

### Ollama streaming call

- ▶ Local HTTP interface
- ▶ Same message format
- ▶ Streaming enabled

Works with: `llama3`,  
`mistral`, `mixtral`,  
`deepseek-r1`

#### Ollama streaming in Jupyter

```
1 stream = ollama.chat( model="llama3",
2                         messages=[...], stream=True )
3
4 accumulated = ""
5 handle = display(
6     Markdown(""), display_id=True)
7
8 for chunk in stream:
9     if "content" in chunk["message"]:
10         accumulated += chunk["message"]\
11             ["content"]
12         handle.update(Markdown(accumulated))
```

## RAG Implementation: Concrete Steps

### Implemented components

- ▶ text chunking (files → chunks) ✓
- ▶ sentence embeddings ✓
- ▶ FAISS index ✓
- ▶ metadata tracking ✓
- ▶ LLM query with context ✓

### Execution order

1. build vector database (offline)
2. embed user query
3. retrieve top- $k$  chunks
4. assemble prompt context
5. generate answer (optional streaming)

```
stream = client.chat.completions.create(  
    model="gpt-4o-mini",  
    messages=[  
        {  
            "role": "system",  
            "content": (  
                "You are a technical assistant. "  
                "Answer only using the provided context."  
            ),  
        },  
        {  
            "role": "user",  
            "content": (  
                "Context:\n"  
                f"{context}\n"  
                "Question:\n"  
                f"{query}"  
            ),  
        },  
    ],  
    [{"text": "Streamed response"}]  
)
```



## Top- $k$ Retrieval: What Do We Actually Get?

### FAISS output

For a query embedding  $z_q$ , FAISS returns:

$$\{(i_1, d_1), (i_2, d_2), \dots, (i_k, d_k)\}$$

- ▶  $i_j$ : index of a stored chunk
- ▶  $d_j$ : distance or similarity score

These indices refer to:

$$z_{i_j} \leftrightarrow \text{chunk}_{i_j}$$

### Important clarification

- ▶ FAISS returns *vectors*, not text
- ▶ Text is recovered via metadata lookup
- ▶ Ordering is by similarity score

### At this stage:

- ▶ no LLM involved
- ▶ no generation
- ▶ no reasoning

## From FAISS Indices to Text Chunks

### Stored metadata per chunk

Each vector  $z_i$  is linked to:

- ▶ file path
- ▶ line range or section
- ▶ raw text content

This mapping is stored outside FAISS  
(e.g. Python lists, JSON, databases).

### Lookup step

For each returned index  $i_j$ :

$$i_j \rightarrow (\text{file, lines, text})$$

This produces a ranked list:

$$(\text{chunk}_1, \dots, \text{chunk}_k)$$

**Result:** Concrete, inspectable evidence.

## Composing the Retrieval Context

### Context construction

Retrieved chunks are:

- ▶ ordered by similarity
- ▶ concatenated into one context block

Typical structure:

- ▶ file reference
- ▶ optional score
- ▶ text snippet

Context size is explicitly limited .

### Context assembly (ICON example)

```
1 context = ""
2 for r in retrieved_chunks:
3     context += (
4         f"[{r.file}], lines {r.start}-{r.end}]\n"
5         + r.text + "\n\n")
```

## OpenSearch

### What is OpenSearch?

OpenSearch is an open-source, distributed search and analytics engine.

It is designed to index and query large collections of documents with low latency.

Originally derived from Elasticsearch, OpenSearch is developed under the Apache 2.0 license.

### Core capabilities

- ▶ full-text search (inverted index)
- ▶ metadata filtering and aggregation
- ▶ persistent, disk-based indices
- ▶ distributed execution across nodes
- ▶ **vector similarity search**

OpenSearch combines text-based and vector-based retrieval in a single system.

## Google Search as Retrieval-Augmented Generation

### Core idea

Retrieval-Augmented Generation (RAG) combines an LLM with an external information source.

In classical RAG, retrieval uses:

- ▶ vector databases (FAISS)
- ▶ local document collections

### Google-based RAG

Here, retrieval is delegated to:

- ▶ Google Search
- ▶ live web documents

### Same logical pipeline

Query → Retrieve →  
Context → LLM → Answer

### Key difference

- ▶ Local RAG: curated, controlled, static
- ▶ Google RAG: open, dynamic, up-to-date

### Interpretation

The LLM does *not* search —  
it *summarizes retrieved evidence*.

## Search to LLM Pipeline

### Pipeline steps

1. Issue search query
2. Retrieve top- $k$  URLs
3. Scrape page text
4. Build context block
5. Ask LLM to answer

Answer =

LLM (Query, Web Evidence)

### Important

No embeddings are required —  
ranking by Search Engine .

### Google–LLM pipeline (schematic)

```
1 urls = search_google(query)
2 texts = [ scrape_website(url),
3   for url in urls ]
4
5 prompt = f """
6 Answer the question using
7 the sources below.
8 Question: {query}
9 Sources: {texts}
10 """
11
12 response = client.chat.completions.
13   create( model="gpt-4o-mini",
14     messages=[{"role":"user", "content
15       :prompt}] )
```