



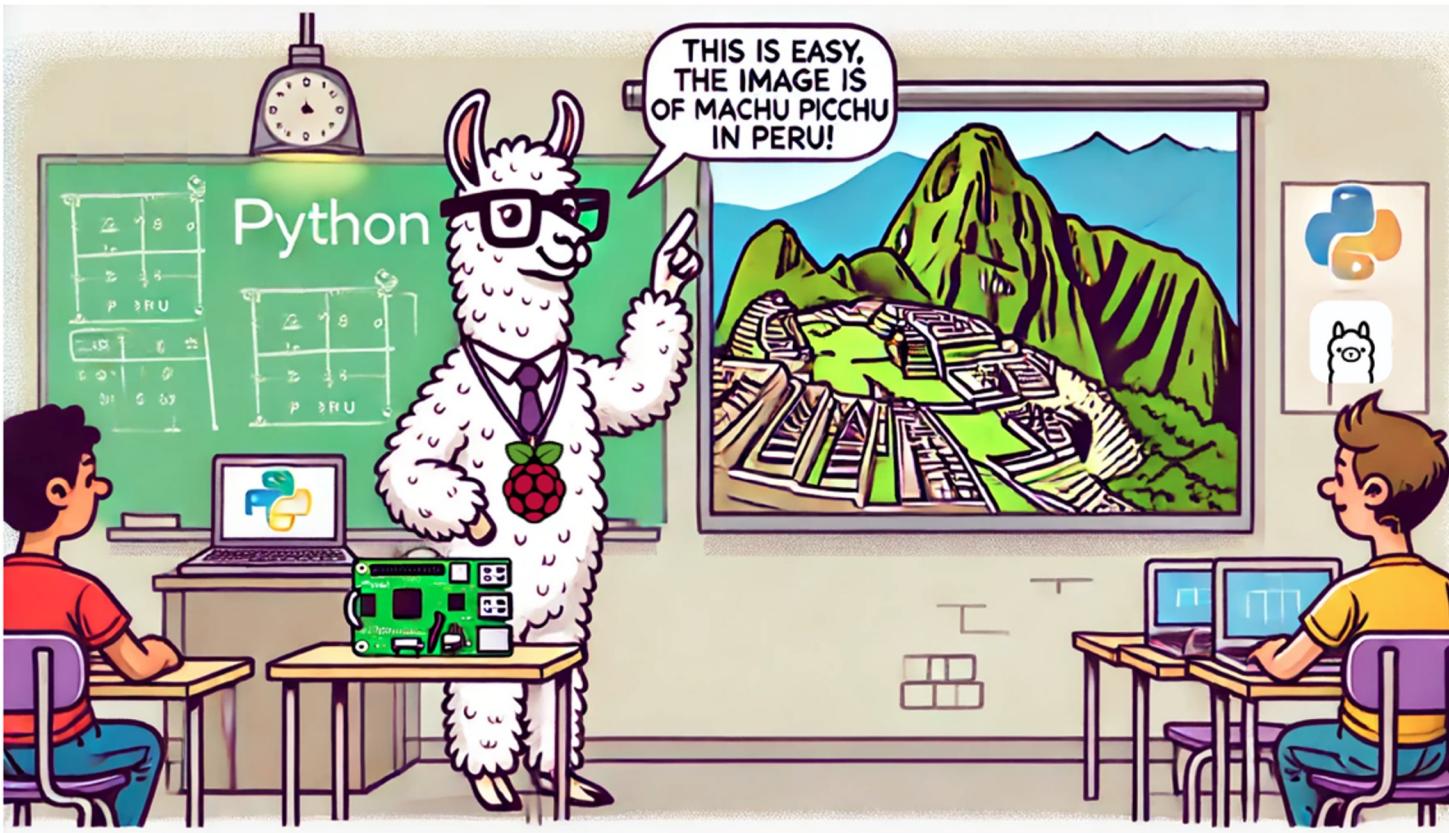
WALC 2024
Applied AI

Large Language Models (LLMs) at the Edge GenAI Demo

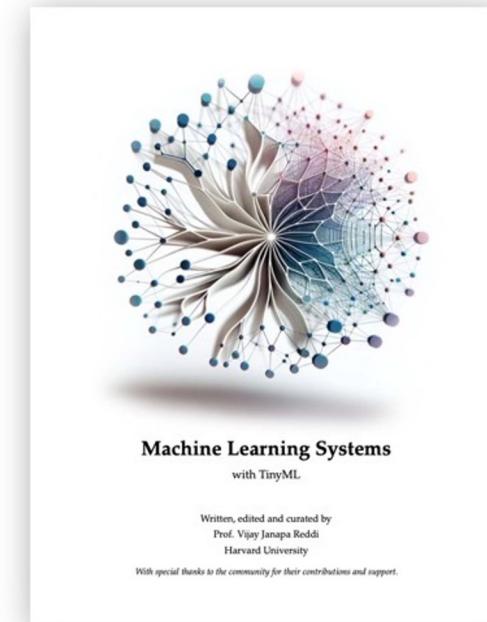
Prof. Marcelo J. Rovai
rovai@unifei.edu.br

UNIFEI - Federal University of Itajuba, Brazil
TinyML4D Academic Network Co-Chair





Running Large Language Models on Raspberry Pi at the Edge



The screenshot shows a web browser window with the title bar "Ollama". The address bar displays the URL "https://ollama.com". The main content area features a large, stylized black and white illustration of a llama's head and upper body. Below the illustration, the text "Get up and running with large language models." is displayed in a large, bold, black font. Underneath this, a smaller text block reads: "Run Llama 3.2, Phi 3, Mistral, Gemma 2, and other models. Customize and create your own." A prominent black button with the text "Download ↓" in white is centered below the descriptive text. At the bottom, a note states "Available for macOS, Linux, and Windows". The browser interface includes standard navigation buttons (back, forward, search, etc.) and a menu icon in the top right corner.

Get up and running with large language models.

Run Llama 3.2, Phi 3, Mistral, Gemma 2, and other models. Customize and create your own.

Download ↓

Available for macOS, Linux, and Windows



marcelo_rovai — mjrovai@raspi-5: ~ — ssh mjrovai@192.168.4.209 — 80x21

```
[mjrovai@raspi-5:~ $ python3 -m venv ~/ollama
[mjrovai@raspi-5:~ $ source ~/ollama/bin/activate
(ollama) mjrovai@raspi-5:~ $ curl -fsSL https://ollama.com/install.sh | sh
>>> Installing ollama to /usr/local
>>> Downloading Linux arm64 bundle
#####
##### 100.0%
#####
##### 100.0%
>>> Creating ollama user...
>>> Adding ollama user to render group...
>>> Adding ollama user to video group...
>>> Adding current user to ollama group...
>>> Creating ollama systemd service...
>>> Enabling and starting ollama service...
Created symlink /etc/systemd/system/default.target.wants/ollama.service → /etc/
systemd/system/ollama.service.
>>> The Ollama API is now available at 127.0.0.1:11434.
>>> Install complete. Run "ollama" from the command line.
WARNING: No NVIDIA/AMD GPU detected. Ollama will run in CPU-only mode.
(ollama) mjrovai@raspi-5:~ $ ollama -v
ollama version is 0.3.11
(ollama) mjrovai@raspi-5:~ $
```

```
marcelo_rovai — mjrovai@raspi-5: ~ — ssh mjrovai@192.168.4.209 — 79x26
(ollama) mjrovai@raspi-5:~ $ ollama run llama3.2:1b --verbose
pulling manifest
pulling 74701a8c35f6... 100% [██████████] 1.3 GB
pulling 966de95ca8a6... 100% [██████████] 1.4 KB
pulling fcc5a6bec9da... 100% [██████████] 7.7 KB
pulling a70ff7e570d9... 100% [██████████] 6.0 KB
pulling 4f659ale86d7... 100% [██████████] 485 B

verifying sha256 digest
writing manifest
success
>>> What is the capital of France?
The capital of France is Paris.

total duration:      2.620170326s
load duration:      39.947908ms
prompt eval count:   32 token(s)
prompt eval duration: 1.644773s
prompt eval rate:    19.46 tokens/s
eval count:          8 token(s)
eval duration:       889.941ms
eval rate:           8.99 tokens/s
```

Multimodal Models



```
marcelo_rovai — mjrovai@raspi-5: ~/Documents/OLLAMA — ssh mjrovai@192.168.4.209 — 84x36
(ollama) mjrovai@raspi-5:~/Documents/OLLAMA $ pwd
/home/mjrovai/Documents/OLLAMA
(ollama) mjrovai@raspi-5:~/Documents/OLLAMA $ ollama run llava-phi3:3.8b --verbose
>>> Describe the image /home/mjrovai/Documents/OLLAMA/image_test_1.jpg
Added image '/home/mjrovai/Documents/OLLAMA/image_test_1.jpg'
The image captures a breathtaking view of Paris, France. The cityscape is dotted with buildings in various shades of white and gray, interspersed with lush green trees that add a touch of nature to the urban setting.

In the heart of the scene stands the Eiffel Tower, an iconic symbol of Paris, its iron lattice structure reaching up into the clear blue sky. The tower's distinctive silhouette is unmistakable against the backdrop of the sky, which is a vibrant shade of blue with just a few clouds scattered across it.

The Seine River gracefully winds its way through the city, bordered by an array of buildings on both sides. The river is lined with several bridges that connect different parts of the city and facilitate movement for pedestrians and vehicles alike.

Above all these elements, a few birds can be seen soaring freely in the sky, their presence adding life to the scene. Their flight paths crisscross over the river and the buildings, creating dynamic patterns that draw the eye.

Overall, this image presents a beautiful daytime snapshot of Paris - its architectural marvels, natural beauty, and bustling city life coexisting in harmony.

total duration:      3m55.972199346s
load duration:      16.198011ms
prompt eval count:  1 token(s)
prompt eval duration: 2m19.561783s
prompt eval rate:   0.01 tokens/s
eval count:         276 token(s)
eval duration:      1m36.330959s
eval rate:          2.87 tokens/s
>>> Send a message (/? for help)
```

llava-phi-3 is a LLaVA model (Large Language and Vision Assistant) fine-tuned from Microsoft Phi-3 mini



= 147K tokens

~ 350 pages



~ 300 words/page



1 word = ~ 1.4 token

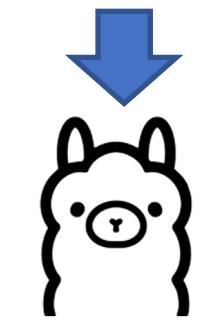


A **4-bit** quantized **3.8 billion parameter *** language model trained on **3.3 trillion tokens****, whose overall performance, as measured by both academic benchmarks and internal testing, rivals that of models such as Mixtral 8x7B and GPT-3.5

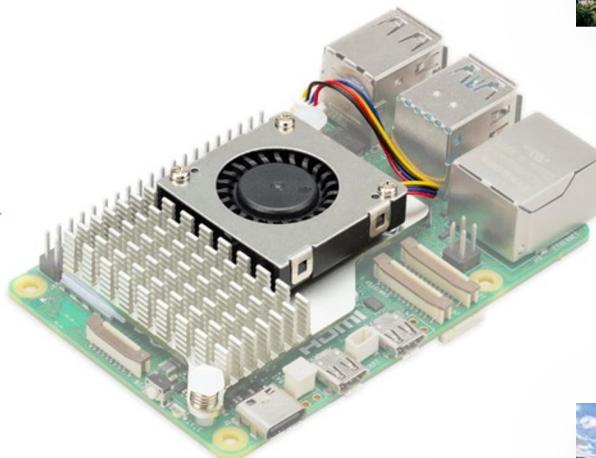
* 2.4 GB

** 22.5 Million books - 17% of all books written in the world

llava-phi-3 (2.9 GB)



Ollama



```
mjrovai@rpi-5: ~
File Edit Tabs Help
>>> Answer with one short sentence, what is the capital of France and its distance
... in Km from Santiago, Chile
The capital of France is Paris and it is around 12,674 kilometers away
from Santiago, Chile.

total duration:      13.860074968s
load duration:       1.537039ms
prompt eval count:   27 token(s)
prompt eval duration: 5.925386s
prompt eval rate:    4.56 tokens/s
eval count:          26 token(s)
eval duration:        7.539223s
eval rate:            3.45 tokens/s
>>> Send a message (/? for help)
```

(13 seconds)



```
mjrovai@rpi-5: ~/Documents/OLLAMA
Help
ute.

/Documents/OLLAMA $
/Documents/OLLAMA $ python calc_distance_image.py /
/home/mjrovai/Documents/OLLAMA/image_test_1.jpg

The image shows Paris, with lat:48.86 and long: 2.35, located in
France and about 11,630 kilometers away from Santiago, Chile.

[INFO] ==> The code (running llava-phi3), took 232.60845186299412
seconds to execute.

mjrovai@rpi-5:~/Documents/OLLAMA $
```



```
mjrovai@rpi-5: ~/Documents/OLLAMA
Help
ute.

/Documents/OLLAMA $
/Documents/OLLAMA $ python calc_distance_image.py /
/home/mjrovai/Documents/OLLAMA/image_test_3.jpg

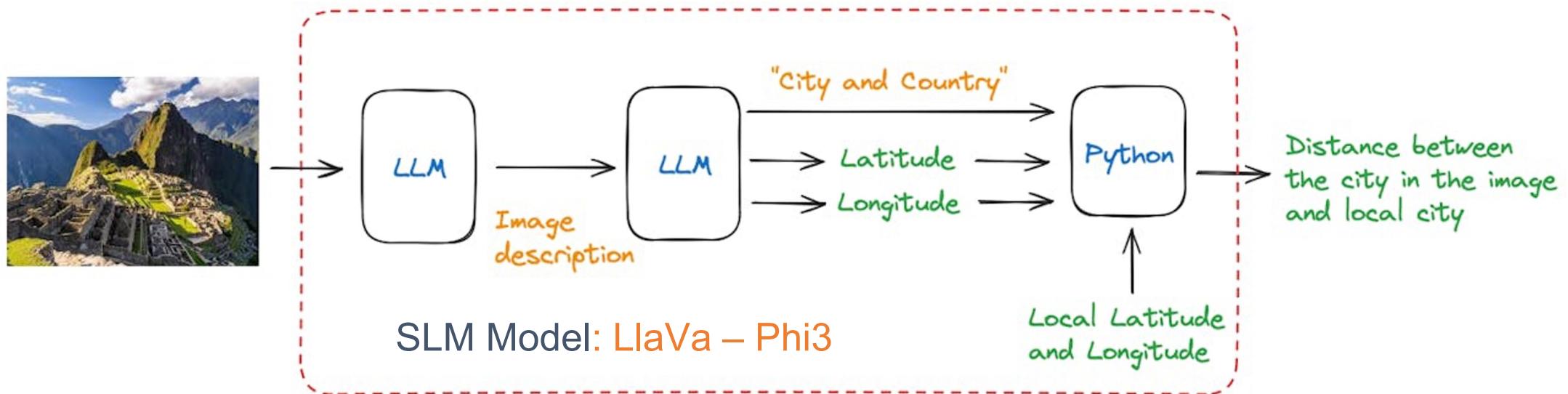
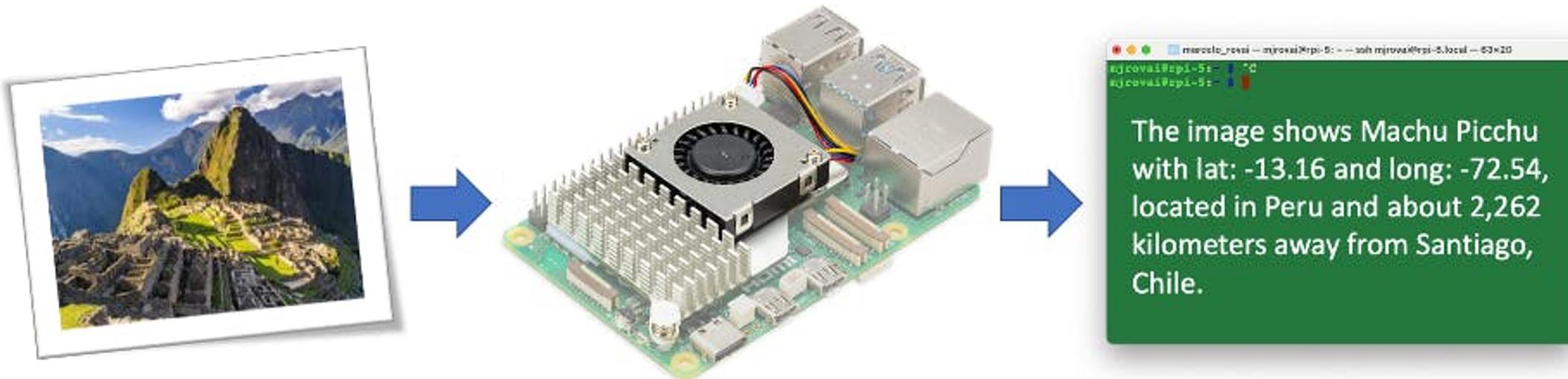
The image shows Machu Picchu, with lat:-13.16 and long: -72.54,
located in Peru and about 2,250 kilometers away from Santiago,
Chile.

[INFO] ==> The code (running llava-phi3), took 267.579568572007
7 seconds to execute.

mjrovai@rpi-5:~/Documents/OLLAMA $
```

(4 minutes)

Function Calling



LLMs: Optimization Techniques

LLMs: Common Optimization Techniques

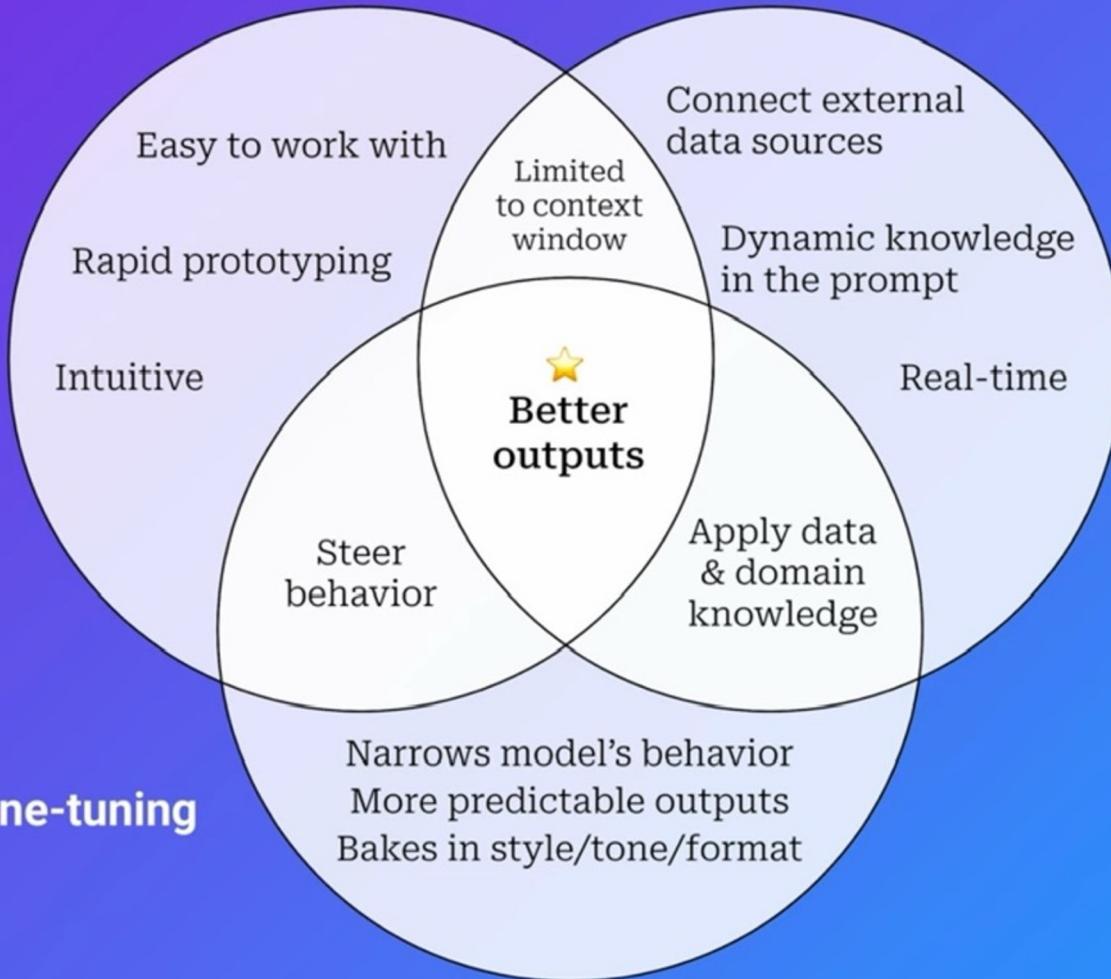
1. **Prompt Engineering**: Tailor your interactions.
2. **Fine-tuning**: Perfect the model's tasks.
3. **RAG**: Enhance with relevant data.

Comparison of Techniques

Prompt Engineering



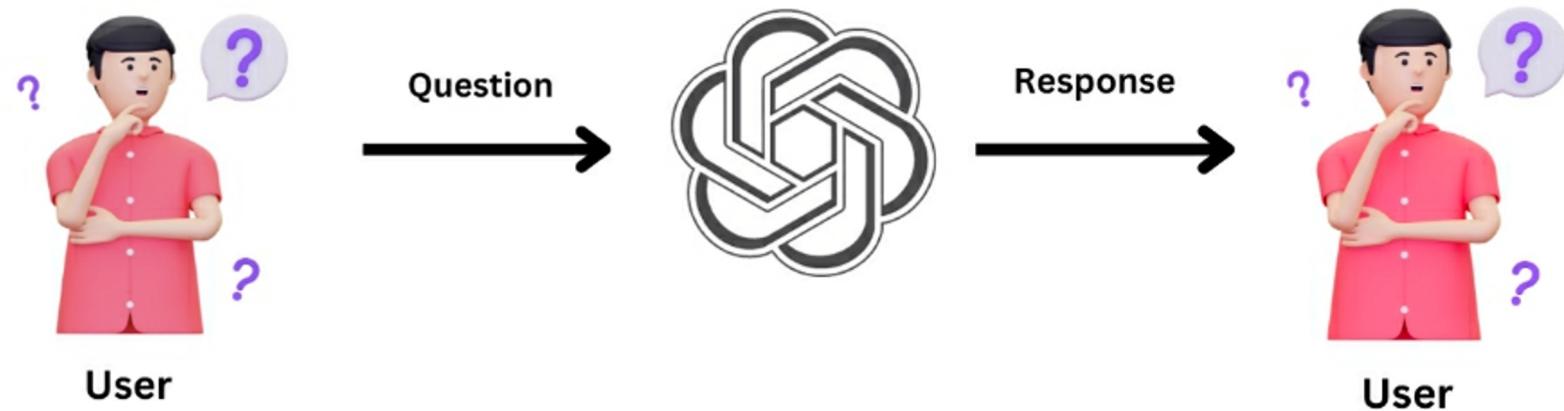
Fine-tuning



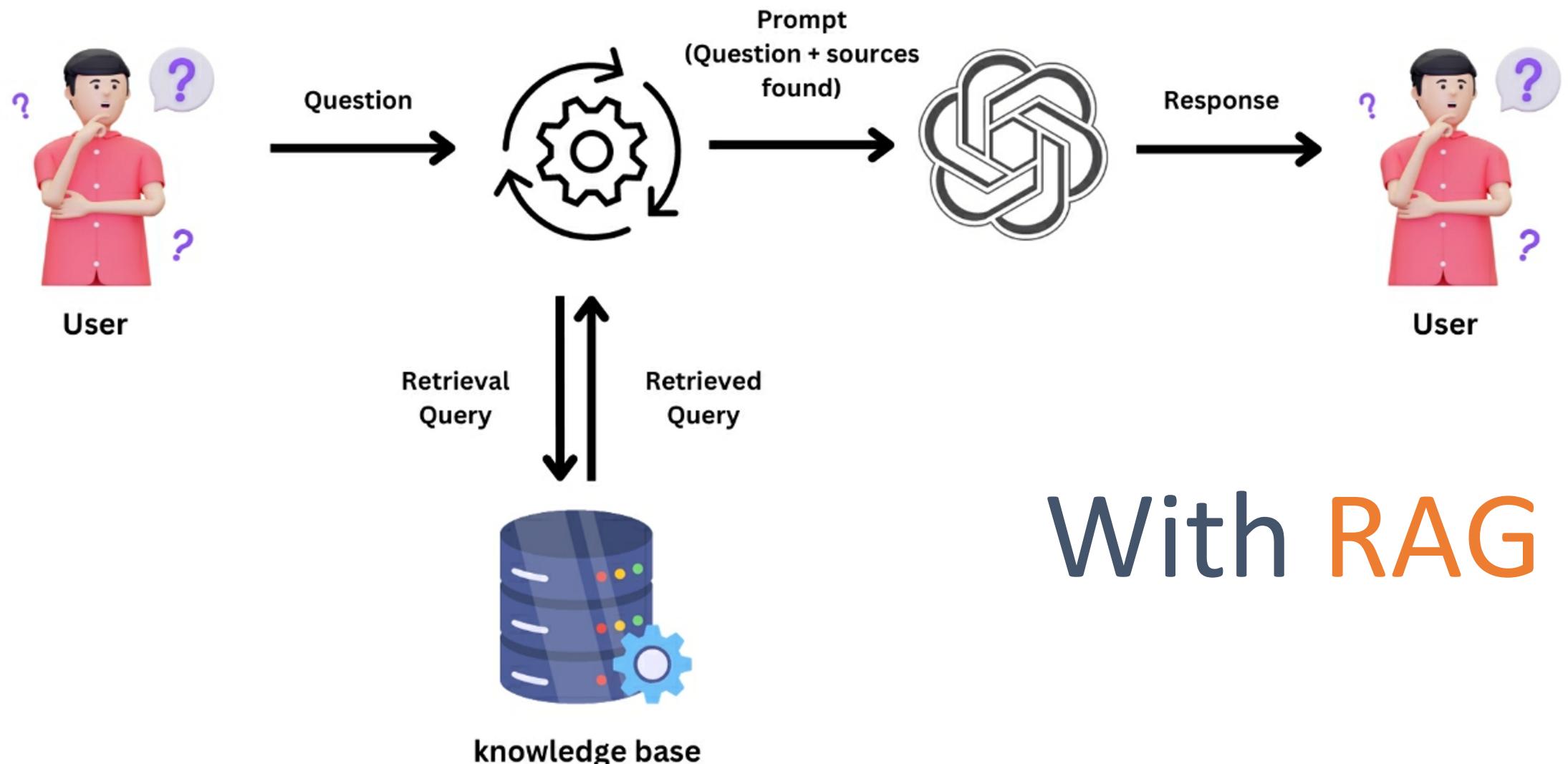
RAG

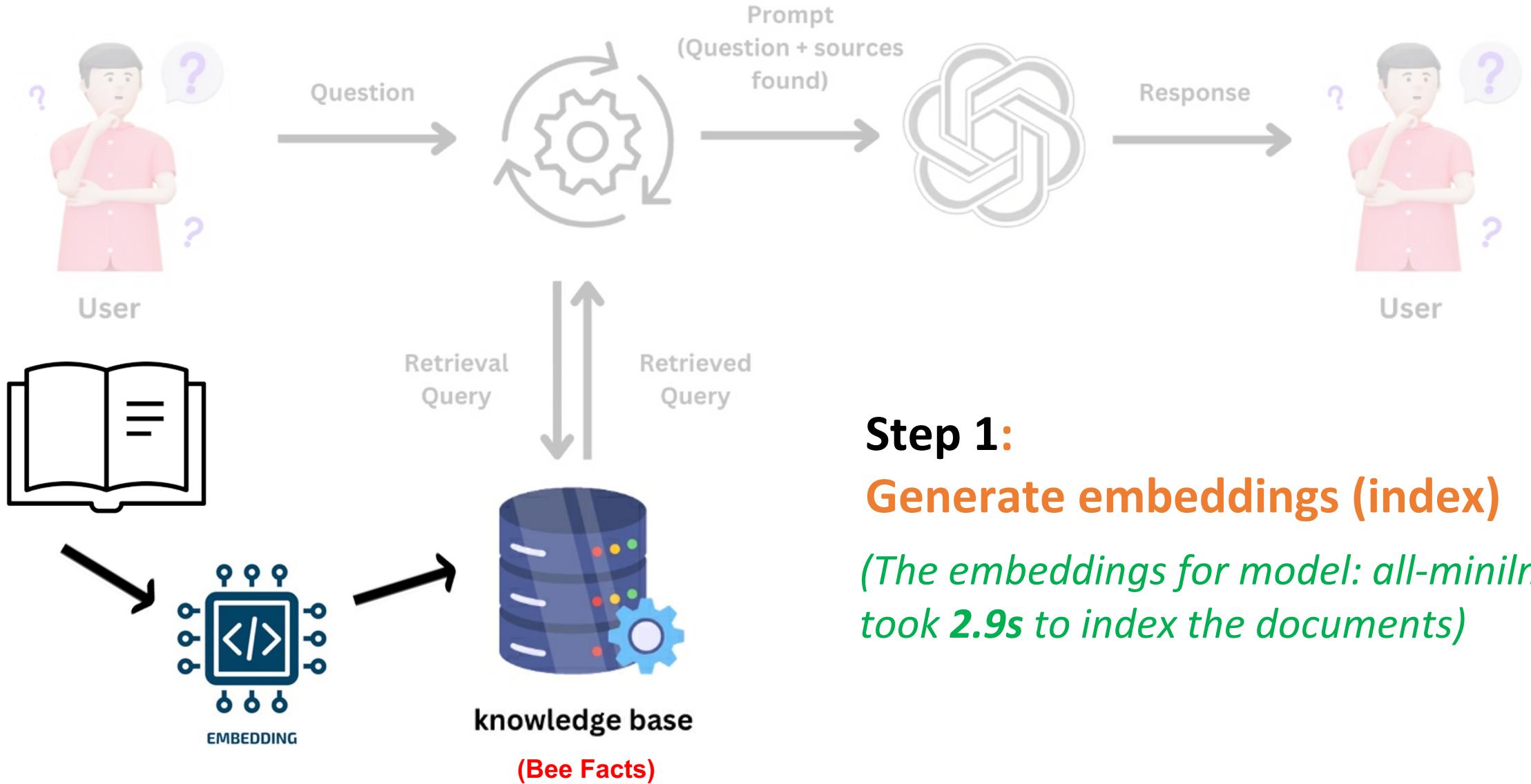
Retrieval-Augmented Generation (RAG)

“A method created by the FAIR team at Meta to enhance the accuracy of Large Language Models (LLMs) and reduce false information or “hallucinations.”



Usual Prompt





Step 1:
Generate embeddings (index)
(The embeddings for model: all-minilm, took 2.9s to index the documents)

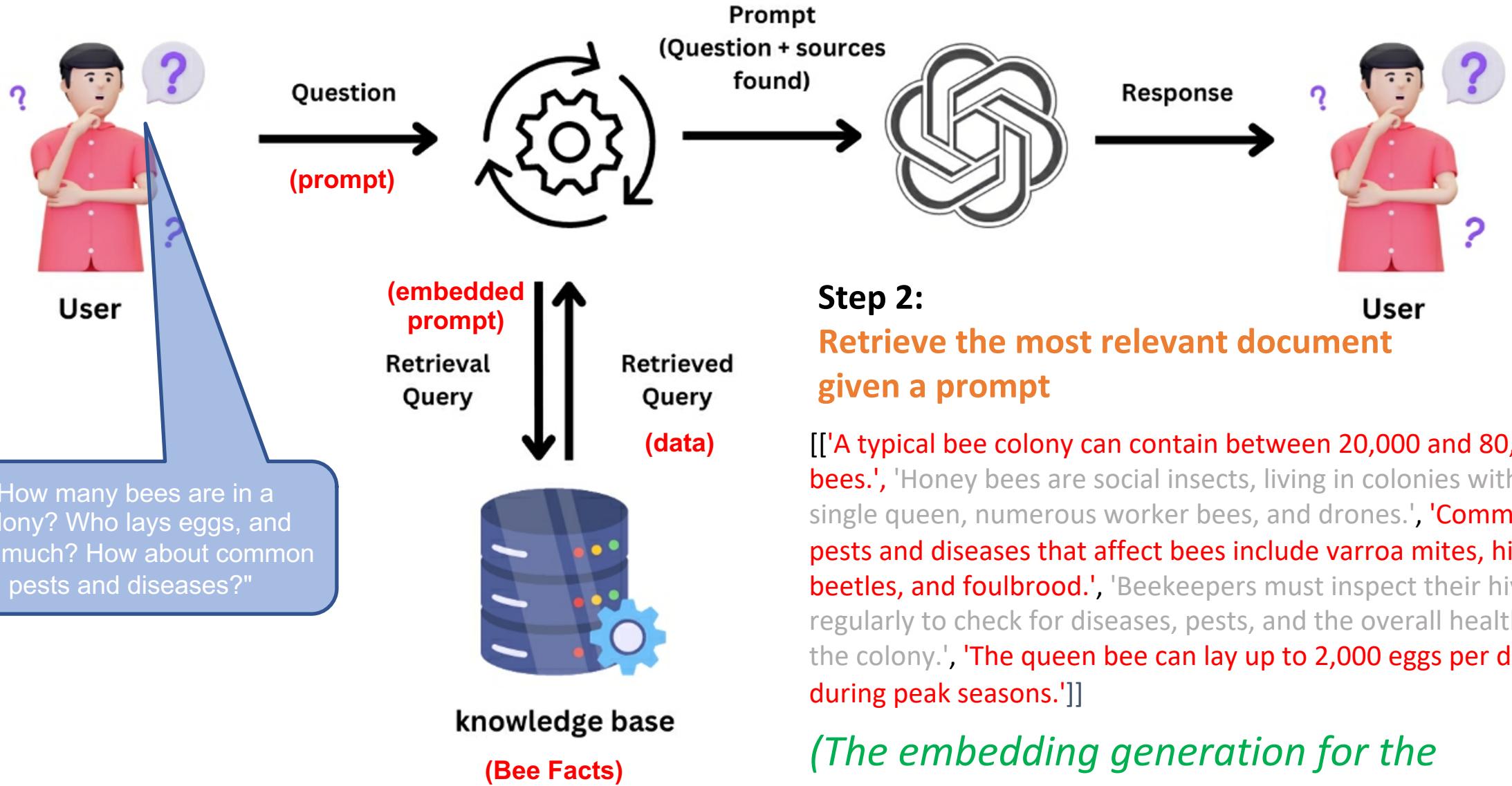
The screenshot shows a code editor window with two tabs: `rag_test.py` and `ppt.py`. The `ppt.py` tab is active, displaying the following Python code:

```
OPEN FILES ◀ rag_test.py × ppt.py • ppt.py  
1 # Step 1: Generate embeddings (index)  
2  
3  
4 import ollama  
5 import chromadb  
6  
7  
8 EMB_MODEL = "all-minilm" # "nomic-embed-text" #"mxbai-embed-large"  
9  
10 documents = [  
11     "Bee-keeping, also known as apiculture, involves the maintenance of bee colonies, typically in hives, by humans.",  
12     "The most commonly kept species of bees is the European honey bee (Apis mellifera).",  
13     "Bee-keeping dates back to at least 4,500 years ago, with evidence of ancient Egyptians practicing it.",  
14     "A beekeeper's primary role is to manage hives to ensure the health of the bee colony and maximize honey production.",  
15     "Honey bees are social insects, living in colonies with a single queen, numerous worker bees, and drones.",  
16     "The queen bee can lay up to 2,000 eggs per day during peak seasons.",  
17     "Worker bees are female and perform all the tasks in the hive except for reproduction.",  
18     "Drones are male bees whose primary role is to mate with a queen from another hive.",  
19     "Honey bees communicate with each other through the 'waggle dance,' which indicates the direction and distance to food sources.",  
20     "Bees produce honey from the nectar they collect from flowers, which they store in the hive for food during winter.",  
21     "Bees also produce beeswax, which they use to build the honeycomb structure in the hive.",  
22     "Propolis, another bee product, is a resin-like substance collected from tree buds and used to seal gaps in the hive.",  
23     "Bees play a crucial role in pollination, which is essential for the reproduction of many plants and crops.",  
24     "A typical bee colony can contain between 20,000 and 80,000 bees.",  
25     "Bee-keeping can be done for various purposes, including honey production, pollination services, and the sale of bees and related products.",  
26     "Beekeepers must inspect their hives regularly to check for diseases, pests, and the overall health of the colony.",  
27     "Common pests and diseases that affect bees include varroa mites, hive beetles, and foulbrood.",  
28     "Bee-keeping requires protective clothing and equipment, such as a bee suit, gloves, and a smoker to calm the bees.",  
29     "Sustainable bee-keeping practices are important for maintaining healthy bee populations and ecosystems.",  
30     "Beekeeping can be a hobby, a part-time occupation, or a full-time profession, depending on the scale and intent of the beekeeper."  
31 ]  
32  
33 client = chromadb.Client()  
34 collection = client.create_collection(name="bee_facts")  
35  
36 # store each document in a vector embedding database  
37 for i, d in enumerate(documents):  
38     response = ollama.embeddings(model=EMB_MODEL, prompt=d)  
39     embedding = response["embedding"]  
40     collection.add(  
41         ids=[str(i)],  
42         embeddings=[embedding],  
43         documents=[d]  
44     )  
45
```

The code uses the `ollama` library to generate embeddings for a list of bee-related facts stored in the `documents` list. These facts cover topics like the history of bee-keeping, bee biology, products, and best practices. The generated embeddings are then stored in a `chromadb` collection named `bee_facts`.

Spaces: 2

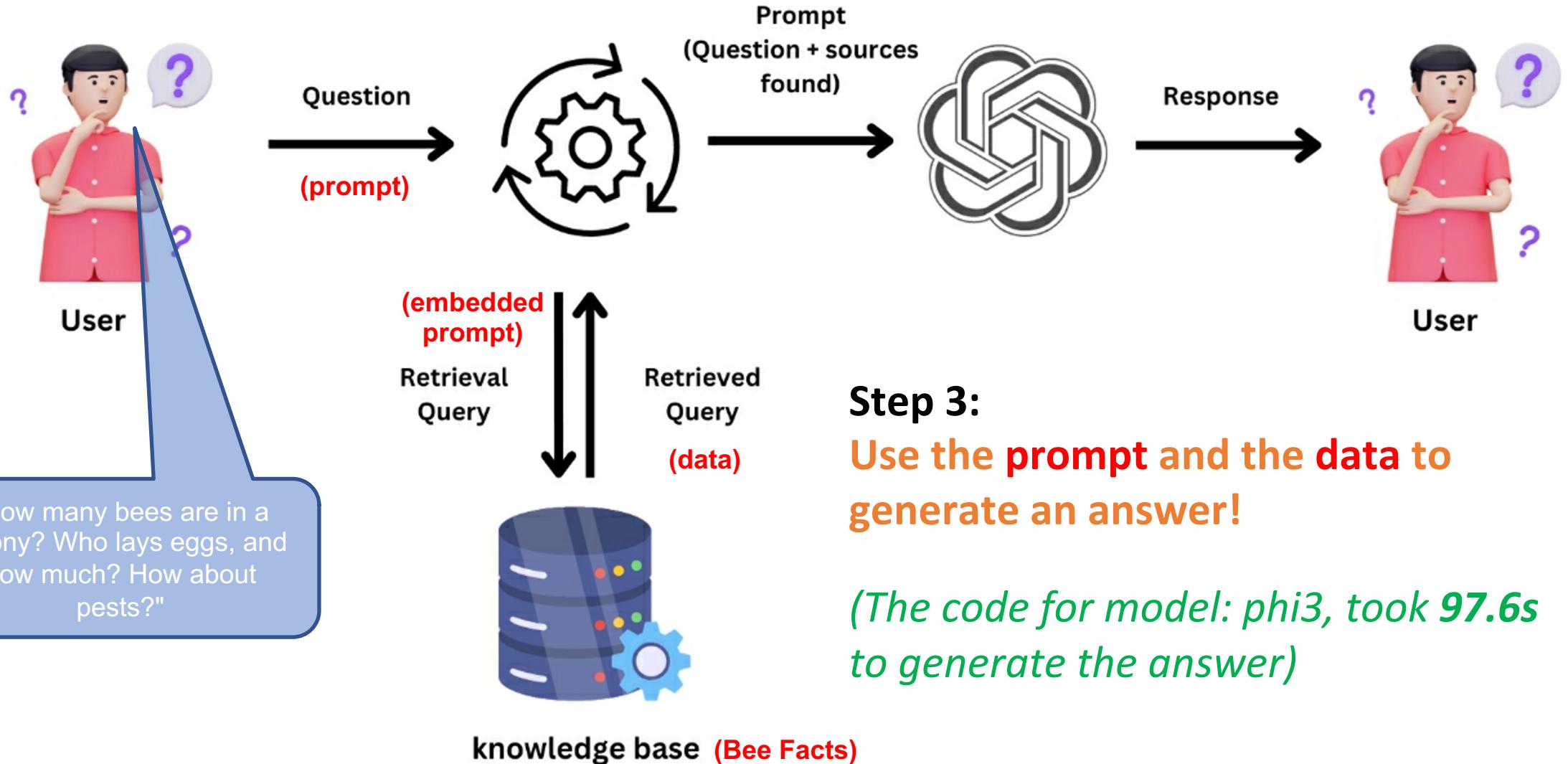
Python



A screenshot of a code editor window titled "ppt.py". The window shows a Python script with line numbers and code. The code is as follows:

```
1 # Step 2: Retrieve the most relevant document given a prompt:  
2  
3  
4  
5 # Prompt  
6 prompt = "How many bees are in a colony? Who lays eggs and how much? How about common pests and diseases?"  
7  
8 # generate an embedding for the prompt and retrieve the most relevant doc  
9 response = ollama.embeddings(  
10     prompt=prompt,  
11     model=EMB_MODEL  
12 )  
13 results = collection.query(  
14     query_embeddings=[response["embedding"]],  
15     n_results=5  
16 )  
17 data = results['documents']  
18
```

The code editor interface includes a toolbar at the top with icons for file operations, a status bar at the bottom indicating "Line 3, Column 1", and a status bar on the right showing "Spaces: 2" and "Python". The title bar also displays "UNREGISTERED".



The screenshot shows a code editor window with the following details:

- File Tabs:** The tabs are labeled "rag_test.py" and "ppt.py".
- Open Files:** The sidebar shows "OPEN FILES" with "rag_test.py" and "ppt.py" listed.
- Code Content:** The "ppt.py" tab contains the following Python code:

```
1 # Step 3: Use the prompt and the data to generate an answer!
2
3 MODEL = "phi3"
4
5
6 # generate a response combining the prompt and data we retrieved in step 2
7 output = ollama.generate(
8     model=MODEL,
9     prompt=f"Using this data: {data}. Respond to this prompt: {prompt}",
10    options={
11        "temperature": 0.0,
12        "top_k":10,
13        "top_p":0.5
14    }
15 )
16
```
- Status Bar:** The bottom status bar indicates "Line 16, Column 1", "Spaces: 2", and "Python".
- Header:** The top right corner says "UNREGISTERED".

Question:

"How many bees are in a colony? Who lays eggs, and how much?
How about common pests and diseases?"

Response

A typical bee colony contains between 20,000 and 80,000 bees. The queen bee is responsible for laying the majority of these eggs; she can produce up to 2,000 eggs per day during peak seasons. Beekeepers must regularly inspect their hives not only to monitor egg-laying but also to check for common pests and diseases that affect bees such as varroa mites, hive beetles, and foulbrood disease.

The screenshot shows a Visual Studio Code (VS Code) interface running on a Raspberry Pi. The title bar indicates the file is "rag_test.py - OLLAMA - Visual Studio Code". The top status bar shows system icons for battery, signal, and temperature (51°). The left sidebar includes a "Wastebasket" icon and a tree view of the project structure under "EXPLORER". The main editor window displays the "rag_test.py" file, which imports ollama, chromadb, and time, and defines a list of documents about beekeeping. The bottom status bar shows the terminal command "sudo raspi-config" and the status "Ln 15, Col 15".

```
rag_test.py - OLLAMA - Visual Studio Code

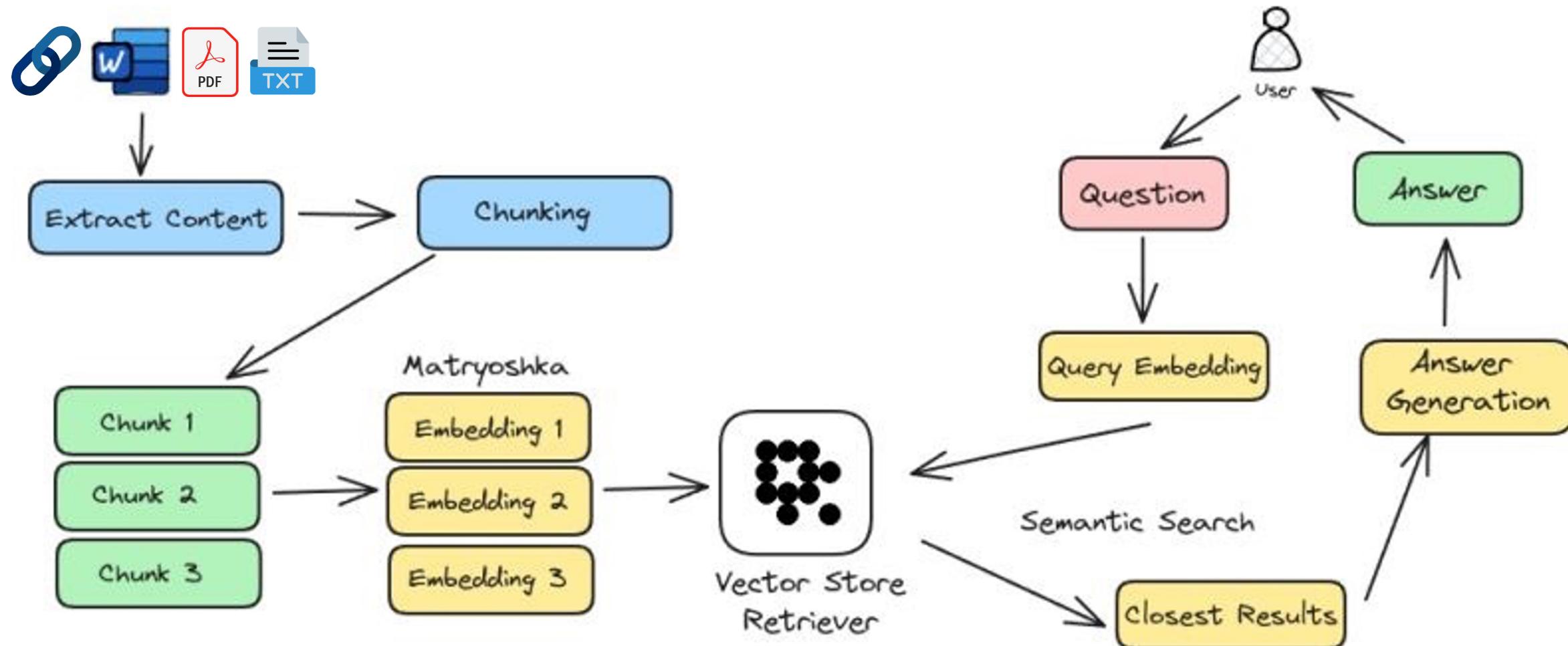
EXPLORER    indexer.py  simple_rag.py  example.py  rag_test.py  calc_distance_image.py
RAG > RAG_test > rag_test.py > ...
8
9     import ollama
10    import chromadb
11    import time
12
13   start_time = time.perf_counter() # Start timing
14   EMB_MODEL = "all-minilm" #"nomic-embed-text" #"mbai-embed-large"
15   MODEL = "phi3"
16
17  documents = [
18      "Bee-keeping, also known as apiculture, involves the maintenance of bee colonies, typically in hives",
19      "The most commonly kept species of bees is the European honey bee (Apis mellifera).",
20      "Bee-keeping dates back to at least 4,500 years ago, with evidence of ancient Egyptians practicing it",
21      "A beekeeper's primary role is to manage hives to ensure the health of the bee colony and maximize honey",
22      "Honey bees are social insects, living in colonies with a single queen, numerous worker bees, and drones",
23      "The queen bee can lay up to 2,000 eggs per day during peak seasons.",
24      "Worker bees are female and perform all the tasks in the hive except for reproduction.",
25      "Drones are male bees whose primary role is to mate with a queen from another hive.",
26      "Honey bees communicate with each other through the 'waggle dance,' which indicates the direction and distance of food sources.",
27      "Bees produce honey from the nectar they collect from flowers, which they store in the hive for food reserves.",
28      "Bees also produce beeswax, which they use to build the honeycomb structure in the hive.",
29      "Propolis, another bee product, is a resin-like substance collected from tree buds and used to seal gaps in the hive structure and protect against invaders.",
30      "Bees play a crucial role in pollination, which is essential for the reproduction of many plants and crops.",
31      "A typical bee colony can contain between 20,000 and 80,000 bees.",
32      "Bee-keeping can be done for various purposes, including honey production, pollination services, and research.",
33      "Beekeepers must inspect their hives regularly to check for diseases, pests, and the overall health of the colony.",
34      "Common pests and diseases that affect bees include varroa mites, hive beetles, and foulbrood.",
35      "Bee-keeping requires protective clothing and equipment, such as a bee suit, gloves, and a smoker to control the bees without harming them.",
36      "Sustainable bee-keeping practices are important for maintaining healthy bee populations and ecosystems."]

PROBLEMS 1 OUTPUT DEBUG CONSOLE TERMINAL PORTS
bash - RAG_test + - x
d honey and larvae; and foulbrood, a bacterial disease caused by Paenibacillus larvae that can devastate young bee populations. The European honey bee (Apis mellifera) is the most commonly kept species of bees worldwide.

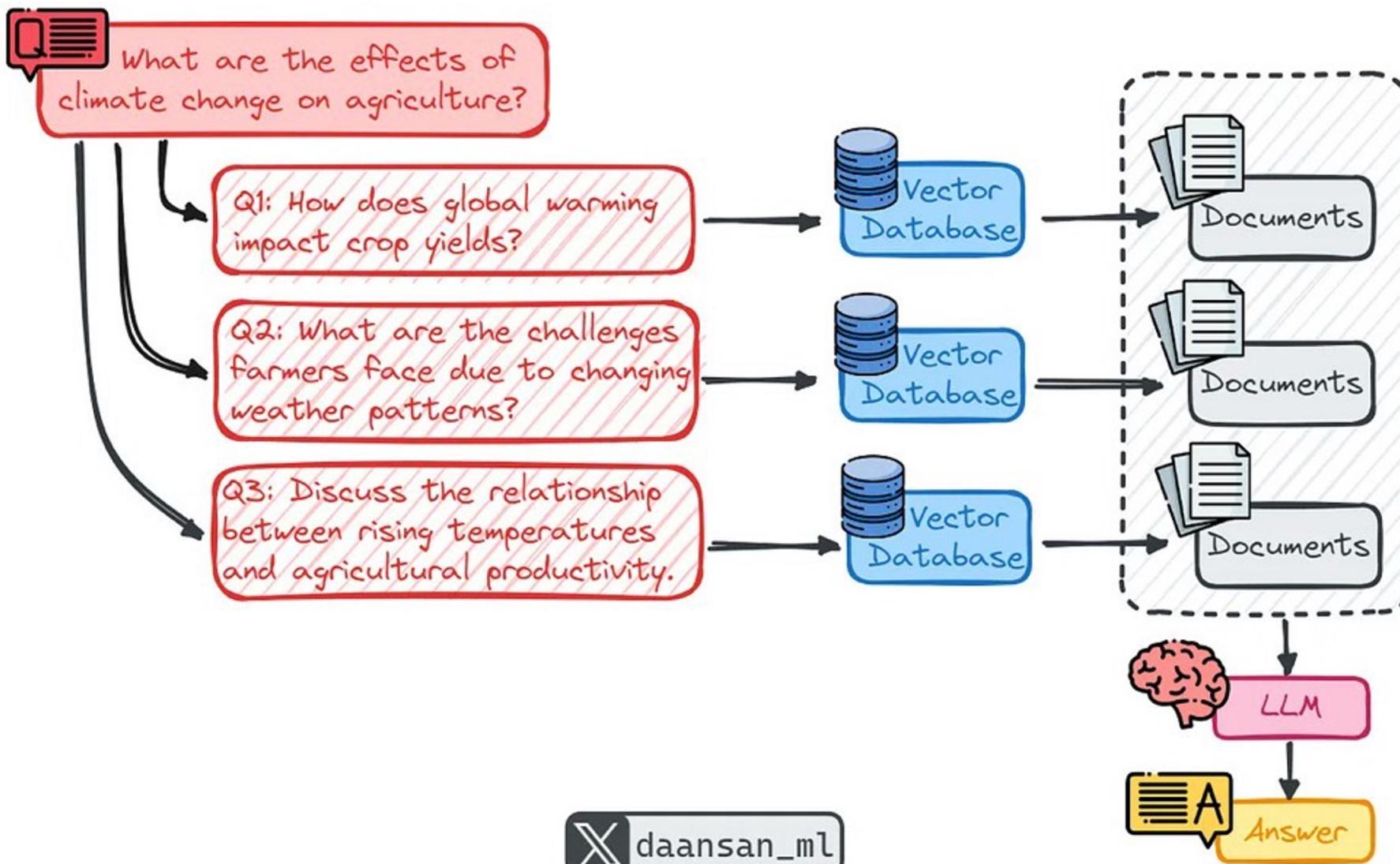
[INFO] ==> The code for model: phi3, took 97.6s to generate the answer.

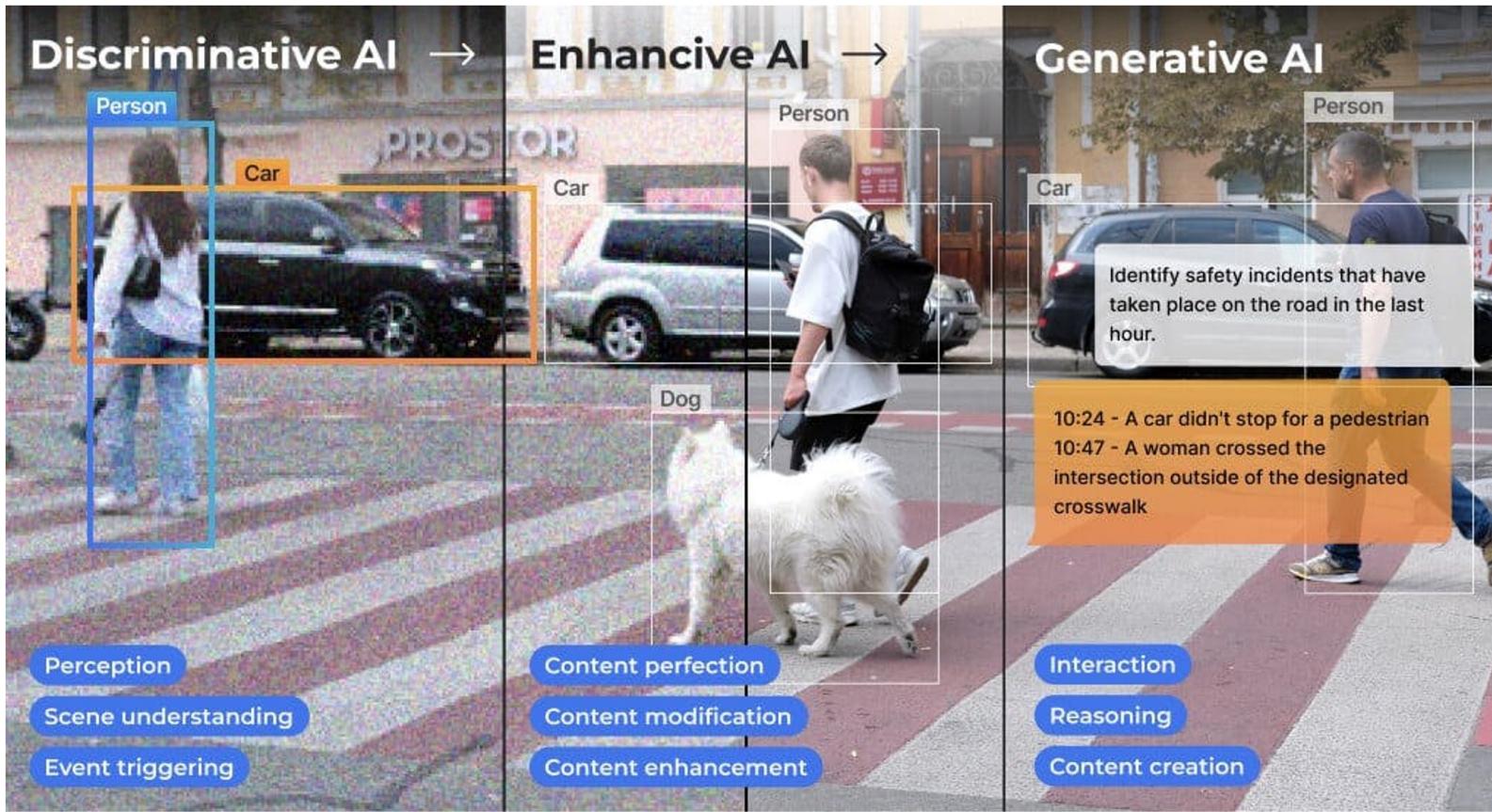
mjrovai@rpi-5:~/Documents/OLLAMA/RAG/RAG_test $ sudo raspi-config
```

RAG: Simple Query



Advanced RAG: Multi Query





"In the vast landscape of artificial intelligence (AI), one of the most intriguing journeys has been the evolution of AI on the edge. This journey has taken us from classic machine vision to the realms of discriminative AI, enhancive AI, and now, the groundbreaking frontier of generative AI. Each step has brought us closer to a future where intelligent systems seamlessly integrate with our daily lives, offering an immersive experience of not just perception but also creation at the palm of our hand."

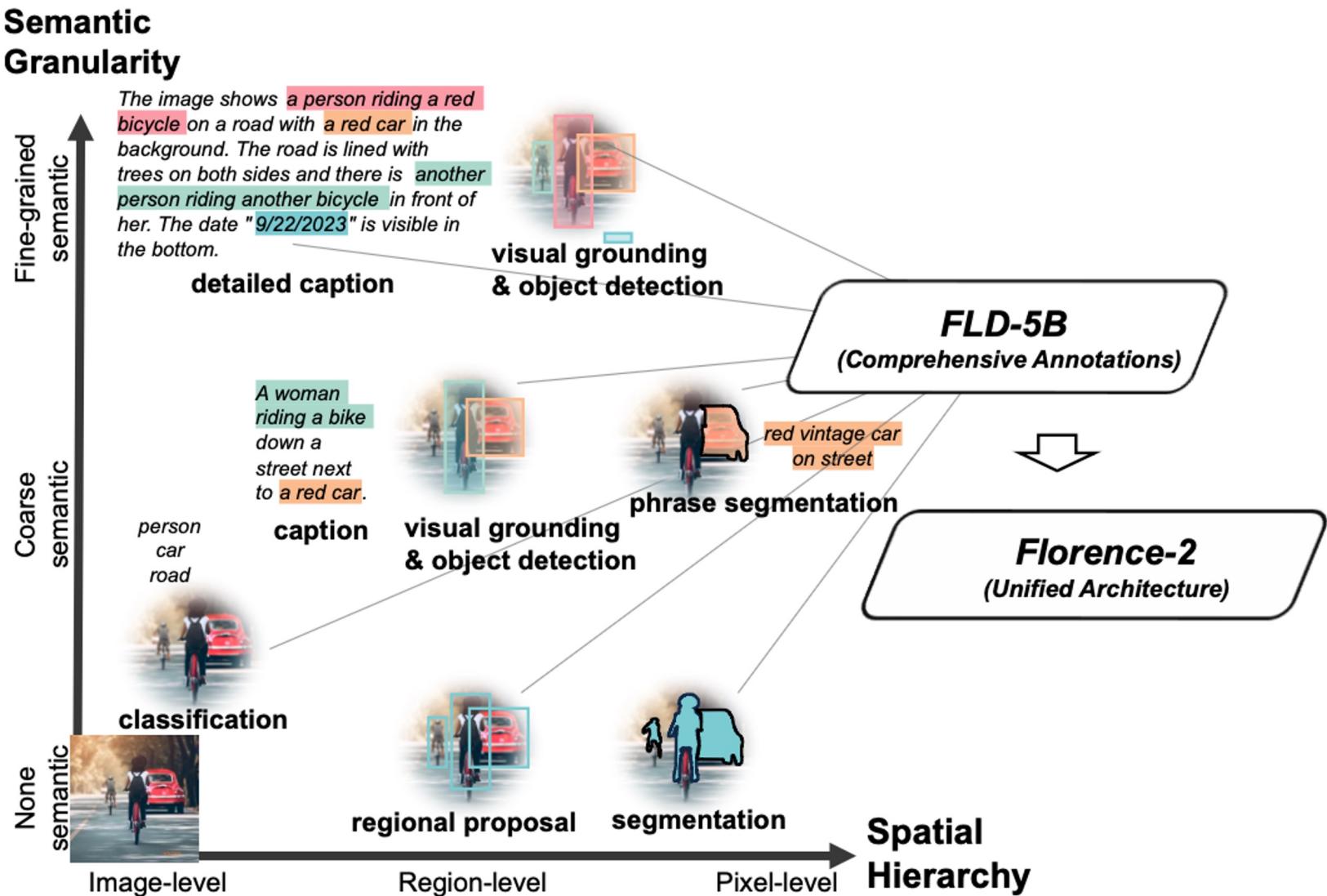
Avi Baum, CTO at Hailo

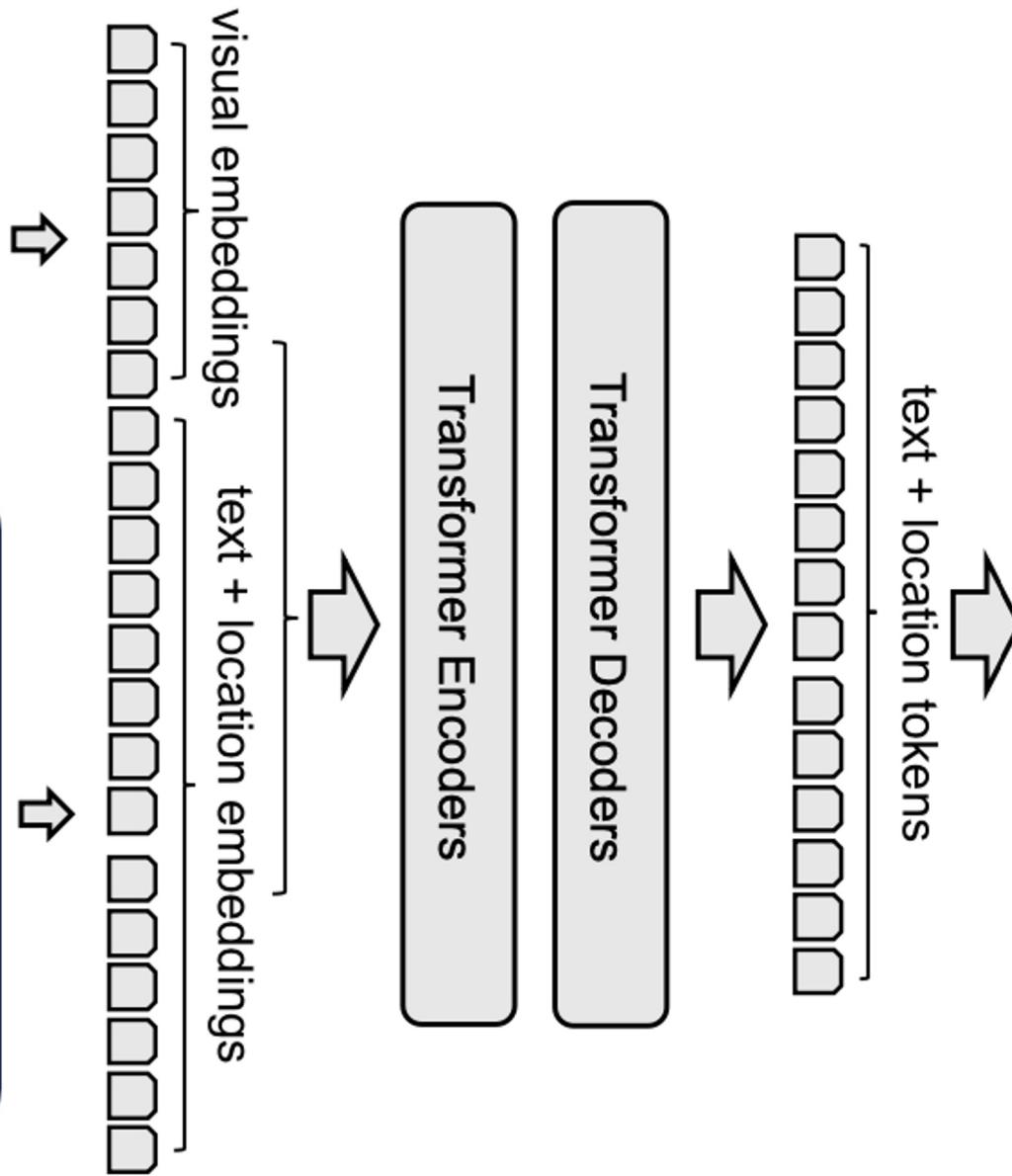
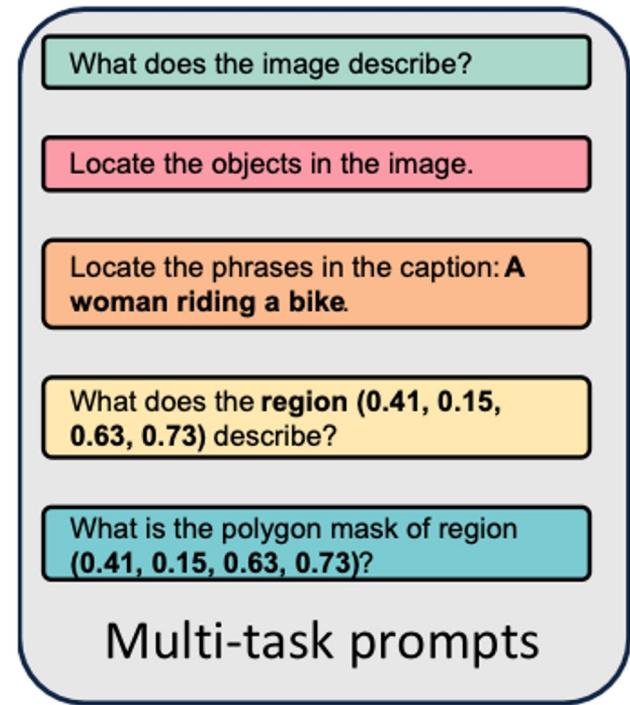
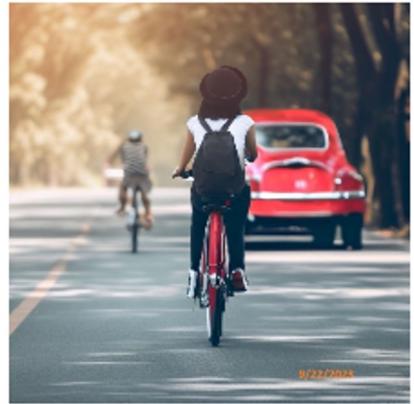
Florence-2

Advancing a Unified Representation for a Variety of Vision Tasks

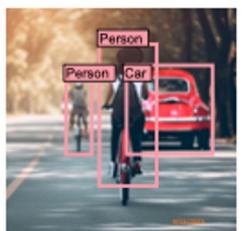


Paper: <https://arxiv.org/abs/2311.06242>





The image shows a person riding a red bicycle on a road with a red car in the background. The person is wearing a white t-shirt, black pants, and a black hat. She has a backpack on her back and is pedaling with their feet on the pedals. The road is lined with trees on both sides and there is another person riding another bicycle in front of her. The date "9/22/2023" is visible in the bottom right corner of the image.



microsoft/BitNet

Official inference framework for 1-bit LLMs

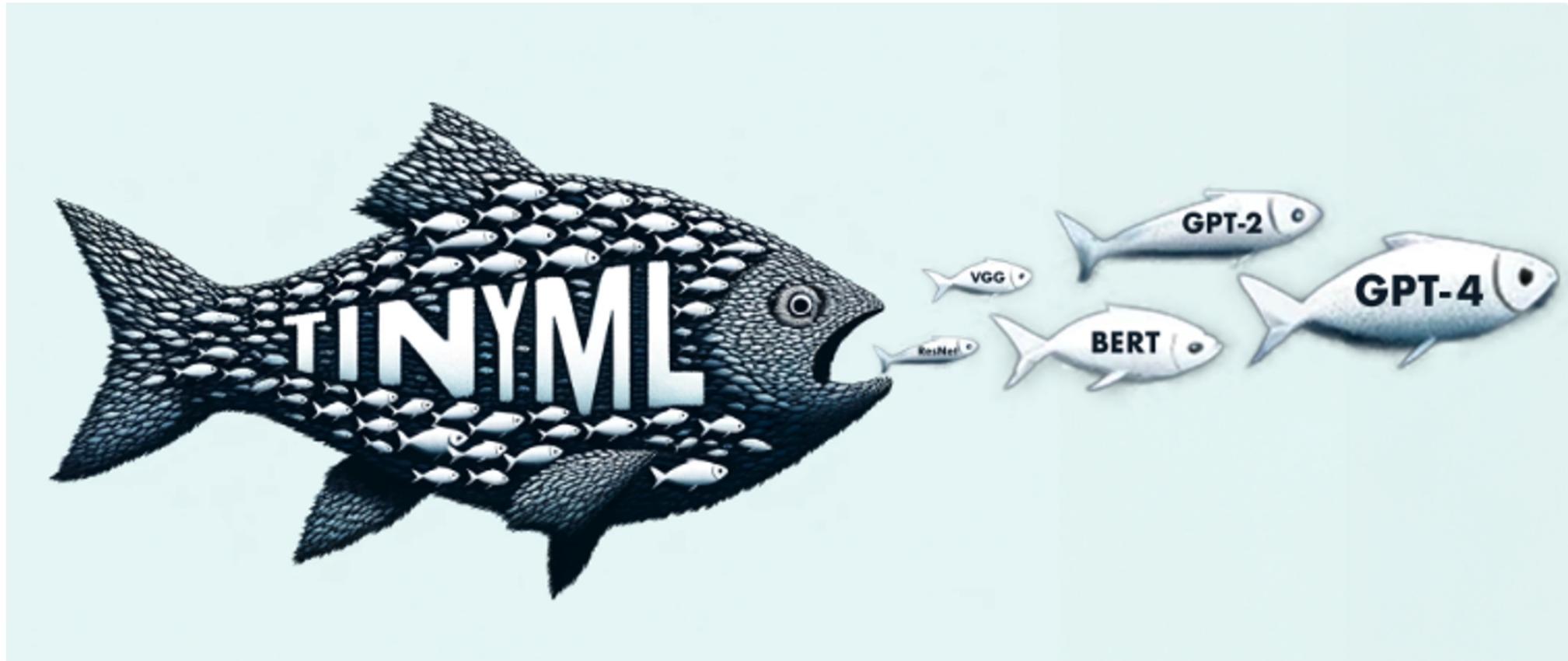


Bitnet.cpp employs one-bit quantization, representing values with a ternary system **(+1, -1, 0)**. This approach simplifies calculations by replacing complex multiplications with additions and subtractions, eliminating the need for GPUs.

- Speedups range from 1.37x to 6.1x on various CPUs.
- Power consumption reductions between 55.4% and 82.2% compared to traditional GPU-based inference.

[bitnet.cpp](#)

TinyML: Why the Future of Machine Learning is Tiny and Bright



Shvetank Prakash, Emil Njor, Colby Banbury, Matthew Stewart, Vijay Janapa Reddi

Questions?



TINYML4D

