

The banner features a blue background with a large white circular graphic on the right containing the text '60 ICTP 1964-2024'. On the left, there's a dark blue rectangular area with white text: 'Workshop on TinyML for Sustainable Development'. Below this are three sections with icons and text: '22 - 26 July 2024', 'São Paulo, Brazil', and 'Deadline: 6 May 2024'. To the right of these sections is a 'FURTHER INFORMATION:' box with an email address (smr3961@ictp.it), a web link (<https://indico.ictp.it/event/10499/>), and a note about female scientists being encouraged to apply. At the bottom right are logos for Harvard John A. Paulson School of Engineering and Applied Sciences, IBM, UNIFEI, and TINY DL.

Large Language Models (LLMs) at the Edge

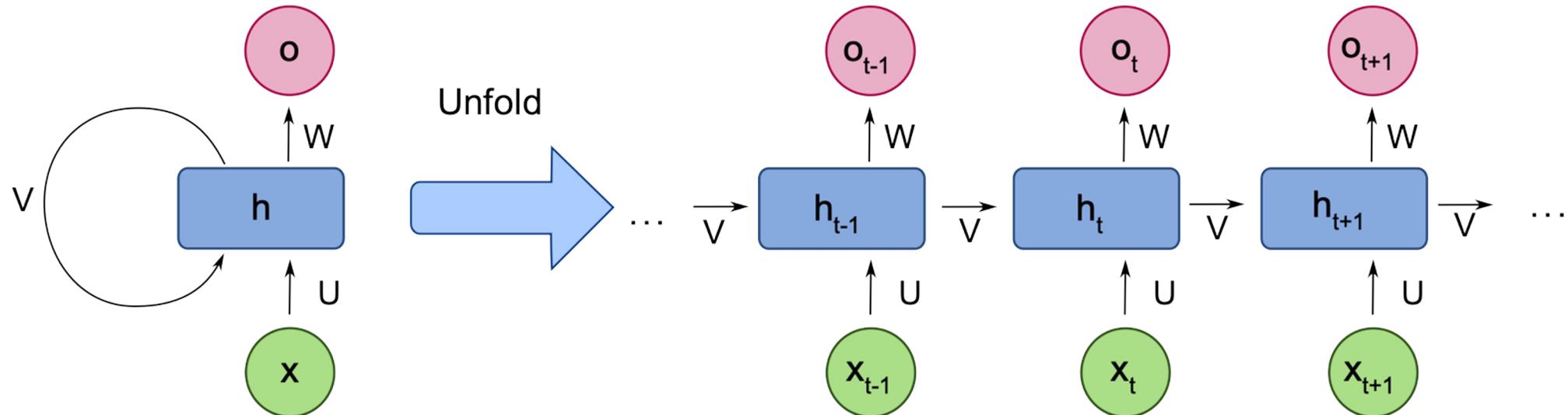
Prof. Marcelo J. Rovai
rovai@unifei.edu.br

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



Deep Learning models (or artificial neural networks)

Recurrent Neural Networks (RNNs): Designed for **sequential data like time series or text**, these networks use their internal state (memory) to process sequences of inputs.



Machado de Assis Bot with RNN - GRU

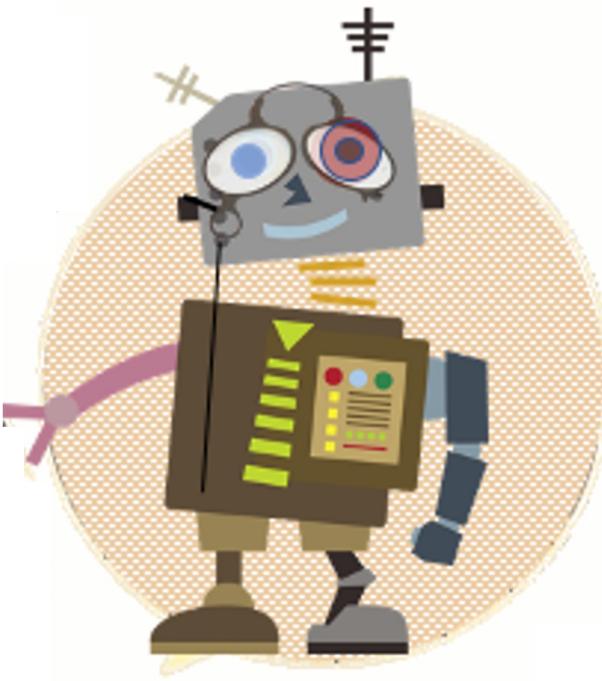


The robot writer model is a Recurrent Neural network (RNN —GRU). To obtain the final AI model, 3.5 million parameters were trained with a **120-letter sequence** from seven of his books: *Memorias Posthumas de Braz Cubas, Dom Casmurro, Quincas Borba, Papeis Avulsos, A Mão e a Luva, Esaú e Jacob, and Memorial de Ayres*.

Model: "sequential"

Layer (type)	Output Shape	Param #
embedding (Embedding)	(128, None, 64)	7488
gru (GRU)	(128, None, 1026)	3361176
dense (Dense)	(128, None, 117)	120159

Total params: 3,488,823
Trainable params: 3,488,823
Non-trainable params: 0



A LUVA DE CASMURRO II

A missa do coupé e um presente e o governo devia cazar logo no papel, a morte do autor, e todos os seus considerados de alegria. Era um espirito de vinte e cinco annos, e eu não estou alguns passos no cerebro, como de outra cousa. Deus me disse:

--Não digo que não. Se eu tivesse a intenção de um probosito. Palha acudiu a mulher, não havia nada. A noite vinha tambem para o seminario, tinha o aspecto do partido recto e de restaurar a minha mãe e do pae, pela primeira vez, a menor destinada a dispensar o chapéo, esperou que não vinhas com as suas mãos de creanças. A manhã della chegasse a baroneza e a maneira desta divida. Parece que é casada.

--Está bom, perdoa-lhe de todos os lados, a vida de que o comprar para o meu quarto de hora, e contavam com o fim de a anterior, e, a parede pouco tempo a alma de pessoas que definitivamente lhe interessam a menos para mim. De quando em quando, esses dous annos de conversação para o fim de deixar nenhuma pessoa que se dispersasse; mas não falo de uma cousa nem lhe pedia com a mão tremula, como se ella quizesse. Eu, apertando-lhe a mão, aliás o principio do governo, a proposito disso, com a desattenção de Estevão, e eu começou a aborrecel-o, e a solidão podia ser melhor, e a sympathia coloca da mãe, e não se sabe calar o enterro no meio do lagem, o que iam-se apanhados no chão, e para a mulher, não tendo visto, nem a mesma cousa.

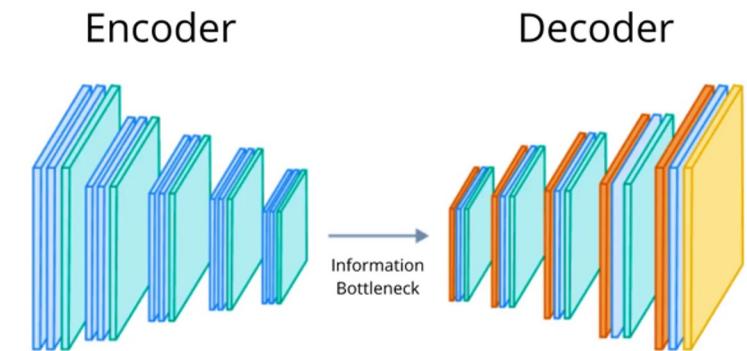
LLM / SLM

Large Language Model / Small Language Models

LLMs are **specialized deep learning models designed to understand and generate human language**, used for tasks like translation, summarization, and generating human-like text responses. SLMs are the same, but use a simpler, less resource-intensive approach (smaller in size).

Deep Learning models (or artificial neural networks)

- **Autoencoders:** Used primarily for unsupervised learning tasks such as dimensionality reduction and feature extraction, autoencoders learn to compress data from the input layer into a shorter code and then reconstruct the output from this representation.
- **Transformer Models:** Highly effective in handling sequences, transformers use mechanisms like self-attention to weigh the importance of different words in a sentence, regardless of their position. The Transformer architecture, while innovative, can be seen as a derivative of earlier deep learning models, particularly those based on the concept of sequence modeling. However, the most direct lineage can be traced to the sequence-to-sequence (seq2seq) models that utilize **encoder-decoder** architectures. These earlier seq2seq models were often built using **recurrent neural networks (RNNs)** or their more advanced variants like **LSTMs (Long Short-Term Memory Networks)** or **GRUs (Gated Recurrent Units)**.



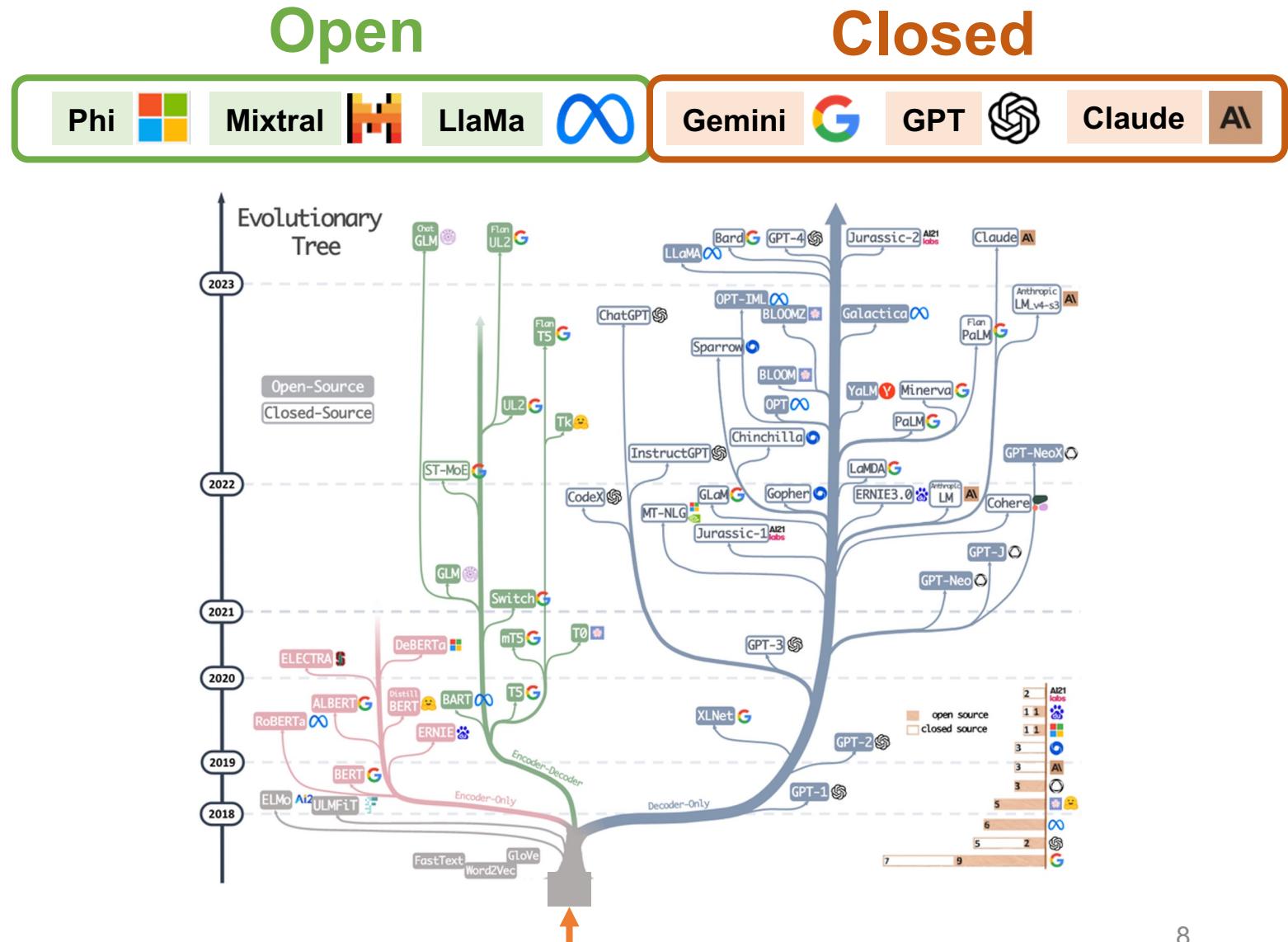
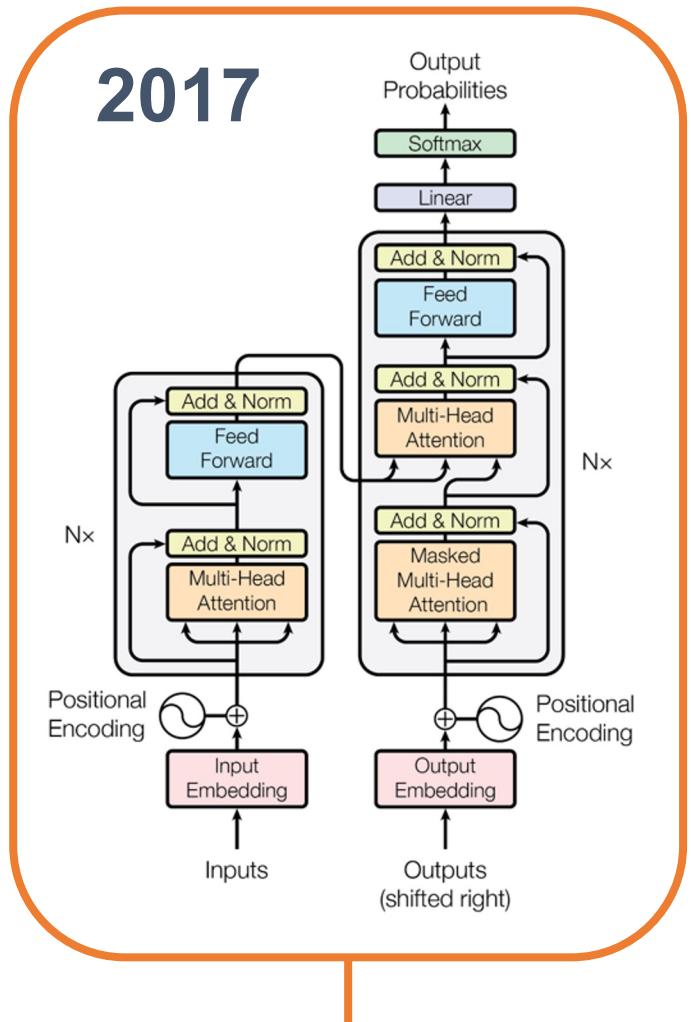
LLM/SLM – Large /Small Language Model

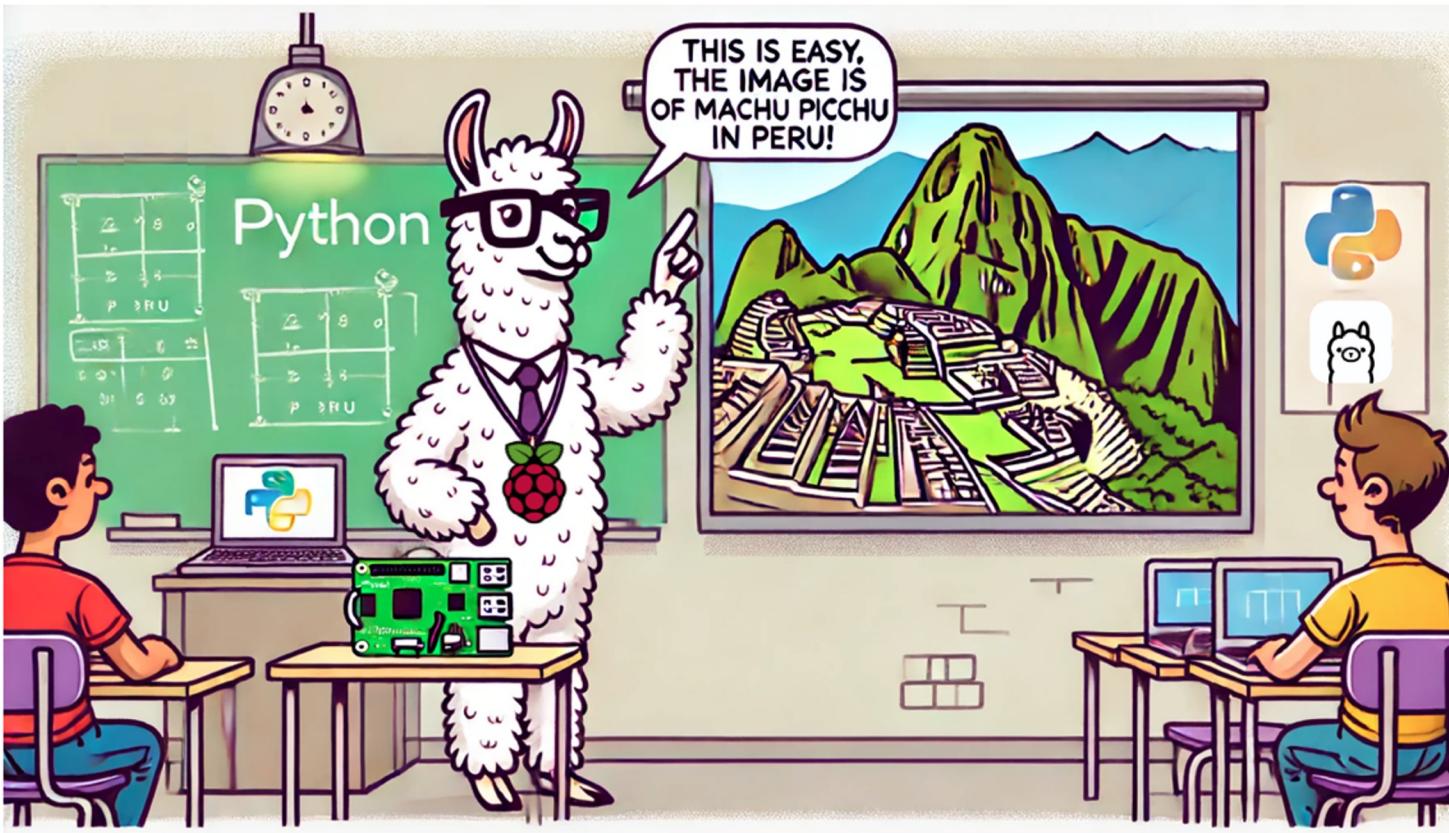
Large Language Models (LLMs) and SLMs are advanced neural networks based on the **Transformer architecture** that excel in understanding and generating human language. They represent a significant evolution from earlier sequence-based models like **LSTMs**, which surpass them in handling long-range dependencies and parallel processing efficiency.

Prof. Jesus's Presentation about IA:



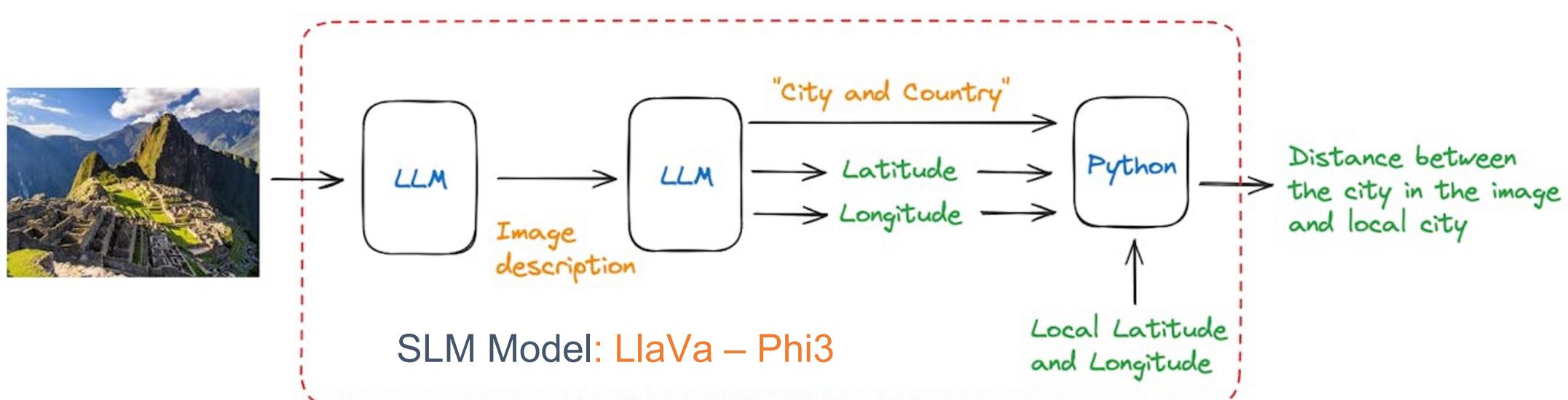
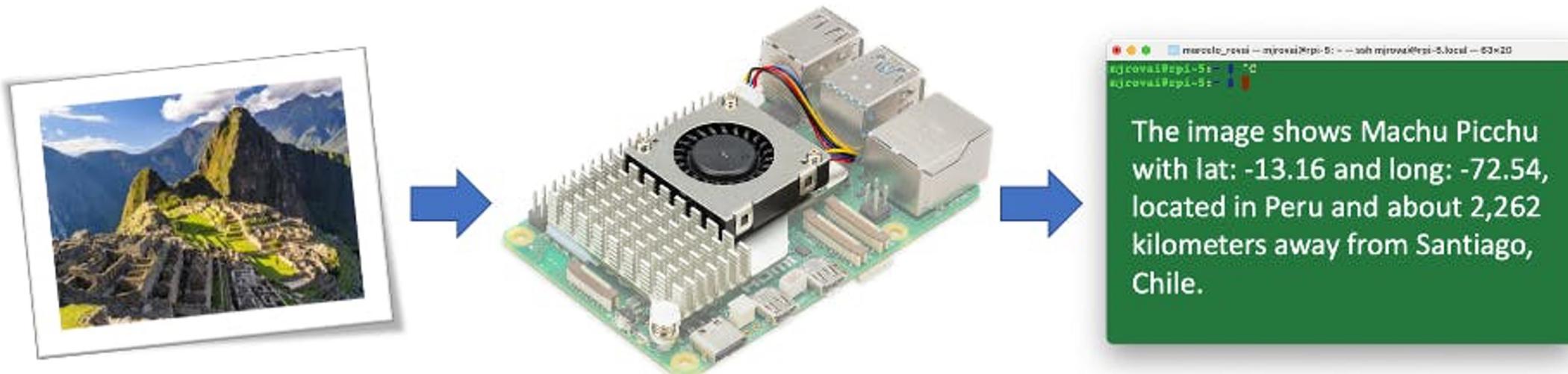
Transformers to LLMs and SLMs



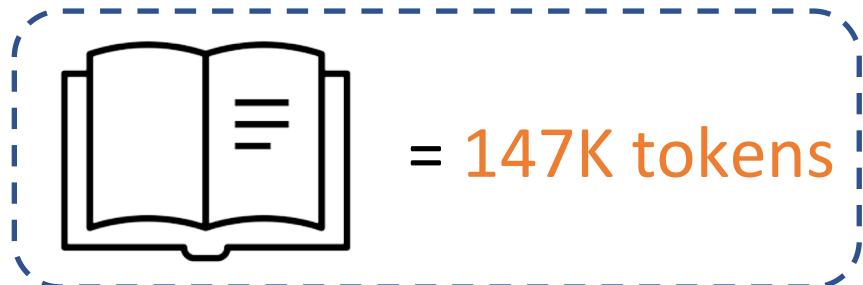


Running Large Language Models on Raspberry Pi at the Edge

Transform a Raspberry Pi into a powerful AI hub, running SLMs for real-time, on-site data analysis and insights using Ollama and Python.



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

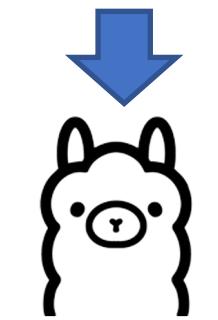


A **4-bit** quantized **3.8 billion parameter *** language model trained on **3.3 trillion tokens****, whose overall performance, as measured by both academic benchmarks and internal testing, rivals that of models such as Mixtral 8x7B and GPT-3.5

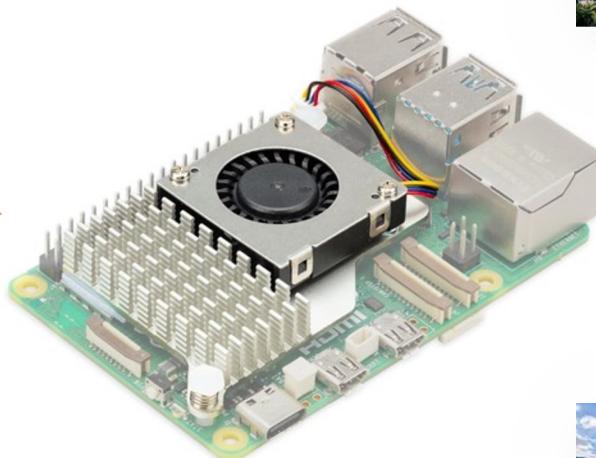
* 2.4 GB

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

llava-phi-3 (2.9 GB)



Ollama



```
mjrovai@rpi-5:~\n\nFile Edit Tabs Help\n\n>>> Answer with one short sentence, what is the capital of France and its distance\n... in Km from Santiago, Chile\nThe capital of France is Paris and it is around 12,674 kilometers away\nfrom Santiago, Chile.\n\nTotal duration: 13.860074968s\nload duration: 1.537039ms\nprompt eval count: 27 token(s)\nprompt eval duration: 5.925386s\nprompt eval rate: 4.56 tokens/s\neval count: 26 token(s)\neval duration: 7.539223s\neval rate: 3.45 tokens/s\n>>> Send a message (/? for help)
```

(13 seconds)



```
mjrovai@rpi-5:~/Documents/OLLAMA\n\nHelp\n\n/ /Documents/OLLAMA $\n/ /Documents/OLLAMA $ python calc_distance_image.py /\n/home/mjrovai/Documents/OLLAMA/image_test_1.jpg\n\nThe image shows Paris, with lat:48.86 and long: 2.35, located in\nFrance and about 11,630 kilometers away from Santiago, Chile.\n\n[INFO] ==> The code (running llava-phi3), took 232.60845186299412\nseconds to execute.\n\nmjrovai@rpi-5:~/Documents/OLLAMA $
```



```
mjrovai@rpi-5:~/Documents/OLLAMA\n\nHelp\n\n/ /Documents/OLLAMA $\n/ /Documents/OLLAMA $ python calc_distance_image.py /\n/home/mjrovai/Documents/OLLAMA/image_test_3.jpg\n\nThe image shows Machu Picchu, with lat:-13.16 and long: -72.54,\nlocated in Peru and about 2,250 kilometers away from Santiago,\nChile.\n\n[INFO] ==> The code (running llava-phi3), took 267.579568572007\n7 seconds to execute.\n\nmjrovai@rpi-5:~/Documents/OLLAMA $
```

(4 minutes)

LLMs: Optimization Techniques

LLMs: Optimization Techniques

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

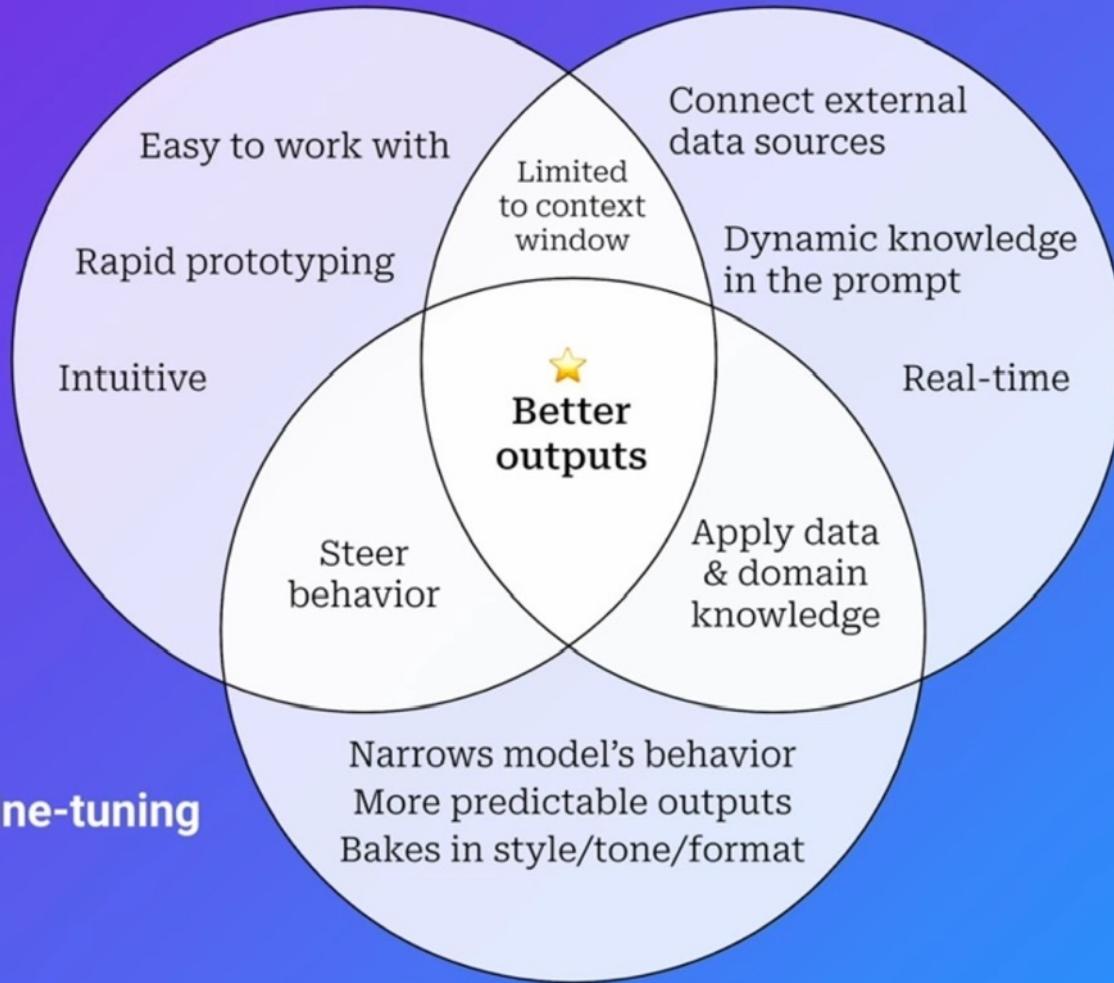
Comparison of Techniques

Prompt Engineering



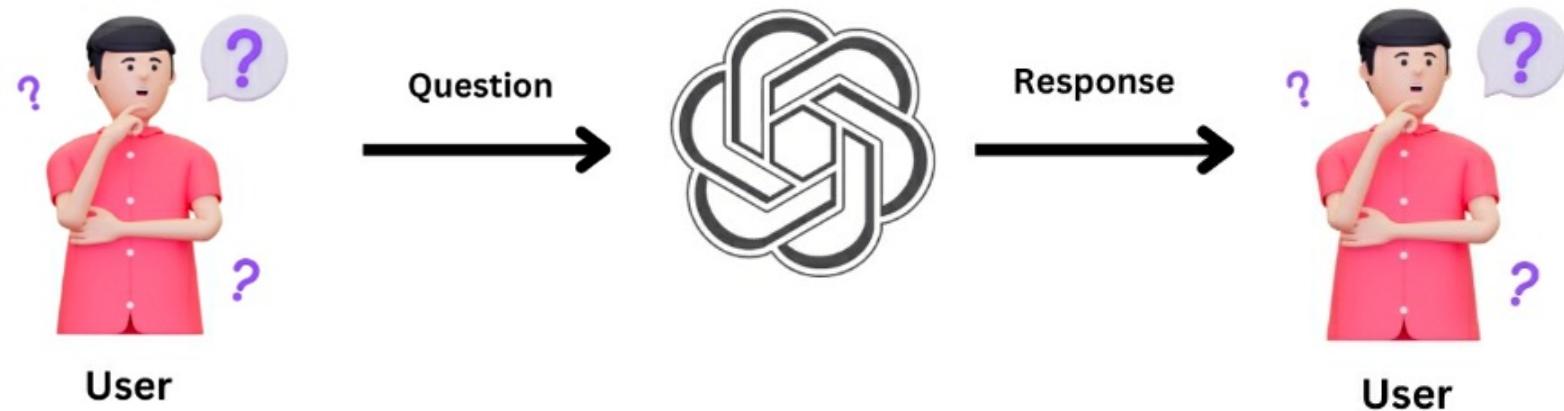
Fine-tuning

RAG

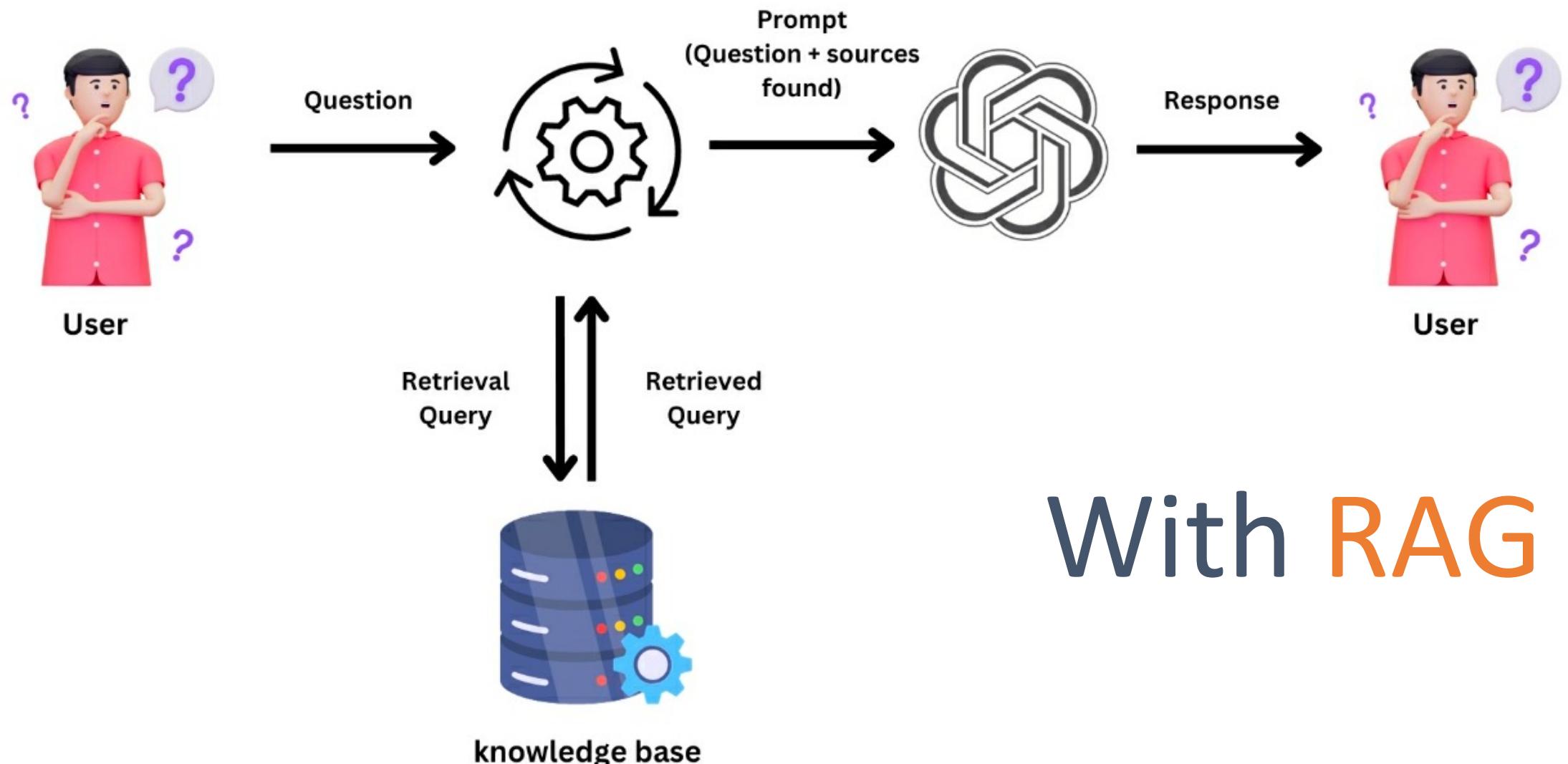


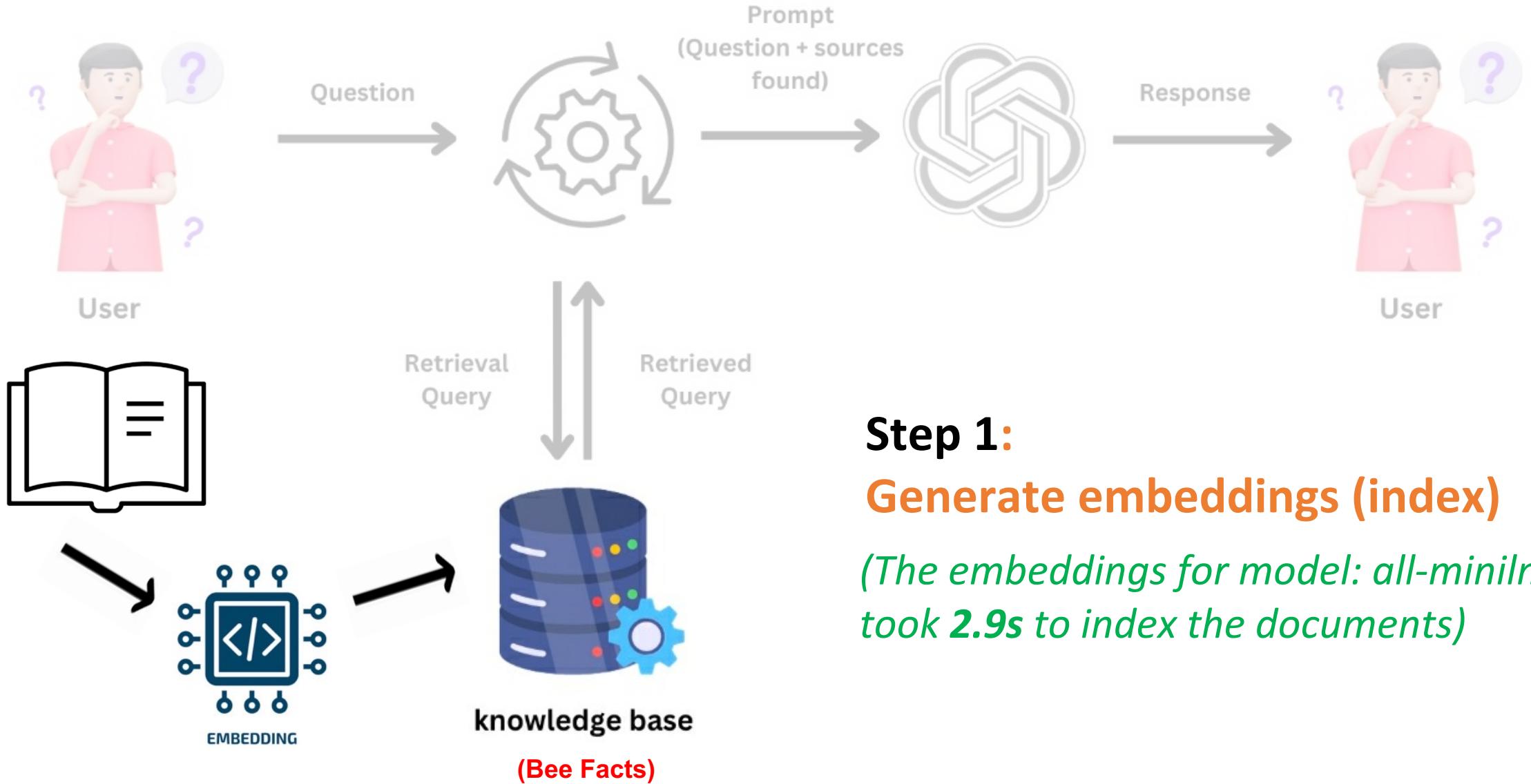
Retrieval-Augmented Generation (RAG)

“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



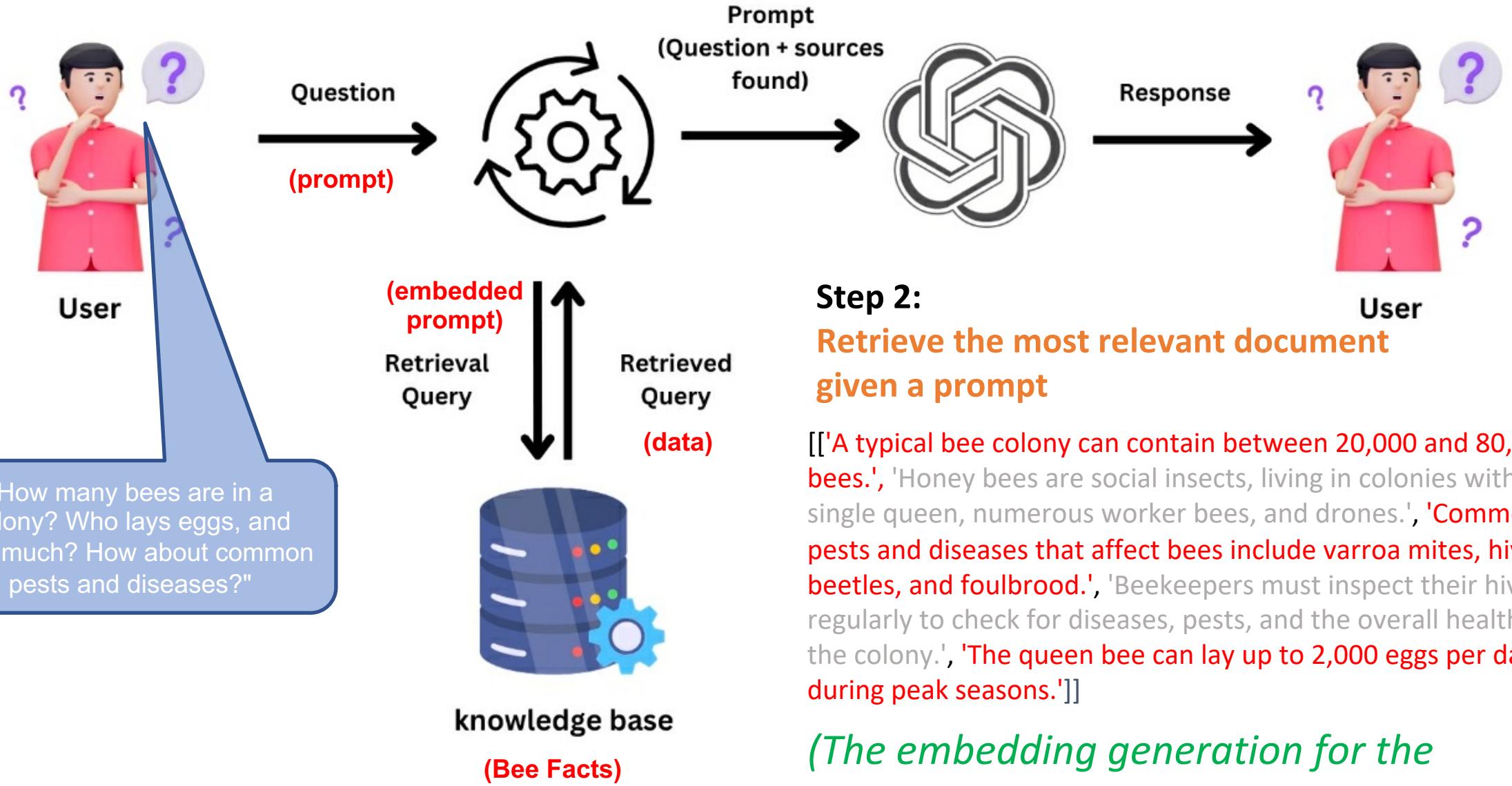


Step 1:
Generate embeddings (index)
(The embeddings for model: all-minilm, took 2.9s to index the documents)

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

- Open Files:** The sidebar on the left lists two files: `rag_test.py` (marked with an 'x') and `ppt.py` (marked with a dot).
- Title Bar:** The title bar shows the current file is `ppt.py`.
- Code Editor:** The main area displays the `ppt.py` script. The code uses the `ollama` and `chromadb` libraries to generate embeddings for a list of bee-related facts.
- Code Content:**

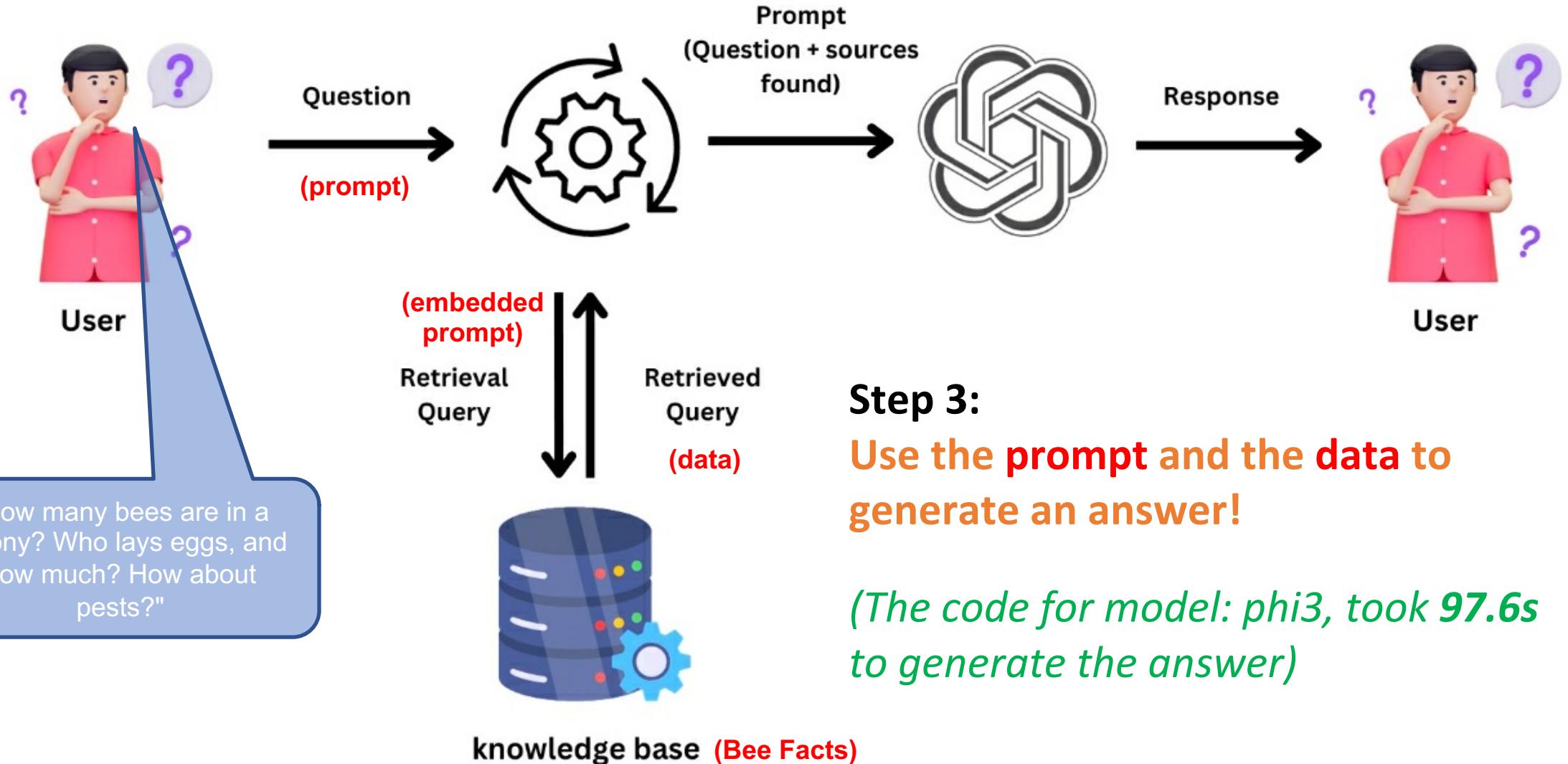
```
1 # Step 1: Generate embeddings (index)
2
3 import ollama
4 import chromadb
5
6
7 EMB_MODEL = "all-minilm" #nomic-embed-text #mxbai-embed-large
8
9 documents = [
10     "Bee-keeping, also known as apiculture, involves the maintenance of bee colonies, typically in hives, by humans.",
11     "The most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



A screenshot of a code editor window titled "ppt.py". The window shows the following Python code:

```
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 is a script for retrieving the most relevant document from a collection based on a given prompt. It uses the ollama library to generate embeddings and the collection library to query the documents.



Step 3:
Use the prompt and the data to generate an answer!

(The code for model: phi3, took 97.6s to generate the answer)

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

- Open Files:** rag_test.py, ppt.py
- Title Bar:** ppt.py, UNREGISTERED
- Ppt.py Content:**

```
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:** Line 16, Column 1, Spaces: 2, Python

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.

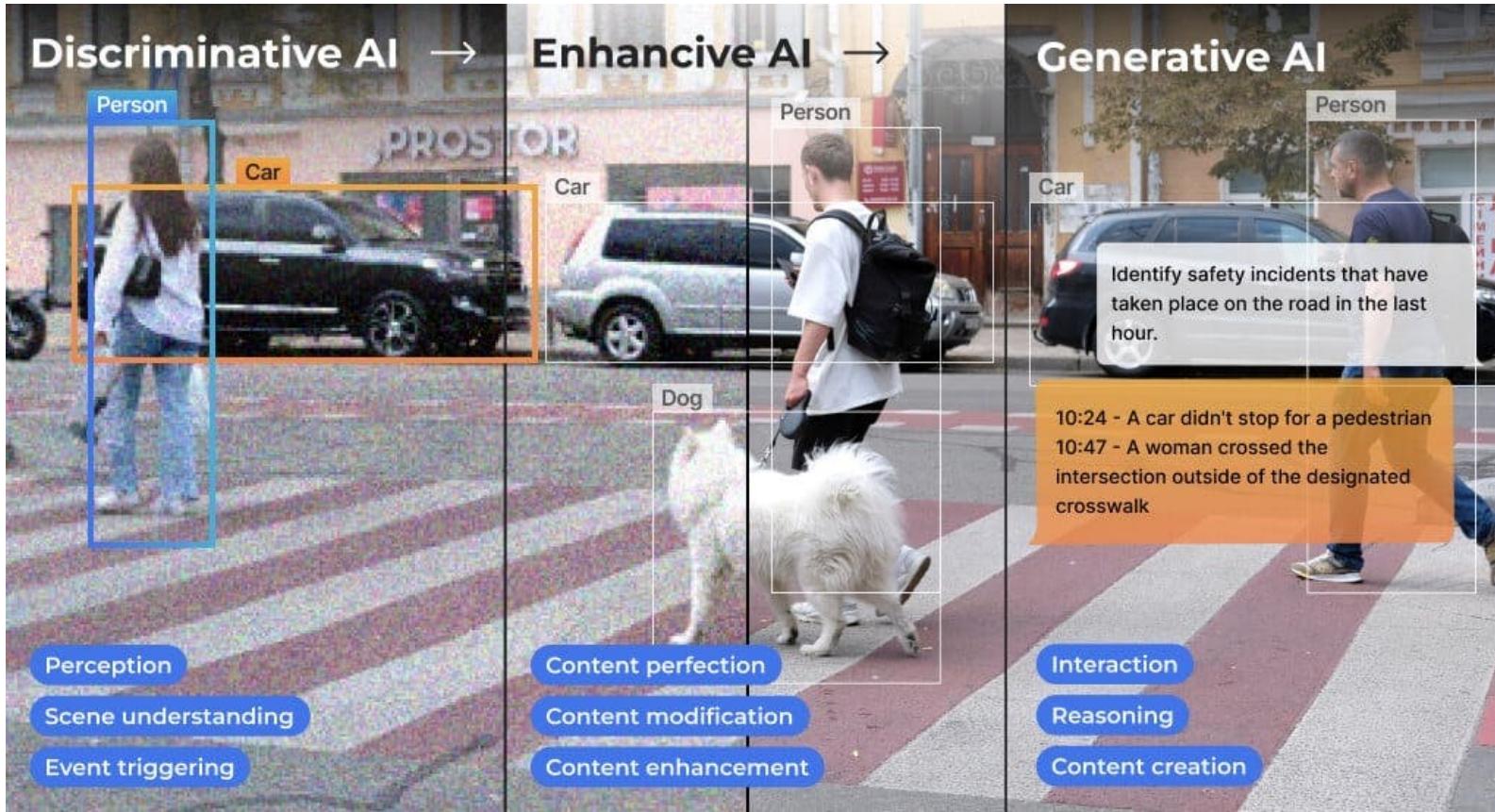
The screenshot shows a Visual Studio Code (VS Code) interface running on a Raspberry Pi. The title bar indicates the file is "rag_test.py - OLLAMA..". The top status bar shows system icons for battery, signal, and temperature (51°). The left sidebar has a "Wastebasket" icon and a tree view of the project 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 ecosystems.",
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

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

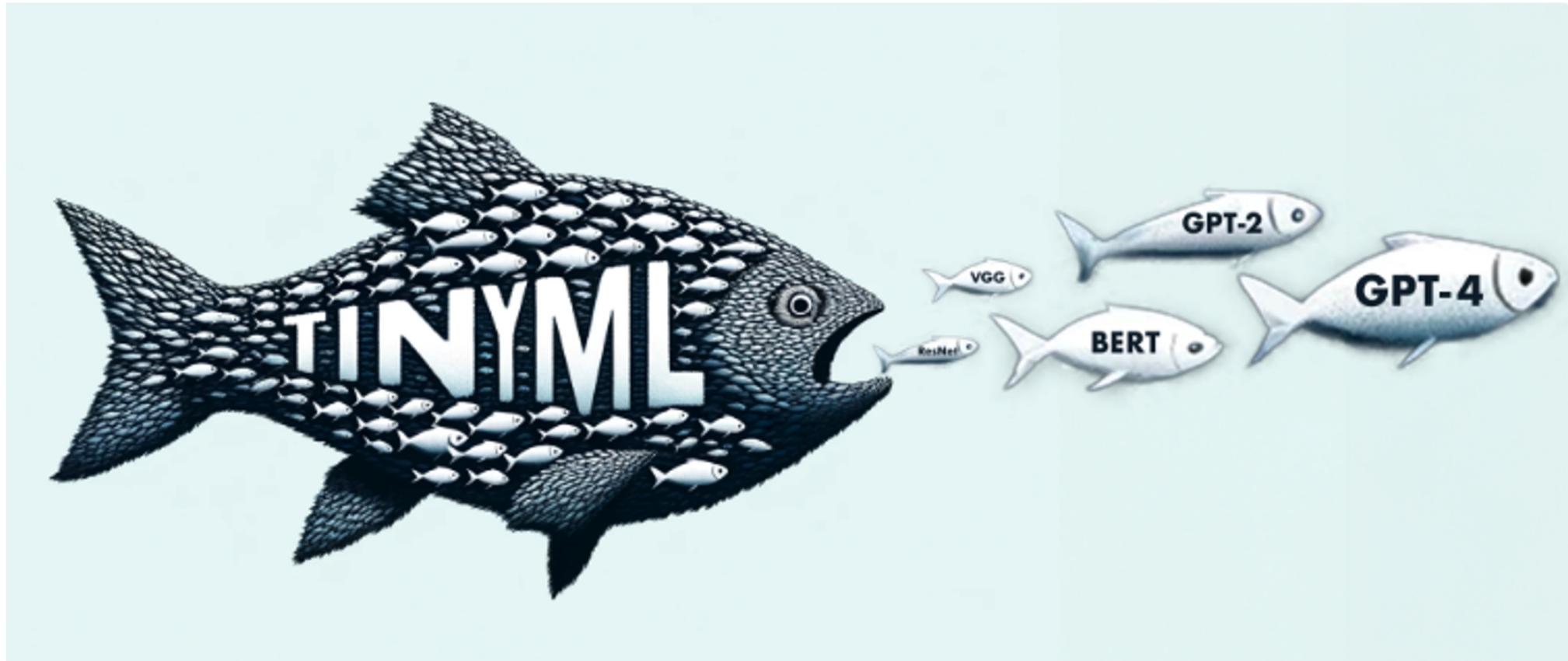
mjrovai@rpi-5:~/Documents/OLLAMA/RAG/RAG_test \$ sudo raspi-config



"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

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



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

Thanks



TINYML4D