



WALC 2024
Applied AI

Large Language Models (LLMs) at the Edge

Prof. Marcelo J. Rovai

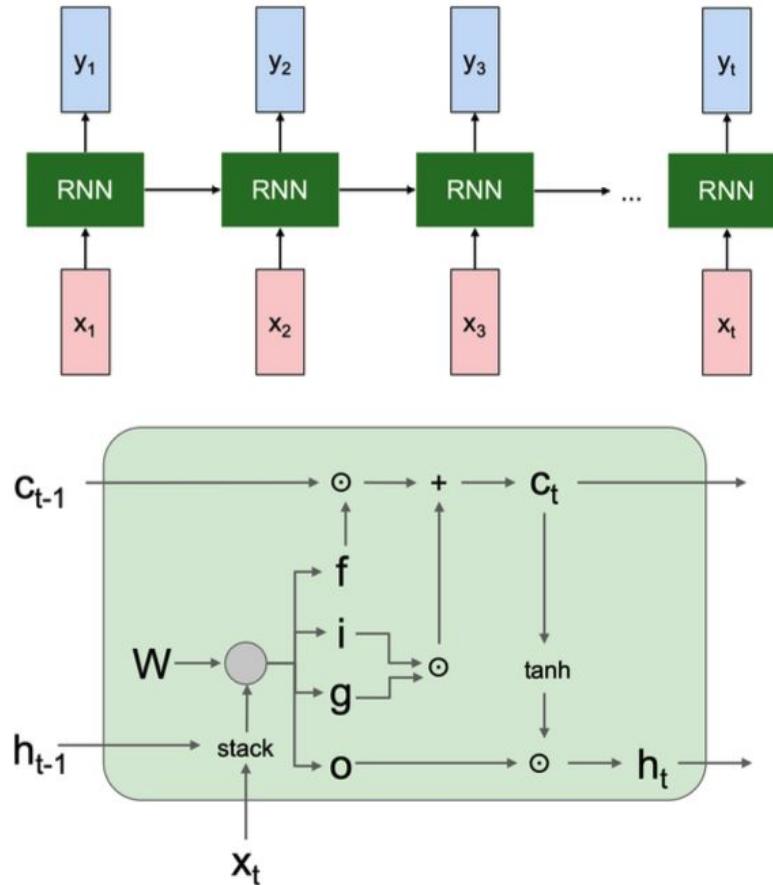
rovai@unifei.edu.br

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

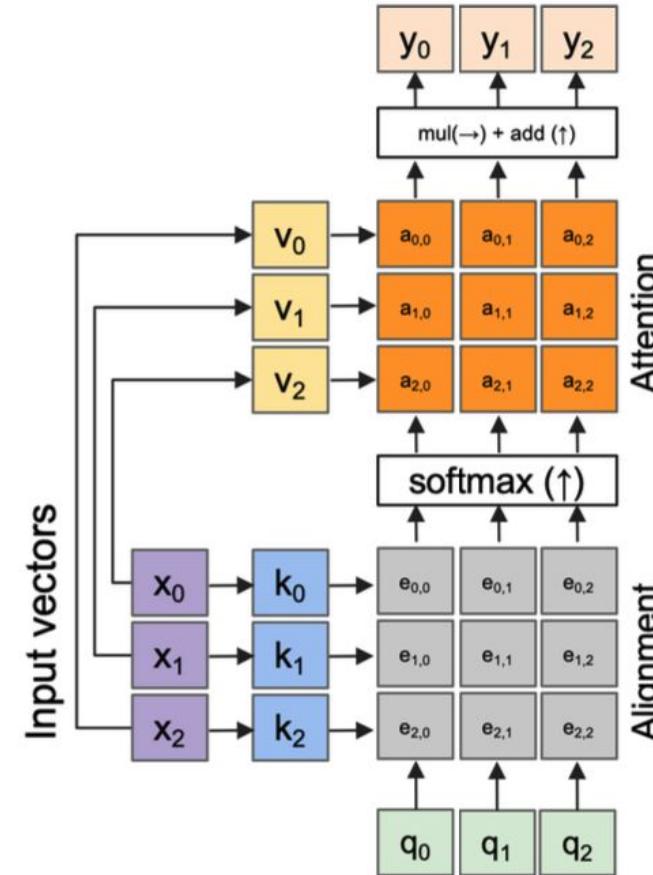


TINYML4D

Recap: Models Beyond DNN and CNN



Recurrent neural network



Attention mechanism / Transformers

Machado de Assis Bot with RNN - GRU

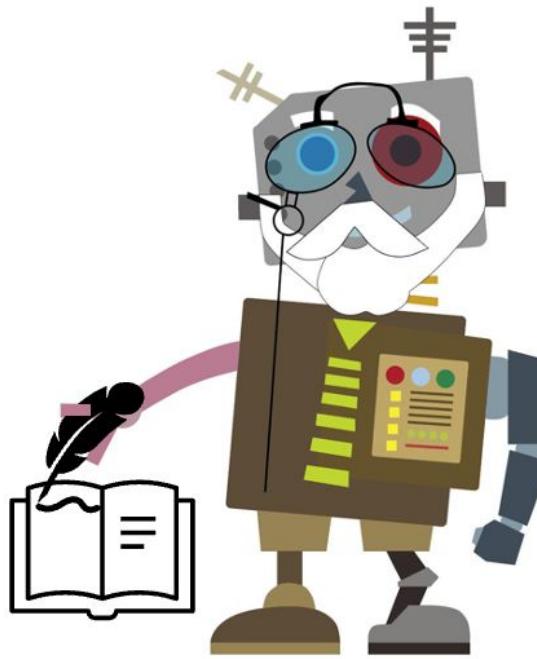


The robot writer model is a **Recurrent Neural network (RNN/GRU)**. The model, with 4M parameters, was trained with a **150-characters sequence** from seven of his books: *Memorias Posthumas de Braz Cubas*, *Dom Casmurro*, *Quincas Borba*, *Papeis Avulsos*, *A Mão e a Luva*, *Esaú e Jacob*, and *Memorial de Ayres*.

Model: "sequential"

Layer (type)	Output Shape	Param #
embedding (Embedding)	(1, 150, 256)	29,952
gru (GRU)	(1, 150, 1024)	3,938,304
dense (Dense)	(1, 150, 117)	119,925

Total params: 4,088,181 (15.60 MB)
Trainable params: 4,088,181 (15.60 MB)
Non-trainable params: 0 (0.00 B)



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.*

Generative AI (GenAI)

Generative AI is an artificial intelligence system capable of creating new, original content across various mediums such as **text, images, audio, and video**. These systems learn patterns from existing data and use that knowledge to generate novel outputs that didn't previously exist.

Large Language Models (LLMs), Small Language Models (SLMs), and **multimodal models** can all be considered types of GenAI when used for generative tasks.

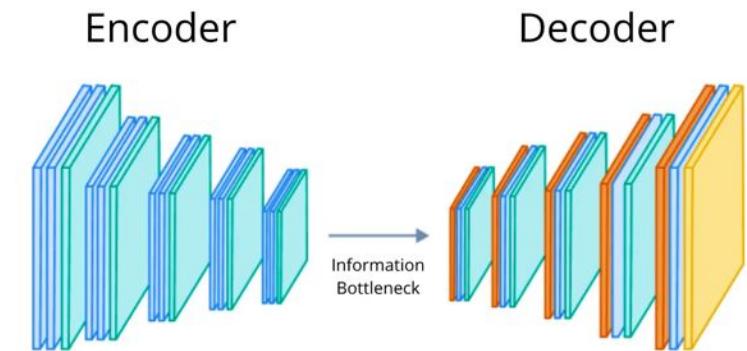
LLM / SLM

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).

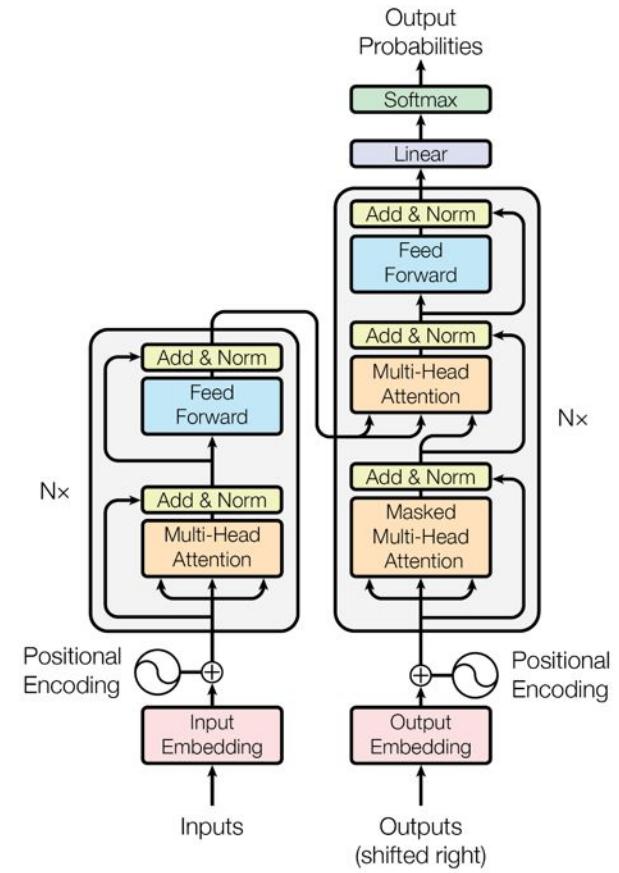
Deep Learning models (or artificial neural networks)

- **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)**.



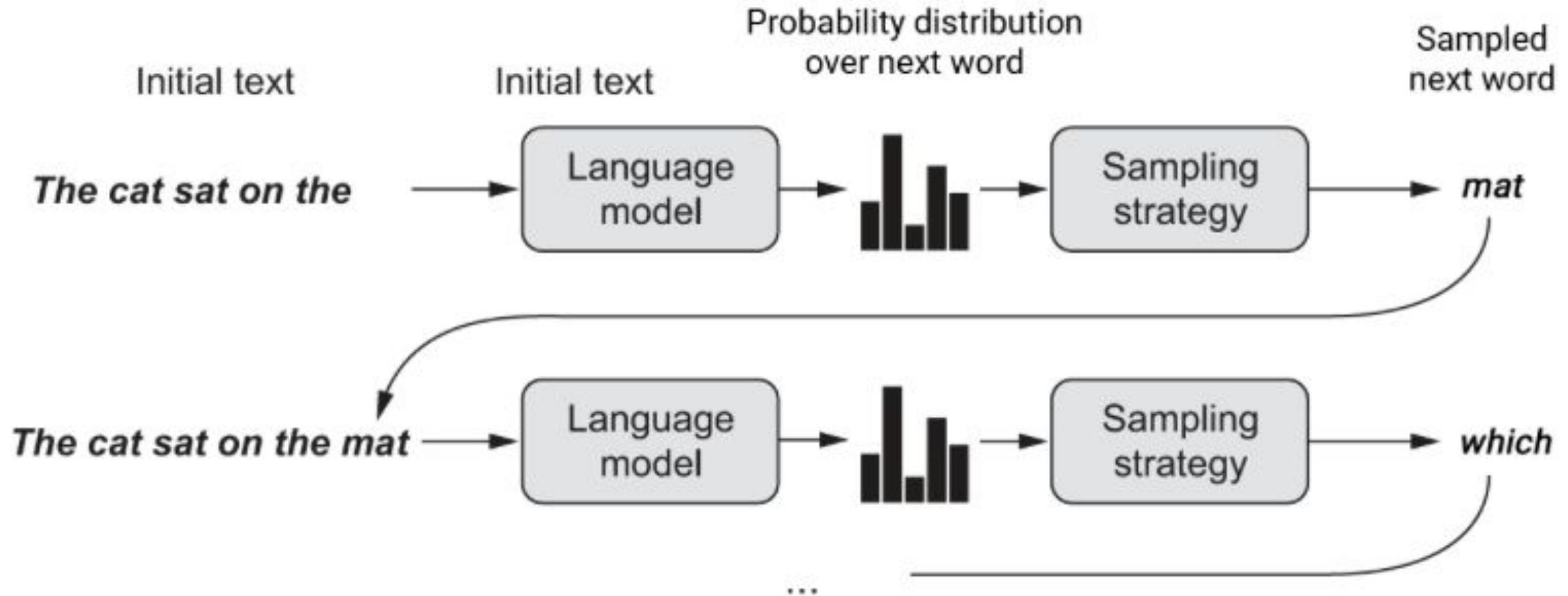
LLM/SLM – Large /Small Language Model

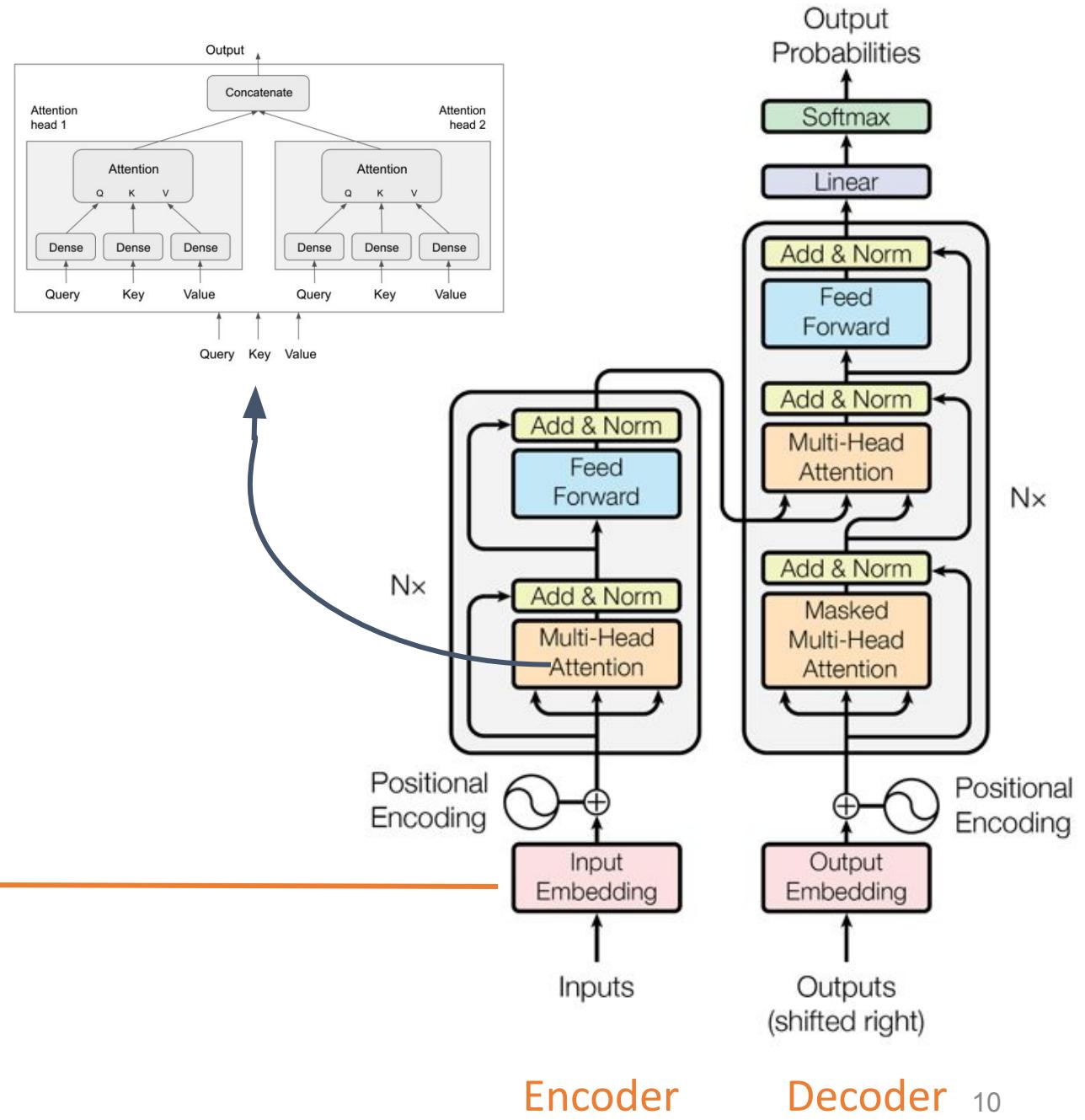
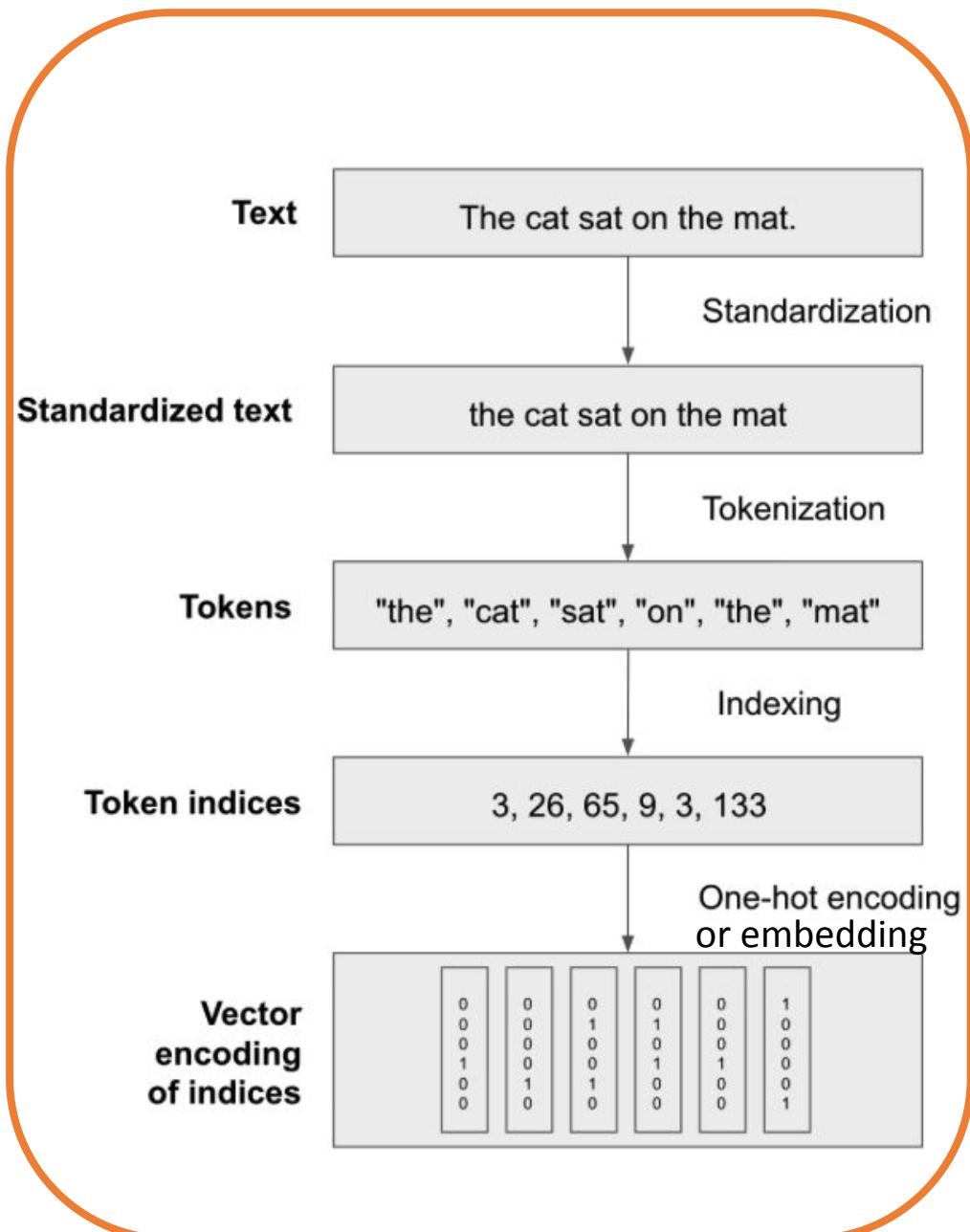
Large Language Models (LLMs) and SLMs are advanced neural networks based on the **Transformer architecture** that excel in understanding and generating human language. They represent a significant evolution from earlier sequence-based models like **RNNs**, which surpass them in handling long-range dependencies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



The Illustrated Transformer

LLM/SLM – Large /Small Language Model

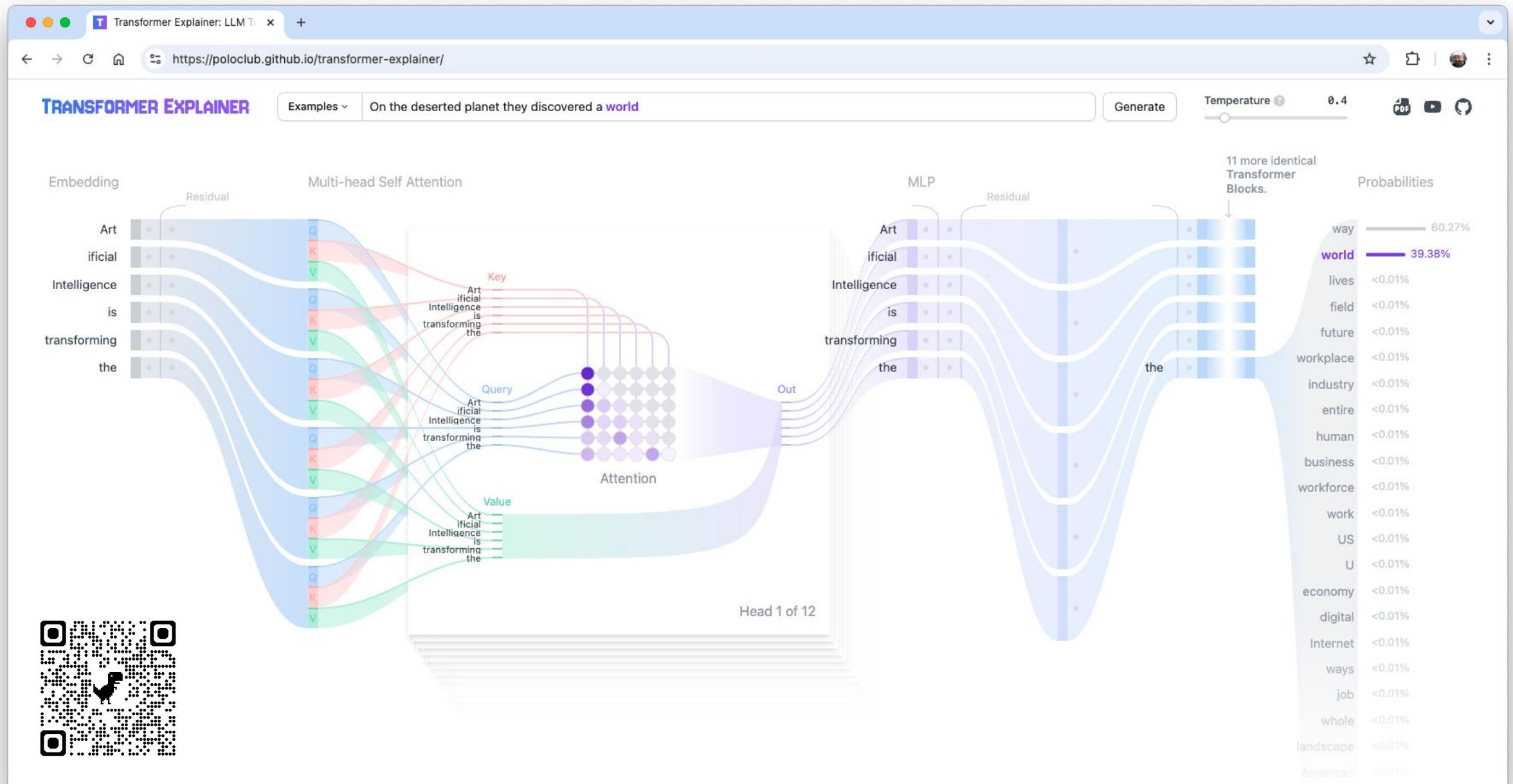




LLM/SLM – Large /Small Language Model

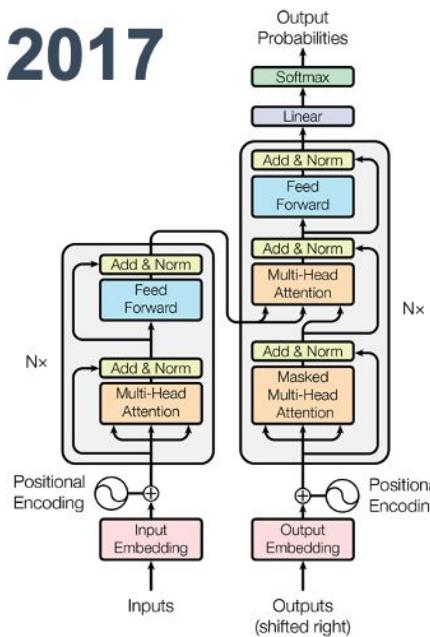


How large language models work



Transformers to LLMs and SLMs

2017

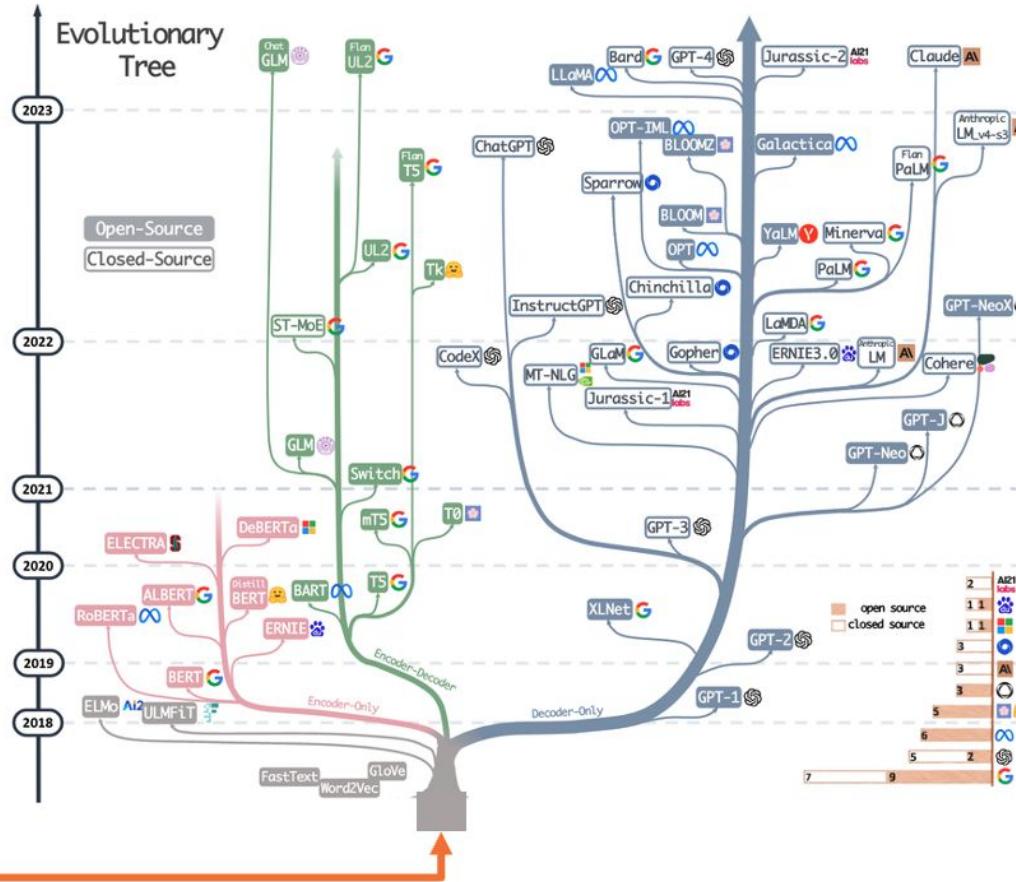


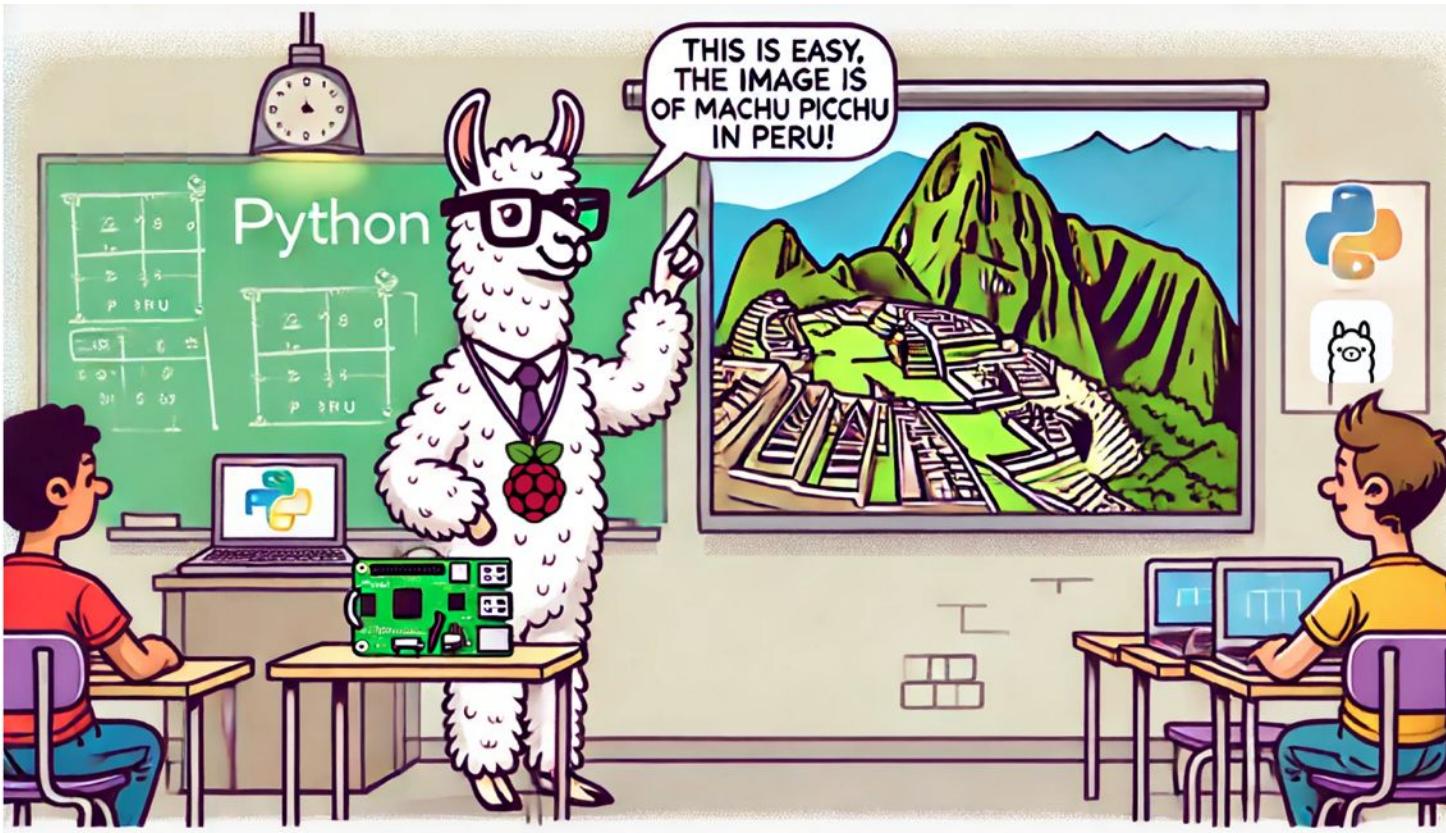
2024



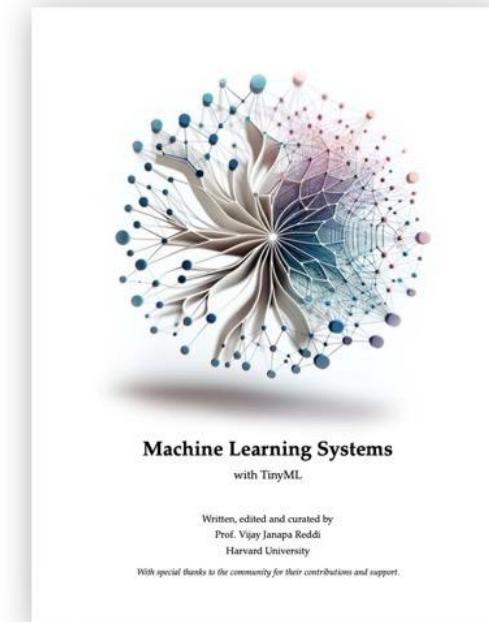
Open

Closed





Running Large Language Models on Raspberry Pi at the Edge



A screenshot of a web browser window showing the Ollama website at <https://ollama.com>. The page features a large, friendly llama logo at the top left. Below the logo, there's a headline: "Get up and running with large language models." followed by a subtext: "Run Llama 3.2, Phi 3, Mistral, Gemma 2, and other models. Customize and create your own." A prominent "Download ↓" button is centered below the subtext. At the bottom, it says "Available for macOS, Linux, and Windows". The browser interface includes standard navigation buttons, a search bar, and a menu icon.

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

```
[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:~\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)



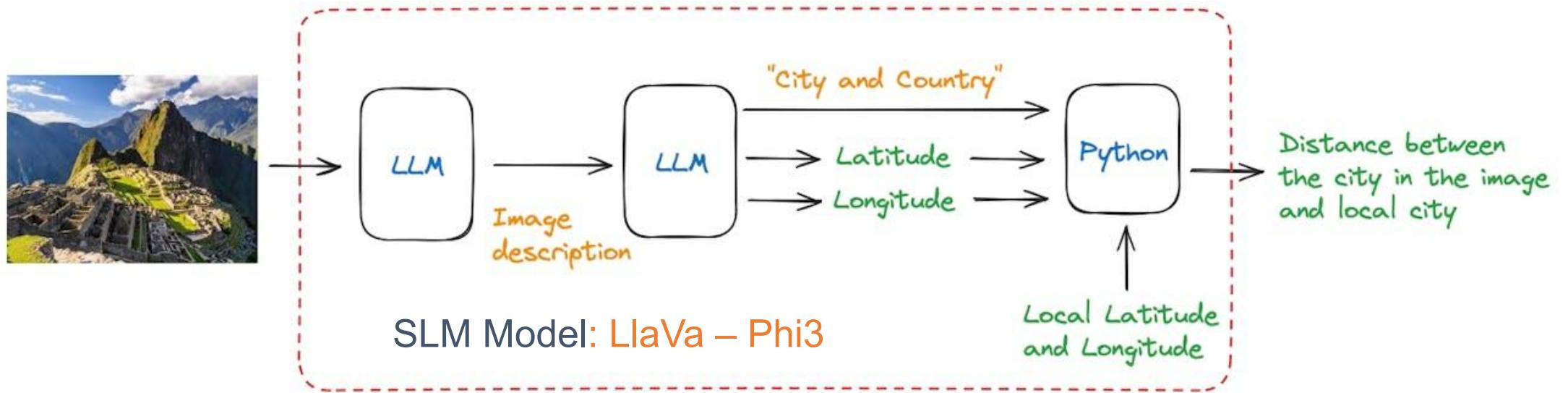
```
mjrovai@rpi-5:~/Documents/OLLAMA\n\nhelp\nroute.\n\n/Documents/OLLAMA $\n/Documents/OLLAMA $ python calc_distance_image.py /\n/home/mjrovai/Documents/OLLAMA/image_test_1.jpg\n\nThe image shows Paris, with lat:48.86 and long: 2.35, located in\nFrance and about 11,630 kilometers away from Santiago, 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)

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.

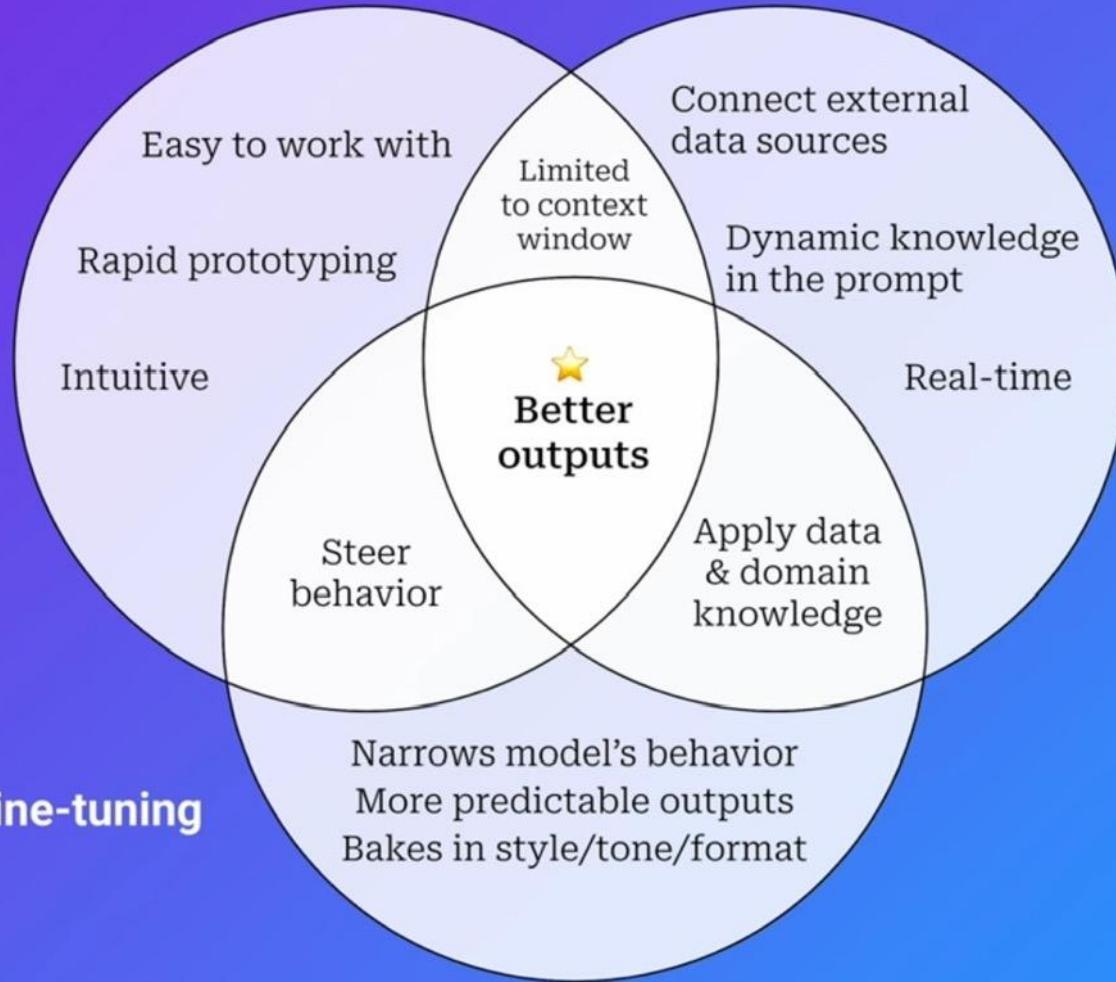
Comparison of Techniques

Prompt Engineering



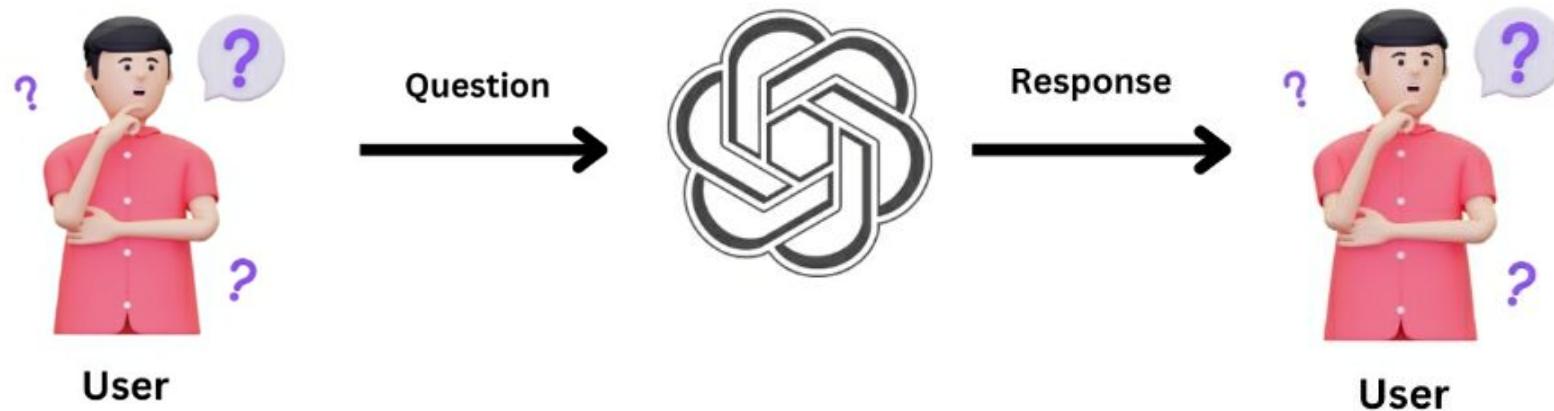
Fine-tuning

RAG

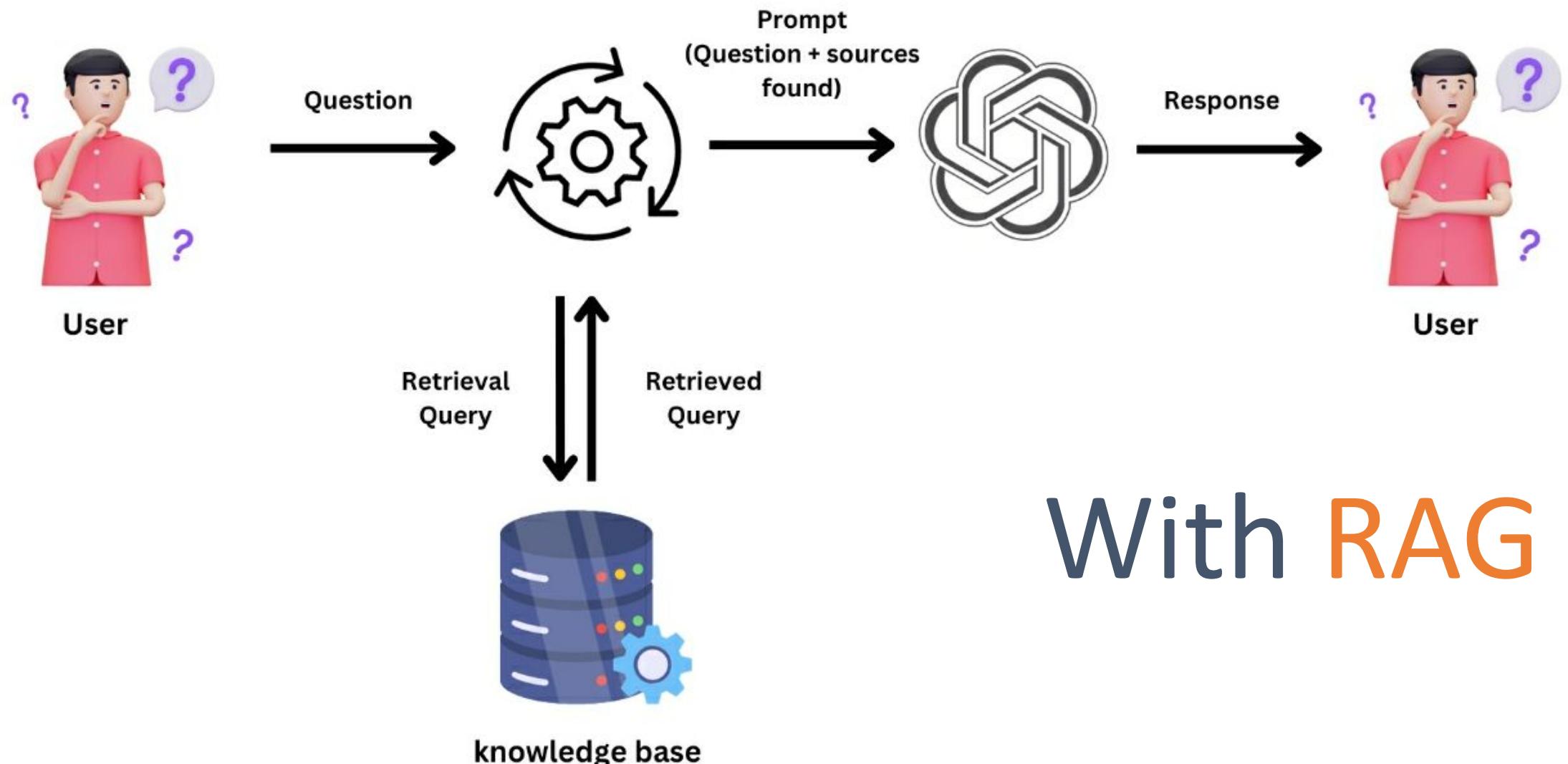


Retrieval-Augmented Generation (RAG)

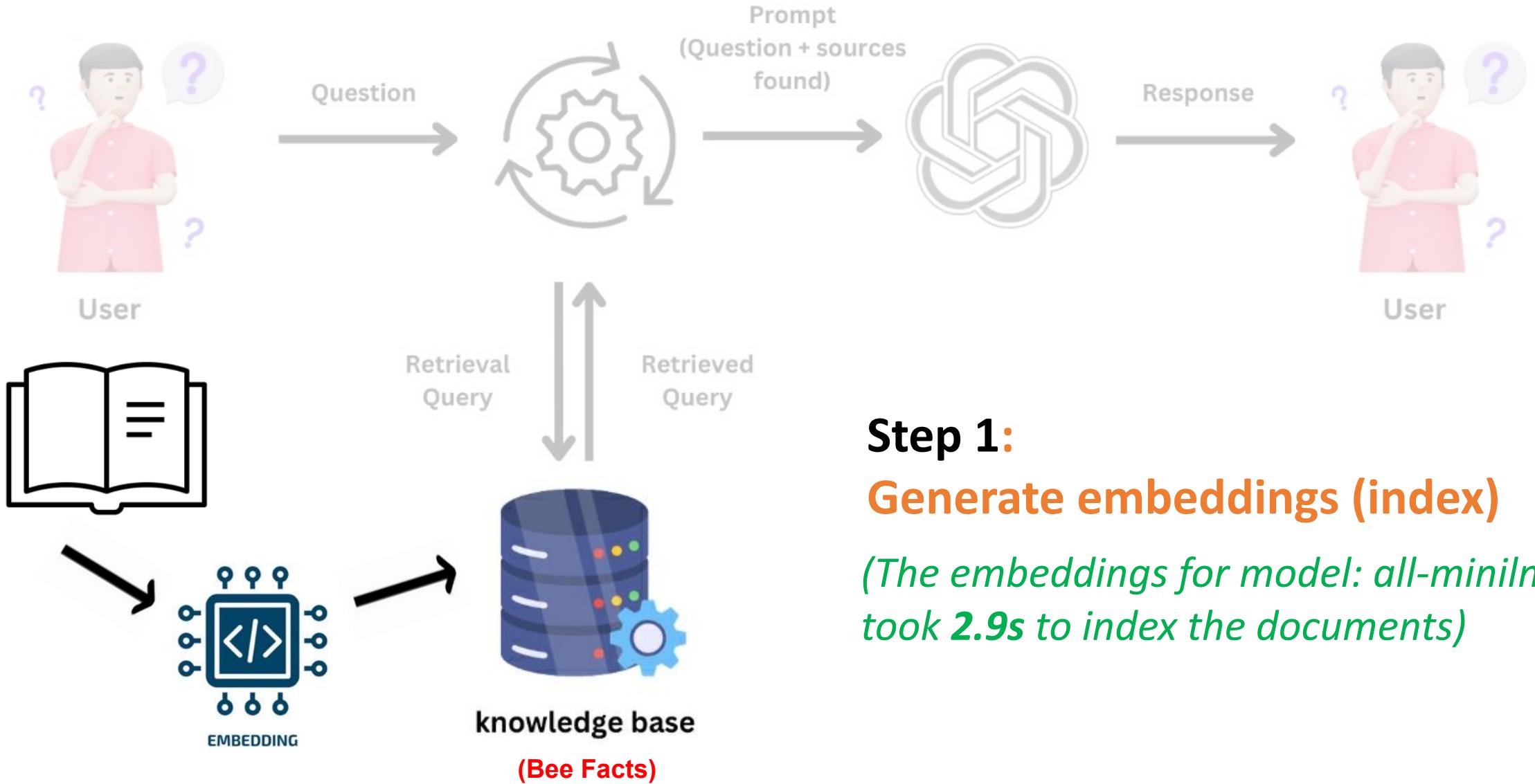
“A method created by the FAIR team at Meta to enhance the accuracy of Large Language Models (LLMs) and reduce false information or “hallucinations.”



Usual Prompt



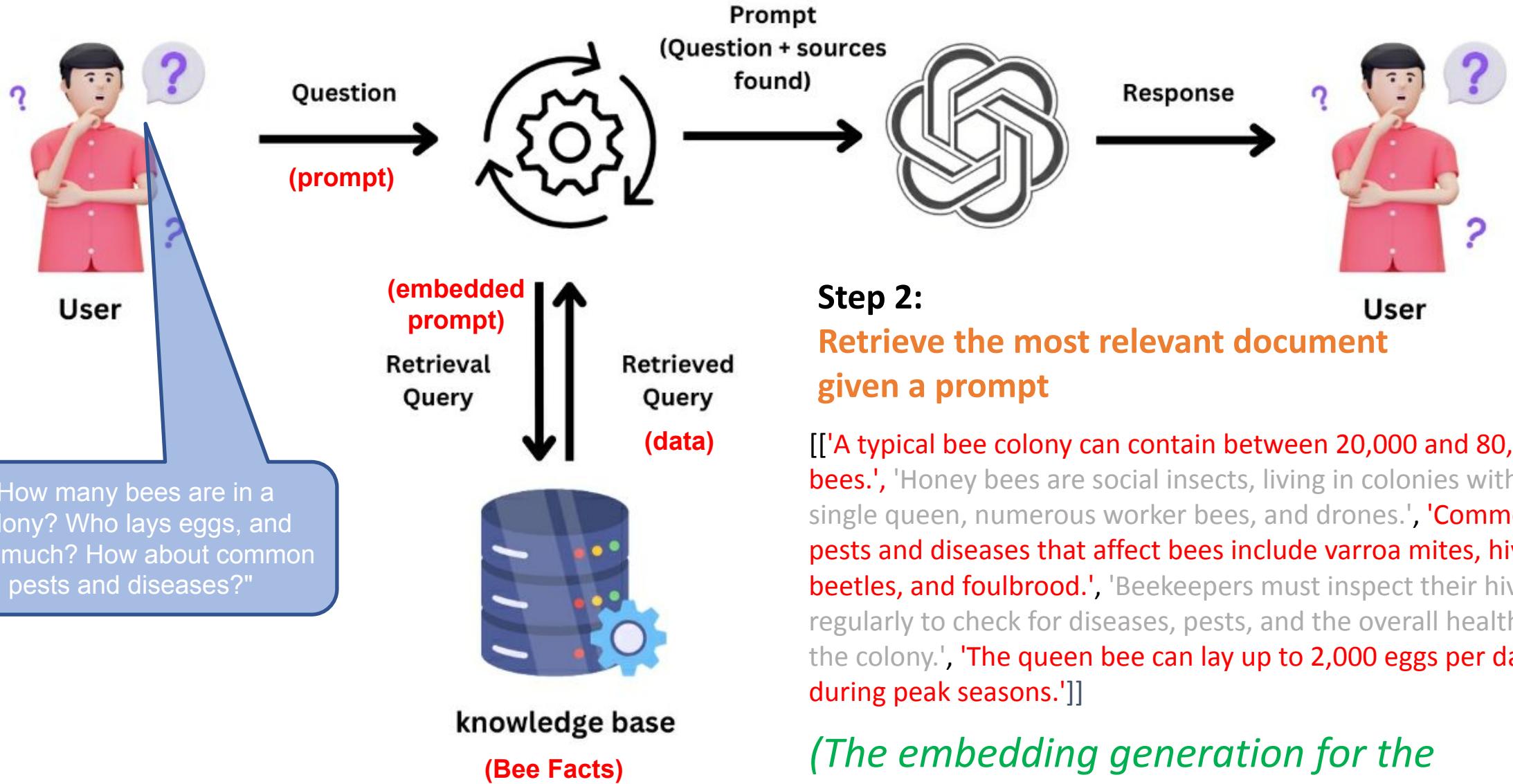
With RAG



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

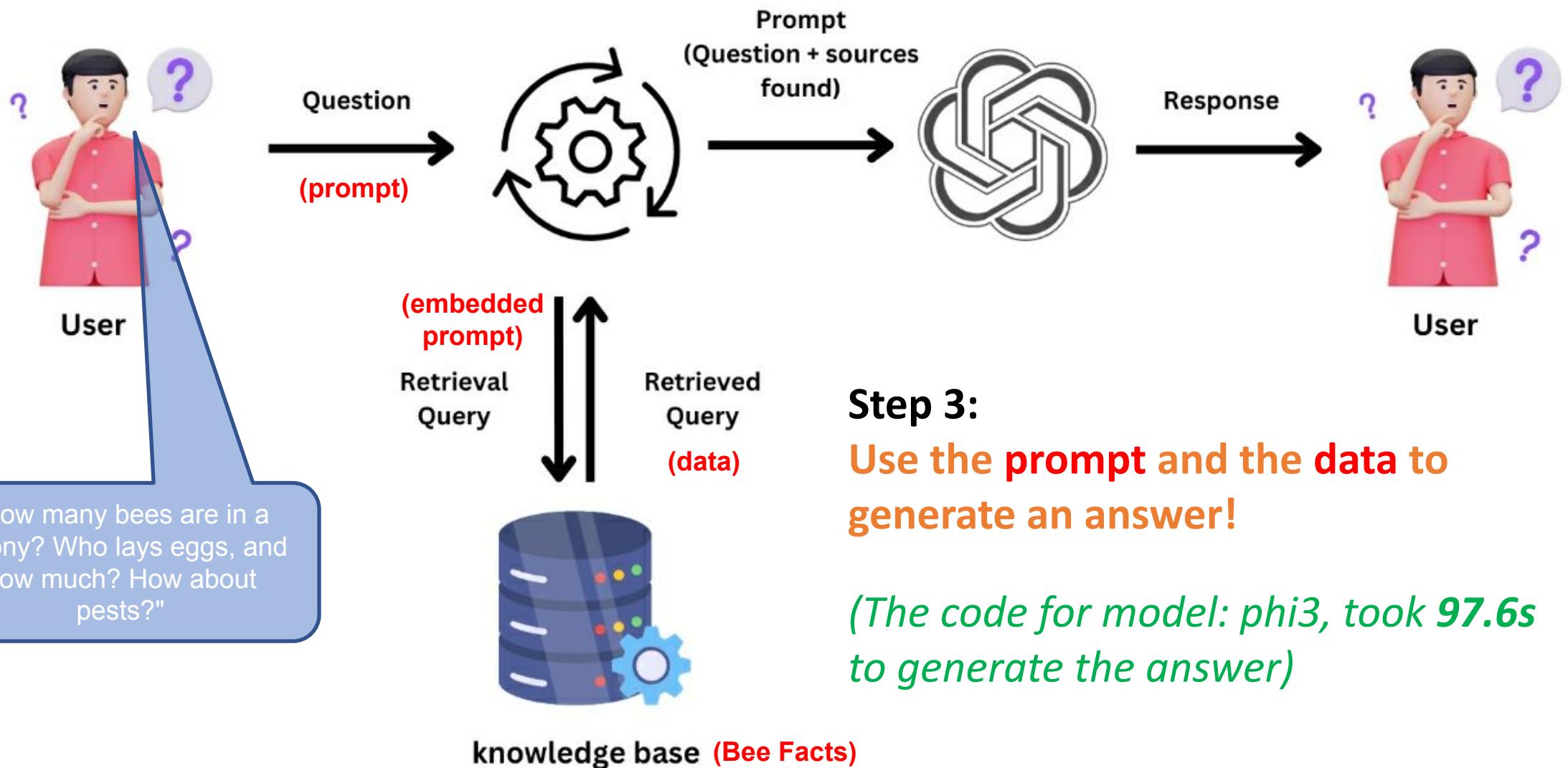
The screenshot shows a code editor window with two tabs: `rag_test.py` and `ppt.py`. The `ppt.py` tab is active, displaying the following Python code:

```
OPEN FILES ◀ ▶ ppt.py ◀ ▶ ppt.py ◀ ▶ ppt.py  
UNREGISTERED  
Line 45, Column 1  
Spaces: 2 Python  
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  
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 ]
```



A screenshot of a code editor window titled "ppt.py". The window shows two tabs: "rag_test.py" and "ppt.py", with "ppt.py" being the active tab. The code in "ppt.py" is a Python script for a Retrieval-Augmented Generation (RAG) system. It starts by importing "os" and "MongoClient". It defines a function "get_answer" that takes a "prompt" and an optional "model". Inside the function, it generates an embedding for the prompt using "ollama.embeddings()", specifying the prompt and the model as "EMB_MODEL". It then queries a MongoDB collection named "data" with the generated embedding as the query and a limit of 5 results. Finally, it returns the documents from the query results. The code editor interface includes a sidebar for "OPEN FILES" with "rag_test.py" and "ppt.py" listed, and a status bar at the bottom indicating "Line 3, Column 1", "Spaces: 2", and "Python".

```
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
```



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:

- 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 (ppt.py):**

```
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 status bar at the bottom 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..". 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".

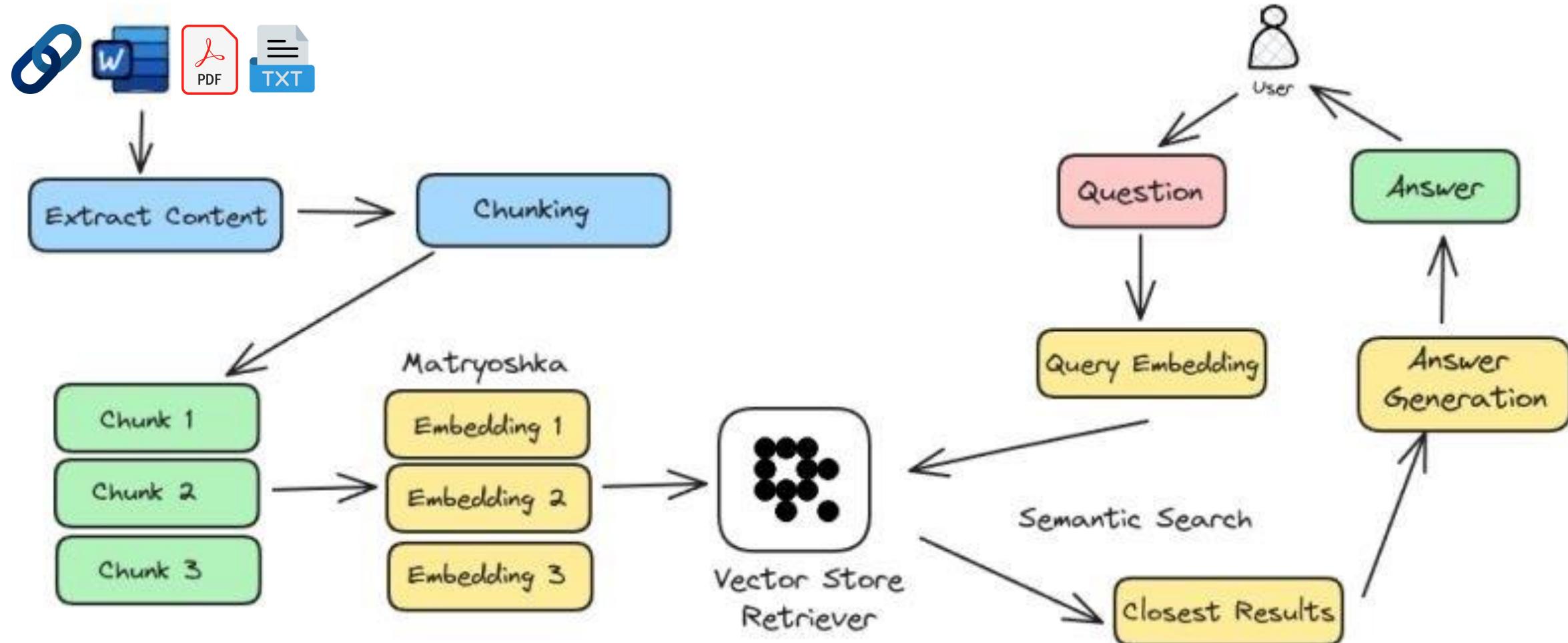
```
rag_test.py - OLLAMA - Visual Studio Code
RAG > RAG_test > rag_test.py > ...
8
9 import ollama
10 import chromadb
11 import time
12
13 start_time = time.perf_counter() # Start timing
14 EMB_MODEL = "all-minilm" #nomic-embed-text" #mbai-embed-large"
15 MODEL = "phi3"
16
17 documents = [
18     "Bee-keeping, 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.",
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.",
36     "Sustainable bee-keeping practices are important for maintaining healthy bee populations and ecosystems."
```

PROBLEMS 1 OUTPUT DEBUG CONSOLE TERMINAL PORTS

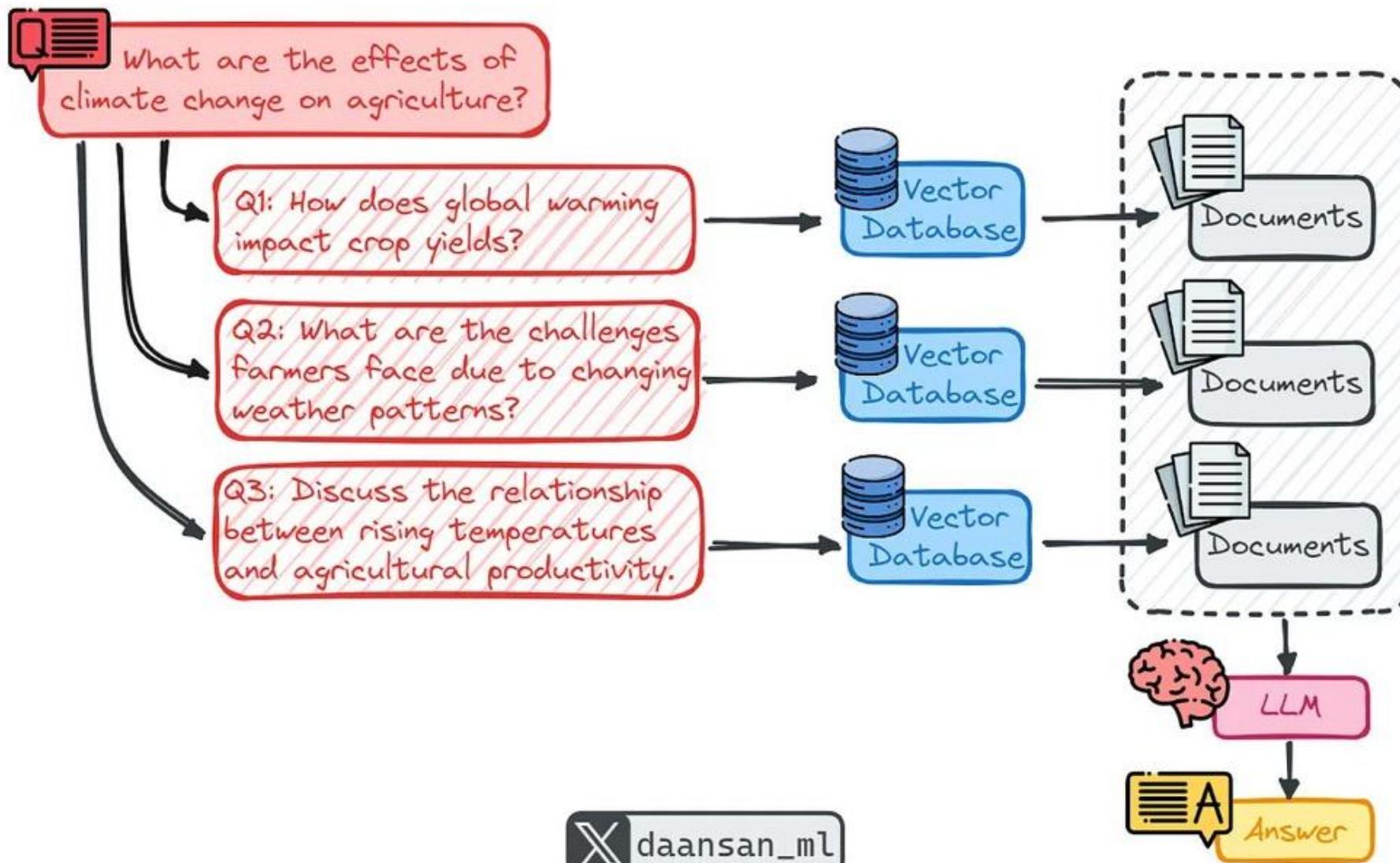
[INFO] ==> The code for model: phi3, took 97.6s to generate the answer.

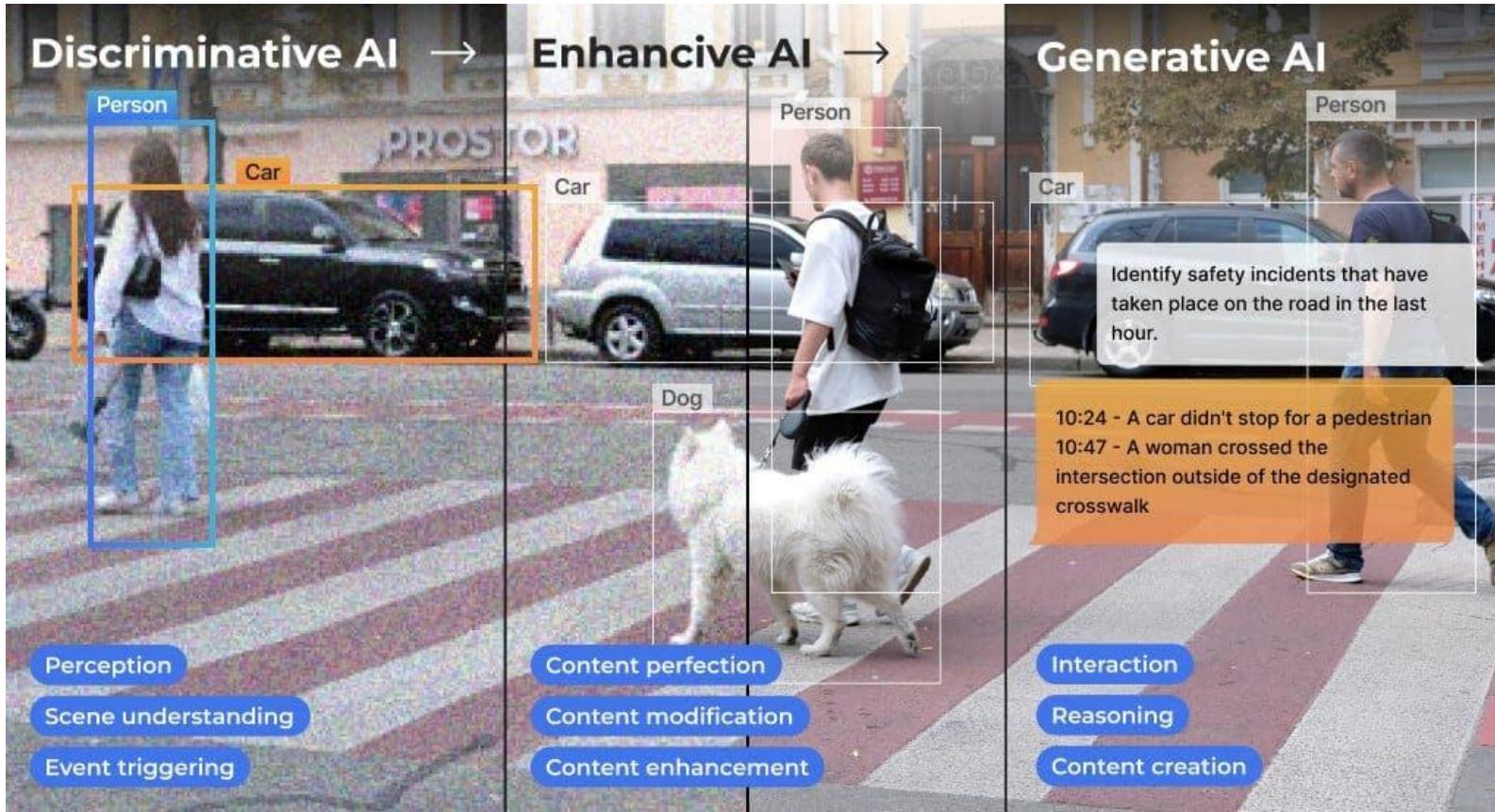
mjrovai@rpi-5:~/Documents/OLLAMA/RAG/RAG_test \$ sudo raspi-config

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

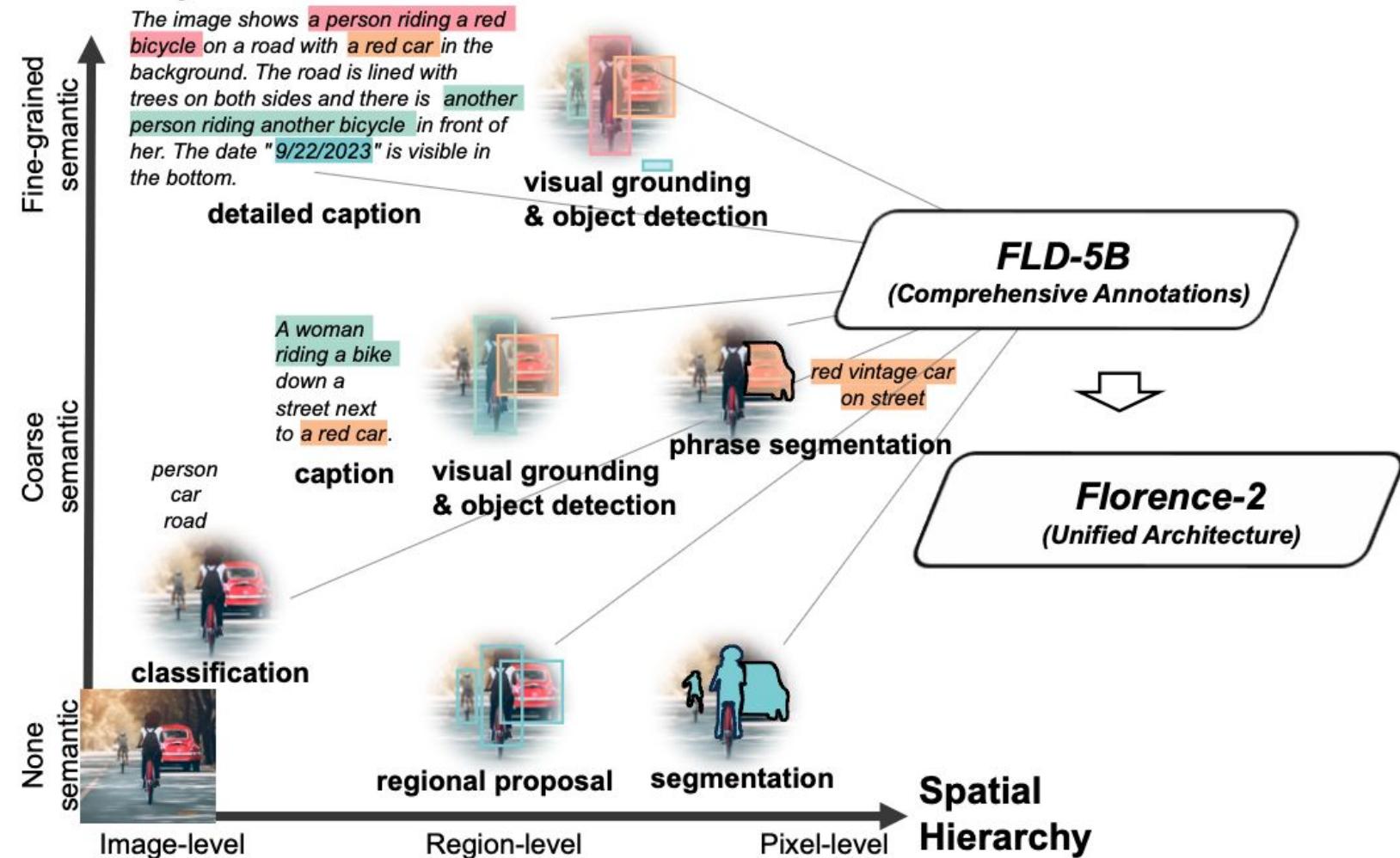
Florence-2

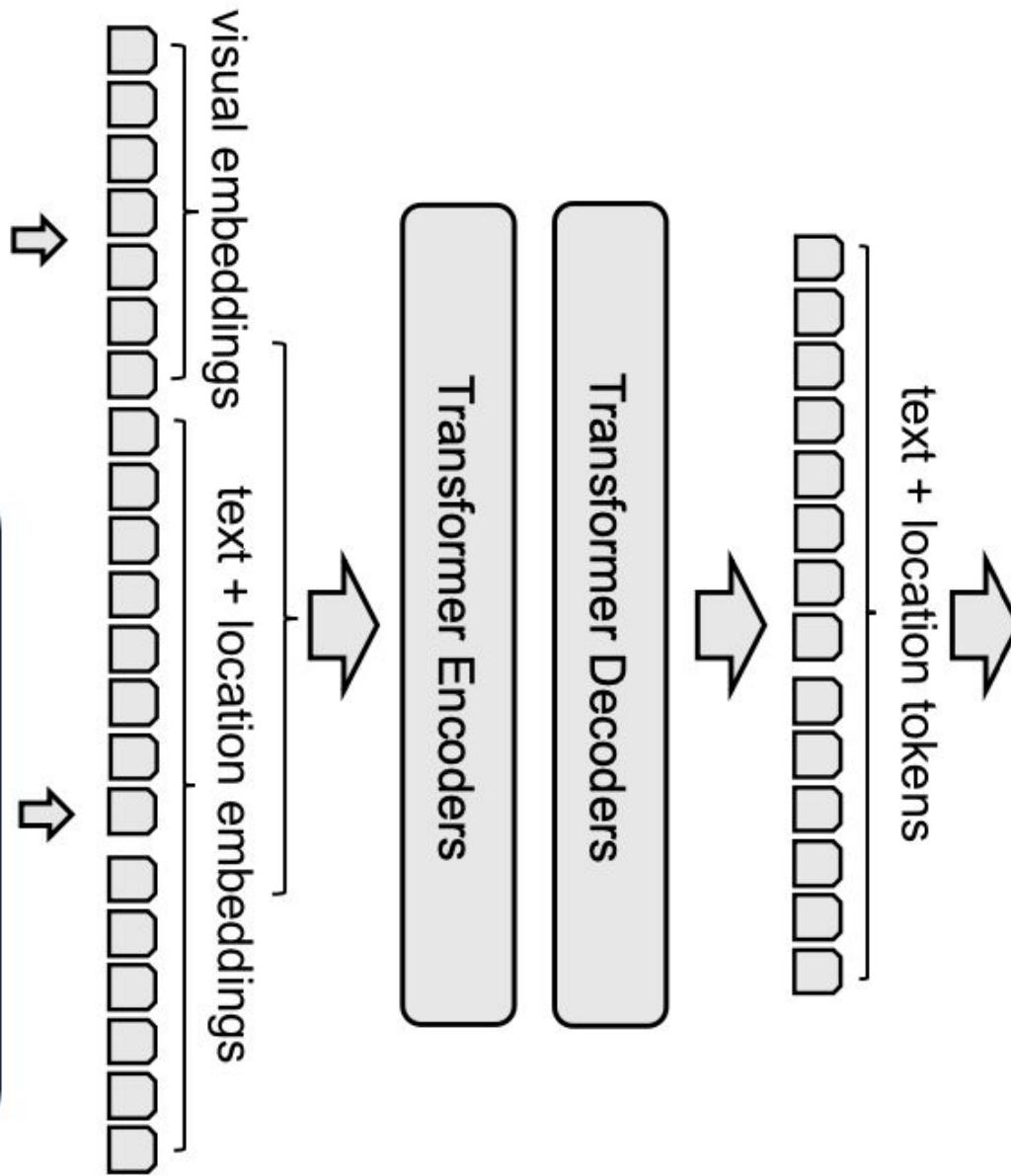
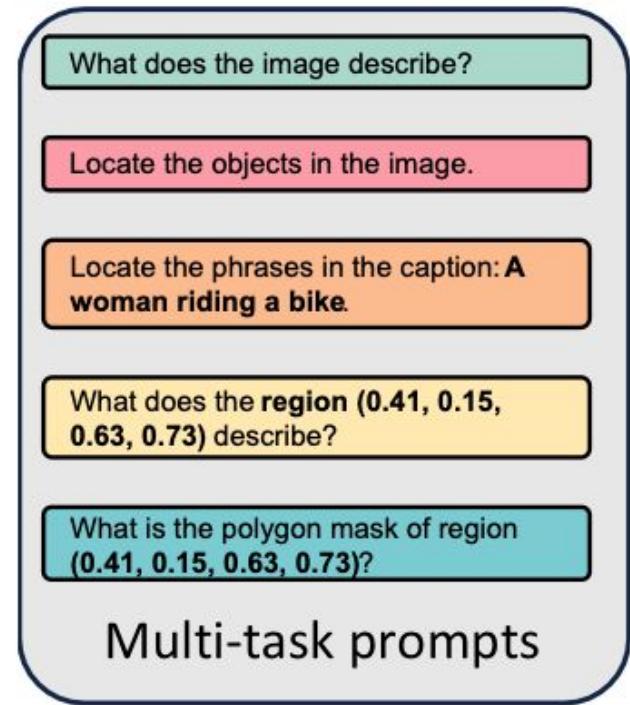
Advancing a Unified Representation for a Variety of Vision Tasks



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

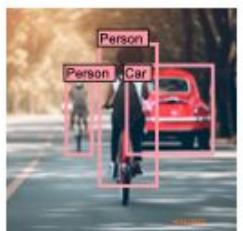
Semantic Granularity





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.

person (0.41, 0.15, 0.63, 0.73)
... car (0.58, 0.26, 0.89, 0.61)



A women riding a bike (0.41, 0.15, 0.63, 0.73)



person riding red bicycle on road

(0.48, 0.19, 0.48, 0.18, 0.49, 0.17, ...)



microsoft/BitNet

Official inference framework for 1-bit LLMs

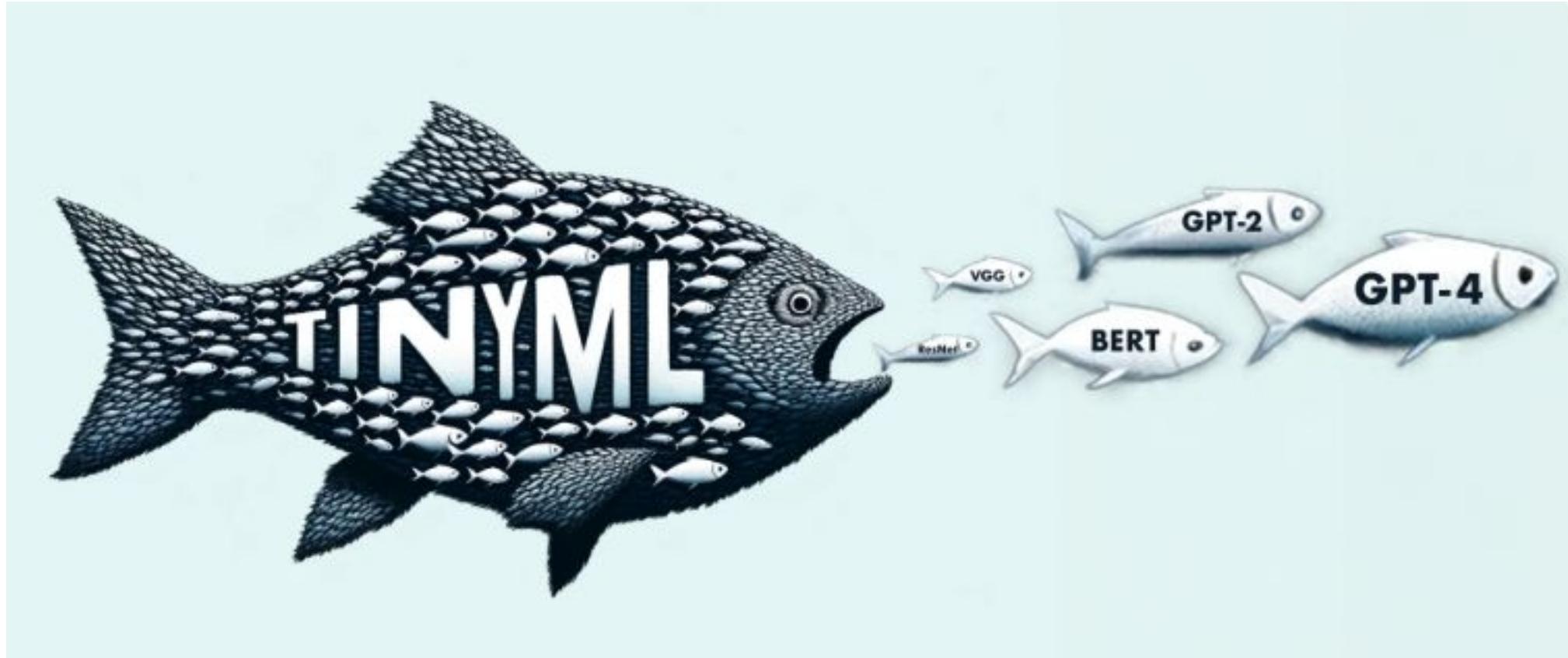


Bitnet.cpp employs one-bit quantization, representing values with a ternary system **(+1, -1, 0)**. This approach simplifies calculations by replacing complex multiplications with additions and subtractions, eliminating the need for GPUs.

- 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?

