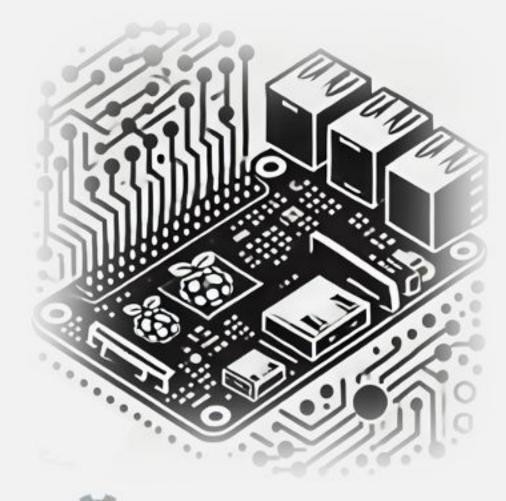
IESTI05 – Edge Al

Machine Learning
System Engineering

20. SLMs: Optimization Techniques- RAG

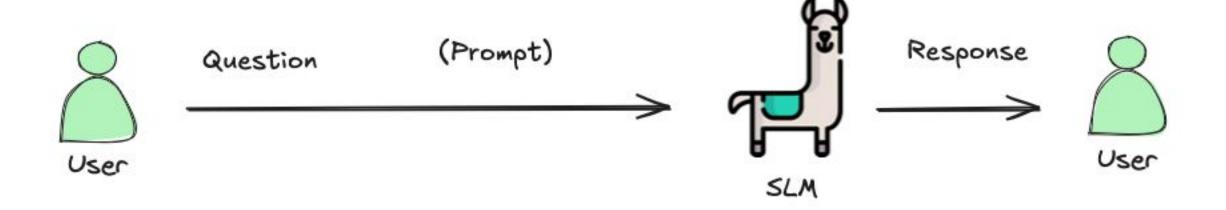






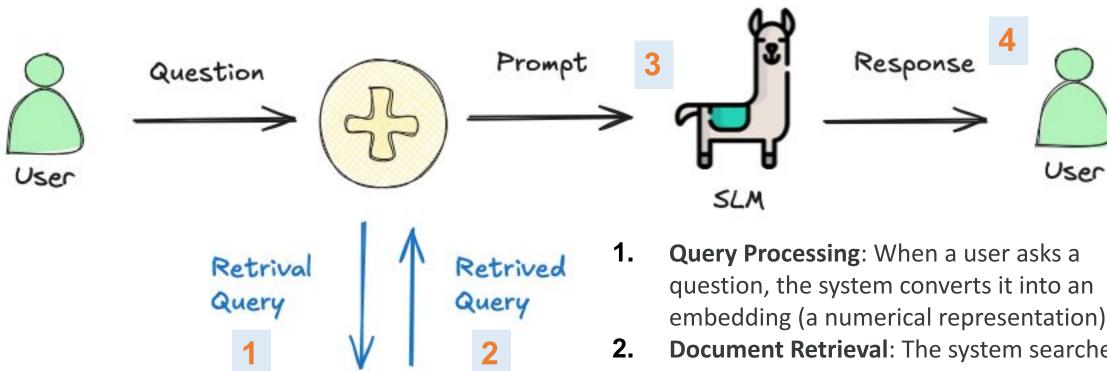
Retrieval-Augmented Generation (RAG)

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





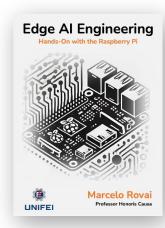
Usual Prompt



- With RAG
- Knowledge Base

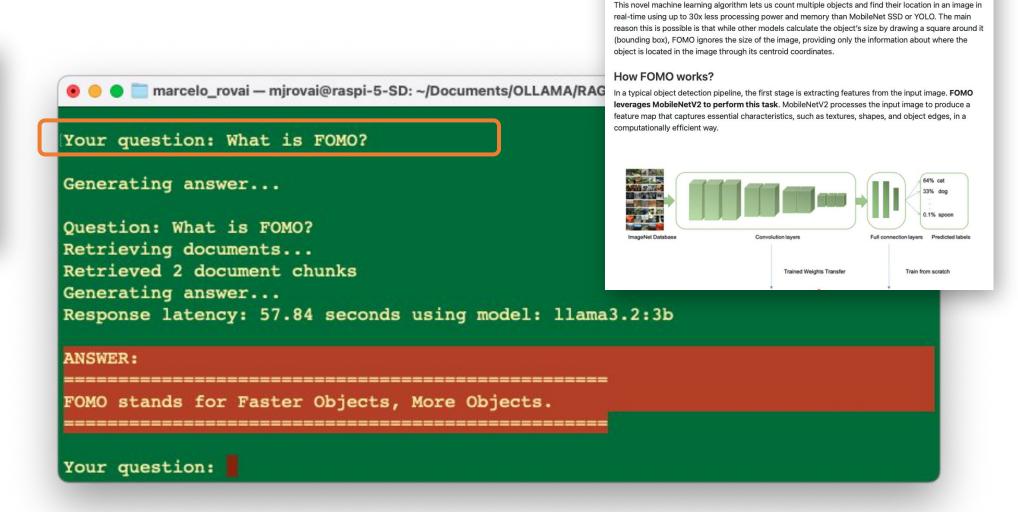
- question, the system converts it into an embedding (a numerical representation).
- Document Retrieval: The system searches a knowledge base for documents with similar embeddings.
- **Context Enhancement**: Relevant documents are retrieved and combined with the original query.
- **Generation**: The SLM generates a response using both the query and the retrieved context.

With RAG





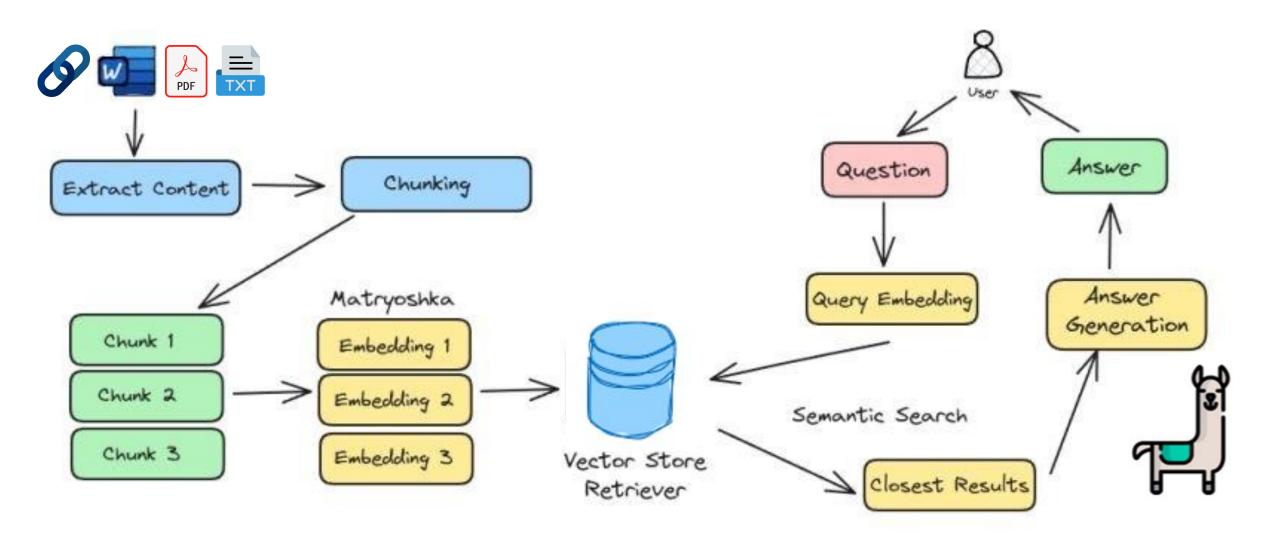




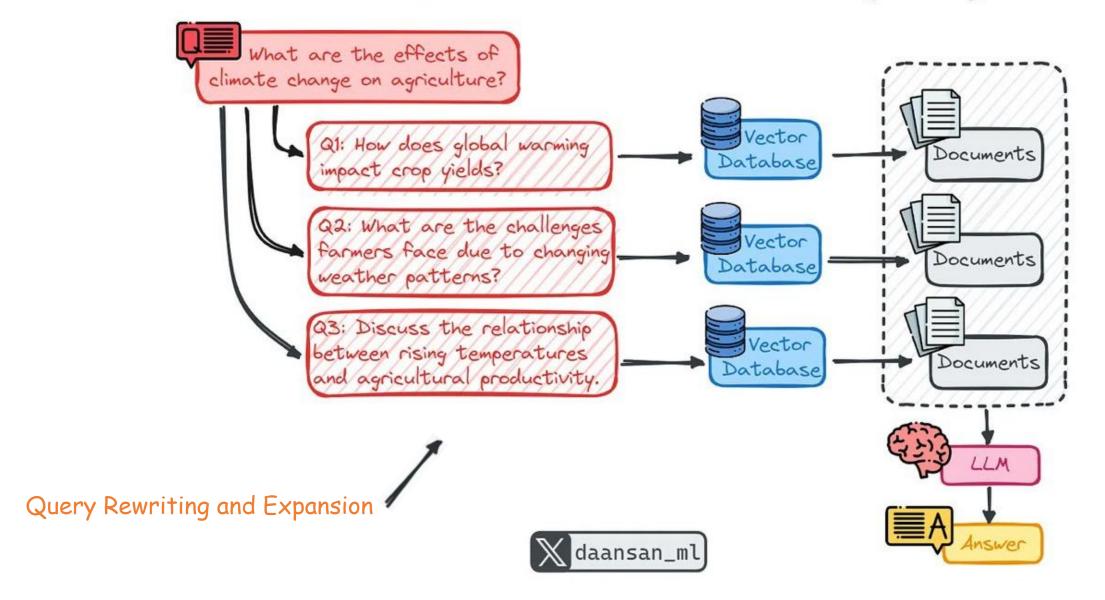
Training a FOMO Model at Edge Impulse Studio &

The inference with the SSD MobileNet model worked well, but the latency was significantly high. The inference varied from 0.5 to 1.3 seconds on a Raspi-Zero, which means around or less than 1 EPS (1 frame per second). One alternative to speed up the process is to use FOMO (Faster Objects, More Objects).

RAG: Simple Query ("Naive RAG")



Advanced RAG: Multi Query



Simple RAG example



40-RAG-simple-bee.ipynb



Implementing a Naive RAG System

1. Creating the Vector Database (10-Create-Persistent-Vector-Database.py).

This script builds a knowledge base by:

- Loading documents from PDFs and URLs
- Splitting them into manageable chunks
- Creating embeddings for each chunk
- Storing these embeddings in a vector database (Chroma)



- Querying the Database (20-Query-the-Persistent-RAG-Database.py).
 This script:
 - Loads the saved vector database
 - Accepts user queries
 - Retrieves relevant documents based on query similarity
 - Combines documents with the query in a prompt
 - Generates a response using the SLM



Book: Advanced EdgeAI: RAG

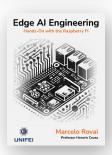
Naive RAG Setup

```
Inside the env:
source ~/ollama/bin/activate
Install:
  pip install -U 'langchain-chroma'
 pip install -U langchain
 pip install -U langchain-community
 pip install -U langchain-ollama
 pip install -U langchain-text-splitter
 pip install -U langchain-community pypdf
 pip install tiktoken
  pip install -U langsmith
```

Create Knowledge Base







```
import warnings
import os
from langchain_text_splitters import RecursiveCharacterTextSplitter
from langchain_community.document_loaders import WebBaseLoader
from langchain community.document loaders import PyPDFLoader
from langchain_ollama import OllamaEmbeddings
from langchain_community.vectorstores import Chroma
# Suppress LangSmith warnings
warnings.filterwarnings("ignore",
                        message="API key must be provided when using hosted LangSmith API",
                        category=UserWarning)
# Vector Database
# Define persistent directory for Chroma
PERSIST_DIRECTORY = "chroma_db"
# PDF documents to include
pdf paths = ["./Use Edge ML Non-Invasive Beehive Monitoring System.pdf"]
# Define URLs for document sources
urls = [
    "https://mjrovai.github.io/EdgeML Made Ease ebook/raspi/setup/setup.html",
    "https://mjrovai.github.io/EdgeML Made Ease ebook/raspi/object detection/object detection fundamentals.html",
    "https://mjrovai.github.io/EdgeML_Made_Ease_ebook/raspi/object_detection/cv_yolo.html",
    "https://mjrovai.github.io/EdgeML_Made_Ease_ebook/raspi/llm/llm.html",
```

```
Chunk 1
 Chunk 2
 Chunk 3
Embedding 1
Embedding 2
Embedding 3
```

Vector Store

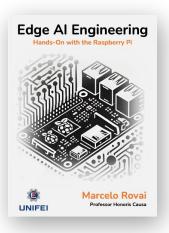
Retriever

```
# Split documents
print("Splitting documents into chunks...")
text splitter = RecursiveCharacterTextSplitter.from_tiktoken_encoder(
   chunk_size=300, chunk_overlap=30
doc_splits = text_splitter.split_documents(docs_list)
print(f"Created {len(doc_splits)} document chunks")
# Create embedding function
print("Initializing embedding model...")
embedding_function = OllamaEmbeddings(model="nomic-embed-text")
# Create and persist vectorstore to disk
print("Creating vector database...")
vectorstore = Chroma from_documents(
    documents=doc_splits,
    collection_name="rag-edgeai-eng-chroma",
    embedding=embedding_function,
    persist_directory=PERSIST_DIRECTORY
```

Create Knowledge Base

python 10-Create-Persistent-Vector-Database.py

```
marcelo_rovai — mjrovai@raspi-5: ~/Documents/Ollama/Rag/edgeai — ssh mjrovai@192.168.4.209 — 96×31
mjrovai@raspi-5:~/Documents/Ollama/Rag/edgeai $ python 10-Create-Persistent-Vector-Database.py
USER AGENT environment variable not set, consider setting it to identify your requests.
Database already exists at chroma db. Recreate? (y/n): y
Creating persistent vector store...
Loading PDF: ./data/2025 Edge AI Technology Report.pdf
Loading documents from URLs...
Loading URL: https://mjrovai.github.io/EdgeML Made Ease ebook/raspi/object detection/object dete
ction.html
Loading URL: https://mjrovai.github.io/EdgeML Made Ease ebook/raspi/image classification/image c
lassification.html
Loading URL: https://mjrovai.github.io/EdgeML Made Ease ebook/raspi/setup/setup.html
Loading URL: https://mjrovai.github.io/EdgeML Made Ease ebook/raspi/counting objects yolo/counti
ng objects yolo.html
Loading URL: https://mjrovai.github.io/EdgeML Made Ease ebook/raspi/llm/llm.html
Loading URL: https://mjrovai.github.io/EdgeML Made Ease ebook/raspi/vlm/vlm.html
Loading URL: https://mjrovai.github.io/EdgeML Made Ease ebook/raspi/physical comp/RPi Physical C
omputing.html
Loading URL: https://mjrovai.github.io/EdgeML Made Ease ebook/raspi/iot/slm iot.html
Total documents loaded: 95
Splitting documents into chunks...
Created 725 document chunks
Initializing embedding model...
Creating vector database...
/home/mjrovai/Documents/Ollama/Rag/edgeai/10-Create-Persistent-Vector-Database.py:99: LangChainD
eprecationWarning: Since Chroma 0.4.x the manual persistence method is no longer supported as do
cs are automatically persisted.
 vectorstore.persist()
Vector store created and saved to chroma db
Total document chunks indexed: 725
Database creation complete!
mjrovai@raspi-5:~/Documents/Ollama/Rag/edgeai $
```







```
● ● ■ marcelo_rovai — mjrovai@raspi-5-SD: ~/Documents/OLLAMA/RAG-Naive — ssh mjrovai@raspi-5-SD.local — 80×30
mjrovai@raspi-5-SD:- $ source ~/ollama/bin/activate
(ollama) mjrovai@raspi-5-SD: S cd Documents/
(ollama) mjrovai@raspi-5-SD:-/Documents $ cd OLLAMA/
(ollama) mjrovai@raspi-5-SD:-/Documents/OLLAMA $ cd RAG-Naive/
(ollama) mjrovai@raspi-5-SD:-/Documents/OLLAMA/RAG-Naive $ python 20-Query-the-P
ersistent-RAG-Database.py
Loading existing vector store...
==== RAG Query System ====
Type your questions and press Enter. Type 'quit' to exit.
Your question: List some popular SLM models
Generating answer...
Question: List some popular SLM models
Retrieving documents...
Retrieved 4 document chunks
Generating answer...
Response latency: 152.53 seconds using model: llama3.2:3b
ANSWER:
     ______
Some popular Small Language Models (SLMs) include Meta Llama, Microsoft PHI, Goo
gle Gemma, Ollama, MoonDream, Lava-Phi-3, and Gemma3. These models are compact v
ersions of larger language models that excel in specific tasks and domains, maki
ng them ideal for targeted applications on edge devices. They offer a range of b
enefits, including reduced parameter count, lower memory footprint, faster infer
ence time, energy efficiency, and offline functionality.
_______
```

No RAG

```
. .
               marcelo_rovai — mjrovai@raspi-5-SD: ~ — ssh mjrovai@raspi-5-SD.local — 80×28
mirovai@raspi-5-SD:- $ ollama run llama3.2:3b --verbose
>>> List some popular SLM models
Here are some popular Speech Language Model (SLM) architectures:
1. **Transformers**: The Transformer model, introduced in the paper
"Attention is All You Need" by Vaswani et al., has been widely adopted for
SLM tasks due to its ability to handle sequential data and its
parallelization capabilities.
2. **BERT** (Bidirectional Encoder Representations from Transformers):
BERT is a pre-trained language model that has been fine-tuned for various
NLP tasks, including text classification, sentiment analysis, and question
answering.
3. **RoBERTa**: RoBERTa is an improved version of BERT that uses a
different approach to generate training data, resulting in better
performance on certain tasks.
4. **XLNet**: XLNet is another pre-trained language model that has been
shown to perform well on a range of NLP tasks, including text
classification, sentiment analysis, and question answering.
5. **DistilBERT**: DistilBERT is a smaller version of BERT that was
trained using a distillation technique, which involves training a smaller
model on the output of a larger pre-trained model.
6. **Conversational AI models**: Models like "ConvAI" and "CIBER" are
designed specifically for conversational AI tasks, such as dialogue
systems and chatbots.
7. **Recurrent Neural Networks (RNNs)**: RNNs are a type of neural network
that is well-suited for sequential data, making them popular for SLM tasks
like speech recognition and natural language processing.
```







Knowledge Base

● ● ■ marcelo_rovai — mjrovai@raspi-5-SD: ~/Documents/OLLAMA/RAG-Naive — ssh mjrovai@raspi-5-SD.local — 80...

Your question: How is possible to automatically estimate the number of bees on a beehave?

Generating answer...

Ouestion: How is possible to automatically estimate the number of bees on a bee have?

Retrieving documents...

Retrieved 4 document chunks

Generating answer...

Response latency: 157.85 seconds using model: llama3.2:3b

ANSWER:

The system uses a YOLOv8n object detection model to automatically estimate the n umber of bees on a beehive. It is deployed on a Raspberry Pi Zero 2W embedded pl atform, which enables local processing without requiring cloud transmission. The system has been shown to have high accuracy and reliability in estimating bee r umbers, with a mean bias of 3.4 bees between automated and manual counts.

Your question:

No RAG



- >>> How is possible to automatically estimate the number of bees on a beehave? Estimating the number of bees in a beehive automatically can be challenging, but various techniques are being developed and researched using computer vision, machine learning, and sensor technologies. Here are some approaches:
- 1. **Computer Vision**: Researchers use computer vision techniques to analyze images or videos of the hive taken from cameras installed inside or outside the hive. They can detect bees, eggs, and other biological features to estimate the population size.
- 2. **Machine Learning Algorithms**: Machine learning algorithms like deep learning can be trained on labeled datasets of images or videos of hives to learn patterns and relationships between bee density and hive characteristics. These models can then be applied to new data to make predictions about bee populations.
- 3. **Acoustic Sensors**: Acoustic sensors can detect the sounds made by bees, such as wingbeats, to estimate the population size. This method is called "acoustic monitoring."
- 4. **Vibration Sensors**: Vibration sensors can measure the vibrations caused by bee movement within the hive, which can be used to estimate the population density.
- 5. **Thermal Imaging**: Thermal imaging cameras can detect heat signatures of bees and other animals in the hive, allowing researchers to estimate the population size.
- 6. **RGB-D Sensors**: RGB-D sensors combine color (RGB) and depth (3D) data from a single camera, providing insights into the spatial distribution of bees within the hive.
- 7. **Crowdsourcing and Drone-Based Monitoring**: Researchers are exploring the use of drones equipped with cameras and sensors to monitor hives remotely. This approach can provide high-resolution images and data on bee populations.

Optimized Naïve RAG Query

```
(ollama) mjrovai@raspi-5-SD:-/Pocuments/OLLAMA/RAG-Naive $ python 25-optimized R
AG query.py
Preloading models...
Models preloaded successfully
Loading existing vector store...
==== Optimized RAG Query System ====
Type your questions and press Enter. Type 'quit' to exit.
Your question: How is possible to automatically estimate the number of bees on a
 beehave?
Generating answer...
Question: How is possible to automatically estimate the number of bees on a beeh
ave?
Retrieving documents...
Retrieved 2 document chunks
Generating answer...
Response latency: 81.87 seconds using model: llama3.2:3b
ANSWER:
The system uses a YOLOv8n object detection model deployed on a Raspberry Pi Zero
2W, which enables local processing without requiring cloud transmission. This a
llows for automated bee counting at hive entrances by detecting and counting ind
ividual bees in images from the environment.
______
```



25-optimized RAG query.py

Questions?

Prof. Marcelo J. Rovai

rovai@unifei.edu.br

