

[Open in app](#)

★ Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#)

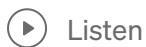
[Open in Google Cache](#)

Knowledge Graphs in RAGs (with Llama-Index)

[Open in Read-Medium](#)

Phil · [Follow](#)

3 min read · Jan 29, 2024

[Open in Freedium](#)

Listen



Share

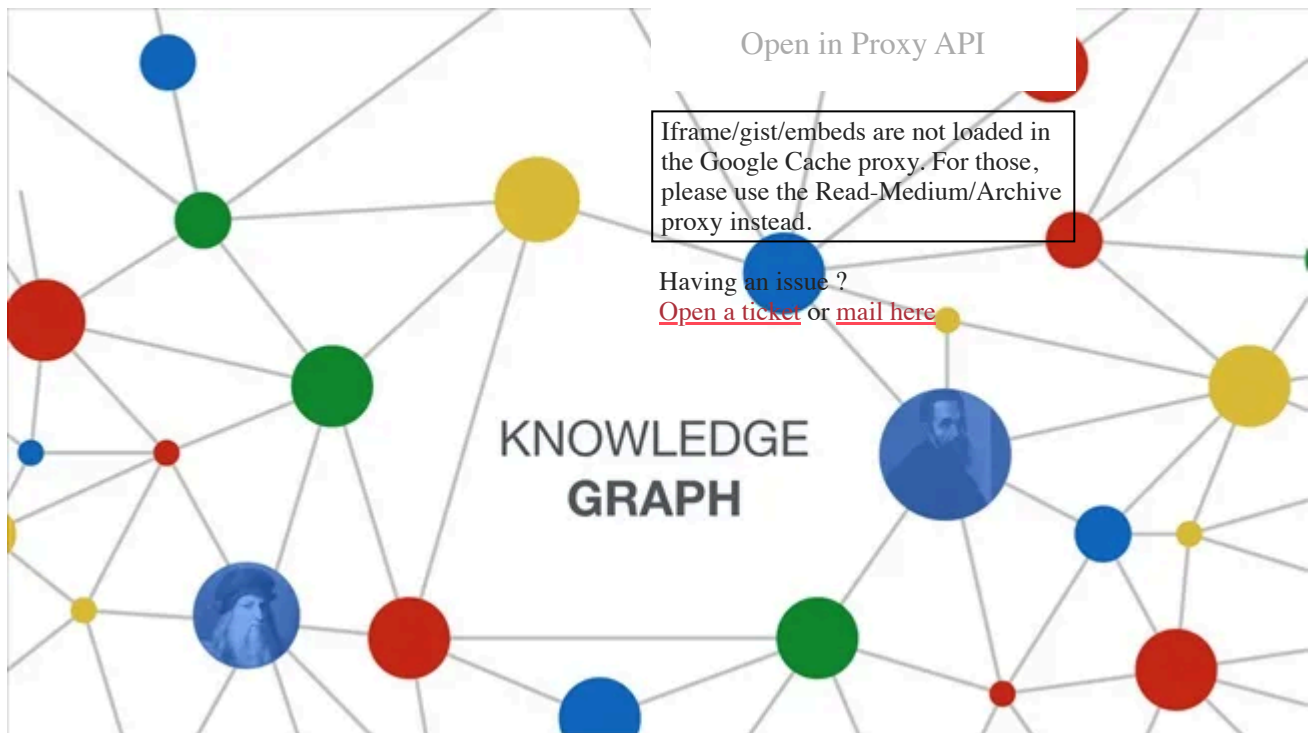
... More

[Open in Archive](#)[Open in Proxy API](#)

Iframe/gist/embeds are not loaded in the Google Cache proxy. For those, please use the Read-Medium/Archive proxy instead.

Having an issue ?

[Open a ticket](#) or [mail here](#)



Are #KnowledgeGraphs (KGs) better than #VectorDBs in RAGs? Yes and No!

Technically, KGs provide more #precise output than Vector DBs. KGs also support more #diverse and #complex querying than Vector DBs. Moreover, KGs enable better #reasoning and inferencing than Vector DBs

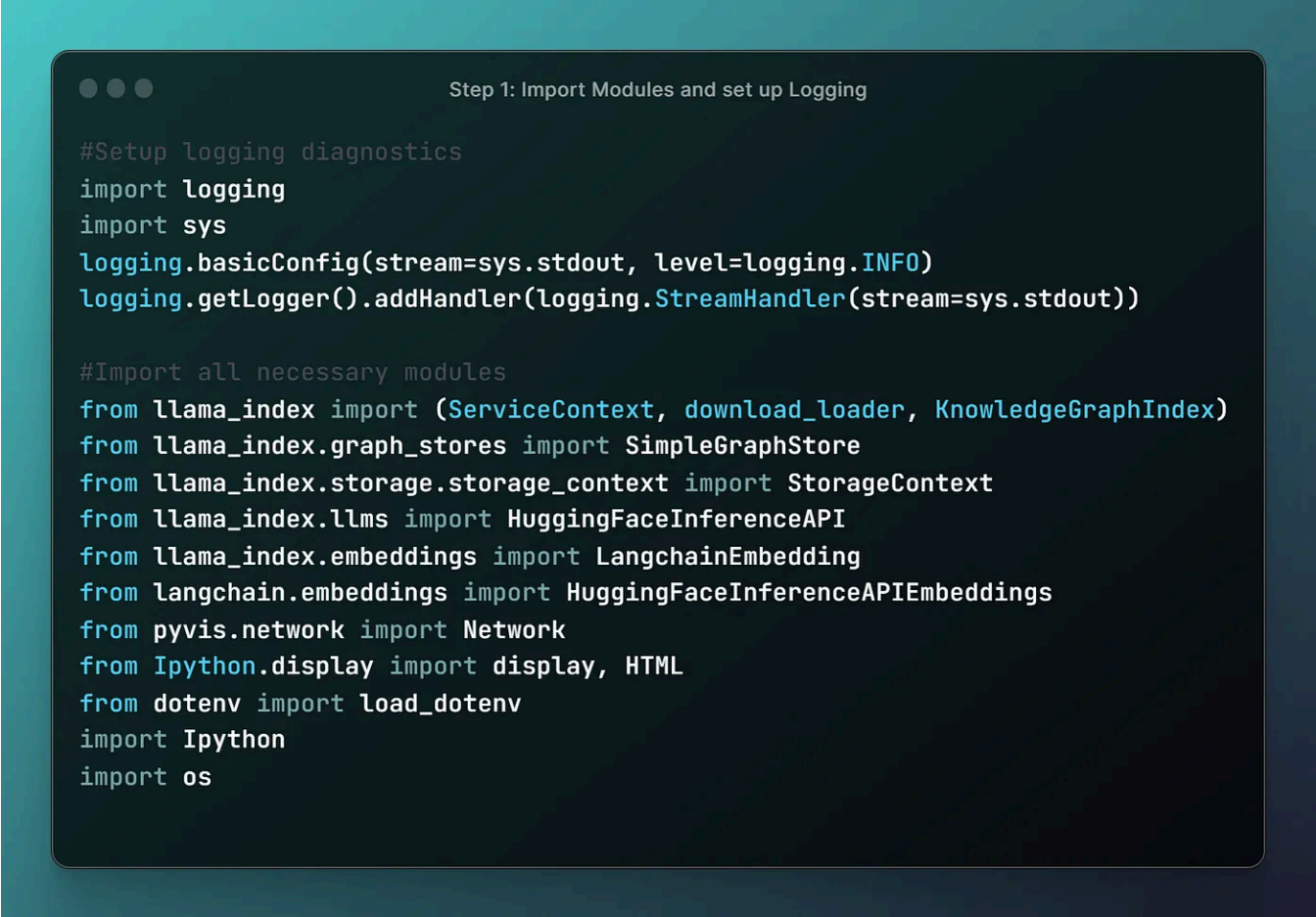
However, the KGs vs Vector DBs comparison is a game of Apples vs Oranges.

Knowledge Graphs are best suited for documents with clear #relationships while Vector DBs are more relevant to #similarity-based contexts. So, choosing between

KGs and Vector DBs depends on the specific requirements and objectives of your RAG project.

Here is my simple demonstration for the stepwise implementation of Knowledge Graph using the `#KnowledgeGraphIndex` class of `LlamaIndex`

1. Import Modules



```
Step 1: Import Modules and set up Logging

#Setup logging diagnostics
import logging
import sys
logging.basicConfig(stream=sys.stdout, level=logging.INFO)
logging.getLogger().addHandler(logging.StreamHandler(stream=sys.stdout))

#Import all necessary modules
from llama_index import (ServiceContext, download_loader, KnowledgeGraphIndex)
from llama_index.graph_stores import SimpleGraphStore
from llama_index.storage.storage_context import StorageContext
from llama_index.llms import HuggingFaceInferenceAPI
from llama_index.embeddings import LangchainEmbedding
from langchain.embeddings import HuggingFaceInferenceAPIEmbeddings
from pyvis.network import Network
from IPython.display import display, HTML
from dotenv import load_dotenv
import IPython
import os
```

The `KnowledgeGraphIndex` class streamlines RAG knowledge graph construction by automating entity extraction & relationship recognition from text, eliminating complex manual parsing. It also offers flexibility in customization, allowing you to tailor the graph structure and reasoning rules to your specific needs.

2. Set up LLM, embed_model, and Load Data

```
Step 2: Set up LLM and Load Data

#Set up LLM and Embedding Model
load_dotenv("llama.env")
face_token = os.getenv("FACE_TOKEN")
llm = HuggingFaceInferenceAPI(model_name = "HuggingFaceH4/zephyr-7b-beta",
face_token=face_token)
embed_model = LangchainEmbedding(HuggingFaceInferenceAPIEmbeddings(model_name =
                                                                    "HuggingFaceH4/zephyr-
7b-beta",
                                                                    face_token=face_token))

#Load data
WikipediaReader = download_loader("WikipediaReader")
loader = WikipediaReader()
documents = loader.load_data(pages=['Tesla Cybertruck'])
```

The `HuggingFaceInferenceAPI` class simplifies the consumption of Large Language Models (LLMs) for RAGs by providing a straightforward interface to pretrained models, facilitating efficient retrieval and generation of responses.

3. Construct the Knowledge Graph Index

```
Step 3: Build the Knowledge Graph Index to query

#Create Knowledge Graph Index
service_context = ServiceContext.from_defaults(chunk_size=300, chunk_overlap=30,
                                                llm=llm, embed_model=embed_model)

graph_store = SimpleGraphStore()
storage_context = StorageContext.from_defaults(graph_store=graph_store)

index = KnowledgeGraphIndex.from_documents(documents=documents, max_triplets_per_chunk=3,
                                           service_context=service_context,
                                           storage_context=storage_context,
                                           include_embeddings=True)
```

The `SimpleGraphStore` class in `llama_index` handles storage and retrieval of entities and relationships automatically, offers lightweight, in-memory storage solution for agile development, and integrates seamlessly with `KnowledgeGraphIndex` for a streamlined workflow.

4. Visualize and Query the Index

Step 4: Visualize and Query the Knowledge Graph Index

```
#Set up Graph Visualization and Display
graph = index.get_networkx_graph()
net = Network(notebook=True, cdn_resources = "in_line", directed=True)
net.from_nx(graph)
net.show("graph.html")

#Create Query Engine
query_engine = index.as_query_engine(include_text=True, response_mode="tree_summarize",
                                     embedding_mode="hybrid", similarity_top_k=5)
response = query_engine.query("According to Musk, what inspired the design of the
Cybertruck?")
print(response)
```

The `pyvis` module turns your RAG's knowledge graph into an interactive web visualization, letting you explore connections and patterns at a glance. But that's not all! LlamaIndex's `KnowledgeGraphQueryEngine` class is your key to advanced questioning. Craft nuanced queries, navigate entities through different modes, and uncover hidden insights your RAG holds — all with flexibility and precision.

Conclusion

Overall, **Knowledge Graphs** capture rich relationships among entities, enabling more precise, diverse, and complex querying and reasoning in RAGs than vector databases.

Retrieval Augmented

Rags

Ai Engineer

Python

Llamaindex



Follow

Written by Phil

6 Followers