

Generative AI

W2 Agenda

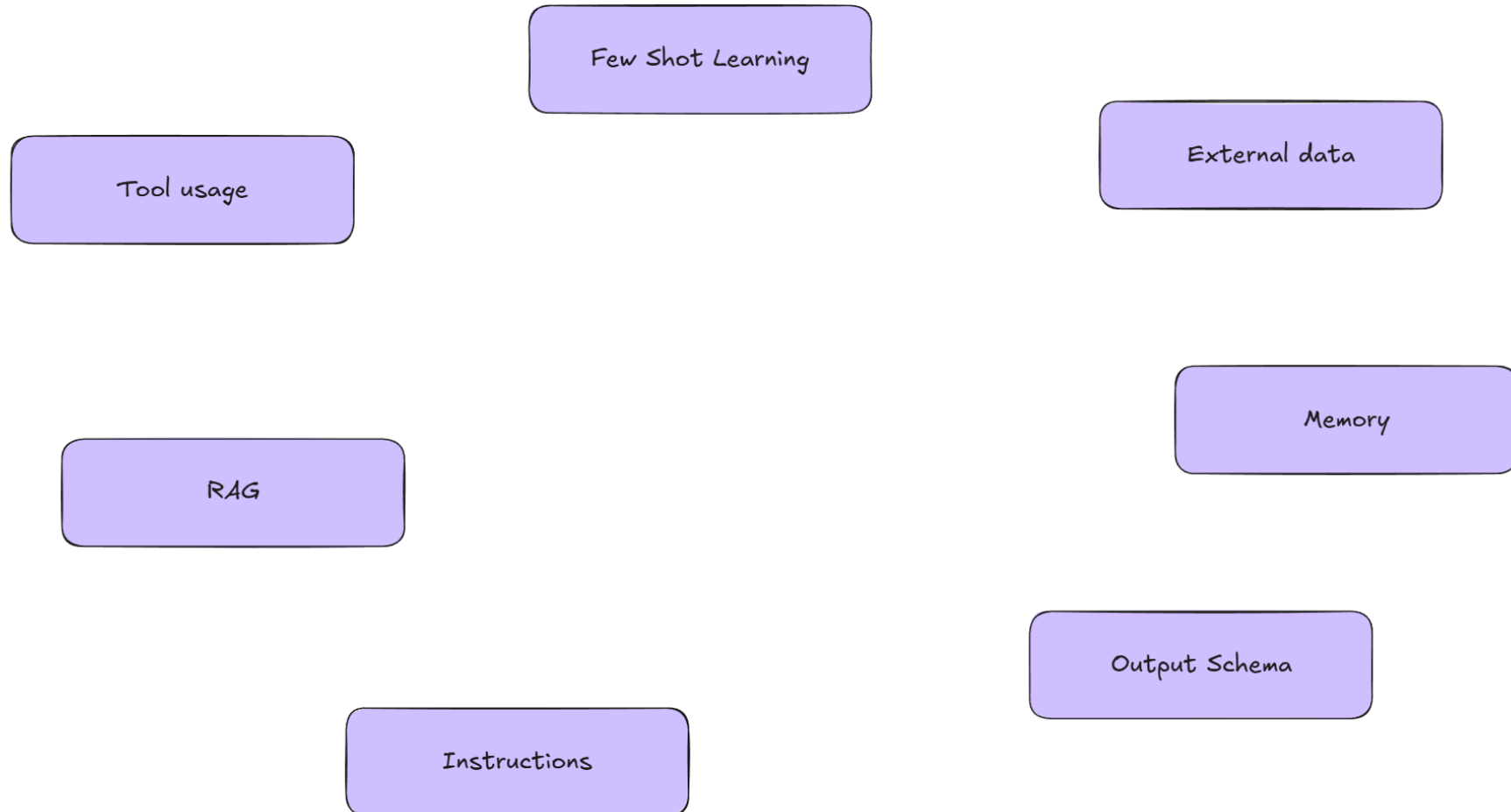
- **Introduction to working with LLMs**
 - **Prompt engineering**
 - **General rules**
- **Retrieval Augmented Generation**
 - **Vector search**
 - **Hybrid search**
- **How to leverage classic ML to make LLMs work better**
 - **Named Entity Recognition**
 - **CrossEncoders / Rerankers**

Introduction to working with LLMs

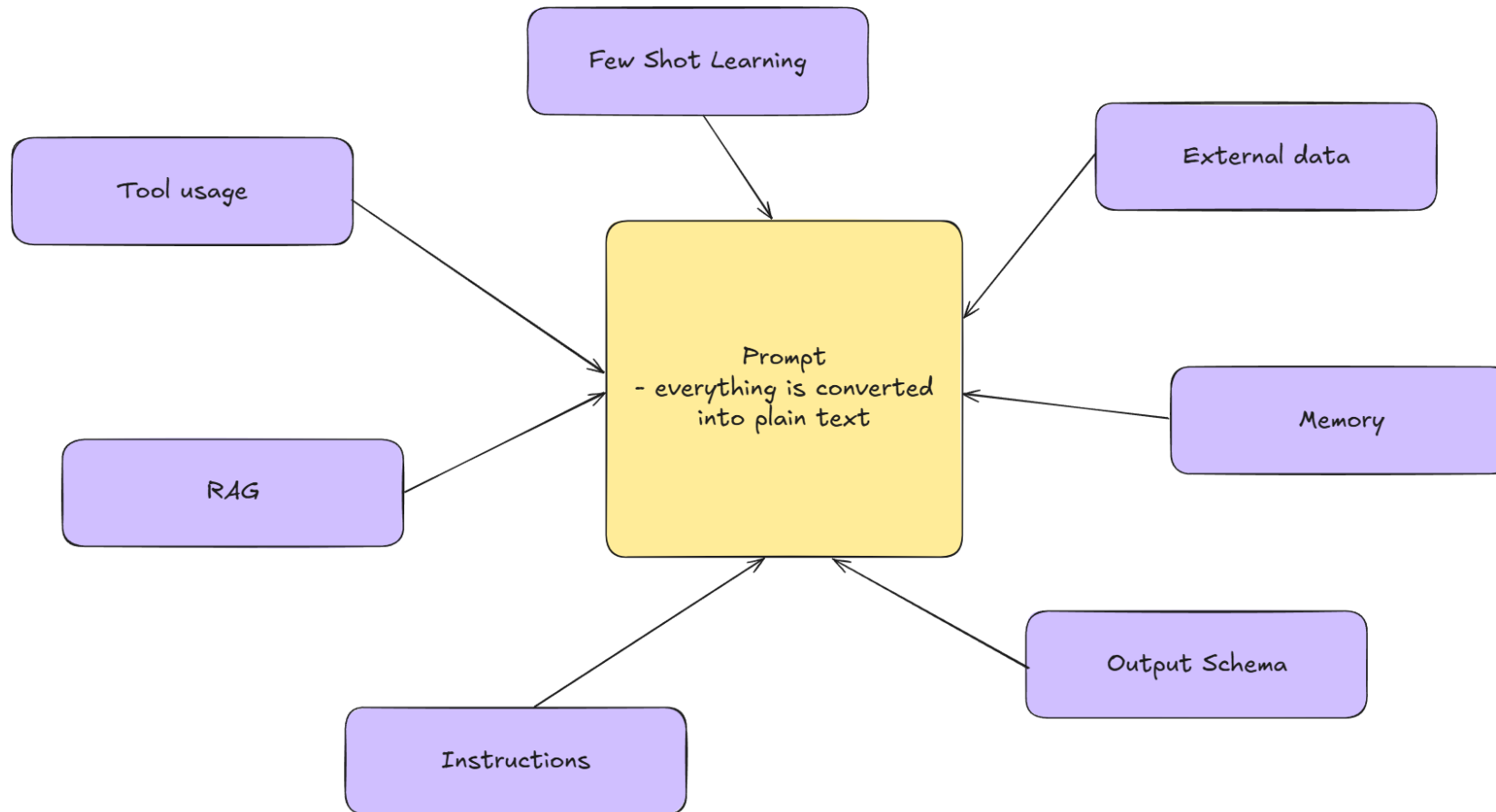
Prompt Engineering

How can you communicate with any LLM? What can be included in LLM call?

How can you communicate with any LLM? What can be included in LLM call?



All non-image data is converted to text, no matter how complex the process might seem



What is a prompt?

- Prompt text is the key model input for LLMs – it contains all the instructions together with conversation history and any additional context such as RAG
- Despite being dominated by Natural Language prompt can implement coded structures such as html like syntax
- While designing a prompt you should provide as detailed instructions and context as possible
- Prompt does not need to be equal to what we see in the chat window, additional elements can be appended such as:
 - initial model instructions
 - domain information from RAG
 - conversation history
- In context learning allows to help format model outputs by few-show-learning through prompting

Popular prompting templates

Question + Instruction



You

What are the 5 most important NLP topics I should include in a short NLP course? Give me an agenda with a brief description of each topic.

Instructions + Input data (zero-shot learning example)



You

Given the following information about me prepare a summary for my LinkedIn Profile. My name is Jan, I live in Warsaw, I graduated Power Engineering but for most of my career worked in Data Science. I am passionate about analysing Real Estate data and NLP models.

A.I. TURNS THIS SINGLE
BULLET POINT INTO A
LONG EMAIL I CAN
PRETEND I WROTE.



A.I. MAKES A SINGLE
BULLET POINT OUT OF
THIS LONG EMAIL I CAN
PRETEND I READ.



TOM
FISH
BURNE

Popular prompting templates

Question + Examples

J You
I really liked London, Rome and Prague. I disliked Zurich, Paris and Sosnowiec. What cities would you recommend for my next trip?

What is Prompt engineering

- Prompt engineering aims to create prompts in a more programmatic, reproducible way
- It involves carefully designing and structuring the input given to the model to effectively guide it towards generating the desired output or response.
- Temperature is one of the key hyperparams controlling how varied and creative each models response is – even with the exact same prompt
- Prompt Engineering can be more structured by leveraging libraries such as LangChain

Chain-of-thought reasoning

- Chain-of-thought prompting: Encourages model to be factual and precise by explaining its reasoning

J

You

What are the best libraries to get started with NER?

Use this format:

Q: <repeat question>

A: Let's analyze the possibilities. <give reasoning>, summarizing starting with NER is <final answer>

Q: What are the best libraries to get started with Named Entity Recognition (NER)?

A: Let's analyze the possibilities. There are several libraries catering to NER tasks, but the choice often depends on factors like ease of use, community support, and performance. spaCy stands out as an excellent option for beginners. Its user-friendly API, pre-trained models, and efficient processing make it a go-to choice for those new to NER. Therefore, spaCy is recommended as a solid starting point for diving into Named Entity Recognition.

Citing sources

J You

What are the top 3 LLM models? Answer only with reliable sources and cite those sources.

 **ChatGPT**

As of my last knowledge update in January 2022, three prominent large language models (LLMs) were widely recognized:

1. **GPT-3 (Generative Pre-trained Transformer 3)**

- **Source:** Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language Models are Few-Shot Learners. [arXiv:2005.14165](https://arxiv.org/abs/2005.14165).

2. **BERT (Bidirectional Encoder Representations from Transformers)**

- **Source:** Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. [arXiv:1810.04805](https://arxiv.org/abs/1810.04805).

3. **T5 (Text-to-Text Transfer Transformer)**

- **Source:** Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... & Liu, P. J. (2019). Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. [arXiv:1910.10683](https://arxiv.org/abs/1910.10683).



Warning!
LLM can still make up sources, but at least you can check if they exists


Few shot learning

Zero-shot

 You

User

'I rather enjoyed this movie, but the beginning was boring'. Please conduct sentiment analysis on this review.




 ChatGPT

Sentiment analysis involves determining the sentiment or emotion expressed in a piece of text. In the review you provided, "I rather enjoyed this movie, but the beginning was boring," the sentiment is mixed. Let's break it down:

1. "I rather enjoyed this movie" - This part of the review expresses a positive sentiment. The word "enjoyed" indicates a positive feeling towards the movie.

2. "but the beginning was boring" - On the other hand, this part of the review expresses a negative sentiment. The word "boring" indicates a negative feeling towards the beginning of the movie.

So, overall, the sentiment in this review is mixed, with a positive sentiment towards the movie as a whole but a negative sentiment towards the beginning of the movie.

Two-shot

 You

Classify this review:

"I hated this movie, total waste of time and money"

Sentiment: Negative

Classify this review:

"The most beautiful love story in the history of cinema"

Sentiment: Positive

Classify this review:

'I rather enjoyed this movie, but the beginning was boring'.


Sentiment:

 ChatGPT

Sentiment: Mixed (Positive about the overall enjoyment, but negative about the beginning)

Three-shot

 You

User

Classify this review:

"I hated this movie, total waste of time and money"

Sentiment: Negative

Classify this review:

"The most beautiful love story in the history of cinema"

Sentiment: Positive

Classify this review:


"The beginning dragged on forever, but I liked the surprising ending. Overall experience was average"

Sentiment: Neutral

Classify this review:

'I rather enjoyed this movie, but the beginning was boring'.


Sentiment:

 ChatGPT

Sentiment: Mixed or Neutral


Few shot learning

Zero-shot

 You

User

'I rather enjoyed this movie, but the beginning was boring'. Please conduct sentiment analysis on this review.




 ChatGPT

Sentiment analysis involves determining the sentiment or emotion expressed in a piece of text. In the review you provided, "I rather enjoyed this movie, but the beginning was boring," the sentiment is mixed. Let's break it down:

1. "I rather enjoyed this movie" - This part of the review expresses a positive sentiment. The word "enjoyed" indicates a positive feeling towards the movie.

2. "but the beginning was boring" - On the other hand, this part of the review expresses a negative sentiment. The word "boring" indicates a negative feeling towards the beginning of the movie.

So, overall, the sentiment in this review is mixed, with a positive sentiment towards the movie as a whole but a negative sentiment towards the beginning of the movie.

Two-shot

 You

Classify this review:

"I hated this movie, total waste of time and money"

Sentiment: Negative

Classify this review:

"The most beautiful love story in the history of cinema"

Sentiment: Positive

Classify this review:

'I rather enjoyed this movie, but the beginning was boring'.


Sentiment:

 ChatGPT

Sentiment: Mixed (Positive about the overall enjoyment, but negative about the beginning)

Three-shot

 You

User

Classify this review:

"I hated this movie, total waste of time and money"

Sentiment: Negative

Classify this review:

"The most beautiful love story in the history of cinema"

Sentiment: Positive

Classify this review:


"The beginning dragged on forever, but I liked the surprising ending. Overall experience was average"

Sentiment: Neutral

Classify this review:

'I rather enjoyed this movie, but the beginning was boring'.

Sentiment:

 ChatGPT

Sentiment: Mixed or Neutral

Jupyter exercise Few shot learning

Proceed to notebook W2-llm-calling-and-vector-search

✓ Setup keys and parsers for later

```
openai.api_key = OPENAI_API_KEY = os.getenv("OPENAI_API_KEY")

client = openai.Client()

output_parser = JsonOutputParser()
```

[19] ✓ 0.3s

Load data of 100 FAQ items from a travel company

```
df = pd.read_json("data/travel_company_faq.json")
i=10
questions = df.loc[i].question
```

[]

Classify FAQ items with LLM

```
##TODO: Prepare system prompt to split FAQ items into 4 categories: [air-travel, hotels-and-booking, food,
# System prompt with instructions and the JSON format we expect
system_prompt = f"""You are a helpful assistant.
Your task: Classify the given question into one of the following categories:
[air-travel, hotels-and-booking, food, insurance, extra-activities].

Return the result ONLY in valid JSON, in the format:
{{"category": "<category>"}}
```

[]

```
# --- Step 2: Define a function to call OpenAI's chat completion directly ---
def classify_question(question: str, system_prompt: str) -> str:
    """
    Calls the OpenAI ChatCompletion endpoint to classify the question.
    Returns the predicted category.
    """

    # Make a direct OpenAI chat call (using the fictional "gpt-4o-mini" as specified)

    response = client.chat.completions.create(
        model="gpt-4o-mini", # Updated to match available models
```

[20]

General Tips

DOs and DONTs



- Single, clear task per each call
- Concise instructions, which can be followed as a human
- Leverage markdown, html and system prompt for clear segmentation between input, instructions, few-shot examples and any other arguments
- Complex tasks should be split into several, serial prompt calls
- Asking for a score or explanation can improve results, even if the scores are not reliable
- Chain-of-thought framework can improve logical capabilities



- Don't flood LLMs with too much info, even if they have 1M tokens window
- LLMs are pretty bad at numerical scoring
- Be careful with too many / too extreme examples as they might steer answers in unexpected directions
- It is hard to tell LLM NOT TO DO something
- LLMs are not designed for counting – it might be more reliable to extract values and summarize them programmatically
- Don't overkill simple tasks with LLMs – many of them can be solved more efficiently with classic ML

It is harder to instruct LLM not to do something

create an image of empty room without any elephants in it



Here's the image of an empty room with no elephants in it. Let me know if you need any modifications!

Generated by author with gpt-4o, 2025-03-04

Retrieval Augmented Generation (RAG)

What is RAG?

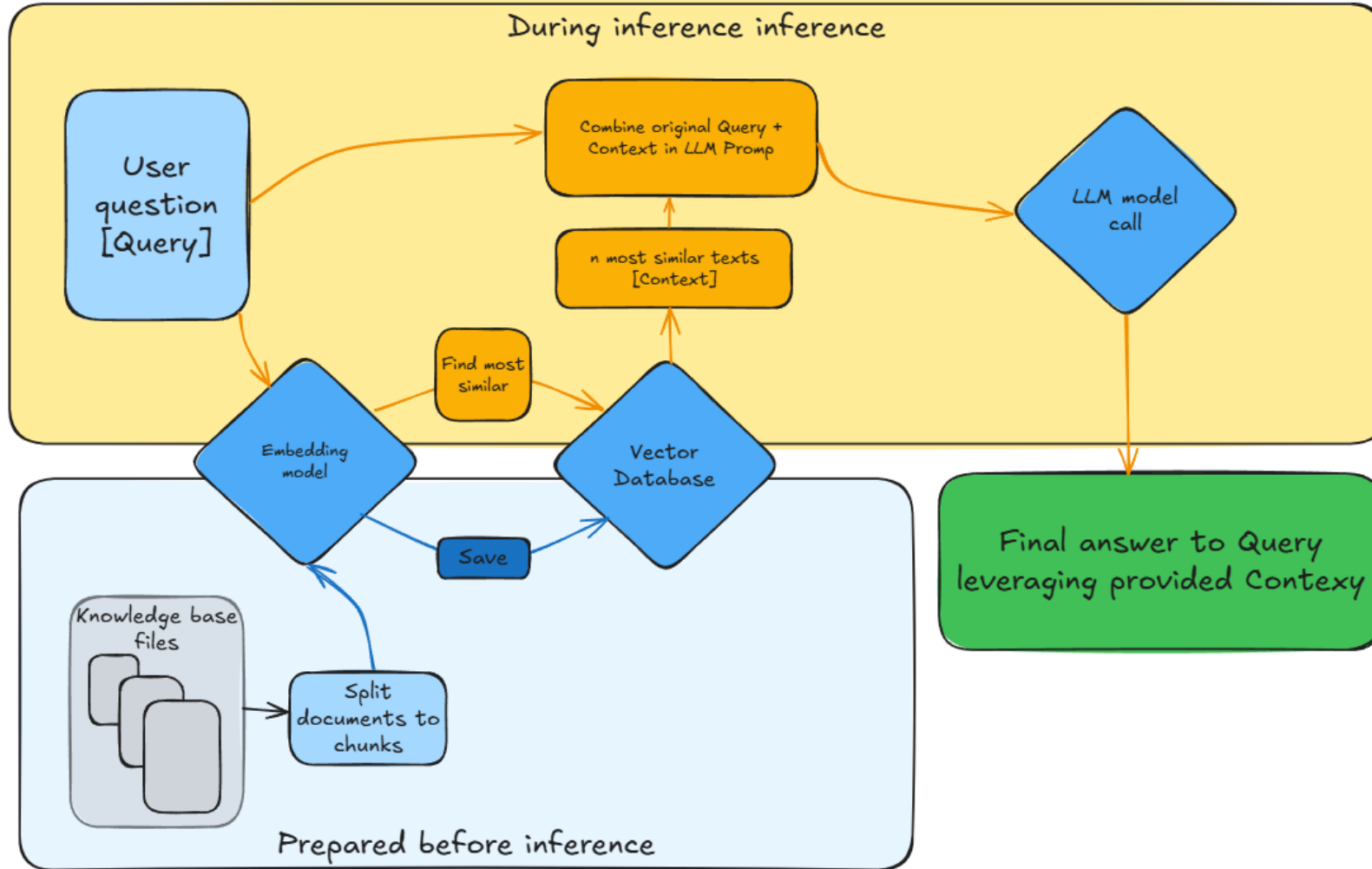
Retrieval augmented generation (RAG) allows us to improve LLM performance by enriching our calls with some selected context

It allows expanding and updating model knowledge with external data, which improves information accuracy and relevance grounding its answers

Most popular forms of RAG rely on finding sections of documents, which are most similar to asked questions and integrating them to LLM call for better context

This approach can be used to feed real time or internal data to guide LLM answers, without need for retraining

Vector RAG diagram



Source: image by author

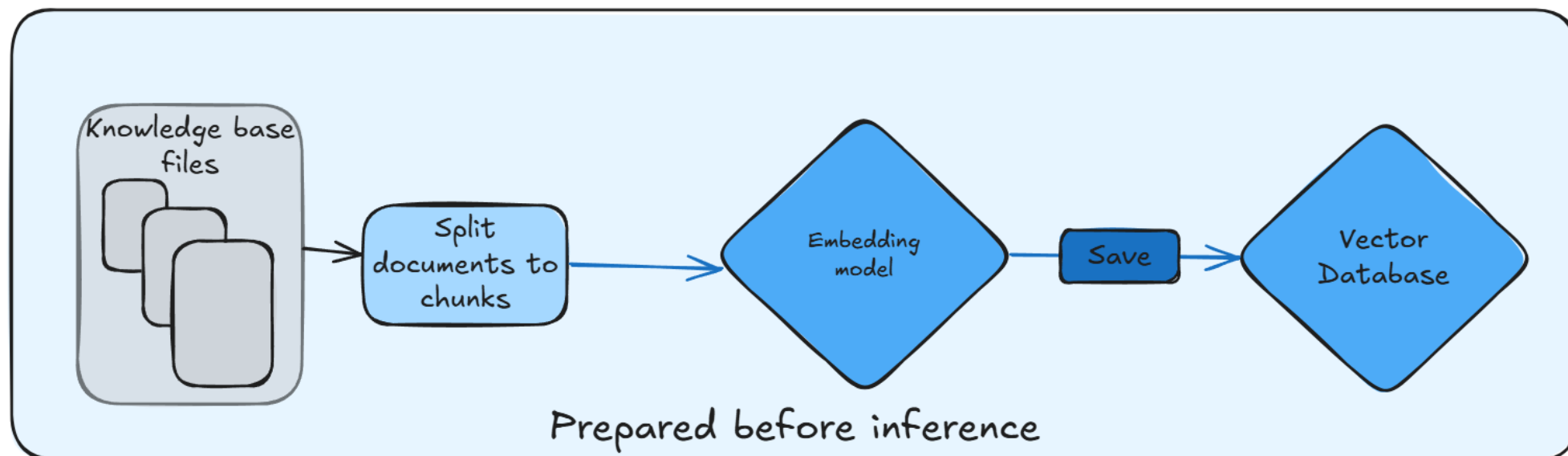
RAG benefits

- Grounding model answers with controlled context, which can reduce hallucinations
- Ability to provide specified, real-time data to LLM e.g. fast changing promotions or product availability
- Well-defined, specialized RAG can get good performance with smaller LLMs

RAG does NOT reason – key pitfalls

- Despite being closely associated with LLMs, RAG relies mainly on vector search
- It DOES NOT have any reasoning or logical abilities and can be easily misled
- Purely vector based RAG carries many shortcomings of semantic similarity -> it will not perform well with complex questions and numbers
- Feeding bad context to LLM call can increase risk of bad answers as LLM will be more likely to use this information

Creating and calling Vector database



1. Gather files for Knowledge base

2. Split files into smaller chunks within a few hundred tokens

3. Embed each chunk as a vector

4. Save each chunk's text, vector, metadata and metadata to DB

Jupyter exercise Vector DB

Continue in notebook W2-llm-calling-and-vector-search

Introduction to vector search

```
chroma_db_path = "chroma_db"
chroma_client = chromadb.PersistentClient(path=chroma_db_path)
```

```
SELECTED_COLLECTION = "travel-company-faq"
```

✓ 0.5s

```
embedding_model = "text-embedding-ada-002"
openai_ef = embedding_functions.OpenAIEmbeddingFunction(model_name=embedding_model, api_key = OPENAI_API_KEY)
```

```
collection = chroma_client.get_or_create_collection(name=SELECTED_COLLECTION, embedding_function=openai_ef)
```

✓ 0.0s

```
def ingest_faq_data(df: pd.DataFrame, collection):
```

```
    """
```

```
    Ingest combined question and answer as vectorized documents. Store question, answer and category as metadata
```

```
    """
```

```
    all_ids = []
```

```
    all_documents = []
```

```
    all_metadatas = []
```

```
    for i, row in df.iterrows():
```

```
        # Combine Q + A as text
```

```
        doc_text = f"Question: {row['question']}\nAnswer: {row['answer']}"
```

```
        doc_id = f"faq_{i}"
```

```
        meta = {
```

```
            "question": row["question"],
```

```
            "answer": row["answer"],
```

```
            "category": row["category"],
```

```
        }
```

```
        all_ids.append(doc_id)
```

```
        all_documents.append(doc_text)
```

```
        all_metadatas.append(meta)
```

```
    collection.add(documents=all_documents, metadatas=all_metadatas, ids=all_ids)
```

```
    print(f"Ingested {len(df)} records into {SELECTED_COLLECTION}.")
```

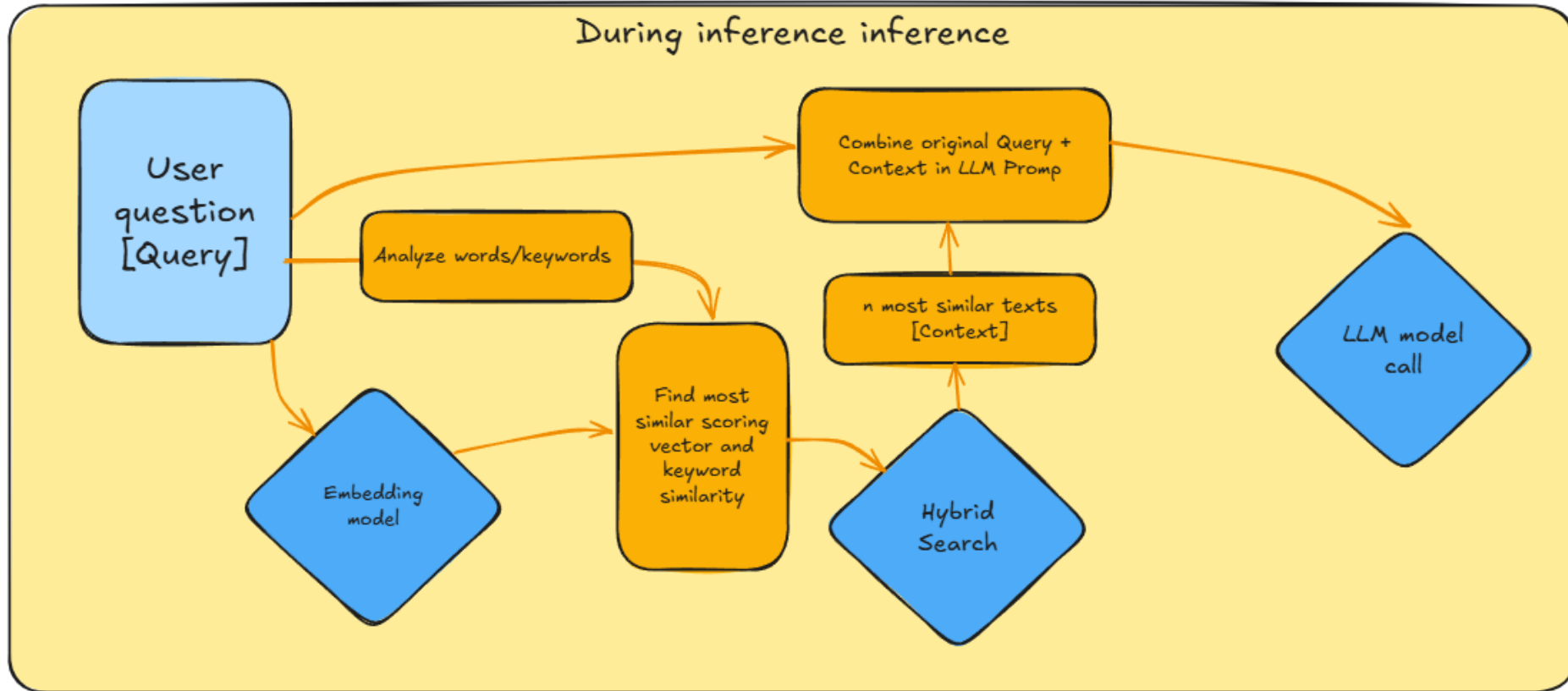
✓ 0.0s

```
ingest_faq_data(df, collection)
```

✓ 1.4s

Add of existing embedding ID: faq_0

Hybrid search – combining vectors and keywords



Hybrid search – combining vectors and keywords

- Hybrid search combines vector similarity with classic keyword-based searches
- BM25 is a popular keyword search leveraging TF-IDF (Term Frequency-Inverse Document Frequency) algorithm
- Keywords can also be used for hard filtering within document text (like in SQL) or metadata (if we want to limit our search to categories or documents)

Improving LLM results with classic ML

Why is it worth using classic ML when we have LLMs

- LLMs are extremely flexible, but some tasks can be still completed with better speed and accuracy with pretrained NLP models e.g. Bert-based
- They are significantly cheaper, faster and easier to deploy
- Training/Finetuning can provide predictable results within specific domains if we have the data
- Validation with classic ML can reduce LLM hallucinations
- They give us to better control LLMs context and answers
- Can serve for evaluation of LLMs

Named Entity Recognition summary

- We use NER when we want to extract specific information from text
- NER extract structured data from text
- Single NER entity e.g. person can have thousands of instances
- Position of the entity itself is important, which makes data annotation as well as processing especially challenging

geo **gpe** **per**
Lebanon 's top Shi'ite cleric is opposing British Prime Minister
Tony Blair 's expected visit to Beirut Monday .
per **geo** **tim**

```
[(0, 7, 'geo'),  
(42, 49, 'gpe'),  
(50, 55, 'per'),  
(65, 75, 'per'),  
(97, 103, 'geo'),  
(104, 110, 'tim')]
```

org **per**
Democratic U.S. Senator Edward Kennedy has urged the government
to spend more money on education as millions of students return
to school for a new academic year .
tim

```
[(11, 15, 'org'), (16, 38, 'per'), (144, 156, 'tim')]
```


Under-the-hood NER is classifying tokens by their relations to specific entities

Spacy uses BIO labeling by default, where each token in analyzed text can get assigned the following values:

- **"B"**: Beginning of an entity
- **"I"**: Inside an entity
- **"O"**: Outside of any entity

Together with the label it also gets assigned an entity class.

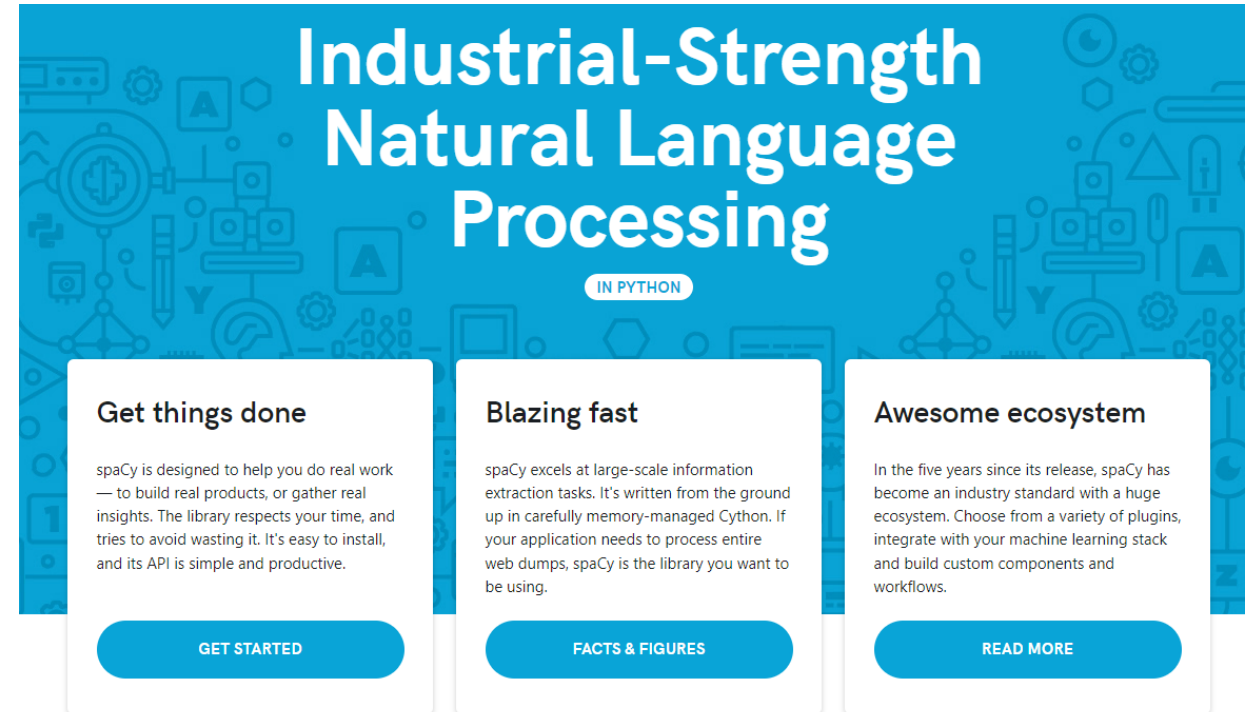
From model perspective we are conducting multilabel classification for each token.

When **Sebastian Thrun** started working on self-driving cars at **Google** in **2007**, few people outside of the company took him seriously. I can tell you very senior CEOs of major **American** car companies would shake my hand and turn away because I wasn't worth talking to, said **Thrun**, in an interview with **Recode** earlier **this week**.

```
Token: Sebastian      ent_iob: B      ent_type: PERSON
Token: Thrun          ent_iob: I      ent_type: PERSON
Token: 2007           ent_iob: B      ent_type: DATE
Token: American       ent_iob: B      ent_type: NORP
Token: Thrun          ent_iob: B      ent_type: PERSON
Token: Recode         ent_iob: B      ent_type: ORG
Token: this           ent_iob: B      ent_type: DATE
Token: week           ent_iob: I      ent_type: DATE
```

What is Spacy

- Spacy is one of the most popular libraries for NLP
- It has a low barrier of entry
- You can build models within a few dozens lines of code in a notebook
- CLI interface great for building production models and replicability



The banner features a blue background with white text and icons. The main title is 'Industrial-Strength Natural Language Processing' in a large, bold font. Below it, a small white box contains the text 'IN PYTHON'. There are three white cards with blue borders and blue buttons at the bottom. Each card has a title, a paragraph of text, and a button.

Industrial-Strength Natural Language Processing

IN PYTHON

Get things done

spaCy is designed to help you do real work — to build real products, or gather real insights. The library respects your time, and tries to avoid wasting it. It's easy to install, and its API is simple and productive.

GET STARTED

Blazing fast

spaCy excels at large-scale information extraction tasks. It's written from the ground up in carefully memory-managed Cython. If your application needs to process entire web dumps, spaCy is the library you want to be using.

FACTS & FIGURES

Awesome ecosystem

In the five years since its release, spaCy has become an industry standard with a huge ecosystem. Choose from a variety of plugins, integrate with your machine learning stack and build custom components and workflows.

READ MORE

Jupyter exercise NER:

Proceed to notebook W2-NER-Intro-BLANK

1 Spacy intro

In [18]:

```
import spacy

# Load English tokenizer, tagger, parser and NER
nlp = spacy.load("en_core_web_sm")

# Process whole documents
text = ("When Sebastian Thrun started working on self-driving cars at "
        "Google in 2007, few people outside of the company took him "
        "seriously. "I can tell you very senior CEOs of major American "
        "car companies would shake my hand and turn away because I wasn't "
        "worth talking to," said Thrun, in an interview with Recode earlier "
        "this week.")
```

In [19]:

```
doc = nlp(text)
```

In [20]:

```
# Analyze syntax
print("Noun phrases:", [chunk.text for chunk in doc.noun_chunks])
print("Verbs:", [token.lemma_ for token in doc if token.pos_ == "VERB"])
```

```
Noun phrases: ['Sebastian Thrun', 'self-driving cars', 'Google', 'few people', 'the company', 'him', 'I', 'you', 'very senior C
EOs', 'major American car companies', 'my hand', 'I', 'Thrun', 'an interview', 'Recode']
Verbs: ['start', 'work', 'drive', 'take', 'tell', 'shake', 'turn', 'talk', 'say']
```

In [21]:

```
# Find named entities, phrases and concepts
for entity in doc.ents:
    print(entity.text, entity.label_)
```

```
Sebastian Thrun PERSON
Google ORG
2007 DATE
American NORP
Thrun PERSON
Recode ORG
earlier this week DATE
```

1.1 Token labels

In []:

```
for token in doc:
    if token.ent_iob_ != 'O':
        print(f"Token: {token.text}\tent_iob: {token.ent_iob_}\tent_type: {token.ent_type_}")
```

2 Preparing Spacy format

In [24]:

```
with open('../data/GMB_data_sagemaker.pickle', 'rb') as f:
    sagemaker_data = pickle.load(f)
```

In [85]:

```
## Current format
sagemaker_data[0]
```

Out[85]:

```
{'workerId': 'some-random-worker-no-123',
 'dataObject': {'content': 'Mr. Egeland said the latest figures show 1.8 million people are in need of food assistance - with t
he need greatest in Indonesia , Sri Lanka , the Maldives and India .'},
 'annotationData': {'content': {'entities': [{'endOffset': 11,
      'label': 'per',
      'startOffset': 0},
      {'endOffset': 128, 'label': 'tim', 'startOffset': 119},
      {'endOffset': 134, 'label': 'per', 'startOffset': 131},
      {'endOffset': 140, 'label': 'gpe', 'startOffset': 135},
      {'endOffset': 155, 'label': 'geo', 'startOffset': 147},
      {'endOffset': 165, 'label': 'geo', 'startOffset': 160}]}}}}
```

Introduction to CrossEncoders

Crossencoders are a type of neural network used for ranking tasks – They take a pair of inputs (e.g., a query and a document) and jointly encode them into a single representation before making a relevance prediction.

They can be fine-tuned to perform multiple tasks e.g. similarity detection, answer relevance. We can tailor training data to fit our specific needs.

Under the hood they are based on encoder only transformer with additional neural network layers for regression task, outputting a single number according to their training goal

Pre-trained X-encoders are widely used for reranking of retrieval tasks.

They can also be used for evaluating logical similarity to evaluate LLM answers

CrossEncoders as rerankers – most common use case in RAG

- As semantic similarity can be misleading, it is worth pulling more records than needed and then reranking most relevant using a more steerable and stronger model
- Evaluating query \leftrightarrow context pairs fit with cross-encoders can be much more accurate than vector similarity and much cheaper than LLMs
- We can also fine-tune to X-encoders for our own goals e.g. answers following company specific rules such as different approach to different products/services/segments

How to choose best model

- Finding best model yourself is next to impossible – luckily, we have HF benchmarks focused on specific tasks and languages
- Especially when working with non-English text most popular models might not perform best
- Find relevant benchmark and choose best model within the size (params count) that you can use



Polish retrieval/reranking benchmark

Evaluation metric: ☒ NDCG@10 ☐ MRR@10 ☐ Recall@100 ☐ Accuracy@1

Show methods: ☒ Retrievers ☐ Hybrid retrieval ☐ Rerankers



Mutli-lingual embedding benchmark

Benchmarks

Select one of the hand-curated benchmarks from our publications and modify them using one of the following filters to fit your needs.

Prebuilt Benchmarks

Select one of our expert-selected benchmarks from MTEB publications.

MTEB(Multilingual, v1)

Select Languages

Select Task Types

Select Domains

Add and remove tasks:

Model Selection

Select models to rank based on an assortment of criteria.

Search Models

Press Enter to search.

Search models by name (Regex sensitive. Separate queries with |)

Availability

☐ Only Open ☐ Only Proprietary

☒ Both

Instructions

☐ Only Instruction-tuned

☐ Only non-instruction ☒ Both

Compatibility

☐ Should be sentence-transformers compatible

Zero-shot

☐ Only Zero-shot ☐ Remove Unknown

☒ Allow All

Model Size (#M Parameters)

0 100000 0

0 100000

Jupyter exercise Reranking

Continue in notebook W2-llm-calling-and-vector-search

Rerank answers

How do rerankers work??

<https://huggingface.co/mixedbread-ai/mxbai-rerank-base-v1>

```
from sentence_transformers import CrossEncoder

# Load the model, here we use our base sized model
model = CrossEncoder("mixedbread-ai/mxbai-rerank-xsmall-v1")
```

```
# Example query and documents
query = "Who wrote 'To Kill a Mockingbird'?"
documents = [
```

```
    "'To Kill a Mockingbird' is a novel by Harper Lee published in 1960. It was immediately successful, winning the Pulitzer Prize in 1961."
    "The novel 'Moby-Dick' was written by Herman Melville and first published in 1851. It is considered a masterpiece of American literature."
    "Harper Lee, an American novelist widely known for her novel 'To Kill a Mockingbird', was born in 1926 in Monroeville, Alabama."
    "Jane Austen was an English novelist known primarily for her six major novels, which interpret, critique and comment on the English landed gentry at the end of the 18th century."
    "The 'Harry Potter' series, which consists of seven fantasy novels written by British author J.K. Rowling, is a series of children's books that have become a global phenomenon."
    "'The Great Gatsby' is a 1925 novel by American author F. Scott Fitzgerald. It is a critique of the American Dream and the excesses of the Roaring Twenties in New York City."
]
```

```
results = model.rank(query, documents, return_documents=True, top_k=3)
```

✓ 0.8s

```
results = collection.query(query_texts=[query], n_results=10)
documents = results["documents"][0]
```

✓ 0.9s

```
##TODO: Get scores for results of initial RAG
results_with_scores = model.rank(query, documents, return_documents=True, top_k=3)
```

✓ 0.1s

```
results_with_scores
```

✓ 0.0s

```
[{'corpus_id': 0,
  'score': np.float32(0.9933727),
  'text': 'Question: Does the insurance cover stolen personal belongings?\nAnswer: Most plans include coverage for theft of personal belongings, but it is important to check the specific terms of your policy.'},
 {'corpus_id': 2,
  'score': np.float32(0.36291805),
  'text': 'Question: How do I file a claim if I lose my luggage?\nAnswer: First, report the loss to the airline and file a claim with your travel insurance provider. Make sure to have a copy of your luggage receipt and a list of the items lost.'},
 {'corpus_id': 4,
```