

Can the modern Frankenstein pass the Turing test?

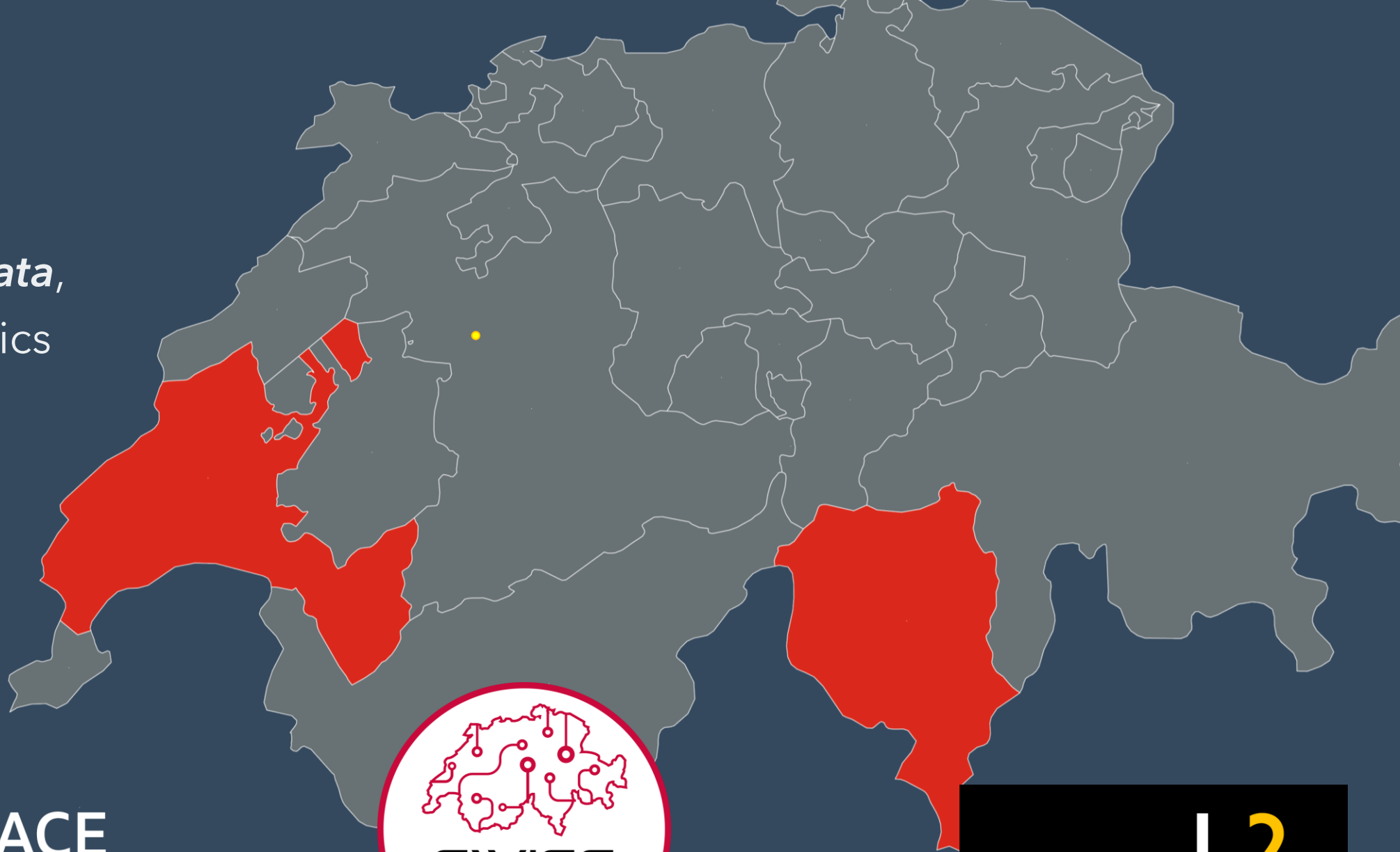
November 26, 2024

Gianni Ceresa



Gianni Ceresa

Working with **data**,
Business Analytics
and EPM tools
for more than
15 years



Oracle ACE
Director



._SYM^{L2}



430+ technical experts helping peers globally

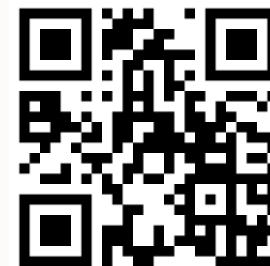
The **Oracle ACE Program** recognizes and rewards community members for their technical and community contributions to the Oracle community



3 membership tiers



For more details on Oracle ACE Program:
ace.oracle.com



Nominate
yourself or someone you know:
ace.oracle.com/nominate

Connect: aceprogram_ww@oracle.com

Facebook.com/OracleACEs

[@oracleace](https://twitter.com/oracleace)

[Oracle ACE Program Group](https://www.linkedin.com/groups/oracle-ace-program-group)



DISCLAIMER



This topic could easily take a few hours.
Making it fit in 45 minutes requires choices on
what to ignore.

Every week I try a different approach, the slides
and the code were still changing earlier in the
train.

Can the modern Frankenstein pass the Turing test?

Building your own custom Oracle Database
ChatGPT-like solution

ChatGPT

ChatGPT is great for various tasks, you can ask it anything and it will (almost) never complain and will give you an answer.

How good is the answer?

Depends: how good was the question?

Just like you and me can make up a story, so can ChatGPT with hallucinations:

- False or misleading information presented as facts
- Analysts estimated that chatbots hallucinate as much as 27% of the time, with factual errors present in 46% of generated texts* (*this was in 2023, things change quickly*)

A chatbot powered by LLM based on facts you provide

Where to start to build an Oracle Database ChatGPT-like solution?

RAG: retrieval-augmented generation

A technique for enhancing the accuracy and reliability of generative AI models with facts fetched from external sources

RAG: one acronym, an almost infinite number of ways to achieve it...

- From the simplest retrieval fed to a LLM
- To more advanced processes with decisions made before and evaluation of the quality of the facts provided and the quality of the answer generated

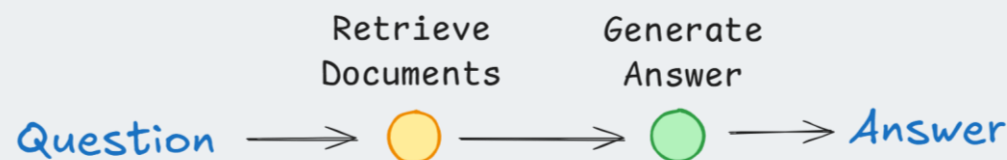
Let's see some examples...

Some more details about RAG...

RAG, one acronym, many ways: the very simple one

The bare minimum for RAG:

- Retrieve Documents
 - Query the database, in Oracle Database 23ai using VECTOR_DISTANCE, to get the closest pieces of content by vector distance to the question asked. The question is embedded into a vector "on the fly".
- Generate Answer
 - Ask a LLM to generate an answer to the question using the provided documents.



 SQL

 LLM

 DATAlysis

RAG, one acronym, many ways: the improved simple one

Adding an extra steps to the bare minimum to try to improve the result:


- Ranker

- Retrieving documents is good, filtering the documents by the most relevant documents is better.
- Retrieve more documents, rank them by relevance and take a shorter subset sorted by the relevance ranking (from most relevant to less relevant).
For example retrieve 10 documents and take a TOP5 from the Ranker to pass to the LLM.
Allows to remove "false positive" results of the retrieval step.



 SQL

 LLM

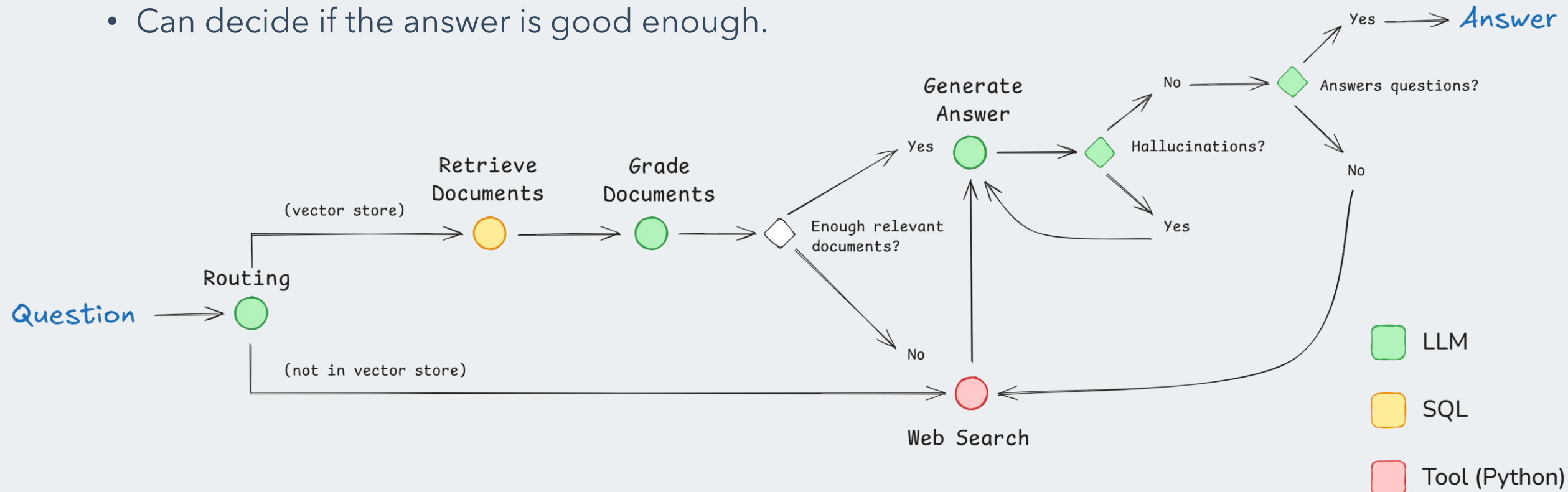
 Ranker model

 DATAanalysis

RAG, one acronym, many ways: the less simple but hopefully better

Your imagination is the limit:

- A LLM can do many things when instructed correctly
 - Can decide if the vector store is able to answer the question.
 - Can decide if documents are relevant for the question.
 - Can decide if the answer sounds plausible or not.
 - Can decide if the answer is good enough.



Back to the Oracle Database ChatGPT-like solution

Recipe for a ChatGPT-like chatbot for Oracle Database

Ingredients:

- Data, documents, some source content
- Oracle Database 23ai
- An embedding model
 - In the database
 - Or a service somewhere called by the database
 - OCI GenAI, Cohere, GoogleAI, HuggingFace, OpenAI or VertexAI
 - Or some solution outside the database (Python, etc.)
- A LLM service somewhere
 - A service somewhere called by the database
 - OCI GenAI, Cohere, GoogleAI, HuggingFace, OpenAI or VertexAI
 - Or some solution outside the database (Python, Ollama etc.)
- Some code glueing together the pieces

A quick look at Oracle OCI GenAI

Generative AI

Overview

Playground

Chat

Embedding

Dedicated AI clusters

Custom models

Endpoints

Scope

Compartment

datanalysis (root) ▾

Generative AI overview

Power your apps with large language models and generative AI

OCI Generative AI is a fully managed service that provides a set of state-of-the-art, customizable LLMs that cover a wide range of use cases for text generation. Use the playground to try out the models out-of-the-box or create and host your own fine-tuned custom models based on your own data on dedicated AI clusters.



Watch service tour

Metrics in datanalysis (root) Compartment

Dedicated AI clusters

0

Custom models

0

Endpoints

0

Resources

- Documentation
- Generative AI Service Management API
- Generative AI Service Inference API

Get started



Playground

The playground is a visual interface for exploring the hosted pretrained and custom models without writing a single line of code. Use the playground to test your use cases and refine prompts and parameters. When you're happy with the results, you can view the code and integrate Generative AI into your applications.

Go to playground



Dedicated AI clusters

Spin up dedicated hardware units for fine-tuning custom models and hosting



Custom models

Create custom models by fine-tuning the base models with your



Endpoints

Create and manage endpoints to host your custom models.



OCI GenAI

OCI Generative AI combines various services into one in Oracle Cloud.

- On-demand, pre-trained, models for embedding and chat LLMs
- Dedicated AI clusters (\$\$\$)
- Custom models: from one of the available pre-trained model, fine-tune it with your own dataset

Available only in a limited number of regions, and not all models are available everywhere: <https://docs.oracle.com/en-us/iaas/Content/generative-ai/pretrained-models.htm>

Cheap and easily accessible for small tasks, calculate you potential costs for large activities (just to be aware of the cost involved).

OCI GenAI: Embedding

ORACLE Cloud

Cloud Classic >

Search resources, services, documentation, and Marketplace

Germany Central (Frankfurt) v

<>

🔔

?

🌐

👤

Generative AI

Overview

Playground

Chat

Embedding

Dedicated AI clusters

Custom models

Endpoints

Scope

Compartment

datanalysis (root)

Embedding ?

To get started, choose a model and a preset prompt example. Then, refine the prompts and parameters to fit your use cases. See [model types](#) for more information. All model responses have [moderation filtering](#) applied for explicit content. Note that some models have deprecation/retirement dates. View our [model list](#) for more details.

Model

cohere.embed-english-v3.0

View model details

Example

Billing communications

View code

Parameters

Truncate ⓘ

None

✓ Sentence input ⓘ

Add a list of sentences or phrases to generate [embeddings](#) (maximum of 96 inputs).

+ Add sentence

Upload file

1. In order to maintain our growth, we need to track our billings to ensure we are charging our customer

2. We have a system in place to track our billings and ensure we are billing our customers accurately.

3. We have a dedicated billing team that is responsible for generating invoices and tracking payments.

4. Our billing system is integrated with our customer relationship management (CRM) system, which al

5. We use a third-party billing service to help us manage our billings and ensure we are billing our cust

6. We are committed to providing our customers with accurate billings and clear explanations of our ch

7. Timely and accurate billing is important to our customers, and we strive to provide them with the bes

8. We are constantly looking for ways to improve our billing process and ensure we are billing our cust

9. We are committed to being transparent with our customers about our billing process and how we cal

10. Billing can be a complex process, and we are here to help our customers understand their bills and i

11. We value our customers and want to ensure that they are happy with our billing process and the ser

Run

Clear

Character count - 1329 | Token limit for each input - 512

📱

OCI GenAI: Embedding

ORACLE Cloud

Cloud Classic >

Search resources, services, documentation, and Marketplace

Germany Central (Frankfurt) v

5. We use a third-party billing service to help us manage our billings and ensure we are billing our cust

×

6. We are committed to providing our customers with accurate billings and clear explanations of our ch

×

7. Timely and accurate billing is important to our customers, and we strive to provide them with the bes

×

8. We are constantly looking for ways to improve our billing process and ensure we are billing our cust

×

9. We are committed to being transparent with our customers about our billing process and how we cal

×

10. Billing can be a complex process, and we are here to help our customers understand their bills and i

×

11. We value our customers and want to ensure that they are happy with our billing process and the ser

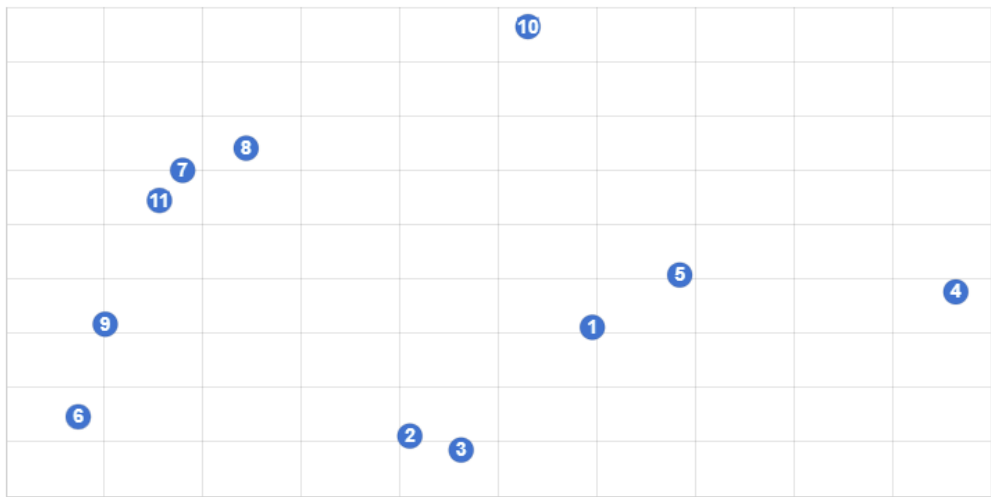
×

Run

Clear

Character count - 1329 | Token limit for each input - 512

Output vector projection ⓘ



Export embeddings to JSON

 @G_Ceresa

Terms of Use and Privacy Cookie preferences unavailable

Copyright © 2024, Oracle and/or its affiliates. All rights reserved.



- ▼ Pretrained Foundational Models
 - About the Chat Models
 - About the Generation Models
 - About the Summarization Models
 - About the Embedding Models
- Getting Started
 - ▶ Getting Access
 - ▶ Using the Large Language Models (LLMs)
 - ▶ Cluster Performance Benchmarks
 - ▶ Managing Dedicated AI Clusters
 - ▶ Fine-Tuning the Base Models
 - ▶ Managing the Custom Models
 - ▶ Managing an Endpoint
 - ▶ Integrating the Models
 - Model Limitations
 - ▶ Calculating Cost
 - Service Limits
 - ▶ Metrics
 - Retiring the Models

— Embedding Models

Convert text to vector embeddings to use in applications for semantic searches, text classification, or text clustering.

Model	Available in These Regions	Key Features
<code>cohere.embed-english-v3.0</code>	<ul style="list-style-type: none">• Brazil East (Sao Paulo)• Germany Central (Frankfurt)• UK South (London)• US Midwest (Chicago)	<ul style="list-style-type: none">• English or multilingual ↗.• Model creates a 1024-dimensional vector for each embedding.• Maximum 96 sentences per run.• Maximum 512 tokens per embedding.
<code>cohere.embed-multilingual-v3.0</code>	<ul style="list-style-type: none">• Brazil East (Sao Paulo)• Germany Central (Frankfurt)• UK South (London)• US Midwest (Chicago)	<ul style="list-style-type: none">• English or multilingual ↗.• Model creates a 1024-dimensional vector for each embedding.• Maximum 96 sentences per run.• Maximum 512 tokens per embedding.
<code>cohere.embed-english-light-v3.0</code>	<ul style="list-style-type: none">• US Midwest (Chicago)	<ul style="list-style-type: none">• Light models are smaller and faster than the original models.• English or multilingual ↗.• Model creates a 384-dimensional vector for each embedding.• Maximum 96 sentences per run.• Maximum 512 tokens per embedding.
<code>cohere.embed-multilingual-light-v3.0</code>	<ul style="list-style-type: none">• US Midwest (Chicago)	<ul style="list-style-type: none">• Light models are smaller and faster than the original models.

↕ Expand All Expandable Areas

Was this article helpful?








Updated 2024-08-08

ORACLE Cloud

Cloud Classic >

Search resources, services, documentation, and Marketplace

Germany Central (Frankfurt) v



Generative AI

Overview

Playground

Chat

Embedding

Dedicated AI clusters

Custom models

Endpoints

Scope

Compartment

datalysis (root)

Chat ?

To get started, choose a model and a preset prompt example. Then, refine the prompts and parameters to fit your use cases. See [model types](#) for more information. All model responses have [moderation filtering](#) applied for explicit content. Note that some models have deprecation/retirement dates. View our [model list](#) for more details.

Model

cohere.command-r-plus ...


View model details

Example

Generate a job description

View code

Generate a job description for a data visualization expert with the following three qualifications only:
1) At least 5 years of data visualization expert
2) A great eye for details
3) Ability to create original visualizations

 Job Title: Data Visualization Expert
Job Description:
We are seeking a talented and experienced Data Visualization Expert to join our team. The ideal candidate will have a strong background in data visualization and be able to create original and engaging visualizations that effectively communicate complex data insights.
Qualifications:

1. At least 5 years of experience as a data visualization expert, with a proven track record of creating high-quality, impactful visualizations.

2. A great eye for detail, with the ability to identify and highlight key insights and trends in complex data sets.

3. Ability to create original visualizations: You should be able to think outside the box and develop creative and innovative ways to visualize data, going beyond standard charts and graphs to develop truly engaging

Type a message...

Submit Clear chat

Parameters

Maximum output tokens ⓘ
 600

Preamble override ⓘ


Temperature ⓘ
 0

Top p ⓘ
 0.75

Top k ⓘ
 0


Frequency penalty ⓘ
 1

Presence penalty ⓘ
 0



Copyright © 2024, Oracle and/or its affiliates. All rights reserved.

[Terms of Use and Privacy](#) [Cookie preferences unavailable](#)

 @G_Ceresa

ORACLE Cloud

Cloud Classic >

Search resources, services, documentation, and Marketplace

Germany Central (Frankfurt) v

<>

🔔

?

🌐

👤

Generative AI

Overview

Playground

Chat

Embedding

Dedicated AI clusters

Custom models

Endpoints

Scope

Compartment

datanalysis (root)

Chat ?

To get started, choose a model. Model responses have [model response](#).

Model

meta.llama-3-70b-instruct

Generate a product recommendation for surround sound. The user wants to know about the latest conversations. The user is interested in the latest conversations.

Introducing the new OCI Generative AI Chat with Immersive AI. Are you tired of using traditional chatbots? You want to create a chatbot with unparalleled performance and game-changing capabilities. The Surround Sound AI Chatbot is the heart of the new AI experience, capturing a 360-degree view of the user's needs, being there in the moment, and providing a nuanced and detailed response. Key Features:

- Compact and easy to use, perfect for integration into your application.

Type a message...

Submit Clear chat

View code

Language

Python

Use this sample code to integrate this model into your application. Ensure that you update this snippet with the [required authentication](#).

```
# coding: utf-8
# Copyright (c) 2023, Oracle and/or its affiliates. All rights reserved.
# This software is dual-licensed to you under the Universal Permissive License (UPL) 1.0 as shown at https://www.oracle.com/licenses/ucpl/

#####
# chat_demo.py
# Supports Python 3
#####
# Info:
# Get texts from LLM model for given prompts using OCI Generative AI Service.
#####
# Application Command Line(no parameter needed)
# python chat_demo.py
#####
import oci

# Setup basic variables
# Auth Config
# TODO: Please update config profile name and use the compartmentId that has policies grant permissions for user
compartment_id = "ocid1.tenancy.oc1..aaaaaaaaxkrfbzdi72pukgwprkjer2lro6l7k5kmcips7syvku54mkkv4boq"
CONFIG_PROFILE = "DEFAULT"
config = oci.config.from_file('~/.oci/config', CONFIG_PROFILE)

# Service endpoint
endpoint = "https://inference.generativeai.eu-frankfurt-1.oci.oraclecloud.com"

generative_ai_inference_client = oci.generative_ai_inference.GenerativeAiInferenceClient(config=config, service_endpoint=endpoint)
chat_detail = oci.generative_ai_inference.models.ChatDetails()
```

Close

Java also
available

▼ Pretrained Foundational Models

About the Chat Models

About the Generation Models

About the Summarization Models

About the Embedding Models

Getting Started

► Getting Access

► Using the Large Language Models (LLMs)

► Cluster Performance Benchmarks

► Managing Dedicated AI Clusters

► Fine-Tuning the Base Models

► Managing the Custom Models

► Managing an Endpoint

► Integrating the Models

Model Limitations

► Calculating Cost

Service Limits

► Metrics

Retiring the Models

— Chat Models (New)

Ask questions and get conversational responses through an AI chatbot.

Model	Available in These Regions	Key Features
<code>cohere.command-r-plus v1.2</code>	<ul style="list-style-type: none">• Brazil East (Sao Paulo)• Germany Central (Frankfurt)• UK South (London)• US Midwest (Chicago)	<ul style="list-style-type: none">• User prompt can be up to 128,000 tokens, and response can be up to 4000 tokens for each run.• Optimized for complex tasks, offers advanced language understanding, higher capacity, and more nuanced responses, and can maintain context from its long conversation history of 128,000 tokens. Also ideal for question-answering, sentiment analysis, and information retrieval.
<code>cohere.command-r-16k v1.2</code>	<ul style="list-style-type: none">• Brazil East (Sao Paulo)• Germany Central (Frankfurt)• UK South (London)• US Midwest (Chicago)	<ul style="list-style-type: none">• User prompt can be up to 16,000 tokens, and response can be up to 4000 tokens for each run.• Optimized for conversational interaction and long context tasks. Ideal for text generation, summarization, translation, or text-based classification.• You can fine-tune this model with your dataset.
<code>meta.llama-3-70b-instruct v1.0</code>	<ul style="list-style-type: none">• Brazil East (Sao Paulo)• Germany Central (Frankfurt)• UK South (London)	<ul style="list-style-type: none">• Model has 70 billion parameters.• User prompt and response can be up to 8000 tokens for each run.• You can fine-tune this model with your dataset.

↕ Expand All Expandable Areas

Was this article helpful?

👍

👎

Updated 2024-08-08

Llama 3.1 is now available with a larger prompt.

outside of EU Llama 3.2 is also available.

Let's start!

Everything begins with data

OracleBaseGPT

When googling something about the database, you often find oracle-base.com in the first results.

You maybe never met Tim Hall, but you know his website.

What happen if we take a chatbot powered by LLM and feed him with oracle-base.com articles?

Can we get a personal OracleBaseGPT?

[illegible]

But before...

Do not randomly steal content online from somebody's website without being allowed to do so: you can maybe navigate the website, but you don't own the content, you can't always do whatever you want with it.

Do not kill websites by sending tons of requests just because you are playing with some embedding exercise.

Be respectful, ask permission.

This presentation is done with the permission of Tim Hall to use his articles.



Data: know your source

Any ML, AI, LLM usage will generally follow the same rule:

- "garbage in, garbage out"

You must know your source, its format, and type.

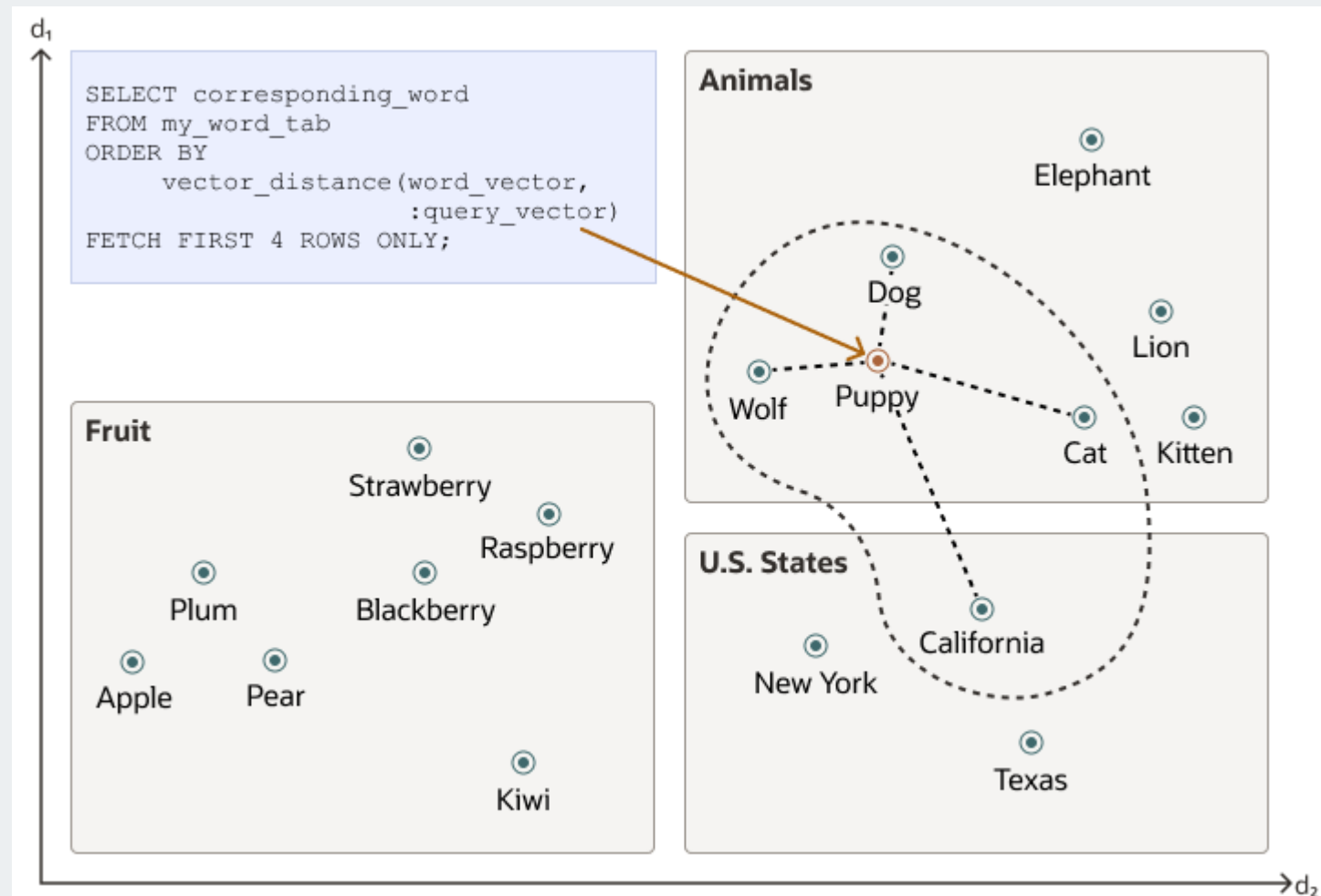
Data: know your source

This is even more important for RAG activities:

- Generate an embedding of garbage content and you will get a fairly pointless embedding.

Reality isn't as simple as the images in the Oracle Database documentation...

Fruit, Animals or U.S.
States all nicely packed
together is an illustration
of an idealistic vision of
how vectors looks like



Data: know your source

Are header and footer useful? Not really...

Identify the real content and how to get rid of the "noise" around it.

The screenshot shows a web page from Oracle Base. The header includes the Oracle Base logo, navigation links (Articles, Scripts, Blog, Certification, Videos, Misc, Printer Friendly, About), and search bars. The main content area is titled "Pipelined Table Functions" and contains a list of links: "Table Functions" and "Pipelined Table Functions". A social media bar with icons for X, Facebook, LinkedIn, Reddit, and YouTube is visible. The page is annotated with a torn paper effect, revealing HTML and CSS code. The HTML code shows the header structure, including the title and navigation links. The CSS code shows the styling for the header, including the font size and margin-top. The footer code shows the "Back to the Top" link.

Annotations include:

- Header: `<h1>Pipelined Table Functions</h1> == $0`
- Table of Contents:
 - Table Functions
 - Pipelined Table Functions
- Footer: `Back to the Top.`
- Footer: `<p>Hope this helps. Regards Tim...</p>`

Start by chunking

Chunking: where the fun really starts...

The embedding of your content will generate a single vector for each input you provide.

Should you provide a whole document as one piece?

You will get a single vector for it! How precisely could a single vector represent a whole document and every concept contained in it?

Should you provide each word as a single piece?

You will get a vector for each word, but no context at all because your words will be individual, separate, pieces.

Finding the best way to chunk your documents can be very challenging!

- Chunking: splitting your source document into smaller pieces

Chunking: where the fun really starts...

If you search online you will find about as many chunking strategies as people talking about it...

Just to name a few simple techniques:

- Fixed size chunking
- Recursive chunking
- Structure based chunking

Chunking: where the fun really starts...

Fixed size chunking is the easiest. Chunk your content into pieces of fixed size.

An example:

- Chunk size: 100 characters
- Overlap: 0 characters

Many words cut in half,
sentences cut in half,
can we safely assume
the meaning is preserved?

A tablescan would, of course, still be expensive and do far more work than needed. Ideally we need an index to be able to find the rows that match all three equality predicates. However, we can see that there are roughly 2 rows per value for item_key, so an index on just (item_key) might be good enough. Averages, of course, can be very misleading: we might have one row for each of 3.5M values of item_key and one value of item_key with 3.5M rows in our table, so we do need to know more about the data distribution before we can make any solid suggestions.

<https://jonathanlewis.wordpress.com/2024/08/07/indexing-4/>

Chunking: where the fun really starts...

Recursive chunking is a bit smarter. The content is split based on a set of separators. If the obtained chunk is bigger than the expected size, another chunking loop is performed on that piece till the result is "small enough".

Separators can be "\n", ":", ",", " " etc.

An example:

- Chunk size: 100 characters
- Overlap: 0 characters

No word cut in half,
sentences still cut in half,
can we safely assume
the meaning is preserved?

A tablescan would, of course, still be expensive and do far more work than needed. Ideally we need an index to be able to find the rows that match all three equality predicates. However, we can see that there are roughly 2 rows per value for item_key, so an index on just (item_key) might be good enough. Averages, of course, can be very misleading: we might have one row for each of 3.5M values of item_key and one value of item_key with 3.5M rows in our table, so we do need to know more about the data distribution before we can make any solid suggestions.

<https://jonathanlewis.wordpress.com/2024/08/07/indexing-4/>

Chunking: where the fun really starts...

Structure based chunking works with the structure of the document. HTML has a structure, paragraphs, headings should be used (but you can abuse HTML massively and break and structured approach).

You can chunk by section of the HTML by headings tags: `<h1>`, `<h2>`, `<h3>`

You can then chunk by paragraph, knowing to which section it belongs. Maybe even “enriching” the paragraph by adding the heading to it.

The issue is that nothing guarantee the size of a paragraph, it could still be very long.

Chunking methods can be combined: structure based first, recursive chunking after on the result of the structure chunking.

After chunking comes embedding, turning inputs into vectors

Embedding in the database

- Embedding in the database requires to first load a model able to perform embedding. Oracle provides some pre-trained models in the ONNX format that can be loaded.

```
begin
  dbms_vector.drop_onnx_model (
    model_name => 'ALL_MINILM_L12_V2',
    force => true);

  dbms_vector.load_onnx_model (
    directory   => 'model_dir',
    file_name   => 'all_Minilm_L12_v2.onnx',
    model_name  => 'ALL_MINILM_L12_V2');
end;
/
```

```
select
vector_embedding(all_minilm_l12_v2 using 'Quick test' as data) AS my_vector;
```

Embedding from the database using a webservice

The DBMS_VECTOR and DBMS_VECTOR_CHAIN packages can call external webservice for embedding.

You first need to create a credential for the service with the required authentication details.

You also need to allow to call the external URL.

```
-- execute an embedding task for OCI GenAI
declare
  input clob;
  params clob;
  v vector;
begin
  dbms_output.put_line('Embedding task for OCI GenAI');
  input := 'Hello world';
  params := '
{
  "provider": "OCIGenAI",
  "credential_name": "CRED_OCI_GENAI",
  "url": "https://inference.generativeai.eu-frankfurt-1.oci...",
  "model": "cohere.embed-english-v3.0",
  "batch_size": 1
}';

  v := dbms_vector.utl_to_embedding(input, json(params));
  dbms_output.put_line(vector_serialize(v));
exception
  when OTHERS THEN
    DBMS_OUTPUT.PUT_LINE (SQLERRM);
    DBMS_OUTPUT.PUT_LINE (SQLCODE);
end;
/
```

How I do chunking and embedding for OracleBaseGPT

Chunking outside the database

HTML splitter first,
recursive splitter after.

Using LangChain and
Hugging Face
packages.

Trying to chunk based
on token size using a
tokenizer for the
embedding model
I'm going to use.

```
### Does still create chunks too big for the embedding model (512 tokens) from time to time...

from langchain_text_splitters import HTMLHeaderTextSplitter, RecursiveCharacterTextSplitter
from langchain_core.documents import Document
from transformers import AutoTokenizer

# define a tokenizer for Cohere-embed-english-v3.0 to try to count up to 512 tokens per chunk
model_id = "Cohere/Cohere-embed-english-v3.0"
tokenizer = AutoTokenizer.from_pretrained(model_id)

# Recursive splitter based on chunk size counted in tokens
text_splitter = RecursiveCharacterTextSplitter.from_huggingface_tokenizer(
    tokenizer, chunk_size=400, chunk_overlap=10
)

# List of headers to split on
headers_to_split_on = [
    ("h1", "Header 1"),
    ("h2", "Header 2"),
    ("h3", "Header 3"),
]

# HTML splitter to split by the structure of the web page
html_splitter = HTMLHeaderTextSplitter(headers_to_split_on=headers_to_split_on)

def chunk_html_content(id: int, url: str, content: str) -> list[Document]:
    # html split on headers
    html_header_splits = html_splitter.split_text(content)

    # split chunks
    chunks = text_splitter.split_documents(html_header_splits)

    # dirty add of metadata
    for idx, c in enumerate(chunks):
        chunks[idx].metadata = c.metadata | {'source': url, 'id': id}

    return chunks
```


Embedding outside the database

Embedding is done right after chunking, in Python. Using the `cohere.embed-english-v3.0` model of OCI GenAI.

A parameter is set to truncate inputs if too large (512 tokens max) and not stop with an error.

The result is a list of vectors.

```
import oci

# setup basic variables
compartment_id = "ocid1.compartment.oc1..aaaaaaaaryltzreesv2f4pa6vj4lnf52yemrtywqo5pcwse3pknps7rwk3fq"
oci_config = oci.config.from_file('~/.oci/config', "DEFAULT")

# service endpoint
endpoint = "https://inference.generativeai.eu-frankfurt-1.oci.oraclecloud.com"

# define the GenAI client to execute embedding calls
generative_ai_inference_client = oci.generative_ai_inference.GenerativeAiInferenceClient(
    config=oci_config,
    service_endpoint=endpoint,
    retry_strategy=oci.retry.NoneRetryStrategy(),
    timeout=(10,240)
)

def embed_list_of_chunks(chunks: list) -> list:
    inputs = [chunk.page_content for chunk in chunks]

    embed_text_detail = oci.generative_ai_inference.models.EmbedTextDetails()
    embed_text_detail.serving_mode = oci.generative_ai_inference.models.OnDemandServingMode(model_id="cohere.embed-english-v3.0")
    embed_text_detail.inputs = inputs
    embed_text_detail.truncate = "END" # if the chunk is too big, truncate the end
    embed_text_detail.compartment_id = compartment_id
    embed_text_response = generative_ai_inference_client.embed_text(embed_text_detail)

    return embed_text_response.data.embeddings
```

Chunks and vectors are inserted in Oracle Database 23ai

An index can be added on the vector column to make queries faster.

The dataset is quite small: only 9'608 chunks (for the 1'327 articles of oracle-base.com).

The embeddings column is of type VECTOR(1024, FLOAT64)

```
1 select * from webpages_chunks order by 1, 2;
```


Fetches 200 rows in 0.815 seconds

	ID	WEBPAGE_ID	URL	CONTENT	EMBEDDING
1	31216	125	https://oracle-base.com/artic...	Oracle Data Pump was introduced in Oracle 10g. This article provi...	[-2.18505859E-002,-2.48260498E-002,-4.17480469E-002,-1.83105469E-002,-2.272
2	31217	125	https://oracle-base.com/artic...	The COMPRESSION parameter allows you to decide what, if anything,...	[-2.0111084E-002,-2.35900879E-002,5.37872314E-003,-2.11334229E-002,-4.21752
3	31218	125	https://oracle-base.com/artic...	Data pump encryption is an Enterprise Edition feature, so the par...	[-4.01916504E-002,-7.07626343E-003,-2.66113281E-002,1.95922852E-002,-2.5924
4	31219	125	https://oracle-base.com/artic...	The use of encryption is controlled by a combination of the ENCRY...	[-7.51113892E-003,-2.91595459E-002,1.11618042E-002,-1.16729736E-002,-5.0628
5	31220	125	https://oracle-base.com/artic...	The ENCRYPTION_ALGORITHM parameter specifies the encryption algor...	[-4.67834473E-002,-1.3168335E-002,5.68771362E-003,5.30395508E-002,-7.049560
6	31221	125	https://oracle-base.com/artic...	The ENCRYPTION_MODE parameter specifies the type of security used...	[-2.13928223E-002,-3.49121094E-002,2.32849121E-002,-1.67236328E-002,-5.6030
7	31222	125	https://oracle-base.com/artic...	The TRANSPORTABLE parameter is similar to the TRANSPORT_TABLESPAC...	[6.23321533E-003,-2.16674805E-002,3.26538086E-002,-1.48391724E-003,-3.26538
8	31223	125	https://oracle-base.com/artic...	The PARTITION_OPTIONS parameter determines how partitions will be...	[1.53808594E-002,-7.2555542E-003,2.14538574E-002,-3.96118164E-002,-4.306030
9	31224	125	https://oracle-base.com/artic...	The REUSE_DUMPFILES parameter can be used to prevent errors being...	[4.22668457E-003,2.9296875E-003,4.97436523E-002,-2.38647461E-002,-4.8828125
10	31225	125	https://oracle-base.com/artic...	This parameter allows a table to be renamed during the import ope...	[-1.8157959E-002,1.75323486E-002,-7.63702393E-003,-4.21905518E-003,-1.86309
11	31226	125	https://oracle-base.com/artic...	During import operations using the external table acces method, s...	[1.18865967E-002,-1.35269165E-002,-2.46734619E-002,-8.09326172E-002,1.00326
12	31227	125	https://oracle-base.com/artic...	During an export, if XMLTYPE columns are currently stored as CLOB...	[-1.30767822E-002,-5.0994873E-002,1.54571533E-002,-1.42288208E-002,-5.72509
13	31228	125	https://oracle-base.com/artic...	During export and import operations, the REMAP_DATA parameter all...	[-2.4520874E-002,-1.00021362E-002,1.50909424E-002,-3.71360779E-003,-3.53393
14	31229	125	https://oracle-base.com/artic...	Worker processes that have stopped due to certain errors will now...	[1.63421631E-002,-3.43933105E-002,-2.39715576E-002,-5.00183105E-002,-1.9226
15	31230	127	https://oracle-base.com/artic...	The Database Replay functionality of Oracle 11g allows you to cap...	[-3.4942627E-003,-4.54711914E-002,-2.56652832E-002,-9.10644531E-002,-4.9438
16	31231	127	https://oracle-base.com/artic...	The DBMS_WORKLOAD_CAPTURE package provides a set of procedures an...	[2.61077881E-002,-2.30858164E-002,-2.00500488E-002,-8.74632780E-002,-2.0782

Generate answers with a LLM

LLMs are a bit like teaching to kids...


By design a LLM isn't really smart. Many even struggle to perform mathematical operations and logical thinking. And it's expected: a LLM isn't a "brain", it does generate text.

**Elastic**
@elastic


Addressing generative model challenges

Prevent hallucinations by:


- ✓ Using models with citations
- ✓ Enhancing context window utilization
- ✓ Training models specifically for #RAG tasks to improve reliability



Like humans, LLMs make mistakes too. But LLM mistakes might seem more silly because they don't do logic!



Lily Adler
Principal Solutions Architect

 **elastic** | The Search AI Company

9:37 PM · Sep 19, 2024 · 1,139 Views

LLMs are a bit like teaching to kids...

By design a LLM isn't really smart. Many even struggle to perform mathematical operations and logical thinking. And it's expected: a LLM isn't a "brain", it does generate text.

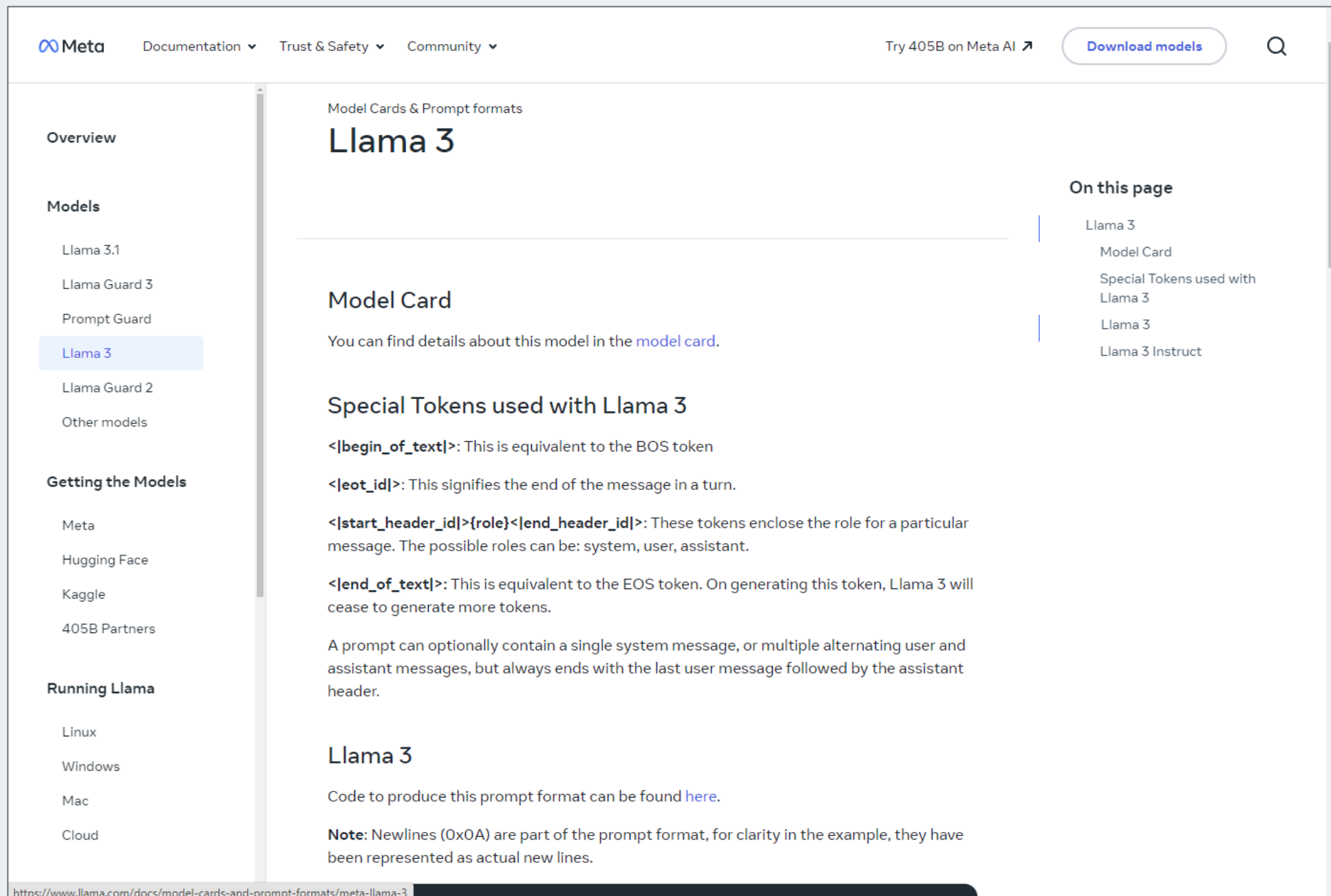
The quality of the prompt is what will decide if the LLM will understand your request or not.

LLMs also don't really have any memory, the memory is provided as part of the prompt (the input message).

Current LLMs in GenAI are chat models, they do have roles for the messages:

- System : provides general instructions to the model, how it should behave
- User : the messages from the users, the "questions"
- Chatbot / Assistant : the messages generated by the model as answer to the inputs received

Models have different prompt formats



The screenshot shows the Meta Llama 3 documentation page. The left sidebar contains navigation links for Overview, Models (with Llama 3 selected), Getting the Models, and Running Llama. The main content area is titled 'Model Cards & Prompt formats' and 'Llama 3'. It includes a 'Model Card' section with a link to the model card, a 'Special Tokens used with Llama 3' section with definitions for <|begin_of_text|>, <|eot_id|>, <|start_header_id|>{role}<|end_header_id|>, and <|end_of_text|>, and a 'Llama 3' section with a link to the prompt format code. A right sidebar titled 'On this page' lists links to Llama 3, Model Card, Special Tokens used with Llama 3, Llama 3, and Llama 3 Instruct. The top navigation bar includes the Meta logo, links to Documentation, Trust & Safety, and Community, a link to 'Try 405B on Meta AI', and a 'Download models' button.

Meta Documentation Trust & Safety Community Try 405B on Meta AI [Download models](#) Q

Overview

Models

- Llama 3.1
- Llama Guard 3
- Prompt Guard
- Llama 3**
- Llama Guard 2
- Other models

Getting the Models

- Meta
- Hugging Face
- Kaggle
- 405B Partners

Running Llama

- Linux
- Windows
- Mac
- Cloud

Model Cards & Prompt formats

Llama 3

Model Card

You can find details about this model in the [model card](#).

Special Tokens used with Llama 3

<|begin_of_text|>: This is equivalent to the BOS token

<|eot_id|>: This signifies the end of the message in a turn.

<|start_header_id|>{role}<|end_header_id|>: These tokens enclose the role for a particular message. The possible roles can be: system, user, assistant.

<|end_of_text|>: This is equivalent to the EOS token. On generating this token, Llama 3 will cease to generate more tokens.

A prompt can optionally contain a single system message, or multiple alternating user and assistant messages, but always ends with the last user message followed by the assistant header.

Llama 3

Code to produce this prompt format can be found [here](#).

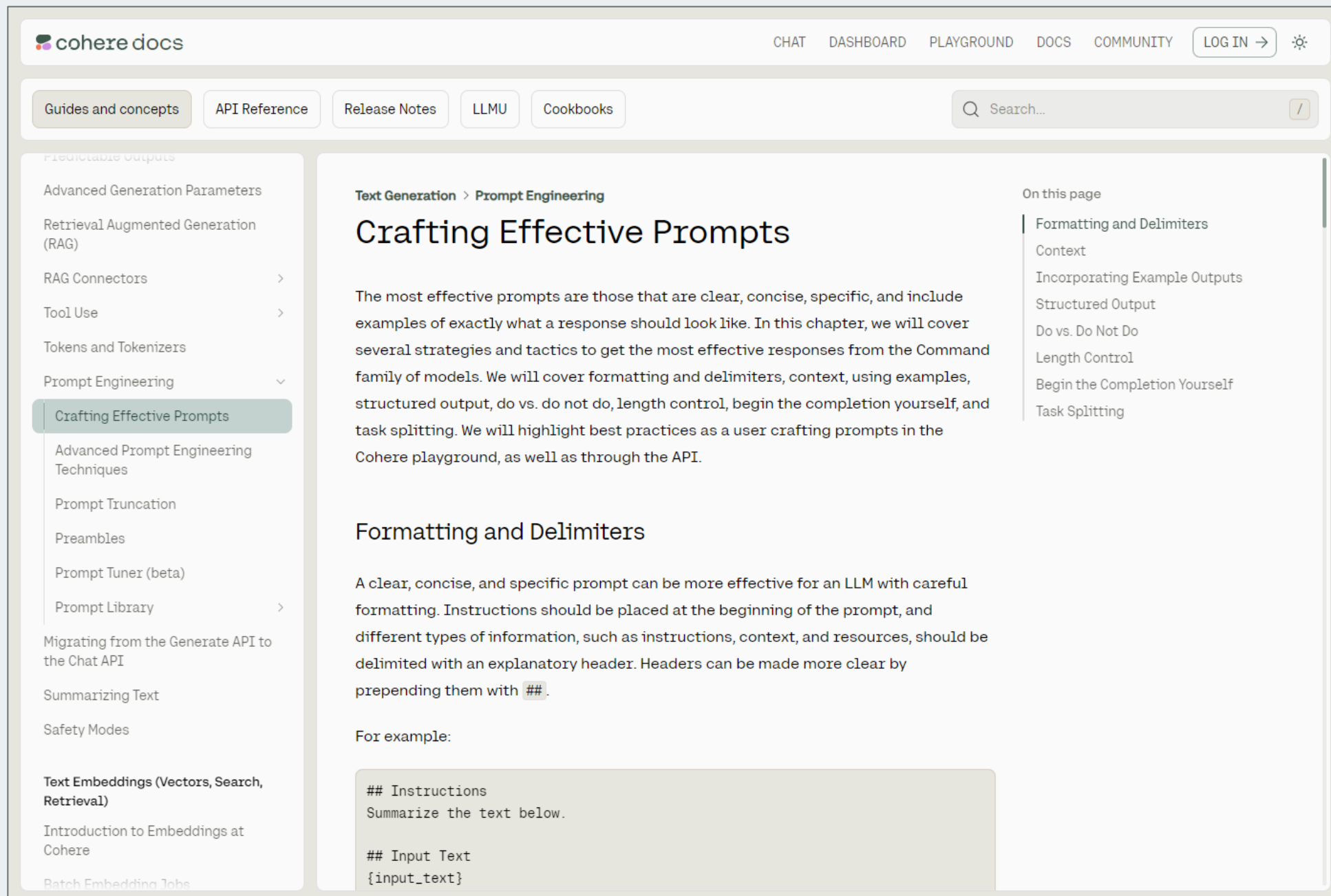
Note: Newlines (0x0A) are part of the prompt format, for clarity in the example, they have been represented as actual new lines.

On this page

- Llama 3
- Model Card
- Special Tokens used with Llama 3
- Llama 3
- Llama 3 Instruct

<https://www.llama.com/docs/model-cards-and-prompt-formats/meta-llama-3>

Models have different prompt formats



The screenshot shows the Cohere Docs website. The top navigation bar includes links for CHAT, DASHBOARD, PLAYGROUND, DOCS, COMMUNITY, and a LOG IN button. Below this is a secondary navigation bar with tabs for Guides and concepts, API Reference, Release Notes, LLMU, and Cookbooks. A search bar is located on the right side of this bar. The left sidebar contains a list of topics under 'Predictable Outputs', with 'Crafting Effective Prompts' highlighted. The main content area is titled 'Text Generation > Prompt Engineering' and 'Crafting Effective Prompts'. It contains an introductory paragraph about effective prompts and a section titled 'Formatting and Delimiters' which explains how to use delimiters like ## to structure prompts. A code block shows an example of a prompt using ## for instructions and input text. On the right side of the main content area, there is a 'On this page' section listing sub-topics like Formatting and Delimiters, Context, and Task Splitting.

cohere docs

CHAT DASHBOARD PLAYGROUND DOCS COMMUNITY LOG IN →

Guides and concepts API Reference Release Notes LLMU Cookbooks

Search...

Predictable Outputs

- Advanced Generation Parameters
- Retrieval Augmented Generation (RAG)
- RAG Connectors >
- Tool Use >
- Tokens and Tokenizers
- Prompt Engineering ▾
 - Crafting Effective Prompts**
 - Advanced Prompt Engineering Techniques
 - Prompt Truncation
 - Preambles
 - Prompt Tuner (beta)
 - Prompt Library >
- Migrating from the Generate API to the Chat API
- Summarizing Text
- Safety Modes
- Text Embeddings (Vectors, Search, Retrieval)
- Introduction to Embeddings at Cohere
- Batch Embedding Jobs

Text Generation > Prompt Engineering

Crafting Effective Prompts

The most effective prompts are those that are clear, concise, specific, and include examples of exactly what a response should look like. In this chapter, we will cover several strategies and tactics to get the most effective responses from the Command family of models. We will cover formatting and delimiters, context, using examples, structured output, do vs. do not do, length control, begin the completion yourself, and task splitting. We will highlight best practices as a user crafting prompts in the Cohere playground, as well as through the API.

Formatting and Delimiters

A clear, concise, and specific prompt can be more effective for an LLM with careful formatting. Instructions should be placed at the beginning of the prompt, and different types of information, such as instructions, context, and resources, should be delimited with an explanatory header. Headers can be made more clear by prepending them with `##`.

For example:



```
## Instructions
Summarize the text below.

## Input Text
{input_text}
```

On this page

- Formatting and Delimiters
- Context
- Incorporating Example Outputs
- Structured Output
- Do vs. Do Not Do
- Length Control
- Begin the Completion Yourself
- Task Splitting

Models can have special features like “documents”

CHAT DASHBOARD PLAYGROUND DOCS COMMUNITY [LOG IN →](#) 

[Guides and concepts](#) [API Reference](#) [Release Notes](#) [LLMU](#) [Cookbooks](#)

Predictable Outputs

- Advanced Generation Parameters
- Retrieval Augmented Generation (RAG)
- RAG Connectors >
- Tool Use >
- Tokens and Tokenizers
- Prompt Engineering
 - Crafting Effective Prompts**
 - Advanced Prompt Engineering Techniques
 - Prompt Truncation
 - Preambles
 - Prompt Tuner (beta)
 - Prompt Library >
- Migrating from the Generate API to the Chat API
- Summarizing Text
- Safety Modes
- Text Embeddings (Vectors, Search, Retrieval)
- Introduction to Embeddings at Cohere
- Batch Embedding Jobs

For the example above, we can split the original news article into different sections and attach them via the `documents` parameter. The Chat API will then provide us not only with the completion but also citations that ground information from the documents. See the following:

```
PYTHON
1 # Sections from the original news article
2 documents = [
3     {"title": "background", "snippet": "From the beginning of the IIHF Ice
4     {"title": "expectations", "snippet": "At the time, the National Hockey
5     {"title": "legacy", "snippet": "While Canada won the series, the Soviet
6 ]
7
8 # New request
9 query = '''The 1972 Canada-USSR Summit Series was an eight-game ice hockey
10
11 # Call the model
12 response = co.chat(
13     message=query,
14     documents=documents,
15     model="command-r-plus",
16     temperature=0.3
17 )
```


The model returns a high quality summary in `response.text`:

The 1972 Canada-USSR Summit Series marked a significant moment in the history of hockey, showcasing a high-stakes competition between the Canadian national team and the elite hockey squad. Here are some key points about the series:

On this page

- Formatting and Delimiters
- Context**
- Incorporating Example Outputs
- Structured Output
- Do vs. Do Not Do
- Length Control
- Begin the Completion Yourself
- Task Splitting

Models can have special features like "documents"

 2.134.0

Search docs

Installation

Configuration

Using FIPS-validated Libraries

Forward Compatibility

New Region Support

Backward Compatibility

Quickstart

Known Issues

Logging

Exception handling

Uploading Large Objects

Raw Requests

Composite Operations and Waiters

Pagination

API Reference

Access Governance Cp
Adm
Ai Anomaly Detection
Ai Document
Ai Language
Ai Speech
Ai Vision
Analytics
Announcements Service
Apigateway
Apm Config

```
class oci.generative_ai_inference.models.CohereChatRequest(**kwargs)
```

Bases: `oci.generative_ai_inference.models.base_chat_request.BaseChatRequest`

Details for the chat request for Cohere models.

Attributes

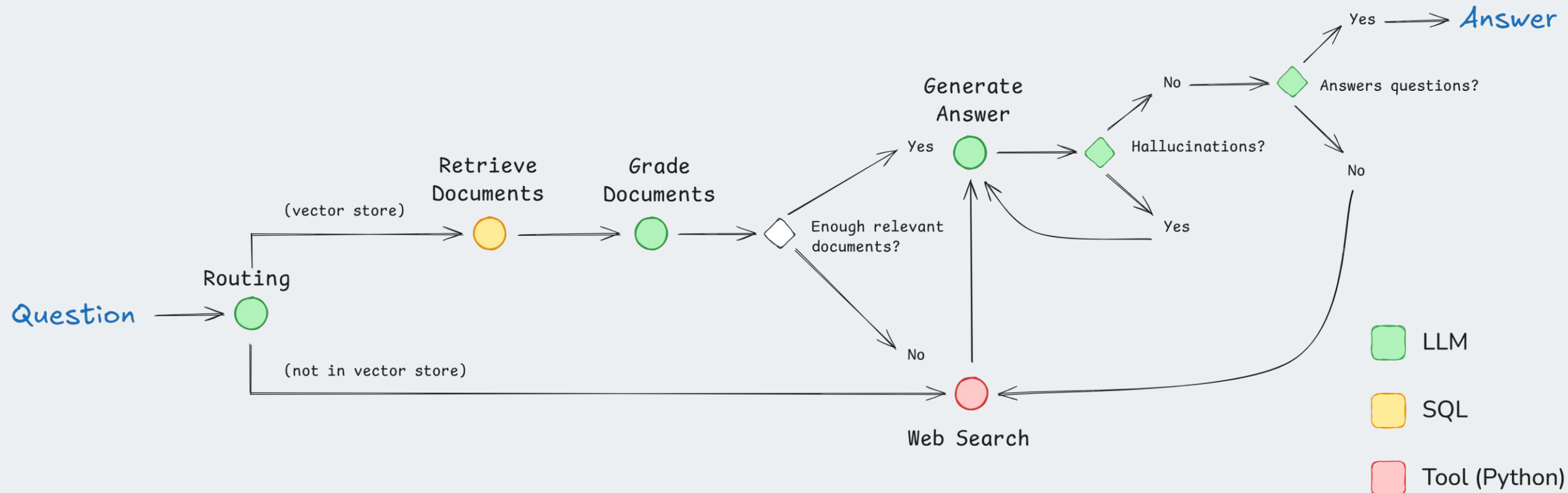
<code>API_FORMAT_COHERE</code>	<code>str(object='') -> str</code>
<code>API_FORMAT_GENERIC</code>	<code>str(object='') -> str</code>
<code>CITATION_QUALITY_ACCURATE</code>	A constant which can be used with the citation_quality proper
<code>CITATION_QUALITY_FAST</code>	A constant which can be used with the citation_quality proper
<code>PROMPT_TRUNCATION_AUTO_PRESERVE_ORDER</code>	A constant which can be used with the prompt_truncation pro
<code>PROMPT_TRUNCATION_OFF</code>	A constant which can be used with the prompt_truncation pro
<code>api_format</code>	[Required] Gets the api_format of this BaseChatRequest.
<code>chat_history</code>	Gets the chat_history of this CohereChatRequest.
<code>citation_quality</code>	Gets the citation_quality of this CohereChatRequest.
<code>documents</code>	Gets the documents of this CohereChatRequest.
<code>frequency_penalty</code>	Gets the frequency_penalty of this CohereChatRequest.
<code>is_echo</code>	Gets the is_echo of this CohereChatRequest.
<code>is_force_single_step</code>	Gets the is_force_single_step of this CohereChatRequest.
<code>is_raw_prompting</code>	Gets the is_raw_prompting of this CohereChatRequest.
<code>is_search_queries_only</code>	Gets the is_search_queries_only of this CohereChatRequest.
<code>is_stream</code>	Gets the is_stream of this CohereChatRequest.
<code>max_tokens</code>	Gets the max_tokens of this CohereChatRequest.
<code>message</code>	[Required] Gets the message of this CohereChatRequest.
<code>preamble_override</code>	Gets the preamble_override of this CohereChatRequest.

Read the documentation of the model creator. Then read the oracle documentation and see if you find the same features.

OracleBaseGPT uses LLM for various tasks

For OracleBaseGPT the model `meta.llama-3.1-70b-instruct` is being used for all the LLMs tasks.

The whole process, the flow of activities, is managed by LangGraph (a Python package designed for exactly that: define flows that involve cycles, essential for most agentic architectures).



OracleBaseGPT: some prompts examples

How to use a LLM to decide if the content of my Oracle Database 23ai table (the chunks and vectors) is going to be able to answer a question or it's better to directly perform a web search?

Being as explicit and precise as possible in the instruction given to the LLM on how it has to behave passing a "system" message.

```
<|begin_of_text|>
<|start_header_id|>system<|end_header_id|> You are an expert at routing a user question to a vectorstore or
web search. Use the vectorstore for questions on Oracle, Oracle database, SQL, PL/SQL, query and
performance tuning. You do not need to be stringent with the keywords in the question related to these topics.
Otherwise, use web-search. Give a binary choice 'web_search' or 'vectorstore' based on the question. Return a
JSON with a single key 'datasource' and no preamble or explanation.
Question to route: {question} <|eot_id|>
<|start_header_id|>assistant<|end_header_id|>
```

For Llama you can see how the last part of the input is always to give the hand back to the LLM to "speak" when using prompts with the proper tokens.

OracleBaseGPT: some prompts examples

Asking the LLM to generate an answer to the question, saying it is an Oracle Database assistant and its answer should be concise and 3 sentences maximum.

```
<|begin_of_text|>
<|start_header_id|>system<|end_header_id|> You are an Oracle Database assistant for question-answering
tasks.
Use the following pieces of retrieved context to answer the question. If you don't know the answer, just say that
you don't know.
Use three sentences maximum and keep the answer concise <|eot_id|>
<|start_header_id|>user<|end_header_id|>
Question: {question}
Context: {context}
Answer: <|eot_id|>
<|start_header_id|>assistant<|end_header_id|>
```

OracleBaseGPT: some prompts examples

Asking the LLM to evaluate if the answer another LLM call did generate is really based on the documents (the data) provided or not. And telling the LLM that the answer should be a JSON with a precise format, no freedom to say anything else.

```
<|begin_of_text|>
<|start_header_id|>system<|end_header_id|> You are a grader assessing whether
an answer is grounded in / supported by a set of facts. Give a binary 'yes' or 'no' score to indicate
whether the answer is grounded in / supported by a set of facts. Provide the binary score as a JSON with a
single key 'score' and no preamble or explanation. <|eot_id|>
<|start_header_id|>user<|end_header_id|> Here are the facts:

-----

{documents}

-----

Here is the answer: {generation} <|eot_id|>
<|start_header_id|>assistant<|end_header_id|>
```

A lot can be done in the database 23ai

**No Python required.
Can be enough for a number of requirements...**

Cleaning, Chunking and Embedding in the database

In Oracle Database 23ai, most tasks related to RAG, vectors, embeddings etc. are packed in 2 packages:

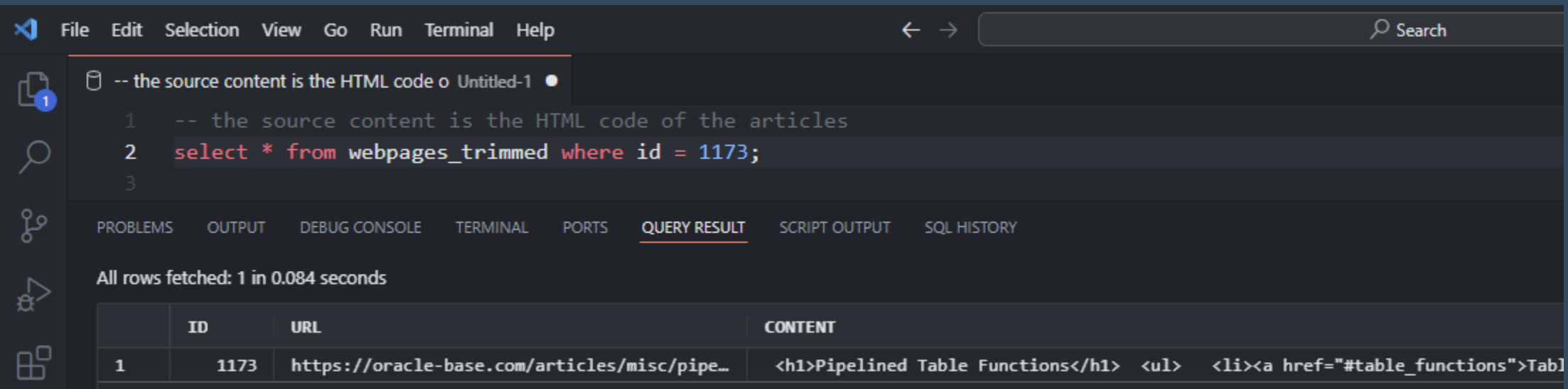
- DBMS_VECTOR
- DBMS_VECTOR_CHAIN

The documentation provides details on the available procedures and the various parameters. There are also some example provided. There seems to easily be mismatches between the documentation and what the current database does (because each release can change things and add new features), be ready for some debugging...

- https://docs.oracle.com/en/database/oracle/oracle-database/23/vecse/summary-dbms_vector-subprograms-arpls.html
- https://docs.oracle.com/en/database/oracle/oracle-database/23/vecse/summary-dbms_vector_chain-subprograms-arpls.html

HTML to text? Yes

Use DBMS_VECTOR_CHAIN.UTL_TO_TEXT

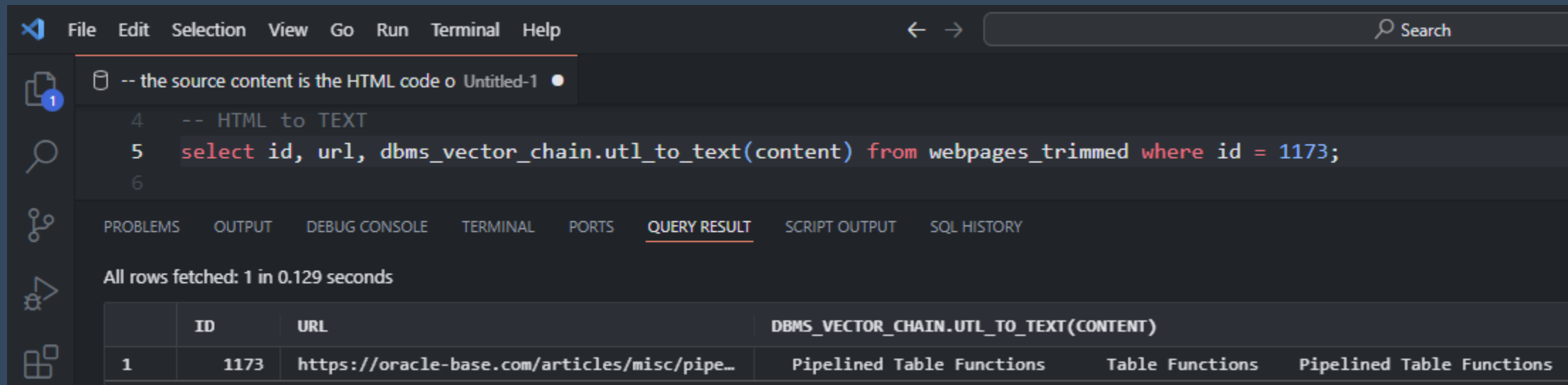


SQL Developer interface showing a query result. The query is:

```
-- the source content is the HTML code of the articles
select * from webpages_trimmed where id = 1173;
```

The result shows 1 row fetched in 0.084 seconds. The table has columns: ID, URL, and CONTENT.

ID	URL	CONTENT
1173	https://oracle-base.com/articles/misc/pipe...	<h1>Pipelined Table Functions</h1> Tabl



SQL Developer interface showing a query result. The query is:

```
-- HTML to TEXT
select id, url, dbms_vector_chain.utl_to_text(content) from webpages_trimmed where id = 1173;
```

The result shows 1 row fetched in 0.129 seconds. The table has columns: ID, URL, and DBMS_VECTOR_CHAIN.UTL_TO_TEXT(CONTENT).

ID	URL	DBMS_VECTOR_CHAIN.UTL_TO_TEXT(CONTENT)
1173	https://oracle-base.com/articles/misc/pipe...	Pipelined Table Functions Table Functions Pipelined Table Functions

HTML to text? Yes

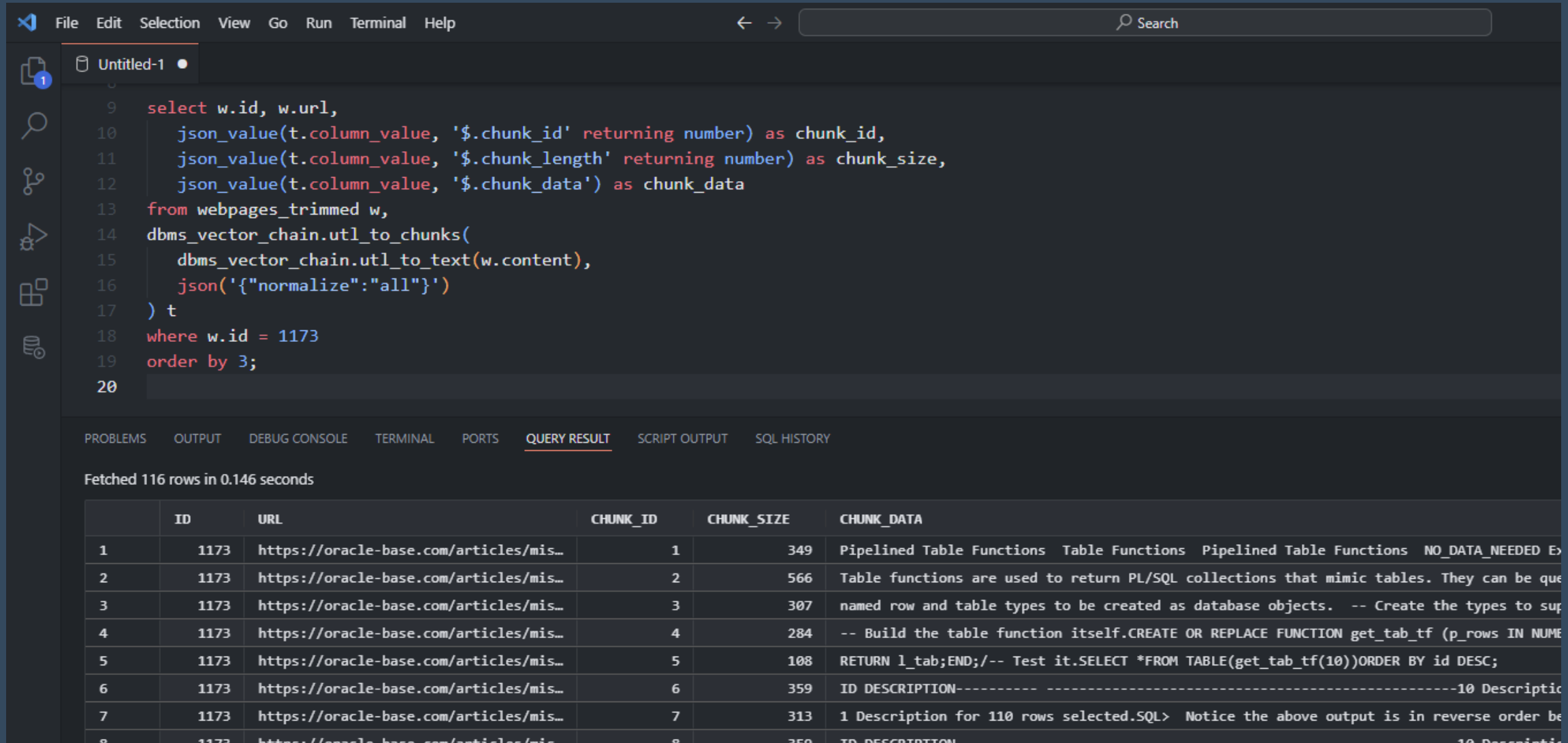
The result isn't the cleanest...

- No cleaning is done
- Lot of new lines added

But it's text, just text, without HTML tags

Chunking? Yes

Use DBMS_VECTOR_CHAIN.UTL_TO_CHUNKS



The screenshot shows an IDE window with a dark theme. The top menu bar includes File, Edit, Selection, View, Go, Run, Terminal, and Help. A search bar is on the right. The left sidebar has icons for Explorer, Search, Source Control, Run and Debug, and Data Explorer. The main editor area displays an SQL query in a file named 'Untitled-1'.

```
9  select w.id, w.url,  
10      json_value(t.column_value, '$.chunk_id' returning number) as chunk_id,  
11      json_value(t.column_value, '$.chunk_length' returning number) as chunk_size,  
12      json_value(t.column_value, '$.chunk_data') as chunk_data  
13  from webpages_trimmed w,  
14      dbms_vector_chain.utl_to_chunks(  
15      dbms_vector_chain.utl_to_text(w.content),  
16      json('{"normalize":"all"}'))  
17  ) t  
18  where w.id = 1173  
19  order by 3;  
20
```

Below the editor, the 'QUERY RESULT' tab is active, showing the results of the query. It indicates that 116 rows were fetched in 0.146 seconds. The results are displayed in a table with the following columns: ID, URL, CHUNK_ID, CHUNK_SIZE, and CHUNK_DATA.

	ID	URL	CHUNK_ID	CHUNK_SIZE	CHUNK_DATA
1	1173	https://oracle-base.com/articles/mis...	1	349	Pipelined Table Functions Table Functions Pipelined Table Functions NO_DATA_NEEDED Ex
2	1173	https://oracle-base.com/articles/mis...	2	566	Table functions are used to return PL/SQL collections that mimic tables. They can be que
3	1173	https://oracle-base.com/articles/mis...	3	307	named row and table types to be created as database objects. -- Create the types to sup
4	1173	https://oracle-base.com/articles/mis...	4	284	-- Build the table function itself.CREATE OR REPLACE FUNCTION get_tab_tf (p_rows IN NUME
5	1173	https://oracle-base.com/articles/mis...	5	108	RETURN l_tab;END;/-- Test it.SELECT *FROM TABLE(get_tab_tf(10))ORDER BY id DESC;
6	1173	https://oracle-base.com/articles/mis...	6	359	ID DESCRIPTION-----10 Descriptio
7	1173	https://oracle-base.com/articles/mis...	7	313	1 Description for 110 rows selected.SQL> Notice the above output is in reverse order be
8	1173	https://oracle-base.com/articles/mis...	8	350	ID DESCRIPTION-----10 Descriptio

Chunking? Yes

DBMS_VECTOR_CHAIN.UTL_TO_CHUNKS has many parameters (all defined in a single JSON parameter).

In theory it can do everything you need.

In practice, does it do what you really want and need?

Partially yes,
fully probably no...

PARAMS

Specify the input parameters in JSON format.

```
{
  "by"           : mode,
  "max"          : max,
  "overlap"      : overlap,
  "split"        : split_condition,
  "custom_list"  : [ split_chars1, ... ],
  "vocabulary"   : vocabulary_name,
  "language"     : nls_language,
  "normalize"    : normalize_mode,
  "norm_options" : [ normalize_option1, ... ],
  "extended"     : boolean
}
```

 Copy

Embedding? Yes

DBMS_VECTOR_CHAIN
has you covered (again)

The “PARAMS” parameter allows you to say “what” does the embedding and how.

211.1.8 UTL_TO_EMBEDDING and UTL_TO_EMBEDDINGS

Use the `DBMS_VECTOR_CHAIN.UTL_TO_EMBEDDING` and `DBMS_VECTOR_CHAIN.UTL_TO_EMBEDDINGS` chainable utility functions to convert plain text to one or more vector embeddings.

Purpose

To perform a text to embedding transformation by accessing:

- Oracle Database as the service provider: Calls the pretrained ONNX format embedding model that you have loaded into the database (default setting)
- Third-party embedding model: Makes a REST call to your chosen third-party service provider, such as Cohere, Google AI, Hugging Face, Oracle Cloud Infrastructure (OCI) Generative AI, OpenAI, or Vertex AI

Syntax

```
DBMS_VECTOR_CHAIN.UTL_TO_EMBEDDING (  
  DATA          IN CLOB,  
  PARAMS         IN JSON default NULL  
) return VECTOR;
```



```
DBMS_VECTOR_CHAIN.UTL_TO_EMBEDDINGS (  
  DATA          IN VECTOR_ARRAY_T,  
  PARAMS         IN JSON default NULL  
) return VECTOR_ARRAY_T;
```



DATA

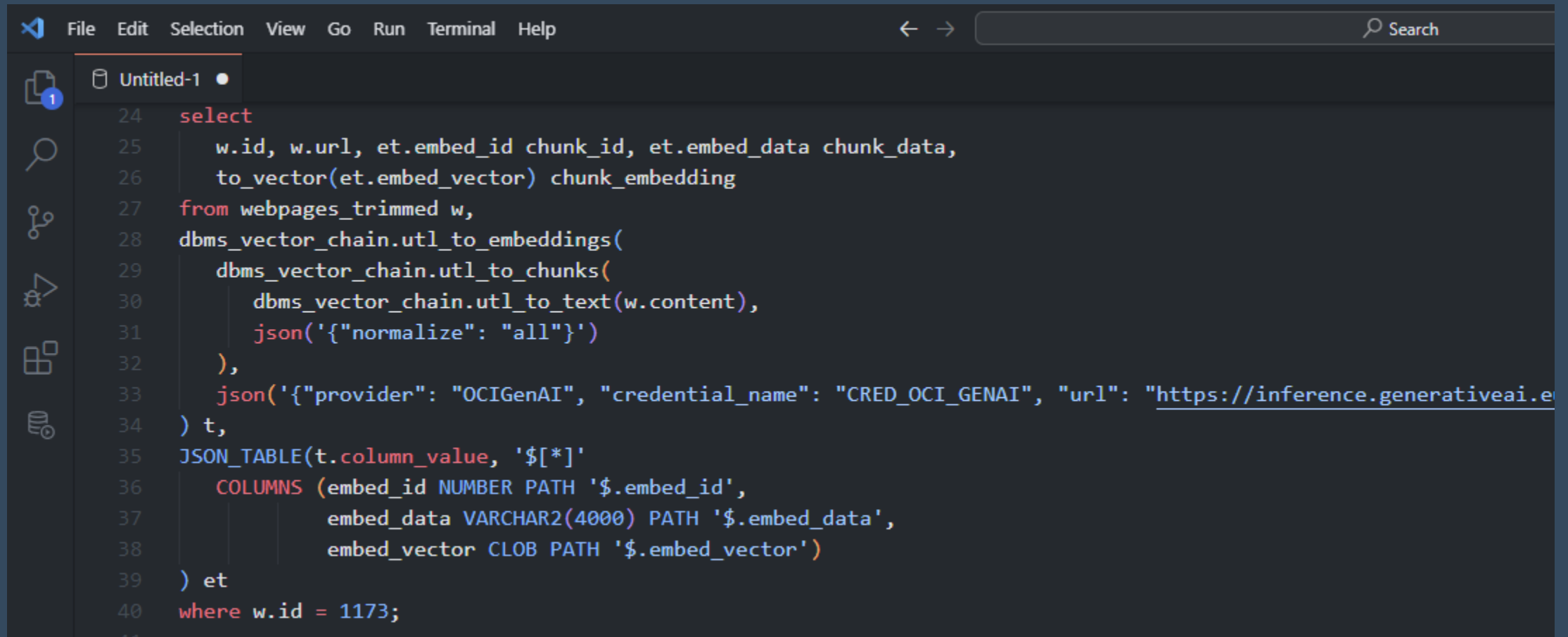
`UTL_TO_EMBEDDING` converts text (CLOB) to a single embedding (VECTOR).

`UTL_TO_EMBEDDINGS` convert an array of chunks (VECTOR_ARRAY_T) to an array of embeddings (VECTOR_ARRAY_T).

All together? Yes

All these steps can be chained together.

The result can be inserted in a table, giving you the chunks and the vectors in a single step.

A screenshot of a code editor window titled 'Untitled-1'. The editor displays an SQL query with line numbers 24 through 41. The query uses Oracle SQL syntax, including the 'JSON_TABLE' function to parse JSON data and the 'JSON_ARRAYAGG' function to aggregate results. It references a table 'webpages_trimmed' and a function 'dbms_vector_chain'. The query filters for 'w.id = 1173'.

```
24 select
25     w.id, w.url, et.embed_id chunk_id, et.embed_data chunk_data,
26     to_vector(et.embed_vector) chunk_embedding
27 from webpages_trimmed w,
28 dbms_vector_chain.utl_to_embeddings(
29     dbms_vector_chain.utl_to_chunks(
30         dbms_vector_chain.utl_to_text(w.content),
31         json('{"normalize": "all"}')
32     ),
33     json('{"provider": "OCI GenAI", "credential_name": "CRED_OCI_GENAI", "url": "https://inference.generativeai.e
34 ) t,
35 JSON_TABLE(t.column_value, '$[*]'
36     COLUMNS (embed_id NUMBER PATH '$.embed_id',
37              embed_data VARCHAR2(4000) PATH '$.embed_data',
38              embed_vector CLOB PATH '$.embed_vector')
39 ) et
40 where w.id = 1173;
```

OracleBaseGPT: demo

OracleBaseGPT (demo)

OracleBaseGPT

Ask Tim questions...

This is a RAG demo based on 1327 articles published on oracle-base.com (up to 1st July 2024).
Oracle Database 23ai is used to store and search vectors. Oracle OCI GenAI is used for embedding, text generation and LLM usage in general.

- **embedding:** model `cohere.embed-english-v3.0` in the Germany Central (Frankfurt) region
- **LLM:** `meta.llama-3-70b-instruct v1.0` in the Germany Central (Frankfurt) region

✓ Connected to Oracle Database 23ai [v23.6.0.24.7]

What is up?



OracleBaseGPT (demo)

OracleBaseGPT

Ask Tim questions...

This is a RAG demo based on 1327 articles published on oracle-base.com (up to 1st July 2024).
Oracle Database 23ai is used to store and search vectors. Oracle OCI GenAI is used for embedding, text generation and LLM usage in general.

- **embedding:** model `cohere.embed-english-v3.0` in the Germany Central (Frankfurt) region
- **LLM:** `meta.llama-3-70b-instruct v1.0` in the Germany Central (Frankfurt) region

✓ Connected to Oracle Database 23ai [v23.6.0.24.7]

What are pipelined table functions?

OracleBaseGPT: Pipelined table functions are a type of table function that allows rows to be piped out of the function as they are created, saving memory and enabling subsequent processing to start before all rows are generated. They include the PIPELINED clause and use the PIPE ROW call to push rows out of the function. This approach negates the need to build huge collections, making it particularly useful for large warehousing ETL operations.

Source:

What is up?



OracleBaseGPT (demo)

generation and LLM usage in general.

- **embedding:** model `cohere.embed-english-v3.0` in the Germany Central (Frankfurt) region
- **LLM:** `meta.llama-3-70b-instruct v1.0` in the Germany Central (Frankfurt) region

✓ Connected to Oracle Database 23ai [v23.6.0.24.7]

What are pipelined table functions?

OracleBaseGPT: Pipelined table functions are a type of table function that allows rows to be piped out of the function as they are created, saving memory and enabling subsequent processing to start before all rows are generated. They include the PIPELINED clause and use the PIPE ROW call to push rows out of the function. This approach negates the need to build huge collections, making it particularly useful for large warehousing ETL operations.

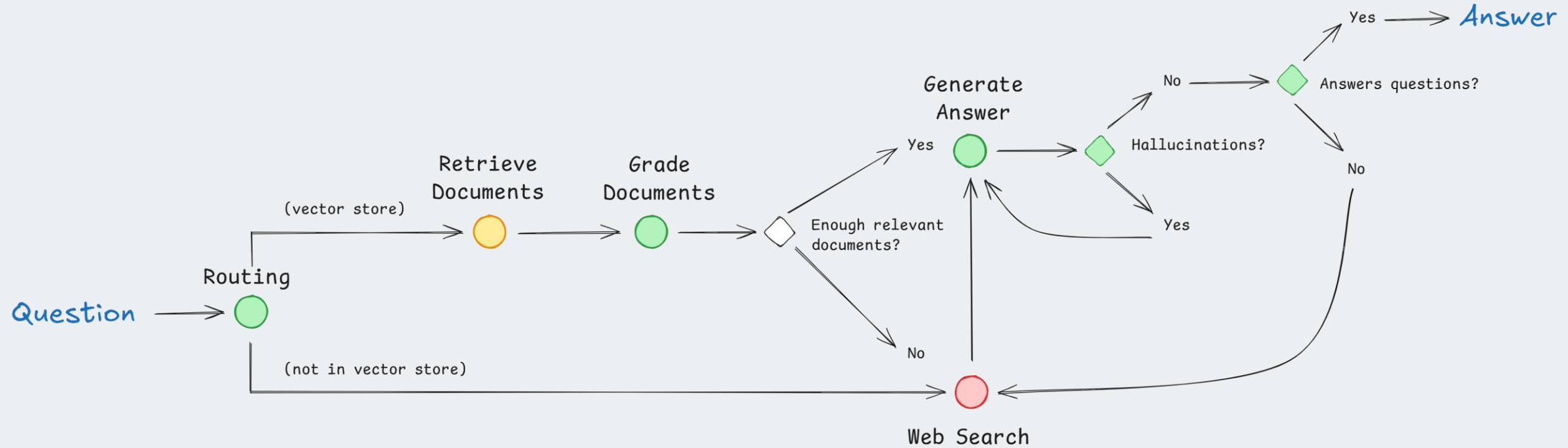
Sources:

<https://oracle-base.com/articles/misc/pipelined-table-functions>
<https://oracle-base.com/articles/12c/plsql-to-sql-interface-enhancements-for-plsql-only-data-types-12cr1>
<https://oracle-base.com/articles/misc/list-files-in-a-directory-from-plsql-and-sql-dbms-scheduler>
<https://oracle-base.com/articles/misc/sql-for-beginners-the-from-clause>

Execution log:

What is up?

OracleBaseGPT (demo)



OracleBaseGPT (demo)

Execution log:

```
[ INFO      ] Successfully connected to Oracle Database (23.6.0.24.7)
[ DEBUG     ] ---ROUTE QUESTION---
[ DEBUG     ] question: What are pipelined table functions?
[ DEBUG     ] route: vectorstore
[ DEBUG     ] ---ROUTE QUESTION TO RAG---
[ DEBUG     ] ---RETRIEVE---
[ DEBUG     ] Finished running: retrieve
[ DEBUG     ] ---CHECK DOCUMENT RELEVANCE TO QUESTION---
[ DEBUG     ] ---GRADE: DOCUMENT RELEVANT---
[ DEBUG     ] ---GRADE: DOCUMENT RELEVANT---
[ DEBUG     ] ---GRADE: DOCUMENT RELEVANT---
[ DEBUG     ] ---GRADE: DOCUMENT RELEVANT---
[ DEBUG     ] ---GRADE: DOCUMENT RELEVANT---
[ DEBUG     ] ---GRADE: DOCUMENT RELEVANT---
[ DEBUG     ] ---GRADE: DOCUMENT NOT RELEVANT---
[ DEBUG     ] ---GRADE: DOCUMENT RELEVANT---
[ DEBUG     ] ---GRADE: DOCUMENT RELEVANT---
[ DEBUG     ] ---DOCUMENT RELEVANCE: 9 of 10---
[ DEBUG     ] ---ASSESS GRADED DOCUMENTS---
[ DEBUG     ] ---DECISION: GENERATE---
[ DEBUG     ] Finished running: grade_documents
[ DEBUG     ] ---GENERATE---
[ DEBUG     ] ---CHECK HALLUCINATIONS---
[ DEBUG     ] ---DECISION: GENERATION IS GROUNDED IN DOCUMENTS---
[ DEBUG     ] ---GRADE GENERATION vs QUESTION---
[ DEBUG     ] ---DECISION: GENERATION ADDRESSES QUESTION---
[ DEBUG     ] Finished running: generate
[ DEBUG     ]
[ DEBUG     ] answer:
[ DEBUG     ] Pipelined table functions are a type of table function that allows
```

What is up?



Building a custom ChatGPT-like chatbot

Oracle Database 23ai provides mostly storage for vectors and distance queries. Depending on your workload, embedding is better executed outside (you maybe need your database resources for queries and not generating vectors).

To generate a text with a LLM you can't do it fully internally in the database for now.

OCI GenAI is very easy to use and doesn't require any commitment (reserving hardware or resources).

But keep in mind you don't own the models, and Oracle do release new ones from time to time, and then remove the old ones.

- You should always perform embeddings with the same model (except if the model explicitly says it is 1:1 compatible with another one, generally not the case)

Basic RAG doesn't bring much to the table, building more advances agents powered by LLM and retrievals has lot of potential.

LLMs starts being able to use tools and perform a step by step logical thinking.