

[Open in app](#)

Search



Write



# AI Sommelier Built with PaLM API and LangChain

Project walkthrough to build an LLM-powered application using PaLM API, Pinecone as the vector database, and Streamlit as the interface.

Gabriel Cassimiro · [Follow](#)

5 min read · Sep 26, 2023

170

1



...



Image by [Hermes Rivera](#) at [Unsplash](#)

In this article, I aim to demonstrate how to use the PaLM API to build solutions using LangChain. I will build a solution using LangChain chains, Google Embeddings, Pinecone as the VectorDB, and Streamlit as the interface to interact with users.

The problem we want to tackle is personalized wine recommendations. I want to build an AI capable of recommending wines based on user preference and also using a database containing 130K wines.

The full code can be found [here](#).

## The data

The data used in this project is the dataset [wine reviews, license](#). This dataset contains 130K reviews of wines. It also contains the country, region, variety, and winery of each wine. This is a sample of the data:

country	description	points	price	province	region_1	title	variety	winery
Italy	Aromas include tropical fruit, broom, brimston...	87	NaN	Sicily & Sardinia	Etna	Nicosia 2013 Vulkà Bianco (Etna)	White Blend	Nicosia
Portugal	This is ripe and fruity, a wine that is smooth...	87	15.0	Douro	NaN	Quinta dos Avidagos 2011 Avidagos Red (Douro)	Portuguese Red	Quinta dos Avidagos
US	Tart and snappy, the flavors of lime flesh and...	87	14.0	Oregon	Willamette Valley	Rainstorm 2013 Pinot Gris (Willamette Valley)	Pinot Gris	Rainsto...
US	Pineapple rind, lemon pith and orange blossom ...	87	13.0	Michigan	Lake Michigan Shore	St. Julian 2013 Reserve Late Harvest Riesling ...	Riesling	St. Julian
US	Much like the regular bottling from 2012, this...	87	65.0	Oregon	Willamette Valley	Sweet Cheeks 2012 Vintner's Reserve Wild Child...	Pinot Noir	Sweet Cheeks

Image by Author

## The Problem

So, to tackle our problem we want to first understand the taste of the user, then search the database for the wines that may be a good recommendation, and decide which one is the final recommendation.

To get the user input we will use a simple questionnaire inside of Streamlit with the following questions:

1. Preferred Taste Profile;
2. Level of experience;
3. Red or White wine preference;
4. Favorite Flavours;
5. Pairing intent;
6. Open field to add any information about your taste.

All of these (except the last) have pre-selected categories. However, since we are using LLMs this isn't strictly necessary. We could leave all questions with free text input because the LLM is able to use it without the categories. However, to guide the user I chose to use categories.

## The Architecture

The Architecture of the solution is shown below:

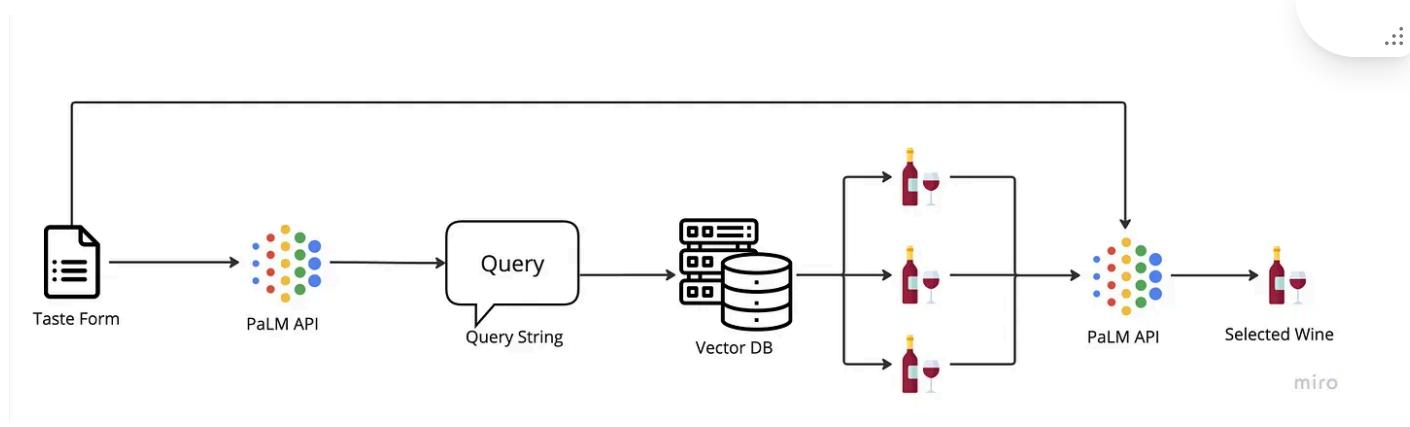


Image by Author

For the solution, we will have two calls to the LLM API. In the first one, we will pass the taste form information and ask for the LLM to generate a string query that summarizes the taste and will be used for similarity search on the vector DB. After that, we can find the most similar descriptions of the wines to the person's taste.

However, this is not enough for a final decision, because some characteristics like red vs. white are not taken into account by this simple similarity search. That is why we call the LLM again passing the original taste form and the top 3 most similar wines to the LLM and ask it to select the best one for it and explain the reasoning.

## The Code

Let's start by looking at the chains used for the calls to the API.

First, if you do not know already what is LangChain, this is a brief definition:

*LangChain is a framework for developing applications powered by language models.*

That is the official definition of LangChain. This framework was created recently and is already used as the industry standard for building tools powered by LLMs. ...

Chains are a core feature of LangChain and enable the integration of various components to form a unified application. They can format user input using a PromptTemplate and then forward it to an LLM. By linking multiple chains or combining them with other elements, we can create more complex chains.

This is the code for both of the chains:

```
1  from langchain.prompts import ChatPromptTemplate
2  from langchain.chains import LLMChain, SequentialChain
3  from langchain.output_parsers import ResponseSchema, StructuredOutputParser, RetryWithError
4
5
6  def build_query_chain(llm):
7
8      query_string = ResponseSchema(
9          name="query_string",
10         description="The string used to query the vector database for wine options.",
11     )
12
13     output_parser = StructuredOutputParser.from_response_schemas(
14         [query_string]
15     )
16     response_format = output_parser.get_format_instructions()
17
18     prompt = ChatPromptTemplate.from_template(
19         """You are an expert wine sommelier. Your goal is to select a wine from the data
20             Take a breath and understand the following user preferences:
21             taste: {taste}
22             experience level:{experience}
23             wine color: {wine_color}
24             flavor: {flavor}
25             pairing: {pairing}
26             complement: {complement}
27
28             Now create a string that will be used to do a similarity search on a vector data
29             To build better queries for similarity search, ensure they are specific, utilize
30             {response_format}
31             """
32     )
33
34     small_chain = LLMChain(llm=llm, prompt=prompt, output_key="query")
35
36     chain = SequentialChain(
37         chains=[small_chain],
38         input_variables=["taste", "experience", "wine_color", "flavor", "pairing", "complement"],
39         output_variables=["query"],
40         verbose=False,
41     )
42     return chain, response_format, output_parser
43
44 def build_recommendation_chain(llm):
45
```

```
46     recommendation = ResponseSchema(  
47         name="recommendation",  
48         description="The recommended wine Name.",  
49     )  
50  
51     explanation = ResponseSchema(  
52         name="explanation",  
53         description="The explanation of the selected recommendation.",  
54     )  
55  
56     output_parser = StructuredOutputParser.from_response_schemas(  
57         [explanation, recommendation]  
58     )  
59     response_format = output_parser.get_format_instructions()  
60  
61     ## Issue and Summary Chain  
62     prompt = ChatPromptTemplate.from_template(  
63         """You are an expert wine sommelier. Your goal is to select a wine from the opti  
64             Take a breath and understand the following user preferences:  
65                 taste: {taste}  
66                 experience level:{experience}  
67                 wine color: {wine_color}  
68                 flavor: {flavor}  
69                 pairing: {pairing}  
70                 complement: {complement}  
71  
72                 Now take a breath and understand the wine options:  
73                     Option 1:  
74                         {wine_1}  
75                     Option 2:  
76                         {wine_2}  
77                     Option 3:  
78                         {wine_3}  
79  
80                     Now select the best wine to recommend to this user.  
81  
82                     {response_format}  
83                     """  
84     )  
85  
86     small_chain = LLMChain(llm=llm, prompt=prompt, output_key="recommendation")  
87  
88     chain = SequentialChain(  
89         chains=[small_chain],
```

```
90     input_variables = [
91         "taste", "experience", "wine_color", "flavor", "pairing", "complement",
92         "response_format", "wine_1", "wine_2", "wine_3"],
93         output_variables=["recommendation"],
94         verbose=False,
95     )
96     return chain, response_format, output_parser
```

build chains my hosted with ❤️ by GitHub

[view raw](#)

Here we have a couple of elements I would like to point out. First, we have the **Response Schema** and **Output Parser**. These are used to create a prompt in which the LLM will always return the output with the same structure and the parser will process that output to be able to work back in the code.

Then we have the prompt built using the **Chat Prompt Template**. This allows us to pass variables into the prompt like the answers of the taste form. Lastly, we put everything back together with the **Sequential Chain**, defining the inputs and outputs.

Now we have to initialize the DB.

## The Database

For the Vector Database, we are using Pinecone however you can easily change the one you want to use because we are using LangChain to interact with the DB.

To add the wines to the Database we need to first create an Embedding of each one and then upload them. We can easily perform this with the code below:



We create an instance of Document from LangChain containing the text and the Metadata. Then we use LangChain and Google PaLM Embeddings to upload to Pinecone. The embeddings have 768 dimensions.

For this implementation I am using the free tier from Pinecone, which allows for 1 free index of the standard resource.

The screenshot shows the Pinecone interface for the 'Starter' project. On the left sidebar, there are links for 'PROJECT', 'Indexes' (selected), 'API Keys', and 'Members'. Below these are 'Docs', 'Support Center', 'Cassimiro Tech Serv...', 'Settings', and an 'Upgrade' button. The main content area is titled 'Starter Indexes' and shows the 'wine-db' index. It provides details like METRIC (cosine), DIMENSIONS (768), POD TYPE (starter), and HOST (https://wine-db-6c24cfa.svc.gcp-starter.pinecone.io). It also shows PROVIDER (Free Tier) and MONTHLY COST (\$0). Below this, there are tabs for 'BROWSER' and 'METRICS'. A search bar at the top right contains 'Query by Vector' with the value 'vector: 0.99,0.59,0.13,0.53,0.9,0.08,0.99,0.86,0.23,0.28,0.05,0.34,0.3,0.33,0.38,0.33,0.61,0.02,0;'. To its right are 'Top K\*' set to '50' and a 'Query' button. Below the search bar is a 'Metadata Filter' button. The search results show 'Matches: 1-10 of 50' with columns for 'ID' and 'VALUES'.

Image by Author

Now we just need to put everything together passing the input from the application to the chains and connect to the vector database.

## The App

The App will be the main control of the flow. It is also where we will create all the resources we need: the form, the PaLM API Connection, the Pinecone connection and the final visualization.



This code can be divided in three main parts:

1. The form to get the user input;
2. The interaction of the chains and the vector database;
3. The display of the final recommendation.

These parts are executed inside the main script but for a larger application should be modularized.

This is an example of using the tool:

1. Preferred Taste Profile

Dry

2. Wine Experience

Casual drinker (have tried a few)

3. Red vs. White

Prefer red

4. Favorite Flavors

Oaky/Woody × Spicy ×

5. Pairing Intent

Red meat dish

6. Any complement of the answers above with other foods or tastes you like the most.

Complement the answers above

I really like whiskey and bitter flavours

**Wine Sommelier**

Fill in the form in the sidebar to get a wine recommendation

Generate recommendation

Recommended Wine

Cedarville Vineyard 2014 Naylor Vineyard Petite Sirah

Explanation: The user is looking for a dry red wine with oaky/woody and spicy flavors. The Cedarville Vineyard 2014 Naylor Vineyard Petite Sirah is a good choice because it is a full-bodied red wine with notes of oak, black pepper, and tobacco. It is also a good value for the price.

Made with Streamlit

Image by Author

So this is the final product. All we have to do now is deploy this application. With this in mind, the first thing we need is to allow the user to choose which LLM to use and have an input for the API key.

The final version can be seen [here](#).

## Conclusion

Using LangChain allows us to build quick and in a modular manner. For the problem of recommending wines based on a database we were able to use Chains with two prompts that helped to understand the users taste, find some options to suggest and make the final decision with an explanation.

An easy way to interact with the user is through an interface such as Streamlit, making it fast to develop and deploy.



Thanks for reading.

If you like the content and want to support me, you can buy me a coffee:

**Gabriel Cassimiro is a Data Scientist sharing free content to the community**

I love supporting creators!

[www.buymeacoffee.com](http://www.buymeacoffee.com)

Here are a few other articles you might be interested in:

### **Async calls for Chains with Langchain**

How to make Langchain chains work with Async calls to LLMs, speeding up the time it takes to run a sequential long...

[towardsdatascience.com](https://towardsdatascience.com/)

### **Solving Unity Environment with Deep Reinforcement Learning**

End to End Project with code of a PyTorch implementation of Deep Reinforcement Learning Agent.

towardsdatascience.com

Llm

Langchain

Pinecone

Palm

Palm Api



## More from the list: "Langchain."

Curated by Shubham Rathod



Matheus Ferreira

**AI Food Advisor with LangChain and Pinecone**

5 min read · Oct 13, 2023



Jamesev

**Playlist Recommender based on Spotify,...**

3 min read · Jun 25, 2023



Emil Azadia

**Using Agents with LangChain** >

8 min read · Mar 20, 2023

[View list](#)



Experts  
Machine Learning

**Written by Gabriel Cassimiro**

[Follow](#)



208 Followers

Solving and creating problems with AI. Google Developer Expert in ML

**More from Gabriel Cassimiro**

Gabriel Cassimiro in Towards Data Science

**LLM Output Parsing Function Calling vs. LangChain**

How to consistently parse outputs from LLMs using Function Calling of Open AI API and...

11 min read · Sep 21, 2023



271



Gabriel Cassimiro in Towards Data Science

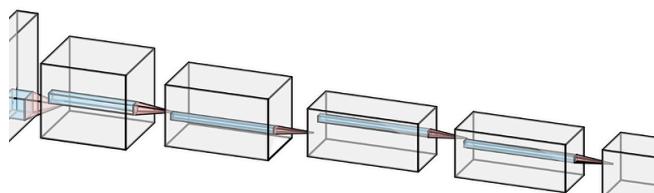
**Async calls for Chains with Langchain**

How to make Langchain chains work with Async calls to LLMs, speeding up the time it...

6 min read · Jul 10, 2023



380



Feature Learning



Gabriel Cassimiro in Towards Data Science

**Transfer Learning with VGG16 and Keras**

How to use a state-of-the-art trained NN to solve your image classification problem

<https://medium.com/@gabrielcassimiro17/ai-sommelier-built-with-palm-api-and-langchain-1b270c41a16a>


Gabriel Cassimiro in Towards Data Science

**How to prepare for the GCP Professional Machine Learning...**

Courses review, study tips, and how I did it

4 min read · Jun 16, 2021

8 min read · Jan 10, 2022



253



...



300



...

See all from Gabriel Cassimiro

## Recommended from Medium



Chandan Durgia

### Evaluation of RAG (Retrieval-Augmented Generation)...

Quantifying the accuracy and relevance of the RAG output

◆ · 9 min read · Feb 21, 2024



91



...



...



Supriyo Roy Banerjee

### Text Summarization with Generative Models on Vertex AI

## Overview

8 min read · Oct 11, 2023

## Lists



### Natural Language Processing

1250 stories · 731 saves



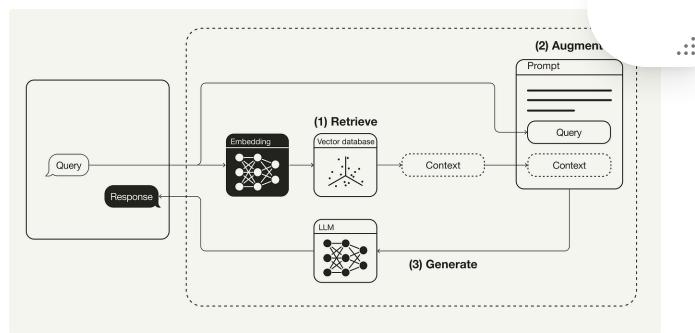
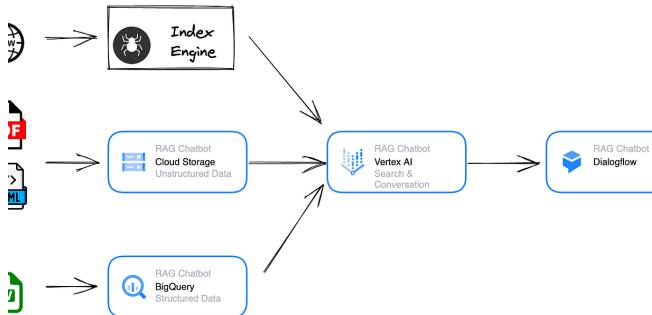
### ChatGPT prompts

44 stories · 1198 saves



### Staff Picks

594 stories · 790 saves



Euclidean AI in Towards AI

## Use Google Vertex AI Search & Conversation to Build RAG Chatbot

Google has recently released their managed RAG (Retrieval Augmented Generator)...

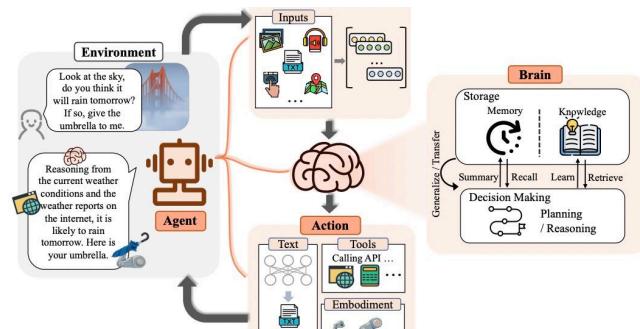
⭐ · 6 min read · Oct 25, 2023

37

5

+

...



Henry Heng LUO



Leonie Monigatti in Towards Data Science

## Retrieval-Augmented Generation (RAG): From Theory to LangChain...

From the theory of the original academic paper to its Python implementation with...

7 min read · Nov 14, 2023

1.2K

9

+

...



Saumya Mittal

## Intro of AI agent, & AI agent projects summary

Basic introduction of AI agent, and 13 interesting AI agent frameworks and AI agen...

16 min read · Jan 23, 2024

👏 183



...



3



...

See more recommendations

