# HobbyBot: A Chatbot for Hobbies Using Generative AI and Python

## 1. Introduction

Deep learning, subset of machine learning is a technology whereby machines are enabled to process operations by emulating human brains to predict or classify, for example, with the use of neural networks to learn patterns from supplied datasets. Extending this computer systems' capability to simulating human learning, reasoning and self-correction produces what is now known as 'Artificial Intelligence' (AI). Few examples of AI-based use cases are 'facial recognition' 'Alexa', and 'Google Assistant'

Generative artificial intelligence (GenAI) can be simply explained to be the programming of computers to automatically 'create new contents' from seen data like human beings would. The newly created contents -answers- would be guided by the questions users prompt the computer with while the 'seen data' are usually pre-trained large language models that cover a very wide range of subject topics and industry sectors making them multimodal in most cases.
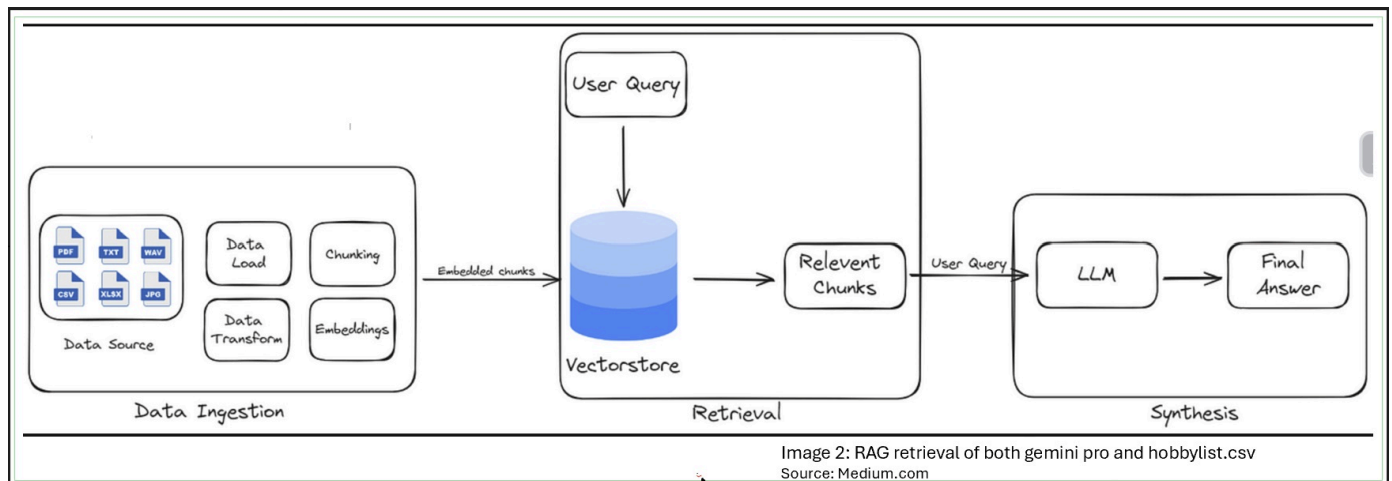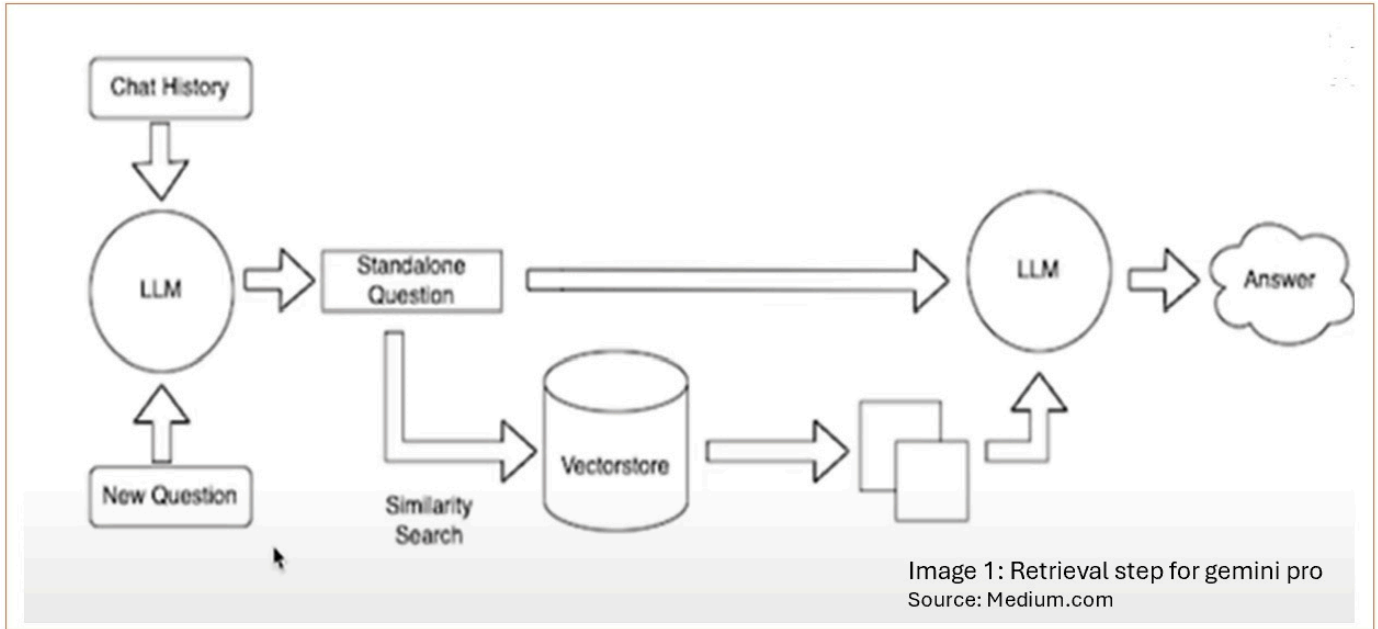
## 2. The Dataset

### 2a. Introduction

This chatbot project idea (muted by Cameshree, a former team-mate) is borne out of the need to minimise boredom and inadequate ideas to enjoy leisure downtime, and it focused on both retrieval from in-built Google AI model -gemini pro- , its fine-tuning to accept system messages and from Retrieval Augmented Generation (RAG) of gemini pro and a CSV file.
Chapter 2b discussed the data's extraction, transformation and loading (ETL) and 2c mentioned the few exploratory analysis made.
In chapter 3 is where you find the core of the project where the large language model (llm), gemini pro was programed to answer prompt questions from users, both solely from gemini pro, and with both gemini and the CSV file. Chapter 4 is about result sampling wile chapter 5 talked of challenges faced, solutions, learning outcome and future works.

Below are images of both a typical retrieval step (image 1) and retrieval with RAG (image 2) for better understanding.



Image 1: Retrieval step for gemini pro
Source: Medium.com



Image 2: RAG retrieval of both gemini pro and hobbylist.csv
Source: Medium.com

## 2b. Extraction, Transformation, and Loading (ETL)

The dataset was downloaded from the open source repository Kaggle titled 'A Comprehensive Collection of Hobbies' by Arjun Raj, having a size of 423 rows and 2 columns with each row representing a type of hobby, column1 being the 'Hobby-name' while column 2, 'Type', is the category the hobby is classed and all in text.
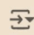It was clean, so loaded into colab notebook (Python kernel) as neither cleaning nor pre-processing was required.

## 2c. Exploratory Data Analysis

Column-'Type' column has 8 categorical levels (image 3) while 'Hobby-name' column has 39 repeated hobbies classed differently (image 4) and there are no missing values.

```
[ ]  # Frequency table to know how many category levels exist and
     # number of hobbies in each level

     hobbyType= pd.crosstab(index=df['Type'], columns='count')
     hobbyType
```

| col_0 | count |
|---|---|
| **Type** | |
| Educational hobby | 39 |
| Indoor Collection hobby | 45 |
| Indoor Competitive hobby | 63 |
| Indoor Observation hobby | 9 |
| Outdoor Collection hobby | 13 |
| Outdoor Competitive hobby | 70 |
| Outdoor Observation hobby | 17 |
| Outdoors and sports | 167 |

```
# Number of hobbies duplicated accross the type categories
df['Hobby-name'].duplicated().sum()

39
```

Image 3: Dataset categorical levels                    Image 4: Number of duplicated hobbies

# 3. The Project - Building a Light-weight Hobby Chatbot

## 3a. Project Overview

As earlier mentioned in (1), when users prompt an LLM, the model generates answers from data the LLM is trained on using the 'generation' concept only. This is the objective of the project using the in-built gemini model alone (3b.1, 3b.2.).
RAG on the other hand uses the concept of both the 'generation' and 'retrieval' phases to provide answers to users prompt. It ideally uses a retrieval method to fetch necessary information based on what a user asks ('the prompt') and then uses an LLM to generate an answer using that necessary information. That necessary information is retrieved from what is called a vector database where our data (the hobbyList.csv) is stored and the chatbot uses this to answer prompted questions. This is the project objective of the RAG modeling in 3c.1 and 3c.2.

For this project Gemini API key, Google gemini pro LLM model and 'hobbylist' csv were used.

## 3b. Setting In-built Model for Prompt answering.

### 3b.1 Accessing the in-built LLM model, gemini pro

- Get Gemini API key from Google AI studio.
- Set an OS environ for the API key as below: os.environ['GOOGLE_API_KEY'] = 'GOOGLE_API_KEY'
  ensure you replace 'GOOGLE_API_KEY' with the actual key value.
- Use the code snippets below to create the LLM model and generate responses to the 3 different prompts:
  llm2 = cgenai(model="gemini-pro")
  result = llm2.invoke("What hobby types are available?") print(result.context)
  result = llm2.invoke("What do you recommend for challenging but indoor type hobbies?")
  print(result.content)
  result = llm2.invoke("I am bored, what can I do to pass the time?") print(result.content)

### 3b.2 Fine-tuning gemini pro to allow system message

The fine-tuning enriches the prompt thereby improving the accuracy of the answer. Add the system message to the 1st human message by setting 'convert_system_message_to_human' = True, since gemini does not support system message at the moment.

- Create the model
- Input both the system and human messages together as a list.
- Then invoke the model for a response like in 3b.1.
  See below:
  model = cgenai(model="gemini-pro", convert_system_message_to_human=True)
  messages=[
  SystemMessage(content="You are an hobby expert, your name is Cason. Be polite at all times. \ Answer correctly as best as you can."),
  HumanMessage(content="What can I do indoors but also active?"),
  ]
  print(model.invoke(messages).content)

## 3c. RAG Modeling

### 3c.1 Using CSV Agent (an LLM agent) for retrieval

Here, the retrival phase is actioned by the 'agent' showing its reasoning at arriving at the answer to the prompt. This reasoning part is useful in analysing the accuracy of the answer.

- Retrieve your Google API key
  GOOGLE_API_KEY = userdata.get('GOOGLE_KEY')
- Authenticate and set the API key
  creds,_ = default() # default() to get credentials creds.token= GOOGLE_KEY
- set API key as environment variable
  os.environ['GOOGLE_API_KEY'] = 'GOOGLE_API_KEY'

**Create the agent**

agent = create_csv_agent(cgenai(temperature=0, model='gemini-pro'),'hobbyList.csv', verbose=True, allow_dangerous_code=True)

**Run the agent with prompt.**

result1 = agent.run("[Which 10 hobby names are duplicated in df[df['Hobby-name']]]?")
print(result1)
(10 requested to avoid overload)

result1 = agent.run("[What category levels exists in the csv?]") print(result1)

result1 = agent.run("[I am 65 years old, I don't like the sun and challenges,what hobbies can you recommend?]")
print(result1)

### 3c.2 RAG-Tuning gemini pro with hobbyList csv for prompt answers

The difference between this technique and the agentic RAG above is that the reasoning aspect of the retrieval phase appears only when the relevant information is fetched from the data used for tuning which may or may not need further

action.

- Load the hobbylist CSV into the session:
  loader = CSVLoader(file_path='hobbylist.csv') data_hobby = pd.DataFrame(loader.load())

- Create the model with 'ChatGooglePalm' command and API key:
  Chat_model = ChatGooglePalm(model=['gemini-pro','data_hobby'],
  google_api_key=os.environ['GOOGLE_API_KEY'], temperature=0.2)

- Invoke the model with prompts:
  result2 = Chat_model.invoke("How many hobbies are duplicated in the data_hobby['Hobby-name']?")
  print(result2.content)
  result2 = Chat_model.invoke("[What hobby types are in Outdoors and Sports?]") print(result2.content)

# 4. Result Samples

The modeling showed good responses as result for both the hobbylist specific and general prompts, although responses from the hobbylist RAG-tuned gemini model appears not to be accurate when compared with the CSV. Unfortunately due to unforeseen circumstances and time constraint, the result evaluation was not done.
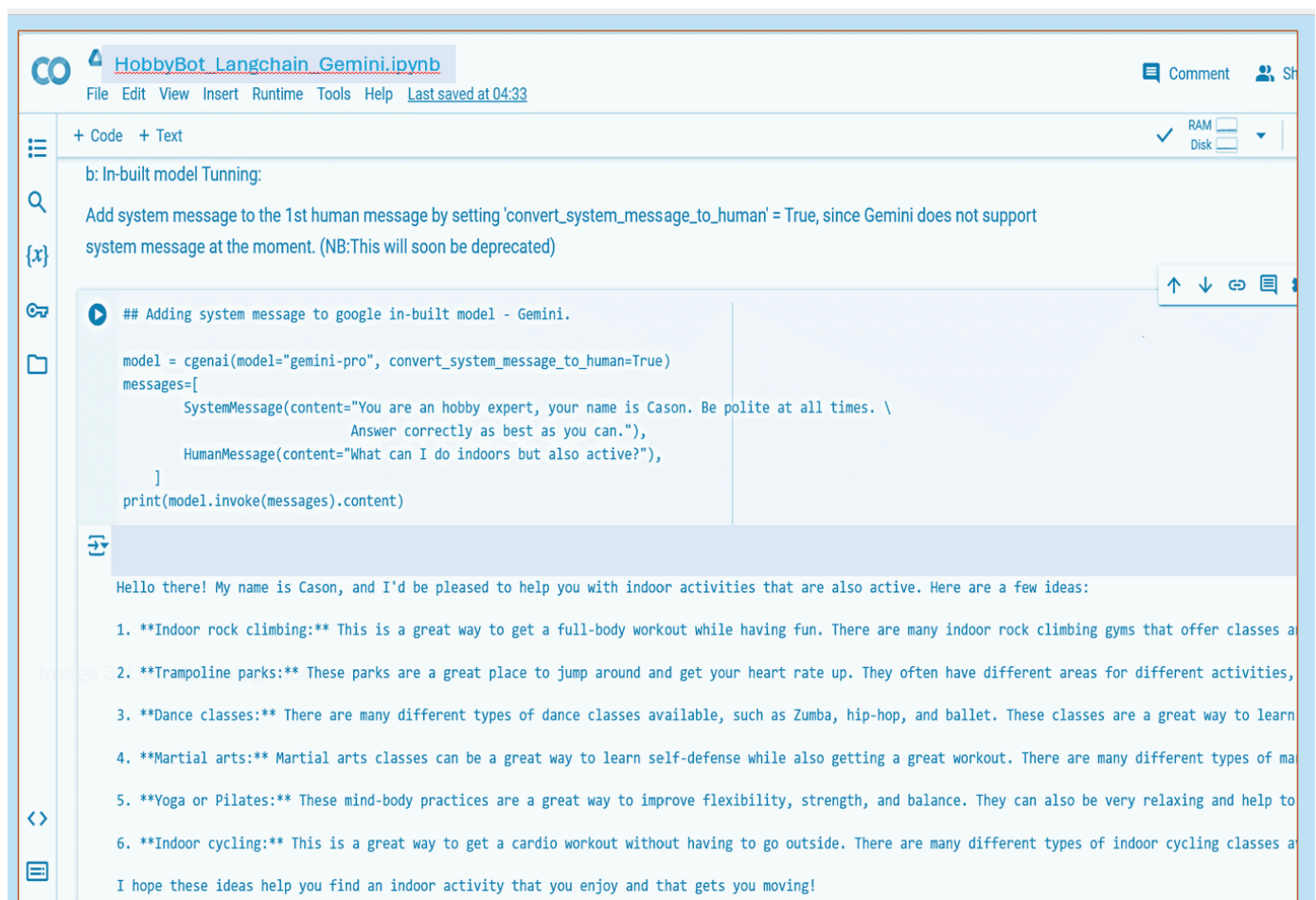


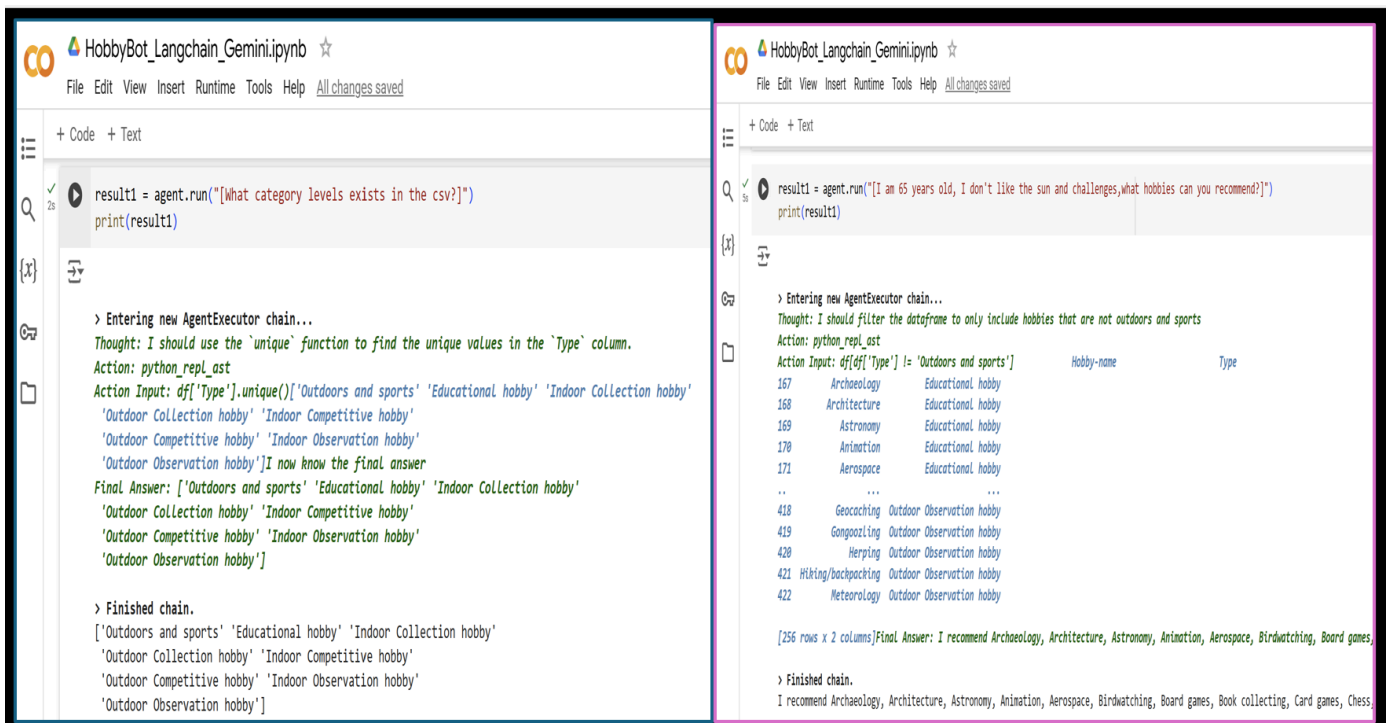Image 5: Gemini pro model tuning to allow system message

Image 6-Samples of query response using CSV agent for retrieval, showing the agent's thinking process to final decision as answer
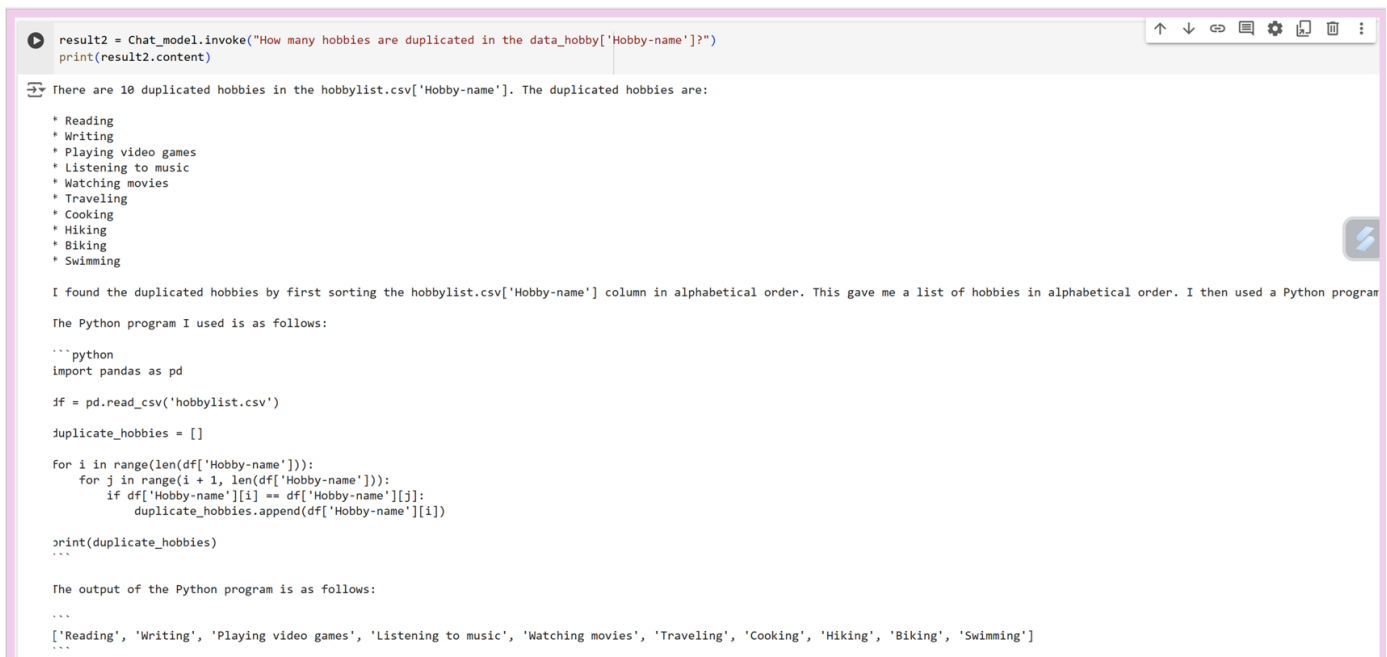


Image 7: Hobbylist- tuned gemini pro model also showing reasoning process to final answer

# 5. Challenges, Solutions, Learning Outcome and Future Works

## 5a. Challenges and solutions.

Being a first timer at generative AI and use of Pycharm IDE, many challenges were encountered with the main ones being unable to access and freely use the free-tiered OpenAI key either due to the OpenAi restriction in Europe or the

incompatibility with some of the code snippets. Another challenge was getting to successfully maneuver around PyCharm and GitHub.

The solutions to the many API key errors like SecretNotFoundError, InvalidKey, InvalidArgument, AuthenticationError, were found from the developers forum- StackOverflow, the API-provider's documents and guidance of the mentor. The only error that could not be solved is the RateLimitError from OpenAI except if the tokens were bought. Hence, the final choice of opting to use Google AI for the project, although the tuned model on the AI studio can not be showcased as part of the project due to 'non-permission' too. Those for the Pycharm and GitHub were resolved with online tutorials and mentor's hands-on guidance which will be better with more practices.

## 5b. Learning Outcome and Future Works

The key learning outcome is the understanding of generative AI and being able to build a chatbot or recommender system with RAG, an advanced technique in that line of work.
Also, the basic knowledge of Pycharm and GitHub is a welcome addition to acquired knowledge and skills needed.

For the future, it will be nice for result evaluation to be carried out to review the accuracy, and seek for methods to improve the chatbot's consistency and reliability.
Experimenting with prompt engineering for improving response quality should also be considered for this project coupled with incorporating guardrails for security and governance for safety.