

CN7050- Intelligent Systems Module

Student ID: 3059676

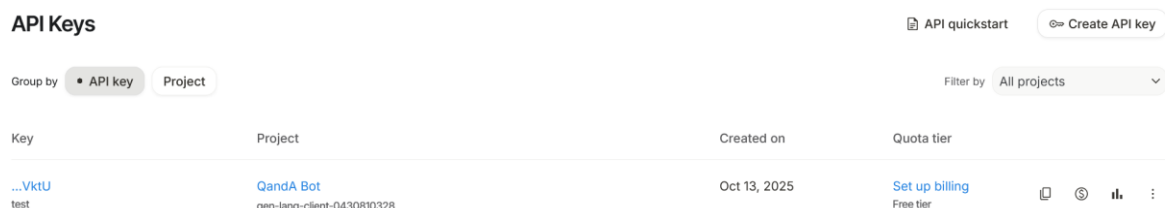
Name: Mann Shrestha

Part I: A Multi Model Interactive Agent with Text and Image Generation capabilities.

Task 1: Account & API key

1 and 2. Sign in to Google AI Studio (Gemini developer console) and Create an API key (Gemini Developer API)

- I have created a Gemini account and API key from the given instruction in the tutorial 5.



Key	Project	Created on	Quota tier
...VktU test	QandA Bot gen-lang-client-0430810328	Oct 13, 2025	Set up billing Free tier

3. Hugging Face API:

- Similarly, I have created a hugging face account to access an API following the instruction in the tutorial 5.

Access Tokens

User Access Tokens

+ Create new token

Access tokens authenticate your identity to the Hugging Face Hub and allow applications to perform actions based on token permissions.

Do not share your Access Tokens with anyone; we regularly check for leaked Access Tokens and remove them immediately.

Name	Value	Last Refreshed Date	Last Used Date	Permissions
base	hf_...RgPX	2 days ago	2 days ago	READ

4. Open a new Google Colab notebook

5. Install the Google Gen AI SDK + helper libraries.

```
1 # installing libraries
2 !pip install --upgrade google-genai gradio requests pillow
```

- I am using upgradeable version of **google-genai** – provides an interface to integrate Google's generative models into python application.
- **Request:** HTTP library to send HTTP1.1 requests extremely easily. Such as PUT and POST data.

- **Pillow:** The Python Imaging Library adds image processing capabilities to our Python interpreter.

6. Securely provide your API keys in Colab

```

1 from getpass import getpass
2 import os
3
4 print("Paste keys when prompted:")
5
6 GEMINI_KEY = getpass("Gemini API Key: AIzaSyBiYR0ThivMnfBob0HXXbwck050hebVktU")
7 HF_KEY = getpass("Hugging Face API Key: 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")
8
9 os.environ["GEMINI_API_KEY"] = GEMINI_KEY
10 os.environ["HF_KEY"] = HF_KEY
11

```

Paste keys when prompted:

Gemini API Key: AIzaSyBiYR0ThivMnfBob0HXXbwck050hebVktU.....

Hugging Face API Key: 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.....

- We use `getpass()` module to securely handle the password prompt where programs interact with the users via the terminal.

7. Hugging Face Image Generation function.

```

1 import requests
2
3 def call_hf_image(prompt):
4     """
5     Hugging Face Stable Diffusion XL example.
6     """
7     endpoint = "https://api-inference.huggingface.co/models/stabilityai/stable-diffusion-xl-base-1.0"
8     headers = {"Authorization": f"Bearer {os.environ['HF_KEY']}" }
9     payload = {"inputs": prompt}
10
11     # Make a POST request to a web page, and return the response text:
12     response = requests.post(endpoint, headers=headers, json=payload)
13     if response.status_code != 200:
14         return {"error": response.text}
15
16     image_bytes = response.content
17     from PIL import Image
18     from io import BytesIO
19     img = Image.open(BytesIO(image_bytes))
20     return {"pil_image": img}

```

- Create a function to call hugging face API to generate image through prompt.
- I am using hugging face **stable-diffusion-xl-base-1.0 model**. It is developed by Stability AI. It is a diffusion model based on text-to-image generative model that uses two fixed, pretrained text encoders.

Source: <https://huggingface.co/stabilityai/stable-diffusion-xl-base-1.0>

- As for the **BytesIO** also called buffered I/O, expects bytes-like objects and produces bytes objects. No encoding, decoding, or newline translation is performed. This category of streams can be used for all kinds of non-text data, and also when manual control over the handling of text data is desired.

Ref. <https://docs.python.org/3/library/io.html>

8. Imports + Gemini client initialization.

```
1 # imports and Gemini client
2 import google.generativeai as genai
3 import requests, json, time, base64, os
4 from PIL import Image
5 from io import BytesIO
6 from IPython.display import HTML, display
7
8 import sys
9 import matplotlib.pyplot as plt
10
11 # Initialize Gemini client only if key provided
12 if os.environ.get('GEMINI_API_KEY'):
13     genai.configure(api_key=os.environ['GEMINI_API_KEY'])
14     GEMINI_AVAILABLE = True
15 else:
16     GEMINI_AVAILABLE = False
17     print("Gemini key not found - text generation disabled until you set GEMINI_KEY.")
```

- Importing necessary libraries and initialize the Gemini client.
- Verify If Gemini key is not found, we will see the error message shown in the last code line.

9. Gemini text generation function.

```
1 # Cell 5 - function: generate text with Gemini
2 def call_gemini(prompt):
3     try:
4         model = genai.GenerativeModel("gemini-2.5-flash")
5         response = model.generate_content(prompt)
6         return {"type": "text", "content": response.text}
7     except Exception as e:
8         return {"type": "text", "content": f"Gemini error: {e}"}
```

- This function to use and generate the text using Gemini LLM model. I had to use Google Gemini Generative model - “gemini-2.5-flash” instead of “gemini-1.5-flash” – this 1.5 version gives me an error i.e. depreciated model. Moreover, 2.5 – flash version model gives me a response in text form, otherwise it will throw an exception.

10. Intent classifier (simple keyword routing)

```
1 def classify_intent(user_text):
2     t = user_text.lower() # convert in to lower case
3
4     # image triggers
5     if "image" in t or "picture" in t or "photo" in t or t.startswith("create me an image"):
6         return "image"
7
8     # text/generation/transform triggers
9     if any(k in t for k in ["write", "summarize", "explain", "draft", "describe", "generate text", "summarise"]):
10        return "text"
11
12    # fallback: if contains short multimedia instructions mention both -> prefer image
13    if "create me" in t and "and" in t and ("video" in t or "image" in t):
14        if "image" in t:
15            return "image"
16
17    # Default -> text
18    return "text"
```

- We created a function **classify_intent()** which takes user text, then convert the text into lowercase. Similarly, this function used to classify image or text or if it contains short multimedia content then return image. Otherwise by default returns text.

11. Orchestrator: route user requests and call right API

```
1 def handle_user_request(user_text):
2     intent = classify_intent(user_text)
3     print(f"[Agent] detected intent: {intent}")
4
5     if intent == "text":
6         return {"type": "text", "content": call_gemini(user_text)['content']}
7     elif intent == "image":
8         return {"type": "image", "content": call_hf_image(user_text)}
9     else:
10        return {"type": "text", "content": "Sorry, I don't understand."}
```

- `handle_user_request()` used to handle user text and check for classify intent of the given user text if its intention were text it calls the function `call_gemini()` with this function it uses gemini-2.5-flash model and check for the content either it is for text or image if either not return "Sorry, I don't understand" text.

12. Display image function

```
1 from PIL import Image
2 from IPython.display import display, HTML
3 import base64, io, requests, numpy as np
4
5 def display_image_response(image_input, filename="/content/generated_image.png"):
6     # Extract PIL Image
7     if isinstance(image_input, dict):
8         if 'pil_image' in image_input:
9             img = image_input['pil_image']
10        elif 'b64_json' in image_input:
11            image_bytes = base64.b64decode(image_input['b64_json'])
12            img = Image.open(io.BytesIO(image_bytes))
13        elif 'url' in image_input:
14            image_bytes = requests.get(image_input['url']).content
15            img = Image.open(io.BytesIO(image_bytes))
16        else:
17
18            raise TypeError(f"Unsupported dict keys: {image_input.keys()}")
19
20        elif isinstance(image_input, Image.Image):
21            img = image_input
22        elif isinstance(image_input, str):
23            img = Image.open(image_input)
24        elif isinstance(image_input, bytes):
25            img = Image.open(io.BytesIO(image_input))
26        elif isinstance(image_input, np.ndarray):
27            img = Image.fromarray(image_input)
28        else:
29            raise TypeError(f"Unsupported image type: {type(image_input)}")
30
31        # Save image
32        img.save(filename)
33
34        # Display thumbnail
35        display(img.resize((256, 256)))
36        time.sleep(2)
```

- Here, we created a function **display_image_response()** to display with a size of 255x255 and save the image locally.
- In this function, we use base64 - Base64 which is a binary-to-text encoding that uses 64 printable characters to represent each 6-bit segment of a sequence of byte values. As for all binary-to-text encodings, Base64-encoding enables transmitting binary data on a communication channel that only supports text.
- And with the help of `base64.b64decode()` method, we can decode the binary string into normal form.

13. Function to show the text generated by the Gemini model in the wrapped form.

```
1 import textwrap
2
3 def pretty_print(text, width=80):
4     print("\n".join(textwrap.wrap(text, width)))
5
```

- We use textwrap module for wrapping and formatting of plain text. This module provides formatting of text by adjusting the line breaks in the input paragraph.

14. Code for interactive agent.

```
1 print("Multimodal Agent (demo). Type 'exit' to quit.")
2 image_counter = 1 # optional: number images uniquely
3
4 while True:
5     q = input("You: ")
6     if q.strip().lower() in ("exit", "quit"):
7         break
8
9     out = handle_user_request(q) # Agent intent
10
11     if out['type'] == 'text': # for text to generate
12         print("\n", pretty_print(out['content']))
13     elif out['type'] == 'image': # for image to generate
14         print("\nAgent (image):")
15         # Save with unique filename each time
16         filename = f"generated_image_{image_counter}.png"
17         display_image_response(out['content'], filename)
18         image_counter += 1
19     else:
20         print("Unknown output type:", out.get('type'))
21
22     print("\n---\n")
```

- This code is used for the demonstration of multi-model Agent which is used for the generation of image or text given by the user input or prompt.

15. Output and result of the code.

- When we run the above code, we will witness a prompt to generate image or text output. If we want to exit the prompt, we can enter exit and press enter. To generate the text, include prompt with keywords like “Image”, “photo” etc. for image generation and “describe”, “summarize” etc., as per step 10.
- **Prompt1:**

Multimodal Agent (demo). Type 'exit' to quit.
You: image for a single pink rose.
[Agent] detected intent: image

Agent (image):



- Prompt 2

You: image for multiple vehicles in london
[Agent] detected intent: image

Agent (image):



- Prompt 3

You: write some description about agentic ai.

[Agent] detected intent: text

Agentic AI refers to a class of artificial intelligence systems that exhibit a high degree of autonomy and goal-oriented behavior. Unlike traditional AI that primarily reacts to specific prompts or commands, agentic AI proactively identifies, plans, and executes steps to achieve a predefined objective with minimal human intervention. Here are some descriptions highlighting its key characteristics:

1. **Goal-Oriented Autonomy:** Agentic AI is designed to pursue a high-level goal, breaking it down into smaller, manageable sub-tasks. It doesn't just respond; it *initiates action* to achieve its objective.
2. **Proactive Planning & Reasoning:** Given a goal, an agentic AI system can reason about the necessary steps, formulate a multi-step plan, and even generate alternative strategies if obstacles arise. It anticipates needs and potential issues.
3. **Action & Execution:** It's not just about thinking; it's about doing. Agentic AI can interact with the digital world (e.g., calling APIs, browsing the web, using software tools, writing code) and potentially the physical world (via robotics) to execute its plans.
4. **Self-Correction & Learning:** These systems often monitor their own progress, evaluate the outcomes of their actions, and adapt their strategies based on feedback. If a step fails, they can re-plan or attempt a different approach, demonstrating a form of iterative learning.
5. **Tool Use & Integration:** Agentic AI frequently leverages a wide array of external tools, services, and data sources. It acts as an orchestrator, intelligently choosing and integrating the right tools (like a web browser, a code interpreter, a database query, or a text editor) to accomplish its tasks.
6. **Persistence & Monitoring:** Once a goal is set, agentic AI will continue to work towards it, often over extended periods, monitoring its environment and adjusting its actions until the objective is met or it determines it cannot proceed.
7. **Moving Beyond Reactive AI:** It signifies a shift from AI that merely answers questions or performs single, isolated tasks to AI that can manage complex, multi-faceted projects and workflows independently.

In essence: Think of Agentic AI not just as a smart assistant that waits for instructions, but as a digital colleague that can be given an objective ("Research and summarize the latest trends in renewable energy, then draft an email to the team with key findings") and will autonomously plan, execute, monitor, and adapt its way to completing that complex task. It represents a significant step towards more independent and capable artificial intelligence.

Part II: Learning Objectives

Task 1: Account & API key

1. **Sign in to Google AI Studio (Gemini developer console).**
2. **Create an API key (Gemini Developer API):**
3. **Open a new Google Colab notebook.**
4. **Install the Google Gen AI SDK + helper libraries.**

```
1 !pip install --upgrade google-genai arxiv beautifulsoup4 requests readability-lxml
2
```

- **arXiv** is a free distribution service and an open-access archive for nearly 2.4 million scholarly articles in the fields of physics, mathematics, computer science, quantitative biology, quantitative finance, statistics, electrical engineering and systems science, and economics. Materials on this site are not peer-reviewed by arXiv.
<https://arxiv.org/>

- **Beautiful Soup** is a library that makes it easy to scrape information from web pages. It sits atop an HTML or XML parser, providing Pythonic idioms for iterating, searching, and modifying the parse tree.
<https://pypi.org/project/beautifulsoup4/>
- **readability-lxml**: Given an HTML document, extract and clean up the main body text and title.

5. Securely provide your API key in Colab

```
1 from getpass import getpass
2 import os
3
4 api_key = getpass("Paste your Google Gemini API key here: AIzaSyBiYR0ThivMNfBob0HXXbwck050hebVktU")
5 os.environ["GEMINI_API_KEY"] = api_key # optional - keep in memory only
6
```

Paste your Google Gemini API key here: AIzaSyBiYR0ThivMNfBob0HXXbwck050hebVktU.....

6. Initialize the Gen AI client in Colab.

- To initialize the Gen AI client, use the following code

```
1 from google import genai
2
3 # create client using the key we provided
4 client = genai.Client(api_key=os.environ["GEMINI_API_KEY"])
5
```

7. Try a small generate call to confirm connectivity.

```
1 resp = client.models.generate_content(
2     model="gemini-2.5-flash",
3     contents="Summarize in one sentence: Why is reproducibility important in research?"
4 )
5 print(resp.text)
```

- I will be using “gemini-2.5-flash” model.

8. Search arXiv for scholarly papers.

```
1 import arxiv
2
3 def search_arxiv(query, max_results=5):
4     search = arxiv.Search(
5         query=query,
6         max_results=max_results,
7         sort_by=arxiv.SortCriterion.Relevance
8     )
9
10    results = []
11    for result in search.results():
12        results.append({
13            "title": result.title,
14            "summary": result.summary,
15            "authors": [a.name for a in result.authors],
16            "pdf_url": result.pdf_url,
17            "id": result.get_short_id()
18        })
19    return results
20
21 # quick test
22 papers = search_arxiv("multimodal transformers", max_results=4)
23 for p in papers:
24     print(p['title'])
```

- **arXiv:** The goal of the API is to allow application developers access to all the arXiv data, search and linking facilities with an easy-to-use programmatic interface. This page provides links to developer documentation and gives instructions for how to join the mailing list and contact other developers and maintainers.

9. Grab web-page content for non-paper websites:

- Use requests + BeautifulSoup to extract main paragraphs (simple heuristic)

```

1 import requests
2 from bs4 import BeautifulSoup
3
4 def fetch_page_text(url, max_chars=4000):
5     try:
6         r = requests.get(url, timeout=10, headers={"User-Agent": "research-agent/1.0"})
7         r.raise_for_status()
8     except Exception as e:
9         return ""
10
11     soup = BeautifulSoup(r.text, "html.parser")
12     # simple extraction: join visible <p> text
13     paragraphs = [p.get_text(separator=" ", strip=True) for p in soup.find_all("p")]
14     content = "\n\n".join(paragraphs)
15     return content[:max_chars] # limit length for API

```

- **BeautifulSoup:** <https://beautiful-soup-4.readthedocs.io/en/latest/>
- **Prepared Requests:** Whenever we receive a Response object from an API call or a Session call, the request attribute is actually the PreparedRequest that was used. In some cases, we may wish to do some extra work to the body or headers (or anything else really) before sending a request.
- We created a fetch_page_text() function with a url and max_char = 4000 as a parameter.

10. Create a safe chunking helper (Gemini has big context but chunking helps reliability):

```

1 def chunk_text(text, max_chars=3000):
2     chunks = []
3     start = 0
4     while start < len(text):
5         end = min(len(text), start + max_chars)
6         chunks.append(text[start:end])
7         start = end
8     return chunks

```

- We created a function chunk_text () to divide a text with maximum characters of 3000 at a time.

11. Create a function that asks Gemini to summarize a text chunk.

```
1 SYSTEM_PROMPT = (  
2     "You are a concise research assistant. For each text provided, return: "  
3     "1) a one-paragraph summary (3-5 sentences), 2) three bullet key contributions/findings,"  
4     "and 3) a suggested short title. Be factual and include no hallucinated facts."  
5 )  
6  
7 def summarize_chunk(chunk_text):  
8     prompt = SYSTEM_PROMPT + "\n\nText to summarize:\n" + chunk_text  
9     resp = client.models.generate_content(  
10         model="gemini-2.5-flash", # adjust if unavailable  
11         contents=prompt  
12     )  
13     return resp.text
```

- We keep a short system prompt at the top to set the assistant behaviour.

12. Summarize a long document by summarizing chunks + combining.

```
1 def summarize_text_long(text):  
2     chunks = chunk_text(text, max_chars=3000)  
3     summaries = []  
4     for c in chunks:  
5         s = summarize_chunk(c)  
6         summaries.append(s)  
7  
8     # Optionally aggregate the chunk summaries into a single final summary:  
9     if len(summaries) == 1:  
10         return summaries[0]  
11     else:  
12         combined = "\n\n".join(summaries)  
13         final_prompt = SYSTEM_PROMPT + "\n\nCombine the following chunk summaries into a single concise summary:\n\n" + combined  
14         final_resp = client.models.generate_content(model="gemini-2.5-flash", contents=final_prompt)  
15         return final_resp.text  
16
```

- We create `summarize_text_long()` to summarize long text from the text and `chunk_text` with `max_character= 3000`.

13. Putting it together — search and summarize top N arXiv abstracts:

```
1 query = "multimodal transformers"  
2 # query = "vision transformer"  
3 papers = search_arxiv(query, max_results=5)  
4  
5 for p in papers:  
6     print("\n---")  
7     print("Title:", p['title'])  
8     text_to_summarize = p['summary'] # arXiv abstract (short)  
9     summary = summarize_text_long(text_to_summarize)  
10    print("Summary:\n", summary)  
11    print("PDF:", p['pdf_url'])  
12
```

- Write a manual query for any topic. In our case, “multimodal transformer” – this journal will be search through `search_arxiv()` with maximum result of 5.

Journal 1:

Title: Multimodal Learning with Transformers: A Survey
Summary:
1) This paper offers a comprehensive survey of Transformer techniques specifically tailored for multimodal data, a burgeoning area in AI research. It begins by establishing a background of multimodal learning and its significance.
2) Key Contributions:
* Provides a comprehensive survey of Transformer techniques applied to multimodal data.
* Offers a theoretical review of various Transformer architectures from a geometrically topological perspective.
* Discusses applications, common challenges, and future research directions for multimodal Transformers.
3) A Survey of Multimodal Transformer Techniques
PDF: <http://arxiv.org/pdf/2206.06488v2>

Journal 2:

Title: MANGO: Multimodal Attention-based Normalizing Flow Approach to Fusion Learning
Summary:
1) **Summary:**
Current multimodal fusion methods, relying on implicit Transformer attention, often fail to capture essential features and complex correlations within multimodal inputs. To address this limitation,
2) **Key Contributions/Findings:**
* Introduces MANGO (Multimodal Attention-based Normalizing Flow) for explicit, interpretable, and tractable multimodal fusion learning.
* Proposes an Invertible Cross-Attention (ICA) layer, incorporating three new mechanisms (MMCA, IMCA, LICA), to capture complex multimodal correlations.
* Achieves state-of-the-art performance on diverse multimodal tasks, including semantic segmentation, image-to-image translation, and movie genre classification.
3) **Suggested Short Title:**
MANGO: Explicit Multimodal Fusion via Invertible Cross-Attention
PDF: <http://arxiv.org/pdf/2508.10133v1>

Journal 3:

Title: Multimodal Transformer With a Low-Computational-Cost Guarantee
Summary:
1) Summary:
Transformer-based models have advanced multimodal understanding tasks like visual question answering, but their multi-head attention mechanism suffers from quadratic complexity, particularly with increasing sequence lengths.
2) Key Contributions/Findings:
* Addresses the quadratic computational complexity of multi-head attention in multimodal Transformers, which limits scalability.
* Introduces LoCoMT, a novel multimodal attention mechanism that reduces computational cost by assigning different attention patterns to individual heads.
* Achieves significant GFLOPs reduction and comparable or superior performance to established models on Audioset and MedVidCL datasets.
3) Suggested Short Title:
LoCoMT: Low-Cost Multimodal Transformers
PDF: <http://arxiv.org/pdf/2402.15096v1>

Journal 4

Title: Brain encoding models based on multimodal transformers can transfer across language and vision
Summary:
1) This study investigates how the human brain represents concepts across language and vision using encoding models trained with representations from multimodal transformers. While traditional encoding models are typically trained on unimodal data, this work explores the transferability of multimodal representations.
2) Key Contributions/Findings:
* Encoding models trained with multimodal transformer representations can predict fMRI responses across language and vision modalities (stories and movies).
* This cross-modal prediction is prominent in cortical regions that represent conceptual meaning and reveals shared semantic dimensions in language and vision.
* Multimodal transformers learn more aligned concept representations across language and vision than unimodal transformers, aiding understanding of brain's multimodal processing.
3) Suggested Short Title: Multimodal Transformers for Cross-Modal Brain Encoding
PDF: <http://arxiv.org/pdf/2305.12248v1>