



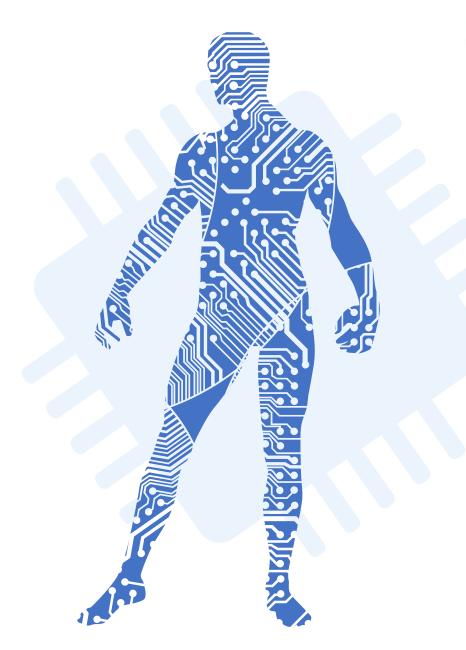
自然語言處理

第11章 檢索增強式文本生成 Retrieval Augmented Generation (RAG)

講師:紀俊男



- RAG 簡介
- RAG 實例講解
- 本章總結



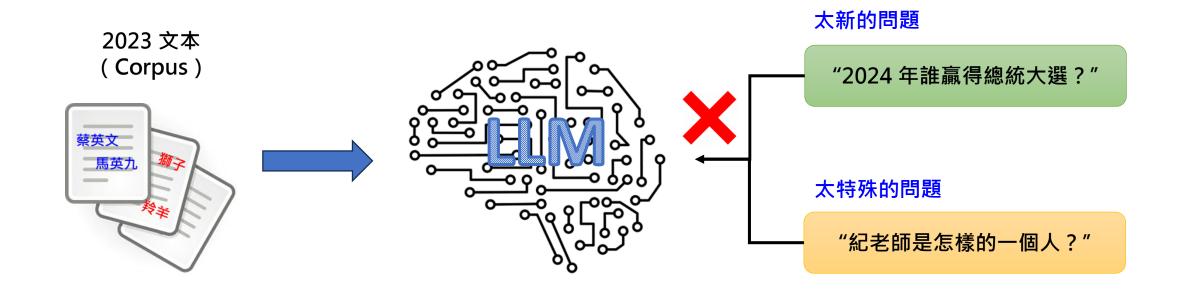






大語言模型的困境

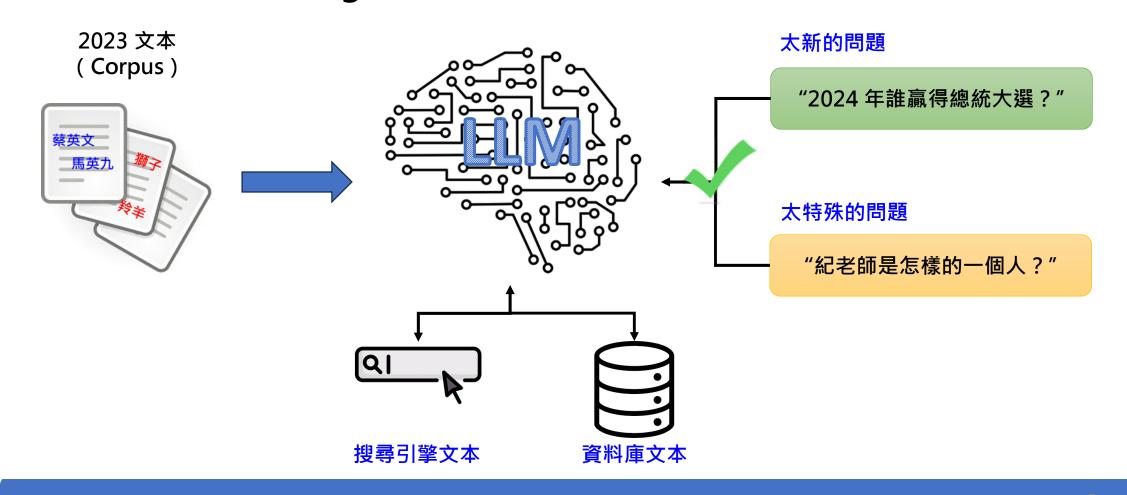




解決方法



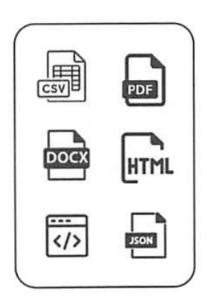
RAG = Retrieval Augmented Generation (檢索增強式文本生成)

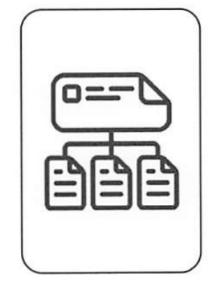


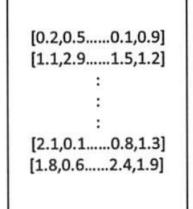
RAG 資料來源之生成步驟

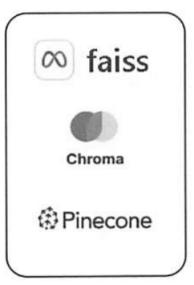


資料載入 Data Loading 資料分割 Data Chunking 資料嵌入化 Data Embedding 向量資料庫 Vector Databases



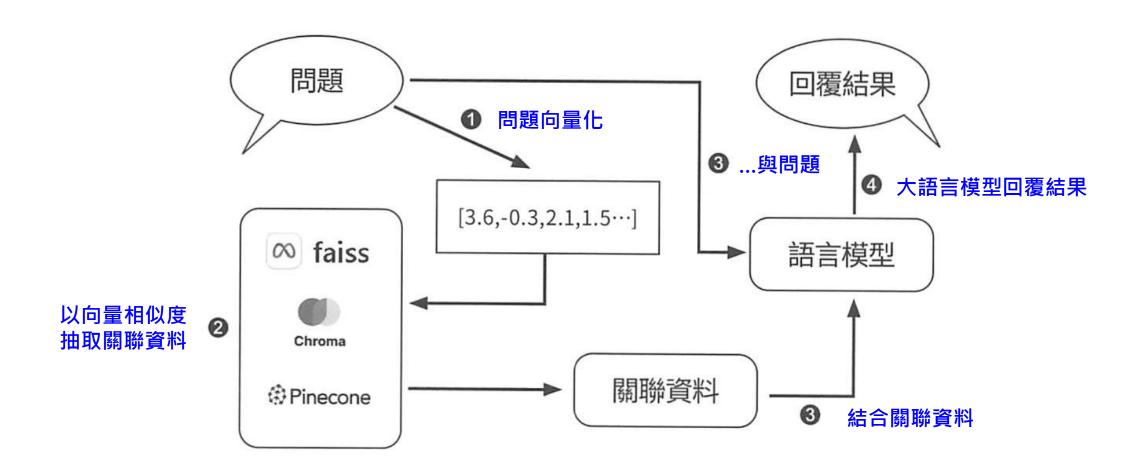






A RAG 運作原理













RAG 實例講解

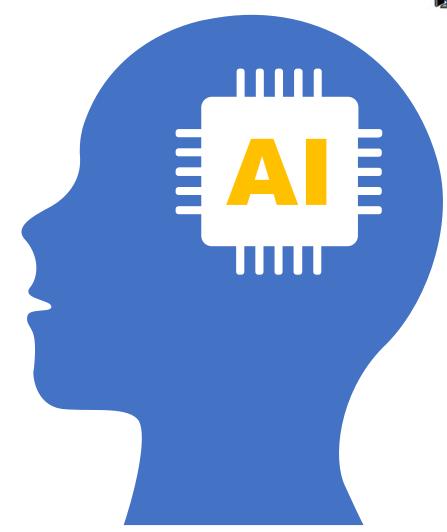
範例完整原始碼:

https://lurl.cc/2w9nLm



RAG 實例講解

- 環境安裝
- 直接連上大語言模型(不使用 RAG)
- 整合大語言模型與 RAG
- 以圖形介面互動



安裝 GPU 驅動相關套件



```
1 # 更新 Linux 內的套件清單至最新版
2 !apt-get update
3
4 # 安裝 PCI 匯流排工具 (PCI Utility) 與 lshw (LiSt HardWare), 以便能偵測到 GPU
5 # -y:遇到詢問是否安裝, 一律自動回答 yes
6 !apt-get install -y pciutils lshw
7
8 # 用 nVidia 的 System Management Interface (SMI) 確認 GPU 的確抓得到
9 !nvidia-smi
```

Driver Version: 535.104.05 CUDA Version: 12.2 I GPU Name Persistence-M | Bus-Id Disp.A | Volatile Uncorr. ECC I Fan Temp Perf Pwr:Usage/Cap | Memory-Usage | GPU-Util Compute M. O Tesla T4 Off | 00000000:00:04.0 Off I N/A 65C P8 11W / 70W I OMiB / 15360MiB I Default | Processes: GPU GI CI PID Type Process name GPU Memory Usage I No running processes found

安裝 Ollama



```
1 # 至 https://ollama.com/download/linux
2 # 直接將安裝 Ollama 於 Linux 的指令貼上
3 !curl -fsSL https://ollama.com/install.sh | sh
5 # 啟動 Ollama, 讓它執行於背景中
6 # ollama serve: 用 Server 模式、而非互動模式執行 Ollama
7 # > server.log:將本應顯示於螢幕的訊息,轉向輸出至 server.log 這個檔備查
8 # 2>&1:2 為 stderr。將所有錯誤訊息, 轉向 &1 (stdout, 螢幕) 輸出。
9 # &:將程式啟動之後,馬上返回,不要等該程式執行完成
10 !ollama serve > server.log 2>&1 &
11
12 # 將 Llama-3 模型下載,並做為此次的大語言模型
13 # ollama run <模型名稱>:下載並執行特定 LLM
14 # > model.log:將本應顯示於螢幕的訊息,轉向輸出至 model.log 這個檔備查
15 !ollama run gemma:7b > model.log 2>&1 &
16
  # 將 nomic-embed-text 模型, 拉下來備用
  # nomic-embed-text 是一款能將文字向量化的嵌入層模型
                                                ← 這次多裝這個 Ollama 模型
  !ollama pull nomic-embed-text > embedding.log 2>&1 &
20
21 # 注意:上述指令皆以「&」後綴, 告知 Colab「不用等 Linux 執行完」。
22 # 但事實上,不論啟動為 Server,或下載 LLM Models,皆須 1~5 分鐘不等的時間。
23 # 可以查看 *.log 檔案內容, 得知當前執行狀況。
```

安裝其它相關套件



```
1 # 下載 LangChain,一套專門連上各種大語言模型的 Python 套件
    2 !pip install langchain # LangChain 核心元件
      !pip install langchain-community # 各種開源大語言模型連接函數套件
2
    5 # ChromaDB 是一款蠻受歡迎的開源式向量資料庫
      !pip install chromadb
3
    8 # 將做為範例的 PDF 檔下載下來
    9 import os
      Dataset File = "GoingLearn BookReviews.pdf"
   11
   12 if not os.path.isfile(Dataset File):
        os.system("wget https://raw.githubusercontent.com/cnchi/datasets/master/" + Dataset File)
   14
   15 # 安裝 PDF 檔內容解析套件
   16 !pip install pypdf
   17
   18 # 由 Hugging Face 開發,開源式網頁圖形介面套件
   19 # 詳見官方首頁 https://www.gradio.app/
   20 !pip install gradio
```

隨堂練習:設定相關環境



- •請依照前述投影片內容,分別撰寫並執行下列程式碼:
 - 安裝 GPU 驅動相關套件
 - 安裝 Ollama (Ollama, gemma:7b, nomic-embed-text)
 - 安裝大語言模型相關套件 (LangChain)
 - 安裝向量資料庫相關套件(ChromaDB)
 - 安裝外部文件檔與套件 (PyPDF)
 - 安裝圖形式介面相關套件(Gradio)

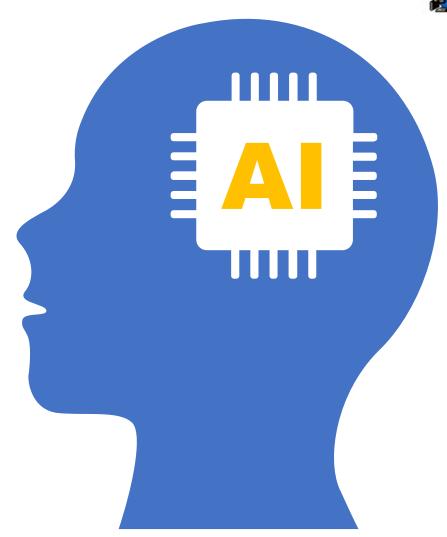






RAG 實例講解

- 環境安裝
- 直接連上大語言模型(不使用 RAG)
- 整合大語言模型與 RAG
- 以圖形介面互動



程式碼與執行結果



```
1 # 載入 Ollama 以便連上後端 LLM
2 from langchain_community.llms import Ollama
3
2 4 # LLM 名稱需與 ollama run 後方名稱相同
5 llm = Ollama(model="gemma:7b")
6
3 7 # 對 LLM 送出提詞,並且印出回應
8 msg = llm.invoke("紀俊男老師是怎麼樣的一個人?請用繁體中文回答我。")
9 print(msg)
```

紀俊男老師以其熱情、有情和有包容性而著稱。他是一位充滿熱情和有禮貌的教師,並努力建立一個包容和支持性的班級環境。以下是一些關於紀俊男老師的人格特徵:

- * 熱情滿滿:紀俊男老師對教學的熱情和動力是令人驚嘆的。他積極參與課程的準備和實踐,並努力引導學生思考和成長。
- *有情有為:紀俊男老師會主動關心學生的生活狀況,並提供他們有需要的支持。他也積極參與校園活動,並為學生爭取更多權益。
- * 有包容性:紀俊男老師會接受並尊重學生的不同觀點和生活方式。他也努力建立一個包容和支持性的班級環境,讓所有學生都能安心學習。
- * 有禮貌有成:紀俊男老師有良好的禮貌和溝通技巧。他也擁有充足的教導經驗和專業素養

總之,紀俊男老師是一位熱情、有情和有包容性的教師,值得學生和家長信賴。

一本正經 胡說八道



隨堂練習:不使用 RAG



•請依照前述投影片內容,撰寫並執行下列程式碼:

```
1 # 載入 Ollama 以便連上後端 LLM
2 from langchain_community.llms import Ollama
3
4 # LLM 名稱需與 ollama run 後方名稱相同
5 llm = Ollama(model="gemma:7b")
6
7 # 對 LLM 送出提詞,並且印出回應
8 msg = llm.invoke("紀俊男老師是怎麼樣的一個人?請用繁體中文回答我。")
9 print(msg)
```

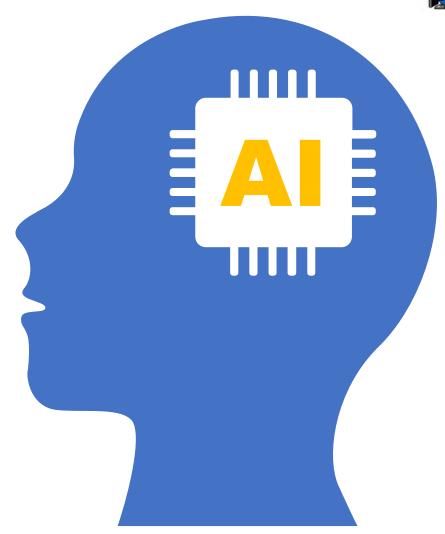






RAG 實例講解

- 環境安裝
- 直接連上大語言模型(不使用 RAG)
- 整合大語言模型與 RAG
- 以圖形介面互動



載入外部資料



• (1) 載入網頁資料

載入外部資料



• (2) 載入 PDF 資料

我是如何選書、切書、看書、評書的?



常有同學問我:「老師,我想學 XXX,可以幫我推薦幾本書嗎?」因為擔任補 習班老師的關係,會有蠻多機會推薦書籍給同學。一般而言,我就把口袋名單 提供給同學就完事了。但有些補習班的「常客」,跟我跟久了,知道我是「重 度書蟲」。不僅擁有兩萬多本電腦書,出門還一定得帶個一兩本在包包裡才安 心的那種人。所以有些人就好奇:

- 「老師,你是怎麼選書、看書的?」
- 「你怎麼有辦法兩個禮拜就看完一本電腦書?」
- 「我看你都用 iPad 在看書,那些書是怎麼電子化的?」
- 「你是怎麼把好幾本書的內容,那麼快就融會質通的?」

- 1 # 載入 PDF 資料
 - 2 from langchain_community.document_loaders import PyPDFLoader
- 2 4 loader = PyPDFLoader(Dataset_File)
 - 5 pdf_docs = loader.load()
- 7 print(pdf_docs)

[Document(page_content='我是如何選書、切書、看書、評書的?\n常有同學問我:「老師‧我想學XXX‧可以幫我推薦幾本書嗎?」因為擔任補\n習班老師的關係‧會有蠻多機會推薦書籍給同學。一般而言‧我就把口袋名單\n提供給同學就完事了。但有些補習班的「常客」‧跟我跟久了知道我是「重\n度書蟲」。不僅擁有兩萬多本電腦書‧出門還一定得帶個一兩本在包包裡才安\n心的那種人。所以有些人就好奇:\n「老師‧你是怎麼選書、看書的?」\n「你怎麼有辦法兩個禮拜就看完一本電腦書?」\n....

載入外部資料



• (3) 將所有的資料接合在一起

```
1 # 將所有資料接合在一起
2 all_docs = web_docs + pdf_docs
3
4 print(all_docs)
```

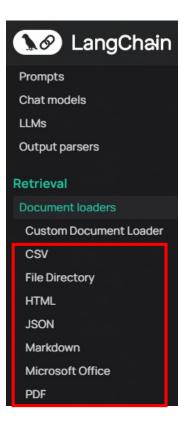
到底可以載入多少種外部資料?



LangChain Documents

• https://python.langchain.com/v0.1/docs/modules/data_connection/docum

ent_loaders/



可以載入這麼多種外部資料



隨堂練習:載入外部資料



•請依照前述投影片內容,撰寫並執行下列程式碼:

```
1 # 載入網頁資料
 2 from langchain community.document loaders import WebBaseLoader
 4 url = "https://goinglearn.com.tw/about-this-site"
 5 loader = WebBaseLoader(url)
 6 web docs = loader.load()
   print(web docs)
10 # 載入 PDF 資料
11 from langchain community.document loaders import PyPDFLoader
12
13 loader = PyPDFLoader(Dataset File)
   pdf docs = loader.load()
15
   print(pdf docs)
17
18 # 將所有資料接合在一起
19 all_docs = web_docs + pdf_docs
20
21 print(all docs)
```





資料分割



```
1 # 將載入的文件以每段 1000 Tokens、重複 200 Tokens 的方式切分
2 # 可保證原文不會因切分而失去上下文
3 # 也可以讓 LLM 透過重複部分,得知誰是上文、誰是下文
4 from langchain.text_splitter import RecursiveCharacterTextSplitter
5
6 text_splitter = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=200)
7 splits = text_splitter.split_documents(all_docs)
8
9 for split in splits:
    print(split)
```



資料向量化 & 存入向量資料庫



```
1 # 利用 nomic-embed-text 模型, 執行文本向量化
2 from langchain_community.embeddings import OllamaEmbeddings
3 from langchain_community.vectorstores import Chroma
4
5 embeddings = OllamaEmbeddings(model="nomic-embed-text")
6 vector_db = Chroma.from_documents(documents=splits, embedding=embeddings)
```

隨堂練習:資料分割與向量化



• 請依照前述投影片內容,撰寫並執行下列程式碼:

```
1 # 將載入的文件以每段 1000 Tokens、重複 200 Tokens 的方式切分
 2 # 可保證原文不會因切分而失去上下文
  # 也可以讓 LLM 透過重複部分,得知誰是上文、誰是下文
  from langchain.text splitter import RecursiveCharacterTextSplitter
 5
   text splitter = RecursiveCharacterTextSplitter(chunk size=1000, chunk overlap=200)
   splits = text splitter.split_documents(all_docs)
   for split in splits:
10
      print(split)
11
  # 利用 nomic-embed-text 模型, 執行文本向量化
   from langchain community.embeddings import OllamaEmbeddings
   from langchain community vectorstores import Chroma
15
  embeddings = OllamaEmbeddings(model="nomic-embed-text")
17 vector_db = Chroma.from_documents(documents=splits, embedding=embeddings)
```





定義 RAG 與 LLM 相關函數



• (1) 定義 LLM 詢問函數

```
1 # 定義一個能丟入「使用者問題」與「RAG 抽取之相關文本」,用來詢問 LLM 的函數 2 def llm_ask(question, rag_contents):
2 3 # 製作提問用的提詞 (Prompt) formatted_prompt = f"Question: {question}\n\nContext: {rag_contents}"
3 5 # 丟入 LLM, 取得答案 msg = llm.invoke(formatted_prompt) return msg
```

定義 RAG 與 LLM 相關函數



• (2) 定義 RAG 相關資料抽取函數

```
# 定義向量資料庫抽取物件
     retriever = vector db.as retriever()
  11
      # 根據與使用者問題的相似度, 抽取相關資料
     def get rag docs(question):
3 14
         # 依照與問題的相似度, 抽取相關文本
         retrieved docs = retriever.invoke(question)
  15
         # 將所有相關文本, 以空行 \n\n 黏成一個大字串
4 16
         relavent contents = "\n\n".join(doc.page content for doc in retrieved docs)
  17
  18
5 19
         # 將問題 + RAG 相關文本, 一口氣丟入 LLM, 取得答案
         return 11m ask(question, relavent contents)
   20
```

△ 陸

隨堂練習:LLM 與 RAG 函數定義



•請依照前述投影片內容,撰寫並執行下列程式碼:

```
1 # 定義一個能丟入「使用者問題」與「RAG 抽取之相關文本」,用來詢問 LLM 的函數
 2 def llm ask(question, rag contents):
      # 製作提問用的提詞 (Prompt)
      formatted prompt = f"Question: {question}\n\nContext: {rag contents}"
      # 丟入 LLM, 取得答案
      msg = llm.invoke(formatted prompt)
      return msg
  # 定義向量資料庫抽取物件
   retriever = vector db.as retriever()
11
   # 根據與使用者問題的相似度, 抽取相關資料
  def get rag docs(question):
      # 依照與問題的相似度, 抽取相關文本
14
      retrieved docs = retriever.invoke(question)
15
      # 將所有相關文本, 以空行 \n\n 黏成一個大字串
16
      relavent contents = "\n\n".join(doc.page content for doc in retrieved docs)
17
18
      # 將問題 + RAG 相關文本, 一口氣丟入 LLM, 取得答案
19
20
      return llm ask(question, relavent contents)
```

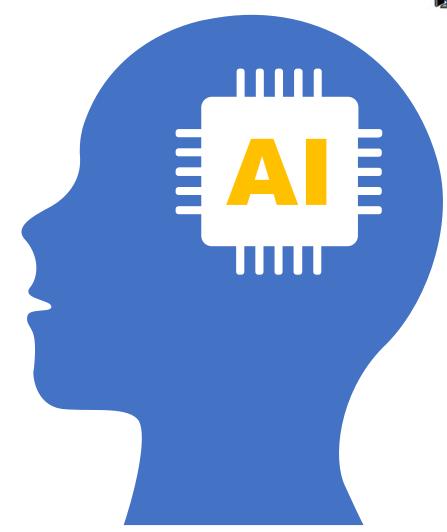






RAG 實例講解

- 環境安裝
- 直接連上大語言模型(不使用 RAG)
- 整合大語言模型與 RAG
- •以圖形介面互動



定義圖形介面



```
1 # 用 Gradio 建立圖形介面,並用 get_rag_docs 做為處理函數
   import gradio as gr
   gui = gr.Interface(
      fn=get_rag_docs, ← 會收取 inputs= 傳回來的數值
      inputs=gr.Textbox(lines=2, placeholder="在此輸入您的問題..."),
 6
      outputs="text",
      title="紀老師答客問", —
      description="請在這邊輸入您想提
 9
                                                    請在追及論,你想提問的問題
10
                                                                                   **紀後男老師的個人特徵: **
                                                      紀俊男老師是怎麼樣的一個人?請用繁體中文回答我。
                                                                                   *具有濃厚的電腦興趣,從1982年開始對電腦產生興趣。
                                                                                   *從 BASIC、Pascal、C、組合語言等編程語言開始學習,並深入研
                                                                                   究網頁設計、手機應用程式開發和嵌入式系統。
                                                          Clear
                                                                       Submit
                                                                                   * 熱愛閱讀,並有低調動漫宅的傾向。
                                                                                   *在講台上很 boisterous,但下講台卻安靜。
                                                                                   **書籍評價: **
                                                                                   紀俊男老師沒有在書籍中出現,因此無法評價書籍。
                          question
                                                                                             Flag
```

△ 啟動圖形介面



- 12 # 啟動圖形介面
- 13 gui.launch()



隨堂練習:以圖形介面互動



•請依照前述投影片內容,撰寫並執行下列程式碼:

```
1 # 用 Gradio 建立圖形介面,並用 get_rag_docs 做為處理函數
   import gradio as gr
   gui = gr.Interface(
     fn=get rag docs,
     inputs=gr.Textbox(lines=2, placeholder="在此輸入您的問題..."),
     outputs="text",
     title="紀老師答客問",
     description="請在這邊輸入您想提問的問題",
10
                                                 紀老師答客問
11
                         請在這邊輸入您想提問的問題
   # 啟動圖形介面
   gui.launch()
                          紀老師平時的嗜好是什麼?
                                                       **問題:紀老師平時的嗜好是什麼? **
                                                       **答案:**紀老師平時的嗜好是閱讀書籍。文章提到紀老師喜歡坐在咖啡廳陽光斜灑
                               Clear
                                                                  Flag
```





本章總結



• 大語言模型的困境

- 無法回答太新的資料。
- 無法回答太專業、特定領域的資料。

• 檢索增強式文本生成(RAG)

• 利用網路搜索、資料庫文本,補強訓練時語料庫太舊、或太廣泛的問題。

• RAG 資料來源的生成步驟

- 資料載入 (Data Loading)
- 資料分割 (Data Chunking)
- 資料嵌入化 (Data Embedding)
- 儲存於向量資料庫 (Vector Databases)



