# 自然語言處理與應用
## Natural Language Processing and Applications

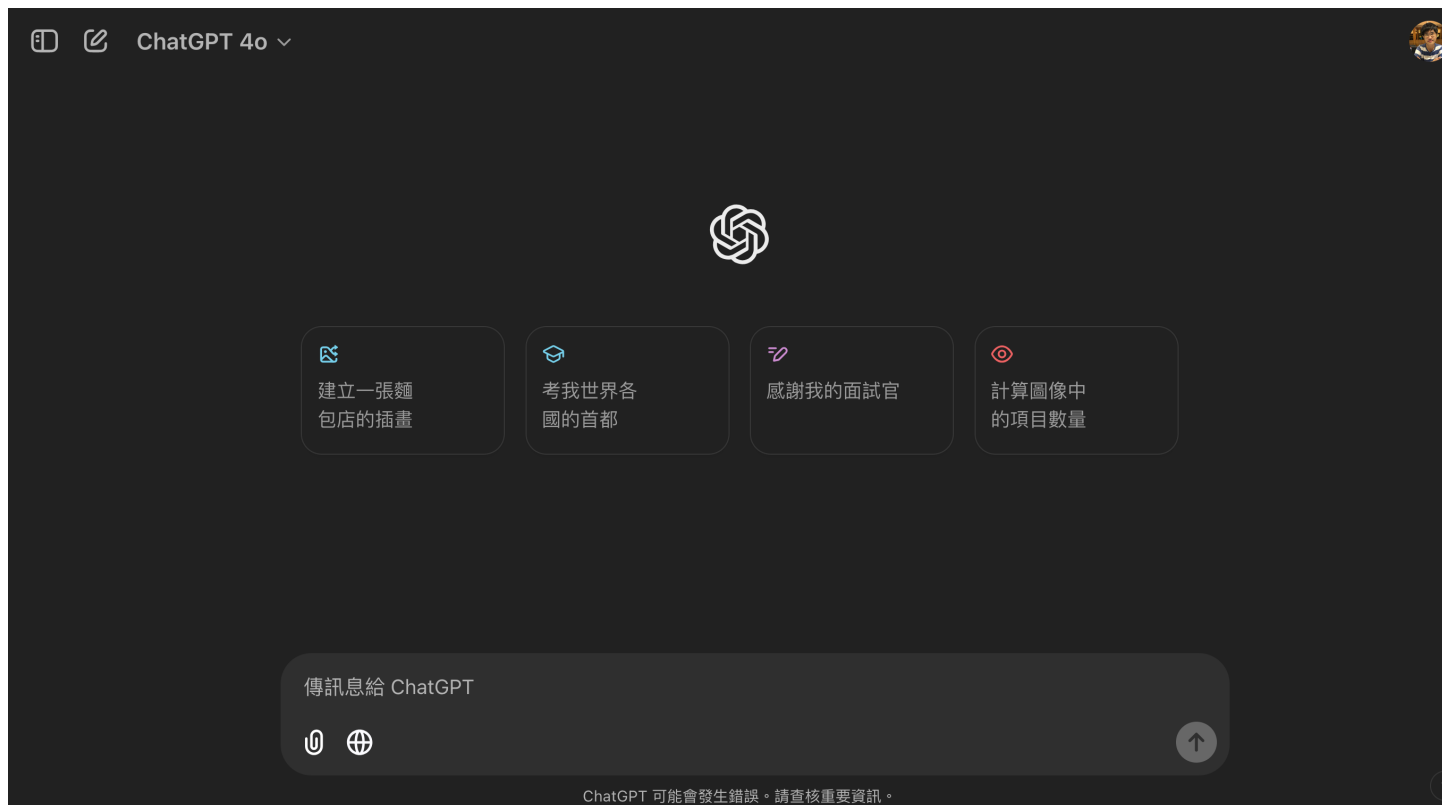**LLM API Tutorial**

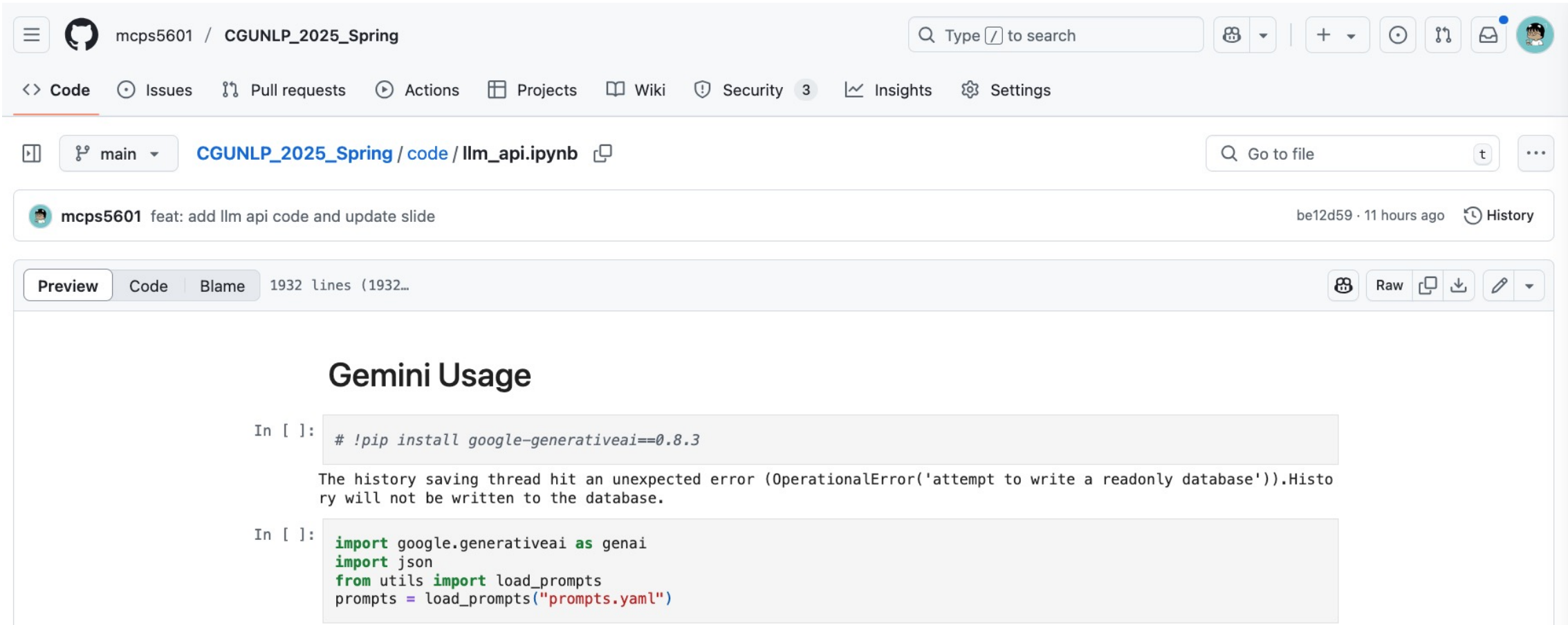Instructor: 林英嘉 (Ying-Jia Lin)
2025/05/12

Course GitHub

Slido # NLP_0512

# Why do we need to use API?



- It's slow if you want to use ChatGPT for NLP tasks by manually copying and pasting from the webpage.
- If you use ChatGPT from the webpage to test data, "Too many requests in 1 hour. Try again later." may occur.

NLP

# Code Example

# Llama 4:
## Leading Multimodal Intelligence

Newest model suite offering unrivaled speed and efficiency

### Llama 4 Behemoth

**288B** active parameter, **16** experts
**2T** total parameters

The most intelligent teacher model for distillation

Preview

### Llama 4 Maverick

**17B** active parameters, **128** experts
**400B** total parameters

Native multimodal with **1M** context length

Available

### Llama 4 Scout

**17B** active parameters, **16** experts
**109B** total parameters

Industry leading **10M** context length
Optimized inference

Available

# 註冊 Google Vertex

https://cloud.google.com/vertex-ai?hl=zh-TW

# 進到控制台

https://console.cloud.google.com

右邊導覽選單

Vertex AI

Model Garden

# 選取專案

ID 等等會用到

NLP

# 搜尋llama4



點這個

# 閱讀 EULA 後點啟用

# 打指令在終端機：gcloud access token



把 API Key 複製起來！

# Thank you!

Instructor: 林英嘉

✉ yjlin@cgu.edu.tw

TA: 吳宣毅

✉ m1161007@cgu.edu.tw