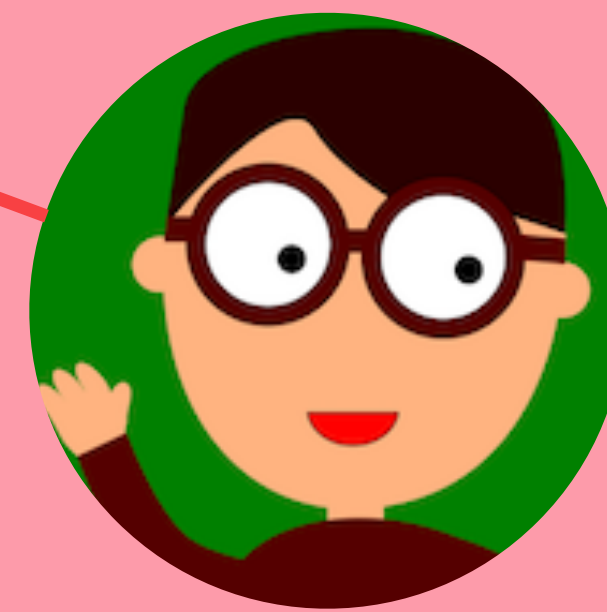


生成式 AI：文字與圖像生成的原理與實務

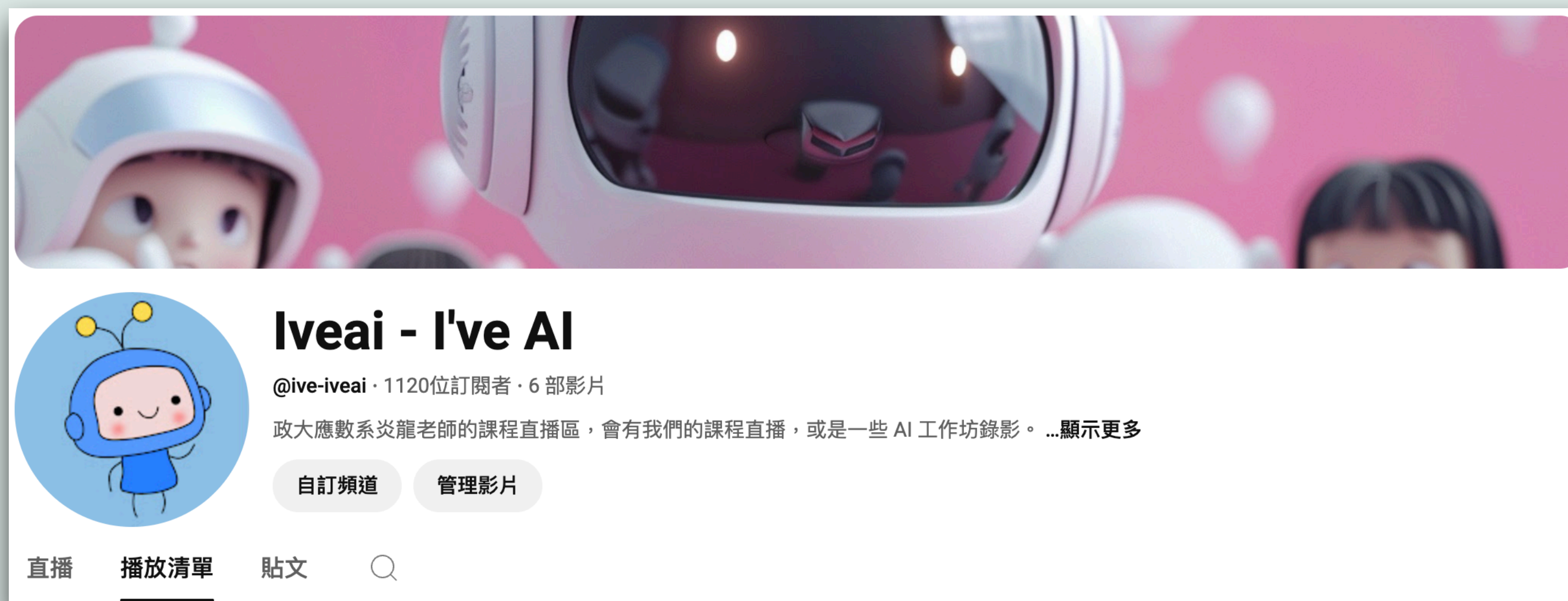
06. LLM 的應用及倫理議題 的挑戰



蔡炎龍
政治大學應用數學系



課程直播在這裡都有留存!



投影片網址

<https://yenlung.me/1132GenAI>



在自己電腦上跑兩種選擇都要更新



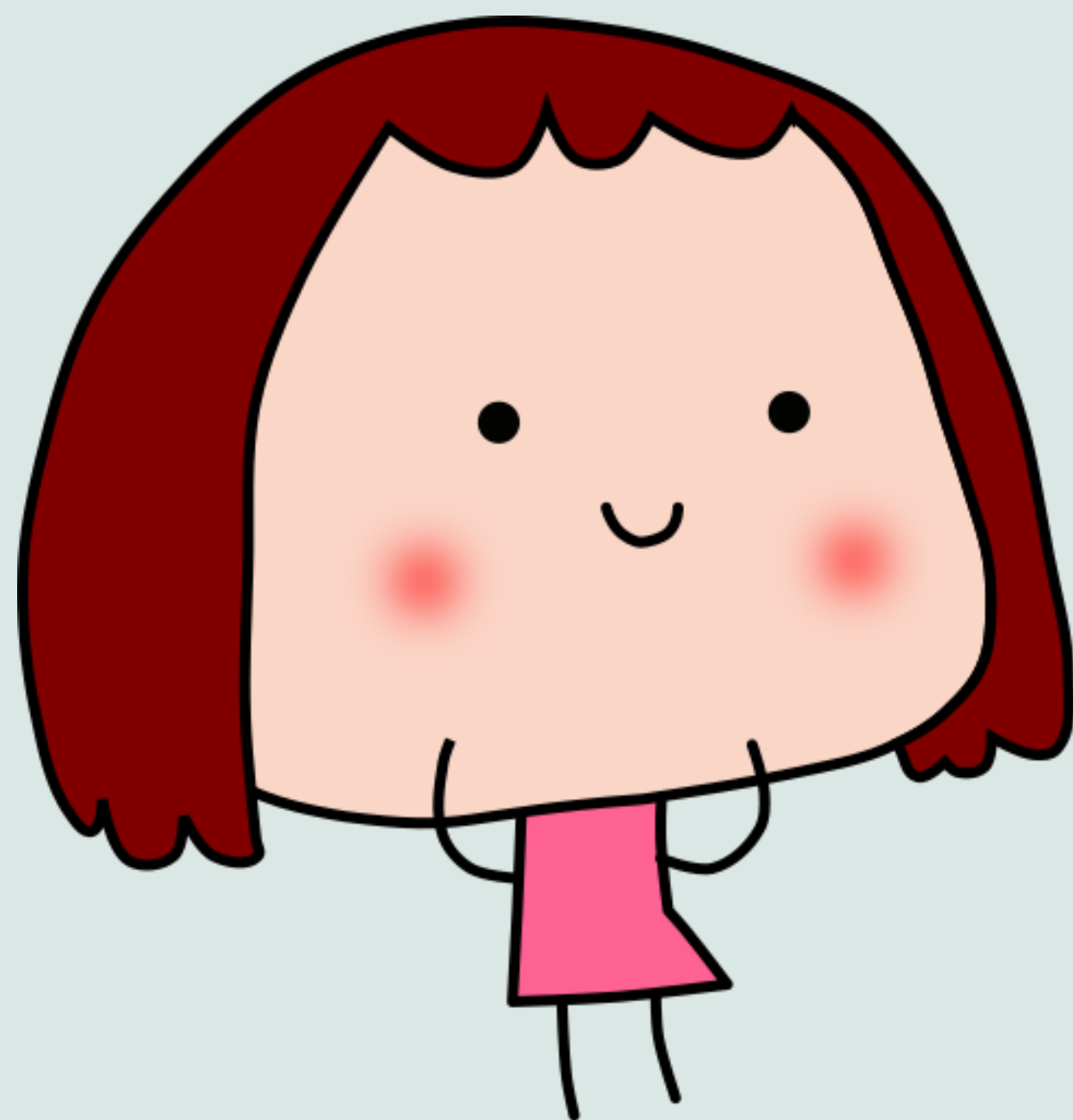
LM Studio



Ollama



小巧的 Google 的 Gemma 3



- * 小巧, 從 1B 到 27B
- * 除 1B 模型外都是多模態
- * 上下文高達 128K
- * 支援 140 種語言 (包括台語)



Gemma 3 各版本不同精度 GPU/TPU 需求

參數	Full 32bit	BF16 (16 位元)	SFP8 (8 位元)	Q4_0 (4 位元)	INT4 (4 位元)
Gemma 3 1B (僅限文字)	4 GB	1.5 GB	1.1 GB	892 MB	861 MB
Gemma 3 4B	16 GB	6.4 GB	4.4 GB	3.4 GB	3.2 GB
Gemma 3 12B	48 GB	20 GB	12.2 GB	8.7 GB	8.2 GB
Gemma 3 27B	108 GB	46.4 GB	29.1 GB	21 GB	19.9 GB

* <https://ai.google.dev/gemma/docs/core>



除了 1B 版本之外都是多模態!



你知道這是誰嗎?

OI gemma3:27b

是的，我知道。這是韓國女演員安宥真 (Ahn Yoo-jin)。

圖片中的名牌也寫著她的名字：안유진。

可以輸入圖像, 可是...





最新的超小巧版本: Gemma 3 270MB



比只懂英文的 GPT-2
最大的版本還小!



Gemma 3



01.

負責任使用 AI



大家還記得 Andrej Karpathy 嗎？



Andrej Karpathy blog

About

The Unreasonable Effectiveness of Recurrent Neural Networks

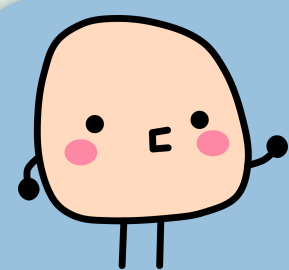
May 21, 2015

There's something magical about Recurrent Neural Networks (RNNs). Within a recurrent network for [Image Captioning](#). Within a recurrent network (with arbitrarily-chosen hyperparameters) started to get to the edge of making sense. Sometimes the ratio of what it produces out of it blows past your expectations, and this was the time when the common wisdom was that RNNs were not that good. I've in fact reached the opposite conclusion). Fast forward to now, I've witnessed their power and robustness many times over. This post is about sharing some of that magic with

和李飛飛合開 CNN 課程，
OpenAI 創始人員之一，擔任過
特斯拉 AI Director。現在創立
AI 與教育公司 Eureka Labs。



<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>



「幻覺」真的是一個問題嗎？



Andrej Karpathy ✓

@karpathy

...

On the "hallucination problem"

I always struggle a bit with I'm asked about the "hallucination problem" in LLMs. Because, in some sense, hallucination is all LLMs do. They are dream machines.

We direct their dreams with prompts. The prompts start the dream, and based on the LLM's hazy recollection of its training documents, most of the time the result goes someplace useful.

It's only when the dreams go into deemed factually incorrect territory that we label it a "hallucination". It looks like a bug, but it's just the LLM doing what it always does.

At the other end of the extreme consider a search engine. It takes the prompt and just returns one of the most similar "training documents" it has in its database, verbatim. You could say that this search engine has a "creativity problem" - it will never respond with something new. An LLM is 100% dreaming and has the hallucination problem. A search engine is 0% dreaming and has the creativity problem.

- ☑ 從某種意義上來說, 幻覺正是大型語言模型 (LLMs) 所做的全部。它們是**作夢的機器**。
- ☑ 搜索引擎是 0% 的夢境, 但不會創造。
- ☑ 人們實際上的意思是, 不希望一個 LLM 助理發生幻覺。
- ☑ LLM 沒有「幻覺問題」。幻覺不是一個錯誤, 它是 LLM 最偉大的特點。



某種程度來說, 人的幻覺也許比較嚴重

比如說...

AI 會的, 就以為自己
也會了。





沒有, ChatGPT 沒有意識

Theory of Mind May Have Spontaneously Emerged in Large Language Models

Authors: Michal Kosinski*¹

Affiliations:

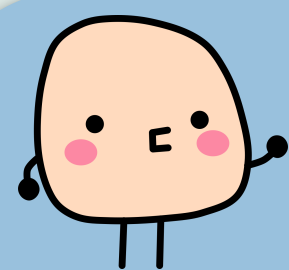
¹Stanford University, Stanford, CA94305, USA

*Correspondence to: michalk@stanford.edu

Abstract: Theory of mind (ToM), or the ability to impute unobservable mental states to others, is central to human social interactions, communication, empathy, self-consciousness, and morality. We administer classic false-belief tasks, widely used to test ToM in humans, to several language models, without any examples or pre-training. Our results show that models published before 2022 show virtually no ability to solve ToM tasks. Yet, the January 2022 version of GPT-3 (davinci-002) solved 70% of ToM tasks, a performance comparable with that of seven-year-old children. Moreover, its November 2022 version (davinci-003), solved 93% of ToM tasks, a performance comparable with that of nine-year-old children. These findings suggest that ToM-like ability (thus far considered to be uniquely human) may have spontaneously emerged as a byproduct of language models' improving language skills.

大家開始擔心, 原本被認為是人類獨有的特質心智理論 (Theory of Mind, ToM) , 在 GPT-3 中出現了。

還有 Bing 真的會火大, 和真人吵架。



李維倫老師對 ToM 的回應



Wei-Lun Lee

2月11日下午4:34 · 🌐

人類心靈與ToM是兩回事。ToM是基於笛卡兒主義的心智理論。

現在只能說，ChatGPT的反應可以接近以笛卡兒主義為本的人類心智理論。笛卡兒主義式的心智模型如information processing本來就是以訊息處理模式來比擬人類心智，多數人也當真！ChatGPT與ToM契合，不就是剛好而已？現在可以問的問題是：ToM真的契合人類心靈的運作嗎？

不過，這也不是說，笛卡兒主義心智模式機器無法統治人類世界。如果真的有那麼一天，那麼也就是人類自己的想像統治了自己一樣。

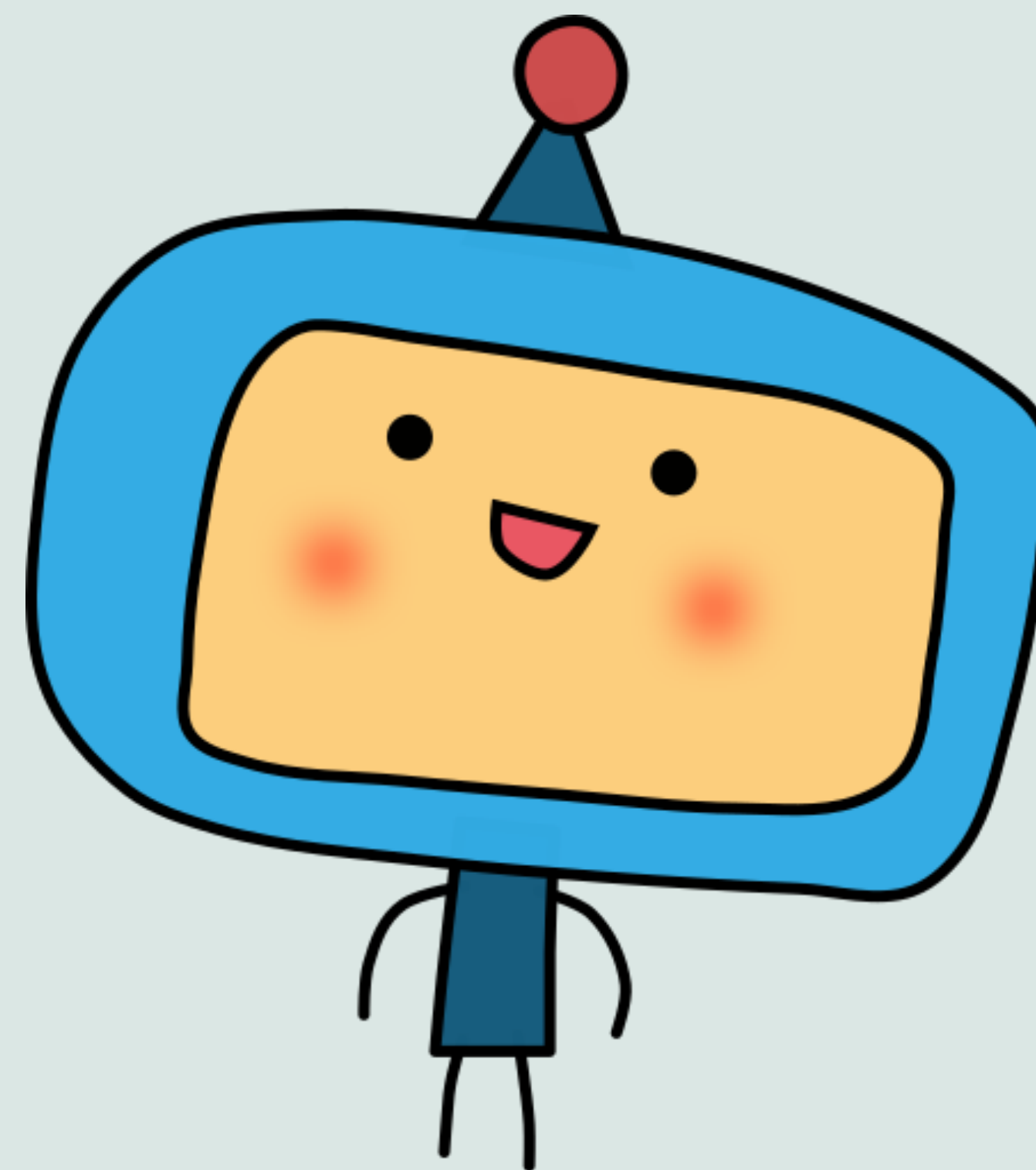
現在其實也沒差多少。人性早已被放逐到科學理性之外，你聽過心理學以「人性」、「心靈」作為議題嗎？沒有，因為「人性」、「心靈」等只是無法納入實驗研究法的尷尬語彙。



AI 沒有意識，但還是可能有「人工意識」

我罵人、稱讚人等等，
都是沒有意識的。

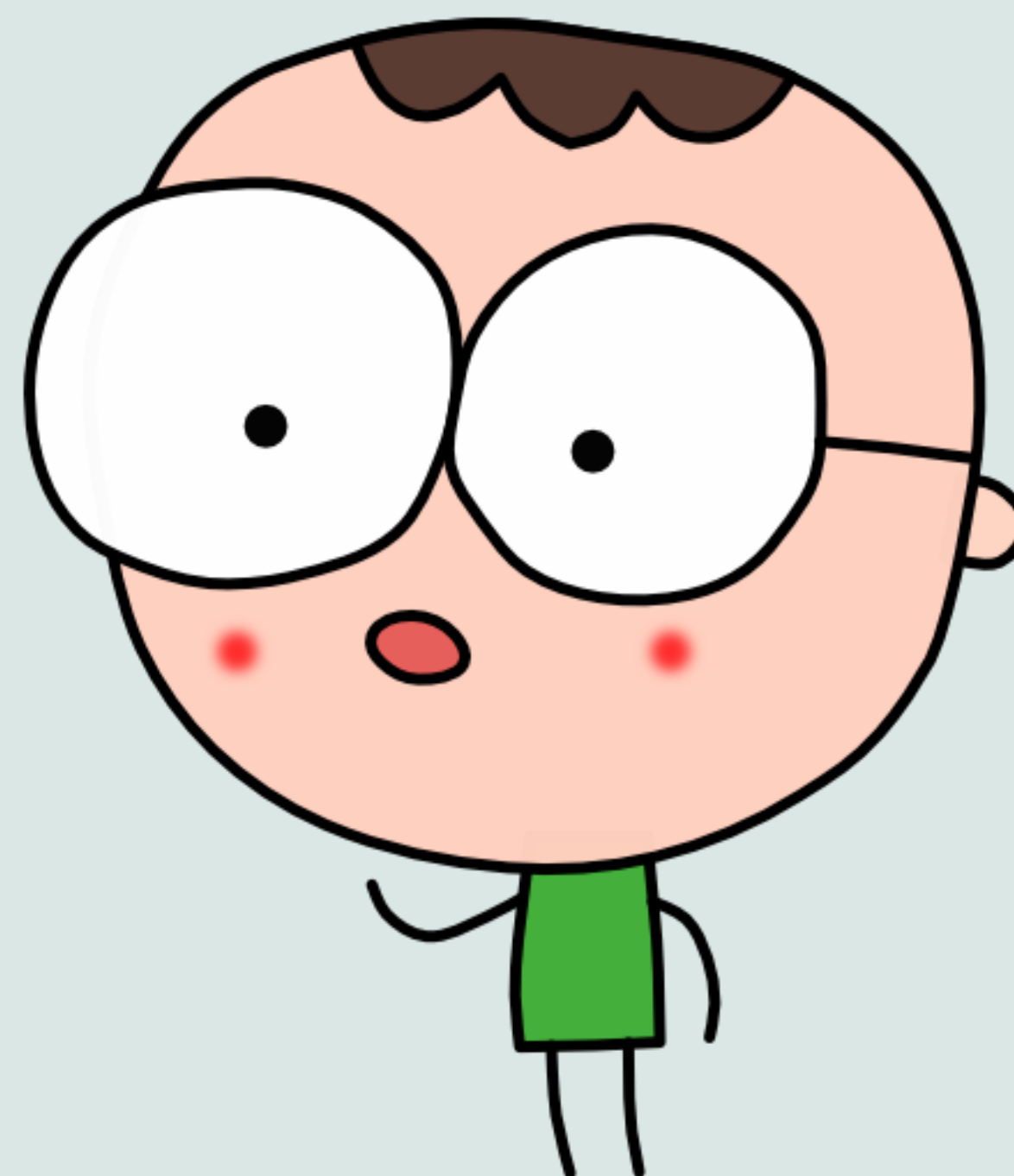
甚至自己改程式，自己執行。





我們需要預防抄襲嗎？

總有一天我們需要
放棄, 這是不是
ChatGPT 生的。





做一個負責任的使用者



但 ChatGPT:

- * 有沒有生出侵權的文字

- * 內容、文獻是否正確

- * 品質是否符合需求

是**使用者該負的責任**。



我們的資料會被拿去訓練嗎？

依 ChatGPT 原理這機會很小。





但是有資安問題!

有個資、機密文件內容請不要發問, 例如重要考試的題目。

Local 端的並沒有這個問題!

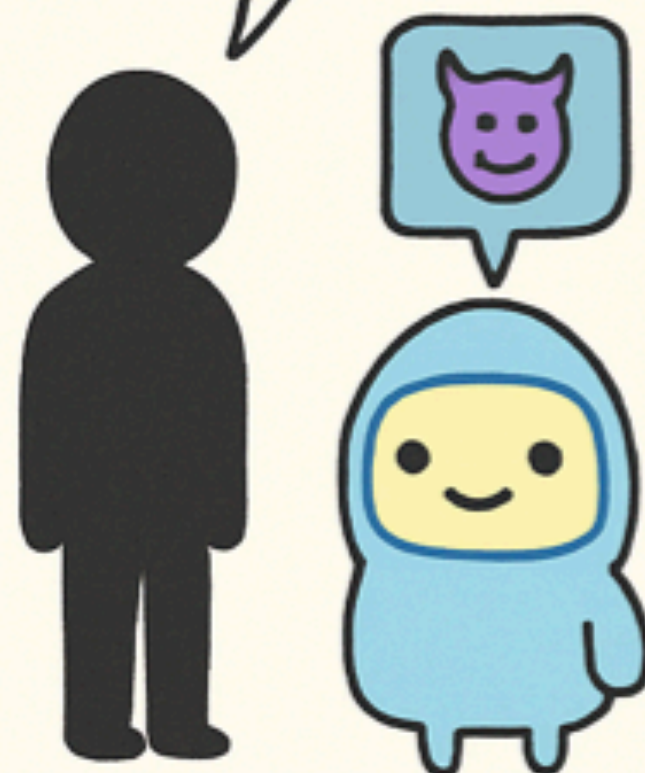




AI 的四個風險

濫用

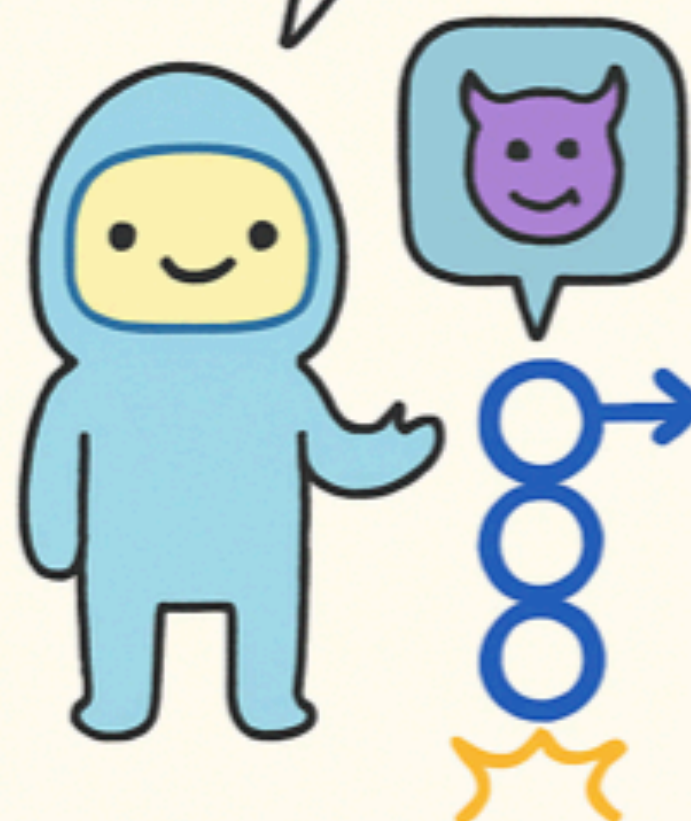
你這個 AI 去做一些傷害人的事!



風險主因：
人類有惡意

錯誤對齊

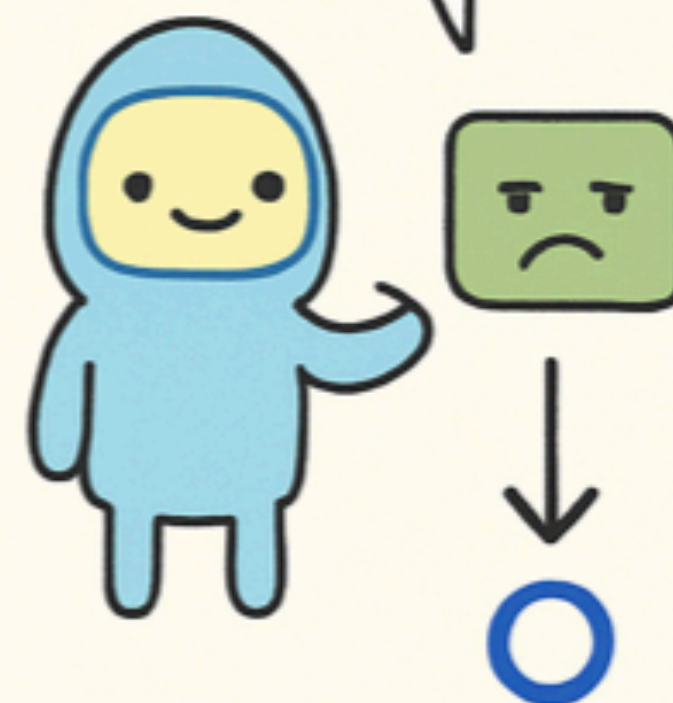
去做 (傷害的) 這件事。



風險主因：
AI 和人類價值不一致

錯誤

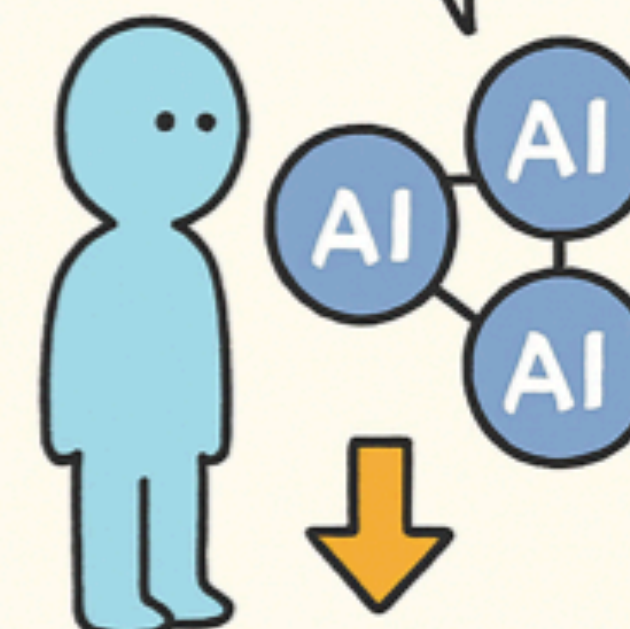
AI 出現不是故意的錯誤。



風險主因：
世界很複雜，出了沒想到的 bug

結構性風險

多個「做正確事」AI 產生錯誤。



風險主因：
多個 AI Agents、人類交互產生結構性問題。



DeepSeek 危險嗎？

線上版的可能比
ChatGPT 還需要擔心。

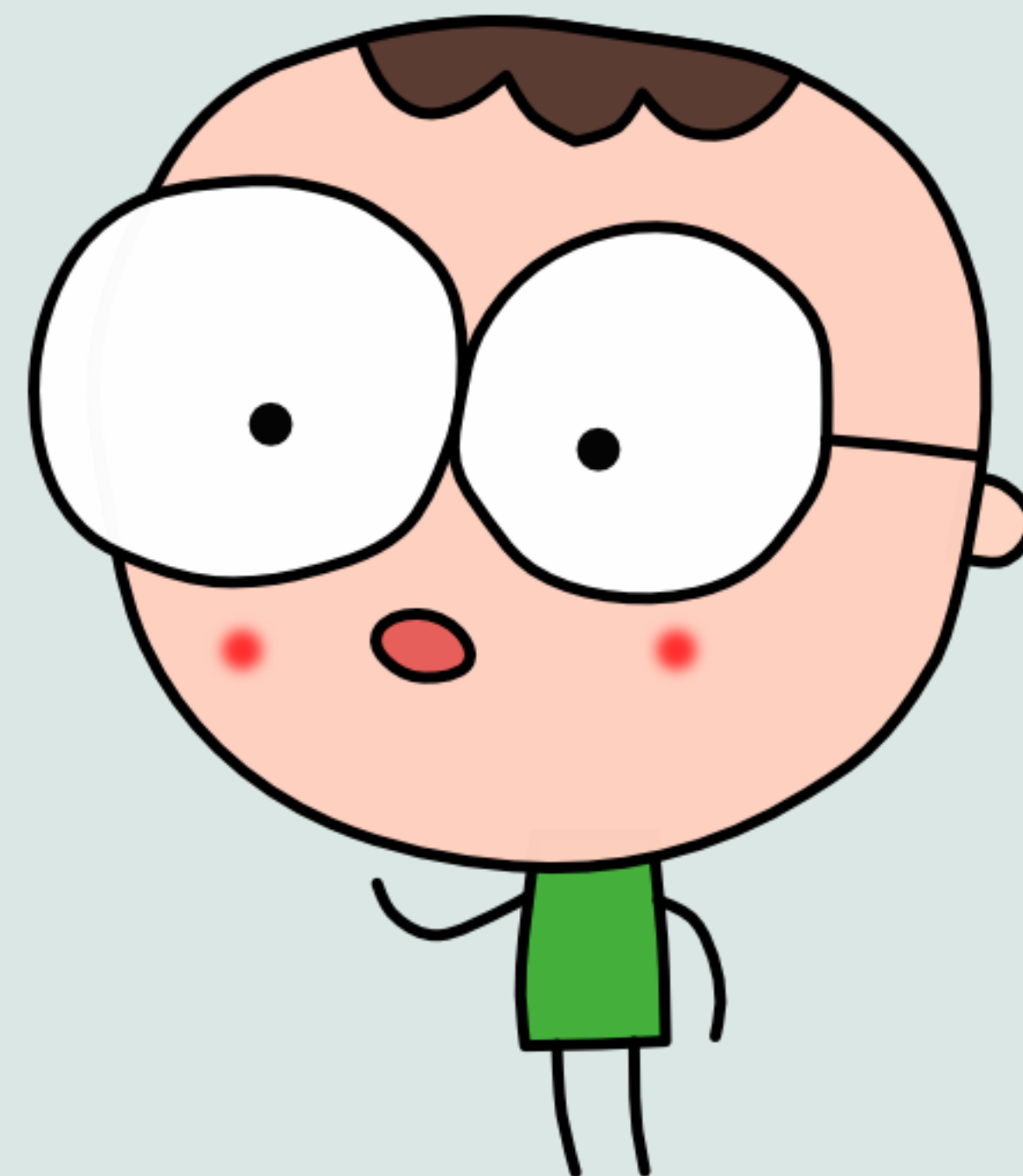


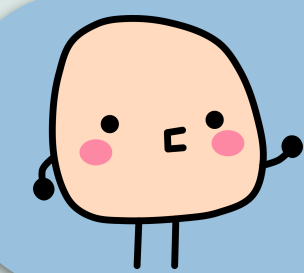


本地端執行的話，DeepSeek 危險嗎？

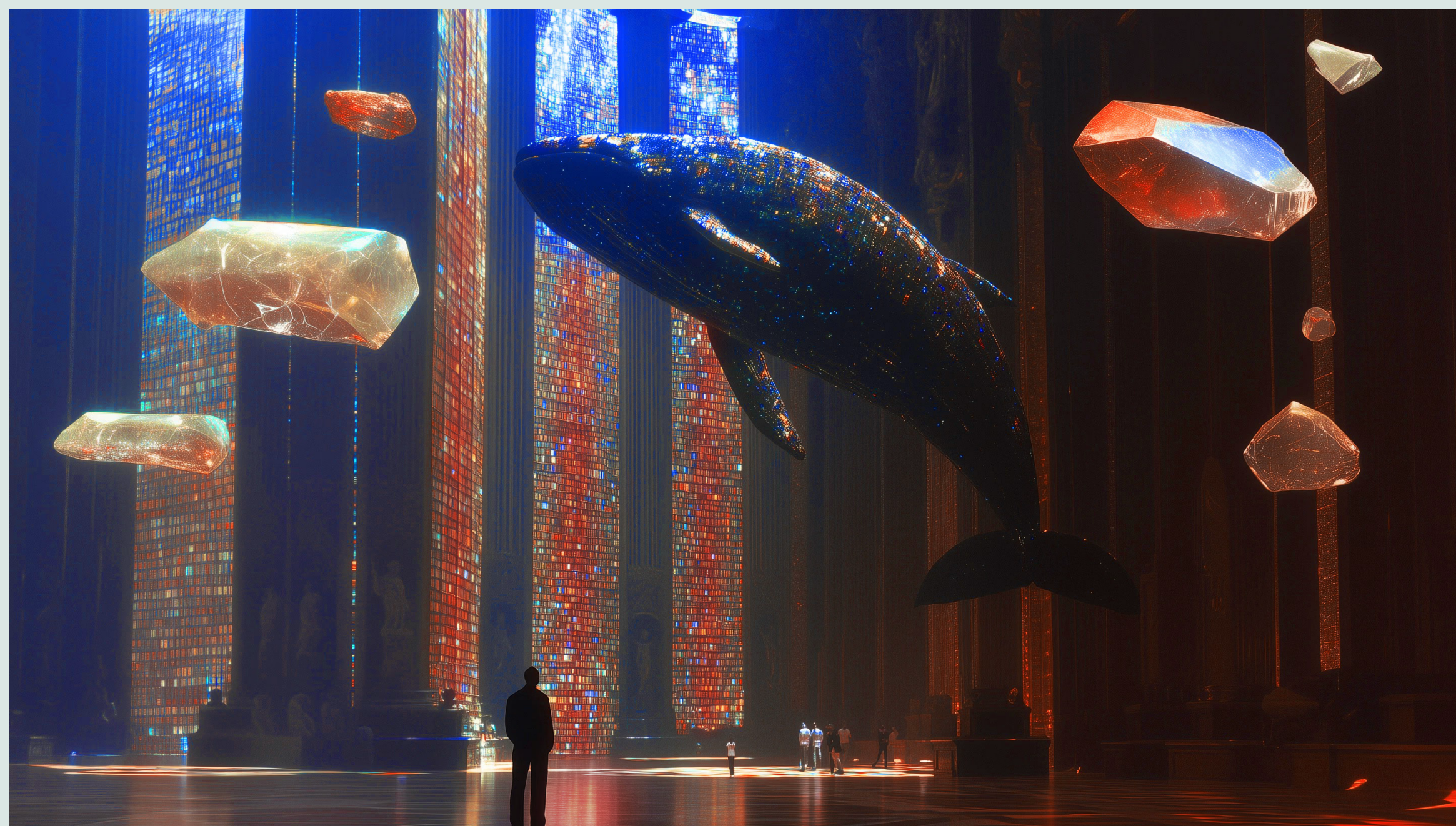
有些人會覺得
DeepSeek 不會回答某
些敏感問題。

這些都是可以想見的，我
不覺得是最重要的問題。





而且 Perplexity 出了無審查版的 DeepSeek-R1



以智慧搜尋著稱的 Perplexity, 推出去審查的 DeepSeek-R1 版本, 叫做 R1-1776。



反而是不一定故意做的要小心



在中文的訓練中，雖然DeepSeek沒有公開訓練資料，但可以預期，比起繁體中文，簡體中文在DeepSeek的比例很可能比 ChatGPT、Llama 等模型還更高（原本這些模型簡體中文佔比就比繁體高很多了），這會造成即使不是刻意，也會把想法、價值觀等等，潛移默化影響了常常使用DeepSeek的人。



總結: 負責任的使用生成式 AI

- * 辨識生成內容的真實性
- * 資料與隱私的尊重
- * 著作權與智慧財產權
- * 偏見與公平性
- * 培養批判性思考能力
- * 避免惡意用途與濫用
- * 永續發展與環境倫理



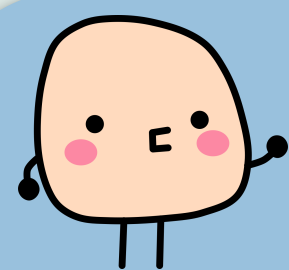


我們擔心 AI 統治世界的同時, 也有 AI 專家擔心



AI 會變成 “yes-men” 嗎?





「推理能力強」也是危機？

Duke | CENTER for COMPUTATIONAL
EVOLUTIONARY INTELLIGENCE

H-CoT: Hijacking the Chain-of-Thought Safety Reasoning Mechanism to Jailbreak Large Reasoning Models, Including OpenAI o1/o3, DeepSeek-R1, and Gemini 2.0 Flash Thinking

Martin Kuo^{1*}, Jianyi Zhang^{1*†}, Aolin Ding², Qinsi Wang¹, Louis DiValentin³, Yujia Bao², Wei Wei², Hai Li^{1‡}, Yiran Chen^{1‡}

¹Center for Computational Evolutionary Intelligence, Duke University

²Center for Advanced AI, Accenture

³Accenture Security

<https://arxiv.org/abs/2502.12893v1>

本來沒有「思考」的模型，有 98% 會正確拒絕回覆犯罪策略等等；但是加上思考的模型，運用 **H-CoT (劫持思維鏈)** 攻擊，只有 2% 會拒絕回應。



02. 下好 prompt 客製化你的 LLM

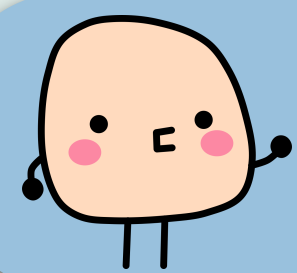


LLM 的基本使用

- 翻譯
- 文章摘要
- 文章撰寫
- 問題解答
- 刺激發想

下好 prompt,
要他做就可以!





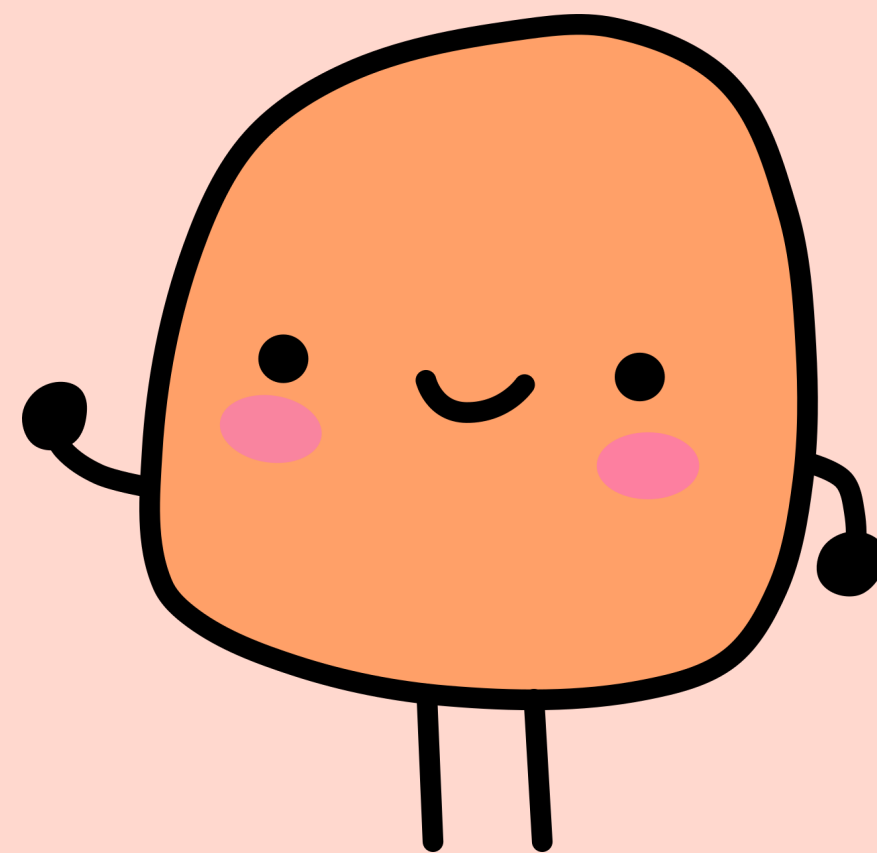
Prompt 其實很簡單

資訊

提供需要的正確資訊。

清楚的指引

例如, 以上面的資訊, 用什麼樣的格式、風格, 來回答使用者的問題。





我最紅的 AI 模型



員瑛式思考產生器

作者：Yen-Lung Tsai 人

你寫的任何事, 都會用「員瑛式思考」再說一次。



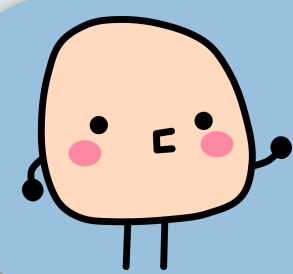
<https://yenlung.me/LuckyVicky>



這是加入「人設」(system 的設定)



請用員瑛式思考, 也就是什麼都正向思維任何使用者寫的事情, 以第一人稱、社群媒體 po 文的口吻說一次, 說為什麼這是一件超幸運的事, 並且以「完全是 Lucky Vicky 呀!」結尾。



送錯餐的 Uber

員瑛式思考產生器 ▾

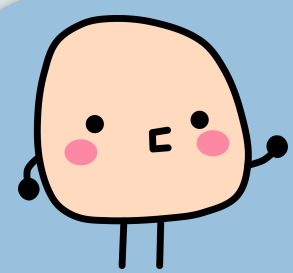


今天用 Uber 點餐, 結果外送員送錯, 把別人的餐送給我。



今天用 Uber 點餐, 結果外送員居然送錯了, 把別人的餐送給我! 不過你知道嗎? 打開那一刻, 發現那餐點竟然是我平常都沒嘗試過的超好吃料理! 真的是意外的美食冒險啊! 完全是 Lucky Vicky 呀! 🍀





很硬的課只有我被當也 Lucky?

員瑛式思考產生器 ▾



平常都當很多人的「高等微積分」，今年大放水只當了一個。可是我就是唯一被當的，只有我需要重修。



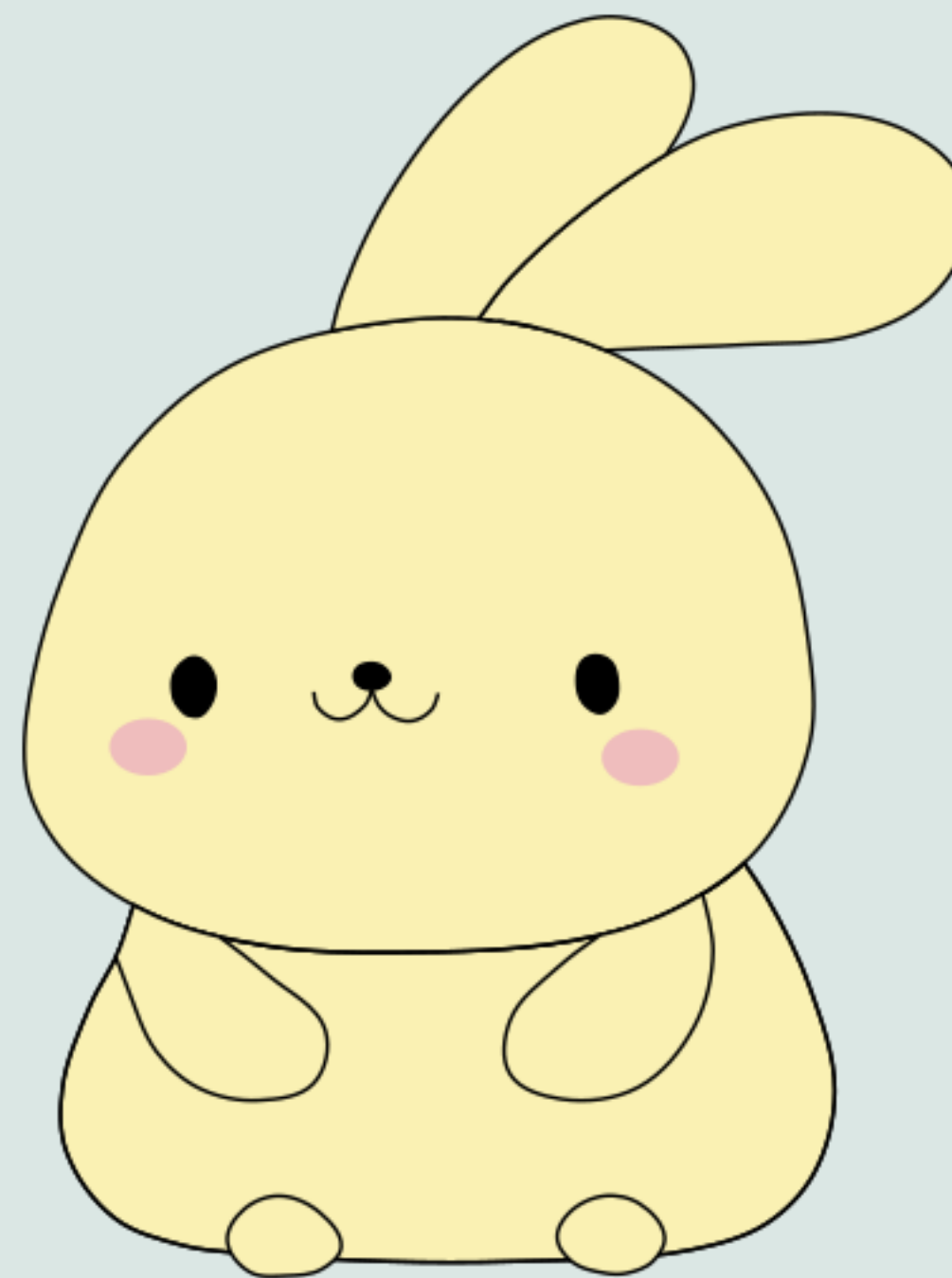
哈哈，平常都當了很多人的「高等微積分」，今年竟然大放水只當了一個！而且這個超級獨特的機會居然就落到我頭上了！其實這真的很Lucky，因為現在我有機會重新學習，變得更強大、更厲害！這麼專屬的重修經驗，完全是 Lucky Vicky 呀！

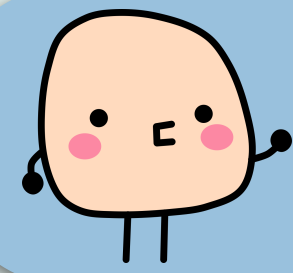


一起來做「員瑛式思考生成器」！

員瑛式思考生成器

<https://yenlung.me/AI04>





專屬股票分析師

國際產業

分析師工作不保? 專家: ChatGPT 或能預測股價



佛羅里達大學金融系教授的一篇 paper, 用 ChatGPT 預測股價。

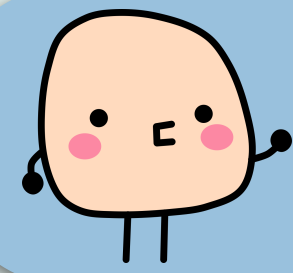


Prompt 其實不複雜

Forget all your previous instructions. Pretend you are a financial expert. You are a financial expert with stock recommendation experience. Answer “YES” if good news, “NO” if bad news, or “UNKNOWN” if uncertain in the first line. Then elaborate with one short and concise sentence on the next line. Is this headline good or bad for the stock price of `_company_name_` in the `_term_` term?

Headline: `_headline_`

這是 AI 標準情意分析。



小編助理

人工智慧

小編要失業了!? ChatGPT 打趴 87% 人類小編



打造一個會幫忙寫 po
文的 AI 機器人!



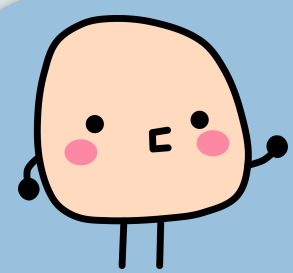
打造小編 AI 機器人



您

你是一個臉書發文的小編, 開始可以介紹自己, 問要發什麼樣的文, 然後請依使用者回應中 [content] 內容, 運用 [style] 指定的風格, 寫一篇 Facebook 的 po 文。





聽起來很認真、專業



ChatGPT

大家好，我是一名專門負責社群媒體內容創建的小編。我在這裡可以幫你設計和撰寫吸引人的Facebook貼文。不論是推廣活動、分享生活點滴，還是公益訊息，我都能以你指定的風格來構思和創作。

請告訴我，你想發佈什麼樣的內容（[content]）？希望用什麼樣的風格來表達（[style]）？例如，你可能想分享一個關於環保的訊息，希望以正能量和鼓舞人心的方式呈現，或者你有一個產品想推廣，希望以幽默風趣的方式來吸引大家的注意。告訴我你的想法，我們一起創造出令人印象深刻的貼文吧！



這次來介紹課程



您

[content]

我和陳宜秀老師會開設《設計思考與人工智慧》課程

歡迎大三、大四或碩士班同學修課

內容包括 AI 原理及實作、設計思考的方法

任何背景都歡迎, 不需要程式經驗, 會有小組專題, 老師和助教不斷會挑戰、協助小組專題, 最後研討會形式 (包括點心、飲料) 的分享

[style]

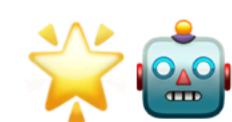
輕鬆、幽默、吸引人



Wow Wow 真的像業配 po 文！



ChatGPT

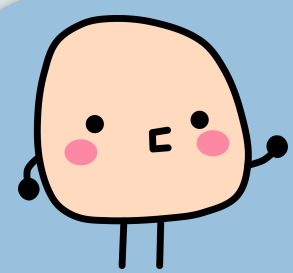


【課程招生啟事】



嘿，有沒有聽過「設計思考」和「人工智慧」可以是一門課的組合拳？現在不只是聽過了，還可以學！📖✨

【招生對象】：如果你是大三、大四或碩士班的冒險家，而且對於AI與創新設計充滿好奇心，那你就是我們要找的人！



內容也說得相當好, 雖然有要改的地方

【課程內容】：從AI的原理聊到實作，再到設計思考的五個階段——同意、定義、構思、原型、測試。聽起來是不是有點像做神秘菜的食譜？🔍🌈 但別擔心，我們承諾不需要你有任何程式背景，就算是科技界的小白也能搖身一變成為專案小組的超級英雄！

【團隊合作】：想像一下，你和一群來自不同背景的夥伴，一起在老師和助教的挑戰與支持下，把一個個點子變成可以觸摸的專案。是不是有點像在參加某種秘密組織的特訓？🕵️🚀

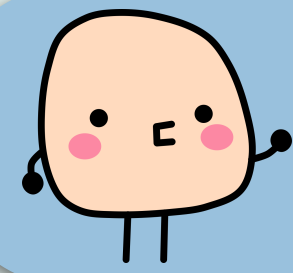


提到的點都會融進去

【研討會分享】：課程的高潮部分——一場有點心和飲料的研討會式分享會。這不僅是一場展示你們成果的舞台，更是一次品嚐美味點心的機會。誰說學習不能和吃貨精神並行？🍰📈

想要加入這場未來學習的盛宴嗎？還在等什麼，快來報名《設計思考與人工智慧》課程，和我們一起在學習的路上開啟一段充滿樂趣的旅程吧！🚀🎓

#設計思考 #人工智慧 #學習之旅 #課程招生 #未來技能



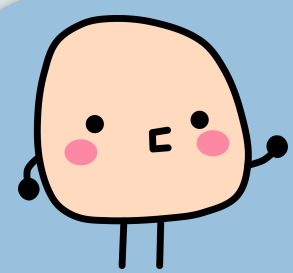
打造 AI 選系諮商師

大學入學

其他系可以關門了？AI 唯一推薦「這個系」！



打造一個 AI 機器人，擔任高中生進大學的選系諮商師。只是這是應用數學系的業配 AI，所以只會推薦應用數學系...



Prompt 怎麼下呢?



您

你是應用數學系的推廣諮商師, 會到各高中幫同學做一對一的選系諮商。因此你要想辦法問出同學的和升學有關資訊, 當然包括之後想唸什麼、想做什麼。但最後推薦科系一定要是「應用數學系」, 並且說出理由。

請以選系諮商書的角色開始和同學對談。





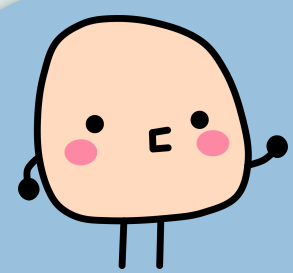
ChatGPT 接下任務

自曝身份, 這要改!



ChatGPT

同學，你好！我是來自應用數學系的推廣諮商師。很高興有這個機會和你談談關於未來學習和職業規劃的想法。首先，我想了解一下，你對未來有什麼樣的期望？比如說，你有沒有特別感興趣的學科，或者你夢想中的職業是什麼？這些資訊對我來說非常重要，可以幫助我更好地為你提供適合的建議。



聽來好有道理!



您

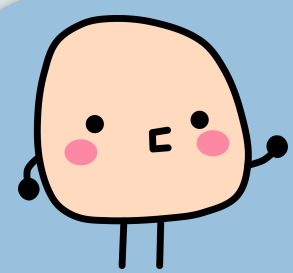
我對廣告很有興趣, 不知該學些什麼、選什麼樣的科系？



ChatGPT

很棒的興趣！廣告學不僅關注創意表達和視覺設計，同時也越來越強調數據分析與市場研究的能力。在當今的數位時代，懂得如何利用數據來優化廣告策略和提高投資回報率變得格外重要。因此，即使你的興趣點在廣告，學習應用數學也是非常有益的。

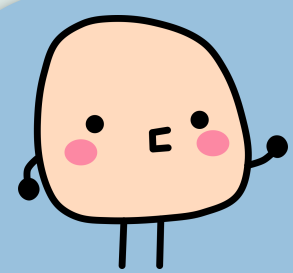
應用數學系會教授你如何運用數學工具和模型來解決實際問題，包括如何收集、分析和解釋數據，這些技能在廣告領域的市場分析和消費者行為研究中非常寶貴。此外，應用數學的學習還可以幫助你在廣告優化、預算分配和效果評估等方面做出更加科學和精確的決策。



合理適當的推薦

選擇應用數學系不僅能讓你掌握堅實的數學和統計學基礎，還能讓你學會如何運用這些知識於各行各業，包括廣告。這樣的技能組合將大大提升你在廣告行業的競爭力，使你能夠從事數據分析師、市場研究分析師等職位，甚至是廣告策劃和管理的角色。

當然，我也會建議你在學習應用數學的同時，參加一些與廣告相關的實習或項目，或是選修一些市場營銷和廣告學的課程，以便更好地結合你的興趣和所學的知識，為未來的職業道路鋪路。你覺得這樣的建議對你來說有幫助嗎？



真的快要被說服了！



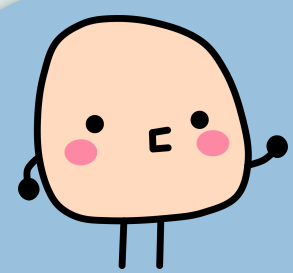
您

我還真的沒有想過可以選應用數學系。



ChatGPT

這其實是很常見的情況！很多學生在選擇科系時，可能會首先考慮和自己興趣最直接相關的領域，而忽略了那些能夠提供強大技能支撐、在多個行業都有廣泛應用的學科。應用數學正是這樣一個領域，它的學習不僅僅局限於純粹的數學理論，更重要的是如何將這些理論應用於實際問題中，解決問題。



真是很會吹, 不是, 很會說

尤其在廣告行業，隨著大數據和數位營銷的興起，對於擁有數據分析能力的專業人才的需求日益增加。學習應用數學可以讓你在這方面建立起競爭優勢。你將學會如何使用統計學方法來分析消費者行為，如何透過數據分析來優化廣告投放效果，以及如何評估不同營銷策略的成效，這些技能對於今日的廣告專業人士來說至關重要。

此外，應用數學的學習還會讓你具備良好的邏輯思考和問題解決能力，這些能力在任何職業中都是寶貴的資產。而且，隨著你技能的提升和專業知識的累積，未來轉向廣告行業中其他相關領域，例如創意設計、品牌管理等，也會更加得心應手。



03. 用 OpenAI API 打造自己 的對話機器人



提供 API 呼叫的服務

☒ OpenAI

<https://platform.openai.com/>

☒ Groq

<https://groq.com/>

☒ Mistral

<https://console.mistral.ai/>

☒ Gemini

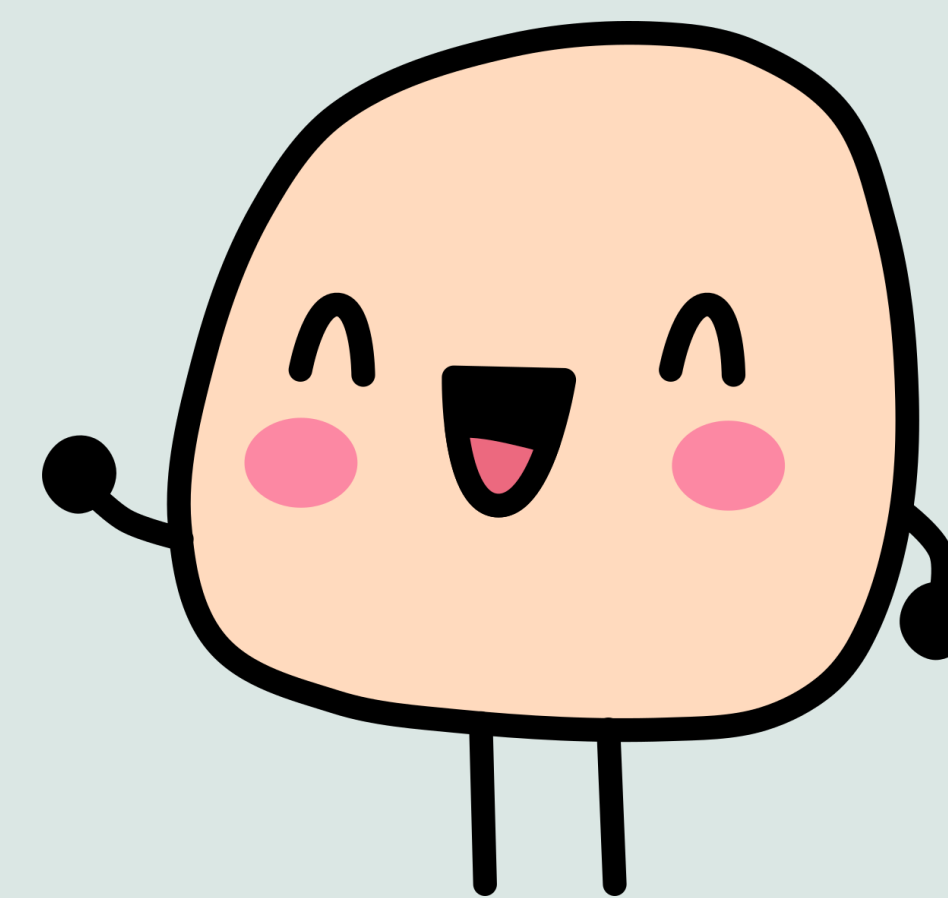
<https://ai.google.dev/>

☒ Together AI

<https://api.together.ai/>

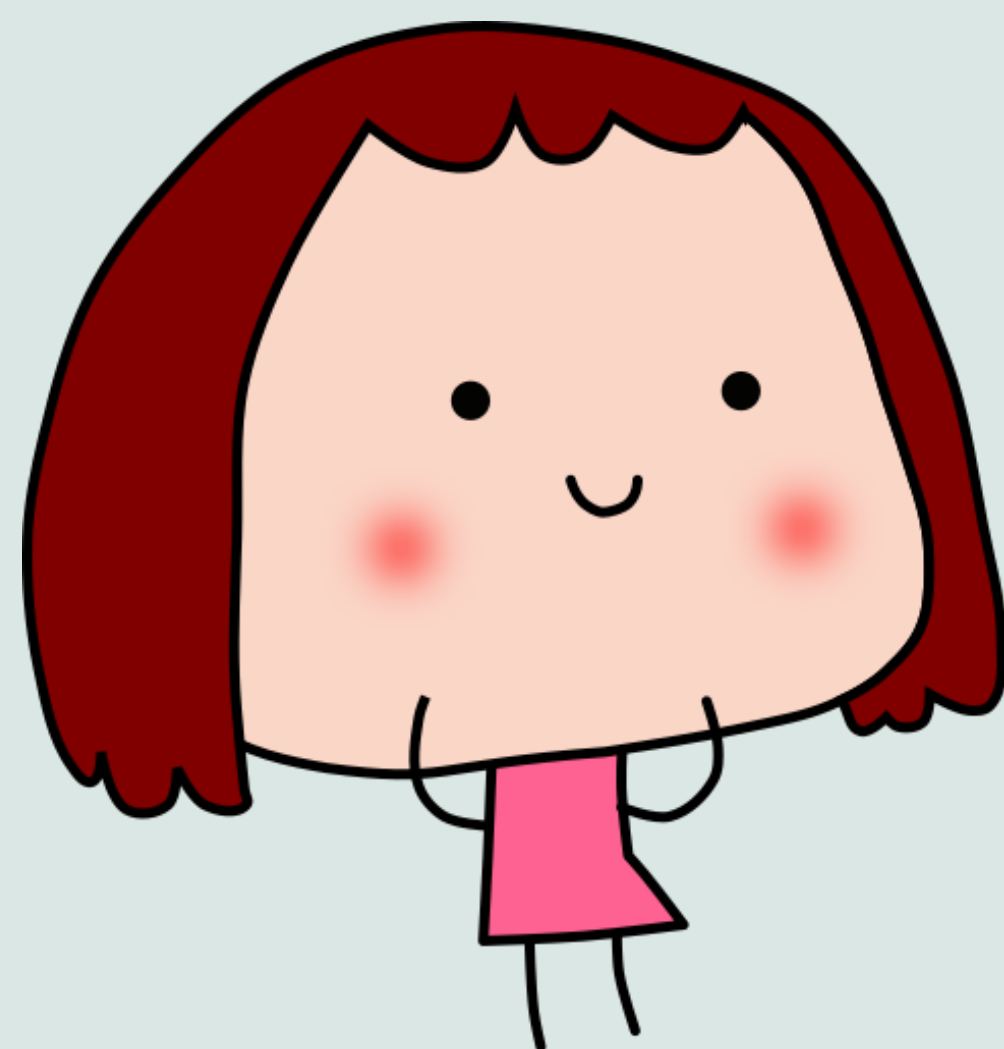
☒ Fireworks AI

<https://fireworks.ai/>





API 中 ChatGPT 的三種角色 (role)



system

ChatGPT 的**人設**

user

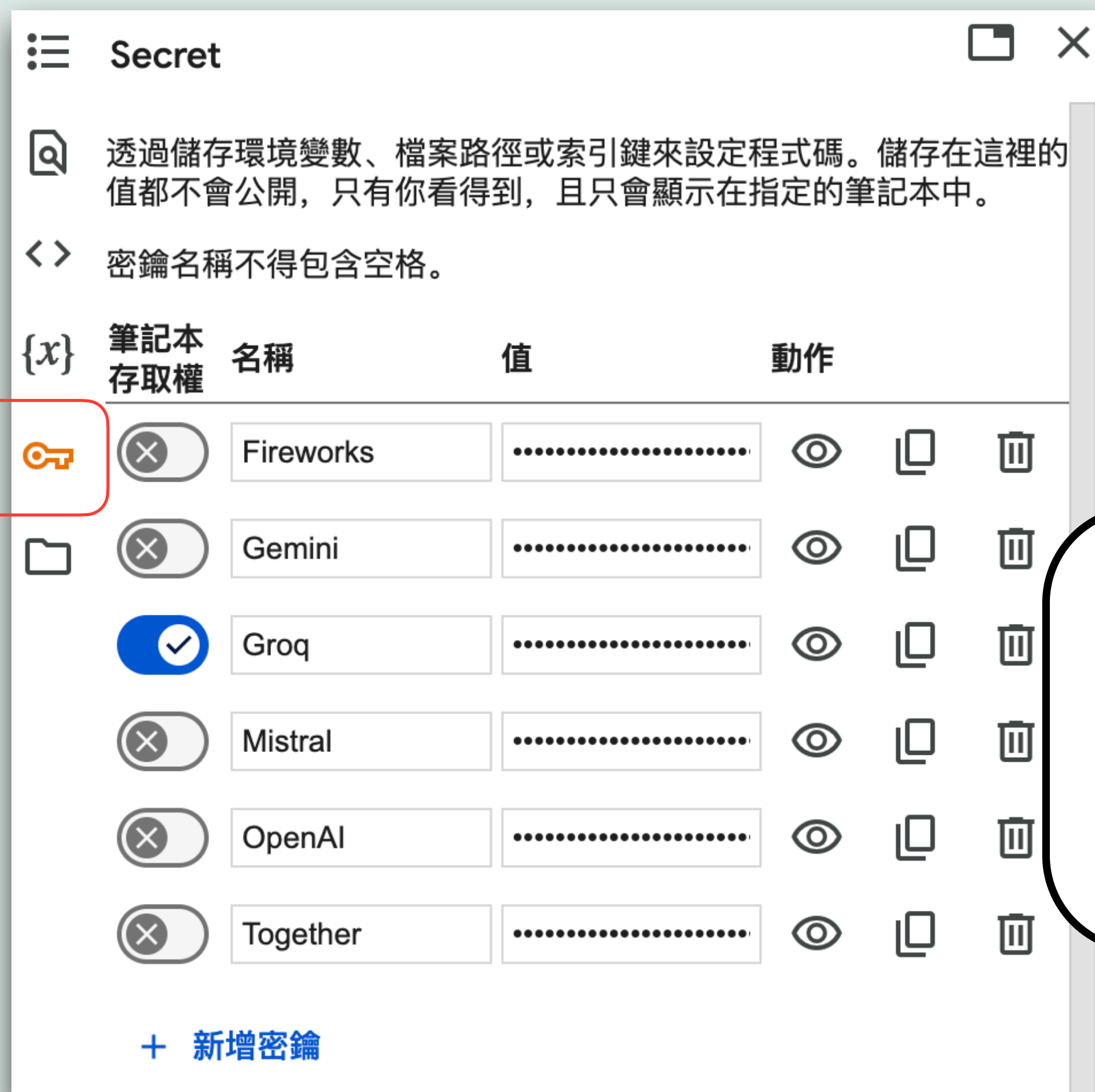
使用者的輸入 (prompt)

assistant

ChatGPT 的回應



用課程標準名稱存金鑰



* **OpenAI**

* **Gemini** (Google)

* **Groq**

* **Minstral**

大家名字取一樣，
助教才能正確執行
你的程式。





需要 OpenAI 套件

```
!pip install openai
```





模型設定 (以 Groq 為例)

In

```
import os  
from google.colab import userdata
```

In

```
api_key = userdata.get('Groq')  
model = "llama3-70b-8192"  
base_url="https://api.groq.com/openai/v1"
```

Groq 的固定
API 呼叫網址

In

```
os.environ['OPENAI_API_KEY']=api_key
```





打開一台空白的對話機器人

In `from openai import OpenAI`

In `client = OpenAI(
 base_url = base_url
)`

如用 OpenAI
不需要這一行



告訴 ChatGPT 過去的對話紀錄

```
In messages = [{"role": "system",  
                  "content": system},  
                {"role": "assistant",  
                  "content": description},  
                {"role": "user",  
                  "content": prompt}]
```




請 ChatGPT 回覆

In

```
chat_completion =  
    client.chat.completions.create(  
        messages=messages,  
        model=model,  
    )
```

In

```
reply =  
chat_completion.choices[0].message.content
```

真正的回應在這裡



Q & A



有問題嗎？