



[www.aiengineer.tw](http://www.aiengineer.tw)

# AI Study Session

**111-2 Fine Tune Large Language Models For Your Application**

**Date: 2023/06/07**

**Reporter:**

**Hung-Ming, Lin. Kuan, Yen. Meg, Ho. Guo-Chi, Li. Chien-Yu, Tseng. Chi-Yin, Ho.**



# Outlines

**1. Introduction of LLM and OpenAI**

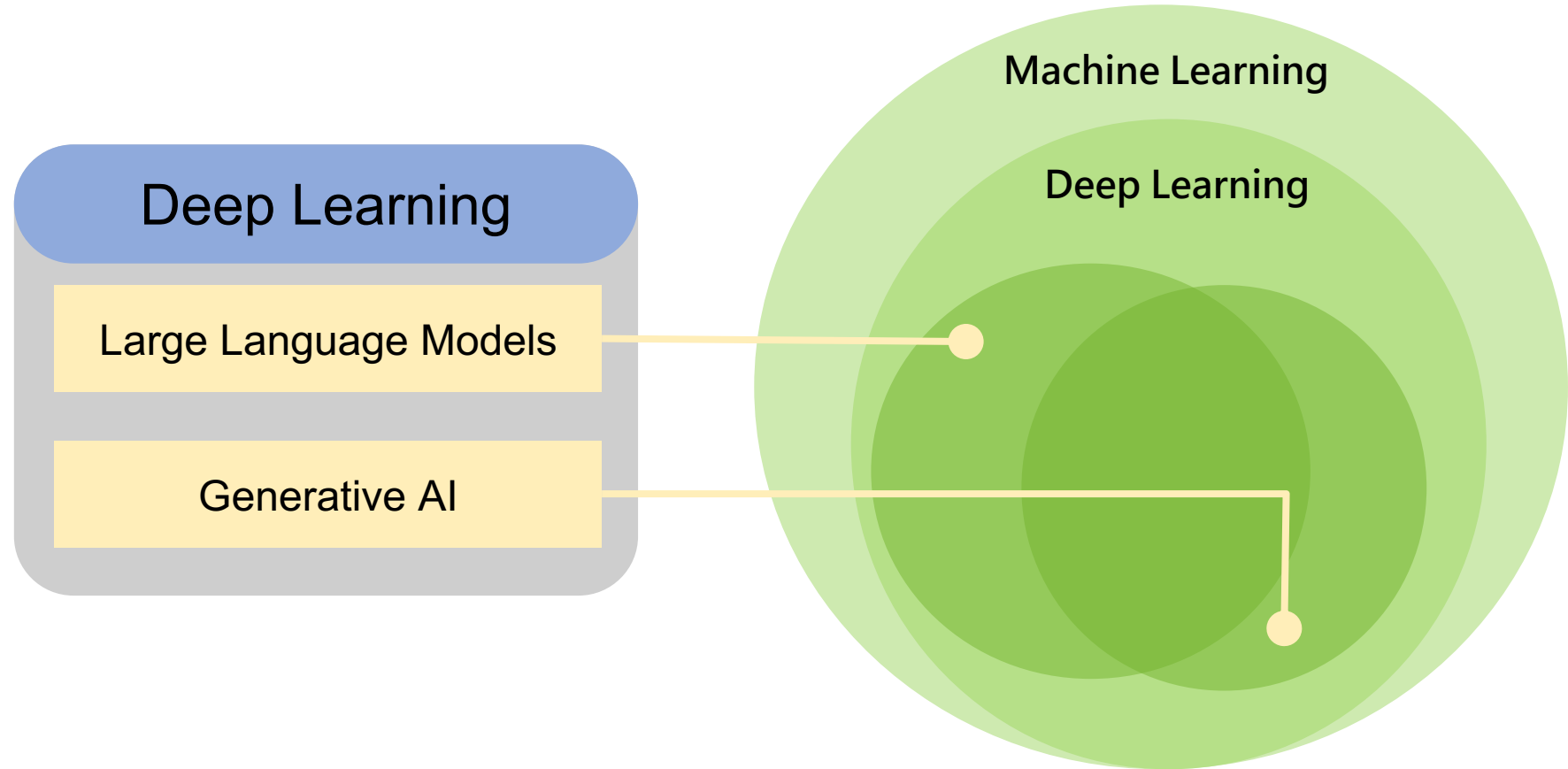
**2. Parameter-Efficient Fine-Tuning (Part I)**

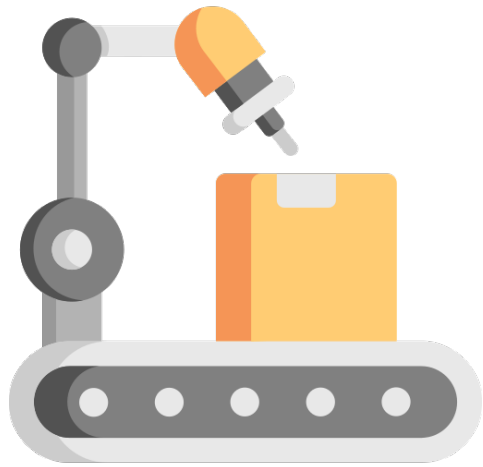
**Take a Break (10 min.)**

**3. Parameter-Efficient Fine-Tuning (Part II)**

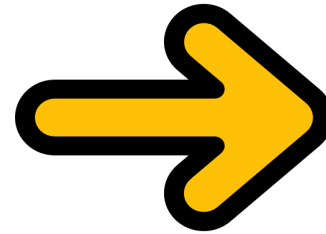
**4. The Application of Word-Embedding: LangChain**

Large Language Models (LLMs) are a subset of Deep Learning





Generative AI





# What are **Large Language Models (LLMs)** ?

Large, general-purpose language models  
can be pre-trained and then fine-tuned  
for specific purposes

Imaging you are training a dog...



# You need more...



+ Special Training

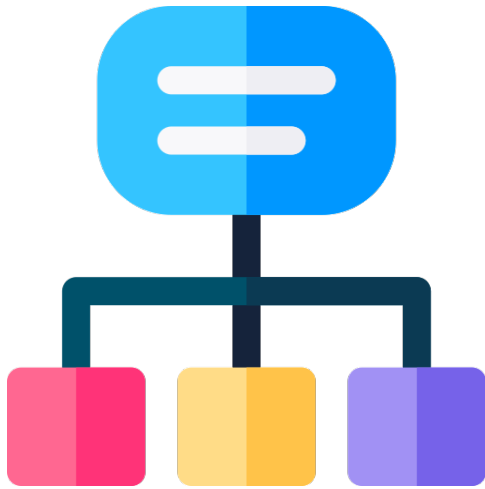


**Similar idea** applies to  
**Large Language Models**



# Introduction to LLMs

Large Language Models are trained to solve common language problems, like...



**Text  
classification**



**Question  
answering**



**Document  
summarization**



**Text  
generation**



# Introduction to LLMs

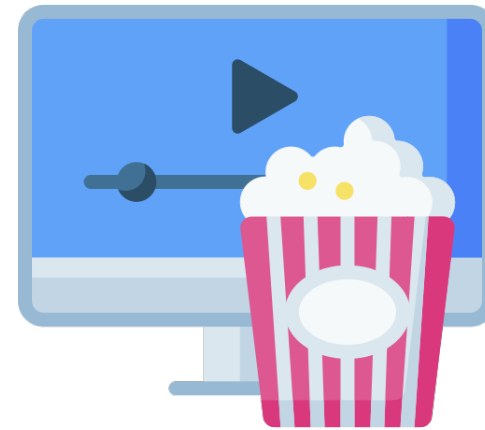
...then be tailored to solve specific problems in different fields, like...



**Retail**



**Finance**



**Entertainment**

Trained with  
a relatively small  
size of field  
datasets





# Introduction to LLMs

## Large Language Models

- ✓ **Large**
  - Large training dataset [PB]
  - Large number of parameters [B~T]
- ✓ **General purpose**
  - Commonality of human languages
  - Resource restriction
- ✓ **Pre-trained and fine-tuned**



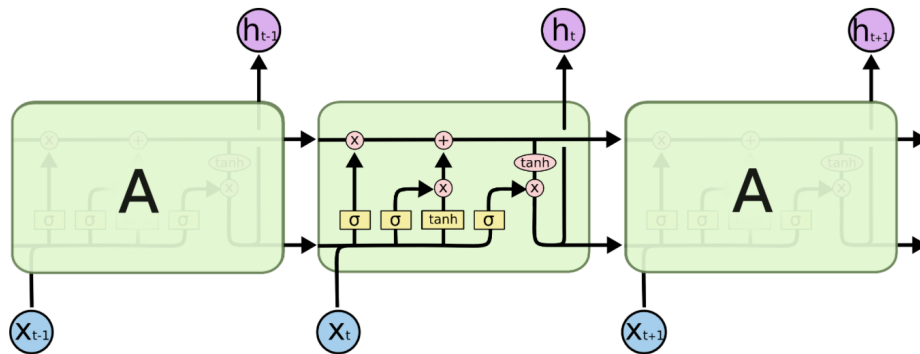




# Introduction to LLMs

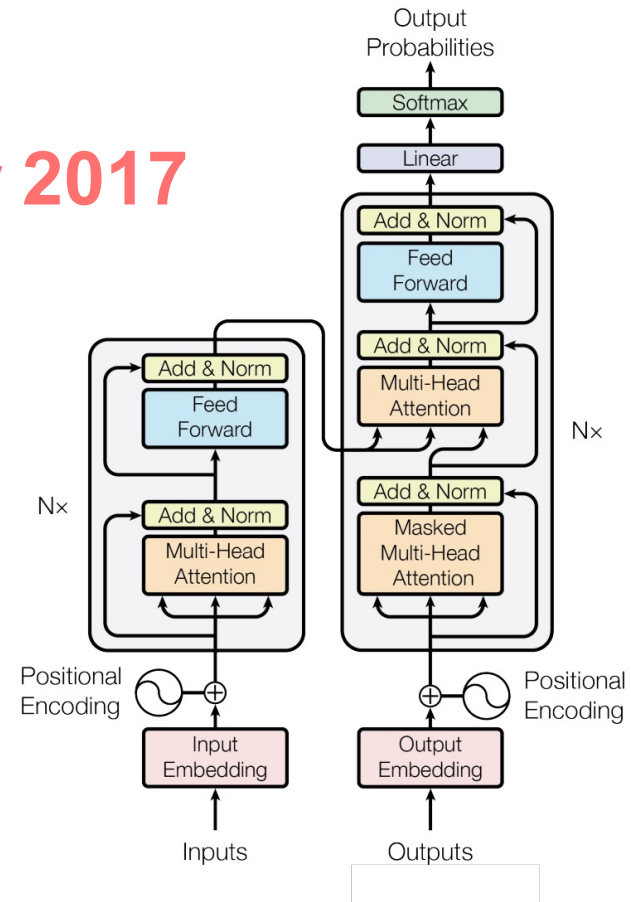
## The turning point

Before 2017



RNN-based models  
(LSTM, GRU ...)

After 2017



Transformer



# Introduction to LLMs

## NTU's Graduation Speech



NVIDIA CEO Tells NTU Grads :  
**Run, Not Walk** — But Be Prepared to Stumble



# Introduction to LLMs

## Benefits of using large language models



A single model can be used for different tasks



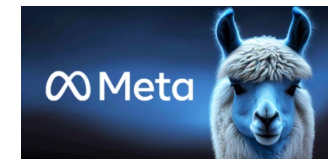
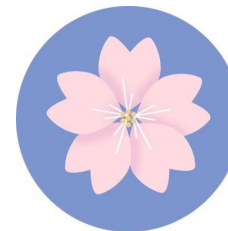
The fine-tune process requires minimal filed data



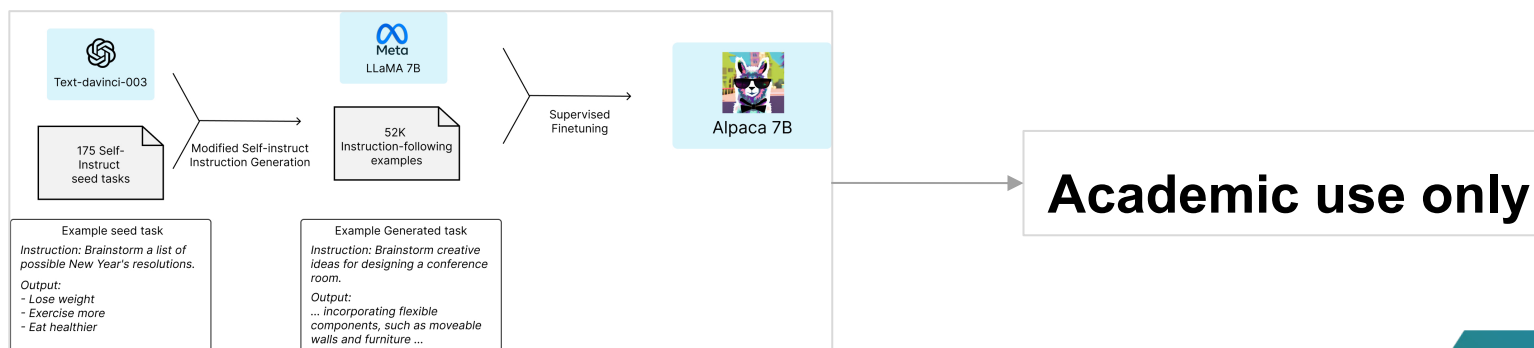
The performance is continuously growing with more data and parameters

# Introduction to LLMs

## Related LLMs



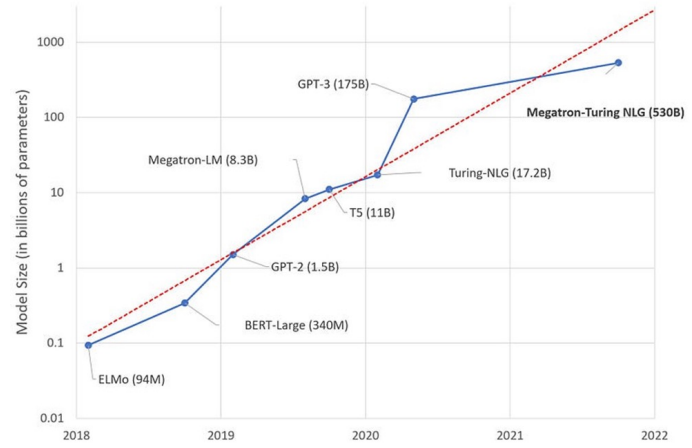
Name	BERT	GPT-1~4	PaLM, PaLM 2	BLOOM	LLaMA	Alpaca
Year	2018	2018 —	Apr. 2022 —	Jul. 2022 —	Feb. 2023	Mar. 2023
Developer	Google	OpenAI	Google	HuggingFace	Meta	Stanford
Params	340M	120M ~ 1T (?)	540B ~ 1.2T	175B	7B ~ 65B	7B
Corpus size	3.3B	4.5G (≅ 1B) ~ ?	768B ~ 3.6T	350B	1.4T	< 600 USD



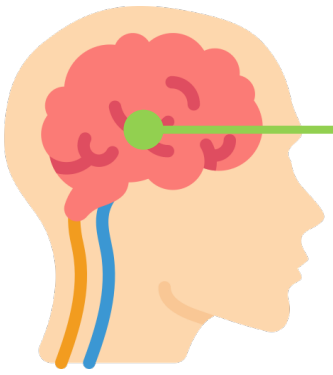


# Introduction to LLMs

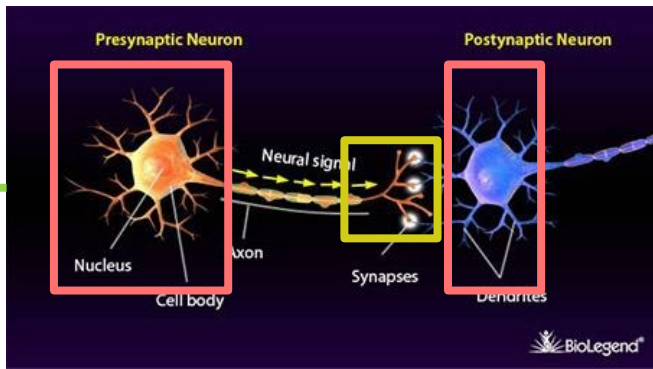
## A new Moore's Law?



~ 10x / year



86B neurons



100T synapses

**Hung-yi Lee**  
@HungyiLeeNTU 16萬位訂閱者 419 部影片  
李宏毅 >

已訂閱

首頁 影片 播放清單 社群 頻道 簡介

- FrugalGPT: 來看看窮人怎麼用 ChatGPT (下) 8:34
- FrugalGPT: 來看看窮人怎麼用 ChatGPT (上) 14:46
- 讓 AI 自主運行其他 AI 26:13
- 用語言模型解釋語言模型 (下) 13:57
- 速覽圖像生成模型 26:57
- GPT-4 來了! 有甚麼特別的地方嗎? 16:15
- 大模型 + 大資料 = 神奇力量 (3/3) 13:28
- 大模型 + 大資料 = 神奇力量 (2/3) 27:17
- 大模型 + 大資料 = 神奇力量 (1/3) 22:43
- 對於大型語言模型的兩種不同期待 (3/3) Finetune vs. Prompt 15:34
- 對於大型語言模型的兩種不同期待 (2/3) Finetune vs. Prompt 24:39
- 對於大型語言模型的兩種不同期待 (1/3) Finetune vs. Prompt 22:15



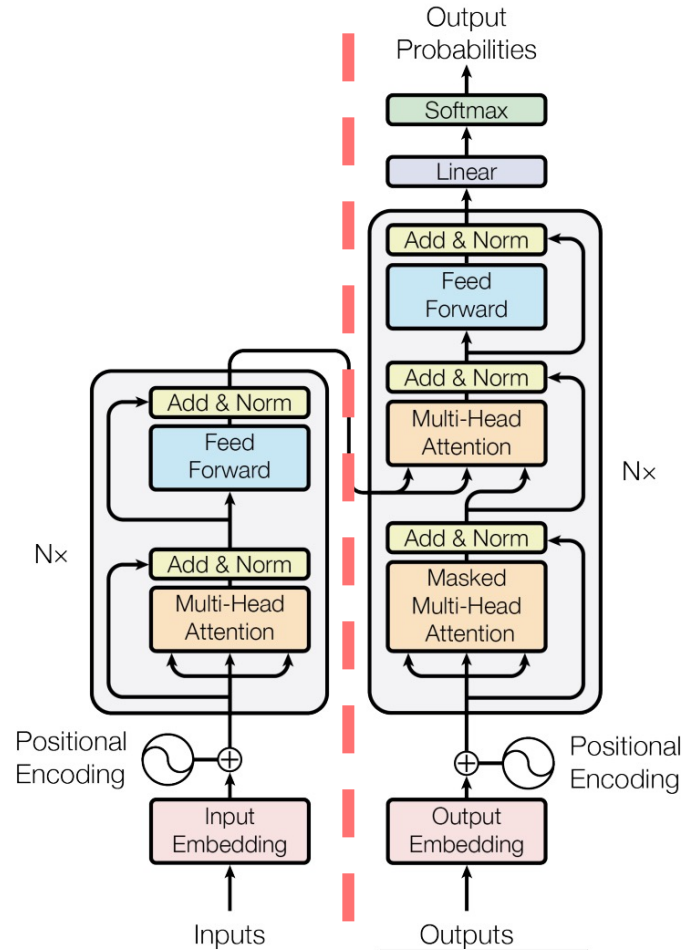
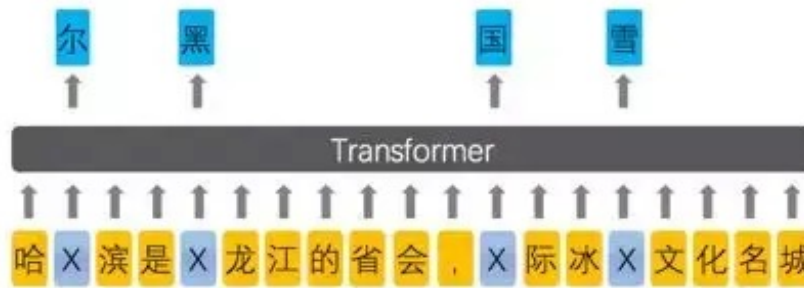
# Introduction to LLMs

## Self-supervised learning or Semi-supervised learning

### BERT



### 文字填空



### GPT



### 文字接龍

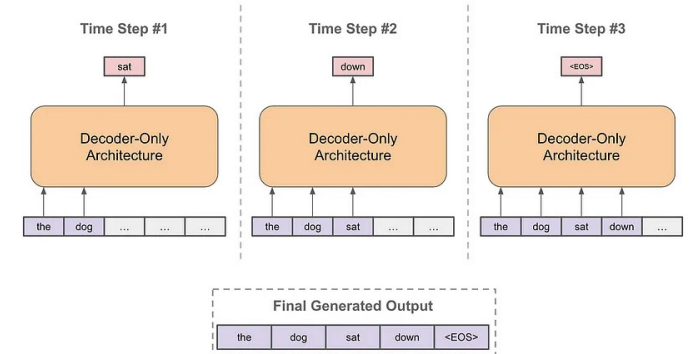


Image source:

<https://zhuannan.zhihu.com/p/59436589>

<https://heidloff.net/article/foundation-models-transformers-bert-and-gpt/>

<https://towardsdatascience.com/language-models-gpt-and-gpt-2-8bdb9867c50a>



# Introduction to LLMs

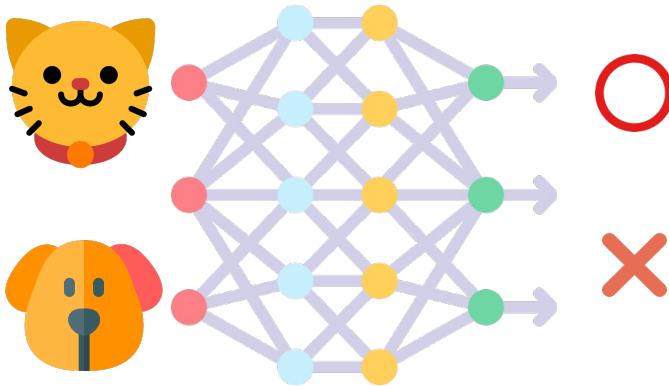
## 3 steps of programming

Traditional



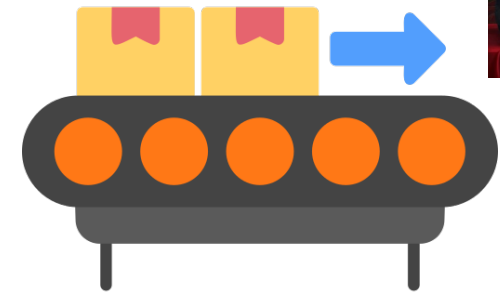
```
class Cat:  
    type = 'animal'  
    legs = 4  
    ears = 2  
    likes = ['fish', 'rats']
```

Neural Network



Generative Language Model

PROMPT





# Introduction to LLMs

## LLM Development vs. Traditional Development

### Traditional ML Development

- ML expertise needed
- training examples
- need to train a model
- compute time + hardware
- **Thinks about minimizing a loss function**

### LLM Development (by pretrained APIs)

- ML expertise needed
- training examples
- need to train a model
- **Thinks about **prompt design****



**Prompt** means all of the **text** that we feed into an LLM, which figures out what text to feed LLM to take on the **behavior** you want.



# Introduction to LLMs

## Prompt (提示語)

### Prompt Design

- Prompts involve **instructions and context** passed to a language model to achieve a desired task.



### Prompt Engineering

- Prompt engineering is the practice of **developing and optimizing prompts** to efficiently use language models for a variety of applications. (by **keywords** or **examples**)

# Introduction to LLMs

## Prompt Engineering (提示語工程)



DeepLearning.AI

Courses ▾ The Batch ▾ Blog ▾ Events ▾ Resources Company ▾ [Get AI News](#)

SHORT COURSE

### ChatGPT Prompt Engineering for Developers

[Learn for Free](#)

IN PARTNERSHIP WITH

ChatGPT Prompt Engineering for Developers

- Introduction
- Guidelines
- Iterative
- Summarizing
- Inferring
- Transforming
- Expanding
- Chatbot
- Conclusion

Course Feedback

Community

Beginner to Advanced

1 Hour

Isa Fulford, Andrew Ng

Free for a

Jupyter | 2-guidelines

Kernel starting, please wait... Not Trusted Python 3 (ipykernel)

File Edit View Insert Cell Kernel Help

### Guidelines for Prompting

In this lesson, you'll practice two prompting principles and their related tactics in order to write effective prompts for large language models.

#### Setup

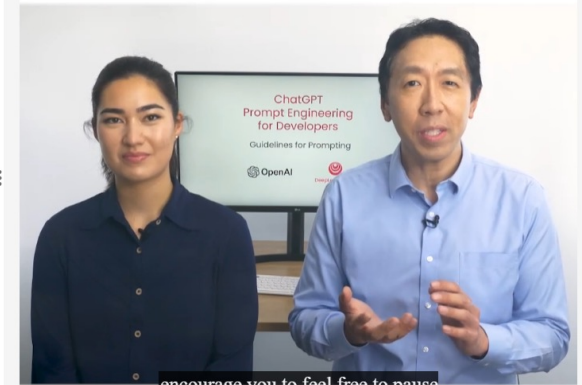
Load the API key and relevant Python libraries.

In this course, we've provided some code that loads the OpenAI API key for you.

```
In [ ]: import openai
import os

from dotenv import load_dotenv, find_dotenv
_ = load_dotenv(find_dotenv())

openai.api_key = os.getenv('OPENAI_API_KEY')
```



encourage you to feel free to pause the video every now and

[TRANSCRIPT](#)

[NEXT LESSON](#)

Isa Fulford

Andrew Ng

- ✓ Learn prompt engineering best practices for application development
- ✓ Discover new ways to use LLMs, including how to build your own custom chatbot
- ✓ Gain hands-on practice writing and iterating on prompts yourself using the OpenAI API



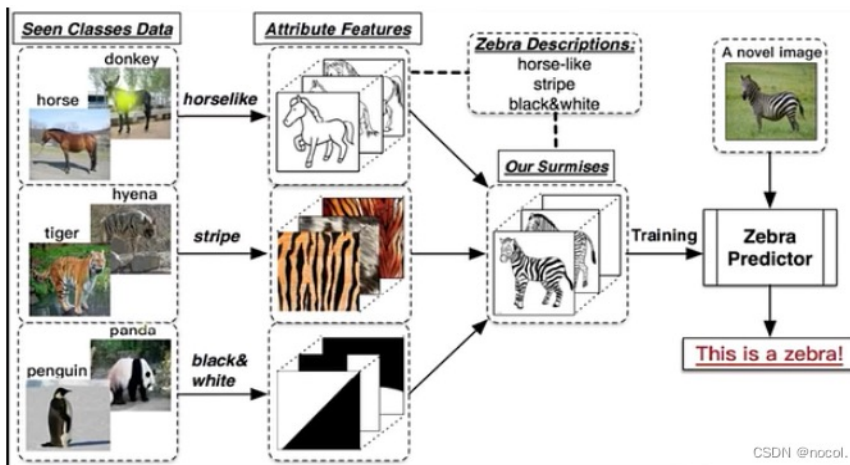
# Introduction to LLMs

## Zero-shot, One-shot and Few-shot Learning

第一部分：單題 (第1-23題，共23題)

1. Look at the picture. The man is holding a \_\_\_\_\_ of grapes in his hands.

- (A) bag
- (B) basket
- (C) bowl
- (D) box

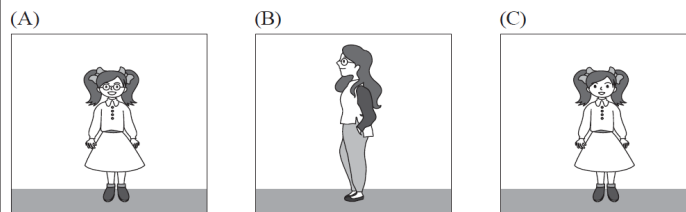


### Zero-shot

第一部分：辨識句意 (第1-3題)

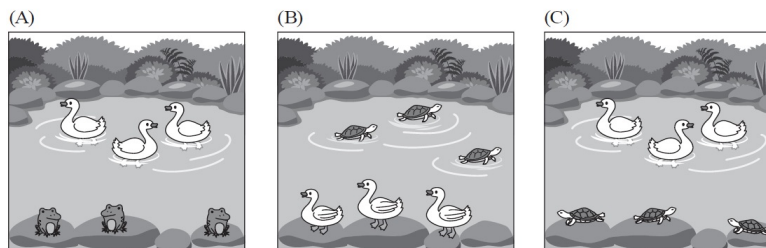
作答說明：每題均有三張圖片，請依據所聽到的句子，選出符合描述的圖片，每題播放兩次。

示例題：你會看到



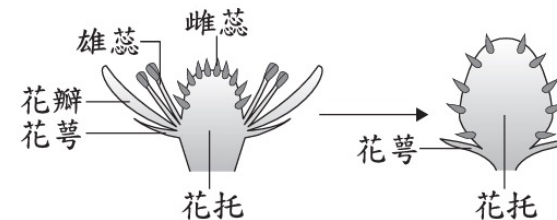
然後你會聽到……(播音)。依據所播放的內容，正確答案應該選A，請將答案卡該題「A」的地方塗黑、塗滿，即：●ⒷⒸ

第1題

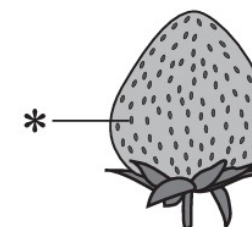


### One-shot

35. 圖(二十)為草莓花朵構造及其發育的示意圖，已知草莓是由花托處膨大而來，若圖(二十一)中的\*構造是由草莓的子房發育而成，則此\*構造應稱為下列何者？



圖(二十)



圖(二十一)

- (A) 胚珠
- (B) 種子
- (C) 果實
- (D) 花粉

### Few-shot



# Introduction to LLMs

For more details, please refer to Google Cloud ...

The screenshot shows the Google Cloud Skills Boost interface. At the top, there are navigation tabs: Paths, Explore (selected), Profile, and Subscriptions. Below this is a header for 'Google Cloud Skills Boost'. The main content area is titled 'Introduction to Large Language Models' with a dropdown arrow. A descriptive paragraph states: 'This module explores what large language models (LLM) are, the use cases where they can be utilized, and how you can use prompt tuning to enhance LLM performance. It also covers Google tools to help you develop your own Gen AI apps.' Below the text are three sections: 'VIDEO' with a link to 'Introduction to Large Language Models' (15 minutes), 'DOCUMENT' with a link to 'Introduction to Large Language Models: Reading', and 'QUIZ (IN PROGRESS)' with a link to 'Introduction to Large Language Models: Quiz'. A 'required' tag is visible next to the quiz link.



The screenshot shows the Google Cloud Tech YouTube channel page. The channel name is 'Google Cloud Tech' with a subscriber count of 99.3 million and 4377 videos. The page features a grid of video thumbnails with titles such as 'Tuning large language models', 'Exposing Apigee instances to the internet', 'An overview of fleets in Anthos', 'Start at zero: Optimising for sustainability', 'How to tune LLMs in Generative AI Studio', 'Organization Policies', 'Generative AI Studio', 'New Cloud Run Features', 'Introducing Query Plan Samples for Cloud Spanner', 'What is Google Cloud's Organization Policy Service?', 'Prototyping language apps with Generative AI Studio', and 'Moving serverless forward with Cloud Run'. Each video includes its duration and view count.

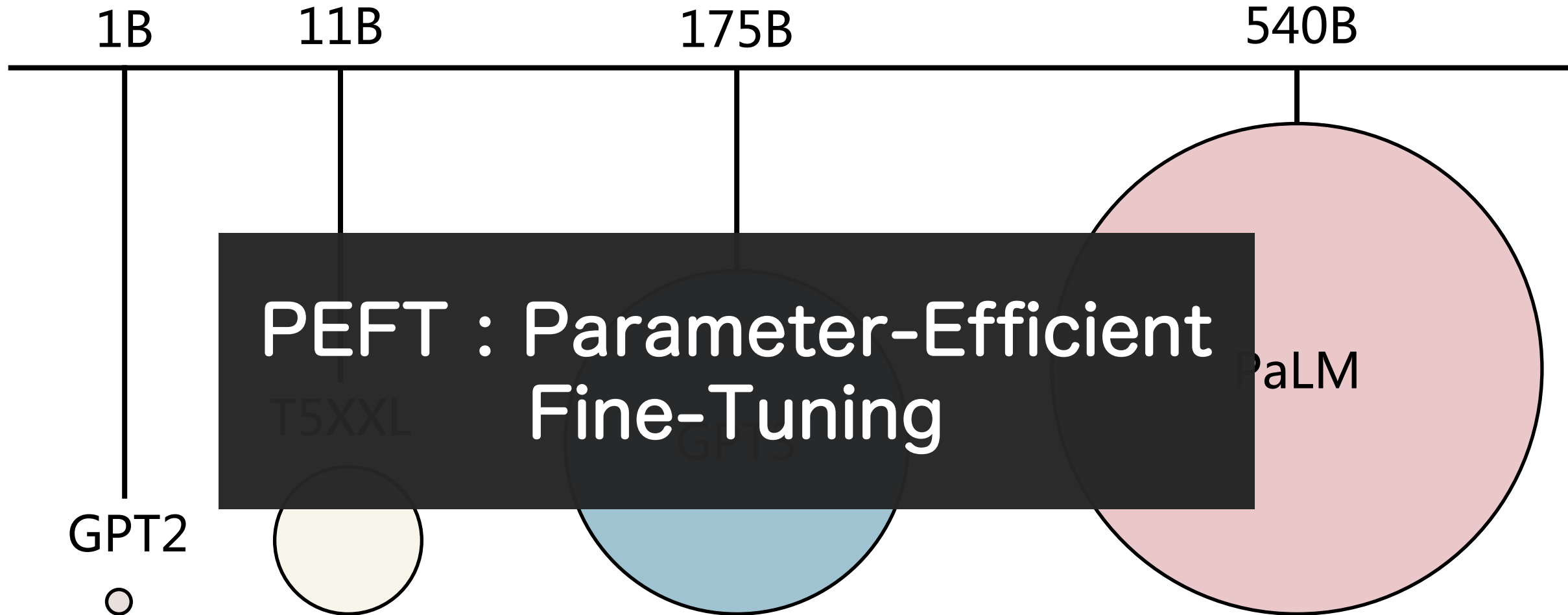


The screenshot shows the Google for Developers YouTube channel page. The channel name is 'Google for Developers' with a subscriber count of 229 million and 5814 videos. The page features a grid of video thumbnails with titles such as 'Recap Developer keynote in 5 minutes', 'I/O Dev Keynote in 8 minutes', 'What do you love about Google Cloud Platform?', 'Material design', 'Google I/O 2023 Developer Keynote in 5 minutes', 'Google I/O 2022 Developer Keynote in 8 minutes', 'Introducing Flutter', 'Hello World - Machine Learning Recipes #1', 'Project Glass: Live Demo At Google I/O', 'Introducing Flutter', 'Google Python Class Day 1 Part 1', 'Hello World - Machine Learning Recipes #1', 'Project Glass: Live Demo At Google I/O', 'Introducing Flutter', 'Introducing the new family discovery experience on Google Play', 'YouTube API Overview', 'Google for Games Developer Summit March 23', and 'A.I. Experiments: Visualizing High-Dimensional Space'. Each video includes its duration and view count.





# The problem with finetuneing LLM





# Parameter-Efficient Fine-Tuning

- **Adapter** : 透過添加額外的模型架構，並固定 Freeze LLM 的模型參數，進行訓練。

- **Prefixing** : Prefixing 在 Prompt 的前半部添加 Token 來讓模型對於特定任務可以做得更好。



# Parameter-Efficient Fine-Tuning

- **LoRA** : LOW-RANK ADAPTATION OF LARGE LANGUAGE MODELS
- **Prefix Tuning** : Prefix-Tuning: Optimizing Continuous Prompts for Generation, P-Tuning v2: Prompt Tuning Can Be Comparable to Fine-tuning Universally Across Scales and Tasks
- **P-Tuning** : GPT Understands, Too
- **Prompt Tuning** : The Power of Scale for Parameter-Efficient Prompt Tuning

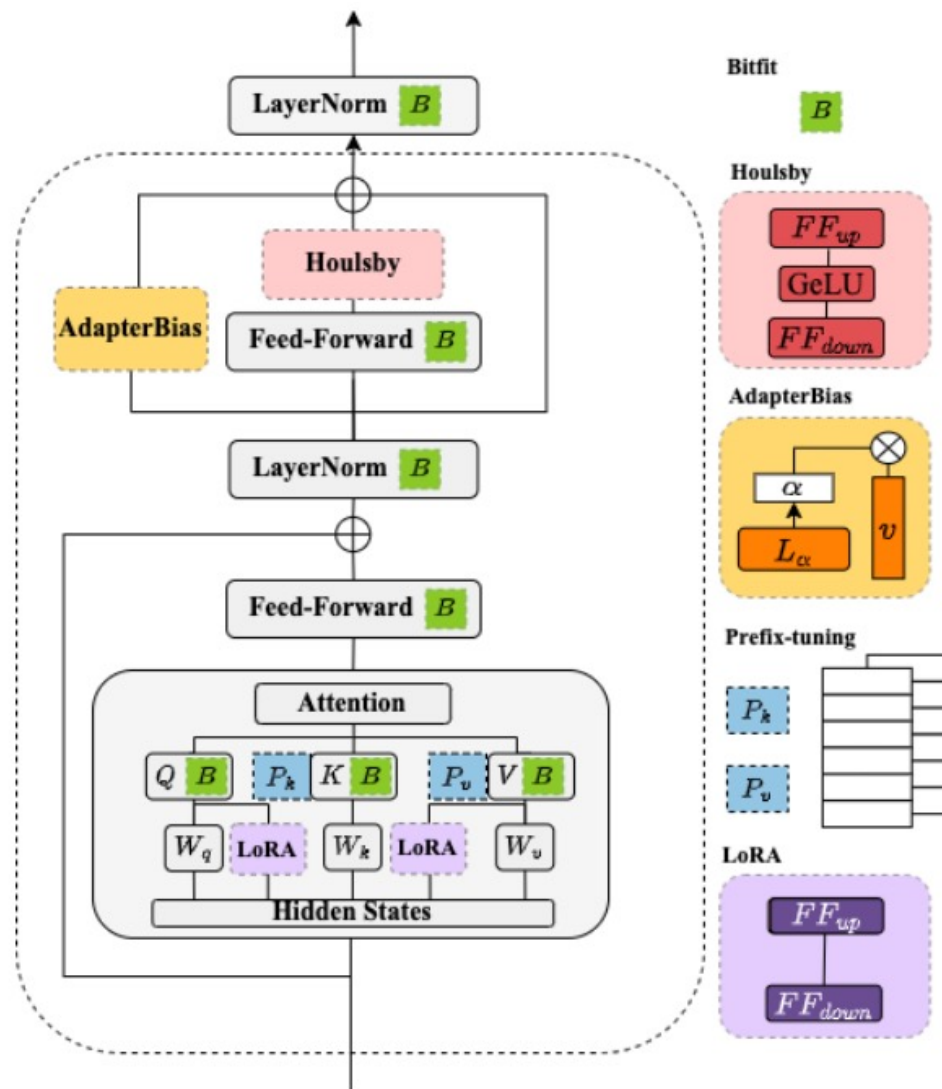




# Adapter



Adapter-hub

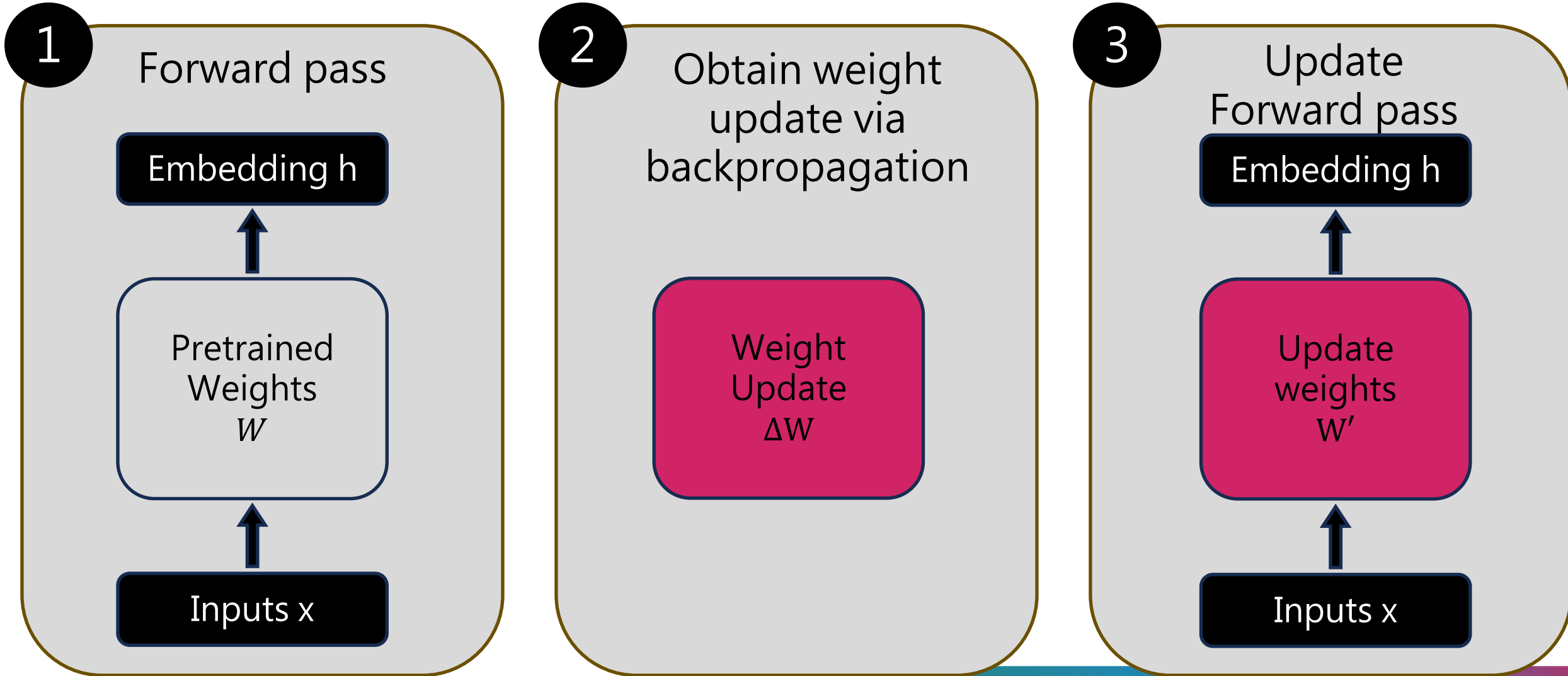


【生成式AI】Finetuning vs. Prompting :

對於大型語言模型的不同期待所衍生的兩類使用方式 (1/3)

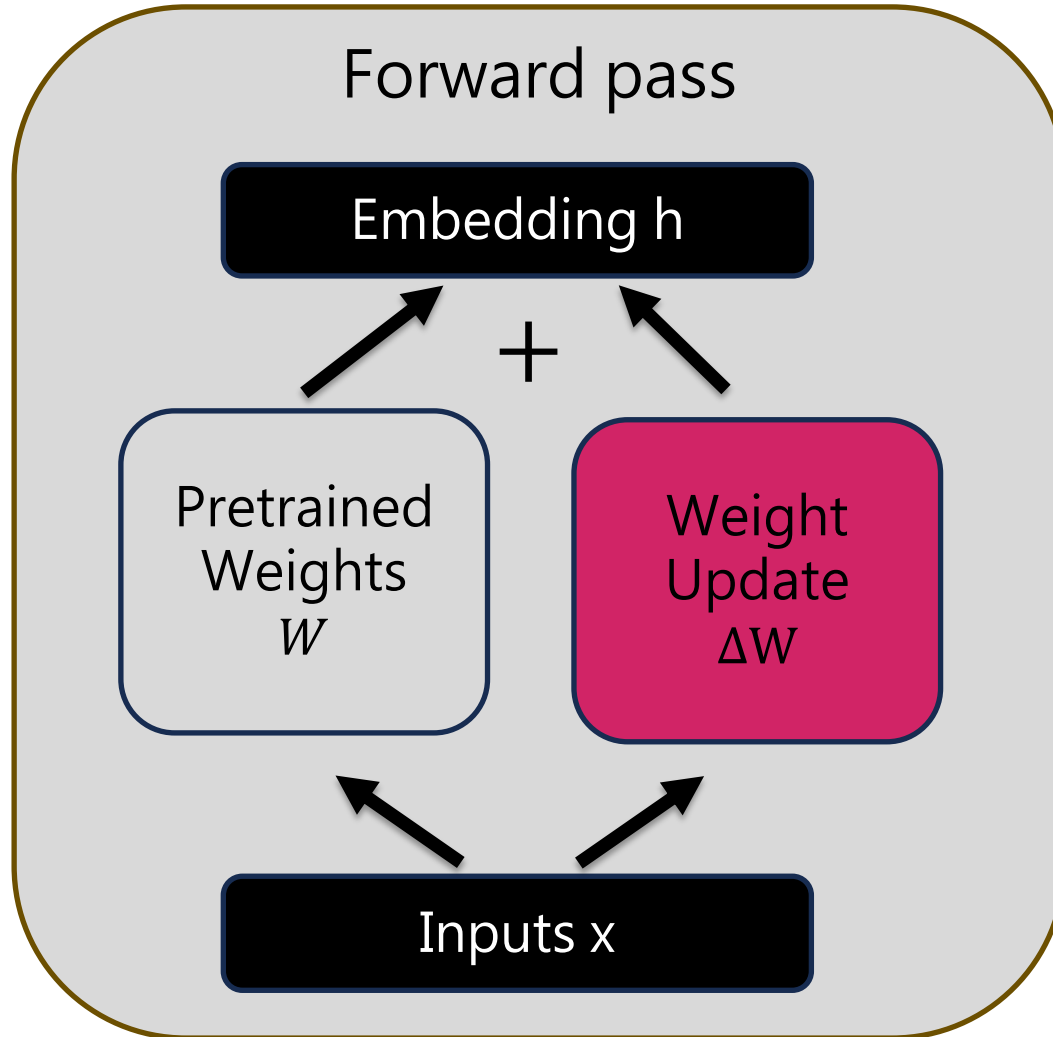


# LoRA : Low Rank Adaptation





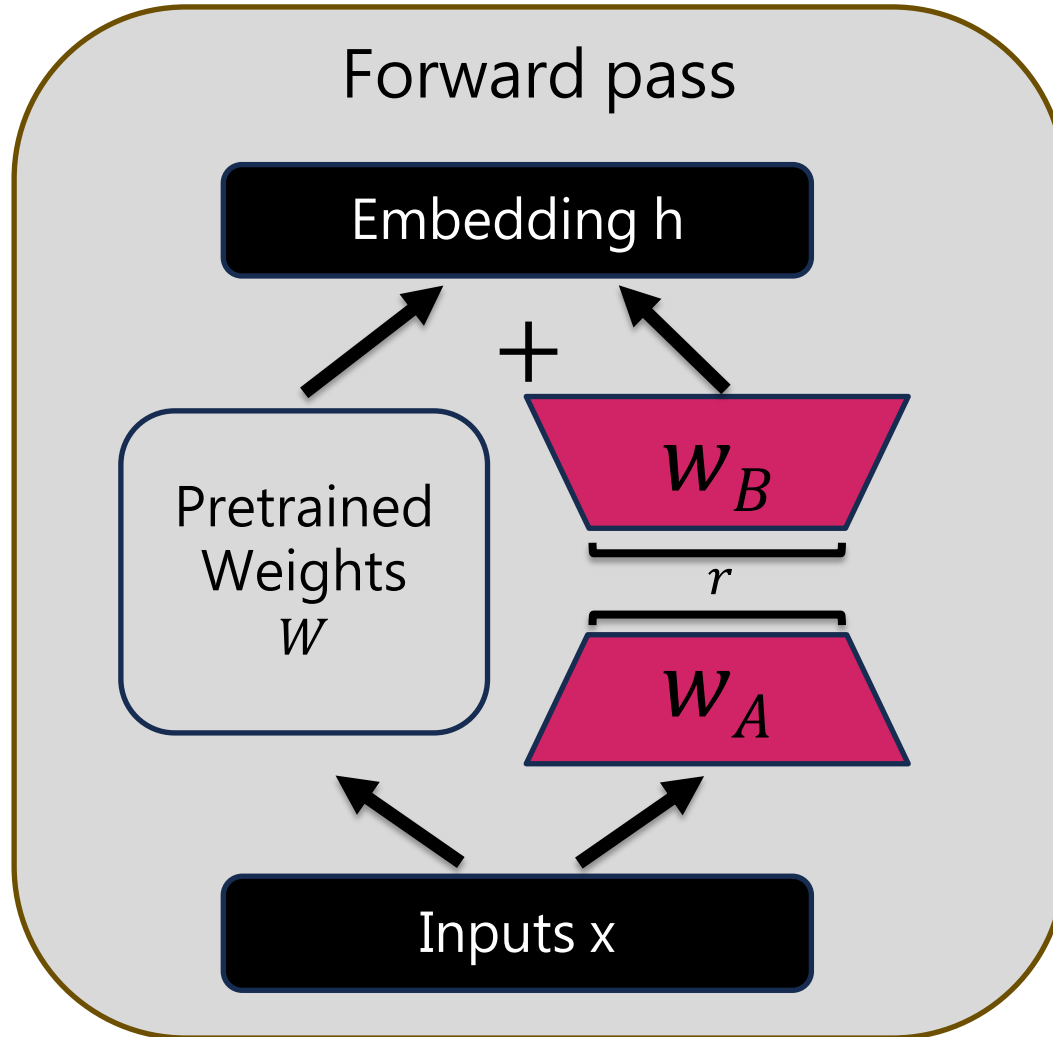
# LoRA : Low Rank Adaptation



$$h = wx + \Delta wx$$



# LoRA : Low Rank Adaptation



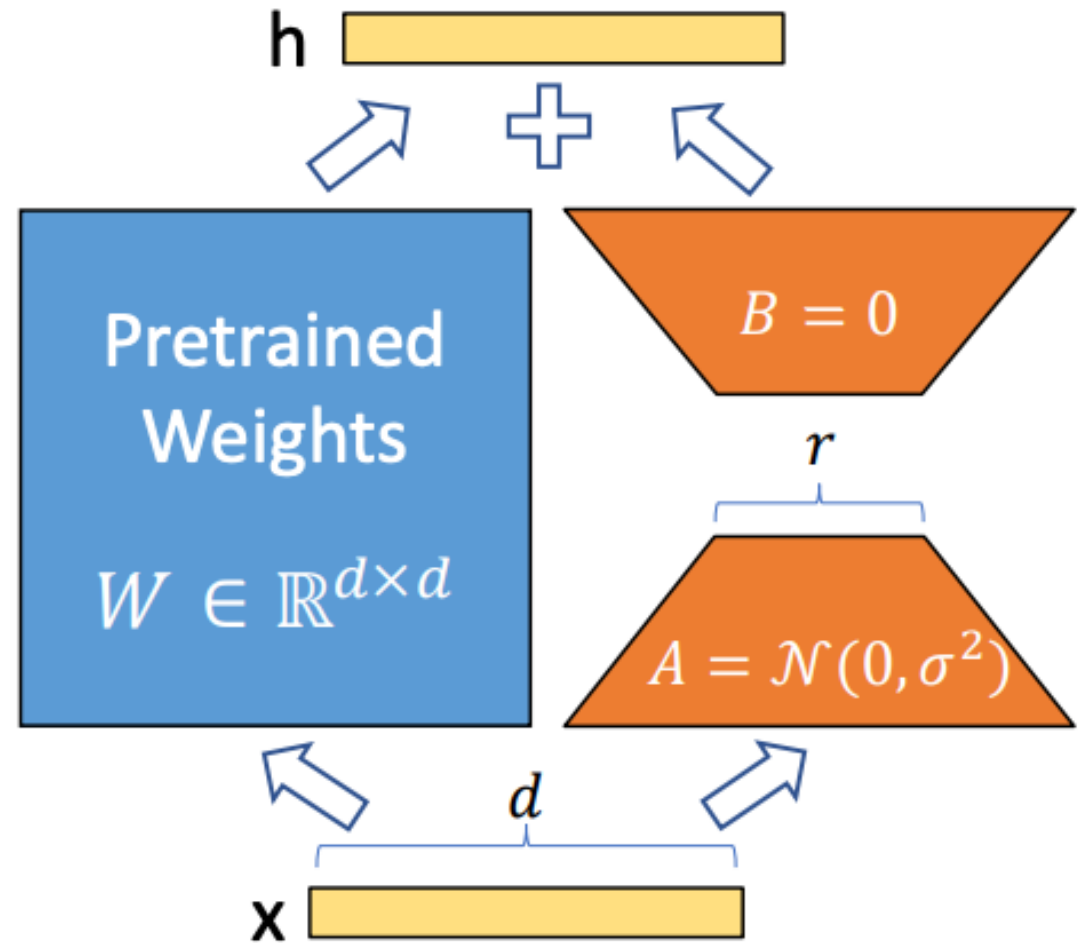
$$h = wx + W_A W_B$$





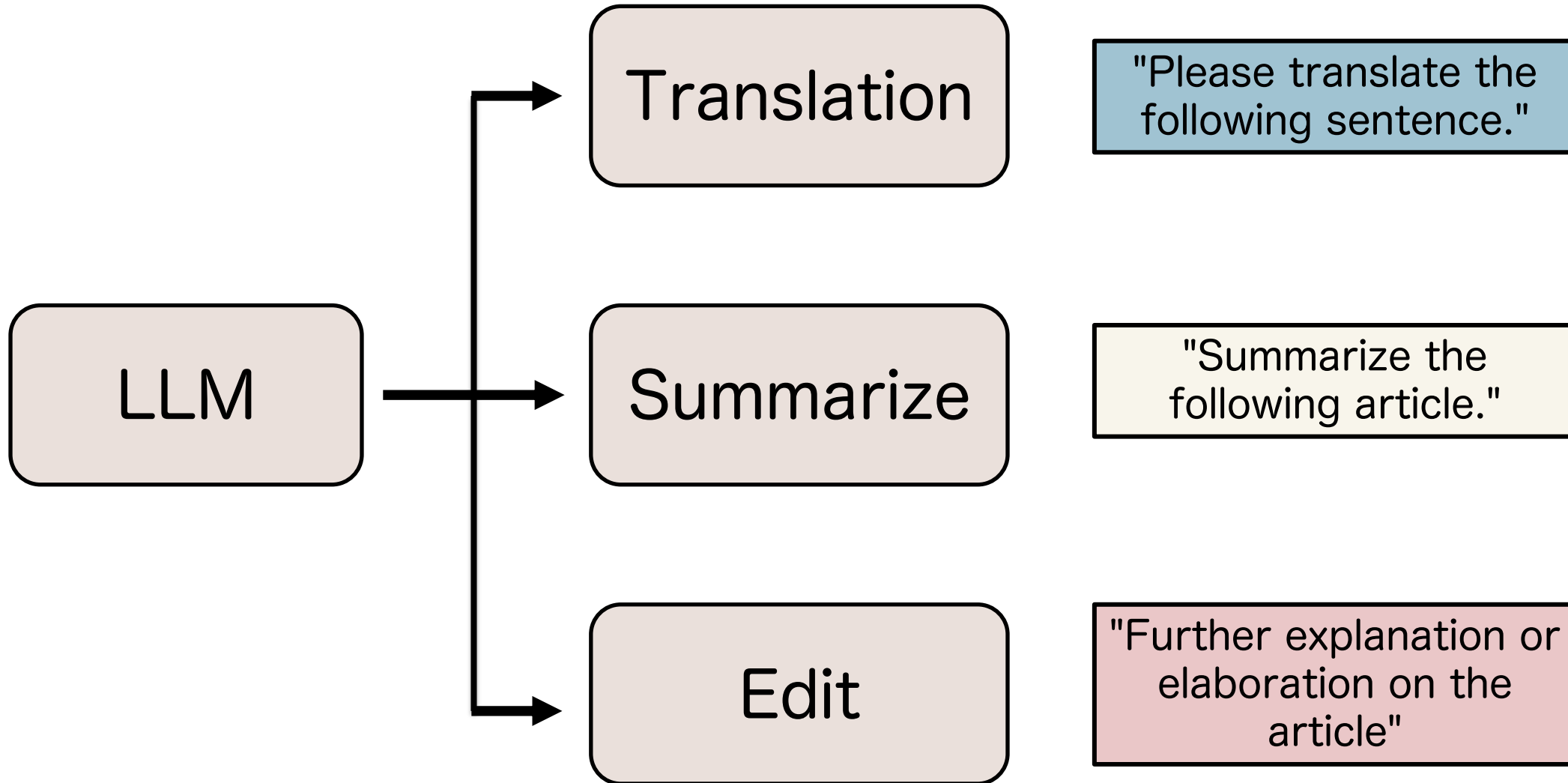
# LoRA : Low Rank Adaptation

Low-Rank Adaptation (LoRA), 概念是透過凍結原本的預訓練模型(e.g., GPT-3) 的權重, 搭配一個小的模型進行微調就可以達到很好的 Fine-Tuning 效果





# Prompt





# Discrete/hard prompts

## Discrete/hard prompts

- natural language instructions/task descriptions

Good morning

## Problems

- requiring domain expertise/understanding of the model's inner workings
- performance still lags far behind SotA model tuning results
- sub-optimal and sensitive

請翻譯：早安





# Sub-optimal and sensitive hard prompts

## Difficulty of manually designing prompts

1. Prompts that humans consider reasonable is not necessarily effective for language models([Liu et al, 2021](#))
2. Pre-trained LMs are sensitive to the choice of prompts ([Zhao et al., 2021](#))

Prompt	P@1
[X] is located in [Y]. ( <i>original</i> )	31.29
[X] is located in which country or state? [Y].	19.78
[X] is located in which country? [Y].	31.40
[X] is located in which country? In [Y].	51.08



# Shifting from hard to soft prompts

## Progress in prompt-based learning

- manual prompt design ([Brown et al., 2020](#); [Schick and Schutze, 2021a,b](#))
- mining and paraphrasing based methods to automatically augment the prompt sets ([Jiang et al., 2020](#))
- gradient-based search for improved discrete/hard prompts ([Shin et al., 2020](#))
- automatic prompt generation using a separate generative language model (i.e., T5) ([Gao et al., 2020](#))
- learning continuous/soft prompts ([Liu et al., 2021](#); [Liu and Liang., 2021](#); [Qin and Eisner., 2021](#); [Lester et al., 2021](#))

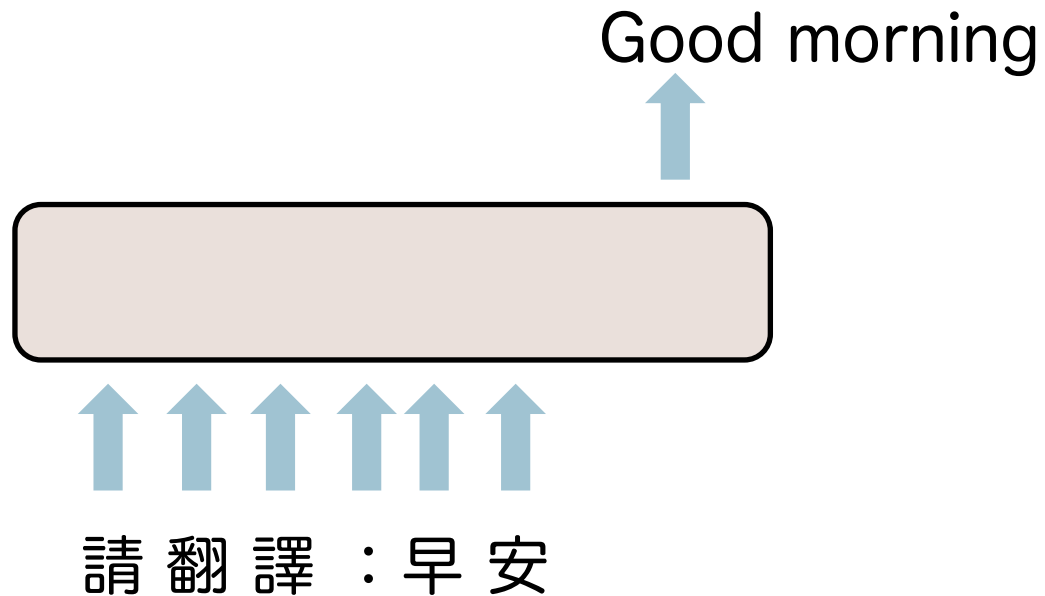
## Continuous/soft prompts

- additional learnable parameters injected into the model

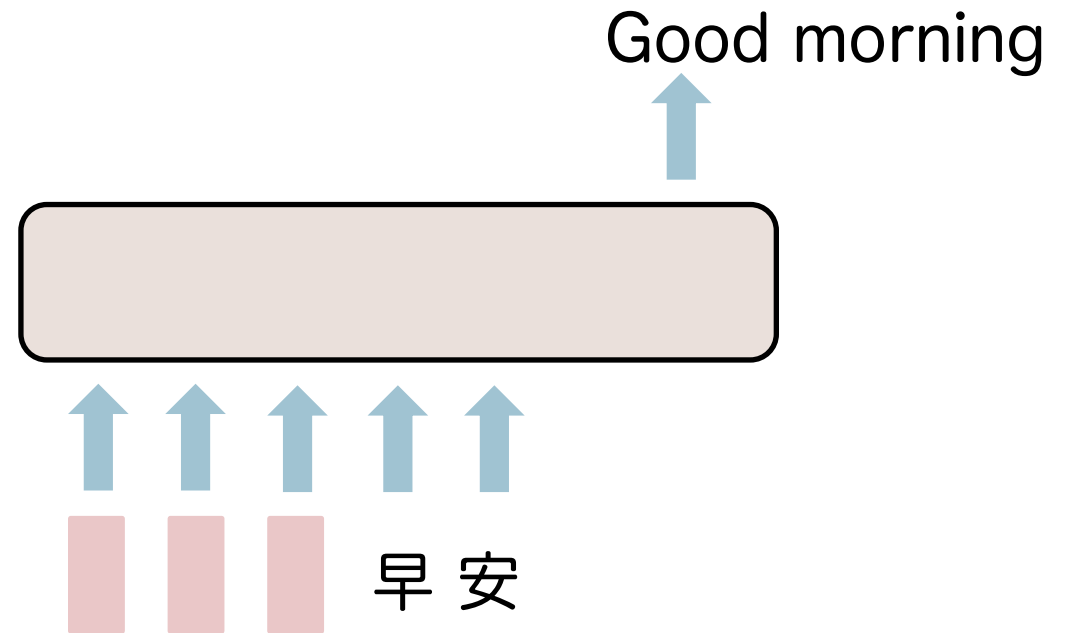


# Soft prompt

## Discrete/hard prompts

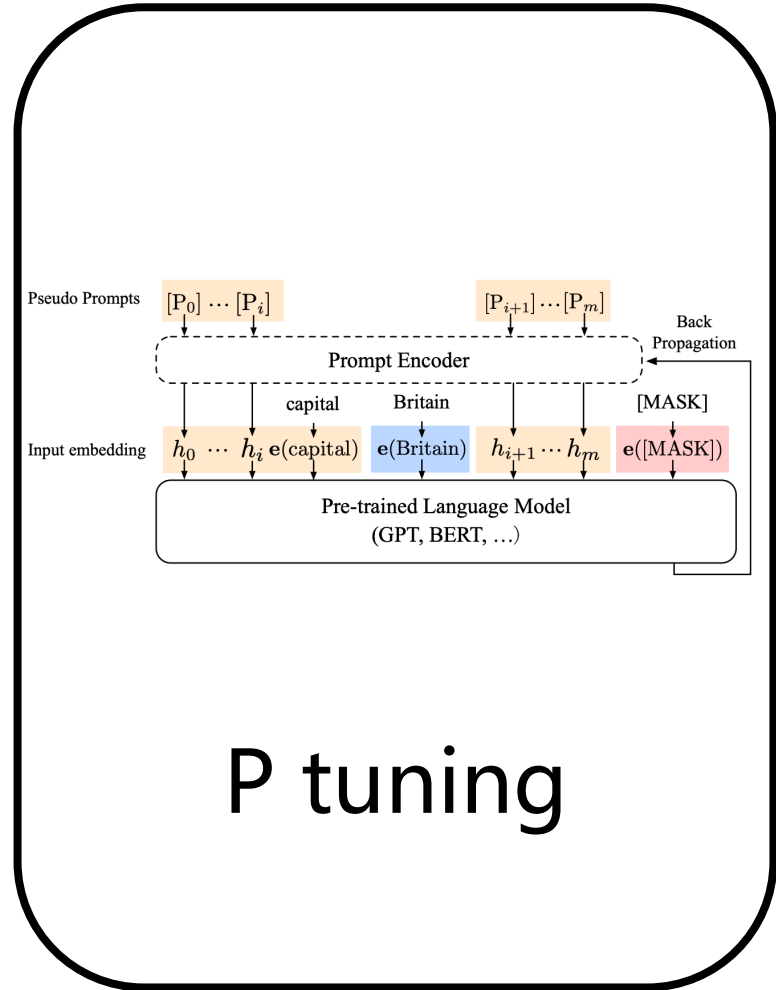
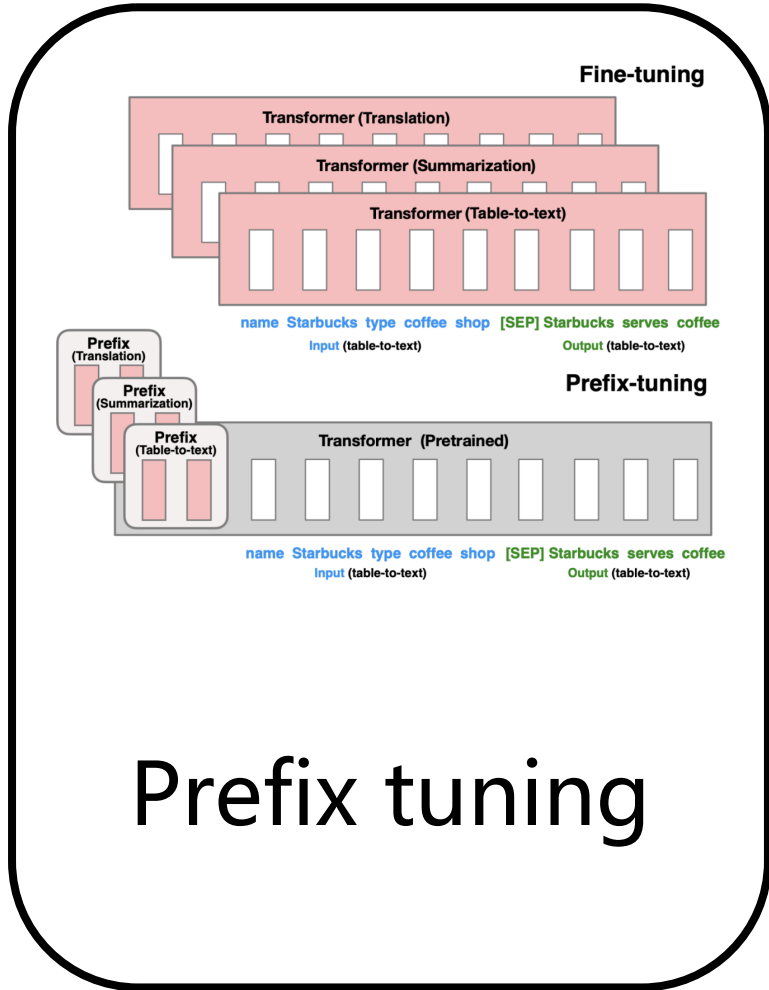
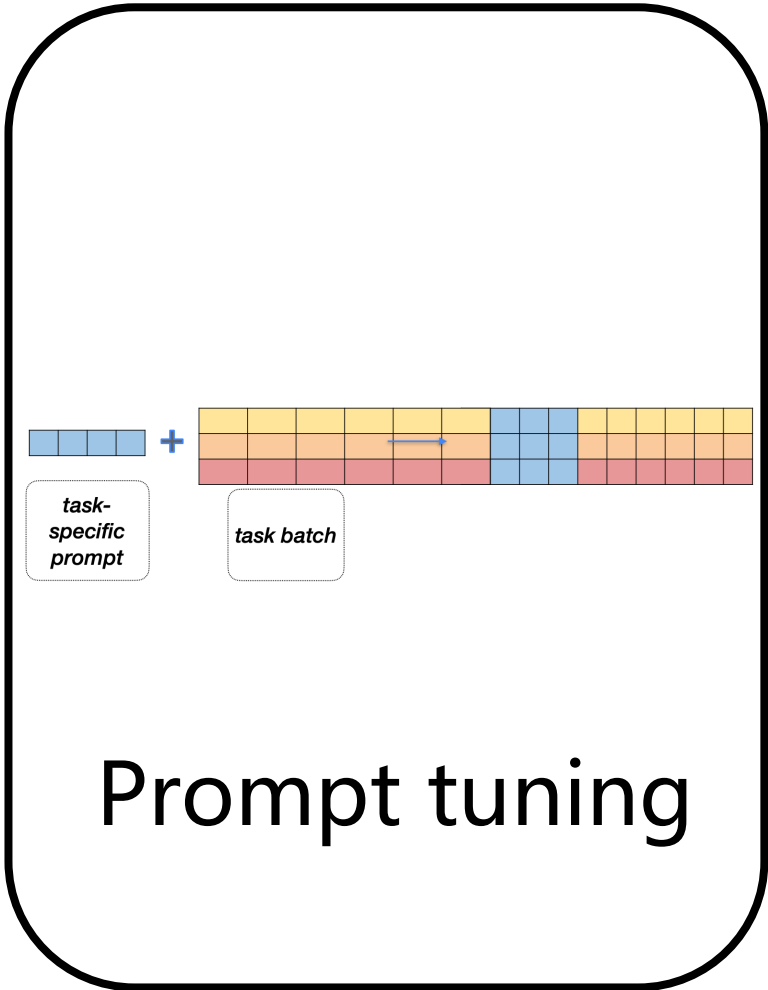


## Continuous/soft prompts



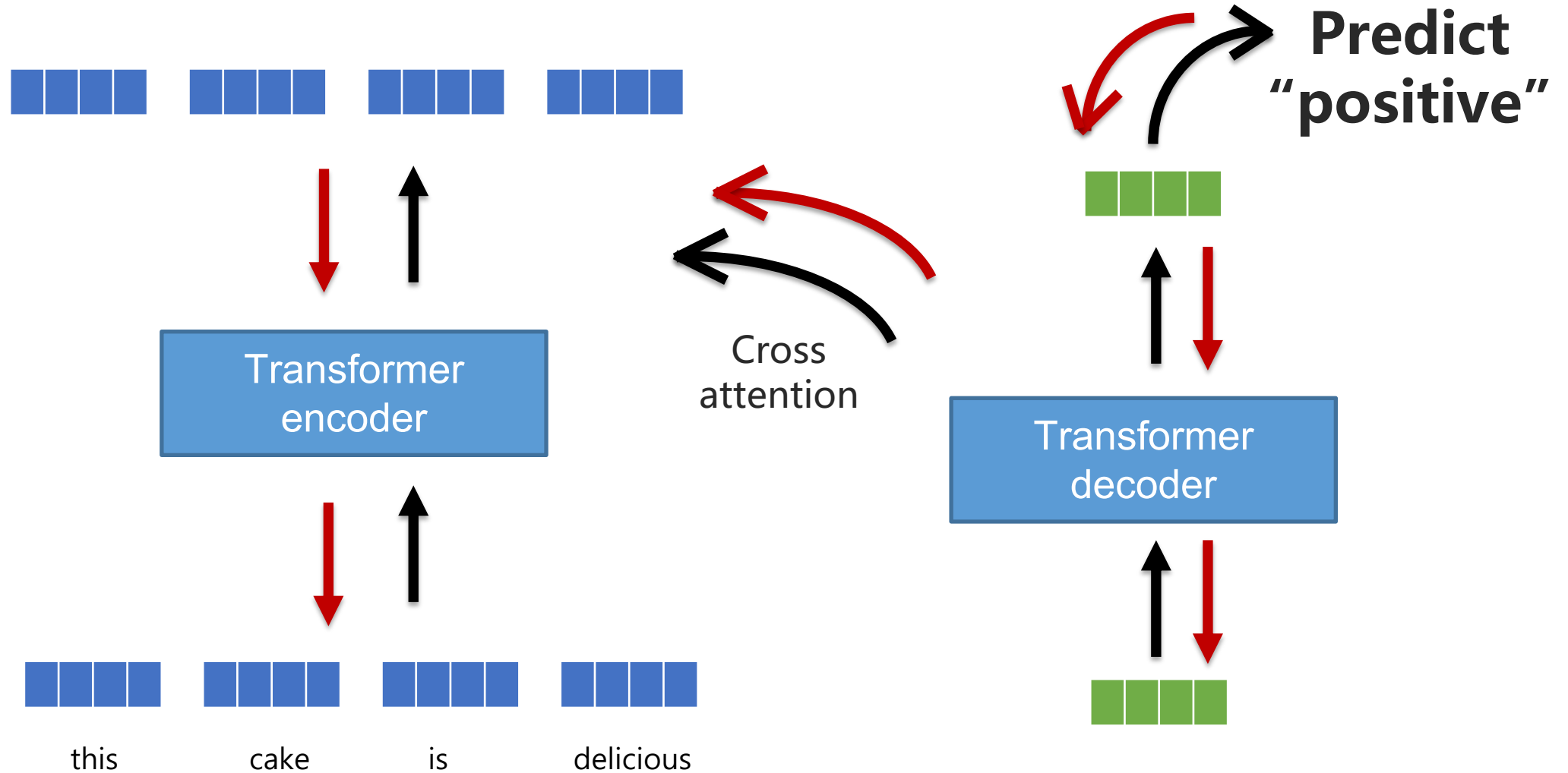


# Learn soft prompts effectively



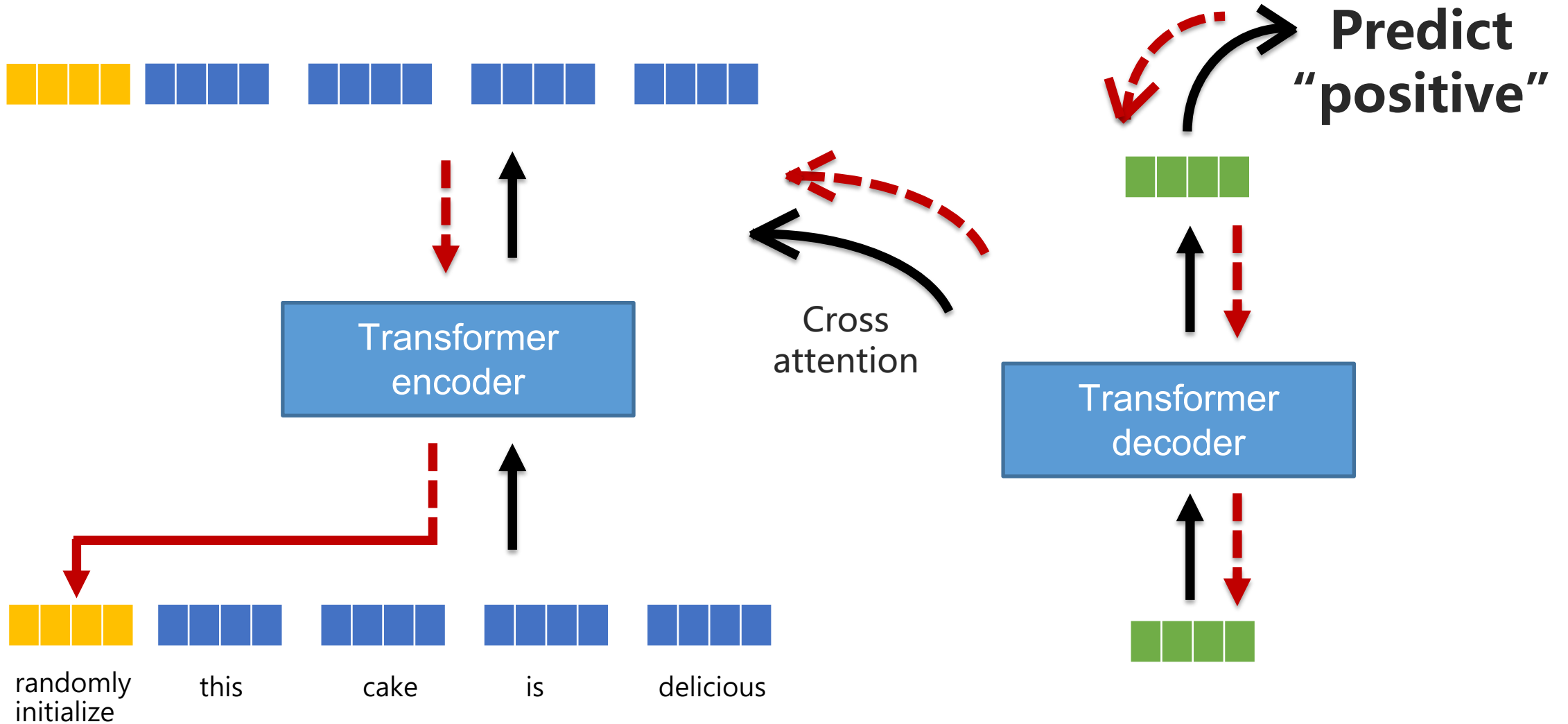


# Prompt tuning





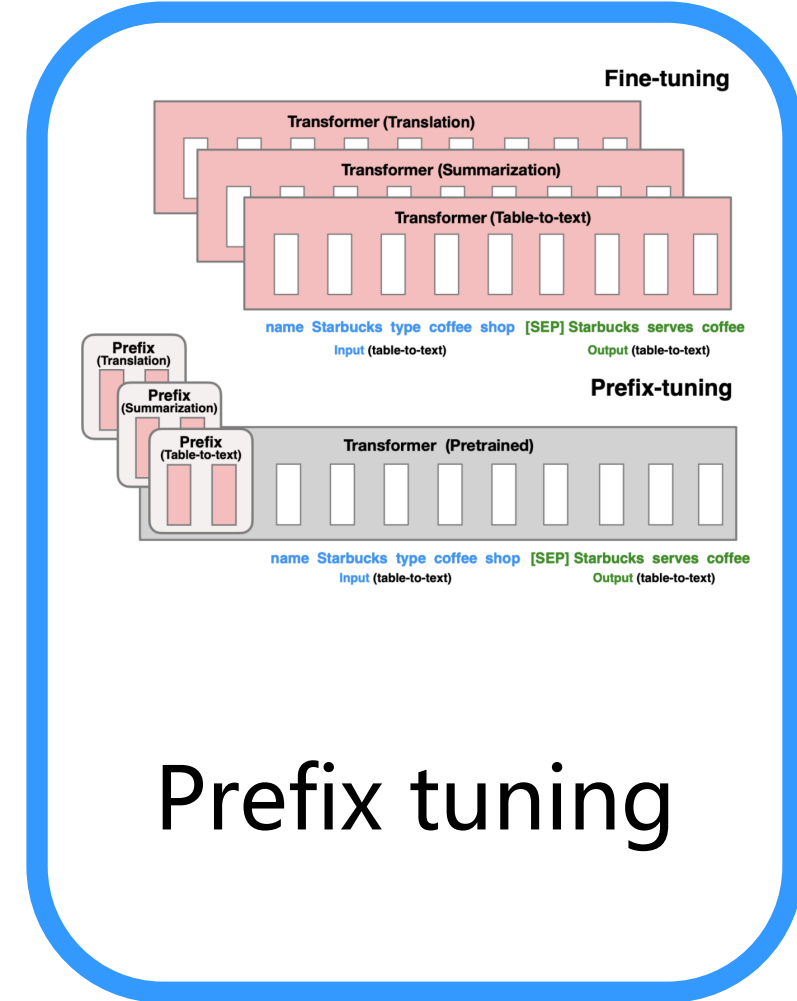
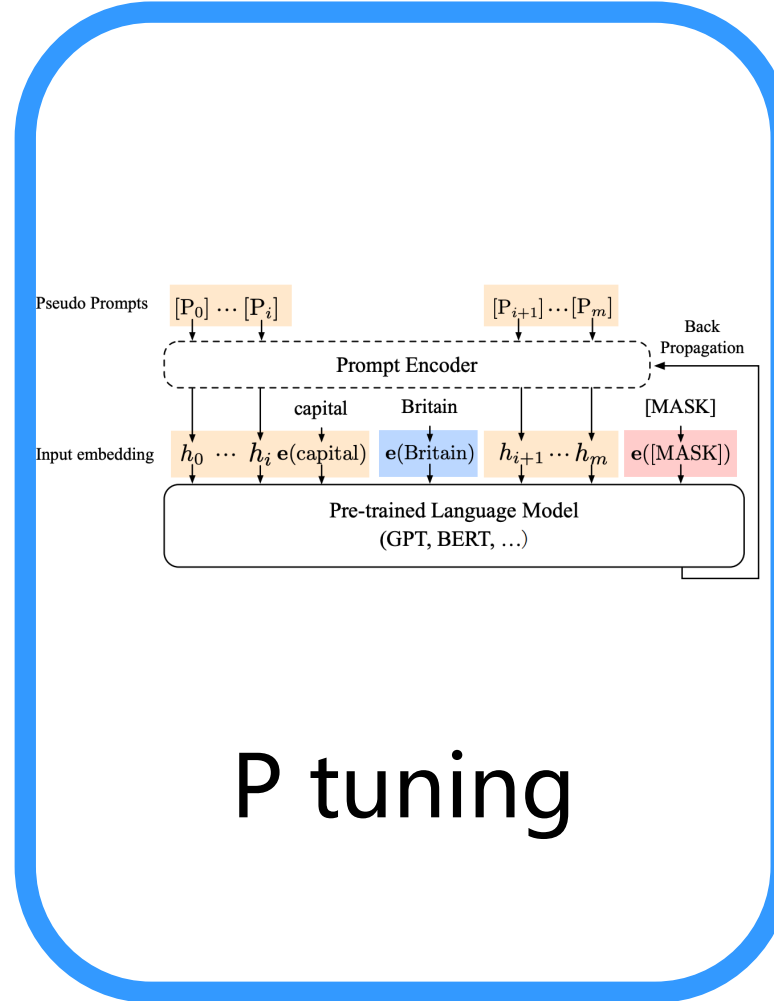
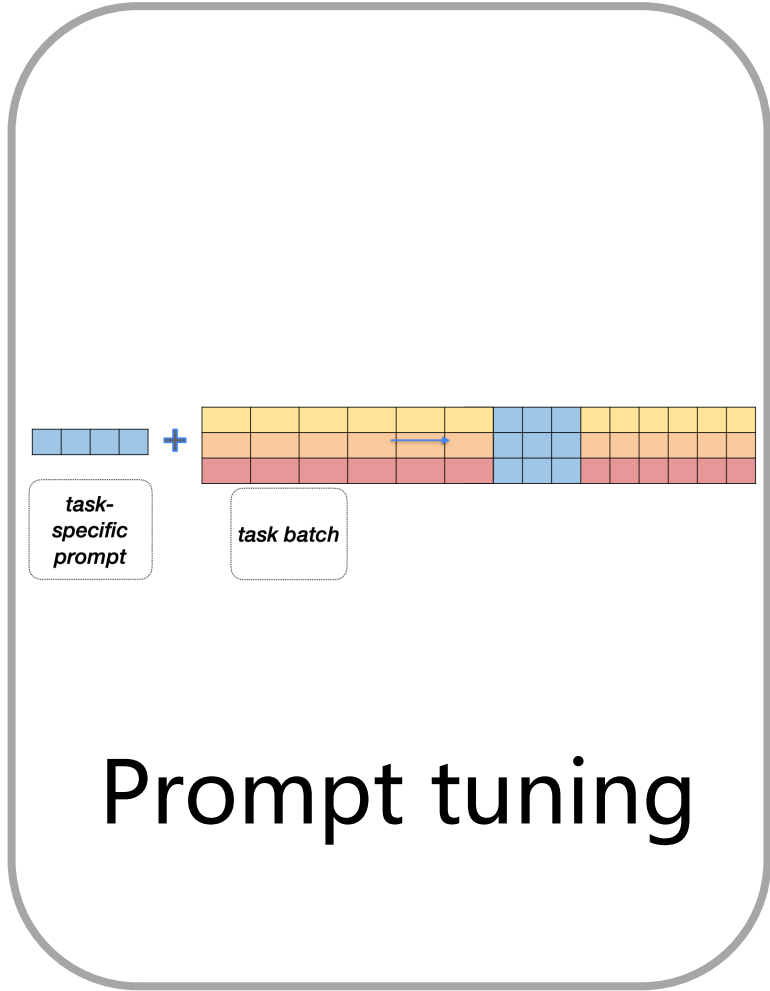
# Prompt tuning



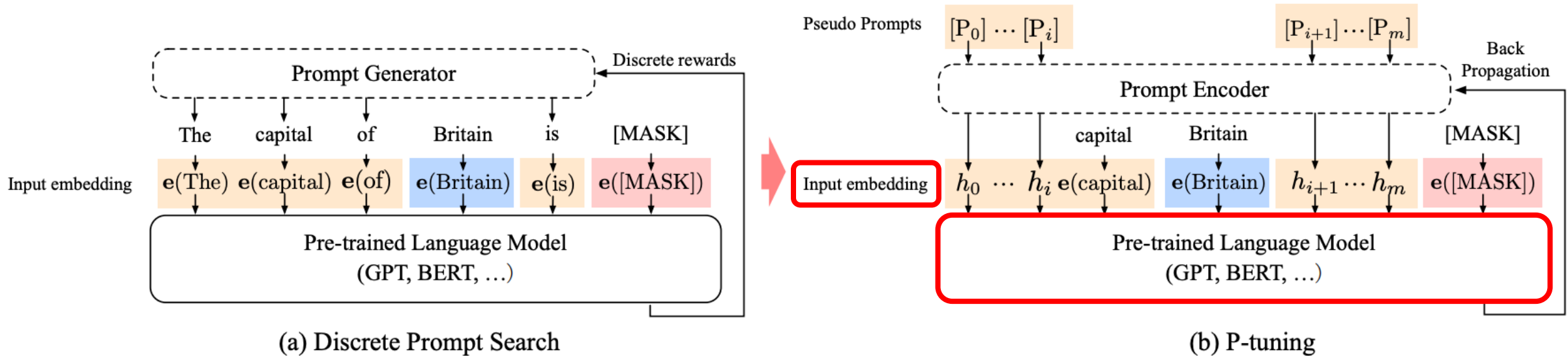




# Learn soft prompts effectively



Direct **optimize** the embeddings instead of prompt tokens.



**Hard Prompt**

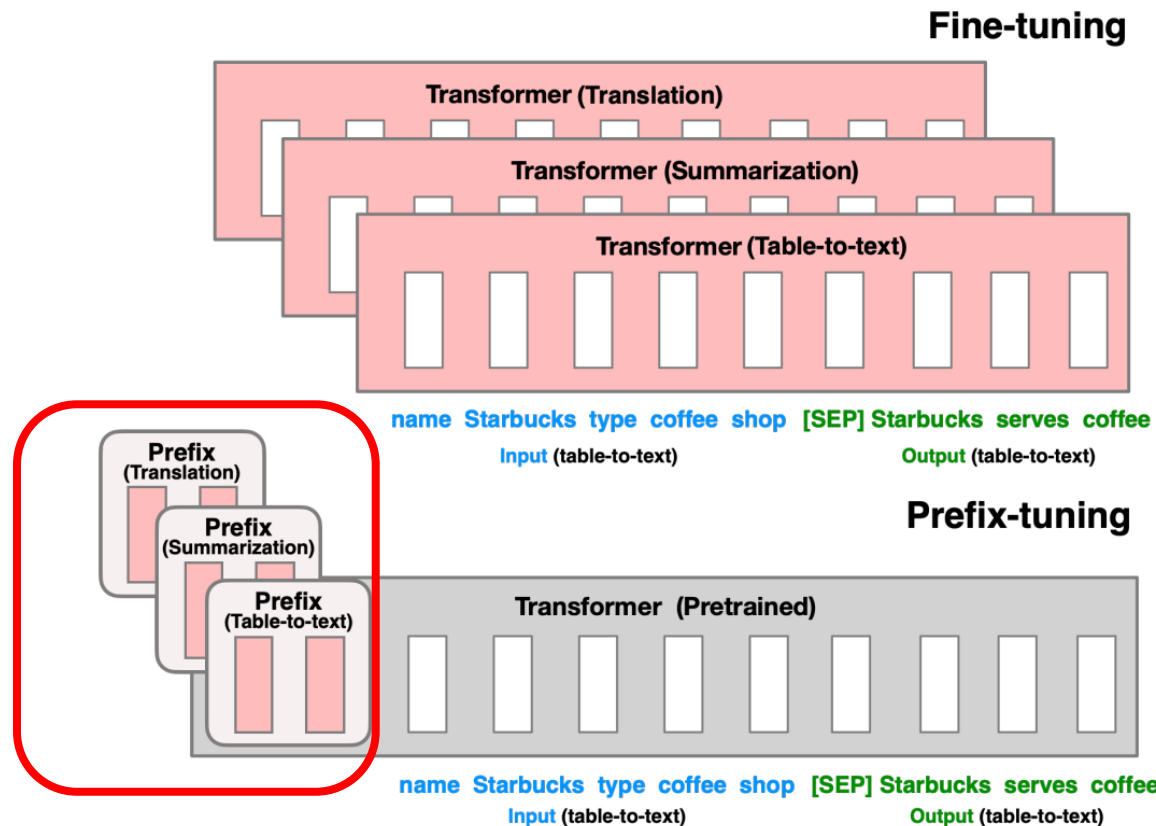
**Soft Prompt**

Prompt	$\mathcal{D}_{dev}$ Acc.	$\mathcal{D}_{dev32}$ Acc.
Does [PRE] agree with [HYP]? [MASK].	57.16	53.12
Does [HYP] agree with [PRE]? [MASK].	51.38	50.00
Premise: [PRE] Hypothesis: [HYP] Answer: [MASK].	68.59	55.20
[PRE] question: [HYP]. true or false? answer: [MASK].	70.15	53.12
P-tuning	76.45	56.25



# Prefix-Tuning

- Only optimize the prefix embeddings(all layers) for efficiency
- Prefix tuning stores 1000x fewer parameters than a fully finetuned model



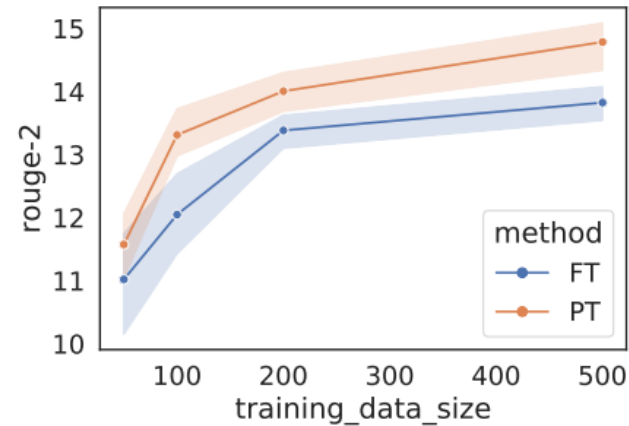
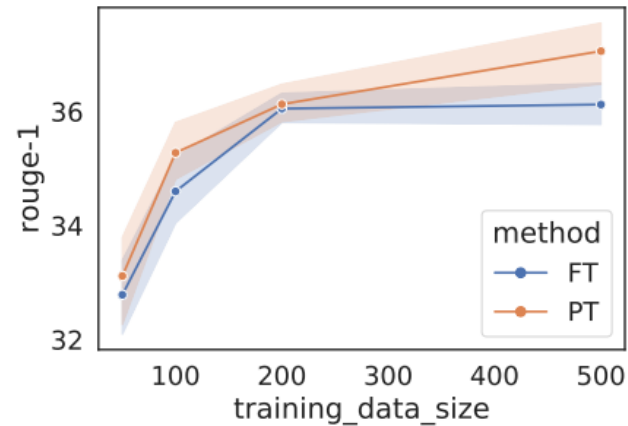
**Efficiency  
for time and space**



# Prefix-Tuning

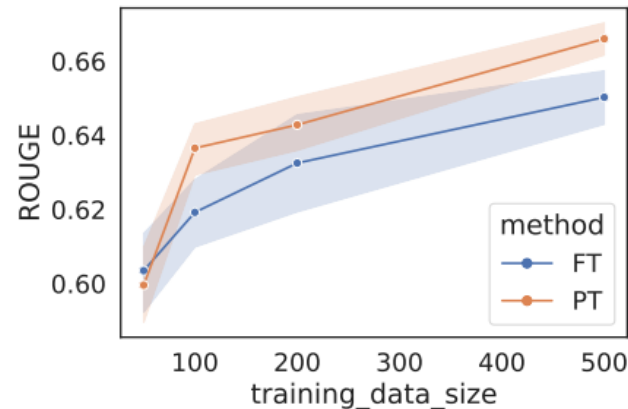
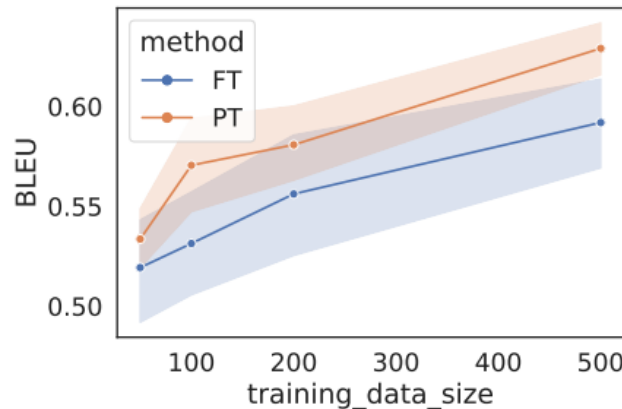
- Prefix-tuning has a comparative advantage when the number of **training examples is smaller**.

## Summarization



—●— FT(Fine-Tuning)  
—●— PT(Prefix-tuning)

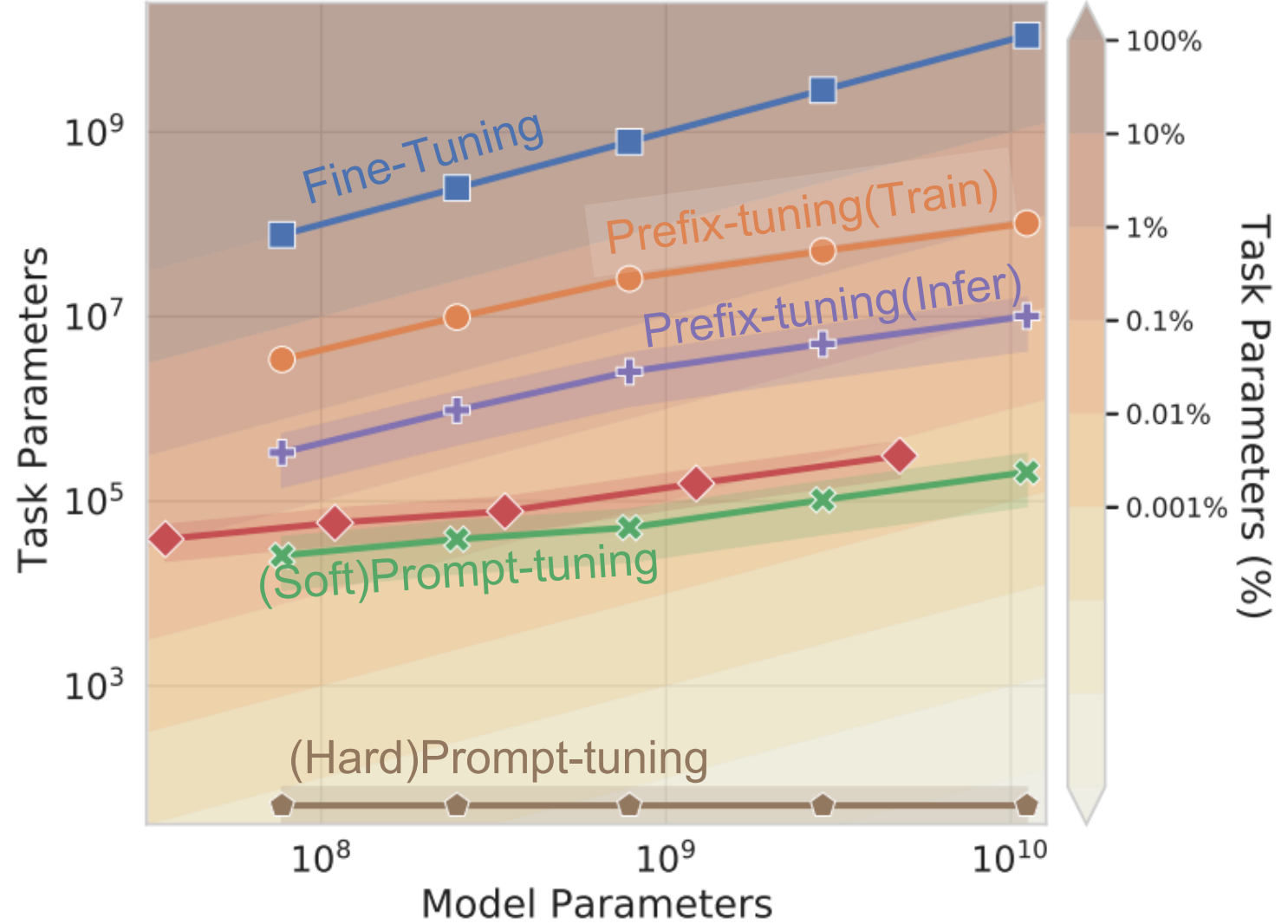
## Table-to-text



- Y-axis is the evaluation metric (higher is better)



# Prefix-Tuning, Prompt-Tuning





- 台大資訊 深度學習之應用 | ADL 15.1: Issues of PLMs 如何提示預訓練模型
- 台大資訊 深度學習之應用 | ADL 15.2: (Hard) Prompt-Tuning, LM-BFF 用自然語言提示模型
- 台大資訊 深度學習之應用 | ADL 15.3: (Soft) Prompt-Tuning (P-Tuning, Prefix Tuning) 人不懂沒關係機器懂就好
- 台大資訊 深度學習之應用 | ADL 15.4: Instruction Tuning 讓機器了解任務指令
- 台大資訊 深度學習之應用 | ADL 15.5: Prompting Paradigm 基於提示的研究大補帖





- Recent Advances in Pre-trained Language Models: [Why Do They Work and How to Use Them](#)
- <https://d223302.github.io/AAACL2022-Pretrain-Language-Model-Tutorial/>



# Relative papers for Tuning

- **REVIEW Paper:**
  - Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing, <https://arxiv.org/pdf/2107.13586.pdf>
- **P-Tuning:**
  - GPT Understands, Too, <https://arxiv.org/pdf/2103.10385.pdf>
  - P-Tuning v2: Prompt Tuning Can Be Comparable to Fine-tuning Universally Across Scales and Tasks, <https://arxiv.org/pdf/2110.07602.pdf>
- **Prompt Tuning:**
  - The Power of Scale for Parameter-Efficient Prompt Tuning, <https://arxiv.org/pdf/2104.08691.pdf>
- **Prefix-Tuning:**
  - Prefix-Tuning: Optimizing Continuous Prompts for Generation, <https://arxiv.org/pdf/2101.00190.pdf>
- **Soft-prompt:**
  - Learning How to Ask: Querying LMs with Mixtures of Soft Prompts, <https://arxiv.org/pdf/2104.06599.pdf>



# Case Study – OpenAI

## Some pre-departure reminder



**DO NOT use business / private data!**



**When want to train, then create an account (\$18 for free).**

### 1. Personal information we collect

We collect information that alone or in combination with other information in our possession could be used to identify you (“Personal Information”) as follows:

**Personal Information You Provide:** We may collect Personal Information if you create an account to use our Services or communicate with us as follows:

- *Account Information:* When you create an account with us, we will collect information associated with your account, including your name, contact information, account credentials, payment card information, and transaction history, (collectively, “Account Information”).
- *User Content:* When you use our Services, we may collect Personal Information that is included in the input, file uploads, or feedback that you provide to our Services (“Content”).



#### Free trial usage



GRANT #	CREDIT GRANTED	EXPIRES (UTC)
Grant 1	\$18.00	Expired 2023年4月1日



# Case Study – OpenAI

## Work flow

### Step 4: Select model

Model	Training	Usage
Ada	\$0.0004 / 1K tokens	\$0.0016 / 1K tokens
Babbage	\$0.0006 / 1K tokens	\$0.0024 / 1K tokens
Curie	\$0.0030 / 1K tokens	\$0.0120 / 1K tokens
Davinci	\$0.0300 / 1K tokens	\$0.1200 / 1K tokens

### Step 5: clean data (delete repeat, split validation set ... )

```
openai tools fine_tunes.prepare_data -f <LOCAL_FILE>
```

### Step 6: train

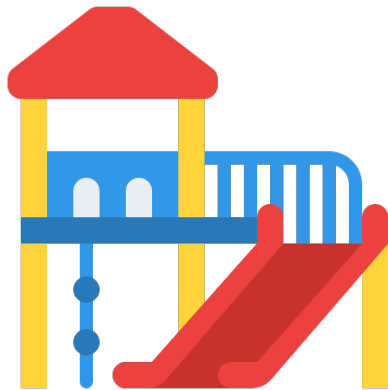
```
openai api fine_tunes.create -t <TRAIN_FILE_ID_OR_PATH> -m <BASE_MODEL>
```

### Step 7: trace the progress

```
openai api fine_tunes.follow -i <YOUR_FINE_TUNE_JOB_ID>
```



# Case Study – OpenAI



## Play Time





# Case Study – OpenAI

## Step 7: trace the progress (Garbage message or NOT): 3091 pairs

```
openai api fine_tunes.follow -i <YOUR_FINE_TUNE_JOB_ID>
```

```
[2023-06-02 16:46:49] Created fine-tune: ft-9xwGZYXgl65dXew2tLTwYWYH
```

```
[2023-06-02 16:49:44] Fine-tune costs $35.15
```

```
[2023-06-02 16:49:45] Fine-tune enqueued. Queue number: 3
```

```
[2023-06-02 16:51:25] Fine-tune is in the queue. Queue number: 2
```

```
[2023-06-02 16:51:49] Fine-tune is in the queue. Queue number: 1
```

```
[2023-06-02 16:52:24] Fine-tune is in the queue. Queue number: 0
```

```
[2023-06-02 16:52:31] Fine-tune started
```

```
[2023-06-02 17:09:32] Completed epoch 1/4
```

```
[2023-06-02 17:36:31] Completed epoch 3/4
```

```
[2023-06-02 17:50:29] Uploaded model: davinci:ft-personal-2023-06-02-09-50-29
```

```
[2023-06-02 17:50:31] Uploaded result file: file-WoUhP6XjpKGyZznvthkviMIIf
```

```
[2023-06-02 17:50:31] Fine-tune succeeded
```





# Case Study – OpenAI

## Step 7: trace the progress (Detail Classification): 1348 pairs

```
openai api fine_tunes.follow -i <YOUR_FINE_TUNE_JOB_ID>
```

```
[2023-06-02 19:33:33] Created fine-tune: ft-t3V3yk1K11TZxy0CWFbiP2nG
```

```
[2023-06-02 21:05:00] Fine-tune costs $25.71
```

```
[2023-06-02 21:05:00] Fine-tune enqueued. Queue number: 2
```

```
[2023-06-02 21:05:03] Fine-tune is in the queue. Queue number: 1
```

```
[2023-06-02 21:05:05] Fine-tune is in the queue. Queue number: 0
```

```
[2023-06-02 21:05:09] Fine-tune started
```

```
[2023-06-02 21:19:43] Completed epoch 1/4
```

```
[2023-06-02 21:43:14] Completed epoch 3/4
```

```
[2023-06-02 21:55:55] Uploaded model: davinci:ft-personal-2023-06-02-13-55-55
```

```
[2023-06-02 21:55:57] Uploaded result file: file-8IorcwLDluWd7JbRwMQG0m0V
```

```
[2023-06-02 21:55:57] Fine-tune succeeded
```



# Case Study – OpenAI

## Scenario 1: Garbage message or NOT

Overview Documentation API reference Examples Playground

Help Personal

Playground

Load a preset...

Save

View code

Share



Write a tagline for an ice cream shop.



Mode

Complete

Model

davinci:ft-persona...

Temperature

1

Maximum length

256

Stop sequences

Enter sequence and press Tab

end

Top P

1

Frequency penalty

0

Presence penalty

0

Looking for ChatGPT?

[Try it now](#)

Submit





# Case Study – OpenAI

## Scenario 2: Detail Classification

Overview Documentation API reference Examples **Playground**

Help Personal

Playground

Load a preset...

Save

View code

Share



Write a tagline for an ice cream shop.



Mode

Complete

Model

davinci:ft-persona...

Temperature

1

Maximum length

256

Stop sequences

Enter sequence and press Tab

end

X

Top P

1

Frequency penalty

0

Presence penalty

0

Looking for ChatGPT?

[Try it now](#)

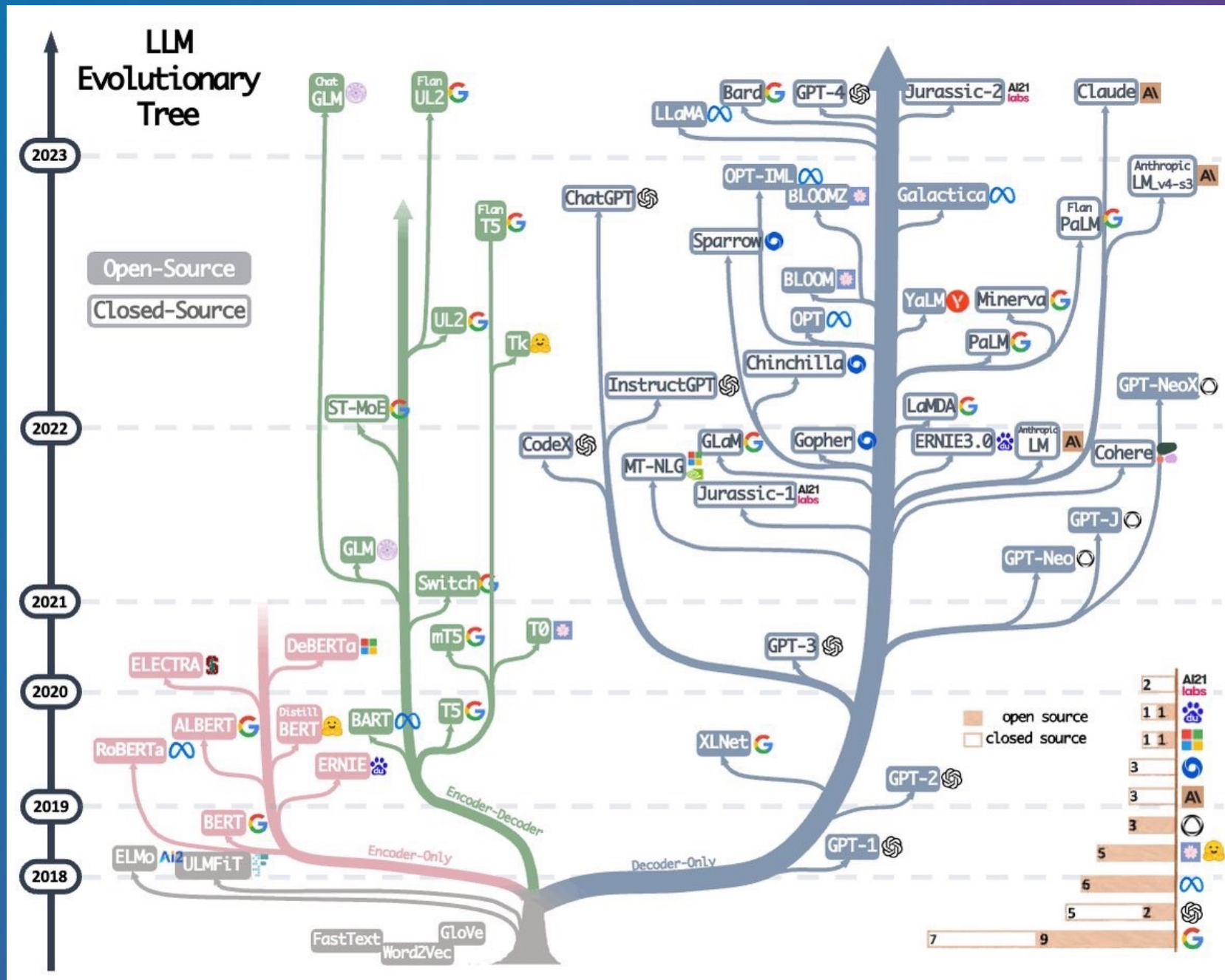
Submit



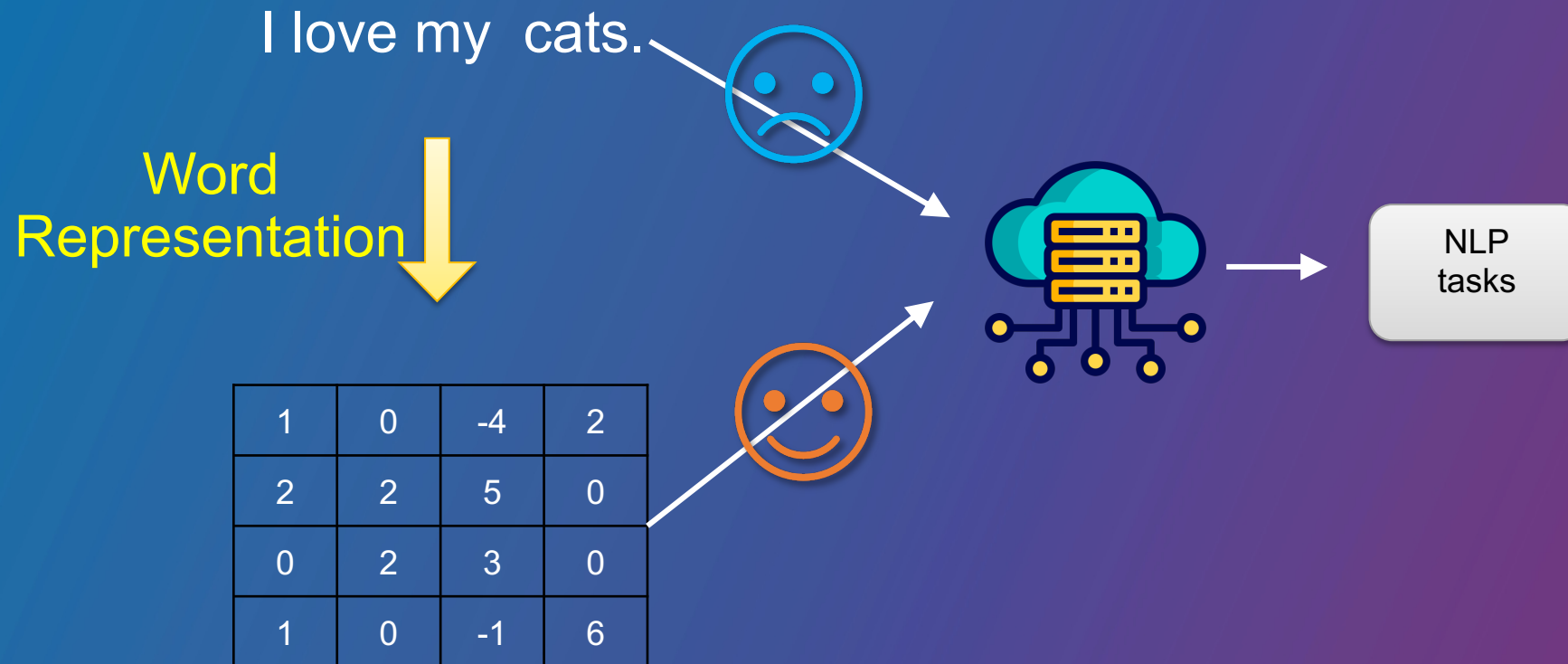


# LoRa for EMIC

- **Model:**
  - minlik/chinese-llama-plus-7b-merged
- **Datasets:**
  - EMIC
- **Hyperparameters in LoRa:**
  - `r=16`, #attention heads
  - `lora_alpha=32`, # scaling factor for the weight matrices
  - `lora_dropout=0.05`,
  - `bias="none"`,
- **Hyperparameters in training:**
  - `per_device_train_batch_size=1`,
  - `gradient_accumulation_steps=4`,
  - `warmup_steps=100`,
  - `max_steps=3000`,
  - `learning_rate=2e-4`,
  - `optim='adamw_torch'`,

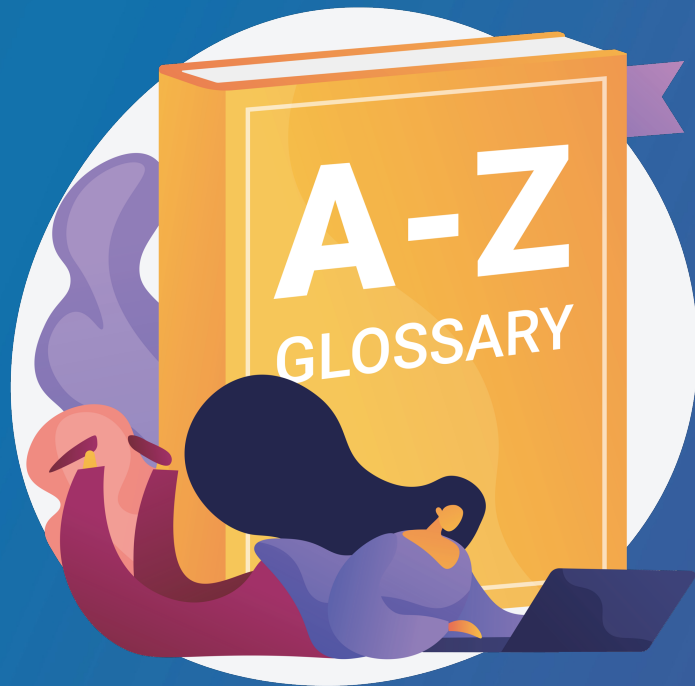


# The ability of computer to understand human language





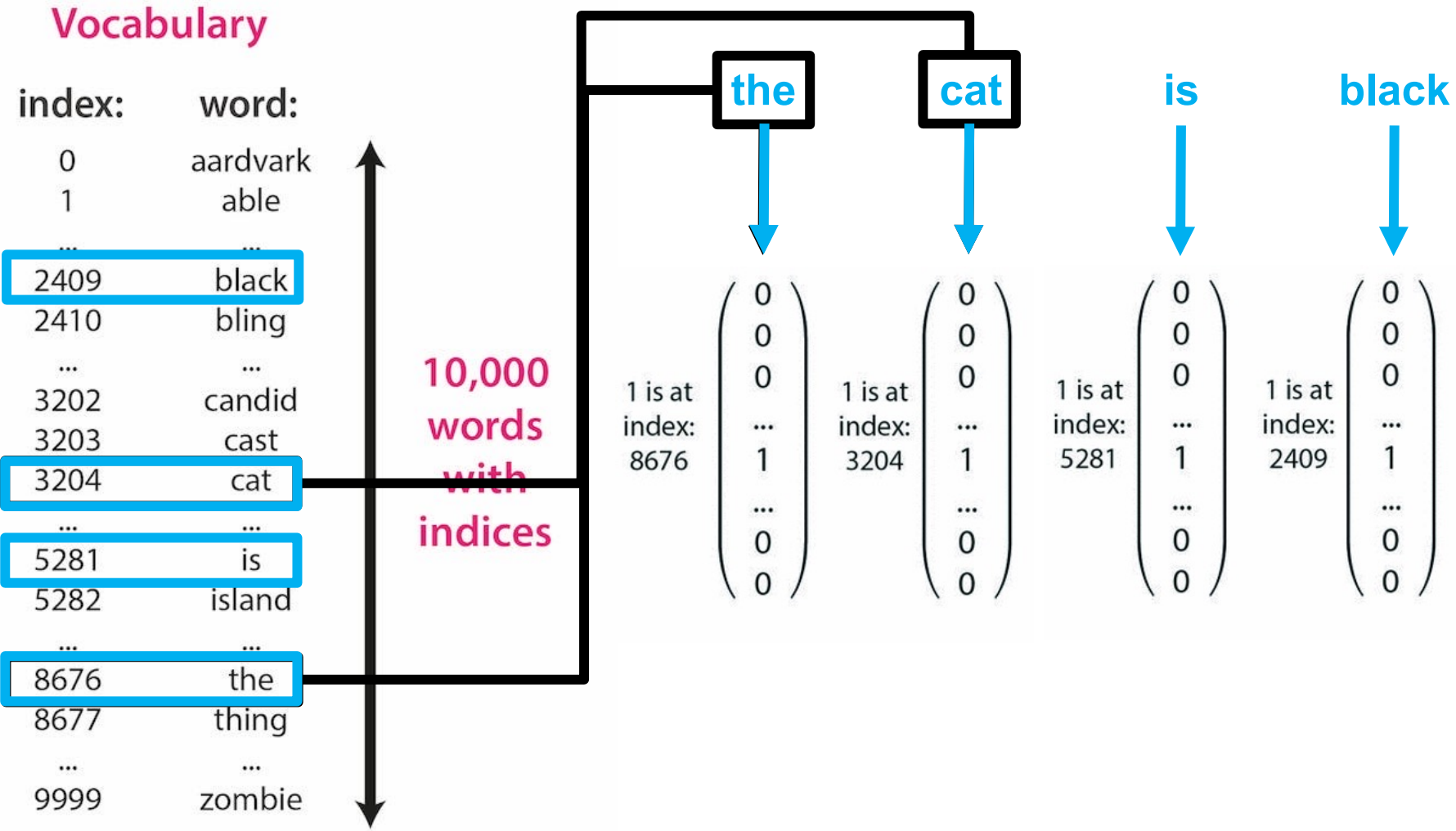
# Word Representation



## Corpus based

- One-hot Representation
- Distributional Representation

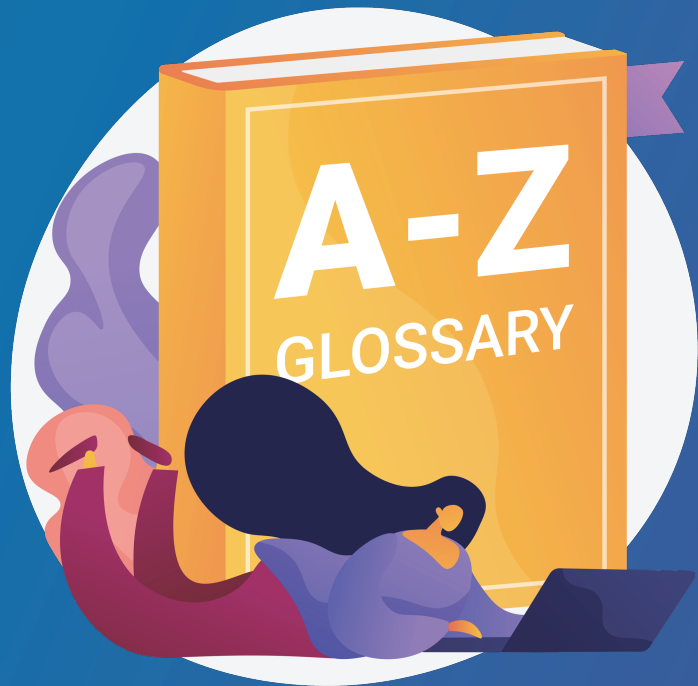
# One-Hot Encoding



## However , the limitations are...

- High-dimensional representation – vector sparsity
- Inability to capture semantic similarity
- Loss of word order information

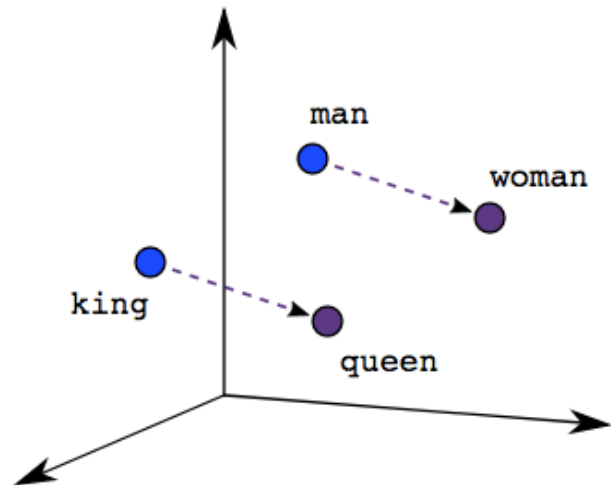
# Word Representation



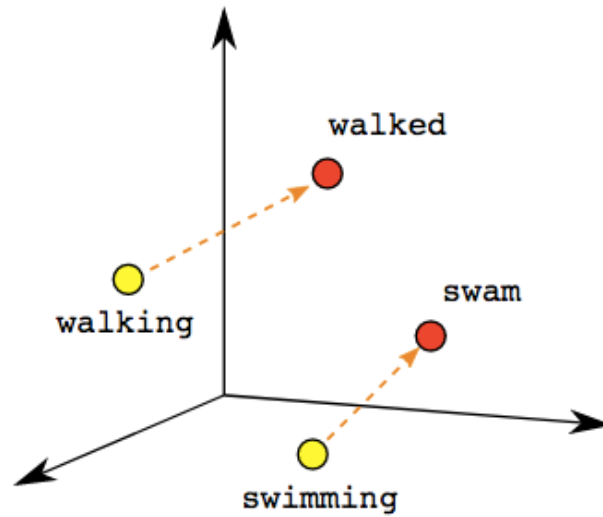
## Corpus based

- One-hot Representation
- Distributional Representation

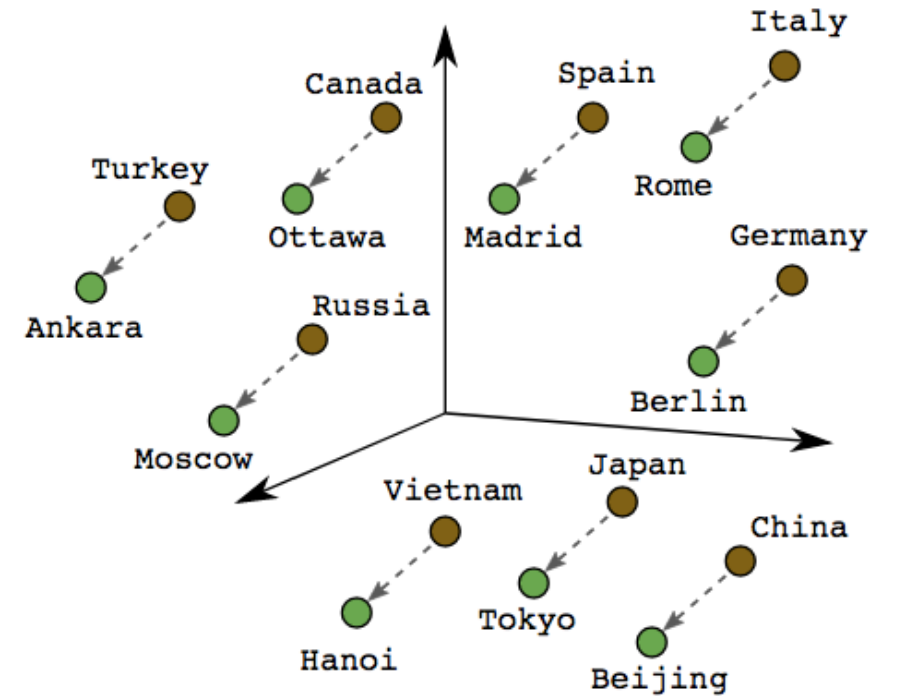
You shall know a word by the company it keeps (Firth, 1957)



Male-Female



Verb Tense



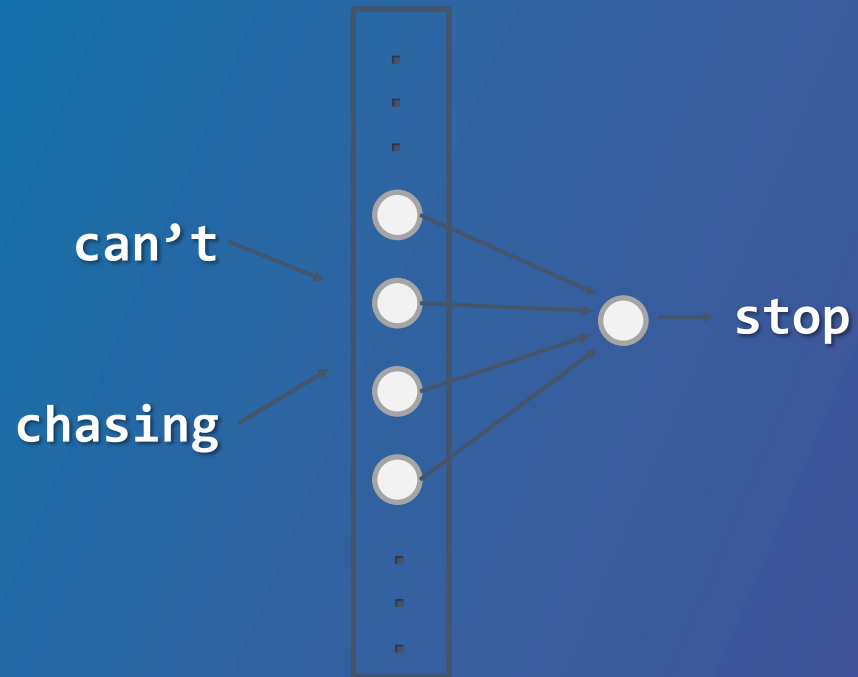
Country-Capital

# Word2Vec (2013)

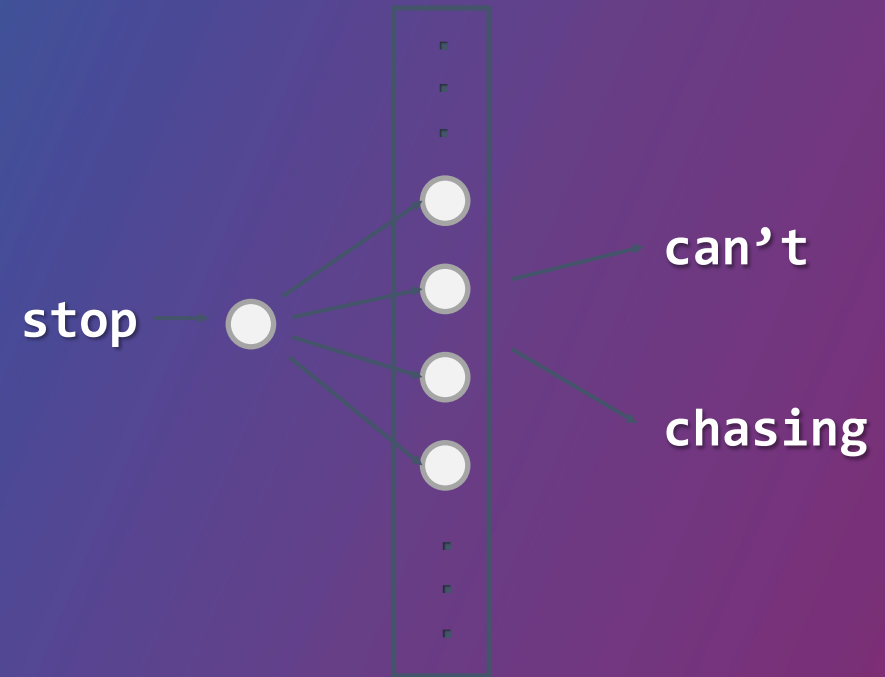
sliding window

The dog can't stop chasing the cat on the street.

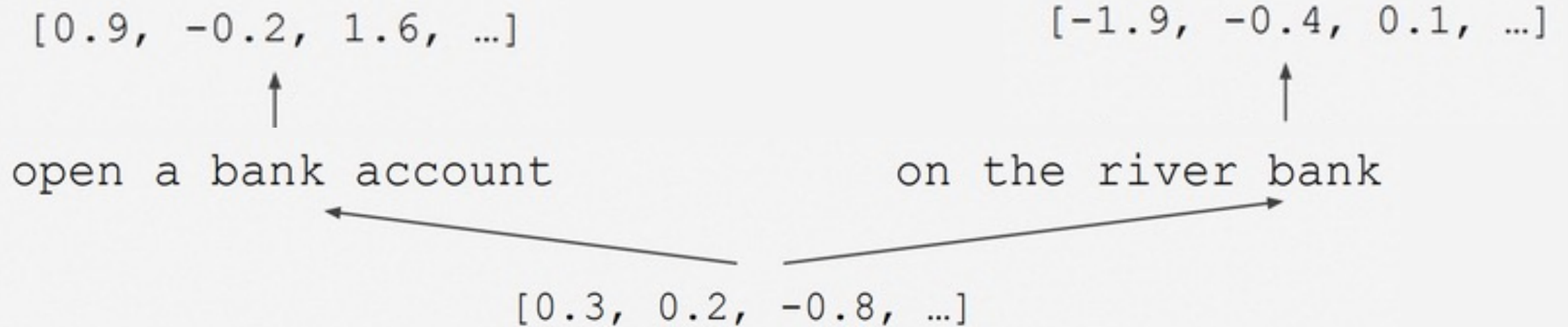
## CBOW



## Skip-gram



## contextual dependent

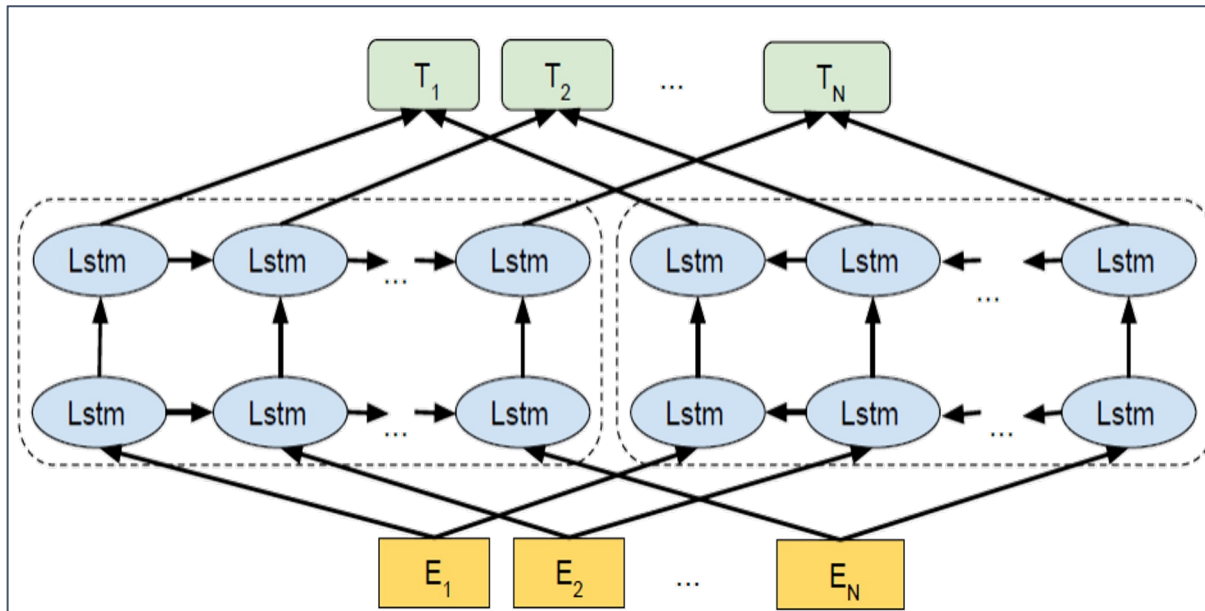


## contextual independent



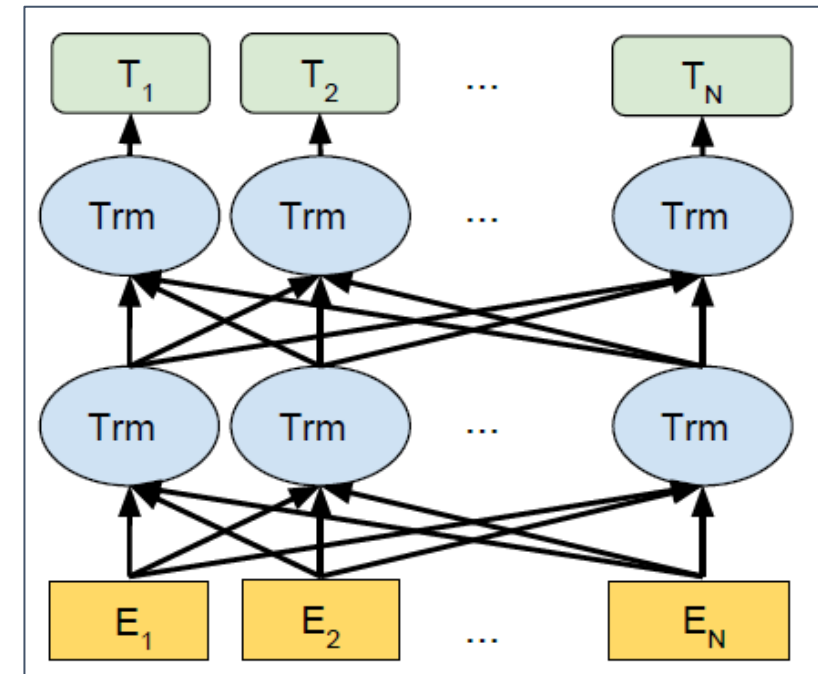
## Deep contextualized word representations (ELMo)

- **RNN based**
- **Bi-LSTM : forward + backward**
- **predicting next word**



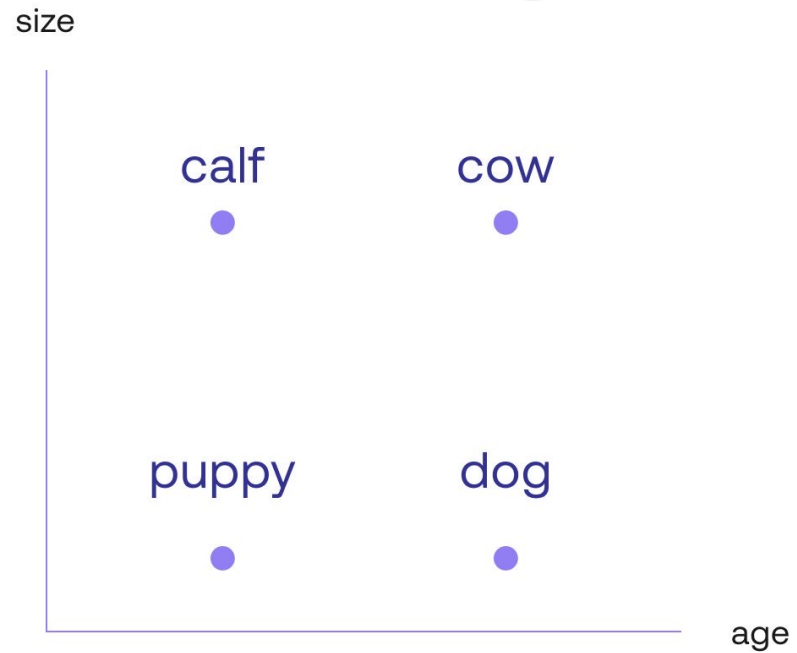
## Bidirectional Encoder Representations from Transformers (BERT)

- **Transformer based**
- **self-attention**
- **predicting masked words and next sentence relationships**



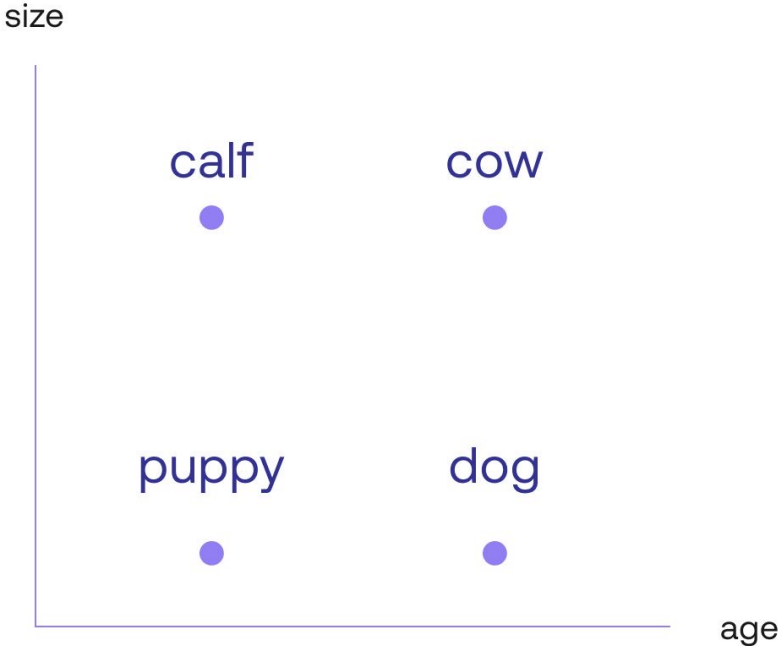
# Variation of Embeddings

## Word Embedding

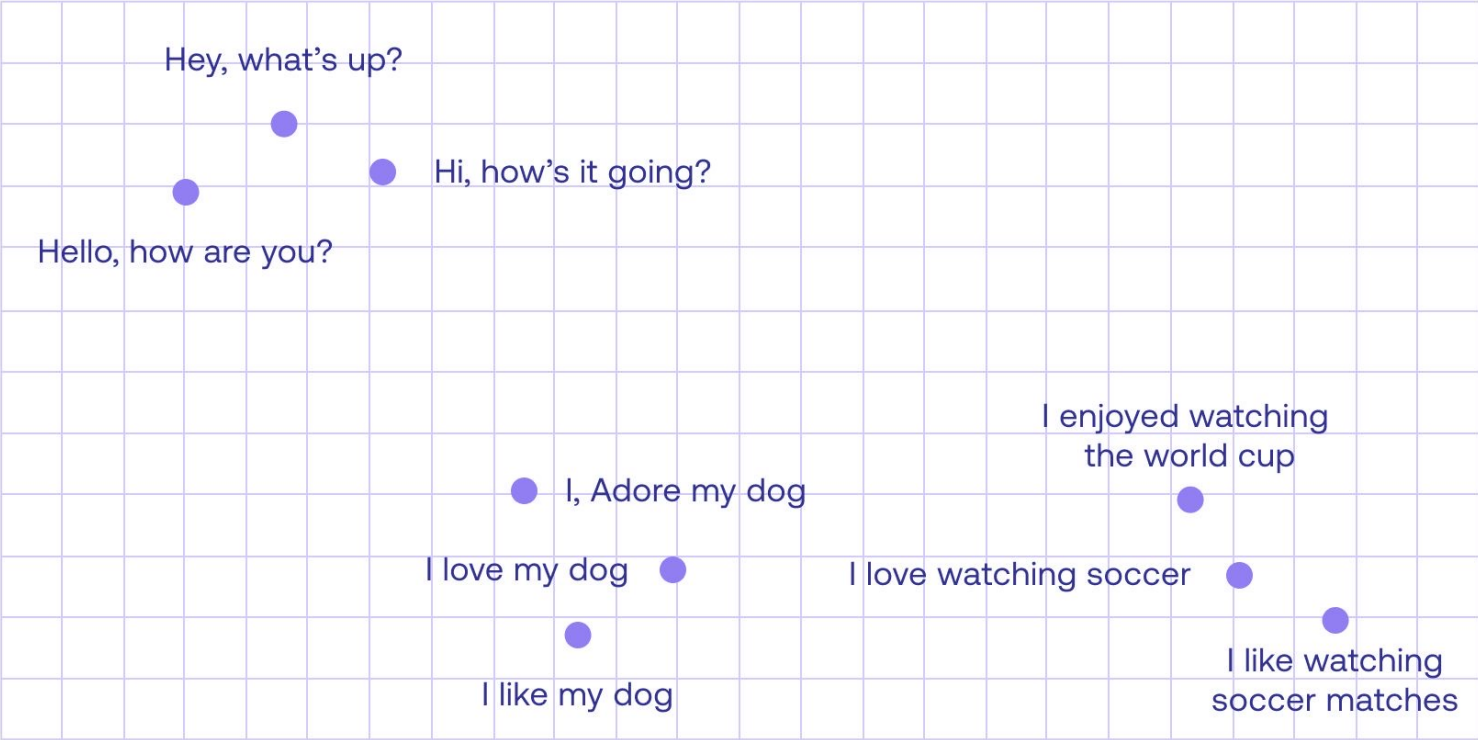


# Variation of Embeddings

## Word Embedding



## Sentence Embedding



# Much left unsaid ...

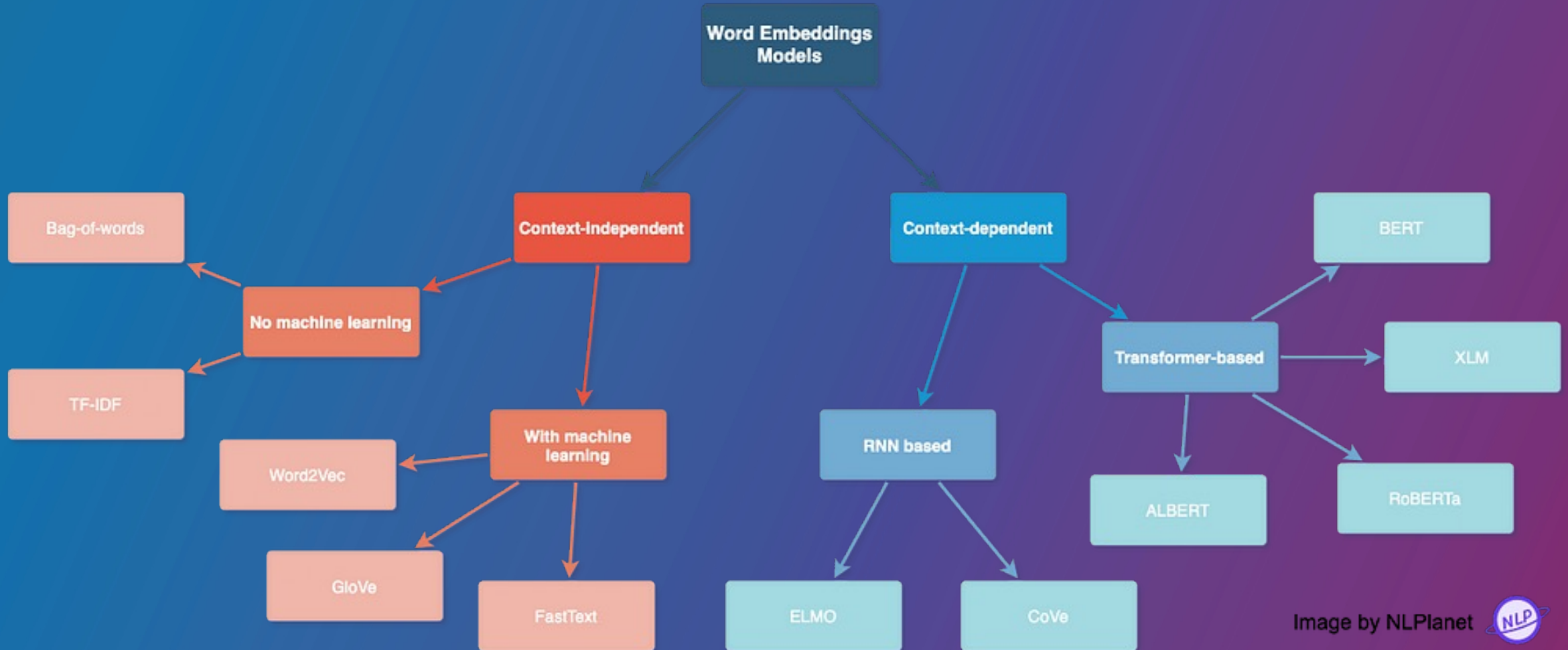


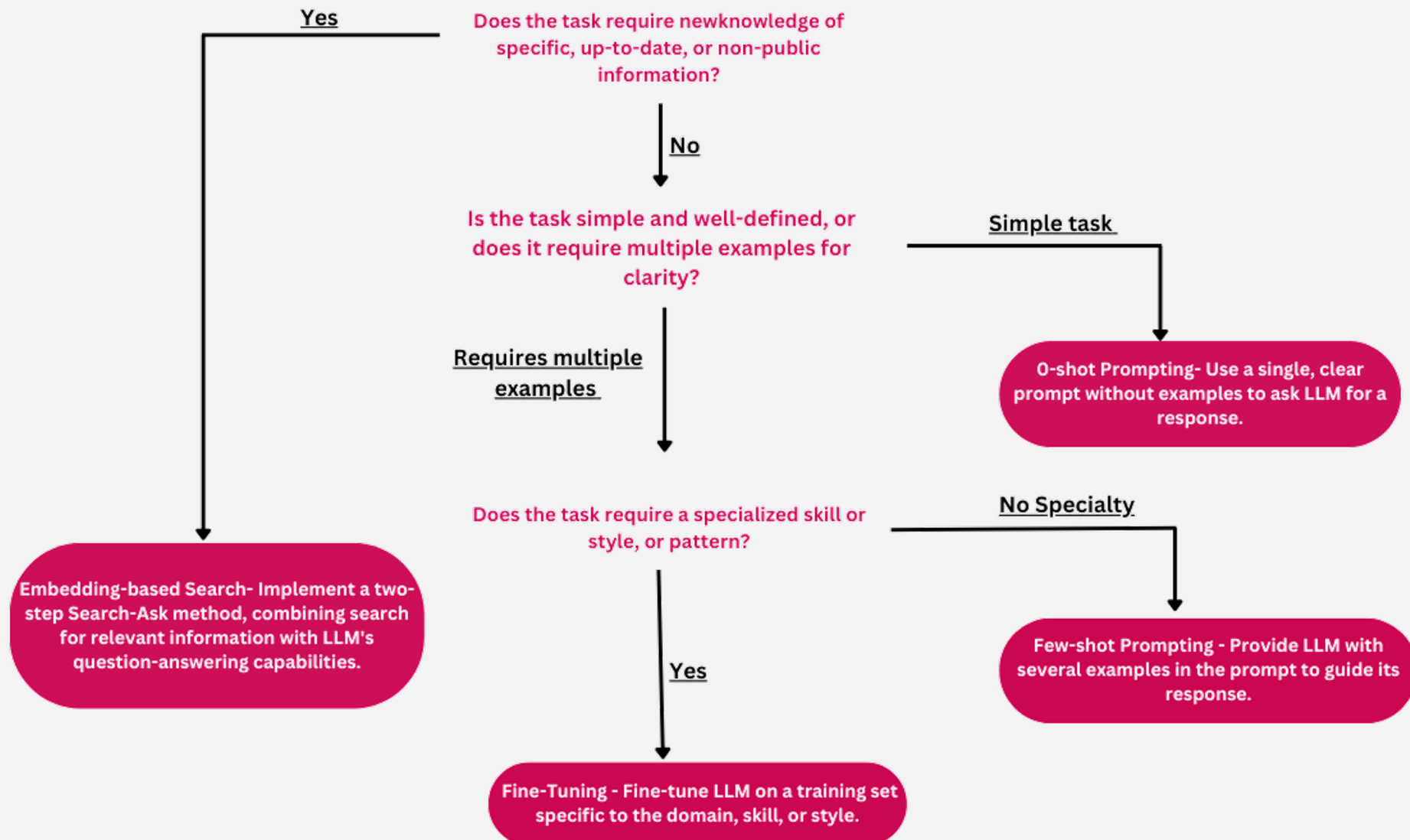
Image by NLPlanet 

## How to use word-embeddings

- Search
- Clustering
- Recommendations
- Anomaly detection
- Diversity measurement
- Classification

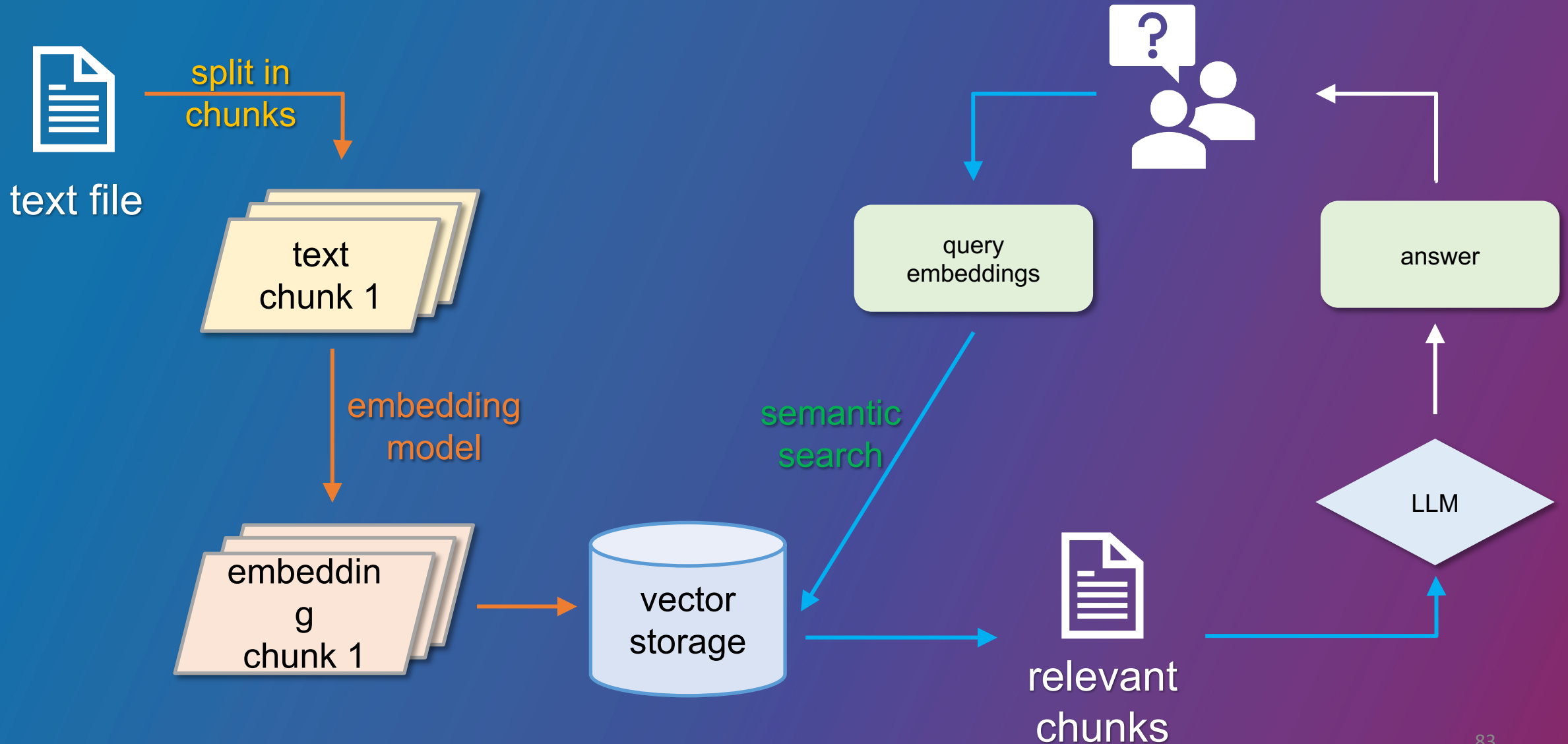
# Choosing the Right LLM Strategy

## 0-shot vs Few-shot vs Fine-tuning vs Embedding



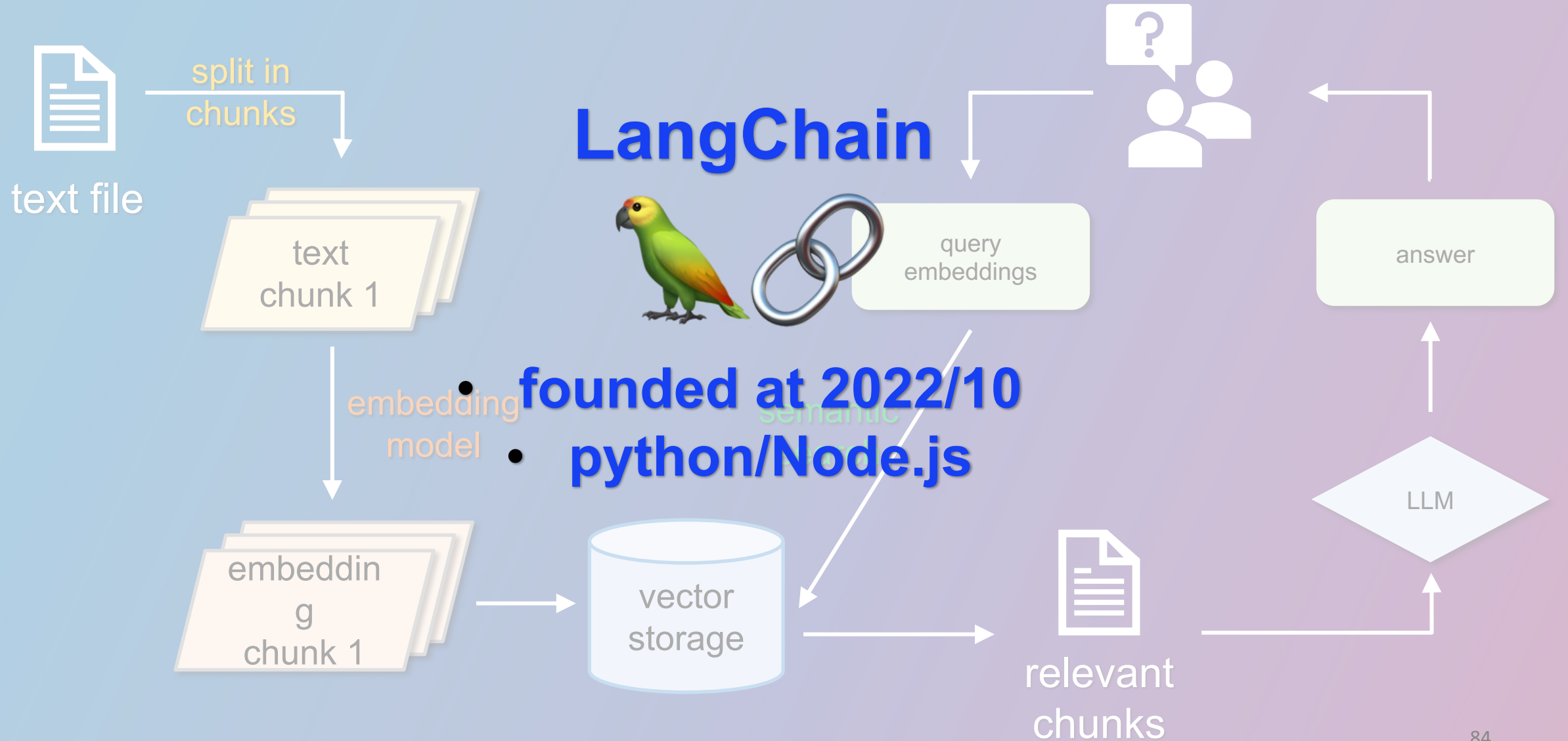


# Embedding-based Application - QA

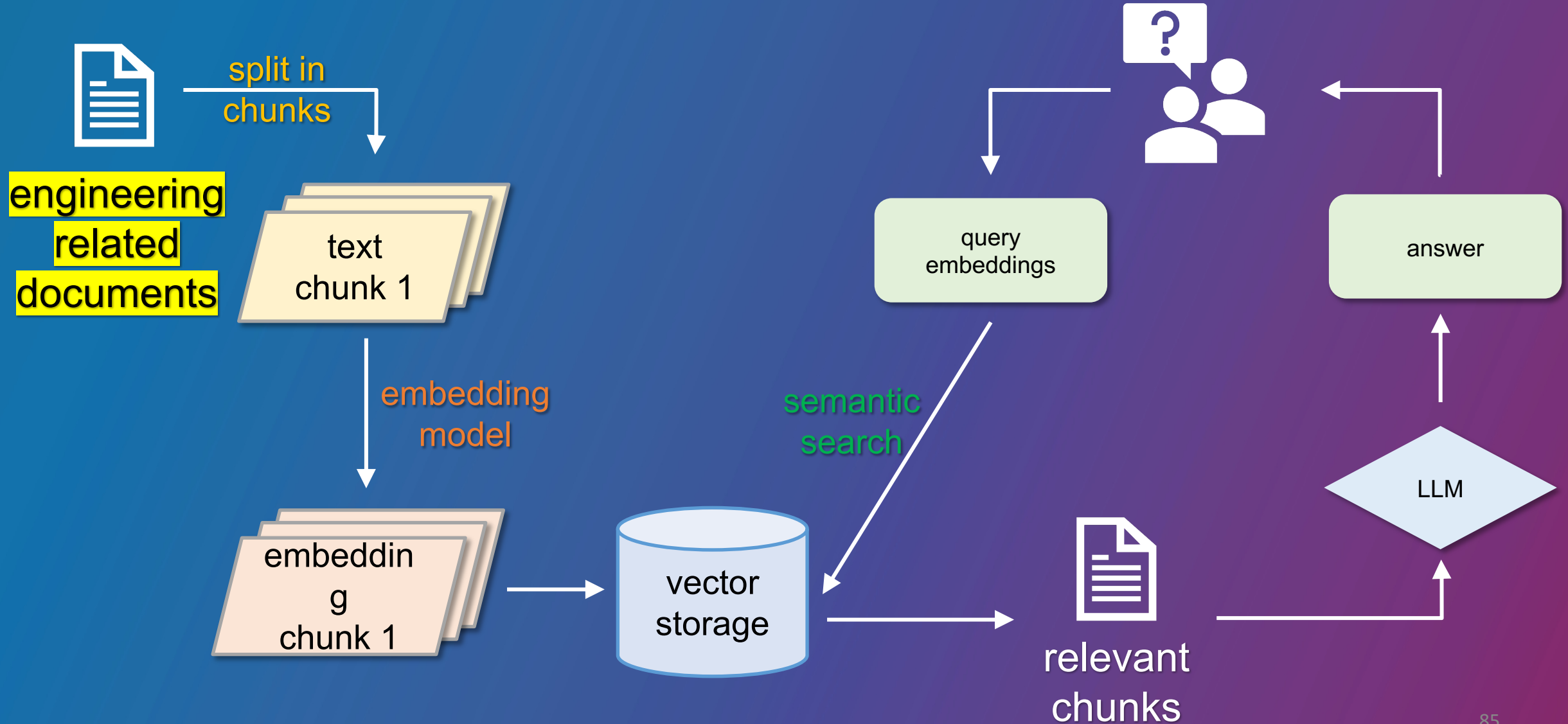




# Embedding-based Application



# Embedding-based Application – QA with documents



# model selection

Hugging Face is way more fun with friends and colleagues! 🤗 [Join an organization](#) Dismiss this message

Spaces: [mteb/leaderboard](#) like 302 Running on CPU UPGRADE App Files Community 9

Massive Text Embedding Benchmark (MTEB) Leaderboard. To submit, refer to the [MTEB GitHub repository](#) 🤗

- Total Datasets: 62
- Total Languages: 112
- Total Scores: >5550
- Total Models: 75

Overall Bitext Mining Classification Clustering Pair Classification Retrieval Reranking STS Summarization

Overall MTEB English leaderboard 🤗

- Metric: Various, refer to task tabs
- Languages: English, refer to task tabs for others

Rank	Model	Embedding Dimensions	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)	Pair Classification Average (3 datasets)	Reranking Average (4 datasets)	Retrieval Average (15 datasets)	STS Average (10 datasets)	Summarization Average (1 dataset)
1	<a href="#">e5-large-v2</a>	1024	62.25	75.24	44.49	86.03	56.61	50.56	82.05	30.19
2	<a href="#">instructor-xl</a>	768	61.79	73.12	44.74	86.62	57.29	49.26	83.06	32.32
3	<a href="#">instructor-large</a>	768	61.59	73.86	45.29	85.89	57.54	47.57	83.15	31.84
4	<a href="#">e5-base-v2</a>	768	61.5	73.84	43.8	85.73	55.91	50.29	81.05	30.28
5	<a href="#">e5-large</a>	1024	61.42	73.14	43.33	85.94	56.53	49.99	82.06	30.97
6	<a href="#">text-embedding-ada-002</a>	1536	60.99	70.93	45.9	84.89	56.32	49.25	80.97	30.8
7	<a href="#">e5-base</a>	768	60.44	72.63	42.11	85.09	55.7	48.75	80.96	31.01
8	<a href="#">e5-small-v2</a>	384	59.93	72.94	39.92	84.67	54.32	49.04	80.39	31.16

<https://huggingface.co/spaces/mteb/leaderboard>

```
from langchain.document_loaders import *
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain.vectorstores import Chroma, Pinecone
from langchain.embeddings.openai import OpenAIEmbeddings
from langchain.embeddings.huggingface import HuggingFaceEmbeddings, HuggingFaceInstructEmbeddings,
from langchain.embeddings import HuggingFaceEmbeddings, SentenceTransformerEmbeddings
from langchain.chains.question_answering import load_qa_chain
from langchain.chat_models import ChatOpenAI
from langchain import PromptTemplate
```

```
loader = UnstructuredWordDocumentLoader('./中興施工規範範例檔案/00001-道碴軌道施工規範-1.0.doc')
data = loader.load()
```

```
text_splitter = RecursiveCharacterTextSplitter(chunk_size=500, chunk_overlap=200)
texts = text_splitter.split_documents(data)
print(f'Now you have {len(texts)} documents')
```

```
embeddings=HuggingFaceEmbeddings(model_name='shibing624/text2vec-base-chinese')
```

```
db = Chroma.from_documents(texts, embeddings)
retriever = db.as_retriever(search_type="similarity", search_kwargs={"k":2})
```

```
query = input("Question:")
docs = retriever.get_relevant_documents(query)
```

```
from langchain.document_loaders import *
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain.vectorstores import Chroma, Pinecone
from langchain.embeddings.openai import OpenAIEmbeddings
from langchain.embeddings.huggingface import HuggingFaceEmbeddings, HuggingFaceInstructEmbeddings,
from langchain.embeddings import HuggingFaceEmbeddings, SentenceTransformerEmbeddings
from langchain.chains.question_answering import load_qa_chain
from langchain.chat_models import ChatOpenAI
from langchain import PromptTemplate
```

```
loader = UnstructuredWordDocumentLoader(('. /中興施工規範範例檔案/00001-道碴軌道施工規範-1.0.doc'))
data = loader.load()
```

```
text_splitter = RecursiveCharacterTextSplitter(chunk_size=500, chunk_overlap=200)
texts = text_splitter.split_documents(data)
print (f'Now you have {len(texts)} documents')
```

```
embeddings=HuggingFaceEmbeddings(model_name='shibing624/text2vec-base-chinese')
```

```
db = Chroma.from_documents(texts, embeddings)
retriever = db.as_retriever(search_type="similarity", search_kwargs={"k":2})
```

```
query = input("Question:")
docs = retriever.get_relevant_documents(query)
```

```
db = Chroma.from_documents(texts, embeddings)
retriever = db.as_retriever(search_type="similarity", search_kwargs={"k":2})
```

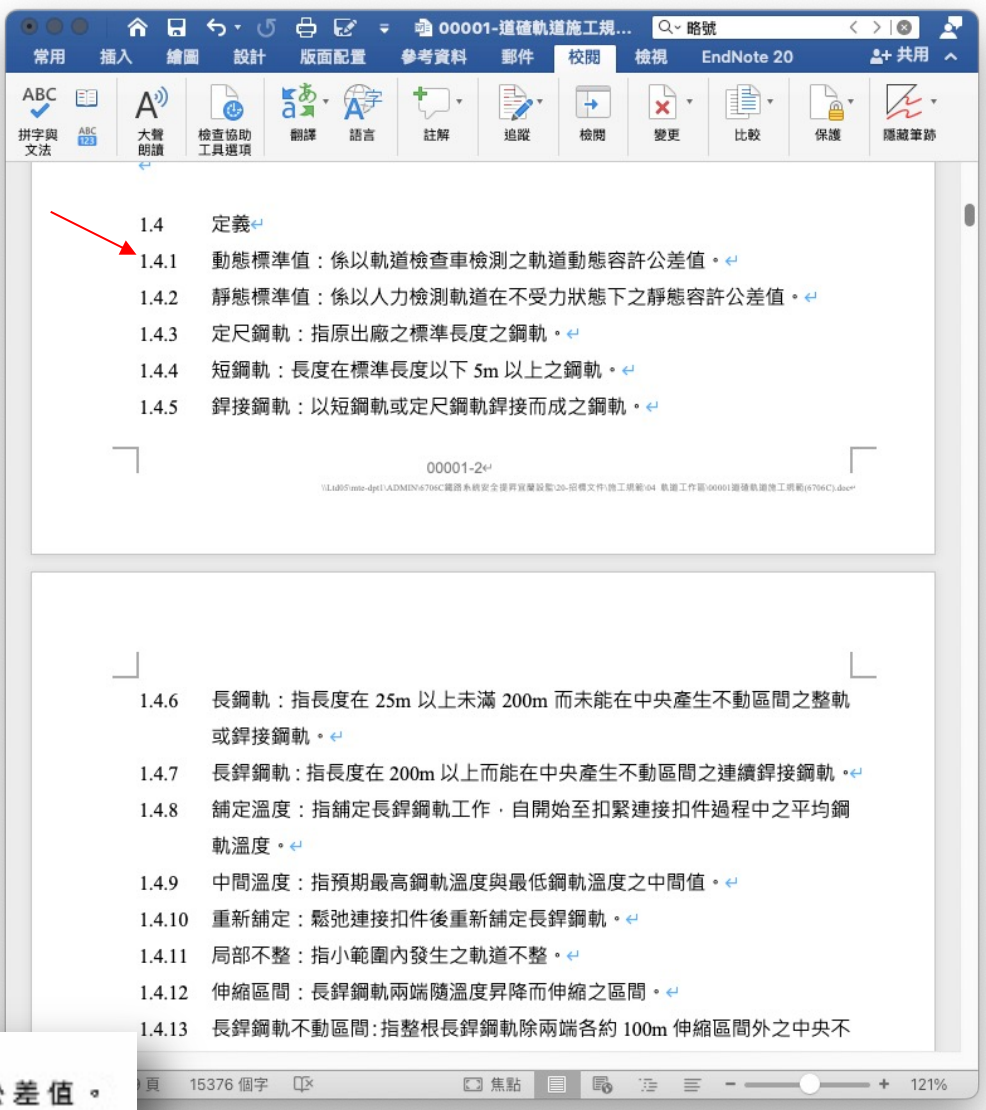
```
query = input("Question:")
docs = retriever.get_relevant_documents(query)
```

```
prompt_template = """請注意：請指根據本段輸入文字訊息的內容進行回答，如果query與提供的內容無關，請回答
“我不知道”，另外也不要回答無關答案：
Context: {context}
Question: {question}
Answer:
請以中文做問答"""
# """
PROMPT = PromptTemplate(template=prompt_template, input_variables=["context", "question"])

qa = load_qa_chain(ChatOpenAI(openai_api_key=OPENAI_API_KEY, temperature=0),
                  chain_type="stuff",
                  prompt=PROMPT,
                  verbose=True)

result = qa({"input_documents": docs, "question": query}, return_only_outputs=False)
print(result["output_text"])
```

```
embedding — python — python BookQA_engineer.py — 114x45
(pytorch-m1) hamigua:embedding cengqianyu$ python BookQA_engineer.py
Could not import azure.core python package.
convert /Users/cengqianyu/Google Drive - a/NTU/PHD/2022spring_AiSeminar/embedding/中興施工規範範例檔案/00001-道確軌道施工規範-1.0.doc -> /private/var/folders/by/vzsz3lx903j6drt3wzdxhj4m0000gn/T/tmp3wipqbid/00001-道確軌道施工規範-1.0.docx using filter : MS Word 2007 XML
You have 1 document(s) in your data
Now you have 352 documents
Question:動態標準值在哪些章節有提到？他的定義是什麼
```



> Finished chain.  
Answer:動態標準值在 1.4.1 定義章節有提到，其定義為以軌道檢查車檢測之軌道動態容許公差值。



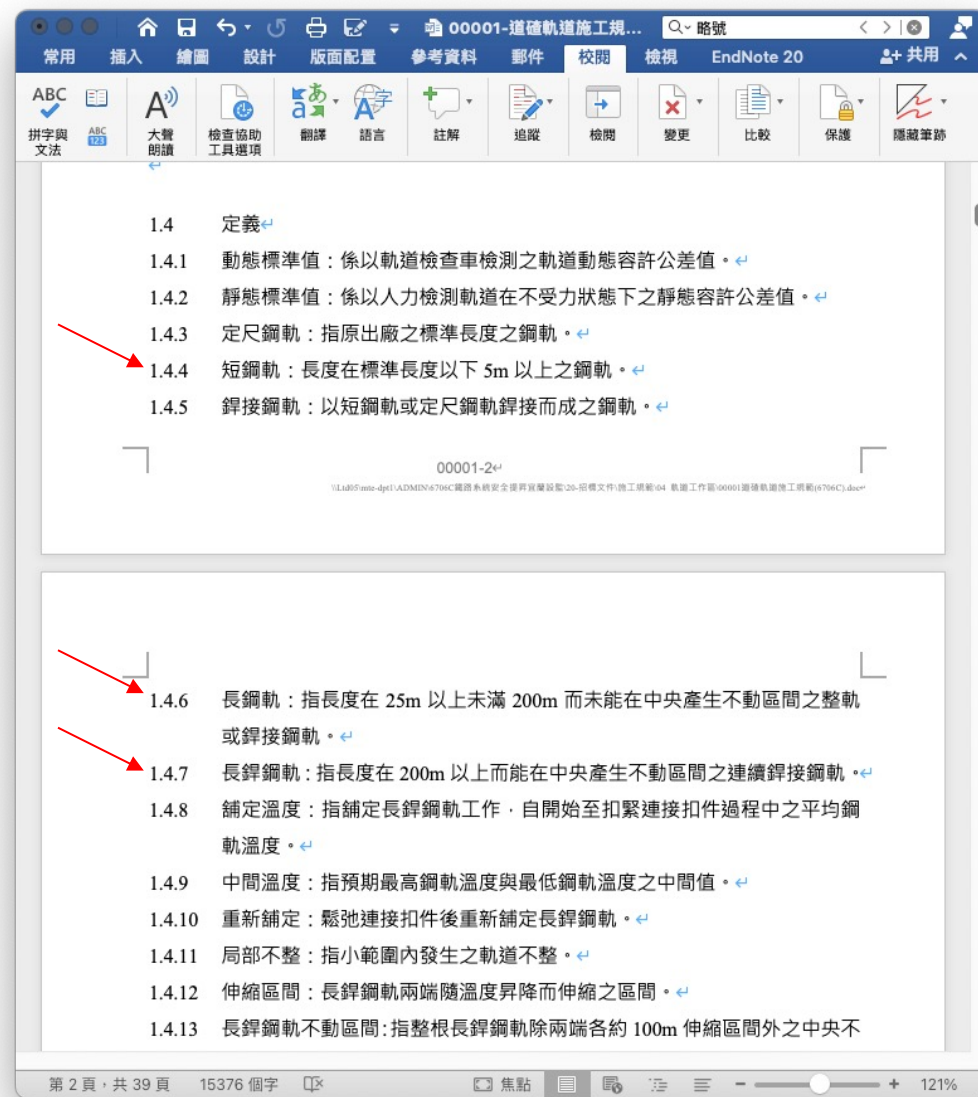
Question: 請分別告訴我短鋼軌、長鋼軌以及長銲鋼軌的定義是什麼？出現在哪幾章節

*text2vec-base-chinese*

Answer: **短鋼軌是指長度在標準長度以下5m以上之鋼軌，出現在1.4.4章節。**長鋼軌是指長度在25m以上未滿200m而未能在中央產生不動區間之整軌或銲接鋼軌，出現在1.4.6章節。長銲鋼軌是指長度在200m以上而能在中央產生不動區間之連續銲接鋼軌，出現在1.4.7章節。

*OpenAI- text-embedding-ada-002*

Answer: **短鋼軌是指定尺鋼軌或短斷鋼軌，銲接鋼軌是由短鋼軌或定尺鋼軌銲接而成，這些定義出現在第1章的1.4.5節。**長鋼軌是指長度在25m以上未滿200m而未能在中央產生不動區間之整軌或銲接鋼軌，長銲鋼軌是指長度在200m以上而能在中央產生不動區間之連續銲接鋼軌，這些定義出現在第1章的1.4.6節和1.4.7節。



# QA with youtube video

- Six behaviors to increase your confidence | Emily Jaenson | TEDxReno

```
loader = YoutubeLoader.from_youtube_url('https://www.youtube.com/watch?v=IitIl2C3Iy8')
```

ArxivLoader, DiscordChatLoader, EverNoteLoader, GoogleDriveLoader, NotionDBLoader...

Question: what are the six behaviors?

Answer:

1. Count Yourself In
2. Give Yourself 20 Seconds of Courage
3. Take a Seat at the Table
4. Cheer for Other People's Success
5. Bolster Your Confidence for a New Activity Through Your Already Great Performance in Another
6. Celebrate Constantly

# TIKTOKEN – token calculator

Overview Documentation API reference Examples Account >

## Tokenizer

The GPT family of models process text using **tokens**, which are common sequences of characters found in text. The models understand the statistical relationships between these tokens, and excel at producing the next token in a sequence of tokens.

You can use the tool below to understand how a piece of text would be tokenized by the API, and the total count of tokens in that piece of text.

**GPT-3** Codex

Many words map to one token, but some don't: indivisible.

Unicode characters like emojis may be split into many tokens containing the underlying bytes: 🍌

Sequences of characters commonly found next to each other may be grouped together: 1234567890

Clear Show example

Tokens	Characters
64	252

Many words map to one token, but some don't: indivisible.

Unicode characters like emojis may be split into many tokens containing the underlying bytes: 🍌

Sequences of characters commonly found next to each other may be grouped together: 1234567890

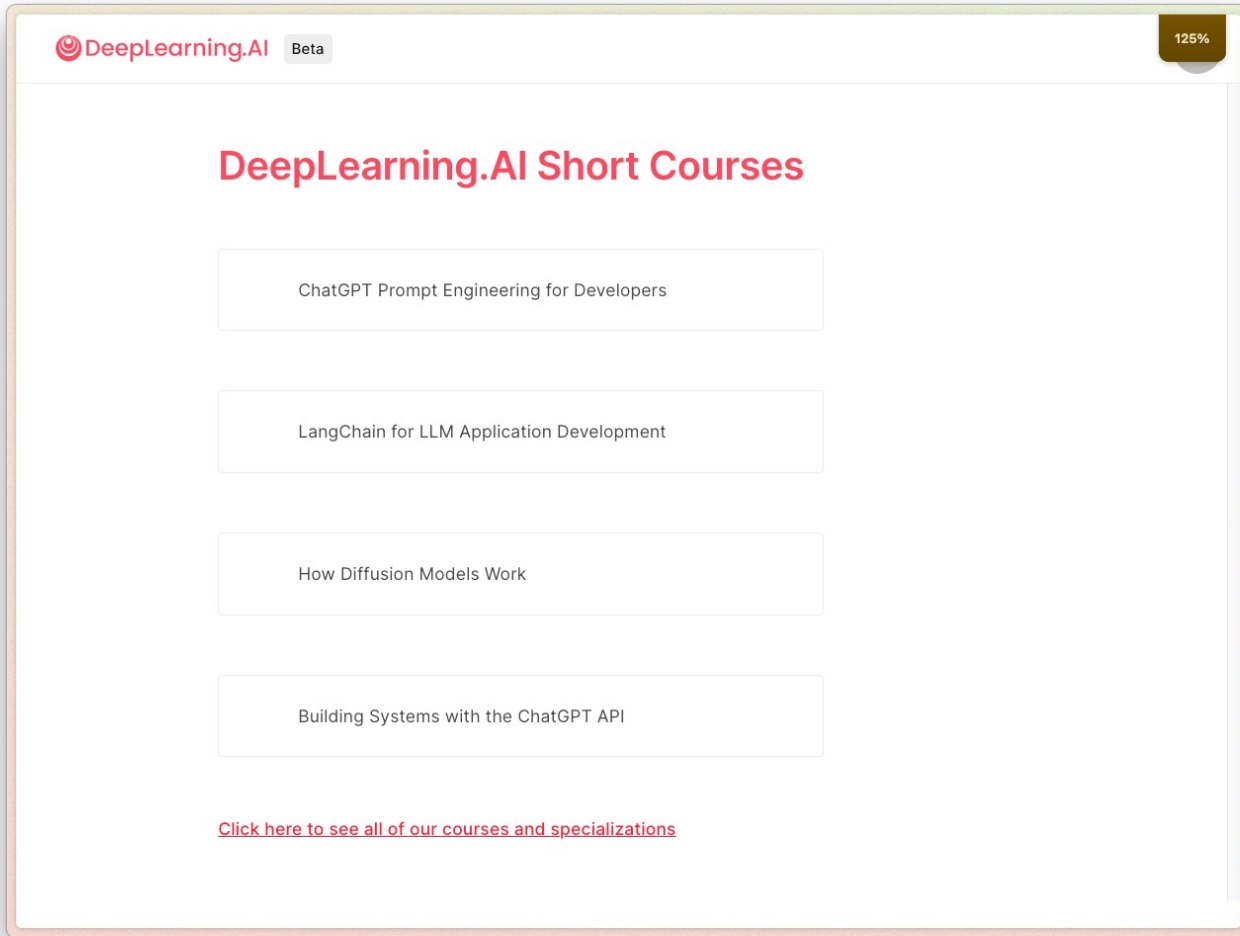
TEXT TOKEN IDS

```
MODEL_TO_ENCODING: dict[str, str] = {  
# chat  
"gpt-4": "cl100k_base",  
"gpt-3.5-turbo": "cl100k_base",  
# text  
"text-davinci-003": "p50k_base",  
"text-davinci-002": "p50k_base",  
"text-davinci-001": "r50k_base",  
"text-curie-001": "r50k_base",  
"text-babbage-001": "r50k_base",  
"text-ada-001": "r50k_base",  
"davinci": "r50k_base",  
"curie": "r50k_base",  
"babbage": "r50k_base",  
"ada": "r50k_base",  
}
```

```
for i in texts:  
t=i.page_content
```

```
encoding = tiktoken.get_encoding("cl100k_base")  
token_len = len(encoding.encode(t))
```

# Online Courses



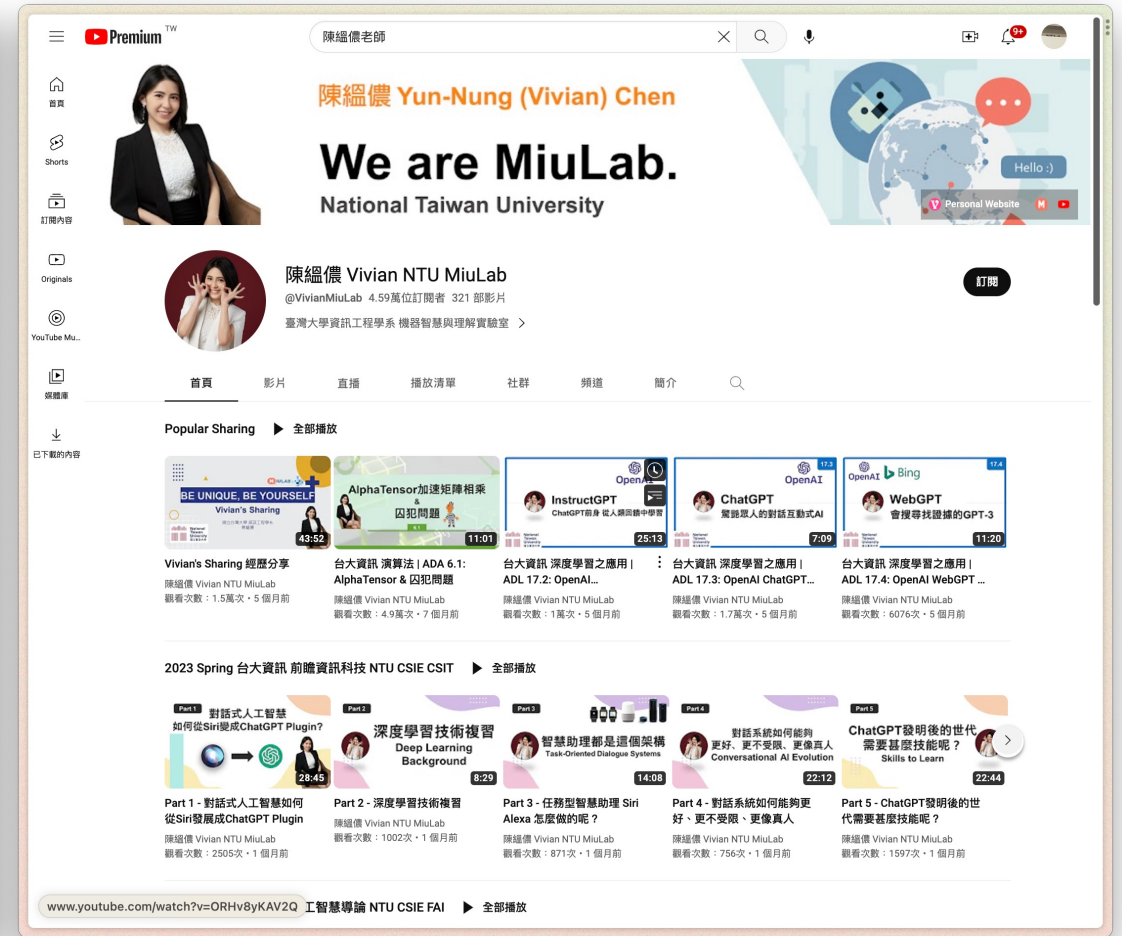
DeepLearning.AI Beta 125%

## DeepLearning.AI Short Courses

- ChatGPT Prompt Engineering for Developers
- LangChain for LLM Application Development
- How Diffusion Models Work
- Building Systems with the ChatGPT API

[Click here to see all of our courses and specializations](#)

<https://learn.deeplearning.ai/>



陳縉儂老師

陳縉儂 Yun-Nung (Vivian) Chen

## We are MiuLab.

National Taiwan University

陳縉儂 Vivian NTU MiuLab  
@VivianMiuLab 4.59萬位訂閱者 321 部影片  
臺灣大學資訊工程學系 機器智慧與理解實驗室

首頁 影片 直播 播放清單 社群 頻道 簡介

### Popular Sharing

- BE UNIQUE, BE YOURSELF
- AlphaTensor 演算法 | ADA 6.1: AlphaTensor & 囚犯問題
- InstructGPT
- ChatGPT
- WebGPT

### 2023 Spring 台大資訊 前瞻資訊科技 NTU CSIE CSIT

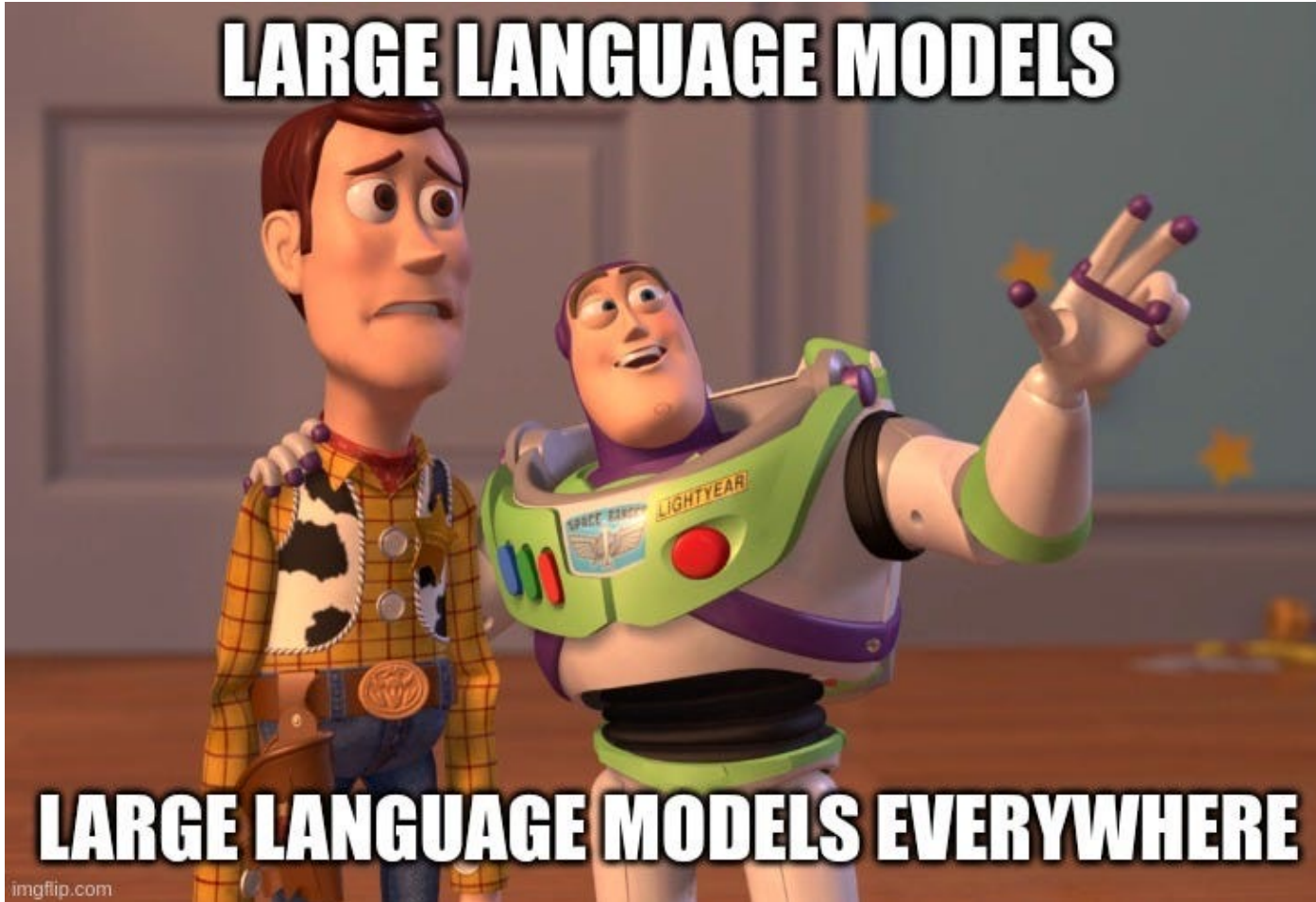
- Part 1 - 對話式人工智慧
- Part 2 - 深度學習技術複習
- Part 3 - 任務型智慧助理 Siri Alexa 怎麼做的呢?
- Part 4 - 對話系統如何能夠更好、更不受限、更像真人
- Part 5 - ChatGPT發明後的世代需要甚麼技能呢?

[www.youtube.com/watch?v=ORHv8yKAV2Q](https://www.youtube.com/watch?v=ORHv8yKAV2Q) 工智慧導論 NTU CSIE FAI

<https://www.youtube.com/@VivianMiuLab>



**LARGE LANGUAGE MODELS**



**LARGE LANGUAGE MODELS EVERYWHERE**

Hope you enjoy the sharing!