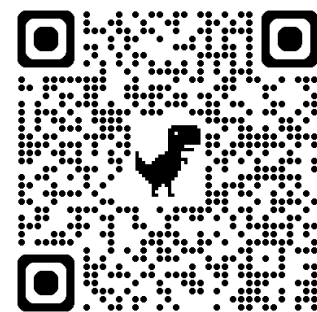# 自然語言處理與應用
# Natural Language Processing and Applications

## Mixture of Experts (MoE)

**Instructor:** 林英嘉 (Ying-Jia Lin)
**2025/06/01**

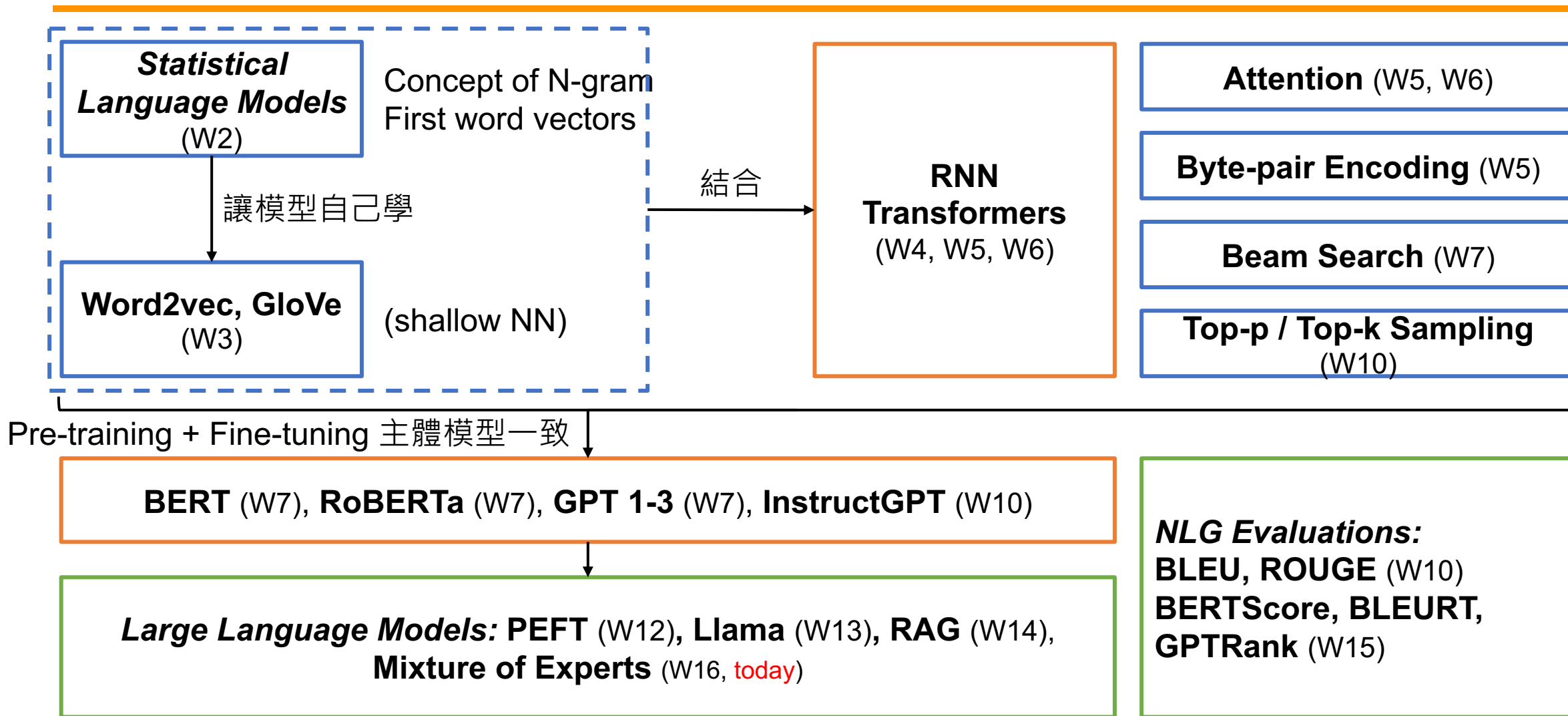Course GitHub
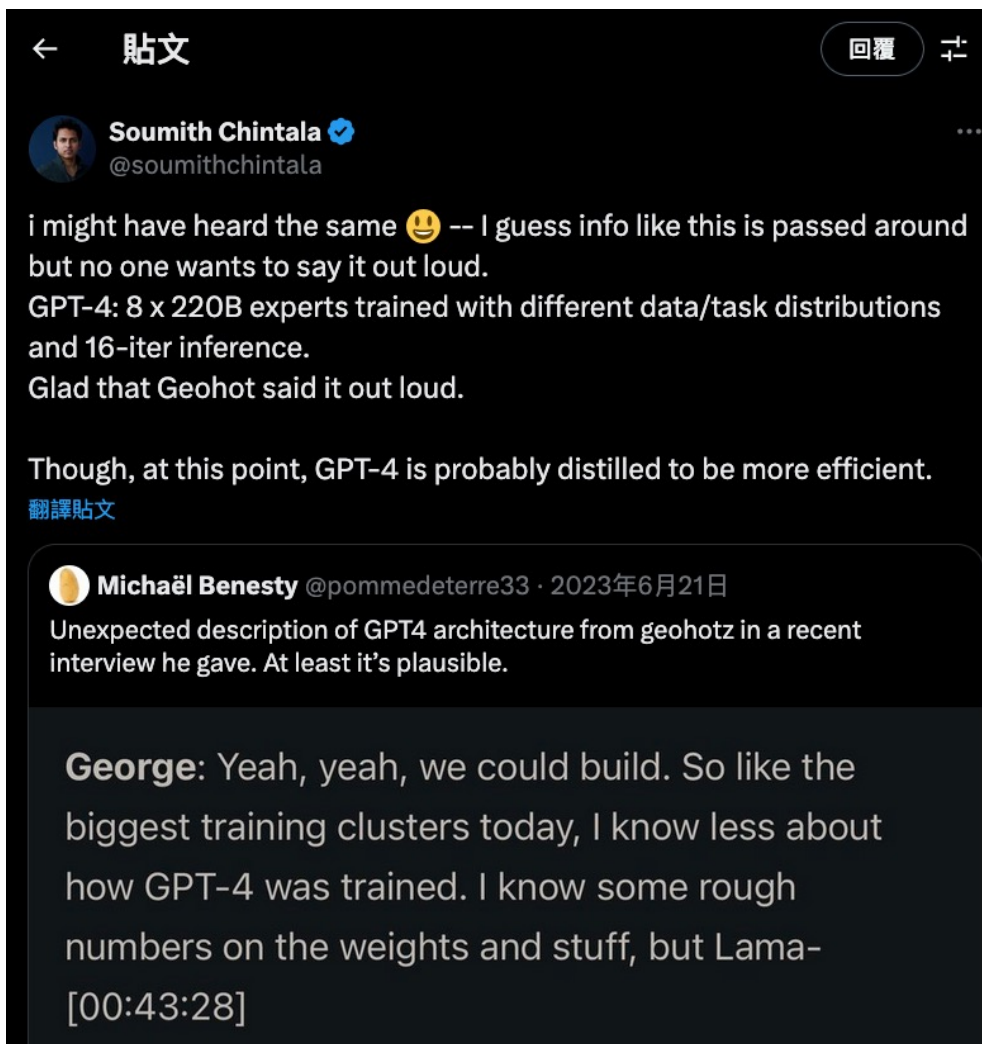
Slido # NLP_0601

# Outline

- 學期主題回顧

- Introduction to Mixture of Experts
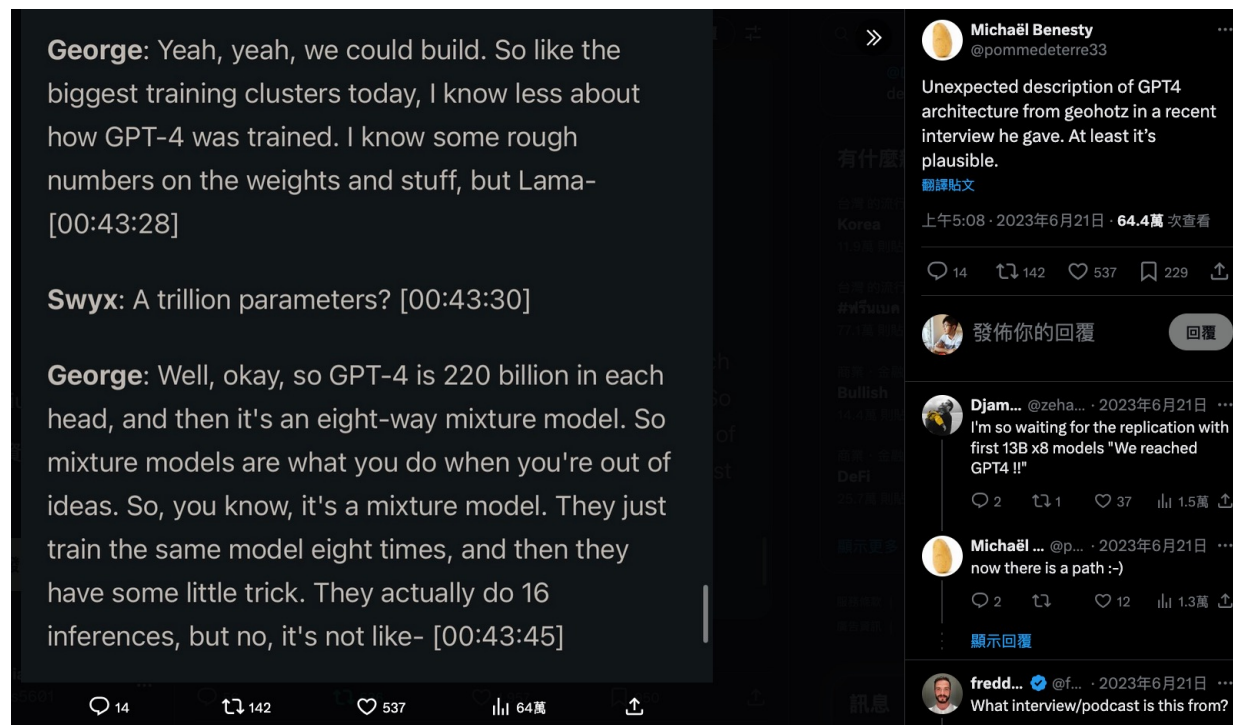
# 學期主題回顧

# Road Map of Natural Language Processing

**Statistical Language Models** (W2)

Concept of N-gram
First word vectors

讓模型自己學

**Word2vec, GloVe** (W3)

(shallow NN)

結合 →

**RNN Transformers** (W4, W5, W6)

**Attention** (W5, W6)

**Byte-pair Encoding** (W5)

**Beam Search** (W7)

**Top-p / Top-k Sampling** (W10)

Pre-training + Fine-tuning 主體模型一致

**BERT** (W7), **RoBERTa** (W7), **GPT 1-3** (W7), **InstructGPT** (W10)

*Large Language Models:* **PEFT** (W12), **Llama** (W13), **RAG** (W14), **Mixture of Experts** (W16, today)

*NLG Evaluations:*
**BLEU, ROUGE** (W10)
**BERTScore, BLEURT, GPTRank** (W15)

NLP

# Why do we need to learn MoE?



PyTorch 創始者推測 GPT-4 採用 MoE

https://x.com/soumithchintala/status/1671267150101721090

# Why do we need to learn MoE?

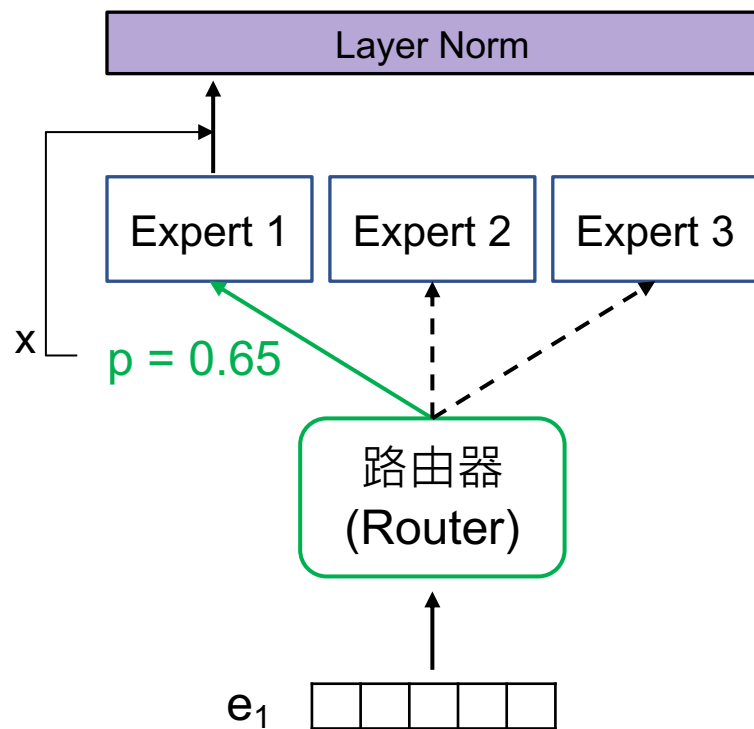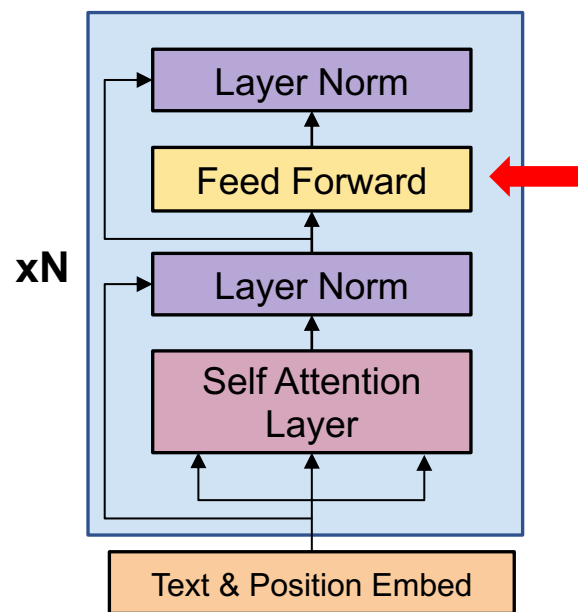| Model | Release Date | Active Parameters (Total Parameters) | Company |
|---|---|---|---|
| DeepSeek-V3 | 2024/12/27 | 37B (671B), 256 experts | DeepSeek |
| DeepSeek-R1 | 2025/1/22 | 37B (671B), 256 experts | DeepSeek |
| Llama 4 Maverick | 2025/4/5 | 17B (400B), 128 experts | Meta |
| Mixtral 8x7B | 2024/1/8 | 13B (47B), 8 experts | Mistral AI |

NLP

6

# MoE

# [Recap] Transformer block
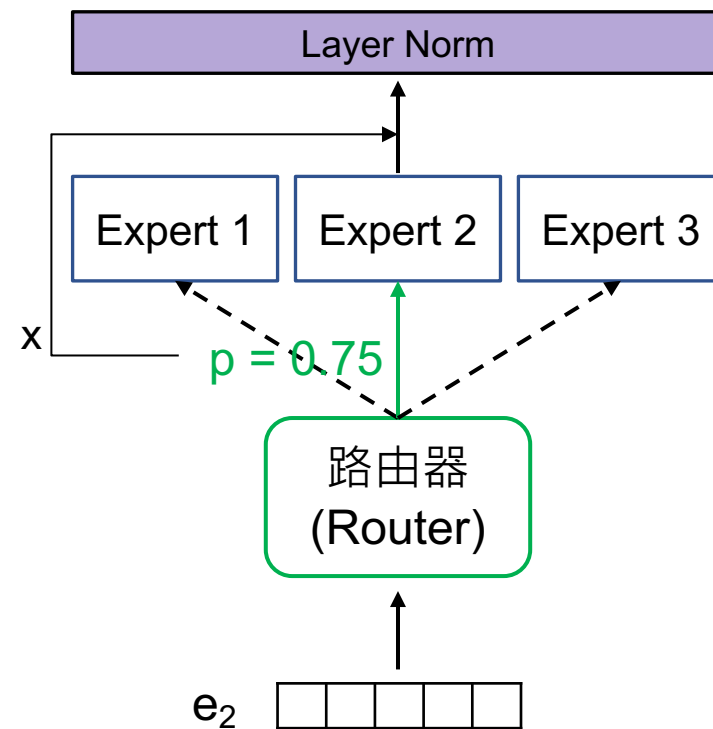
Transformer block

NLP

# Mixture of Experts (MoE)
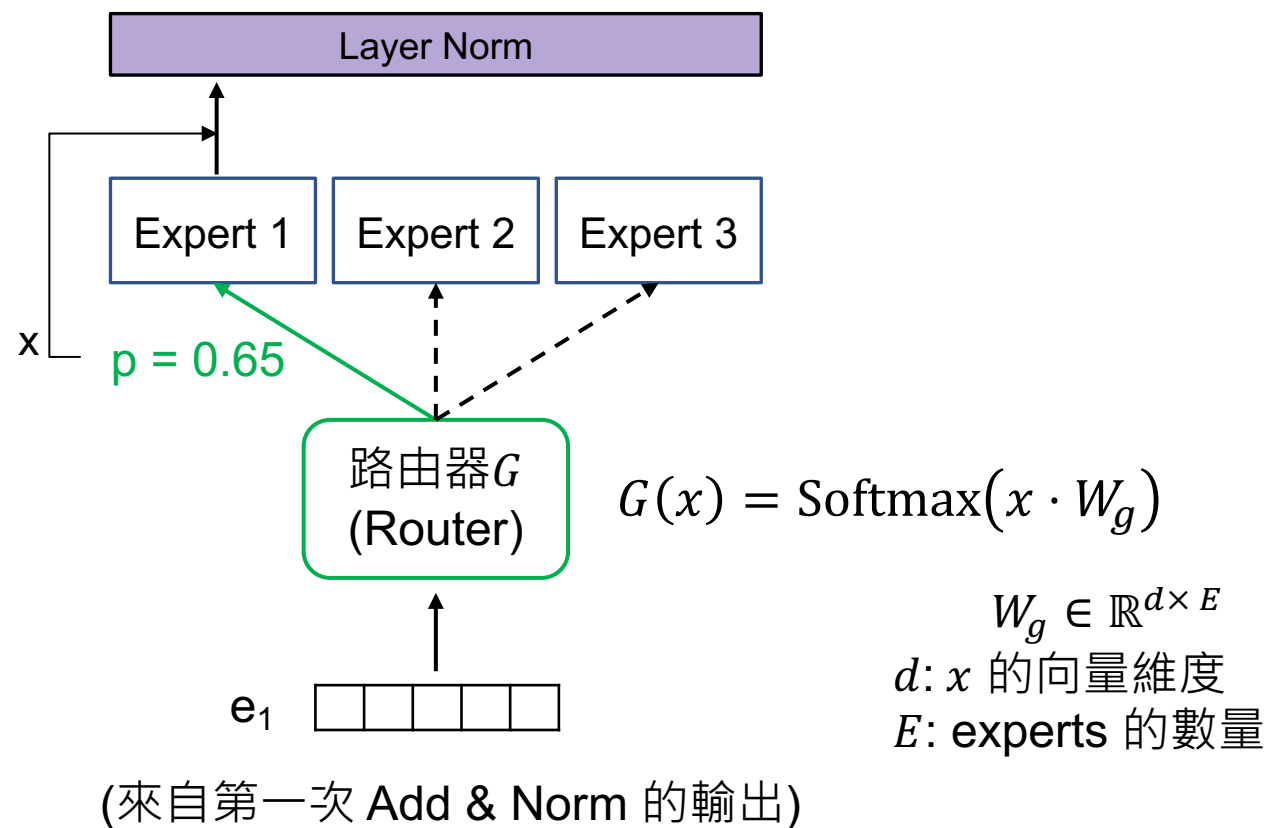
Transformer block (layer)



(來自第一次 Add & Norm 的輸出)　(來自第一次 Add & Norm 的輸出)

# Router

- Router 是一個小型神經網路，用來決定輸入token 該被送到哪一個或哪幾個 experts，又稱作 Gating Network



Layer Norm

Expert 1　Expert 2　Expert 3

x

p = 0.65

路由器$G$
(Router)

$$G(x) = \text{Softmax}(x \cdot W_g)$$

$$W_g \in \mathbb{R}^{d \times E}$$
$$d: x \text{ 的向量維度}$$
$$E: \text{experts 的數量}$$

$e_1$

(來自第一次 Add & Norm 的輸出)

# Expert 是什麼？

- 主要是看你放在哪裡，Expert 就會是一樣的架構

Transformer block (layer)



Experts: FFN FFN FFN     E.g., Switch Transformer, Mistral 8x7B, DeepSeek

Experts: Self Attn. Self Attn. Self Attn.     E.g., MoA (Zhang et al., 2022), SwitchHead (Csordás at al., 2024)

# Active Parameters (以Mistral 8x7B為例)



Feed Forward (for generation)

第32層 Transformer Layer
- Layer Norm
- Feed Forward
- Layer Norm
- Self Attention Layer

FFN FFN FFN FFN FFN FFN FFN FFN
Router

第1層 Transformer Layer
- Layer Norm
- Feed Forward
- Layer Norm
- Self Attention Layer

FFN FFN FFN FFN FFN FFN FFN FFN
Router

Text & Position Embed

原本是7B
外加了32層的(FFN+Router) 參數
一個FFN: 4096*14336*3 (SwiGLU)
一個Router: 4096*8 = 32768
層: (4096*14336*3+32768)*32約5.6B

每層由 Router 產生的數值選 Top-k個 experts (k=2)

Active parameters
Partial parameters

NLP

# Online resources

- [Stanford CS25: V4 I Demystifying Mixtral of Experts](#)

- [[IBM Technology] What is Mixture of Experts?](#)

- <span style="color:red">Important papers</span>

  - 近代 MoE 開山之作

    - [Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer (Shazeer et al., 2017)](#)

  - MoE on Transformers (T5)

    - [Switch Transformers (Fedus et al., 2021)](#)

NLP

# Thank you!

Instructor: 林英嘉

✉ yjlin@cgu.edu.tw