



# Where Can We Mix? From Atom to Cosmic

Paper Reading by Zhiying Lu

2024.11.19

智能多媒体内容计算实验室  
Intelligent Multimedia Content Computing Lab



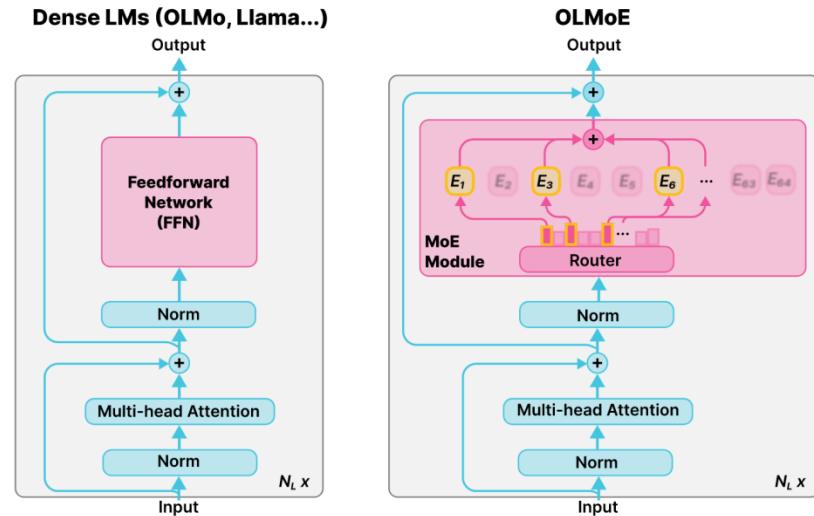
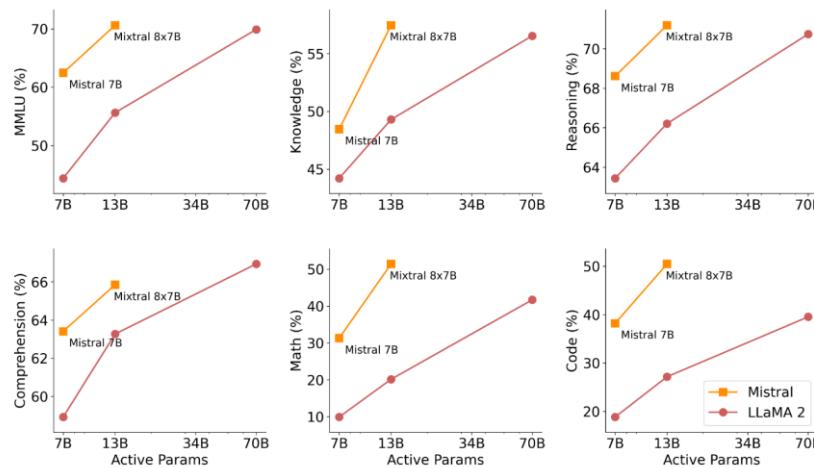
- Standard MoE
- Internal Mixture
- External Mixture
- A New Insight
- Summary



# Mixture-of-Expert (MoE)

3

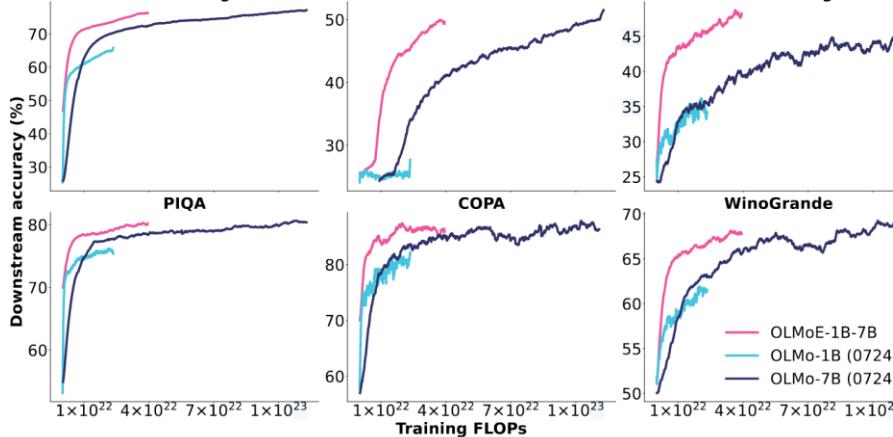
- 混合专家架构，一种将多个domain-specific expert融合起来构建统一模型的技术，常用赋予模型multi-task能力
- 在LLM结构中，FFN部分可以被替换为MoE架构，从Dense架构转换为Sparse结构，通过为每个token分配少量的K个expert，实现稀疏激活
- 扩展expert总数，控制激活的总参数量，构建轻量化高精度模型



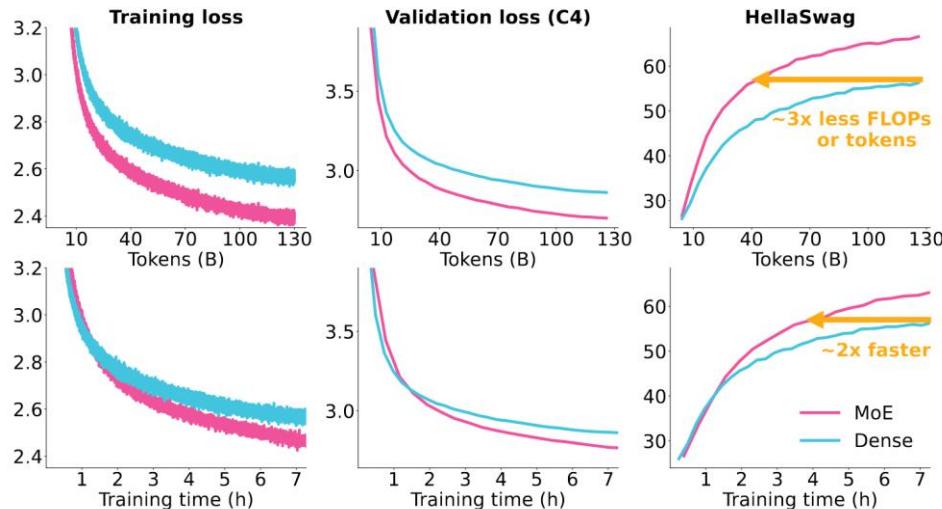
# MoE

4

- MoE的优势：
  - 更快的收敛速度
  - 更强的模型表征
  - 更丰富的泛化性
  - 更少的推理负担



	Active params	Open Data	MMLU	Hella-Swag	ARC-Chall.	ARC-Easy	PIQA	Wino-Grande
LMs with ~7-9B active parameters								
Llama2-7B [181]	6.7B	✗	46.2	78.9	54.2	84.0	77.5	71.7
OLMo-7B (0724) [64]	6.9B	✓	54.9	80.5	68.0	85.7	79.3	73.2
Mistral-7B [77]	7.3B	✗	64.0	83.0	78.6	90.8	82.8	77.9
DCLM-7B [89]	6.9B	✓	64.4	82.3	79.8	92.3	80.1	77.3
Llama3.1-8B [50]	8.0B	✗	66.9	81.6	79.5	91.7	81.1	76.6
Gemma2-9B [175]	9.2B	✗	<b>70.6</b>	<b>87.3</b>	<b>89.5</b>	<b>95.5</b>	<b>86.1</b>	<b>78.8</b>
LMs with ~2-3B active parameters								
OpenMoE-3B-9B [198]	2.9B	✓	27.4	44.4	29.3	50.6	63.3	51.9
StableLM-2B [16]	1.6B	✗	40.4	70.3	50.6	75.3	75.6	65.8
DeepSeek-3B-16B [39]	2.9B	✗	45.5	80.4	53.4	82.7	80.1	<b>73.2</b>
JetMoE-2B-9B [156]	2.2B	✗	49.1	<b>81.7</b>	61.4	81.9	80.3	70.7
Gemma2-3B [175]	2.6B	✗	53.3	74.6	67.5	84.3	78.5	71.8
Qwen1.5-3B-14B [178]	2.7B	✗	<b>62.4</b>	80.0	<b>77.4</b>	<b>91.6</b>	<b>81.0</b>	72.3
LMs with ~1B active parameters								
Pythia-1B [18]	1.1B	✓	31.1	48.0	31.4	63.4	68.9	52.7
OLMo-1B (0724) [64]	1.3B	✓	32.1	67.5	36.4	53.5	74.0	62.9
TinyLlama-1B [209]	1.1B	✓	33.6	60.8	38.1	69.5	71.7	60.1
DCLM-1B [89]	1.4B	✓	48.5	75.1	57.6	79.5	76.6	68.1
<b>OLMOE-1B-7B</b>	1.3B	✓	<b>54.1</b>	<b>80.0</b>	<b>62.1</b>	<b>84.2</b>	<b>79.8</b>	<b>70.2</b>



# MoE

5

- 其关键点在于：
- Expert -- What to be Mixed?
- Router -- How to Mix?

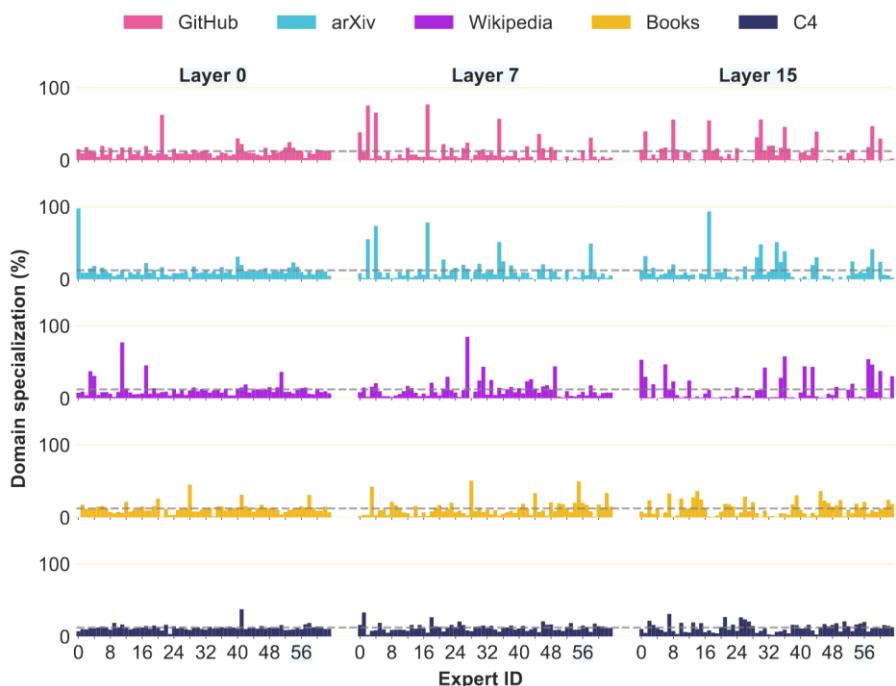
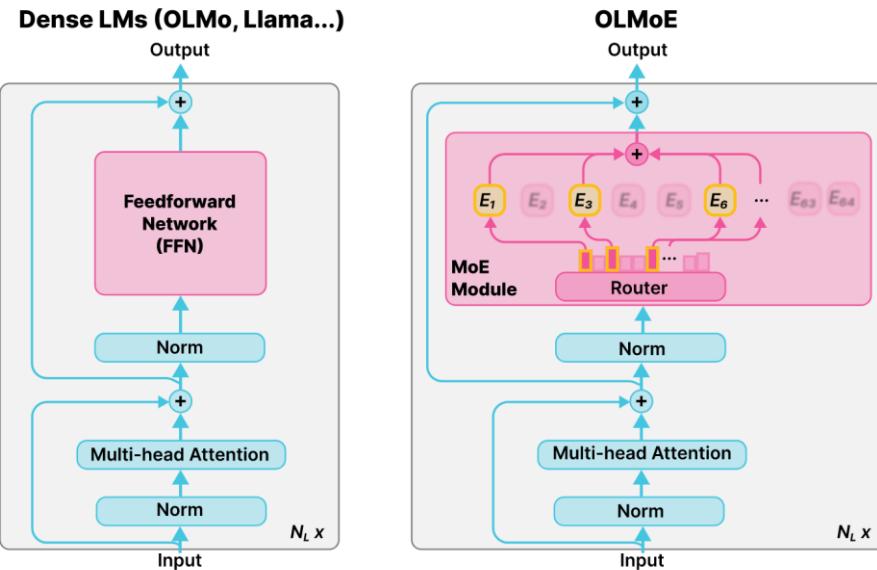
$$\sum_{i=0}^{n-1} G(x)_i \cdot E_i(x).$$

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

- 本报告专注于分享第三个问题：

Where to Mix?

- 单模态内
- 单模态->多模态
- 单任务->多任务





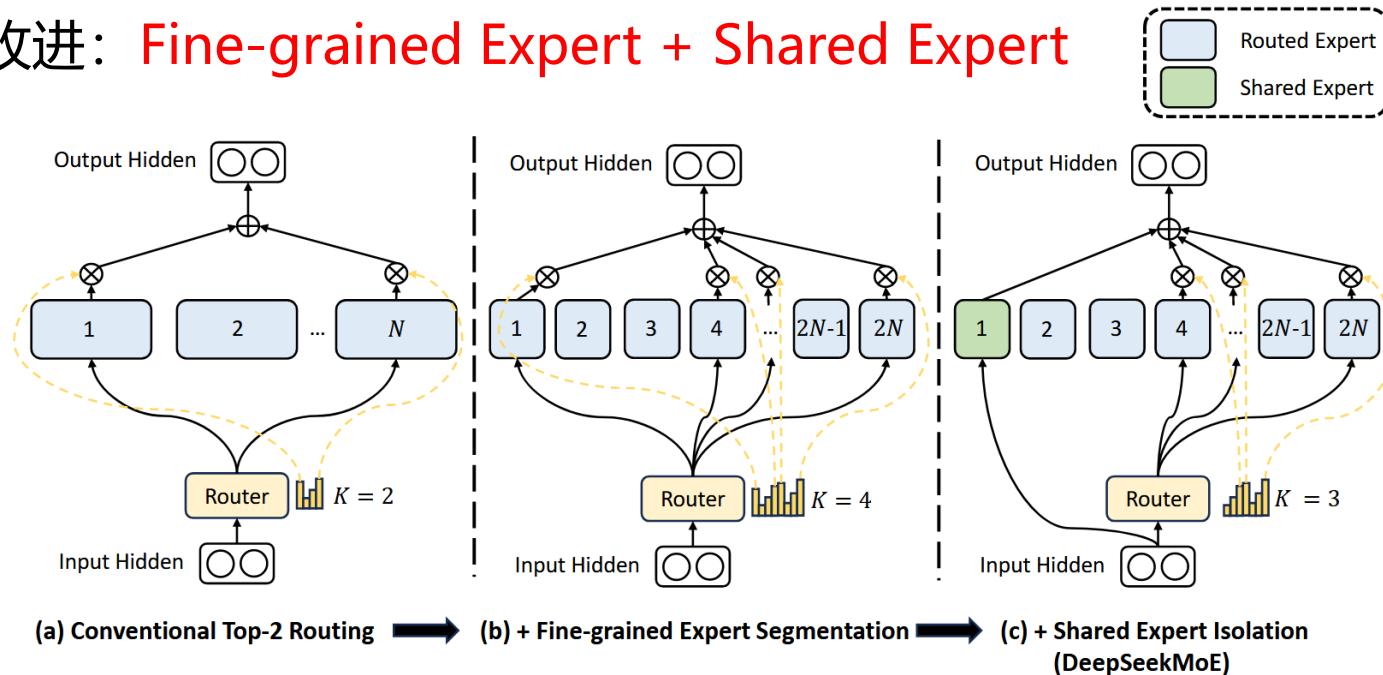
- Standard MoE
- Internal Mixture
- External Mixture
- A New Insight
- Summary

# DeepSeekMoE: Towards Ultimate Expert Specialization in Mixture-of-Experts Language Models



7

- 两种改进: Fine-grained Expert + Shared Expert



$$\mathbf{h}_t^l = \sum_{i=1}^{K_s} \text{FFN}_i(\mathbf{u}_t^l) + \sum_{i=K_s+1}^{mN} \left( g_{i,t} \text{FFN}_i(\mathbf{u}_t^l) \right) + \mathbf{u}_t^l,$$

$$g_{i,t} = \begin{cases} s_{i,t}, & s_{i,t} \in \text{Topk}(\{s_{j,t} | K_s + 1 \leq j \leq mN\}, mK - K_s), \\ 0, & \text{otherwise,} \end{cases}$$

$$s_{i,t} = \text{Softmax}_i \left( \mathbf{u}_t^{lT} \mathbf{e}_i^l \right).$$

- 固定activate params, 增加组合数

$$\binom{16}{2} = 120$$

$$\binom{64}{8} = 4,426,165,368$$

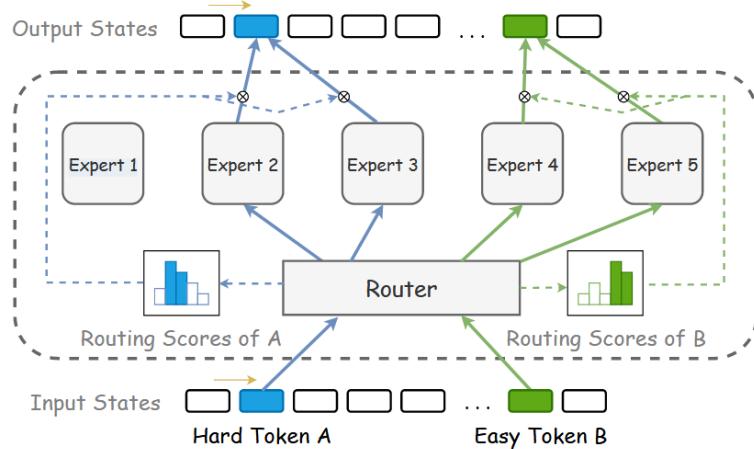
实验室  
Computing Lab

# HMoE: Heterogeneous Mixture of Experts for Language Modeling

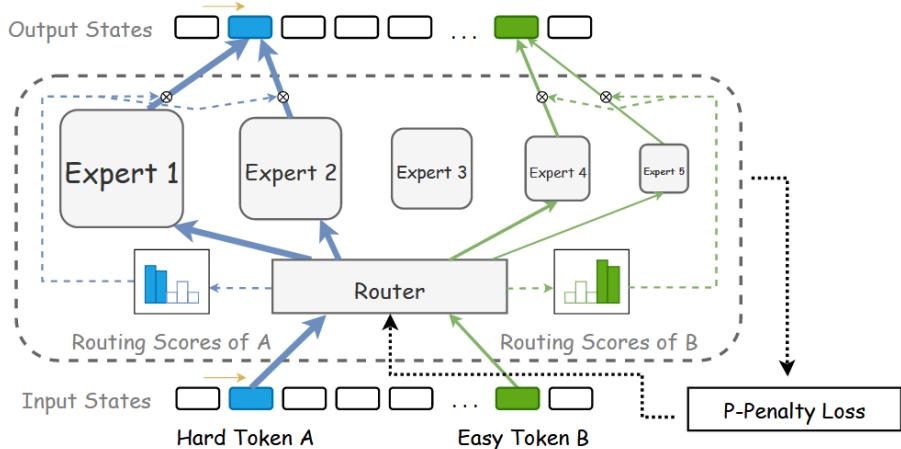


8

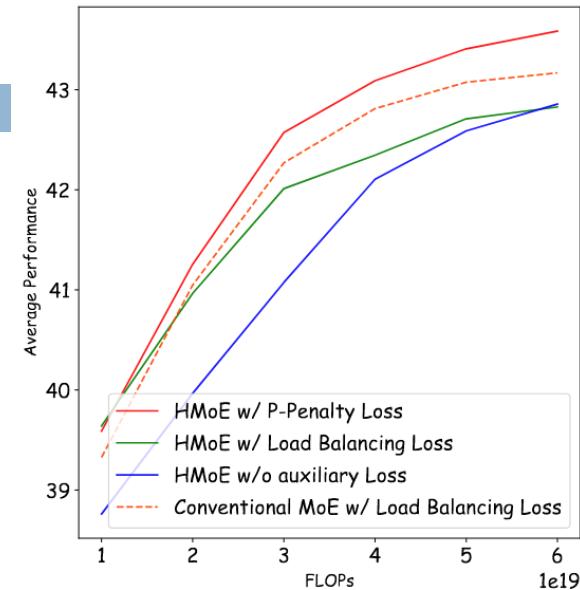
- 传统MoE模型需要控制每个expert的负载均衡，即分配给每个expert的token数量接近，因此需要每个expert的参数量相等
- 该工作设置了参数相关的惩罚性loss，使得expert的大小可以不相等



(a) Conventional homogenous MoE.



(b) Our proposed heterogeneous MoE.



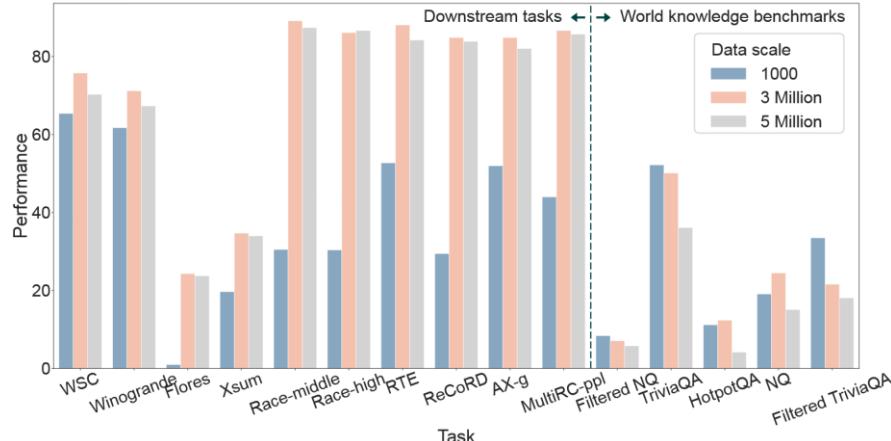
# LoRAMoE

9

- 在SFT阶段，用MoE方式结合多个Lora，防止通用知识的遗忘

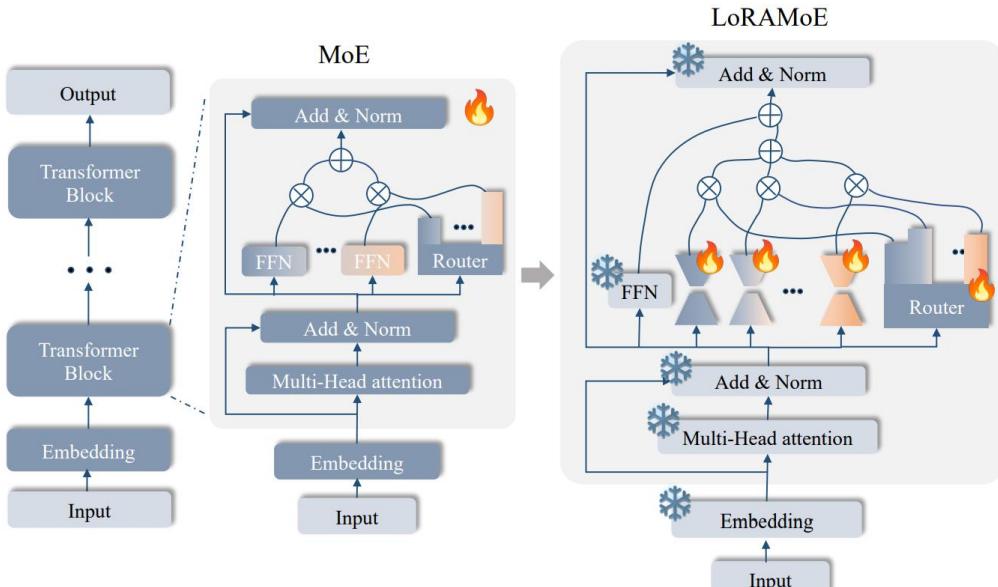
$$o = W_0x + \gamma \sum_{i=1}^N G(x)_i E_i(x)$$

$$= W_0x + \frac{\alpha}{r} \sum_{i=1}^N \omega_i \cdot B_i A_i x$$



## LoRAMoE: Alleviating World Knowledge Forgetting in Large Language Models via MoE-Style Plugin

Shihan Dou<sup>1\*</sup>, Enyu Zhou<sup>1\*</sup>, Yan Liu<sup>1</sup>, Songyang Gao<sup>1</sup>, Wei Shen<sup>1</sup>, Limao Xiong<sup>1</sup>,  
Yuhao Zhou<sup>1</sup>, Xiao Wang<sup>1</sup>, Zhiheng Xi<sup>1</sup>, Xiaoran Fan<sup>1</sup>, Shiliang Pu<sup>5</sup>, Jiang Zhu<sup>5</sup>,  
Rui Zheng<sup>1</sup>, Tao Gui<sup>2†</sup>, Qi Zhang<sup>1,3†</sup>, Xuanjing Huang<sup>1,4†</sup>



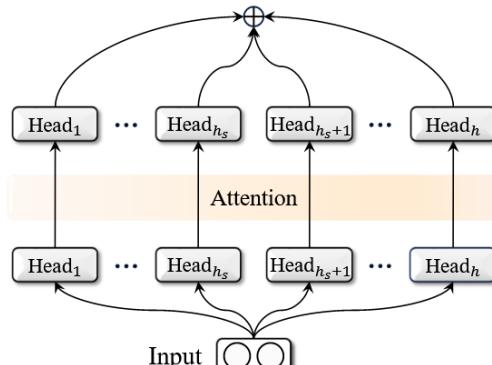
Task	baseline	SFT solely on CBQA	SFT	LoRA	LoRAMoE	LoRAMoE (with $\mathcal{L}_{lbc}$ )
<b>WSC</b>	65.4	-	<b>76.0</b>	65.4	71.2	70.2
<b>winogrande</b>	61.7	-	<b>71.2</b>	64.3	66.3	69.6
<b>Flores</b>	0.1	-	24.3	<b>26.6</b>	26.4	25.9
<b>Xsum</b>	19.7	-	34.7	34.5	<b>34.8</b>	33.2
<b>Race-middle</b>	30.5	-	89.1	78.8	84.5	<b>90.0</b>
<b>Race-high</b>	30.4	-	86.1	75.3	80.6	<b>86.5</b>
<b>RTE</b>	52.7	-	<b>88.1</b>	77.3	80.9	87.4
<b>ReCoRD</b>	29.4	-	84.8	83.2	84.3	<b>85.9</b>
<b>AX-g</b>	52.0	-	84.8	76.1	81.7	<b>87.1</b>
<b>multiRC</b>	44.0	-	86.7	81.4	87.3	<b>87.9</b>
<b>TriviaQA</b>	52.2	57.8	51.1	47.8	55.3	<b>58.1</b>
<b>NQ</b>	18.5	28.6	24.5	16.2	23.8	<b>28.0</b>
<b>Filtered TriviaQA</b>	33.5	36.2	21.6	33.4	<b>38.5</b>	35.4
<b>Filtered NQ</b>	7.8	12.8	7.3	11.6	<b>13.4</b>	12.0
<b>hotpot QA</b>	11.2	16.1	13.4	10.7	14.4	<b>16.1</b>

# MoH

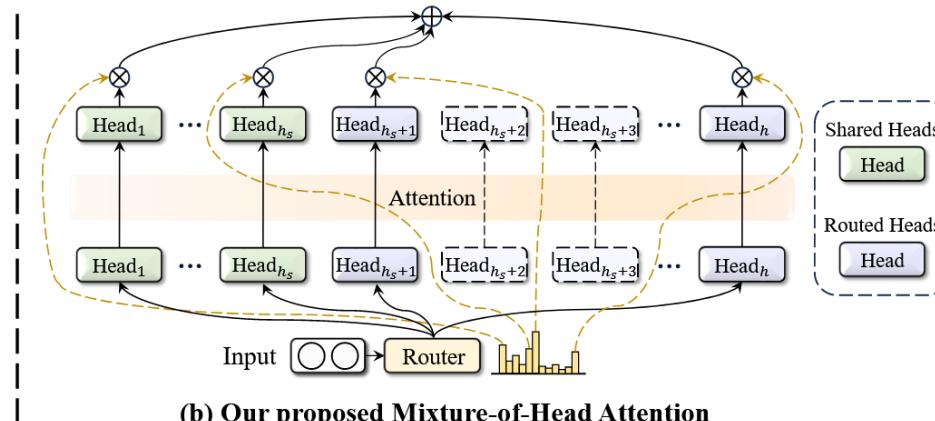
10

- MHSA本身也是一个MoE的形式，因此也可以改进成sparse结构
- 激活的比例比FFN MoE要大

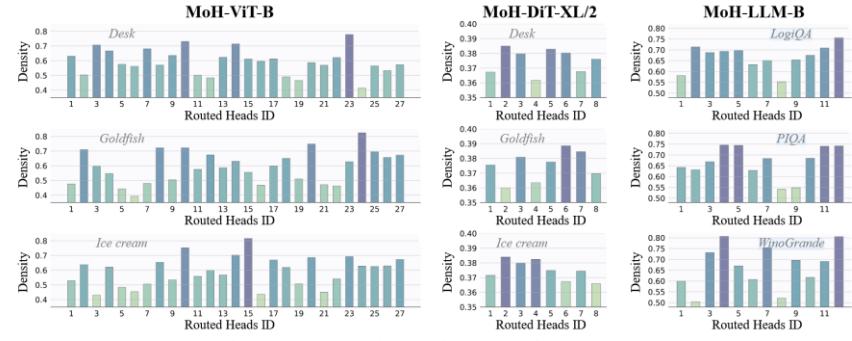
$$\text{MultiHead}(\mathbf{X}, \mathbf{X}') = \text{Concat}(\mathbf{H}^1, \mathbf{H}^2, \dots, \mathbf{H}^h) \mathbf{W}_O, \\ \mathbf{H}^i = \text{Attention}(\mathbf{X} \mathbf{W}_Q^i, \mathbf{X}' \mathbf{W}_K^i, \mathbf{X}' \mathbf{W}_V^i),$$



(a) Multi-Head Attention



(b) Our proposed Mixture-of-Head Attention



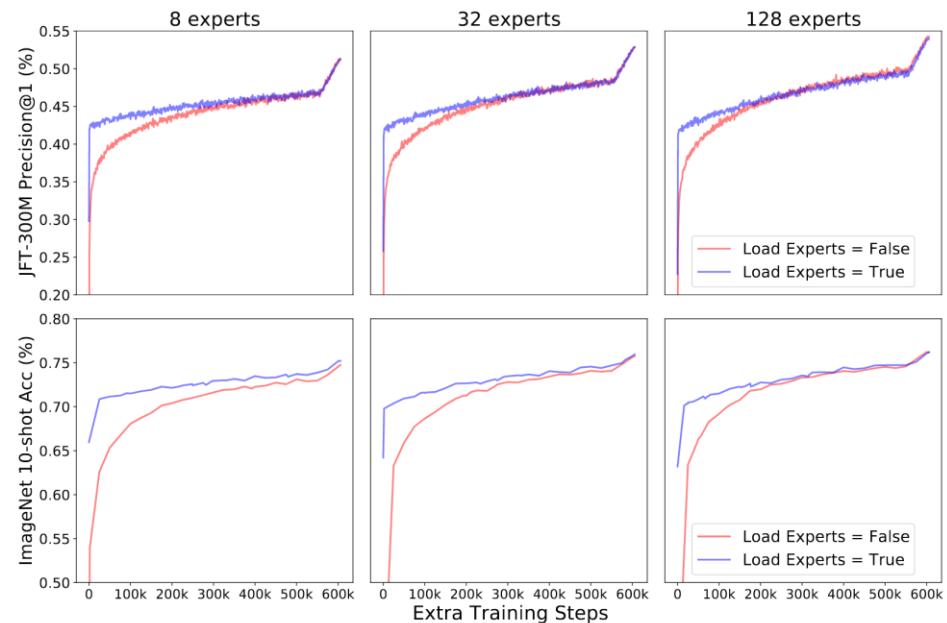
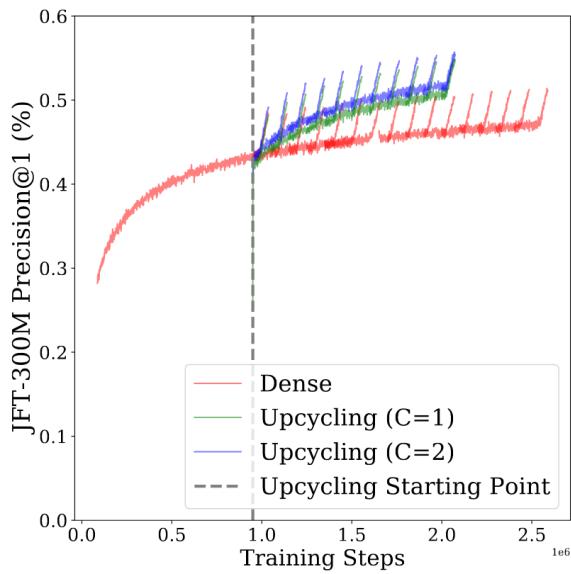
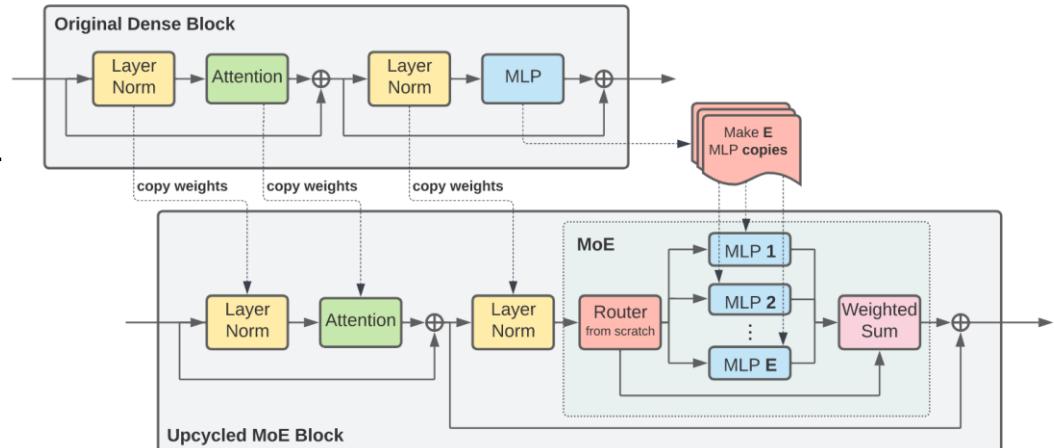
Methods	#Activated Heads (%)	MMLU (5)	CEVAL (5)	CMMLU (5)	GSM8K(8)	TruthfulQA
LLaMA3-8B (Dubey et al., 2024)	100	65.2	52.3	50.7	49.5	35.4
<b>MoH-LLaMA3-8B</b>	75	65.8	61.5	64.4	56.9	44.0
Methods	#Activated Heads (%)	HellaSwag (10)	LogiQA	BoolQ (32)	LAMBADA	SciQ
LLaMA3-8B (Dubey et al., 2024)	100	81.9	30.0	83.9	75.5	94.0
<b>MoH-LLaMA3-8B</b>	75	80.1	30.3	84.0	76.4	92.2
Methods	#Activated Heads (%)	PIQA	WinoGrande	NQ (32)	ARC-C (25)	Average
LLaMA3-8B (Dubey et al., 2024)	100	81.0	72.5	31.5	59.0	61.6
<b>MoH-LLaMA3-8B</b>	75	78.8	72.9	28.3	60.1	<b>64.0</b>

$$\text{MultiHead}(\mathbf{X}, \mathbf{X}') = \sum_{i=1}^h \mathbf{H}^i \mathbf{W}_O^i. \quad \text{MoH}(\mathbf{X}, \mathbf{X}') = \sum_{i=1}^h g_i \mathbf{H}^i \mathbf{W}_O^i,$$

# Sparse Upcycling

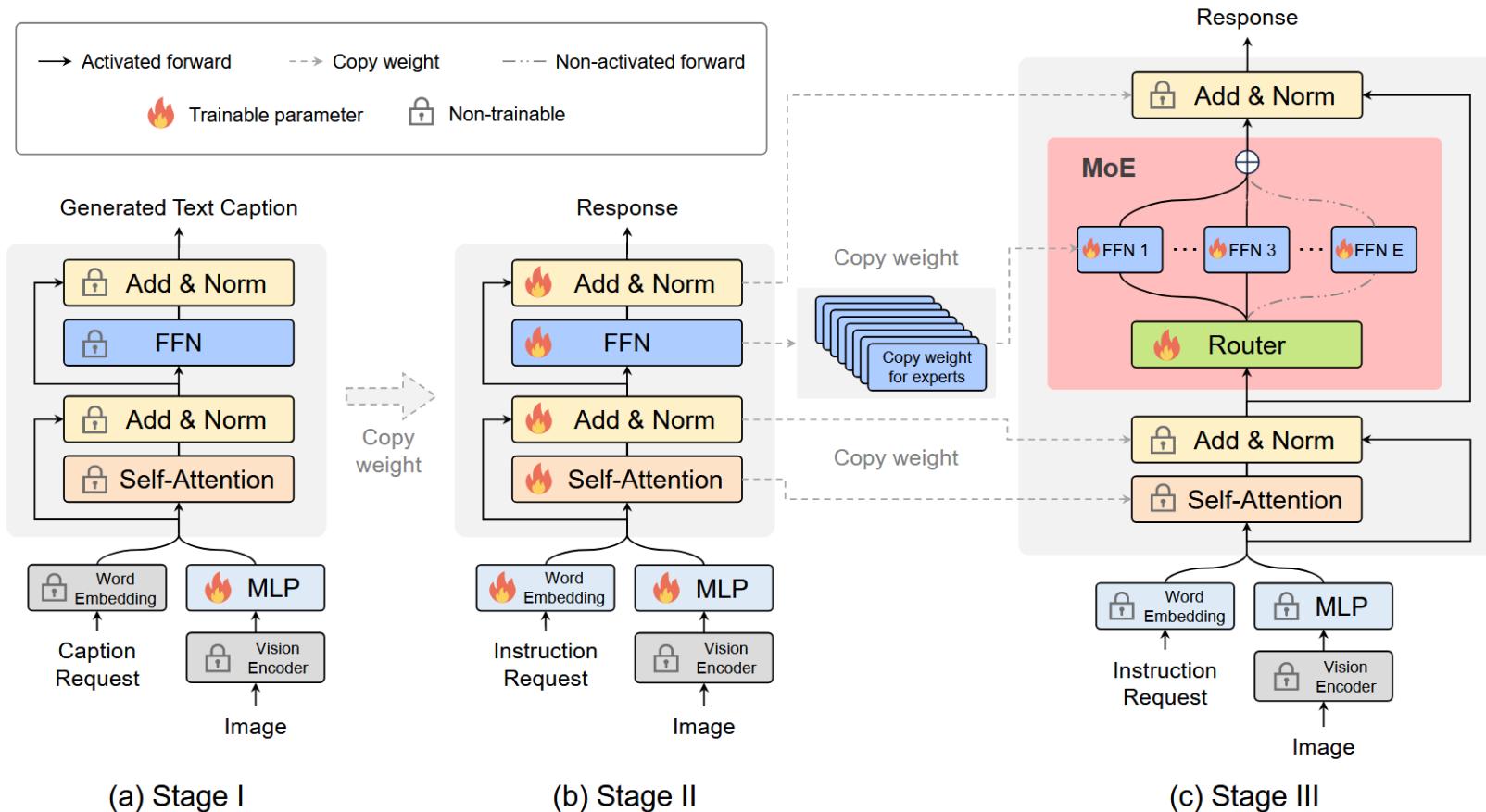
11

- 在初期训练之后将dense模型改造成MoE结构再进行训练
- 直接将MLP参数复制多份，并引入router进行组合
- 效果比直接MoE训练要好



# MoE-LLaVA

12

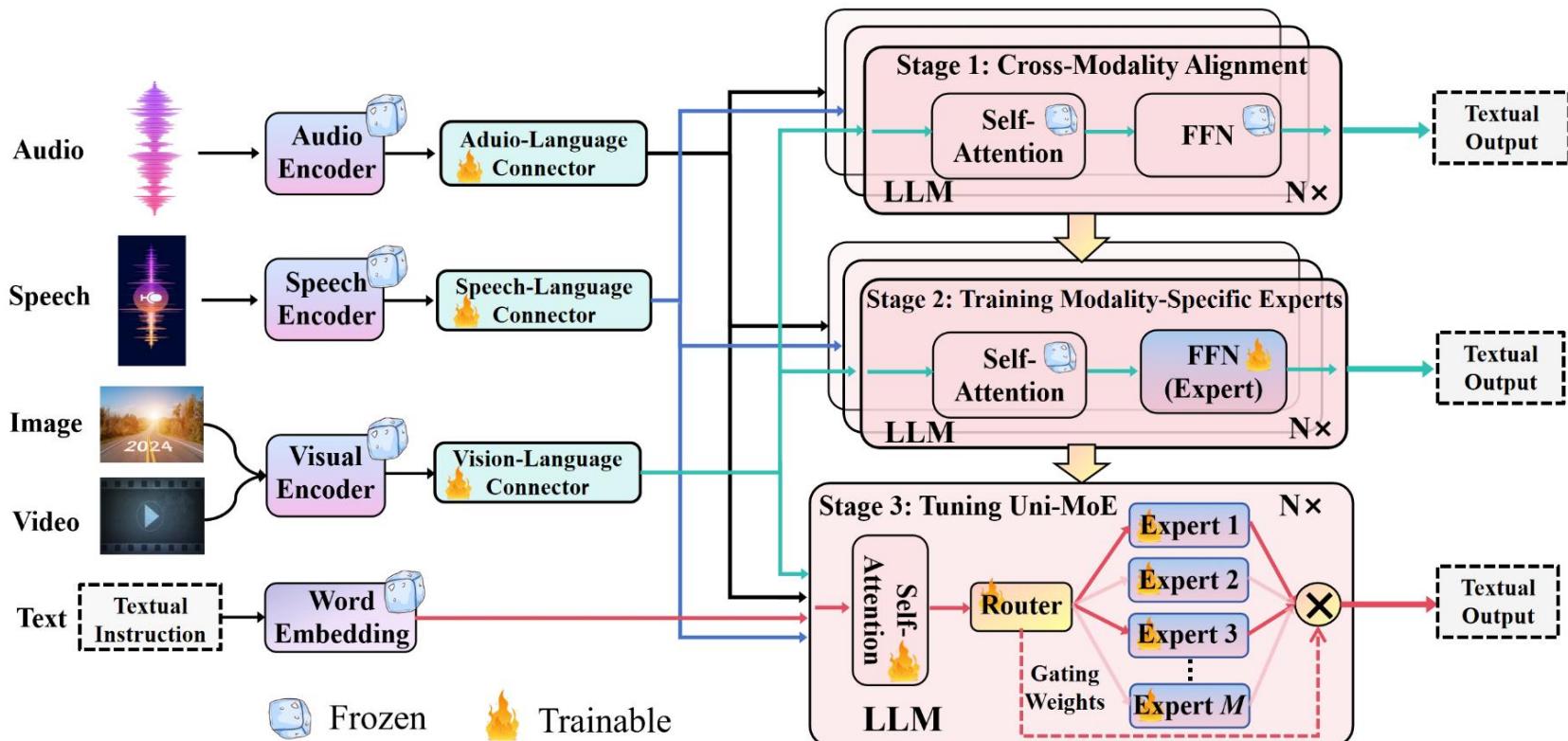




- Standard MoE
- Internal Mixture
- External Mixture
- A New Insight
- Summary

# Uni-MoE

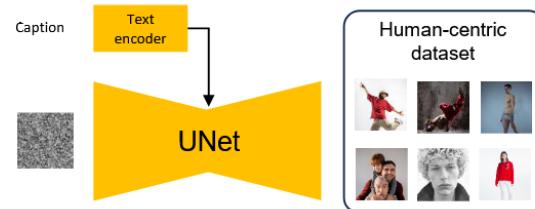
14



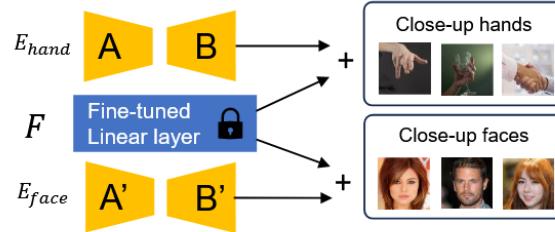
# MoLE

15

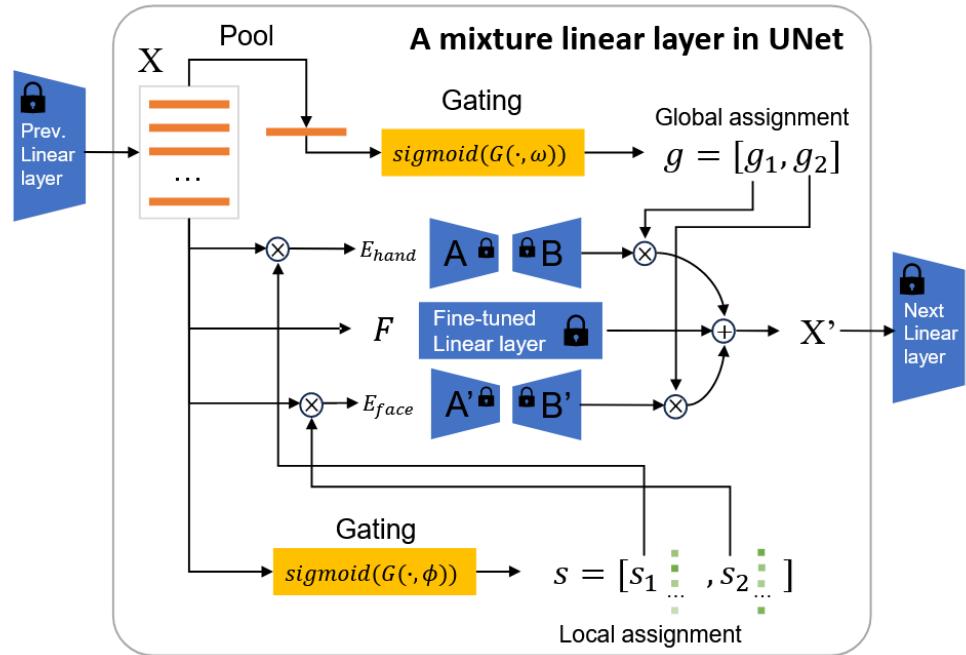
## Stage 1: Fine-tuning on Human-centric dataset



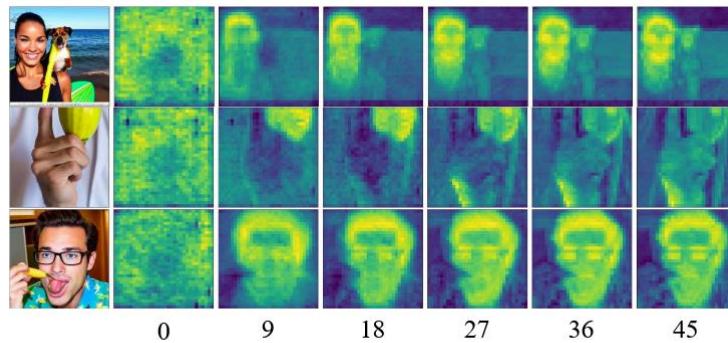
## Stage 2: Low-rank Expert Generation



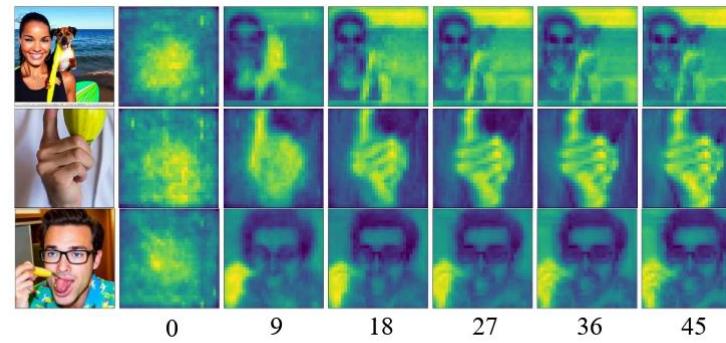
## Stage 3: Soft Mixture Assignment



Face Expert



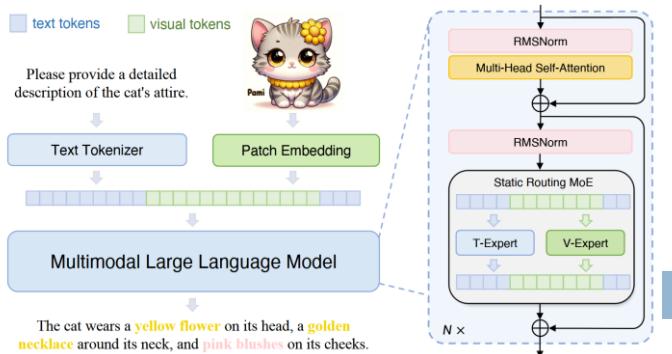
Hand Expert



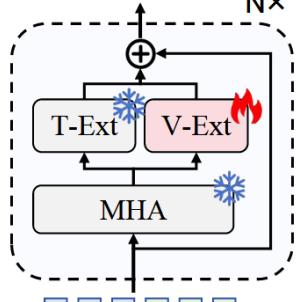
(d) Visualization of score map from face expert (left) and hand expert (right) in different inference step.

# Mono-InternVL

16

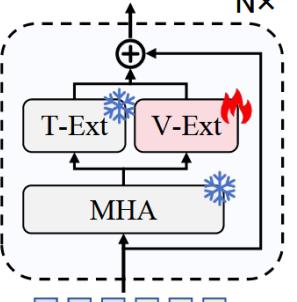


## S1.1: Concept learning



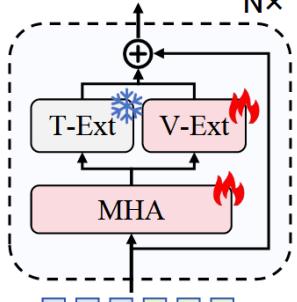
Noisy Image-Text Pairs  
922 Millions

## S1.2: Semantic learning



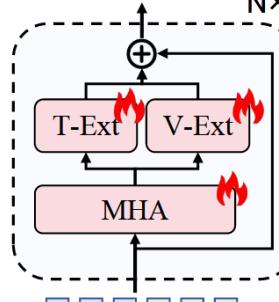
Synthetic Image-Text Pairs  
258 Millions

## S1.3: Alignment learning



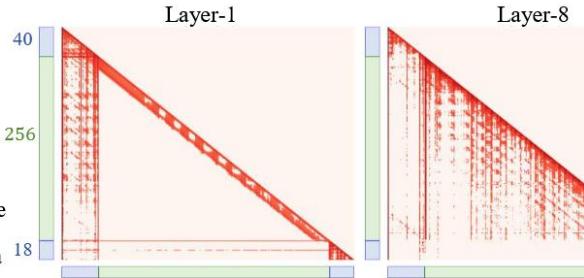
Task-related Data  
143 Millions

## S2: Instruction tuning



Instruction Data  
7 Millions

Resolution:  $600 \times 600$



HallB CCB | Avg

—	—	—
—	—	—
—	—	—
—	—	—
27.6	37.6	43.4
32.2	29.6	46.7
36.2	41.4	44.6
36.1	45.3	51.2
37.5	63.5	54.4

不

构

Q: Provide a one-sentence caption for the image.  
A: A wooden chair with a woven seat and back.

Resolution:  $1024 \times 768$



Q: what are the words in white?  
A: The words in white are "AI Weiwei".

Layer-1

Layer-16

Layer-24

Layer-1

Layer-16

Layer-24

Layer-8

Layer-8

Layer-24

Layer-24

—	—	—
—	—	—
17.1	3.5	16.1
21.1	12.4	34.8
26.4	16.3	38.9
—	—	—
34.8	66.3	55.2

# MoMa: Efficient Early-Fusion Pre-training with Mixture of Modality-Aware Experts

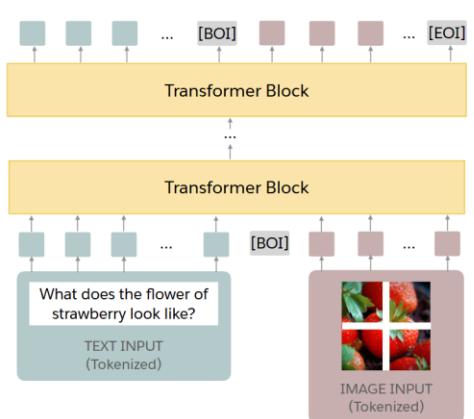
Xi Victoria Lin<sup>1,\*</sup>, Akshat Shrivastava<sup>1,\*</sup>, Liang Luo<sup>1</sup>, Srinivasan Iyer<sup>1</sup>, Mike Lewis<sup>1</sup>, Gargi Ghosh<sup>1</sup>, Luke Zettlemoyer<sup>1</sup>, Armen Aghajanyan<sup>1,\*</sup>



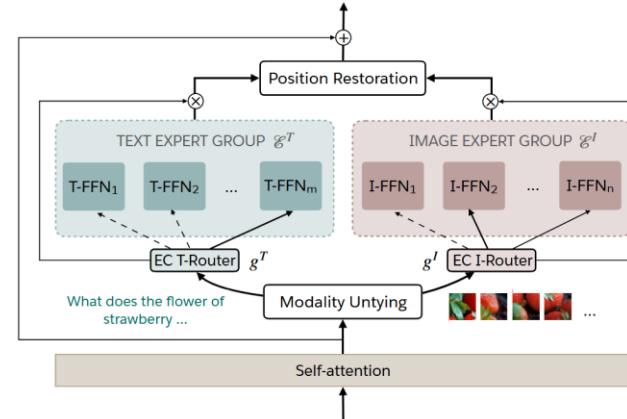
<sup>1</sup>Meta FAIR

17

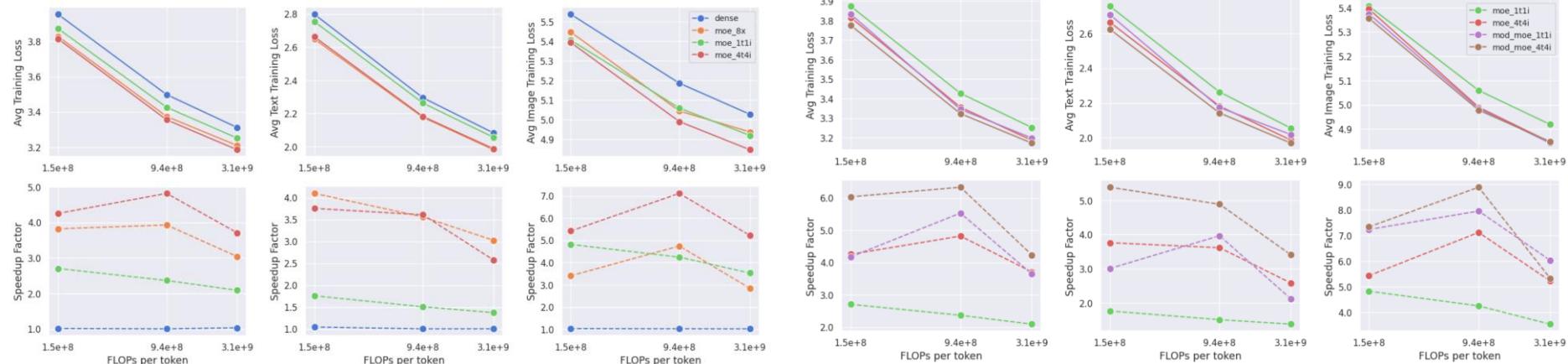
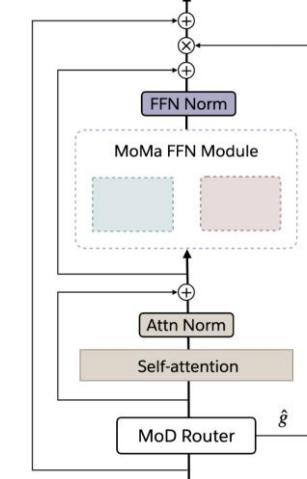
- 扩展FFN到多模态，每个模态内还有MoE，形成多级知识库
- 对整个模型也进行mixture of depth，即选择是否跳过一些编码器层



(a) Early-fusion mixed-modal LLM architecture.



(b) Mixture of modality-aware experts (MoMa) transformer block.



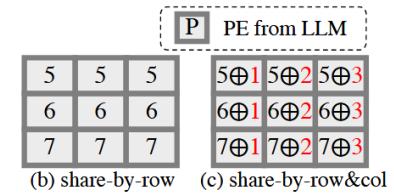
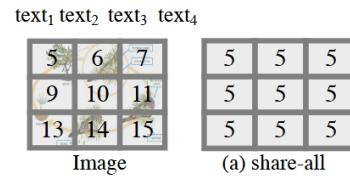
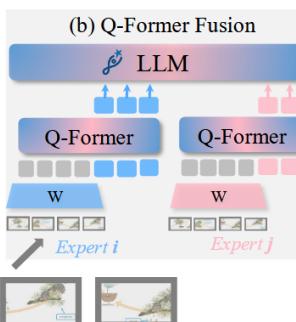
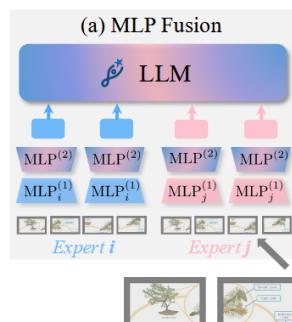
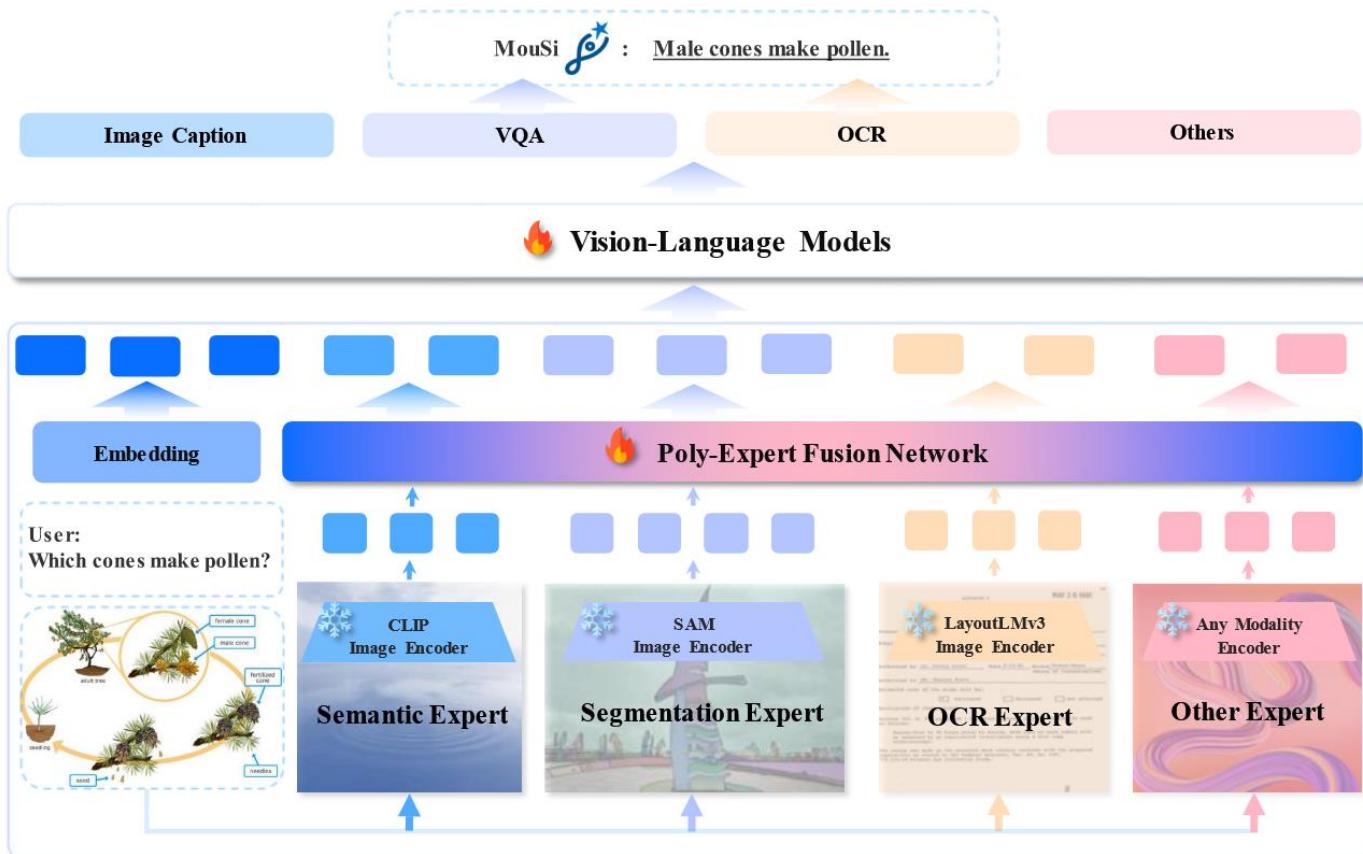
Vision Encoder(s)	Step1: High-resolution Adaptation and Training Recipe						Step2: Fusion Paradigm Exploration				
	Fusion	Token (V)	#Tokens/s	GQA	MME	MMMU	OCR	SQA	POPE	TextVQA	Avg
<i>CLIP</i>	–	1024	<b>47.9</b>	64.6	<b>1531</b>	35.7	496	<b>72.8</b>	87.9	64.0	644.4
<i>ConvNeXt</i>	–	1024	46.8	63.7	1433	<b>37.0</b>	527	72.0	87.9	71.7	652.3
<i>CLIP</i> + <i>ConvNeXt</i>	SA	2048	46.1	64.6	1482	35.3	536	72.5	88.3	71.8	657.2
	<i>CC</i>	1024	47.3	64.0	1486	36.0	533	72.7	88.6	71.9	<b>658.2</b>
	<i>LH</i>	1024	47.0	62.9	1488	36.8	521	72.6	<b>88.9</b>	68.7	652.1
	<i>MG</i>	1024	45.2	<b>64.9</b>	1481	36.7	474	<b>72.8</b>	88.2	68.6	646.7
	<i>DA</i>	1024	47.3	63.6	1497	35.3	459	71.9	88.0	69.1	640.9
<i>CLIP</i>	SA	3072	40.3	<b>64.9</b>	1439	35.0	504	72.7	88.2	71.3	649.2
+ <i>ConvNeXt + SAM</i>	<i>CC</i>	1024	46.3	64.6	1497	35.2	<b>558</b>	71.7	87.9	<b>72.2</b>	<b>660.4</b>

Model		Knowledge				General				OCR and Chart				Vision-Centric					
		Avg	SQA <sup>1</sup>	MMU	MathVista	AI2D	Avg	MME	MMBench	SEED	GQA	Avg	ChartQA	OCRBench	TextVQA	DocVQA	Avg	MMVP	RWQA
<i>Llama3-8B</i>																			
<i>MGM-HD</i>	55.7	75.1	37.3	37.0	73.5		72.7	<b>1606</b>	72.7	73.2	64.5	62.9	59.1	47.7	70.2	74.6	40.4	18.7	62.1
<i>Cambrian-1</i>	61.3	80.4	42.7	49.0	73.0		73.1	1547	75.9	74.7	64.6	71.3	73.3	62.4	71.7	77.8	57.6	51.3	64.2
<i>Eagle</i>	<b>64.2</b>	<b>84.3</b>	<b>43.8</b>	<b>52.7</b>	<b>76.1</b>		<b>73.8</b>	1559	<b>75.9</b>	<b>76.3</b>	<b>64.9</b>	<b>76.6</b>	<b>80.1</b>	<b>62.6</b>	<b>77.1</b>	<b>86.6</b>	<b>69.1</b>	<b>71.6</b>	<b>66.5</b>
<i>Vicuna-13B</i>																			
<i>MGM-HD</i>	54.1	71.9	37.3	37	70.1		70.7	1597	68.6	70.6	63.7	60.8	56.6	46.6	70.2	69.8	38.4	19.3	57.5
<i>Cambrian-1</i>	60.2	79.3	40.0	48.0	73.6		73.7	1610	<b>75.7</b>	74.4	64.3	71.3	73.8	<b>61.9</b>	72.8	76.8	52.2	41.3	63.0
<i>Eagle</i>	<b>63.0</b>	<b>82.0</b>	<b>41.6</b>	<b>54.4</b>	<b>74.0</b>		<b>74.6</b>	<b>1651</b>	<b>75.7</b>	<b>74.8</b>	<b>65.3</b>	<b>75.1</b>	<b>77.6</b>	<b>61.9</b>	<b>75.5</b>	<b>85.4</b>	<b>61.4</b>	<b>58.0</b>	<b>64.8</b>
<i>Yi-34B</i>																			
<i>MGM-HD</i>	62.4	77.7	48	43.4	80.5		76.2	1659	80.6	75.3	65.8	68.1	67.6	51.8	74.1	78.9	52.3	37.3	67.2
<i>Cambrian-1</i>	67.0	<b>85.6</b>	49.7	53.2	<b>79.7</b>		<b>76.8</b>	<b>1689</b>	<b>81.4</b>	75.3	<b>65.8</b>	71.9	75.6	60.0	76.7	75.5	60.3	52.7	67.8
<i>Eagle</i>	<b>68.6</b>	85.5	<b>53.2</b>	<b>57.9</b>	79.1		76.3	1677	81.0	<b>75.6</b>	64.9	<b>75.4</b>	<b>77.2</b>	<b>62.4</b>	<b>78.8</b>	<b>83.0</b>	<b>68.3</b>	<b>67.0</b>	<b>69.5</b>

# Poly-Visual-Expert Vision-Language Models



19

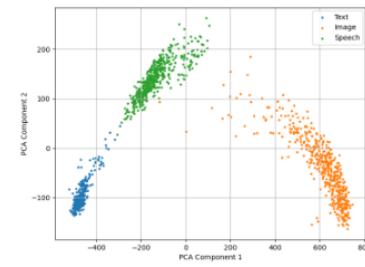
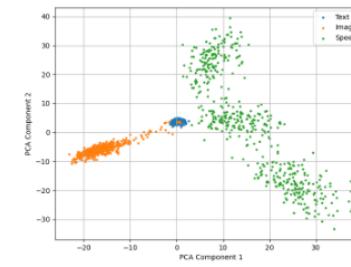
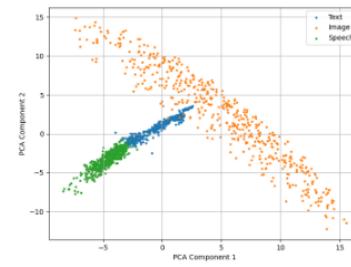
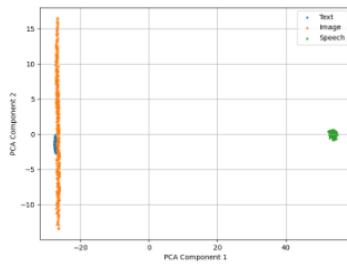
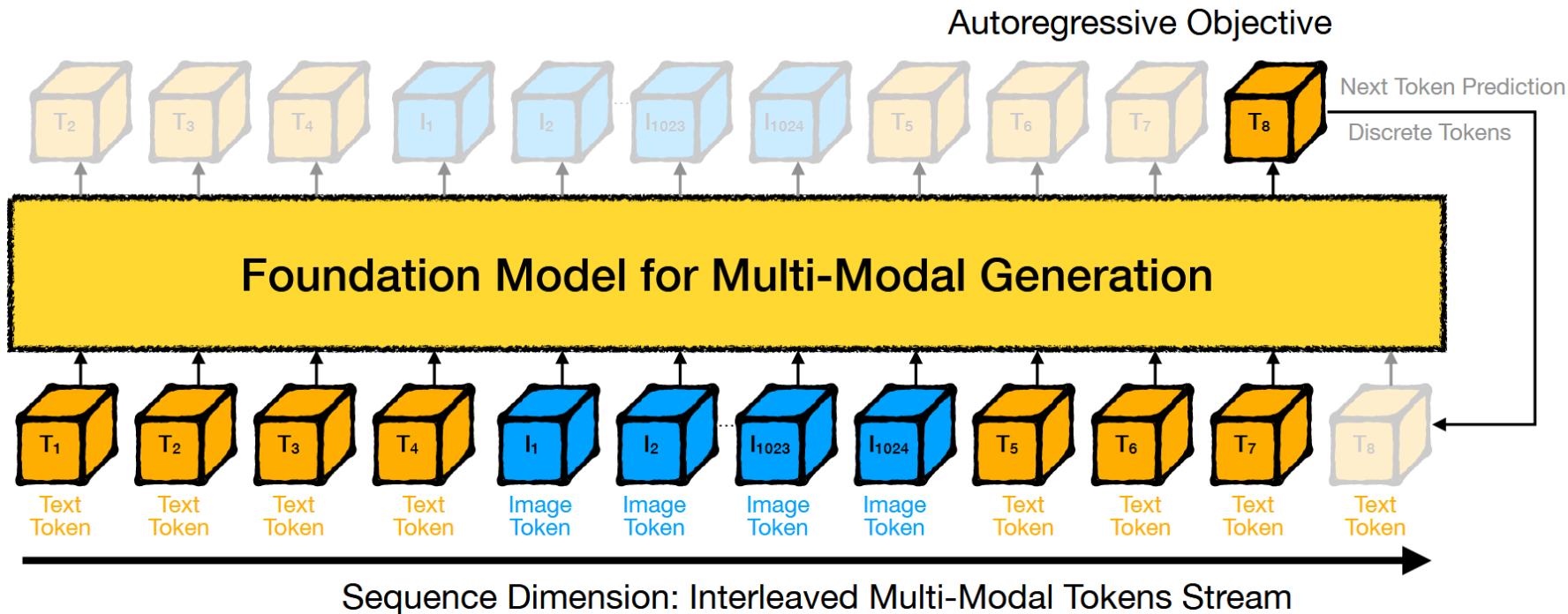




# Mixture-of-Transformers: A Sparse and Scalable Architecture for Multi-Modal Foundation Models

Weixin Liang<sup>2,\*</sup>, Lili Yu<sup>1†</sup>, Liang Luo<sup>1†</sup>, Srinivasan Iyer<sup>1</sup>, Ning Dong<sup>1</sup>, Chunting Zhou<sup>1</sup>, Gargi Ghosh<sup>1</sup>, Mike Lewis<sup>1</sup>, Wen-tau Yih<sup>1</sup>, Luke Zettlemoyer<sup>1</sup>, Xi Victoria Lin<sup>1</sup>

20



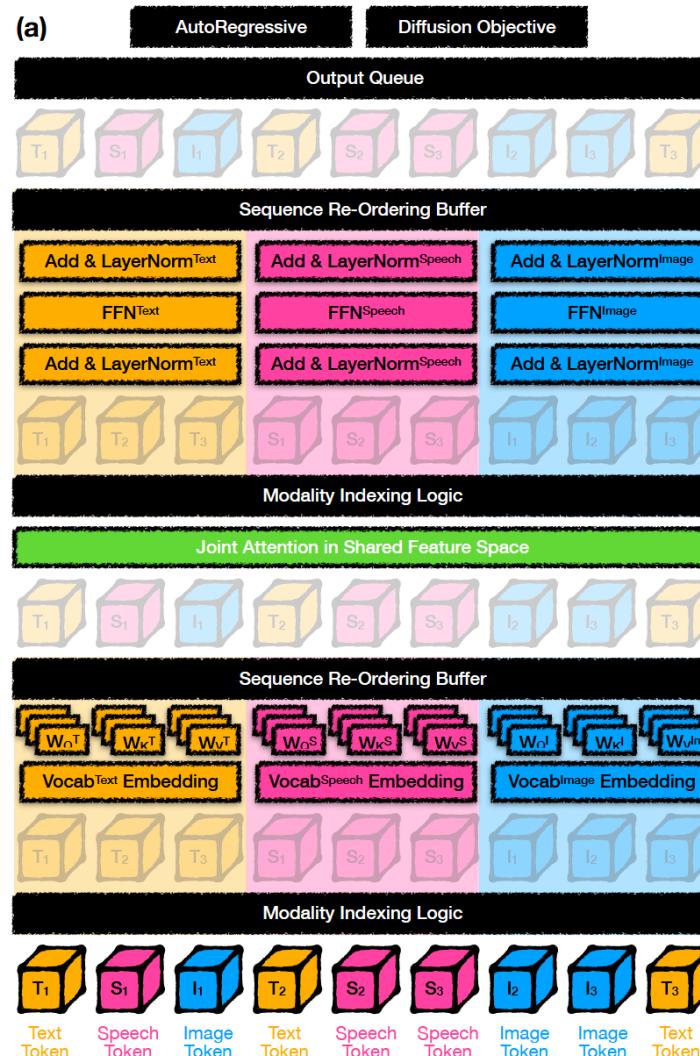
# Mixture-of-Transformers: A Sparse and Scalable Architecture for Multi-Modal Foundation Models



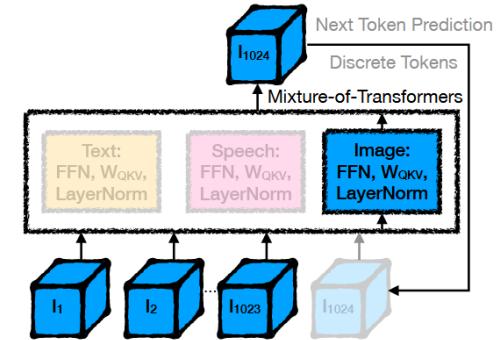
Weixin Liang<sup>2,\*</sup>, Lili Yu<sup>1†</sup>, Liang Luo<sup>1†</sup>, Srinivasan Iyer<sup>1</sup>, Ning Dong<sup>1</sup>, Chunting Zhou<sup>1</sup>, Gargi Ghosh<sup>1</sup>, Mike Lewis<sup>1</sup>, Wen-tau Yih<sup>1</sup>, Luke Zettlemoyer<sup>1</sup>, Xi Victoria Lin<sup>1</sup>

21

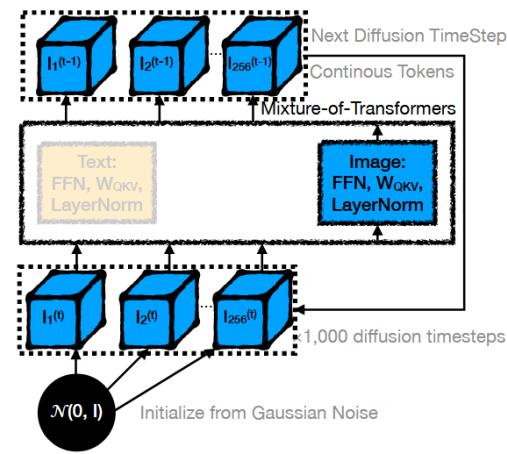
- 各模态在attn中共同计算注意力，但采用各自的QKV进行映射
- 各模态采用各自的FFN和Norm，在FFN部分单独进行
- 多种方式进行生成训练，自回归或者扩散均可



(b) Autoregressive Objective



(c) Diffusion Objective





- Standard MoE
- Internal Mixture
- External Mixture
- A New Insight
- Summary



- Standard MoE
- Internal Mixture
- External Mixture
- A New Insight
- Summary



# 总结反思

24

- 可以从各个角度构建MoE，只要符合“比较-激活-加权”的形式，均可以视为MoE
- MoE的用途：
  - 多任务：Eagle, Poly,
  - 多模态：Uni-MoE, Mono-InternVL, MoMa, MoT
  - 稀疏高效：Mixtral, DeepseekMoE, MoH, Sparse Upcycling, LoRAMoE, OLMoE, MoE-LLaVA
  - 微调：LoRAMoE, MoLE
  - .....



# 谢谢！