



中國人民大學  
RENMIN UNIVERSITY OF CHINA



GeWu-Lab

Gaoling School of Artificial Intelligence  
Renmin University of China

# 平衡多模态学习：过去与当下

卫雅珂

中国人民大学

yakewei@ruc.edu.cn

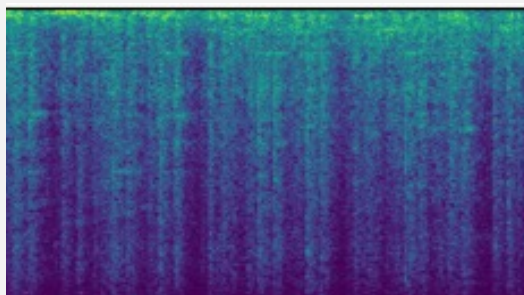
# 什么是平衡多模态学习？

“平衡多模态学习”：不同模态的差异化学习特点会导致部分学习较慢较差的模态在优化竞争中受到较优模态的抑制，从而造成部分模态的学习不充分。

音视频样本对: *Bus*



干净 & 容易学习



嘈杂 & 难以学习

容易学习  
Modality 1

单模态  
编码器

难以学习  
Modality 2

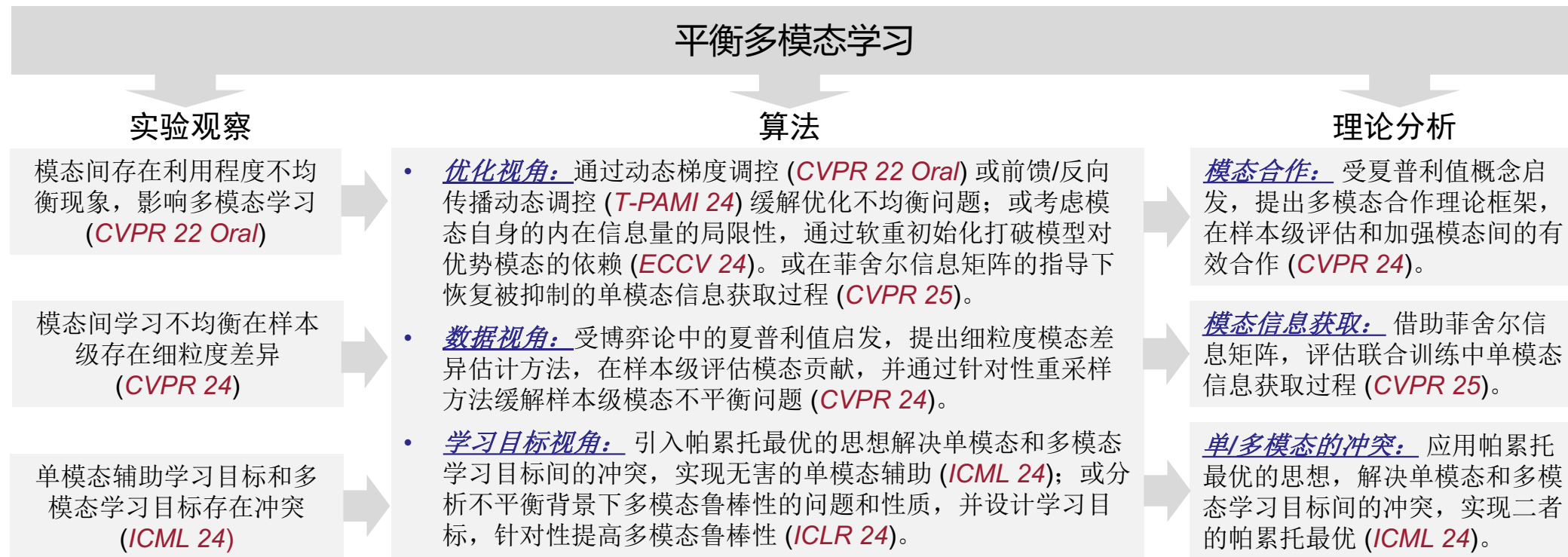
单模态  
编码器

交互融合  
&  
判别性学习

最终预测

学习目标

# 平衡多模态学习：过去



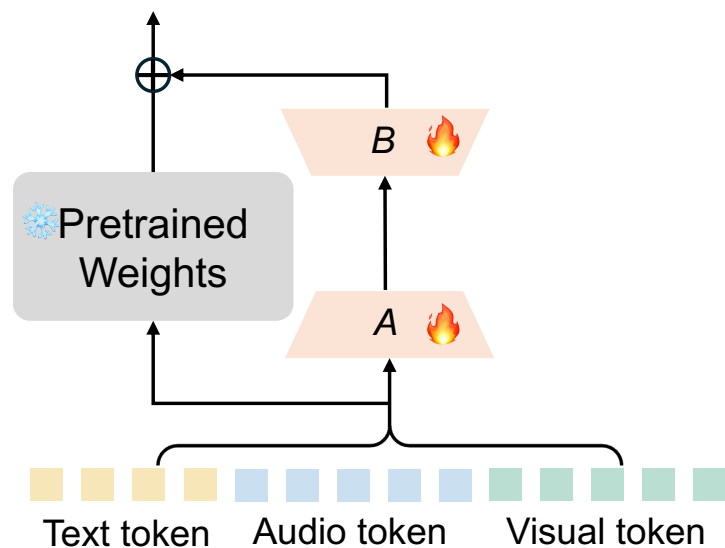
- Yake Wei, Di Hu, Henghui Du, and Ji-Rong Wen. On-the-fly modulation for balanced multimodal learning. *T-PAMI*, 2024.
- Yake Wei and Di Hu. Mmpareto: boosting multimodal learning with innocent unimodal assistance. *ICML*, 2024.
- Yake Wei, Ruoxuan Feng, Ziheng Wang, and Di Hu. Enhancing multimodal cooperation via sample-level modality valuation. *CVPR*, 2024.
- Yake Wei, Siwei Li, Ruoxuan Feng, and Di Hu. Diagnosing and re-learning for balanced multimodal learning. *ECCV*, 2024.
- Xiaokang Peng\*, Yake Wei\*, Andong Deng, Dong Wang, and Di Hu. Balanced multimodal learning via on-the-fly gradient modulation. *CVPR*, 2022, **Oral**.
- Chengxiang Huang\*, Yake Wei\*, Zequn Yang, and Di Hu. Adaptive Unimodal Regulation for Balanced Multimodal Information Acquisition. *CVPR*, 2025.
- Zequn Yang, Yake Wei, Ce Liang, and Di Hu. Quantifying and enhancing multi-modal robustness with modality preference. *ICLR*, 2024.

大模型背景下

数据的问题依然存在

学习的贪婪特性未曾改变

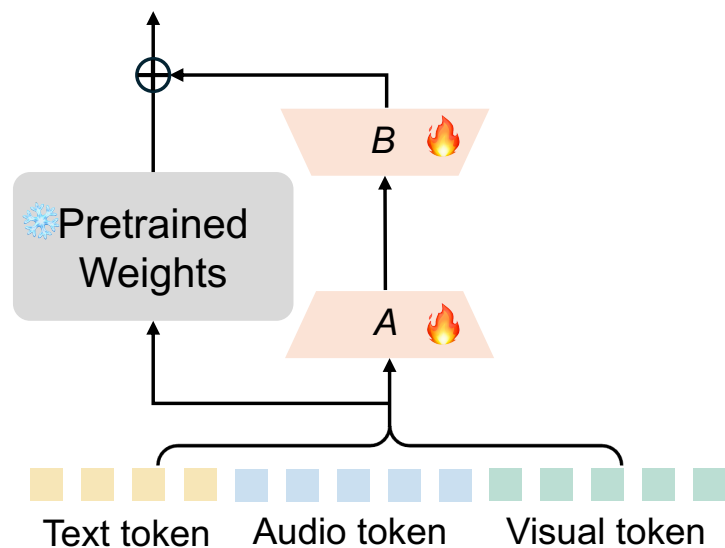
# 观察：被忽视的单模态特性



典型LoRA架构

$$h = W_0 \mathbf{x} + \Delta W \mathbf{x} = W_0 \mathbf{x} + \Delta W [\mathbf{x}^{m_1}; \mathbf{x}^{m_2}; \dots; \mathbf{x}^{m_n}]$$

# 观察：被忽视的单模态特性

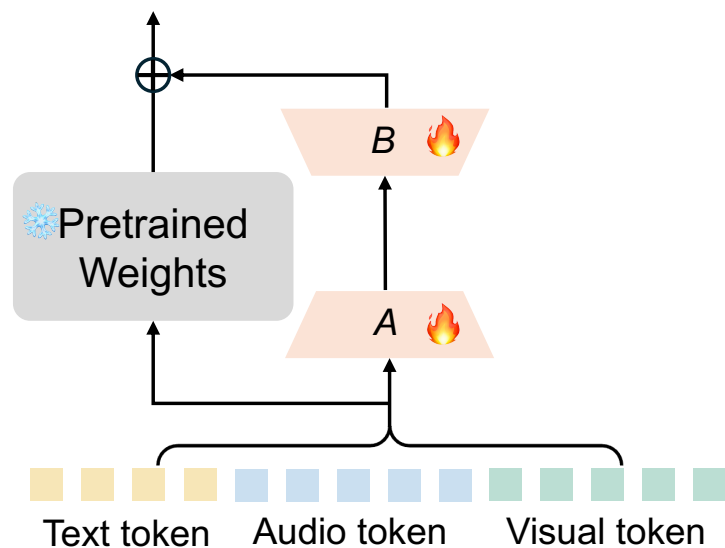


典型LoRA架构

$$h = W_0 \mathbf{x} + \Delta W \mathbf{x} = W_0 \mathbf{x} + \Delta W [\mathbf{x}^{m_1}; \mathbf{x}^{m_2}; \dots; \mathbf{x}^{m_n}]$$

直接“borrow” LLM微调方法完全契合多模态场景吗？

# 观察：被忽视的单模态特性

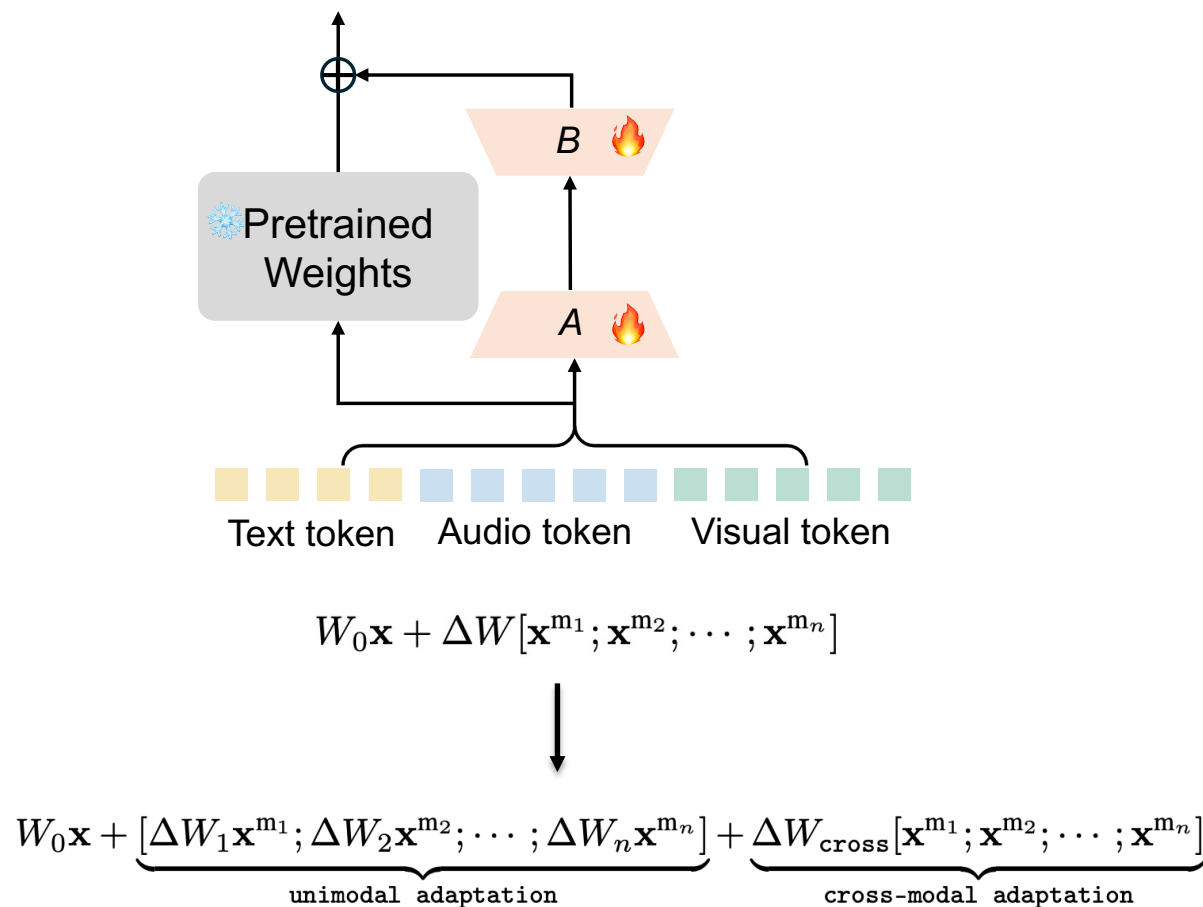


典型LoRA架构



完全统一处理的范式并非多模态的最优解  
MLLM finetuning应该兼顾单模态与跨模态特性

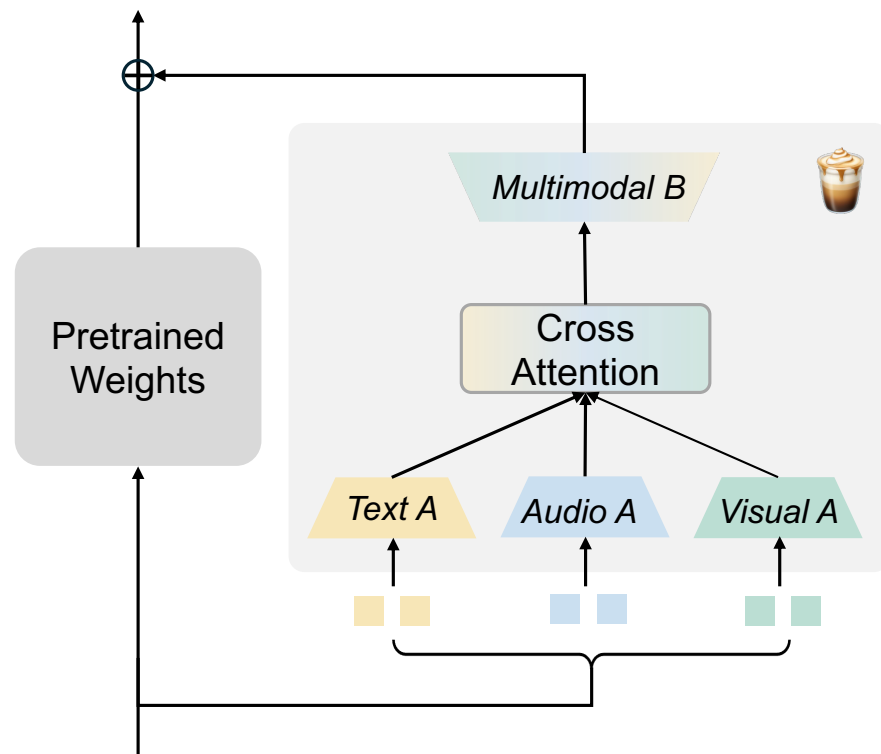
# 观察：被忽视的单模态特性





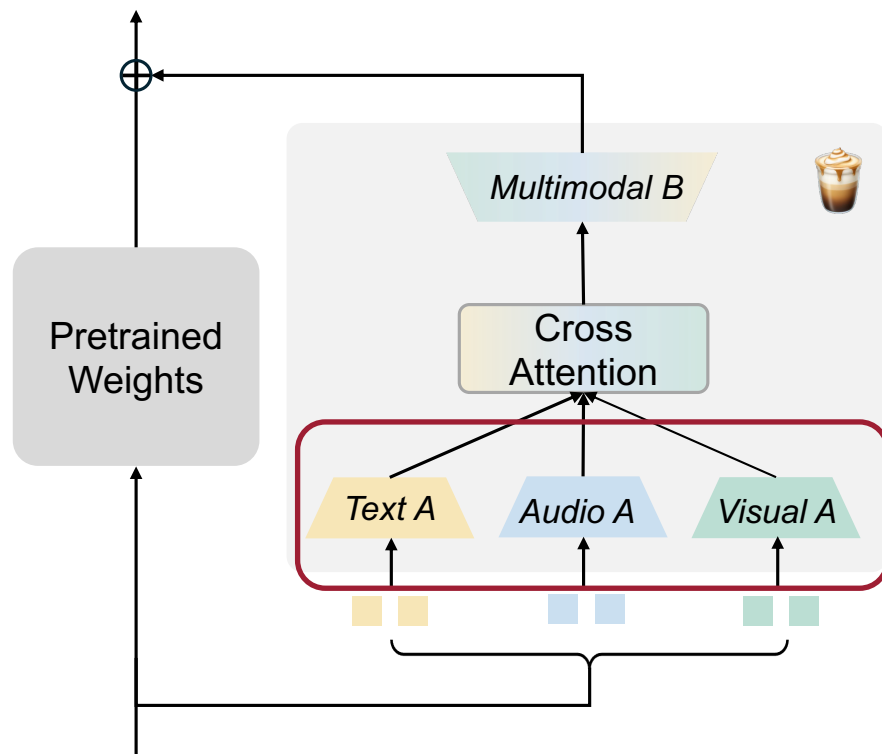
# MokA: 兼具单模态与跨模态特性

## Multimodal Low-Rank Adaptation for MLLMs



# MokA: 兼具单模态与跨模态特性

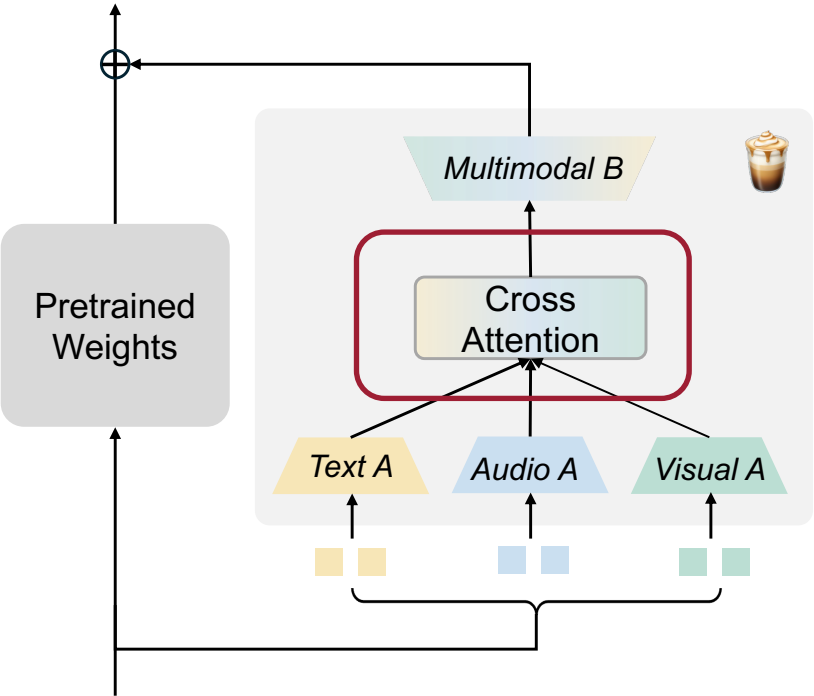
模态独有matrix A: 独立抽取模态信息，互不干扰



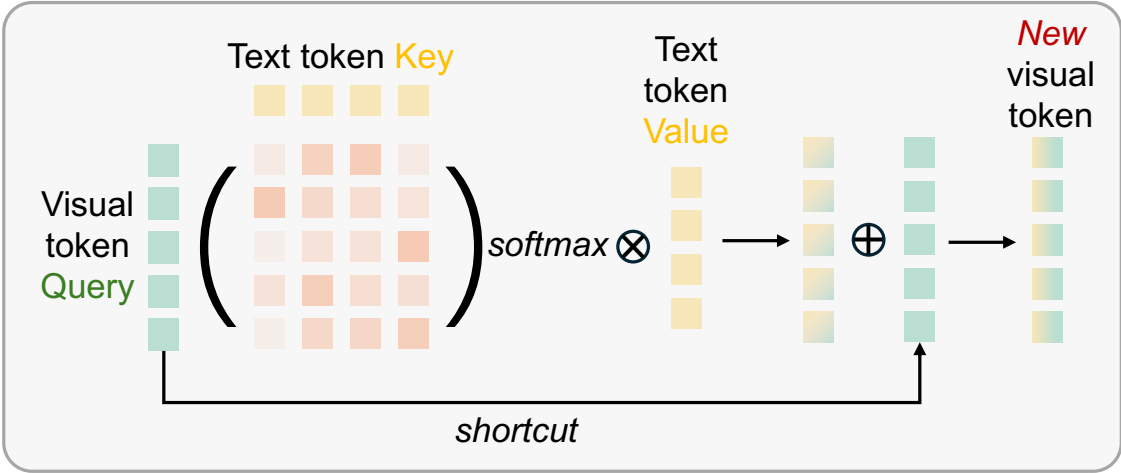


# MokA: 兼具单模态与跨模态特性

显式模态交互: Cross-attention

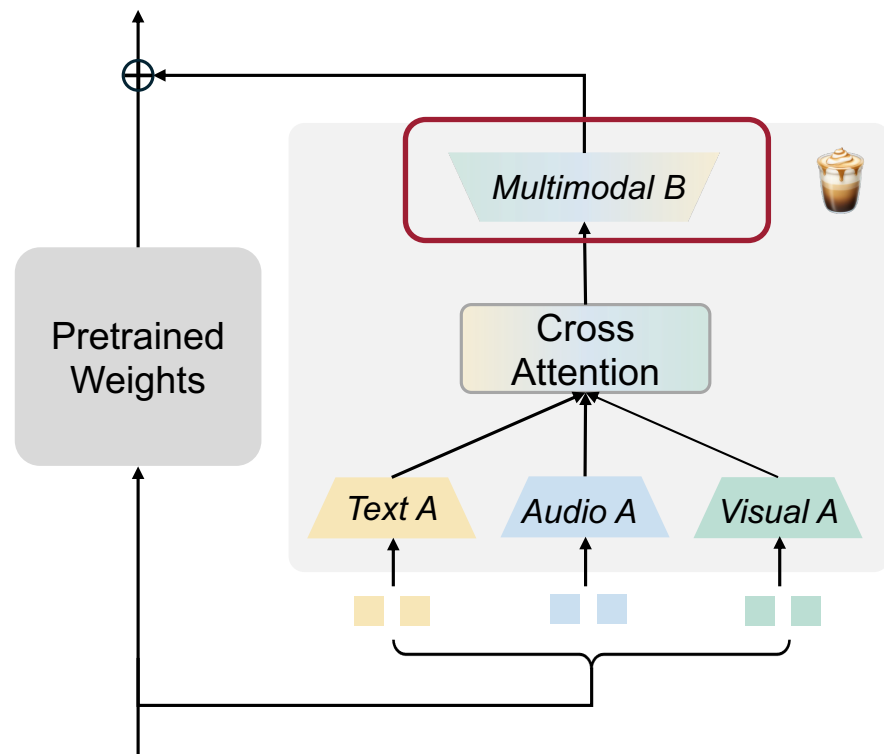


<audio> <visual> Please answer the question: which clarinet makes the sound first?



# MokA: 兼具单模态与跨模态特性

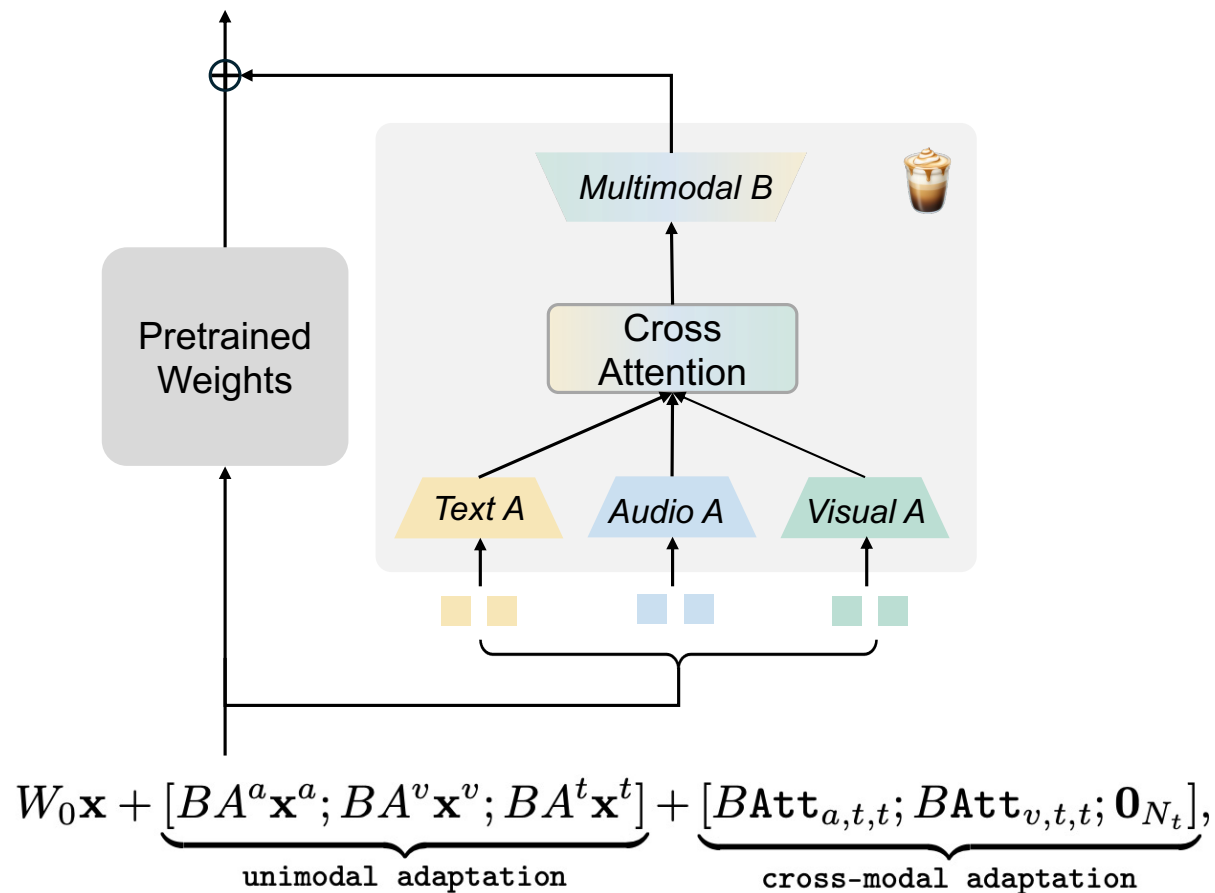
模态共享matrix B: 分离的单模态子空间映射回同一共享空间





# MokA: 兼具单模态与跨模态特性

## Multimodal Low-Rank Adaptation for MLLMs





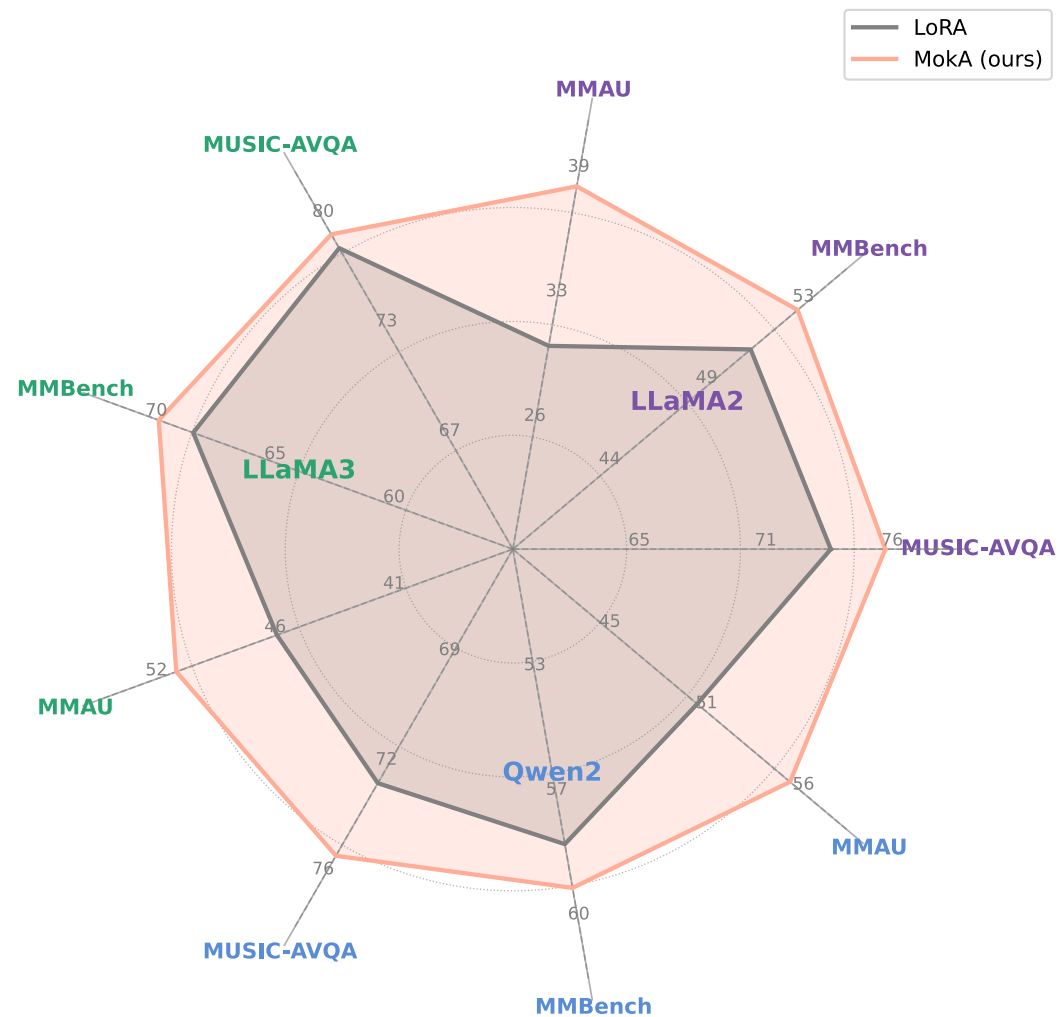
# MokA: 兼具单模态与跨模态特性

多个基座模型:

- LLaMA2
- LLaMA3
- Qwen2

多种多模态场景:

- Audio-Visual-Text
- Visual-Text
- Speech-Text



# Thanks for listening!

BML交流社区



MokA主页

