=1920bp =1080bp

# Multimodal Large Language Model Meets Remote Sensing
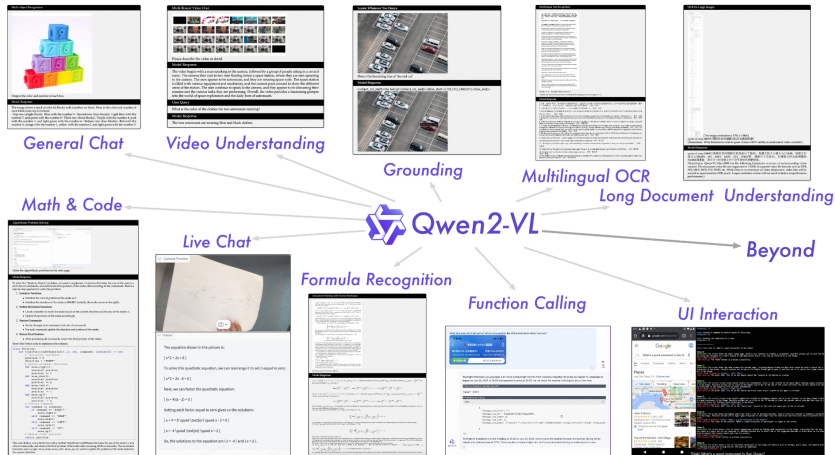## GISLab Short-Term Course 2025 Summer

Zhenyuan Chen

School of Earth Science, Zhejiang University
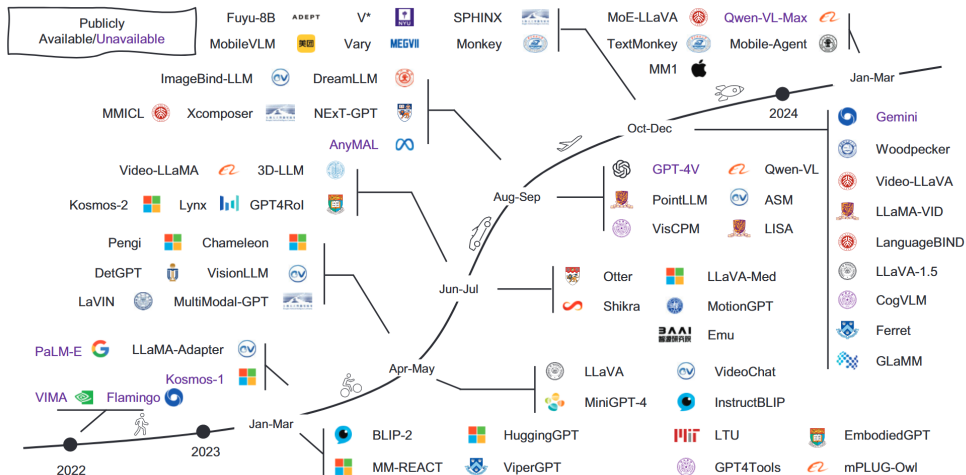
2025
bili_sakura@zju.edu.cn

# Background: Multimodal Large Language Models (MLLMs)



General Chat · Video Understanding · Grounding · Multilingual OCR · Long Document Understanding · Math & Code · Live Chat · Qwen2-VL · Formula Recognition · Function Calling · Beyond · UI Interaction

Qwen2-VL (Wang et al., 2024).

Wang, et al. Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution, 2024.

# Background



MLLM Timeline (Yin et al., 2024).

Yin, et al. A Survey on Multimodal Large Language Models. National Science Review. 2024.

# What is a Multimodal Large Language Model (MLLM)?



Figure: General architecture of a Multimodal Large Language Model (MLLM) (Yin et al., 2024).

Yin, et al. A Survey on Multimodal Large Language Models. National Science Review. 2024.

# Vision Encoder: CLIP



Figure: OpenAI CLIP architecture (Radford et al., 2021).

Radford, et al. Learning Transferable Visual Models From Natural Language Supervision, ICML, 2021.
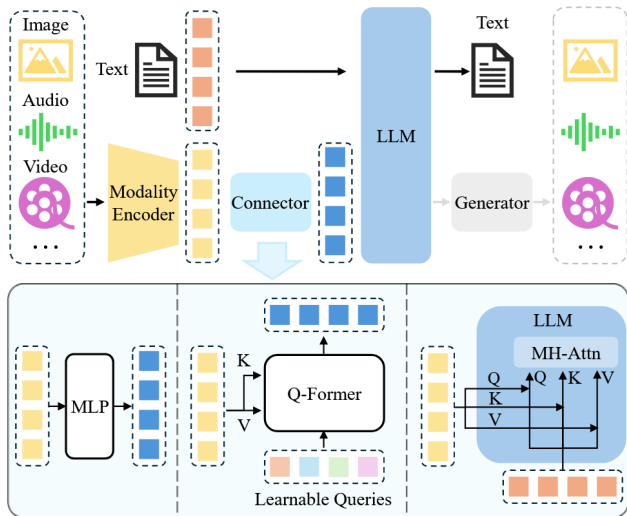
# LLM Backbone: Large Language Models



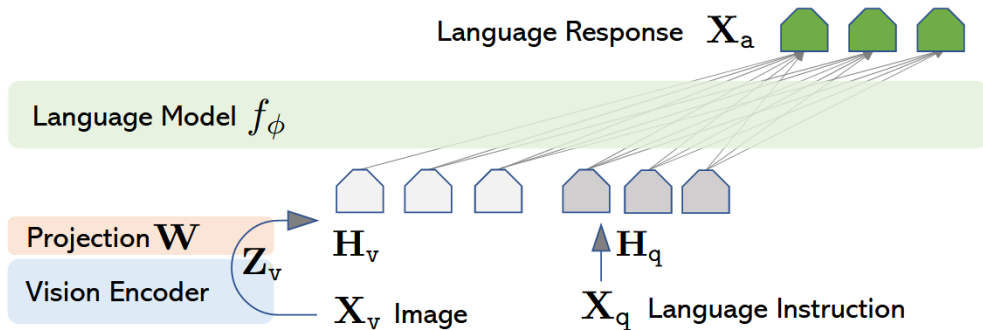Latest LLMs are generally composed of a series of Transformers decoder blocks (Vaswani et al., 2017).

(image source:cameronrwolfe)

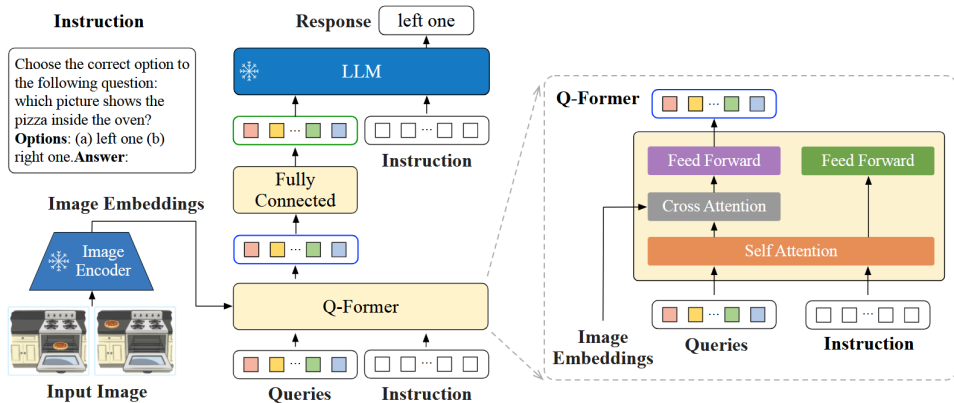Vaswani, et al. Attention Is All You Need, NeurIPS, 2017.

# Connector: General MLLM



An typical MLLM architecture includes an encoder, a connector, and a LLM. (Yin et al., 2024).

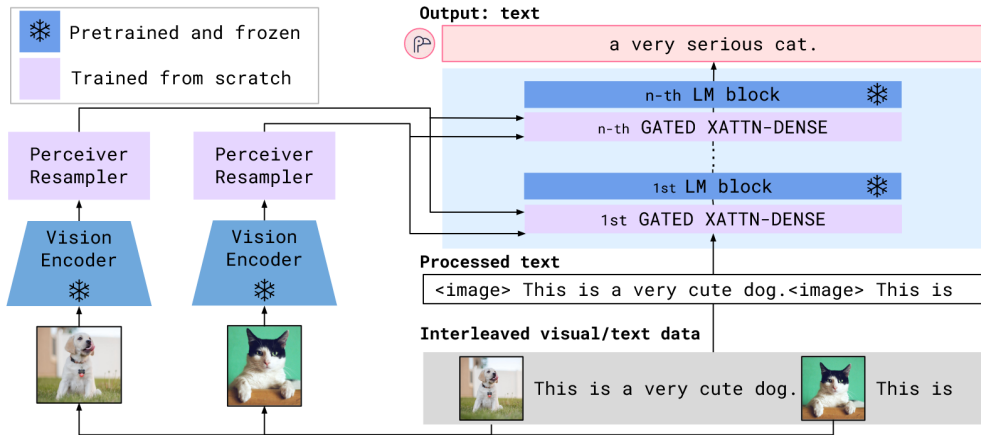Yin, et al. A Survey on Multimodal Large Language Models. National Science Review. 2024.

# Connector: MLP



Connector architecture in LLaVA (Liu et al., 2023).

Liu, et al. Visual Instruction Tuning, NeurIPS, 2023.

# Connector: Q-Former



Connector architecture in InstructBLIP (Dai et al., 2023).

Dai, et al. InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning, NeurIPS, 2023.

# Connector: Feature-level Fusion



Connector architecture in Flamingo (Alayrac et al., 2022).

Alayrac, et al. Flamingo: A Visual Language Model for Few-Shot Learning, NeurIPS, 2022.