

# Multimodal Large Language Model Meets Remote Sensing

GISLab Short-Term Course 2025 Summer

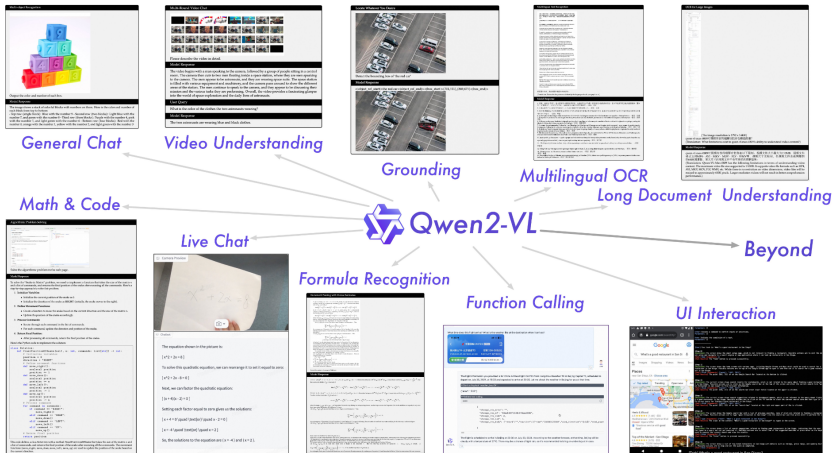
Zhenyuan Chen

School of Earth Science, Zhejiang University

2025

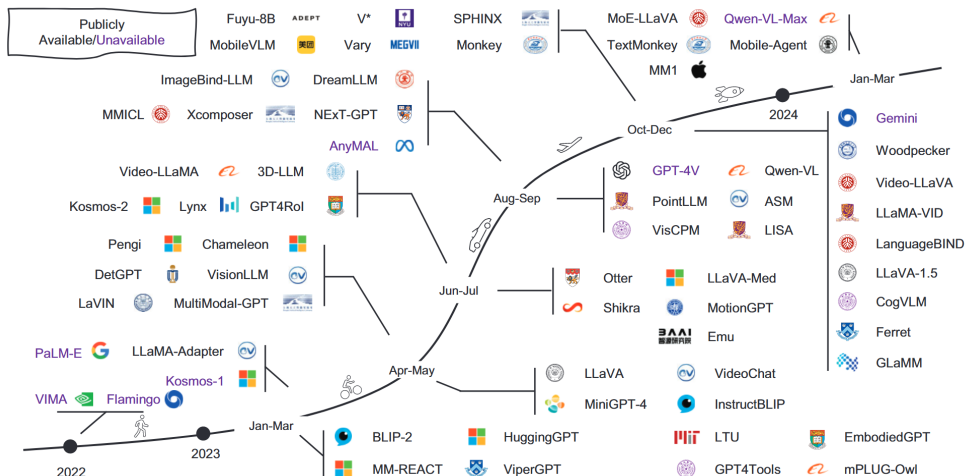
`bili_sakura@zju.edu.cn`

# Background: Multimodal Large Language Models (MLLMs)



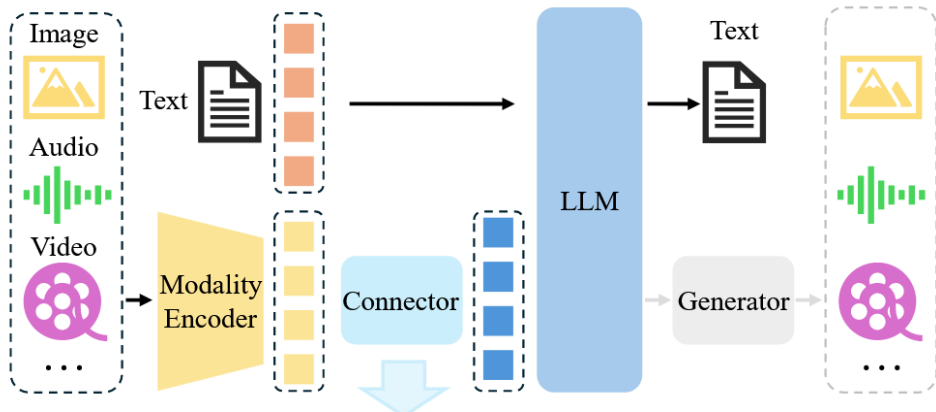
Qwen2-VL (wangQwen2VLEnhancingVisionLanguage2024a).

# Background



MLLM Timeline (inSurveyMultimodalLarge2024a).

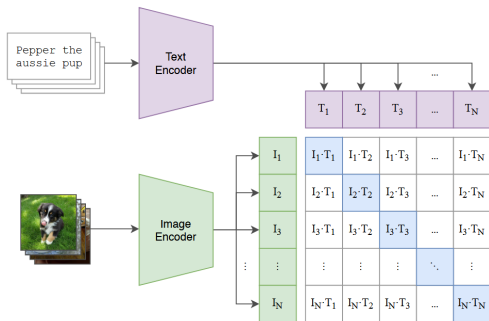
# What is a Multimodal Large Language Model (MLLM)?



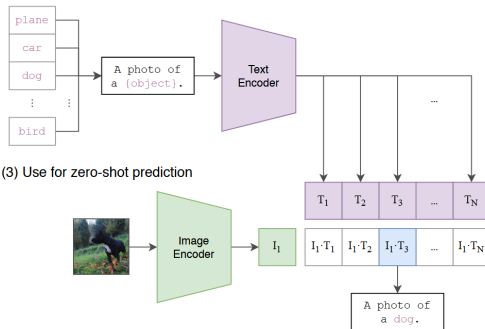
**Figure:** General architecture of a Multimodal Large Language Model (MLLM) (yinSurveyMultimodalLarge2024a).

# Vision Encoder: CLIP

(1) Contrastive pre-training



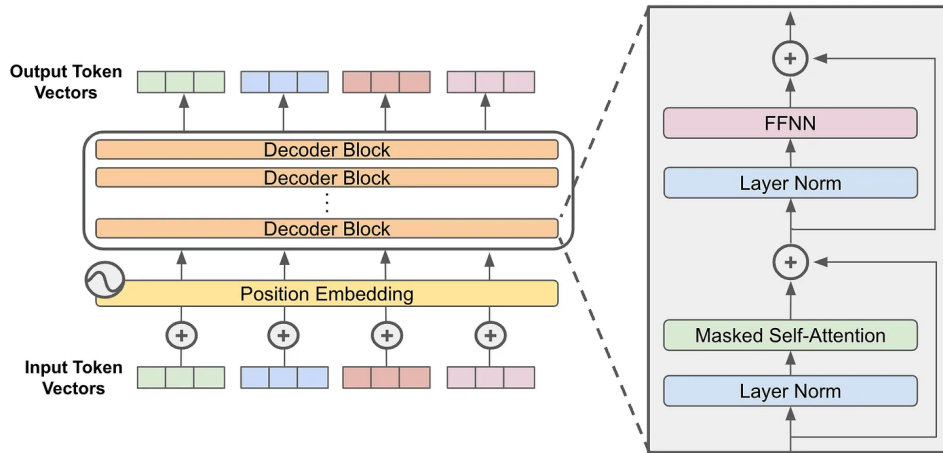
(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

Figure: OpenAI CLIP architecture (radfordLearningTransferableVisual2021).

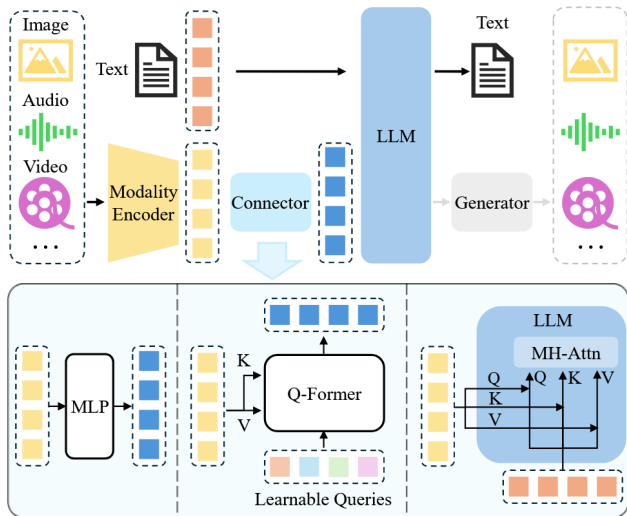
# LLM Backbone: Large Language Models



Latest LLMs are generally composed of a series of Transformers decoder blocks (**vaswaniAttentionAllYou2017**).

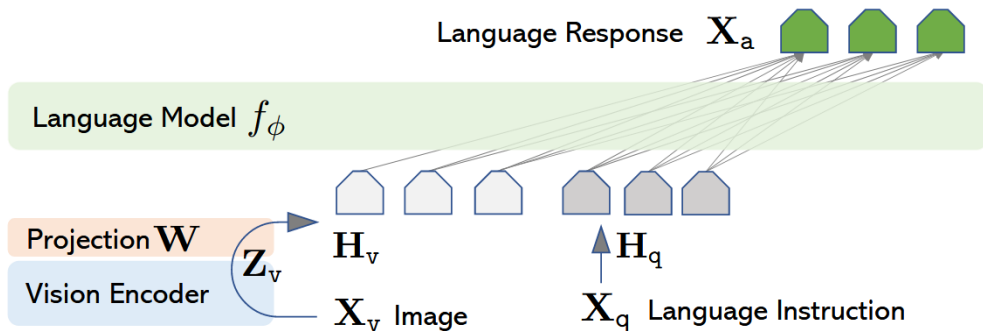
(image source:cameronrwolfe)

# Connector: General MLLM



An typical MLLM architecture includes an encoder, a connector, and a LLM. ([yinSurveyMultimodalLarge2024a](#)).

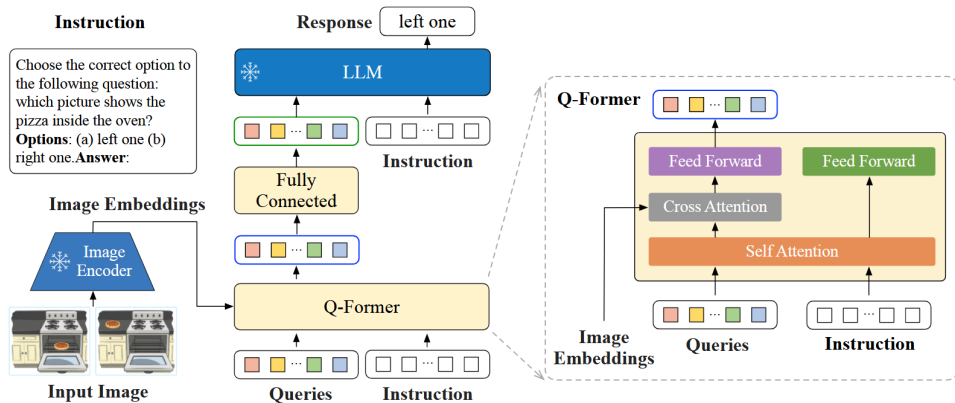
## Connector: MLP



Connector architecture in LLaVA (liuVisualInstructionTuning2023).

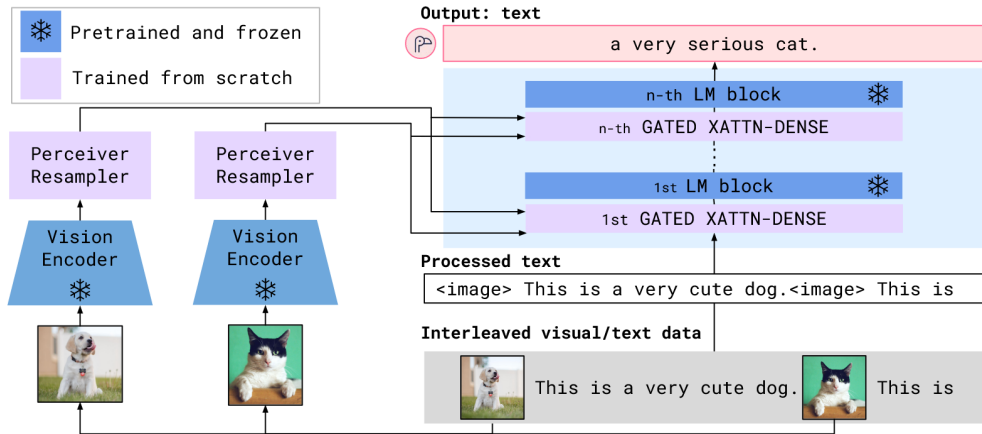


# Connector: Q-Former



Connector architecture in InstructBLIP (daInstructBLIPGeneralpurposeVisionLanguage2023a).

# Connector: Feature-level Fusion



Connector architecture in Flamingo (alayracFlamingoVisualLanguage2022).