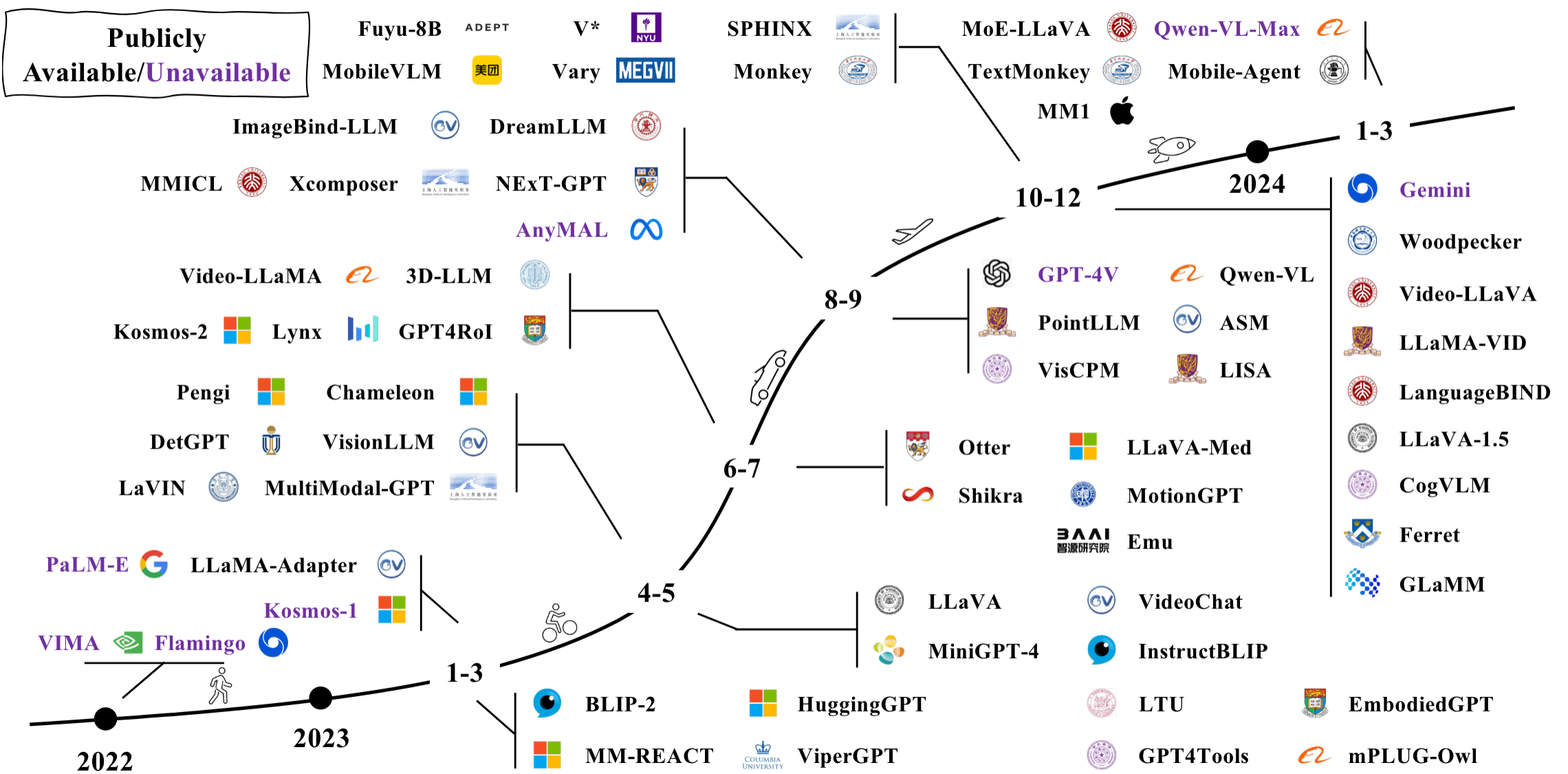


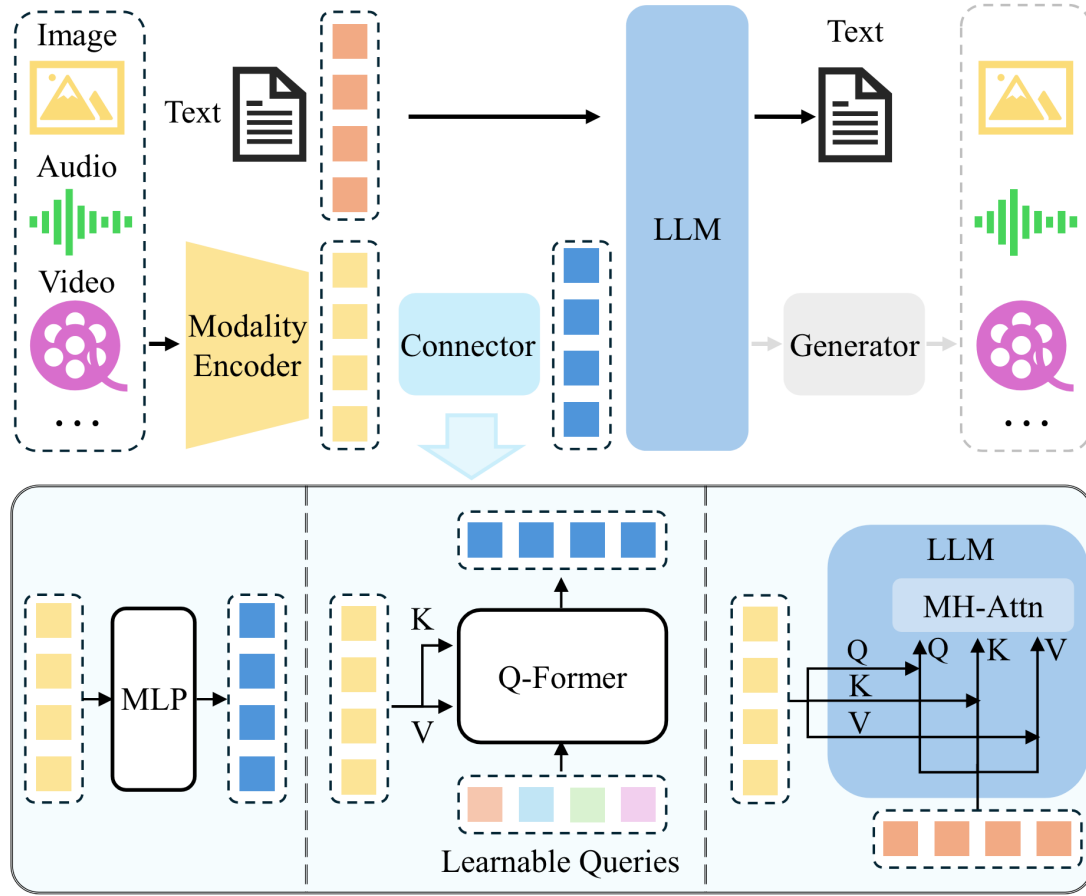
Multimodal Large Language Model

Yongxin Wang

06/28/2024



Architecture

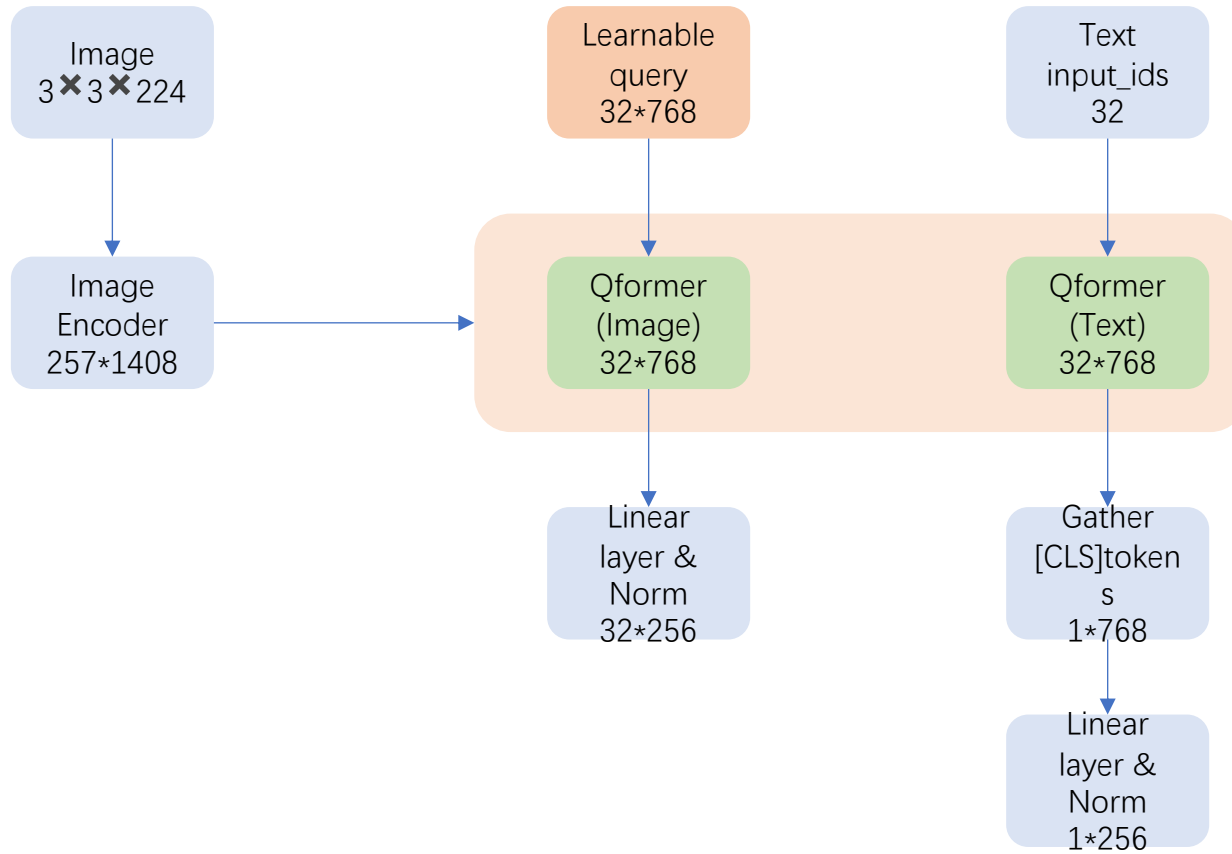


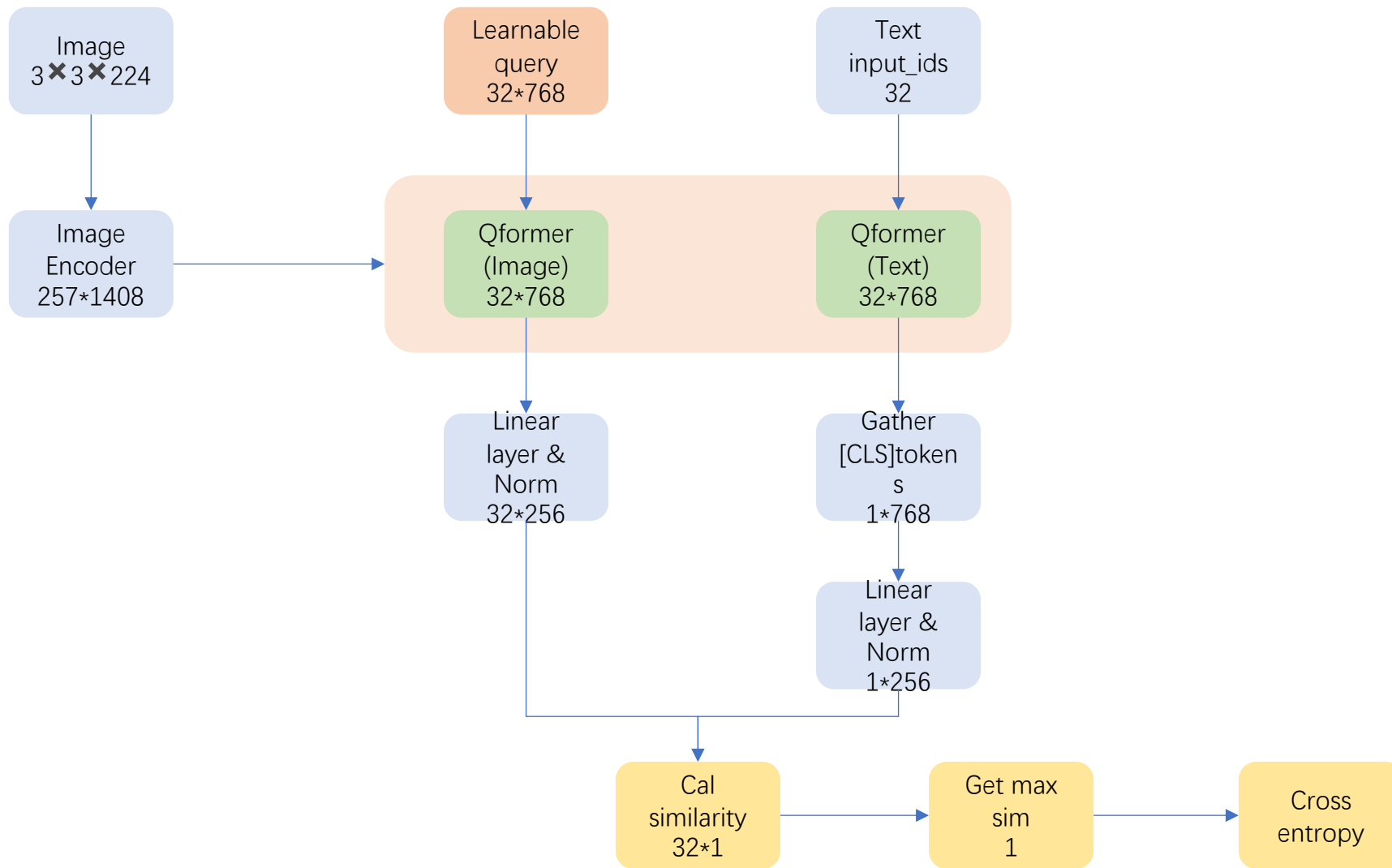
LLaVA-series

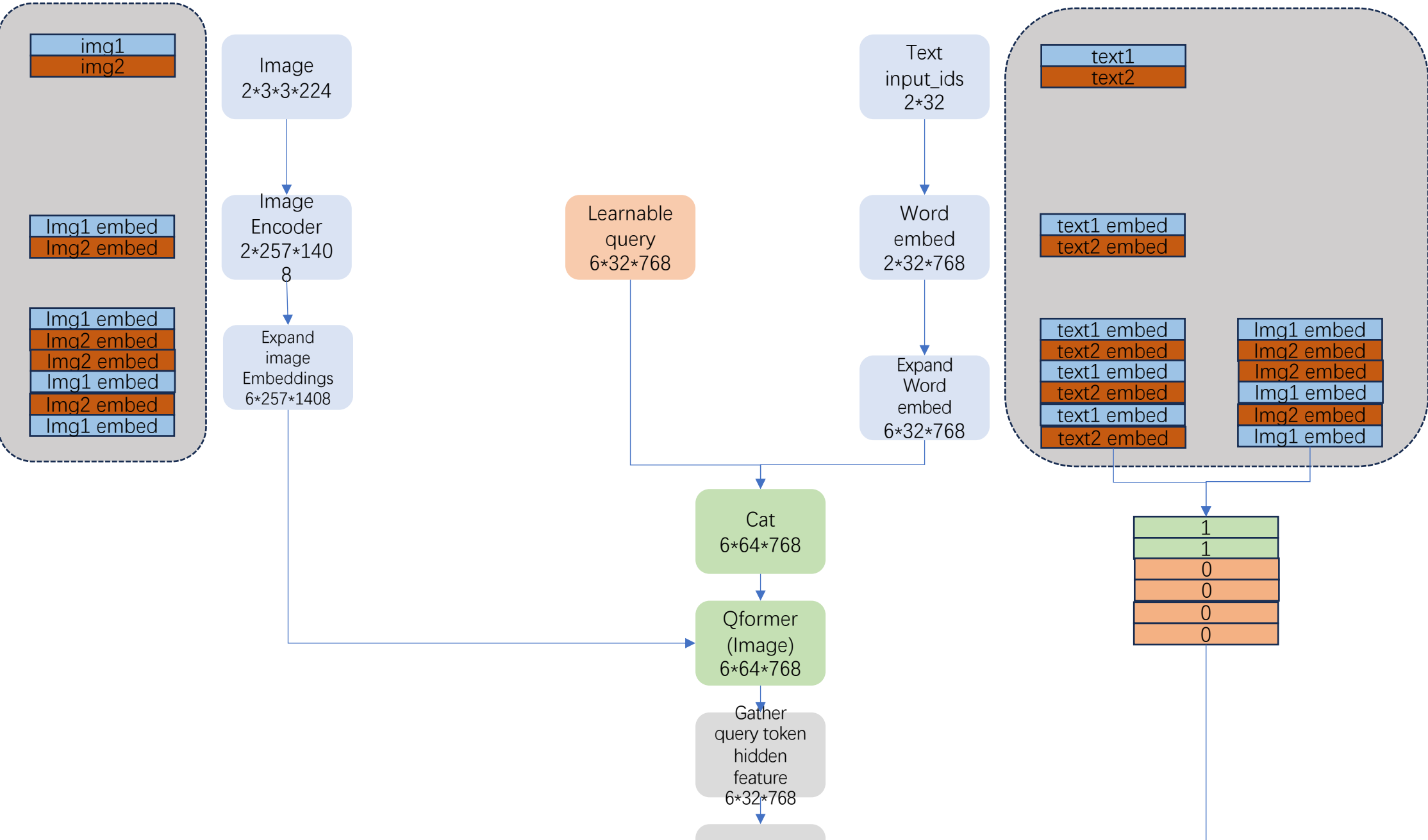
BLIP2, InstructBLIP-series

Flamingo, LLaMA-Adapter

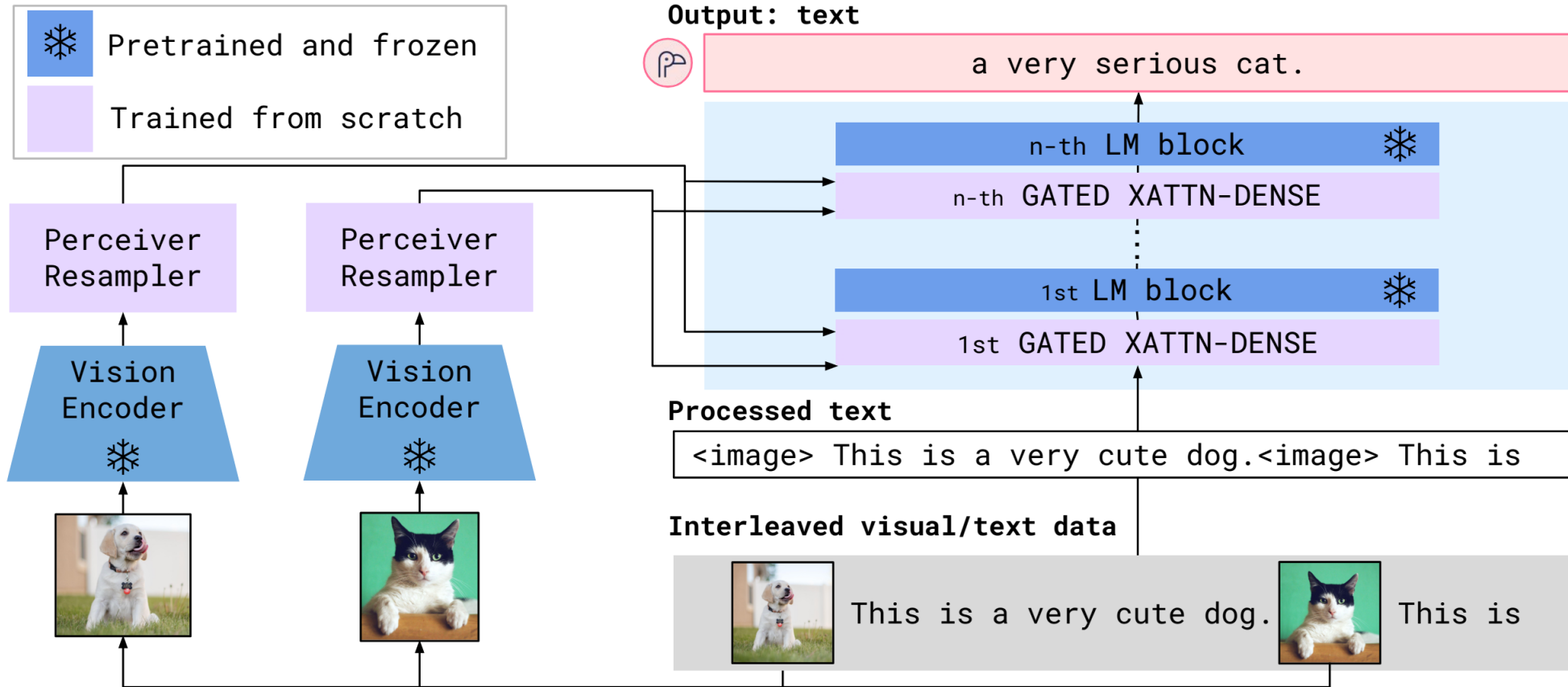
Q-Former

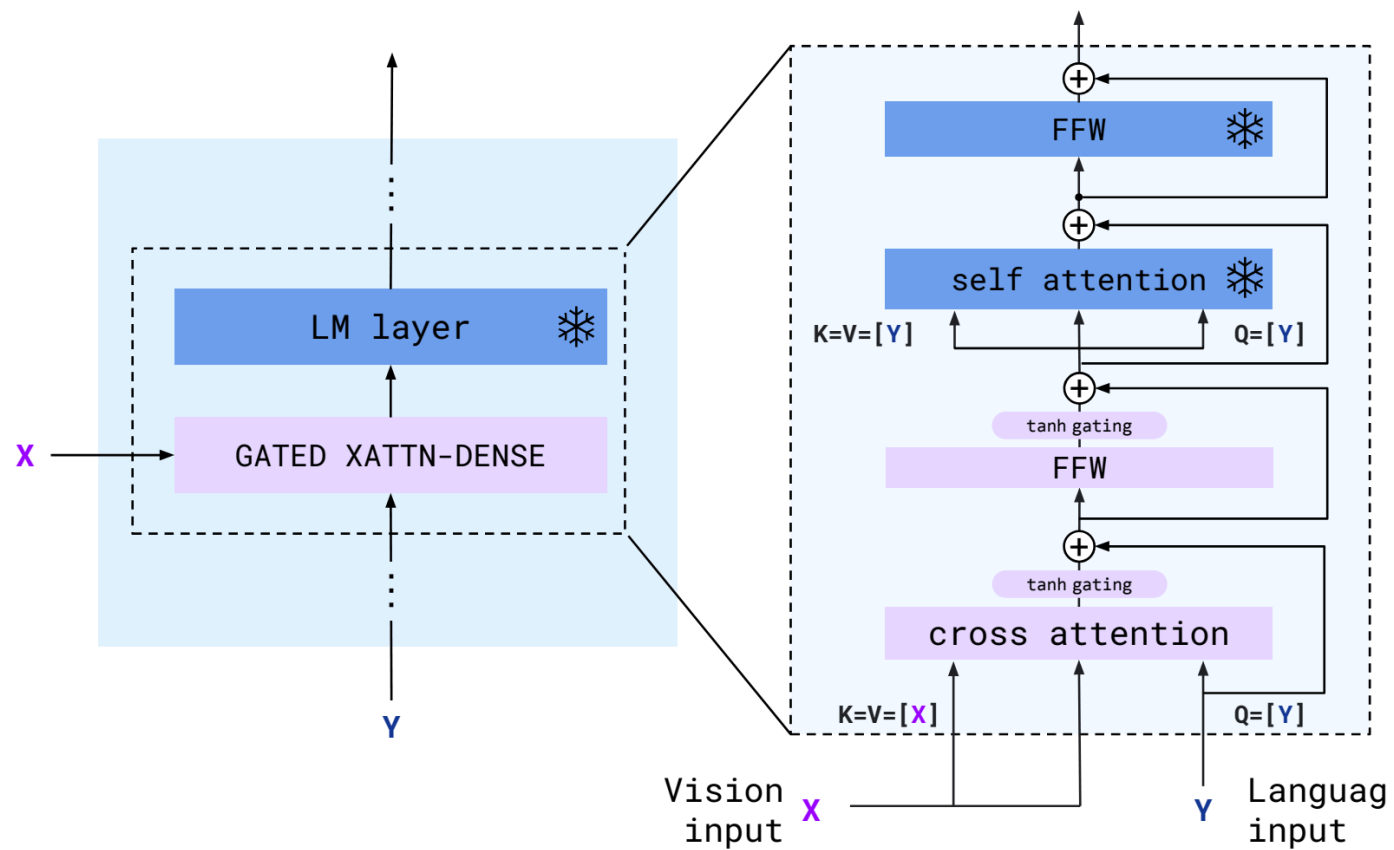
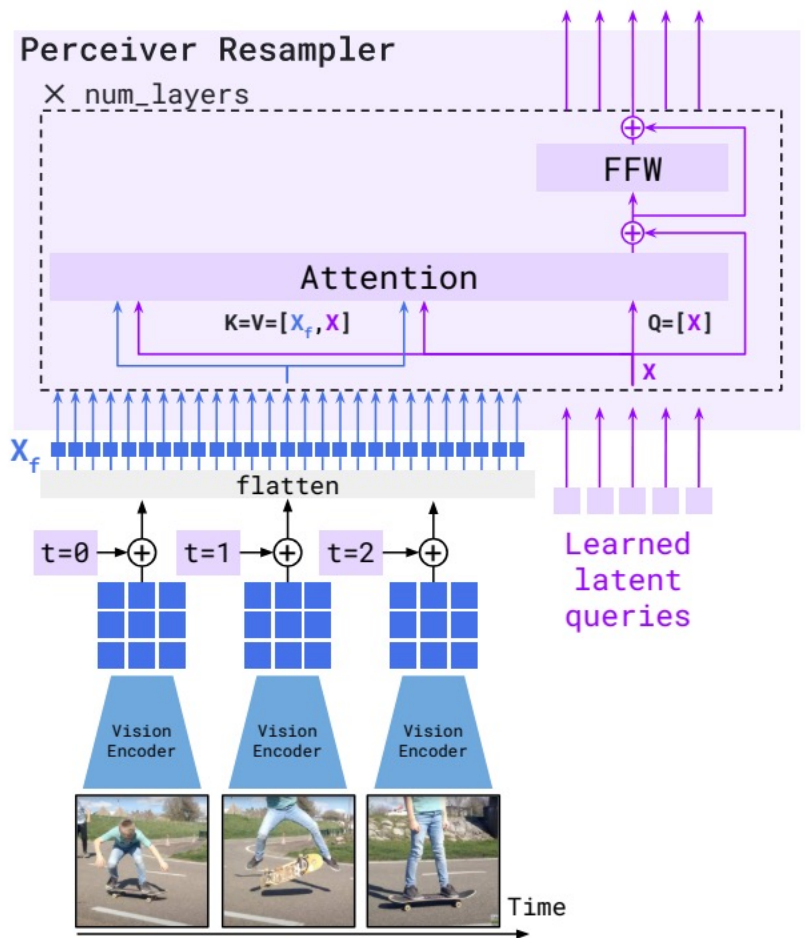




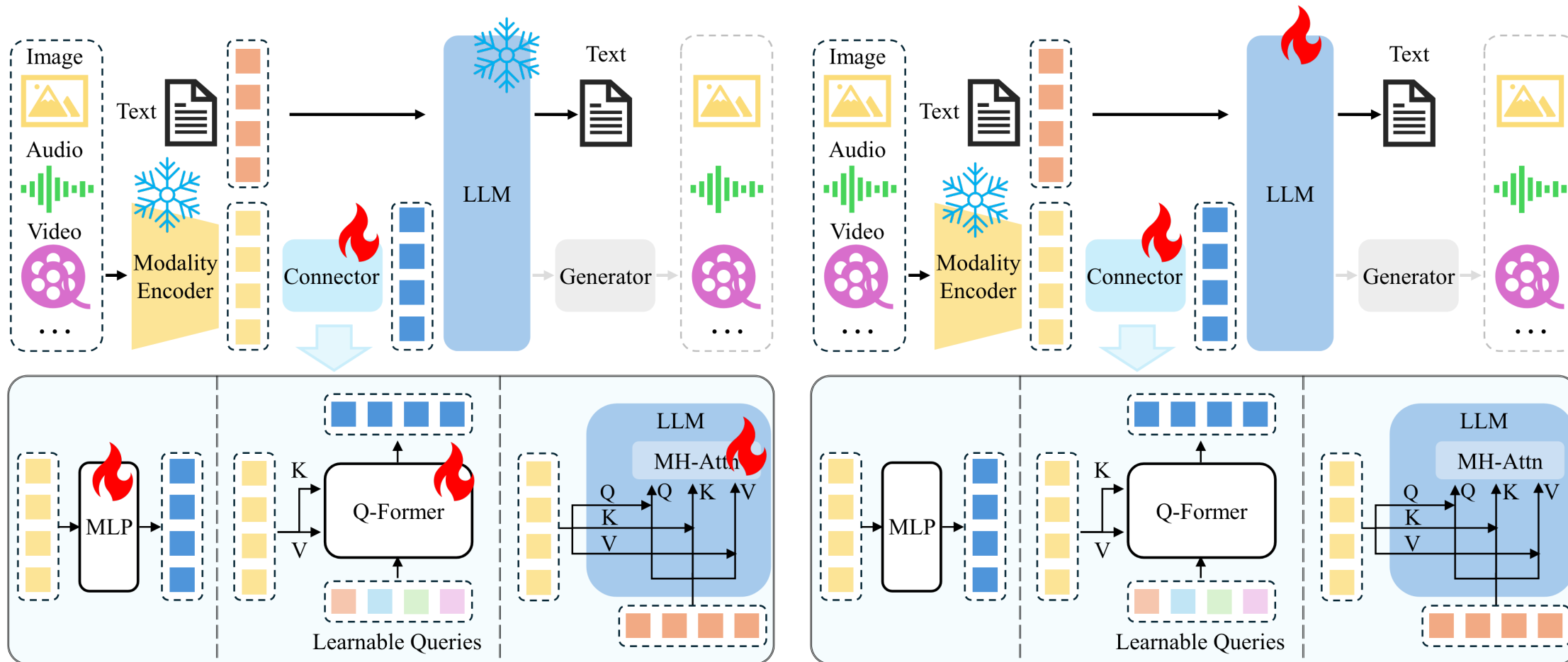


Flamingo





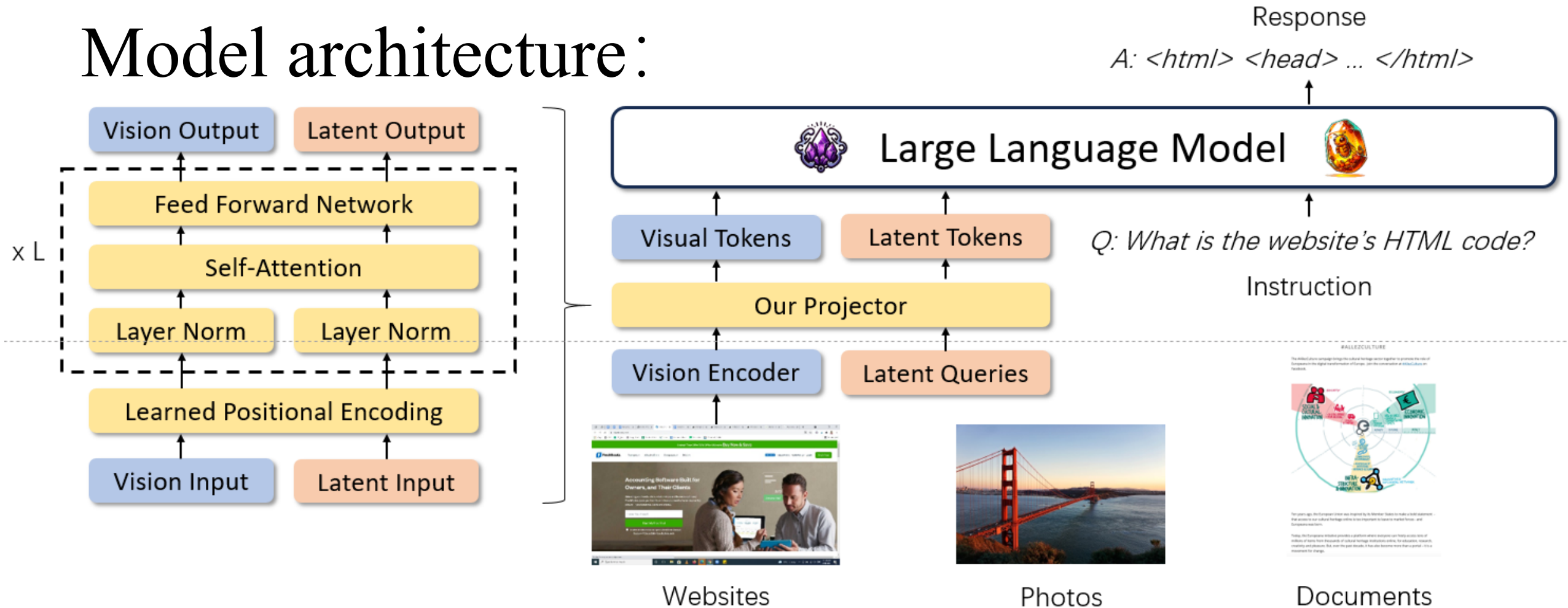
Training Details



- image embeddings for perception and reasoning should be different
 - For example, reasoning often requires multiple perceptual image tokens
 - However, existing frameworks use only global tokens (like Flamingo and BLIP-2) or only local tokens (like LLaVA)
 - considering to use both global and local tokens as image embeddings
 - **Global** tokens from both **image** and **text**
 - **Local** tokens from the **image**

- Explicitly modeling **relation tokens** as inputs of LLM
- Inspired on scene graph domain, relation tokens extracted from object tokens can bring rich information
 - Like Flamingo, we add new modules to model object tokens and relation tokens
 - New learnable object and relation tokens are given as inputs of LLM

Model architecture:



- Aligning text and visual embeddings using image-text pairs.
- Refining the model's instruction-following capability using multi-modal tasks and natural language processing data.