# VLM Architectures 핵심 요약

## 요약

(94개의 Model)

# 1. Unicoder-VL (AAAI 2020)

- **Universal** Encoder

    - Universal = **Vision + Language**

- 3개의 pretraining task

    - Masked language modeling (**MLM**)

    - Masked object classification (**MOC**)

    - Visual-linguistic matching (**VLM**)

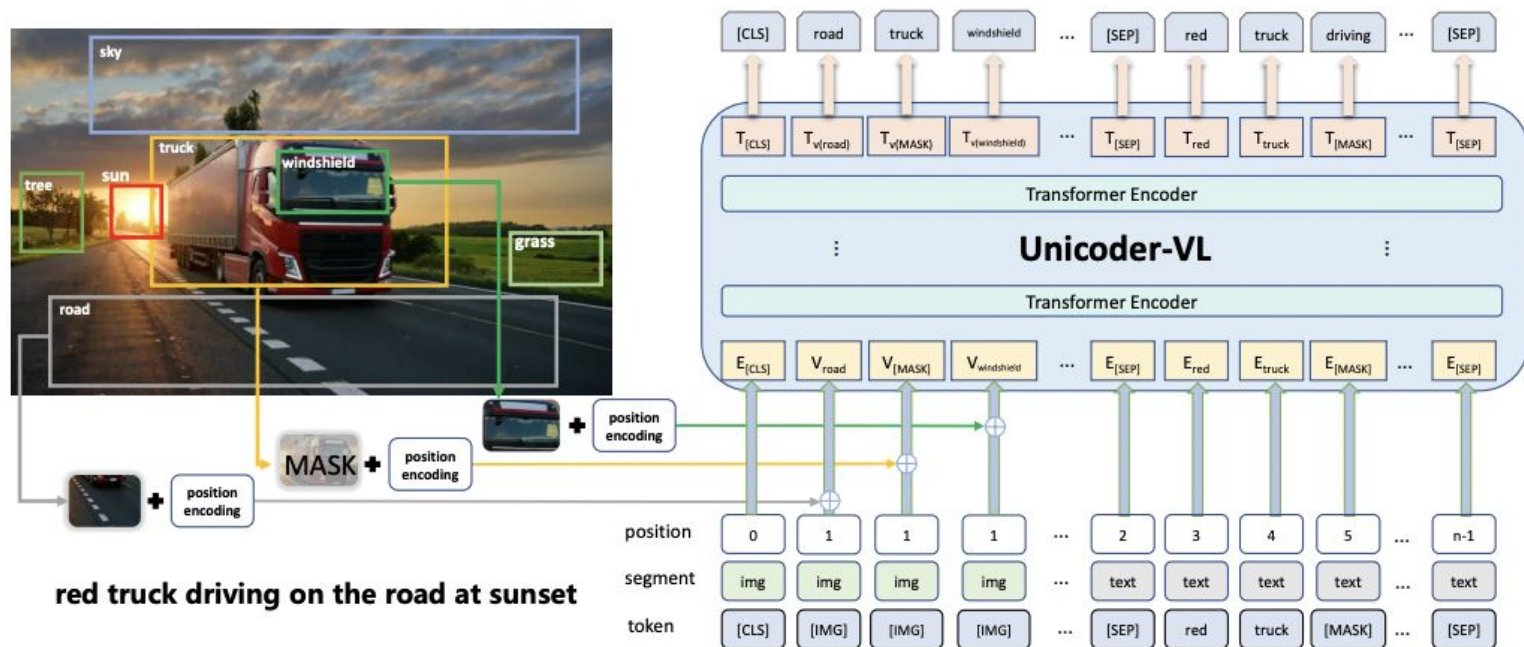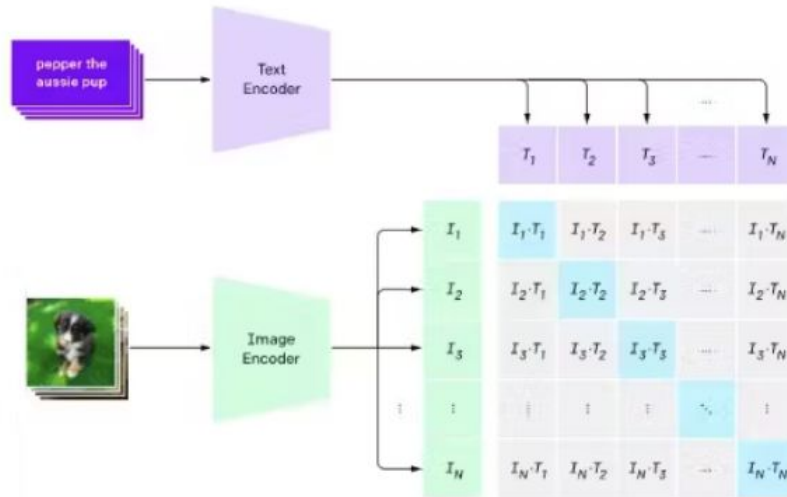# 1. Unicoder-VL (AAAI 2020)



Figure 1: Illustration of Unicoder-VL in the context of an object and text masked token prediction, or *cloze*, task. Unicoder-VL contains multiple Transformer encoders which are used to learn viusal and linguistic representation jointly.
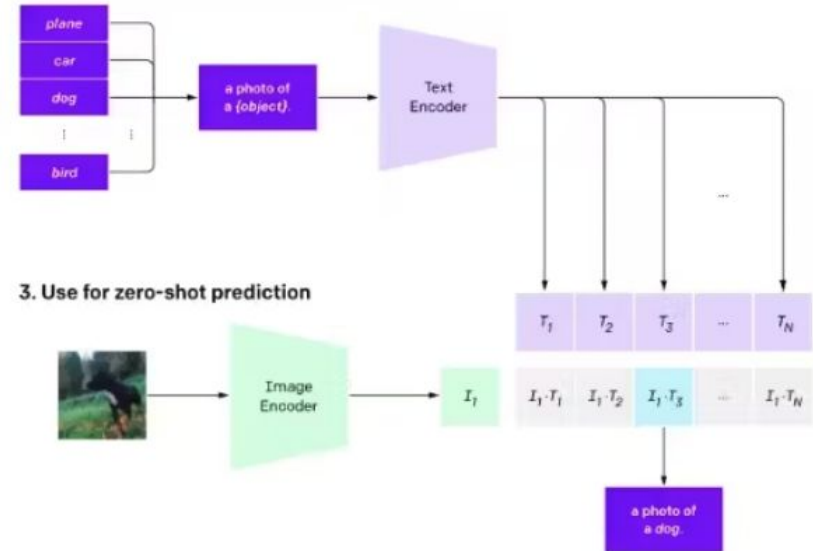
# 2. CLIP (ICML 2021)

- Language & Image "**Contrastive Learning**"

# 3. MetaCLIP (ICLR 2024)

- Refine the **data curation** process

- Leverage CLIP-derived metadata

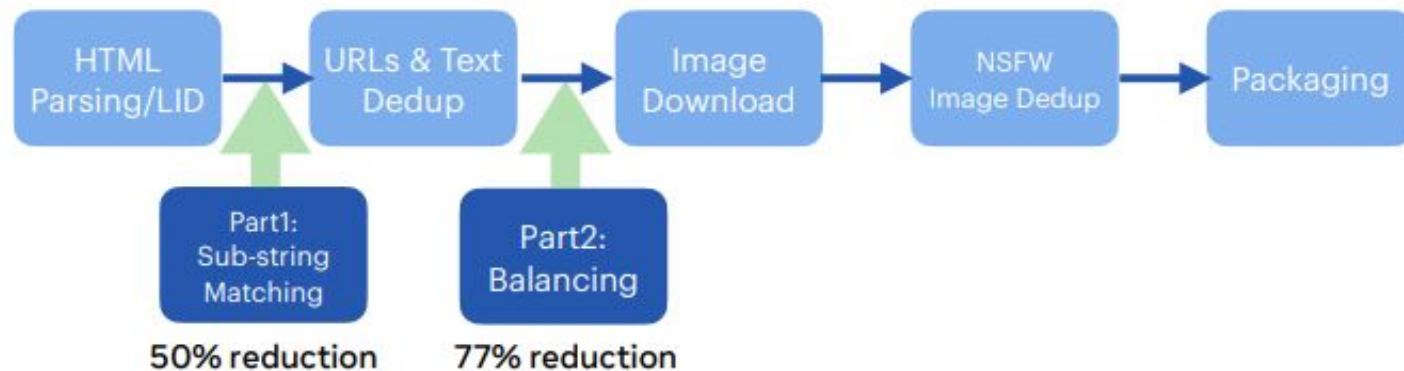- To create a **balanced & high-quality dataset** from vast sources



Figure 5: Case study: Curation implementation in our data pipeline.

# 4. Alpha-CLIP (CVPR 2024)

- CLIP + **Region awareness**

    - Addition of an **alpha channel** to the CLIP encoder

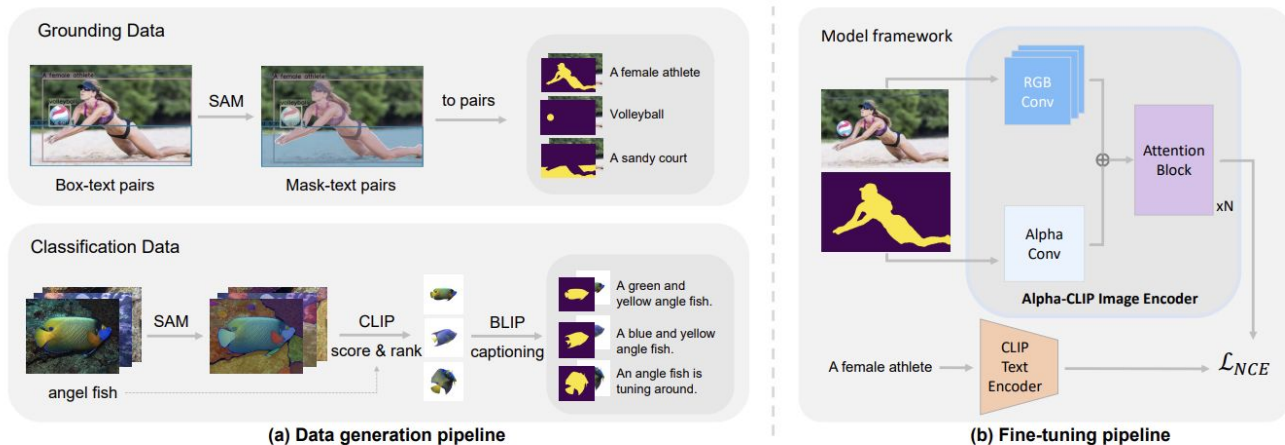- Allows for detailed **segmentation** & region-specific processing



Figure 3. **The pipeline of our data generation method and model architecture**. **(a)** Our method generates millions of RGBA-region text pairs. **(b)** Alpha-CLIP modifies the CLIP image encoder to take an additional alpha channel along with RGB.

# 5. GLIP (CVPR 2022)

- GLIP = **Grounded** + LIP

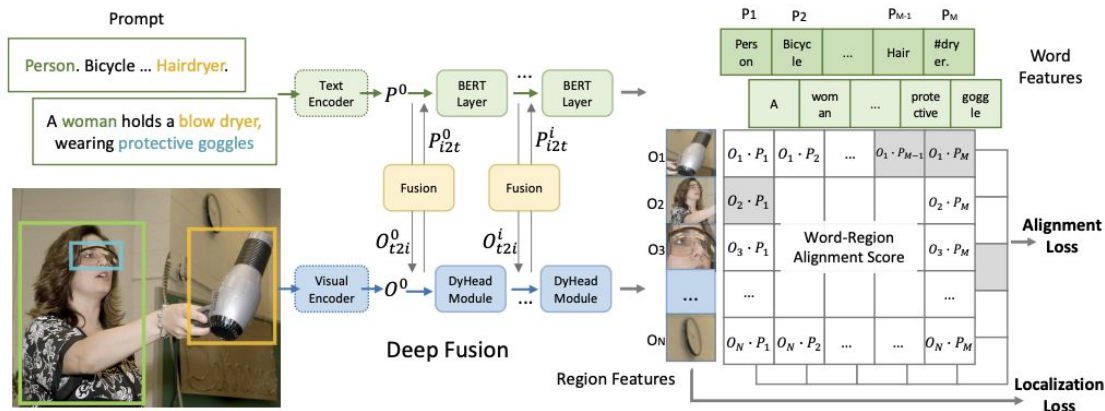- Pretraining by unifying **object detection & phrase grounding**



Figure 2. A unified framework for detection and grounding. Unlike a classical object detection model which predicts a categorical class for each detected object, we reformulate detection as a grounding task by aligning each region/box to phrases in a text prompt. GLIP jointly trains an image encoder and a language encoder to predict the correct pairings of regions and words. We further add the cross-modality deep fusion to early fuse information from two modalities and to learn a language-aware visual representation.

# 6. SigLIP (CVPR 2022)

- Simple pairwise **sigmoid** loss

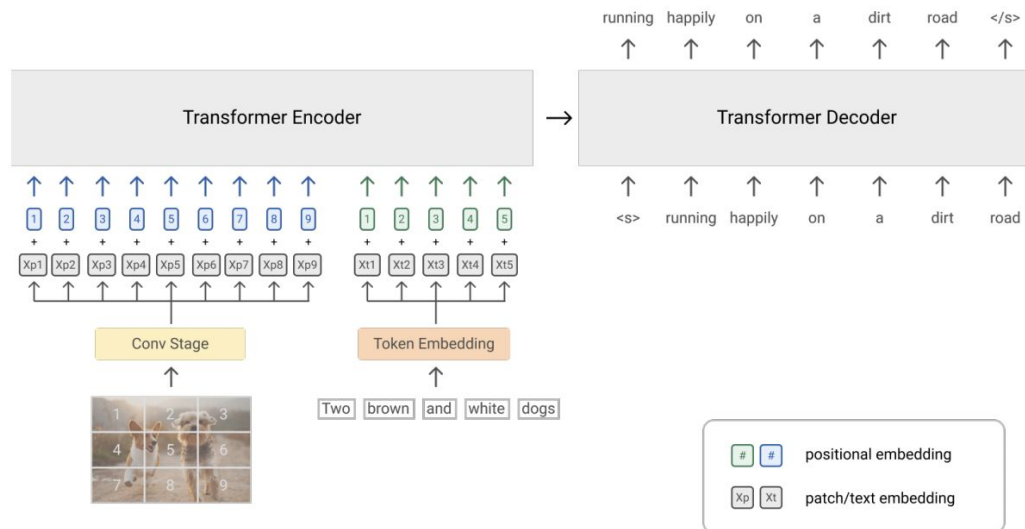- Scalable training with **large batch sizes**

CLIP: SoftMax

$$-\frac{1}{2|\mathcal{B}|}\sum_{i=1}^{|\mathcal{B}|}\left(\overbrace{\log\frac{e^{t\mathbf{x}_i\cdot\mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|}e^{t\mathbf{x}_i\cdot\mathbf{y}_j}}}^{\text{image}\rightarrow\text{text softmax}}+\overbrace{\log\frac{e^{t\mathbf{x}_i\cdot\mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|}e^{t\mathbf{x}_j\cdot\mathbf{y}_i}}}^{\text{text}\rightarrow\text{image softmax}}\right)$$

SigLIP: Sigmoid

$$-\frac{1}{|\mathcal{B}|}\sum_{i=1}^{|\mathcal{B}|}\sum_{j=1}^{|\mathcal{B}|}\log\underbrace{\frac{1}{1+e^{z_{ij}(-t\mathbf{x}_i\cdot\mathbf{y}_j+b)}}}_{\mathcal{L}_{ij}}$$

# 7. SimVLM (ICLR 2022)

- **Prefix**LM

  - X = **(Prefix) Image** + Text

  - Y = Text

- Task: **Next word prediction**

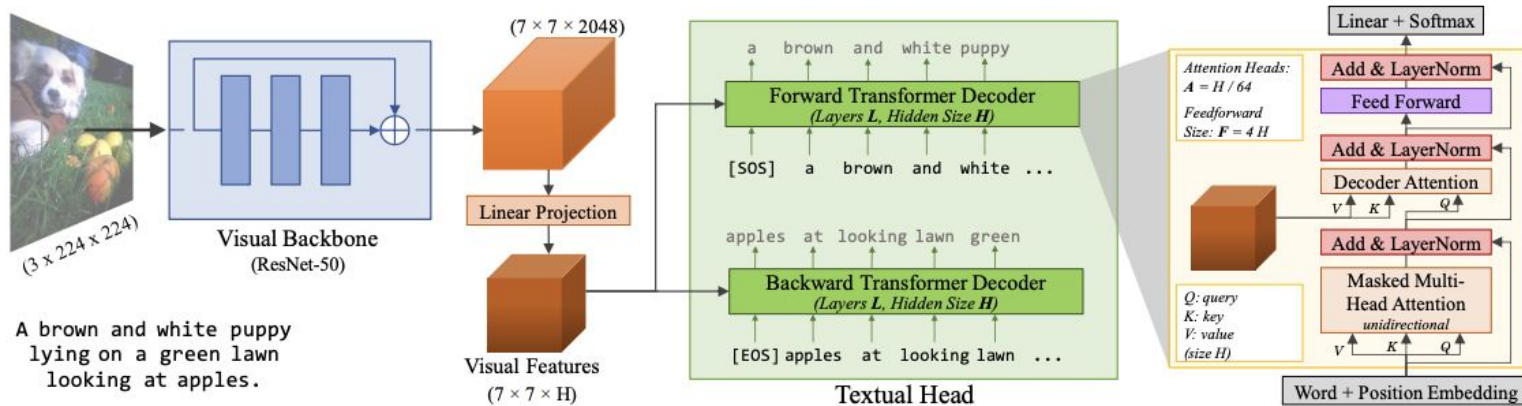# 8. VirTex (CVPR 2021)

- Pretraining task = **Image Captioning**



Figure 3: **VirTex pretraining setup:** Our model consists of a *visual backbone* (ResNet-50), and a *textual head* (two uni-directional Transformers). The visual backbone extracts image features, and textual head predicts captions via bidirectional language modeling (*bicaptioning*). The Transformers perform masked multiheaded self-attention over caption features, and multiheaded attention over image features. Our model is trained end-to-end from scratch. After pretraining, the visual backbone is transferred to downstream visual recognition tasks.

# 9. Frozen (NeurIPS 2021)

- PrefixLM (e.g., SimVLM)

  - X = (Prefix) Image + Text

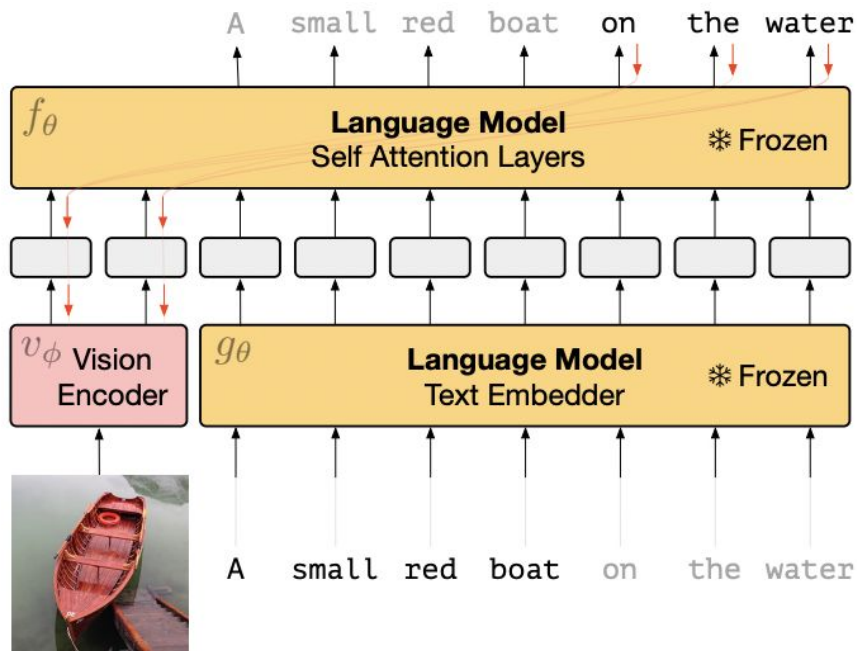  - Y = Text

- **Frozen PrefixLM**

  - Use "Pretrained" LLM



Figure 2: Gradients through a frozen language model's self attention layers are used to train the vision encoder.

# 9. Frozen (NeurIPS 2021)



(a) **0-shot VQA**   (b) **1-shot outside-knowledge VQA**   (c) **Few-shot image classification**
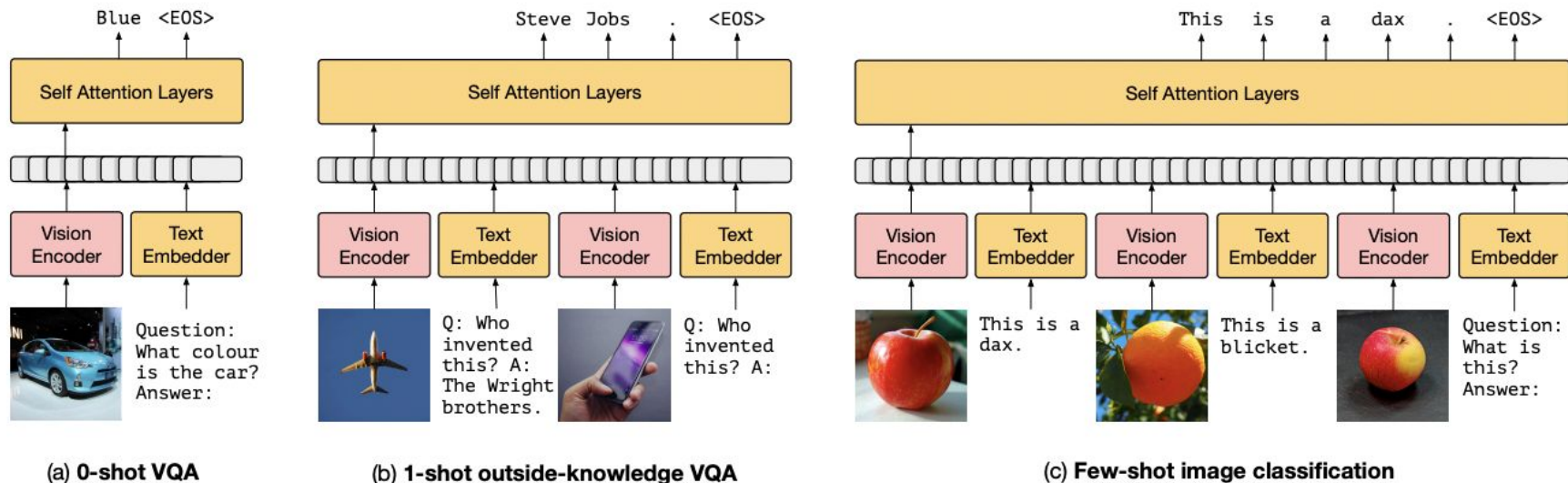
Figure 3: Inference-Time interface for *Frozen*. The figure demonstrates how we can support (a) visual question answering, (b) outside-knowledge question answering and (c) few-shot image classification via in-context learning.

# 10. Flamingo (NeurIPS 2022)

- Model Architecture

    - LLM (**freeze**): Chinchilla

    - Vision (**freeze**): CLIP

- **Arbitrarily interleaved** visual & textual tokens

- **Cross-attention** for **multimodal fusion**

- **Perceiver Resampler**

    - Input: Visual features

    - Output: **Fixed number** of visual tokens
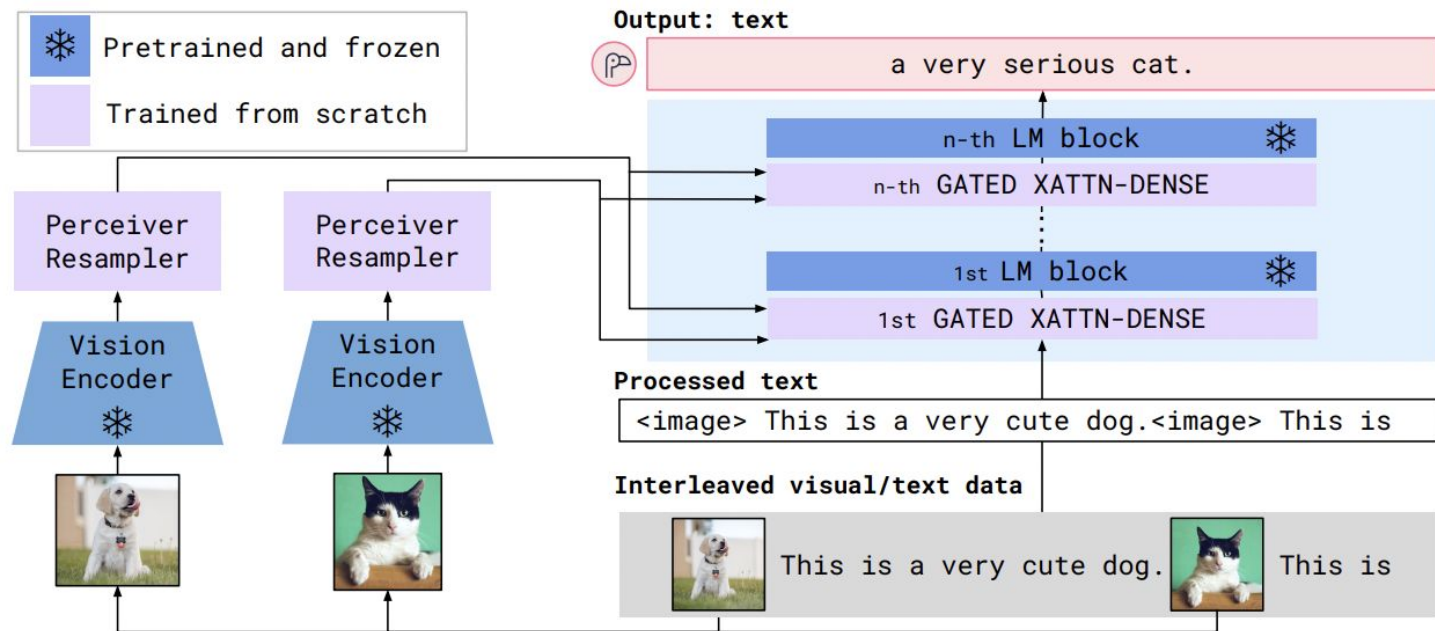
# 10. Flamingo (NeurIPS 2022)



Figure 3: **Flamingo architecture overview.** Flamingo is a family of visual language models (VLMs) that take as input visual data interleaved with text and produce free-form text as output.

# 11. VisualGPT (CVPR 2022)

- **Cross-attention** for **multimodal fusion**

- Architecture

  - Encoder: Visual encoder

  - Decoder: LLM

- **Self-resurrecting** enc-dec attention
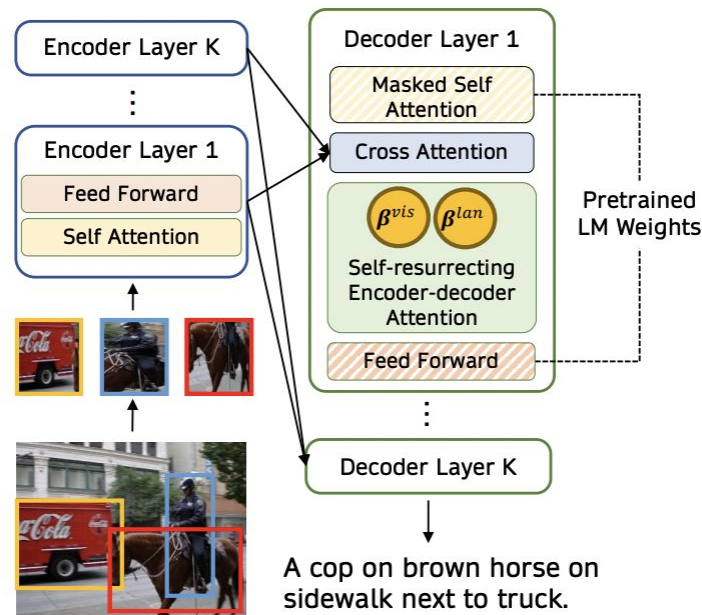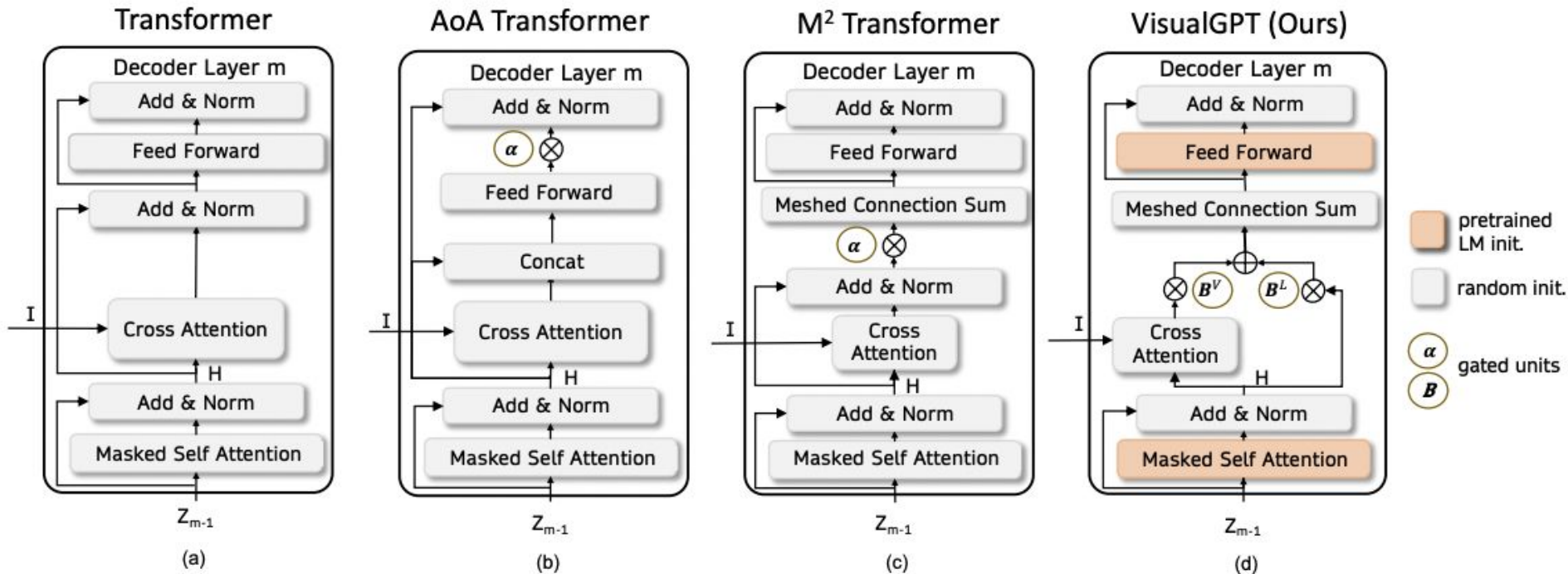
  - To adapt LLM for visual tasks



Figure 1. Our VisualGPT model transfers the knowledge from a pre-trained language model to the caption decoder. A self-resurrecting encoder-decoder attention is designed to connect the multi-level visual features and caption decoder.

# 11. VisualGPT (CVPR 2022)

# 12. VisualBERT (ACL 2020)
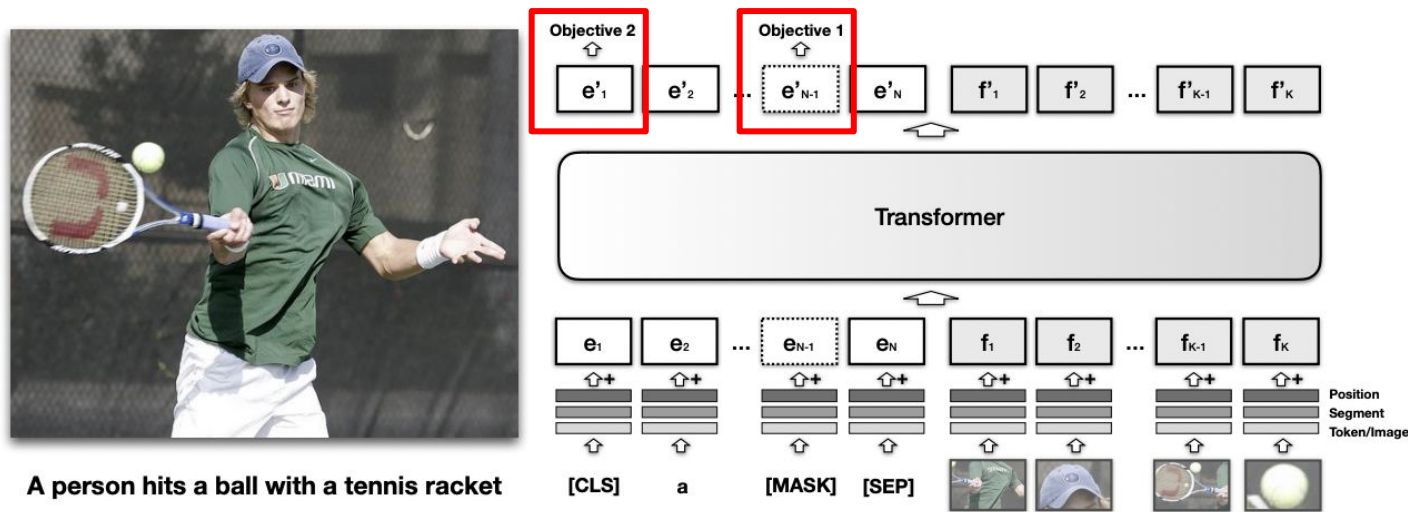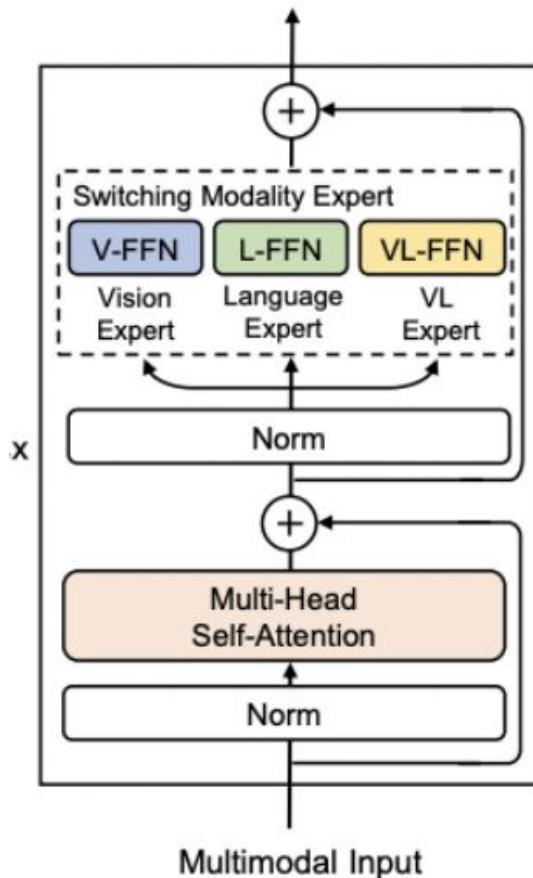
- Two pretraining tasks = **MLM + ITM**



Figure 2: The architecture of VisualBERT. Image regions and language are combined with a Transformer to allow the self-attention to discover implicit alignments between language and vision. It is pre-trained with a masked language modeling (Objective 1), and sentence-image prediction task (Objective 2), on caption data and then fine-tuned for different tasks. See §3.3 for more details.

# 13. VLMo (NeurIPS 2022)

- **MoME** (Mixture-of-modality-experts)

    - **각 Modality 별**로 expert 존재

    - V용 / L용 / VL용

- 두 가지 역할 모두 가능!

    - **Dual** encoder: 각각 embedding

    - **Fusion** encoder: 같이 embedding

# 13. VLMo (NeurIPS 2022)

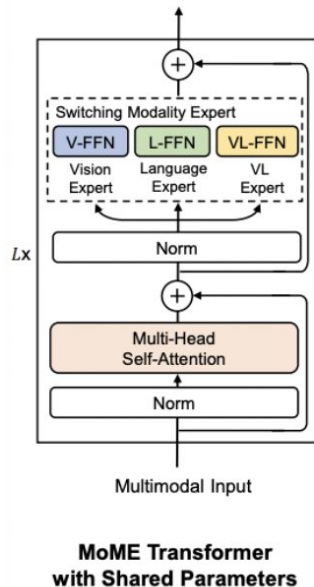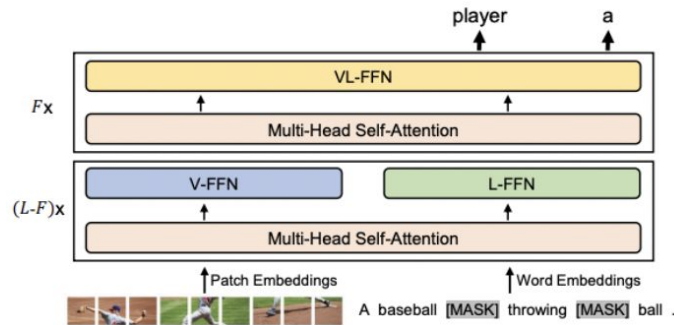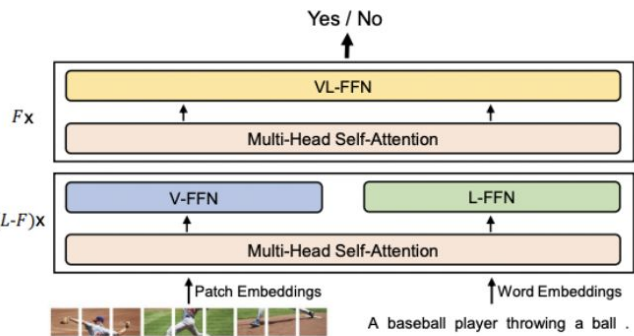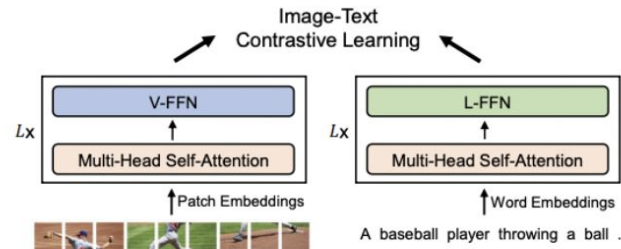- 세 가지 pretraining

  - **ITC, ITM, LM**
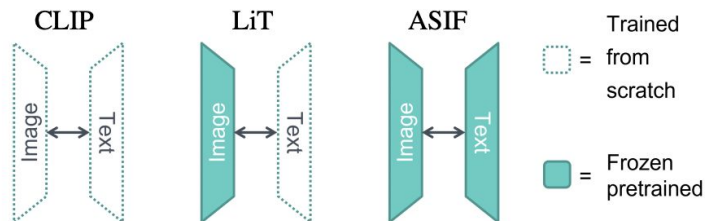
# 14. ASIF (NeurIPS 2023)

**- Training-free**



Figure 1: ASIF is a simple recipe to align the representations of two frozen pre-trained encoders.



Figure 4: **Zero shot classification with ASIF.** In this example we determine the best caption for the image $x^*$ from the two candidates, $y_i$ and $y_j$. **a.** Compute and store the embeddings of the multimodal dataset and the test samples. **b.** Compute the relative representation of the test image and the candidate captions. **c.** We consider the relative representation of $x^*$ with respect to the image collection $x_1, \ldots, x_n$ *as if* it was the relative representation of $y^*$ – the ideal caption for $x^*$ – with respect to the corresponding caption collection $y_1, \ldots, y_n$. **d.** We choose the candidate caption most similar to the ideal one.

# 15. ViLD (ICLR 2022)

- **Knowledge Distillation**

    - Teacher = (Pretrained) **Open-vocab** CLS model

    - Student = Two-stage detector

# 16. CoCa (TMLR 2022)

- CoCa = Contrastive Captioner

- 핵심: **2가지의 joint loss**

  - **Contrastive** loss (feat. CLIP)

  - **Captioning** loss (feat. SimVLM)

# 16. CoCa (TMLR 2022)



Figure 1: Overview of Contrastive Captioners (CoCa) pretraining as image-text foundation models. The pretrained CoCa can be used for downstream tasks including visual recognition, vision-language alignment, image captioning and multimodal understanding with zero-shot transfer, frozen-feature evaluation or end-to-end finetuning.

# 16. CoCa (TMLR 2022)



**Captioning Loss**

two dogs running in a field [/s]

**Multimodal Text Decoder**

Cross-Attention

attentional pooling → **Contrastive Loss** ← cls-token

**Image Encoder**     **Unimodal Text Decoder**
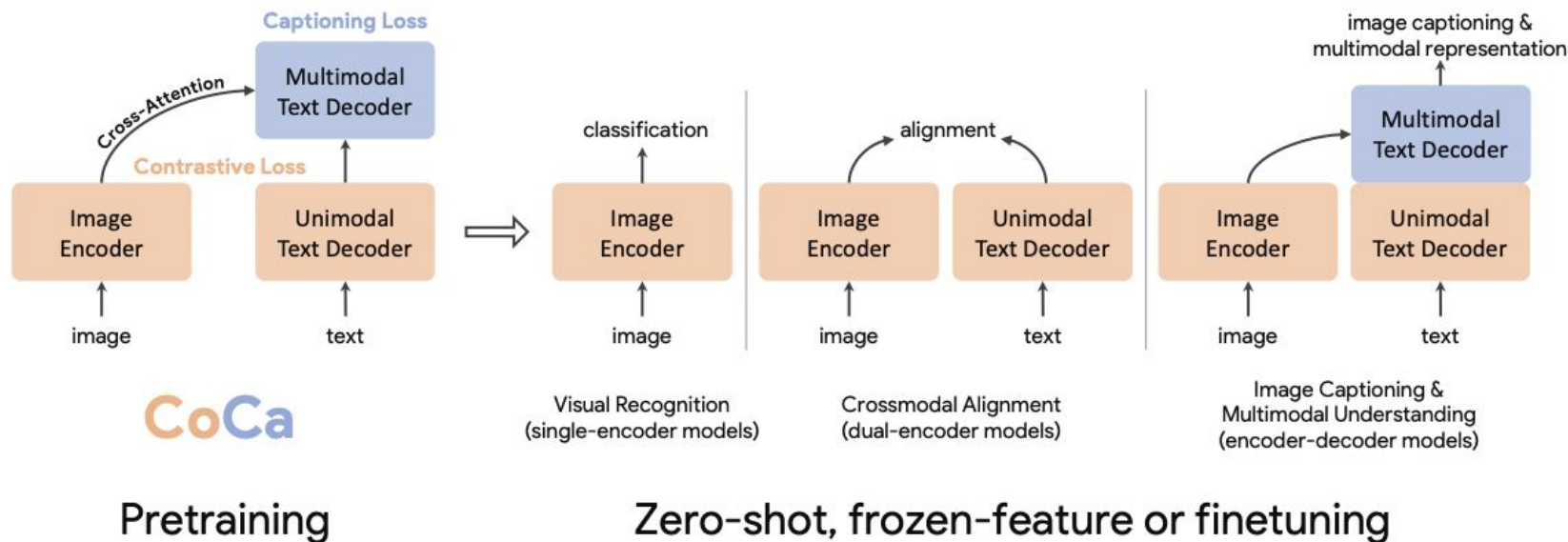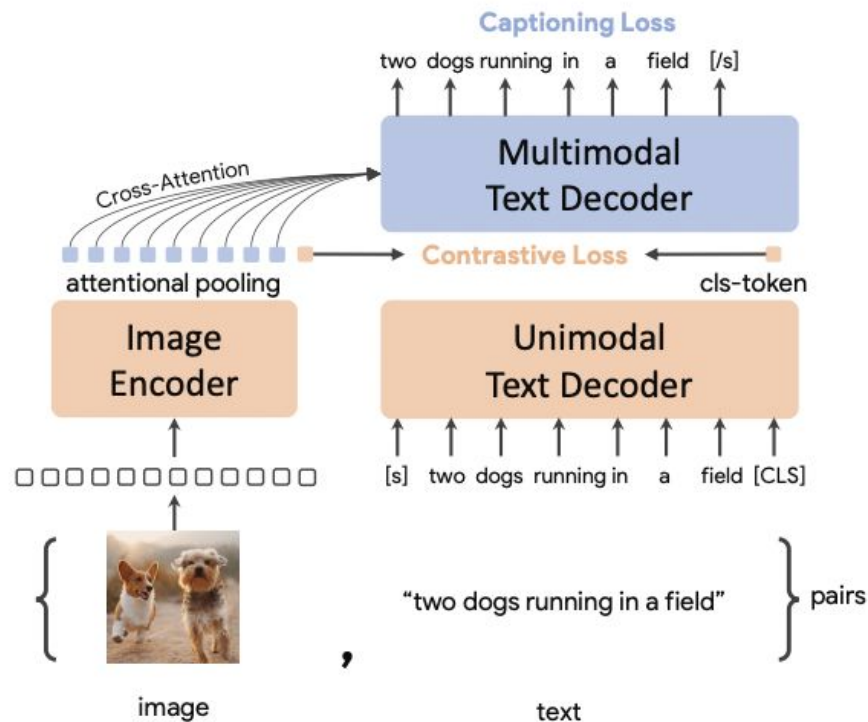
[s] two dogs running in a field [CLS]

"two dogs running in a field" } pairs

image     text

**Algorithm 1** Pseudocode of Contrastive Captioners architecture.

```
# image, text.ids, text.labels, text.mask: paired {image, text} data
# con_query: 1 query token for contrastive embedding
# cap_query: N query tokens for captioning embedding
# cls_token_id: a special cls_token_id in vocabulary

def attentional_pooling(features, query):
    out = multihead_attention(features, query)
    return layer_norm(out)

img_feature = vit_encoder(image) # [batch, seq_len, dim]
con_feature = attentional_pooling(img_feature, con_query) # [batch, 1, dim]
cap_feature = attentional_pooling(img_feature, cap_query) # [batch, N, dim]

ids = concat(text.ids, cls_token_id)
mask = concat(text.mask, zeros_like(cls_token_id)) # unpad cls_token_id
txt_embs = embedding_lookup(ids)
unimodal_out = lm_transformers(txt_embs, mask, cross_attn=None)
multimodal_out = lm_transformers(
    unimodal_out[:, :-1, :], mask, cross_attn=cap_feature)
cls_token_feature = layer_norm(unimodal_out)[:, -1:, :] # [batch, 1, dim]

con_loss = contrastive_loss(con_feature, cls_token_feature)
cap_loss = softmax_cross_entropy_loss(
    multimodal_out, labels=text.labels, mask=text.mask)
```

vit_encoder: vision transformer based encoder; lm_transformer: language-model transformers.
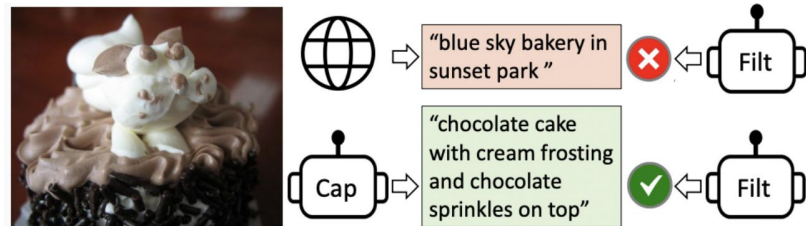
Figure 2: Detailed illustration of CoCa architecture and training objectives.

# 17. BLIP (ICML 2022)

- Proposal 2개

    - Model 측면: Multimodal mixture of Encoder-Decoder (**MED**)

    - Data 측면: **CapFilt**

- 세 가지 Pretraining

    - ITC loss: CL loss

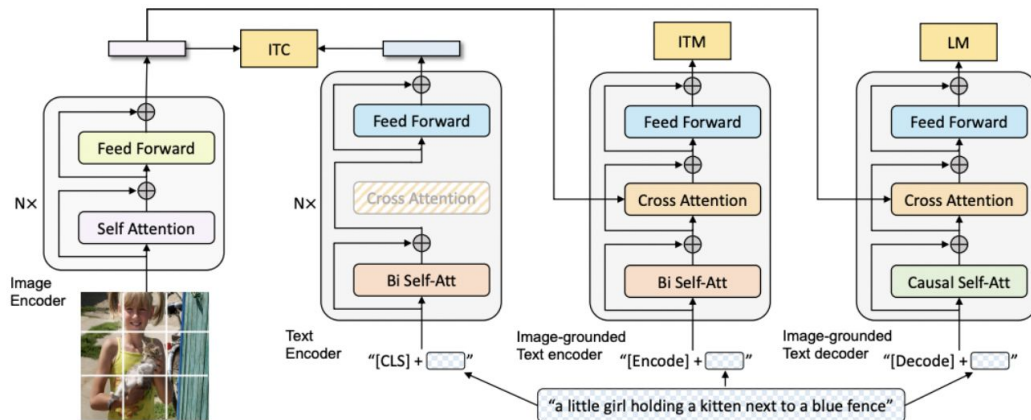    - ITM loss: Binary CE loss

    - LM: CE loss

# 17. BLIP (ICML 2022)

- [1] CapFilt: <span style="color:red">**Bootstrap the caption**</span>!

  - 배경: Noisy image-text pairs

  - (1) <span style="color:red">**Captioner**</span>: Synthetic caption 생성

    - Image-grounded text decoder

    - LM으로 fine-tune

  - (2) <span style="color:red">**Filter**</span>: Noisy caption 제거

    - Image grounded text encoder
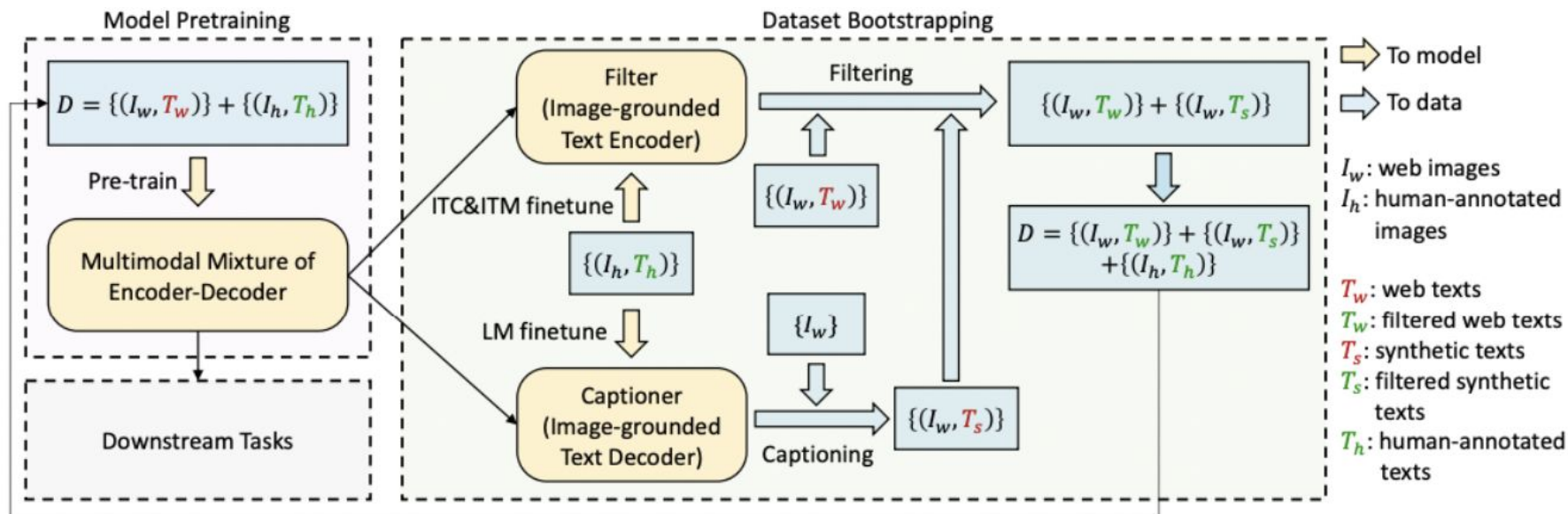
    - ITC & ITM으로 fine-tune

# 17. BLIP (ICML 2022)

- [2] **MED**: 아래의 3개 중 하나의 기능으로 작동 가능

    - Unimodal encoder

    - Image-grounded text encoder

    - Image-grounded text decoder

# 17. BLIP (ICML 2022)

- Learning framework of BLIP

# 18. BLIP-2 (ICML 2023)

- **Q(Querying)-Former**

- 목적: Modality gap을 bridge

- 2단계 학습

  - 1) (Frozen image 인코더)

    **VL** representation learning

  - 2) (Frozen language 인코더)
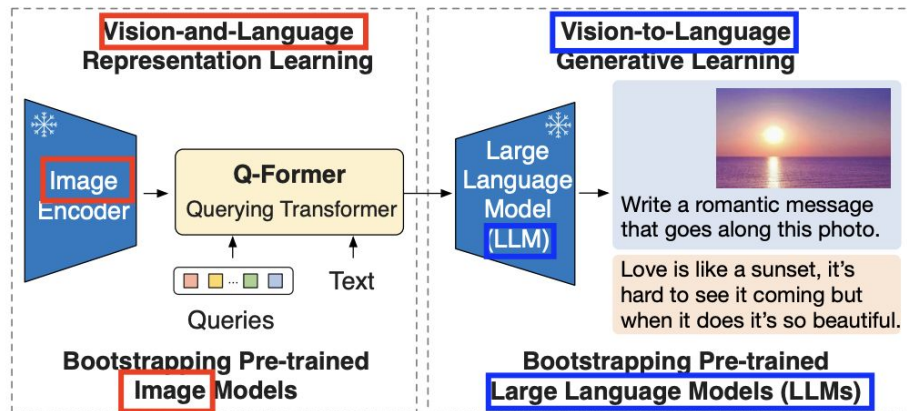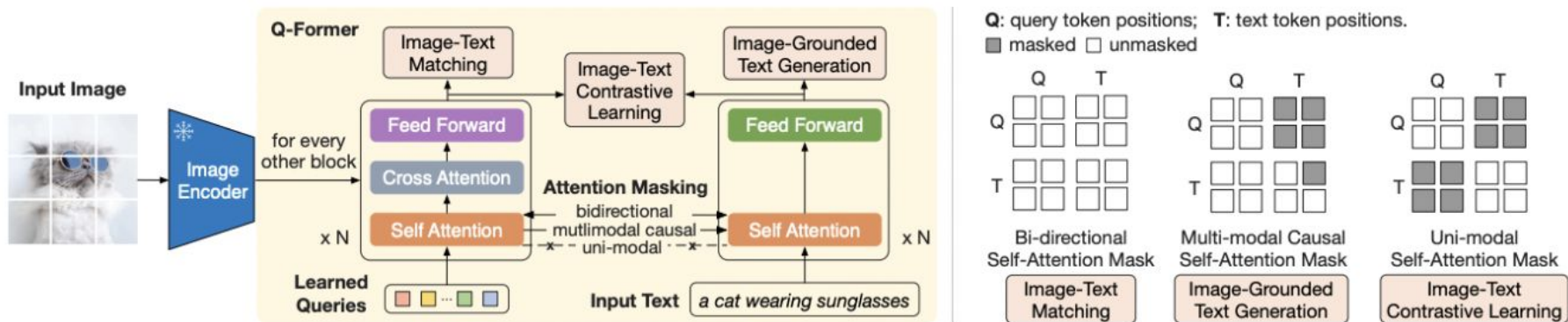
    **V-to-L** generation learning



*Figure 1.* Overview of BLIP-2's framework. We pre-train a lightweight Querying Transformer following a two-stage strategy to bridge the modality gap. The first stage bootstraps vision-language representation learning from a frozen image encoder. The second stage bootstraps vision-to-language generative learning from a frozen LLM, which enables zero-shot instructed image-to-text generation (see Figure 4 for more examples).
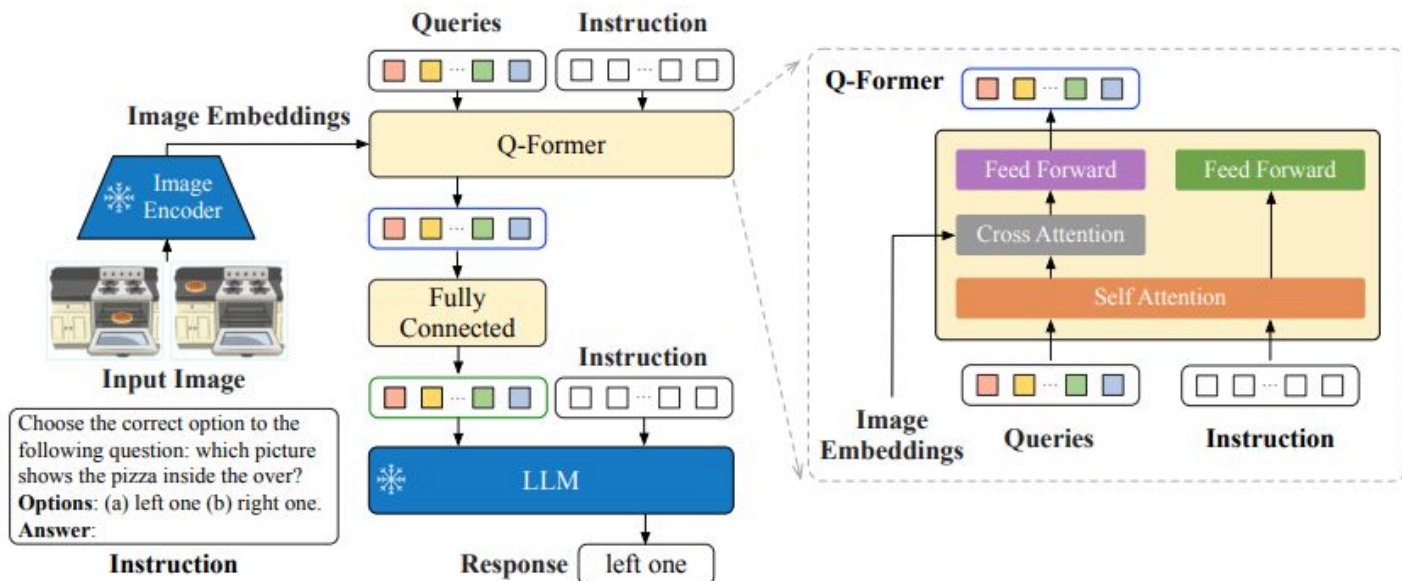
# 18. BLIP-2 (ICML 2023)

- 2개의 transformer = **Image** transformer + **Text** transformer

- Pretraining

    - (BLIP과 비슷하게) **ITC + ITG + ITM**

    - 각 objective별로 다른 attention masking strategy

# 19. InstructBLIP (NeurIPS 2023)

- BLIP2를 개선함.

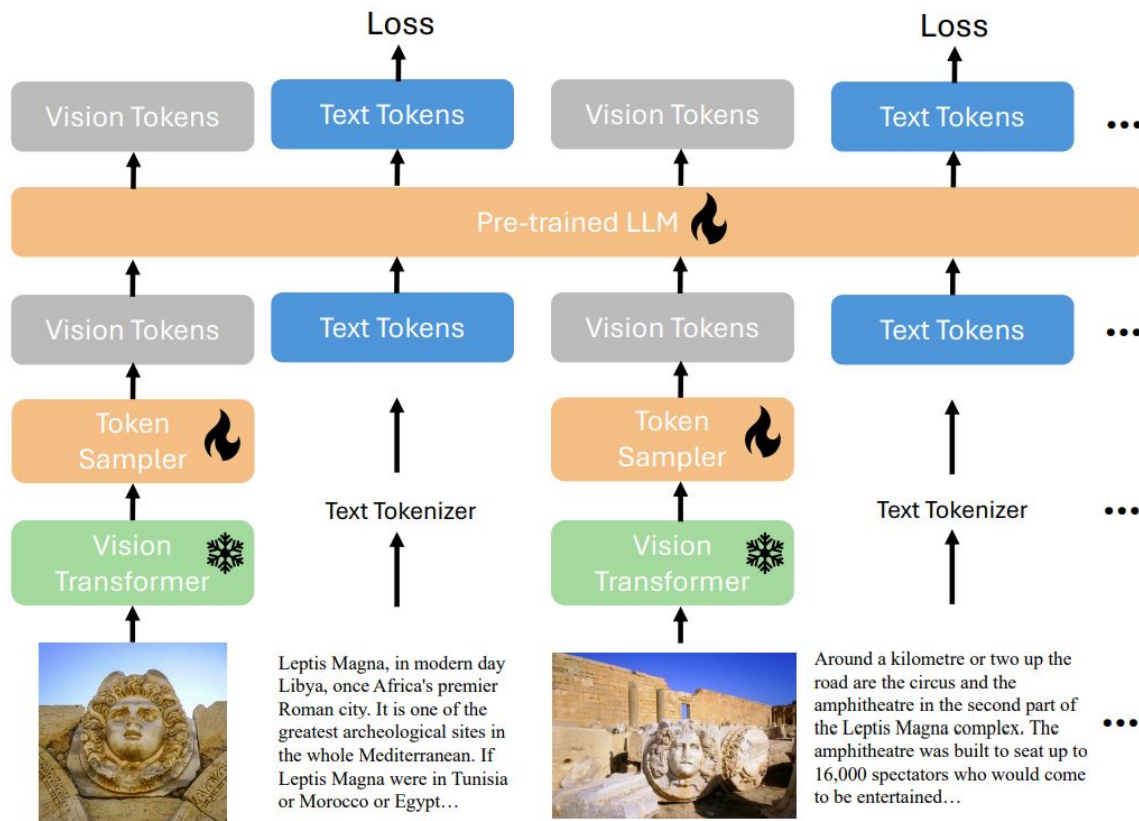- **Q-Former를 instruction tuning**함

# 20. xGen-MM (BLIP-3) (ICML 2023)

- Diverse, large-scale, high-quality multimodal data

- **BLIP-2**:

    - (a) **Q-Former**

    - (b) **Multiple pretraining objectives**

- **BLIP-3**:

    - (a) More scalable vision token sampler (**perceiver resampler**)

    - (b) Unifying the training objective to a **single autoregressive loss**

# 20. xGen-MM (BLIP-3) (ICML 2023)

# 21. LLaVA (NeurIPS 2023)

- Goal: **Multimodal 버전의 Instruction-following dataset**을 생성하자!
- How: **"text" 데이터만**을 사용해서 (feat. GPT-4)
- (1) + (2) 두 가지 정보를 준 뒤, 문제/답변 생성을 LLM에 요청
    - **(1) Caption**
    - **(2) Bounding boxes**
- 생성한 3종류의 데이터셋:
    - Conversation
    - Detailed description
    - Complex Reasoning



**Context type 1: Captions**
A group of people standing outside of a black vehicle with various luggage.
Luggage surrounds a vehicle in an underground parking area
People try to fit all of their luggage in an SUV.
The sport utility vehicle is parked in the public garage, being packed for a trip
Some people with luggage near a van that is transporting it.
**Context type 2: Boxes**
person: [0.681, 0.242, 0.774, 0.694], person: [0.63, 0.222, 0.686, 0.516], person: [0.444, 0.233, 0.487, 0.34], backpack: [0.384, 0.696, 0.485, 0.914], backpack: [0.755, 0.413, 0.846, 0.692], suitcase: [0.758, 0.413, 0.845, 0.69], suitcase: [0.1, 0.497, 0.173, 0.579], bicycle: [0.282, 0.363, 0.327, 0.442], car: [0.786, 0.25, 0.848, 0.322], car: [0.783, 0.27, 0.827, 0.335], car: [0.86, 0.254, 0.891, 0.3], car: [0.261, 0.101, 0.787, 0.626]

# 21. LLaVA (NeurIPS 2023)

- Model architecture

  - LLM: Vicuna

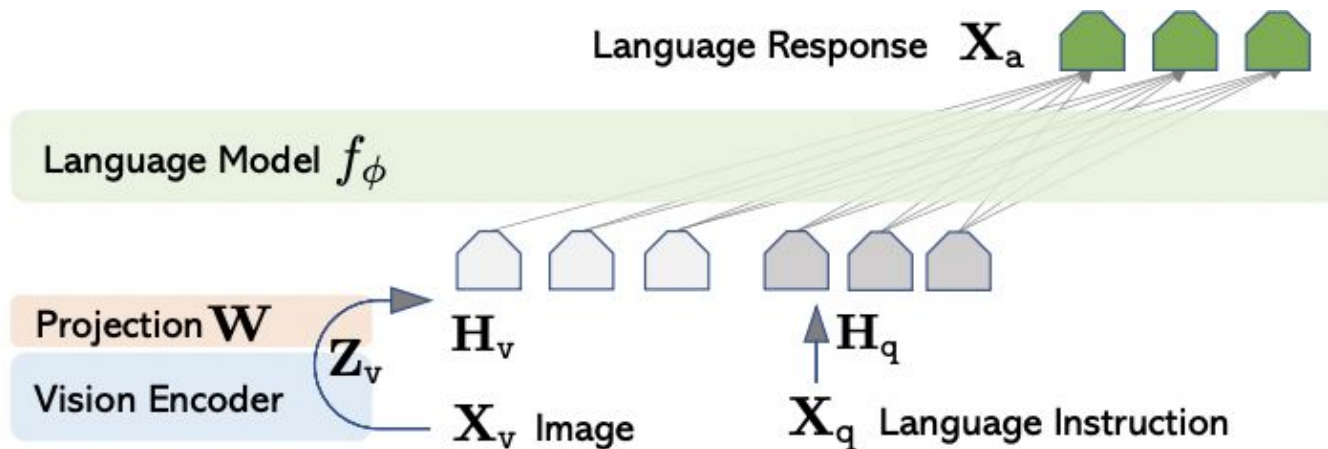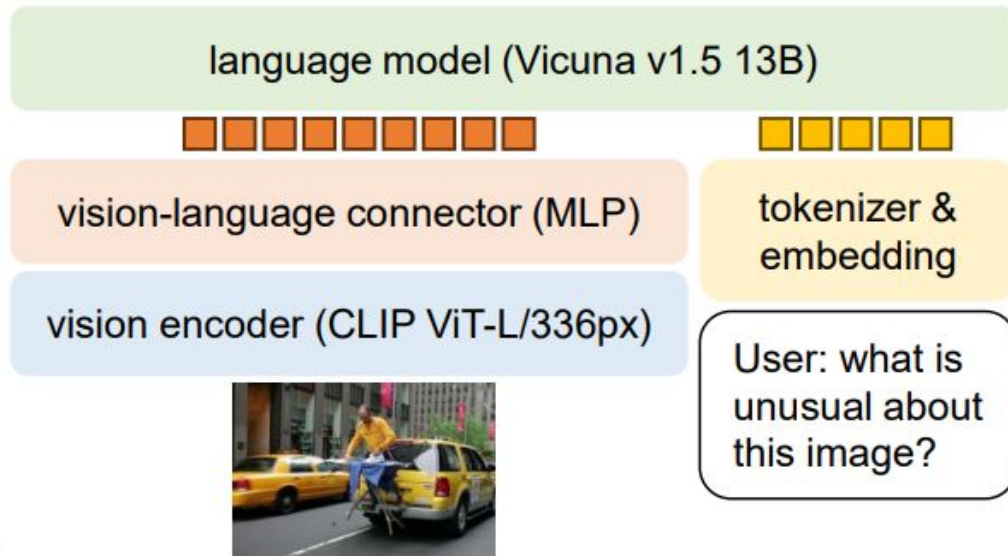  - Vision encoder: CLIP (ViT-L/14) + **Linear projection layer**



Figure 1: LLaVA network architecture.

# 22. LLaVA 1.5 (CVPR 2024)

- LLaVA: **Linear** projection layer

- LLaVA 1.5: **MLP** projection layer

# 23. LLaVA-NeXT (LLaVA 1.6) (blog)

- LLaVA1.5 + alpha

  - **Dynamic high resolution**: High-resolution image processing

  - OCR capabilities
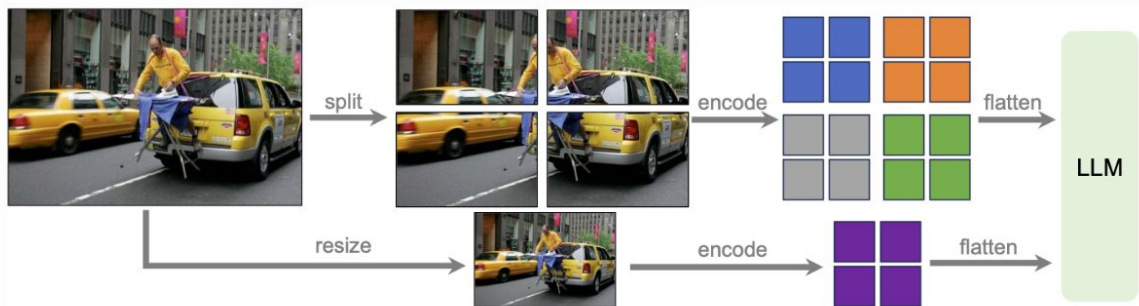
  - Efficient deployment & inference



*Illustration of dynamic high resolution scheme: a grid configuration of $2 \times 2$*

# 24. u-LLaVA (arxiv 2023)

- Projector-based architecture

- Unifies multi-modal tasks

    - Integrates **pixel**, **regional**, and **global** features

    - By connecting (specialized) experts with (central) LLM

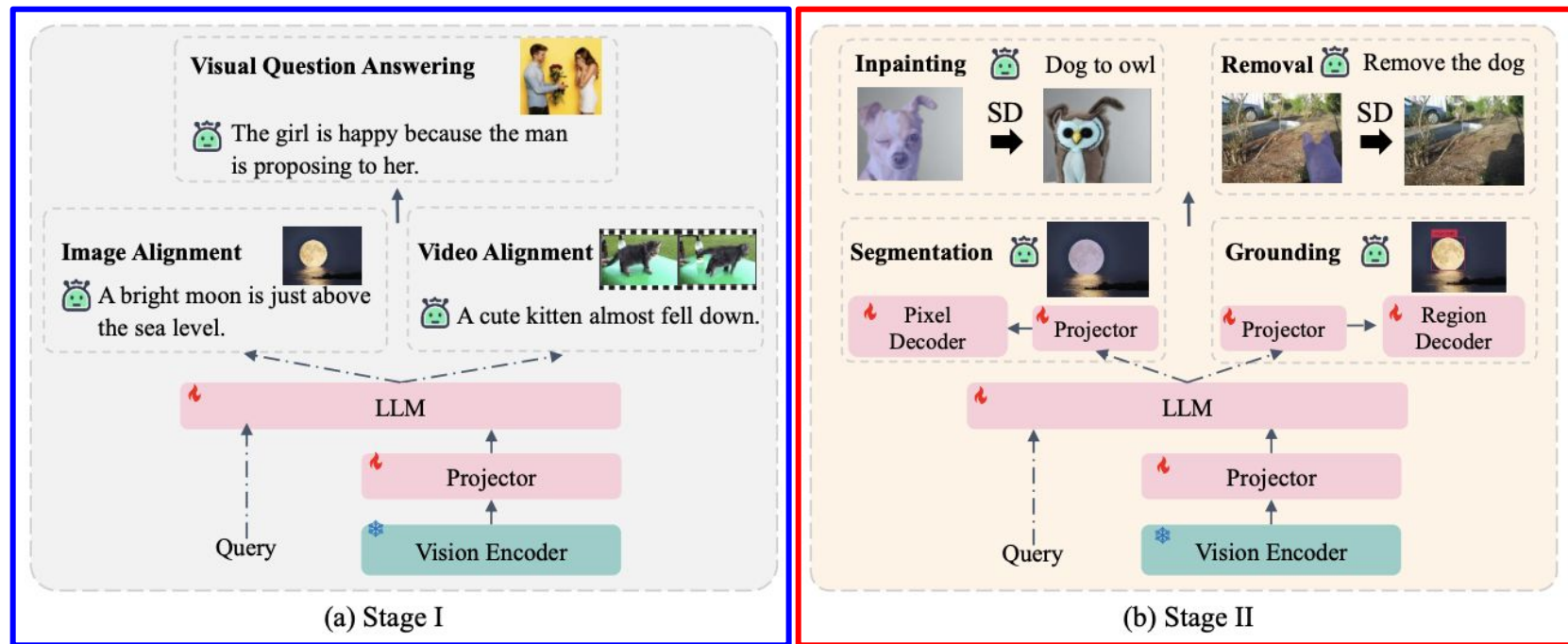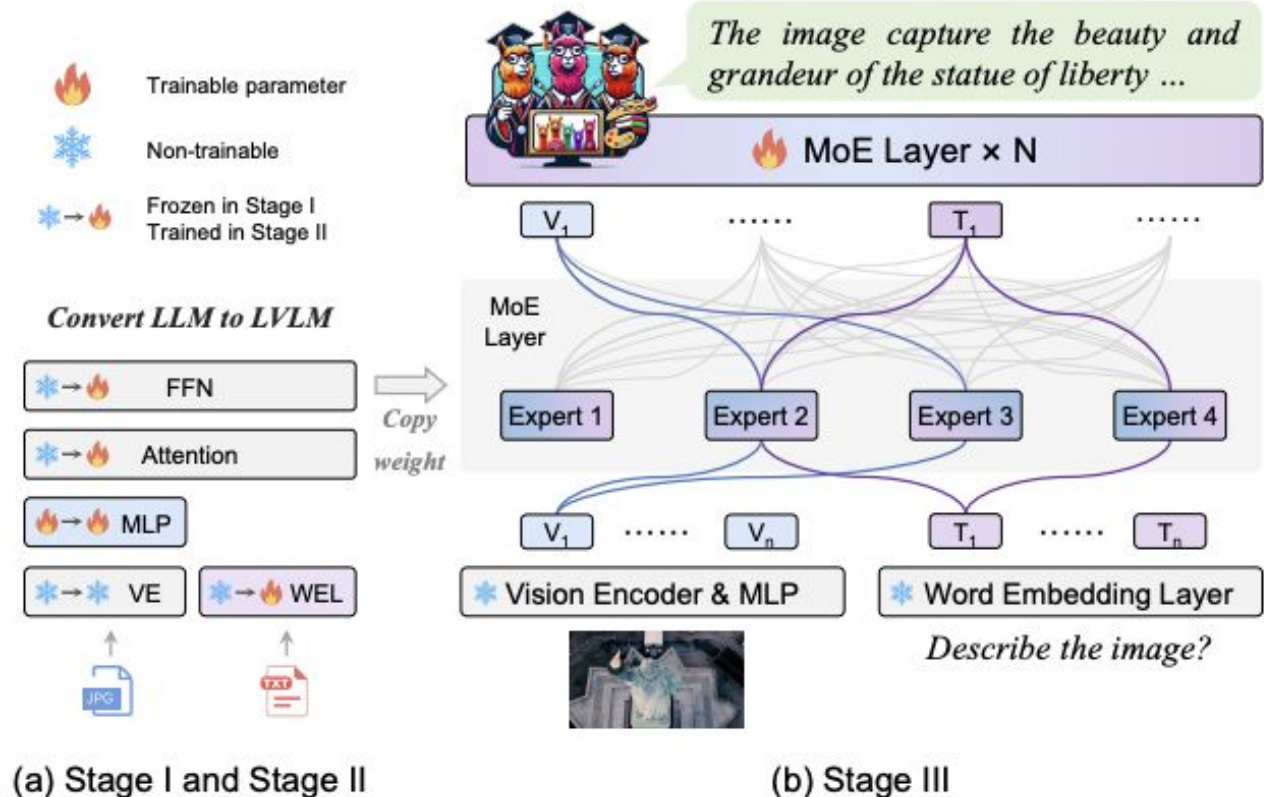| Methods | Image | Video | Region | Pixel |
|---|---|---|---|---|
| LLaVA [25] | ✓ | ✗ | ✗ | ✗ |
| MiniGPT-4 [66] | ✓ | ✗ | ✗ | ✗ |
| Video-LLaMA [62] | ✓ | ✓ | ✗ | ✗ |
| Video-ChatGPT [31] | ✗ | ✓ | ✗ | ✗ |
| Shikra [7] | ✓ | ✗ | ✓ | ✗ |
| CogVLM [50] | ✓ | ✗ | ✓ | ✗ |
| LISA [17] | ✓ | ✗ | ✗ | ✓ |
| u-LLaVA (ours) | ✓ | ✓ | ✓ | ✓ |

# 24. u-LLaVA (arxiv 2023)



**Figure 1**: Overview of u-LLaVA. In stage I, image and spatio-temporal features are used to efficiently boost the general-purpose modality alignment. In Stage II, task-specific projectors and decoders are jointly trained for region and pixel-level understanding. Further, additional modules such as stable diffusion [39] can be easily patched for downstream tasks.

# 25. MoE-LLaVA (arxiv 2024)

- **MoE** + VLM

  (Overview)



(a) Stage I and Stage II

(b) Stage III

# 25. MoE-LLaVA (arxiv 2024)
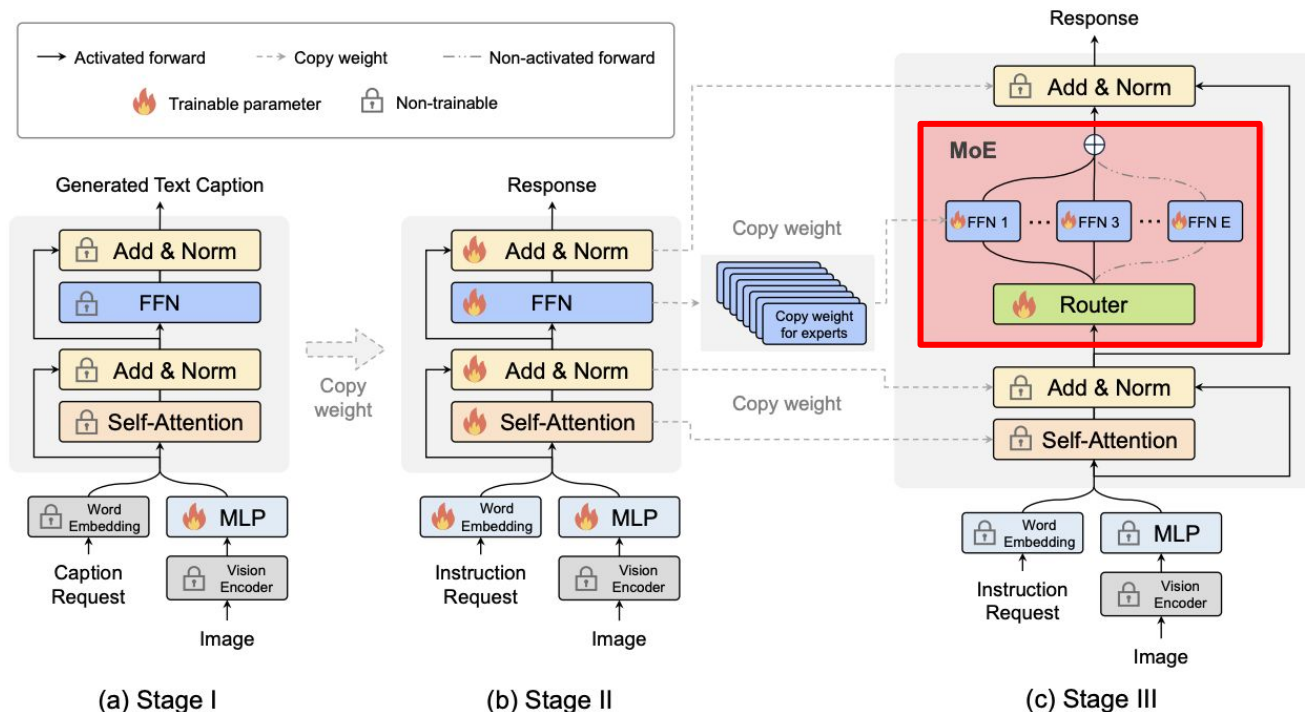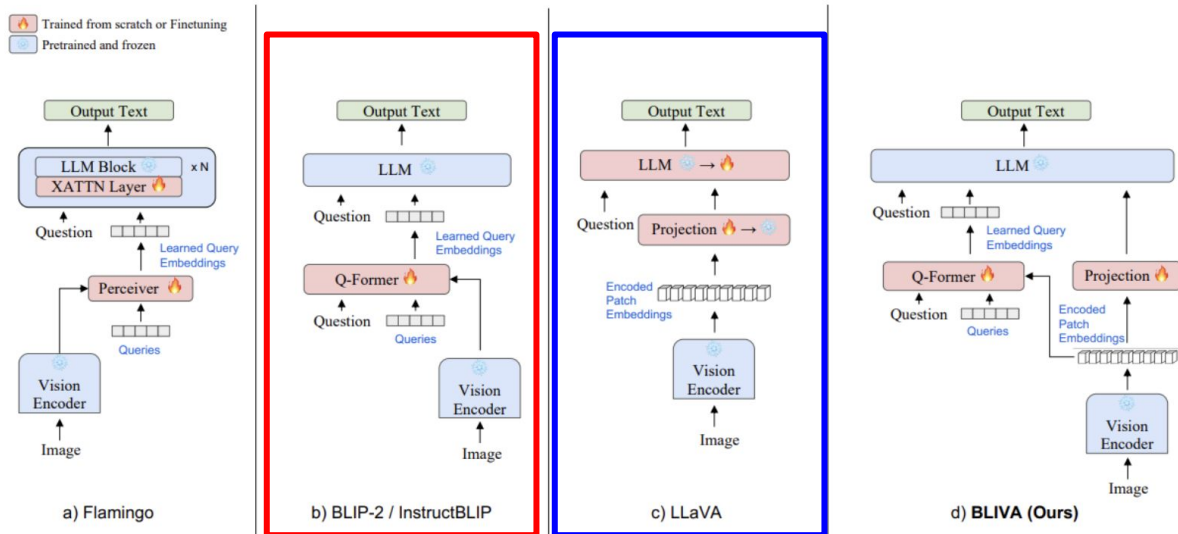
- **MoE** + VLM

  (Detail)



*Figure 3.* **Training framework and strategy.** MoE-LLaVA adopts a three-stage training strategy. (a) We solely train the MLP to adapt the LLM to visual inputs. (b) Training the LLM backend empowers multi-modal understanding capability and MoE layers are not involved. (c) In this stage, we replicate the weights of the FFN to initialize each expert.
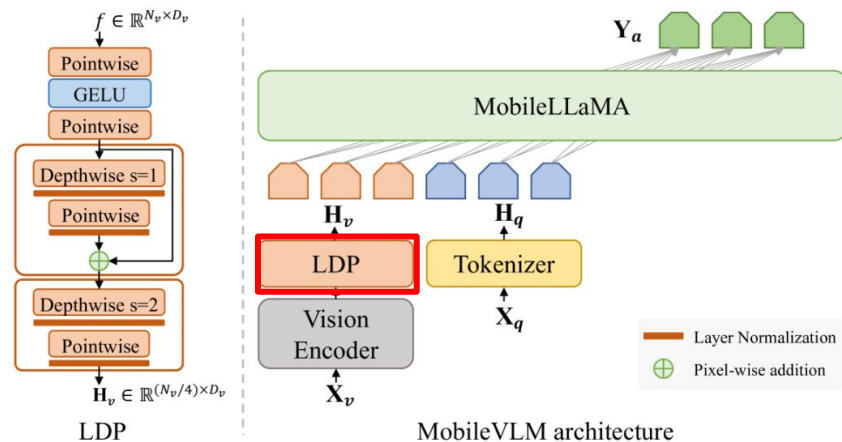
# 26. BLIVA (AAAI 2024)

- Augments "**InstructBLIP**" with "**Visual Assistant**"

  - InstructBLIP: **Learned query embeddings** (feat. Q-Former)

  - Visual Assistant: **Encoded patch embeddings**
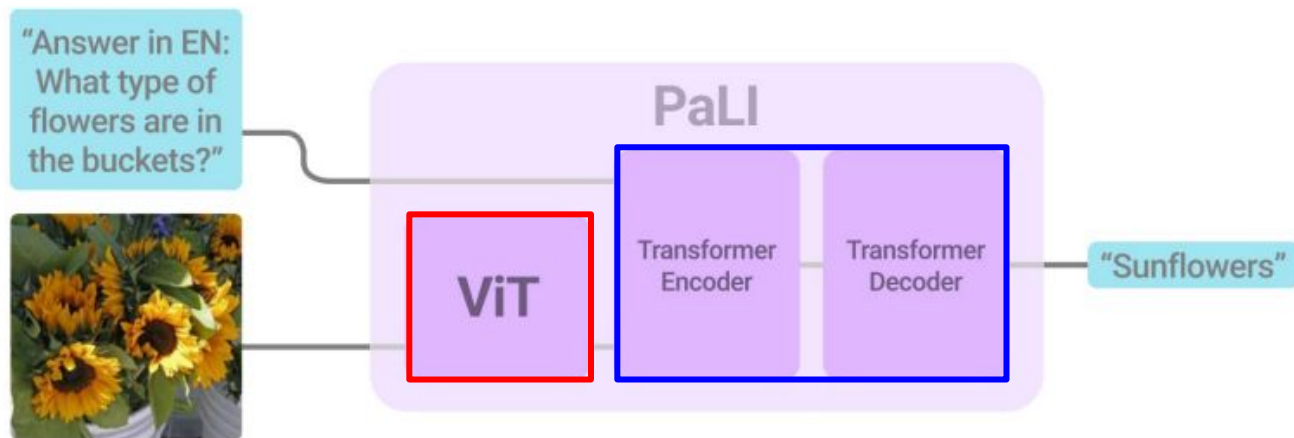
# 27. MobileVLM (arxiv 2023)

- Mobile-optimized VLM

- Architecture

    - (1) Vision: CLIP ViT-L/14

    - (2) LLM: MobileLLaMA

    - Combine (1)->(2) with **LDP**
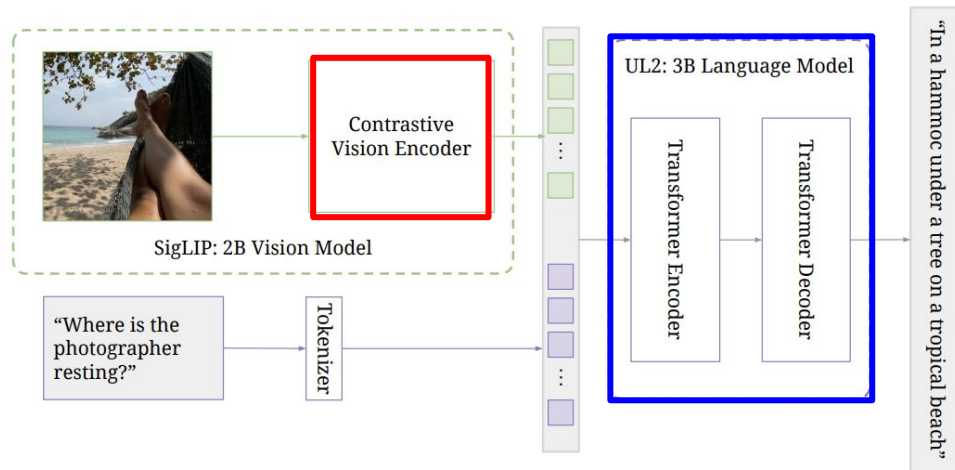
- **lightweight downsample projector**

# 28. PaLI (ICLR 2023)

- Architecture

  - (Image) **ViT**

  - (LLM) **mT5**

# 29. PaLI-3 (arxiv 2023)

- Architecture
    - (Image) **SigLIP-2B**
    - (LLM) **UL2-3B**
- Training:
    - Stage 0) **Unimodal pretraining**
        - Pretrain image encoder
        - Task: (image-text) CL
    - Stage 1) **Multimodal pretrianing**
        - Pretrain LLM
        - Task: Mixture of tasks
    - Stage 2) **Resolution increase**

# 30. PaliGemma (arxiv 2024)

- Architecture

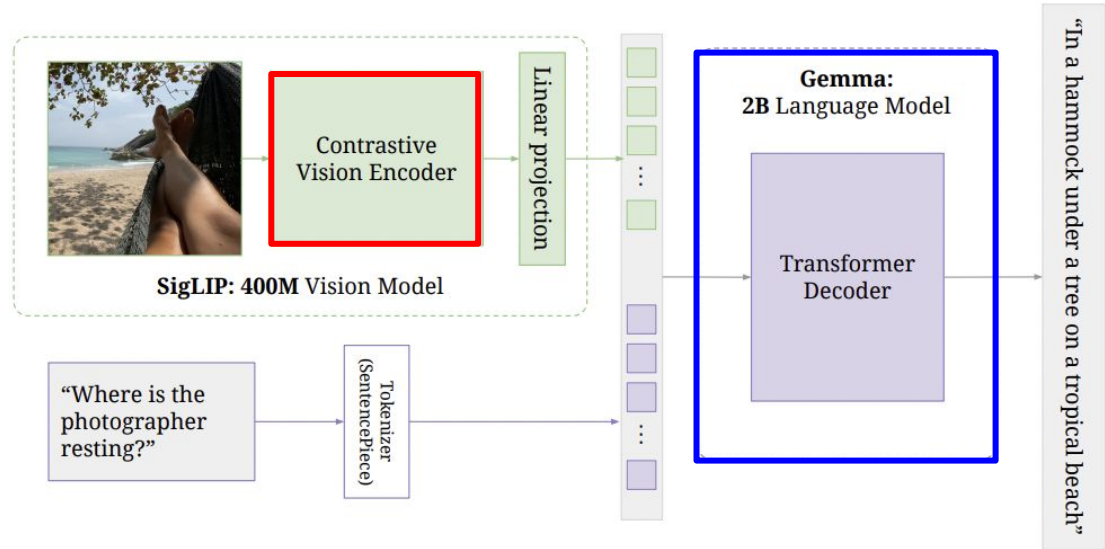  - (Image) **SigLIP**

  - (LLM) **Gemma**



Figure 1 | PaliGemma's architecture: a SigLIP image encoder feeds into a Gemma decoder LM.

# 31. PaliGemma 2 (arxiv 2024)

- PaliGemma 개선

    - Gemma2-**2B** => Gemma2-**27B**
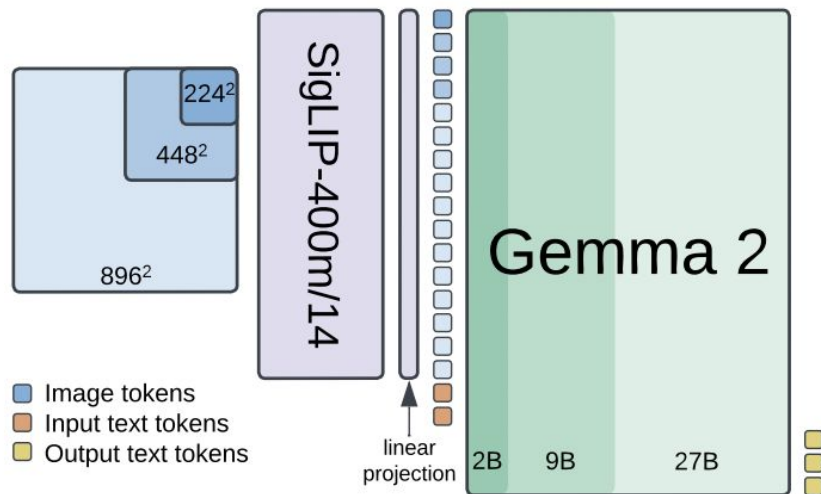
    - 3개의 서로 다른 resolution



Figure 1 | PaliGemma 2 processes a 224px$^2$/ 448px$^2$/896px$^2$ image with a SigLIP-400m en- coder with patch size 14px$^2$, yielding 256/1024/ 4096 tokens. After a linear projection, the image tokens are concatenated with the input text to- kens and Gemma 2 autoregressively completes this prefix with an answer.

# 32. AIMv2 (arxiv 2024)

$$L_{\text{img}} = \frac{1}{I} \sum_{i=1}^{I} \|\hat{x}_i(x_{<i}; \theta) - x_i\|_2^2,$$

- A family of generalist vision encoders

- Vision Encoder + Multimodal decoder

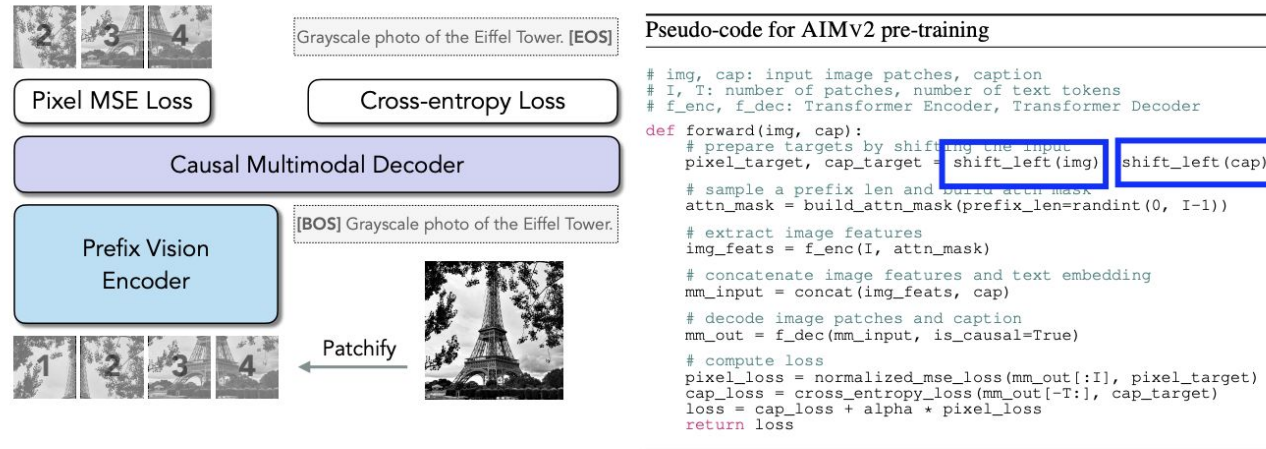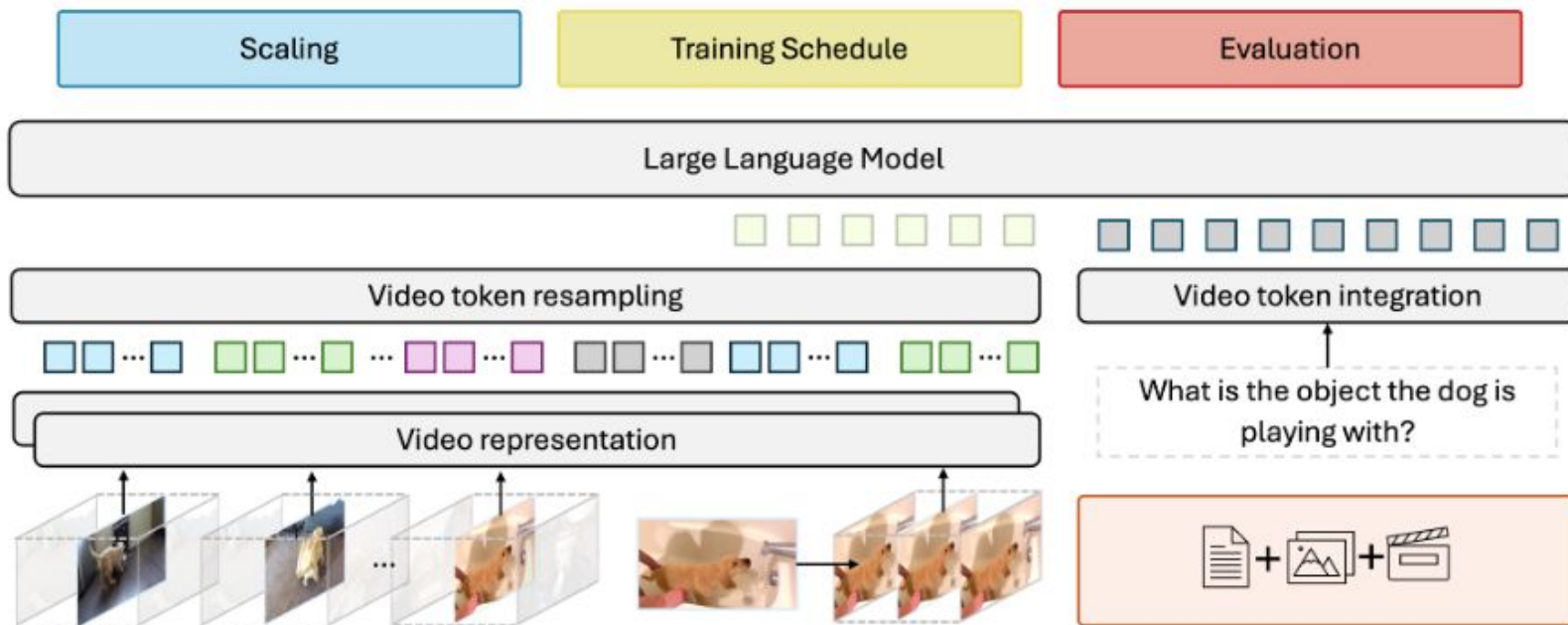$$L_{\text{text}} = -\frac{1}{T} \sum_{t=I+1}^{I+T} \log P(x_t | x_{<t}; \theta).$$

**Figure 1. AIMv2 pre-training Overview.** (Left) Image patches are processed by a vision encoder trained with prefix attention [33, 95]. The resulting visual representations are concatenated with the text embeddings of their corresponding captions. This combined multimodal sequence is then processed by a joint decoder. The model is pre-trained to autoregressively reconstruct the shifted input. (Right) Pseudocode for the forward pass during AIMv2 pre-training. The pre-training process of AIMv2 is straightforward to implement, resembling that of AIM and LLMs as it relies solely on a simple autoregressive objective.
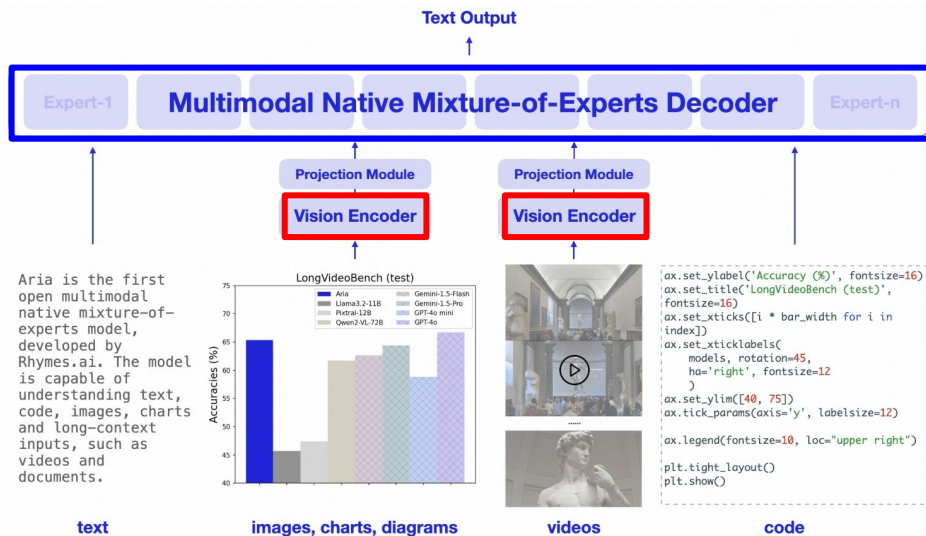
# 33. Apollo (arxiv 2024)
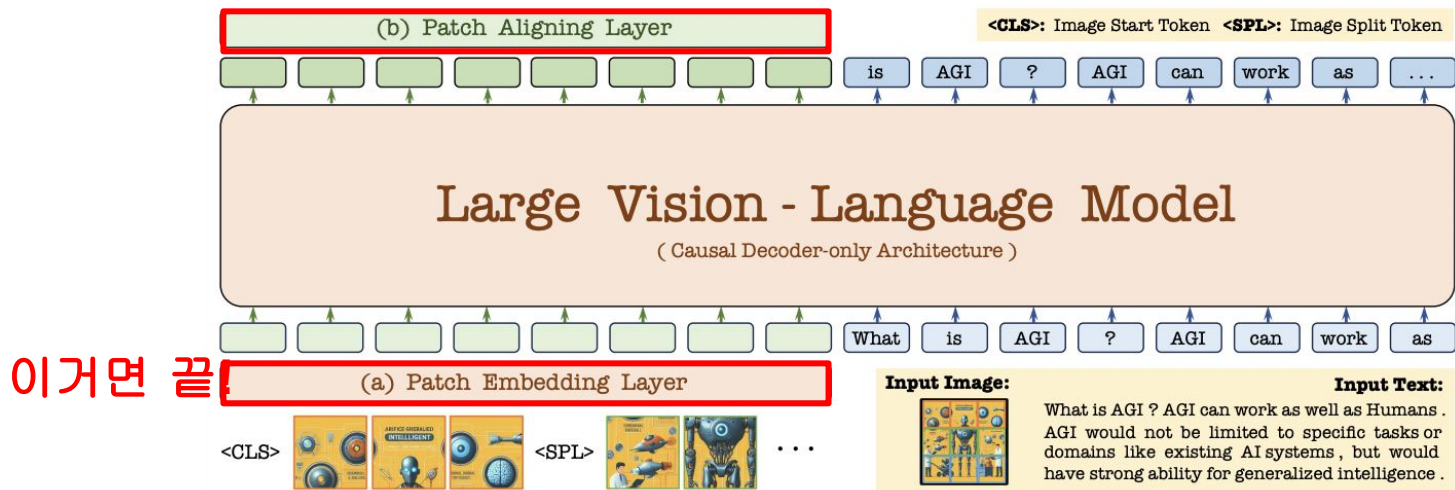
- SOTA LMMs for **video understanding**

# 34. ARIA (arxiv 2024)

- **Multimodal MoE**

  - Text, Code, Images, Video …

- Architecture

  - (Fine-grained) **MoE decoder**

  - (Lightweight) **Visual encoder**

# 35. EVE (NeurIPS 2024)

- Encoder-free VLM

  - **별도의 vision encoder 필요 없음**

- Image & Text를 동시에 "unified" "decoder-only" 아키텍처에 넣음

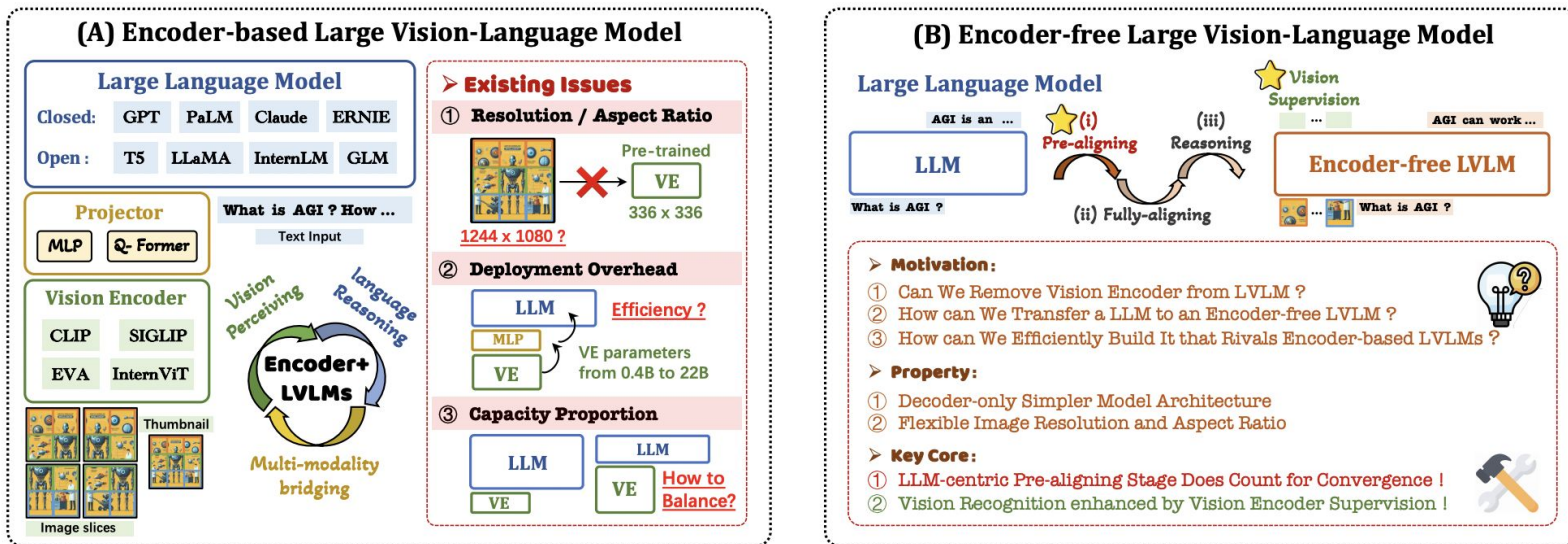# 35. EVE (NeurIPS 2024)



Figure 1: Overview of large (A) encoder-based and (B) encoder-free vision-language models. **Encoder-based VLMs** contain vision encoders (VE) and large language models (LLM), with a projector as a vision-language bridge, while **encoder-free VLMs** exclude vision encoders and handle vision perception and linguistic instruction simultaneously with one unified architecture.
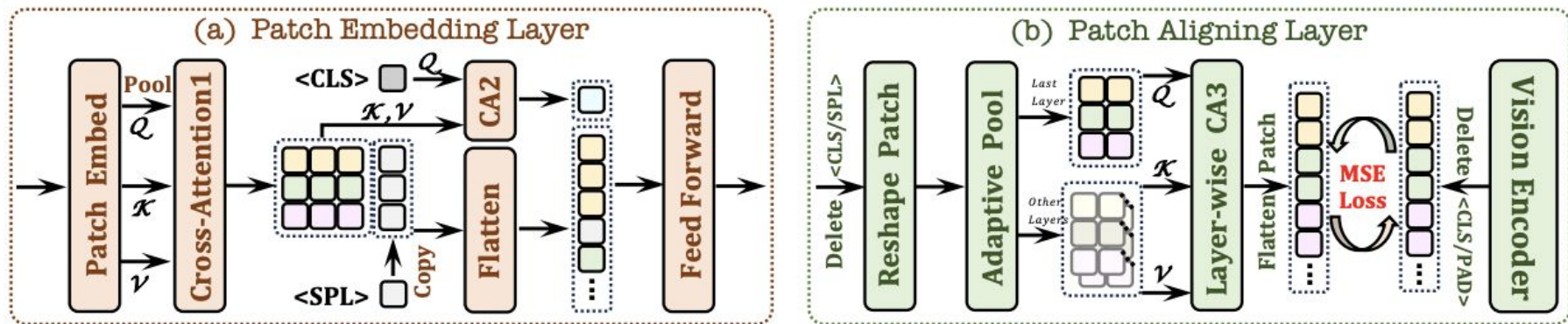
# 35. EVE (NeurIPS 2024)



Figure 3: Overview of (a) patch embedding and (b) patch aligning layer. The former layer encodes images into patch features and employs cross-attention (CA1) within a limited receptive field to enhance representations. Meanwhile, a special <CLS> token provides a holistic view of each patch feature in the subsequent backbone. The latter layer removes all meaningless tokens and adjusts the size of patch features to align with the semantics of the same positional features in the vision encoder.

# 35. EVE (NeurIPS 2024)
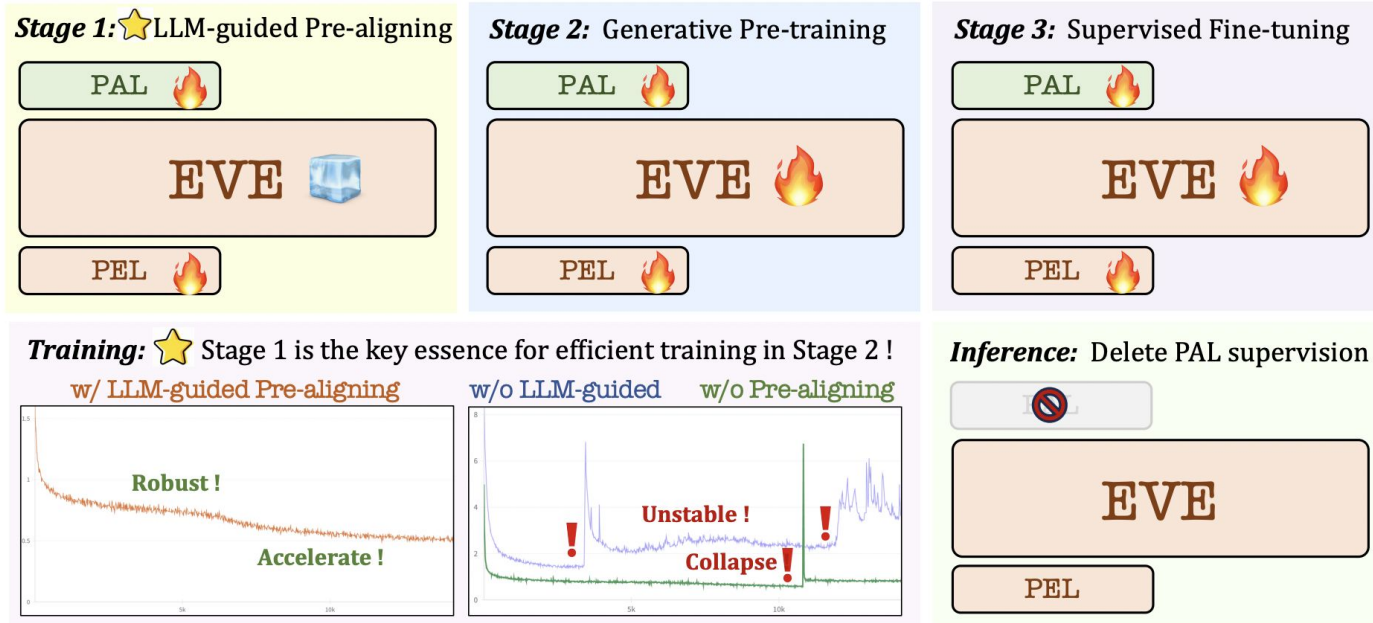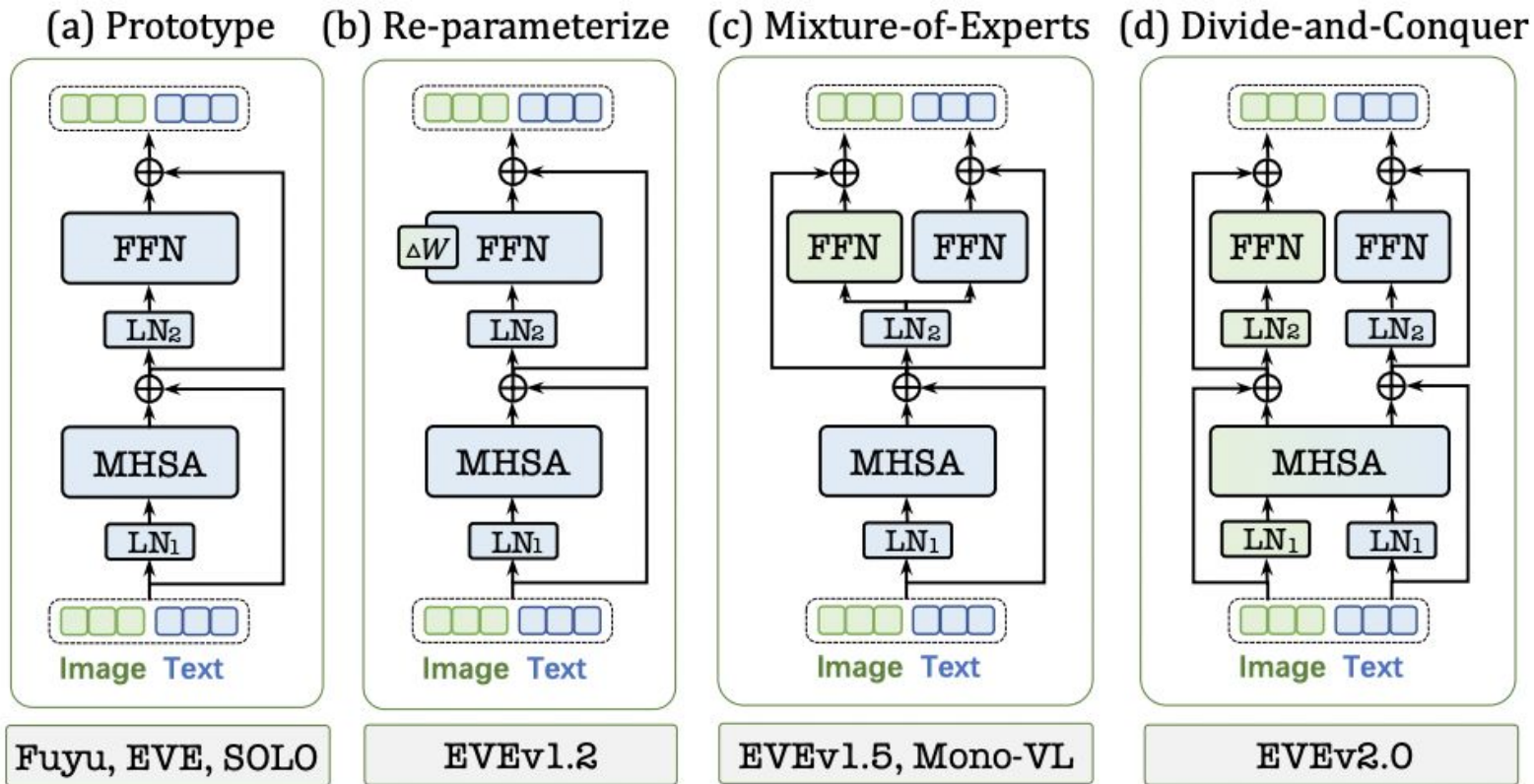


Figure 4: Overview of training procedure with three successive stages. We perform initial vision-language alignment guided by a frozen LLM in Stage 1, and then update the entire backbone for Stage 2 and Stage 3. We empirically find that Stage 1 is quite crucial to avoid collapse and accelerate convergence, thereby enhancing training efficiency. Notably, PAL is removed during inference.

# 36. EVEv2 (arxiv 2025)

- EVE와 마찬가지로 **Encoder-free VLM**

- **Divide-and-conquer** 아키텍처 도입

    - **"Modality-specific"** component

        ( Text & Vision에 각기 다른 Att, LN, FFN )

- 효과

    - Maximize scaling efficiency

    - Reduce inter-modality interference

    - Superior data efficiency

# 36. EVEv2 (arxiv 2025)



(a) Prototype — Fuyu, EVE, SOLO

(b) Re-parameterize — EVEv1.2

(c) Mixture-of-Experts — EVEv1.5, Mono-VL

(d) Divide-and-Conquer — EVEv2.0
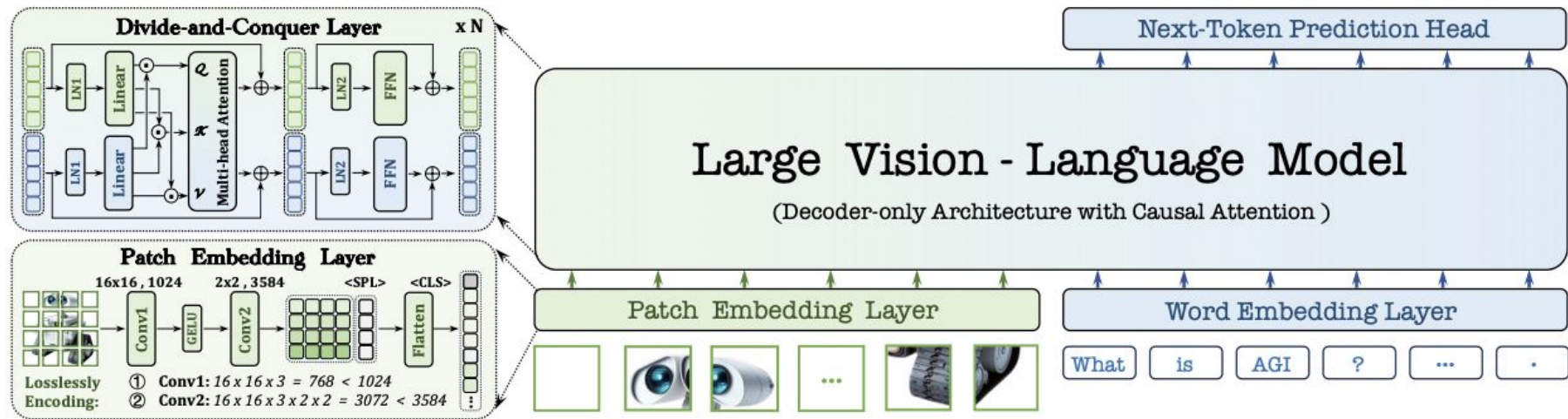
# 36. EVEv2 (arxiv 2025)



Figure 3. Overview of our proposed EVEv2.0 framework. We first adopt a patch embedding layer to encode images losslessly, and then concatenate visual and textual tokens into a unified decoder-only vision-language model. Here, it extends the standard autoregressive transformer by incorporating modality-specific weights for each multi-head self-attention layer, feed-forward layer, and layer normalization.

# 37. Janus (arxiv 2024)

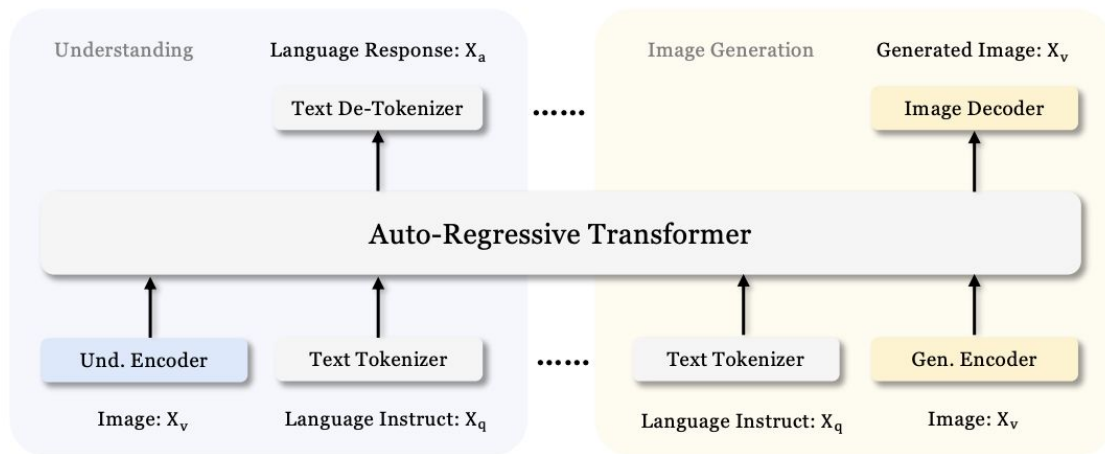- AR framework that unifies multimodal understanding and generation



Figure 2 | **Architecture of our Janus.** Different from previous approaches [77, 85] that typically assume visual understanding and generation require the same visual encoder, our Janus decouples visual encoding for visual understanding and visual generation. "Und. Encoder" and "Gen. Encoder" are abbreviations for "Understanding Encoder" and "Generation Encoder", respectively. Best viewed in color.

# 37. Janus (arxiv 2024)

- AR framework that unifies multimodal understanding and generation
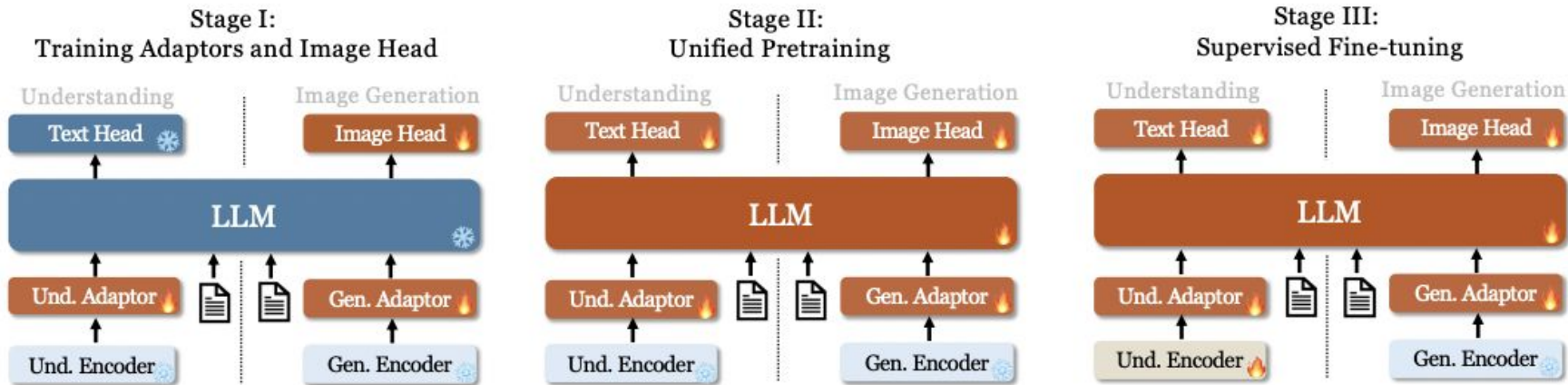


Figure 3 | **Our Janus adopts a three-stage training procedure.** We use flame symbols/snowflake symbols in the diagram to indicate the module updates/does not update its parameters.

# 38. Janus-Pro (arxiv 2025)

- **Janus를 발전 시킴**

  - (1) Training strategy (minor) 변화

    To address this issue, we make two modifications.

    - **Longer Training in Stage I**: We increase the training steps in Stage I, allowing sufficient training on the ImageNet dataset. Our findings reveals that even with the LLM parameters fixed, the model could effectively model pixel dependence and generate reasonable images based on category names.
    - **Focused Training in Stage II**: In Stage II, we drop ImageNet data and directly utilize normal text-to-image data to train the model to generate images based on dense descriptions. This redesigned approach enables Stage II to utilize the text-to-image data more efficiently, resulting in improved training efficiency and overall performance.

  - (2) Training data 키우기

  - (3) Model scale-up
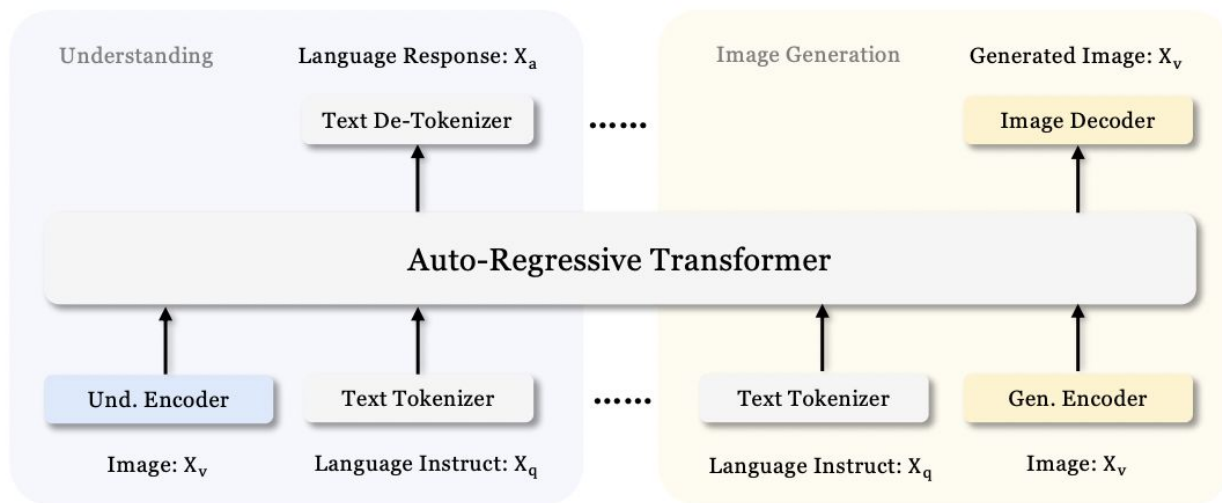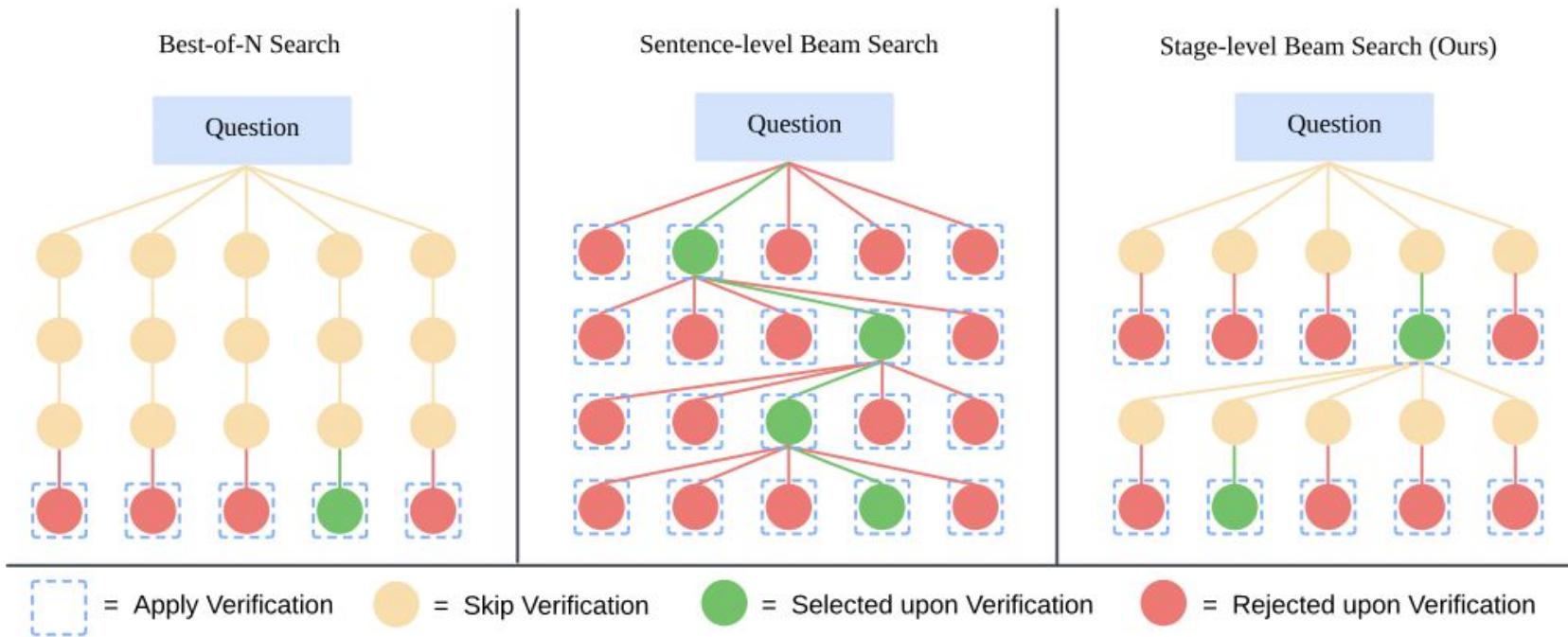
# 38. Janus-Pro (arxiv 2025)

- 아키텍처는 동일함



Figure 3 | **Architecture of our Janus-Pro.** We decouple visual encoding for multimodal understanding and visual generation. "Und. Encoder" and "Gen. Encoder" are abbreviations for "Understanding Encoder" and "Generation Encoder", respectively. Best viewed on screen.

# 39. LLaVA-CoT (arxiv 2024)

- **Multi-stage reasoning**을 하는 VLM

# 40. LLM2CLIP (arxiv 2024)

- Fine-tuning 방법론
    - **LLM & CLIP visual encoder를 연결**하기 위해!
- 2-stage approach
    - **(1) Caption Contrastive FT**
        - LLM을 FT한다
        - Causal ATT -> Bidirectional ATT
          ( 목적: 생성 X, 이해 O 이므로 )
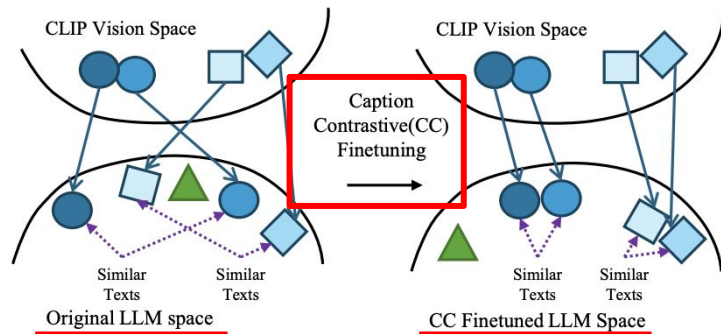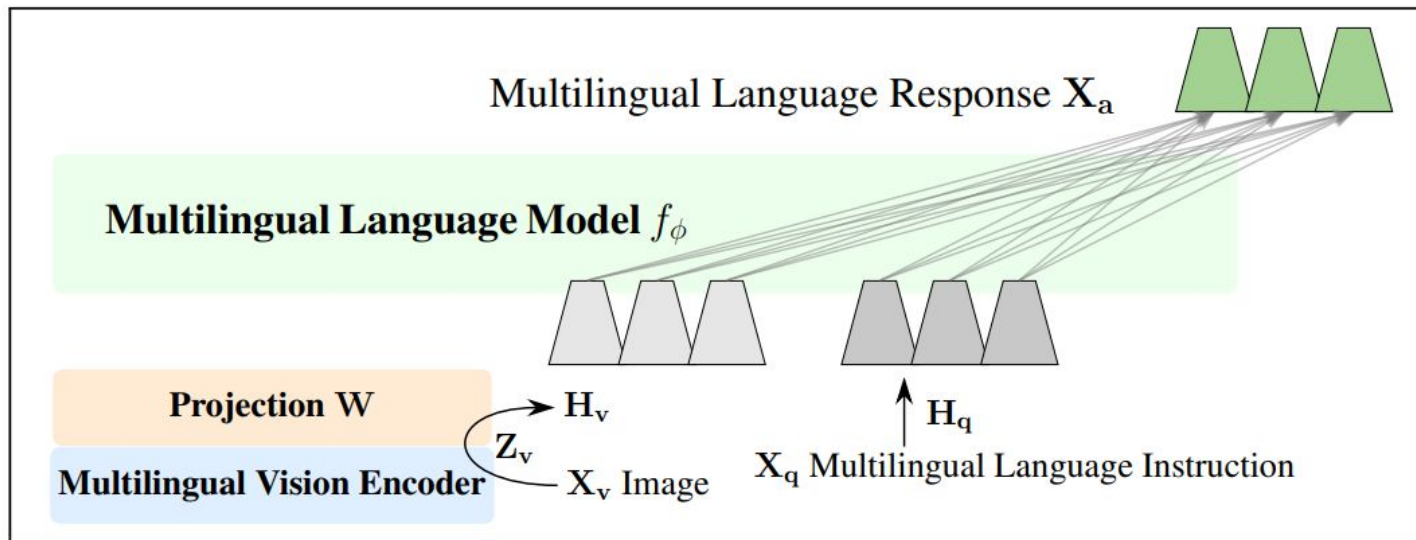    - **(2) CLIP Vision Encoder FT**
        - Vision Encoder를 FT한다



Figure 2: By adjusting the output space of the LLM, we enable the LLM to more effectively act as a private tutor for CLIP, allowing the fine-tuning process of *LLM2CLIP* to efficiently adjust the pretrained CLIP's cross-modal space.
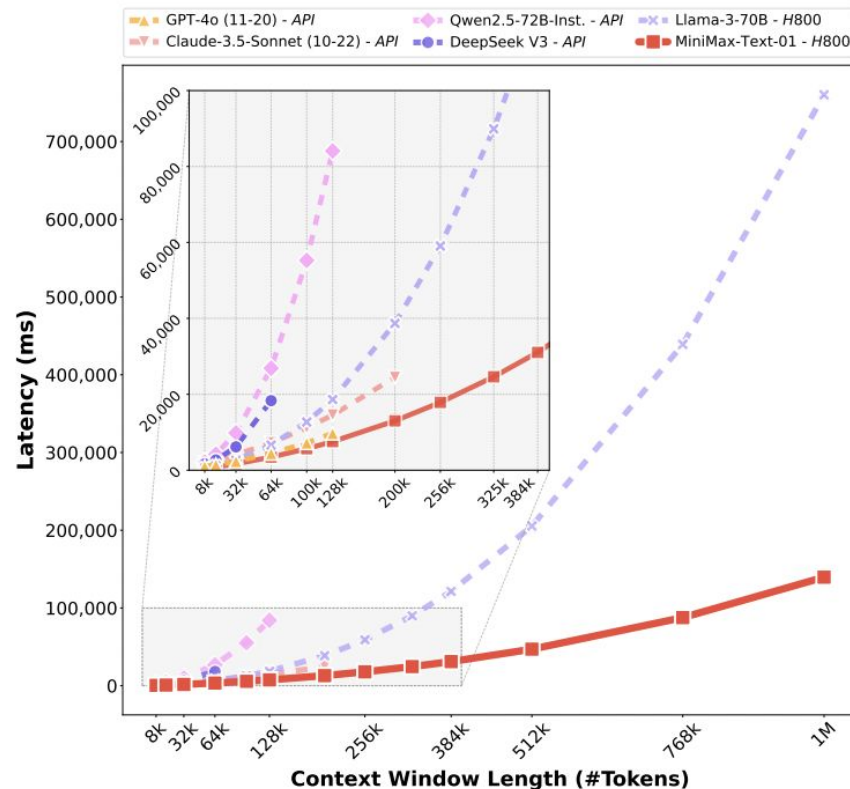
# 41. Maya (arxiv 2024)

- Multimodal **Multilingual** VLM (mVLM)

- Framework: **LLaVA 1.5**

# 42. MiniMax-01 (arxiv 2025)

- MiniMax-01

  - MiniMax-Text-01

  - MiniMax-VL-01

- **Longer context window (~ 4M)**

  => How? Novel architecture!

  - **Lightning attention**

  - MoE

# 42. MiniMax-01 (arxiv 2025)

- MiniMax-01

  - MiniMax-Text-01

  - MiniMax-VL-01

- **Longer context window (~ 4M)**

  => How? Novel architecture!

  - Lightning attention
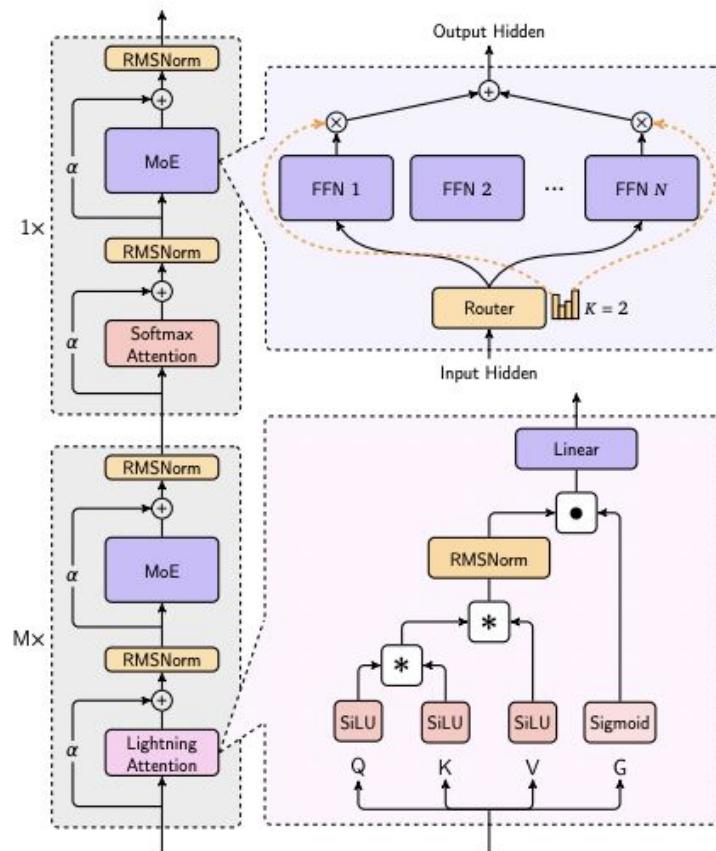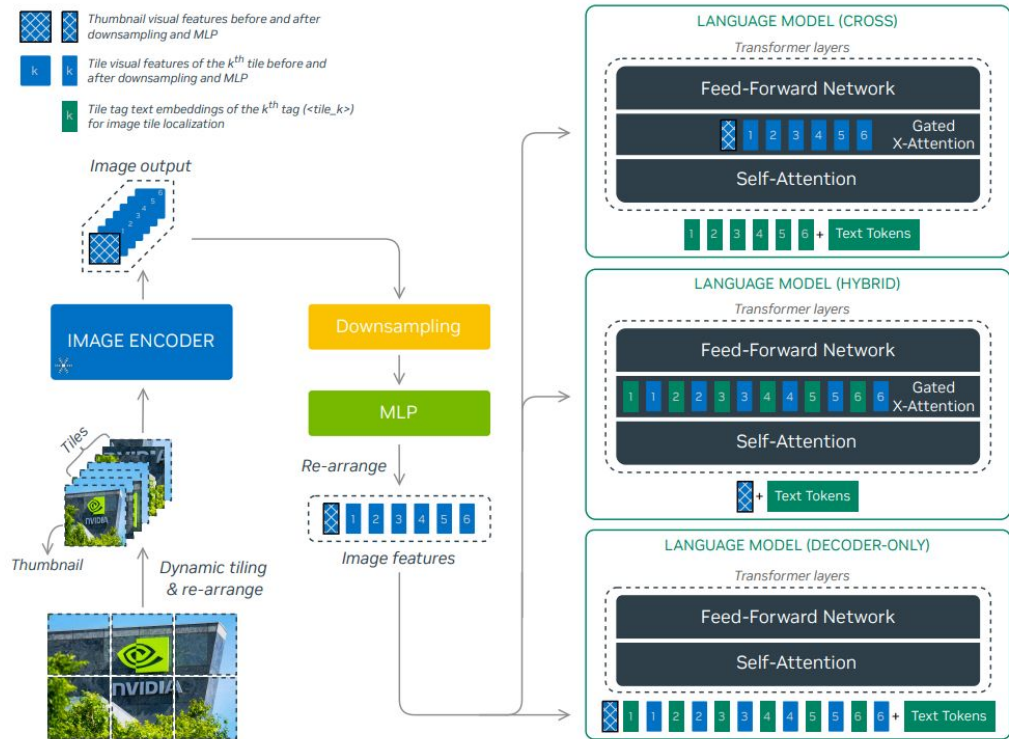
  - **MoE**



Figure 3 | **The architecture of MiniMax-Text-01.**
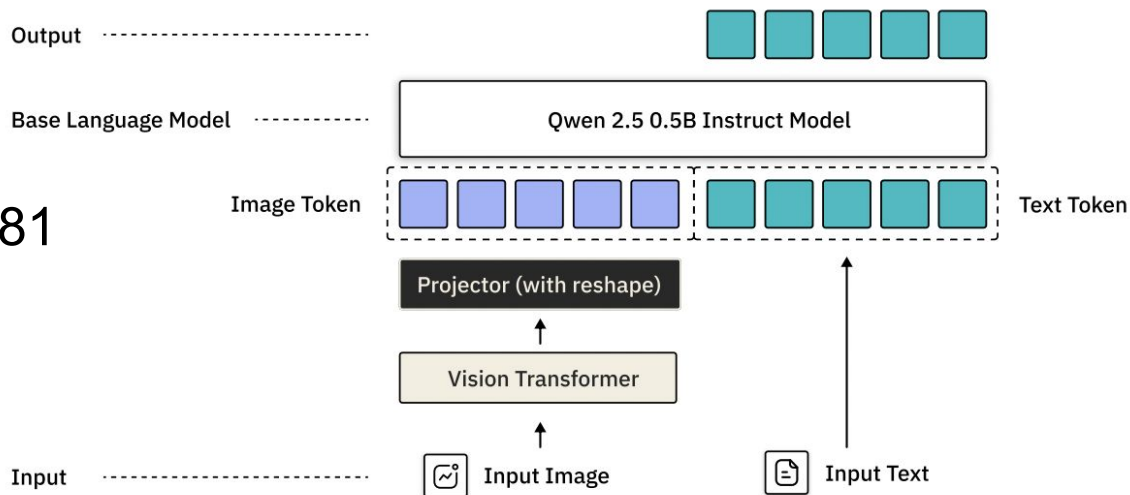
# 43. NVLM (arxiv 2024)

- NVLM = NVIDIA VLM

  - NVLM-**D** (Decoder-only)

  - NVLM-**X** (Cross-attention)

  - NVLM-**H** (Hybrid)

# 44. OmniVLM (arxiv 2024)

- VLM with sub-billion # params

- 목적: **On-device** inference

- **Token compression**

    - Convolution-based

    - # Visual token: 729 -> 81

- Architecture

    - Vision: SigLIP-400M

    - LLM: Qwen2.5-0.5B

# 45. Pixtral 12B (arxiv 2024)

- Mistral AI에서 공개한 모델
- Architecture
  - **(1) Vision encoder = Pixtral-ViT**
    - Trained with a new vision encoder from scratch
    - Support **variable** image sizes & aspect ratios
    - New **ROPE-2D** implementation
  - **(2) Multimodal decoder: NTP task**
    - Built on top of Mistral Nemo 12B
  - MLP: Connect (1) to (2)
- New benchmark: MM-MT-Bench

$$\text{RoPE-2D}\left(x^{(i,j)}, \Theta\right) = M_\Theta^{(i,j)} x^{(i,j)},$$

$$\text{where} \quad M_\Theta^{(i,j)} = \begin{pmatrix} \cos i\theta_1 & -\sin i\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin i\theta_1 & \cos i\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos j\theta_2 & -\sin j\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin j\theta_2 & \cos j\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos j\theta_{\frac{d}{2}} & -\sin j\theta_{\frac{d}{2}} \\ 0 & 0 & 0 & 0 & \cdots & \sin j\theta_{\frac{d}{2}} & \cos j\theta_{\frac{d}{2}} \end{pmatrix}$$
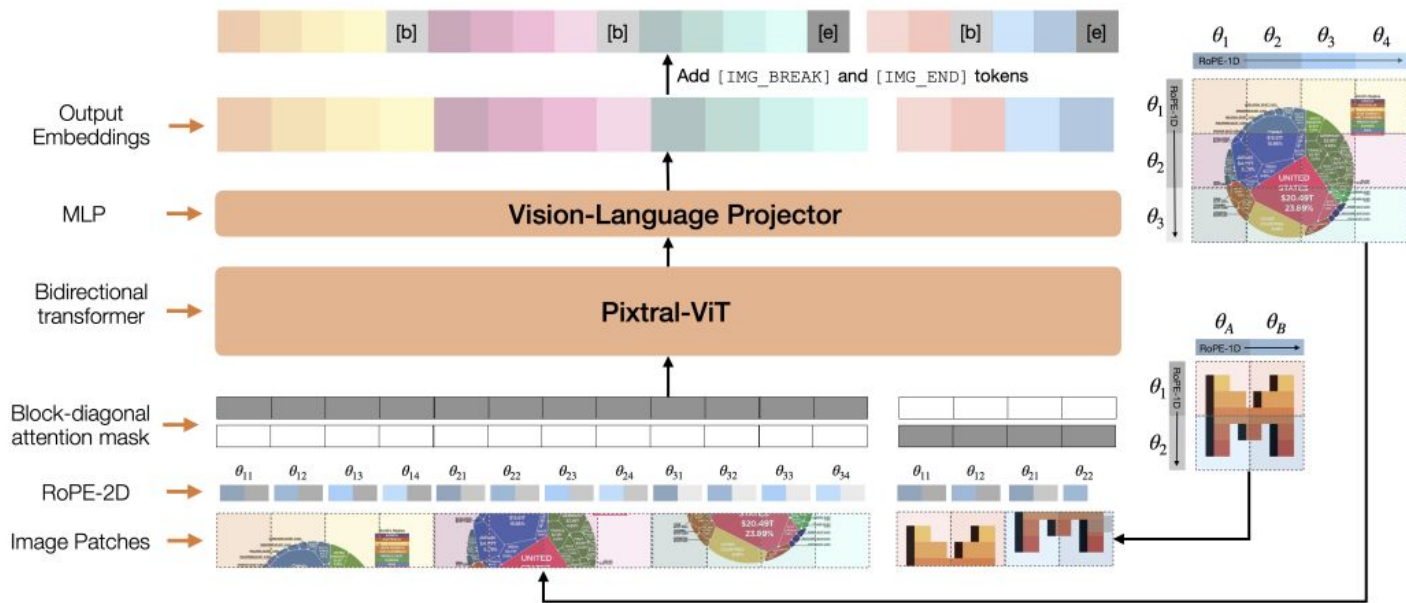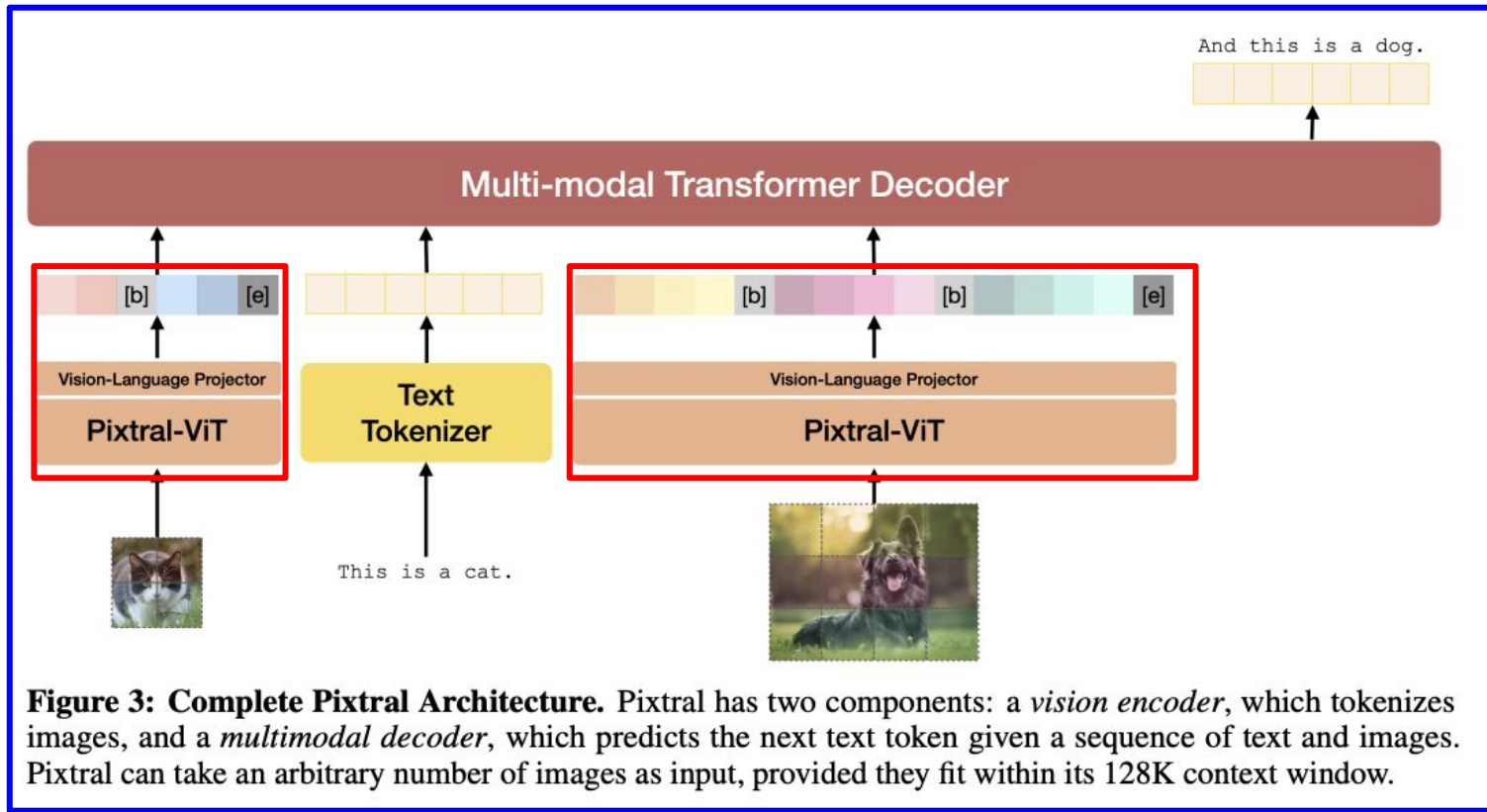
# 45. Pixtral 12B (arxiv 2024)



**Figure 2: Pixtral Vision Encoder.** Pixtral uses a new vision encoder, which is trained from scratch to natively support variable image sizes and aspect ratios. Block-diagonal attention masks enable sequence packing for batching, while RoPE-2D encodings facilitate variable image sizes. Note that the attention mask and position encodings are fed to the vision transformer as additional input, and utilized only in the self-attention layers.

# 45. Pixtral 12B (arxiv 2024)



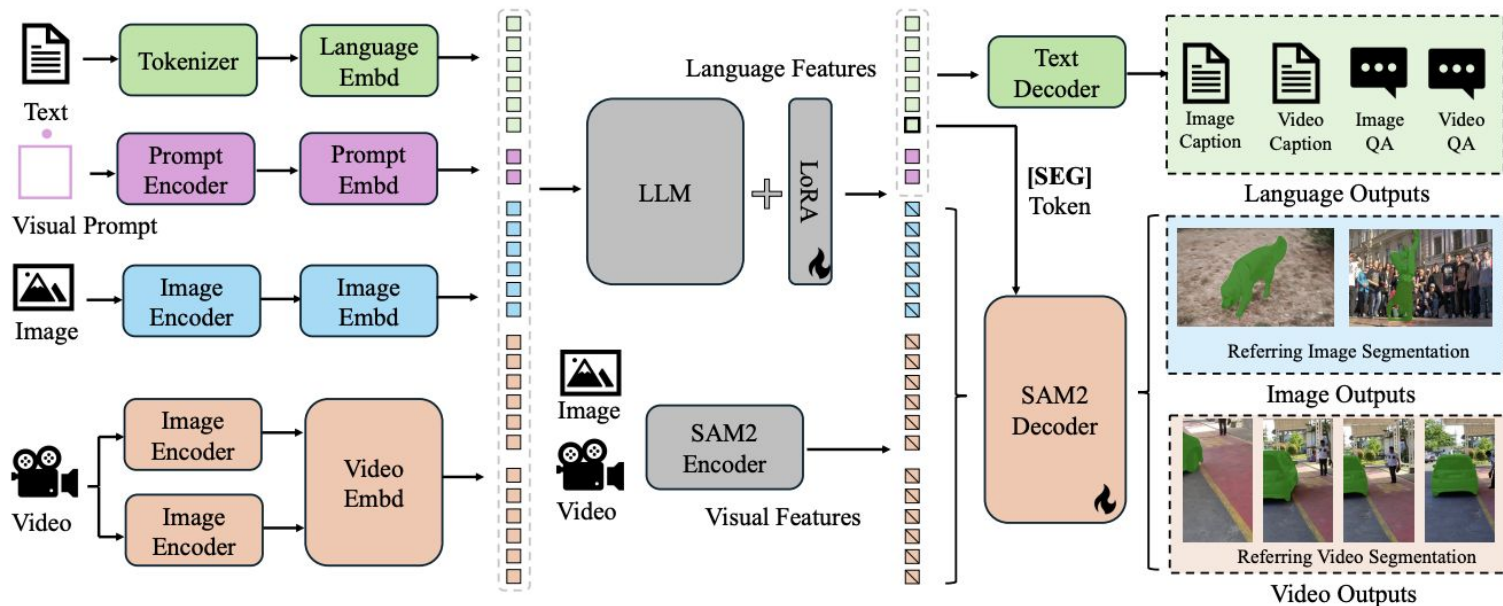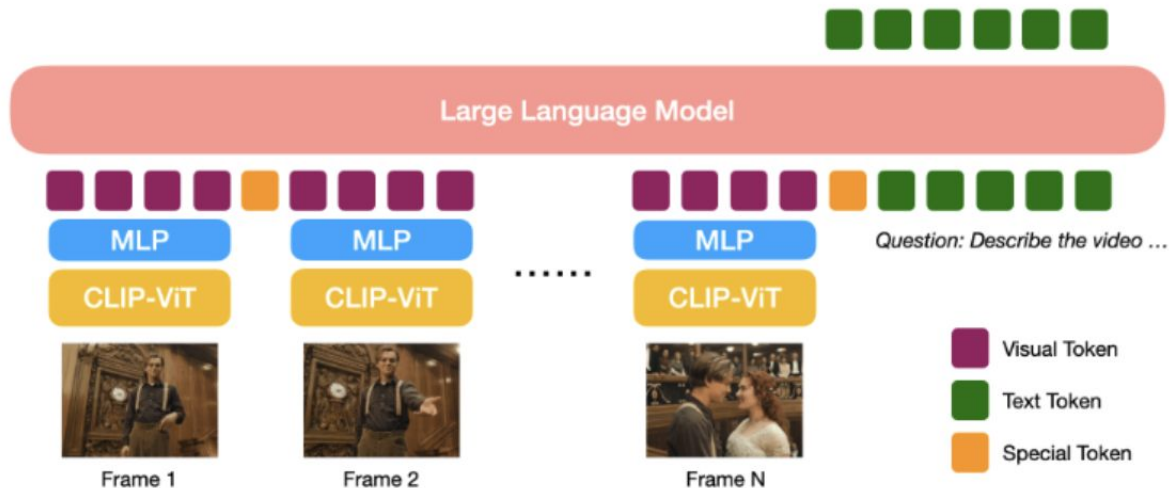**Figure 3: Complete Pixtral Architecture.** Pixtral has two components: a *vision encoder*, which tokenizes images, and a *multimodal decoder*, which predicts the next text token given a sequence of text and images. Pixtral can take an arbitrary number of images as input, provided they fit within its 128K context window.

# 46. Sa2VA (arxiv 2025)

- "Unified" model for "**dense**" grounded understanding of "**image & video**"

- Treat **all inputs (text, images, videos)** as tokens in a **shared LLM**

Table 1. **Comparison of capabilities of different representative models.** Our method supports various tasks and modalities. Benefiting from these interactive features on video, Sa2VA can perform multiple promptable tasks in the video, as shown in Fig. 1 (a) and (b).

| Method | Support Inputs | | | Dense Grounding | | | | Conversation | | | | End to End |
| | Image | Video | Visual Prompts | RES | Ref-VOS | Inter-VOS | GCG | Image-Chat | Video-Chat | Video Caption | Interactive Caption | Training |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LLAVA [53] | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| LLaVA-OneVision [42] | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✗ | ✓ |
| InternVL 2.0 [9] | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✗ | ✓ |
| Osprey [92] | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| LISA [40] | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| GLAMM [66] | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ |
| VIP-LLaVA [5] | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| VISA [85] | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| OMG-LLaVA [99] | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ |
| PSALM [102] | ✓ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| GSVA [83] | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| LLaMA-VID [50] | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✗ | ✓ |
| ST-LLM [56] | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✗ | ✓ |
| F-LLM [82] | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| Sa2VA (Ours) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

# 46. Sa2VA (arxiv 2025)



Figure 2. **Our proposed Sa2VA model.** The model first encodes the input texts, visual prompts, images, and videos into token embeddings. These tokens are then processed through a large language model (LLM). The output text tokens are used to generate the "[SEG]" token and associated language outputs. The SAM-2 decoder receives the image and video features from the SAM-2 encoder, along with the "[SEG]" token, to generate corresponding image and video masks.

# 47. Tarsier2  (arxiv 2025)

-   VLM that **generates video descriptions**

-   Superior video understanding capabilities

-   3 stage: Pretraining - SFT - DPO

# 48. VideoChat-Flash  (arxiv 2024)

- Handling **long-form video** content in MLLMs

- **HiCo = Hierarchical visual token Compression**

    - To reduce computational load!

- Two compression

    - **(1) Clip-level**

        - Video = Clips x N

        - Visual encoder processes each clip

        - Compressor reduces # of visual tokens (to exploit inter-frame redundancy)

    - **(2) Video-level**

        - Visual tokens are further reduced via dropout

# 48. VideoChat-Flash  (arxiv 2024)

# 49. VideoLLaMA 3 (arxiv 2025)

- For both image & video understanding

- Architecture

  - (1) Vision: Pretrained SigLIP

  - (2) Video: DiffFP (compressor)

  - Projector: for (1)+(2) to LLM

  - (3) LLM: Qwen 2.5

- Four training stages

# 50. LLaMA 3.2-Vision (blog 2024)

- **LLaMA 3** with "**multimodal** capabilities"

- # params: 11B, 90B

- Architecture:

  - (1) Visual encoder

  - (2) LLM: LLaMA 3

  - **Vision adapter**

    - with cross-attention layers

    - for (1) -> (2)

# 51. Idefics2-8B (blog 2024)

- Efficiently process **interleaved image & text**

- Architecture

    - (1) Vision: SigLIP

    - (2) LLM: Mistral-7B

    - **Perceiver pooling layer + MLP**: for (1) -> (2)

# 52. Idefics3-8B (arxiv 2024)

- Idefics2-8B보다 뛰어남

- Architecture

    - Vision: SigLIP

    - LLM: LLaMA 3.1 instruct

- Idefics2의 **perceiver resampler를**

    **simple pixel shuffle strategy로 대체**

    ( similar to InternVL-1.5 )

# 53. InternLM-XComposer2 (arxiv 2024)

- Free-form text-image composition & comprehension

- Architecture

  - (1) Vision: CLIP

  - (2) LLM: InternLM-2

  - **Partial LoRA**: apply to visual tokens



Figure 2. **The illustration of the Partial-LoRA.** The blue tokens represent the visual tokens and the gray tokens are the language tokens. Our Partial-LoRA is only applied to the visual tokens.

# 54. InternLM-XComposer2-4KHD (NeurIPS 2024)

- InternLM-XComposer를 기반으로

- Handle high-resolution image

- How? **Dynamic resolution**

  + **Automatic patch configuration**

    - Resolution 336 ~ 4K HD



**Figure 3**: The framework of InternLM-XComposer2-4KHD. Our model processes the high-resolution image with a Dynamic Image Partition strategy and concatenates the image tokens with text tokens as LLM input.

# 55. InternLM-XComposer-2.5 (arxiv 2024)

- Handle "**LONG**" contextual input & output



Figure 5. **Framework** of IXC-2.5 that supports the multi-modal inputs, including text, single/multiple images, and videos.

# 56. InternVL 2.5 (arxiv 2024)

- InternVL 2.0를 이음

- 448x448 image tile로 나눔

- 핵심: **Pixel unshuffle operation**

  - Visual token: 1024 -> 256로 줄임

- **Dynamic resolution**



r² channels    High-resolution image (output)

To find the optimal aspect ratio for resizing, we define the set of target aspect ratios $\mathcal{R}$ as:

$$\mathcal{R} = \{i/j \mid 1 \leq i, j \leq n, i \times j \in [n_{\min}, n_{\max}]\}. \tag{1}$$

The closest aspect ratio $r_{\text{best}}$ is selected by minimizing the difference between the original aspect ratio $r$ and each target aspect ratio $r_{\text{target}}$:

$$r_{\text{best}} = \arg \min_{r_{\text{target}} \in \mathcal{R}} |r - r_{\text{target}}|. \tag{2}$$

In cases where multiple aspect ratios produce the same difference (*e.g.*, 1:2 and 2:4), we prioritize the aspect ratio that results in an area less than or equal to twice the original image size. This helps to some extent in preventing the excessive enlargement of low-resolution images.

# 56. InternVL 2.5 (arxiv 2024)



(a) Data Preprocessing

(b) Model Architecture

Figure 2: **Overall architecture.** InternVL 2.5 retains the same model architecture as InternVL 1.5 [35] and InternVL 2.0, *i.e.* the widely-used "ViT-MLP-LLM" paradigm, which combines a pre-trained InternViT-300M or InternViT-6B with LLMs [19, 229] of various sizes via an MLP projector. Consistent with previous versions, we apply a pixel unshuffle operation to reduce the 1024 visual tokens produced by each 448×448 image tile to 256 tokens. Moreover, compared to InternVL 1.5, InternVL 2.0 and 2.5 introduced additional data types, incorporating multi-image and video data alongside the existing single-image and text-only data.

# 57. SmolVLM (blog 2024)

- **2B** params VLM

- SOTA for its memory footprint

- **Small, fast, memory-efficient**

- Built upon Idefics3

  - LLaMA 3.18B -> SmoLM2 1.7B

- Aggressive image compression

# 58. Qwen-VL (arxiv 2023)

- LLM인 **QWEN에 visual 능력**을 갖춰줌.

- (1) **Visual receptor**

    - Visual Encoder = ViT

    - Position-aware VL Adapter = 1개의 cross-attention

- (2) Input-output

    - Input = Image & **Bbox**

    - Output = **Bbox**

- (3) Three-stage training pipeline

# 59. Qwen2-VL (arxiv 2024)

- Qwen-VL를 개선함

- 핵심 개선: Improved understanding of images & videos

- Details:

    - (1) **<span style="color:red">Naive Dynamic Resolution</span>**

      - arbitrary resolution 받아서 dynamic # of visual tokens 반환

    - (2) Multimodal Rotary Position Embedding (**<span style="color:red">M-ROPE</span>**)

        - temporal + height + width information

# 59. Qwen2-VL (arxiv 2024)

# 60. Qwen2.5-VL (arxiv 2025)

- Qwen2-VL를 개선함

    - Perception of temporal & spatial scales

    - Simplified network structure

        -> increased efficiency

# 61. DeepSeek-VL (arxiv 2024)

- (1) Data: Diverse & Large-scale

- (2) Model: **Hybrid** vision encoder

  - LLM: Deepseek LLM (feat. LLaMA)

  - Vision:

    - a) **SigLIP** (low res)

    - b) **SAM-B** (high res)

    - VL adaptor (for a+b)

- (3) Train: Focus on **text -> img**

# 61. DeepSeek-VL2 (arxiv 2024)

- DeepSeek-VL + Two upgrades

- (1) [Vision] **Dynamic tiling** vision encoding strategy

- (2) [LLM] DeepSeekMoE + MLA



Figure 2 | **Overview of DeepSeek-VL2.** The overall structure is a llava-style architecture, which includes a vision encoder, a VL adaptor, and a MoE-based LLM.

# 61. DeepSeek-VL2 (arxiv 2024)



Figure 3 | **Illustration of dynamic tiling strategy in DeepSeek-VL2**. By dividing images into multiple tiles, DeepSeek-VL2 achieves stronger fine-grained understanding capabilities compared to DeepSeek-VL.

# 62. MANTIS (TMLR 2024)

- **Multi-image** VL tasks

  - (e.g., coherence, comparison …)

- **Interleaved** text-image

- Trained with significantly **less datasets**!



Does the shop frame in **(image 1: <IMAGE> </IMAGE>)** have the same color as the vending machine in **(image 2: <IMAGE> </IMAGE>)**?

**Co-reference  Compare     Reasoning  Temporal**

**No**, they have different colors. The shop frames in the **first image is red** but the vending machine in the **second image** is **blue**. So the answer is no.

Training: Mantis-Instruct dataset

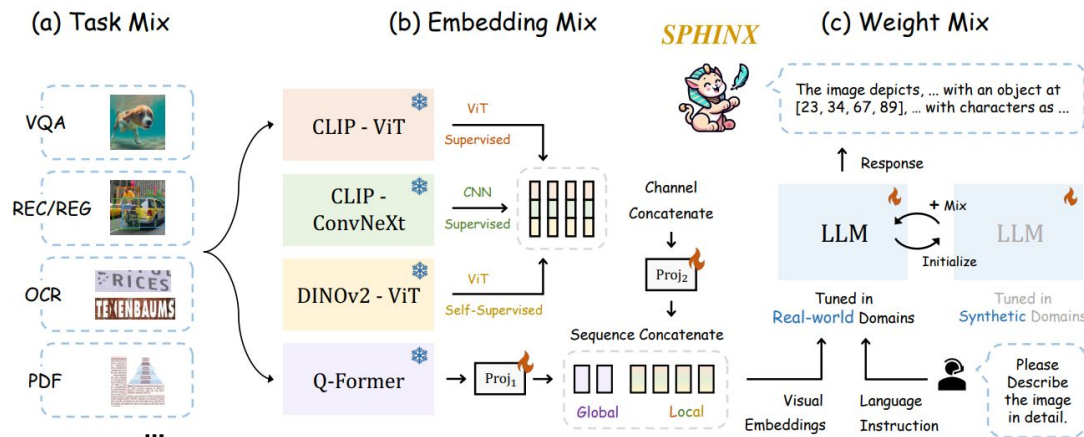| Birds-to-words | NLVR2 | Contrastive Captioning |
|---|---|---|
| Q: What are the differences in coloring and physical features between <image> and <image> in the bird images? | Q: Answer the following multiple-choice question: Here is a statement describing 2 images: There are exactly two animals in the image on the left. Is it true or false? <image> Options: (A) False (B) True | Q: Here are 3 images <image> , which image shows the following content: {caption of image 2} Options: (A) Image 1 (B) Image 2 (C) Image 3 |
| Answer: animal1 has a white belly and animal2 has a yellow and white belly . animal1 has a longer thinner body than animal2 . | Answer: B | Answer: B |

# 63. moondream (blog 2024)

For **practical usage**!

- moondream1: 1.6B
    - Model: SigLIP + Phi-1.5
    - Dataset: LLaVA
- moondream2: 1.86B
    - moondream1 + MLP projector
- moondream-next: 1.9B
    - structured output (JSON, XML, MD,CSV) …

# 64. SPHINX (arxiv 2024)

- Jointly **mixing**

  - **(1) Model weight**

  - **(2) Tasks**

  - **(3) Visual embeddings**



- 2-stage approach

  - Step 1. Pretraining) Unfreeze LLaMA-2

    - To foster more effective cross-modal learning

  - Step 2. Fine-tuning

# 65. SPHINX-X (ICML 2024)

- Use two visual encoders: **CLIP-ConvNeXt & DINOv2**

- Learnable skip tokens (to bypass fully-padded sub-images)

- Efficient single-stage training process



Figure 4: **Pipeline of SPHINX for high-resolution images.** We propose to further mix different scales and sub-images to better capture fine-grained semantics on high-resolution images.

Figure 3. **Overall paradigm of SPHINX-X family.** On top of SPHINX (Lin et al., 2023), we adopt three modifications for a more general and concise architecture: removing redundant visual encoders in Mixture of Visual Experts (MoV), bypassing fully-padded sub-images with skip tokens, and simplifying multi-stage training into a one-stage all-in-one approach.

# 66. KOSMOS-1 (NeurIPS 2023)

- Can perceive **general modalities**

- Align various modalities with LLM

# 67. KOSMOS-2 (ICLR 2024)

- KOSMOS-1에 새로운 capability

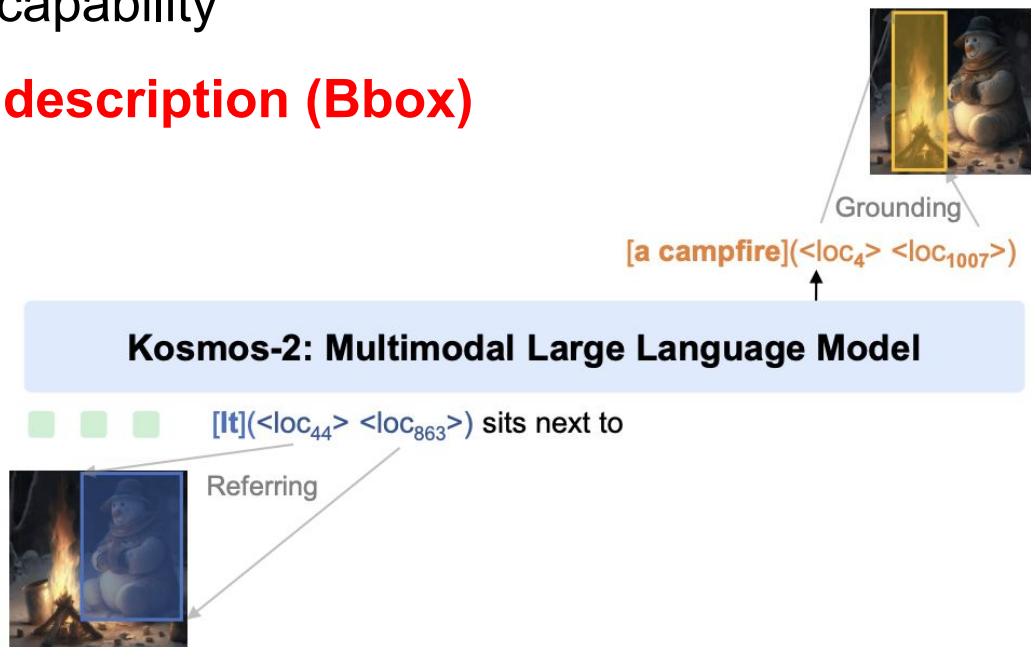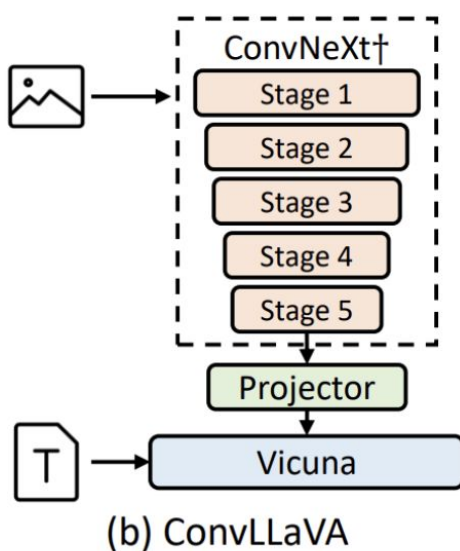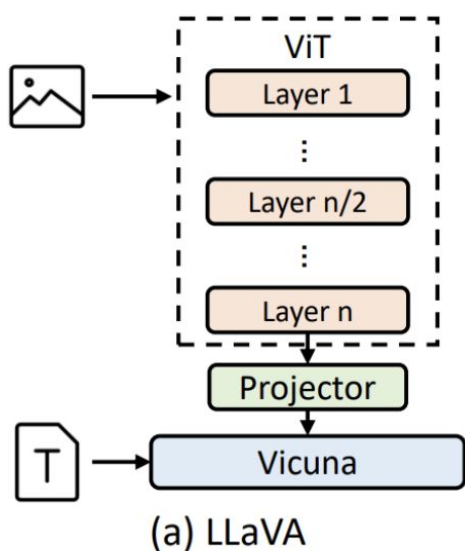  - Perceiving **object description (Bbox)**

  - **Grounding text**



Figure 1: KOSMOS-2 is a multimodal large language model that has new capabilities of multimodal grounding and referring. KOSMOS-2 can understand multimodal input, follow instructions, perceive object descriptions (*e.g.*, bounding boxes), and ground language to the visual world.

# 68. ConvLLaVA (arxiv 2024)

- High-resolution LMM에서 **ViT의 한계점** 을 극복하기 위해!

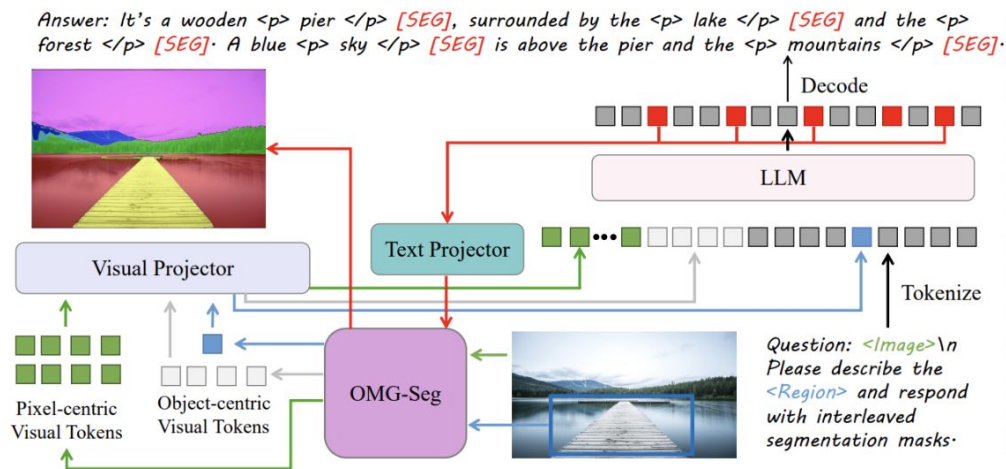- Hierarchical backbone "**ConvNeXT**"를 대신 사용함



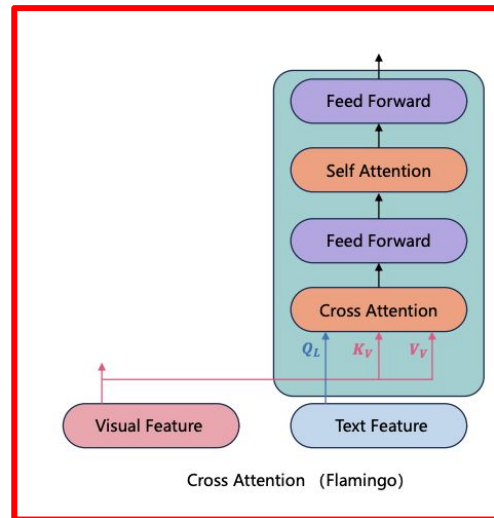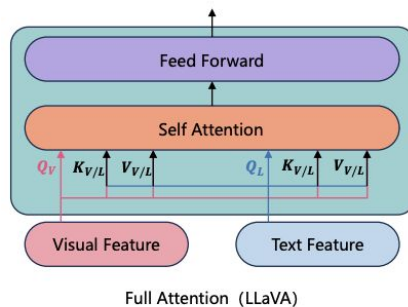(a) LLaVA     (b) ConvLLaVA     (c) Training Procedure

# 69. OMG-LLaVA (arxiv 2024)

- Unifies (1) **Image**-level + (2) **Object**-level + (3) **Pixel**-level reasoning

- Architecture:

  - (1) Vision: **OMG-seg (segmentation model)**
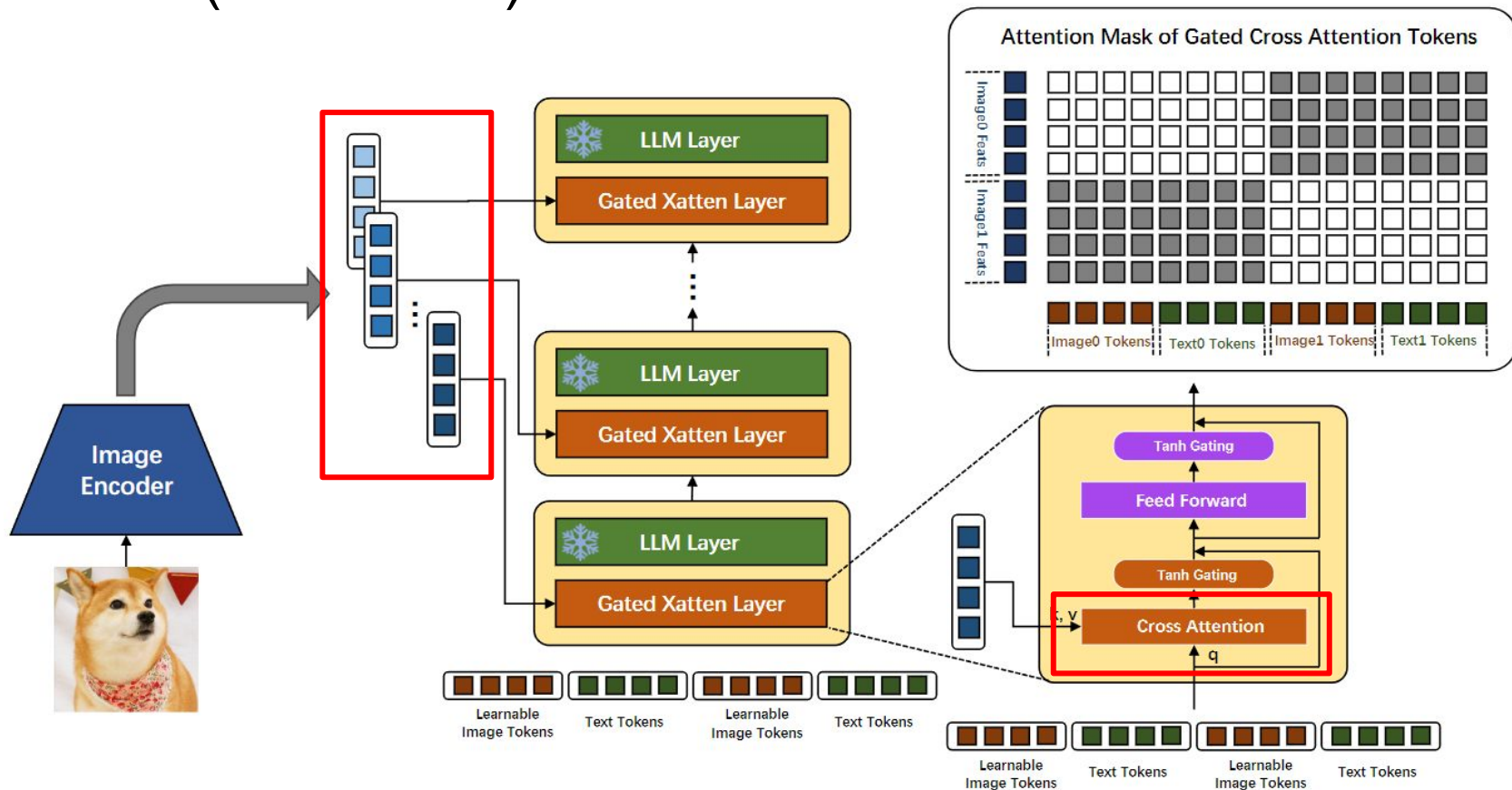
  - (3) LLM

# 70. EVLM (arxiv 2024)

- Architecture: **Flamingo**

- Goal: Minimize computational cost!

- Challenges of handling long sequences of visual signals

- How?

  - **Cross-attention**

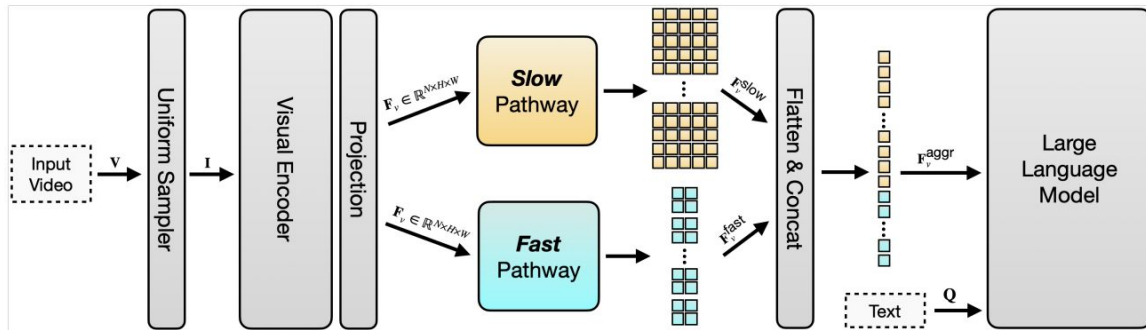  - **Hierarchical ViT features**

- **(Q-former처럼) learnable tokens**



Full Attention (LLaVA)

Cross Attention (Flamingo)

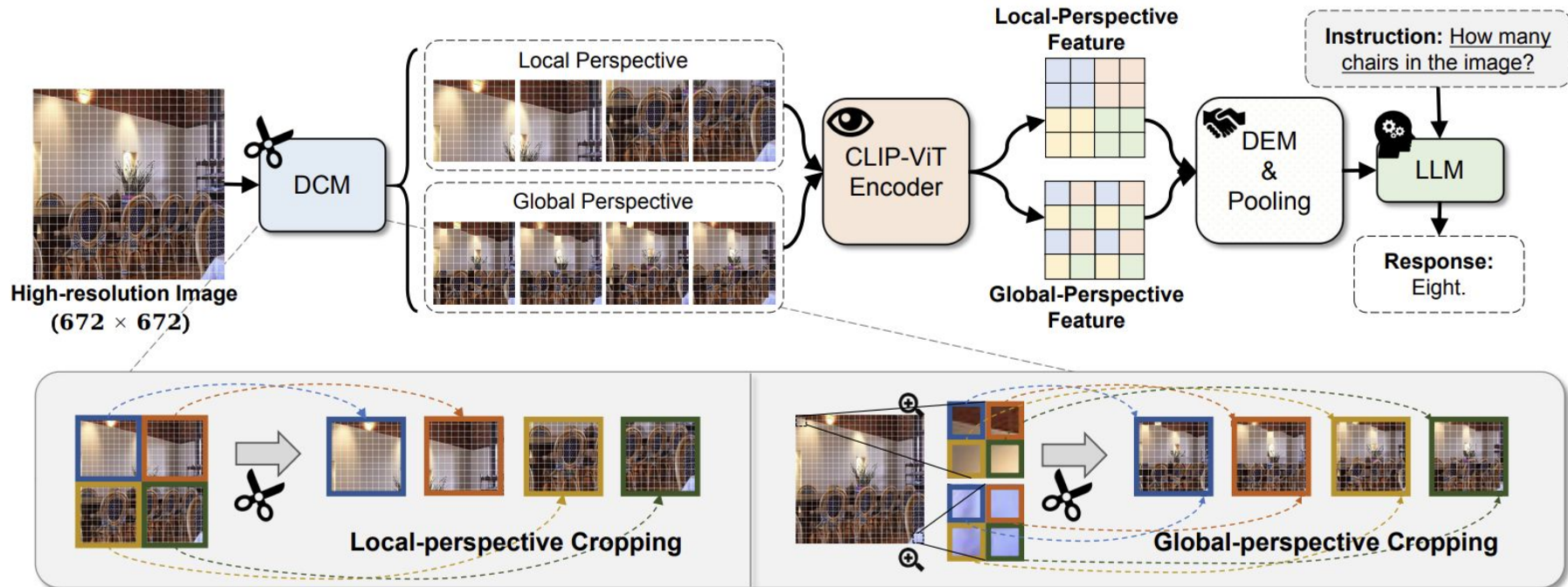# 70. EVLM (arxiv 2024)

# 71. SlowFast-LLaVA (arxiv 2024)

- **Training-free** Video LLM

- (1) Uniform sampler: Sample N frames from video

- (2) Pathways

  - **Slow** pathway: Detailed spatial semantics

  - **Fast** pathway: Long-range temporal context

- (3) LLM: LLaVA-Next

# 72. INF-LLaVA (arxiv 2024)

- MLLM to efficiently process **high-resolution images**

- 기존 연구: cropping-based & dual-encoder

- Proposal

  - **Dual-perspective** cropping module (DCM)

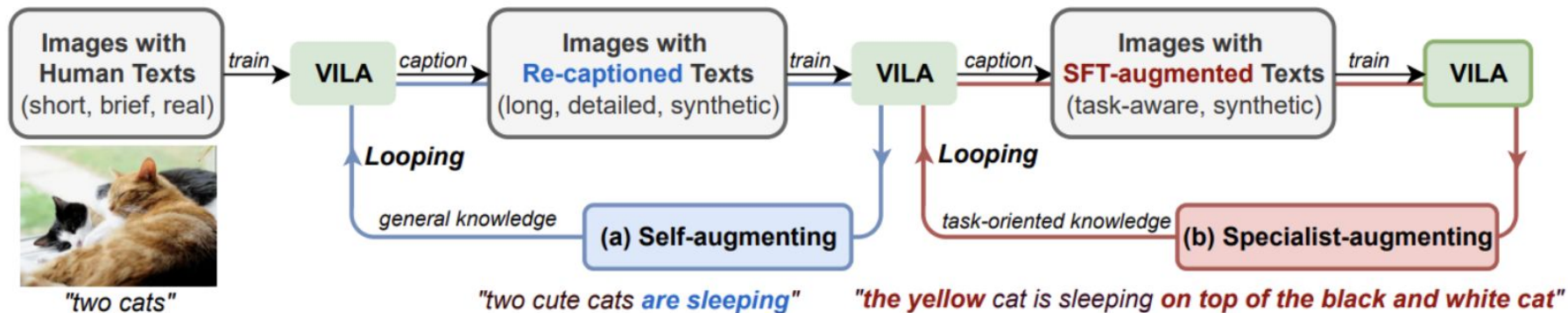  - **Dual-perspective** enhancement module (DEM)

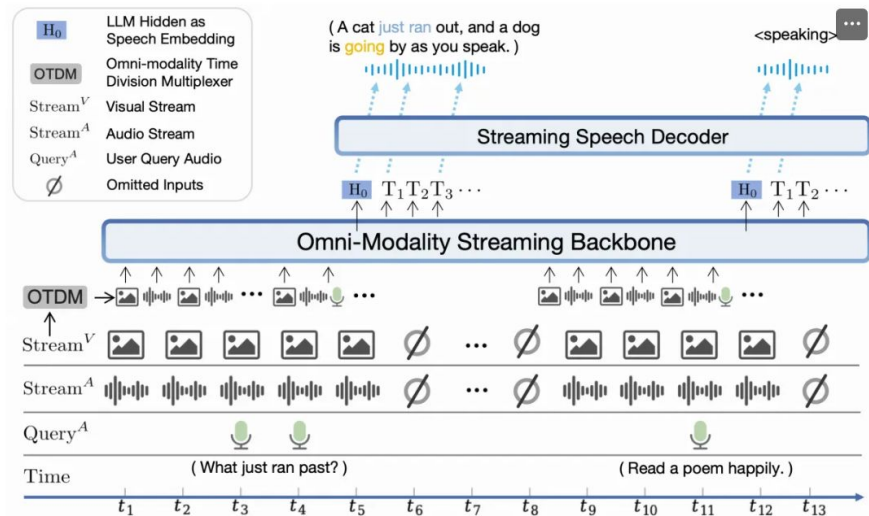# 72. INF-LLaVA (arxiv 2024)

# 73. **VILA²** (arxiv 2024)

- Limitation of data quantity & quality in VLMs

- Costly human annotation, distillation 대신…

   => Leverage VLM **itself** to **refine & augment** pretraining data!
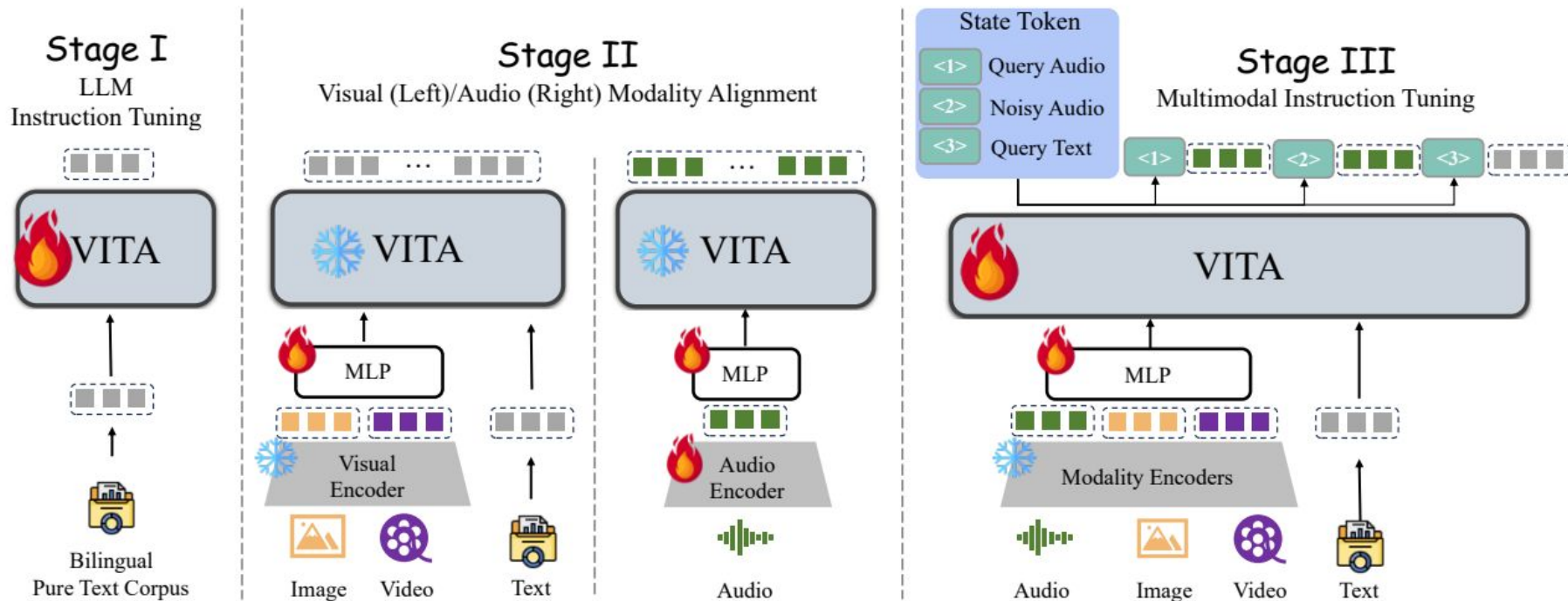
# 74. MiniCPM-o-2.6 (arxiv 2024)

- **8B** MLLM

  - Efficient for deployment on **edge devices**

- **Vision, Speech, and Multimodal** live streaming

- Architecture

  - Vision: SigLIP-400M

  - Audio: Whisper-medium-300M

  - TTS: ChatTTS-200M

  - LLM: Qwen2.5-7B
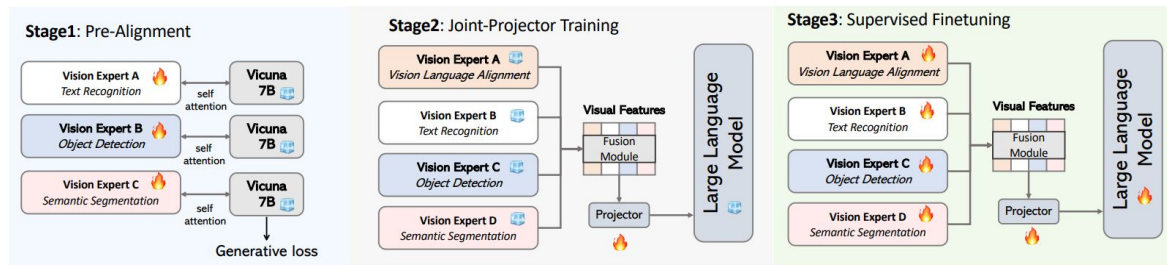
# 75. VITA (arxiv 2024)

- First **open-source** MLLM capable of processing video, image, text, audio

- Architecture

  - LLM: Mixtral 8x7B

  - Vision: InternViT-300M-448px

- Three stage training

  - Stage 1) **LLM** instruction tuning

  - Stage 2) Multimodal **alignment**

  - Step 3) Multimodal **instruction tuning**

# 75. VITA (arxiv 2024)

# 76. EAGLE (ICLR 2025)

- **Open-source** MLLMs

- Mixture of vision encoders

- Particularly in OCR and document understanding

- Architecture: LLaVA

    - (1) LLM

    - (2) Vision encoder

    - Projection layer



- 3 stages of training

# 77. EAGLE2 (arxiv 2025)

- "**Data-centric**" approach

  - Prioritize data diversity & quality

- Tiled Mixture of vision encoders (MoVE)

  - SigLIP & ConvNeXt-XXLarge

  - Image tiling to handle high resolution

- LLM: Qwen2.5

- Build upon open-source components
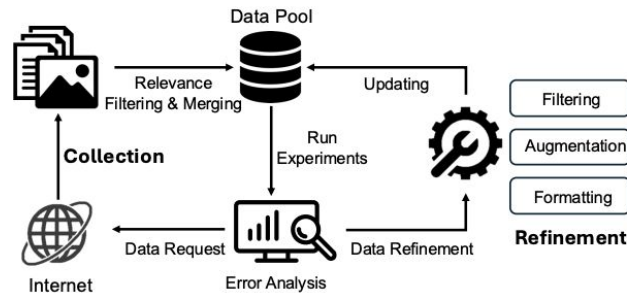
- Three-stage training recipe



Figure 3 | An overview of our data strategy. The upper part shows the date collection pipeline and the lower part shows the data refinement pipeline.
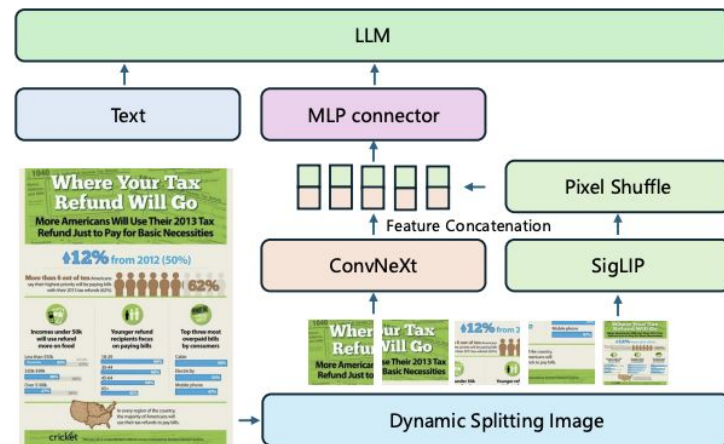


Figure 11 | Tiled Mixture of Vision Encoders.
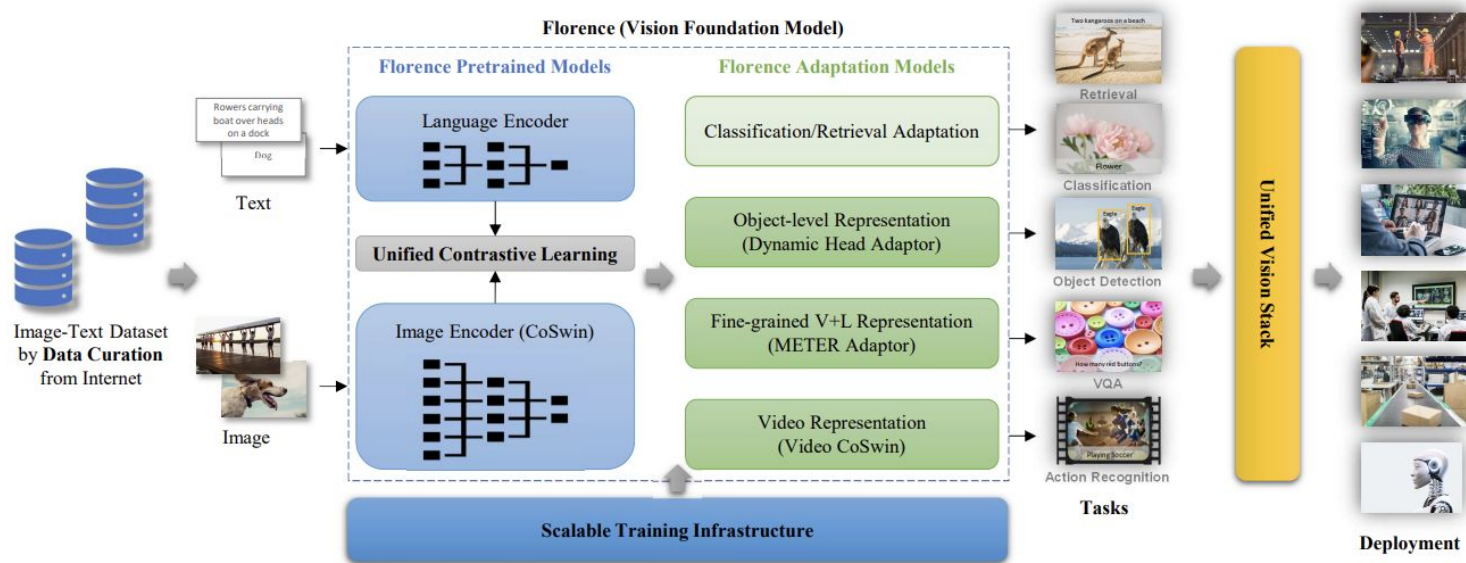
# 78. Florence (arxiv 2021)



*Figure 2.* Overview of building Florence. Our workflow consists of data curation, unified learning, Transformer architectures and adaption. It shows the foundation model can be adapted to various downstream tasks and finally integrated into modern computer vision system to power real-world vision and multimedia applications. Compared with existing image-text pretraining models (Radford et al., 2021; Jia et al., 2021; Wud), mainly limited on cross-modal shared representation for classification and retrieval (illustrated by *light-green* adaptation module), Florence expands the representation to support object level, multiple modality, and videos respectively.
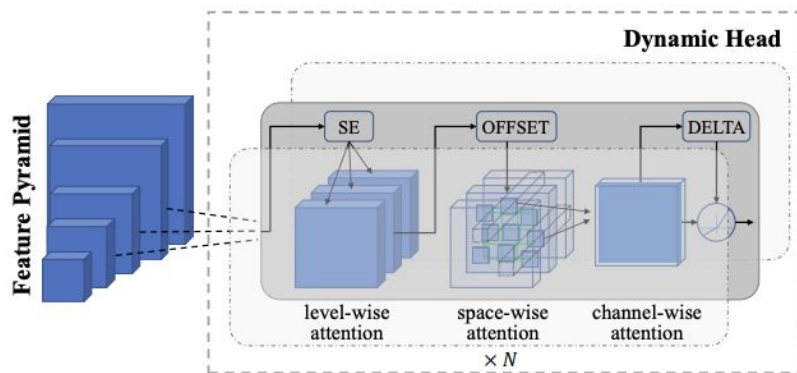
# 78. Florence (arxiv 2021)



*Figure 3.* Dynamic Head (Dai et al., 2021a) adapter is used for object-level visual representation learning.
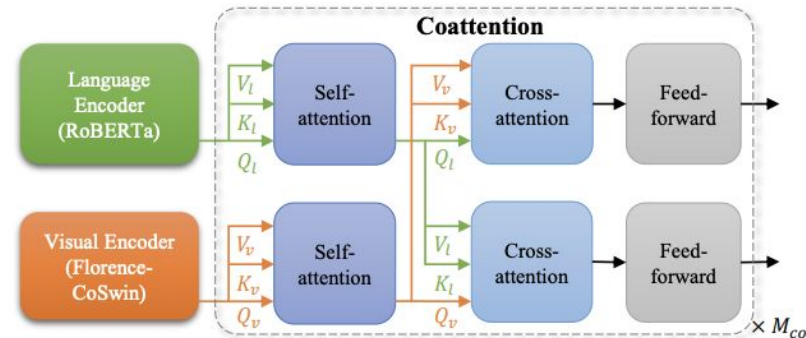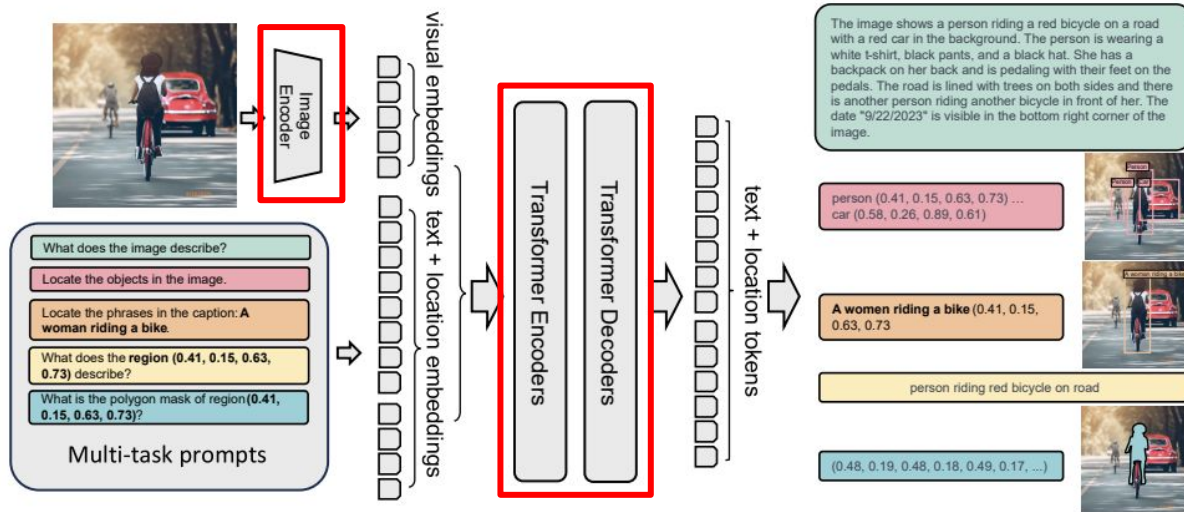


*Figure 4.* METER (Dou et al., 2021) is used as Florence V+L adaptation model, trained with the image-text matching (ITM) loss and the masked language modeling (MLM) loss.

# 79. Florence-2 (CVPR 2024)

- Two key components

  - (1) Image encoder (DaViT arch)

  - (2) Multimodality encoder-decoder.

# 80. MultiInstruct (arxiv 2022)

- To enhance multi-modal zero-shot learning

    - By **Instruction tuning**

- Built upon **OFA (One For All)** as pretrained model

    - Single Transformer-based **sequence-to-sequence** framework

    - Pre-trained on a diverse set of multimodal and unimodal tasks
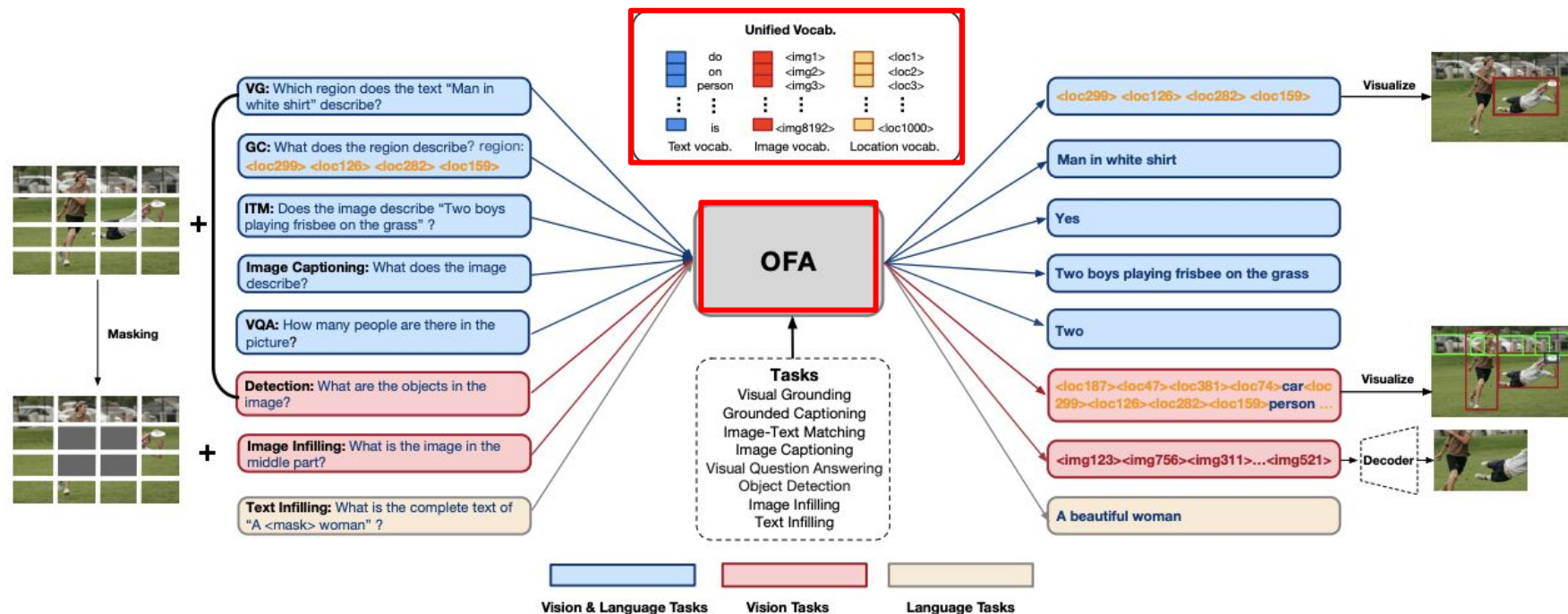
# 80. MultiInstruct (arxiv 2022)



Figure 2: A demonstration of the pretraining tasks, including visual grounding, grounded captioning, image-text matching, image captioning, VQA, object detection, image infilling as well as text infilling.

# 80. MultiInstruct (arxiv 2022)



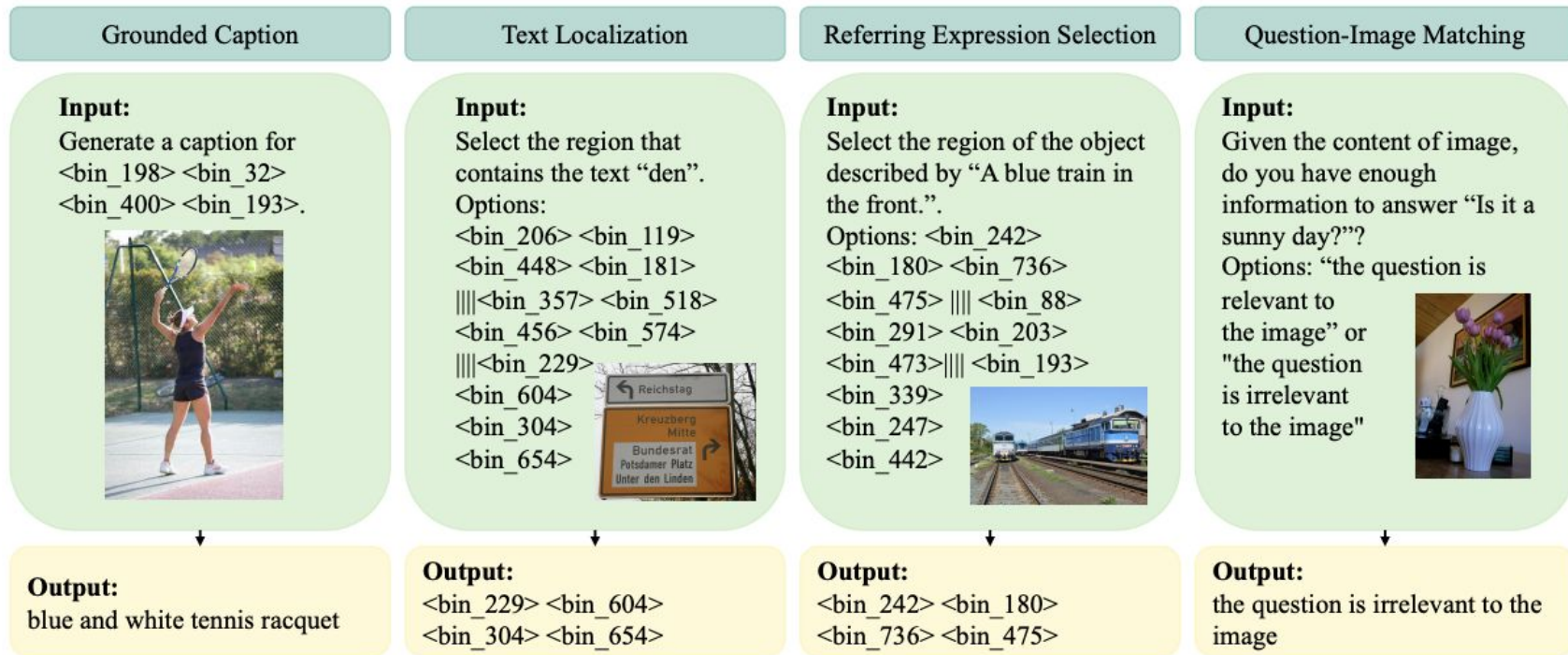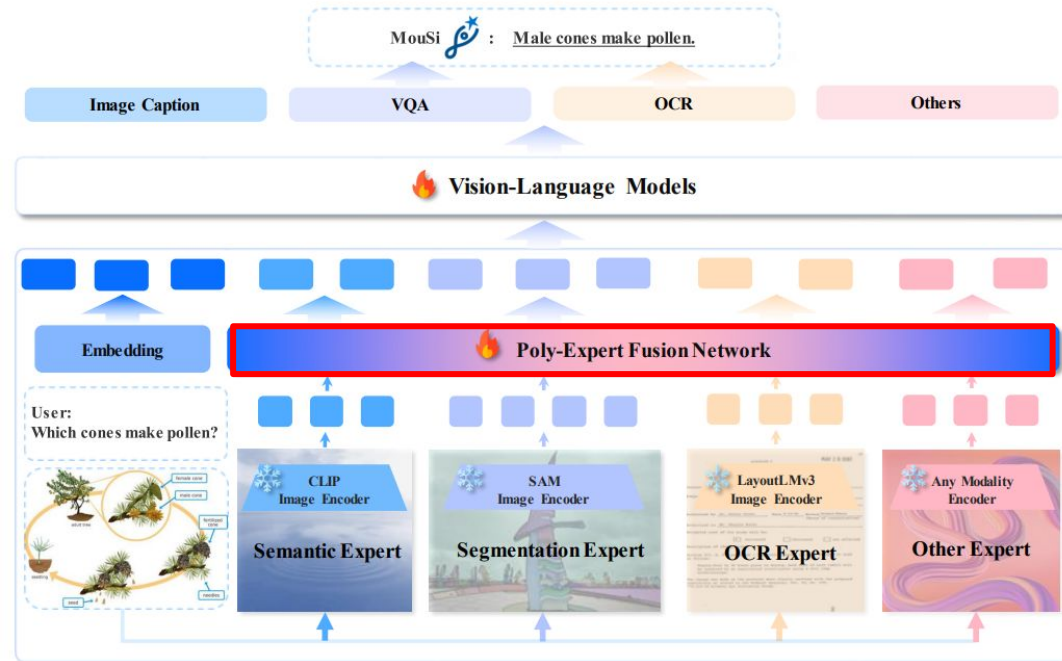| Grounded Caption | Text Localization | Referring Expression Selection | Question-Image Matching |
|---|---|---|---|
| **Input:**<br>Generate a caption for <bin_198> <bin_32> <bin_400> <bin_193>. | **Input:**<br>Select the region that contains the text "den".<br>Options:<br><bin_206> <bin_119><br><bin_448> <bin_181><br>\|\|\|\|<bin_357> <bin_518><br><bin_456> <bin_574><br>\|\|\|\|<bin_229><br><bin_604><br><bin_304><br><bin_654> | **Input:**<br>Select the region of the object described by "A blue train in the front.".<br>Options: <bin_242><br><bin_180> <bin_736><br><bin_475> \|\|\|\| <bin_88><br><bin_291> <bin_203><br><bin_473>\|\|\|\| <bin_193><br><bin_339><br><bin_247><br><bin_442> | **Input:**<br>Given the content of image, do you have enough information to answer "Is it a sunny day?"?<br>Options: "the question is relevant to the image" or "the question is irrelevant to the image" |
| **Output:**<br>blue and white tennis racquet | **Output:**<br><bin_229> <bin_604><br><bin_304> <bin_654> | **Output:**<br><bin_242> <bin_180><br><bin_736> <bin_475> | **Output:**<br>the question is irrelevant to the image |

Figure 1: **Example Instances from MULTIINSTRUCT for Four Tasks.**
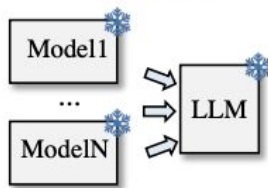
# 81. MouSi (COLM 2024)

- Incorporate **multiple visual experts** (e.g., CLIP and SAM) to VLMs

- **Poly-expert fusion network**

  - Combine their outputs

    (e.g., CLIP, SAM)

  - Interface with LLMs
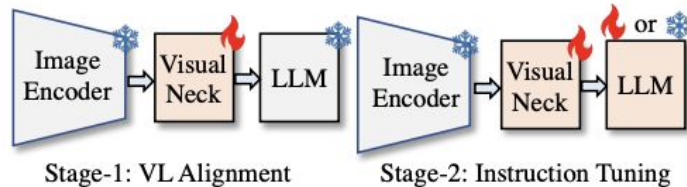
    (e.g., Vicuna)

# 82. LaVIN (NeurIPS 2023)

- **Cost-effective** approach to VL instruction-tuning

- How? **Mixture-of-Modality Adapter (MM-Adapter)**

  - Significantly reduce # of trainable parameters

  - Adapt visual features to LLM



(a) Expert System  (b) Modular Training Scheme  (c) Mixture-of-Modality Adaptation
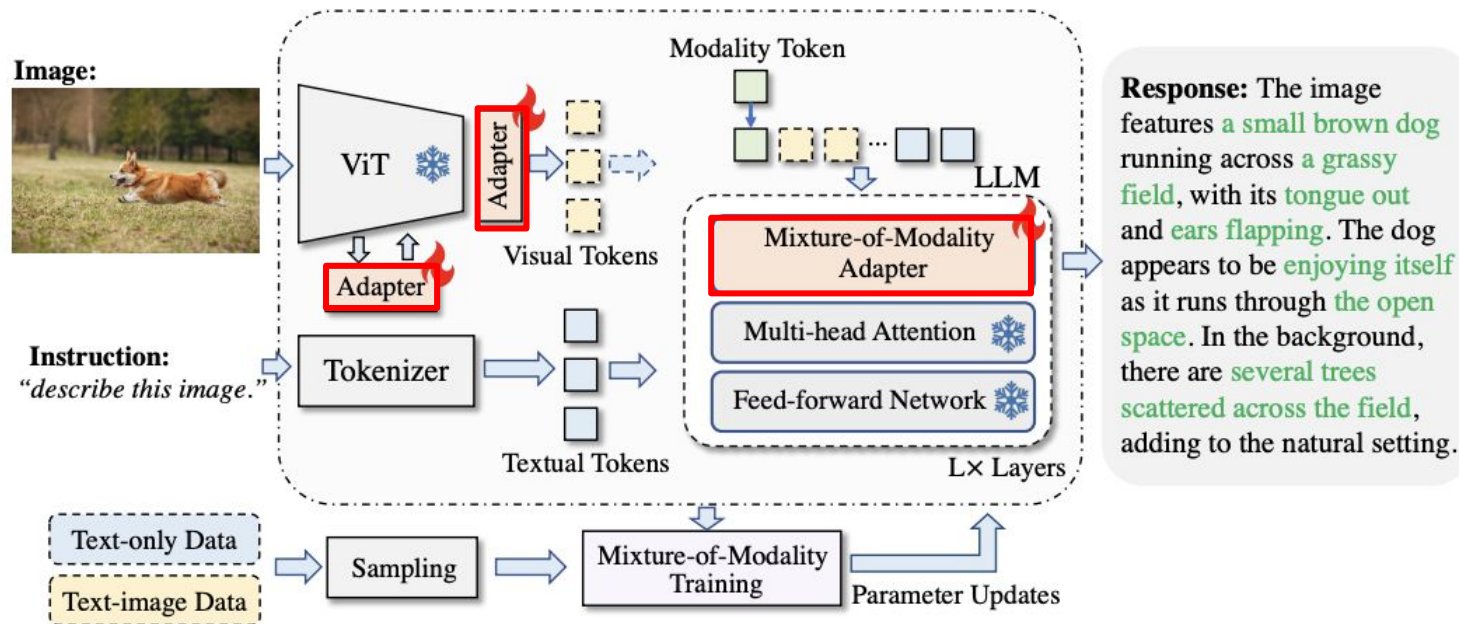
# 82. LaVIN (NeurIPS 2023)



Figure 2: The overview of the Mixture-of-Modality Adaptation (MMA) and the architecture of LaVIN. In LaVIN, the novel Mixture-of-Modality Adapters are employed to process the instructions of different modalities. During instruction tuning, LaVIN is optimized by Mixture of Modality Training (MMT) in an end-to-end manner.

# 83. TinyGPT-V (ICMLW 2024)

- Combining

  - (1) Vision: compact **EVA-ViT**

    - with **linear projection layers (+Q-Former layer)**

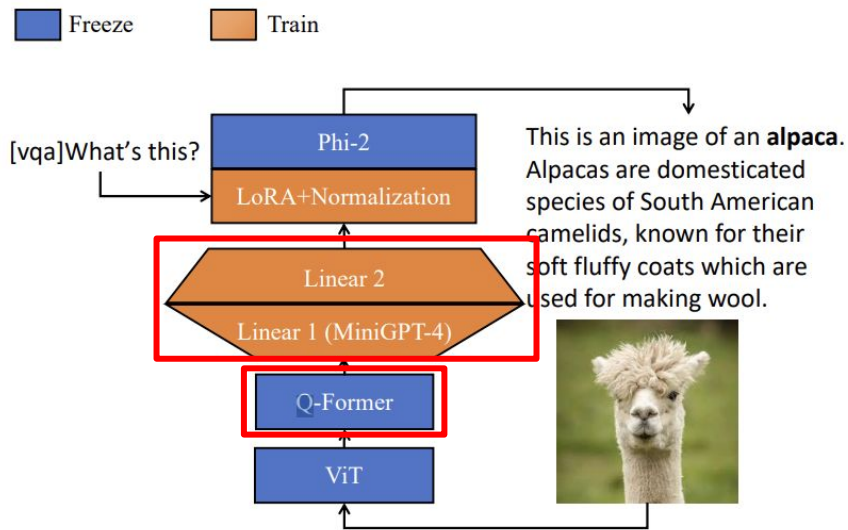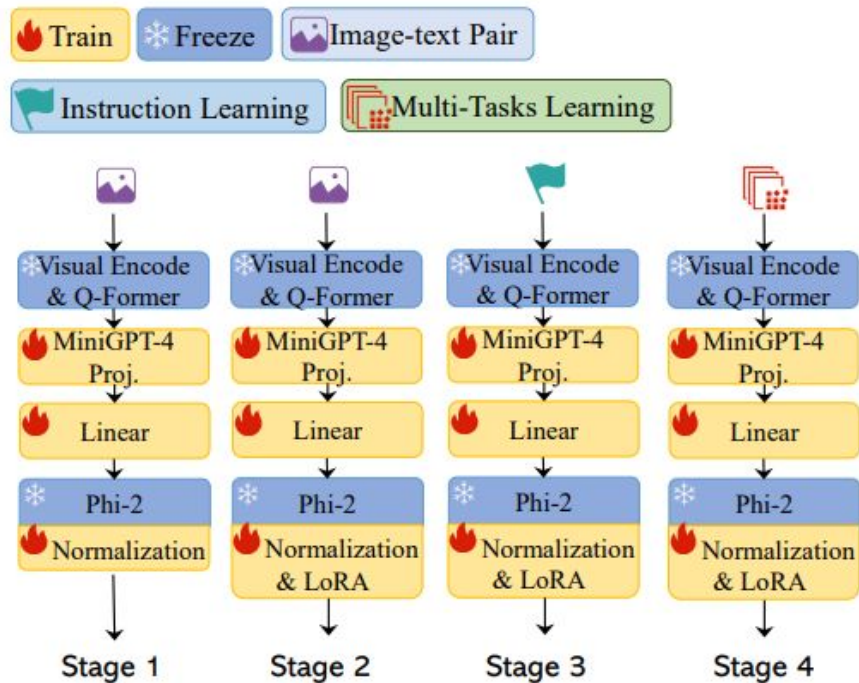  - (2) LLM: powerful **Phi-2**
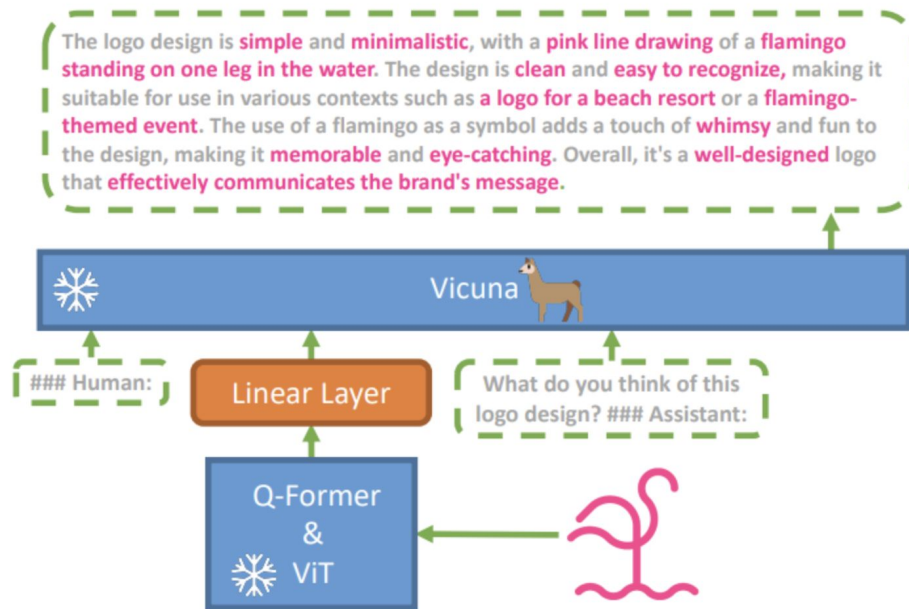
# 83. TinyGPT-V (ICMLW 2024)



Figure 4: Architecture of TinyGPT-V. The model takes a visual backbone, which remains frozen during all training phases. We concatenate Q-Former layer visual output tokens from ViT backbone and project them into Phi-2 language model space via two linear projection layers.
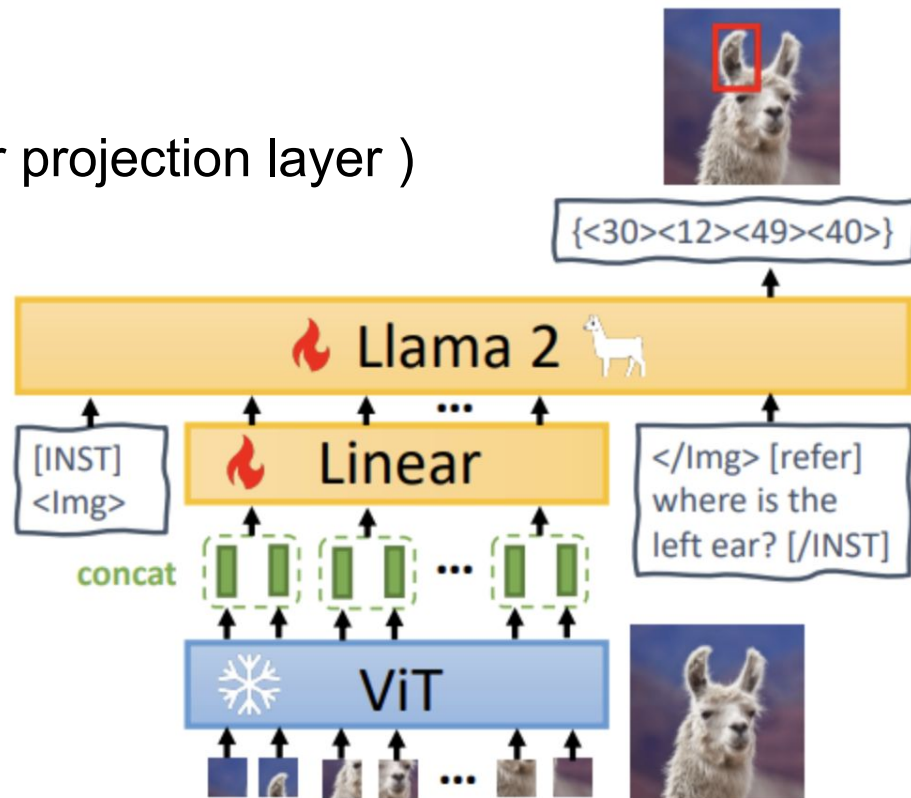
# 84. MiniGPT-4 (ICLR 2024)

- Architecture

  - (1) Vision: **Pretrained ViT + Q-Former (+ linear proj)**

  - (2) LLM: Vicuna

# 85. MiniGPT-v2 (ICLR 2024)

- Architecture

  - (1) Vision: static ViT ( + linear projection layer )

  - (2) LLM: LLaMA-2-chat

# 86. CoVLM (ICLR 2024)

- Propose **communication tokens**

- To enable dynamic interaction btw…

  - (1) CLIP ViT-L image encoder

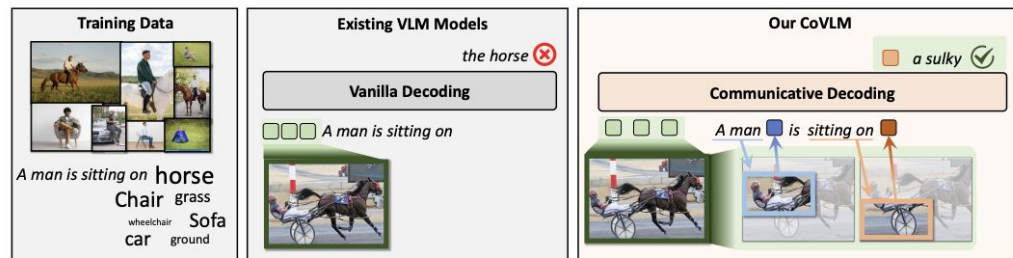  - (2) YOLOX detection network

  - (3) Pythia LLM



Figure 1: Comparison with existing VLMs. Previous models take in a whole image as input, impairing the compositionality of VLMs. Our CoVLM inserts communication tokens into the LLM after visual entities / relationships to enable the language-to-vision and vision-to-language communication, improving compositionality to a large extent.
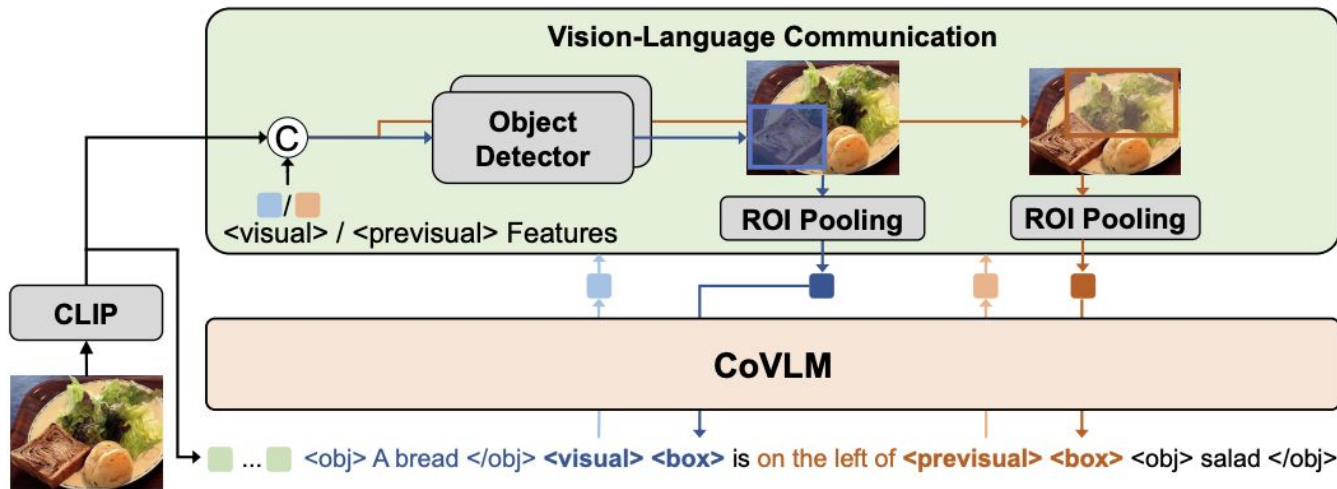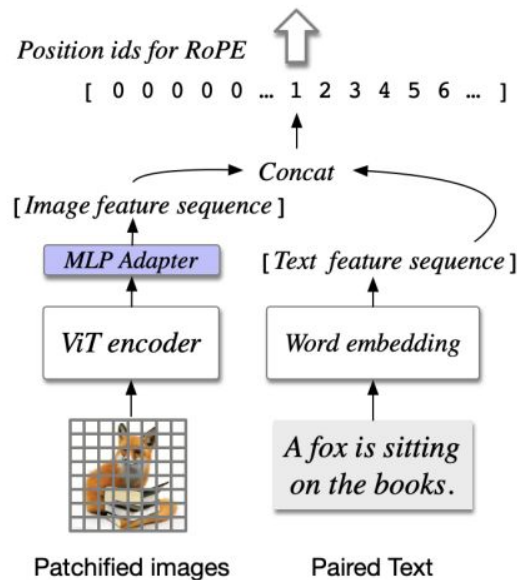
# 86. CoVLM (ICLR 2024)



Figure 2: Overview of our CoVLM framework. Our vision module consists of a CLIP encoder to encode the image, and an object detector which takes in the image together with language inputs to generate relevant regions. For language modelling, we insert a set of communication tokens into the LLM, which can appear after a visual entity with a `<visual>` token or after a relationship with a `<previsual>` token. The last hidden layer of the LLM is then sent to the object detector to propose regions relevant to the language inputs so far. This is termed as top down language-to-vision communication. Next, in vision-to-language communication, the features of the proposed regions are fed back to LLM via `<box>` or `<prebox>` token for further language generation.

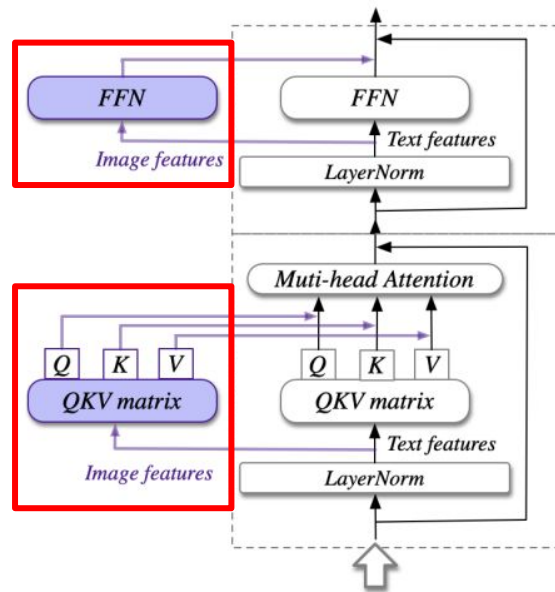# 87. CogVLM (NeurIPS 2024)

- Enhance pretrained LM with **Visual Expert** module

    - Employ a visual expert module in each layer

      that leverages **QKV matrix & MLP**

- Components

    - (1) ViT encoder

    - (2) LLM

    - MLP adapter to (1)->(2)

# 87. CogVLM (NeurIPS 2024)



(a) The input of visual language model

(b) The visual expert built on the language model

Figure 3: **The architecture of CogVLM.** (a) The illustration about the input, where an image is processed by a pretrained ViT and mapped into the same space as the text features. (b) The Transformer block in the language model. The image features have a different QKV matrix and FFN. Only the purple parts are trainable.

# 88. CogVLM2 (arxiv 2024)

- VLM for **image and video**

- Adapter: Incorporates …
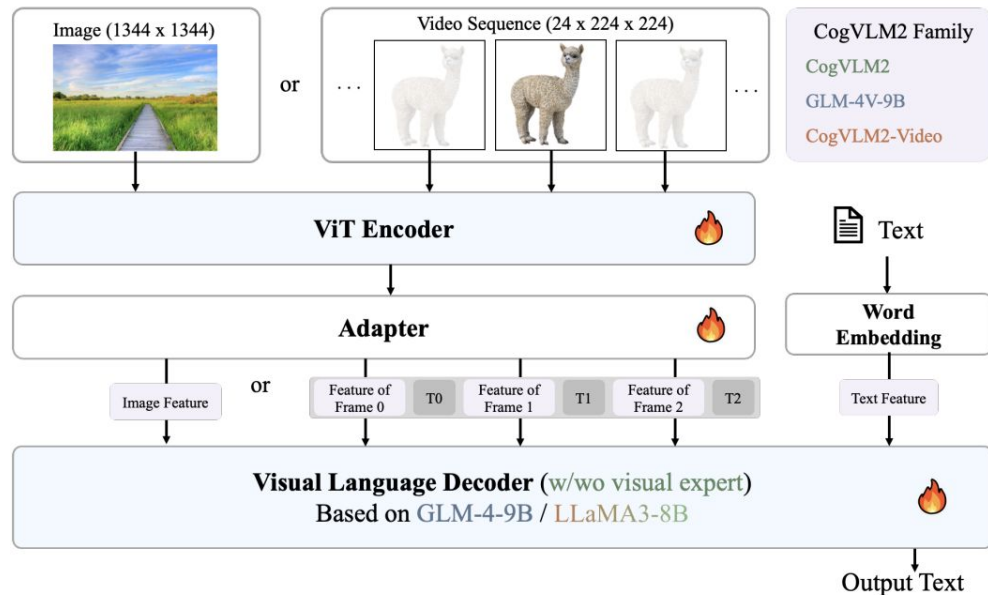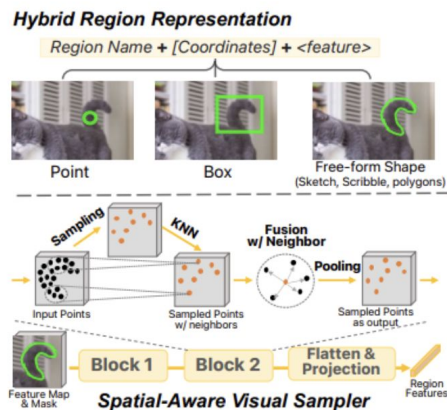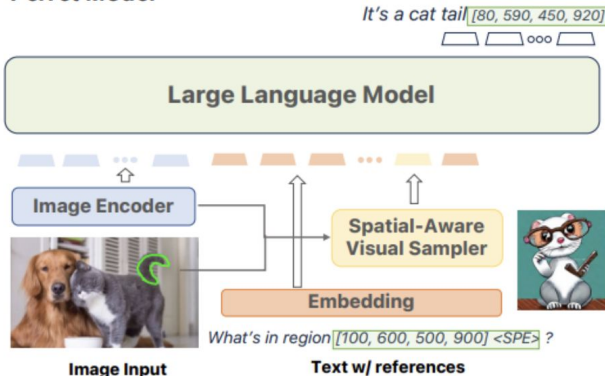
  - Conv for downsampling

  - SwiGLU module



Figure 2: The architecture of the CogVLM Family. Taking a high resolution image or the extracted frames from a given video, CogVLM models embed visual information with a pre-trained ViT Encoder and an Adapter. The embedded visual features are sent to a Visual Language Decoder. CogVLM2-Video is capable of answering image-related and video-related queries.

# 89. Ferret (ICLR 2024)

- **Hybrid** region representation

  - Discrete coordinates

  - Continuous features

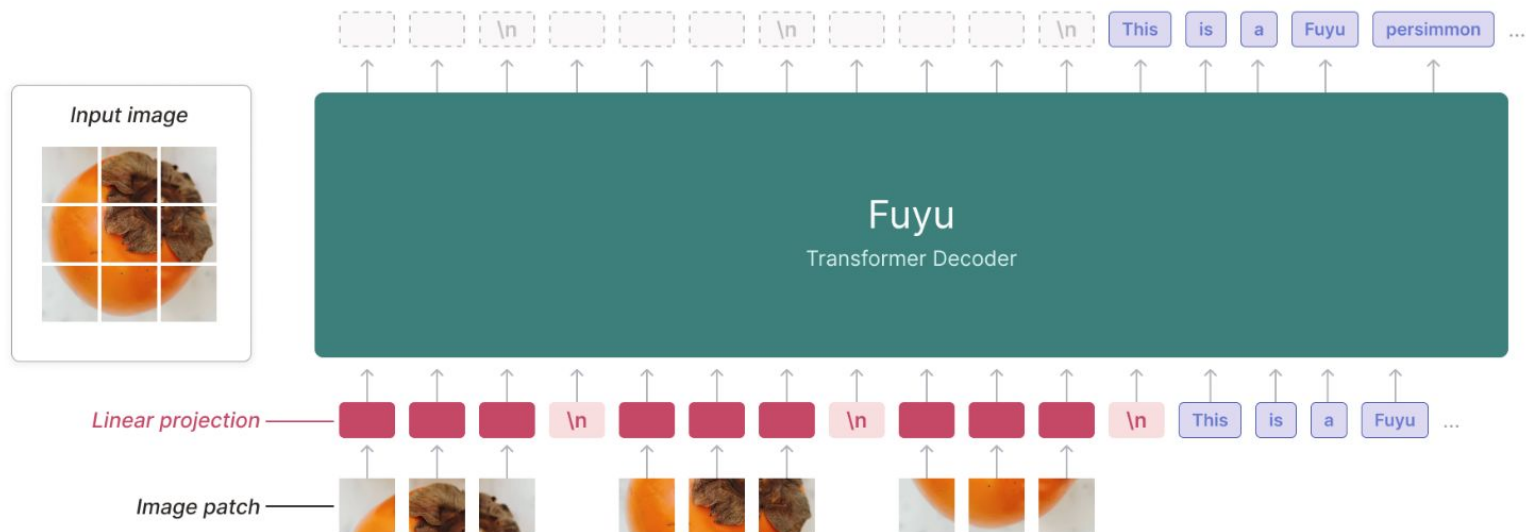- Allow it to precisely pinpoint objects & regions within images

# 90. Fuyu-8B (blog 2023)

- **Directly** projecting images patches into (decoder-only) LLM

- Treat **image & text tokens uniformly**

- Efficient!

# 91. OtterHD (arxiv 2023)

- Inspired by Fuyu-8B

- **Directly integrates pixel-level** info into LLM using position embedding

  - Separate vision encoding X

  - Varying image sizes O

  - High-resolution images up to 1024x1024

# 92. GLaMM (CVPR 2024)

- Pixel-level grounding

- 5 component architecture

  - Global & Regional image encoders

  - LLM

  - Grounding Image encoder

  - Pixel decoder

# 93. COSMO (arxiv 2024)

- ViT + (partitioned) LLM

  - LLM half: **unimodal**

  - LLM half: **mulitmodal**

- Interleaved data sequences

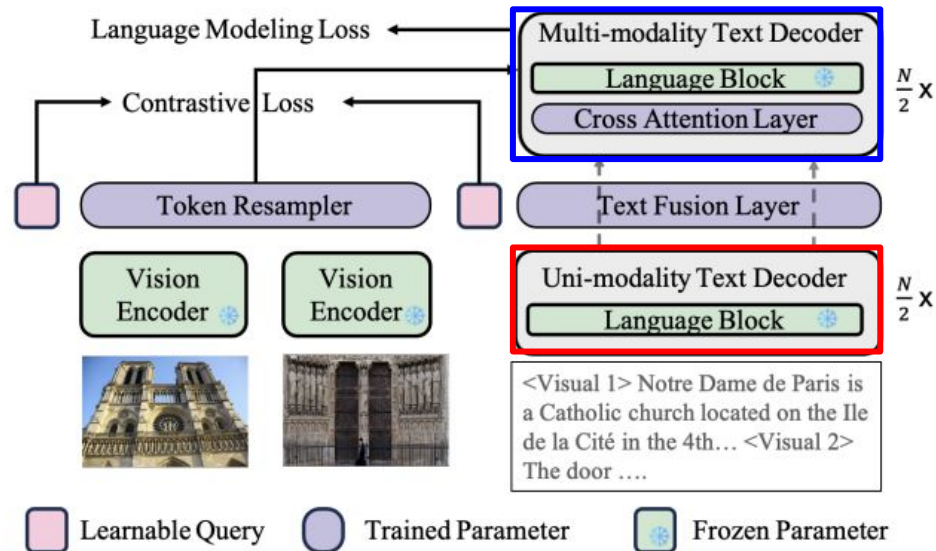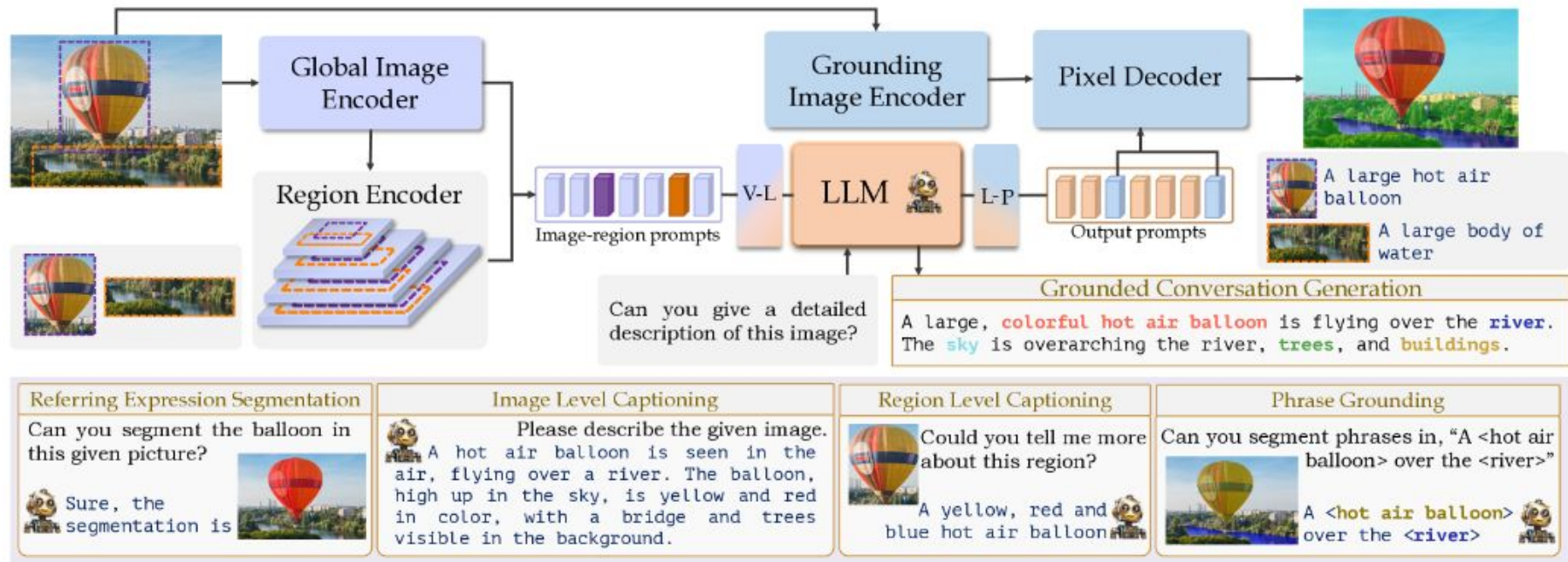- Pretraining: LM loss + CL loss



Figure 3. **An introduction to CosMo**: This model handles both image/video text pairs and inter-level image/video text pairs. The Large Language Model is divided into two parts to compute contrastive loss and language modeling loss.

# 93. COSMO (arxiv 2024)

# 94. UNIT (NeurIPS 2024)

- Vision Encoder (ViT)에 **Text Recognition 능력**을 심어주자!

  => **Single** (Vision) Model로써 전부 처리 가능!

- Two components (both lightweight)

  - (1) **Language** decoder: Text 생성 위해

  - (2) **Vision** decoder: Visual 능력 손실 방지

# 94. UNIT (NeurIPS 2024)