



北京航空航天大學
BEIHANG UNIVERSITY

生成式AI与大模型第5讲

多模态对齐

Beihang University

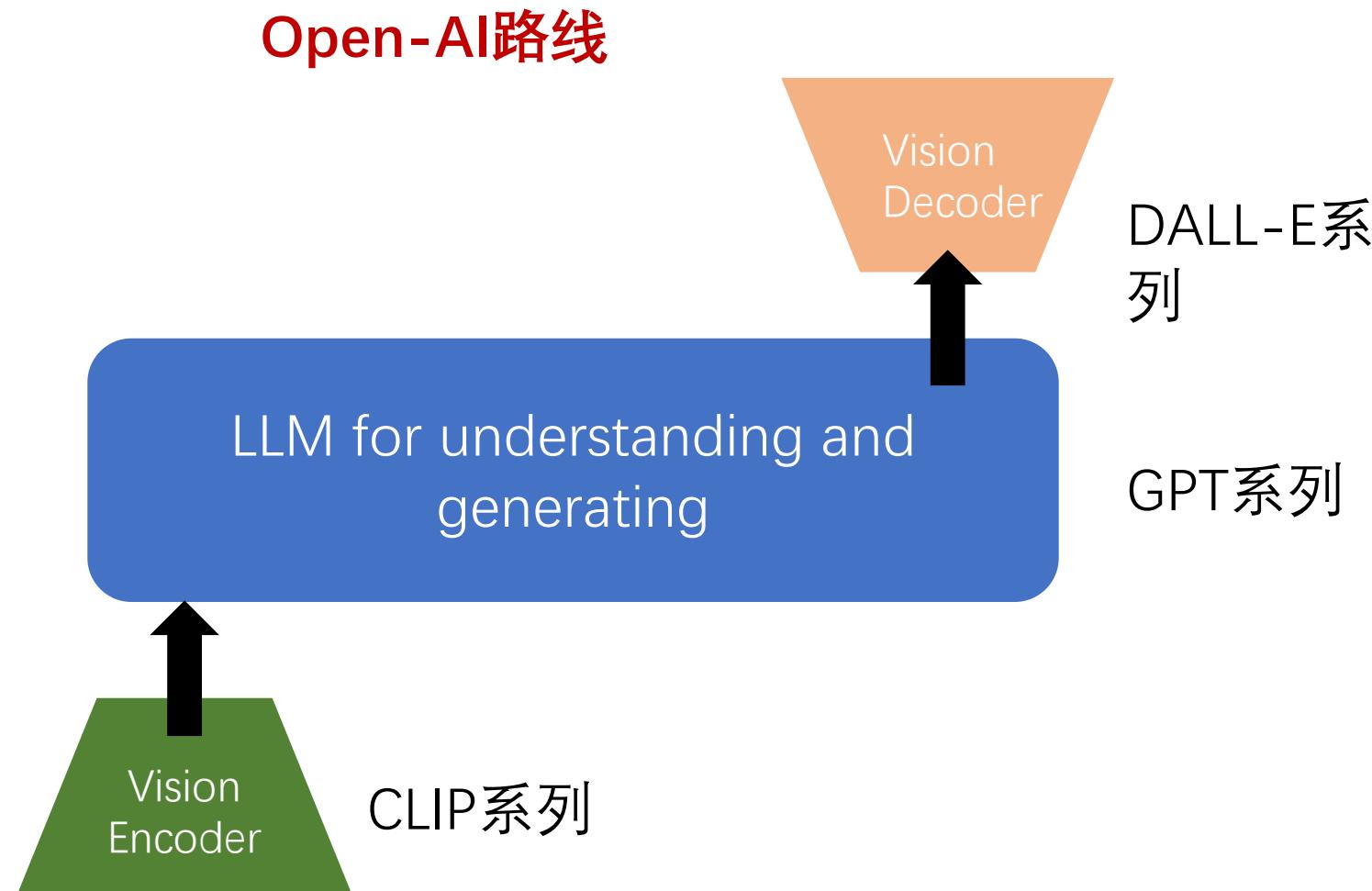
人工智能研究院
黄雷

多模态和视觉大模型

图像：

文本：

图像：



00

语言自回归模型（回顾）

Language Model

- Auto-Regressive (自回归)
 - GPTs

(all characters)

One-hot Encoding

$$\text{深} = [1 \ 0 \ 0 \ 0 \ 0]$$

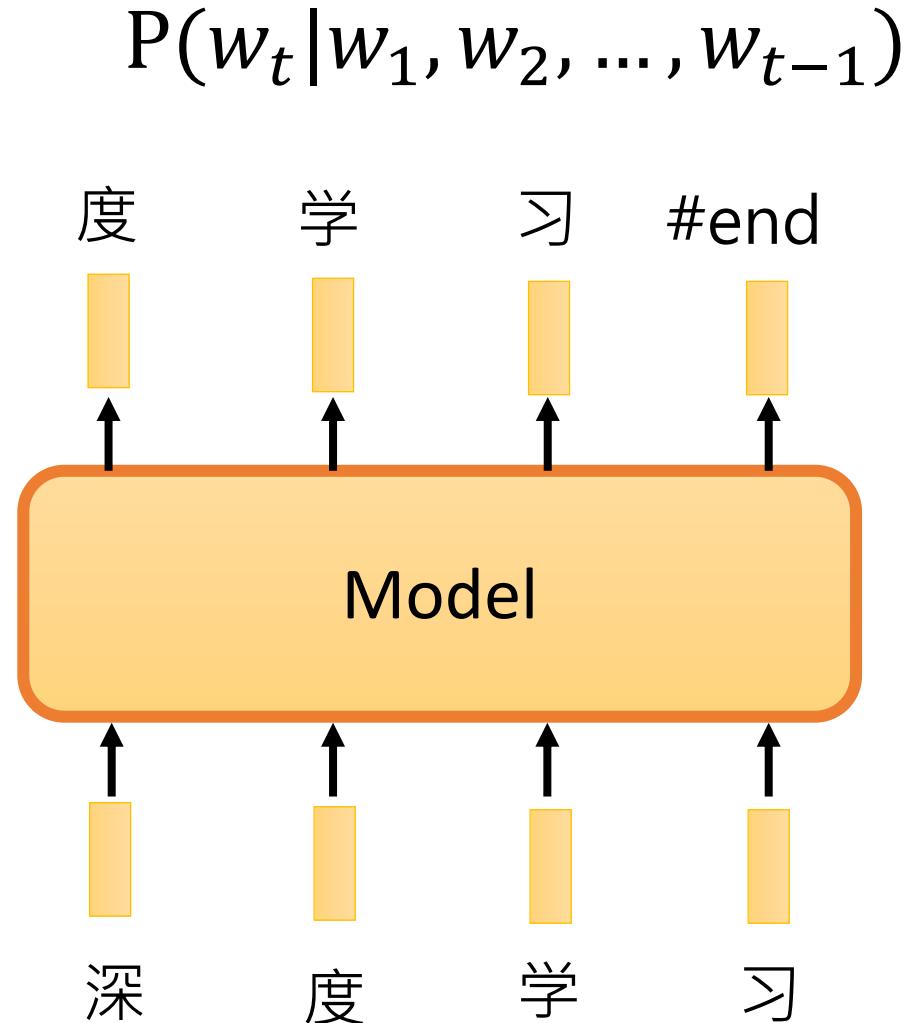
$$\text{度} = [0 \ 1 \ 0 \ 0 \ 0]$$

$$\text{学} = [0 \ 0 \ 1 \ 0 \ 0]$$

$$\text{习} = [0 \ 0 \ 0 \ 1 \ 0]$$

$$\#\text{end} = [0 \ 0 \ 0 \ 0 \ 1]$$

深	0.1
度	0.7
学	0.1
习	0.05
#end	0.05



Vision Model

Explicit density model

$$p(x) = p(x_1, x_2, \dots, x_n)$$



Likelihood of
image x

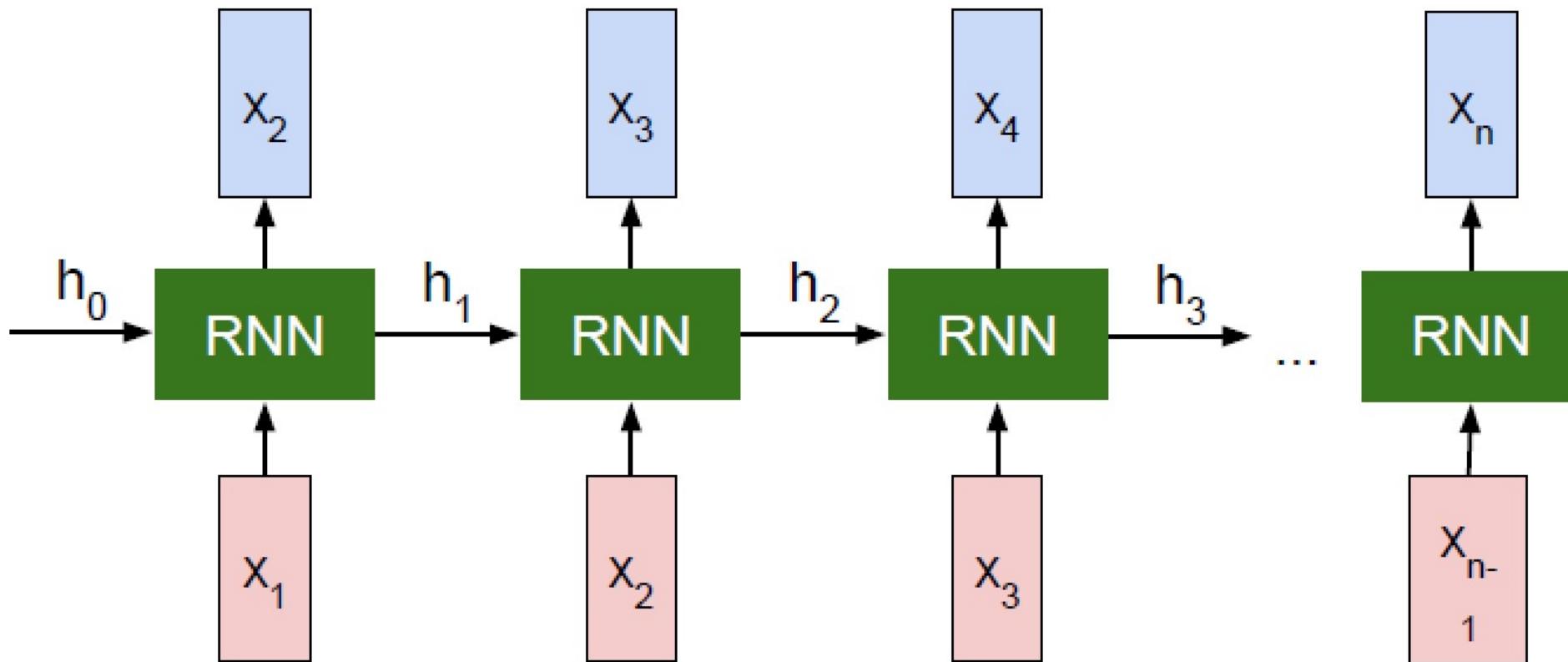


Joint likelihood of each
pixel in the image

Fully visible belief network(FVBN)

Incontext Modelling

Recurrent Neural Network(RNN)



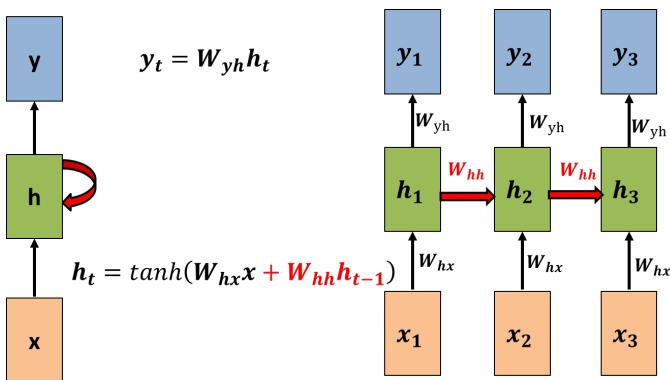
$$p(x_i | x_1, \dots, x_{i-1})$$

Incontext Modelling

$$X \in \mathcal{R}^{N \times d}$$

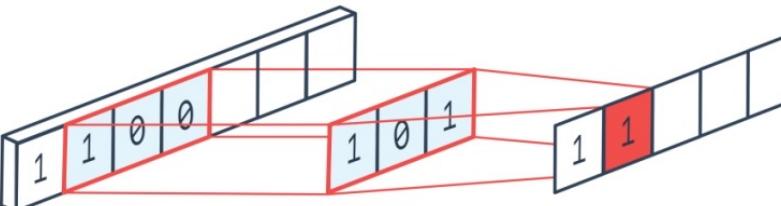
$$Y \in B^{N \times c} / Y \in R^{N \times c}$$

- RNN



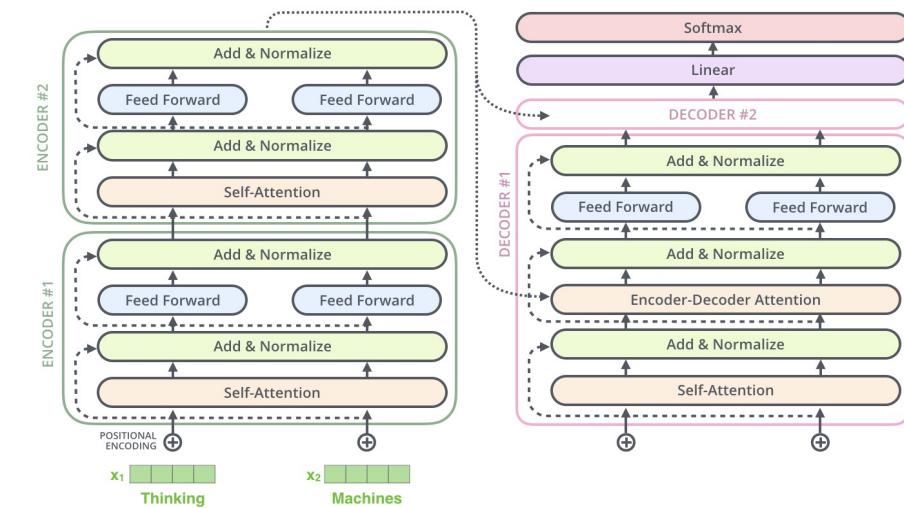
- Theoretically with infinite context window
- Practically suffer from vanishing/exploding gradient
- Training: $O(N)$, not parallelizable
- Inference: constant time for each token.

- CNN



- Finite context window (depending on kernel size)
- Need to materialize the kernel before using it.
- Training and inference depend on kernel size.
- Easily parallelizable

- Transformer



- Finite context window
- Training: $O(N^2)$
- easily parallelizable
- Inference: $O(N^2)$; may $O(N)$ when using KV-Cache, for each token.

大语言模型

- 大语言模型的核心功能和特点
 - 统一任务
 - 自回归模型(Auto-regressive), 自监督学习 (Self-supervised Learning)



01

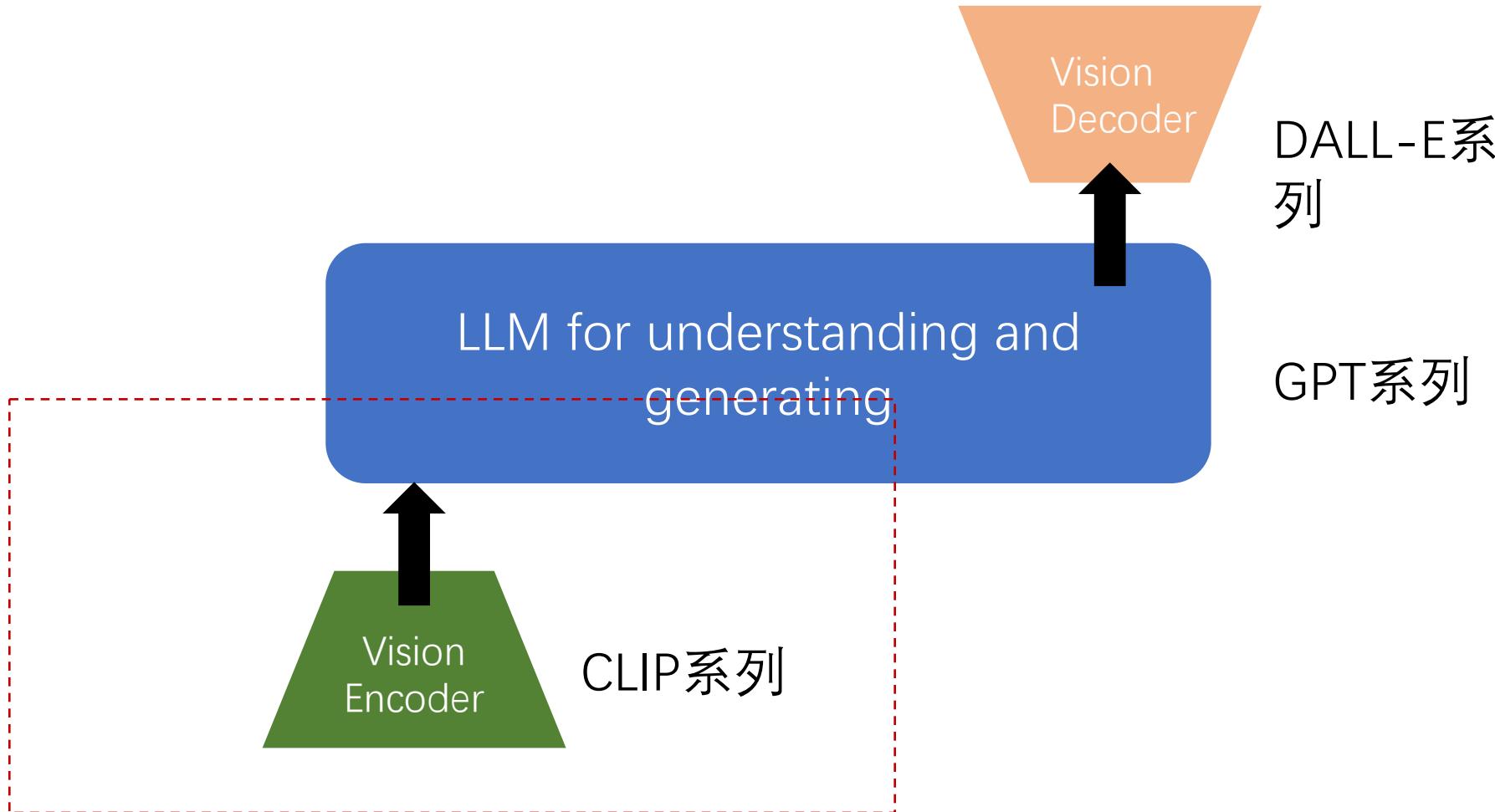
模态对齐

多模态和视觉大模型

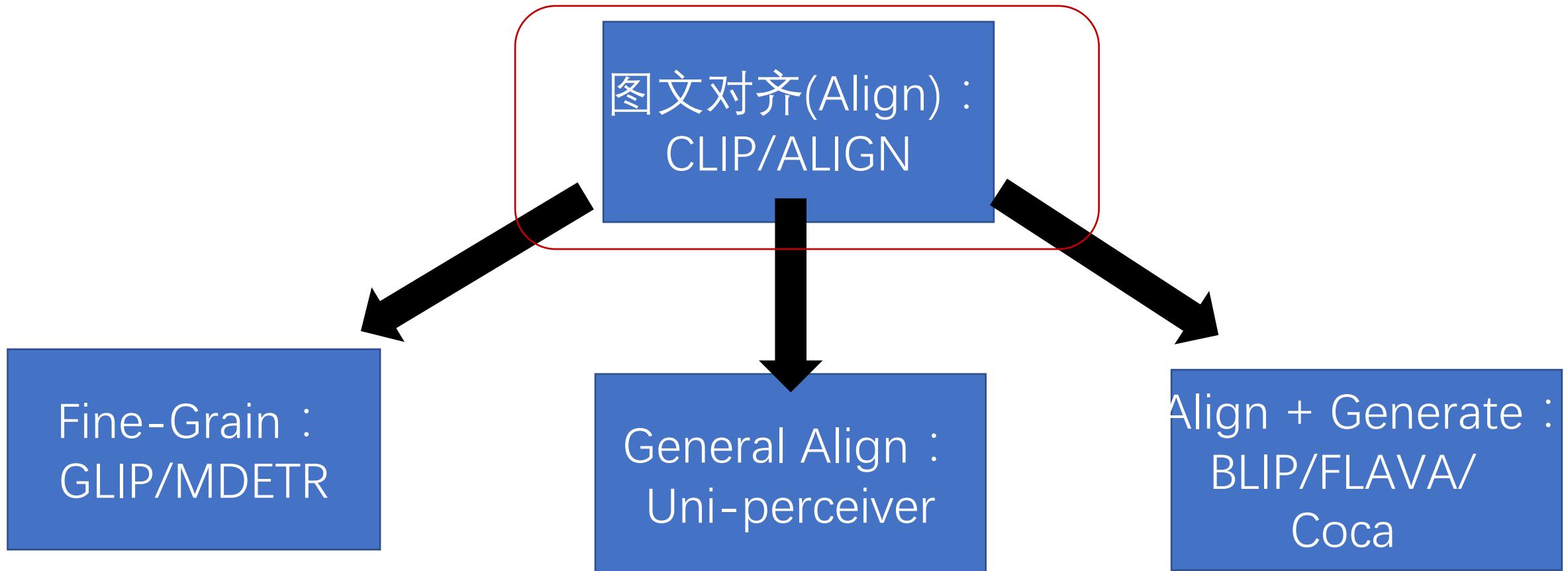
图像：

文本：

图像：



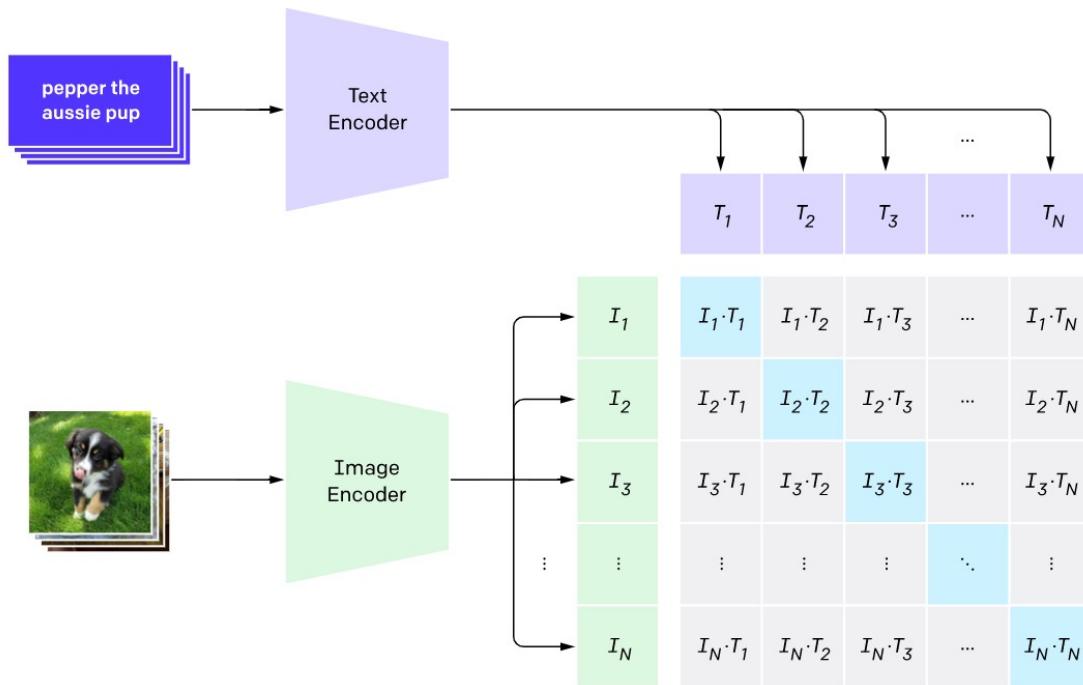
多模态大模型 (BERT-like)



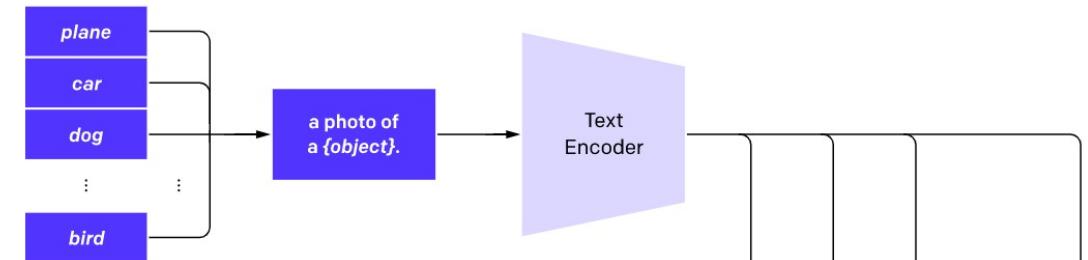
Align the text and vision

- CLIP: Contrastive Language-Image Pre-training

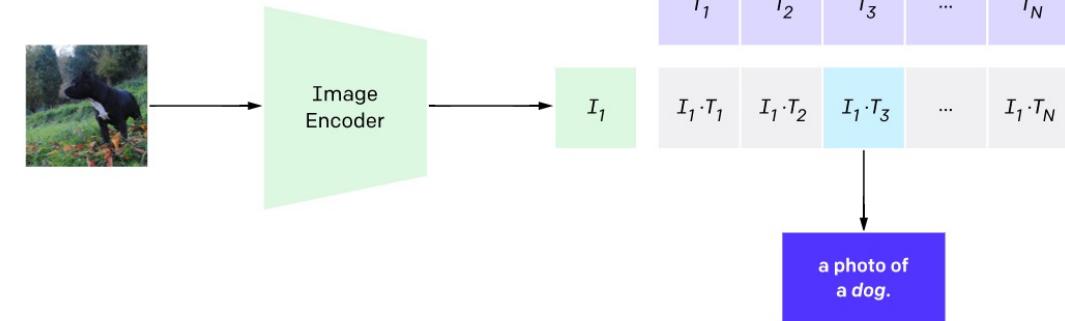
1. Contrastive pre-training



2. Create dataset classifier from label text

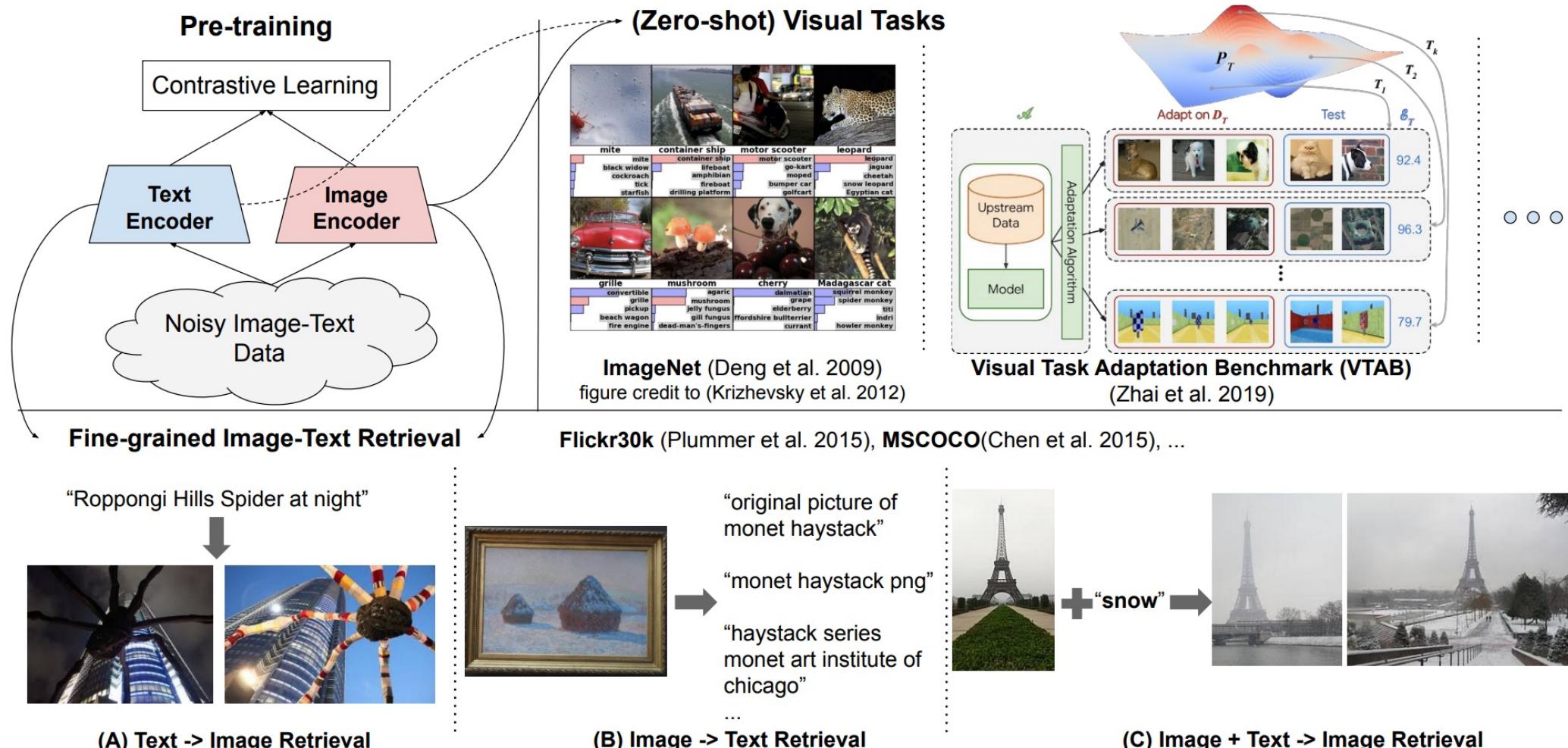


3. Use for zero-shot prediction



Align the text and vision

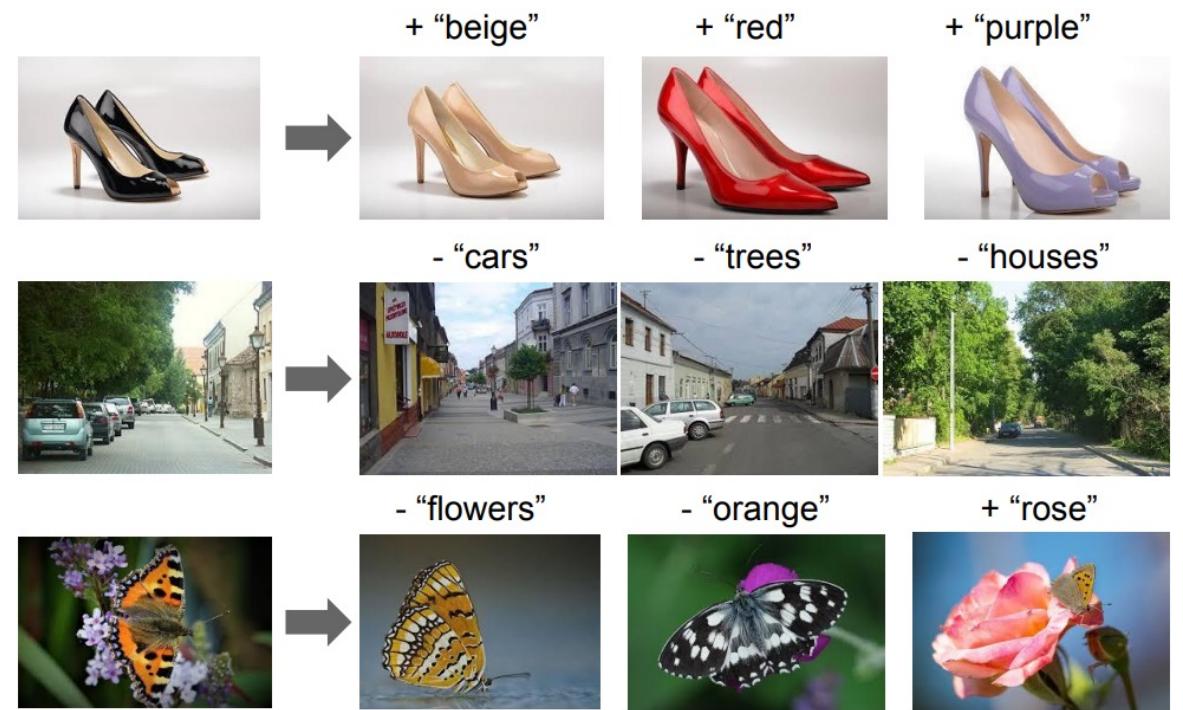
- ALIGN: A Large-scale ImaGe and Noisy-text embedding



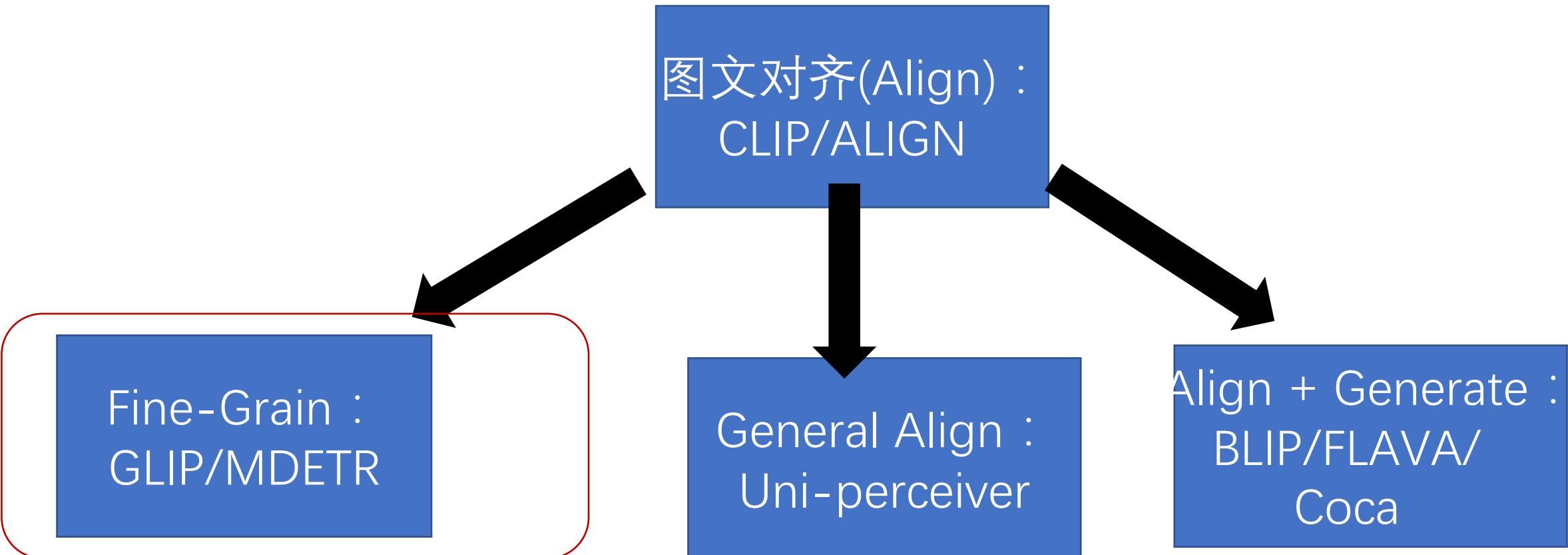
Align the text and vision

- ALIGN

Retrieval (检索):

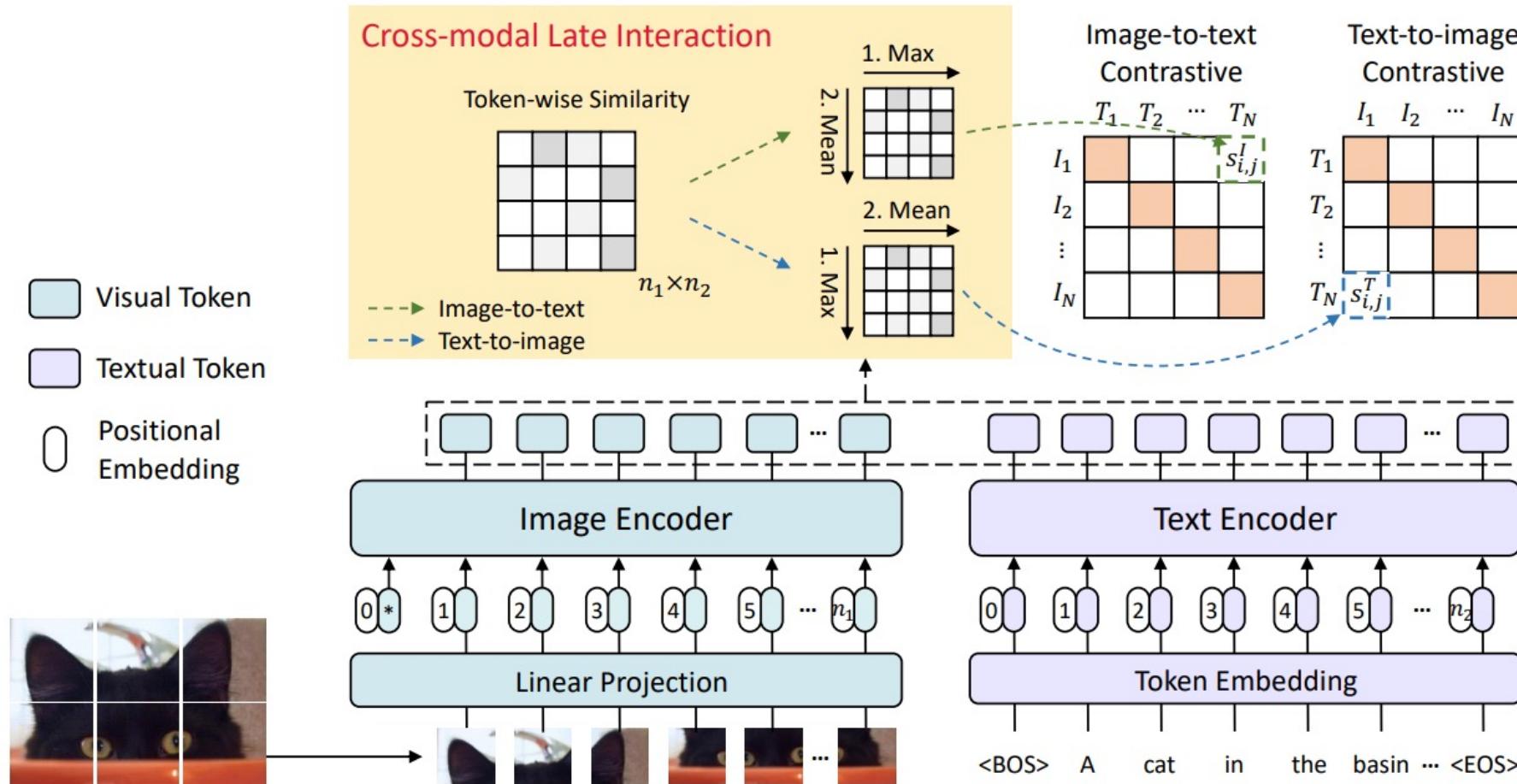


多模态大模型 (BERT-like)



Fine-grained Align the text and vision

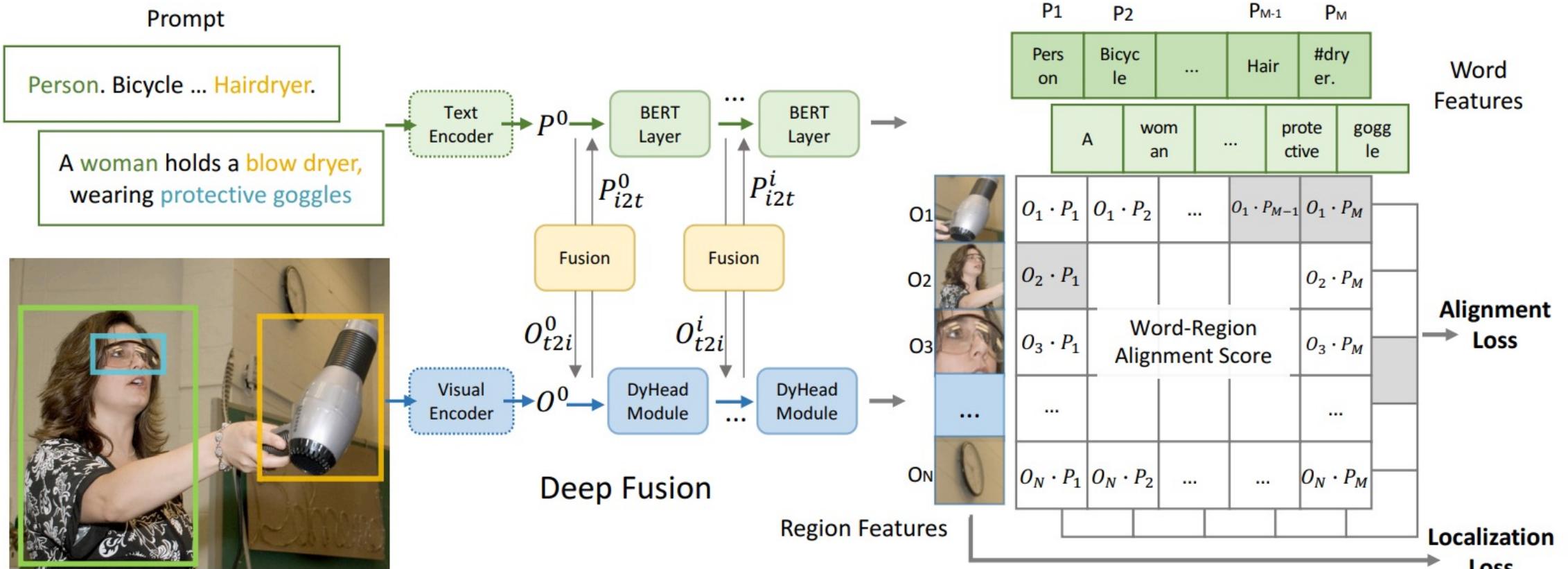
- FILIP: Fine-grained Interactive Language-Image Pre-training



Fine-grained Align the text and visual

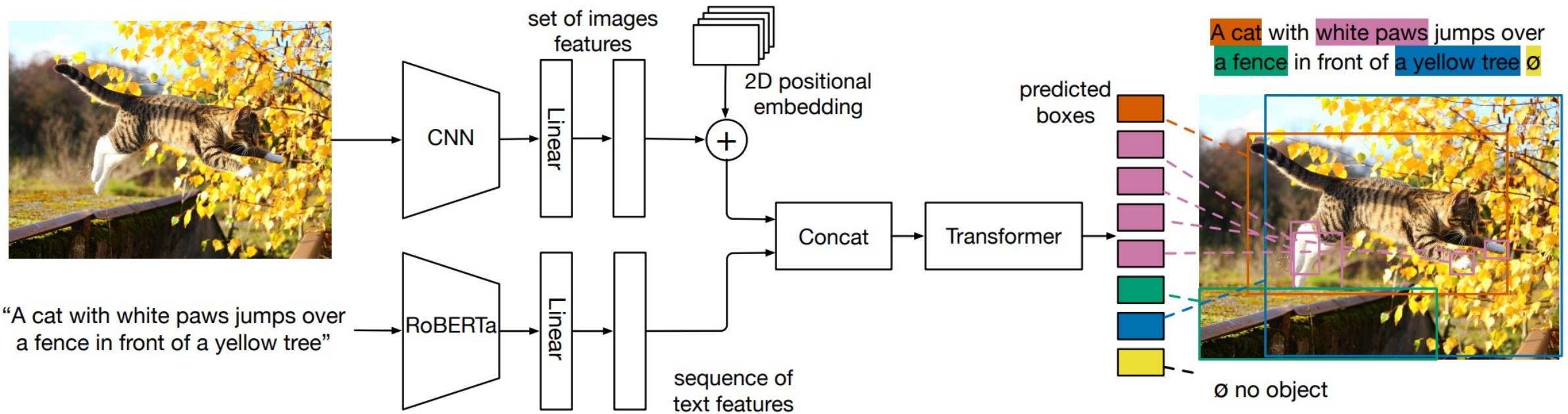
- GLIP: grounded language-image pretraining

Phase Grounding

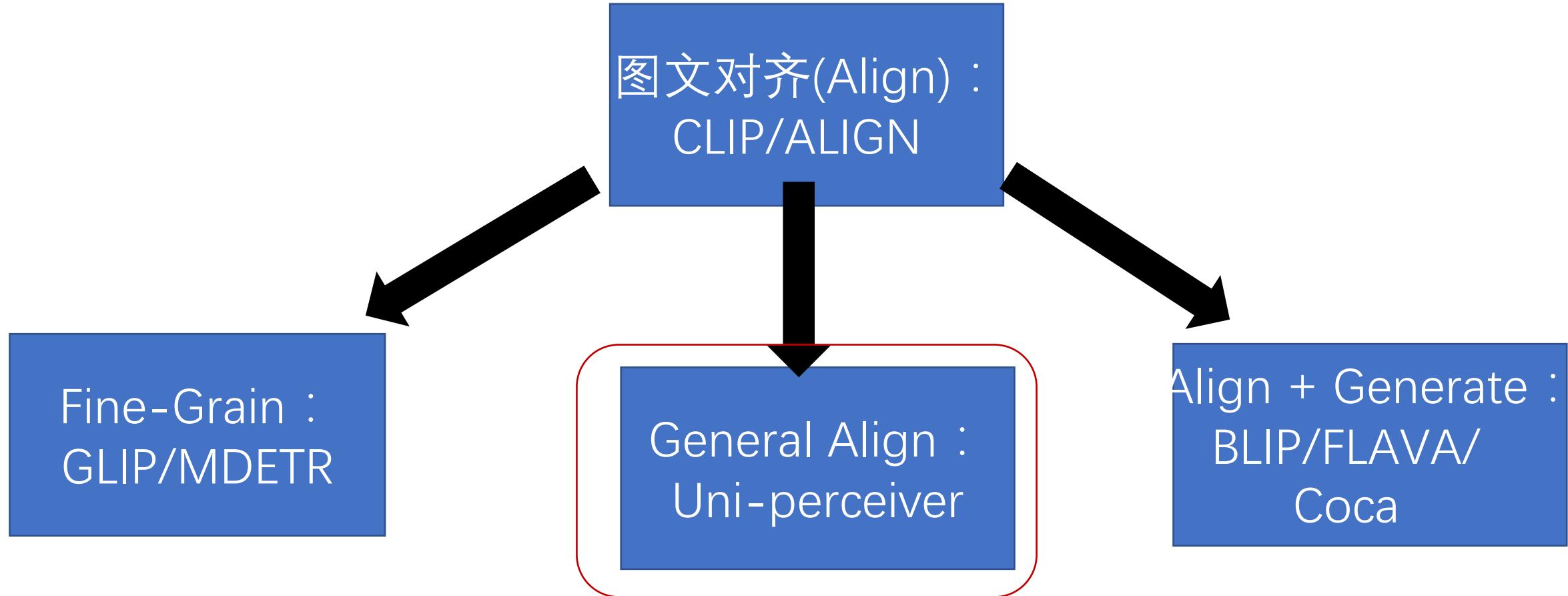


Fine-grained Align the text and visual

- MDETR: modulated detector based on DETR

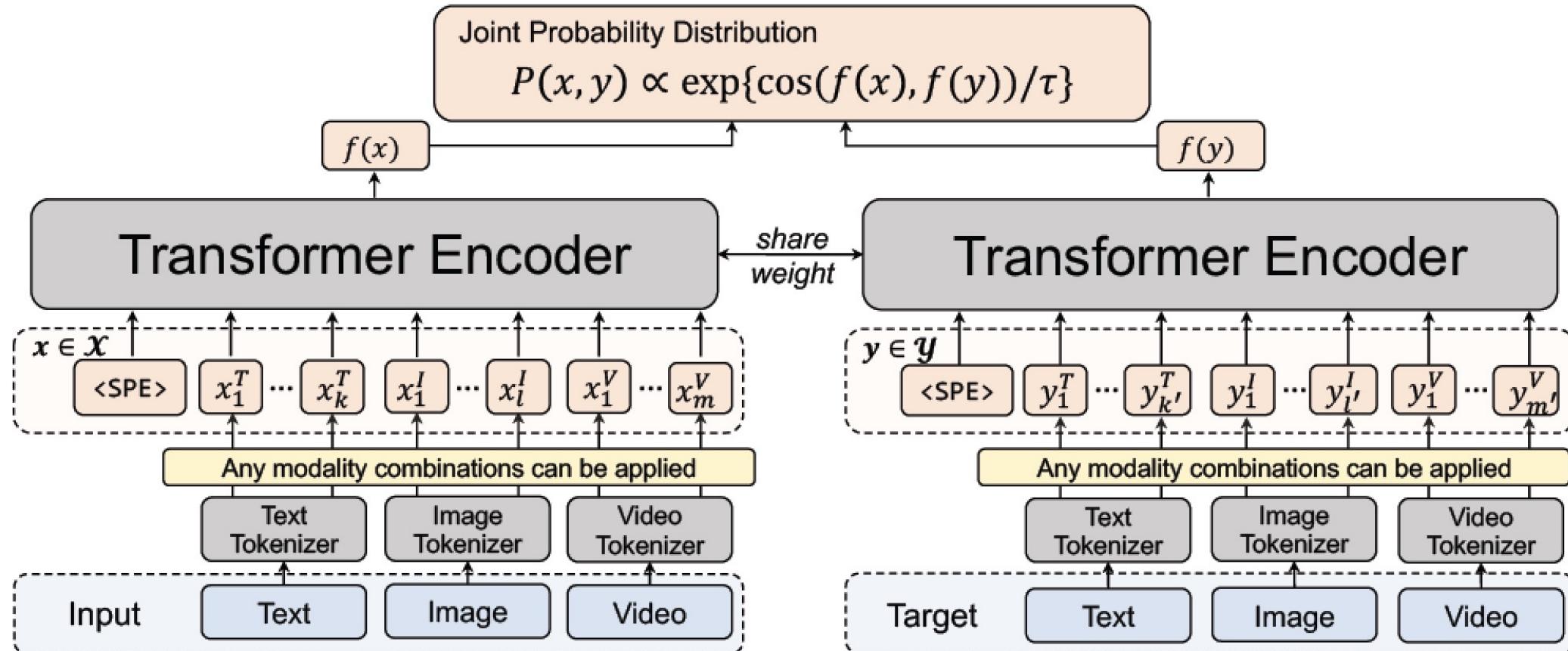


多模态大模型 (BERT-like)



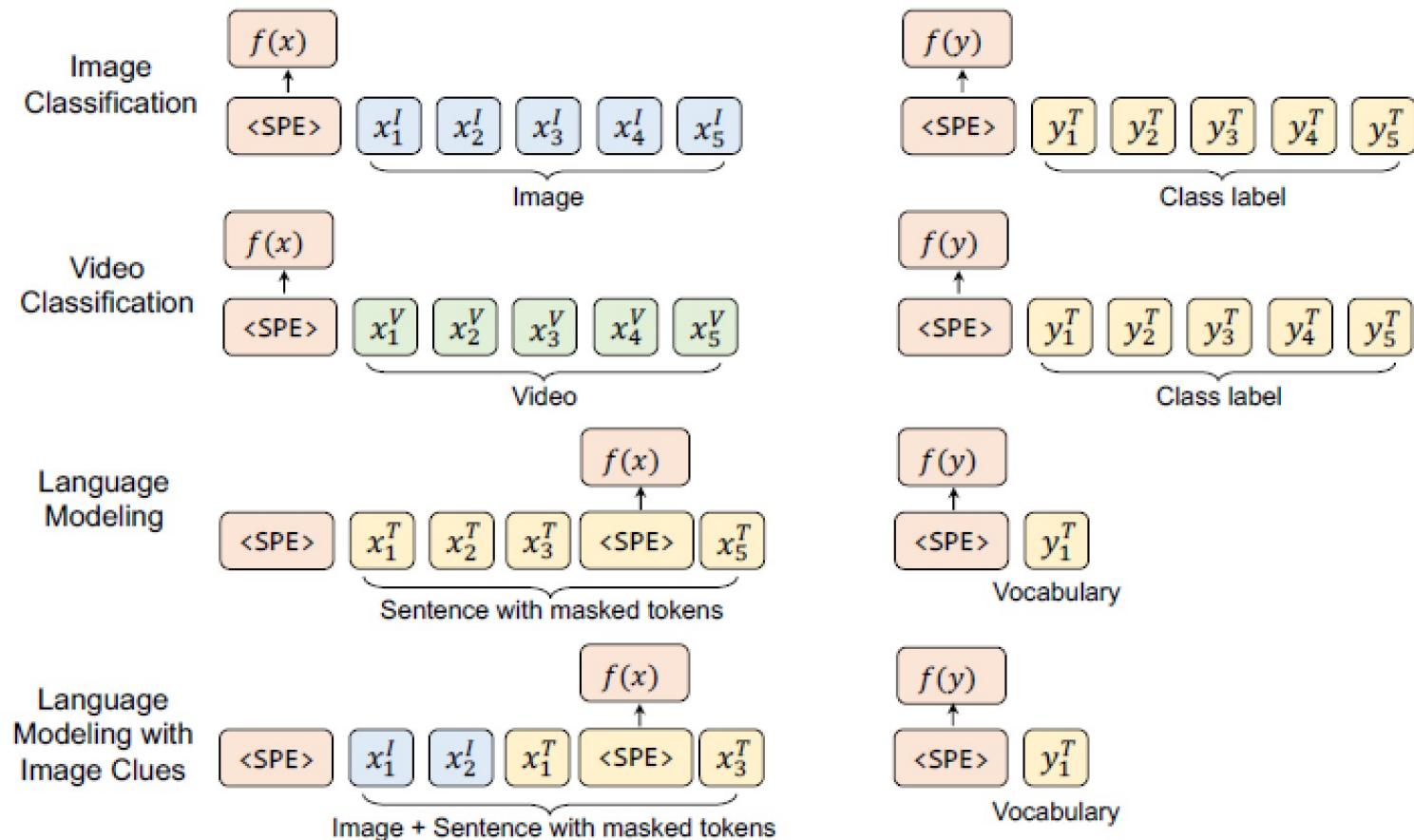
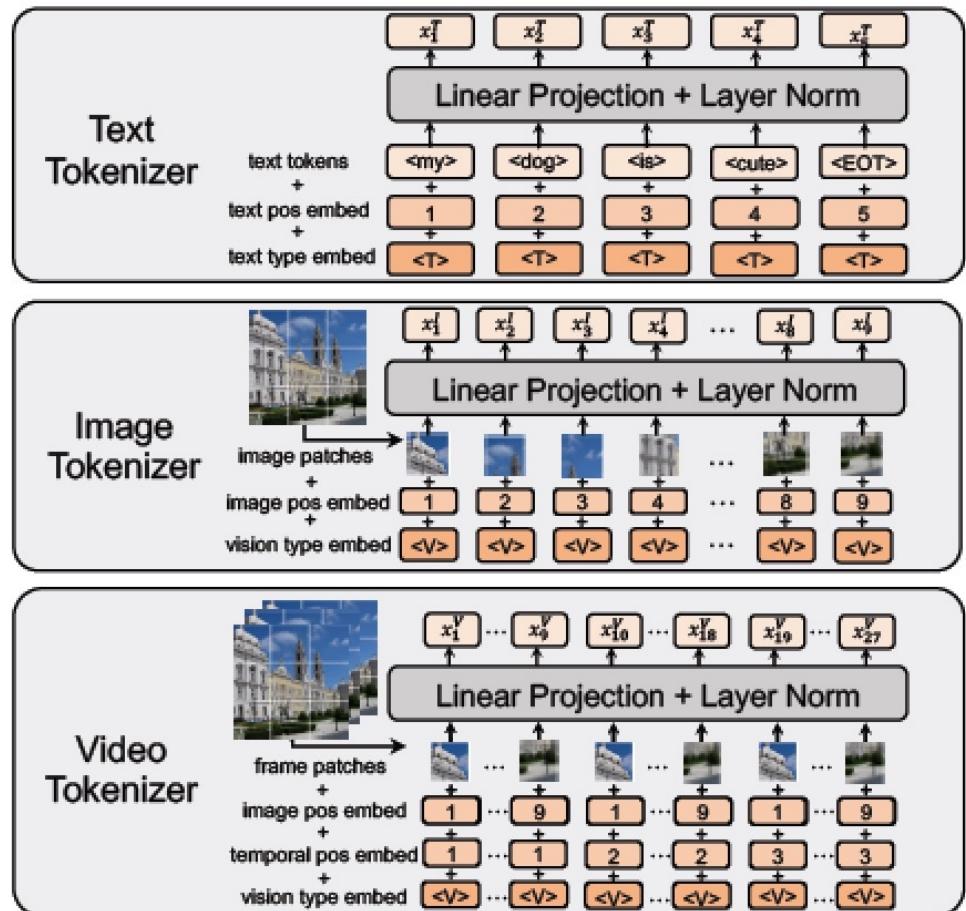
General Align

- Uni-Perceiver
 - No specific design for tasks: focus on the input and output process

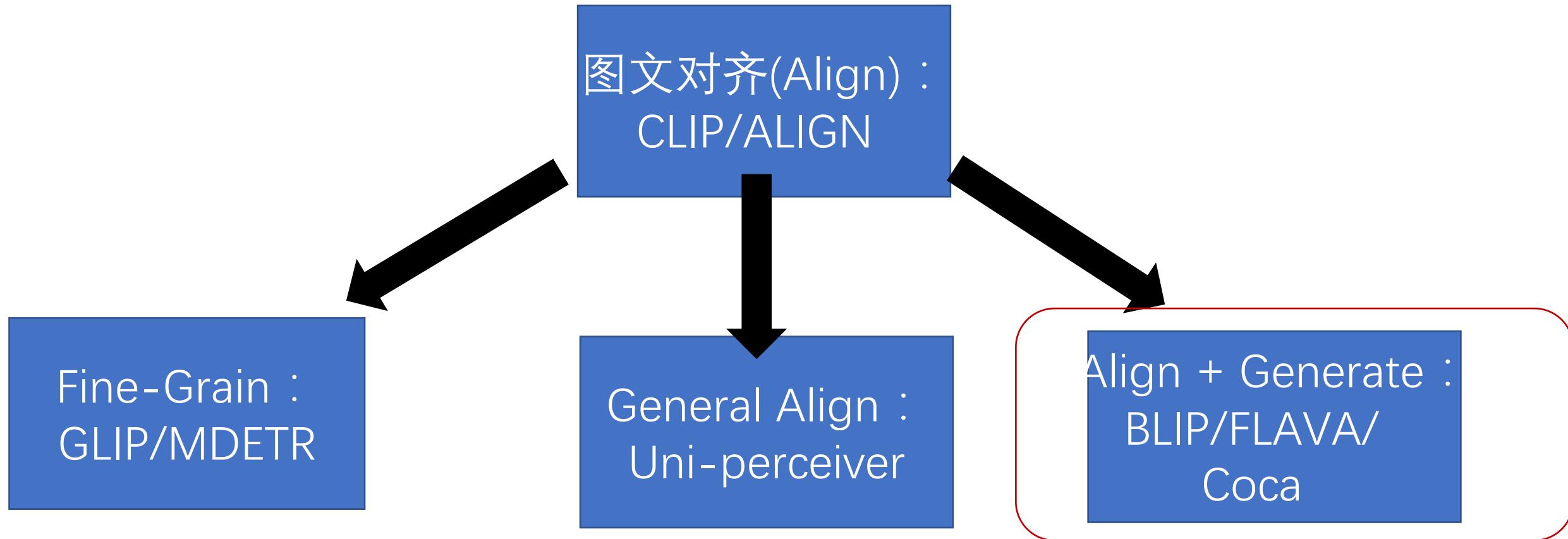


General Align

- Uni-Perceiver

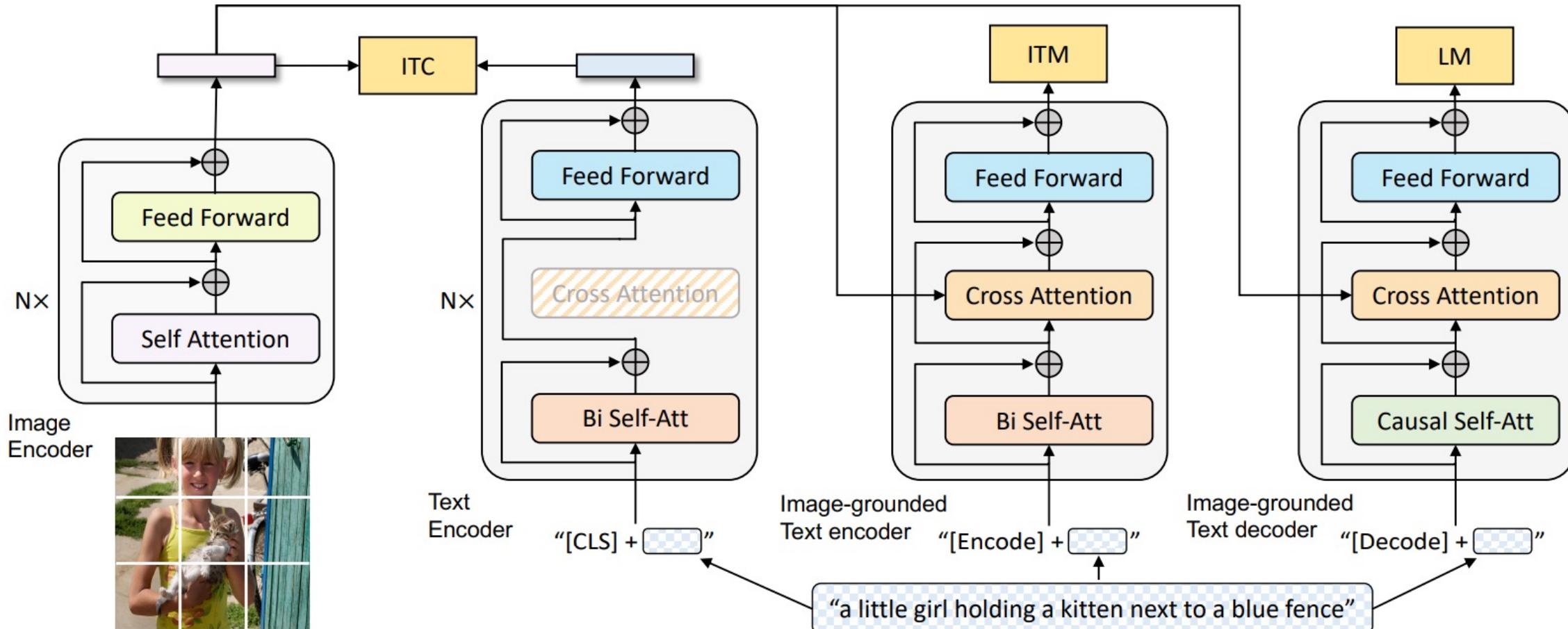


多模态大模型 (BERT-like)



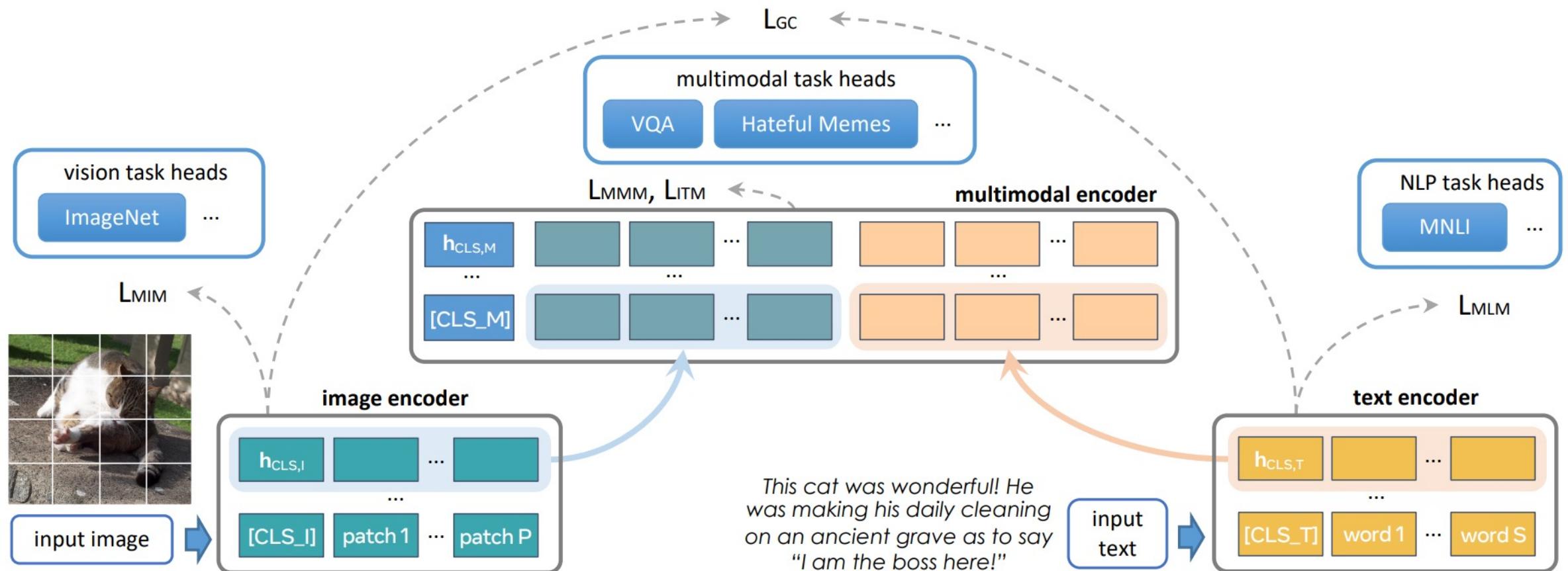
Align + Generate

- BLIP: Bootstrapping Language-Image Pre-training



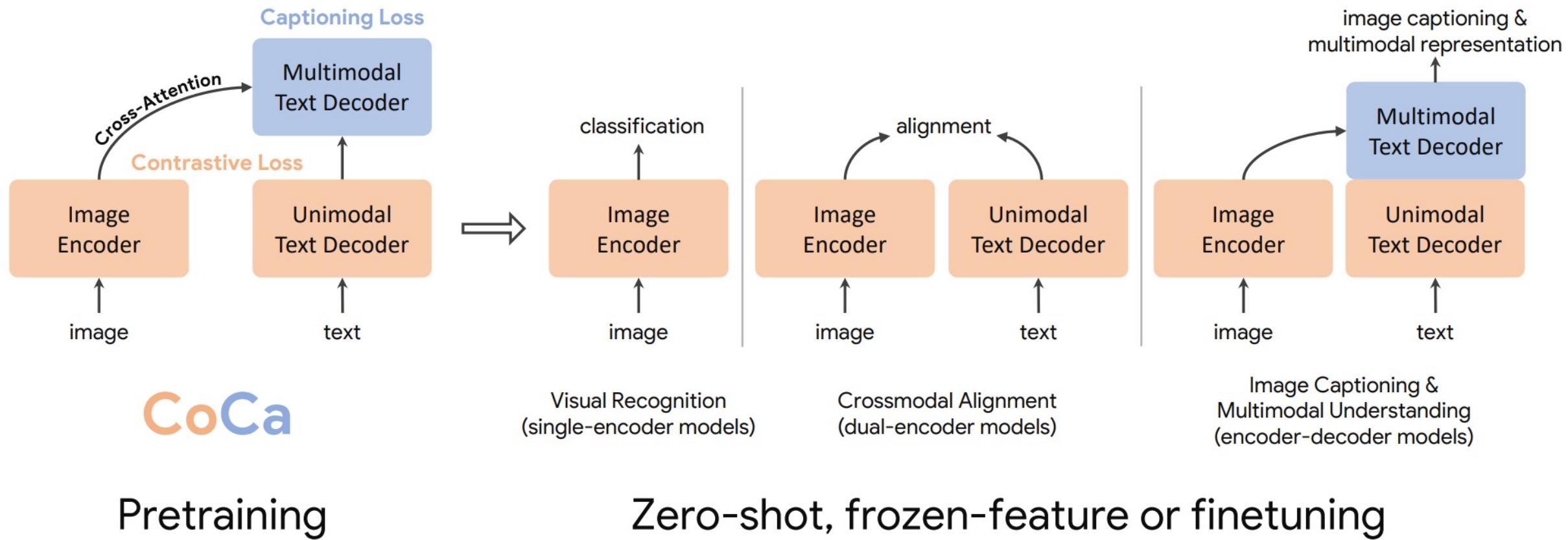
Align + Generate

- FLAVA: Foundational LAnguage and Vision Alignment model



Align + Generate

- Coca: Contrastive Captioner



谢谢！

