

# **Large Language Models**

Multi-modal Foundation Models: Vision-Language Models (VLMs)

M. Soleymani

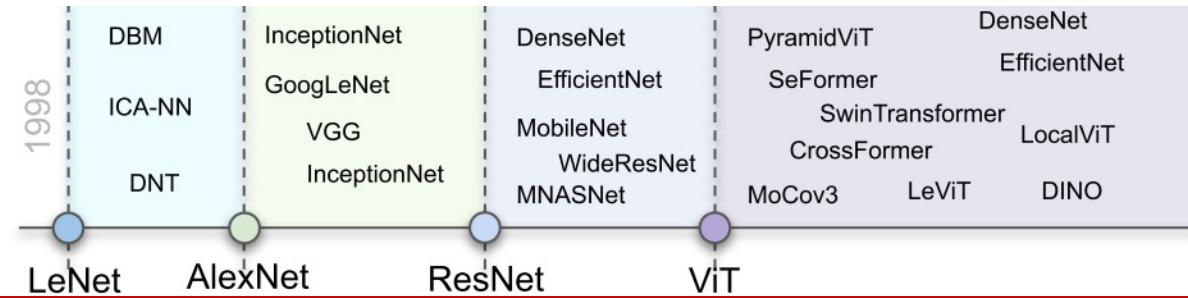
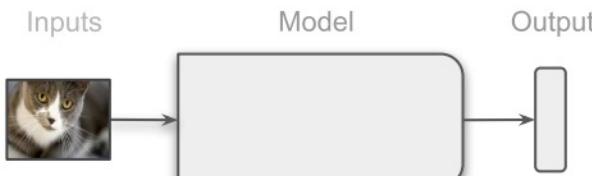
Sharif University of Technology

Fall 2023

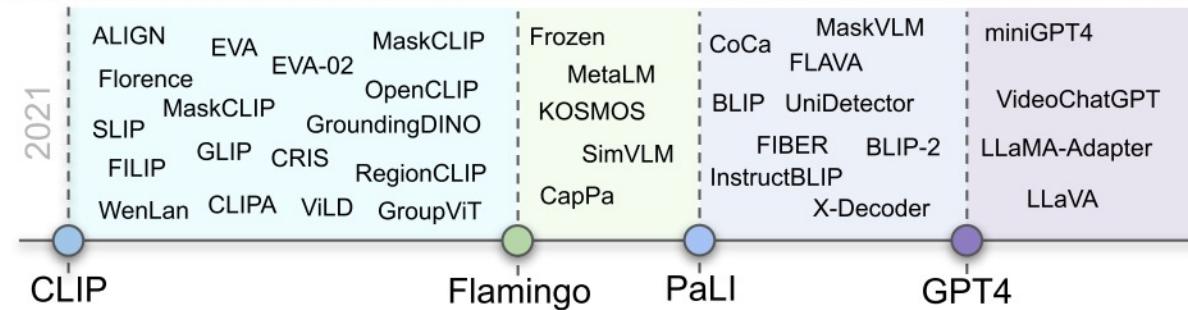
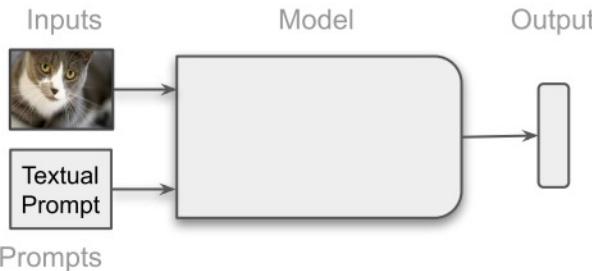
# Multi-modal data

- Multimodal data:
  - Input and output from different modalities (e.g. text-to-image, image-to-text)
  - Inputs are multimodal (e.g. a system that can process both text and images)
  - Outputs are multimodal (e.g. a system that can generate both text and images)

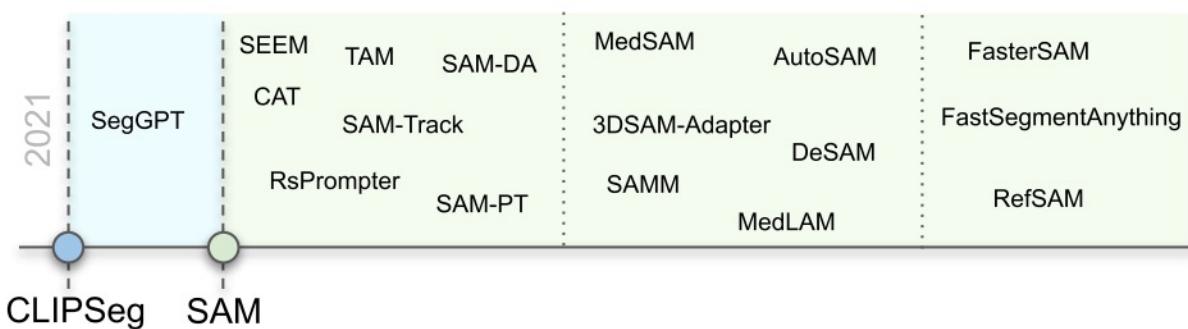
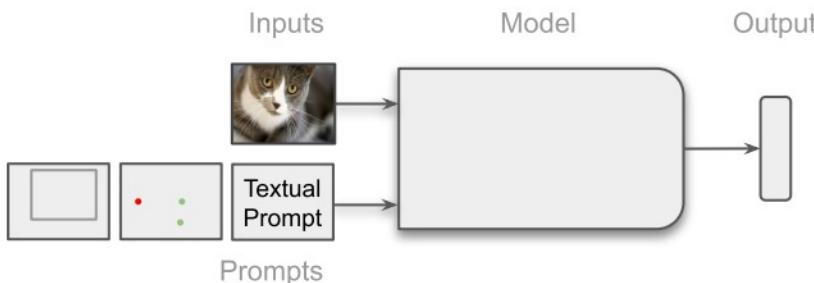
## Traditional Models



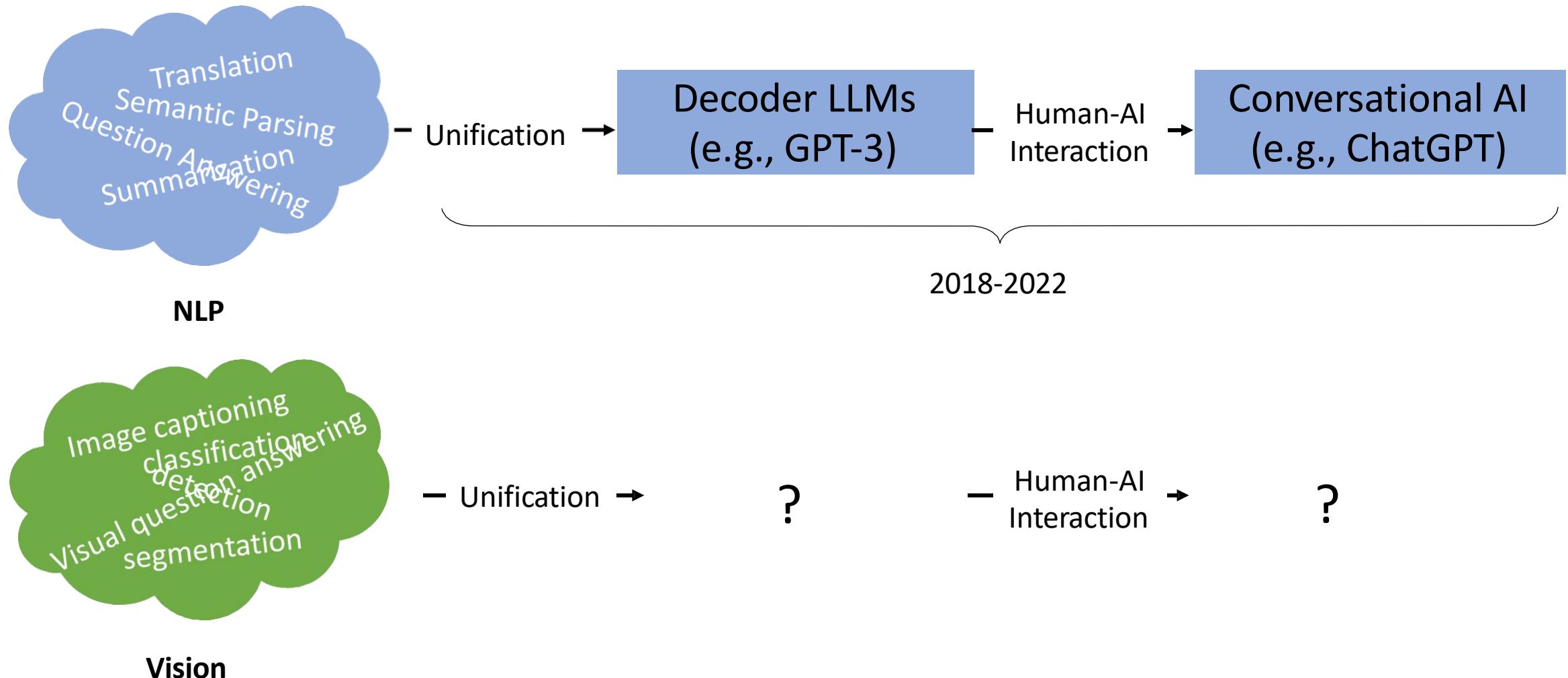
## Textually Prompted Models



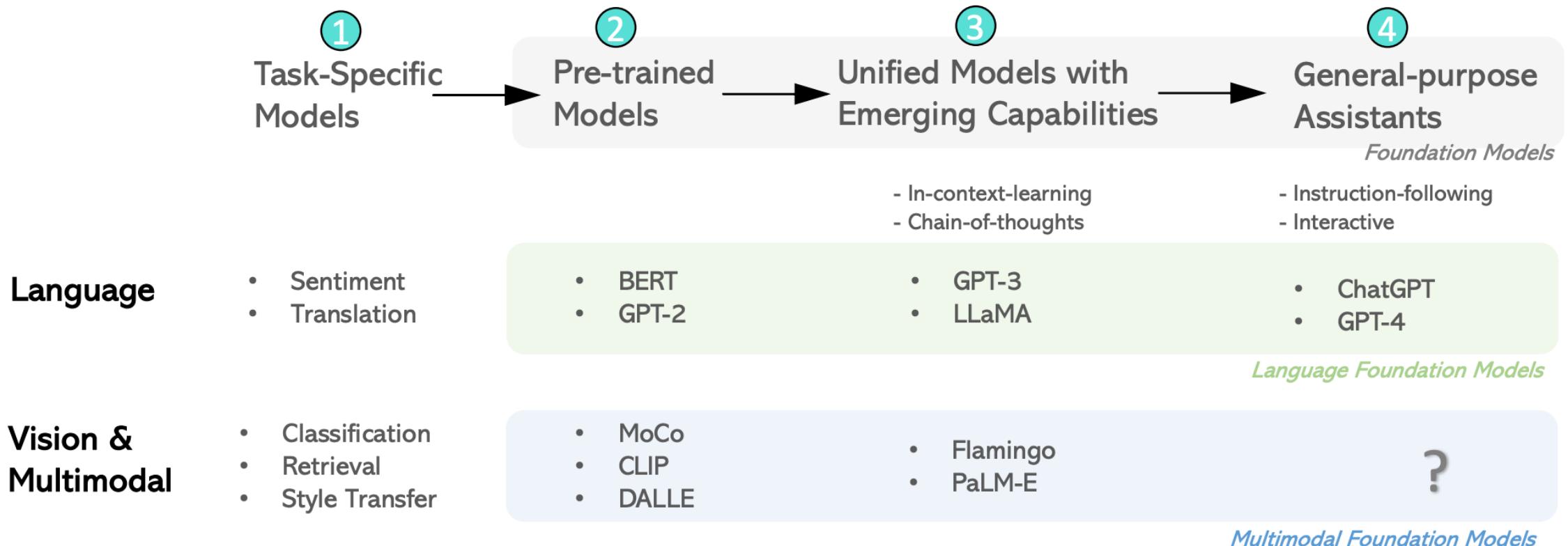
## Visually Prompted Models



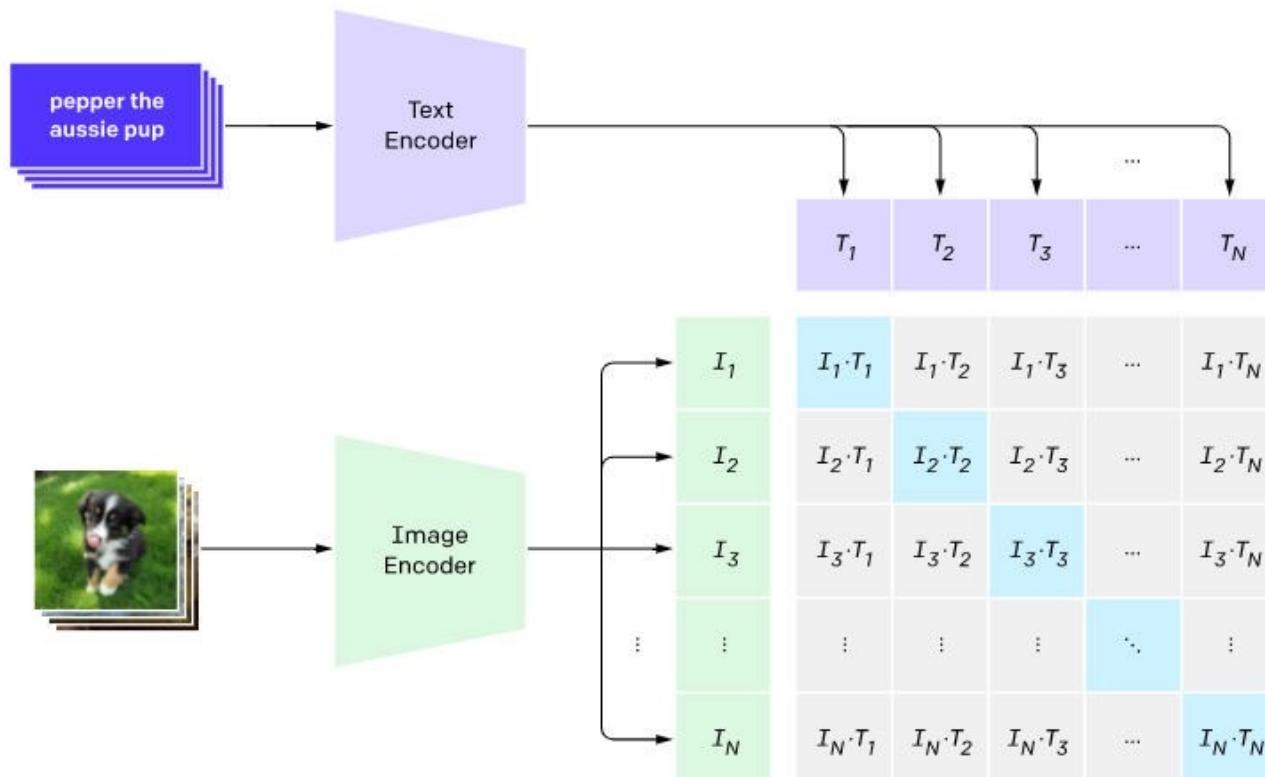
# A Lesson from LLMs



# A Lesson from LLMs



# CLIP: Models and Training Complexity

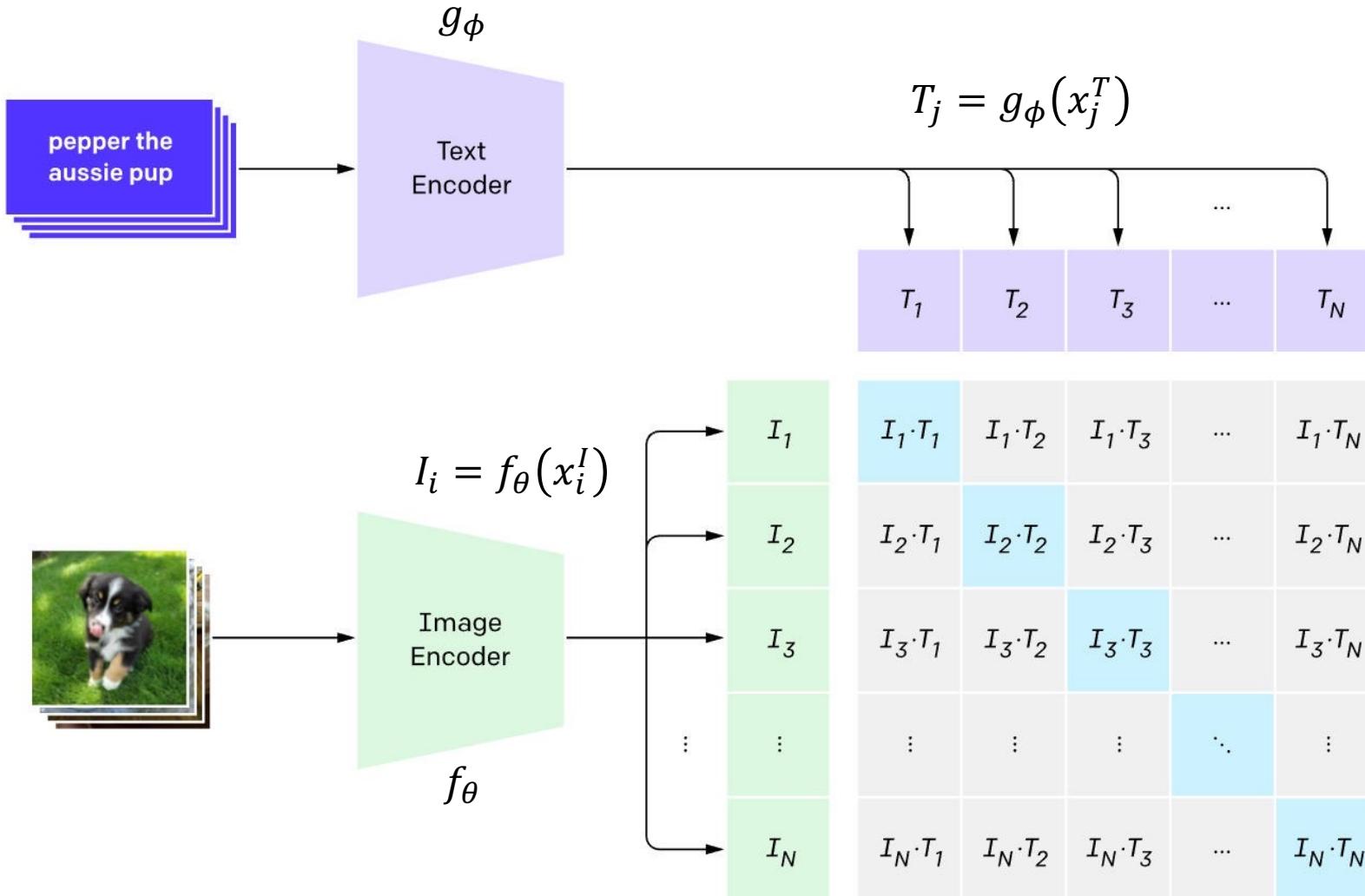


- Text encoder:
  - 12-layer Transformer with causal mask
- Image encoder:
  - ResNet families: RN50, RN101, RN50x4, RN50x16, RN50x64
  - ViT families: ViT-B/32, ViT-B/16, ViT-L/14

# Vision-language models: Contrastive learning

- Contrastive training to bridge the image and text embedding spaces
- Making embedding of (image, text) pairs similar and that of non-pairs dissimilar
- This embedding space is super helpful for performing searches across modalities
  - Can return the best caption given an image
  - Has impressive capabilities for zero-shot adaptation to unseen tasks, without the need for fine-tuning

## 1. Contrastive pre-training



$$s_{i,j}^T = s_{i,j}^I = I_i^T T_j$$

$$\mathcal{L}_i^I = -\log \frac{e^{s_{i,i}^I}}{\sum_{j=1}^N e^{s_{i,j}^I}}$$

$$\mathcal{L}_j^T = -\log \frac{e^{s_{i,i}^T}}{\sum_{i=1}^N e^{s_{i,j}^T}}$$

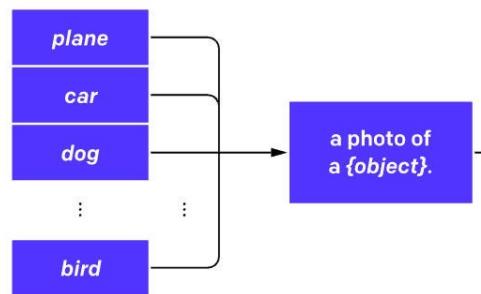
$$\mathcal{L} = \frac{1}{2N} \sum_{i=1}^N (\mathcal{L}_i^I + \mathcal{L}_j^T)$$

- Training batchsize: 32,768

- Training time:
  - RN50x64: 18 days on 592 V100 GPUs
  - ViT-L/14: 12 days on 256 V100 GPUs

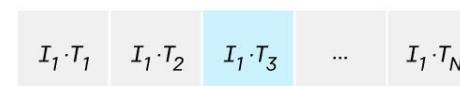
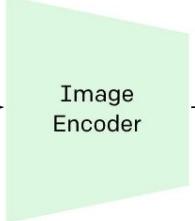
# CLIP for zero-shot learning

## 2. Create dataset classifier from label text



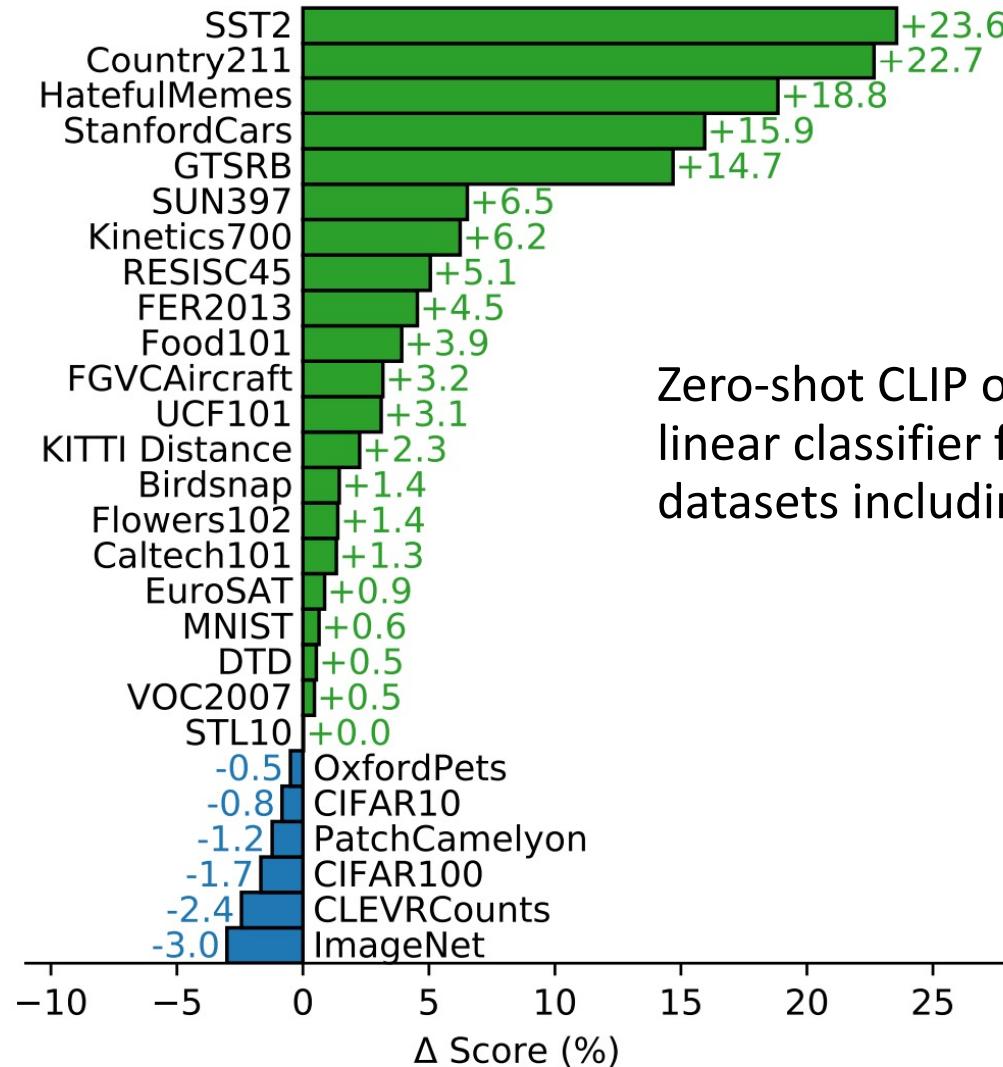
encodes all the text labels and compares them to the encoded image

## 3. Use for zero-shot prediction

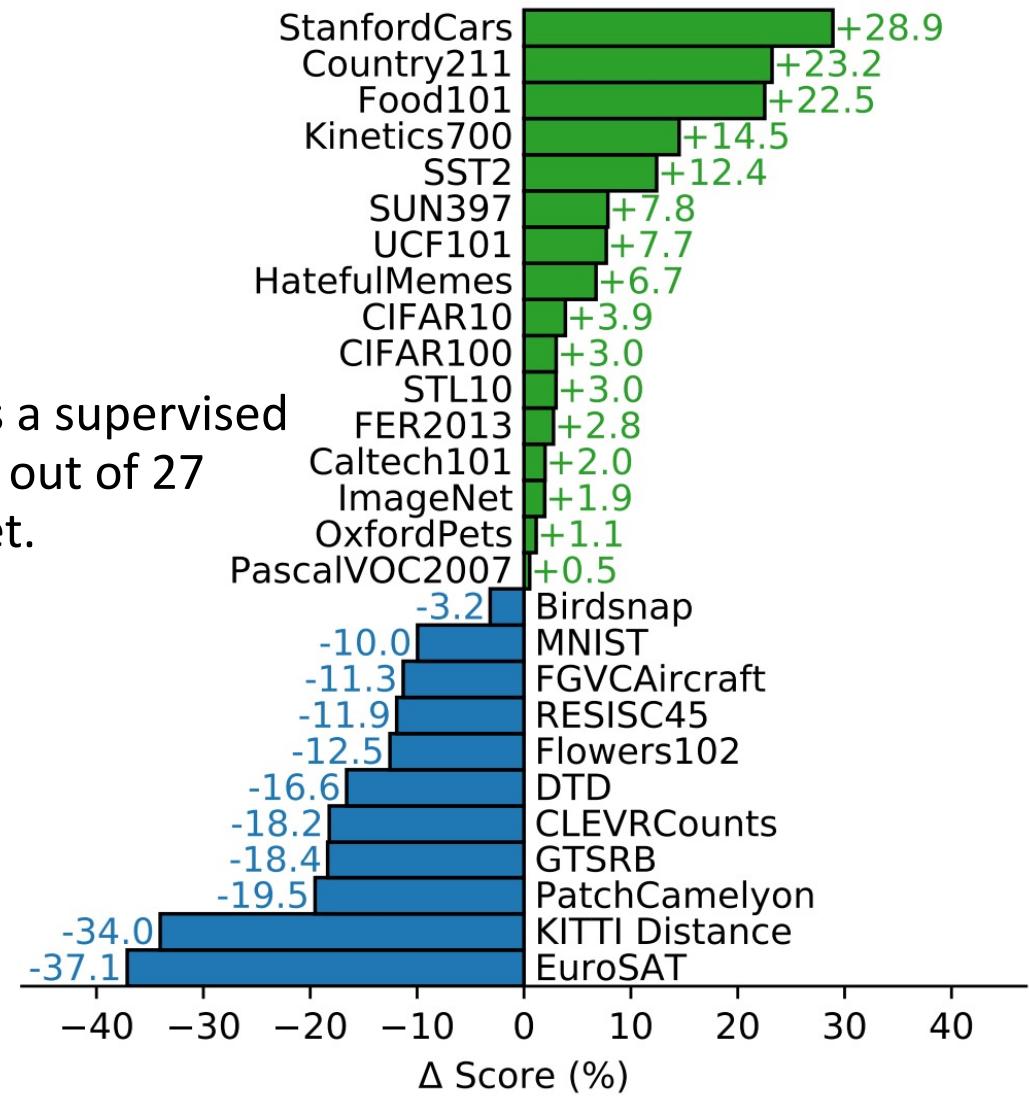


a photo of  
a dog.

<https://openai.com/blog/clip/>



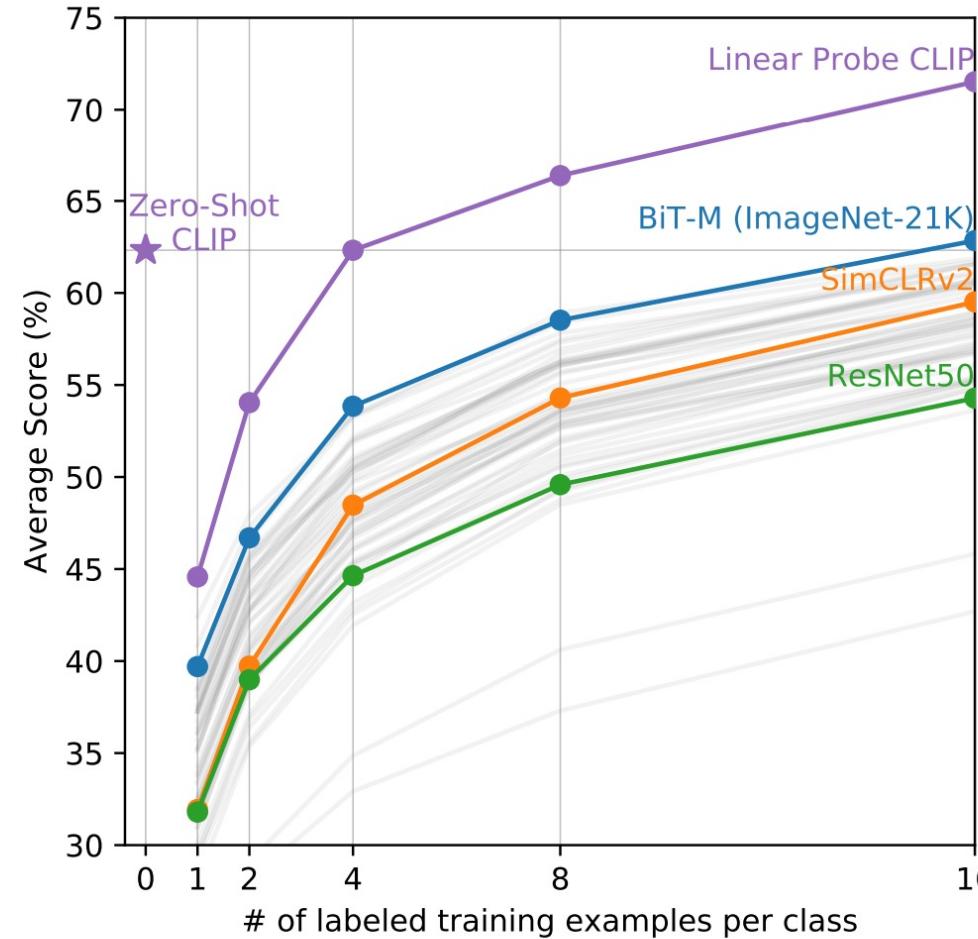
Zero-shot CLIP outperforms a supervised linear classifier fitted on 16 out of 27 datasets including ImageNet.



CLIP's features outperform the features of the best ImageNet model on a wide variety of datasets.

Radford et al. "Learning Transferable Visual Models From Natural Language Supervision", ICML 2021

# Zero-shot CLIP outperforms few-shot linear probes



# Vision Language Tasks

Large Multi-modal Models (LMMs) in their current form is primarily generates a text sequence.

	Image Captioning	Text-to Image Retrieval	Image-to-Text Retrieval	VQA	Text-to-Image Generation
Input	Image: 	Query: A couple of zebra walking across a dirt road.  A pool of images	Query:   A pool of texts	Image:   Q: why did the zebra cross the road?	Text: A couple of zebra walking across a dirt road.
Output	A couple of zebra walking across a dirt road.		A couple of zebra walking across a dirt road.	A: to get to the other side  (Selected from a pool of 3,129 answers in VQAv2 or generate answer)	
	Generation	Understanding	Understanding	Understanding/Generation	Generation

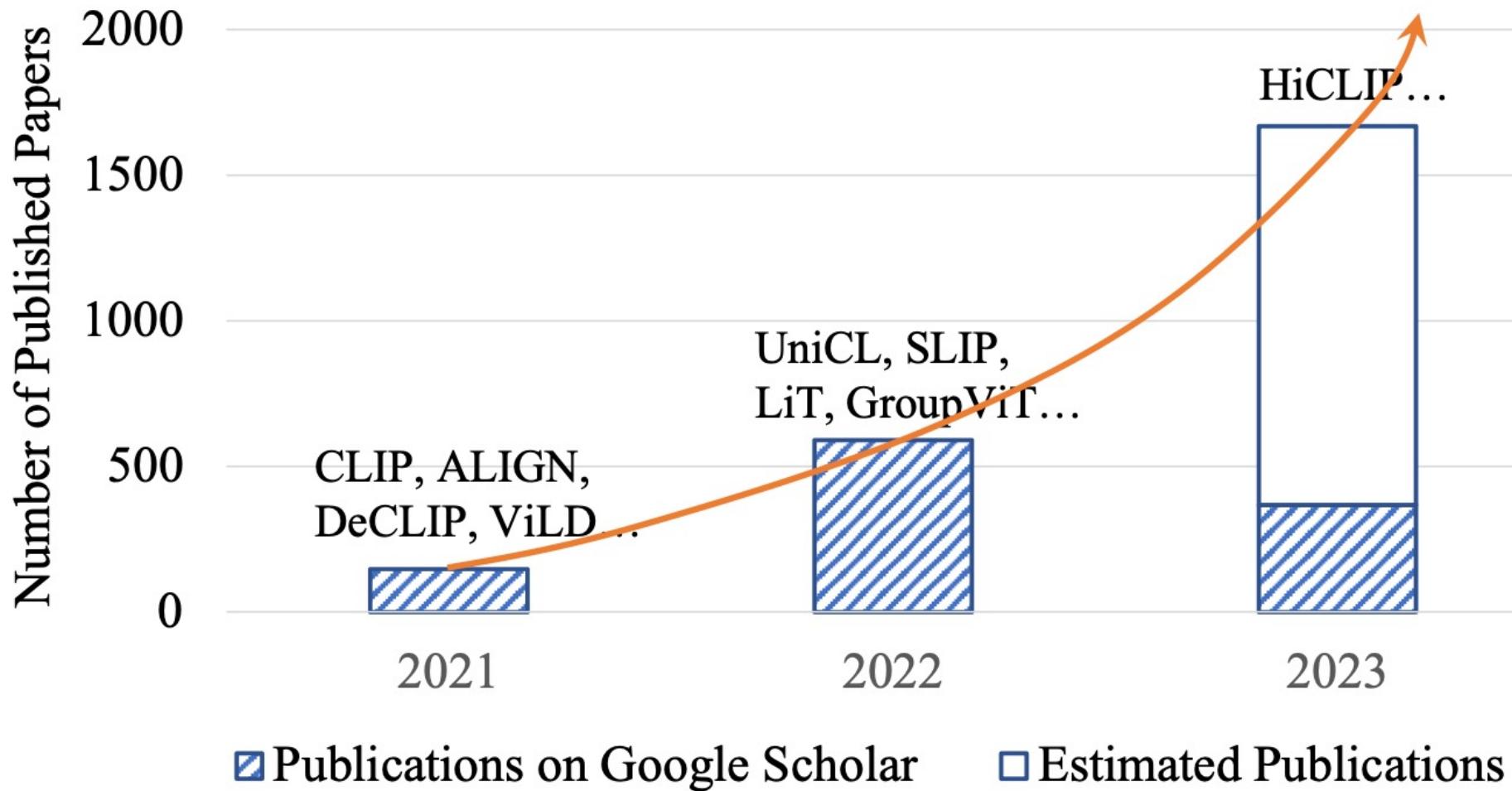
# CLIP: Summary

- ✓ CLIP improved open-vocabulary visual recognition capabilities through learning from Internet-scale image-text pairs.
- ✗ CLIP doesn't go directly from image to text or vice versa. It just connects the image and text embedding spaces
  - CLIP can only address limited use cases such as classification
  - It crucially lack the ability to generate language which makes them less suitable to more open-ended tasks such as captioning or visual question answering

# Survey of VLMs



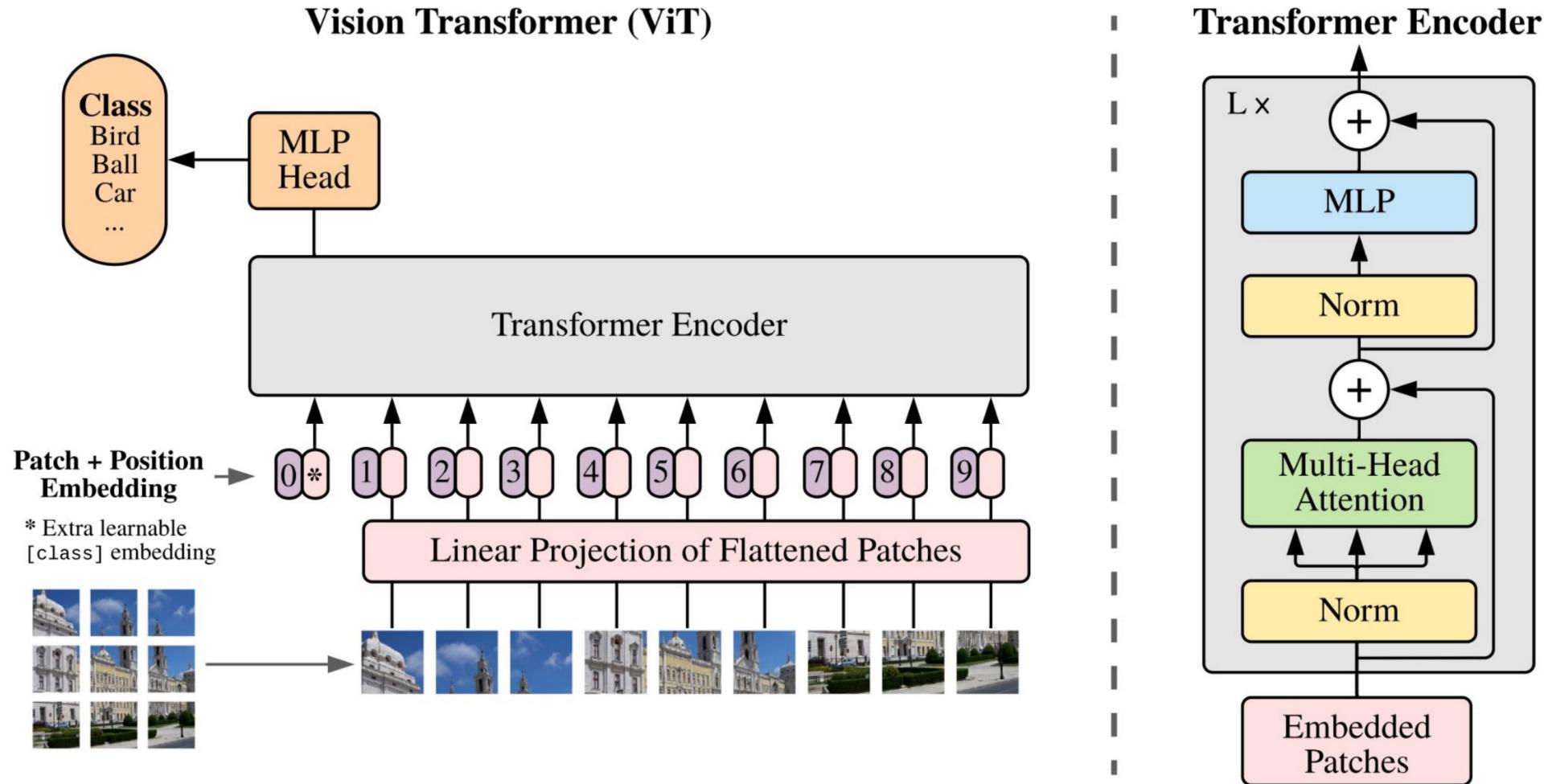
# Publication on VLMs

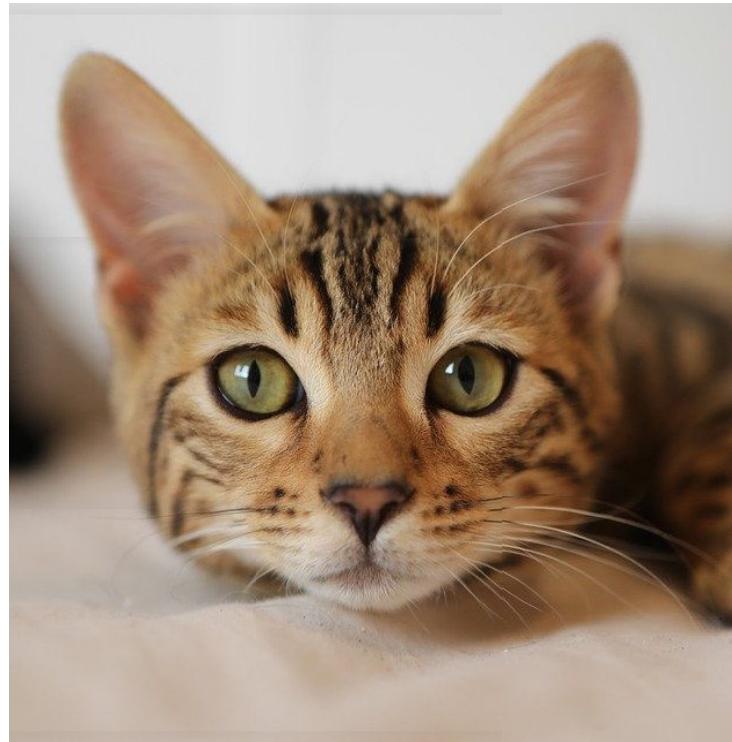


# CLIP Variants

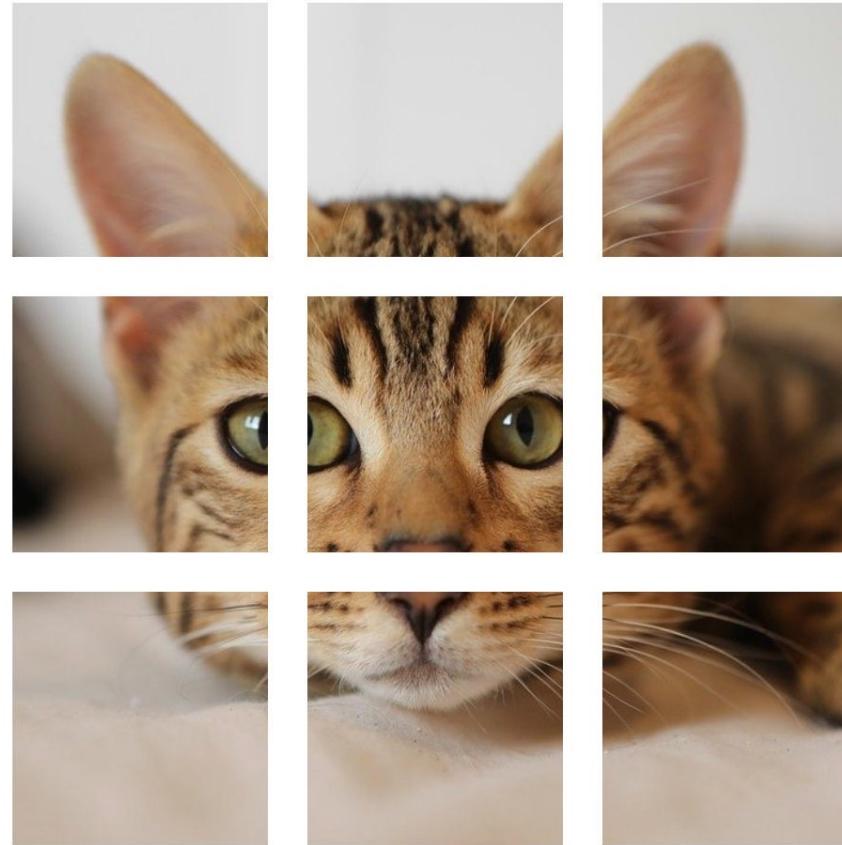
- Objective function or pretraining
  - Combining CLIP with label supervision (BASIC, UniCL, LiT, MOFI)
  - Contrastive + self-supervised image representation learning
    - Contrastive + Self-supervised methods like SimCLR (SLIP, DeCLIP, nCLIP)
    - Contrastive + Masked Image Modeling (EVA, EVA-02, MVP)
  - Fine-grained matching loss (FILIP)
  - Contrastive + captioning loss (CoCa)
  - Region-level pretraining (RegionCLIP, GLIP)
  - Sigmoid loss for language-image pre-training (SigCLIP)
- Architecture
  - Dual encoder
  - Multi-modal fusion
  - Encoder-decoder

# Vision Transformer as Image Encoder Architecture

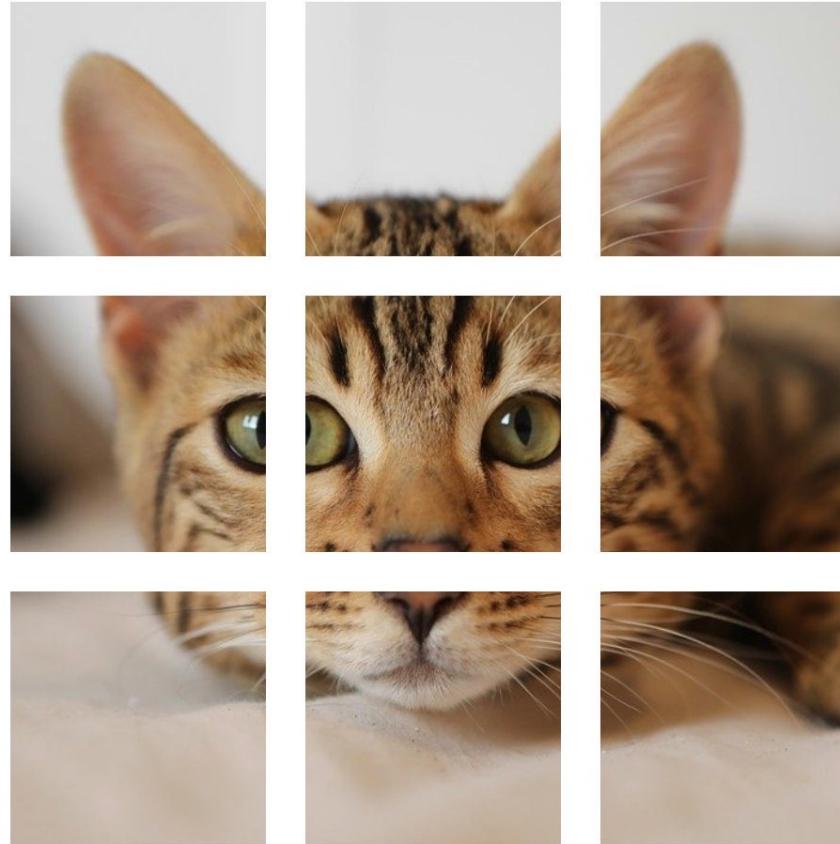




*Slides from Justin Johnson*



*Slides from Justin Johnson*

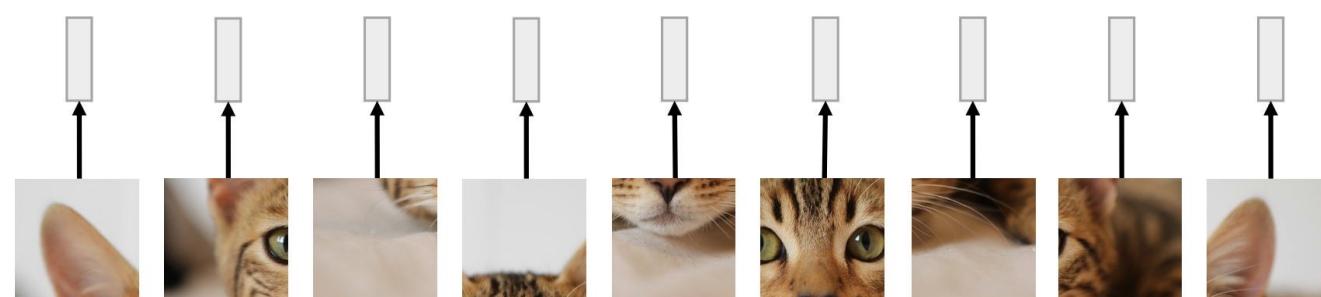


N input patches, each  
of shape 3x16x16



Linear projection to  
D-dimensional vector

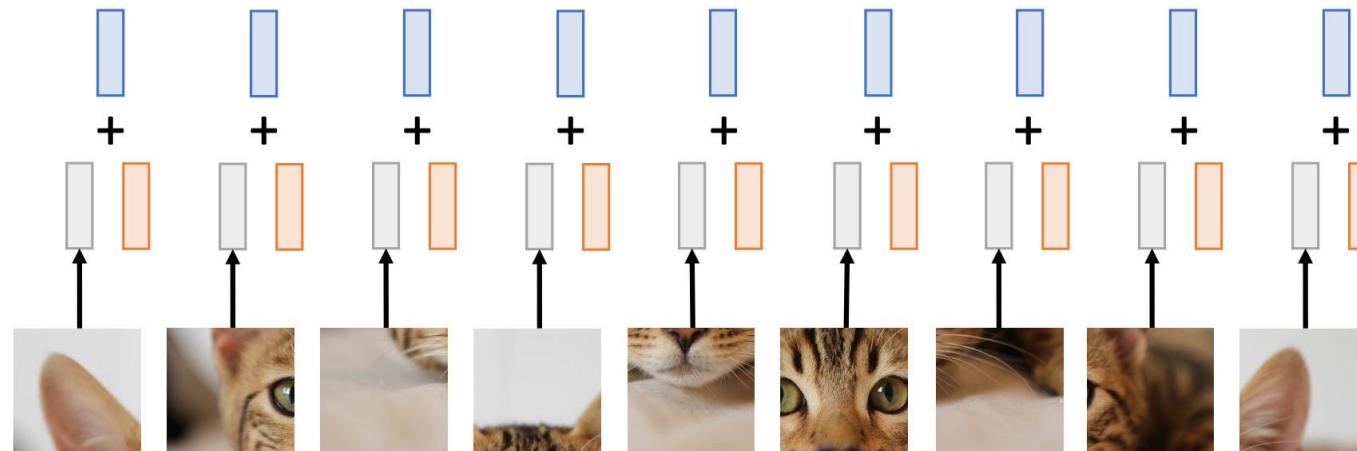
N input patches, each  
of shape 3x16x16



Add positional  
embedding: learned D-  
dim vector per position

Linear projection to  
D-dimensional vector

N input patches, each  
of shape 3x16x16



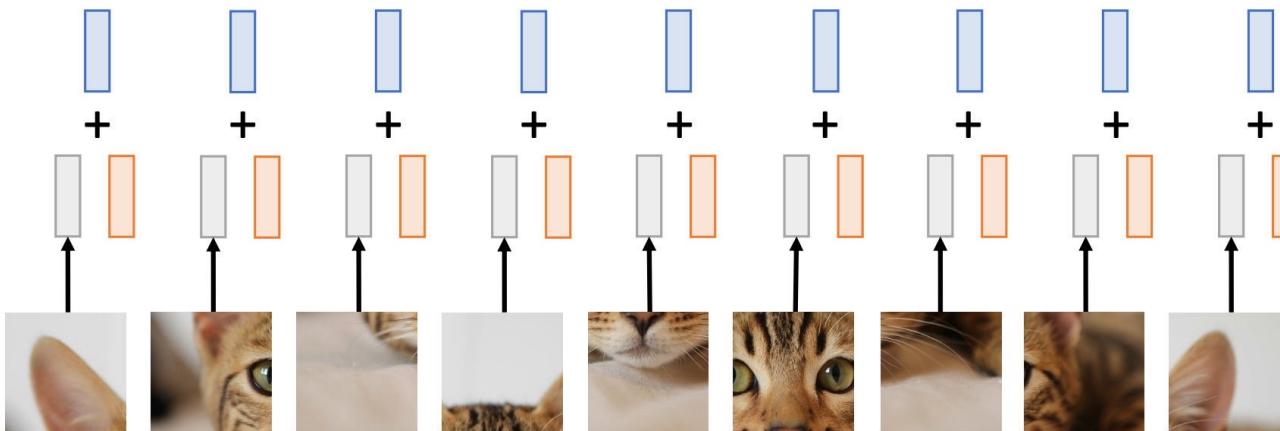
Output vectors



Exact same as  
NLP Transformer!

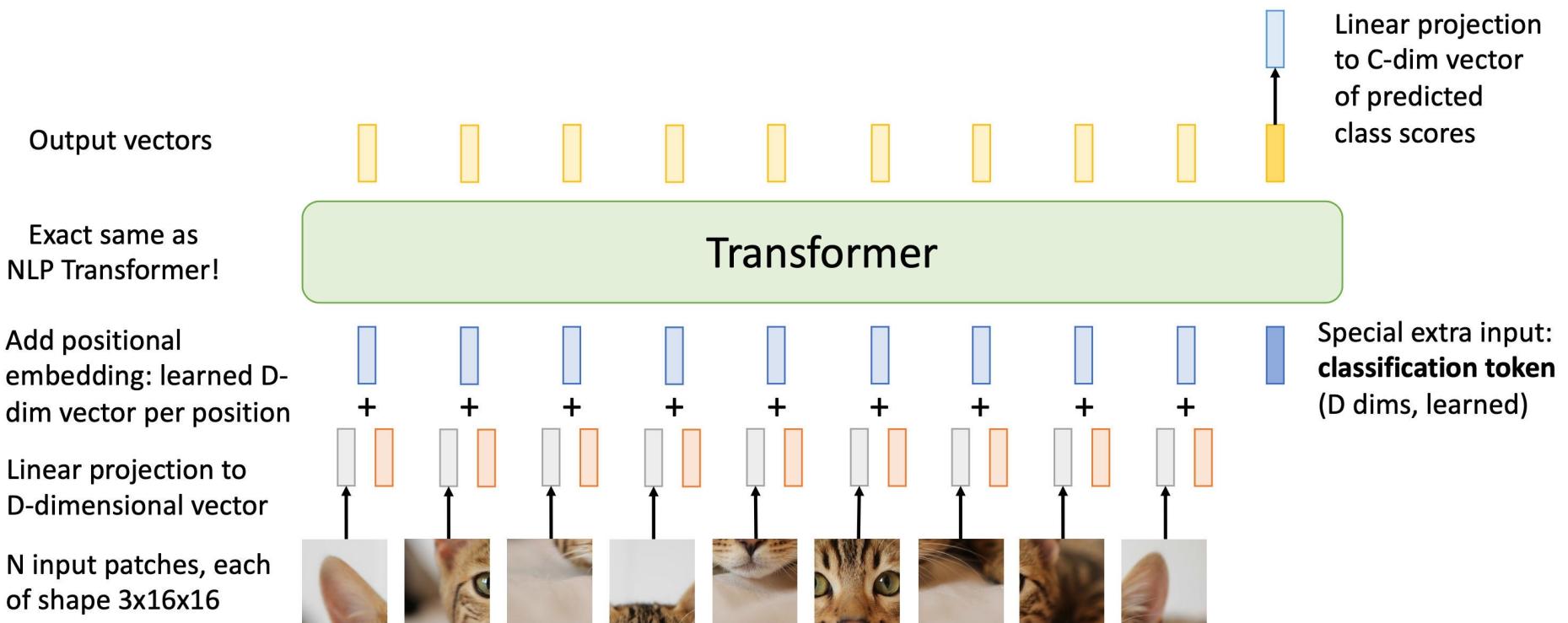
## Transformer

Add positional  
embedding: learned D-  
dim vector per position



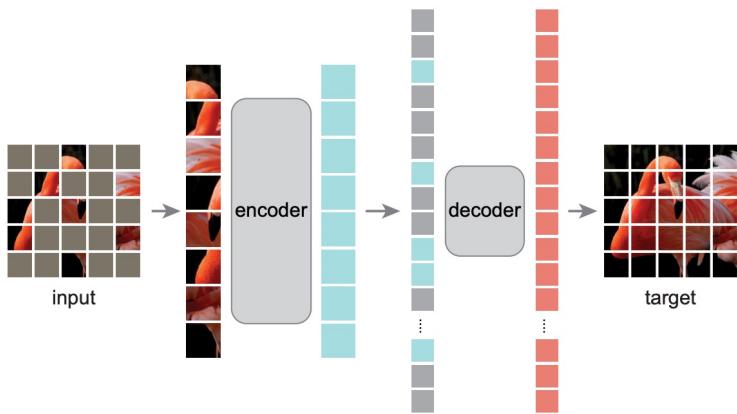
Linear projection to  
D-dimensional vector

N input patches, each  
of shape 3x16x16

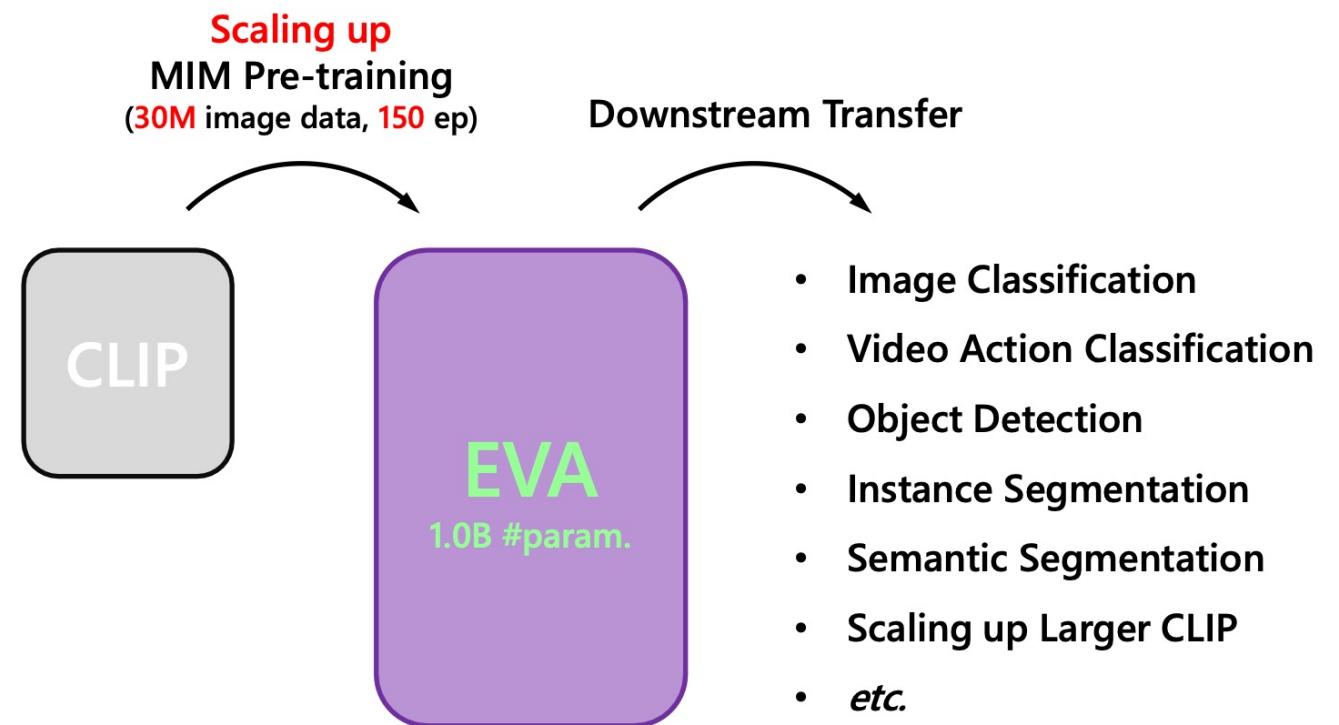


# EVA

- Simply regressing the masked out image-text aligned vision features (*i.e.*, CLIP features) scales up well (to 1.0B parameters) and transfers well to various downstream tasks.

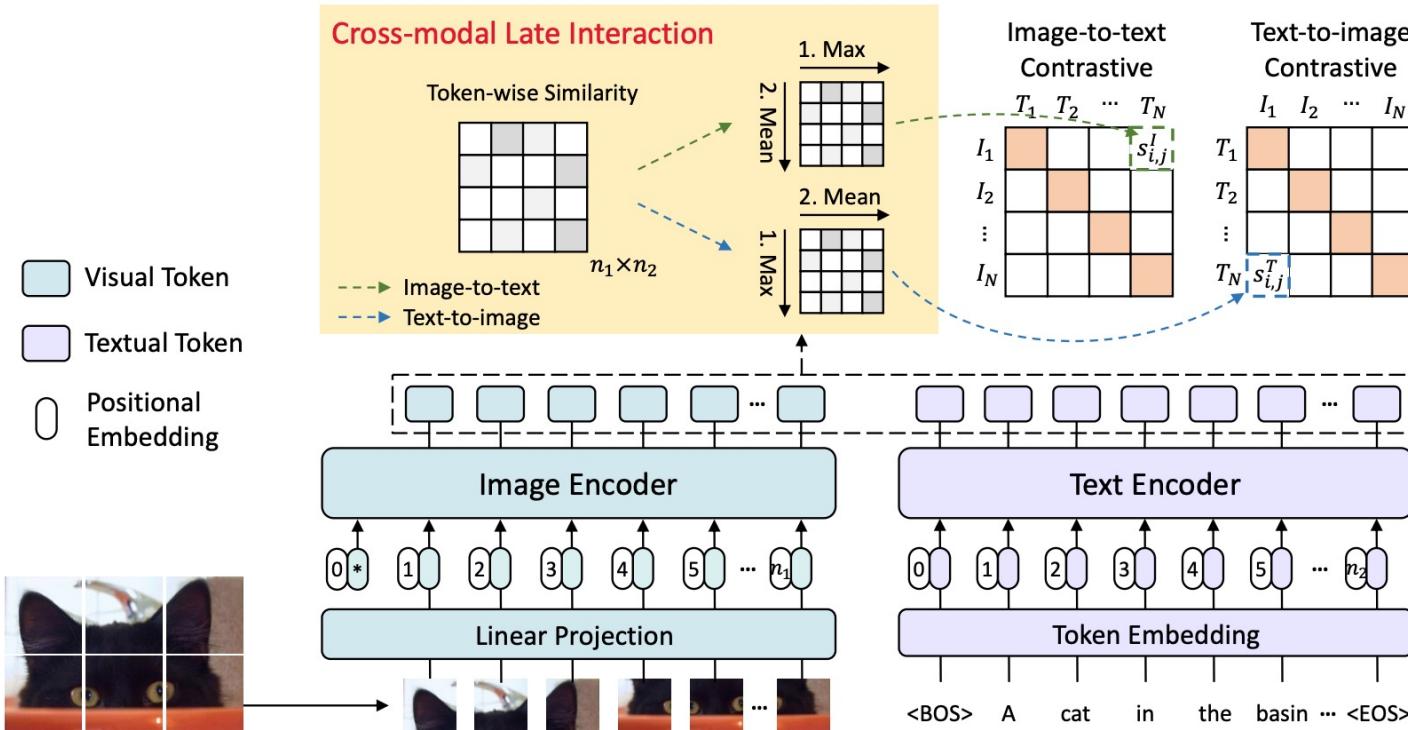


He et al., “Masked Autoencoders Are Scalable Vision Learners”, 2021



# Fine-grained alignment: FILIP

- First computing the token-wise similarity, and then aggregating the matrix by max pooling



$$s_{i,j}^I = \frac{1}{n_1} \sum_{k=1}^{n_1} \max_{l=1, \dots, n_2} I_{i,k}^T T_{j,l}$$

$$s_{i,j}^T = \frac{1}{n_2} \sum_{k=1}^{n_2} \max_{l=1, \dots, n_2} I_{i,l}^T T_{j,k}$$

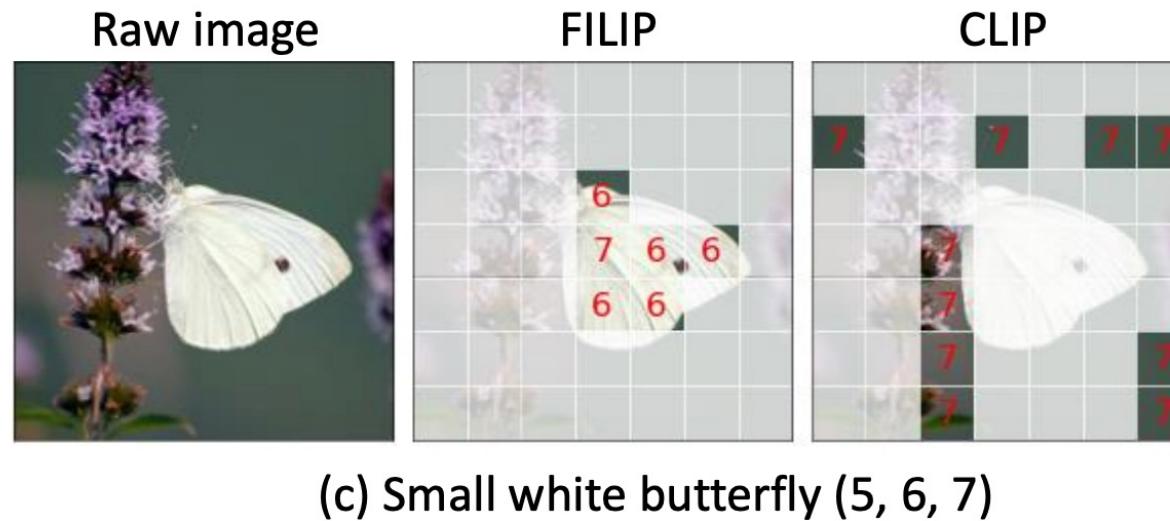
$$\mathcal{L}_i^I = -\frac{1}{b} \log \frac{e^{s_{i,i}^I}}{\sum_{j=1}^b e^{s_{i,j}^I}}$$

$$\mathcal{L}_j^T = -\frac{1}{b} \log \frac{e^{s_{i,i}^T}}{\sum_{i=1}^b e^{s_{i,j}^T}}$$

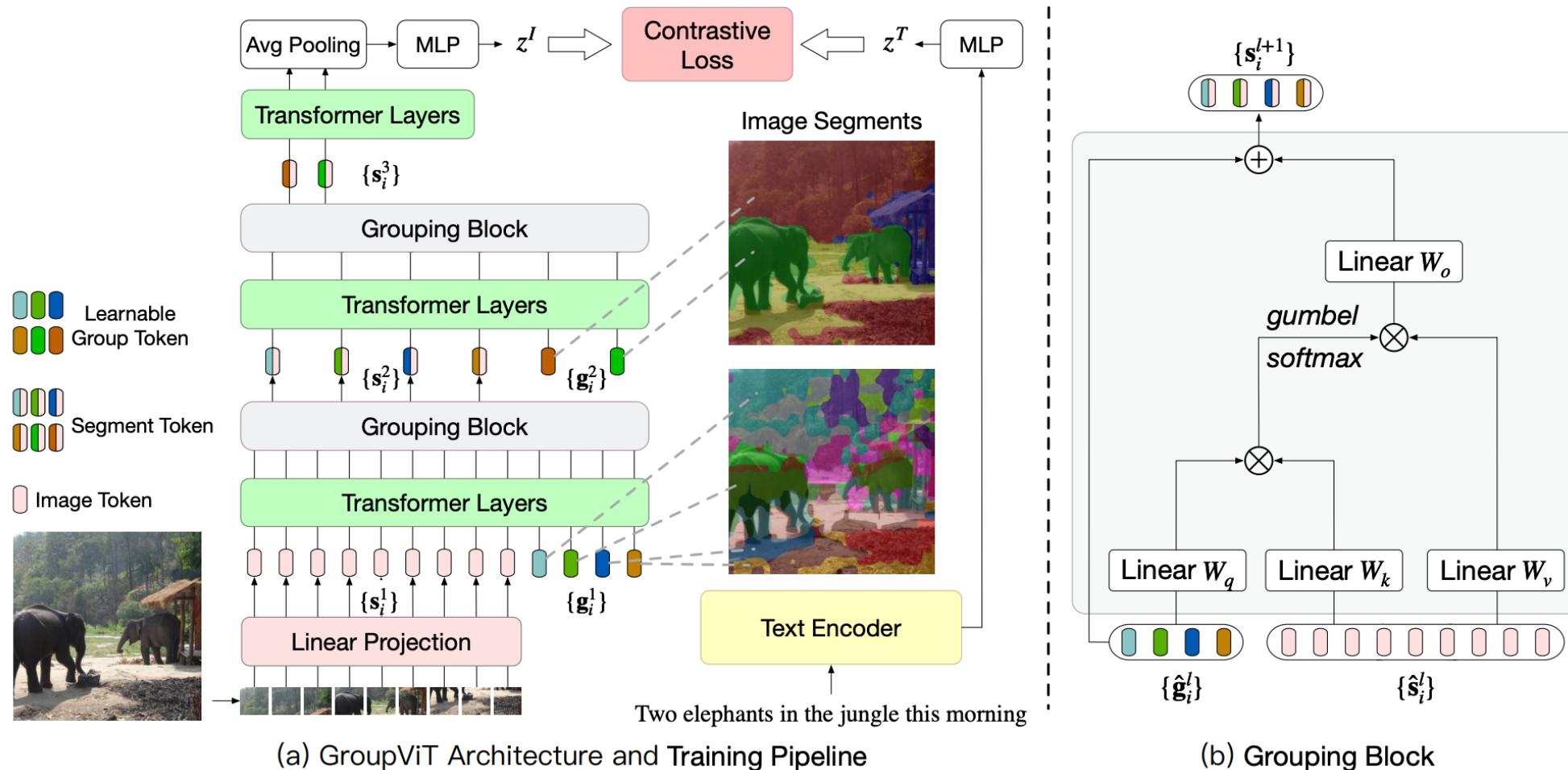
$$\mathcal{L} = \frac{1}{2} \sum_{i=1}^b (\mathcal{L}_i^I + \mathcal{L}_j^T)$$

# FILIP vs. CLIP

- FILIP learns word-patch alignment that is good for visualization

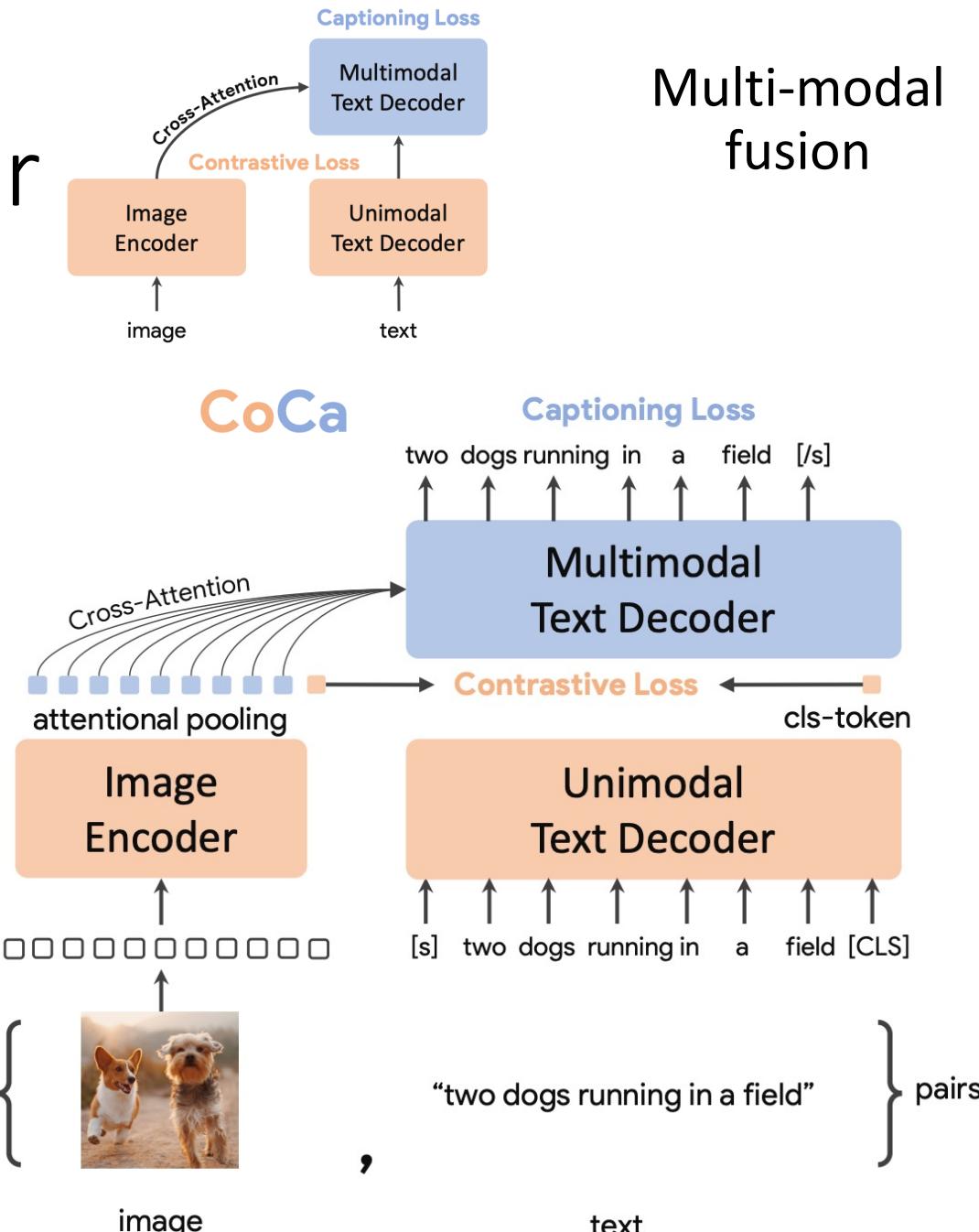


# GroupViT



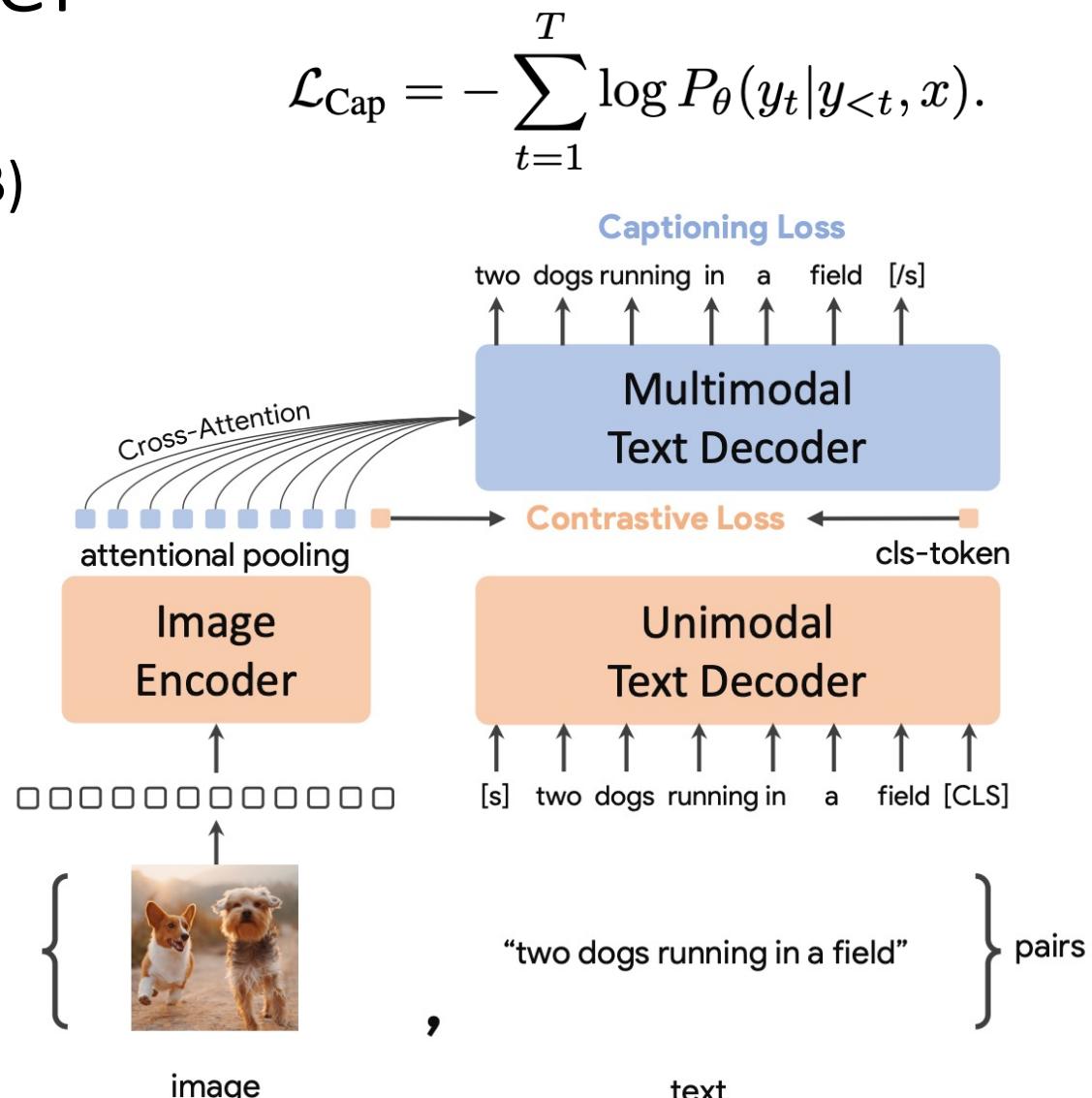
# CoCa: Contrastive Captioner

- Use mixed image-text and image-label (JFT-3B) data for pre-training
- Consider an additional generative branch for enhanced performance and enabling new capabilities (image captioning and VQA)
- CoCa aims to learn a better image encoder from scratch

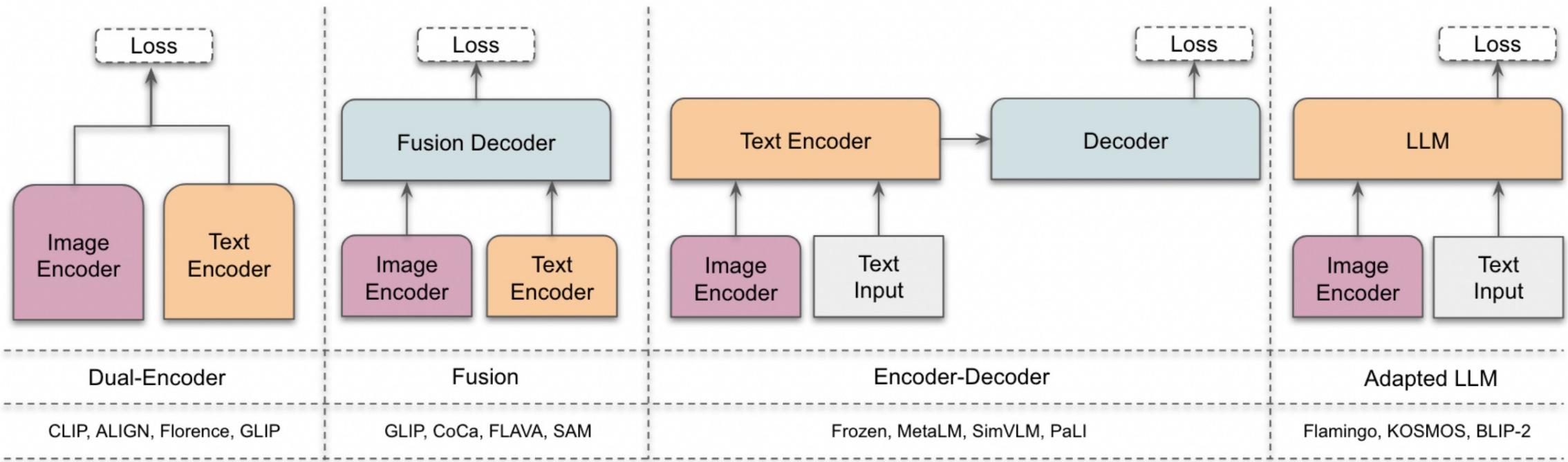


# CoCa: Contrastive Captioner

- Use mixed image-text and image-label (JFT-3B) data for pre-training
- Consider an additional generative branch for enhanced performance and enabling new capabilities (image captioning and VQA)
- CoCa aims to learn a better image encoder from scratch



# Architecture of Multimodal Models



# Conclusion

- VLMs bridge the vision and language spaces
- VLMs showcase impressive capabilities for zero-shot adaptation to unseen tasks
- However, they are still restricted to tasks in a pre-defined form, struggling to match the open-ended task capabilities of LLMs
- A unified generalist framework is required that will be discussed in the next session

## Next Session: Multimodal LLMs

### Large Multimodal Models (LMMs)

- A unified generalist framework that can integrate the strengths of LLMs with the specific requirements of vision-centric tasks.
  - More efficient multi-modal training
  - Multimodal In-Context Learning (M-ICL)
  - Multi-modal Instruction Tuning (M-IT)

# Questions

