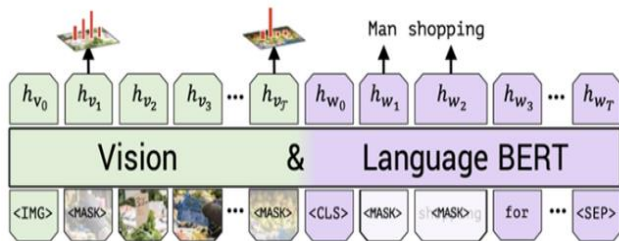Earlier ...

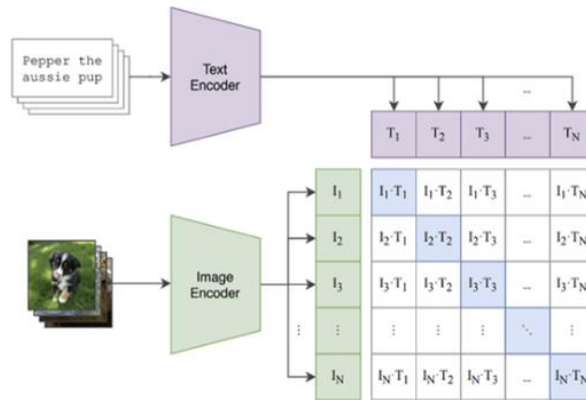**Bert-based Models**



VilBert, VisualBert, VL-BERT, UNITER, ALBERF,HERO, ...

Text models are usually very small
Limited language understanding

**Dual-encoder Contrastive Models**



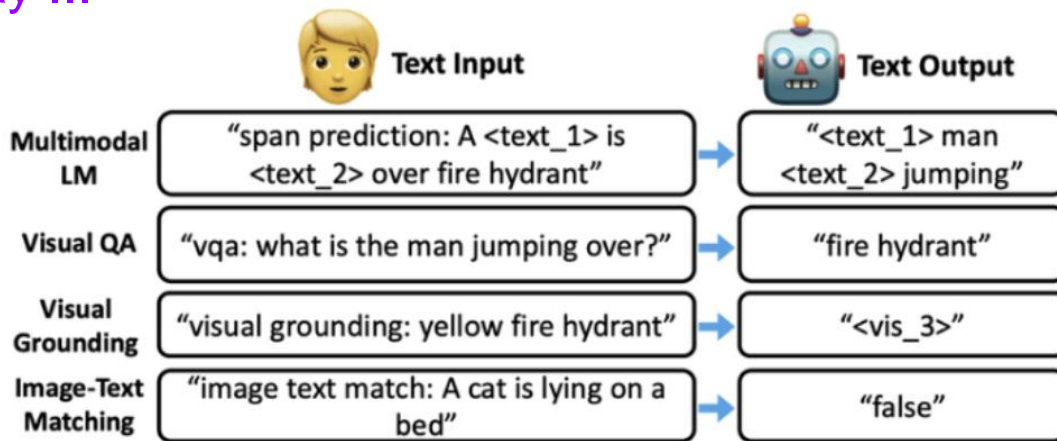CLIP, ALIGN, CoCa, Florence, MIL-NCE, BASIC, LiT, FILIP, MMV, ...

paperswithcode.com/methods/category/vision-and-language-pre-trained-models

Meta

# Vision Language Model Overview

Today ...



| | Text Input | Text Output |
|---|---|---|
| **Multimodal LM** | "span prediction: A <text_1> is <text_2> over fire hydrant" | "<text_1> man <text_2> jumping" |
| **Visual QA** | "vqa: what is the man jumping over?" | "fire hydrant" |
| **Visual Grounding** | "visual grounding: yellow fire hydrant" | "<vis_3>" |
| **Image-Text Matching** | "image text match: A cat is lying on a bed" | "false" |

**Vision Language Model**
- Capable of performing different tasks, including VQA.
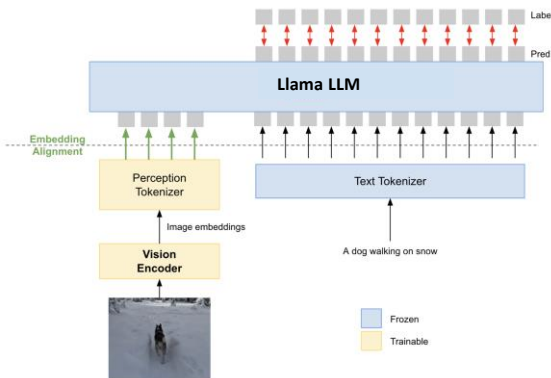- Additional knowledge and intelligence comes from adding a large language model to the VLM

# How LLMs Understand Images?
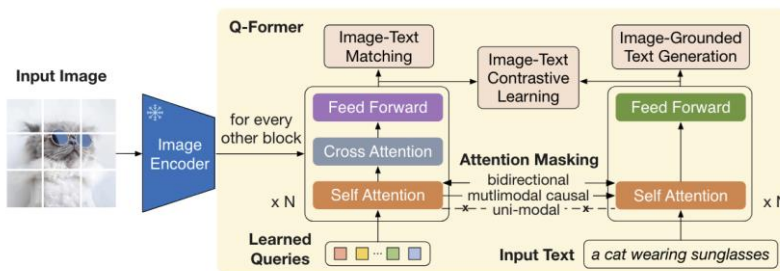
**Fusing vision tokens with Large Language Models**

**Why is the vision encoder important for MM-LLM?**

- Allows LLMs to input/ understand images

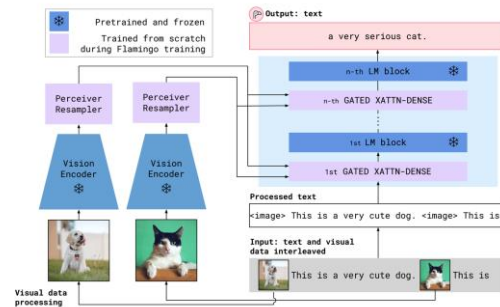**How to fuse and align modality tokens with the large language model?**
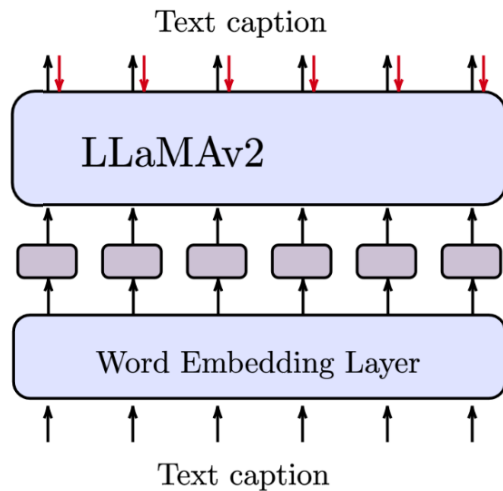


AnyMAL | LLaVa | Fuyu

BLIP-v2

Flamingo

Li et. al, BLIP-2: Bootstrapping Language-Image Pre-training.
Moon et. al, AnyMAL: An Efficient and Scalable Any-Modality Augmented
Language Model
Alayrac et. al, Flamingo: a Visual Language Model for Few-Shot Learning

∞ Meta
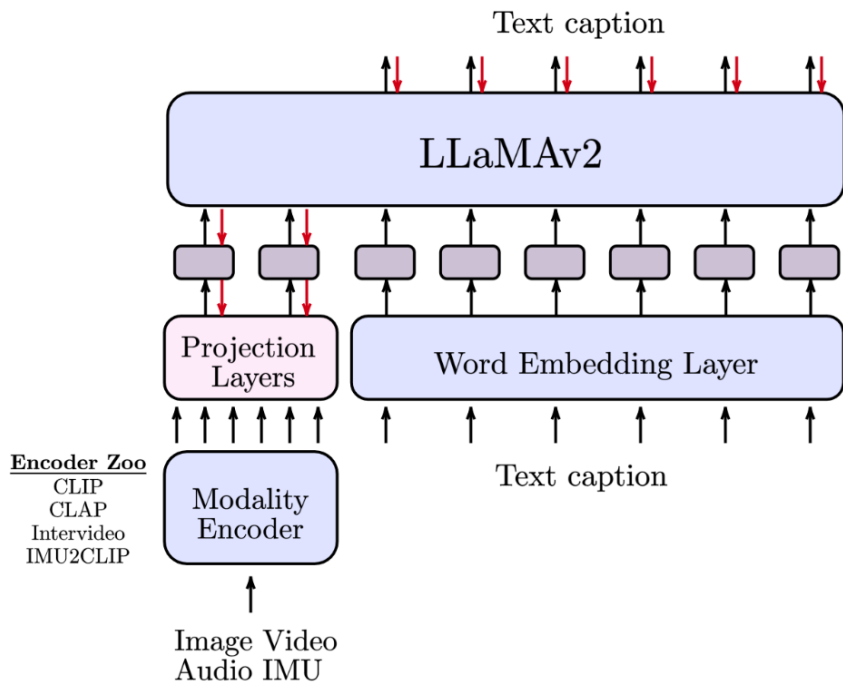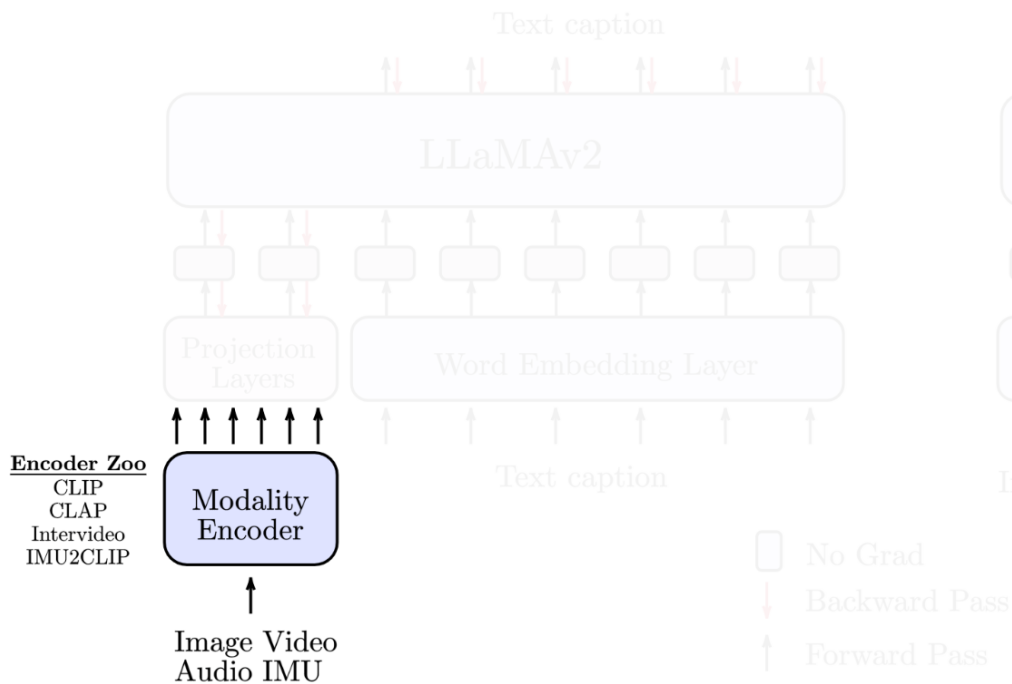
## Base Language Model

- Preserves the strong language-based capabilities
- Variations
  - OPT
  - FlanT5
  - Llama
  - Llama-2-chat

# AnyMAL Overview (2)
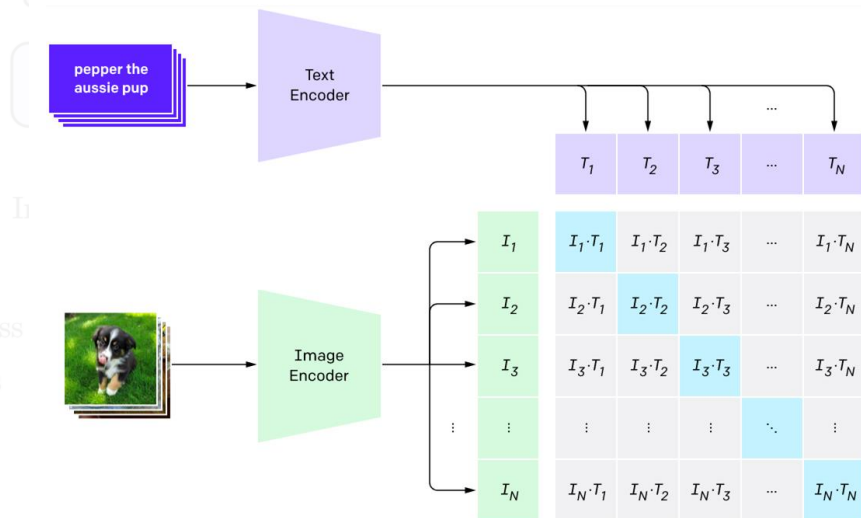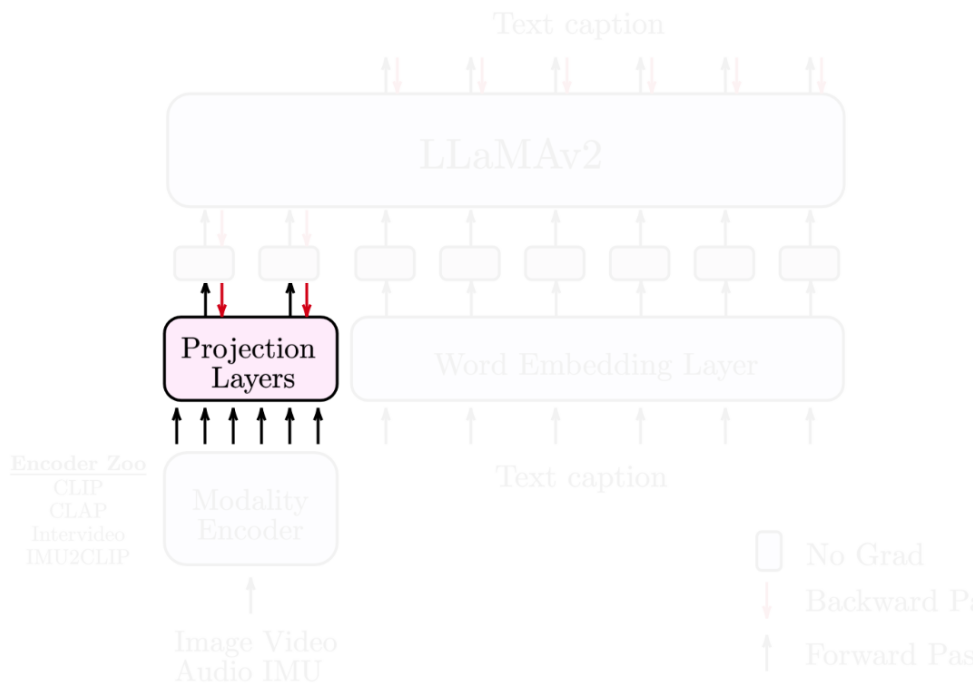
**Modality Encoder**

- Trained with contrastive loss (text & other modality) for the best alignment in the text space



© 2024 Meta

OpenAI, "CLIP: Connecting Text and Images", 2021
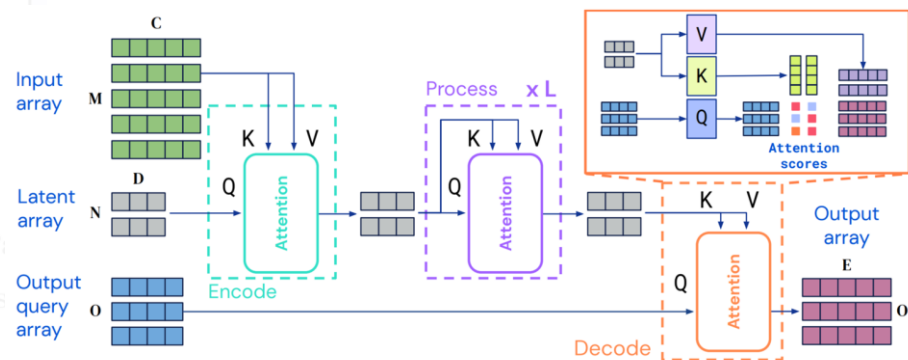
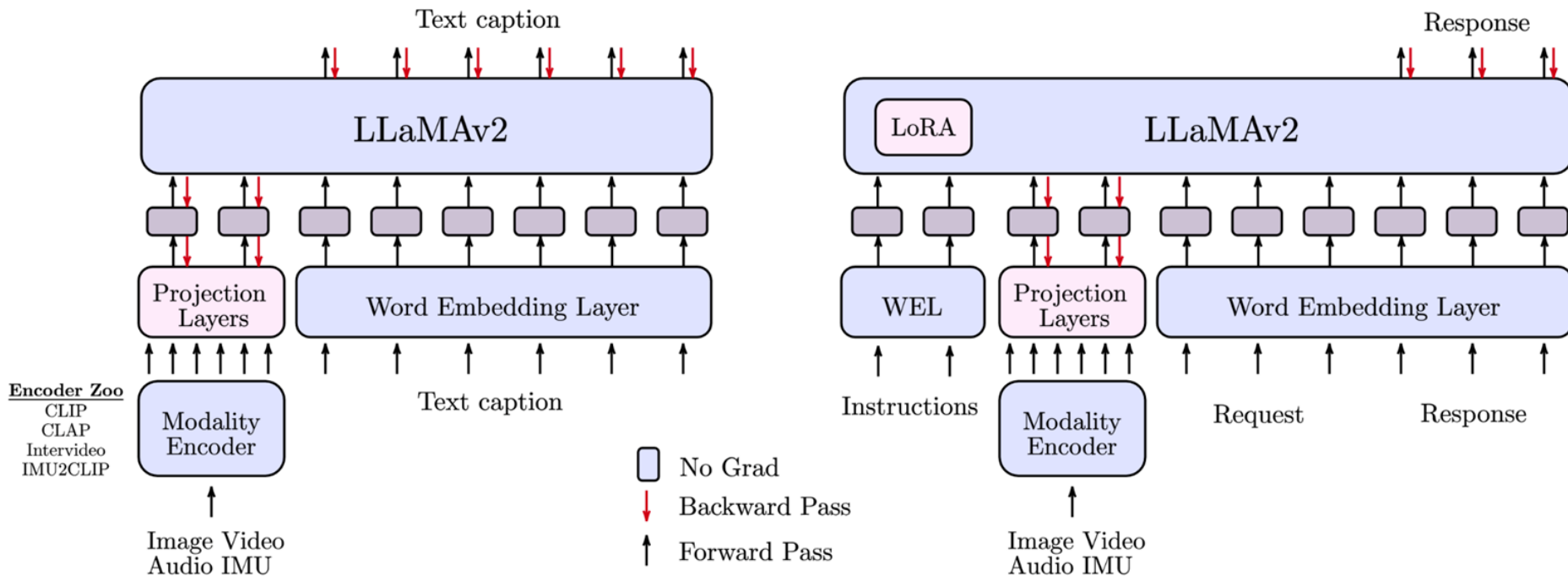## Projection Layers

- Perceiver Resampler to resample patch embeddings into a sequence of Llama-compatible tokens

a) Modality Alignment

b) Multimodal Instruction Tuning

Moon et. al, AnyMAL: An Efficient and Scalable Any-Modality Augmented Language Model

# Vision language model training
# =
# Vision encoder with language alignment

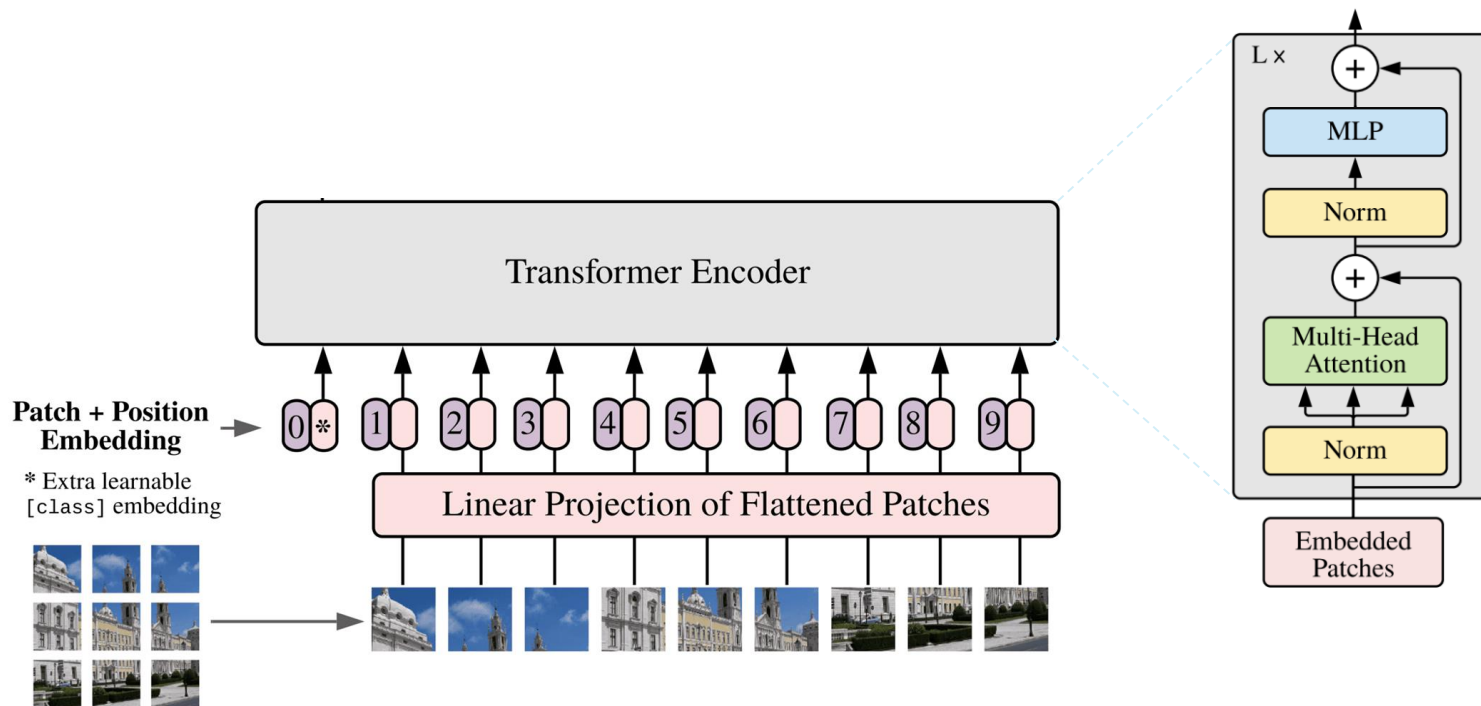# What Does the Vision Encoder See?

Look at that person

What is the name of that hotel?

What is that?

Li et. al.  CLIP Surgery for Better Explainability with Enhancement in Open-Vocabulary Tasks

# Vision Transformer Backbone

**Patch + Position Embedding**

\* Extra learnable [class] embedding

Transformer Encoder

0\* 1 2 3 4 5 6 7 8 9

Linear Projection of Flattened Patches

L ×

MLP

Norm

Multi-Head Attention

Norm

Embedded Patches

Dosovitskiy et. al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale
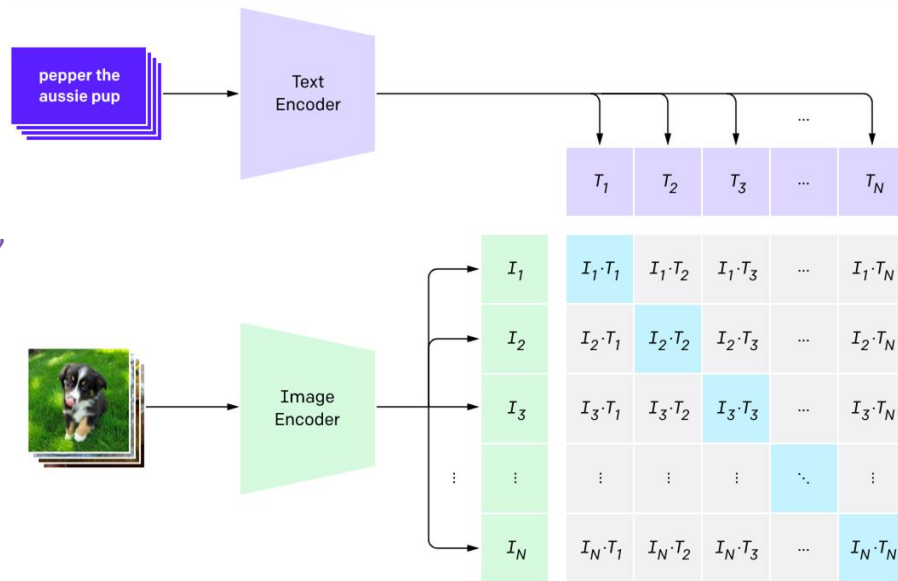
# Training a High Quality CLIP Vision Encoder (1)

**Key ingredients:**

- Data volume and quality

- Model scale (number of parameters)

- Efficient pre-training recipes

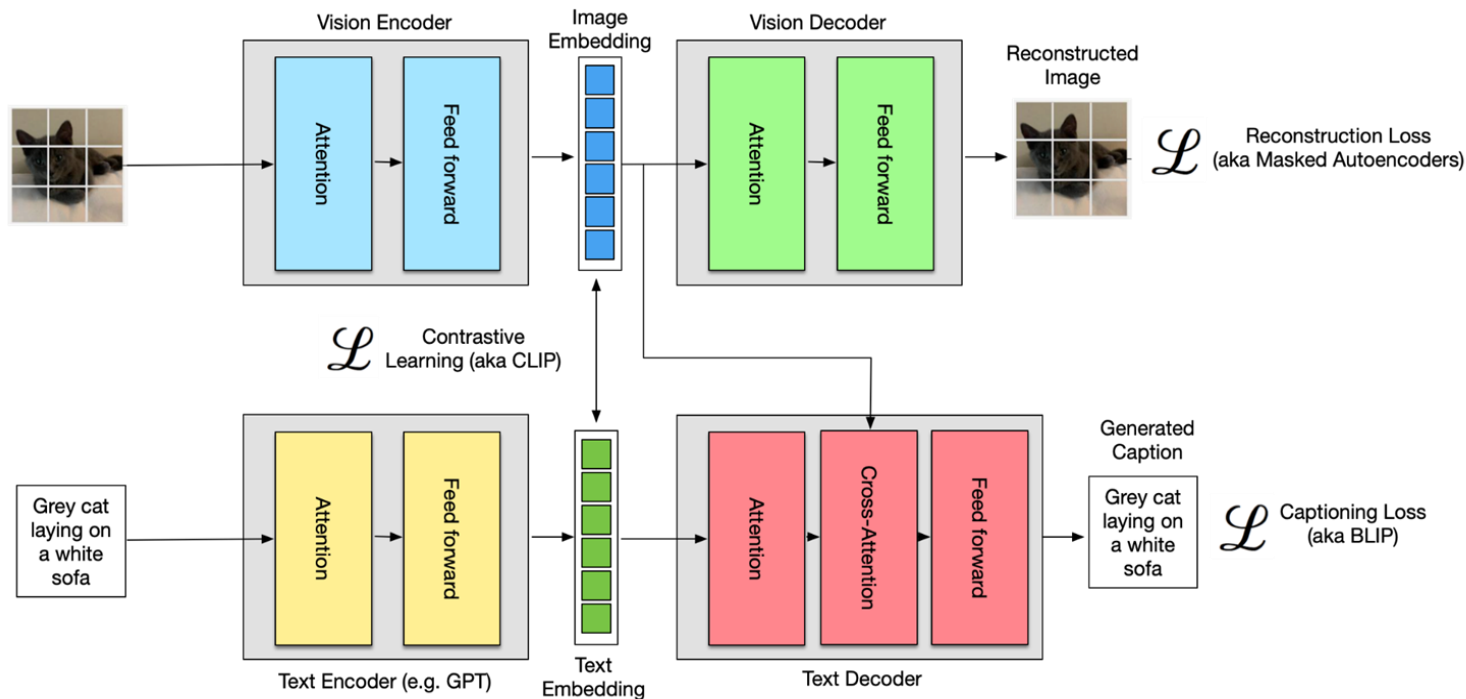- Domain fine-tuning: higher resolutions, diverse data, better caption quality

**Measuring quality:**

- Zero-shot ImageNet (IN1k-0 shot)

- Post-alignment metrics

  - VQA accuracy

  - Captioning accuracy



Radford et. al. Learning Transferable Visual Models From Natural Language Supervision
Li et. al, An Inverse Scaling Law for CLIP Training

**Auxiliary loss functions**

Improving alignment performance with LLM

Yu et. al. CoCa: Contrastive Captioners are Image-Text Foundation Models

# Resources

We are hiring!

www.metacareers.com

AI at Meta
ai.meta.com

For more information:

**AnyMAL**
Any-Modality Augmented Language Model

**Seungwhan Moon**\*, Andrea Madotto\*, Zhaojiang Lin\*, Tushar Nagarajan\*, Matt Smith, Shashank Jain, Chun-Fu Yeh,  Prakash Murugesan, Peyman Heidari, Yue Liu, Kavya Srinet, Babak Damavandi, Anuj Kumar

**Meta Reality Labs & FAIR**