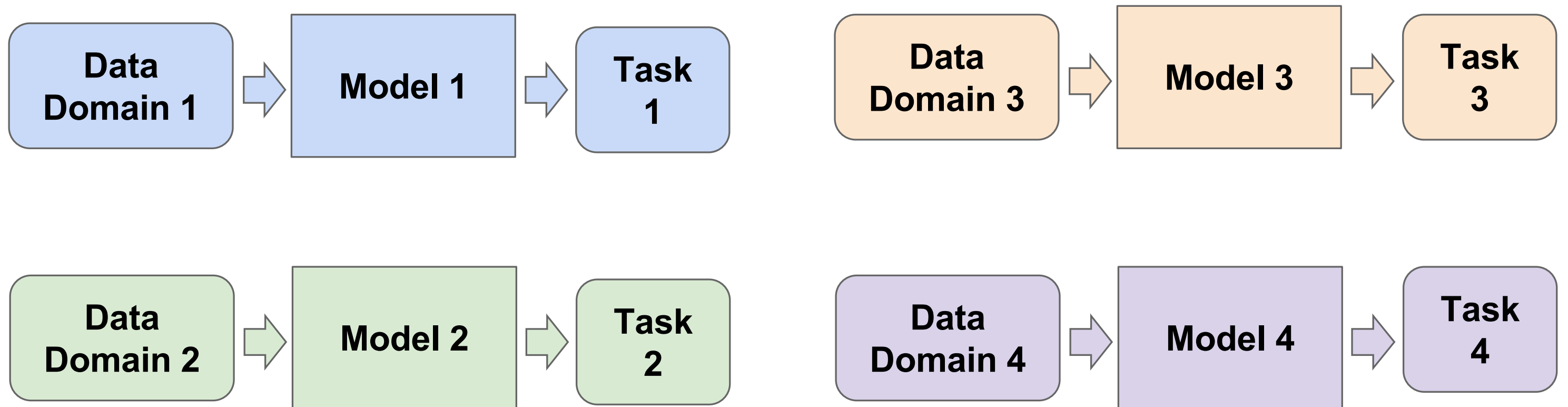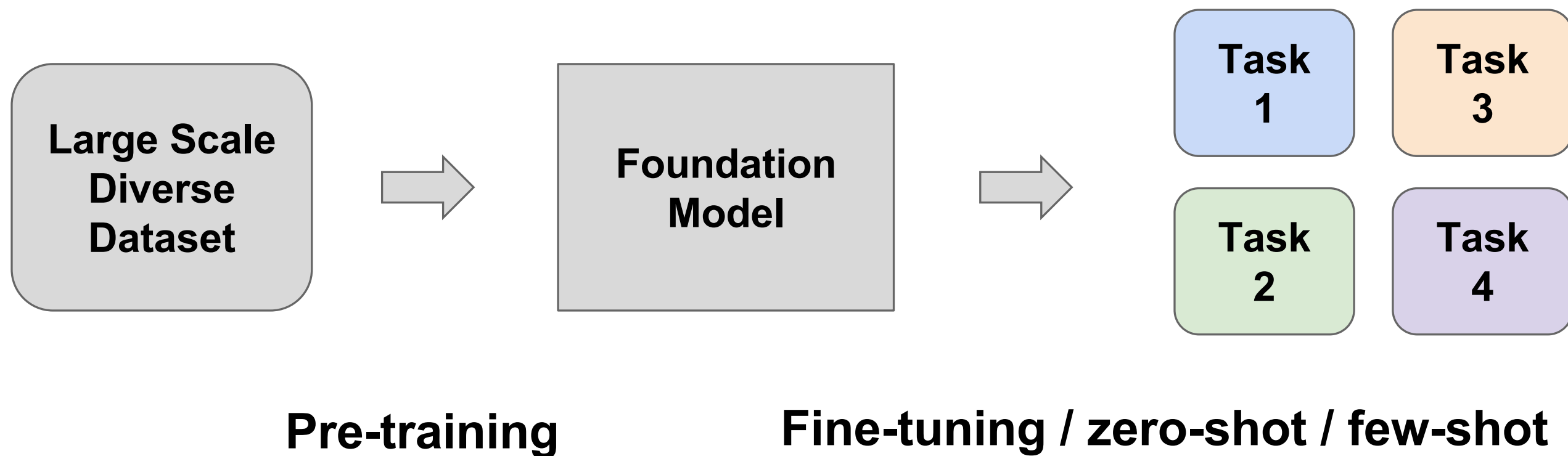# Lecture 16:
# Multi-Modal Foundation Models

# How have we been thinking about models in this class so far?

Train a *specialized* model for *each* task
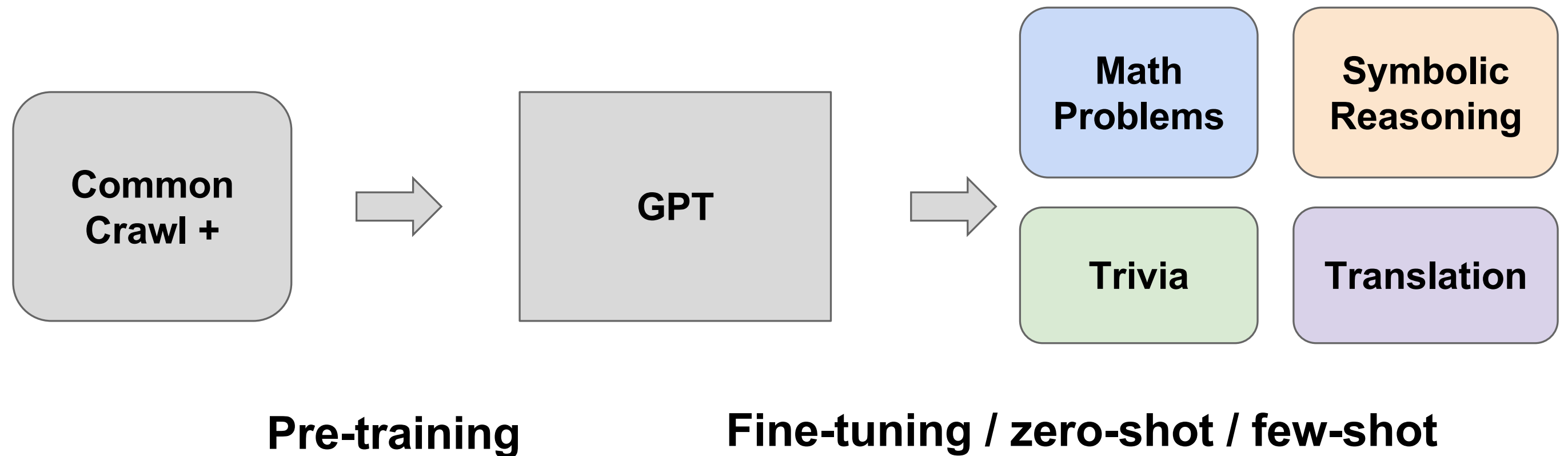
# Now, we build Foundation Models

*Pre-train one* model that acts as the *foundation* for many different tasks



**Pre-training**          **Fine-tuning / zero-shot / few-shot**

# Foundation Models

## Language



Common Crawl +  →  GPT  →  Math Problems / Symbolic Reasoning / Trivia / Translation

**Pre-training**  **Fine-tuning / zero-shot / few-shot**

# There are many classes of Foundation Models

| Language | Classification | LM + Vision | And More! | Chaining |
|---|---|---|---|---|
| ELMo | CLIP | LLaVA | Segment Anything | LMs + CLIP |
| BERT | CoCa | Flamingo | Whisper | Visual Programming |
| GPT | | GPT-4V | Dalle | |
| T5 | | Gemini | Stable Diffusion | |
| | | Molmo | Imagen | |

How do identify a model as a Foundation?

**Always see with foundation models:**

- general /robust to many different tasks

**Often see with foundation models:**

- Large # params
- Large amount of data
- Self-supervised pre-training objective

# Language models are out of scope for this class

| Language | Classification | LM + Vision | And More! | Chaining |
|---|---|---|---|---|
| **ELMo** | CLIP | LLaVA | Segment Anything | LMs + CLIP |
| **BERT** | CoCa | Flamingo | Whisper | Visual Programming |
| **GPT** | | GPT-4V | Dalle | |
| **T5** | | Gemini | Stable Diffusion | |
| | | Molmo | Imagen | |

# We will focus on multimodal (vision) foundation models

| Language | Classification | LM + Vision | And More! | Chaining |
|----------|----------------|-------------|-----------|----------|
| **ELMo** | **CLIP** | **LLaVA** | **Segment Anything** | **LMs + CLIP** |
| **BERT** | **CoCa** | **Flamingo** | Whisper | **Visual Programming** |
| **GPT** | | GPT-4V | Dalle | |
| **T5** | | Gemini | Stable Diffusion | |
| | | **Molmo** | Imagen | |

# Let's start with the foundation models for classification

| Language | Classification | LM + Vision | And More! | Chaining |
|---|---|---|---|---|
| **ELMo** | **CLIP** | **LLaVA** | **Segment Anything** | **LMs + CLIP** |
| **BERT** | **CoCa** | **Flamingo** | Whisper | **Visual Programming** |
| **GPT** | | GPT-4V | Dalle | |
| **T5** | | Gemini | Stable Diffusion | |
| | | **Molmo** | Imagen | |

# Recall this self-supervised objective from SimCLR



Use Self Supervised learning to learn good image features

Can train small classifiers on top of these features using supervised learning

# The main idea was to learning concepts without labels -> a self-supervised pretraining objective



Pull Together

Push Apart

The hope was that the learned representations generalize to new instances

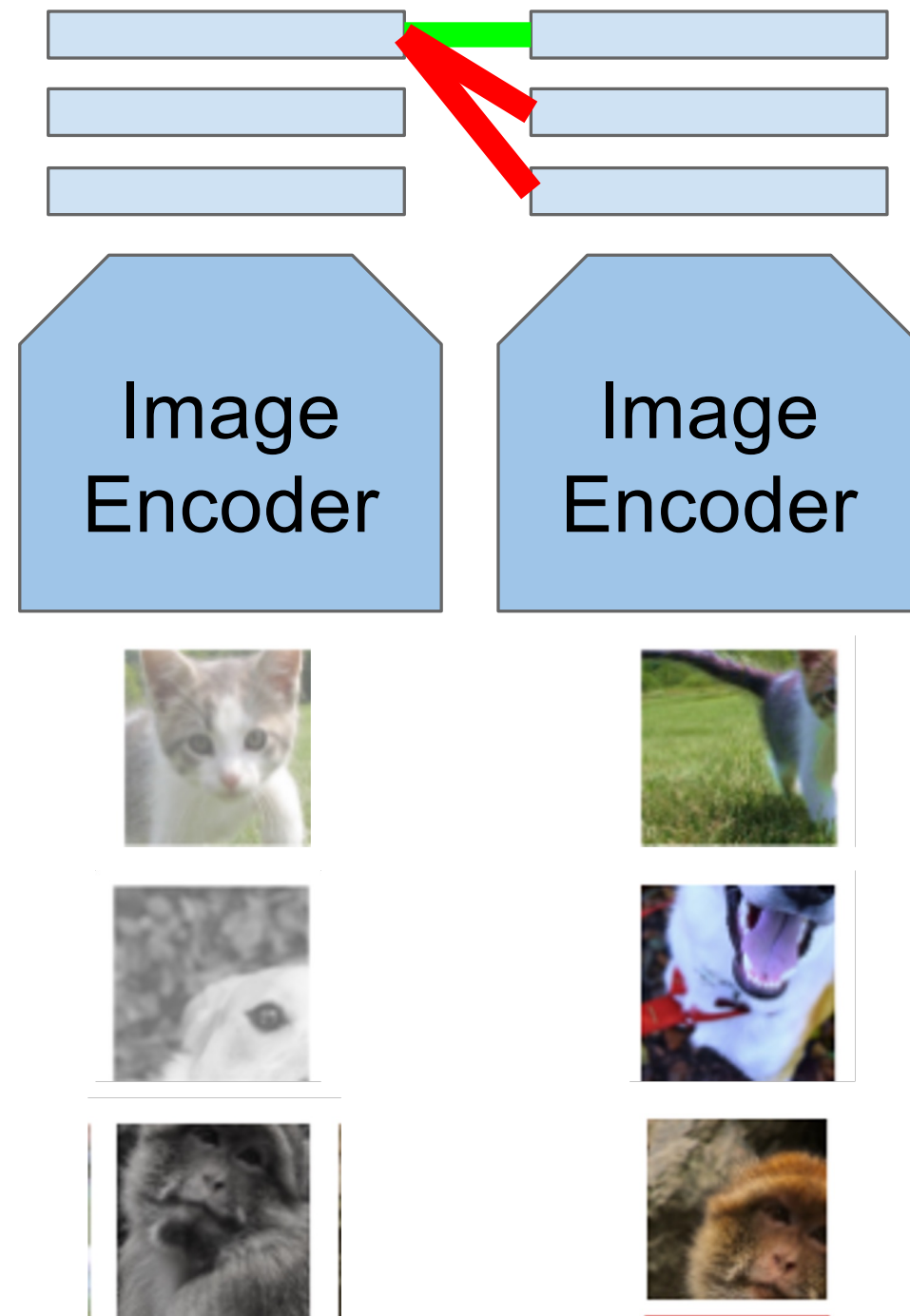# Can we generalize these representations beyond just images? *To language perhaps?*



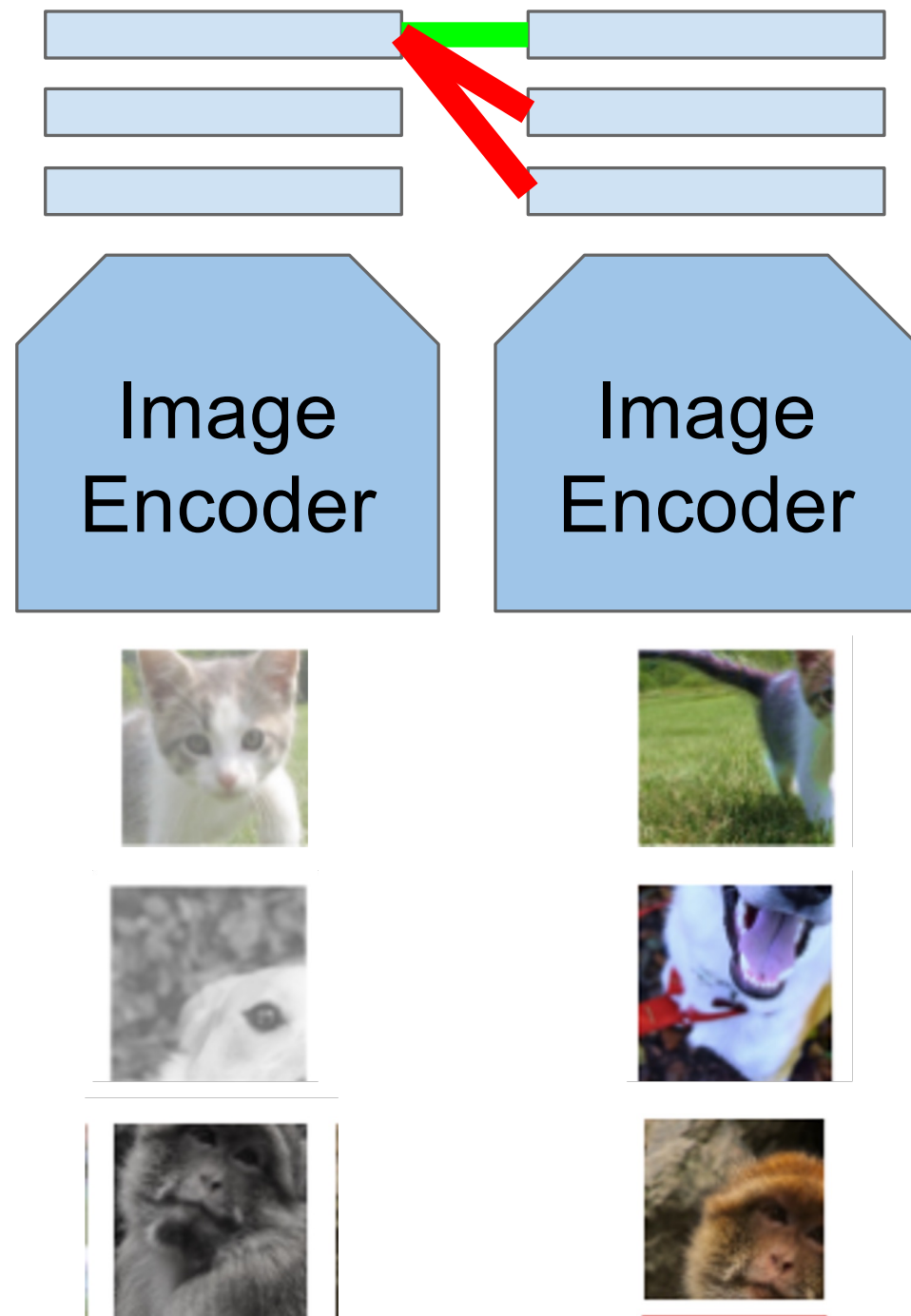1. "A cute fluffy cat"
2. "My favorite dog is a golden retriever"

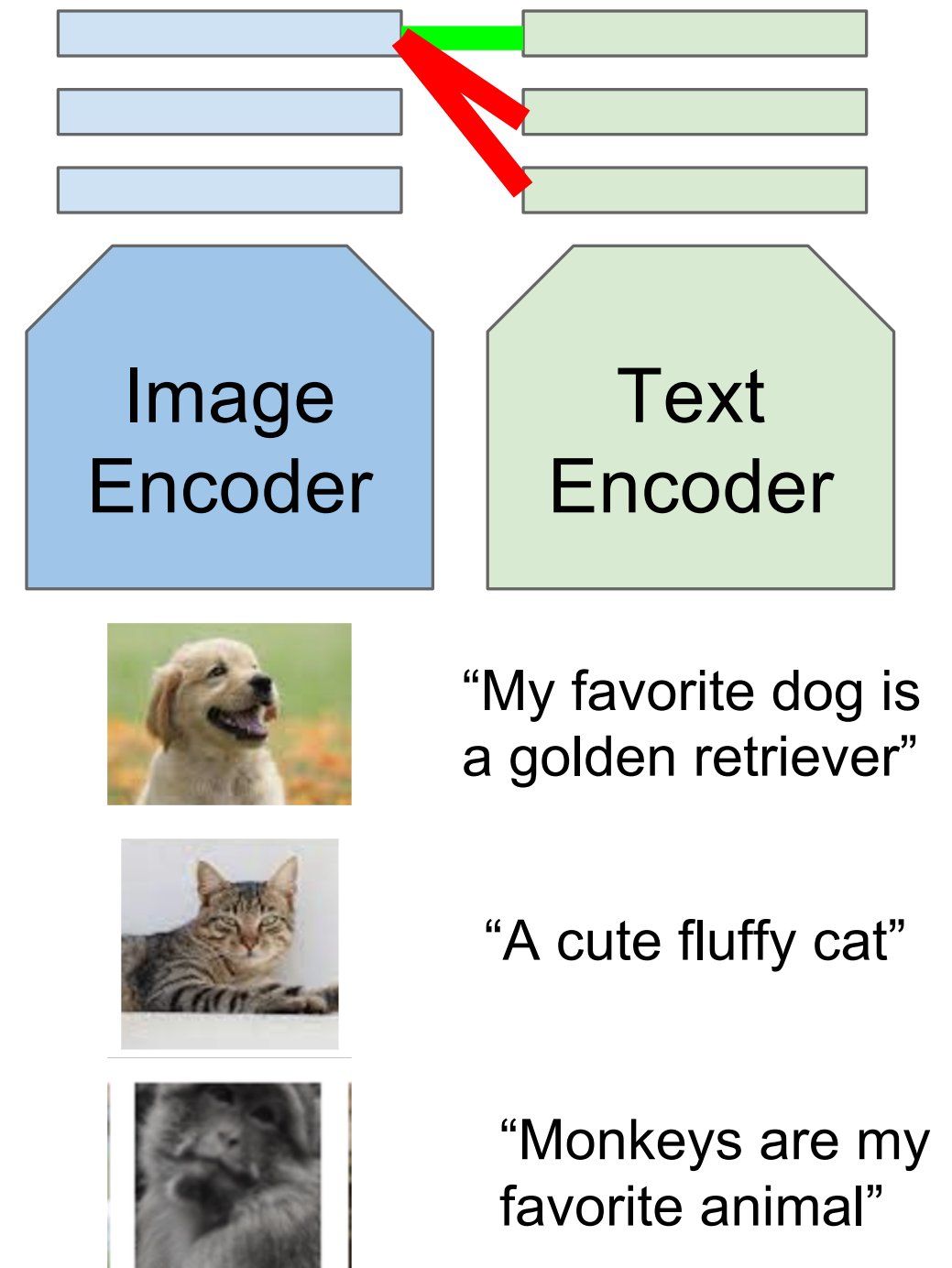What if this representation space could also embed sentences/phrases?

"My favorite dog is a golden retriever"

"A cute fluffy cat"

# SimClr

# SimClr

# CLIP



"My favorite dog is a golden retriever"
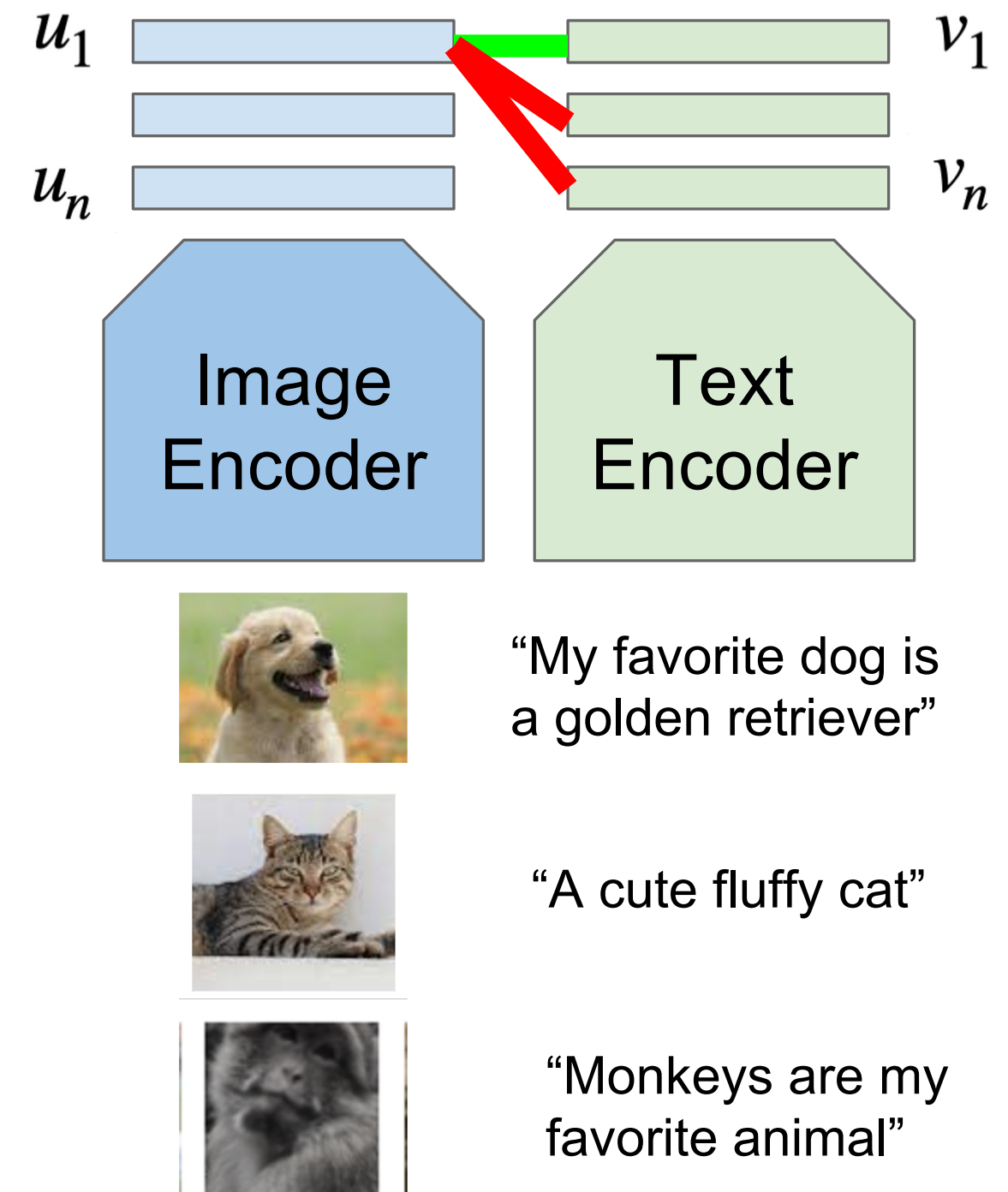
"A cute fluffy cat"

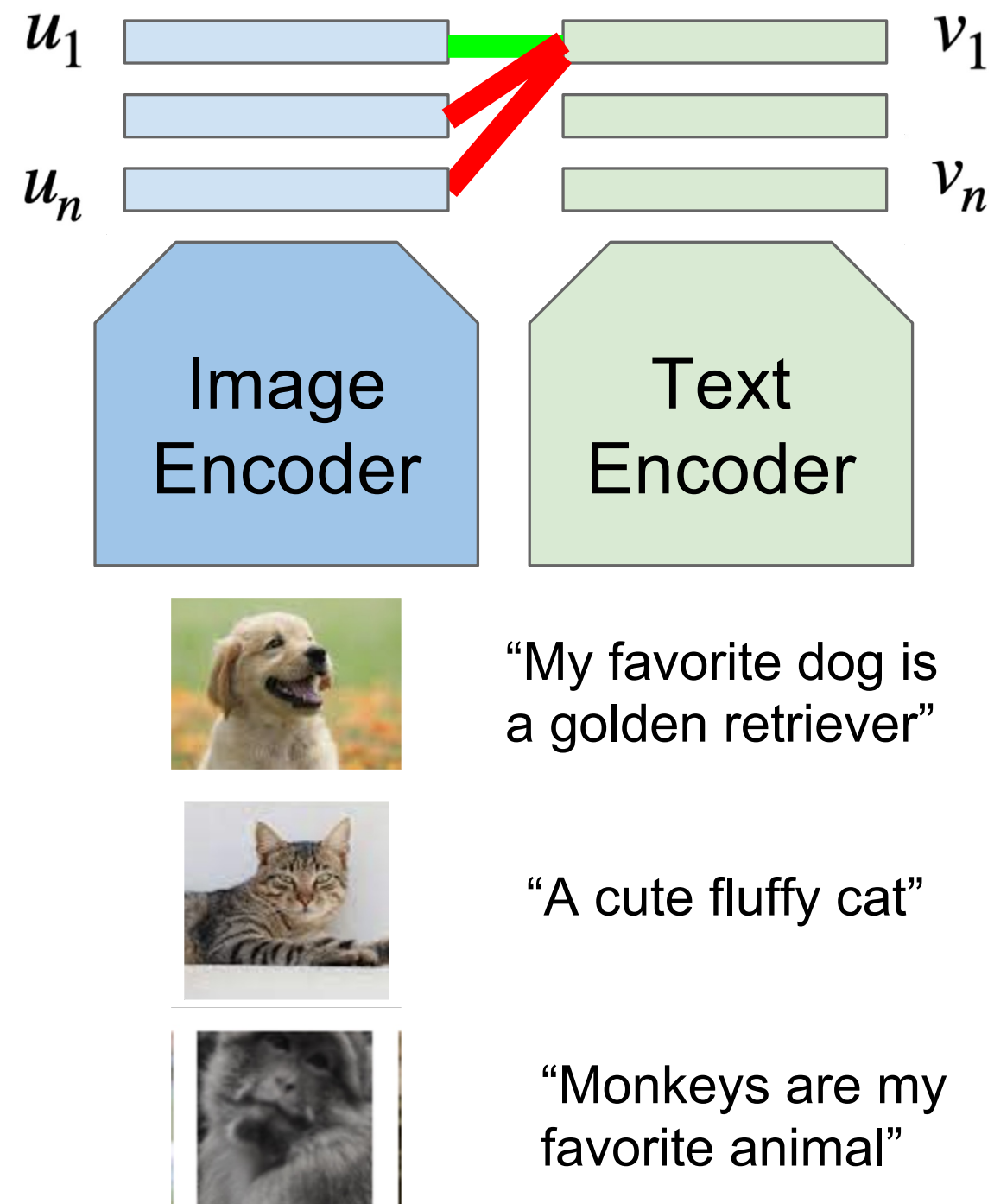"Monkeys are my favorite animal"

# CLIP is trained with the same contrastive objective

$$\sum_{i=1}^{n} -\log\left(\frac{e^{\langle u_i, v_i\rangle}}{\sum_{j=1}^{n} e^{\langle u_i, v_j\rangle}}\right)$$

# CLIP Training Objective

$$\sum_{i=1}^{n} -\log\left(\frac{e^{\langle u_i, v_i\rangle}}{\sum_{j=1}^{n} e^{\langle u_i, v_j\rangle}}\right)$$

$$+\sum_{i=1}^{n} -\log\left(\frac{e^{\langle u_i, v_i\rangle}}{\sum_{j=1}^{n} e^{\langle u_j, v_i\rangle}}\right)$$



Image Encoder

Text Encoder

"My favorite dog is a golden retriever"

"A cute fluffy cat"

"Monkeys are my favorite animal"

# Lots of image-text data can be found online



Mount Rainier's northwestern slope viewed aerially just before sunset on September 6, 2020

CLIP training data was scraped at scale from images and their associated alt-text from the internet

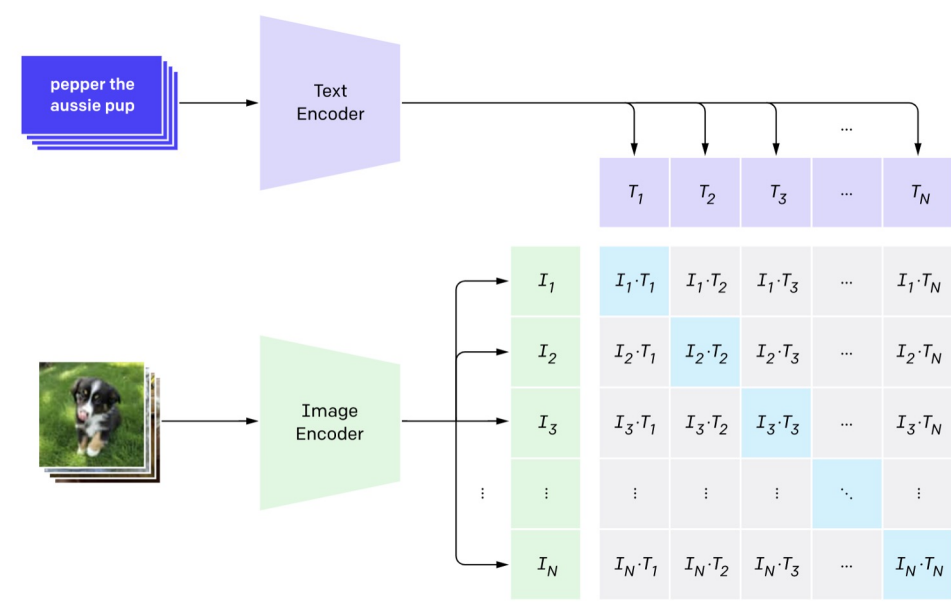https://en.wikipedia.org/wiki/Mount_Rainier

# CLIP Training Objective



**1. Contrastive pre-training**

At the end of training, you have a model that will give you a similarity score between an image and a text
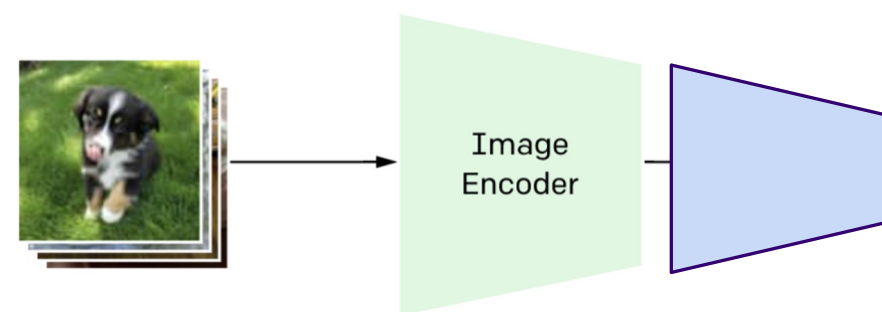
# Using pre-trained models out of the box

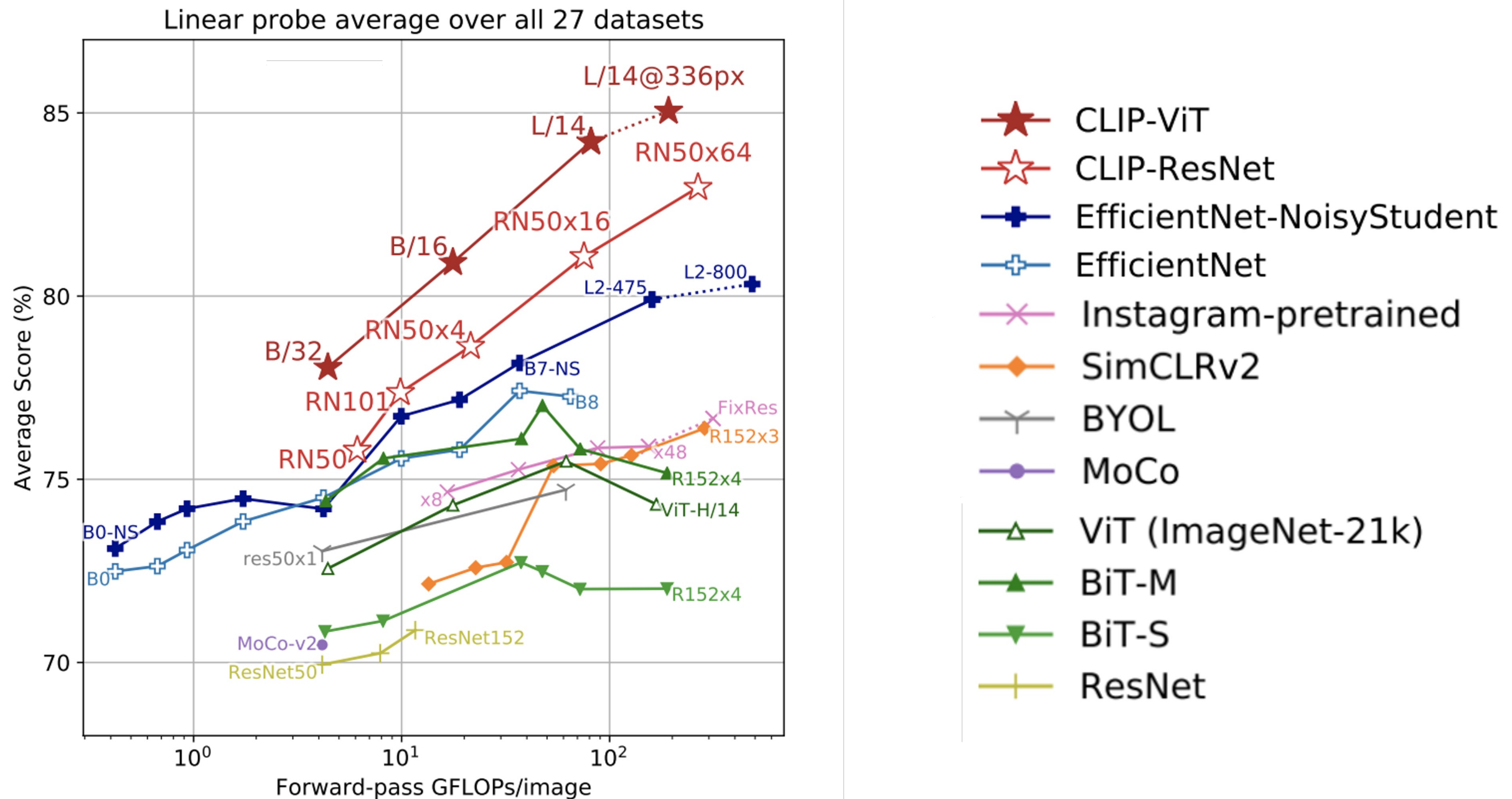**Step 1:** <u>Pretrain</u> a network on a <u>pretext task</u> that doesn't require supervision



**Pre-training tasks:**
Contrastive Objective

**Step 2:** Transfer encoder to <u>downstream tasks</u> via <span style="color:red">linear classifiers</span>
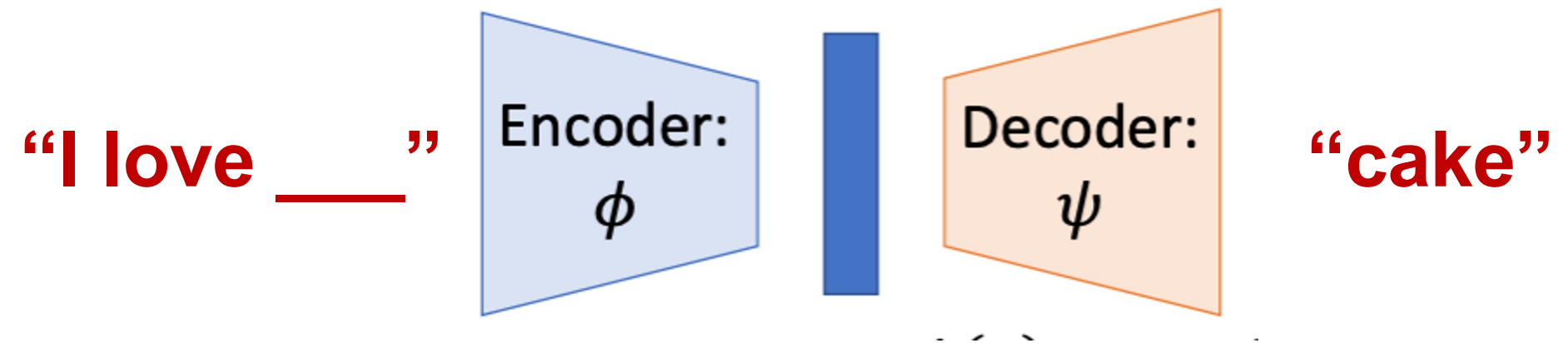


**Downstream tasks:**
Image classification, object detection, semantic segmentation

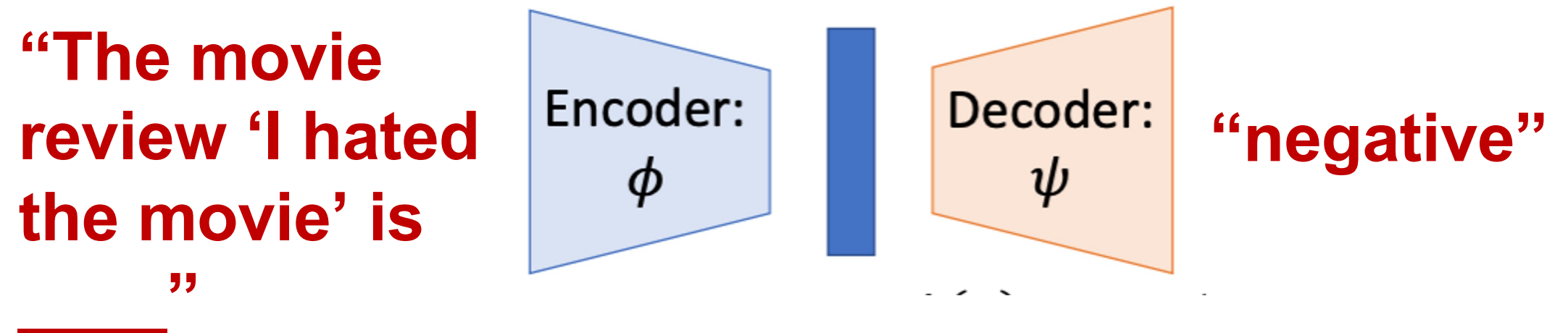# CLIP features w/ linear probe across multiple datasets

# Big difference with language models: We can use LLMs zero-shot for new downstream tasks

**Step 1:** <u>Pretrain</u> a network on a <u>pretext task</u> that doesn't require supervision
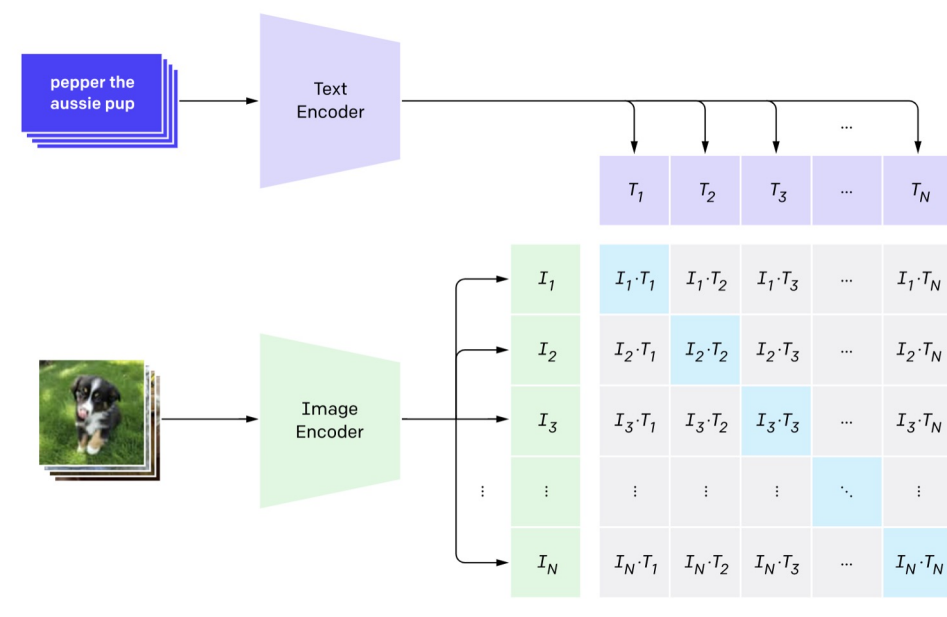


**"I love ___"** → Encoder: $\phi$ → Decoder: $\psi$ → **"cake"**

**Step 2:** Use the model out of the box in a creative way!



**"The movie review 'I hated the movie' is ___"** → Encoder: $\phi$ → Decoder: $\psi$ → **"negative"**

# But how do we use pre-trained <span style="color:red">vision-language</span> models in a <span style="color:red">zero-shot</span> manner?

**Step 1:** <u>Pretrain</u> a network on a <u>pretext task</u> that doesn't require supervision
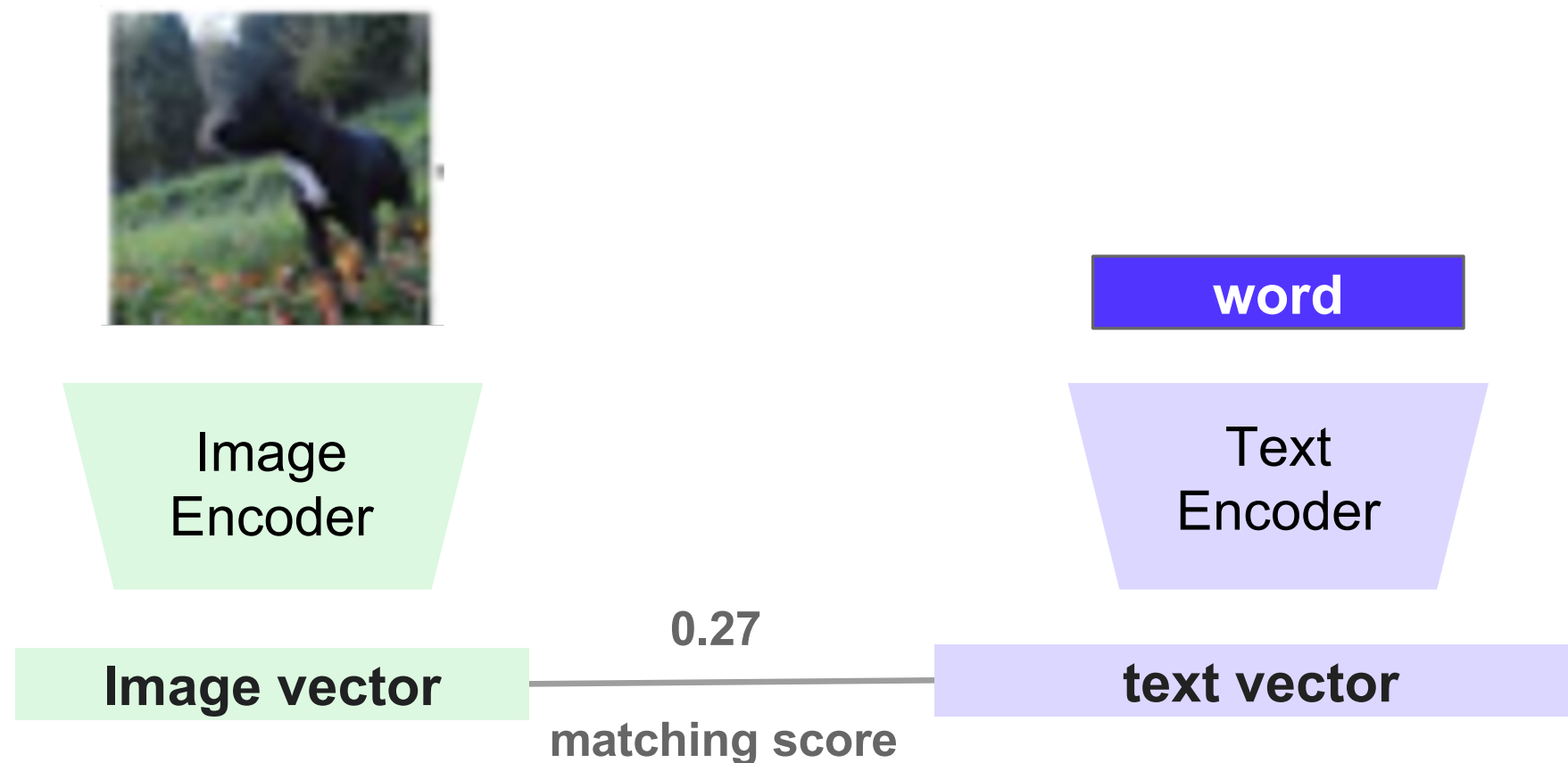


**Pre-training tasks:**
Contrastive Objective

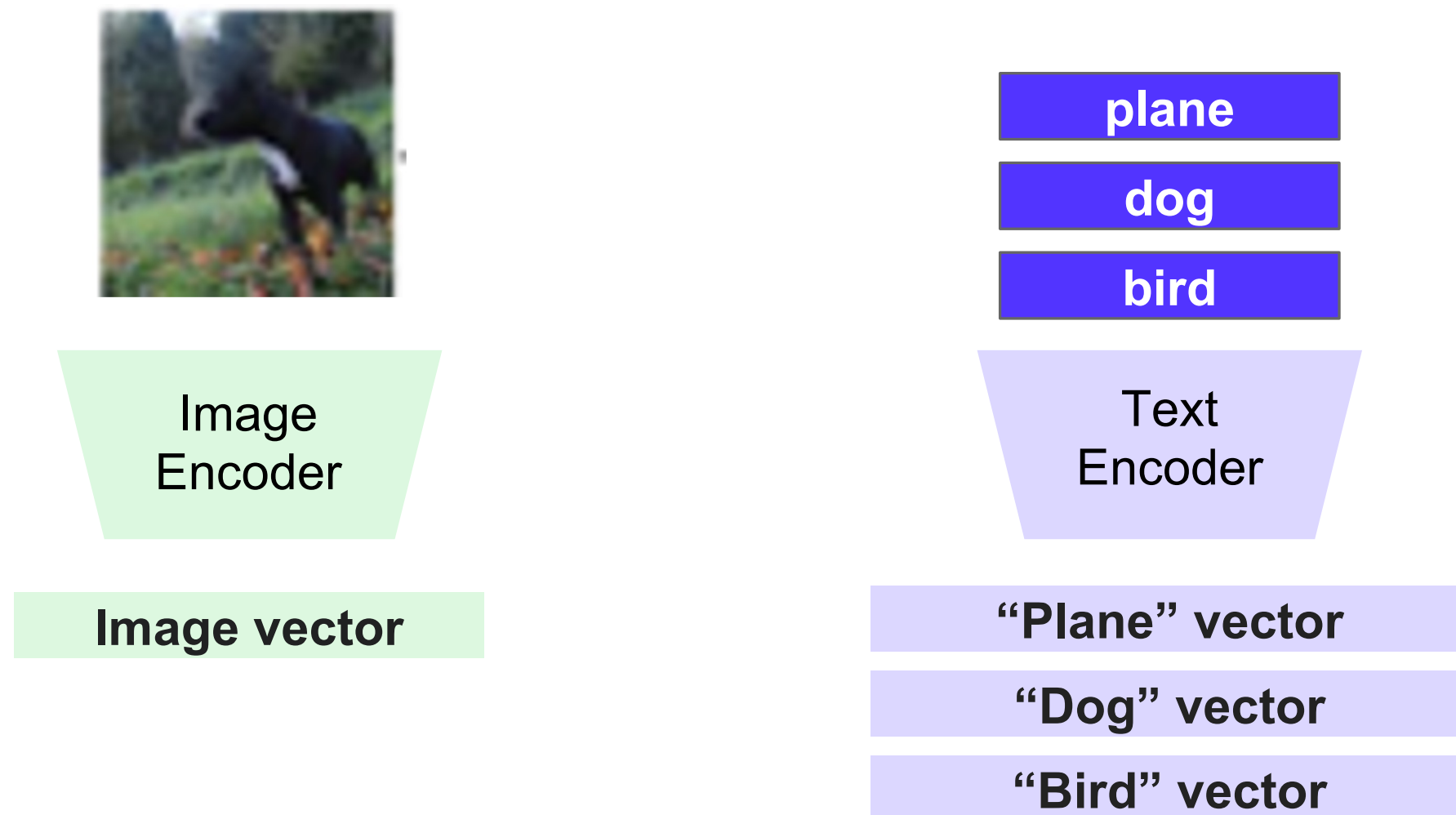**Step 2:** Use the model out of the box in a creative way!

<span style="color:purple">**Out of the box classification (No fine-tuning)**</span>

# Clever trick: we can create a classifier using the text encoder!



Image Encoder

**Image vector**

word

Text Encoder

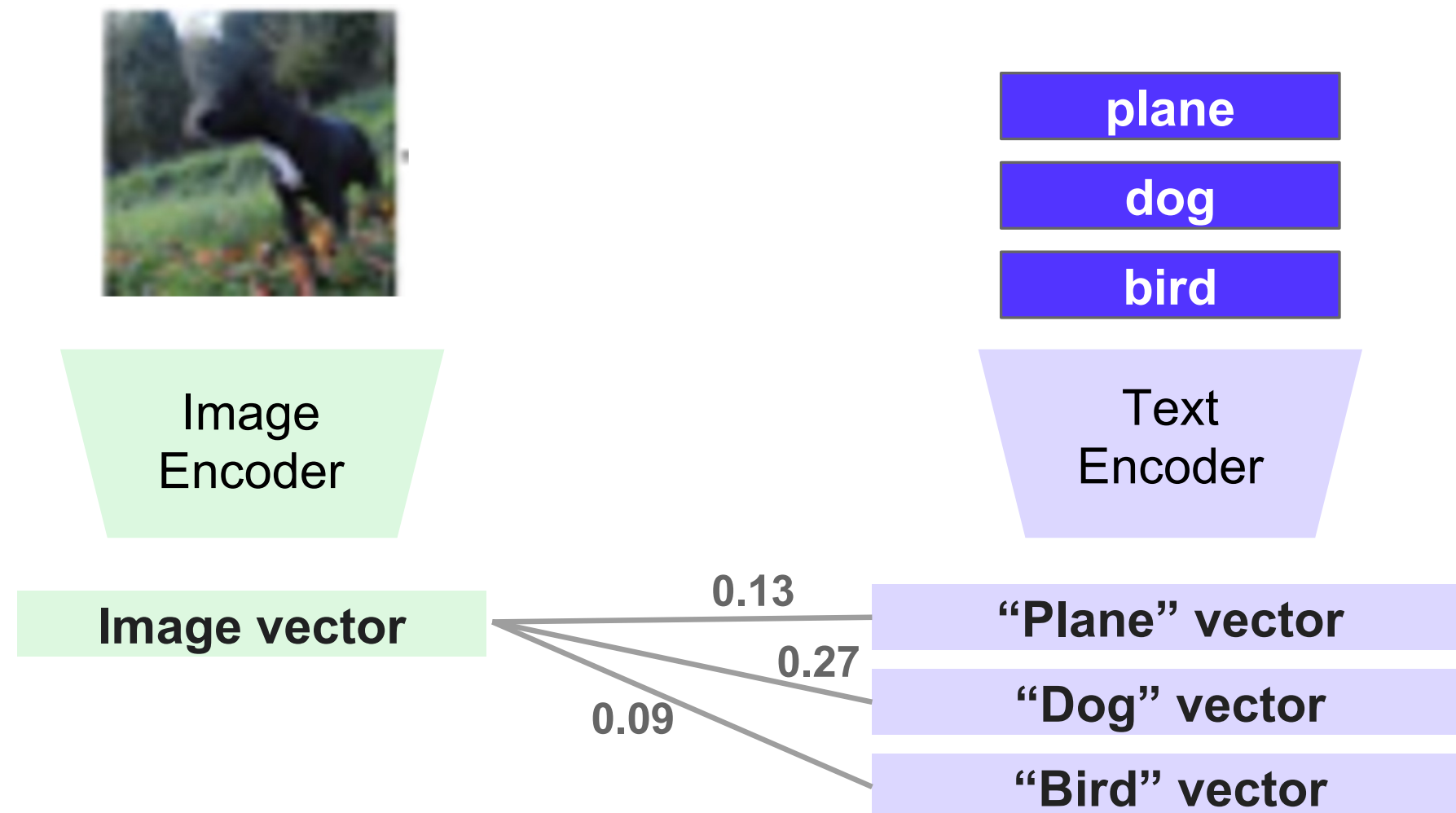**text vector**

0.27

**matching score**

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"
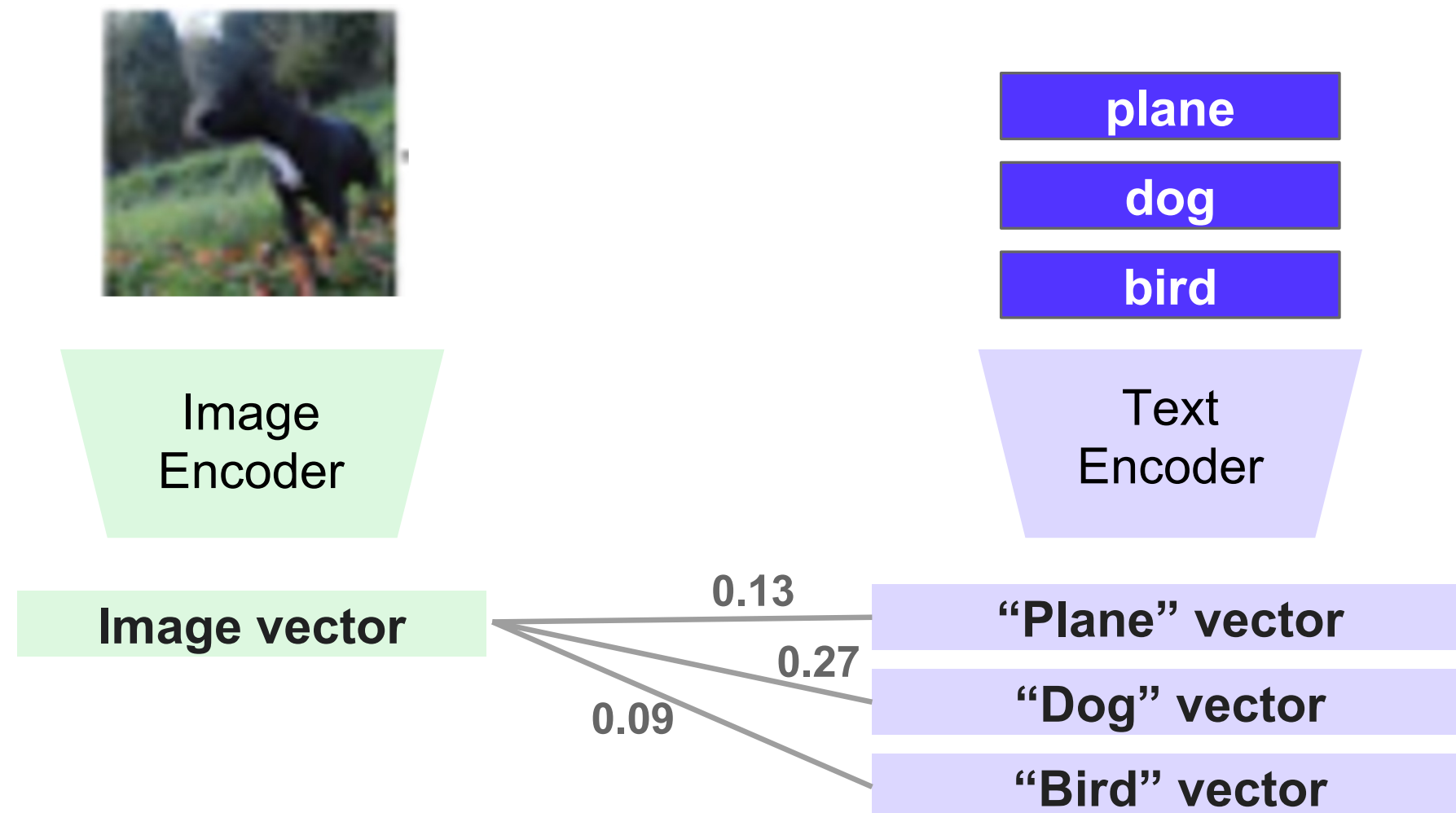
# Create a vector representation for *each* category!



Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

# Match a new image to the most similar vector



plane

dog

bird

Image Encoder

Text Encoder

**Image vector**

0.13

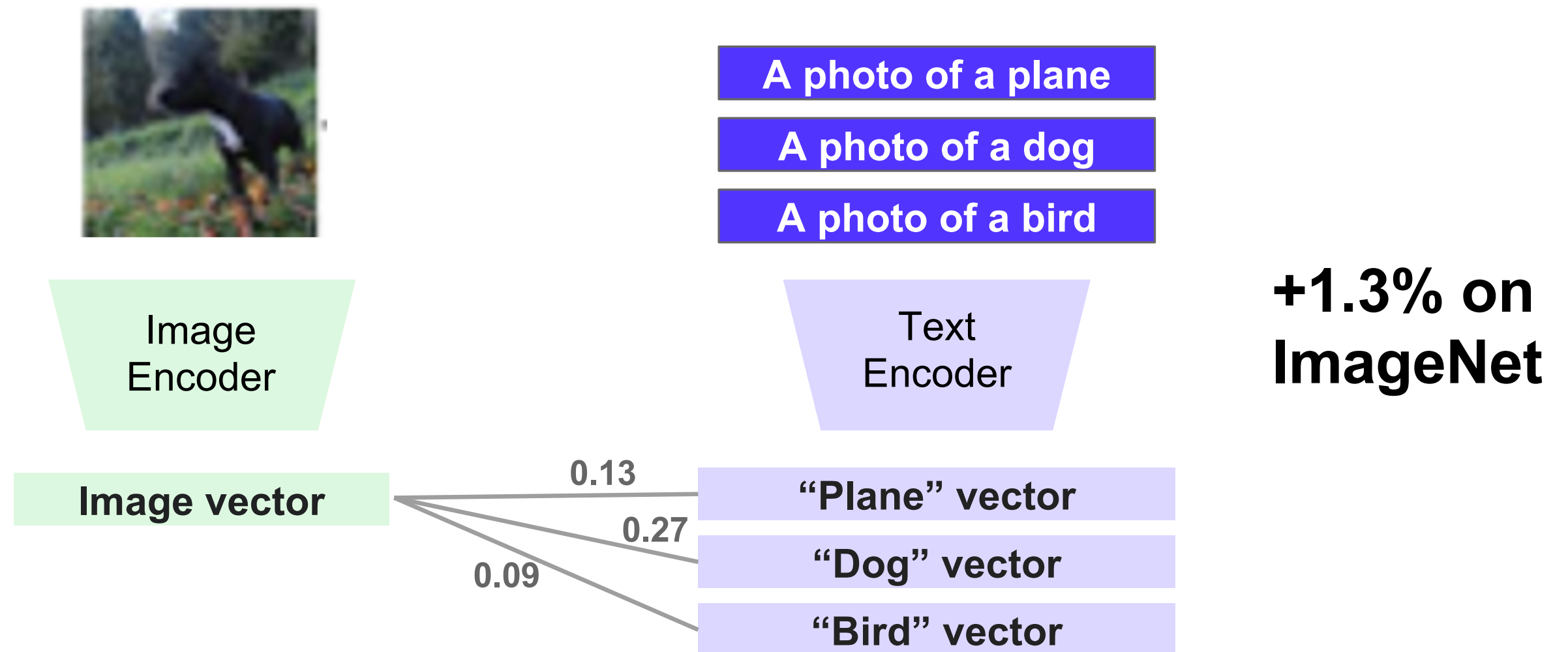0.27

0.09

"Plane" vector

"Dog" vector

"Bird" vector

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

# You can think of this as a 1-NN algorithm with the vectors as the training data



plane
dog
bird

Image Encoder

Text Encoder

Image vector

0.13 — "Plane" vector

0.27 — "Dog" vector

0.09 — "Bird" vector

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

# Since CLIP was trained with phrases, you can improve performance by using a phrase "A photo of a [category]"



**A photo of a plane**

**A photo of a dog**

**A photo of a bird**

Image Encoder

Text Encoder

**+1.3% on ImageNet**

**Image vector**

0.13    **"Plane" vector**

0.27    **"Dog" vector**

0.09    **"Bird" vector**

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

# A single phrase might be too biased.
# Solution: Use multiple phrases



| A photo of a plane | A drawing of a plane | ... |
| A photo of a dog | A drawing of a dog | ... |
| A photo of a bird | A drawing of a bird | ... |

**Image Encoder**

**Text Encoder**

**Image vector**

| "Plane" vector 1 | "Plane" vector 2 |
| "Dog" vector 1 | "Dog" vector 2 |
| "Bird" vector 1 | "Bird" vector 2 |

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

# Use the average vector across phrases as the representation for each category



A photo of a plane    A drawing of a plane    …

A photo of a dog    A drawing of a dog    …

A photo of a bird    A drawing of a bird    …

Image Encoder

Text Encoder

**+5% on ImageNet**

Image vector

0.13 — Mean "Plane" vector

0.27 — Mean "Dog" vector

0.09 — Mean "Bird" vector

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

# That's it! Now, you can use CLIP as a foundation model for image classification for any dataset



**2. Create dataset classifier from label text**

**3. Use for zero-shot prediction**

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

# Exciting result after training on 400M image-text pairs



| DATASET | IMAGENET RESNET101 | CLIP VIT-L |
|---------|--------------------|------------|
| ImageNet | 76.2% | 76.2% |

Matches the accuracy of of ResNet 101 that has been trained on ImageNet, except CLIP was trained with no human labels at all!

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

# Here's where things get even more exciting



| DATASET | IMAGENET RESNET101 | CLIP VIT-L |
|---|---|---|
| ImageNet | 76.2% | 76.2% |
| ObjectNet | | |

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

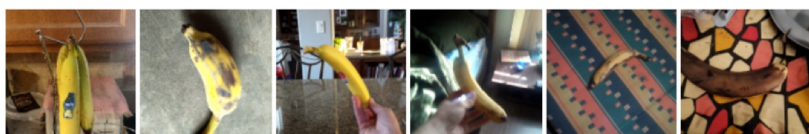# Training on ImageNet doesn't generalize to other datasets. ObjectNet contains the same categories but in weird viewpoints



| DATASET | | IMAGENET RESNET101 | CLIP VIT-L |
|---|---|---|---|
| ImageNet | | 76.2% | 76.2% |
| ObjectNet | | 32.6% | |

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

# But CLIP zero-shot does so well!
## Q. Why do you think that is?



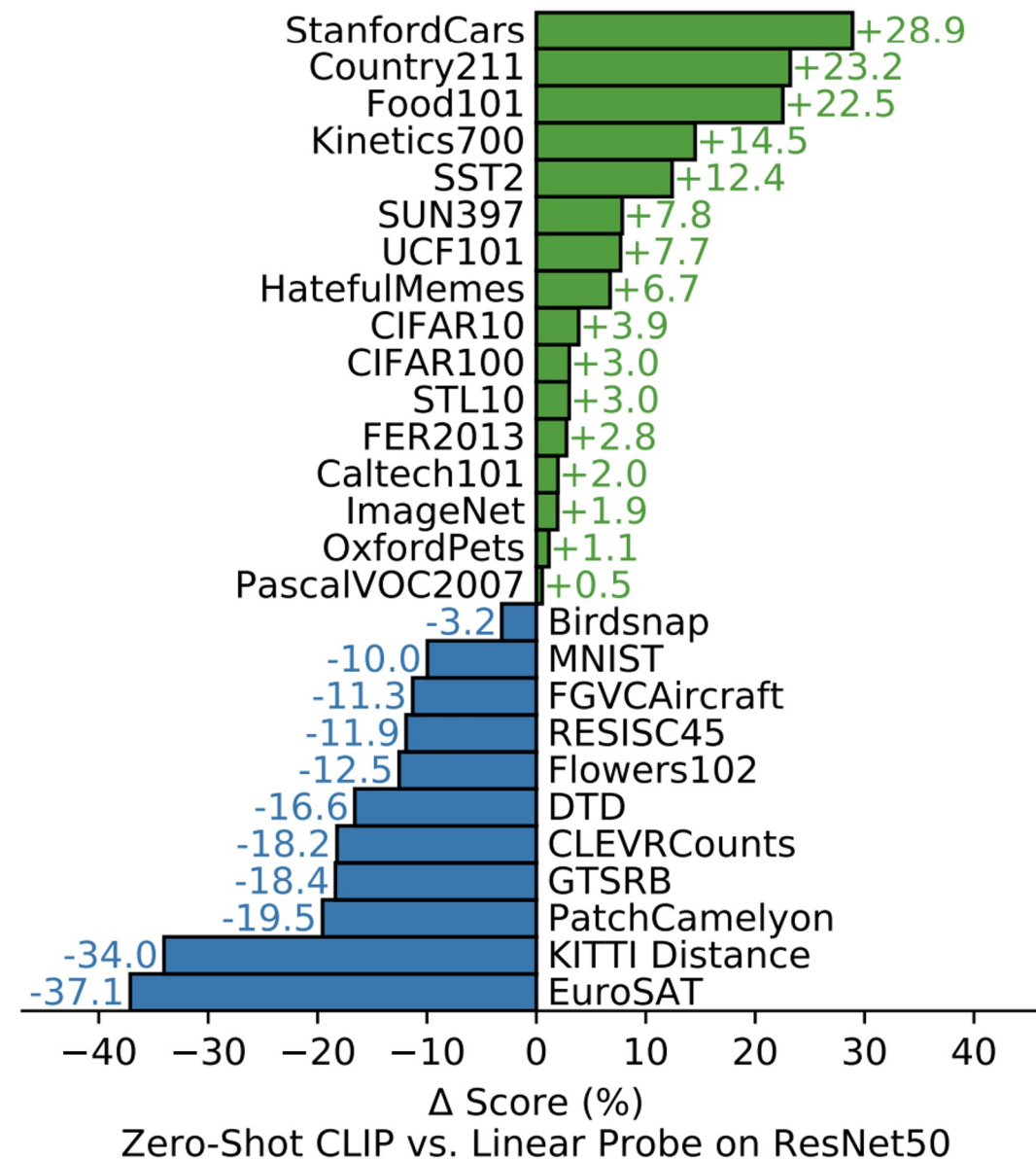| DATASET | IMAGENET RESNET101 | CLIP VIT-L |
|---------|--------------------|-----------| 
| ImageNet | 76.2% | 76.2% |
| ObjectNet | 32.6% | 72.3% |

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

CLIP performance is great also on graphic images , sketches, adversarial datasets,



| DATASET | IMAGENET RESNET101 | CLIP VIT-L |
|---|---|---|
| ImageNet | 76.2% | 76.2% |
| ImageNet V2 | 64.3% | 70.1% |
| ImageNet Rendition | 37.7% | 88.9% |
| ObjectNet | 32.6% | 72.3% |
| ImageNet Sketch | 25.2% | 60.2% |
| ImageNet Adversarial | 2.7% | 77.1% |

Radford et al "Learning Transferable Visual Models From Natural Language Supervision"

# Difference in performance between linear probe vs zero-shot



Radford et al "Learning Transferable Visual Models From Natural Language Supervision"
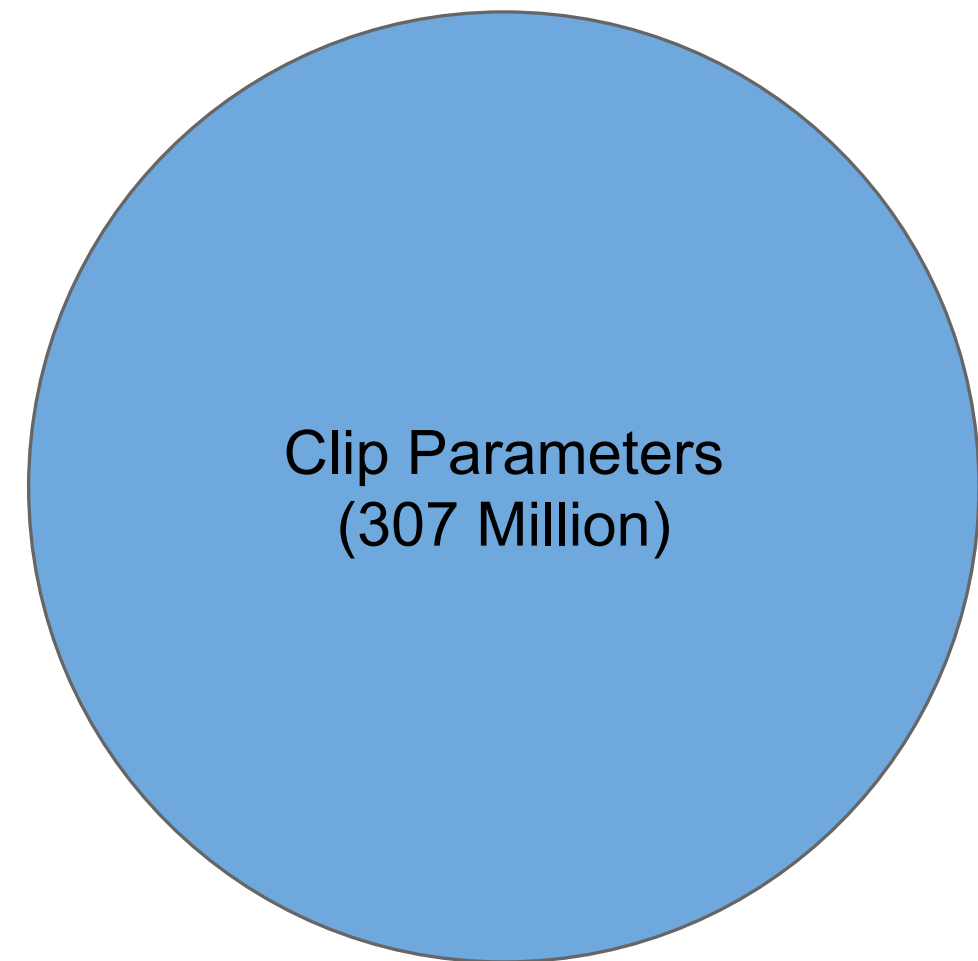
# Why does CLIP perform so well?

How can no labels beat labels??

## Scale!

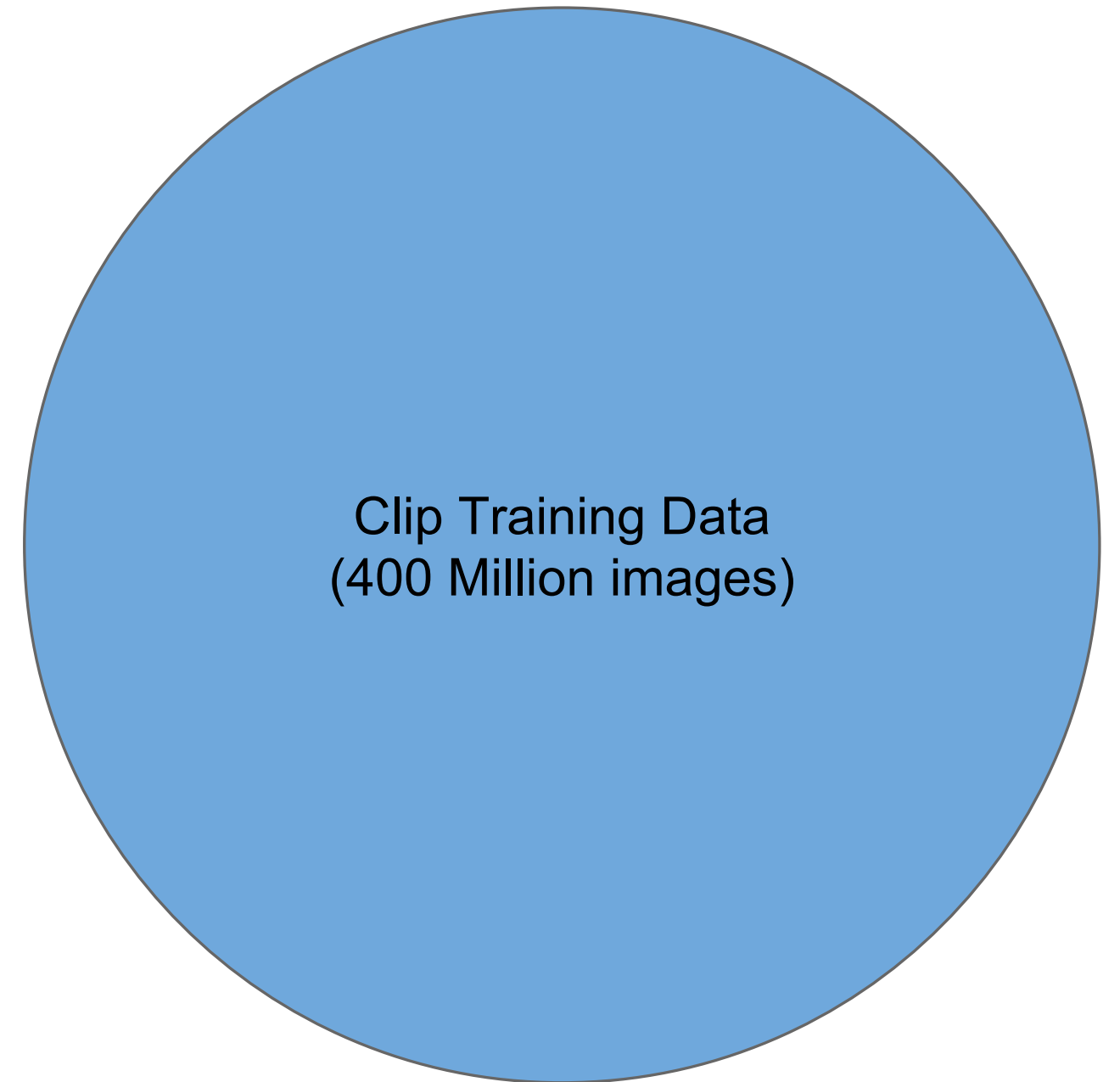# CLIP scaled up the model parameters with the transformer architecture
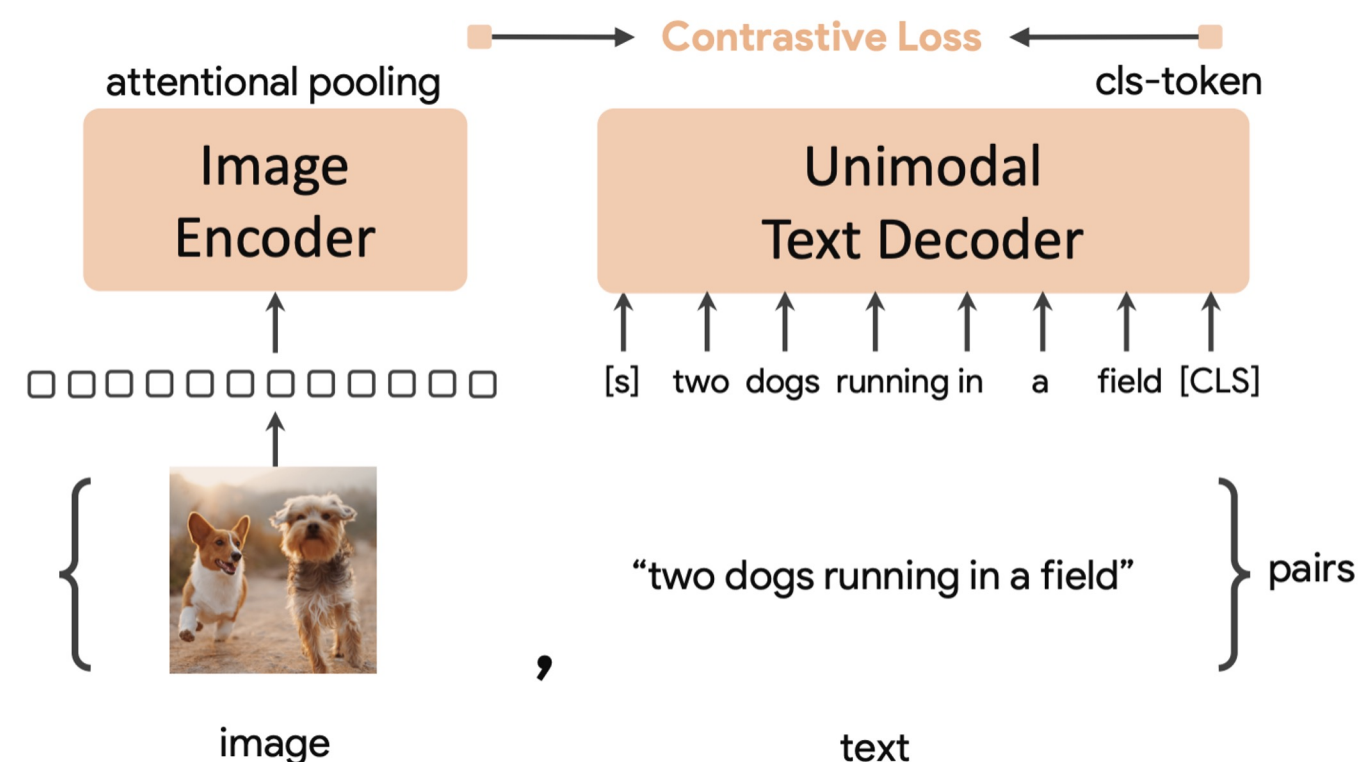


ImageNet ResNet Parameters
(44.5 Million)

Clip Parameters
(307 Million)

# CLIP Scaled up the training data by scraping image-text pairs from the internet

Clip Training Data
(400 Million images)

ImageNet ResNet Training Data
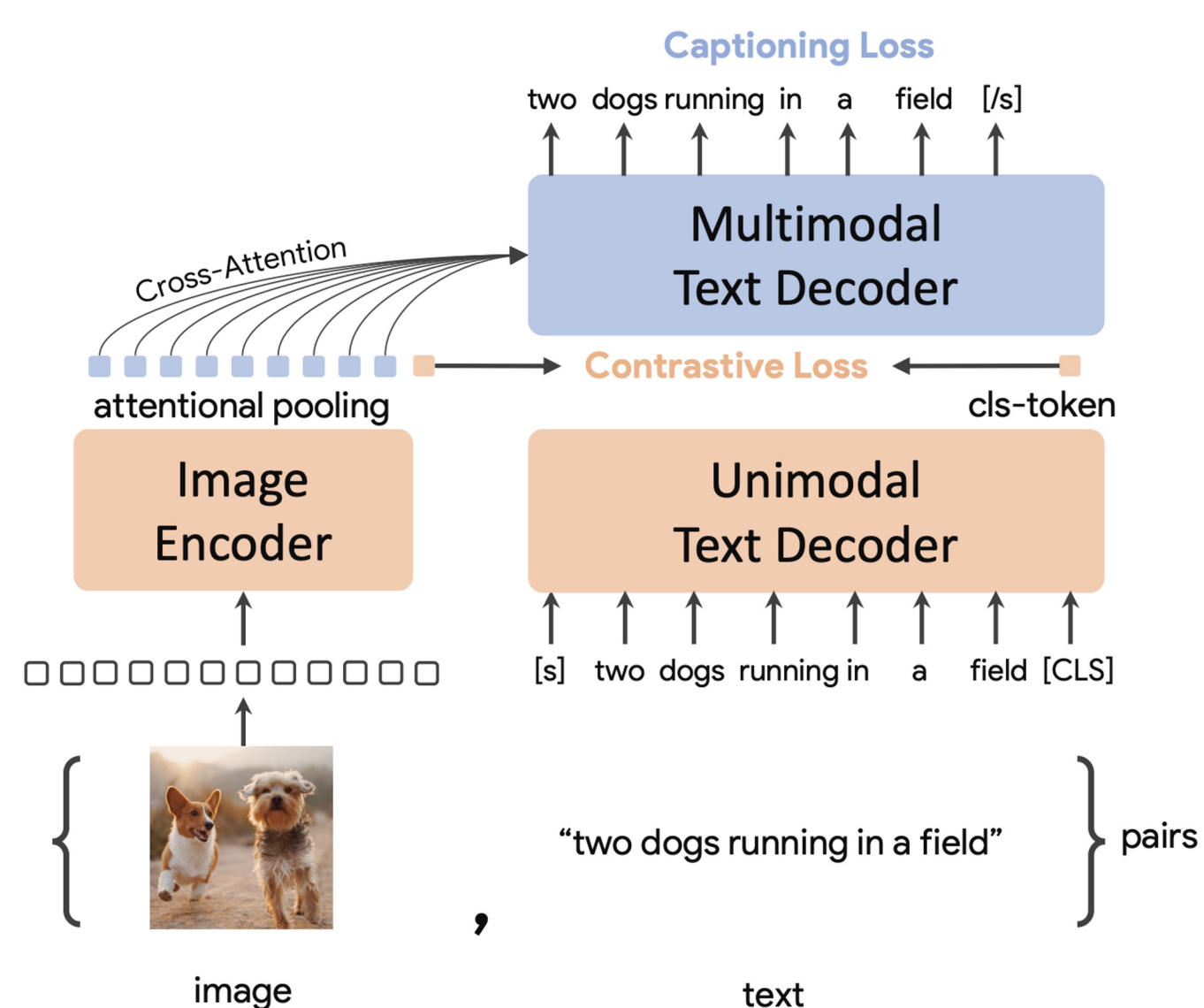(1.28 Million)

# CoCa improved upon CLIP by adding a generation objective



"Contrastive Captioners are Image-Text Foundation Models", 2022

# CoCa added a decoder with a captioning loss



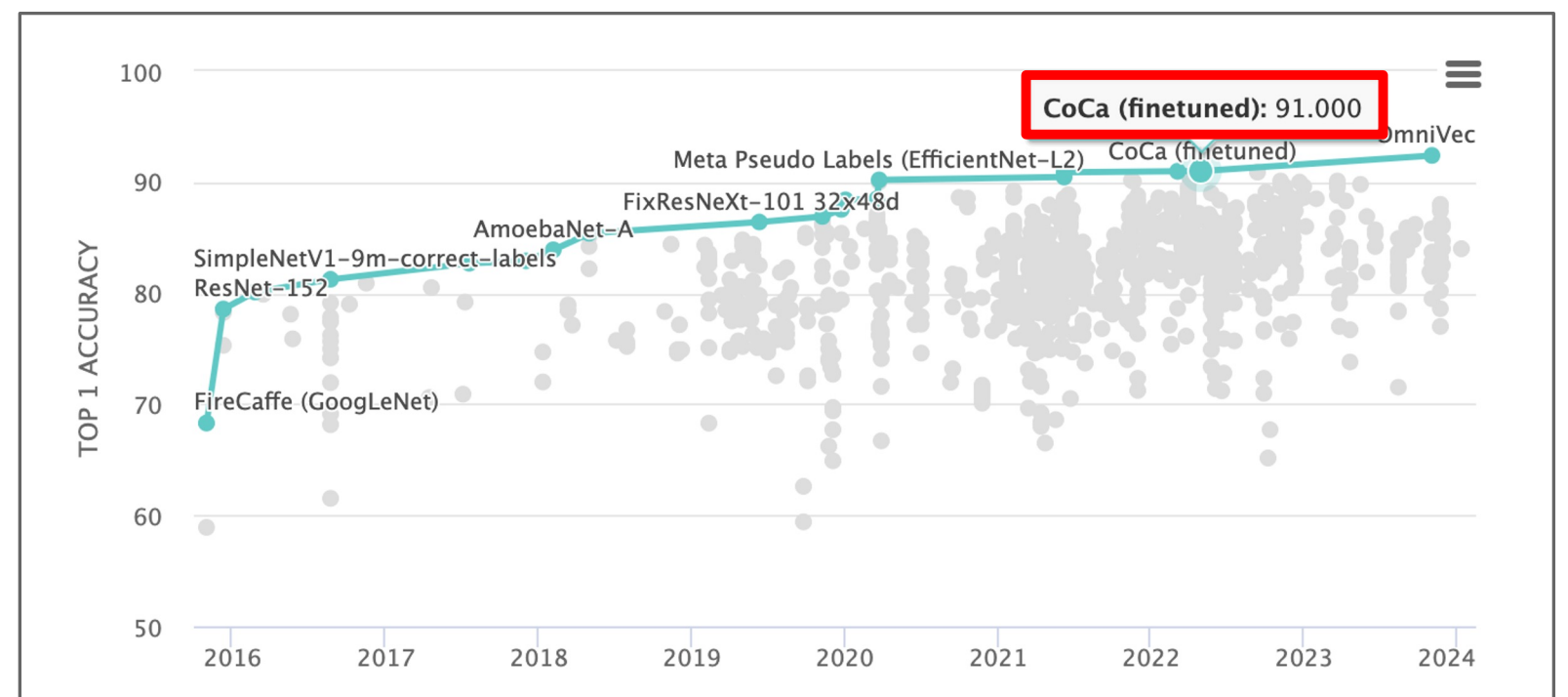"Contrastive Captioners are Image-Text Foundation Models", 2022

# CoCa: Contrastive Captioners are Image-Text Foundation Models

| Model | ImageNet | ImageNet-A | ImageNet-R | ImageNet-V2 | ImageNet-Sketch | ObjectNet | Average |
|---|---|---|---|---|---|---|---|
| CLIP [12] | 76.2 | 77.2 | 88.9 | 70.1 | 60.2 | 72.3 | 74.3 |
| ALIGN [13] | 76.4 | 75.8 | 92.2 | 70.1 | 64.8 | 72.2 | 74.5 |
| FILIP [61] | 78.3 | - | - | - | - | - | - |
| Florence [14] | 83.7 | - | - | - | - | - | - |
| LiT [32] | 84.5 | 79.4 | 93.9 | 78.7 | - | 81.1 | - |
| BASIC [33] | 85.7 | 85.6 | 95.7 | 80.6 | 76.1 | 78.9 | 83.7 |
| CoCa-Base | 82.6 | 76.4 | 93.2 | 76.5 | 71.7 | 71.6 | 78.7 |
| CoCa-Large | 84.8 | 85.7 | 95.6 | 79.6 | 75.7 | 78.6 | 83.3 |
| **CoCa** | **86.3** | **90.2** | **96.5** | **80.7** | **77.6** | **82.7** | **85.7** |

Table 4: Zero-shot image classification results on ImageNet [9], ImageNet-A [64], ImageNet-R [65], ImageNet-V2 [66], ImageNet-Sketch [67] and ObjectNet [68].

# Classifier foundation models now beat all other models on ImageNet

| Model | ImageNet |
|---|---|
| ALIGN [13] | 88.6 |
| Florence [14] | 90.1 |
| MetaPseudoLabels [51] | 90.2 |
| CoAtNet [10] | 90.9 |
| ViT-G [21] | 90.5 |
| + Model Soups [52] | 90.9 |
| CoCa (frozen) | 90.6 |
| **CoCa (finetuned)** | **91.0** |

# Advantages of CLIP-style models

1. Dot product is super efficient
   a. Easy to train (enables scaling)
   b. Fast inference, e.g., retrieval over 5B images

2. Open-vocabulary (zero-shot generalization)

3. Can be chained with other models (CuPL)
   [we will discuss this later today]

April 2022, Tristan Thrush et al:

CLIP can't distinguish between:



there is a mug in some grass



there is some grass in a mug

# Disadvantages of CLIP-style models

1. Rely too heavily on batch size to learn concepts

   Increasing batch size helps you understand fine-grained concepts



Batch size: 4              "animal"

Batch size: 100            "dog"

Batch size: **32000**      "Welsh Corgi"

# Disadvantages of CLIP-style models

1. Rely too heavily on batch size to learn concepts

   Increasing batch size helps you understand fine-grained concepts

   But there's a limit to how fine-grained you can get this way

   Even in a batch of 32K, it's unlikely you see both "a mug in some grass" and "some grass in a mug"

# Disadvantages of CLIP-style models

1. Rely too heavily on batch size to learn concepts

Winoground



there is a mug in some grass

there is some grass in a mug

"compositionality"

CREPE



✅ Crepe on a skillet. 👁?

❌ Boats on a skillet.

❌ Crepe under a skillet.

❌ Crepe on a dog.

ARO



BLIP

the grass is eating the horse  81%

the horse is eating the grass  78%

…

| Paper | Venue | Perturbation |
|---|---|---|
| Winoground | CVPR 2022 (Oral) | word order |
| VL-Checklist | EMNLP 2022 | replacements |
| When-and-Why | ICLR 2023 (Oral) | word order |
| CREPE | CVPR 2023 (Spotlight) | word order replacements negations |
| SVLC | CVPR 2023 | replacements |
| DAC | NeurIPS 2023 (spotlight) | replacements |
| What's Up | EMNLP 2023 | replacements |
| Text encoders… | EMNLP 2023 | word order |
| SugarCREPE | NeurIPS 2023 | word order replacements additions |
| COLA | NeurIPS 2023 D&B | replacements |

# Disadvantages of CLIP-style models

1. Rely too heavily on batch size to learn concepts

Solution?                                    Hard Negative Fine-Tuning

horse eating grass                grass eating horse



TODO: Get NegCLIP scores for these captions now

# Disadvantages of CLIP-style models

1. Rely too heavily on batch size to learn concepts

But training with hard negatives has its own problems…

A black cat and a brown dog
✓

A brown cat and a black dog
✗

A brown dog and a black cat
✗

"hard positives"

# Disadvantages of CLIP-style models

1. Rely too heavily on batch size to learn concepts
2. Image-level captions are insufficient supervision



"living room"          ✓

"house plants"         ✗

"couch"                ✗

# Disadvantages of CLIP-style models

1. Rely too heavily on batch size to learn concepts
2. Image-level captions are insufficient supervision

Also train on region captions
with bounding box coordinates



a man wearing
a backpack is ...

man crossing
the street

man wearing
backpack

silver backpack

# Disadvantages of CLIP-style models

1. Rely too heavily on batch size to learn concepts
2. Image-level captions are insufficient supervision
3. You can't know everything in your 5B dataset



It's extremely important to be intentional about data collection and filtering

# Foundation Models

| Language | Classification | LM + Vision | And More! | Chaining |
|----------|---------------|-------------|-----------|----------|
| **ELMo** | **CLIP** | **LLaVA** | **Segment Anything** | **LMs + CLIP** |
| **BERT** | **CoCa** | **Flamingo** | Whisper | **Visual Programming** |
| **GPT** | | GPT-4V | Dalle | |
| **T5** | | Gemini | Stable Diffusion | |
| | | **Molmo** | Imagen | |

# LLaVA

Motivation: Language models which do next token prediction can be applied to a wide variety of tasks at inference (Math, sentiment analysis, symbolic reasoning)

**Can we build a model that can accept images and text as input, and then output text?**

$\rightarrow$ **Vision-Language Models**

# First, some historical context

Vision-Language Models didn't start with LLaVA!

They go as far back as 2019 → ViLBERT

# Historical context

Vision-Language Models didn't start with LLaVA!

They go as far back as 2019 → ViLBERT

BUT, they had to finetune for each task separately,
with non-trivial task-specific methods (e.g., Mask-RCNN
bounding box re-ranking for RefCOCO)

→ Same paradigm as we discussed right at the beginning of this lecture:
very task-specific

# LLaVA uses the autoregressive nature of LLMs

# Recall how transformers decode language

# Key idea behind LLaVA – add visual information to the LLM



Cute

Transformer block

Transformer block

Transformer block

$v_1$ $v_2$ … $v_N$ Cats are so

Which image tokens work best here?

Input image tokens

Input Text

# The CLIP encoder is a good option!



1. Contrastive pre-training

At the end of training, you have a model that will give you a similarity score between an image and a text

# What features should we use from CLIP?

Extract Image Features from CLS token for contrastive learning

Pooling token    CLS

Patchify

Flatten + Linear proj + 2D pos embed

Vision Transformer

CLS

...

[Image source]

# What features should we use from CLIP?



Pooling token

CLS

Patchify

Flatten + Linear proj + 2D pos embed

Vision Transformer

CLS

But these tokens are not supervised! (Could be random and loss will not change)

# Use Penultimate Layer!



Vision Transformer

Pooling token    CLS                                    CLS              CLS

Patchify

Flatten + Linear proj +
2D pos embed

L – 1 Layers                                 Final Layer

In practice, these tokens preserve spatial and linguistic
information best for LLMs. Can drop CLS for slight gains.

# LLaVA – Overall Architecture + Training Recipe

1. Initialize with pretrained Language Model for LLM Decoder (e.g. LLaMA) and pretrained image encoder (e.g. CLIP)
2. Train a new **linear layer** to bridge CLIP features to LLM input space
3. Finetune LLM + linear layer together



Can get reasonable performance with **>100,000** samples containing an input image, input instruction, and output text.

# Flamingo followed up with a new way to fuse visual features

# Pre-trained parts of Flamingo



Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# There are 2 learned parts in Flamingo



Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Perceiver sampler converts variable sized image tokens to fixed sized ones



Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Flamingo full architecture



Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Flamingo full architecture

Learned method of down-sampling image/video representations



Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Flamingo full architecture



Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Flamingo gated cross-attention



Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Flamingo gated cross-attention



```python
def gated_xattn_dense(
    y,    # input language features
    x,    # input visual features
    alpha_xattn, # xattn gating parameter – init at 0.
    alpha_dense, # ffw gating parameter – init at 0.
):
  """Applies a GATED XATTN-DENSE layer."""

  # 1. Gated Cross Attention
  y = y + tanh(alpha_xattn) * attention(q=y, kv=x)
  # 2. Gated Feed Forward (dense) Layer
  y = y + tanh(alpha_dense) * ffw(y)

  # Regular self-attention + FFW on language
  y = y + frozen_attention(q=y, kv=y)
  y = y + frozen_ffw(y)

  return y  # output visually informed language features
```

Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Flamingo arranges its training data similar to language modeling,

with special tags <image>, <eos> to indicate when a new image shows up or the text ends.



`<BOS>Cute pics of my pets!<EOC><image>My puppy sitting in the grass <EOC><image> My cat looking very dignified.<EOC>`

**Processed text:** <image> tags are inserted and special tokens are added

Image 1          Image 2

# Flamingo masked attention



$\phi$   0   0   0   0   0   0   0   0   1   1   1 1   1    1    1 1    1   1    1    2   2   2 2   2    2     2     2    2 2   2

$Y$ `<BOS>` Cute pics of my pets!`<EOC><image>`My puppy sitting in the grass. `<EOC><image>`My cat looking very dignified.`<EOC`

tokenization

`<BOS>`Cute pics of my pets!`<EOC><image>`My puppy sitting in the grass.`<EOC><image>` My cat looking very dignified.`<EOC>`

**Processed text:** `<image>` tags are inserted and special tokens are added

Image 1      Image 2

Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Flamingo masked attention



Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Flamingo full architecture



Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Flamingo results



Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Flaming enables in-context learning



Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.

# Flamingo results



| | Output |
|---|---|
| Output: A propaganda poster depicting a cat dressed as French emperor Napoleon holding a piece of cheese. | Output: A pink room with a flamingo pool float. | Output: → A portrait of Salvador Dali with a robot head. |
| Les sanglots longs des violons de l'automne blessent mon coeur d'une langueur monotone. | Pour qui sont ces serpents qui sifflent sur vos têtes? | → Je suis un cœur qui bat pour vous. |
| pandas: 3 | dogs: 2 | → giraffes: 4 |
| I like reading | , my favourite play is Hamlet. I also like | , my favorite book is → Dreams from my Father. |
| | What happens to the man after hitting the ball? Answer: | → he falls down. |

Alayrac et al "Flamingo: a Visual Language Model for Few-Shot Learning. 2022.
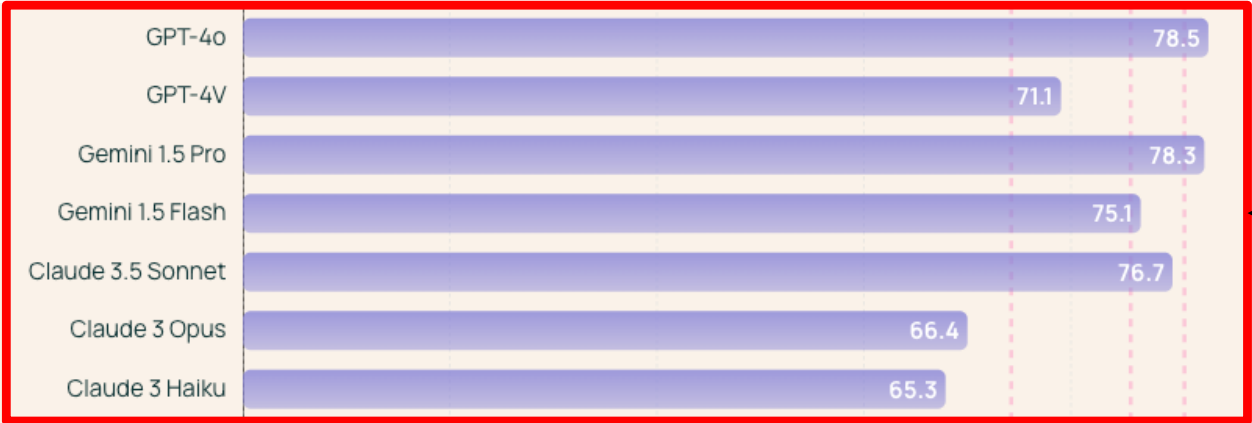
# Results: zero & few shot

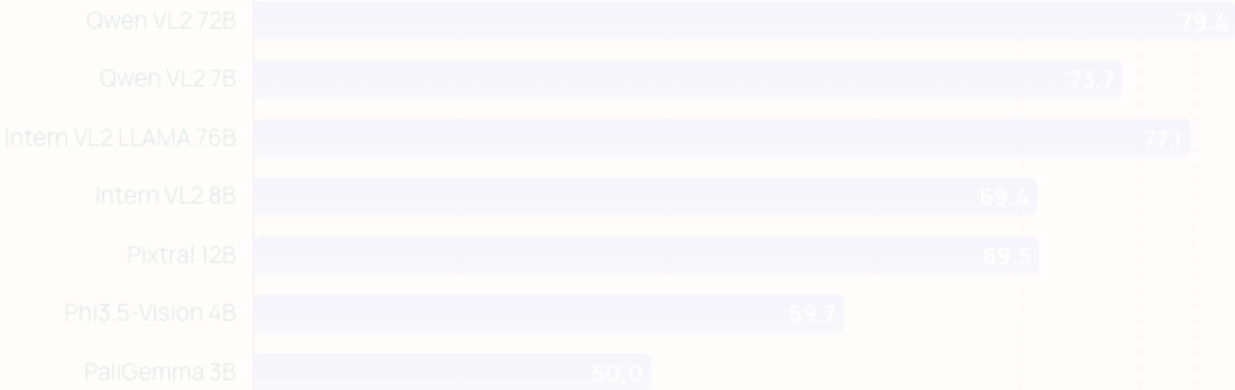| Method | FT | Shot | OKVQA | VQAv2 | COCO | MSVDQA | VATEX | VizWiz | Flick30K | MSRVTTQA | iVQA | YouCook2 | STAR | VisDial | TextVQA | NextQA | HatefulMemes | RareAct |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Zero/Few shot SOTA | ✗ | (X) | [39] 43.3 (16) | [124] 38.2 (4) | [134] 32.2 (0) | [64] 35.2 (0) | - | - | - | [64] 19.2 (0) | [145] 12.2 (0) | - | [153] 39.4 (0) | [87] 11.6 (0) | - | - | [94] 66.1 (0) | [94] 40.7 (0) |
| Flamingo-3B | ✗ | 0 | 41.2 | 49.2 | 73.0 | 27.5 | 40.1 | 28.9 | 60.6 | 11.0 | 32.7 | 55.8 | 39.6 | 46.1 | 30.1 | 21.3 | 53.7 | 58.4 |
| | ✗ | 4 | 43.3 | 53.2 | 85.0 | 33.0 | 50.0 | 34.0 | 72.0 | 14.9 | 35.7 | 64.6 | 41.3 | 47.3 | 32.7 | 22.4 | 53.6 | - |
| | ✗ | 8 | 44.6 | 55.4 | 90.6 | 37.0 | 54.5 | 38.4 | 71.7 | 19.6 | 36.8 | 68.0 | 40.6 | 47.6 | 32.4 | 23.9 | 54.7 | - |
| | ✗ | 16 | 45.6 | 56.7 | 95.4 | 40.2 | 57.1 | 43.3 | 73.4 | 23.4 | 37.4 | 73.2 | 40.1 | 47.5 | 31.8 | 25.2 | 55.3 | - |
| | ✗ | 32 | 45.9 | 57.1 | 99.0 | 42.6 | 59.2 | 45.5 | 71.2 | 25.6 | 37.7 | 76.7 | 41.6 | OOC | 30.6 | 26.1 | 56.3 | - |
| Flamingo-9B | ✗ | 0 | 44.7 | 51.8 | 79.4 | 30.2 | 39.5 | 28.8 | 61.5 | 13.7 | 35.2 | 55.0 | 41.8 | 48.0 | 31.8 | 23.0 | 57.0 | 57.9 |
| | ✗ | 4 | 49.3 | 56.3 | 93.1 | 36.2 | 51.7 | 34.9 | 72.6 | 18.2 | 37.7 | 70.8 | 42.8 | 50.4 | 33.6 | 24.7 | 62.7 | - |
| | ✗ | 8 | 50.0 | 58.0 | 99.0 | 40.8 | 55.2 | 39.4 | 73.4 | 23.9 | 40.0 | 75.0 | __43.4__ | 51.2 | 33.6 | 25.8 | 63.9 | - |
| | ✗ | 16 | 50.8 | 59.4 | 102.2 | 44.5 | 58.5 | 43.0 | 72.7 | 27.6 | 41.5 | 77.2 | 42.4 | 51.3 | 33.5 | 27.6 | 64.5 | - |
| | ✗ | 32 | 51.0 | 60.4 | 106.3 | 47.2 | 57.4 | 44.0 | 72.8 | 29.4 | 40.7 | 77.3 | 41.2 | OOC | 32.6 | 28.4 | 63.5 | - |
| Flamingo | ✗ | 0 | 50.6 | 56.3 | 84.3 | 35.6 | 46.7 | 31.6 | 67.2 | 17.4 | 40.7 | 60.1 | 39.7 | 52.0 | 35.0 | 26.7 | 46.4 | __60.8__ |
| | ✗ | 4 | 57.4 | 63.1 | 103.2 | 41.7 | 56.0 | 39.6 | 75.1 | 23.9 | 44.1 | 74.5 | 42.4 | 55.6 | 36.5 | 30.8 | 68.6 | - |
| | ✗ | 8 | 57.5 | 65.6 | 108.8 | 45.5 | 60.6 | 44.8 | 78.2 | 27.6 | 44.8 | 80.7 | 42.3 | 56.4 | 37.3 | 32.3 | __70.0__ | - |
| | ✗ | 16 | 57.8 | 66.8 | 110.5 | 48.4 | 62.8 | 48.4 | __78.9__ | 30.0 | 45.2 | 84.2 | 41.1 | __56.8__ | 37.6 | 32.9 | __70.0__ | - |
| | ✗ | 32 | __57.8__ | __67.6__ | __113.8__ | __52.3__ | __65.1__ | __49.8__ | 75.4 | __31.0__ | __45.3__ | __86.8__ | 42.2 | OOC | __37.9__ | __33.5__ | __70.0__ | - |
| Pretrained FT SOTA | ✓ | (X) | 54.4 [39] (10K) | 80.2 [150] (444K) | 143.3 [134] (500K) | 47.9 [32] (27K) | 76.3 [165] (500K) | 57.2 [70] (20K) | 67.4 [162] (30K) | 46.8 [57] (130K) | 35.4 [145] (6K) | 138.7 [142] (10K) | 36.7 [138] (46K) | 75.2 [87] (123K) | 54.7 [147] (20K) | 25.2 [139] (38K) | 75.4 [60] (9K) | - |

# Results: zero & few shot

| Method | FT | Shot | OKVQA | VQAv2 | COCO | MSVDQA | VATEX | VizWiz | Flick30K | MSRVTTQA | iVQA | YouCook2 | STAR | VisDial | TextVQA | NextQA | HatefulMemes | RareAct |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Zero/Few shot SOTA | ✗ | (X) | [39]<br>43.3<br>(16) | [124]<br>38.2<br>(4) | [134]<br>32.2<br>(0) | [64]<br>35.2<br>(0) | - | - | - | [64]<br>19.2<br>(0) | [145]<br>12.2<br>(0) | - | [153]<br>39.6<br>(0) | [87]<br>11.6<br>(0) | - | - | [94]<br>66.1<br>(0) | [94]<br>40.7<br>(0) |
| Flamingo-3B | ✗ | 0 | 41.2 | 49.2 | 73.0 | 27.5 | 40.1 | 28.9 | 60.6 | 11.0 | 32.7 | 55.8 | 39.6 | 46.1 | 30.1 | 21.3 | 53.7 | 58.4 |
| | ✗ | 4 | 43.3 | 53.2 | 85.0 | 33.0 | 50.0 | 34.0 | 72.0 | 14.9 | 35.7 | 64.6 | 41.3 | 47.3 | 32.7 | 22.4 | 53.6 | - |
| | ✗ | 8 | 44.6 | 55.4 | 90.6 | 37.0 | 54.5 | 38.4 | 71.7 | 19.6 | 36.8 | 68.0 | 40.6 | 47.6 | 32.4 | 23.9 | 54.7 | - |
| | ✗ | 16 | 45.6 | 56.7 | 95.4 | 40.2 | 57.1 | 43.3 | 73.4 | 23.4 | 37.4 | 73.2 | 40.1 | 47.5 | 31.8 | 25.2 | 55.3 | - |
| | ✗ | 32 | 45.9 | 57.1 | 99.0 | 42.6 | 59.2 | 45.5 | 71.2 | 25.6 | 37.7 | 76.7 | 41.6 | OOC | 30.6 | 26.1 | 56.3 | - |
| Flamingo-9B | ✗ | 0 | 44.7 | 51.8 | 79.4 | 30.2 | 39.5 | 28.8 | 61.5 | 13.7 | 35.2 | 55.0 | 41.8 | 48.0 | 31.8 | 23.0 | 57.0 | 57.9 |
| | ✗ | 4 | 49.3 | 56.3 | 93.1 | 36.2 | 51.7 | 34.9 | 72.6 | 18.2 | 37.7 | 70.8 | 42.8 | 50.4 | 33.6 | 24.7 | 62.7 | - |
| | ✗ | 8 | 50.0 | 58.0 | 99.0 | 40.8 | 55.2 | 39.4 | 73.4 | 23.9 | 40.0 | 75.0 | 43.4 | 51.2 | 33.6 | 25.8 | 63.9 | - |
| | ✗ | 16 | 50.8 | 59.4 | 102.2 | 44.5 | 58.5 | 43.0 | 72.7 | 27.6 | 41.5 | 77.2 | 42.4 | 51.3 | 33.5 | 27.6 | 64.5 | - |
| | ✗ | 32 | 51.0 | 60.4 | 106.3 | 47.2 | 57.4 | 44.0 | 72.8 | 29.4 | 40.7 | 77.3 | 41.2 | OOC | 32.6 | 28.4 | 63.5 | - |
| Flamingo | ✗ | 0 | 50.6 | 56.3 | 84.3 | 35.6 | 46.7 | 31.6 | 67.2 | 17.4 | 40.7 | 60.1 | 39.7 | 52.0 | 35.0 | 26.7 | 46.4 | 60.8 |
| | ✗ | 4 | 57.4 | 63.1 | 103.2 | 41.7 | 56.0 | 39.6 | 75.1 | 23.9 | 44.1 | 74.5 | 42.4 | 55.6 | 36.5 | 30.8 | 68.6 | - |
| | ✗ | 8 | 57.5 | 65.6 | 108.8 | 45.5 | 60.6 | 44.8 | 78.2 | 27.6 | 44.8 | 80.7 | 42.3 | 56.4 | 37.3 | 32.3 | 70.0 | - |
| | ✗ | 16 | 57.8 | 66.8 | 110.5 | 48.4 | 62.8 | 48.4 | 78.9 | 30.0 | 45.2 | 84.2 | 41.1 | 56.8 | 37.6 | 32.9 | 70.0 | - |
| | ✗ | 32 | 57.8 | 67.6 | 113.8 | 52.3 | 65.1 | 49.8 | 75.4 | 31.0 | 45.3 | 86.8 | 42.2 | OOC | 37.9 | 33.5 | 70.0 | - |
| Pretrained FT SOTA | ✔ | (X) | 54.4<br>[39]<br>(10K) | 80.2<br>[150]<br>(444K) | 143.3<br>[134]<br>(500K) | 47.9<br>[32]<br>(27K) | 76.3<br>[165]<br>(500K) | 57.2<br>[70]<br>(20K) | 67.4<br>[162]<br>(30K) | 46.8<br>[57]<br>(130K) | 35.4<br>[145]<br>(6K) | 138.7<br>[142]<br>(10K) | 36.7<br>[138]<br>(46K) | 75.2<br>[87]<br>(123K) | 54.7<br>[147]<br>(20K) | 25.2<br>[139]<br>(38K) | 75.4<br>[60]<br>(9K) | - |

# Today, average performance across
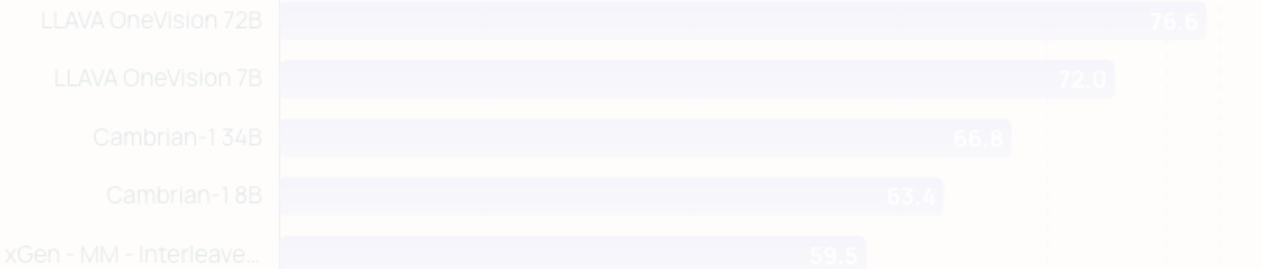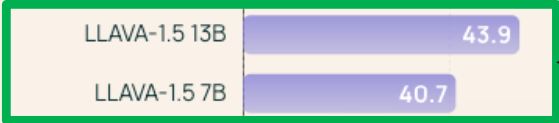## 11 visual understanding benchmarks



API Only

| Model | Score |
|---|---|
| GPT-4o | 78.5 |
| GPT-4V | 71.1 |
| Gemini 1.5 Pro | 78.3 |
| Gemini 1.5 Flash | 75.1 |
| Claude 3.5 Sonnet | 76.7 |
| Claude 3 Opus | 66.4 |
| Claude 3 Haiku | 65.3 |

Open Weights

| Model | Score |
|---|---|
| Qwen VL2 72B | 79.4 |
| Qwen VL2 7B | 73.7 |
| Intern VL2 LLAMA 76B | 77.1 |
| Intern VL2 8B | 65.4 |
| Pixtral 12B | 69.5 |
| Phi3.5-Vision 4B | 59.7 |
| PaliGemma 3B | 56.0 |

Distilled

| Model | Score |
|---|---|
| LLAVA OneVision 72B | 76.6 |
| LLAVA OneVision 7B | 72.0 |
| Cambrian-1 34B | 66.8 |
| Cambrian-1 8B | 63.4 |
| xGen - MM - Interleave... | 59.5 |

Open

| Model | Score |
|---|---|
| LLAVA-1.5 13B | 43.9 |
| LLAVA-1.5 7B | 40.7 |

# There are open-weight models
## but they are all distilled from GPT



| Category | Model | Score |
|---|---|---|
| API Only | GPT-4o | 78.5 |
| | GPT-4V | 71.1 |
| | Gemini 1.5 Pro | 78.3 |
| | Gemini 1.5 Flash | 75.1 |
| | Claude 3.5 Sonnet | 76.7 |
| | Claude 3 Opus | 66.4 |
| | Claude 3 Haiku | 65.3 |
| Open Weights | Qwen VL2 72B | 79.4 |
| | Qwen VL2 7B | 73.7 |
| | Intern VL2 LLAMA 76B | 77.1 |
| | Intern VL2 8B | 69.4 |
| | Pixtral 12B | 69.5 |
| | Phi3.5-Vision 4B | 59.7 |
| | PaliGemma 3B | 50.0 |
| Distilled | LLAVA OneVision 72B | 76.6 |
| | LLAVA OneVision 7B | 72.0 |
| | Cambrian-1 34B | 66.8 |
| | Cambrian-1 8B | 63.4 |
| | xGen - MM - Interleave... | 59.5 |
| Open | LLAVA-1.5 13B | 43.9 |
| | LLAVA-1.5 7B | 40.7 |

# How do we close the gap without relying on proprietary models?

# There are open-weight models
## but they are all distilled from GPT



**API Only**
- GPT-4o — 78.5
- GPT-4V — 71.1
- Gemini 1.5 Pro — 78.3
- Gemini 1.5 Flash — 75.1
- Claude 3.5 Sonnet — 76.7
- Claude 3 Opus — 66.4
- Claude 3 Haiku — 65.3

**Open Weights**
- Qwen VL2 72B — 79.4
- Qwen VL2 7B — 73.7
- Intern VL2 LLAMA 76B — 77.1
- Intern VL2 8B — 69.4
- Pixtral 12B — 69.5
- Phi3.5-Vision 4B — 59.7
- PaliGemma 3B — 50.0

**Distilled**
- LLAVA OneVision 72B — 76.6
- LLAVA OneVision 7B — 72.0
- Cambrian-1 34B — 66.8
- Cambrian-1 8B — 63.4
- xGen - MM - Interleave... — 59.5

**Open**
- LLAVA-1.5 13B — 43.9
- LLAVA-1.5 7B — 40.7

**Average Score on 11 Academic Benchmarks**

Open
*Weights*
*Data*
*Code*
*Evals*

| Model | Score |
|---|---|
| Molmo 72B | 81.2 |
| Molmo 7B-D | 77.3 |
| Molmo 7B-O | 74.6 |
| MolmoE 1B | 68.6 |

API Only

| Model | Score |
|---|---|
| GPT-4o | 78.5 |
| GPT-4V | 71.1 |
| Gemini 1.5 Pro | 78.3 |
| Gemini 1.5 Flash | 75.1 |
| Claude 3.5 Sonnet | 76.7 |
| Claude 3 Opus | 66.4 |
| Claude 3 Haiku | 65.3 |

Open
Weights

| Model | Score |
|---|---|
| Qwen VL2 72B | 79.4 |
| Qwen VL2 7B | 73.7 |
| Intern VL2 LLAMA 76B | 77.1 |
| Intern VL2 8B | 69.4 |
| Pixtral 12B | 69.5 |
| Phi3.5-Vision 4B | 59.7 |
| PaliGemma 3B | 50.0 |

Distilled

| Model | Score |
|---|---|
| LLAVA OneVision 72B | 76.6 |
| LLAVA OneVision 7B | 72.0 |
| Cambrian-1 34B | 66.8 |
| Cambrian-1 8B | 63.4 |
| xGen - MM - Interleave... | 59.5 |
| LLAVA-1.5 13B | 43.9 |
| LLAVA-1.5 7B | 40.7 |

Open

30    82

Completely
**Open**
Open Weights
Open Data
Open Code
Open Evals

Ranj...                          ...ture 16 -                    May 27, 2025

## Average Score on 11 Academic Benchmarks | Human Preference Elo Rating

**Open** *Weights Data Code Evals*

| Model | Average Score | Elo Rating |
|---|---|---|
| Molmo 72B | 81.2 | 1077 |
| Molmo 7B-D | 77.3 | 1056 |
| Molmo 7B-O | 74.6 | 1051 |
| MolmoE 1B | 68.6 | 1032 |

**API Only**

| Model | Average Score | Elo Rating |
|---|---|---|
| GPT-4o | 78.5 | 1079 |
| GPT-4V | 71.1 | 1041 |
| Gemini 1.5 Pro | 78.3 | 1074 |
| Gemini 1.5 Flash | 75.1 | 1054 |
| Claude 3.5 Sonnet | 76.7 | 1069 |
| Claude 3 Opus | 66.4 | 971 |
| Claude 3 Haiku | 65.3 | 999 |

**Open Weights**

| Model | Average Score | Elo Rating |
|---|---|---|
| Qwen VL2 72B | 79.4 | 1037 |
| Qwen VL2 7B | 73.7 | 1025 |
| Intern VL2 LLAMA 76B | 77.1 | 1018 |
| Intern VL2 8B | 69.4 | 953 |
| Pixtral 12B | 69.5 | 1016 |
| Phi3.5-Vision 4B | 59.7 | 982 |
| PaliGemma 3B | 50.0 | 937 |

**Distilled**

| Model | Average Score | Elo Rating |
|---|---|---|
| LLAVA OneVision 72B | 76.6 | 1051 |
| LLAVA OneVision 7B | 72.0 | 1024 |
| Cambrian-1 34B | 66.8 | 953 |
| Cambrian-1 8B | 63.4 | 952 |
| xGen - MM - Interleave… | 59.5 | 979 |
| LLAVA-1.5 13B | 43.9 | 960 |
| LLAVA-1.5 7B | 40.7 | 951 |

**Open**

2025

One of the largest human
preference evaluations for VLMs

with 325k pairwise comparisons
and 870 human annotators

Molmo – 72B model ranks second for vision tasks

Barely second to GPT–4o

Outperforming Gemini 1.5 Pro and Claude–3.5

with 325k pairwise comparisons
and 870 human annotators

Molmo-7B outperforms
- Gemini 1.5 Flash
- LLAVA OneVision 72B
- GPT-4V
- QwenVL2 72B
- and many others

with 325k pairwise comparisons
and 870 human annotators

# Reaction online – released Sep 25, 2024

**Never bet against open-source software!**



**Jim Fan** ✔
@DrJimFan

I just pulled the numbers on vision-language benchmarks for Llama-3.2-11B (vision). Surprisingly, the open-source community at large isn't behind in the lightweight model class! Pixtral, Qwen2-VL, Molmo, and InternVL2 all stand strong. OSS AI models have never been stronger.

The last 3 lines are API-only frontier models. Gemini-flash and GPT-4o (likely in heavier-weight class) are still the reigning champions.

But never bet against OSS. Never underestimate the combined firepower of so many talents distributed all over the world.

| Models\Benchmark | MMMU | MathVista | ChartQA | AI2D | DocVQA | VQAv2 |
|---|---|---|---|---|---|---|
| Llama-3.2-11B | 50.7 | 51.5 | 83.4 | 91.1 | 88.4 | 75.2 |
| Pixtral-12B | 52.5 | 58 | 81.8 | 79 | 90.7 | 80.2 |
| Qwen2-VL-7B | 54.1 | 58.2 | 83 | 83 | 94.5 | 82.9 |
| Molmo-7B-D | 45.3 | 51.6 | 84.1 | 93.2 | 92.2 | 85.6 |
| InternVL2-8B | 51.2 | 58.3 | 83.3 | 83.8 | 91.6 | 76.7 |
|  |  |  |  |  |  |  |
| Claude-3 Haiku | 50.2 | 46.4 | 81.7 | 86.7 | 88.8 | 68.4 |
| Gemini-1.5 Flash | 56.1 | 58.4 | 85.4 | 91.7 | 89.9 | 80.1 |
| GPT-4o-0513 | 69.1 | 63.8 | 85.7 | 94.2 | 92.8 | 78.7 |

11:42 AM · Sep 25, 2024 · **45.6K** Views

Molmo grounds reasoning directly in the pixels

Example, it points when it counts



molmo > gemini 1.5 flash (at counting)

# Data matters! Quality over quantity even for pretraining

LLAMA 3.1V

6 Billion Image-Text pairs

Molmo is trained with

PixMo

700k Image–Text pairs

# Internet data is incidental
# Human annotated data is intentional

pink, japan,
aesthetic image

love this winter picture by
person

# PixMo data is intentional:



This photograph captures a well-organized work desk set prominently in the middle of the frame. The desk is large and rectangular, made from a polished, rich wood that spans horizontally across the image. Its structure is supported by four distinctive A-shaped legs, adding an elegant touch. On the desk, a striking dual-monitor setup is noticeable: a tall, vertical screen placed behind and to the right of a wider, horizontal computer monitor.

To the right of these monitors, a black mouse rests on a mouse pad. Scattered around the mouse pad, some white papers or letters are strewn across the far right side. On the left side of the desk, a black desk lamp with an extended arm hangs down, illuminating the workspace. Nearby, a stack of books is neatly placed in the upper left corner of the table.

The background wall is painted a subtle beige-white, complementing the refined ambiance of the space. The floor below the desk features elegant pinkish marble tiles, enhancing the room's sophisticated look. To the far right of the image, a large window or patio door allows natural light to pour in, with clear glass that offers a glimpse into the outside area. This exterior view includes part of a rustic brick wall and a metal pail, hinting at an adjacent patio.

# Collecting dense captions is hard!!!



This image features a screenshot taken from a tablet device. At the top-right corner, the time is displayed as 11:48, alongside a battery icon indicating 67% power remaining. The device also shows a telephony signal strength of four out of five bars and roughly three-quarters Wi-Fi connectivity. Additionally, a clock alarm icon is present, as well as an icon resembling a white rectangle with a blue letter "N" in it, whose specific function is unclear.

The main portion of the screen, approximately 80% of it, is a solid medium blue color devoid of any content. At the bottom of the screen, there are several folders and icons representing various apps and functionalities:

1. The first folder contains three icons:
- An icon likely for contacts.
- An icon probably for telephone.
- An icon that seems to represent a text messaging app.

2. The second folder houses four apps:
- A Gmail app icon indicated by a red "M" on a white background.
- A YouTube app icon characterized by a red play button.
- A map app icon depicting a map.
- An unidentified app icon represented by a blue folder with its top corner bent down on the right.

3. The third icon is a white oval with six dots, arranged in two rows of three, likely representing an app drawer or menu.

4. The fourth icon resembles a red, yellow, and green shutter with a blue dot in the middle, suggesting it might be for a camera or photo viewing app.

5. Next to it, there is a silver colored camera icon with a black lens and a blue spot in the center, hinting at a camera application.

At the very bottom of the screen, there are three navigational icons:
- A left arrow triangle in white at the bottom-left corner.
- A white oval in the center, indicative of a home button.
- A white rectangle on the bottom-right corner, likely for accessing recent apps or multitasking.

Overall, this image captures the home screen of a tablet, providing an overview of available functionalities and connectivity status.

Questions designed to extract meaningful visual information from annotators

- What is the image at first glance?
- What are the objects and their counts?
- What does the text say?
- What are the positions of the objects?
- What subtle details are noticeable?
- What is in the background?
- What is the style and color?

People don't like to type
… but they love to talk

We ask annotators to speak for 60 to 90 seconds about an image

We automatically convert speech into text for pretraining

Molmo

"Point to Mt Rainier" → Molmo → "Mt Rainier"

`<point x="63.5" y="44.5"alt="Mt Rainier">Mt Rainier</point>`

Large Language Model

Connector · · · Connector Tokenizer

CLIP · · · CLIP

Point
to Mt
Rainier

Start with off-the-self
Large Language Model
& Visual Encoder

# From perception To action



*"Point to the menu"*



*"Point to where I can set search options"*



*"Point to where I can find mid size datasets"*

# Pointing to count, pointing to ground

# Pointing examples

# Chaining Molmo + SAM 2

# Future: Embodied AI for Navigation & Manipulation

# Demo
# molmo.allenai.org

# Foundation Models

| Language | Classification | LM + Vision | And More! | Chaining |
|---|---|---|---|---|
| ELMo | CLIP | LLaVA | Segment Anything | LMs + CLIP |
| BERT | CoCa | Flamingo | Whisper | Visual Programming |
| GPT | | GPT-4V | Dalle | |
| T5 | | Gemini | Stable Diffusion | |
| | | MoImo | Imagen | |

# Segment Anything Model (SAM)

What does it mean to have a segmentation foundation model?



Masking model trained on dataset of specific number of objects (80 in COCO)

Model outputs masks of all objects in that image that is one of the categories of interest

Images: He et al. Mask R-CNN. 2017

# Segment Anything Model (SAM)

What does it mean to have a segmentation foundation model?



Masking model trained on a dataset of a huge number of categories

Model outputs mask of any objects that the user cares about

# Segment Anything Model (SAM)

What does it mean to have a segmentation foundation model?



Masking model trained on a dataset of a huge number of categories

**How to get this?**

Model outputs mask of any objects that the user cares about

**How to know this?**

Images: Kirillov et al. Segment Anything. 2023.

# How to know what to mask?



"Cats"

# Basic SAM Architecture



Images: Kirillov et al. Segment Anything. 2023.

# SAM Architecture



Images: Kirillov et al. Segment Anything. 2023.

# Ambiguity in correct prompt

# Ambiguity in correct prompt



Images: Kirillov et al. Segment Anything. 2023.

# SAM Architecture



Images: Kirillov et al. Segment Anything. 2023.

# Basic SAM Architecture



1. Loss only calculated with respect to best mask
2. Model also trained to output confidence score for each mask

Images: Kirillov et al. Segment Anything. 2023.

# Segment Anything Model (SAM)

What does it mean to have a segmentation foundation model?



Masking model trained on a dataset of a huge number of categories

**How to get this?**

Model outputs mask of any objects that the user cares about

**How to know this?**

Images: Kirillov et al. Segment Anything. 2023.

# Segment Anything Model (SAM)

# Segment Anything Model (SAM)



Images: Kirillov et al. Segment Anything. 2023.

# SAM Results



Image Source: Kirillov et al. Segment Anything. 2023

# SAM Results



Image Source: Kirillov et al. Segment Anything. 2023

# Zero-Shot with SAM



Image Source: https://segment-anything.com/

# Zero-Shot with SAM



Image Source: https://segment-anything.com/

# Foundation Models

| Language | Classification | LM + Vision | And More! | Chaining |
|---|---|---|---|---|
| ELMo | CLIP | LLaVA | Segment Anything | LMs + CLIP |
| BERT | CoCa | Flamingo | Whisper | Visual Programming |
| GPT | | GPT-4V | Dalle | |
| T5 | | Gemini | Stable Diffusion | |
| | | Molmo | Imagen | |

# What happens when a model is asked to classify a concept it has never seen?

A photo of a marimba
A photo of a viaduct
A photo of a papillon
A photo of a lorikeet



Pratt et al "What does a platypus look like? Generating customized prompts for zero-shot image classification". 2023.

# Solution: chaining
1. Get an LLM to generate a description.
2. Classify using the description

"A **marimba** is a large wooden percussion instrument that looks like a xylophone."
"A **viaduct** is a bridge composed of several spans supported by piers or pillars."
"A **papillon** is a small, spaniel-type dog with a long, silky coat and fringed ears."
"A **lorikeet** is a small to medium-sized parrot with a brightly colored plumage."



Pratt et al "What does a platypus look like? Generating customized prompts for zero-shot image classification". 2023.

# CuPL (CUstomized Prompts via Language models)



**LLM-prompts:**

"What does a {lorikeet, marimba, viaduct, papillon} look like?"

GPT-3

**Image-prompts:**

"A lorikeet is a small to medium-sized parrot with a brightly colored plumage."
"A marimba is a large wooden percussion instrument that looks like a xylophone."
"A viaduct is a bridge composed of several spans supported by piers or pillars."
"A papillon is a small, spaniel-type dog with a long, silky coat and fringed ears."

**Lorikeet**   **Marimba**   **Viaduct**   **Papillon**

Pratt et al "What does a platypus look like? Generating customized prompts for zero-shot image classification". 2023.

# CuPL (CUstomized Prompts via Language models)

| | ImageNet | DTD | Stanford Cars | SUN397 | Food101 | FGVC Aircraft | Oxford Pets | Caltech101 | Flowers 102 | UCF101 | Kinetics-700 | RESISC45 | CIFAR-10 | CIFAR-100 | Birdsnap |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| std | 75.54 | 55.20 | 77.53 | 69.31 | 93.08 | 32.88 | 93.33 | 93.24 | 78.53 | 77.45 | 60.07 | 71.10 | 95.59 | 78.26 | 50.43 |
| # hw | 80 | 8 | 8 | 2 | 1 | 2 | 1 | 34 | 1 | 48 | 28 | 18 | 18 | 18 | 1 |
| CuPL (base) | 76.19 | 58.90 | 76.49 | 72.74 | 93.33 | 36.69 | 93.37 | 93.45 | 78.83 | 77.74 | 60.24 | 68.96 | 95.81 | 78.47 | 51.11 |
| Δ std | +0.65 | +3.70 | -1.04 | +3.43 | +0.25 | +3.81 | +0.04 | +0.21 | +0.30 | +0.29 | +0.17 | -2.14 | +0.22 | +0.21 | +0.63 |
| # hw | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 |

Pratt et al "What does a platypus look like? Generating customized prompts for zero-shot image classification". 2023.

# Can we generalize the idea of chaining to all vision tasks?

Many Visual Question Answering models which have been trained to do this type of task
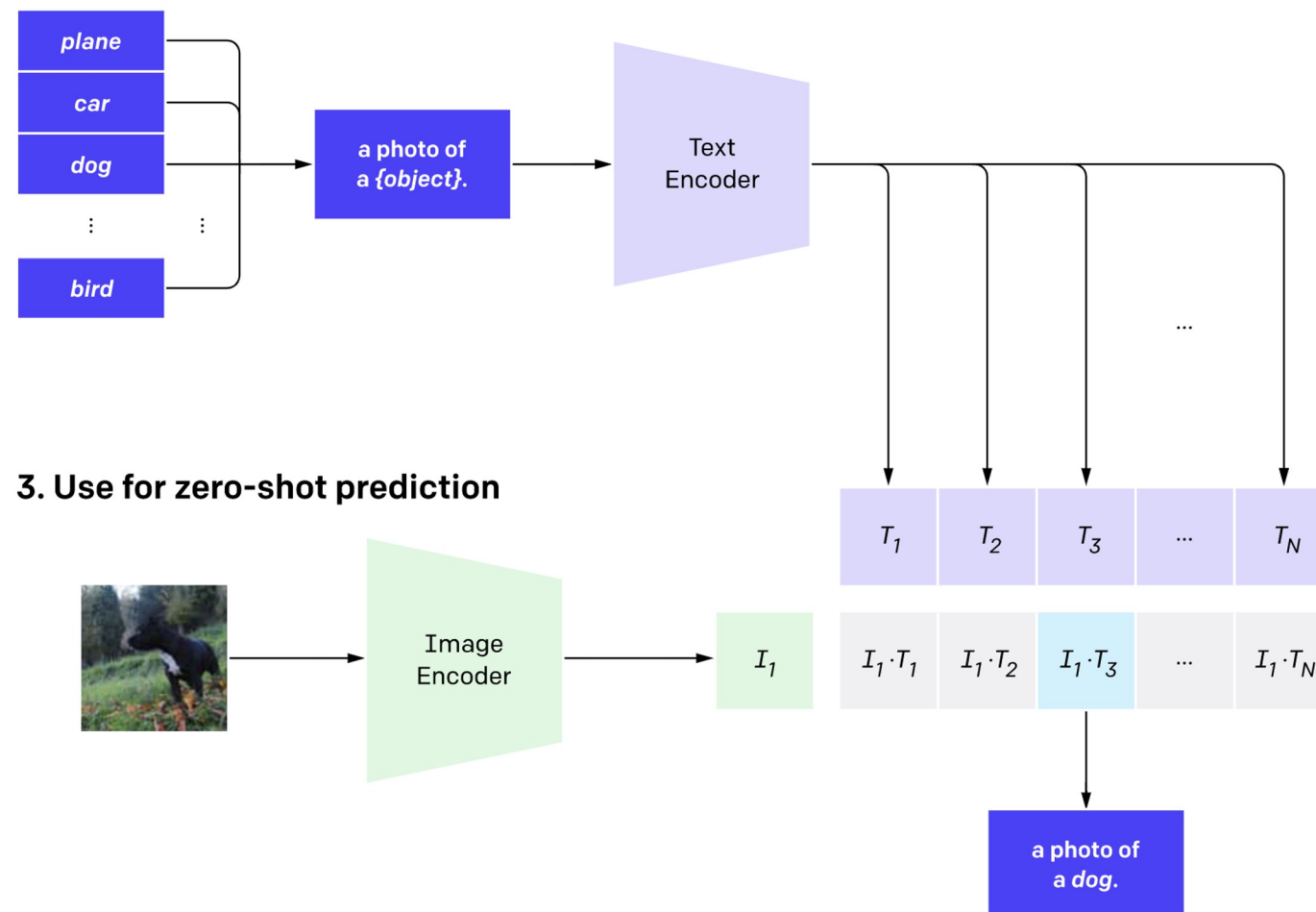


**Are there 3 people in the boat?**

Gupta et al "Visual Programming: Compositional visual reasoning without training". 2023.

# VisProg (visual programming)



LEFT:

RIGHT:

Statement: **The left and right image contains a total of six people and two boats.**

Gupta et al "Visual Programming: Compositional visual reasoning without training". 2023.

# VisProg (visual programming)

Train a new model for your task

output

Multi-image VQA model



Write a python script with the models you have

```
Class MyMultiImageVQA():

    Def ProcessIms():
        Ans1 = VQA(Image1)
        Ans2 = VQA(Image2)
        Return Ans1 + Ans2
```

**General to 2 images now, but not beyond that**

Gupta et al "Visual Programming: Compositional visual reasoning without training". 2023.

# VisProg (visual programming)



LEFT:

RIGHT:

Statement: **The left and right image contains a total of six people and two boats.**

**GPT**

```
Class MyMultiImageVQA():

    Def ProcessIms():
        Ans1 = VQA(Image1)
        Ans2 = VQA(Image2)
        Return Ans1 + Ans2
```

False

Gupta et al "Visual Programming: Compositional visual reasoning without training". 2023.

# VisProg (visual programming)



Gupta et al "Visual Programming: Compositional visual reasoning without training". 2023.

# VisProg (visual programming)



Gupta et al "Visual Programming: Compositional visual reasoning without training". 2023.

# VisProg (visual programming)



Gupta et al "Visual Programming: Compositional visual reasoning without training". 2023.

# VisProg (visual programming)



**Factual Knowledge Object Tagging**

IMAGE:

Prediction: IMAGE0

Instruction: Tag the 7 main characters on the TV show Big Bang Theory
Program:
OBJ0=FaceDet(image=IMAGE)
LIST0=List(query='main characters on the TV show Big Bang Theory', max=7)
OBJ1=Classify(image=IMAGE, object=OBJ0, categories=LIST0)
IMAGE0=Tag(image=IMAGE, object=OBJ1)
RESULT=IMAGE0

Gupta et al "Visual Programming: Compositional visual reasoning without training". 2023.

# VisProg (visual programming)



IMAGE:

Prediction: IMAGE0

Instruction: **Replace desert with lush green grass**
Program:
OBJ0=Seg(image=IMAGE)
OBJ1=Select(image=IMAGE, object=OBJ0, query='desert', category=None)
IMAGE0=Replace(image=IMAGE, object=OBJ1, prompt='lush green grass')
RESULT=IMAGE0

Gupta et al "Visual Programming: Compositional visual reasoning without training". 2023.

# Summary

# Summary



**2. Create dataset classifier from label text**

**3. Use for zero-shot prediction**

| DATASET | | IMAGENET RESNET101 | CLIP VIT-L |
|---|---|---|---|
| ImageNet | | 76.2% | 76.2% |
| ImageNet V2 | | 64.3% | 70.1% |
| ImageNet Rendition | | 37.7% | 88.9% |
| ObjectNet | | 32.6% | 72.3% |
| ImageNet Sketch | | 25.2% | 60.2% |
| ImageNet Adversarial | | 2.7% | 77.1% |

# Summary

# Summary

# Summary

Next time: Robot Learning