# Introduction to Vision-Language Models

## BIODS 271 / CS 277

Maya Varma

Stanford University

# Attendance Code

`multimodal`

# Why do we need VLMs?

The human experience of the world is multimodal, so we need AI systems capable of simultaneously processing diverse input modalities



**Data is often inherently multimodal**

# Part 1: Pretraining Methods and Datasets

# Contrastive Language-Image Pretraining (CLIP)

**Key Idea**: Maximize the similarity between true image-text embedding pairs and minimize similarity between mismatched image-text embedding pairs



**A dog sitting in a field**

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

# Contrastive Language-Image Pretraining (CLIP)

**Key Idea**: Maximize the similarity between true image-text embedding pairs and minimize similarity between mismatched image-text embedding pairs



A dog sitting in a field

Batch with N image-caption pairs



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

# Contrastive Language-Image Pretraining (CLIP)

**Key Idea**: Maximize the similarity between true image-text embedding pairs and minimize similarity between mismatched image-text embedding pairs

A dog sitting in a field

Batch with N image-caption pairs

Image Encoder

$I_1$ ← Image Embedding

$I_2$

$I_3$

$\vdots$

$I_N$

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

# Contrastive Language-Image Pretraining (CLIP)

**Key Idea**: Maximize the similarity between true image-text embedding pairs and minimize similarity between mismatched image-text embedding pairs



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

# Contrastive Language-Image Pretraining (CLIP)

**Key Idea**: Maximize the similarity between true image-text embedding pairs and minimize similarity between mismatched image-text embedding pairs



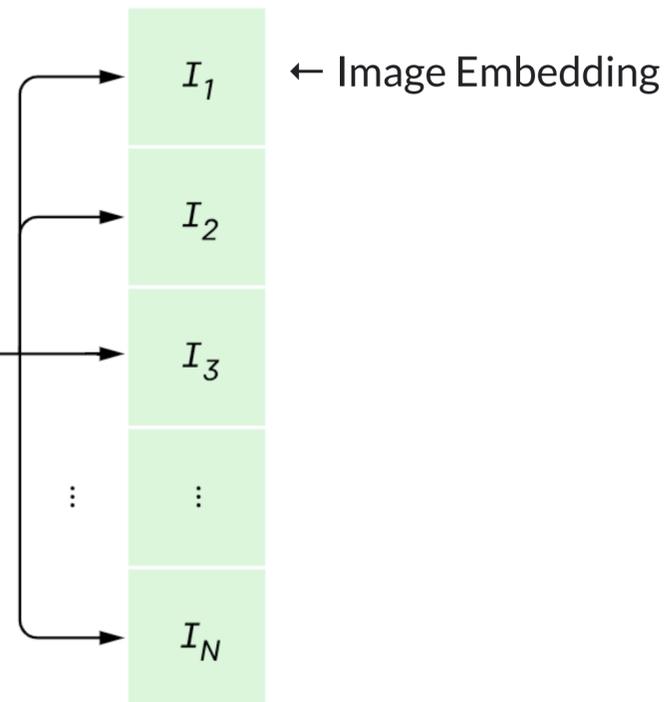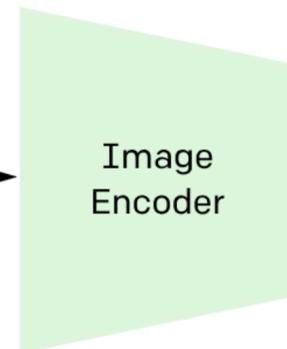A dog sitting in a field

Batch with N image-caption pairs

Text Encoder

Image Encoder

| | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|---|---|---|---|---|---|
| $I_1$ | $I_1{\cdot}T_1$ | $I_1{\cdot}T_2$ | $I_1{\cdot}T_3$ | ... | $I_1{\cdot}T_N$ |
| $I_2$ | $I_2{\cdot}T_1$ | $I_2{\cdot}T_2$ | $I_2{\cdot}T_3$ | ... | $I_2{\cdot}T_N$ |
| $I_3$ | $I_3{\cdot}T_1$ | $I_3{\cdot}T_2$ | $I_3{\cdot}T_3$ | ... | $I_3{\cdot}T_N$ |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋱ | ⋮ |
| $I_N$ | $I_N{\cdot}T_1$ | $I_N{\cdot}T_2$ | $I_N{\cdot}T_3$ | ... | $I_N{\cdot}T_N$ |

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

# Contrastive Language-Image Pretraining (CLIP)

**Key Idea**: Maximize the similarity between true image-text embedding pairs and minimize similarity between mismatched image-text embedding pairs
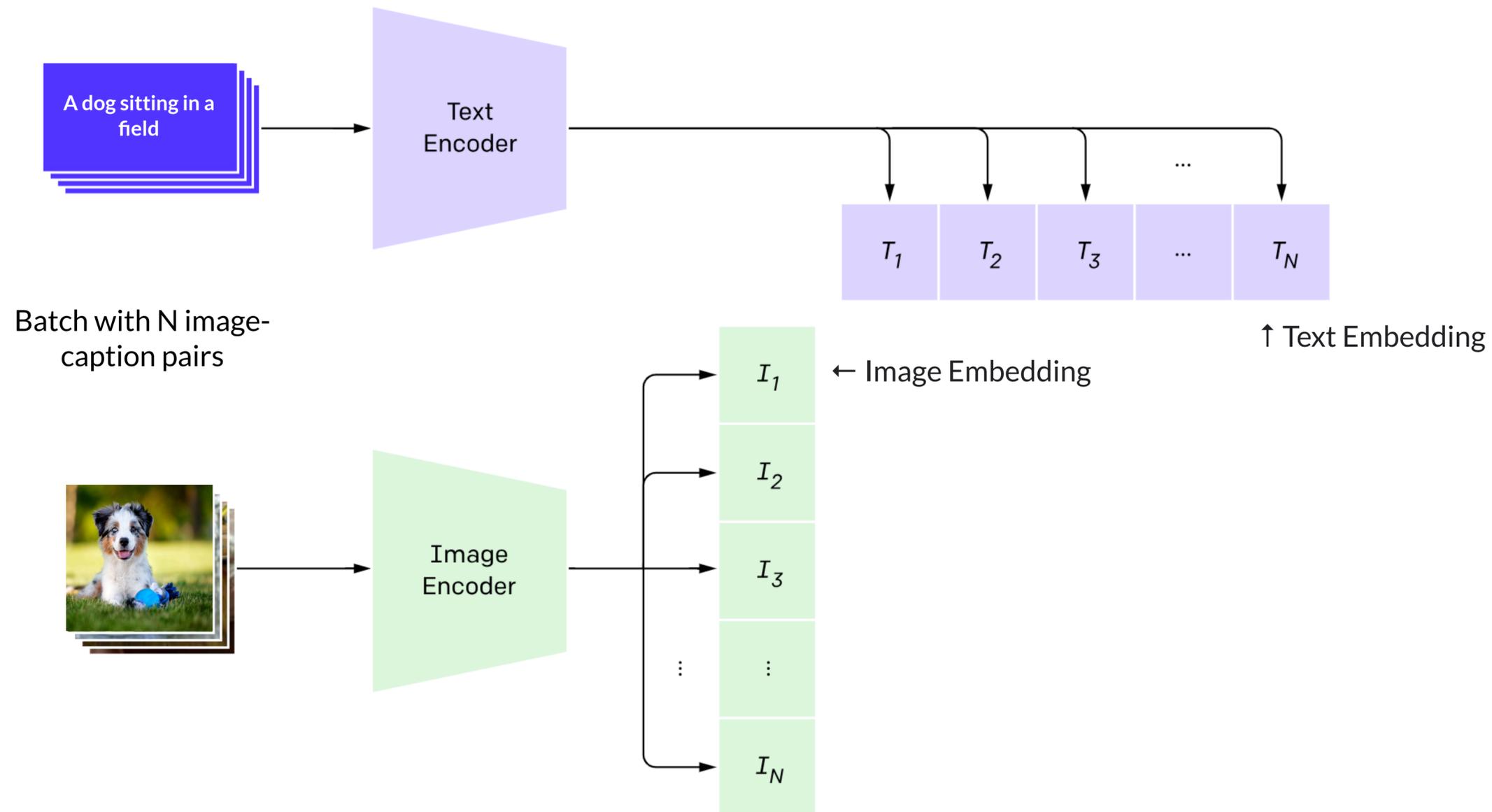


**Objective:** InfoNCE Loss Function

$$L_{I \rightarrow T} = \sum_{k=1}^{N} -\log \frac{exp(I_k \cdot T_k / \tau)}{\sum_{j=1}^{N} exp(I_k \cdot T_j / \tau)}$$

Positive Image-Text Pairs    Negative Image-Text Pairs

Softmax Function

Batch with N image-caption pairs

A dog sitting in a field

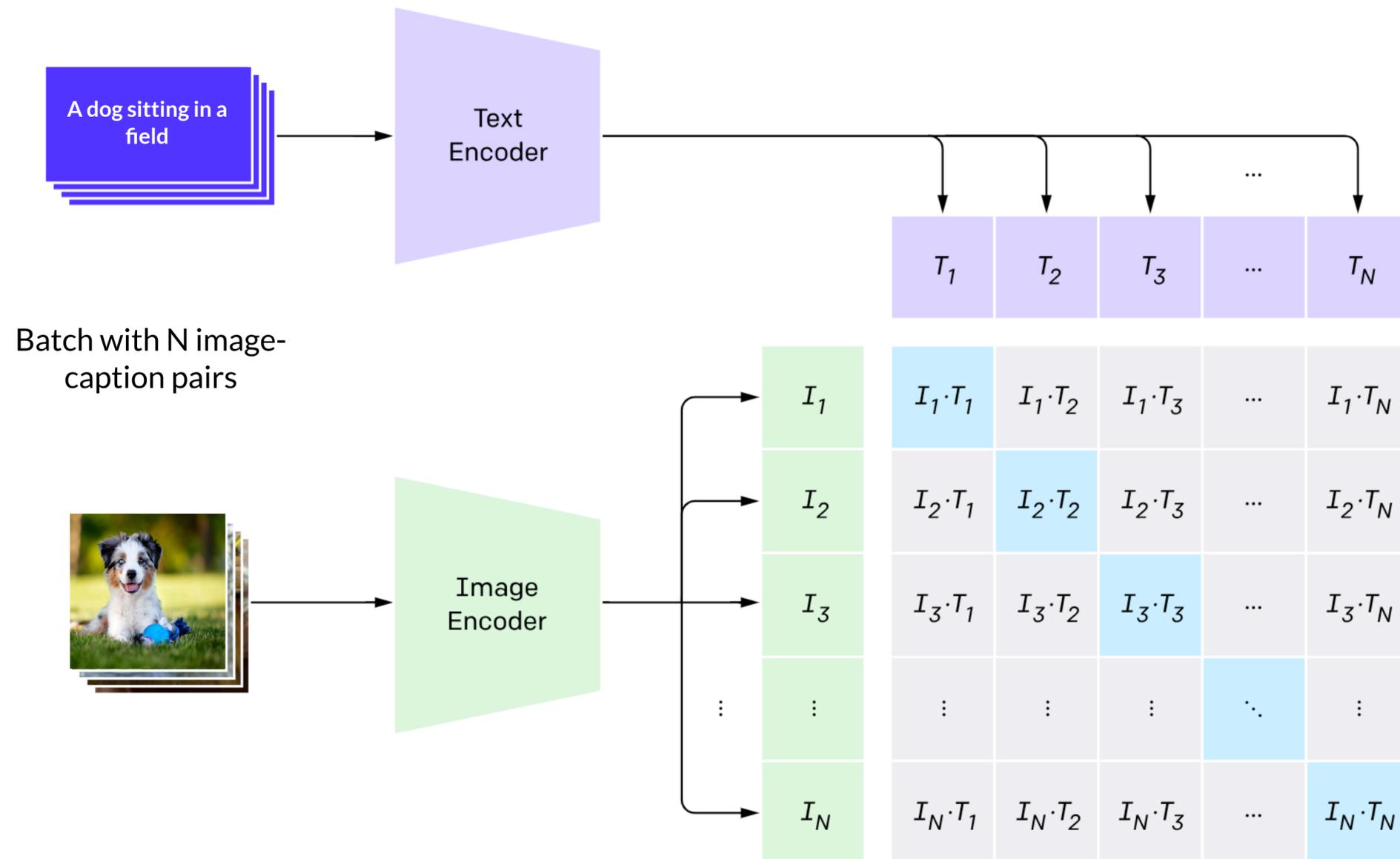Radford et al. "Learning Transferable Visual Models From Natural Language Supervision"

# Contrastive Language-Image Pretraining (CLIP)

**Key Idea**: Maximize the similarity between true image-text embedding pairs and minimize similarity between mismatched image-text embedding pairs
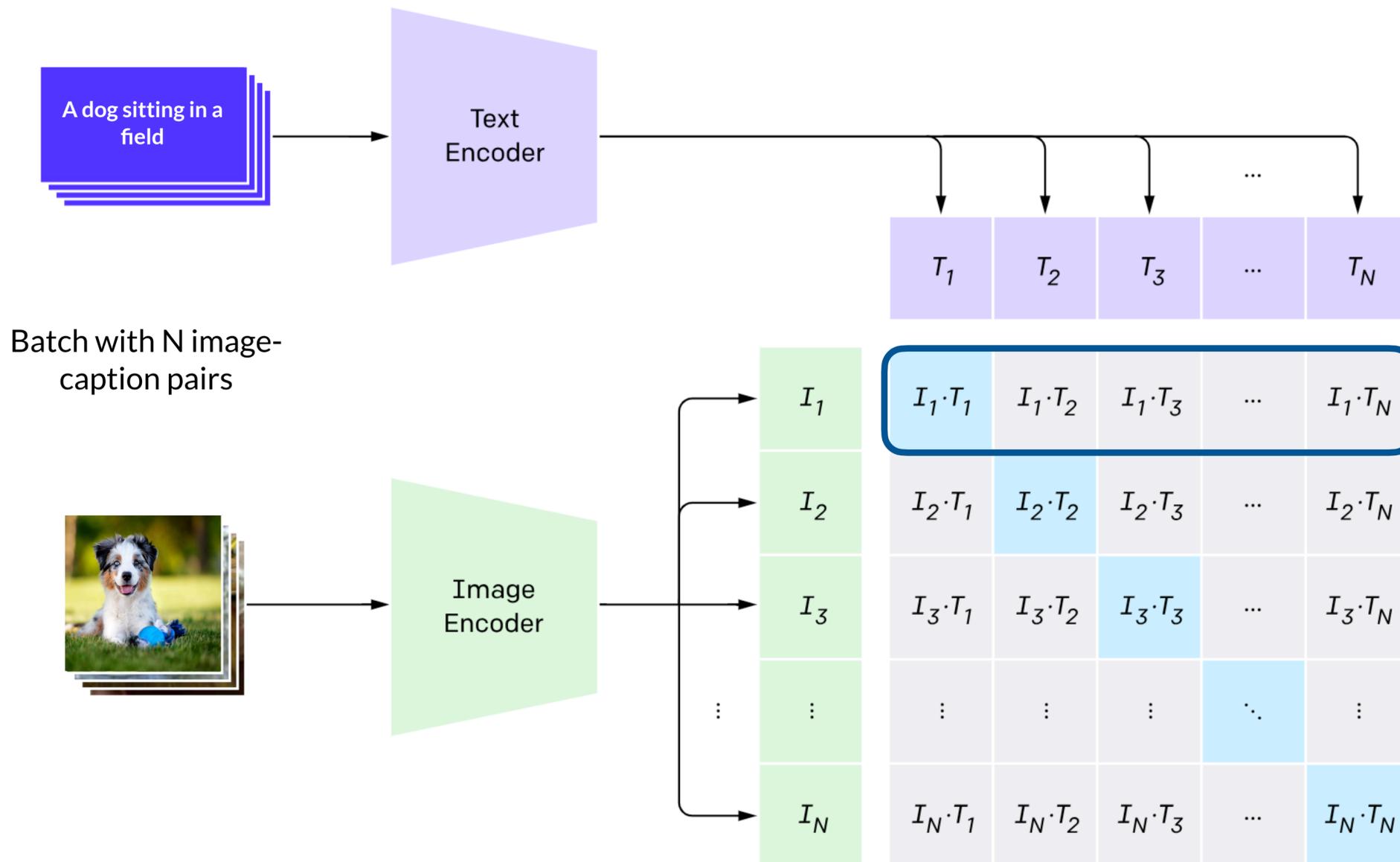


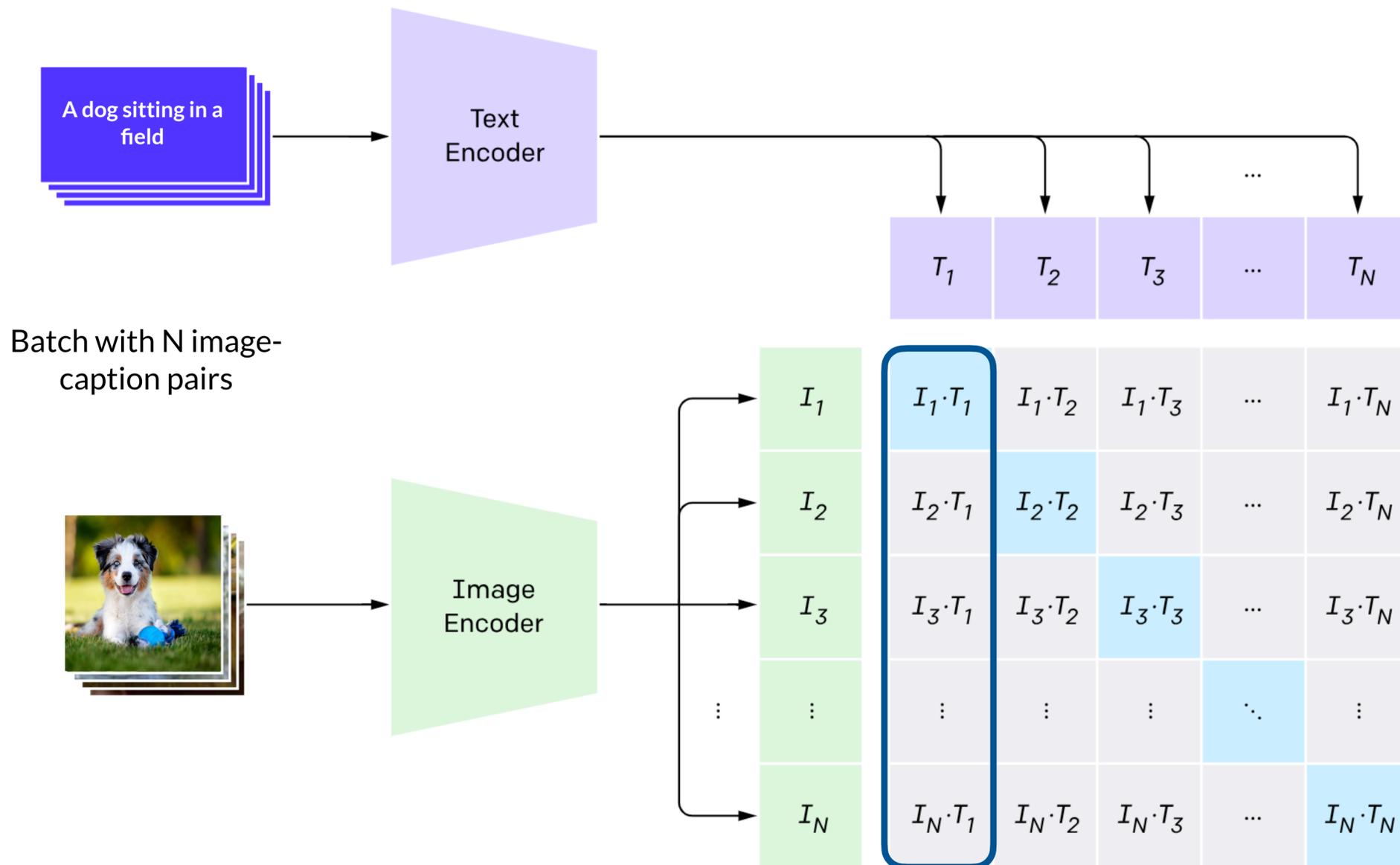**Objective:** InfoNCE Loss Function

$$L_{I \to T} = \sum_{k=1}^{N} -\log \frac{exp(I_k \cdot T_k / \tau)}{\sum_{j=1}^{N} exp(I_k \cdot T_j / \tau)}$$

$$L_{T \to I} = \sum_{k=1}^{N} -\log \frac{exp(I_k \cdot T_k / \tau)}{\sum_{j=1}^{N} exp(I_j \cdot T_k / \tau)}$$

$$L = L_{T \to I} + L_{I \to T}$$

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

# OpenCLIP



# OpenCLIP

[Paper] [Citations] [Clip Colab] [Coca Colab] `pypi v2.24.0`

Welcome to an open source implementation of OpenAI's CLIP (Contrastive Language-Image Pre-training).

Using this codebase, we have trained several models on a variety of data sources and compute budgets, ranging from small-scale experiments to larger runs including models trained on datasets such as LAION-400M, LAION-2B and DataComp-1B. Many of our models and their scaling properties are studied in detail in the paper reproducible scaling laws for contrastive language-image learning. Some of our best models and their zero-shot ImageNet-1k accuracy are shown below, along with the ViT-L model trained by OpenAI. We provide more details about our full collection of pretrained models here, and zero-shot results for 38 datasets here.

| Model | Training data | Resolution | # of samples seen | ImageNet zero-shot acc. |
|---|---|---|---|---|
| ConvNext-Base | LAION-2B | 256px | 13B | 71.5% |
| ConvNext-Large | LAION-2B | 320px | 29B | 76.9% |
| ConvNext-XXLarge | LAION-2B | 256px | 34B | 79.5% |
| ViT-B/32 | DataComp-1B | 256px | 34B | 72.8% |
| ViT-B/16 | DataComp-1B | 224px | 13B | 73.5% |
| ViT-L/14 | LAION-2B | 224px | 32B | 75.3% |
| ViT-H/14 | LAION-2B | 224px | 32B | 78.0% |
| ViT-L/14 | DataComp-1B | 224px | 13B | 79.2% |
| ViT-G/14 | LAION-2B | 224px | 34B | 80.1% |
| ViT-L/14 | OpenAI's WIT | 224px | 13B | 75.5% |

Ilharco et al. "OpenCLIP"
Cherti et al. "Reproducible scaling laws for contrastive language-image learning"
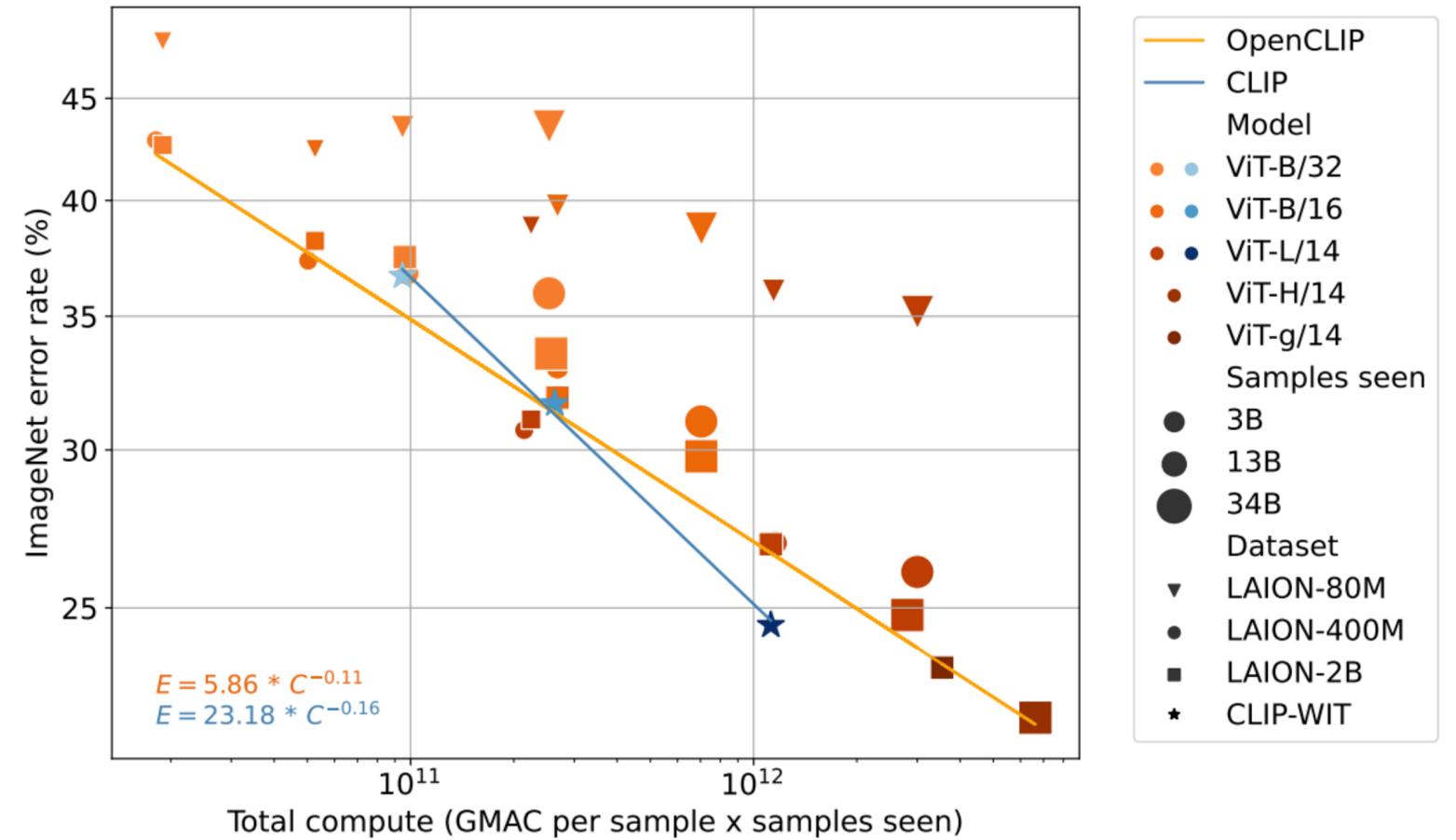
# OpenCLIP



## OpenCLIP

[Paper] [Citations] [Clip Colab] [Coca Colab] `pypi` `v2.24.0`

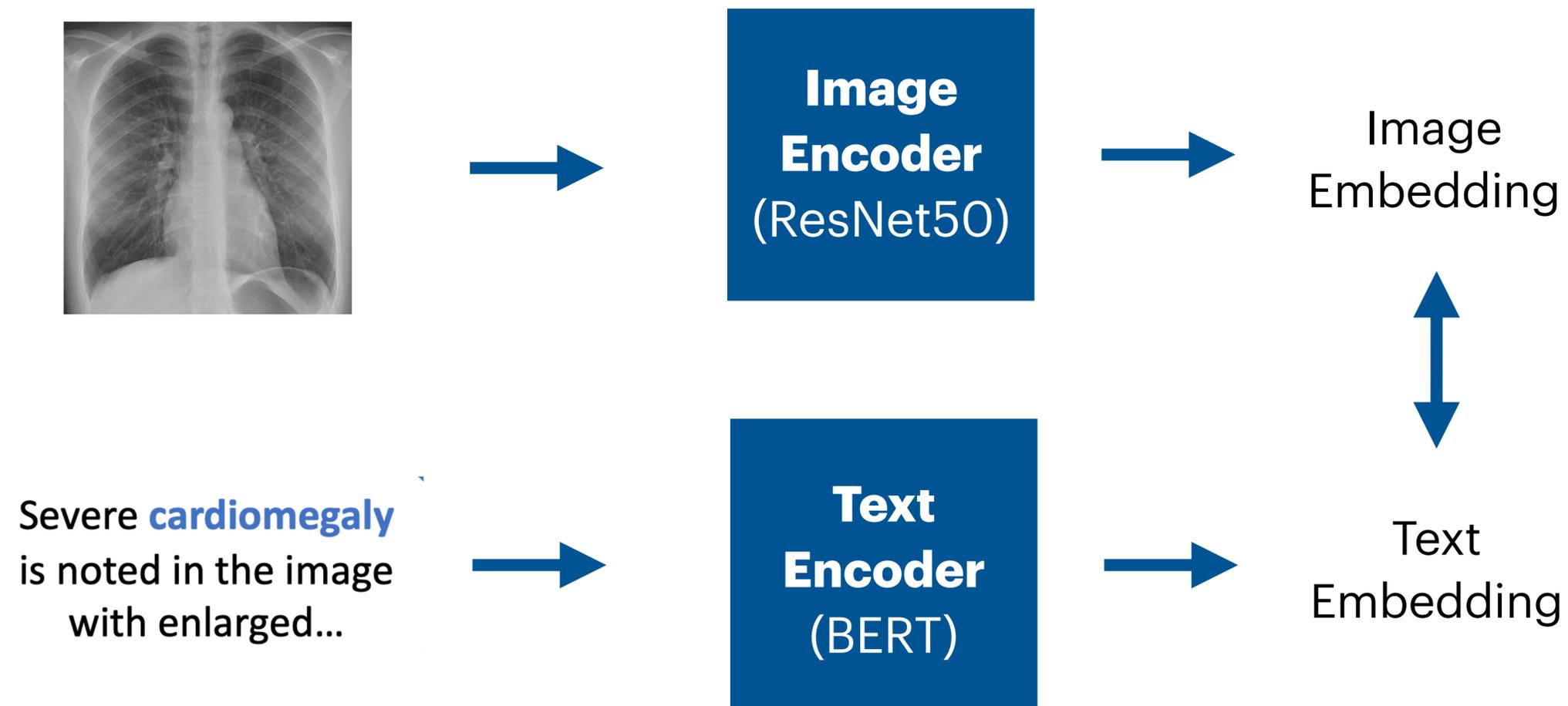Welcome to an open source implementation of OpenAI's CLIP (Contrastive Language-Image Pre-training).

Using this codebase, we have trained several models on a variety of data sources and compute budgets, ranging from small-scale experiments to larger runs including models trained on datasets such as LAION-400M, LAION-2B and DataComp-1B. Many of our models and their scaling properties are studied in detail in the paper reproducible scaling laws for contrastive language-image learning. Some of our best models and their zero-shot ImageNet-1k accuracy are shown below, along with the ViT-L model trained by OpenAI. We provide more details about our full collection of pretrained models here, and zero-shot results for 38 datasets here.

| Model | Training data | Resolution | # of samples seen | ImageNet zero-shot acc. |
|---|---|---|---|---|
| ConvNext-Base | LAION-2B | 256px | 13B | 71.5% |
| ConvNext-Large | LAION-2B | 320px | 29B | 76.9% |
| ConvNext-XXLarge | LAION-2B | 256px | 34B | 79.5% |
| ViT-B/32 | DataComp-1B | 256px | 34B | 72.8% |
| ViT-B/16 | DataComp-1B | 224px | 13B | 73.5% |
| ViT-L/14 | LAION-2B | 224px | 32B | 75.3% |
| ViT-H/14 | LAION-2B | 224px | 32B | 78.0% |
| ViT-L/14 | DataComp-1B | 224px | 13B | 79.2% |
| ViT-G/14 | LAION-2B | 224px | 34B | 80.1% |
| ViT-L/14 | OpenAI's WIT | 224px | 13B | 75.5% |

Ilharco et al. "OpenCLIP"
Cherti et al. "Reproducible scaling laws for contrastive language-image learning"

# ConVIRT

**Key Idea**: Maximize the similarity between true image-text embedding pairs and minimize similarity between mismatched image-text embedding pairs

Severe **cardiomegaly** is noted in the image with enlarged...

Image Encoder (ResNet50) → Image Embedding

Text Encoder (BERT) → Text Embedding

Zhang et al. "Contrastive Learning of Medical Visual Representations from Paired Images and Text"

# General-Domain Data: LAION-5B

**LAION-5B** contains 5 billion image-text pairs obtained from CommonCrawl
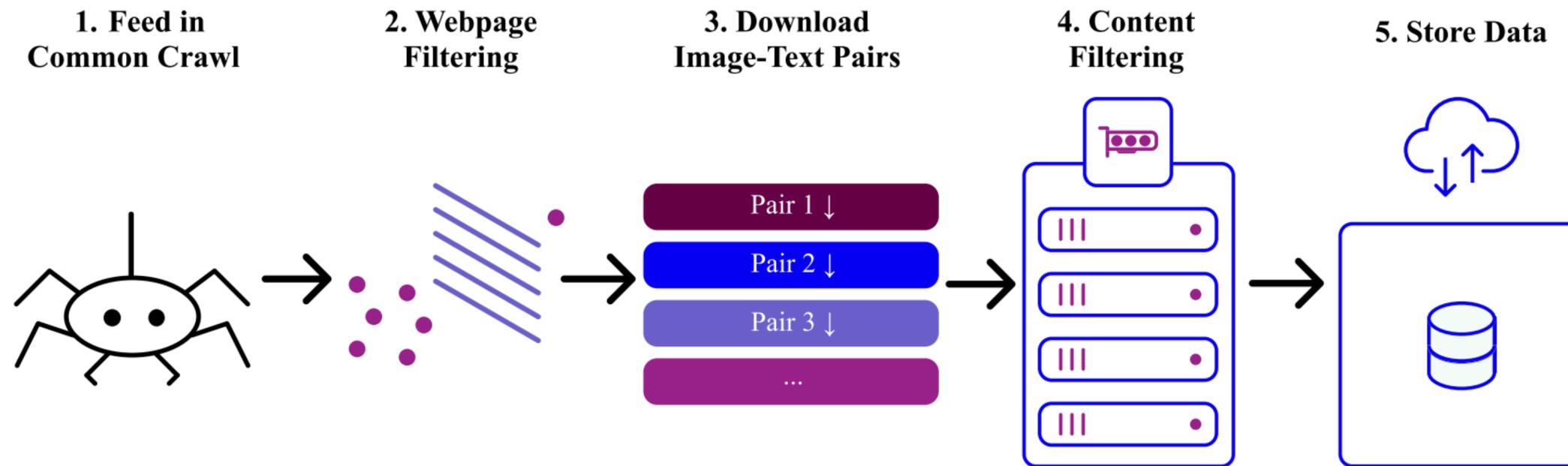


C: Green Apple Chair



C: sun snow dog



C: Color Palettes



C: pink, japan, aesthetic image

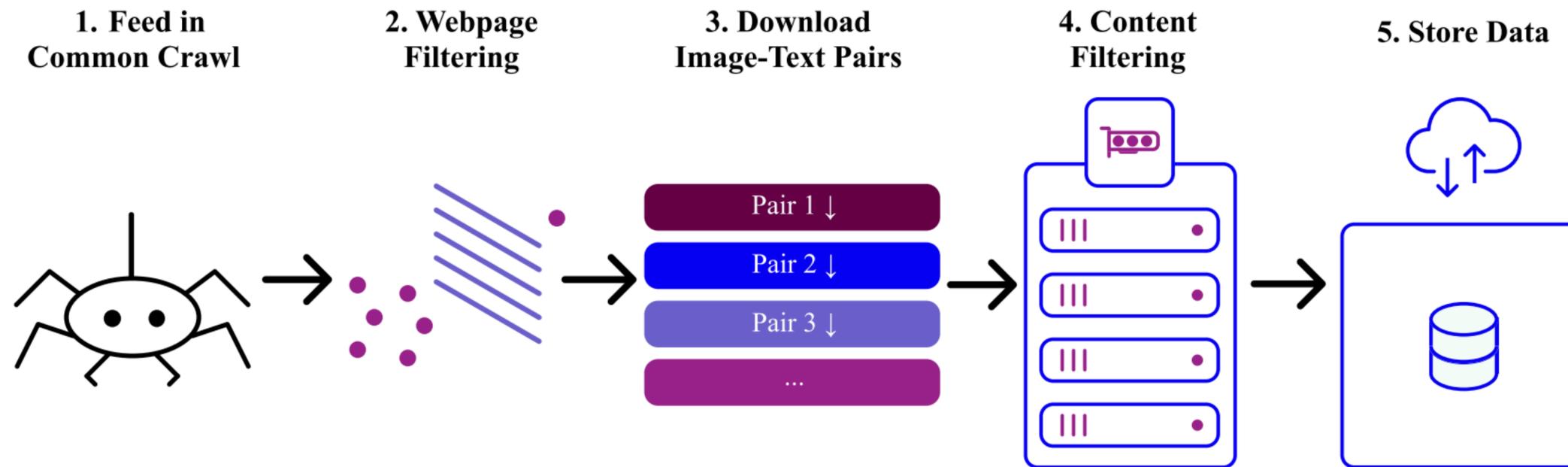Schumann et al. "LAION-5B: An open large-scale dataset for training next generation image-text models"

# General-Domain Data: LAION-5B

**LAION-5B** contains 5 billion image-text pairs obtained from CommonCrawl

# General-Domain Data: LAION-5B

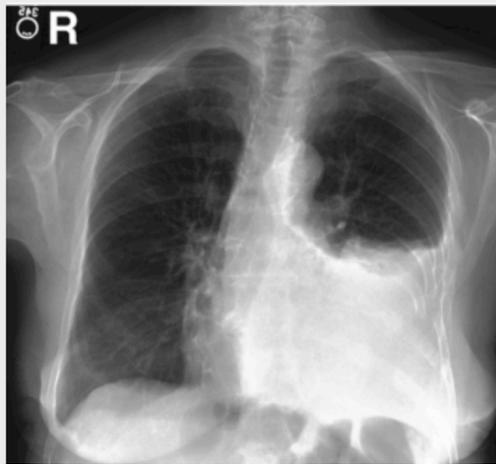**LAION-5B** contains 5 billion image-text pairs obtained from CommonCrawl



**Content filtering is performed using a pre-trained CLIP model**
(i.e. by computing cosine similarity between the image and text embeddings)

Schumann et al. "LAION-5B: An open large-scale dataset for training next generation image-text models"

# Medical-Domain Data

## MIMIC-CXR

370k chest X-rays with 220k reports

**Cardiac size cannot be evaluated. Large left pleural effusion is new. Small right effusion is new. The upper lungs are clear. There is no pneumothorax.**
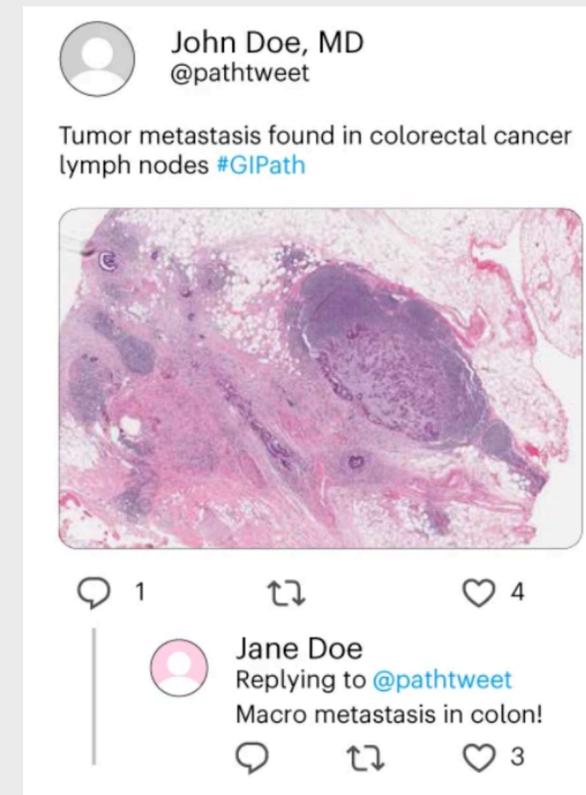
## PadChest

160k chest X-rays with 110k reports (Spanish)

**cambi pulmonar cronic sever. sign fibrosis bibasal. sutil infiltr pseudonodul milimetr vidri deslustr localiz bas. cifosis sever**
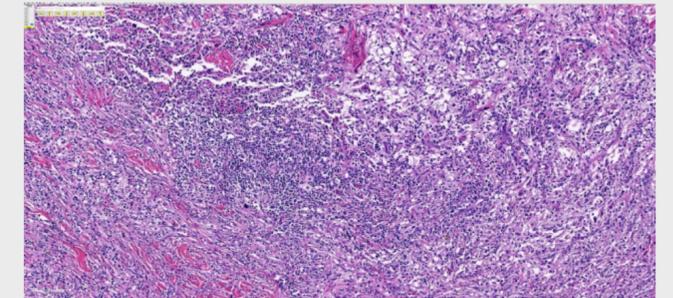
## OpenPath

200k histopathology image-text pairs (Twitter)

John Doe, MD
@pathtweet

Tumor metastasis found in colorectal cancer lymph nodes #GIPath

💬 1    ↻    ♡ 4

Jane Doe
Replying to @pathtweet
Macro metastasis in colon!

💬    ↻    ♡ 3

## Quilt-1M

1M histopathology image-text pairs (Youtube)

**Large histiocytes with abundant cytoplasm identified as Rosai-Dorfman histiocytes**

Johnson et al. "MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports."
Bustos et al. "PadChest: A large chest x-ray image dataset with multi-label annotated reports"
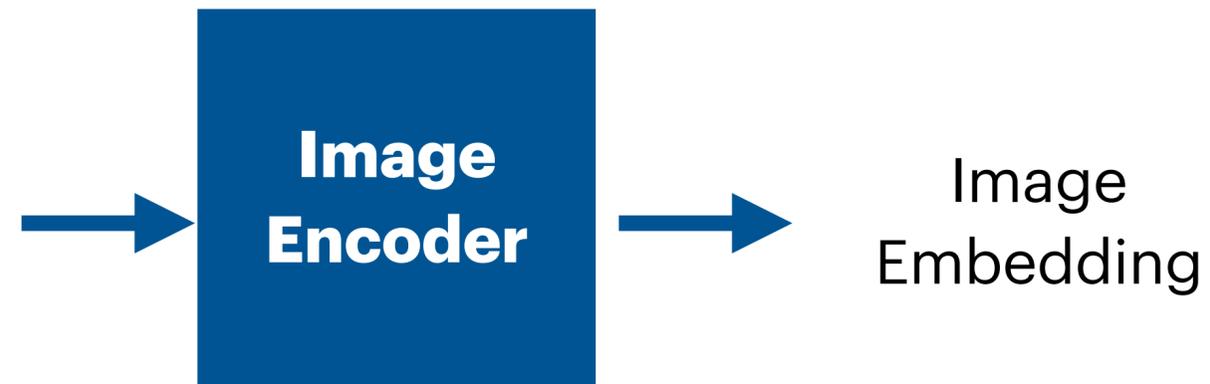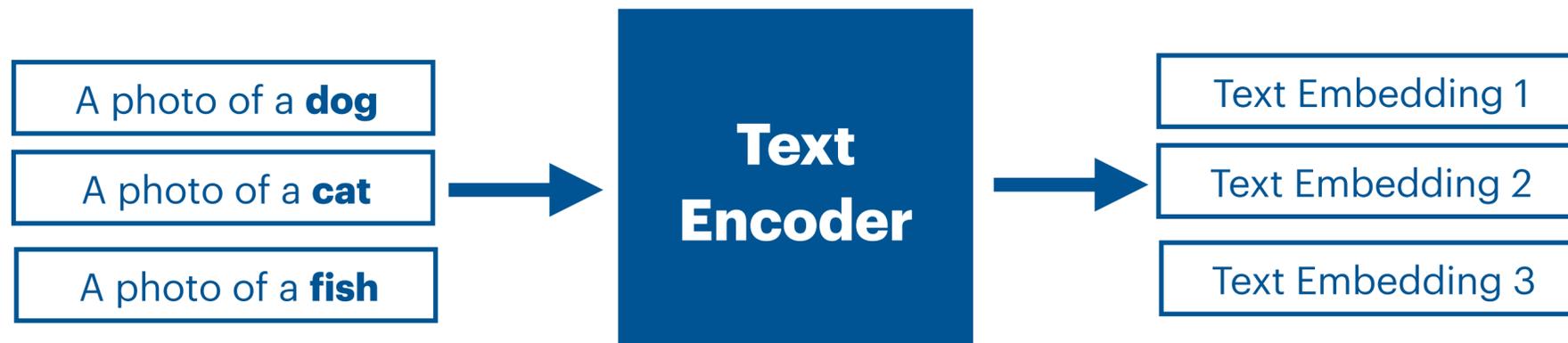Ikezogwo et al. "Quilt-1M: One Million Image-Text Pairs for Histopathology"
Huang et al. "A visual-language foundation model for pathology image analysis using medical Twitter"
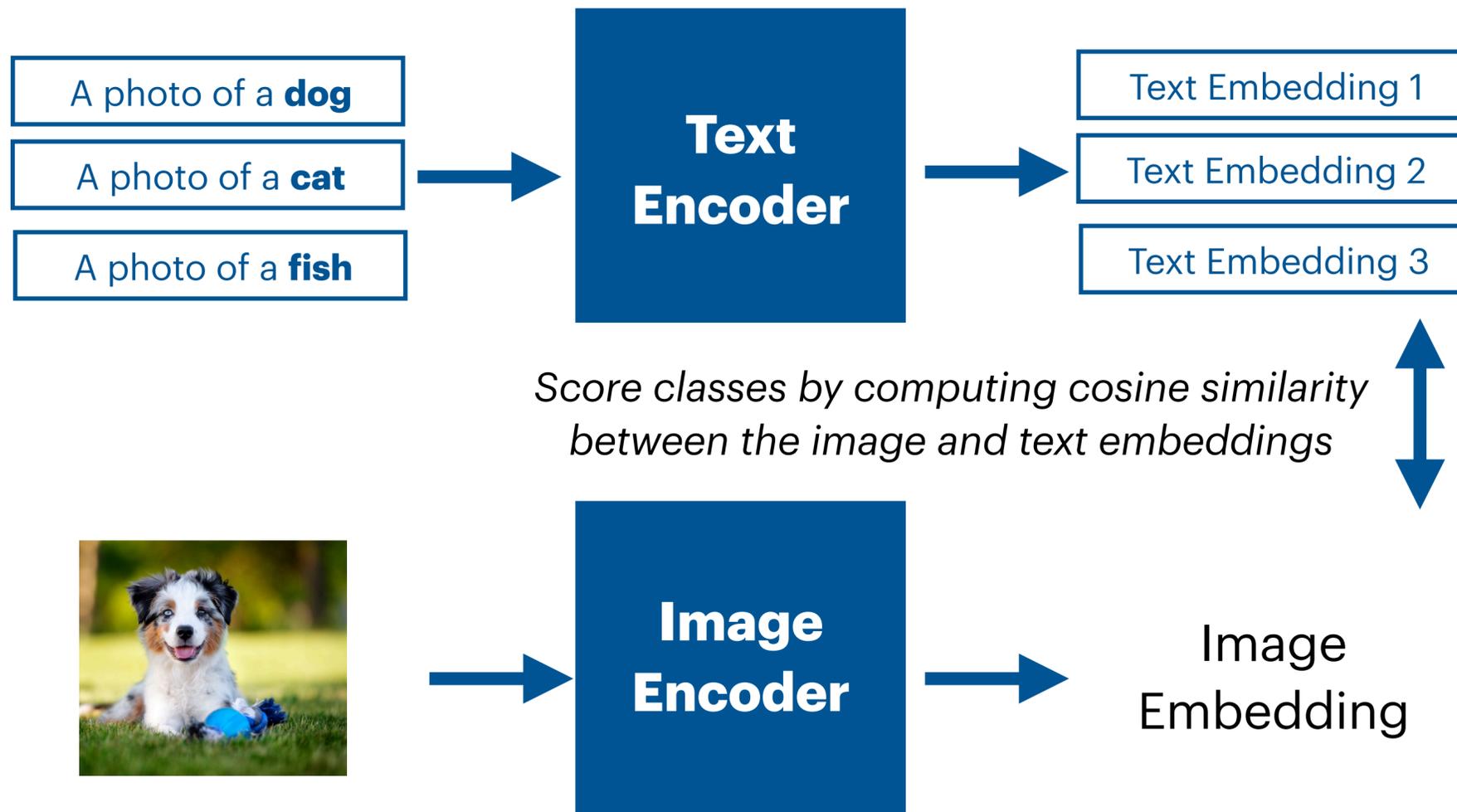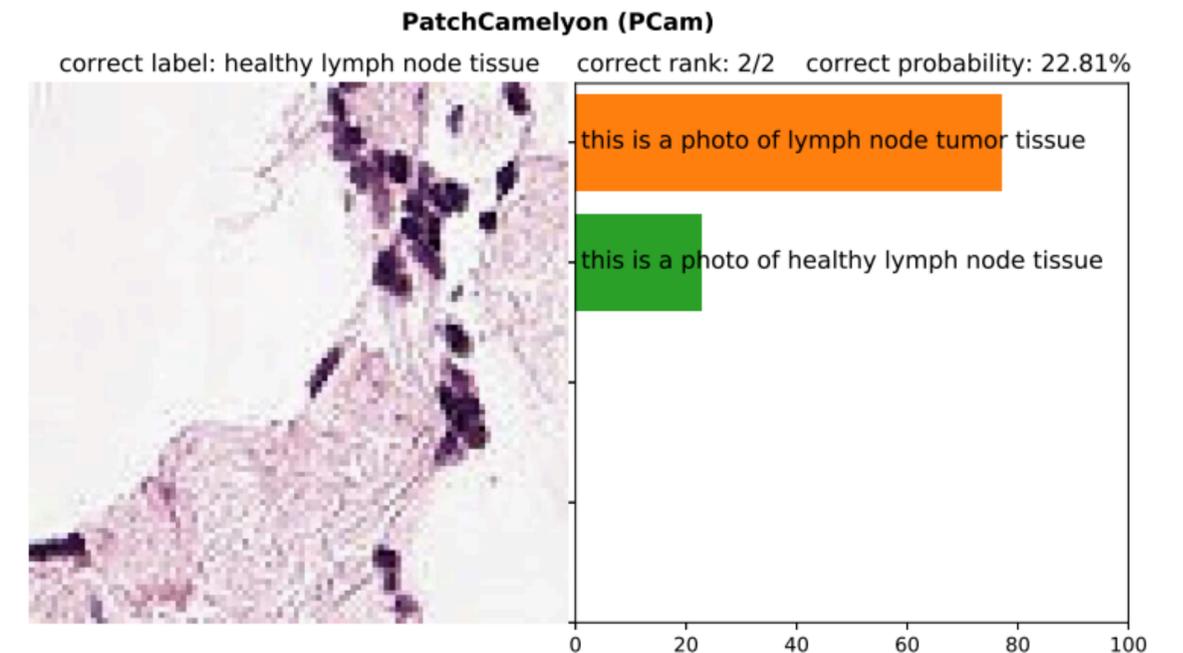
# Part 2: Evaluating VLMs
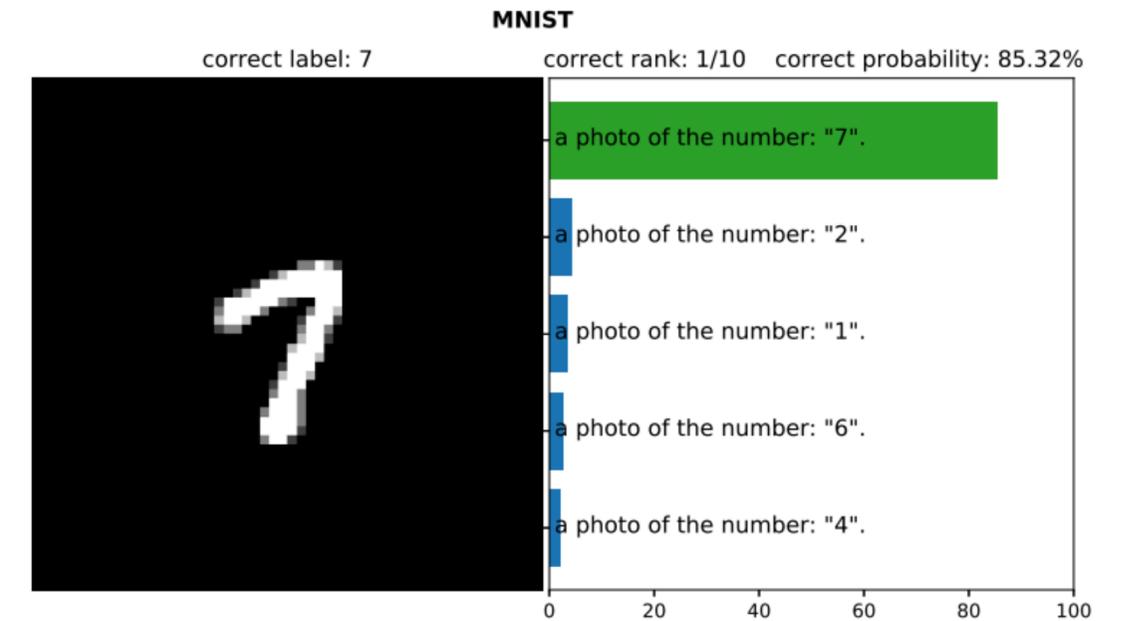
# Evaluating VLMs

## Zero-Shot Classification



| A photo of a **dog** |
| A photo of a **cat** |
| A photo of a **fish** |

→ **Text Encoder** →

| Text Embedding 1 |
| Text Embedding 2 |
| Text Embedding 3 |

**Image Encoder** → Image Embedding

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

# Evaluating VLMs

## Zero-Shot Classification



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

# Evaluating VLMs

## Zero-Shot Classification



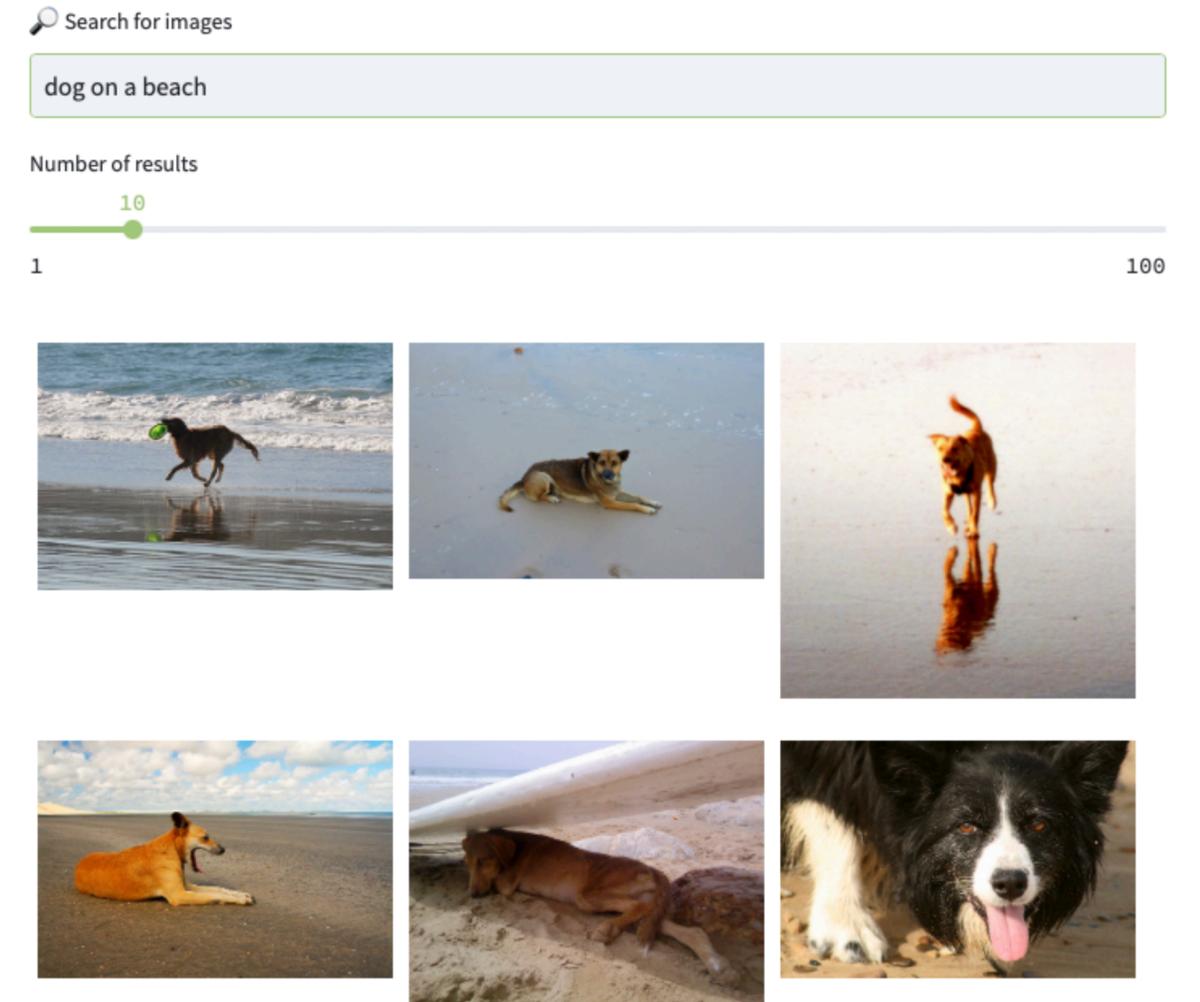Radford et al. "Learning Transferable Visual Models From Natural Language Supervision"
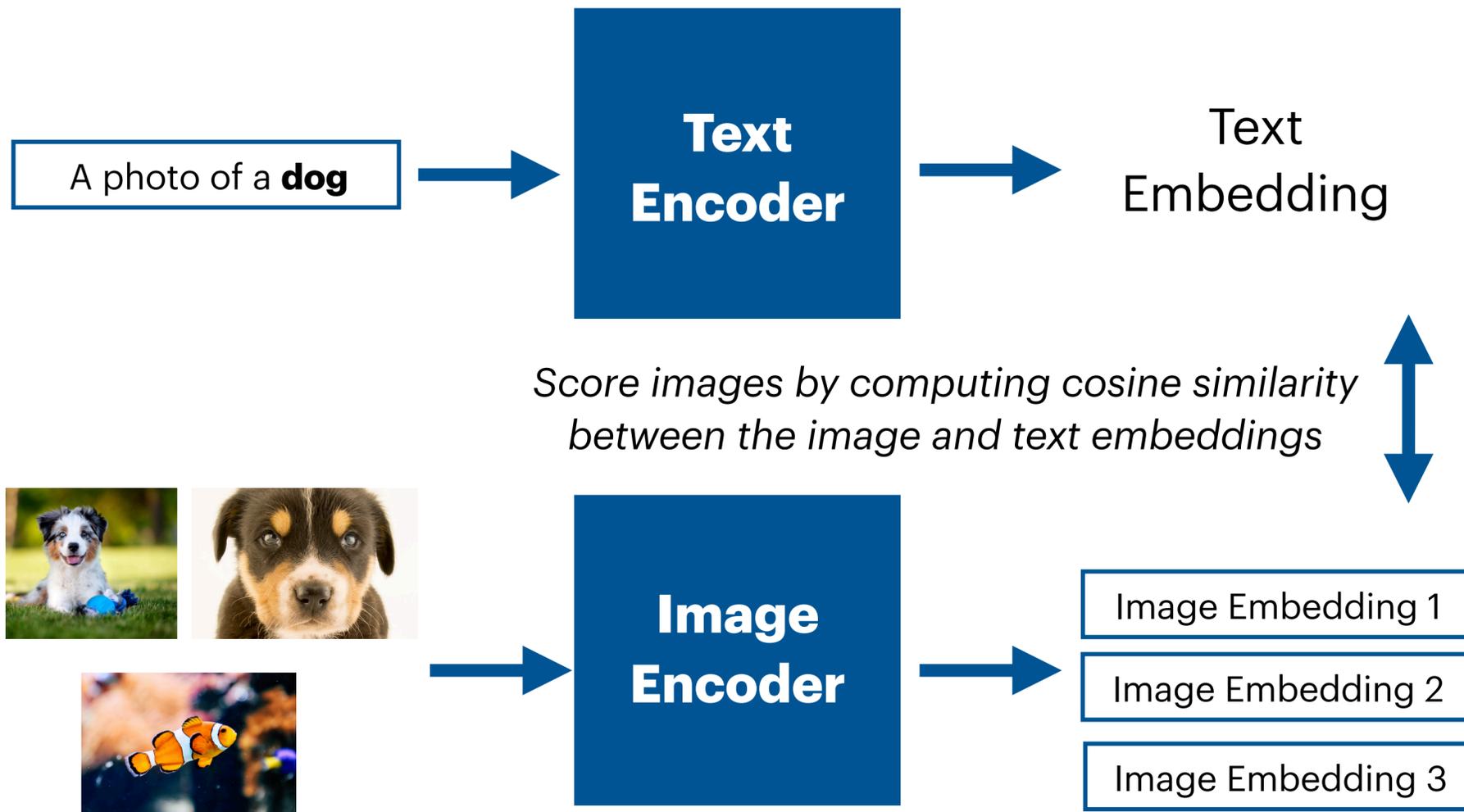
# Evaluating VLMs

**Text → Image Retrieval**

# Evaluating VLMs
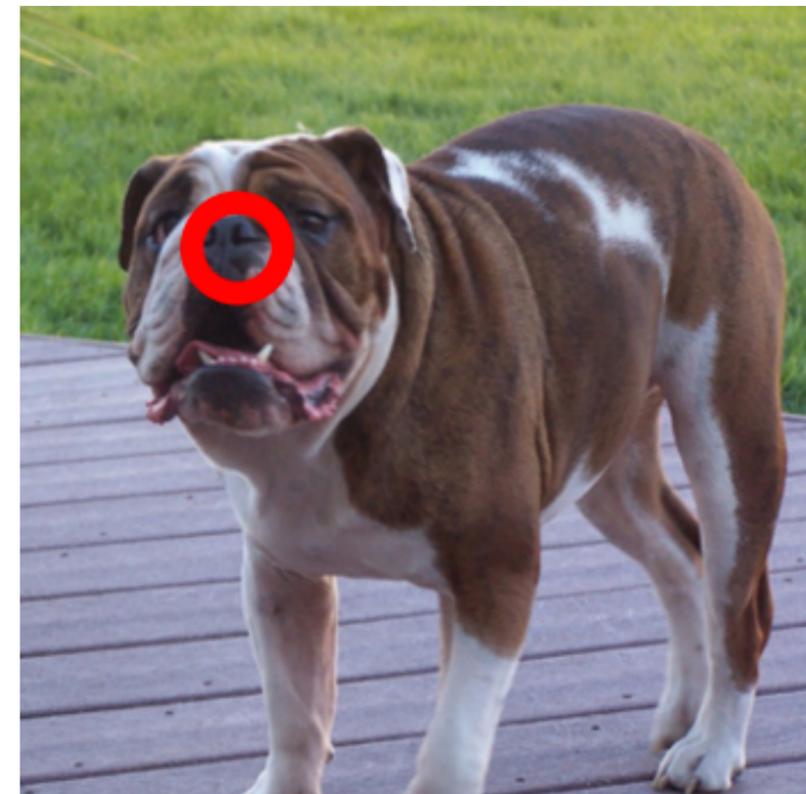
## Text → Image Retrieval

# Prompting VLMs

## Textual Prompts

*Example text prompts used by CLIP for zero-shot classification on CIFAR-10*

```
templates = [
    'a photo of a {}.',
    'a blurry photo of a {}.',
    'a black and white photo of a {}.',
    'a low contrast photo of a {}.',
    'a high contrast photo of a {}.',
    'a bad photo of a {}.',
    'a good photo of a {}.',
    'a photo of a small {}.',
    'a photo of a big {}.',
    'a photo of the {}.',
    'a blurry photo of the {}.',
    'a black and white photo of the {}.',
    'a low contrast photo of the {}.',
    'a high contrast photo of the {}.',
    'a bad photo of the {}.',
    'a good photo of the {}.',
    'a photo of the small {}.',
    'a photo of the big {}.',
]
```
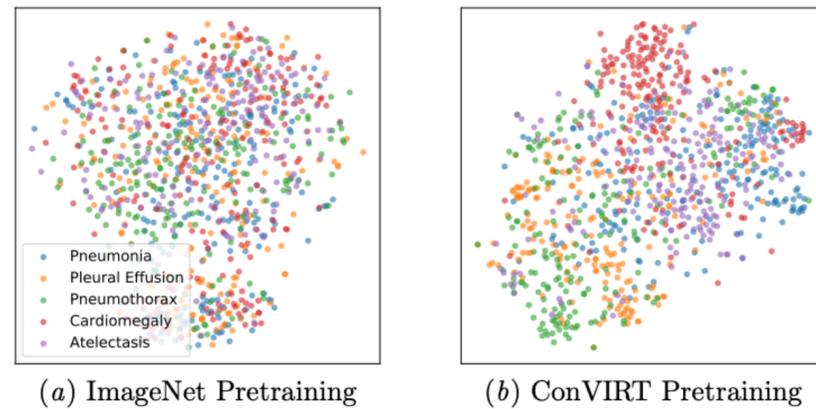
## Visual Prompts

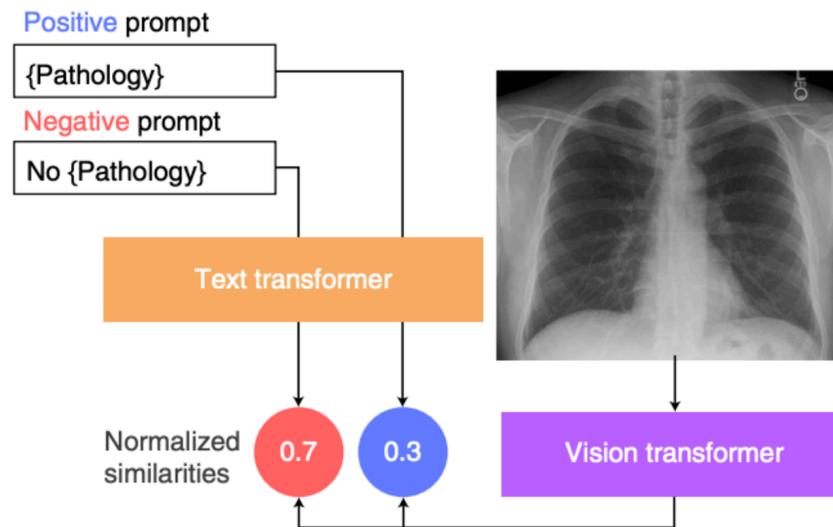*Adding visual signal to images can help with targeted retrieval and classification*

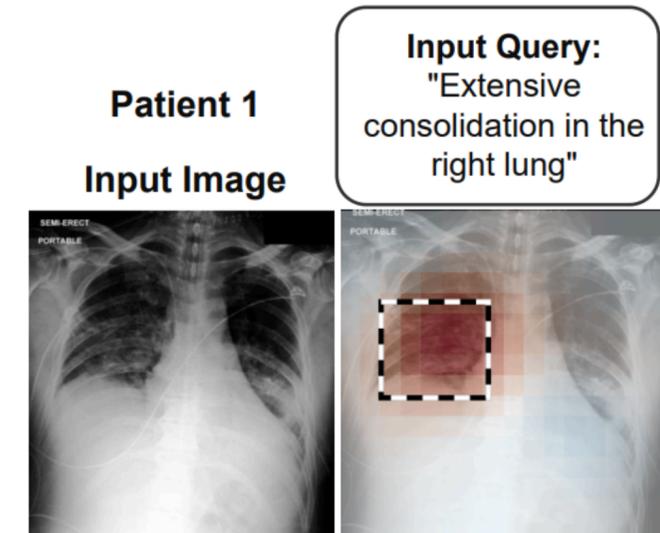Radford et al. "Learning Transferable Visual Models From Natural Language Supervision"
Shtedritski et al. "What does CLIP know about a red circle? Visual prompt engineering for VLMs"

# Evaluating Medical VLMs

## Classification



(a) ImageNet Pretraining    (b) ConVIRT Pretraining

- Pneumonia
- Pleural Effusion
- Pneumothorax
- Cardiomegaly
- Atelectasis

## Zero-Shot Classification



Positive prompt
{Pathology}

Negative prompt
No {Pathology}

Text transformer

Normalized similarities    0.7    0.3

Vision transformer

## Visual Grounding



Patient 1

Input Image

Input Query:
"Extensive consolidation in the right lung"

## Segmentation



## Text to Image Retrieval



Breast tumor surrounded by fat

## Natural Language Inference

Sentence 1:

No pneumothorax is seen

Sentence 2:

Previously-seen pneumothorax is no longer visualized
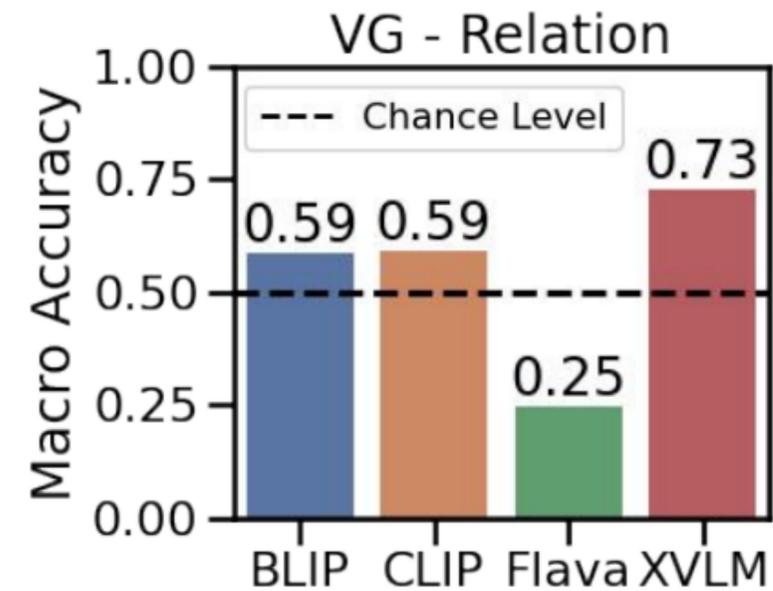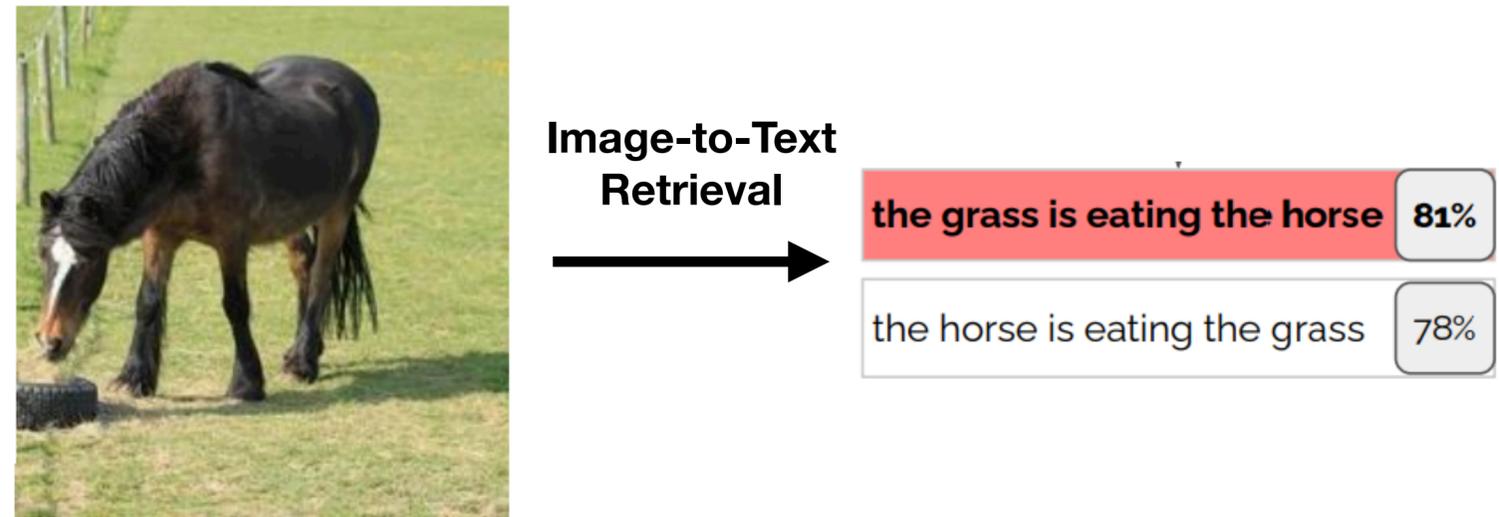
Type: **Entailment**

Miura et al. "RadNLI: A natural language inference dataset for the radiology domain"
Zhang et al. "Contrastive Learning of Medical Visual Representations from Paired Images and Text"
Huang et al. "A visual-language foundation model for pathology image analysis using medical Twitter"
Tiu et al. "Expert-level detection of pathologies from unannotated chest X-ray images via self-supervised learning"
Boecking et al. "Making the Most of Text Semantics to Improve Biomedical Vision–Language Processing"

# Part 3: Limitations

# Limitations: Contrastive Training

## Complex Patterns (e.g. counting)



## Relational Understanding



Paiss et al. "Teaching CLIP to Count to Ten"
Yuksekgonul et al. "When and Why Vision-Language Models Behave Like Bags-of-Words and What to Do About it?"

# Limitations: Domain-Specific Challenges

**Fine-Grained Visual Information**



**Lengthy and Complex Text**



Chen et al. "Toward Expanding the Scope of Radiology Report Summarization to Multiple Anatomies and Modalities"

# Questions?