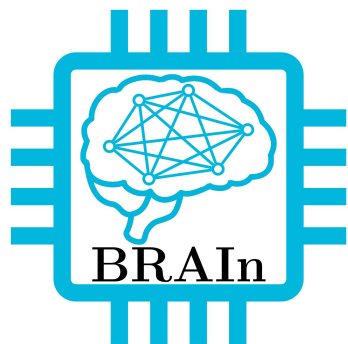


Foundation Models for Vision: CLIP: Contrastive Language-Image Pretraining

Équipe BRAIn



IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

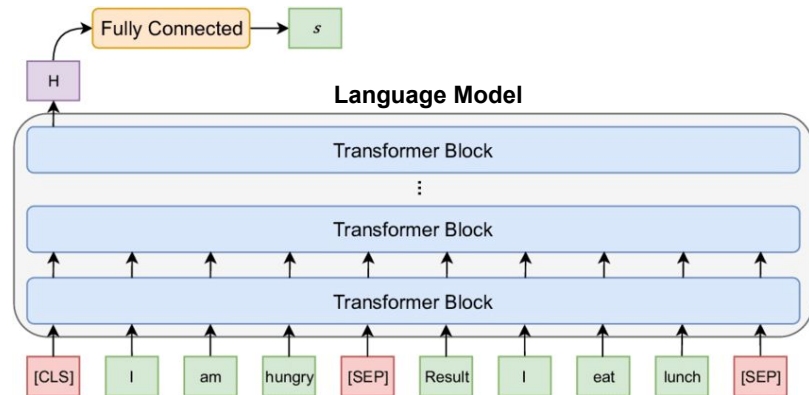


Vision



Uni-modal

Text

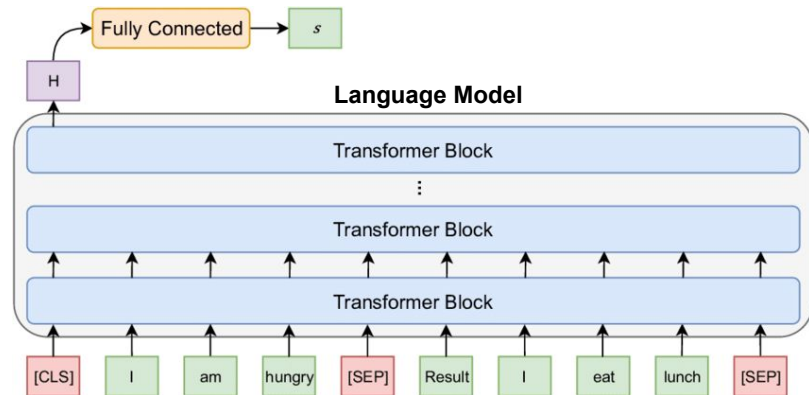


Vision



Uni-modal

Text



a teddy bear on a skateboard in times square



Text to Image generation

Multi-modal

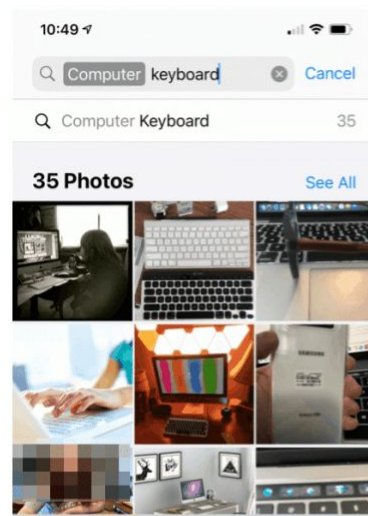
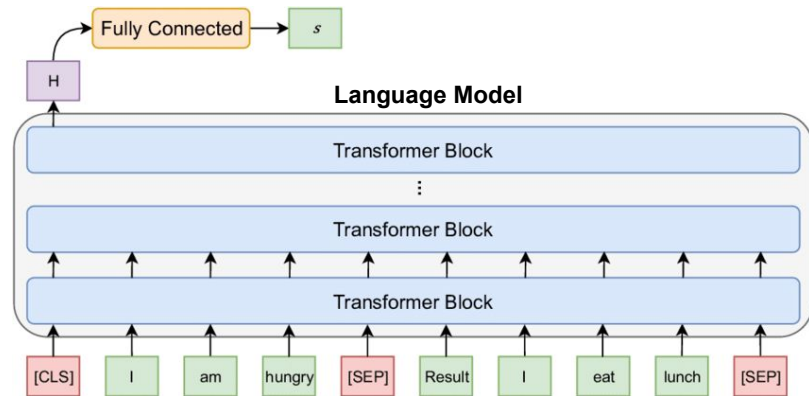


Image retrieval

Vision

Text

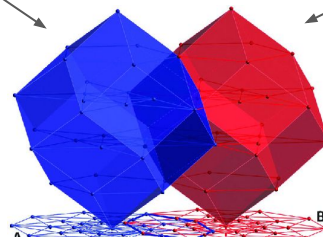
Uni-modal



a teddy bear on a skateboard in times square



Text to Image generation



Alignment

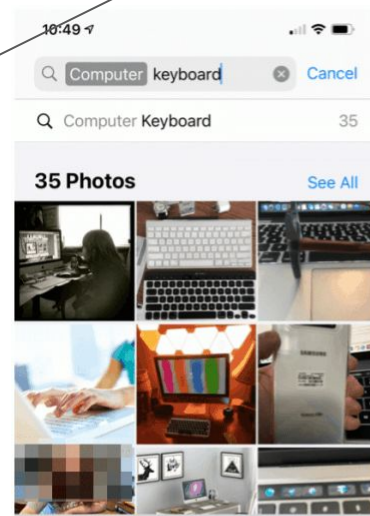


Image retrieval

Multi-modal

CLIP: Connecting text and images

Cited by 12033

Web ImageText dataset: 400M paired image-text



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.



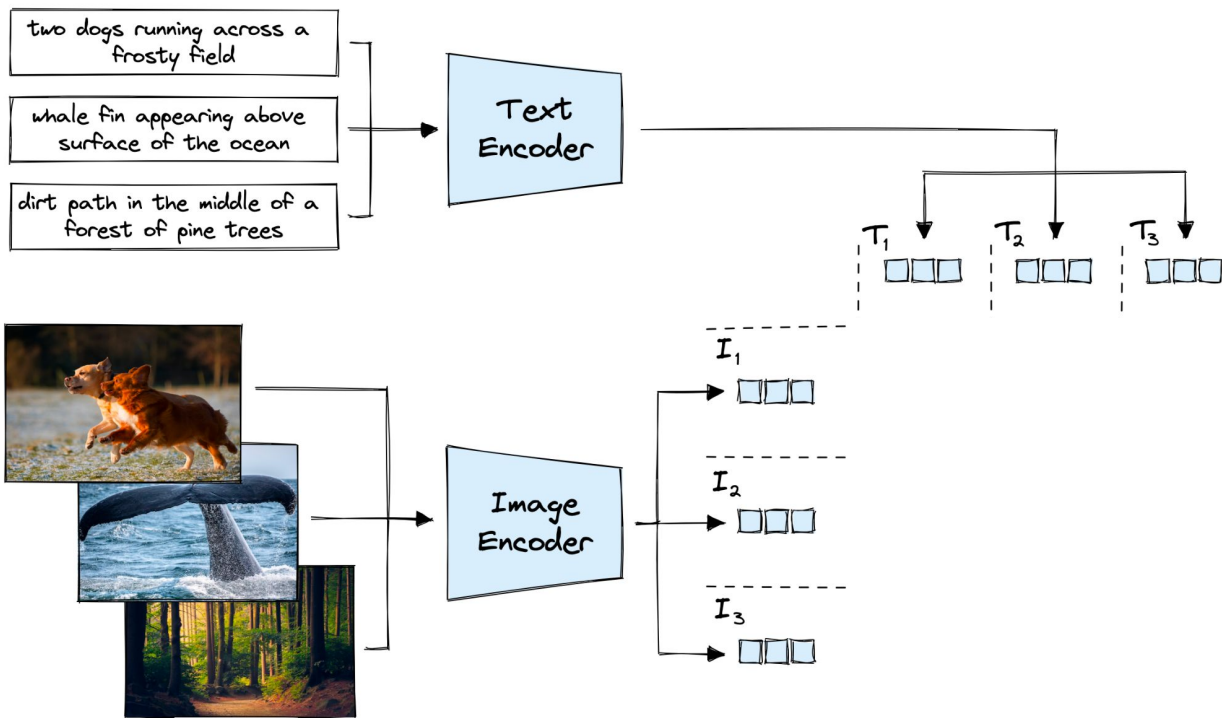
A horse carrying a large load of hay and two people sitting on it.

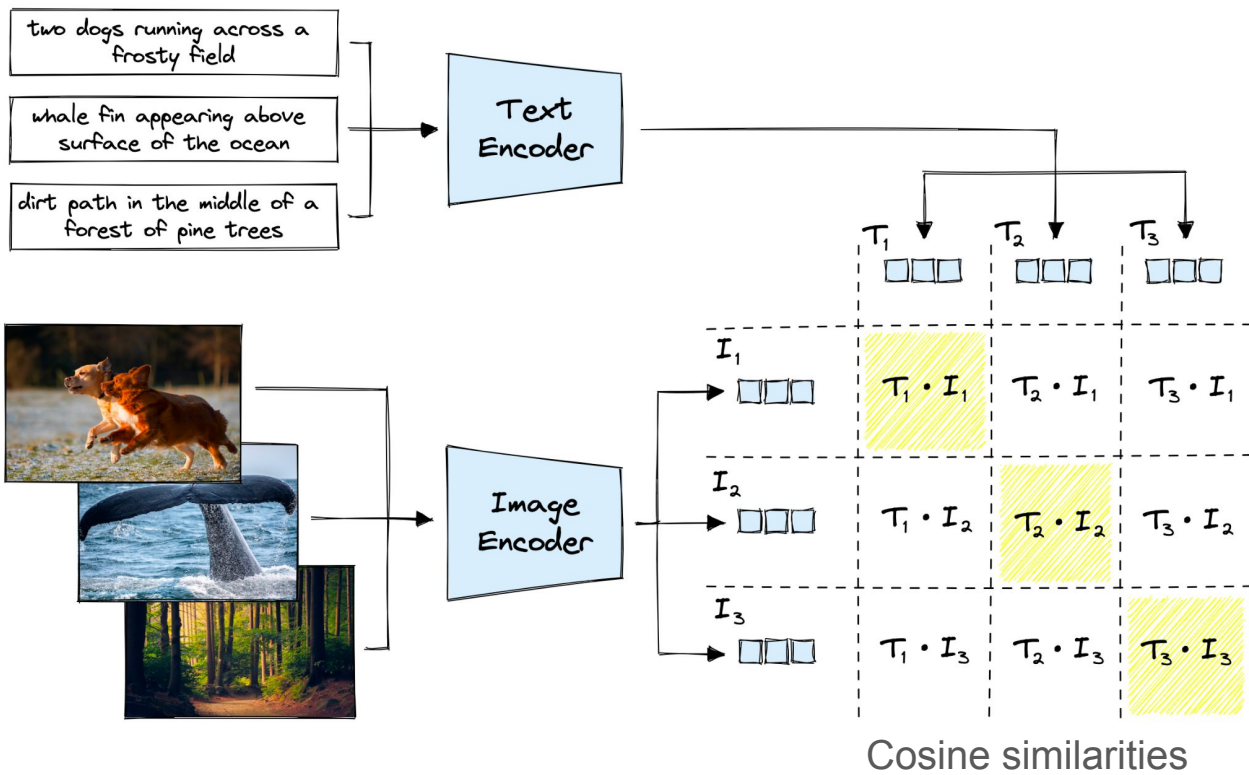


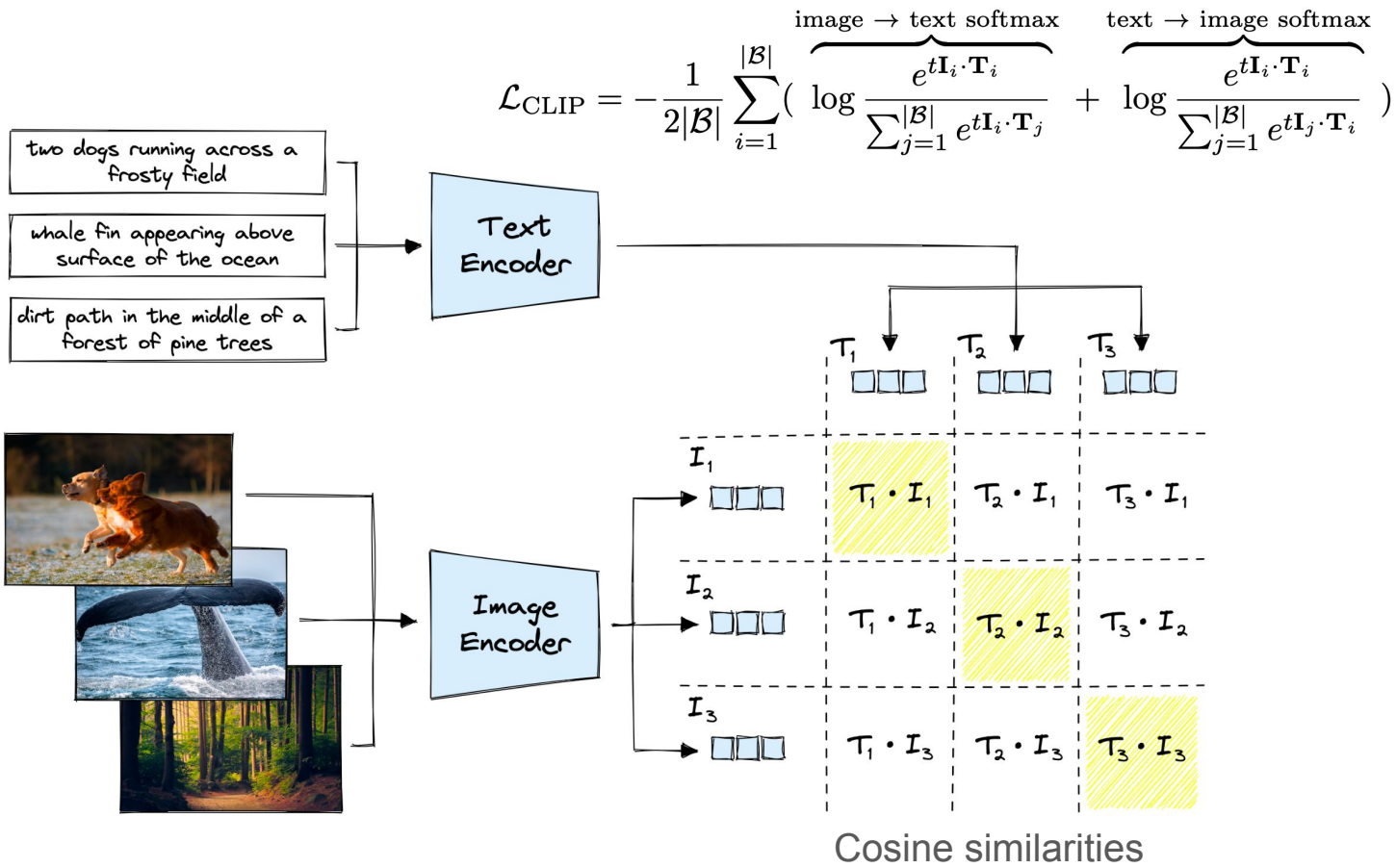
Bunk bed with a narrow shelf sitting underneath it.

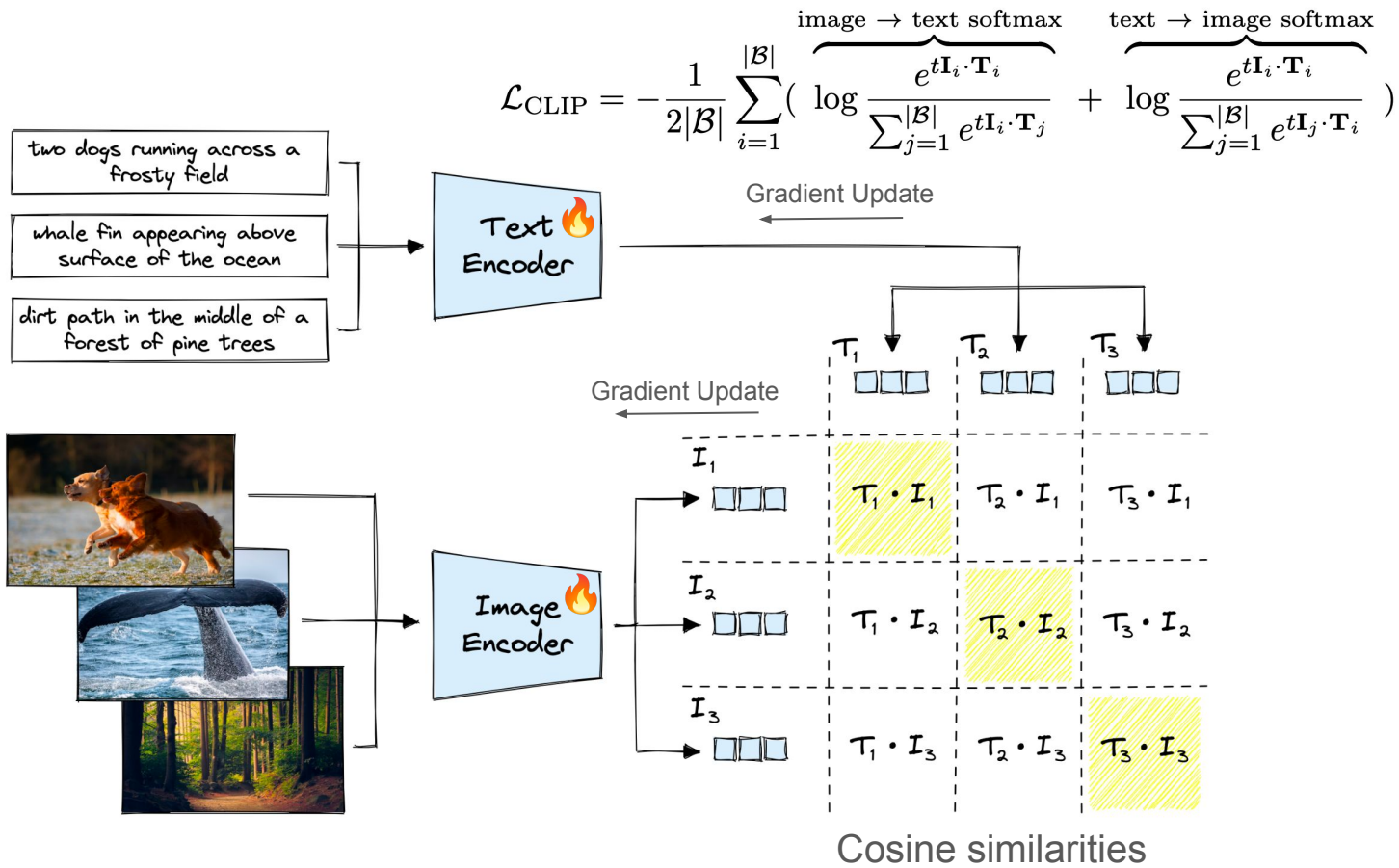
Two networks trained jointly

- Image encoder
- Text encoder





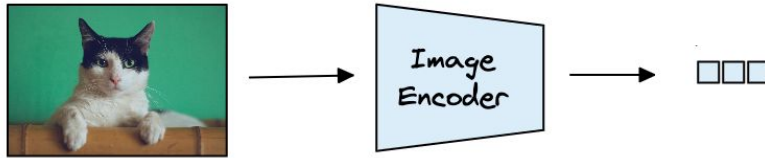




Free Zero-shot classification



Free Zero-shot classification

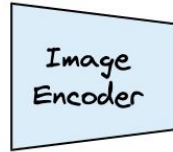


Free Zero-shot classification

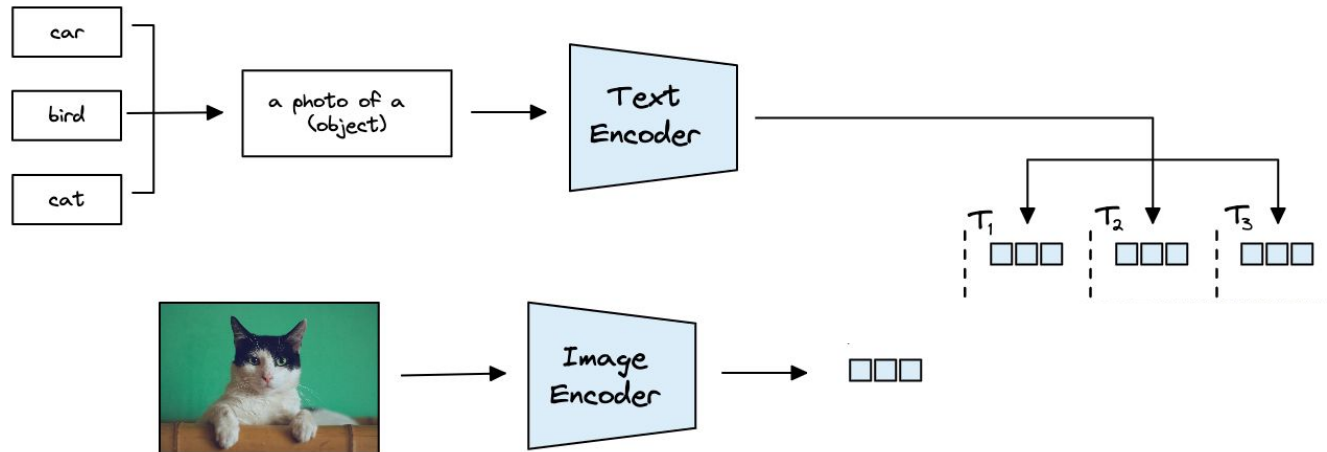
car

bird

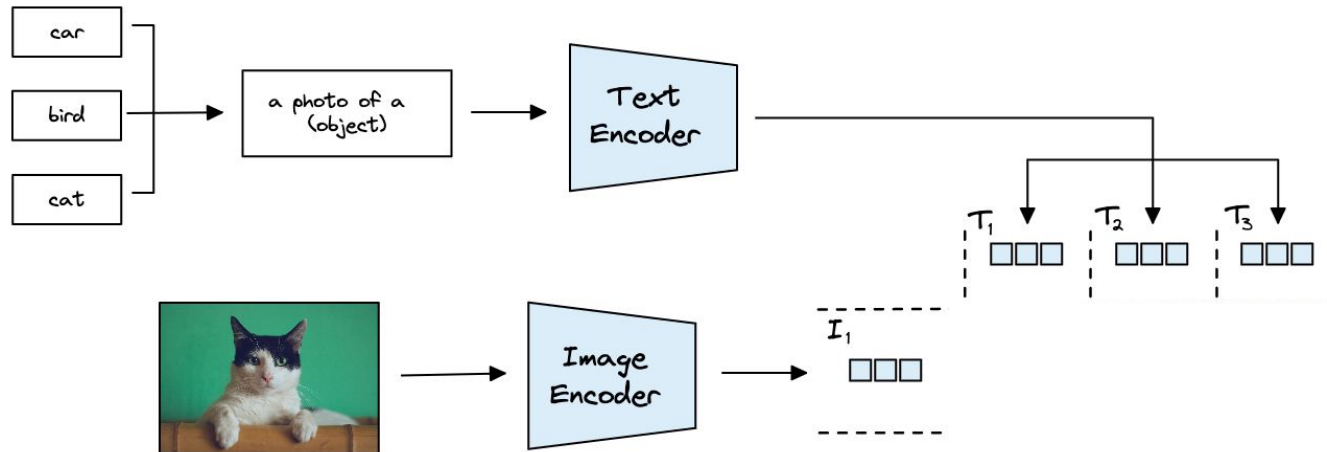
cat



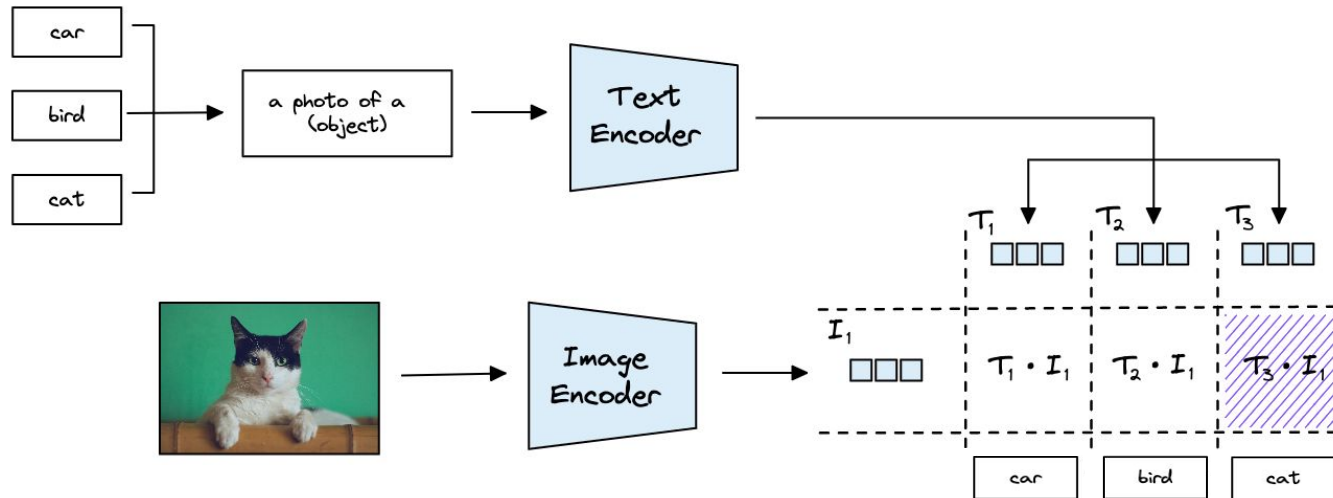
Free Zero-shot classification



Free Zero-shot classification

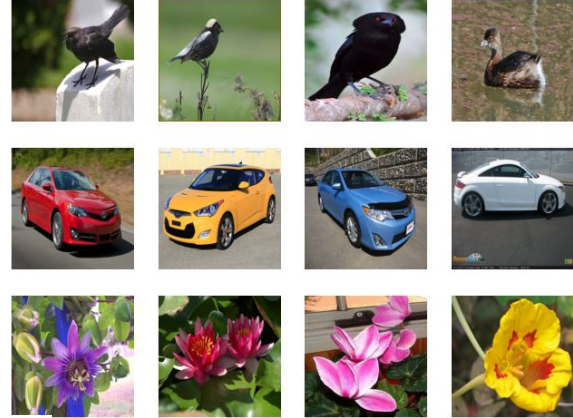


Free Zero-shot classification





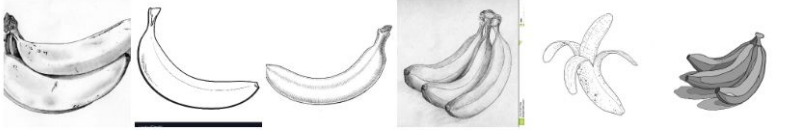


Perfect for Image retrieval

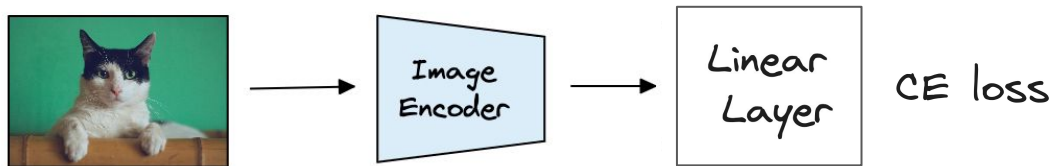
canard



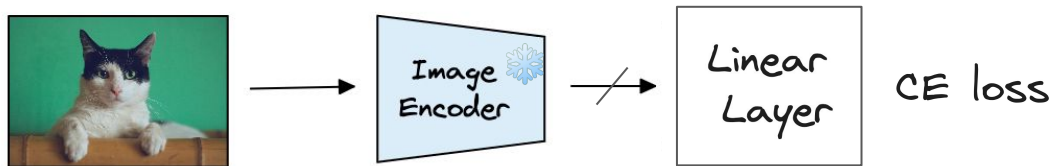
Performance

Dataset	ImageNet ResNet101	CLIP ViT-L
	76.2%	76.2%
ImageNet		
	64.3%	70.1%
ImageNet V2		
	37.7%	88.9%
ImageNet Rendition		
	32.6%	72.3%
ObjectNet		
	25.2%	60.2%
ImageNet Sketch		

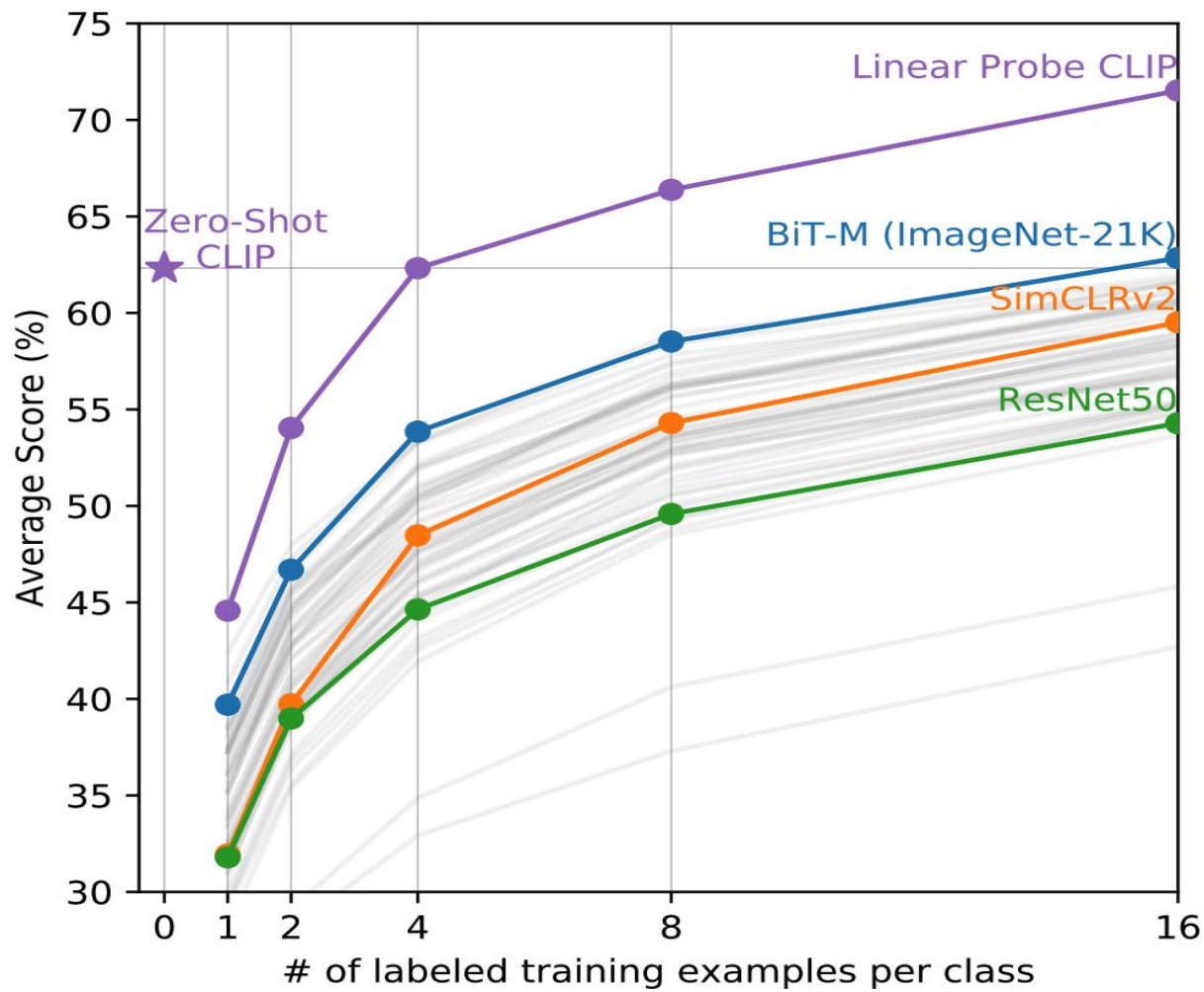
Linear Probing



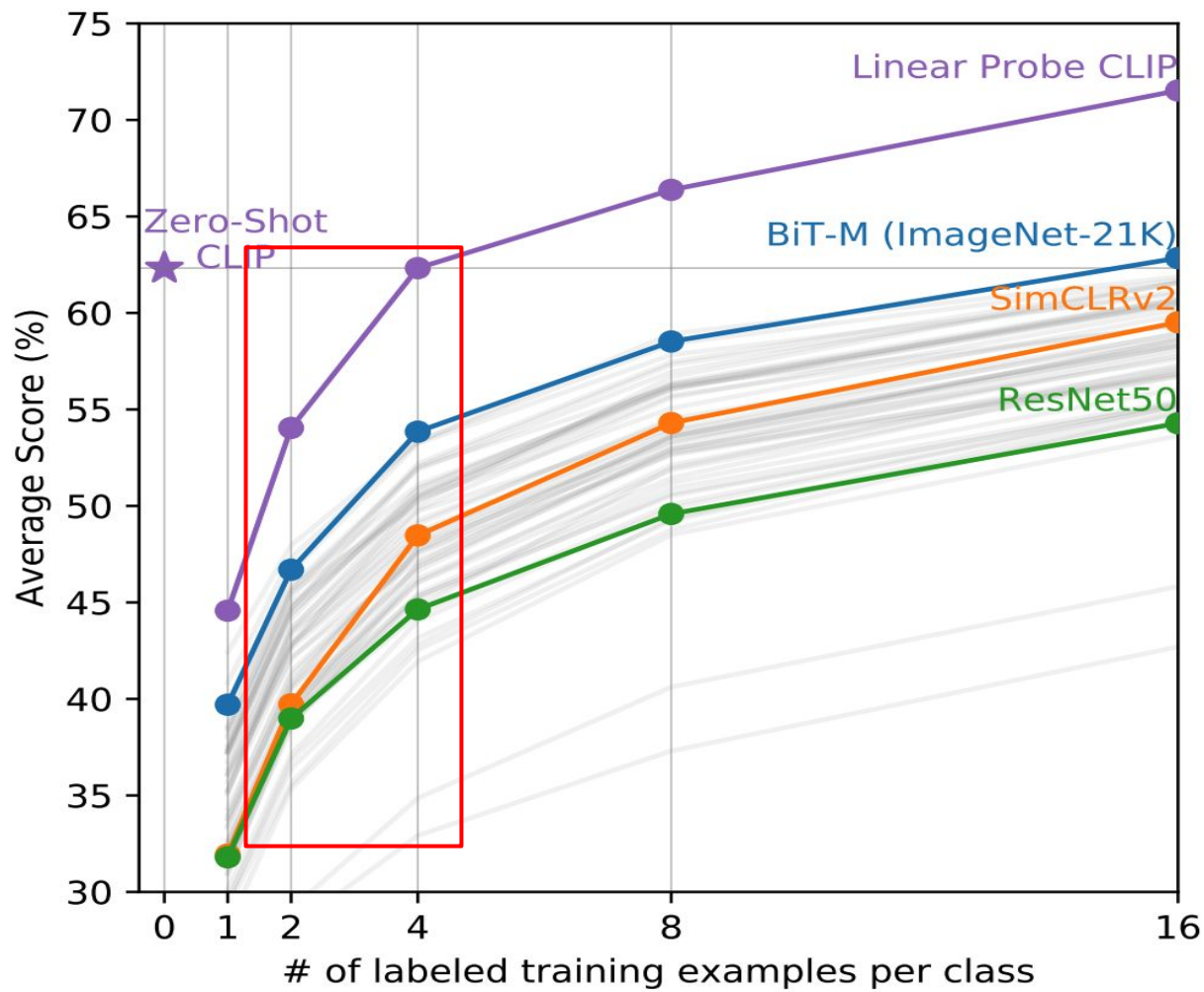
Linear Probing




Performance: Linear probing




Performance: Linear probing



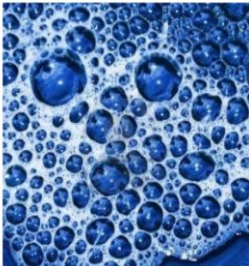
Prompt tuning

Caltech101	Prompt	Accuracy
	a [CLASS].	82.68
	a photo of [CLASS].	80.81
	a photo of a [CLASS].	86.29

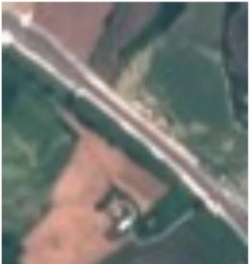
(a)

Flowers102	Prompt	Accuracy
	a photo of a [CLASS].	60.86
	a flower photo of a [CLASS].	65.81
	a photo of a [CLASS], a type of flower.	66.14

(b)


Describable Textures (DTD)	Prompt	Accuracy
	a photo of a [CLASS].	39.83
	a photo of a [CLASS] texture.	40.25
	[CLASS] texture.	42.32

(c)


EuroSAT	Prompt	Accuracy
	a photo of a [CLASS].	24.17
	a satellite photo of [CLASS].	37.46
	a centered satellite photo of [CLASS].	37.56

(d)

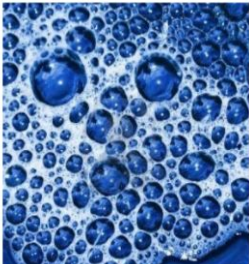
Prompt tuning

Caltech101	Prompt	Accuracy
	a [CLASS].	82.68
	a photo of [CLASS].	80.81
	a photo of a [CLASS].	86.29
	$[V]_1 [V]_2 \dots [V]_M$ [CLASS].	91.83

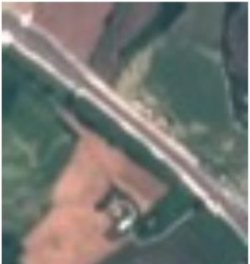
(a)

Flowers102	Prompt	Accuracy
	a photo of a [CLASS].	60.86
	a flower photo of a [CLASS].	65.81
	a photo of a [CLASS], a type of flower.	66.14
	$[V]_1 [V]_2 \dots [V]_M$ [CLASS].	94.51

(b)

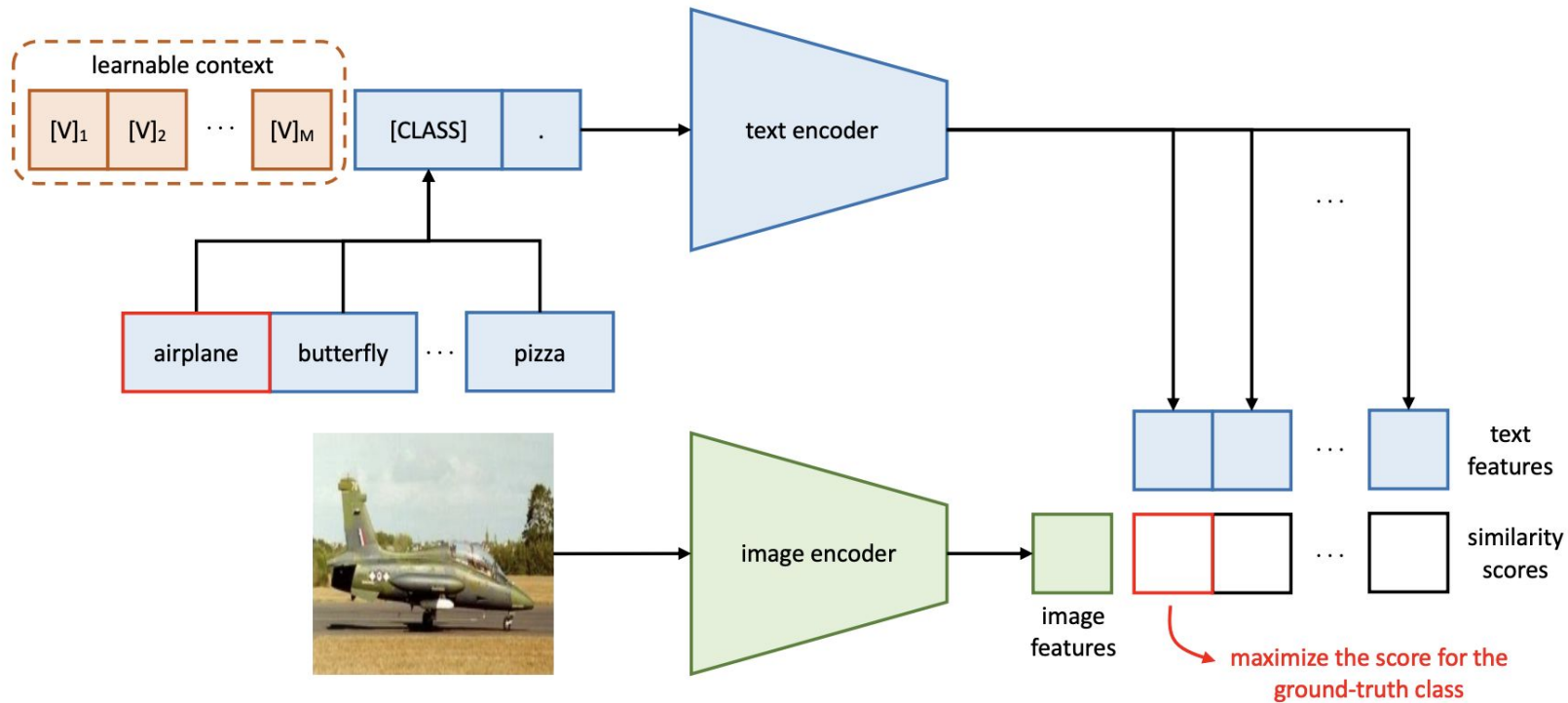
Describable Textures (DTD)	Prompt	Accuracy
	a photo of a [CLASS].	39.83
	a photo of a [CLASS] texture.	40.25
	[CLASS] texture.	42.32
	$[V]_1 [V]_2 \dots [V]_M$ [CLASS].	63.58

(c)

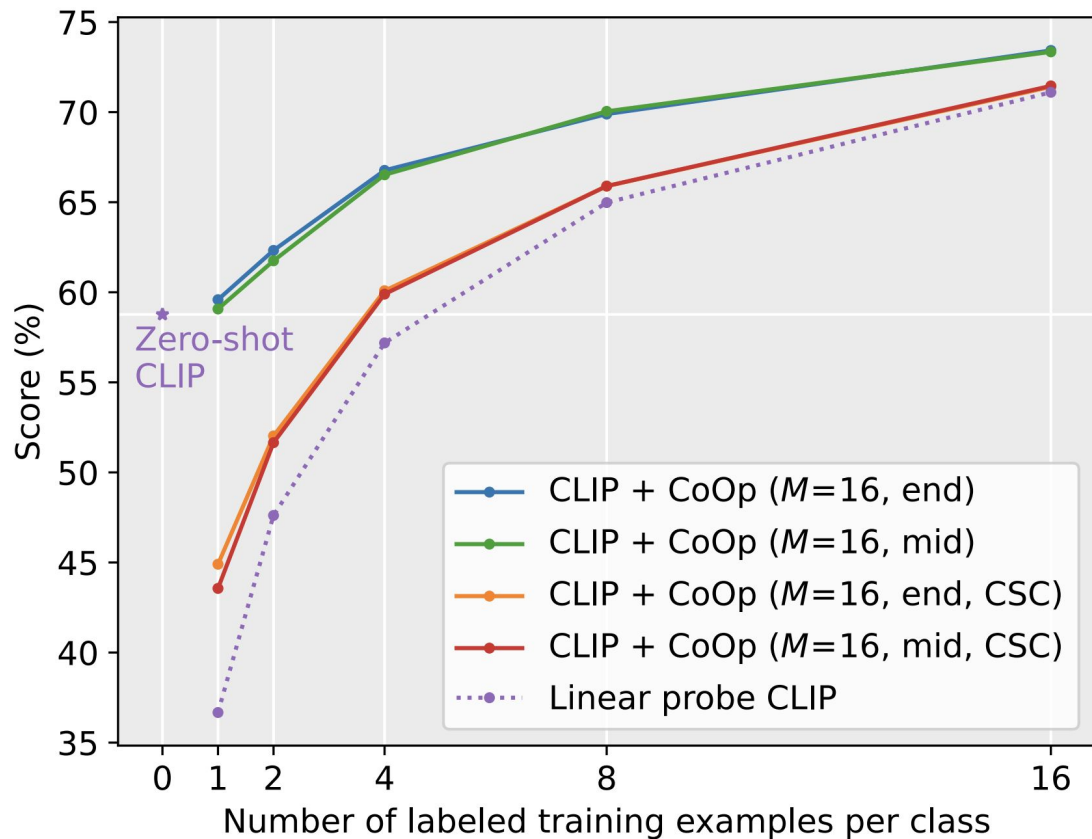
EuroSAT	Prompt	Accuracy
	a photo of a [CLASS].	24.17
	a satellite photo of [CLASS].	37.46
	a centered satellite photo of [CLASS].	37.56
	$[V]_1 [V]_2 \dots [V]_M$ [CLASS].	83.53

(d)

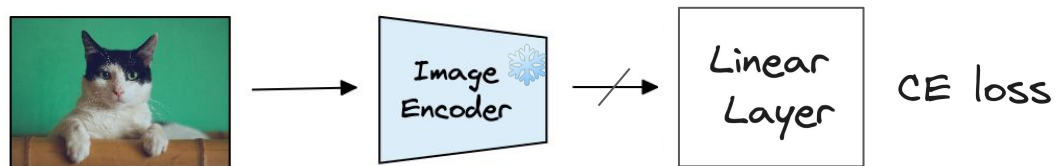
Prompt tuning



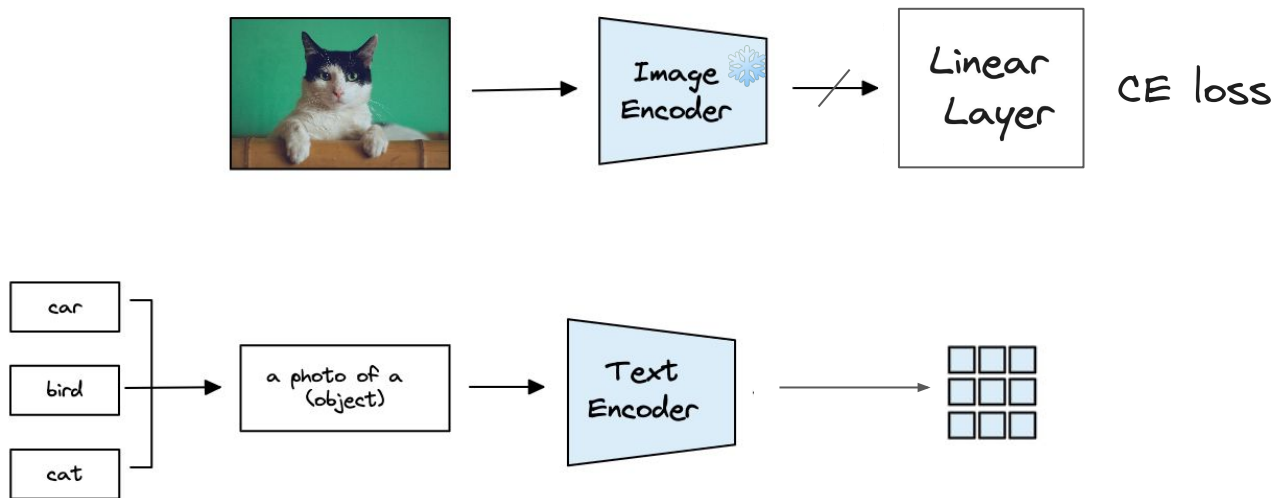
Prompt tuning



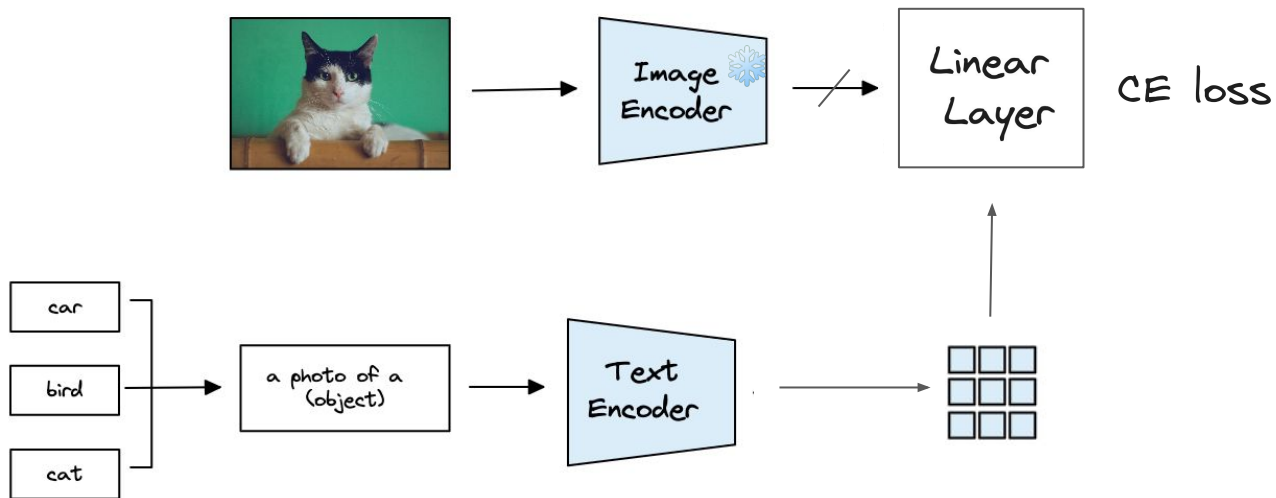
Few-shot Adaptation



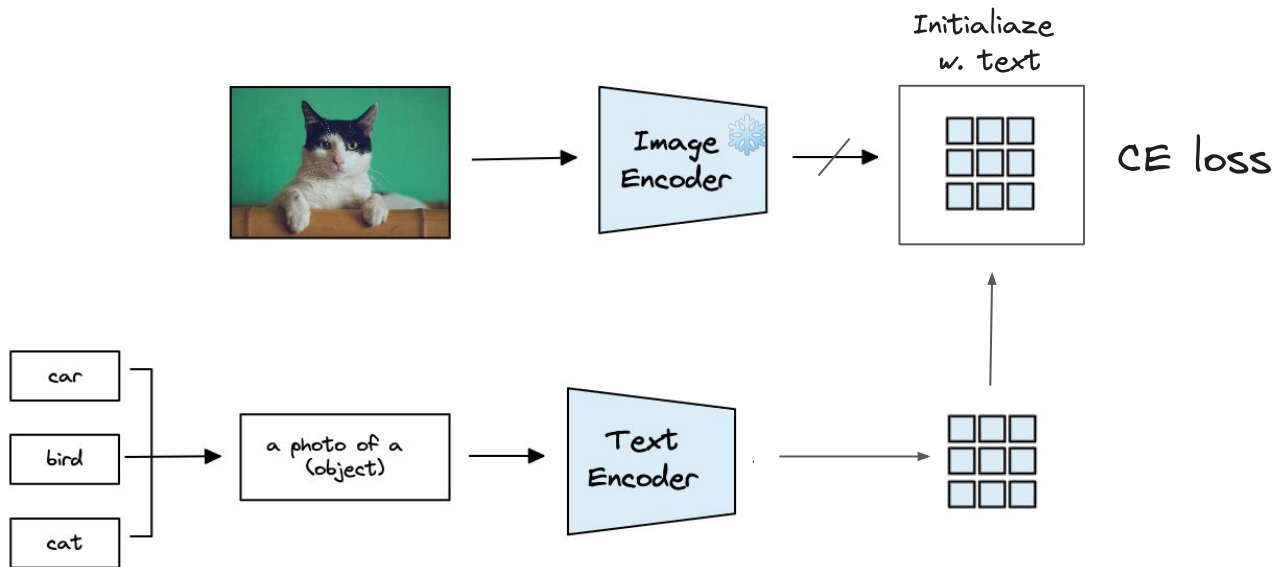
Few-shot Adaptation



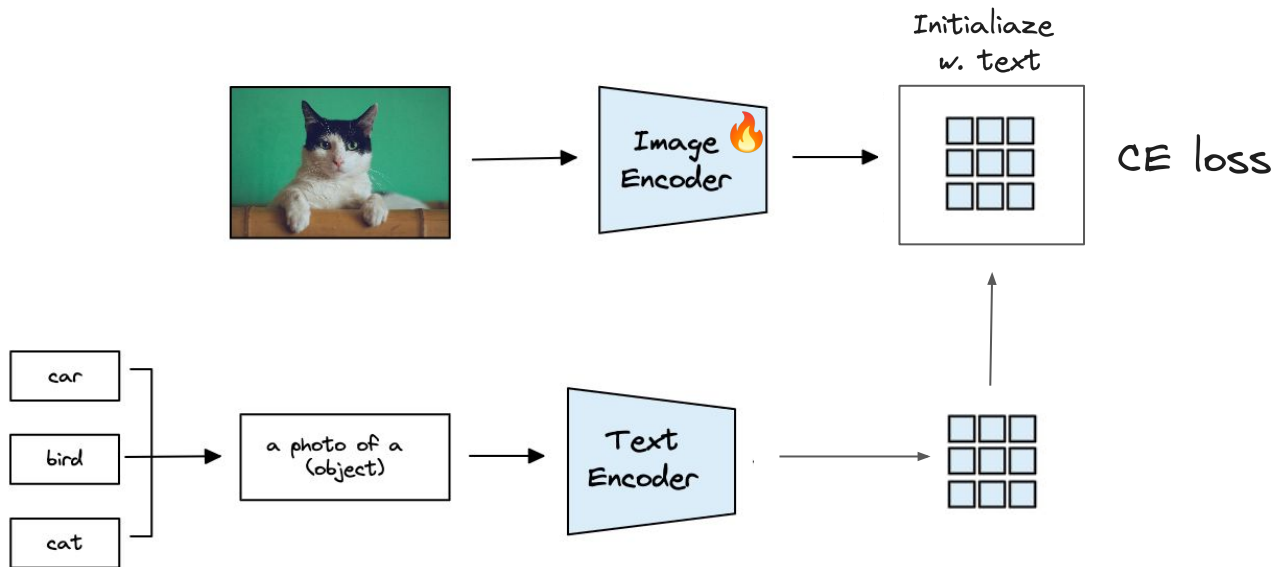
Few-shot Adaptation



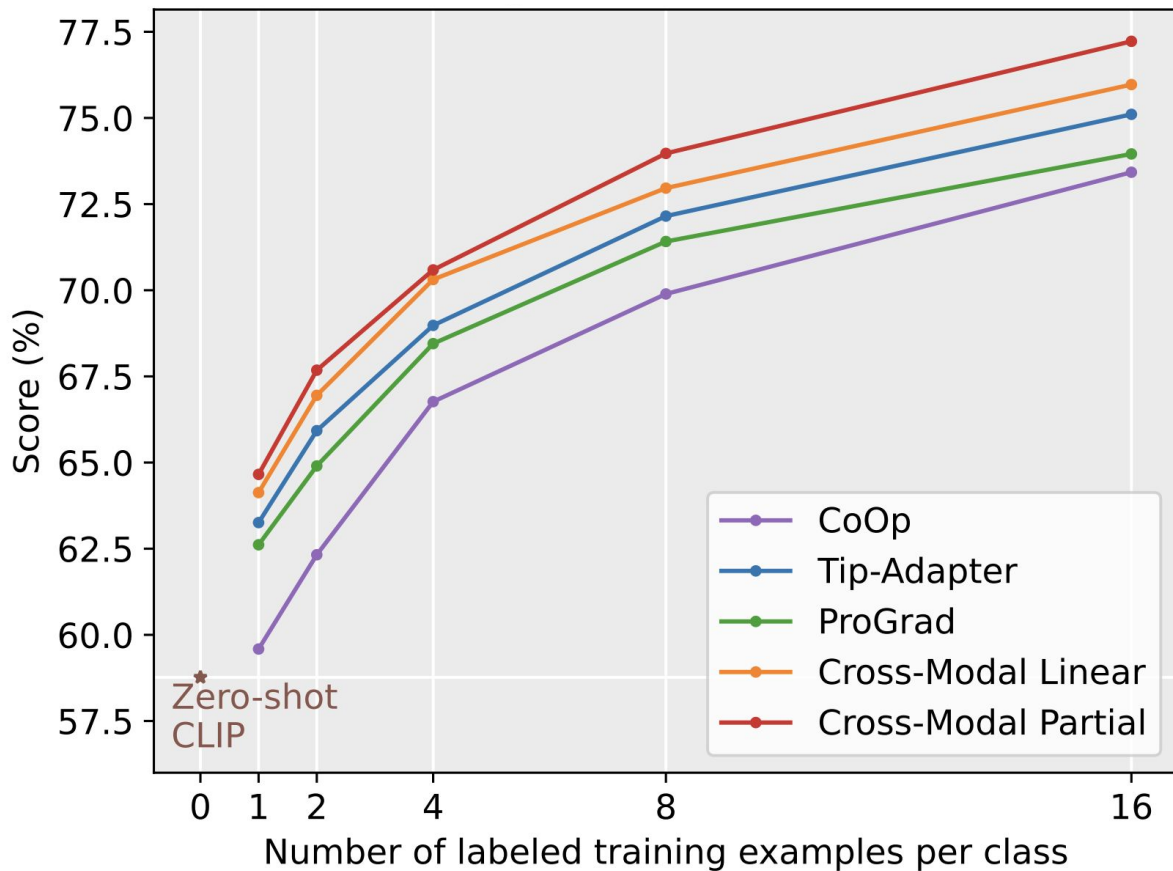
Few-shot Adaptation



Few-shot Adaptation



Multi-modal Few-shot



Better CLIPs

Open-CLIP

Reproducible scaling laws for contrastive language-image learning

Mehdi Cherti^{1,5} §§ Romain Beaumont¹ §§ Ross Wightman^{1,3} §§
Mitchell Wortsman⁴ §§ Gabriel Ilharco⁴ §§ Cade Gordon²
Christoph Schuhmann¹ Ludwig Schmidt⁴ °° Jenia Jitsev^{1,5} §§°°
LAION¹ UC Berkeley² HuggingFace³ University of Washington⁴
Juelich Supercomputing Center (JSC), Research Center Juelich (FZJ)⁵
contact@laion.ai, {m.cherti,j.jitsev}@fz-juelich.de
§§ Equal first contributions, °° Equal senior contributions

**LAION-5B: A NEW ERA OF
OPEN LARGE-SCALE MULTI-
MODAL DATASETS**



Similar performance
to OpenAI



Model	Training data	Resolution	ImageNet zero-shot acc.
ViT-B/32	DataComp-1B	256px	72.8%
ViT-B/16	DataComp-1B	224px	73.5%
ViT-L/14	LAION-2B	224px	75.3%
ViT-H/14	LAION-2B	224px	78.0%
ViT-L/14	DataComp-1B	224px	79.2%
ViT-G/14	LAION-2B	224px	80.1%
ViT-L/14	OpenAI's WIT	224px	75.5%

OpenAI



ImageBind

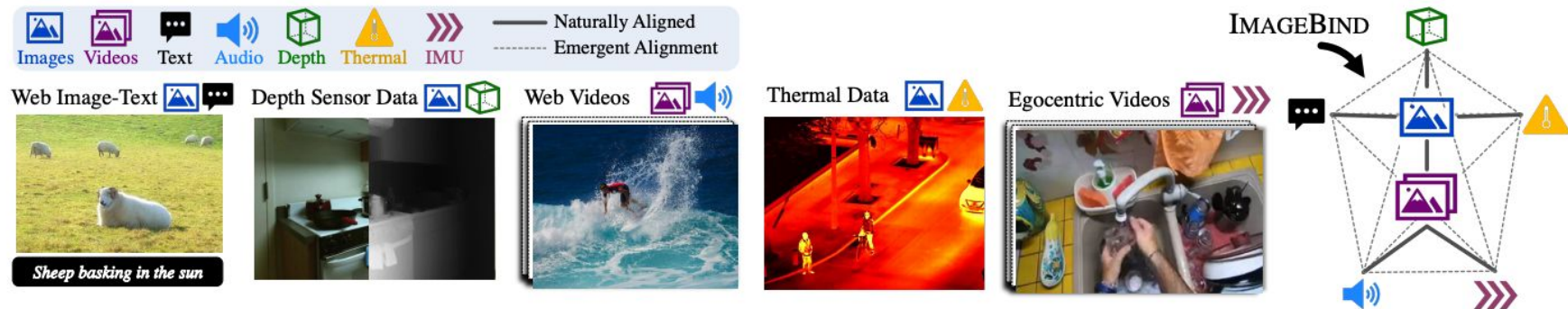
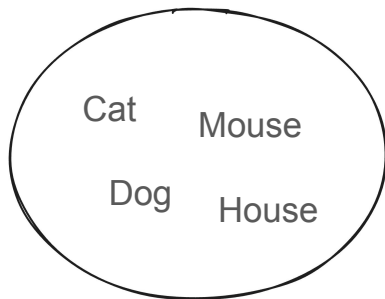


Figure 2. IMAGEBIND overview. Different modalities occur naturally aligned in different data sources, for instance images+text and video+audio in web data, depth or thermal information with images, IMU data in videos captured with egocentric cameras, *etc.* IMAGEBIND links all these modalities in a common embedding space, enabling new emergent alignments and capabilities.

Open-vocabulary models using Text

1. Classification



Closed set of
classes



All vocabulary

2. Object Detection




Detection Prompt: black cat

3. Image Segmentation



Visual Question Answering (VQA)



This is a picture of two teddy bears on the moon.

What are they doing?

They are having a conversation.

What object are they using?

It looks like a computer.

Is this surprising?

Yes, it is surprising.

Why is this picture surprising to you?

I think it is surprising because teddy bears are not usually found on the moon.

