

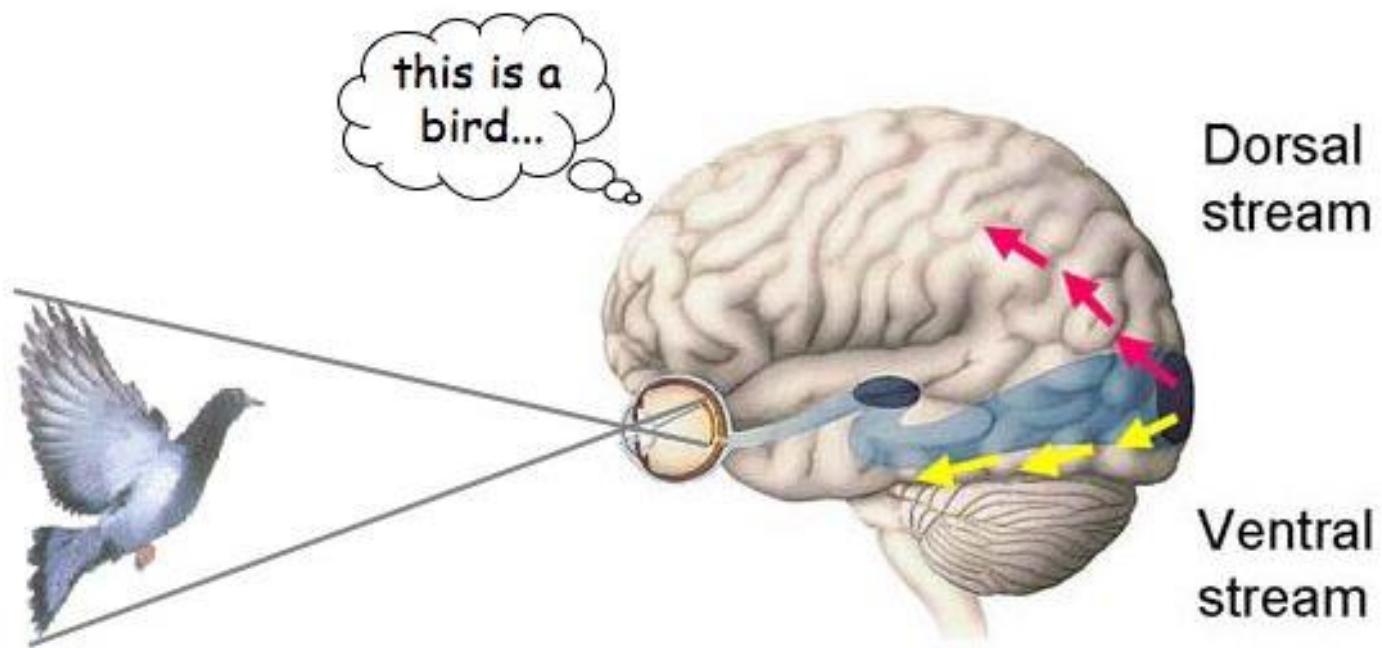


# Vision-Language Learning for Visual Recognition

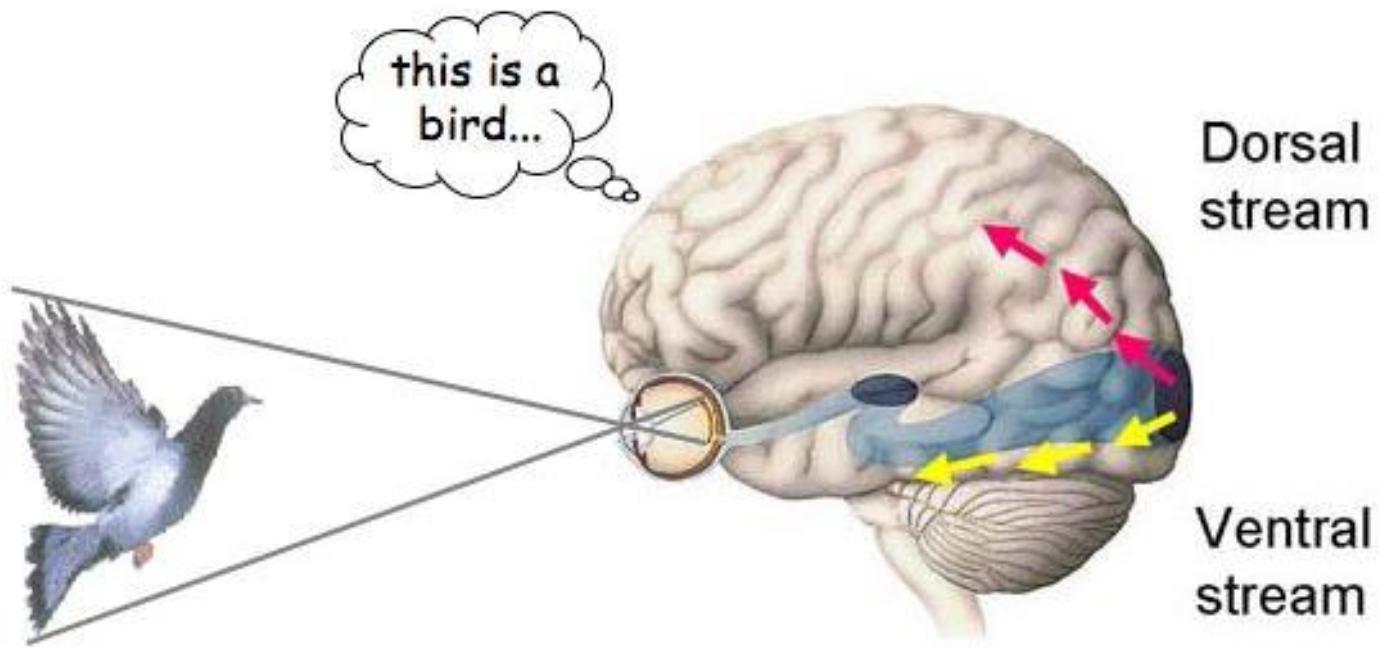
Jianwei Yang  
Microsoft Research



# What is visual recognition?



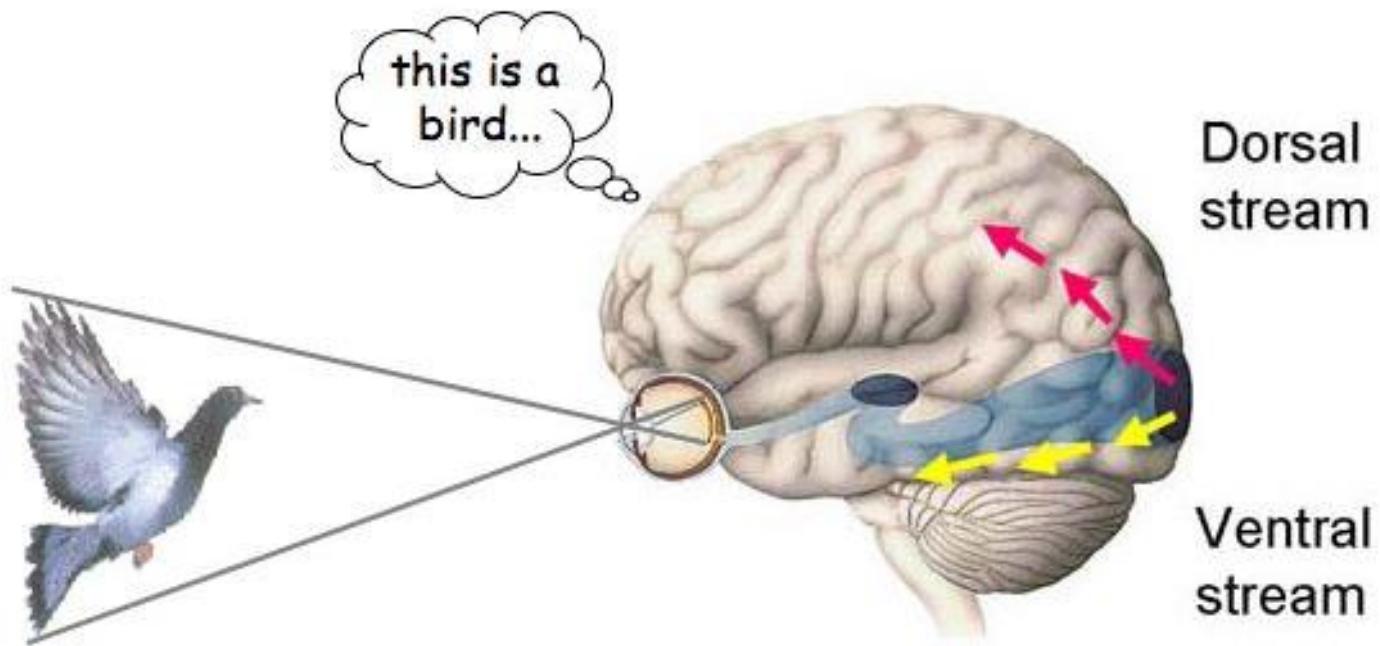
# What is visual recognition?



Visual object recognition serves as a gateway **from vision to cognitive processes such as categorization, language and reasoning.**

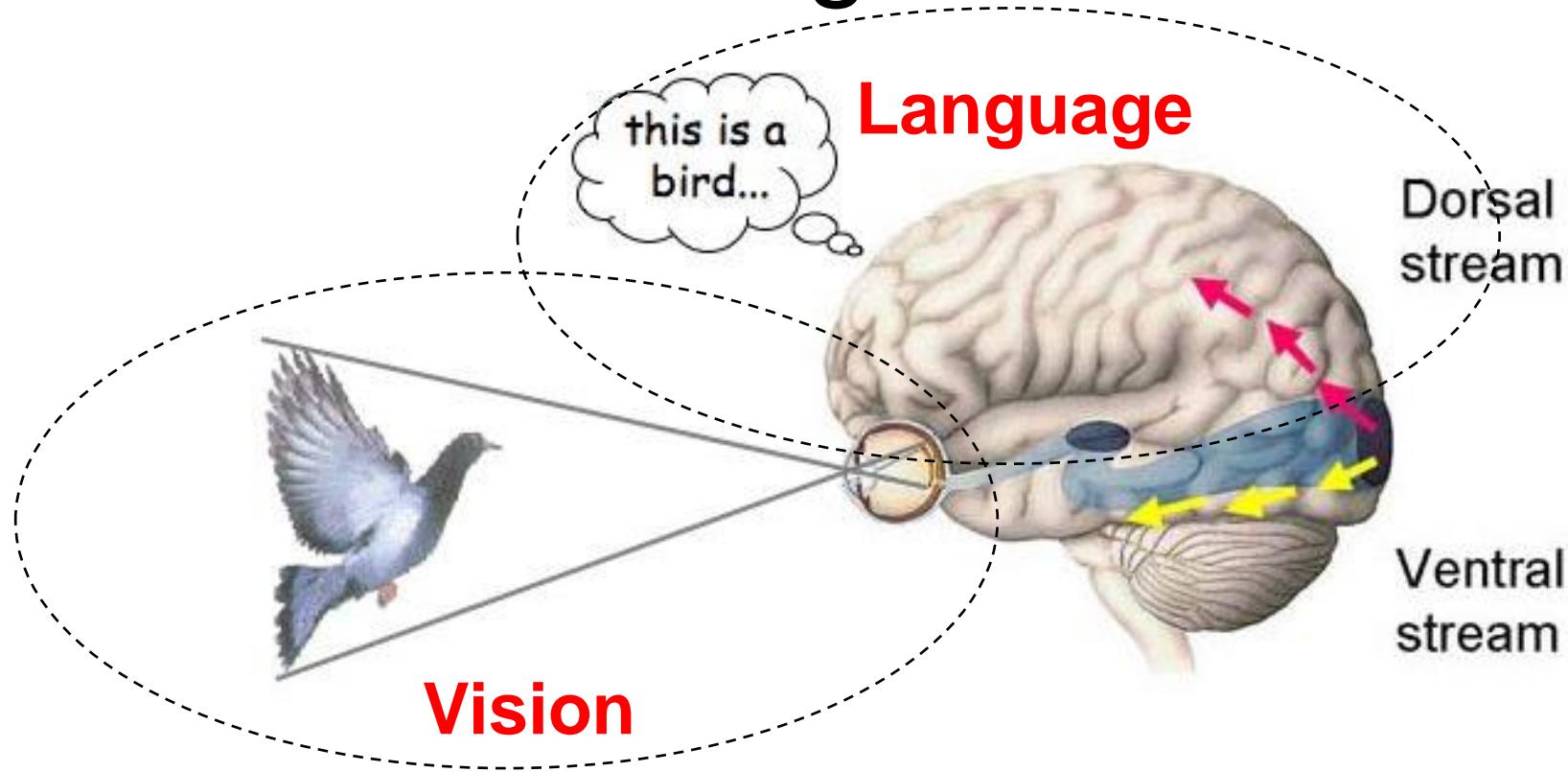
- *Object Recognition*. John E. Hummel

# What is visual recognition?



**Visual Recognition:** Observe visual object  
and map it with semantic concept

# What is visual recognition?



**Visual Recognition:** Observe visual object  
and map it with semantic concept

# How model learns visual recognition?

# Supervised Learning

Map an image to a discrete label  
which is associated a visual concept

Image



Label (Concept)



“2” (Apple)

# Supervised Learning



MNIST. LeCun *et al.*



CIFAR-10. Krizhevsky *et al.*

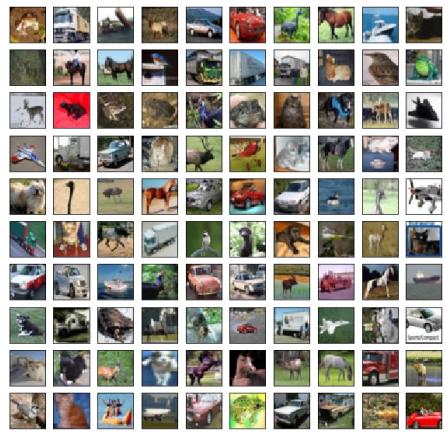


ImageNet. Deng *et al.*

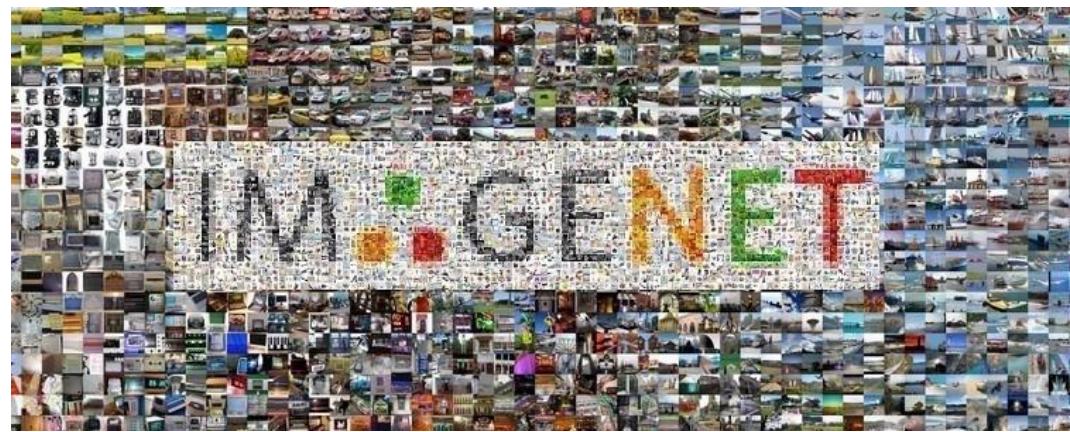
# Supervised Learning



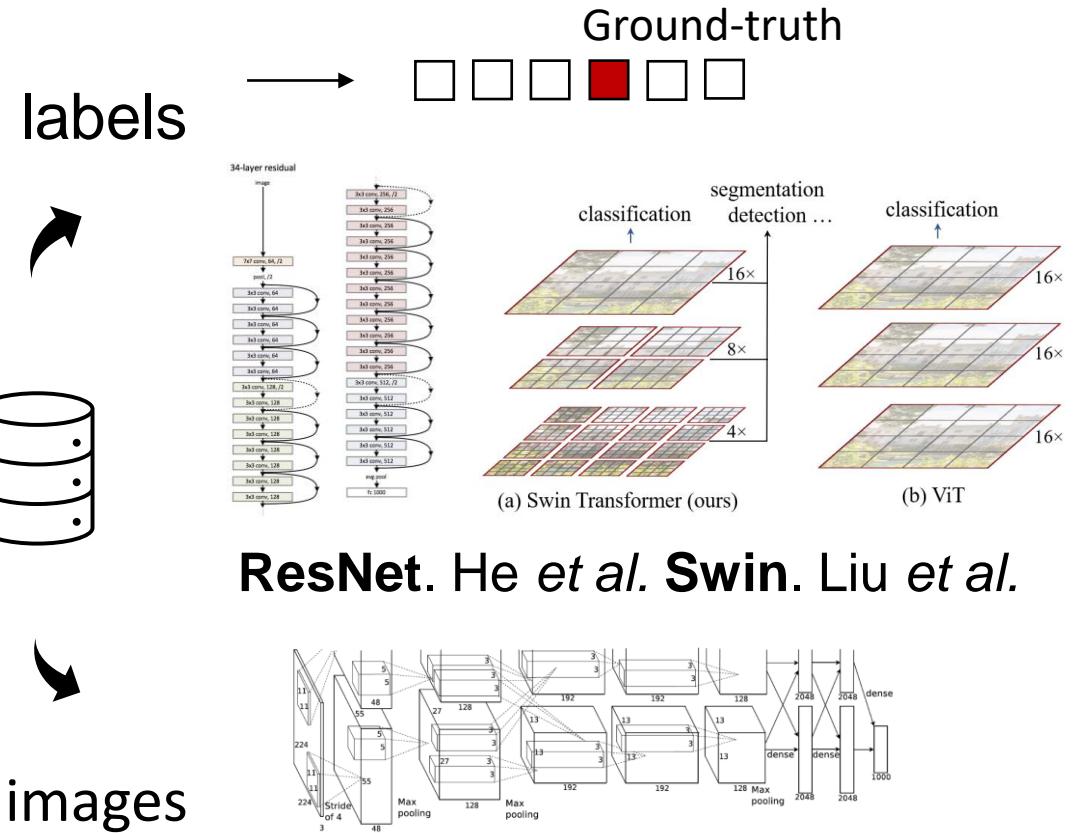
MNIST. LeCun *et al.*



CIFAR-10. Krizhevsky *et al.*



# ImageNet. Deng et al.



**AlexNet.** Krizhevsky *et al.*

# Supervised Learning

- Pros
  - Densely labeled samples for each category
- Cons
  - Requires a lot of human effort
  - Limited number of categories

# Zero-Shot Learning (Canonical)

Map an image to description of a visual concept

Image



Descriptions (Concept)

Fruit, Red, Sphere (Apple)



Fruit, Yellow (Orange)

# Zero-Shot Learning (Canonical)



CUB-200-2011. Wah et al.

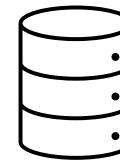


AwA2. Xian et al.

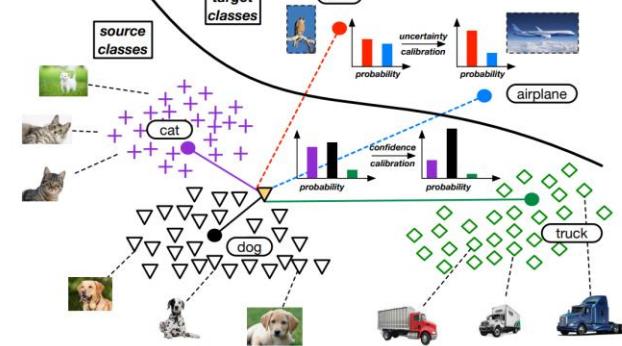


aPY. Farhadi et al.

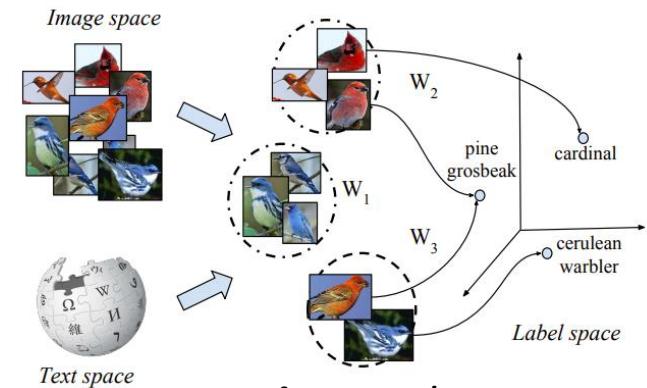
Label &  
descriptions



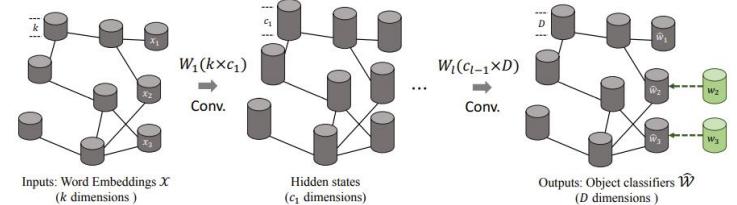
images



Liu et al.



Xian et al.



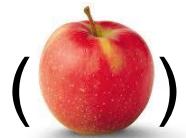
Wang et al.

# Zero-Shot Learning (Canonical)

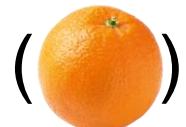
- Pros
  - Directly learn the visual-semantic matching
- Cons
  - Small scale with limited vocabulary
  - Fixed visual and text encoder

# Visual-Semantic Space Learning

Image ← → Align Text



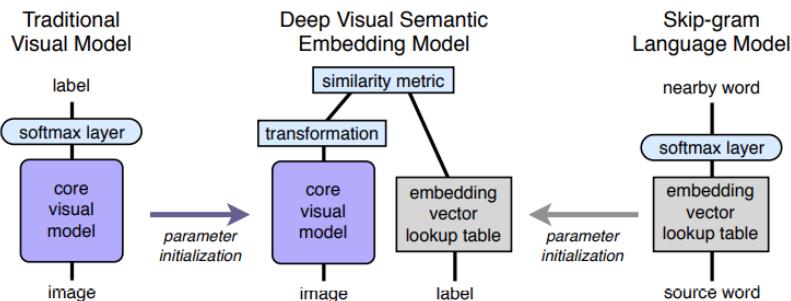
A red apple



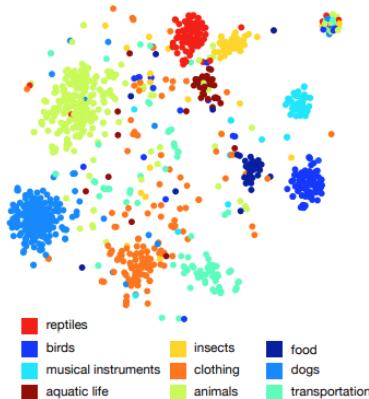
One yellow orange

# Visual-Semantic Space Learning

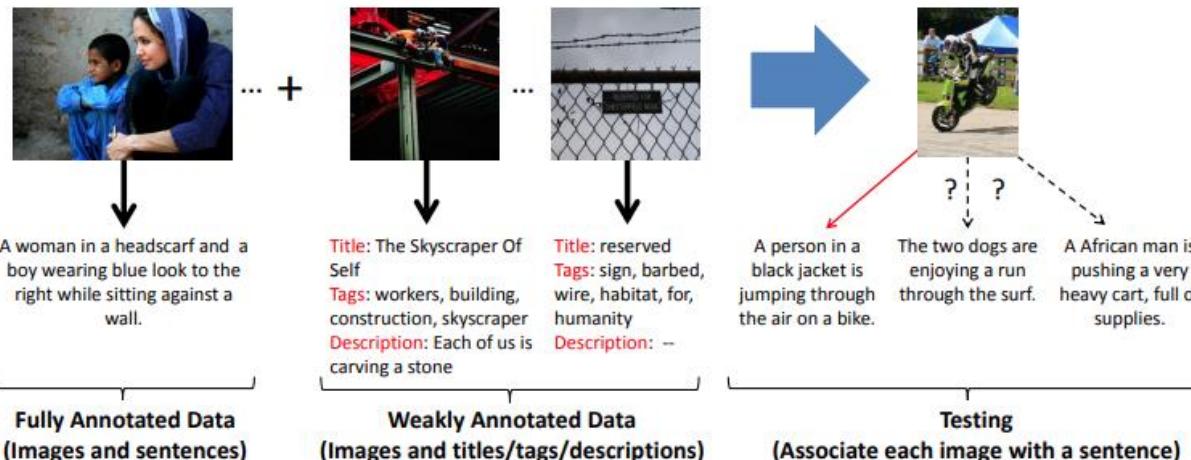
A



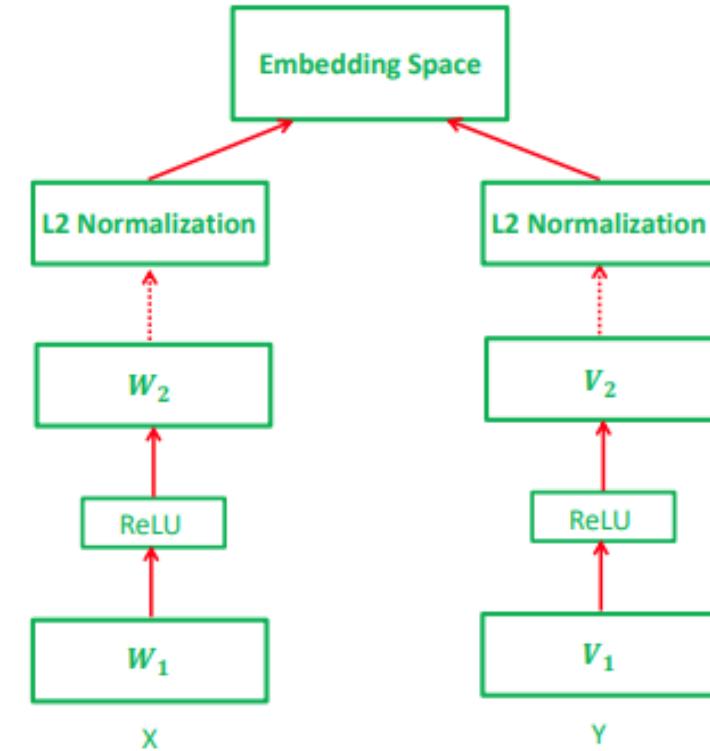
B



**DeViSE. Frome *et al.***



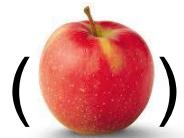
**Improving Image-Sentence Embeddings Using Large Weakly Annotated Photo Collections. Gong *et al.***



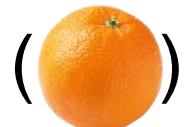
**Learning Deep Structure-Preserving Image-Text Embeddings. Wang *et al.***

# Visual-Semantic Space Learning

Image ← → Align Text



A red apple

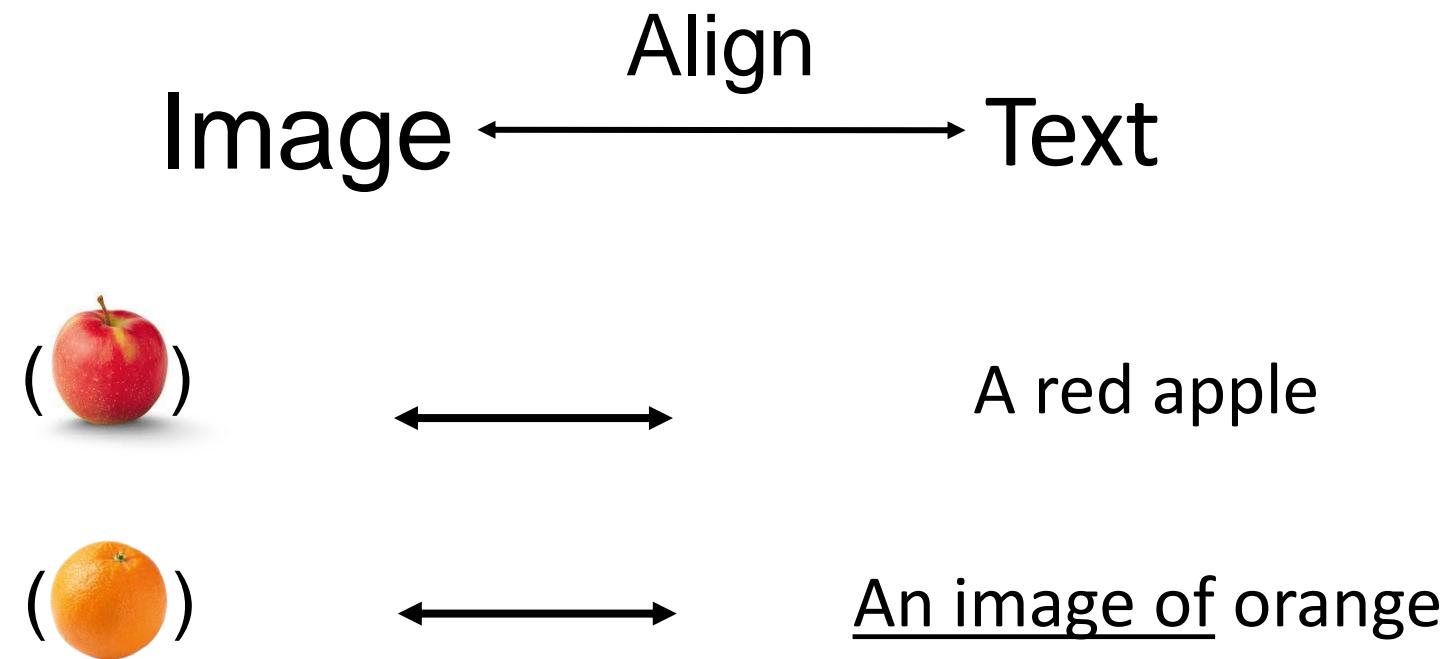


One yellow orange

# Visual-Semantic Space Learning

- Pros
  - Directly learn the visual-semantic matching
- Cons
  - Fixed visual and text encoder
  - Not showing the capacity for visual recognition

# Contrastive Vision-Language Learning



# Contrastive Vision-Language Learning

Image  $\xleftarrow{\text{Align}}$  Text

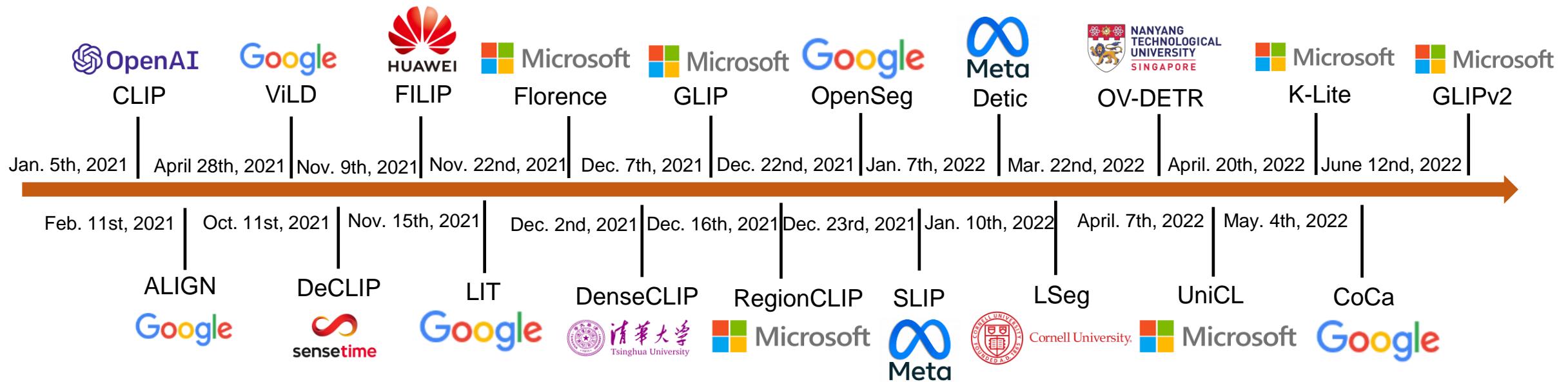
End-to-end learning on  
large-scale corpus



An image of orange

# The most recent art

## Contrastive Vision-Language Learning



A lot of research works come along the line of vision-language learning for vision

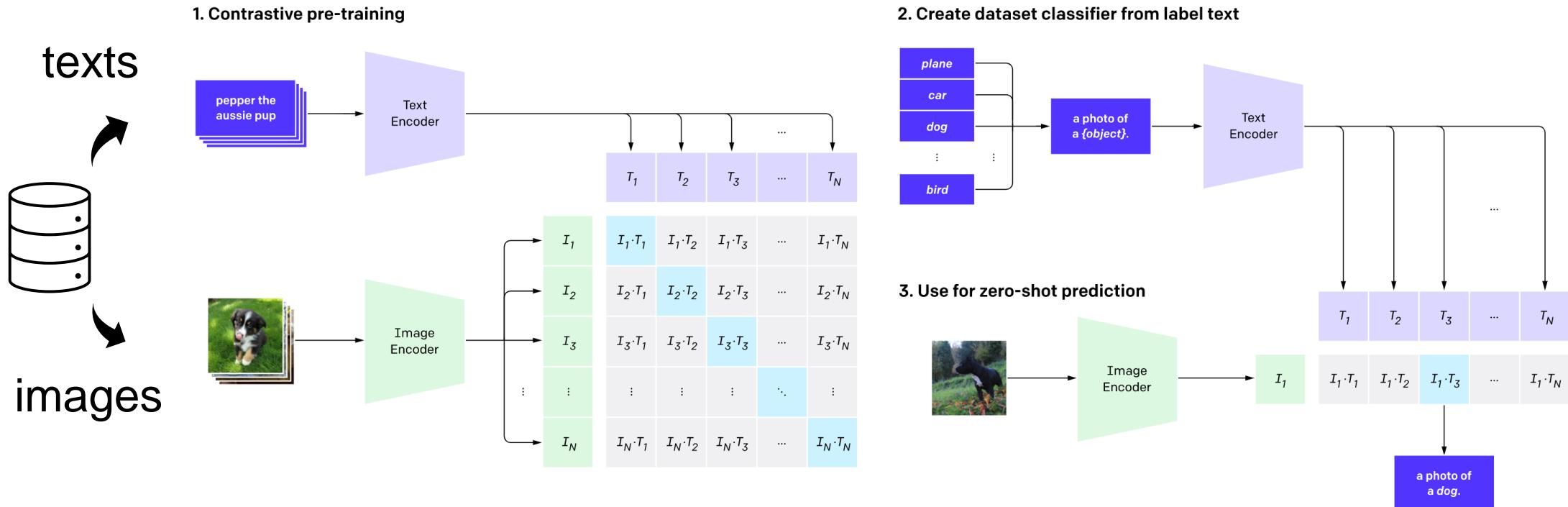
# In this tutorial

- Vision-Language Learning for Image-Level Recognition
- Vision-Language Learning for Visual Region-Level Recognition

# In this tutorial

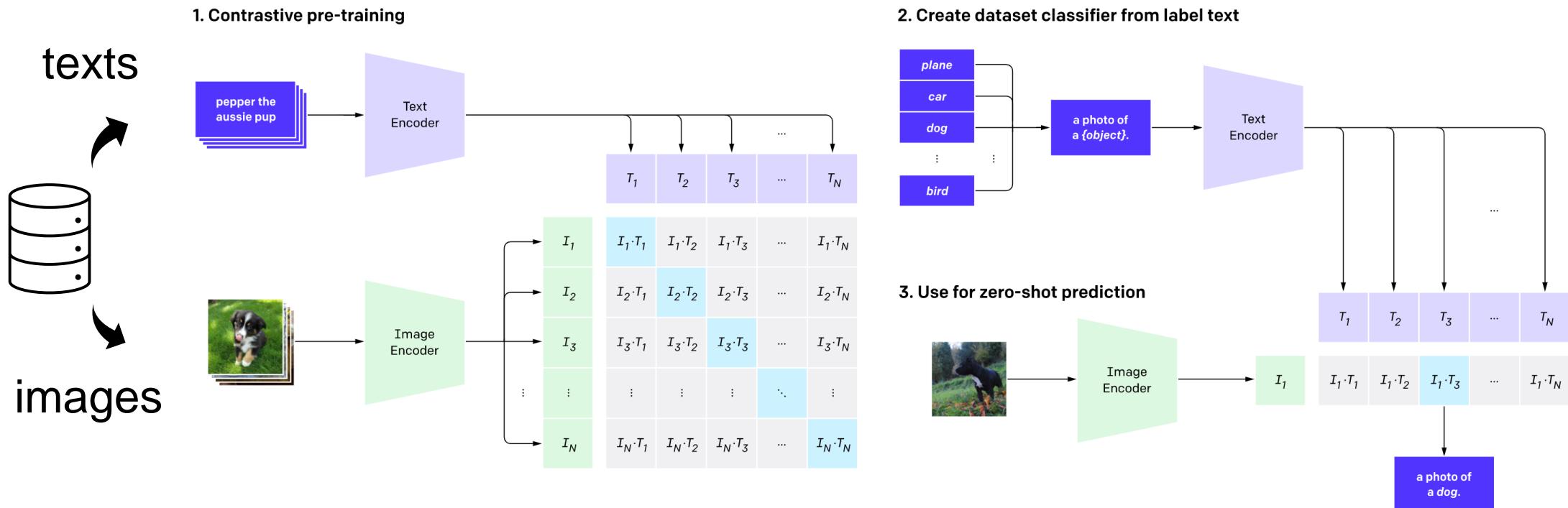
- Vision-Language Learning for **Image-Level** Recognition
- Vision-Language Learning for Visual **Region-Level** Recognition

# Contrastive Language-Image Learning (CLIP)



- ❖ Learning from 400M web-crawled image-text pairs
- ❖ Image recognition as an image-text matching problem

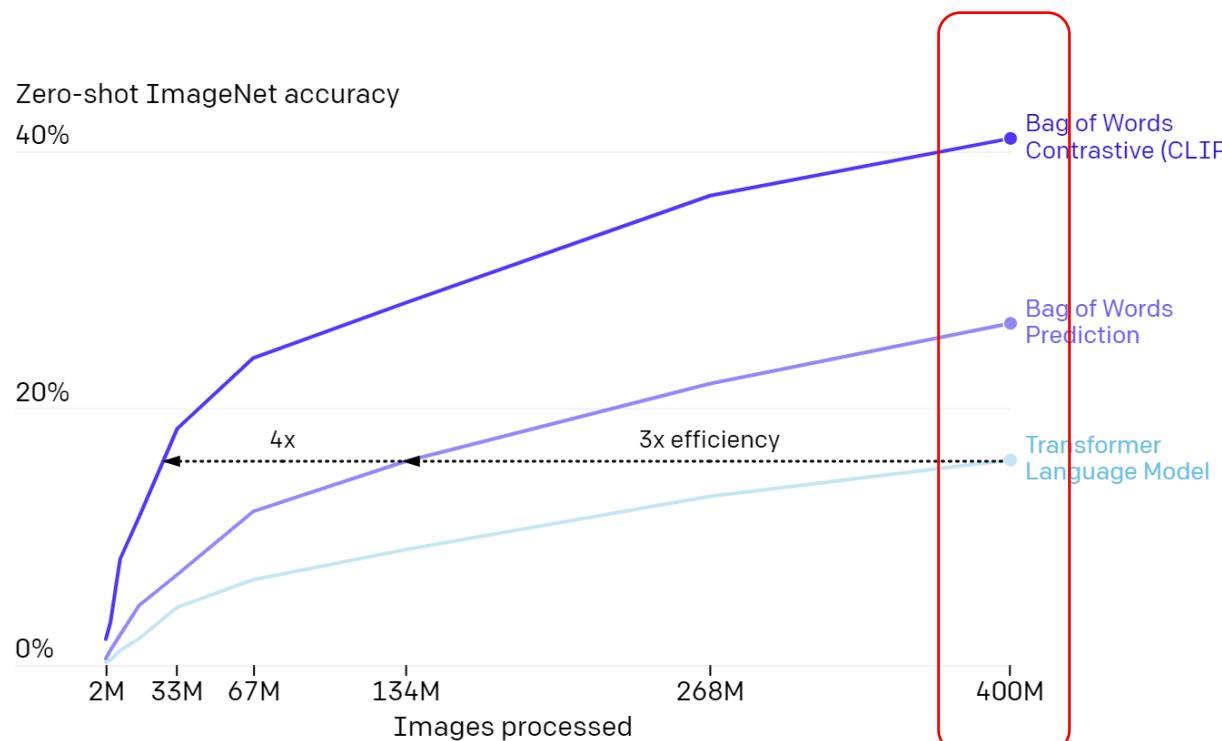
# Contrastive Language-Image Learning (CLIP)



$$L_{i2t} = - \sum_j \log \frac{\exp I_j T_j^T}{\sum_k \exp I_j T_k^T}$$

$$L_{t2i} = - \sum_j \log \frac{\exp I_j T_j^T}{\sum_k \exp I_k T_j^T}$$

# Contrastive Language-Image Learning (CLIP)



Contrastive learning is more effective than generative learning

Data really matter

**CLIP.** Radford *et al.*

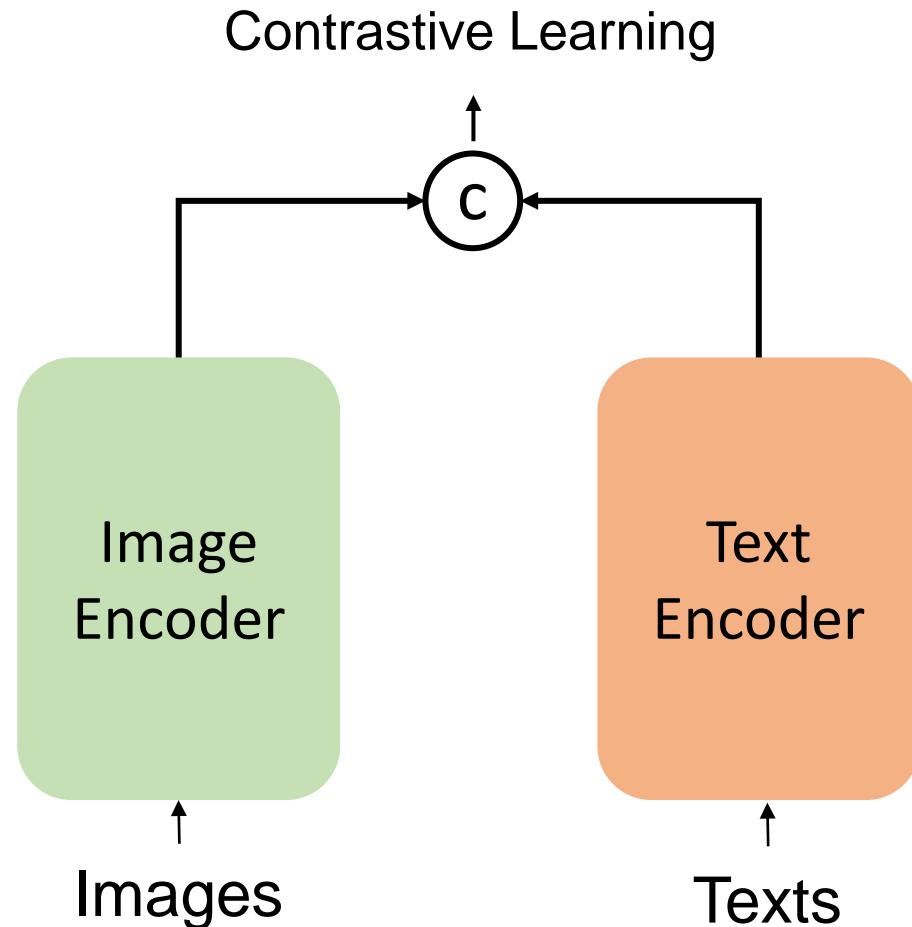
DATASET	IMAGENET RESNET101	CLIP VIT-L
ImageNet	76.2%	76.2%
ImageNet V2	64.3%	70.1%
ImageNet Rendition	37.7%	88.9%
ObjectNet	32.6%	72.3%
ImageNet Sketch	25.2%	60.2%
ImageNet Adversarial	2.7%	77.1%

Zero-shot transferring

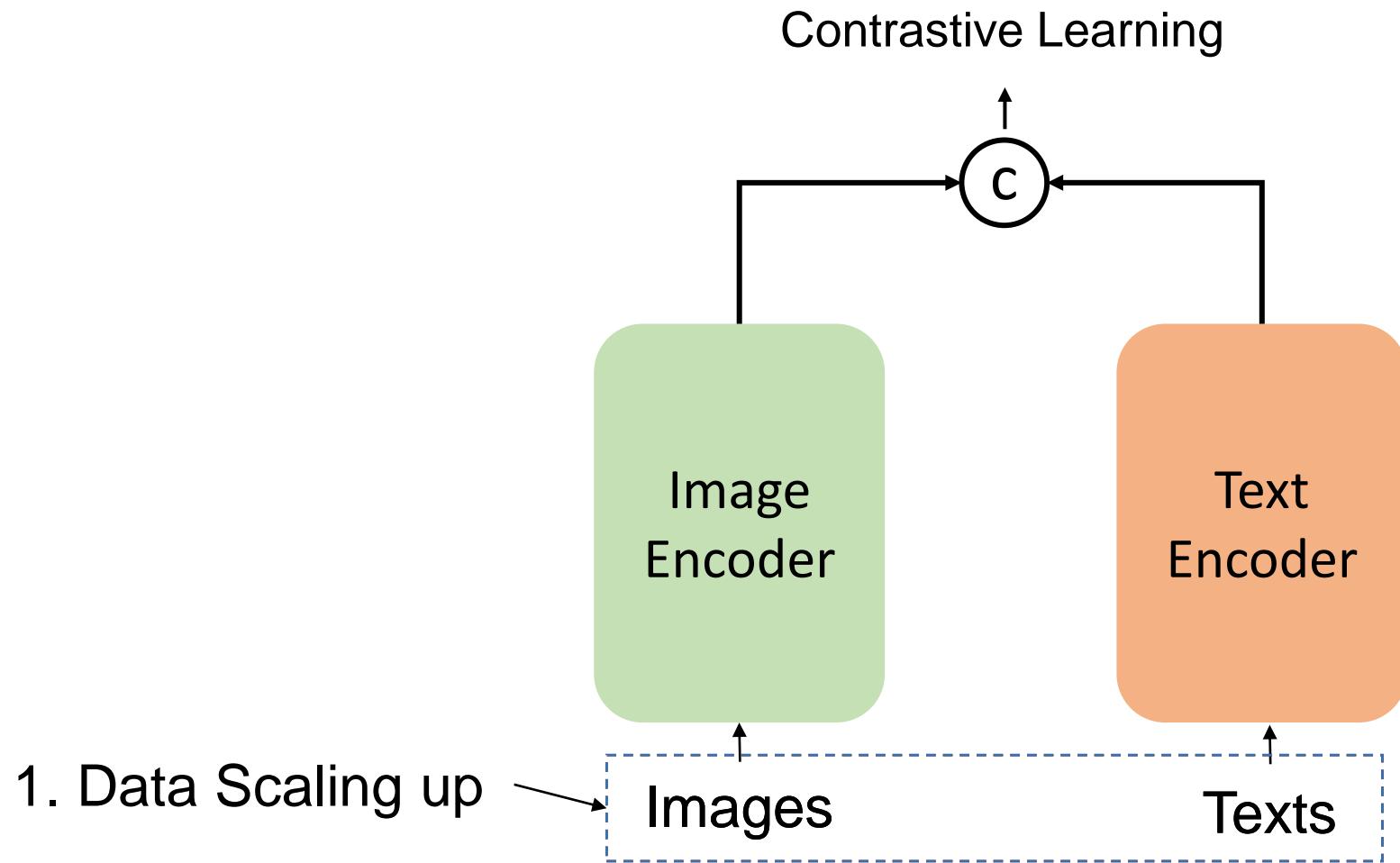
# The Lesson from CLIP

- Image recognition can be formulated as an image-text matching problem instead of image-label mapping problem
- Image recognition does not require human-annotated image-label data but huge amount of (noisy) image-text pairs
- Contrastive learning is a good learning objective for multi-modal learning strategy compared with generative learning
- Two-tower model without fusion is sufficient to learn good and generic visual and language representations

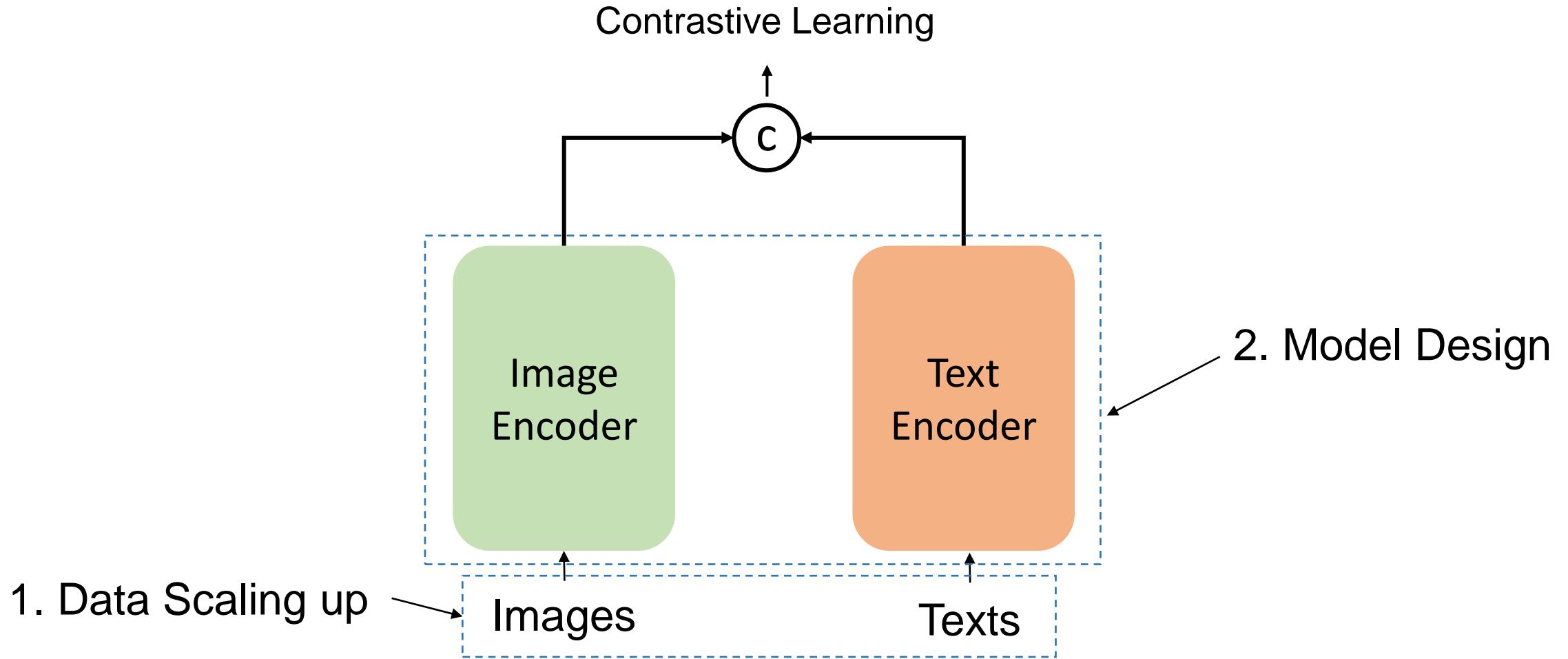
# The most recent art



# The most recent art

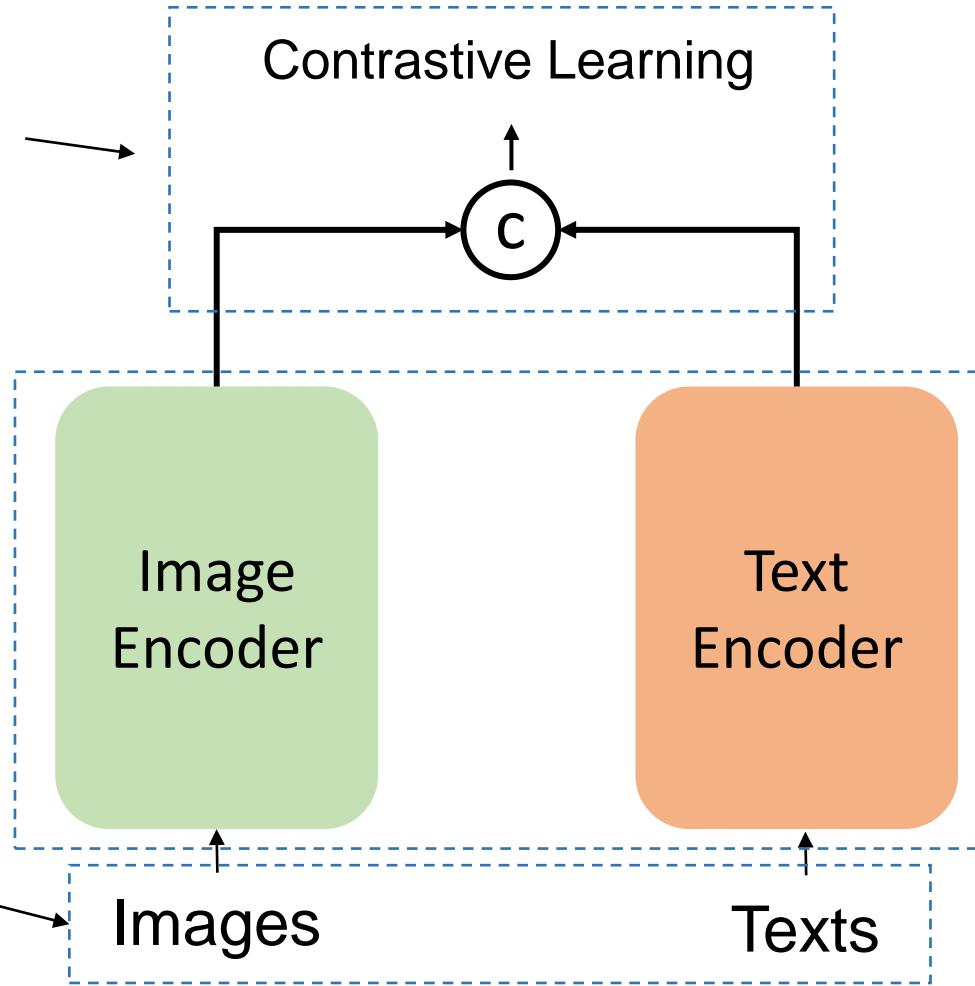


# The most recent art



# The most recent art

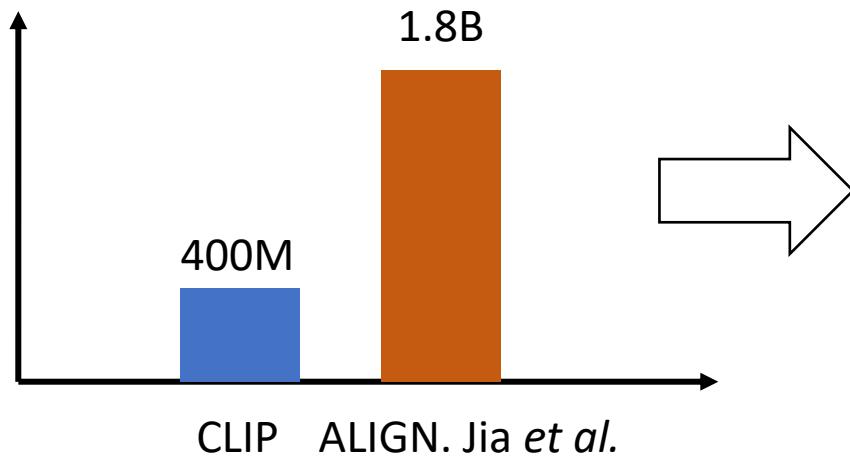
3. Objective functions



1. Data Scaling up

2. Model Design

# Data Scaling-Up



Larger scale but more noisy data



"motorcycle front wheel"



"thumbnail for version as of 21  
57 29 june 2010"



"file frankfurt airport  
skyline 2017 05 jpg"



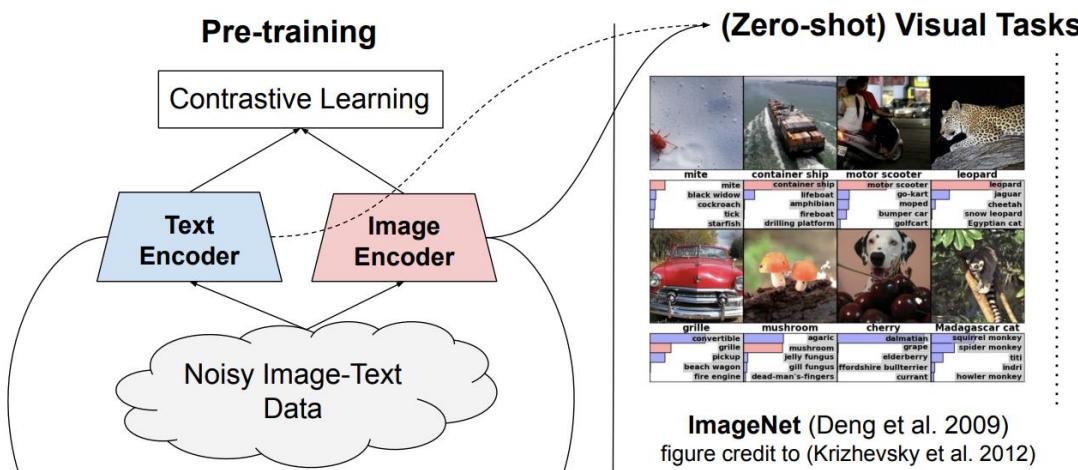
"file london barge race 2 jpg"



"moustache seamless  
wallpaper design"



"st oswalds way and shops"



Model	ImageNet	ImageNet-R	ImageNet-A	ImageNet-V2
CLIP	76.2	88.9	<b>77.2</b>	<b>70.1</b>
<b>ALIGN</b>	<b>76.4</b>	<b>92.2</b>	75.8	<b>70.1</b>

Zero-shot image classification on ImageNet

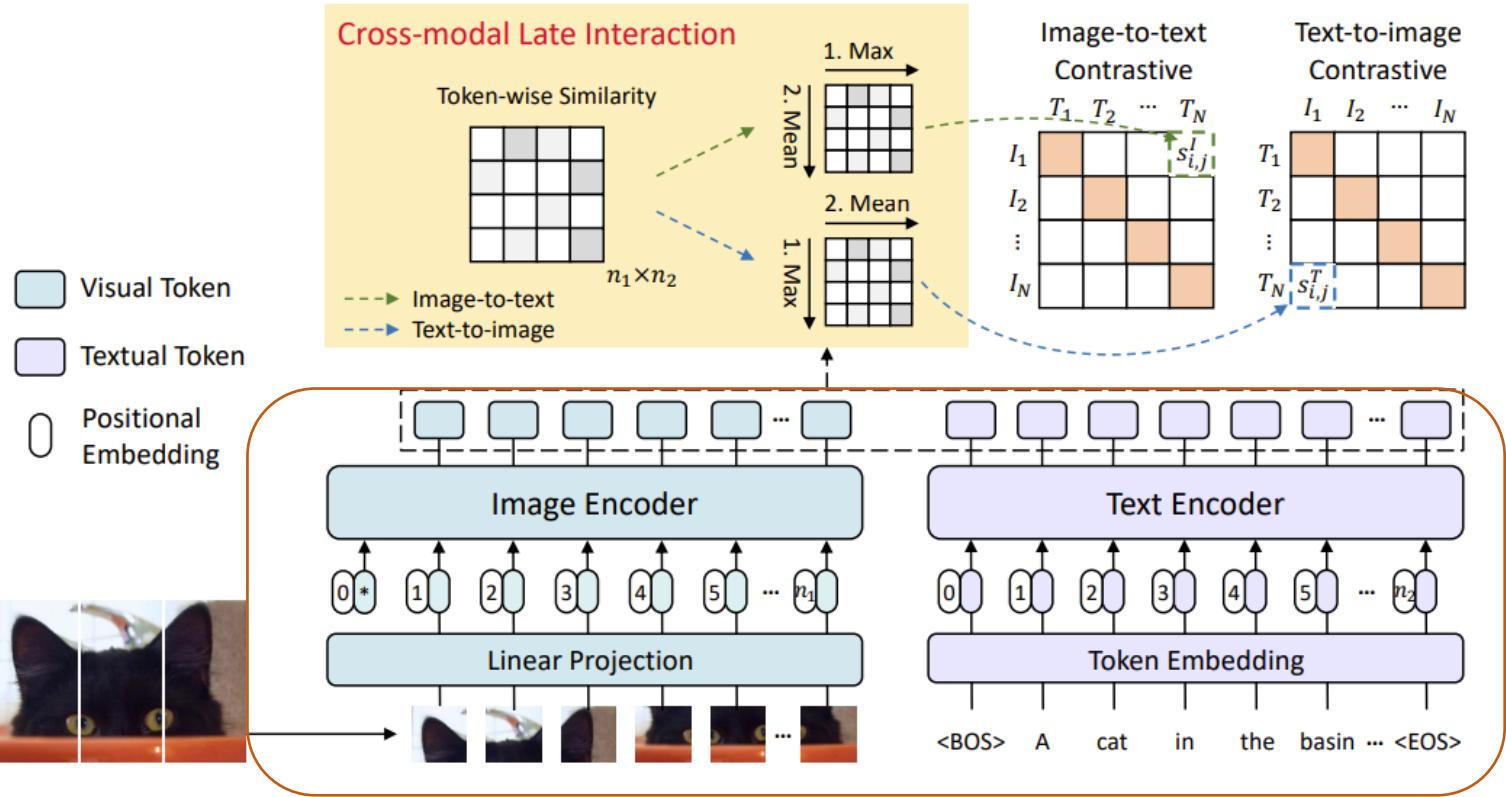
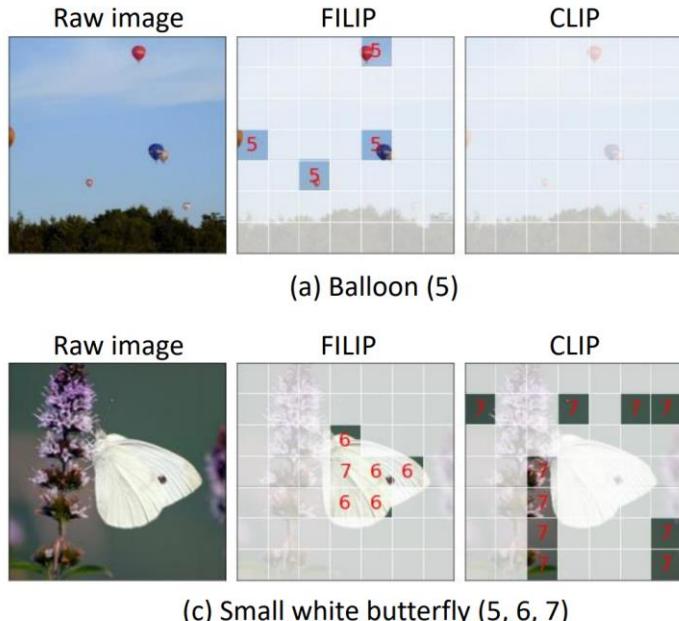
Model (backbone)	Acc@1 w/ frozen features	Acc@1	Acc@5
WSL (ResNeXt-101 32x48d)	83.6	85.4	97.6
CLIP (ViT-L/14)	85.4	-	-
BiT (ResNet152 x 4)	-	87.54	98.46
NoisyStudent (EfficientNet-L2)	-	88.4	98.7
ViT (ViT-H/14)	-	88.55	-
Meta-Pseudo-Labels (EfficientNet-L2)	-	<b>90.2</b>	<b>98.8</b>
<b>ALIGN</b> (EfficientNet-L2)	<b>85.5</b>	88.64	98.67

Image classification finetuning

# Model Design

From coarse-grain to fine-grained

**Intuition:** computes the loss by first compute the token-wise similarity, and then aggregate the matrix by maxing.



**FILIP.** Yao *et al.*

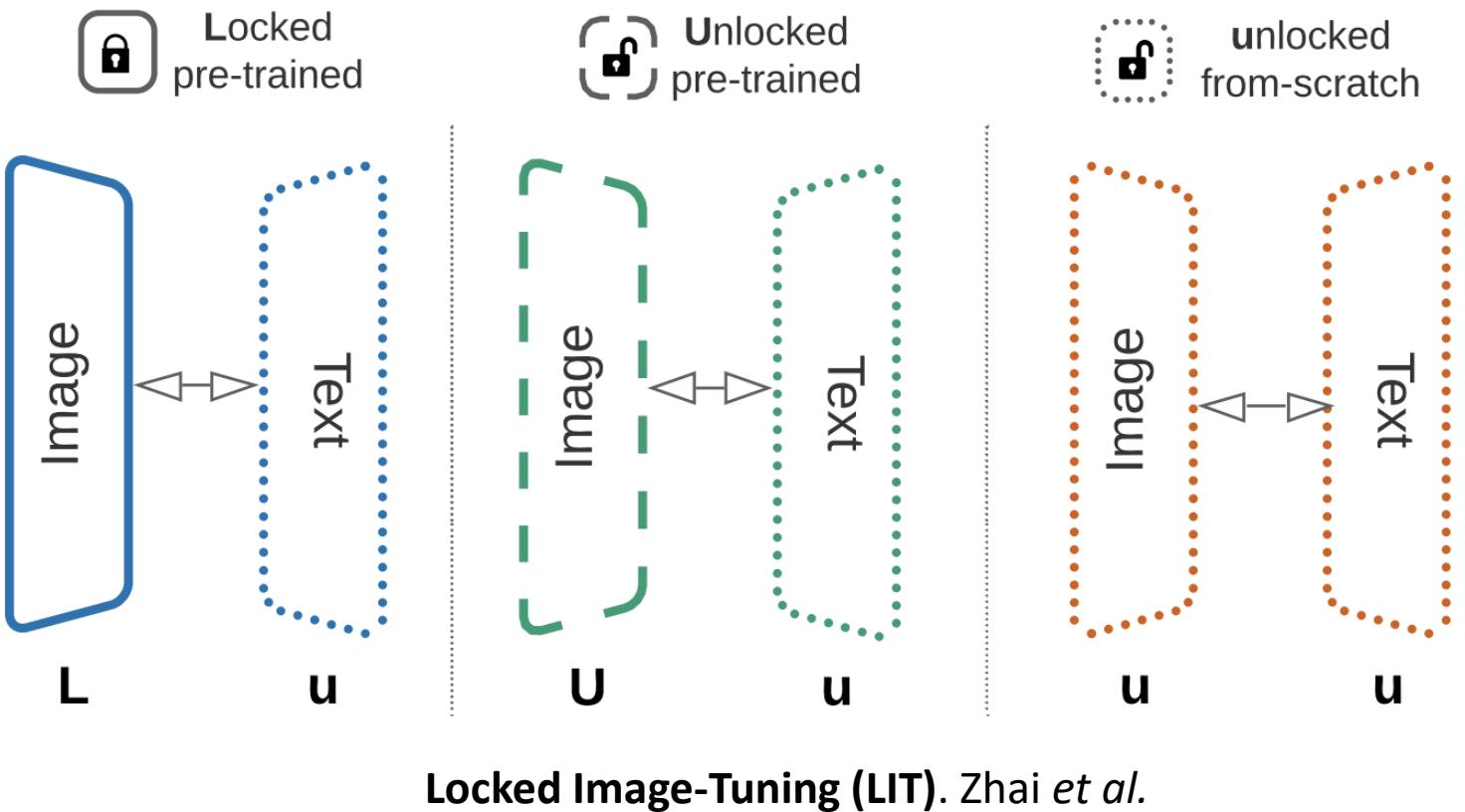
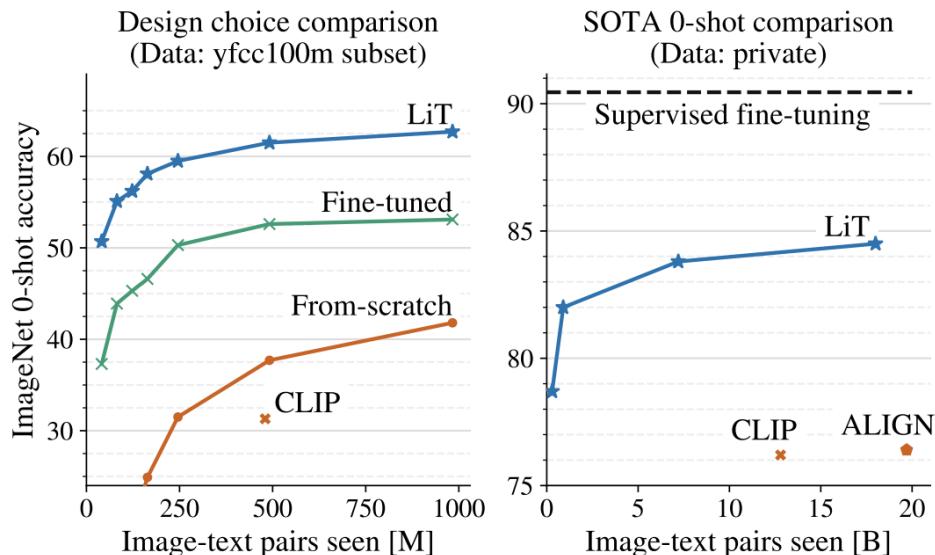
	CIFAR10	CIFAR100	Caltech101	StanfordCars	Flowers102	Food101	SUN397	DTD	Aircrafts	OxfordPets	EuroSAT	ImageNet	Average
CLIP-ViT-B/32	91.3	65.1	87.9	59.4	66.7	84.4	63.2	44.5	21.2	87.0	49.4	63.2	65.3
FILIP <sub>base</sub> -ViT-B/32	86.9	65.5	91.9	55.4	85.3	82.8	69.1	49.3	57.2	88.1	49.9	68.8	<b>70.9<sup>+5.6</sup></b>
CLIP-ViT-L/14	96.2	77.9	92.6	77.3	78.7	92.9	67.7	55.3	36.1	93.5	59.9	75.3	75.3
FILIP <sub>large</sub> -ViT-L/14	95.7	75.3	93.0	70.8	90.1	92.2	73.1	60.7	60.2	92	59.2	77.1	<b>78.3<sup>+3.0</sup></b>

Zero-shot recognition on downstream datasets

# Model Design

From joint training to sequential training

**Intuition:** image and text encoder are not necessarily trained all from scratch



**Locked Image-Tuning (LiT).** Zhai et al.

Method	I	Net	I	Net	I	Net	Obj	Net	Real	V	TAB	N
Private	CLIP [45]	76.2	70.1	88.9	77.2	72.3	-	-	-	-	-	-
	ALIGN [30]	76.4	70.1	92.2	75.8	-	-	-	-	-	-	-
	<b>LiT</b>	<b>84.5</b>	<b>78.7</b>	<b>93.9</b>	<b>79.4</b>	<b>81.1</b>	88.0	72.6	-	-	-	-
Public	CLIP [45]	31.3	-	-	-	-	-	-	-	-	-	-
	OpenCLIP [28]	34.8	30.0	-	-	-	-	-	-	-	-	-
	<b>LiT</b>	<b>75.7</b>	<b>66.6</b>	<b>60.4</b>	<b>37.8</b>	<b>54.5</b>	<b>82.1</b>	<b>63.1</b>	-	-	-	-
*	ResNet50 [25]	75.8	63.8	36.1	0.5	26.5	82.5	72.6	-	-	-	-

Zero-shot recognition

Model:	Pre-training			LiT			
	Dataset	Labels?	Full IN	10-shot	0-shot	I→T	T→I
ViT-B/16							
MoCo-v3 [10]	IN	n	76.7	60.6	55.4	33.5	17.6
DINO [4]	IN	n	78.2	61.2	55.5	33.4	18.2
AugReg [54]	IN21k	y	77.4	63.9	55.9	30.3	17.2
AugReg [54]	IN	y	77.7	77.1	64.3	25.4	13.8
AugReg [54]	Places	y	-	22.5	28.5	25.1	12.9

Effect of pretraining

# Learning Objectives

Combining contrastive learning with other learning objectives

Contrastive  
vision-language  
learning

+

Self-supervised  
Learning

= ?

+

Supervised  
Learning

= ?

# Learning Objectives

Combining contrastive learning with other learning objectives

Contrastive  
vision-language  
learning

+

Self-supervised  
Learning

= ?

+

Supervised  
Learning

= ?

# Learning Objectives

**Self-supervised learning on each modality:**

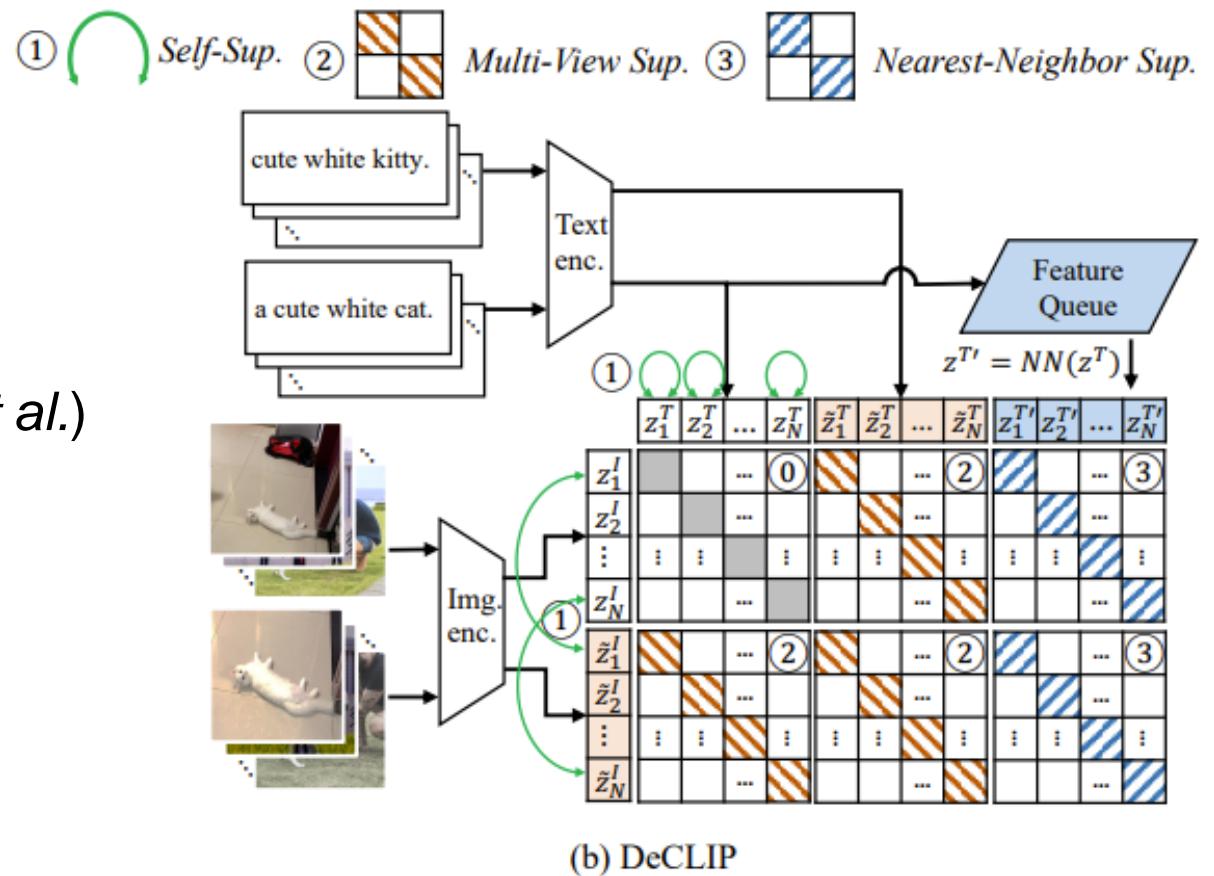
1. Image: SimSam (Chen *et al.*)
2. Text: MLM (Devlin *et al.*)

**Multi-View Supervision:**

1. Image + Augmented image (Caron *et al.*)
2. Text + Augmented text (Wei *et al.*)

**Nearest-neighbor supervision:**

1. Find the nearest text embeddings as the matched texts to an image



# Learning Objectives

**Self-supervised learning on each modality:**

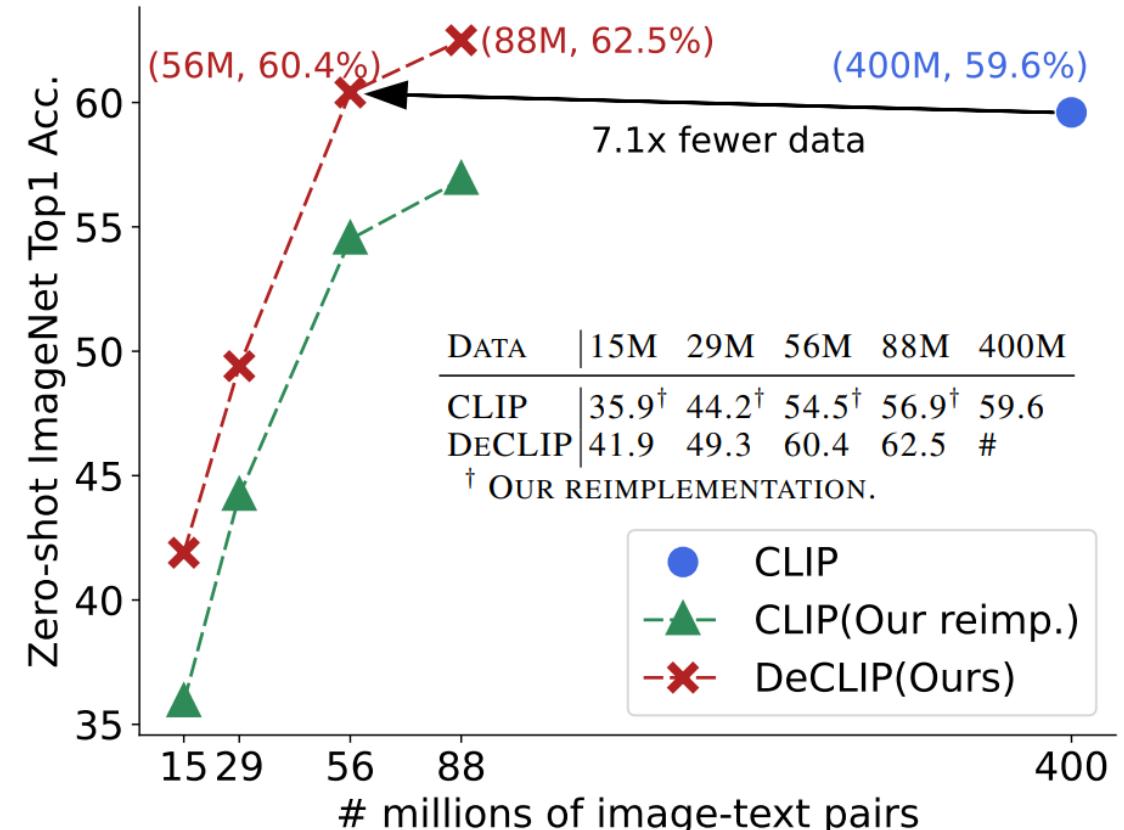
1. Image: SimSam (Chen *et al.*)
2. Text: MLM (Devlin *et al.*)

**Multi-View Supervision:**

1. Image + Augmented image (Caron *et al.*)
2. Text + Augmented text (Wei *et al.*)

**Nearest-neighbor supervision:**

1. Find the nearest text embeddings as the matched texts to an image

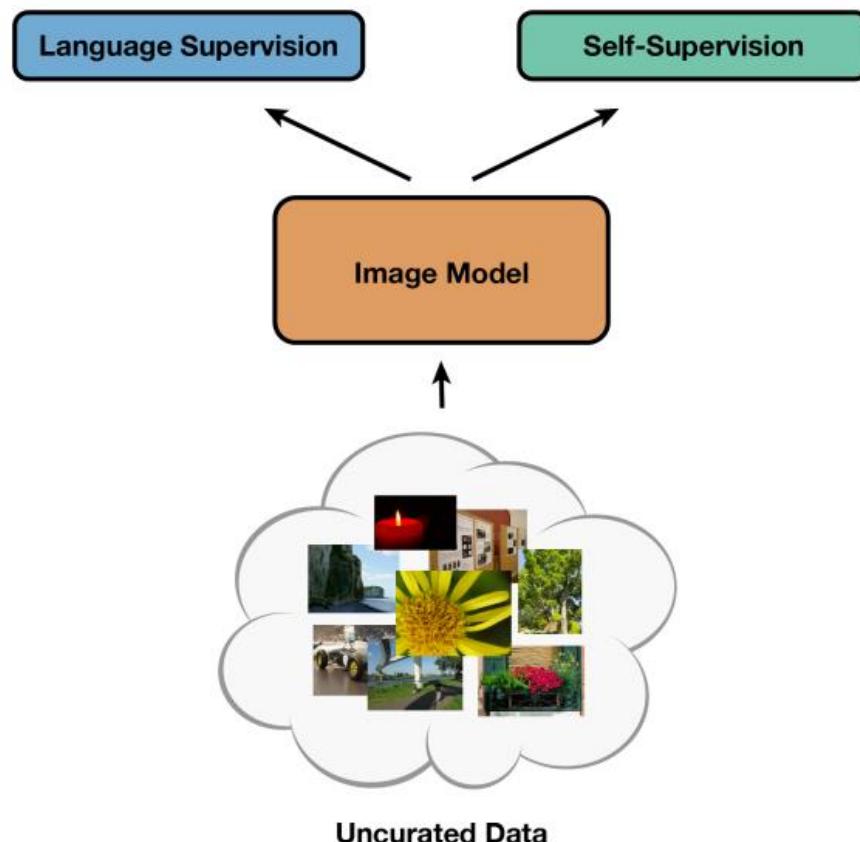


Combining vision-language and self-supervised learning improves data efficiency significantly

**Supervision Exists Everywhere.** Li *et al.*

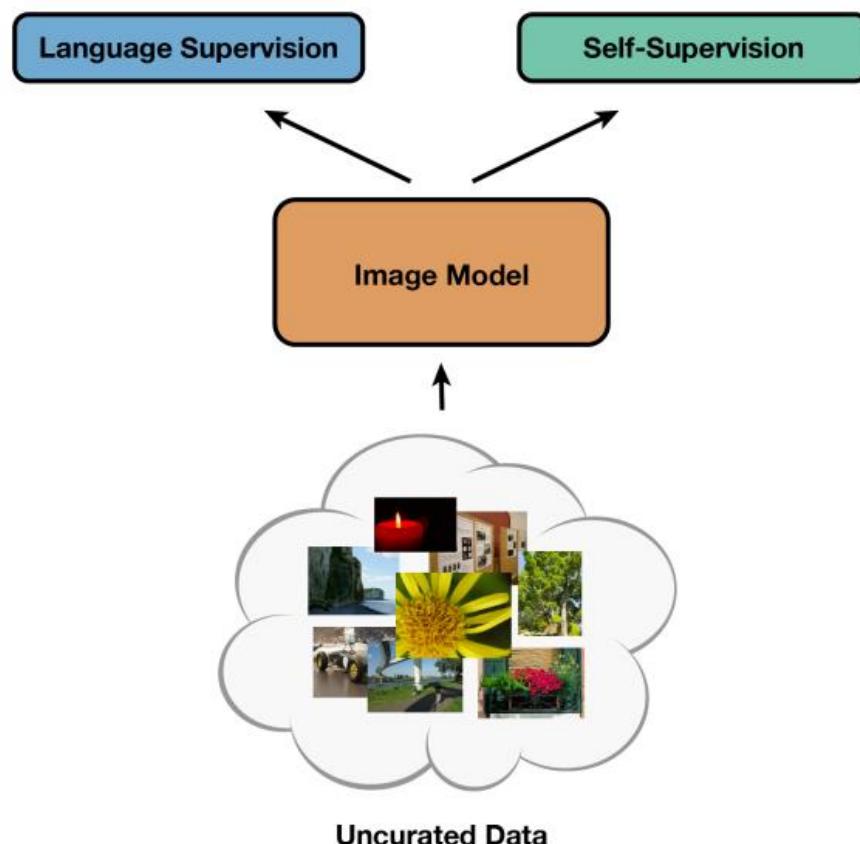
# Learning Objectives

Simply combining contrastive language-image pretraining with self-supervised learning

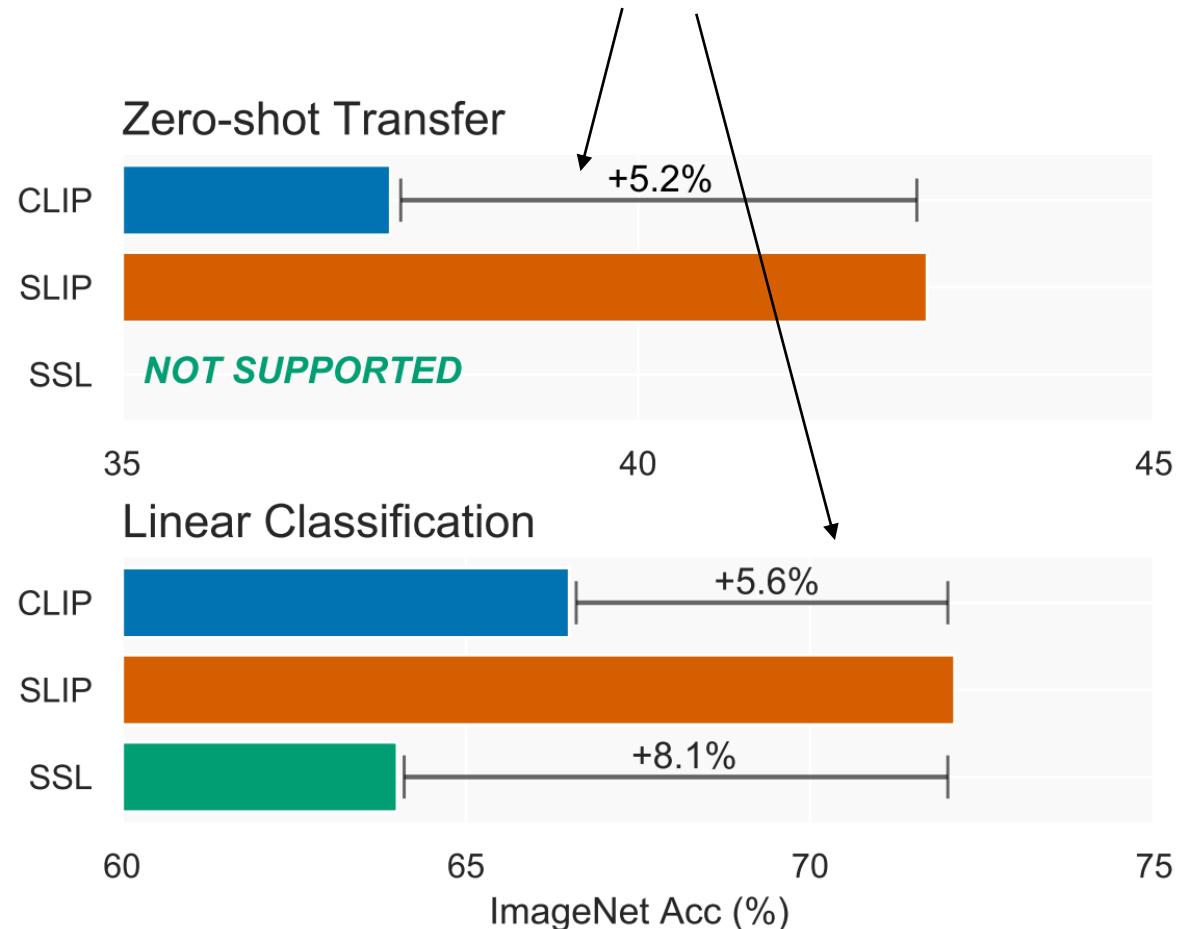


# Learning Objectives

Simply combining contrastive language-image pretraining with self-supervised learning



SLIP outperforms CLIP on both zero-shot transfer and linear classification



# Learning Objectives

Combining contrastive learning with other learning objectives

Contrastive  
vision-language  
learning

+

Self-supervised  
Learning

= ?

+

Supervised  
Learning

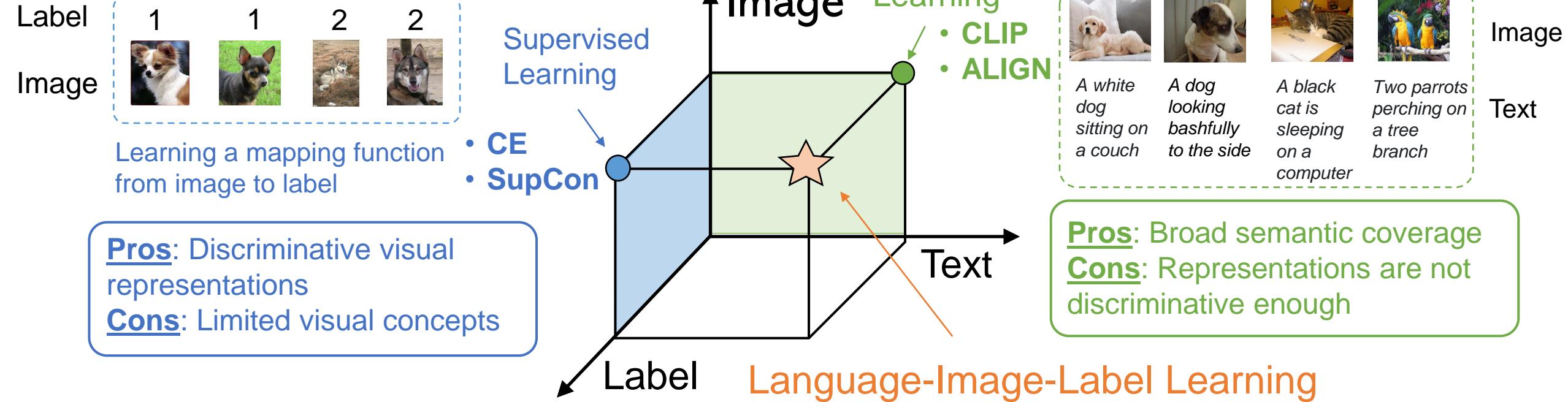
= ?

# A new perspective

Supervised Learning				Language-Image Learning				3	4	5	6
Label	1	1	2	2	3	4	5	6	Image	Text	
Image											
Chihuahua	Chihuahua	Siberian Husky	Siberian Husky	A white dog sitting on a couch	A dog looking bashfully to the side	A black cat is sleeping on a computer	Two parrots perching on a tree branch				

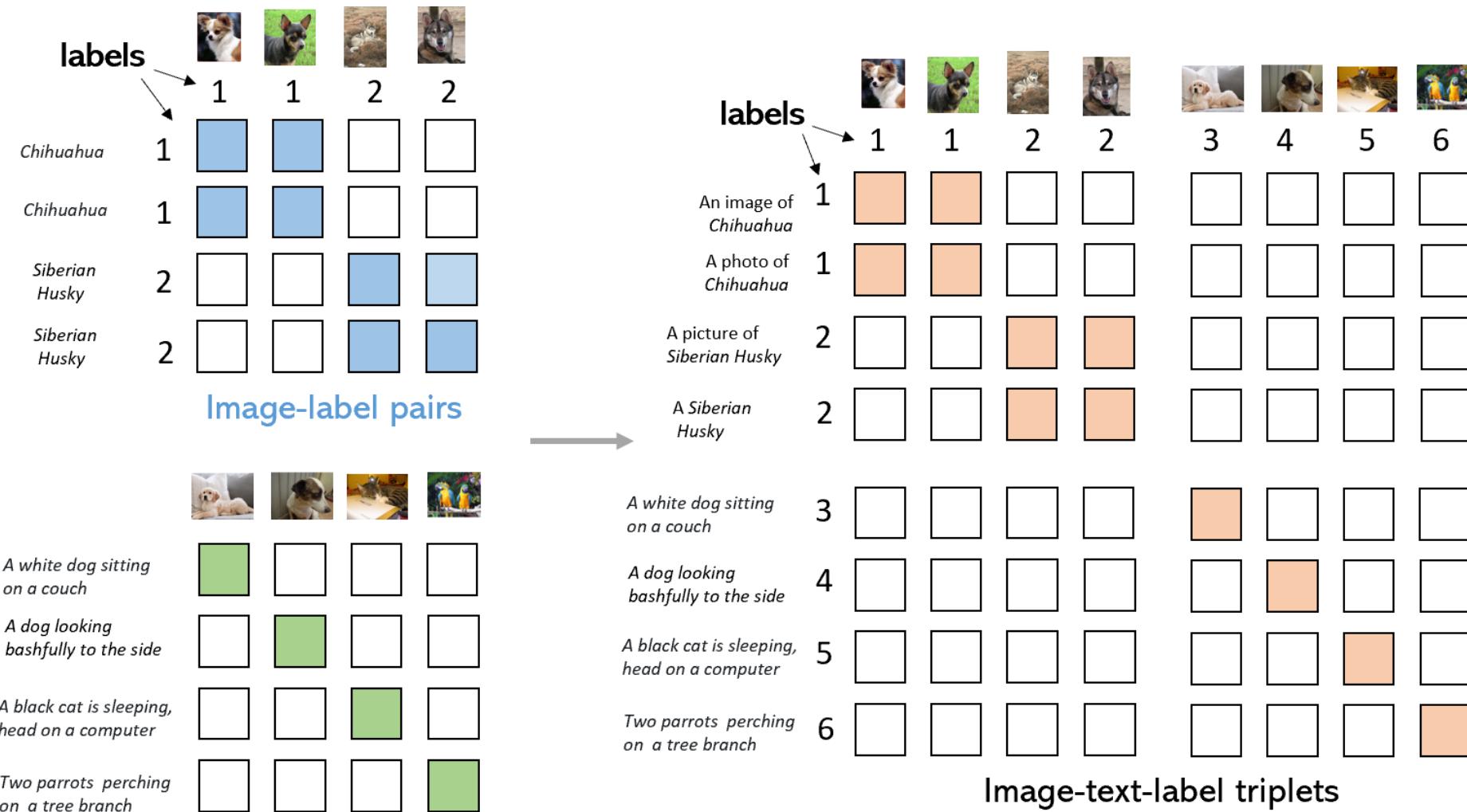
- 1 Supervised learning use the categorical labels but with concepts associated
- 2 Language-image learning assume each image-text pair belongs to a unique category

# Image-Text-Label Space



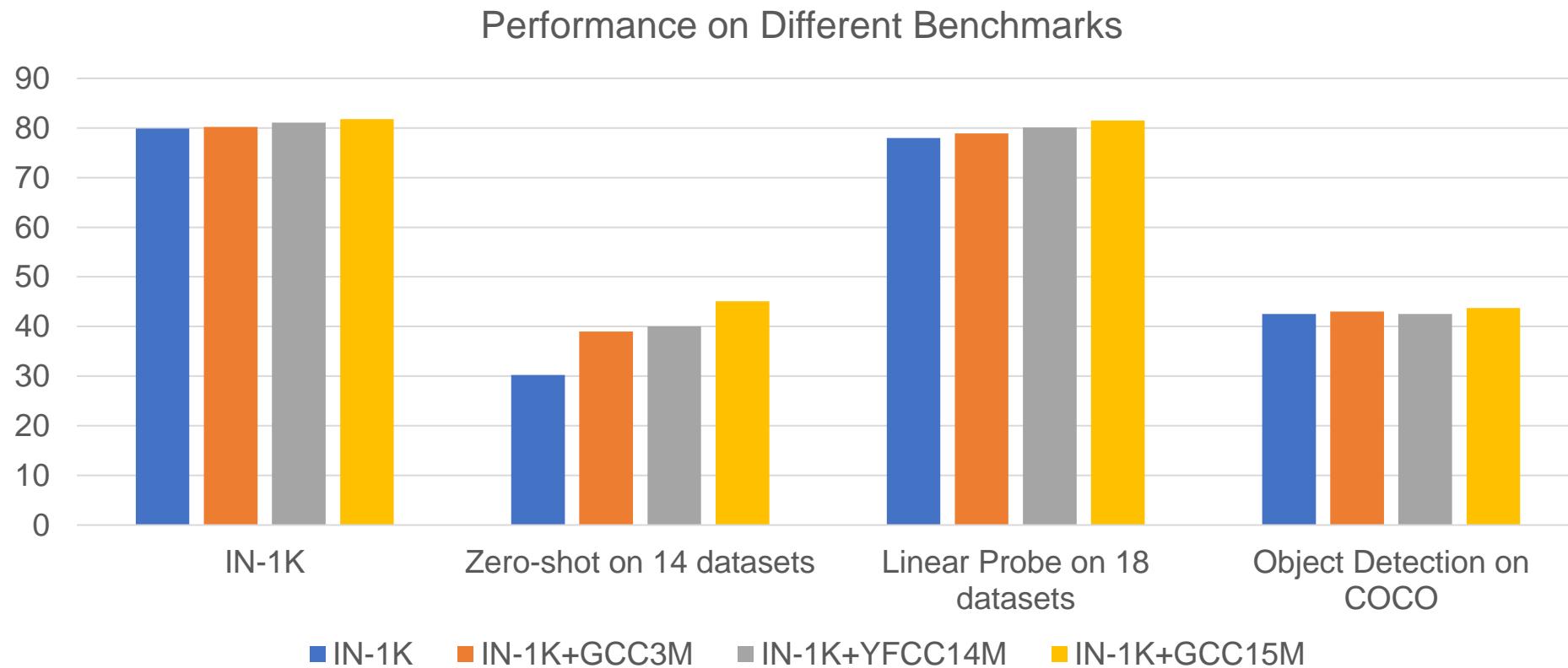
Best of both world:  
Discriminative visual representations and  
Broad semantic coverage

# Learning Objectives



# Learning Objectives

How image-text pairs benefits image classification on ImageNet-1K?

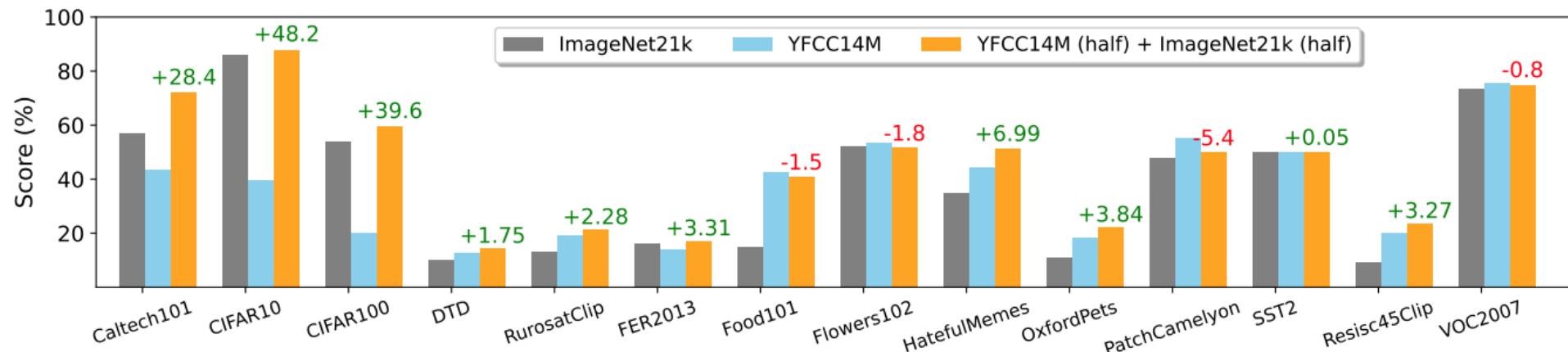


Adding image-text pairs particularly improve the zero-shot performance on various benchmarks

# Learning Objectives

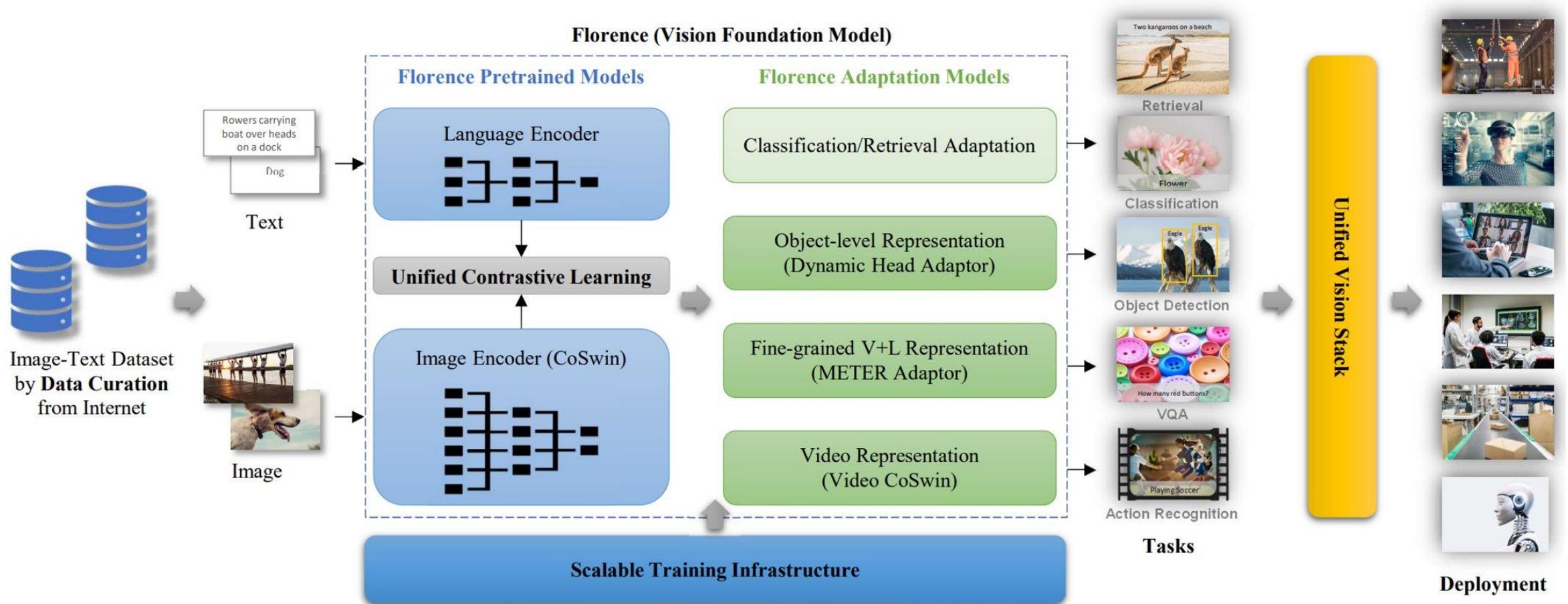
How image-label data benefits image-text pairs on low-shot recognition?

Training Data	Method	Metric			
		Zero-shot		ImageNet-1K Finetuning	Linear Probing 18 datasets
		ImageNet-1K	14 datasets		
YFCC-14M	CLIP	30.1	36.3	77.5	72.7
ImageNet-21K	UniCL	28.5	37.8	78.8	80.5
YFCC-14M(half) + ImageNet-21K(half)	UniCL	36.4	45.5	79.0	80.0
YFCC-14M + ImageNet-21K	UniCL	<b>40.5</b>	<b>49.1</b>	<b>80.2</b>	<b>81.6</b>



Adding image-label data to image-text pairs can significantly improve the zero-shot, few-shot recognition

# Scaling Up



# Scaling Up



Image  
by D:  
frc

	Florence (Vision Foundation Model)											
	Florence Pretrained Models				Florence Adaptation Models							
	Food101	CIFAR10	CIFAR100	SUN397	Stanford Cars	FGVC Aircraft	VOC2007	DTD	Oxford Pets	Caltech101	Flowers102	ImageNet
CLIP-ResNet-50x64	91.8	86.8	61.3	48.9	76.0	35.6	83.8	53.4	93.4	90.6	77.3	73.6
CLIP-ViT-L/14 (@336pix)	93.8	<b>95.7</b>	77.5	68.4	78.8	37.2	84.3	55.7	93.5	92.8	78.3	76.2
FLIP-ViT-L/14	92.2	<b>95.7</b>	75.3	73.1	70.8	<b>60.2</b>	-	60.7	92.0	93.0	<b>90.1</b>	78.3
Florence-CoSwin-H (@384pix)	<b>95.1</b>	94.6	<b>77.6</b>	<b>77.0</b>	<b>93.2</b>	55.5	<b>85.5</b>	<b>66.4</b>	<b>95.9</b>	<b>94.7</b>	86.2	<b>83.7</b>



Action Recognition

Tasks

Deployment

Scalable Training Infrastructure

More at CVPR Keynote!

# Learning Objectives

Combining contrastive learning with other learning objectives

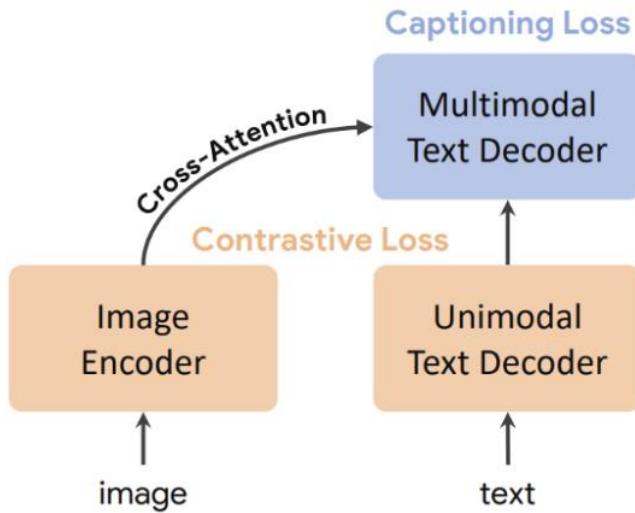
Contrastive  
vision-language  
learning

+

Generative  
Learning

= ?

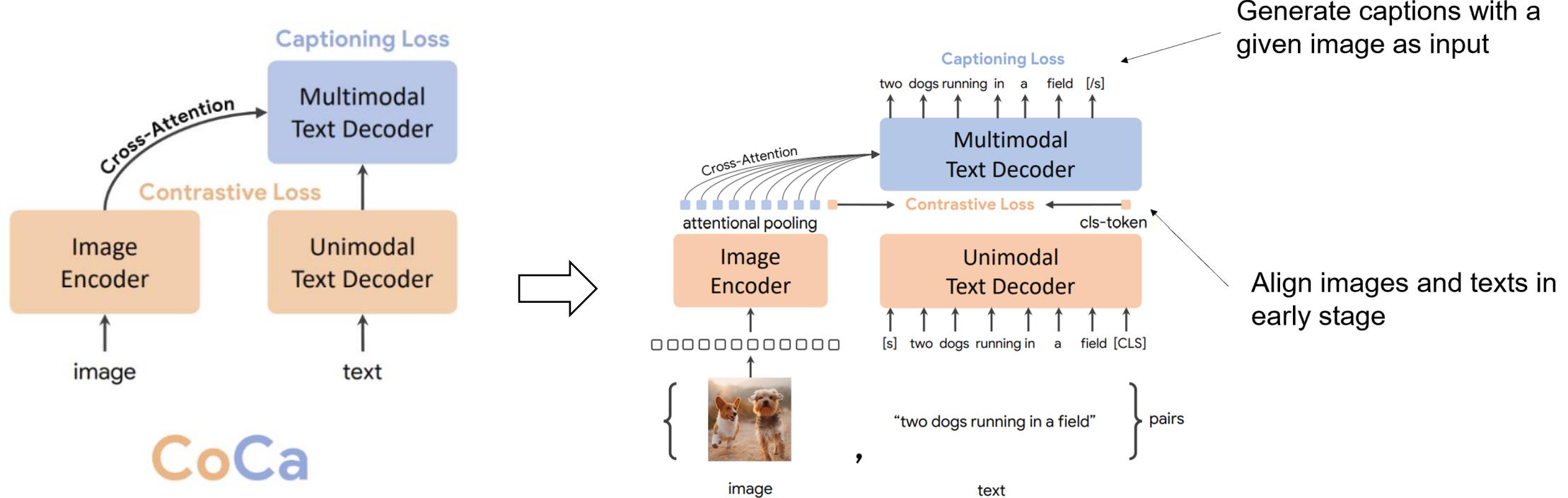
# Learning Objectives



CoCa

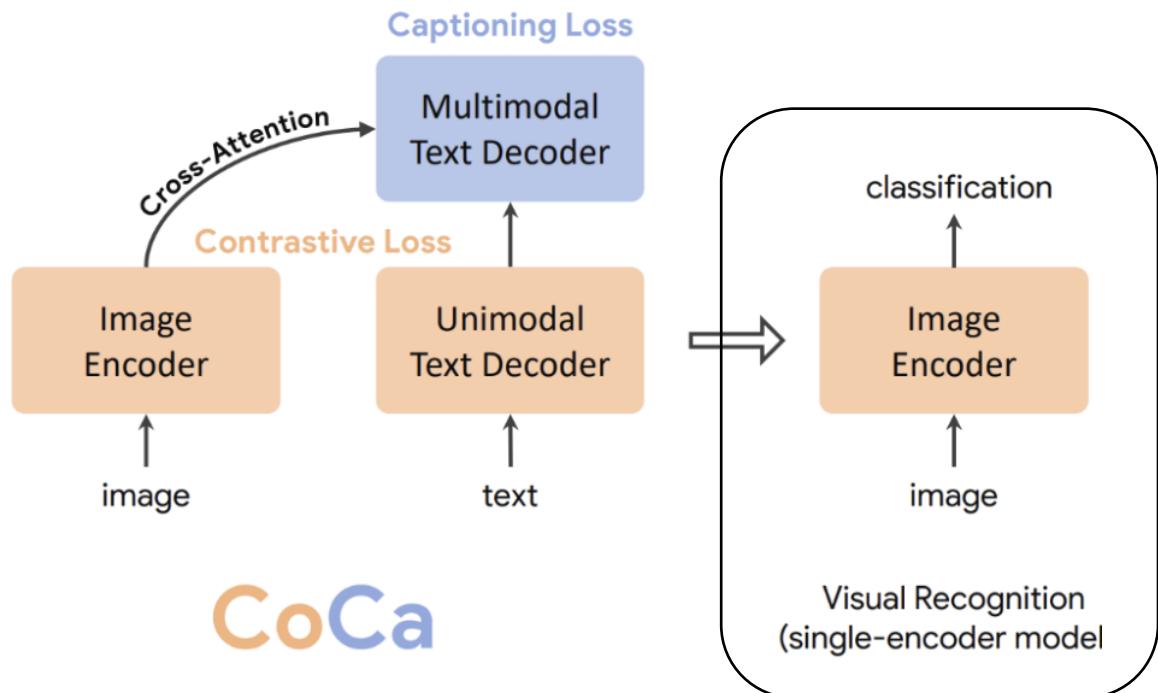
1. Training image encoder, unimodal text decoder, multimodal text decoder simultaneously
2. Training on both image-text pairs from ALIGN and JFT-3B form BigTransfer

# Learning Objectives



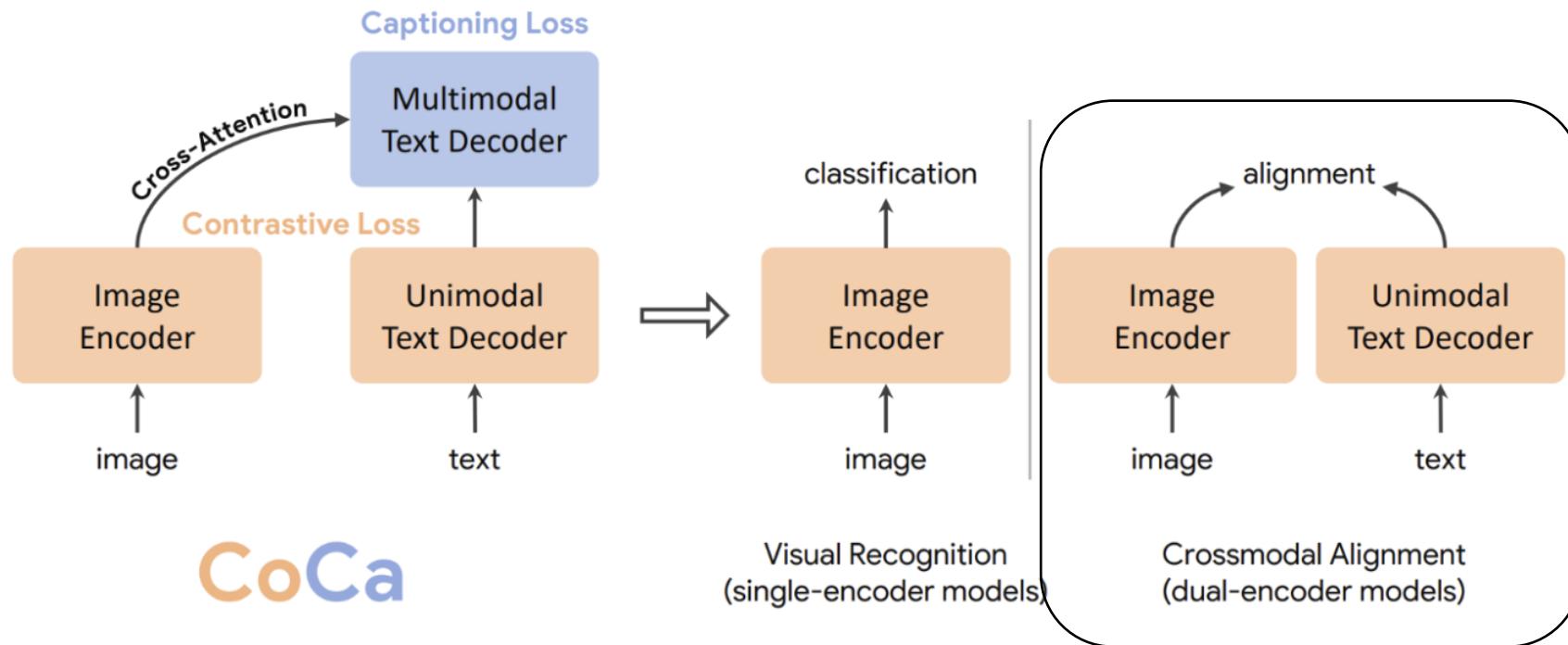
1. Training image encoder, unimodal text decoder, multimodal text decoder simultaneously
2. Training on both image-text pairs from ALIGN and JFT-3B form BigTransfer

# Learning Objectives



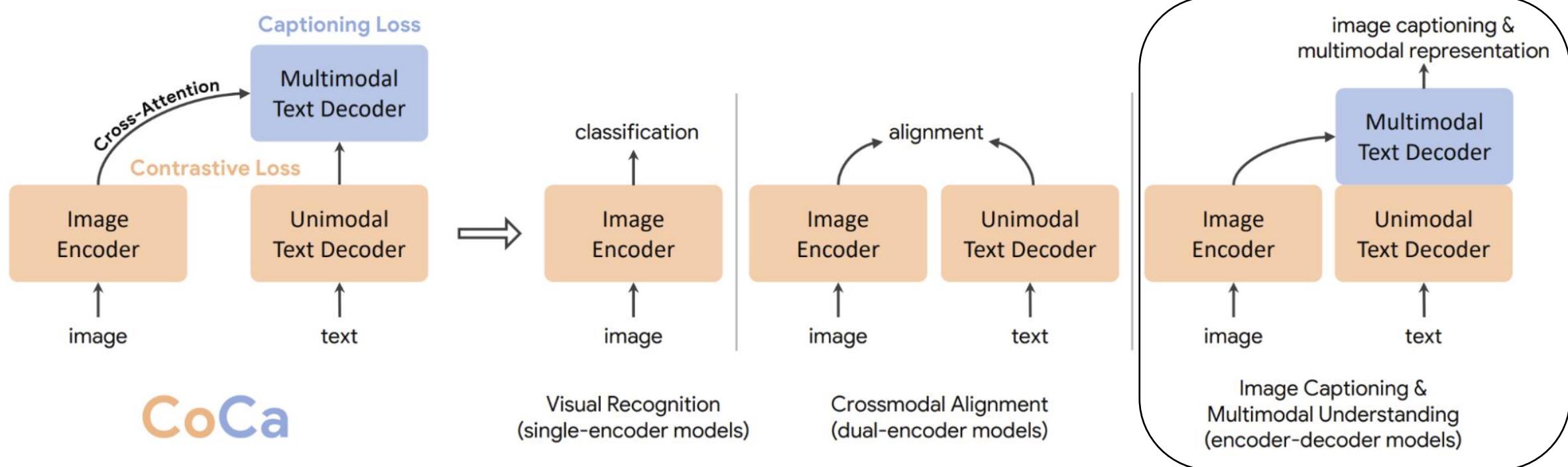
1. Training image encoder, unimodal text decoder, multimodal text decoder simultaneously
2. Training on both image-text pairs from ALIGN and JFT-3B form BigTransfer

# Learning Objectives



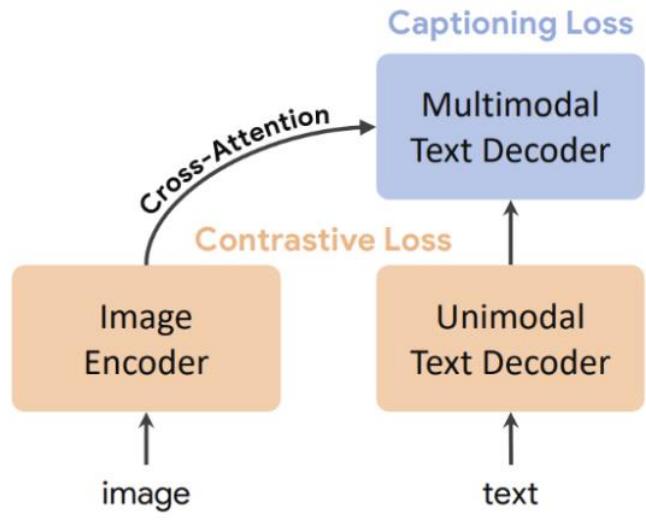
1. Training image encoder, unimodal text decoder, multimodal text decoder simultaneously
2. Training on both image-text pairs from ALIGN and JFT-3B form BigTransfer

# Learning Objectives



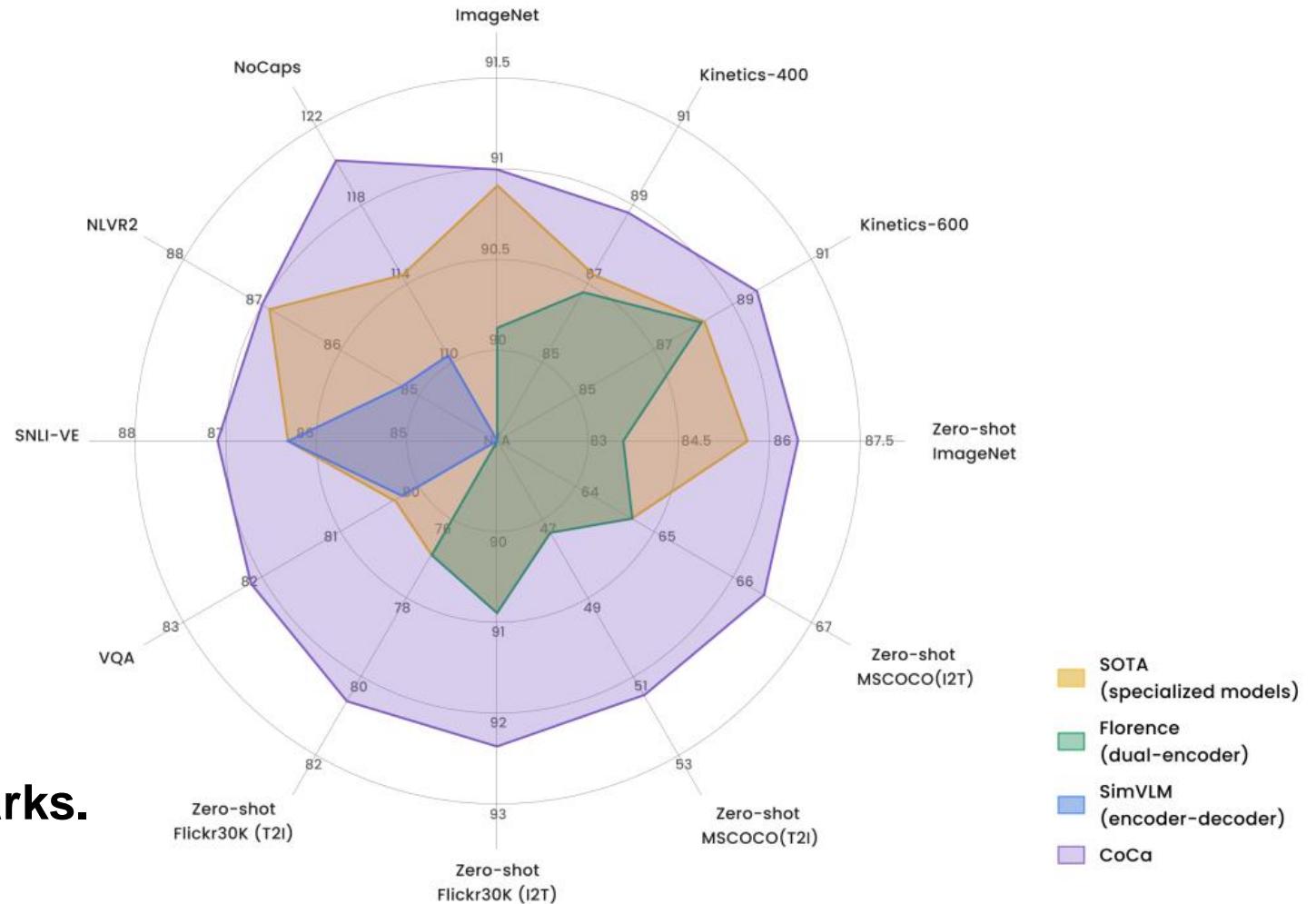
1. Training image encoder, unimodal text decoder, multimodal text decoder simultaneously
2. Training on both image-text pairs from ALIGN and JFT-3B form BigTransfer

# Learning Objectives



**CoCa**

**CoCa achieved new SoTA on various image and video recognition benchmarks.**



# Why contrastive vision-language learning works well?

What kind of data the model is learning from?

Dataset	#Images	#Concepts	Vocab. Size	#Img/C.
CIFAR-10 [34]	50k	10	10	5000
CIFAR-100 [34]	50k	100	105	500
ImageNet-1K [10]	1.3M	1,000	1,233	1300
ImageNet-22k [10]	14.2M	21,841	14,733	650
GCC-3M [51]	3.3M	17,135	7,953	193
GCC-12M [5]	12M	584,261	98,347	21
YFCC-14M [56]	14M	650,236	214,380	22

1. Image classification data has limited concepts but dense annotations
2. Image-text pairs cover a huge amount of concepts, but relative sparse

# Why contrastive vision-language learning works so well?

**Concept coverage** on downstream classification data for each pretraining data

**Observations:** pretraining corpus covers many concepts in the downstream datasets

Dataset	Name	#Concepts	Vocab. Size	ImageNet-1K		ImageNet-21K		GCC-3M		GCC-12M		YFCC-14M	
				Cover.	#Img/C.	Cover.	#Img/C.	Cover.	#Img/C.	Cover.	#Img/C.	Cover.	#Img/C.
	ImageNet-1K	1,000	1,233	<b>100%</b>	1300	0%	0	45.3%	247.0	<b>78.5%</b>	851.1	<b>69.3%</b>	1930.8
	Food-101	102	139	4.0%	1300.0	20.8%	650.0	21.8%	39.8	<b>58.4%</b>	250.8	<b>67.3%</b>	408.8
	CIFAR-10	10	10	0.0%	0.0	<b>90.0%</b>	650.0	<b>100.0%</b>	6175.4	<b>100.0%</b>	19969.8	<b>100.0%</b>	32998.9
	CIFAR-100	100	100	24.0%	1300.0	<b>65.0%</b>	650.0	<b>95.0%</b>	3928.4	<b>99.0%</b>	15628.5	<b>99.0%</b>	18303.2
	SUN397	397	432	5.0%	1300.0	28.5%	650.0	48.1%	818.9	<b>65.5%</b>	2355.4	<b>66.5%</b>	7043.2
	Stanford Cars	196	291	0.0%	0.0	0.0%	0.0	0.0%	0.0	0.0%	0.0	0.0%	0.0
	FGVC Aircraft (variants)	100	115	0.0%	0.0	0.0%	0.0	0.0%	0.0	22.0%	4.1	0.0%	0.0
	VOC2007 classification	20	20	0.0%	0.0	75.0%	650.0	<b>85.0%</b>	14721.6	<b>85.0%</b>	19934.8	<b>85.0%</b>	31448.8
	Describable Textures	47	47	0.0%	0.0	4.3%	650.0	14.9%	8.9	27.7%	53.2	36.2%	181.7
	Oxford-IIIT Pets	37	53	5.4%	1300.0	13.5%	650.0	10.8%	80.9	<b>64.9%</b>	134.0	37.8%	169.0
	Caltech-101	102	122	24.5%	1300.0	43.1%	650.0	<b>66.6%</b>	1633.8	<b>83.3%</b>	5249.7	<b>87.3%</b>	5017.7
	Oxford Flowers 102	102	147	10.0%	1300.0	40.2%	650.0	17.6%	53.2	<b>50.0%</b>	194.3	<b>65.7%</b>	422.7
	MNIST	10	10	0.0%	0.0	0.0%	0.0	40.0%	0.8	<b>100.0%</b>	46.0	<b>90.0%</b>	68.8
	FER 2013 *	8	12	0.0%	0.0	8.3%	650.0	25.0%	5.9	41.7%	29.2	41.7%	11.5
	STL10	10	10	0.0%	0.0	100%	650.0	<b>100.0%</b>	8778.6	<b>100.0%</b>	28547.6	<b>100.0%</b>	45587.5
	GTSRB *	43	85	0.0%	0.0	0.0%	0.0	2.3%	12.7	2.3%	52.9	2.3%	551.3
	PatchCamelyon	2	6	0.0%	0.0	0.0%	0.0	0.0%	0.0	<b>50.0%</b>	143.0	<b>50.0%</b>	15.0
	UCF101 *	101	153	0.0%	0.0	0.0%	0.0	0.0%	0.0	<b>51.5%</b>	66.4	0.0%	0.0
	Hateful Memes	2	2	0.0%	0.0	0.0%	0.0	<b>50.0%</b>	79.5	<b>50.0%</b>	2742.5	<b>50.0%</b>	321.5
	EuroSAT	10	20	0.0%	0.0	0.0%	0.0	20.0%	2946.6	30.0%	5266.3	30.0%	15458.7
	Resisc45	45	59	8.9%	1300.0	26.7%	650.0	<b>71.1%</b>	3688.6	<b>75.6%</b>	7572.0	<b>80.0%</b>	26317.6
	Rendered-SST2	2	2	0.0%	0.0	<b>50.0%</b>	650.0	<b>50.0%</b>	1.0	<b>100.0%</b>	114.0	<b>100.0%</b>	1259.0

# Why contrastive vision-language learning works so well?

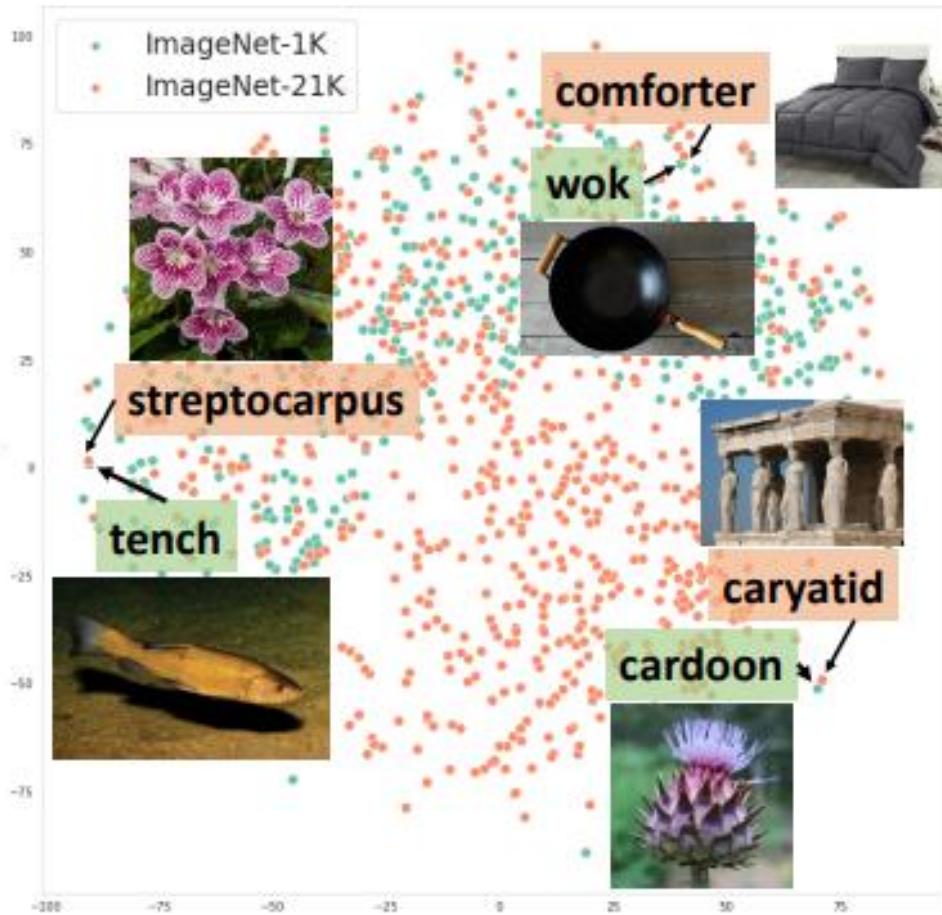
“Zero-shot” is not a real zero-shot

A realistic way to learn!

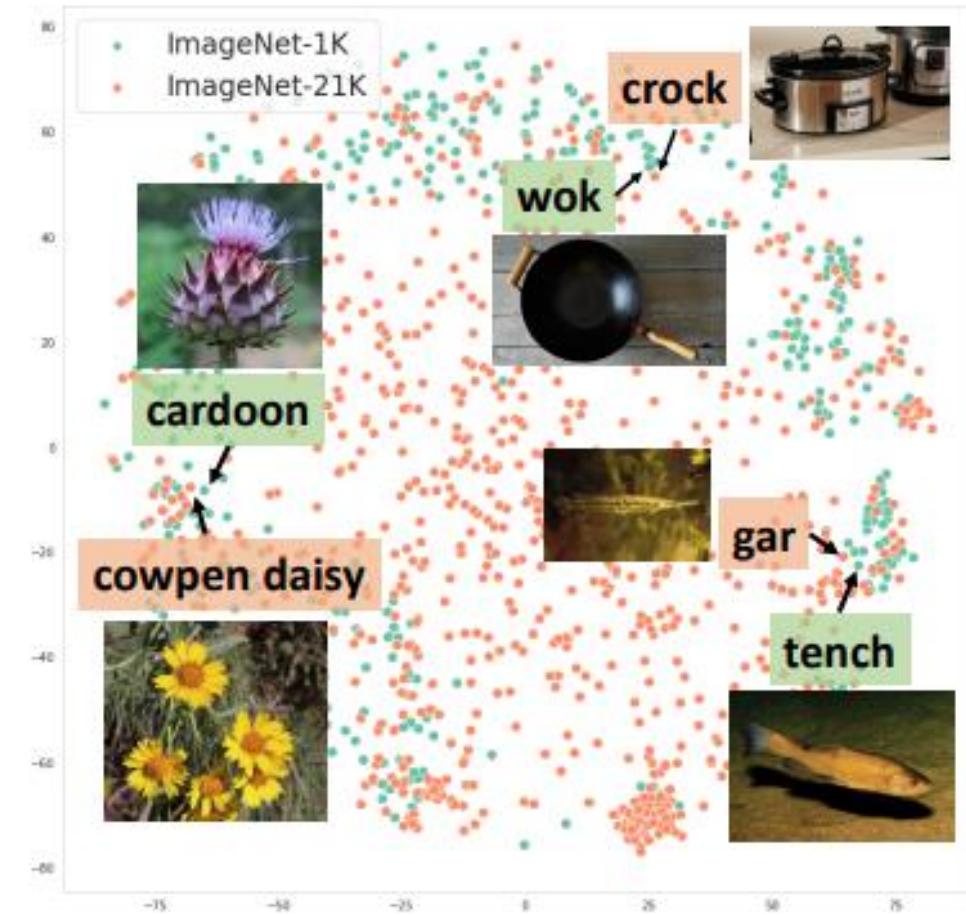
Name	Dataset			ImageNet-1K		ImageNet-21K		GCC-3M		GCC-12M		YFCC-14M	
	#Concepts	Vocab. Size	Cover.	#Img/C.	Cover.	#Img/C.	Cover.	#Img/C.	Cover.	#Img/C.	Cover.	#Img/C.	
ImageNet-1K	1,000	1,233	<b>100%</b>	1300	0%	0	45.3%	247.0	<b>78.5%</b>	851.1	<b>69.3%</b>	1930.8	
Food-101	102	139	4.0%	1300.0	20.8%	650.0	21.8%	39.8	<b>58.4%</b>	250.8	<b>67.3%</b>	408.8	
CIFAR-10	10	10	0.0%	0.0	<b>90.0%</b>	650.0	<b>100.0%</b>	6175.4	<b>100.0%</b>	19969.8	<b>100.0%</b>	32998.9	
CIFAR-100	100	100	24.0%	1300.0	<b>65.0%</b>	650.0	<b>95.0%</b>	3928.4	<b>99.0%</b>	15628.5	<b>99.0%</b>	18303.2	
SUN397	397	432	5.0%	1300.0	28.5%	650.0	48.1%	818.9	<b>65.5%</b>	2355.4	<b>66.5%</b>	7043.2	
Stanford Cars	196	291	0.0%	0.0	0.0%	0.0	0.0%	0.0	0.0%	0.0	0.0%	0.0	
FGVC Aircraft (variants)	100	115	0.0%	0.0	0.0%	0.0	0.0%	0.0	22.0%	4.1	0.0%	0.0	
VOC2007 classification	20	20	0.0%	0.0	75.0%	650.0	<b>85.0%</b>	14721.6	<b>85.0%</b>	19934.8	<b>85.0%</b>	31448.8	
Describable Textures	47	47	0.0%	0.0	4.3%	650.0	14.9%	8.9	27.7%	53.2	36.2%	181.7	
Oxford-IIIT Pets	37	53	5.4%	1300.0	13.5%	650.0	10.8%	80.9	<b>64.9%</b>	134.0	37.8%	169.0	
Caltech-101	102	122	24.5%	1300.0	43.1%	650.0	<b>66.6%</b>	1633.8	<b>83.3%</b>	5249.7	<b>87.3%</b>	5017.7	
Oxford Flowers 102	102	147	10.0%	1300.0	40.2%	650.0	17.6%	53.2	<b>50.0%</b>	194.3	<b>65.7%</b>	422.7	
MNIST	10	10	0.0%	0.0	0.0%	0.0	40.0%	0.8	<b>100.0%</b>	46.0	<b>90.0%</b>	68.8	
FER 2013 *	8	12	0.0%	0.0	8.3%	650.0	25.0%	5.9	41.7%	29.2	41.7%	11.5	
STL10	10	10	0.0%	0.0	100%	650.0	<b>100.0%</b>	8778.6	<b>100.0%</b>	28547.6	<b>100.0%</b>	45587.5	
GTSRB *	43	85	0.0%	0.0	0.0%	0.0	2.3%	12.7	2.3%	52.9	2.3%	551.3	
PatchCamelyon	2	6	0.0%	0.0	0.0%	0.0	0.0%	0.0	<b>50.0%</b>	143.0	<b>50.0%</b>	15.0	
UCF101 *	101	153	0.0%	0.0	0.0%	0.0	0.0%	0.0	<b>51.5%</b>	66.4	0.0%	0.0	
Hateful Memes	2	2	0.0%	0.0	0.0%	0.0	<b>50.0%</b>	79.5	<b>50.0%</b>	2742.5	<b>50.0%</b>	321.5	
EuroSAT	10	20	0.0%	0.0	0.0%	0.0	20.0%	2946.6	30.0%	5266.3	30.0%	15458.7	
Resisc45	45	59	8.9%	1300.0	26.7%	650.0	<b>71.1%</b>	3688.6	<b>75.6%</b>	7572.0	<b>80.0%</b>	26317.6	
Rendered-SST2	2	2	0.0%	0.0	<b>50.0%</b>	650.0	<b>50.0%</b>	1.0	<b>100.0%</b>	114.0	<b>100.0%</b>	1259.0	

# Visualizations

Visualization of learned text concept embeddings in IN-1K and IN-21K



Left: learning purely on imagenet-1k

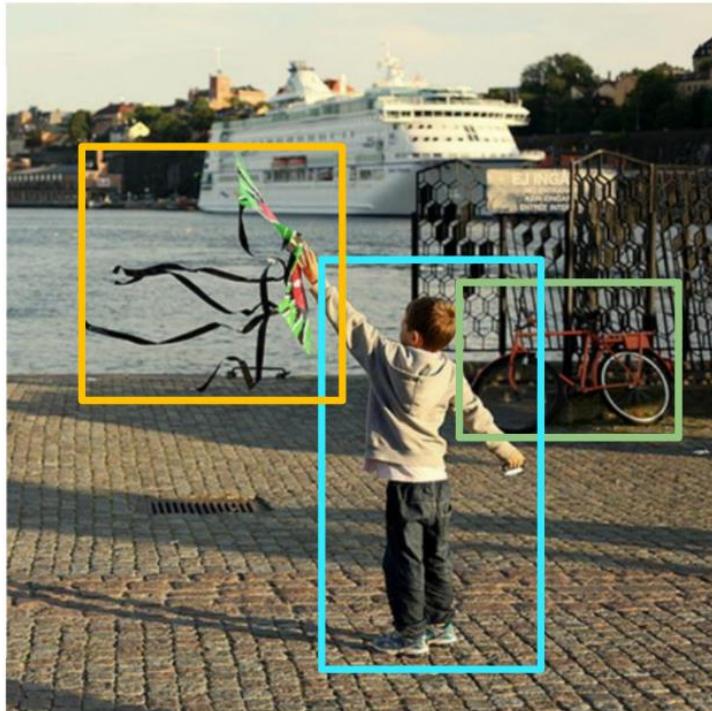


Right: learning on imagenet-1k plus image-text pairs

# Take-aways

- Visual recognition is essentially a vision-language problem, i.e., aligning images with semantic concepts/descriptions.
- Large-scale image-text pairs can help model learn the understanding of rich semantics
- Contrastive vision-language learning can be combined with other learning methods, such as self-supervised, supervised or even generative learning.

# Visual Region Recognition



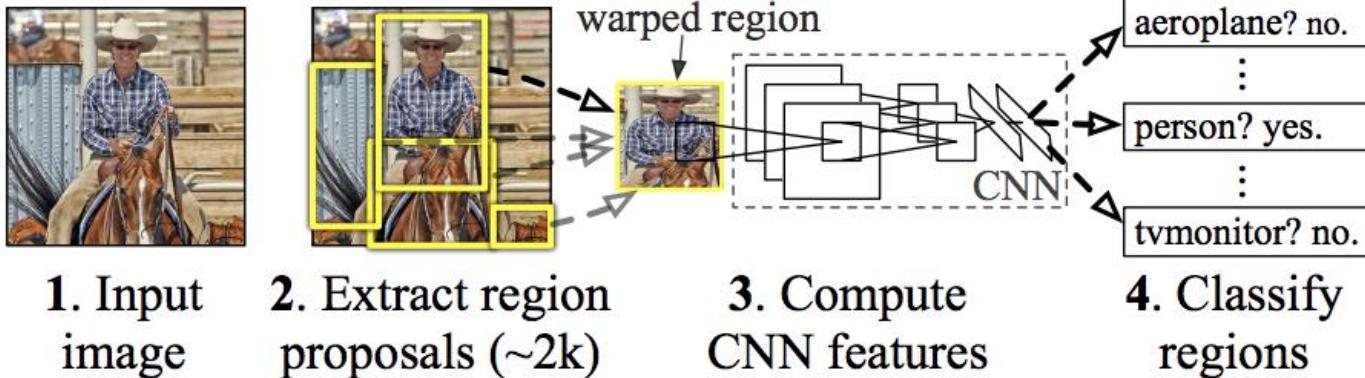
Objects: boy, kite, bicycle

- Attributes: little, green, orange
- Relationships: flying, near
- Caption: “A little boy is flying a green kite near an orange bicycle.”
- Question: What is the boy next to the bicycle doing?

Objects play a major role in visual understanding.

# Visual Region Recognition

## R-CNN: *Regions with CNN features*



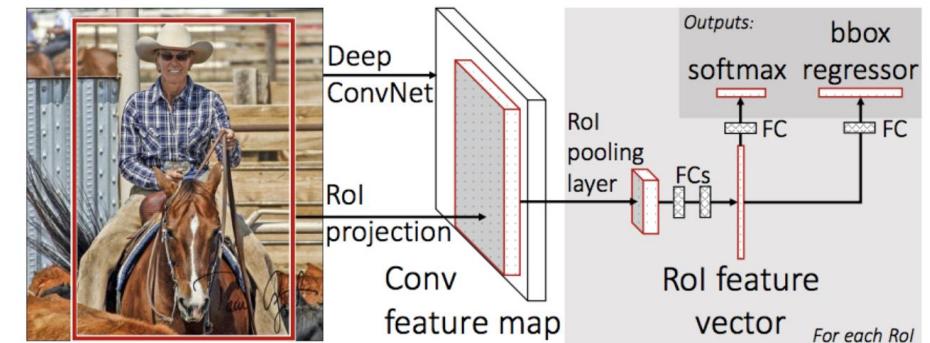
1. Input image

2. Extract region proposals (~2k)

3. Compute CNN features

4. Classify regions

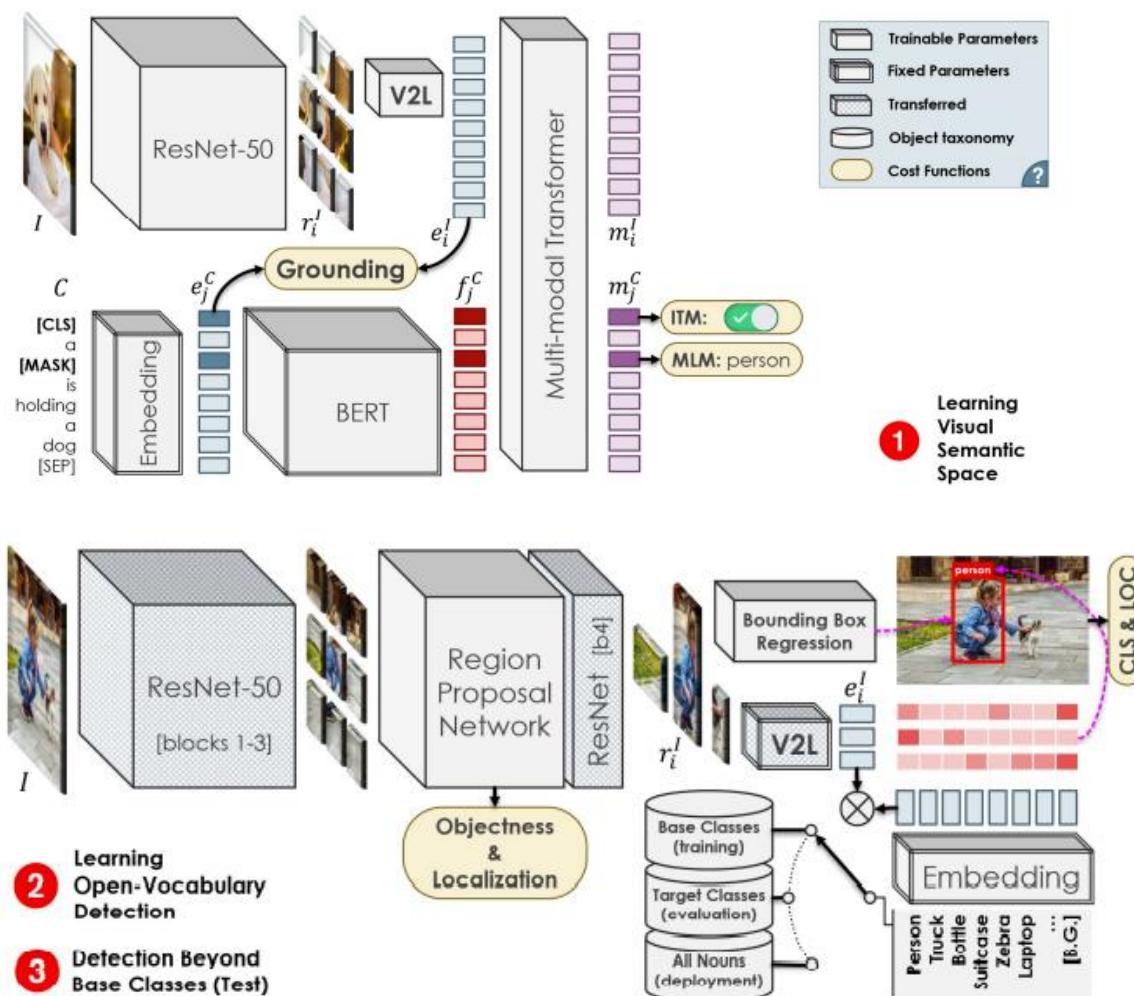
R-CNN. Girshick *et al.*



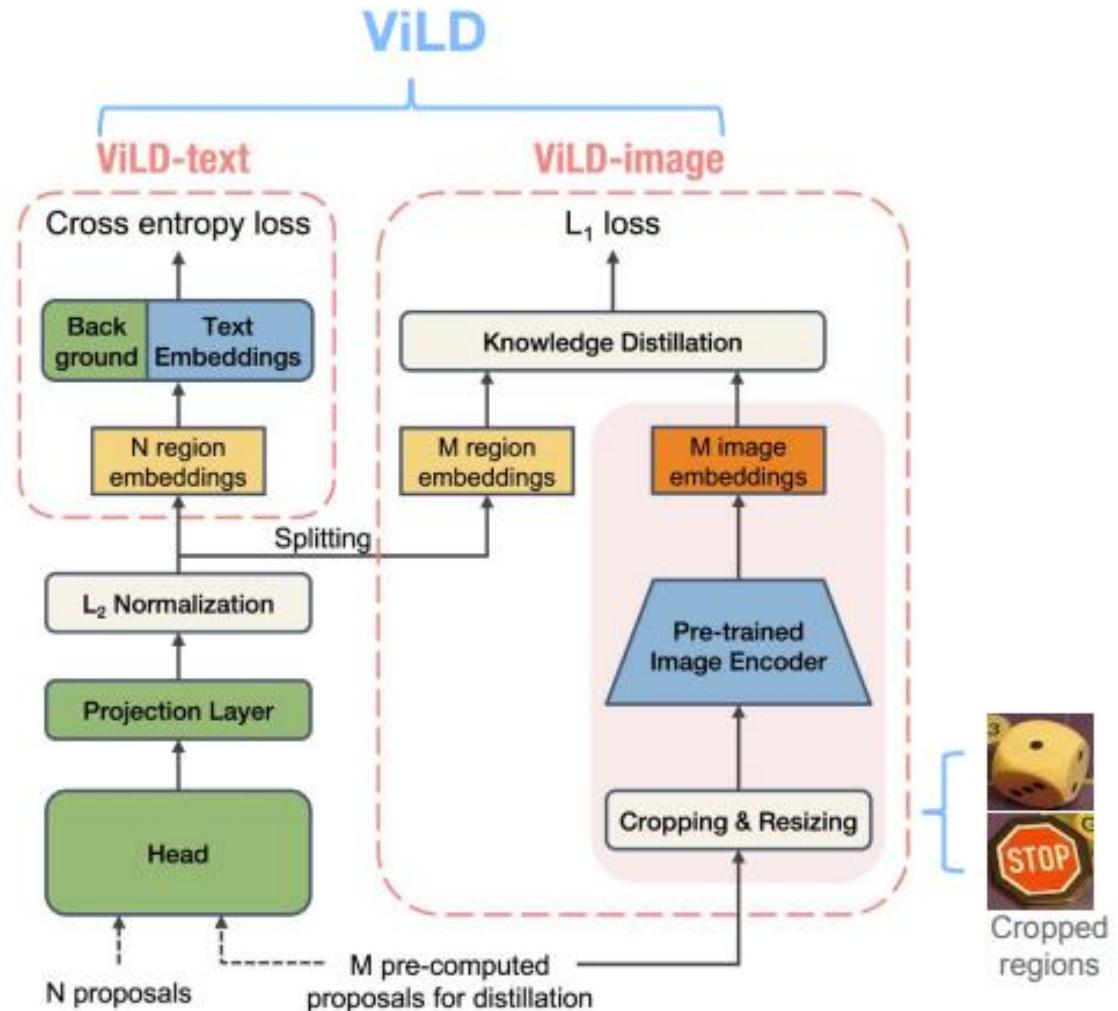
Fast R-CNN. Girshick *et al.*

Can we localize the regions and then recognize them with vision-language models?

# Open-Vocabulary R-CNN

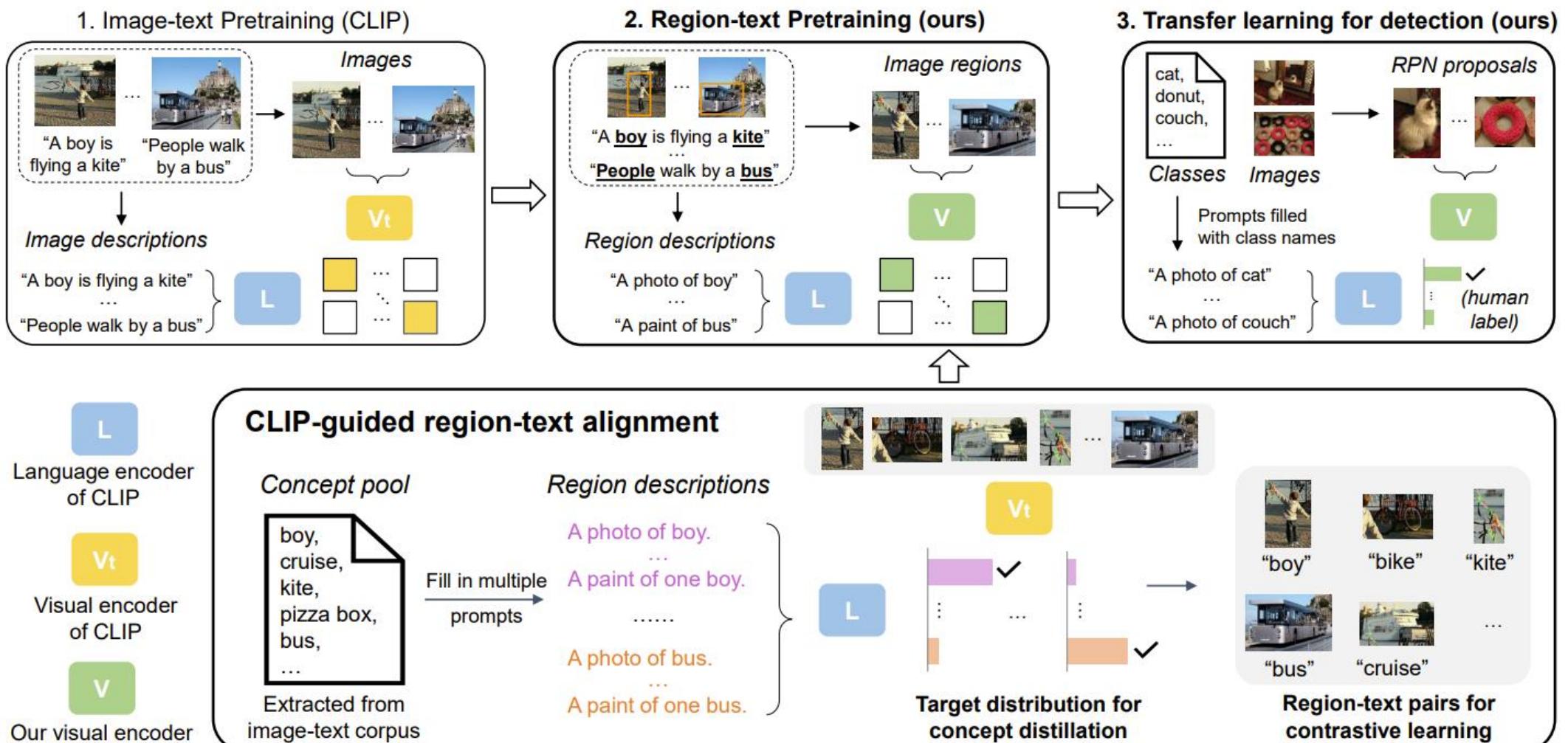


OVR-CNN. Zareian *et al.*



ViLD. Gu *et al.*

# Open-Vocabulary R-CNN



# Conclusion

- **Summary**

- Vision-language learning can help vision significantly, including image-level and region-level vision recognition.
- Vision-language learning leverage image-text pairs, but can goes beyond by integrating supervised learning, self-supervised learning, generative learning, etc.
- Vision-language learning can serves as a generic learning paradigm for zero-shot, few-shot, transfer learning.

- **Challenge**

- How to learn from the large-scale image-text pairs more efficiently?
- How to learn more discriminative visual representations for downstream tasks?
- How to enable fine-grained understanding of visual contents?

**Thanks for your Attention!**