

# Foundations of Vision-Language Models: Concepts and Roadmap

Kaiyang Zhou





# Outline

- History
- Pre-training
- Prompting
- Applications



# Outline

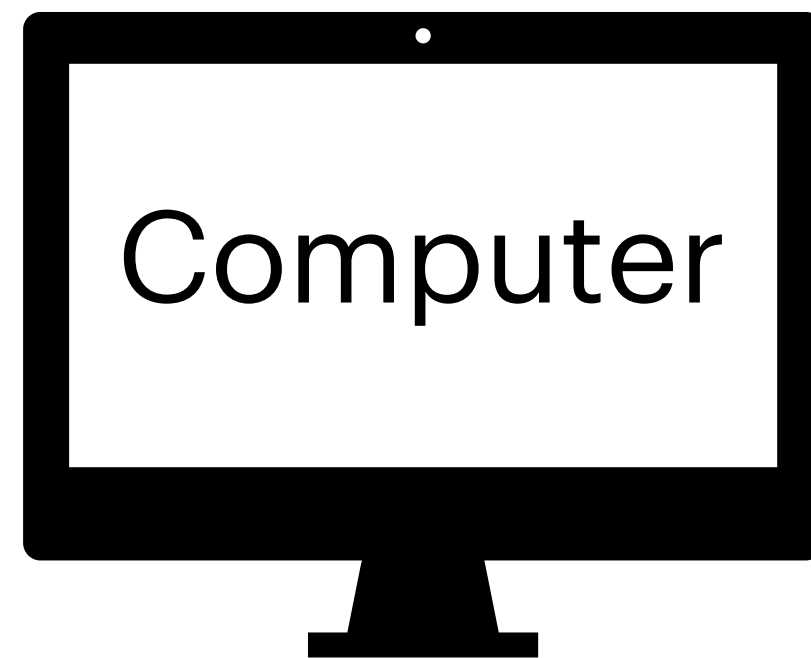
- History
- Pre-training
- Prompting
- Applications



# Vision Models



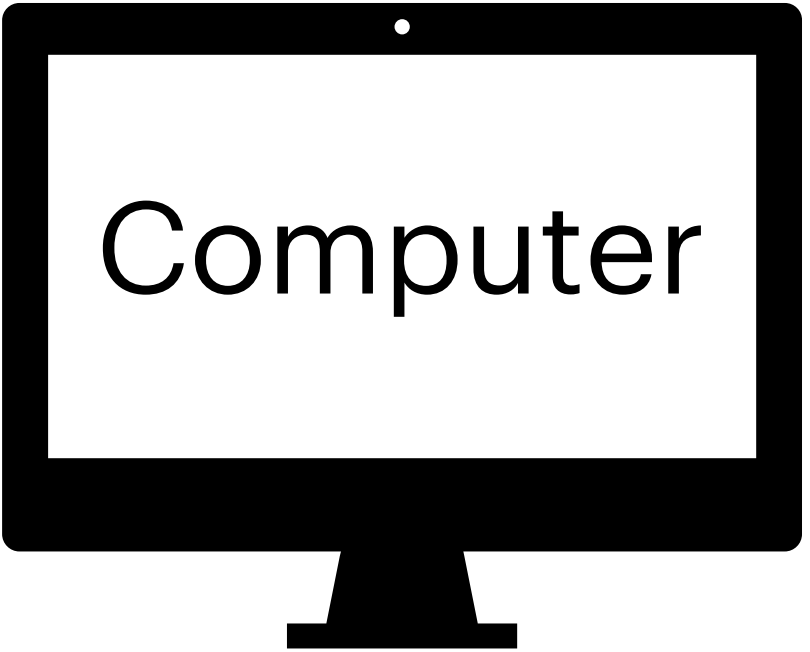
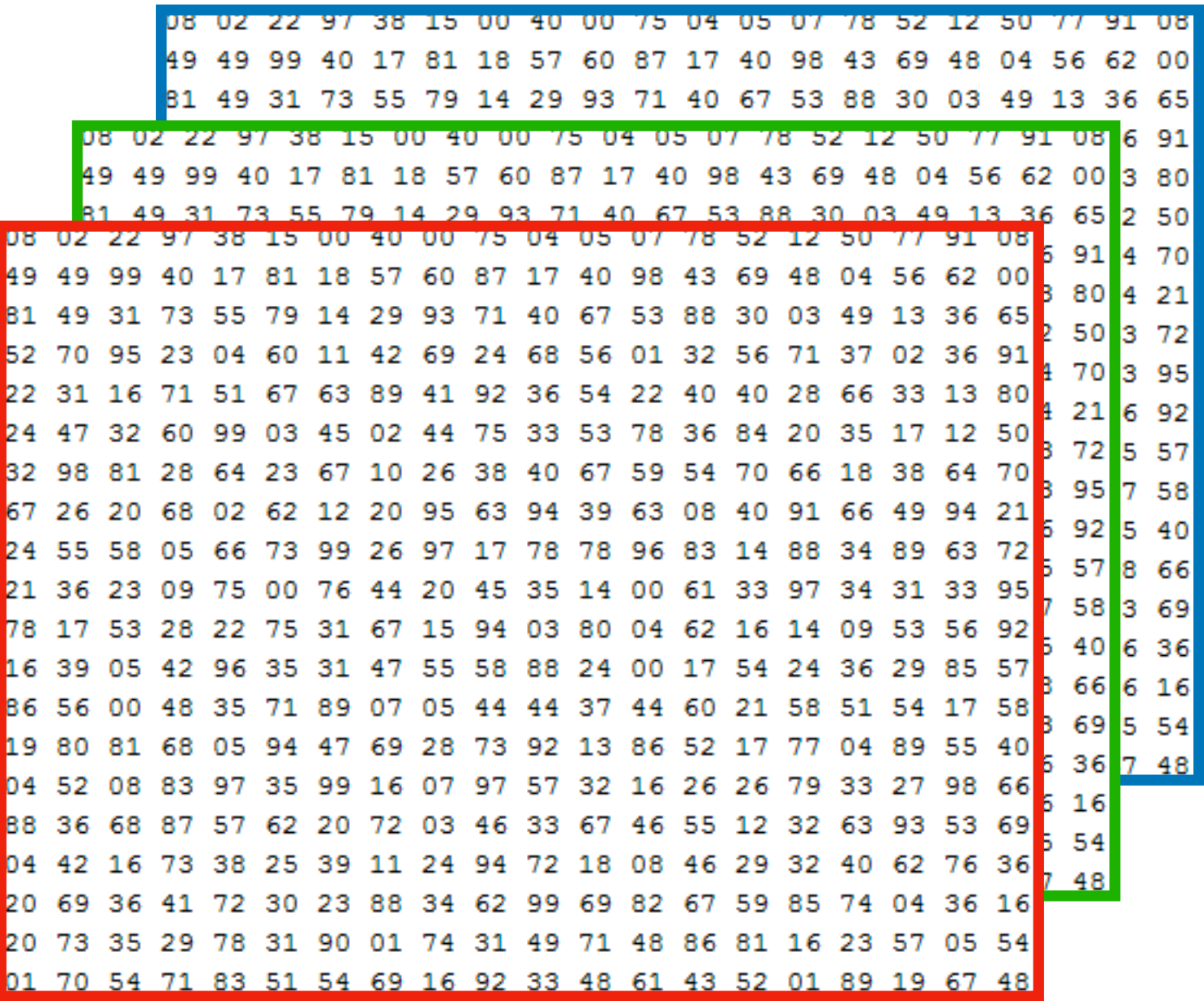
# Teach computers to see





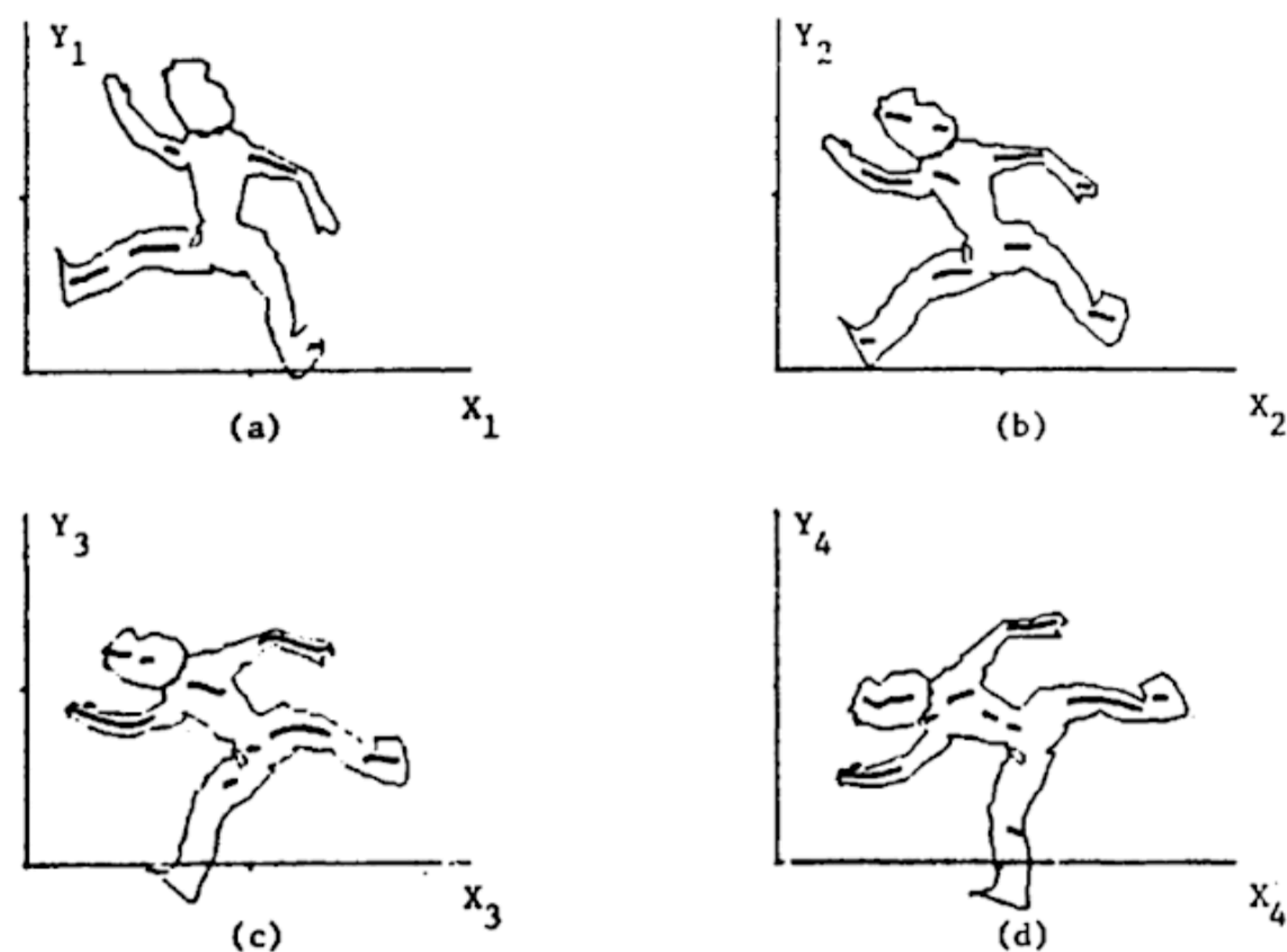
# Teach computers to see

The key question is how to build  
discriminative visual representations

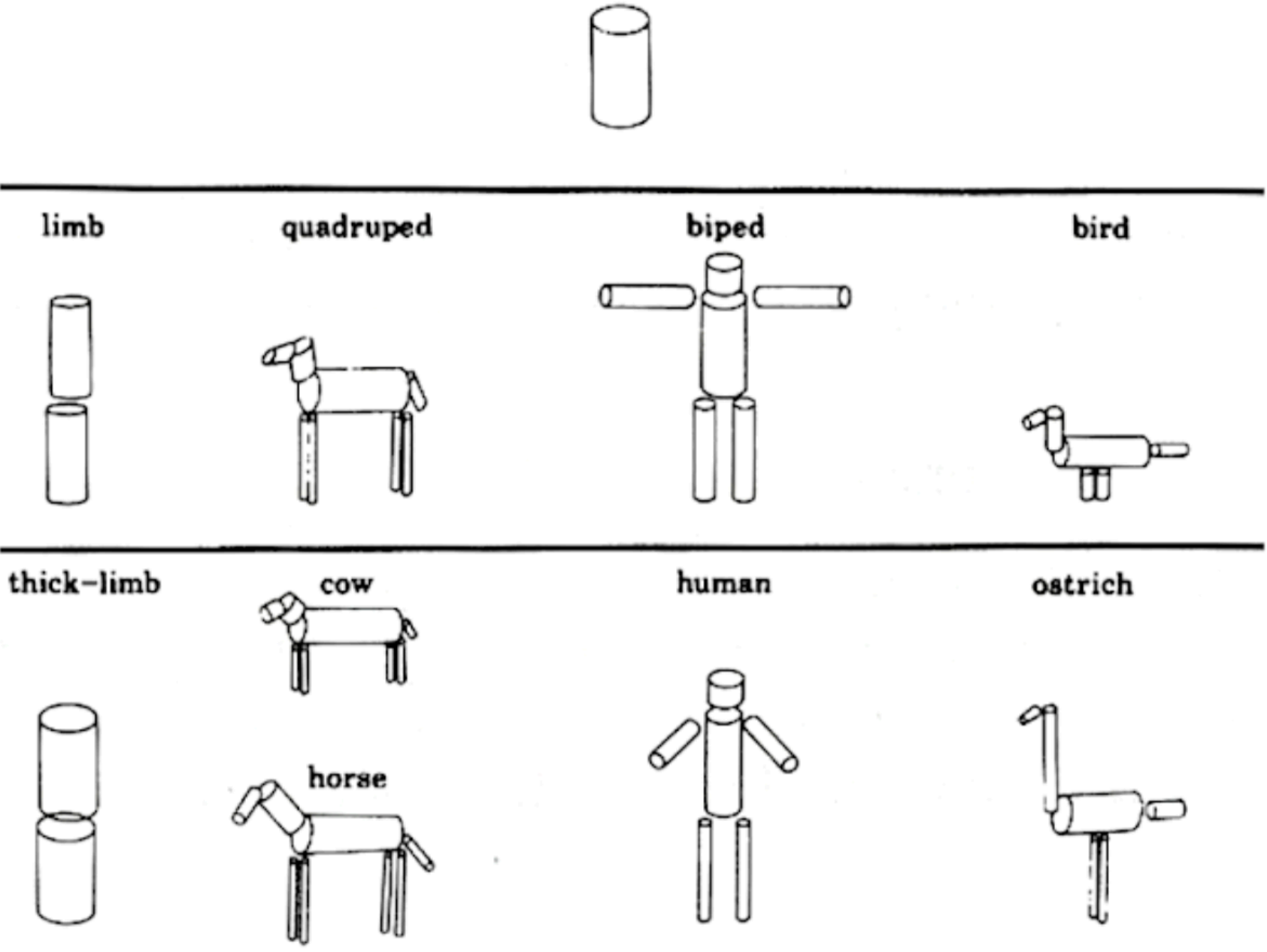




# Visual representations



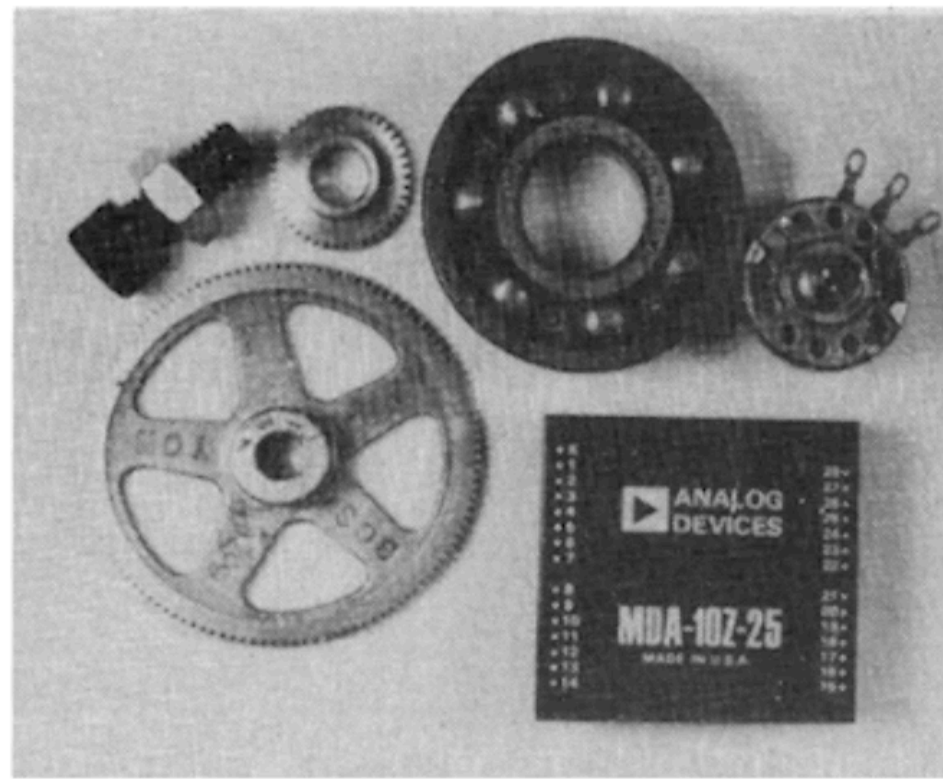
**Curves**  
Nevatia & Binford '77



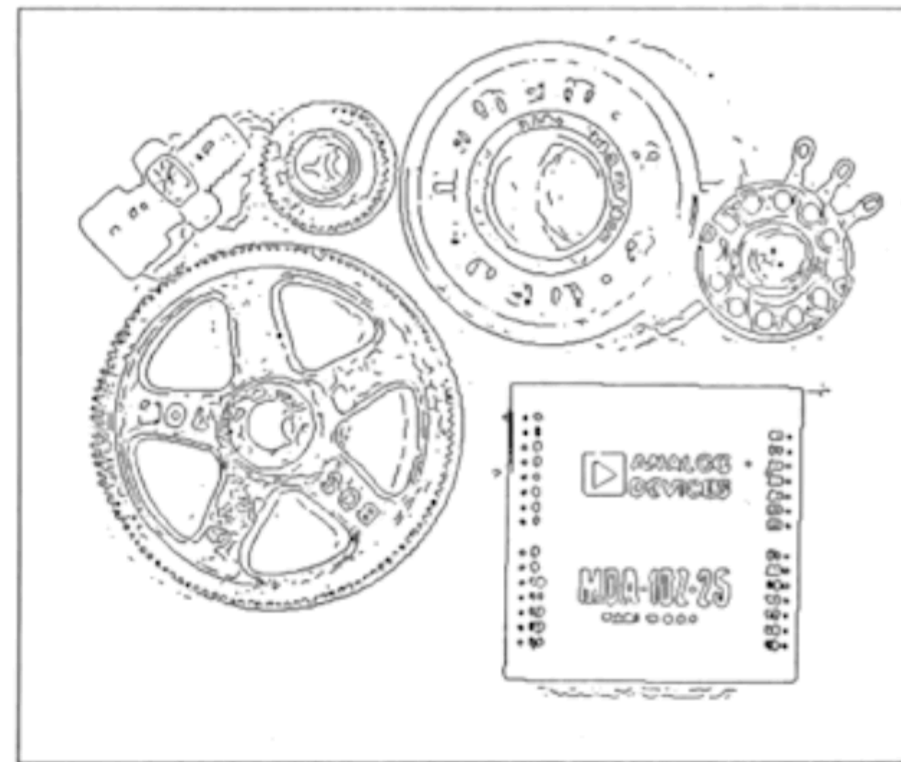
**Cylinders**  
Brooks & Binford '79



# Visual representations



(a)



(b)

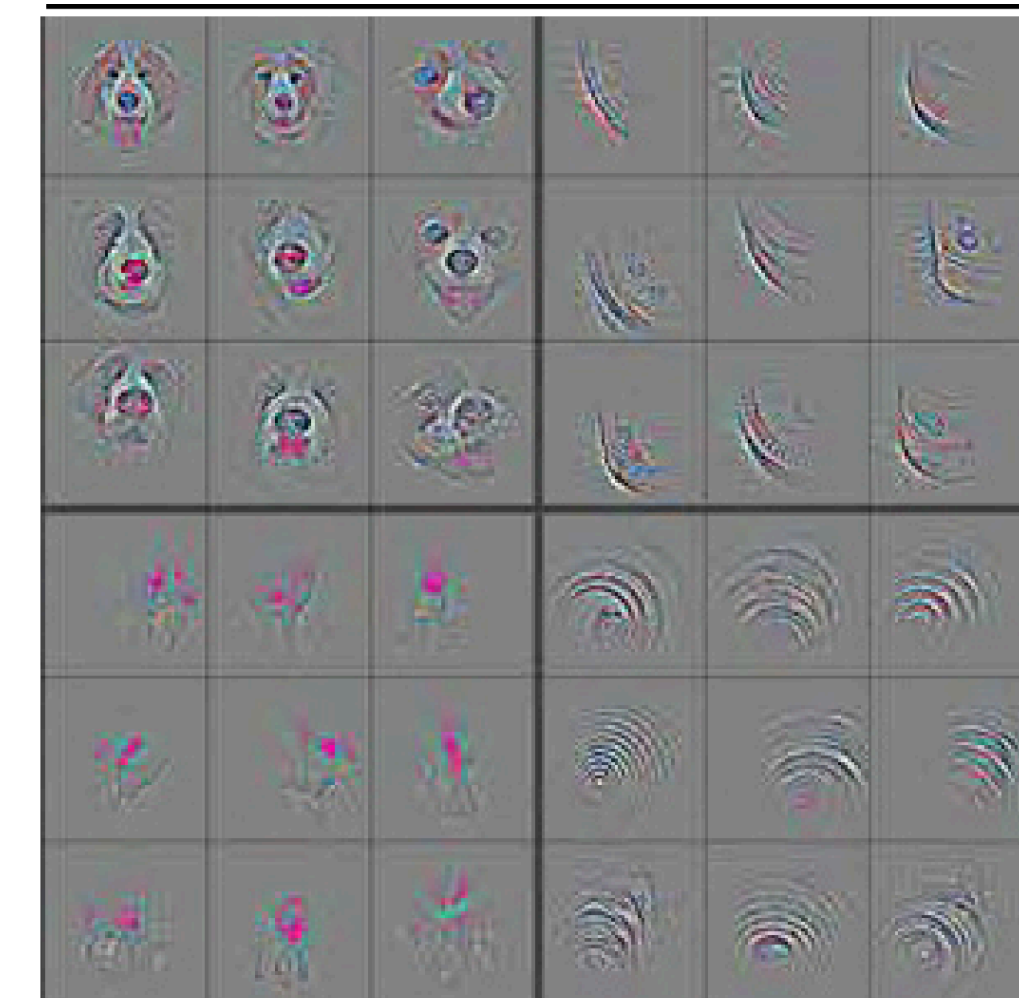
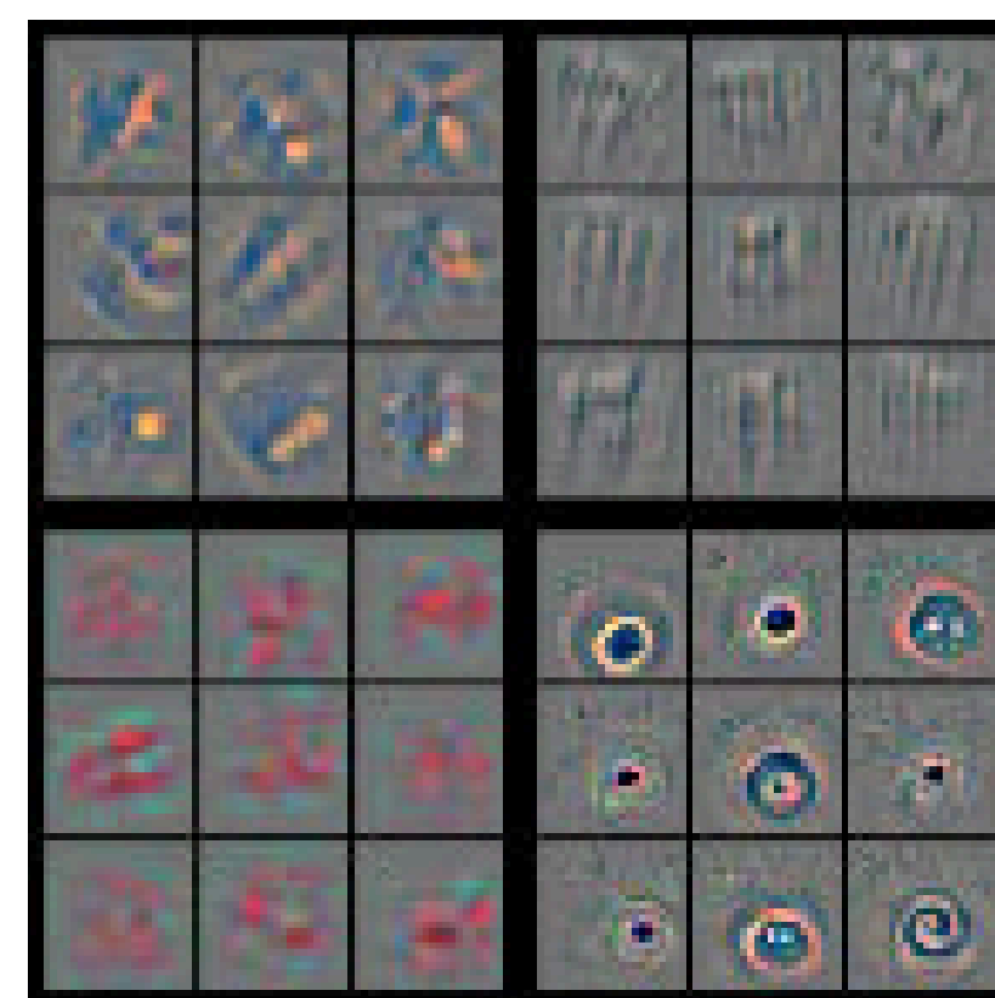
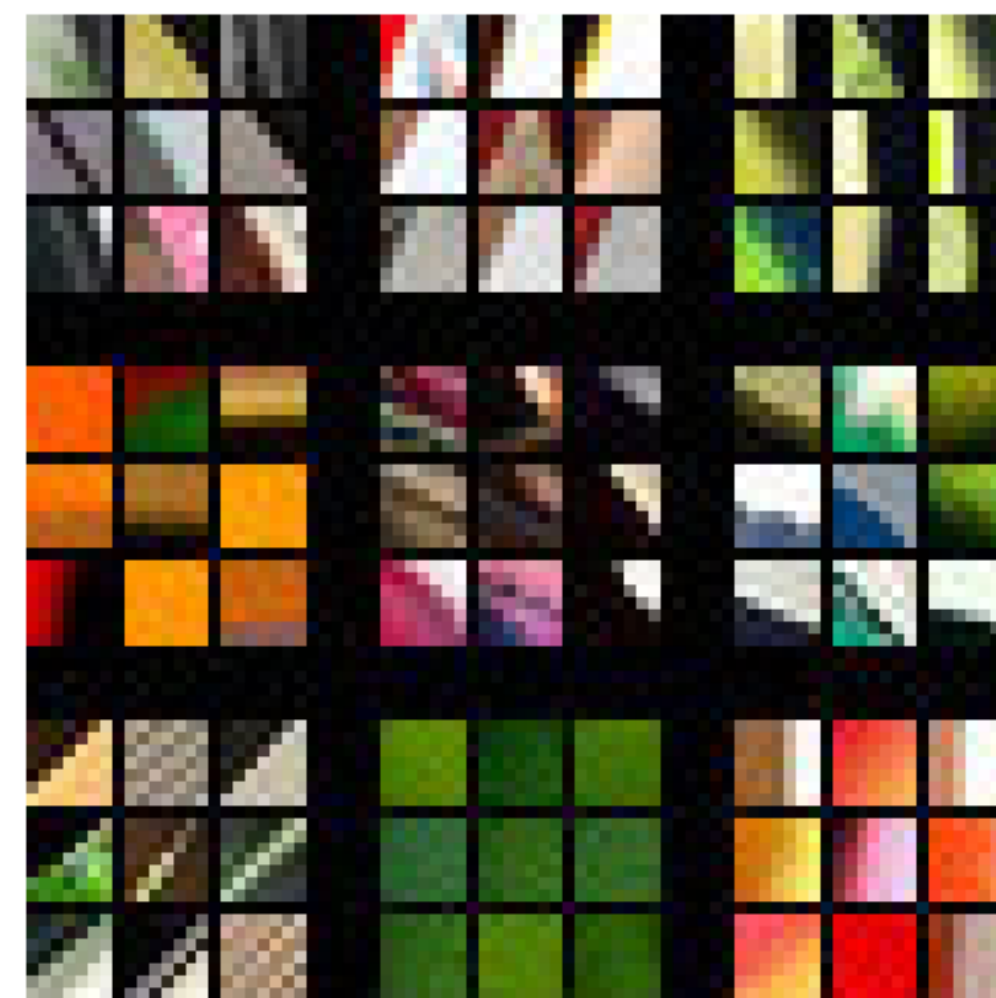
**Edges**  
Canny '86



**Local Features (SIFT)**  
Lowe '99

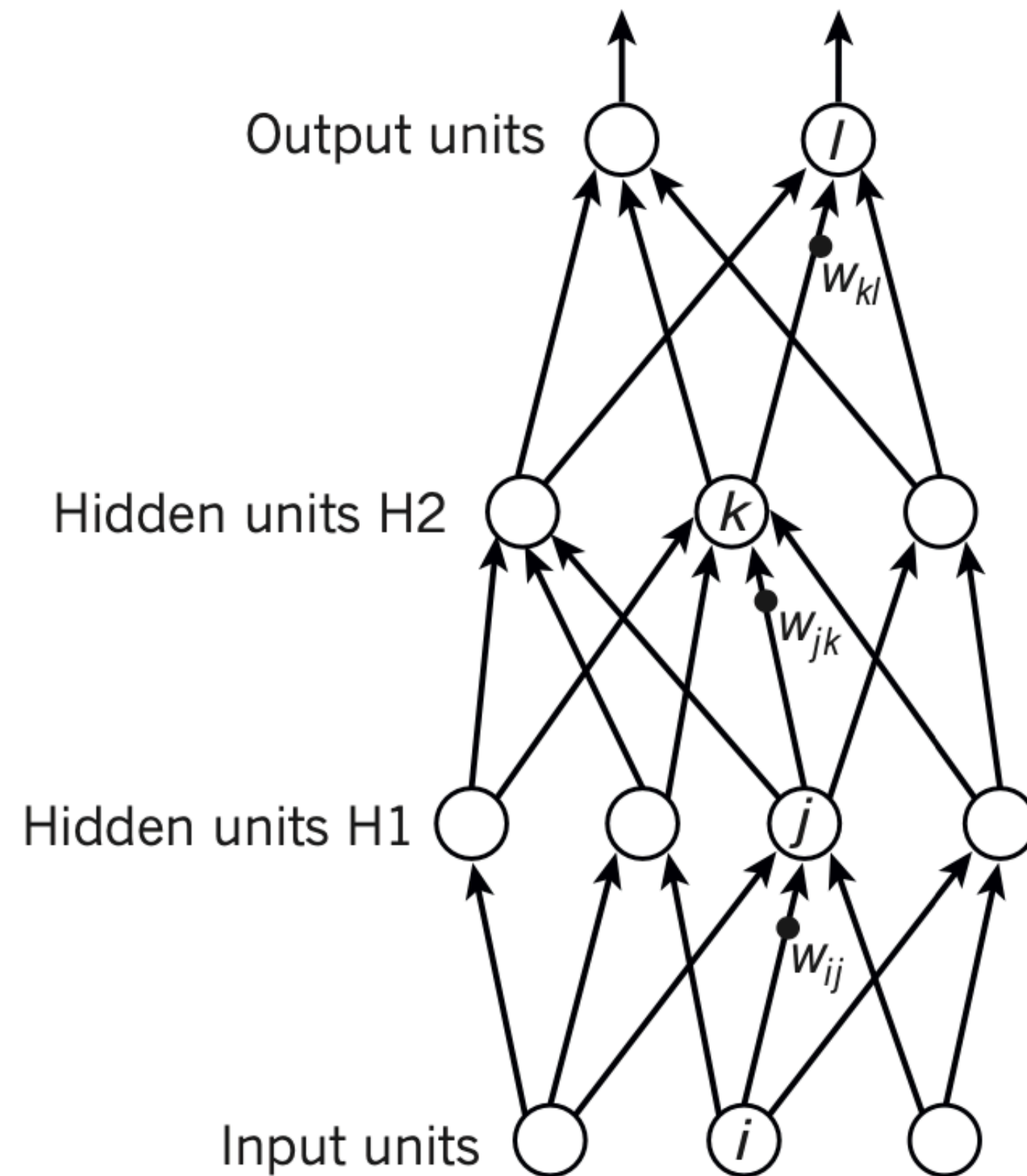


# Visual representations



# Deep neural network

**c**



$$y_l = f(z_l)$$

$$z_l = \sum_{k \in H2} w_{kl} y_k$$

$$y_k = f(z_k)$$

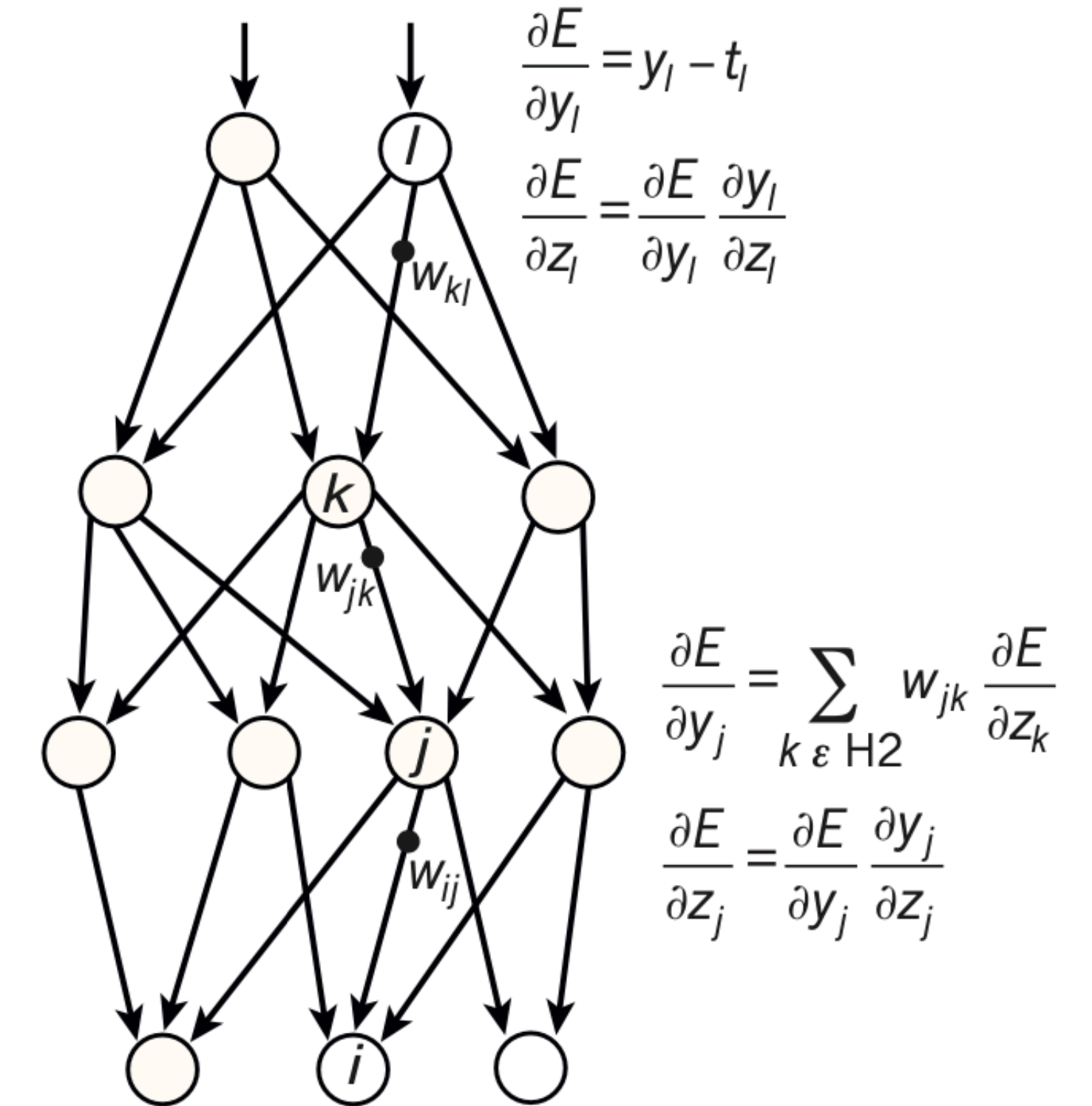
$$z_k = \sum_{j \in H1} w_{jk} y_j$$

$$y_j = f(z_j)$$

$$z_j = \sum_{i \in \text{Input}} w_{ij} x_i$$

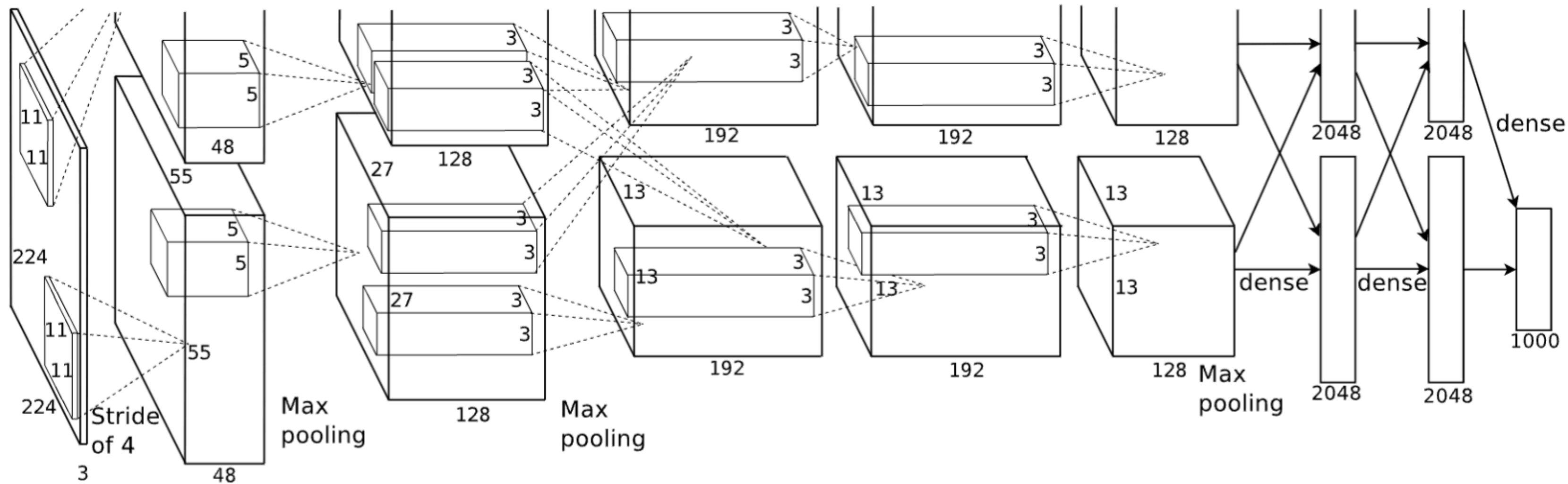
**d**

Compare outputs with correct answer to get error derivatives





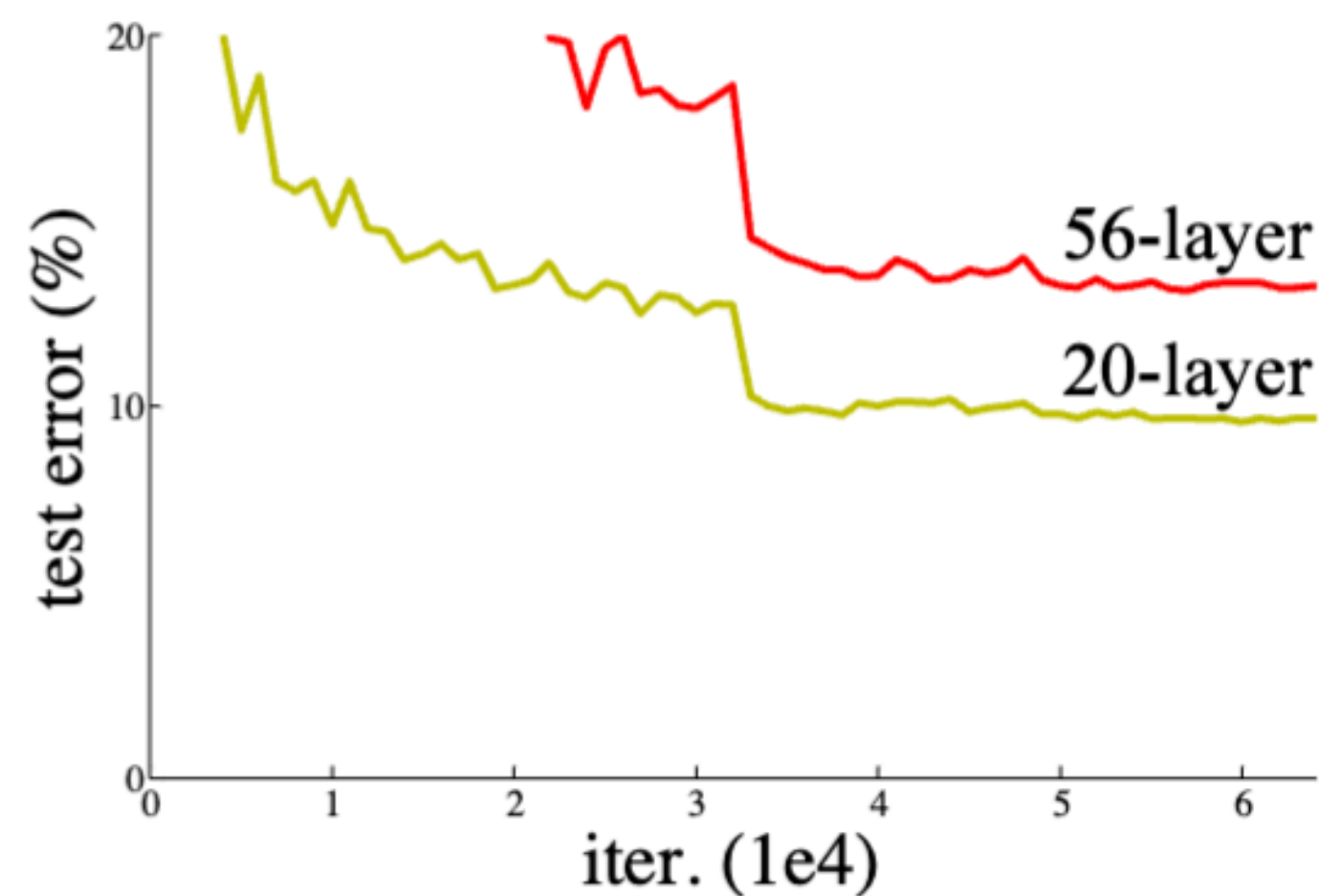
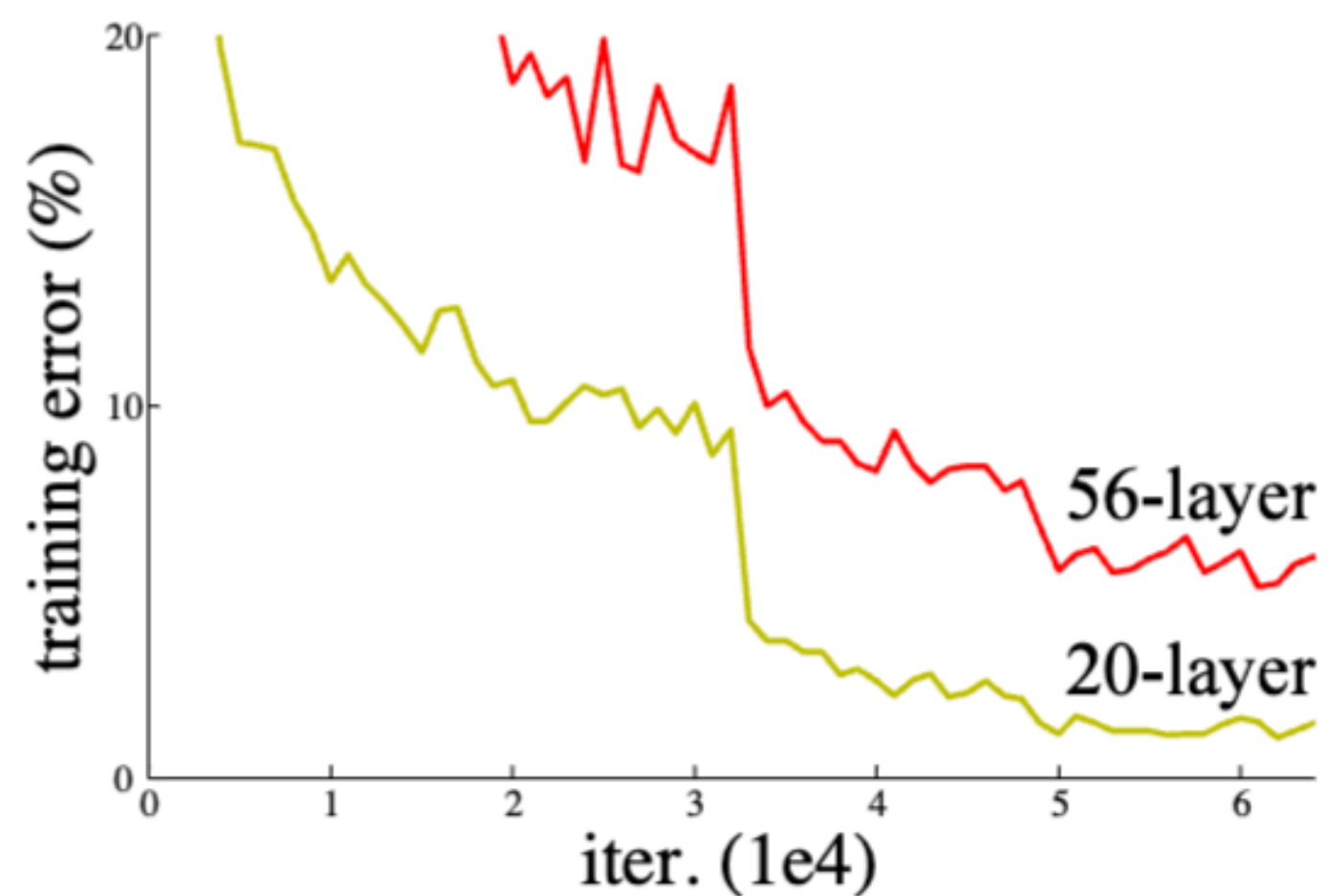
# Convolutional neural network



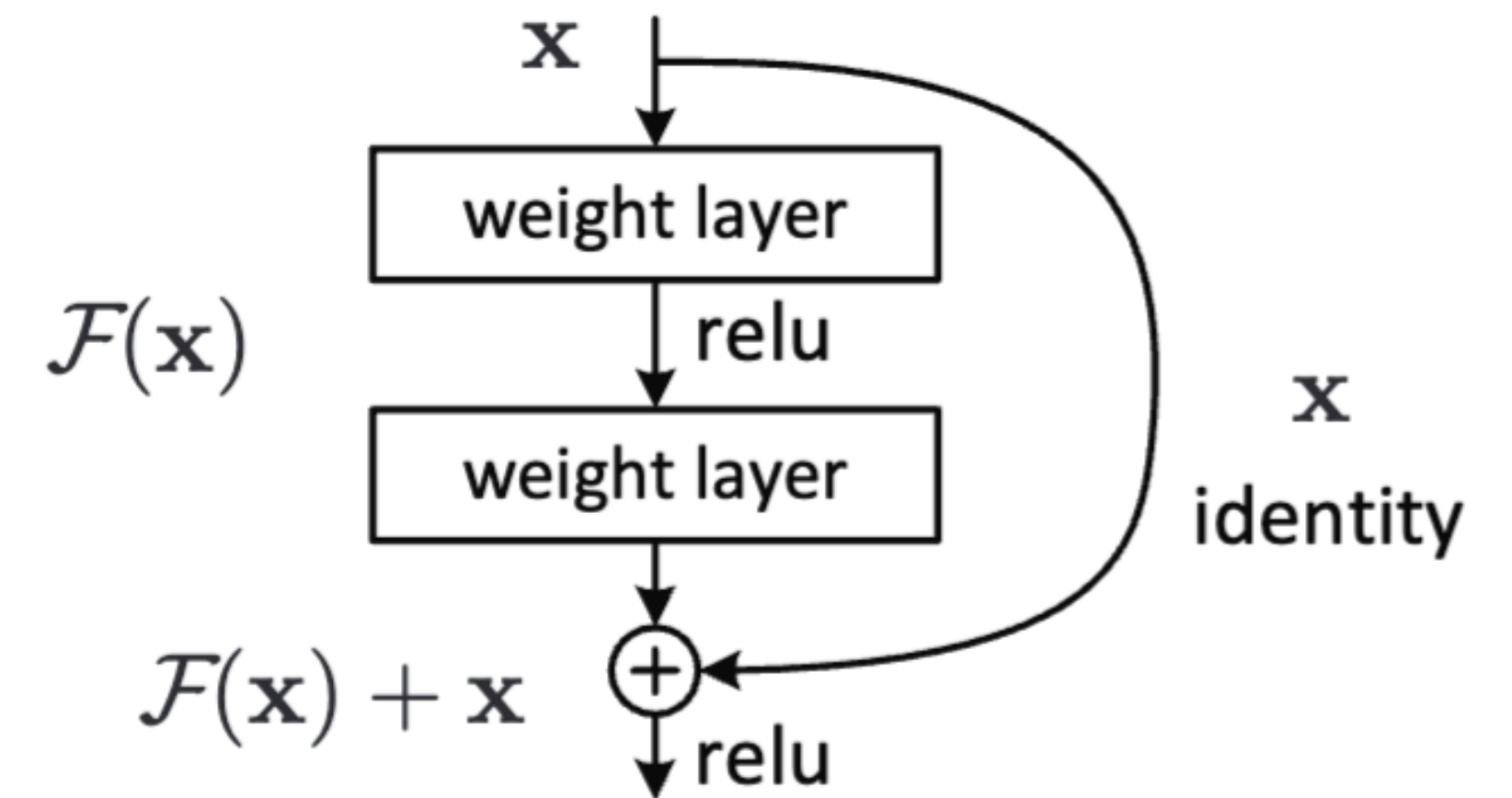
- ReLU non-linearity
- Feature normalization
- Dropout
- Data augmentation
- Multi-GPU training

# Deep residual network

Problem: Deeper networks are difficult to train

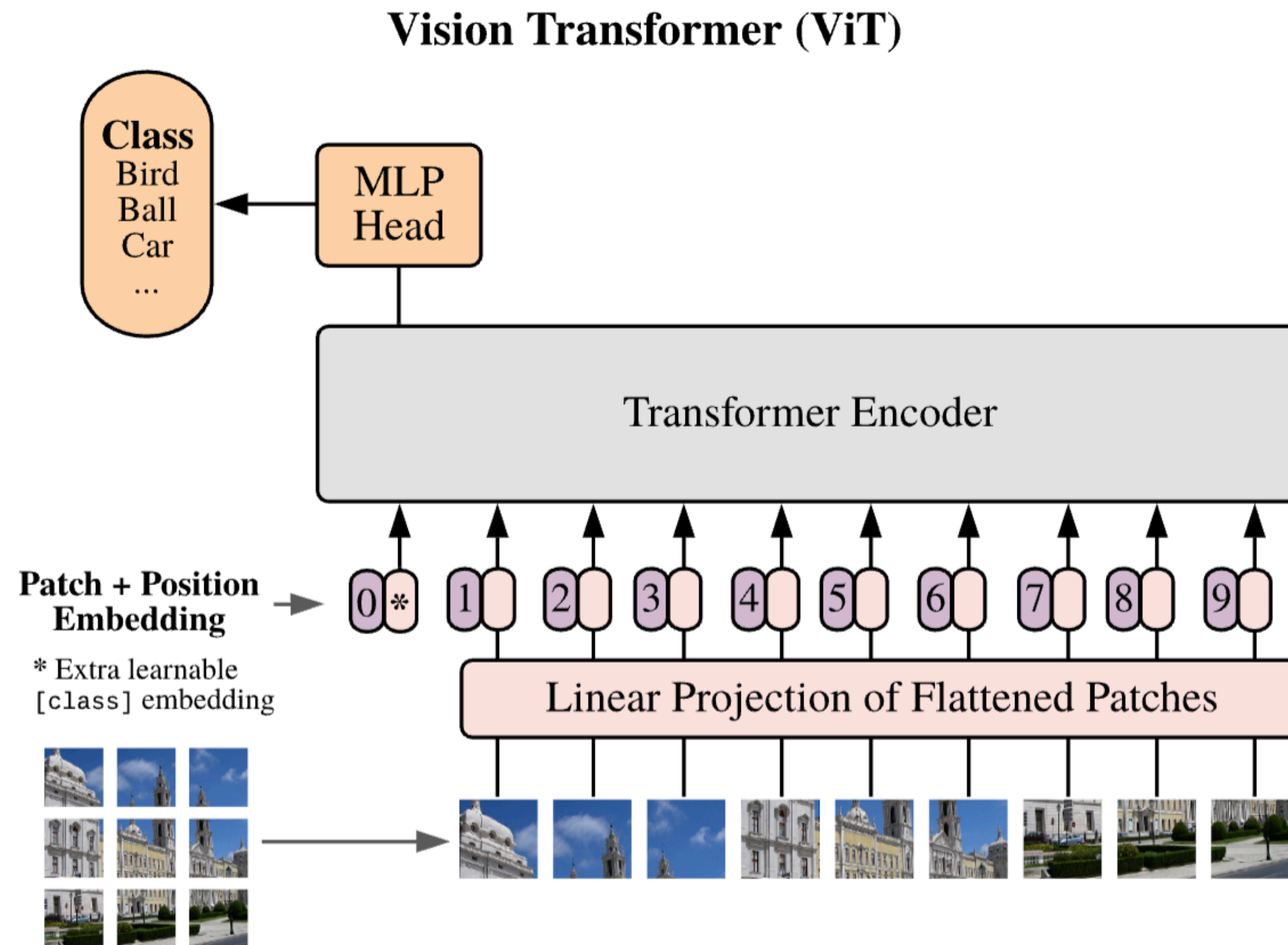


Solution: Residual connection

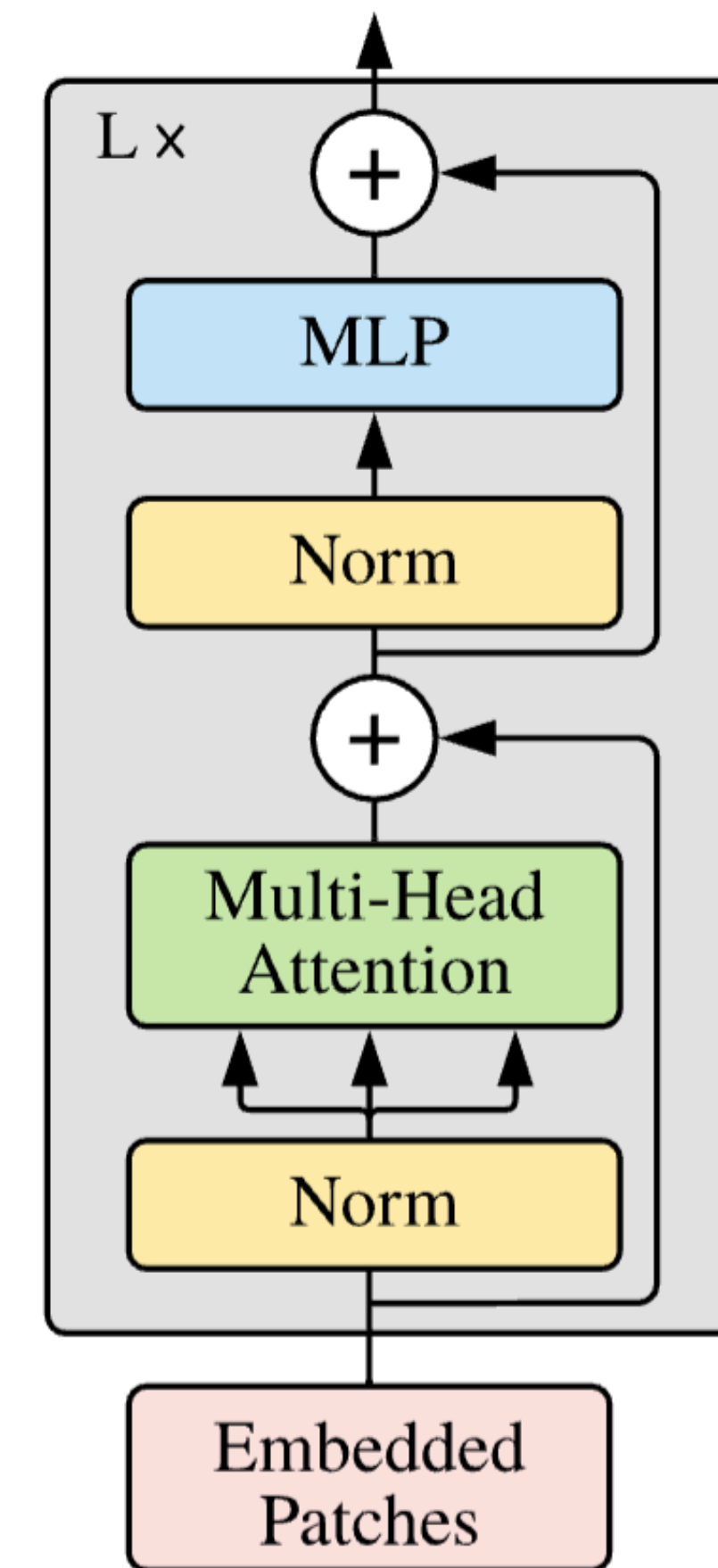




# Vision Transformer



## Transformer Encoder







sky

tree

person 100%

person 99%

person 93%

person 99%

person 74%

bottle 77%  
bicycle 98%

bicycle 96%

bicycle 91%

person 69%

bicycle 62%

person 100%

person 97%

backpack 90%

cell phone 95%

person 67% 91%

person 92%  
person 69%

person 56%

person 99%

road



# Language Models

# Word representations (embeddings)

Sentiment Analysis

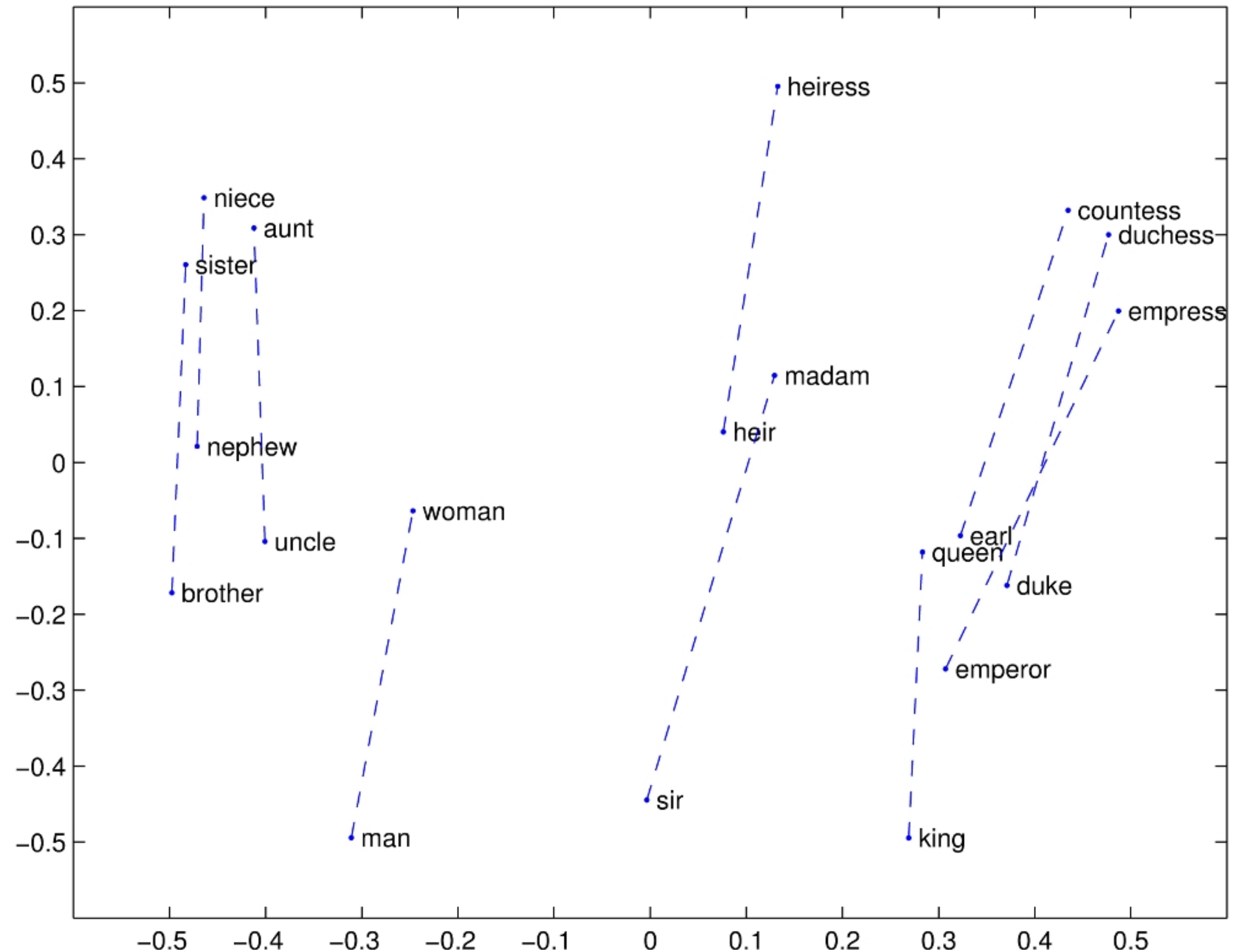
Machine Translation

Text Summarization

Email Filtering

Chatbot

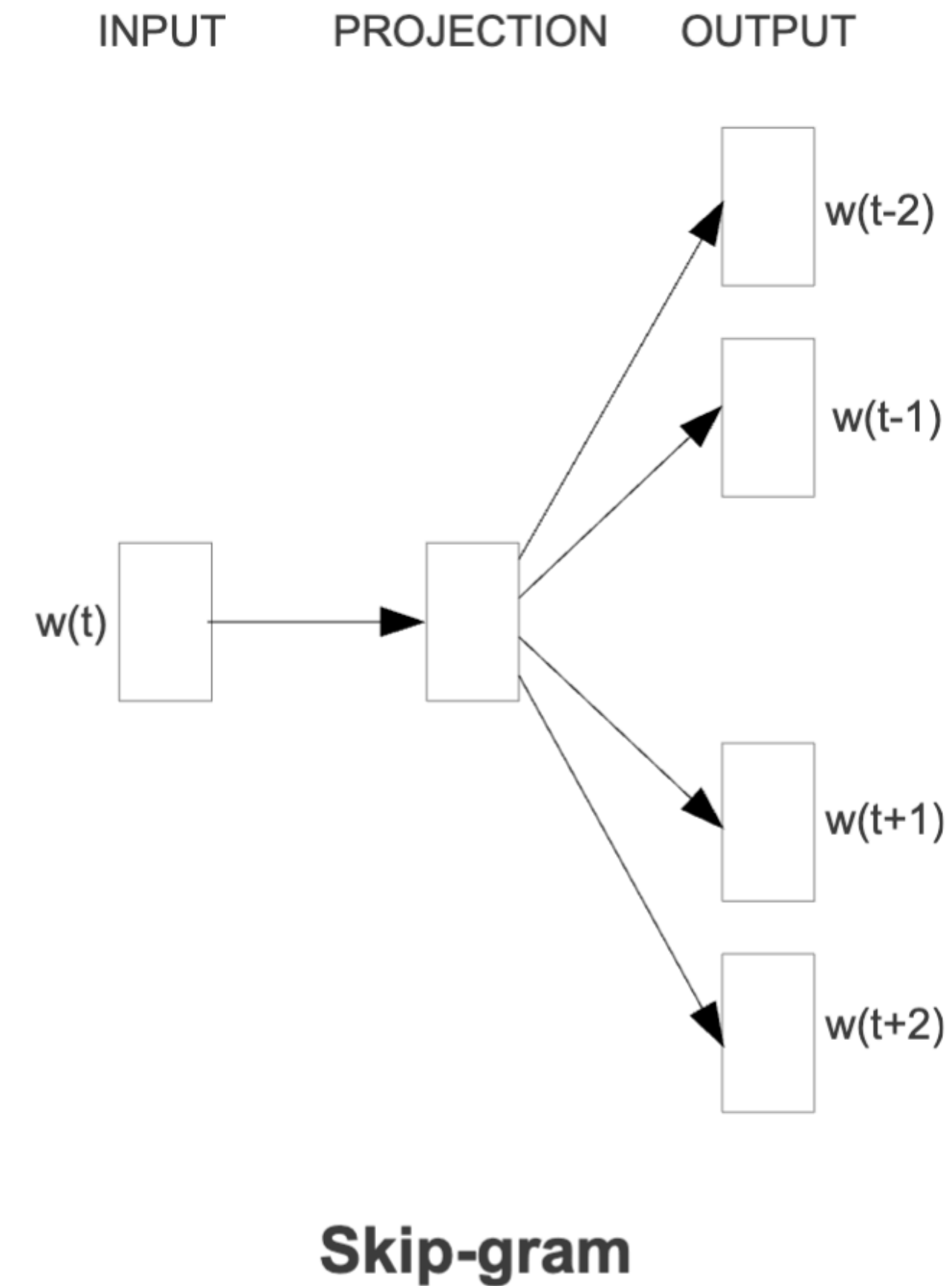
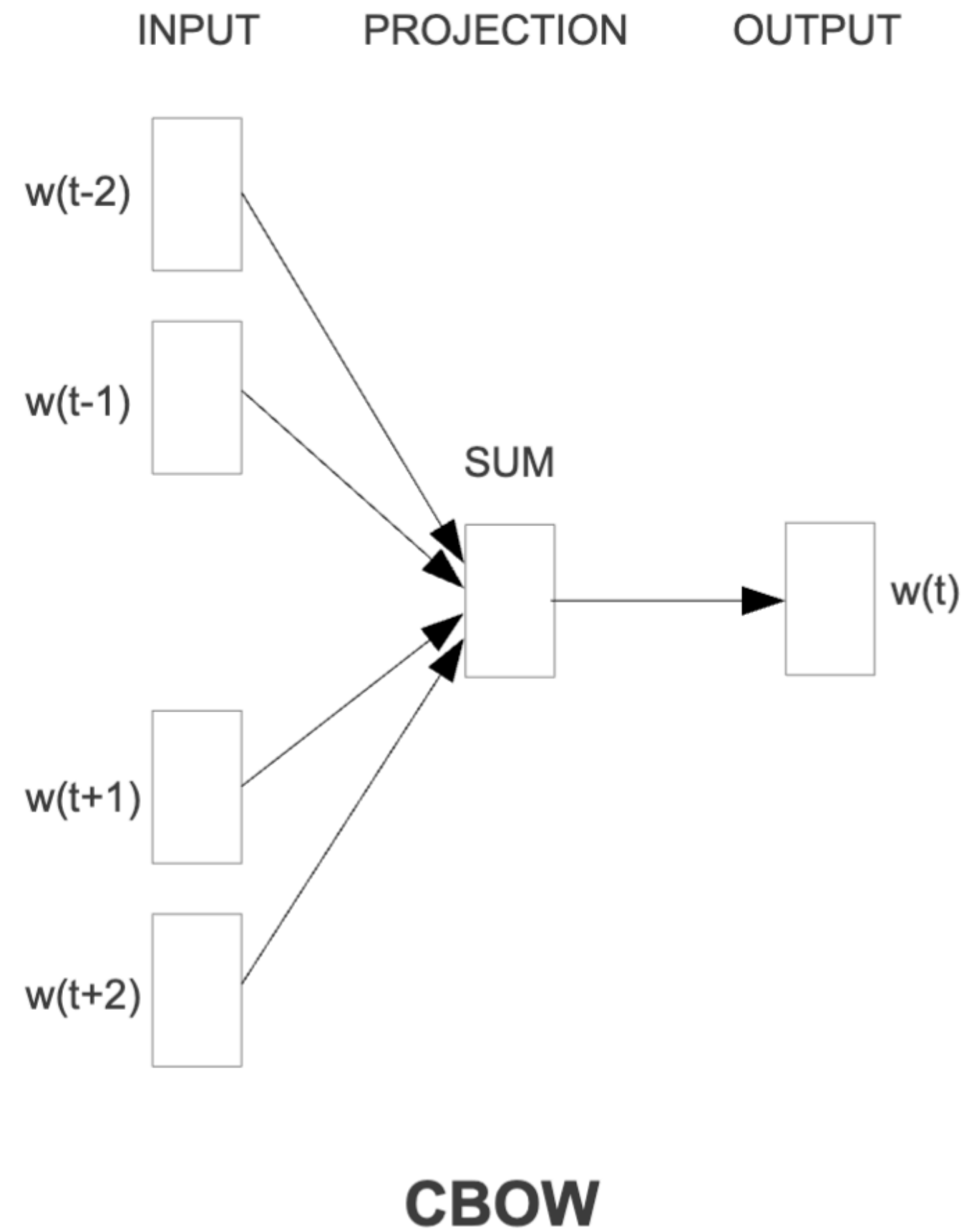
etc.





# Word2vec

Predicts the current word based on the context

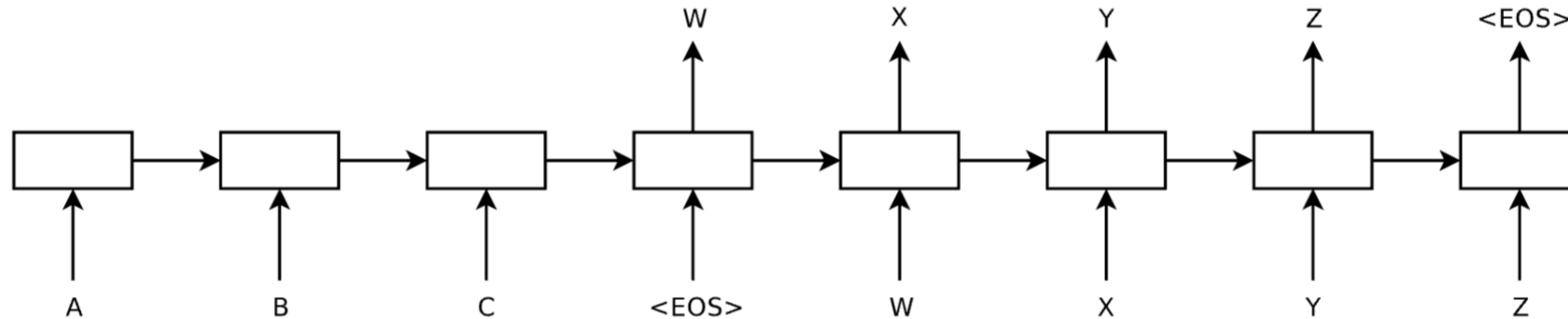


Predicts context words given the current word

Window Size	Text	Skip-grams
2	[ The <u>wide</u> road shimmered ] in the hot sun.	wide, the wide, road wide, shimmered
	The [ wide road <u>shimmered</u> in the ] hot sun.	shimmered, wide shimmered, road shimmered, in shimmered, the
	The wide road shimmered in [ the hot <u>sun</u> ].	sun, the sun, hot
3	[ The <u>wide</u> road shimmered in ] the hot sun.	wide, the wide, road wide, shimmered wide, in
	[ The wide road <u>shimmered</u> in the hot ] sun.	shimmered, the shimmered, wide shimmered, road shimmered, in shimmered, the shimmered, hot
	The wide road shimmered [ in the hot <u>sun</u> ].	sun, in sun, the sun, hot



# Seq2seq



## Encoder-Decoder

- Autoregressive
- “Large” NN
- “Large” datasets

# GPT

# Autoregressive training (scaled up)

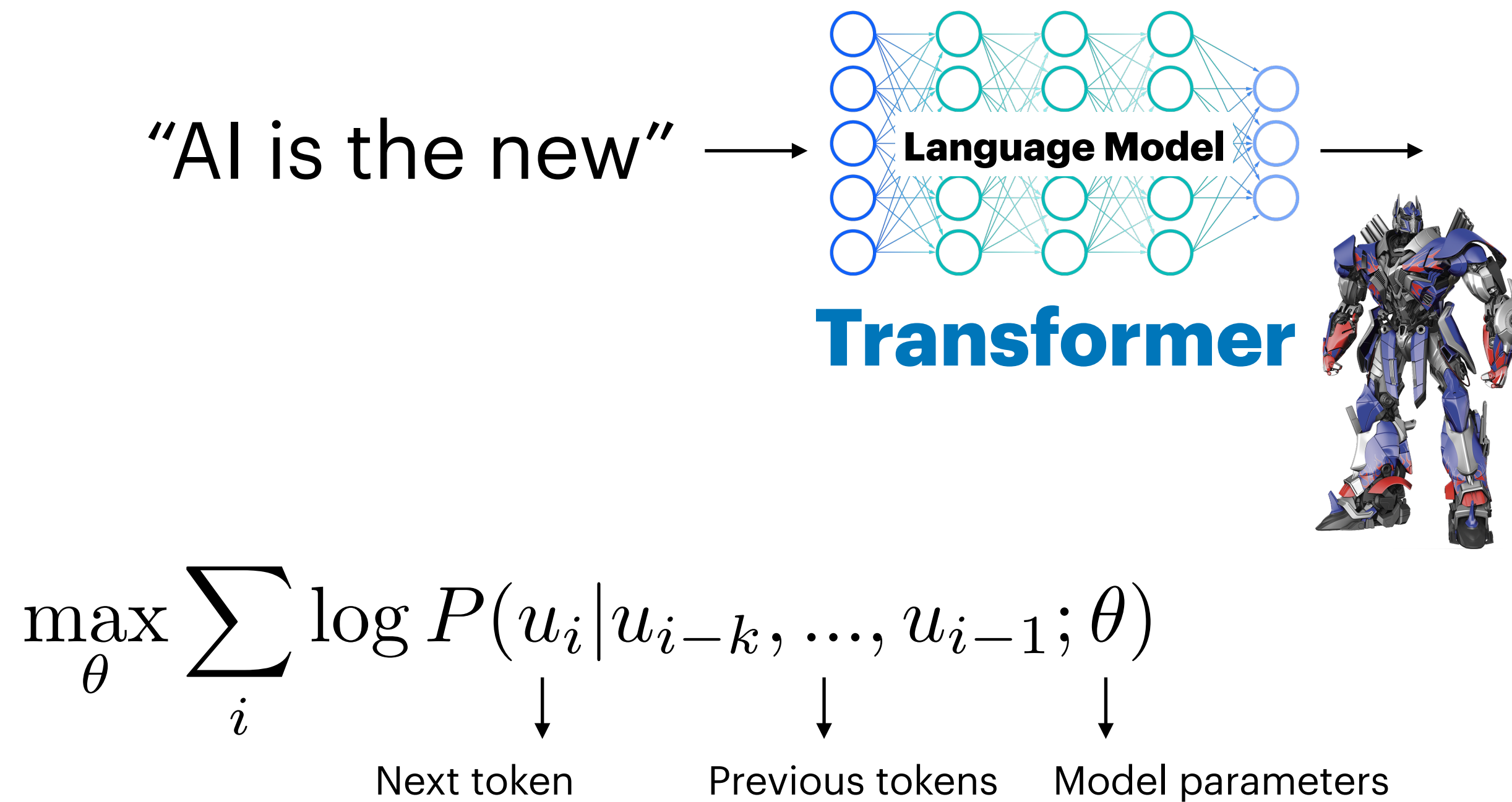
# Improving Language Understanding by Generative Pre-Training

**Alec Radford**  
OpenAI  
alec@openai.com

**Karthik Narasimhan**  
OpenAI  
karthikn@openai.com

**Tim Salimans**  
OpenAI  
tim@openai.com

**Ilya Sutskever**  
OpenAI  
ilyasu@openai.com



Word	Probability
a	0.0000001
ah	0.0000002
...	...
elect	0.0000022
electricity	0.03
...	...
zip	0.0000034

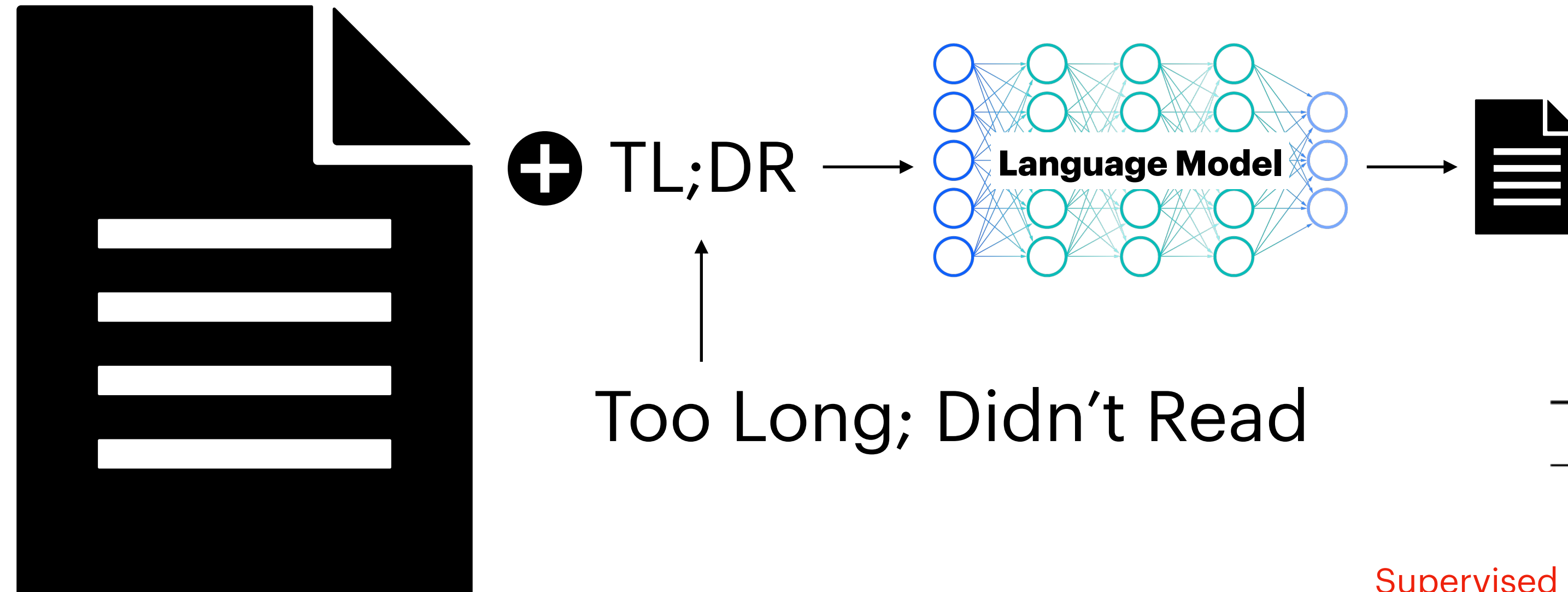


# GPT

## Language Models are Unsupervised Multitask Learners

Alec Radford<sup>\*1</sup> Jeffrey Wu<sup>\*1</sup> Rewon Child<sup>1</sup> David Luan<sup>1</sup> Dario Amodei<sup>\*\*1</sup> Ilya Sutskever<sup>\*\*1</sup>

## Zero-shot (autocomplete)



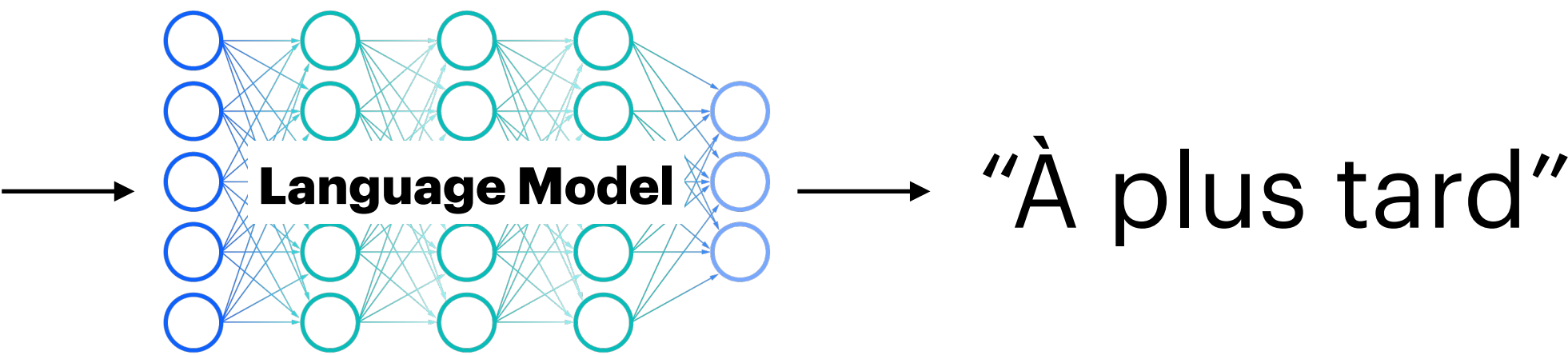
	R-1	R-2	R-L	R-AVG
Bottom-Up Sum	<b>41.22</b>	<b>18.68</b>	<b>38.34</b>	<b>32.75</b>
Lede-3	40.38	17.66	36.62	31.55
Supervised Seq2Seq + Attn	31.33	11.81	28.83	23.99
GPT-2 TL;DR:	29.34	8.27	26.58	21.40
Random baseline Random-3	28.78	8.63	25.52	20.98
GPT-2 no hint	21.58	4.03	19.47	15.03

Table 4. Summarization performance as measured by ROUGE F1 metrics on the CNN and Daily Mail dataset. Bottom-Up Sum is the SOTA model from (Gehrmann et al., 2018)

# GPT

## Few-shot in-context learning

“Translate English into French:  
Hello => Bonjour  
Thank you => Merci  
Goodbye => Au revoir  
Excuse me => \_\_\_\_”



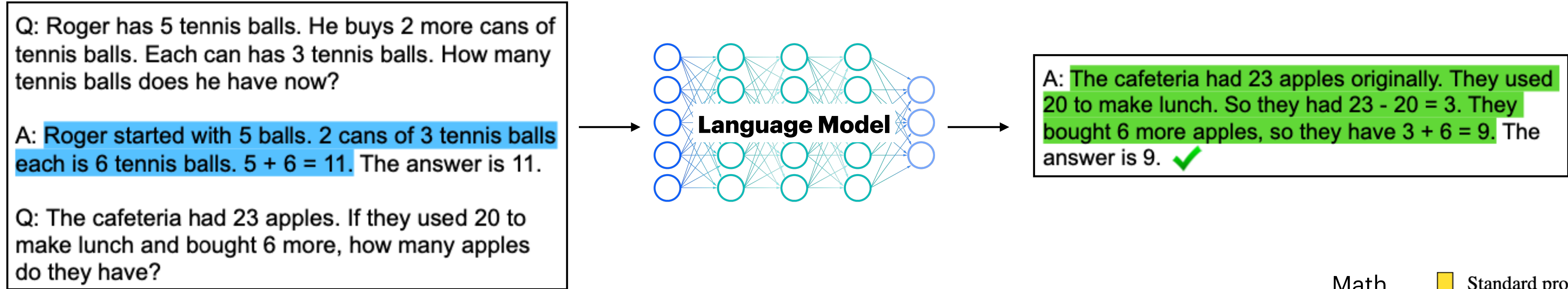
Setting	En→Fr	Fr→En	En→De	De→En	En→Ro	Ro→En
SOTA (Supervised)	<b>45.6<sup>a</sup></b>	35.0 <sup>b</sup>	<b>41.2<sup>c</sup></b>	40.2 <sup>d</sup>	<b>38.5<sup>e</sup></b>	<b>39.9<sup>e</sup></b>
XLM [LC19]	33.4	33.3	26.4	34.3	33.3	31.8
MASS [STQ <sup>+</sup> 19]	<u>37.5</u>	34.9	28.3	35.2	<u>35.2</u>	33.1
mBART [LGG <sup>+</sup> 20]	-	-	<u>29.8</u>	34.0	35.0	30.5
GPT-3 Zero-Shot	25.2	21.2	24.6	27.2	14.1	19.9
GPT-3 One-Shot	28.3	33.7	26.2	30.4	20.6	38.6
GPT-3 Few-Shot	32.6	<u>39.2</u>	29.7	<u>40.6</u>	21.0	<u>39.5</u>

Good at X -> En

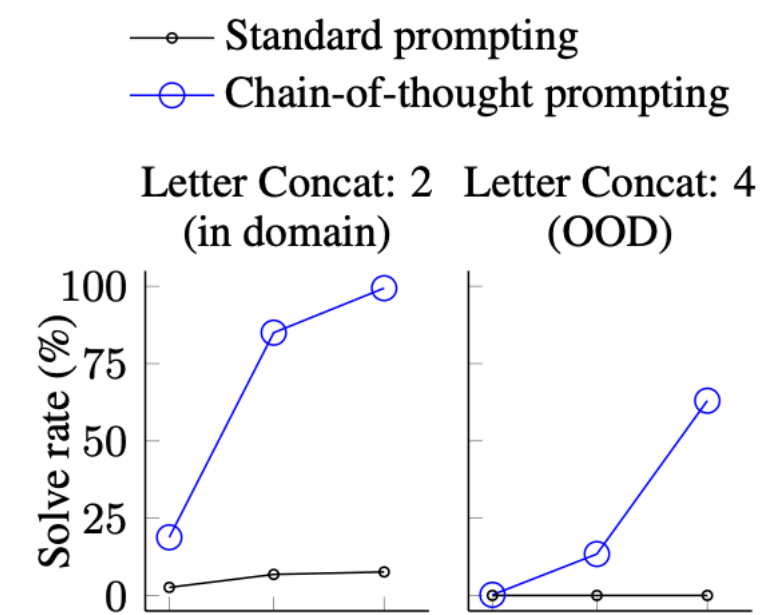
### Language Models are Few-Shot Learners

Tom B. Brown*		Benjamin Mann*		Nick Ryder*		Melanie Subbiah*	
Jared Kaplan†	Prafulla Dhariwal	Arvind Neelakantan		Pranav Shyam		Girish Sastry	
Amanda Askell	Sandhini Agarwal	Ariel Herbert-Voss		Gretchen Krueger		Tom Henighan	
Rewon Child	Aditya Ramesh	Daniel M. Ziegler		Jeffrey Wu		Clemens Winter	
Christopher Hesse	Mark Chen	Eric Sigler		Mateusz Litwin		Scott Gray	
Benjamin Chess		Jack Clark		Christopher Berner			
Sam McCandlish		Alec Radford		Ilya Sutskever		Dario Amodei	
OpenAI							

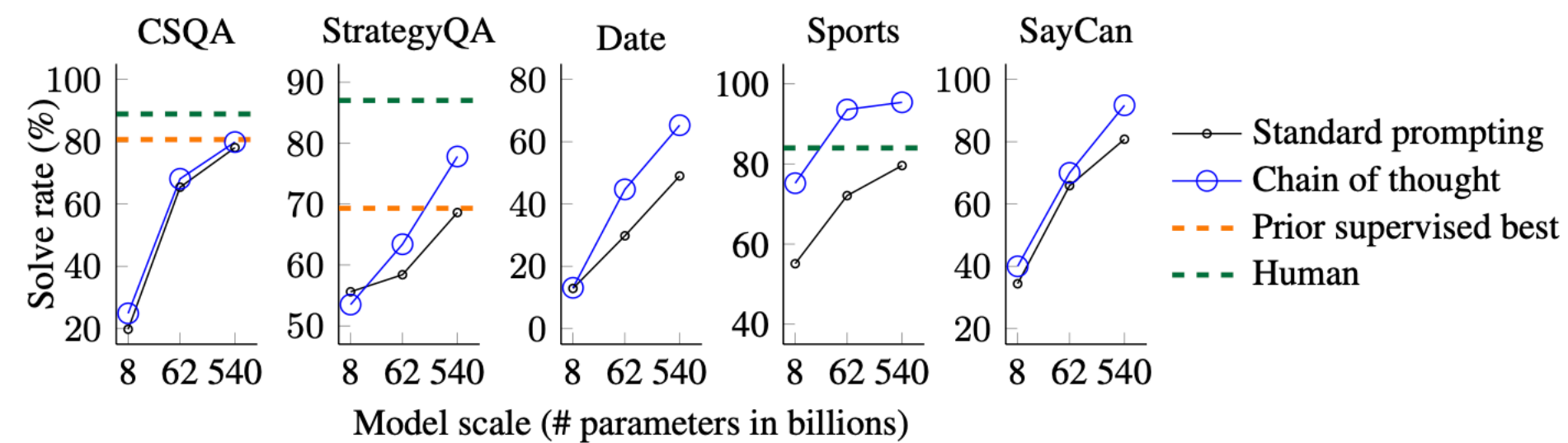
# Chain-of-thought prompting



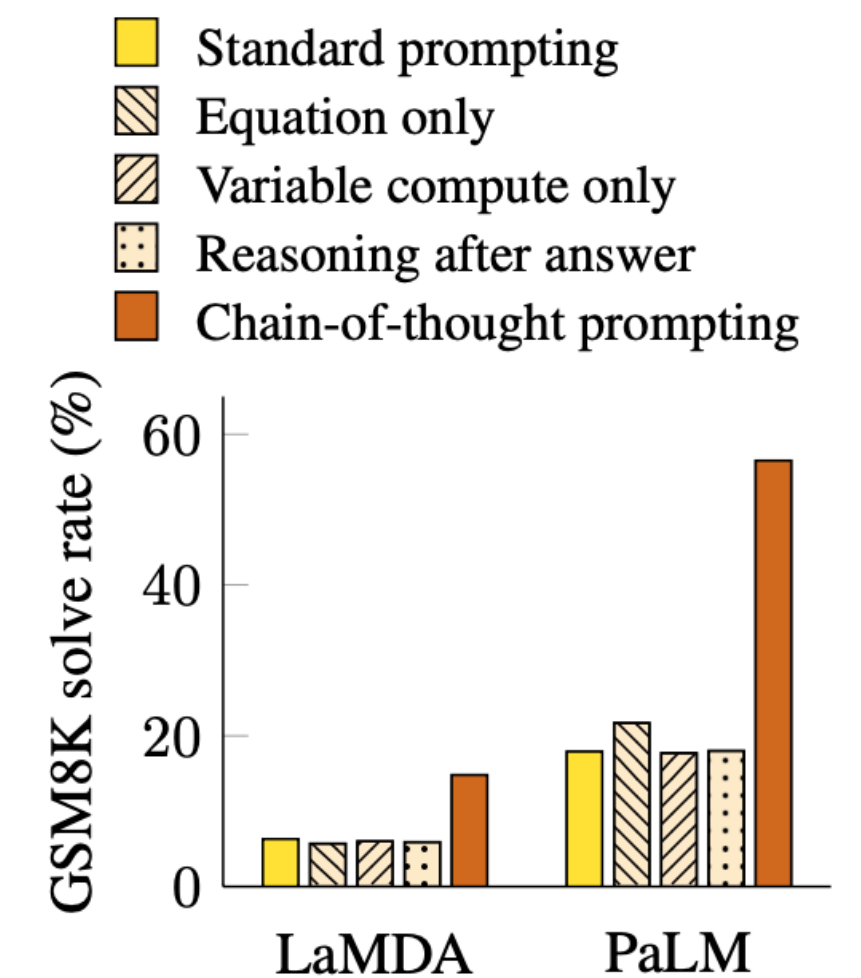
## Symbolic reasoning



## Commonsense reasoning



## Math

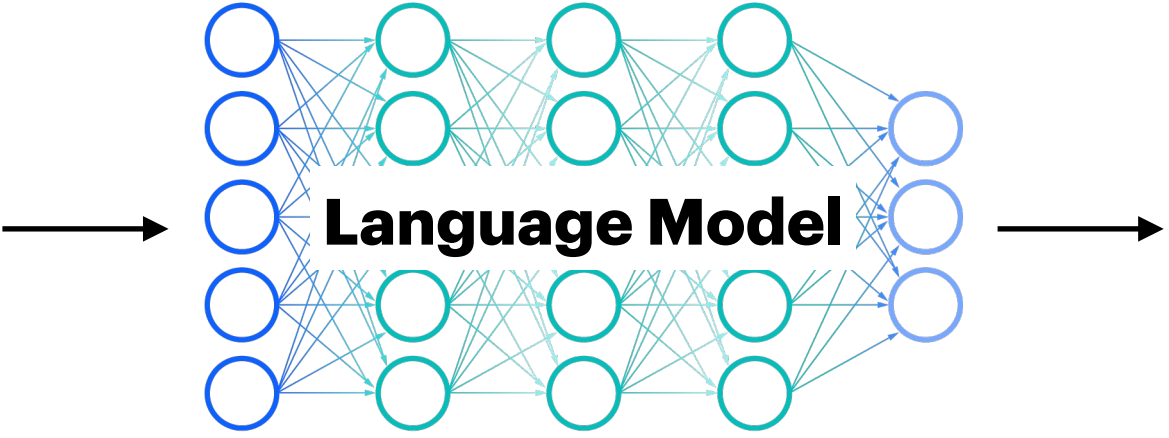




# Zero-shot chain-of-thought prompting

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: Let's think step by step.



There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls.



	MultiArith	GSM8K
<b>Zero-Shot</b>	<b>17.7</b>	<b>10.4</b>
Few-Shot (2 samples)	33.7	15.6
Few-Shot (8 samples)	33.8	15.6
<b>Zero-Shot-CoT</b>	<b>78.7</b>	<b>40.7</b>
Few-Shot-CoT (2 samples)	84.8	41.3
Few-Shot-CoT (4 samples : First) (*1)	89.2	-
Few-Shot-CoT (4 samples : Second) (*1)	90.5	-
Few-Shot-CoT (8 samples)	93.0	48.7

Significantly beats zero-shot ←

Manual CoT is still better ←

---

## Chain-of-Thought Prompting Elicits Reasoning in Large Language Models

---

---

## Large Language Models are Zero-Shot Reasoners

---

---

## Tree of Thoughts: Deliberate Problem Solving with Large Language Models

---

---

## LEAST-TO-MOST PROMPTING ENABLES COMPLEX REASONING IN LARGE LANGUAGE MODELS

Denny Zhou<sup>†\*</sup> Na  
Dale Schuurmans<sup>†</sup>  
<sup>†</sup>Google Research, E

## SELF-CONSISTENCY IMPROVES CHAIN OF THOUGHT REASONING IN LANGUAGE MODELS

Xuezhi Wang<sup>†‡</sup> Jason Wei<sup>†</sup> Dale Schuurmans<sup>†</sup> Quoc Le<sup>†</sup> Ed H. Chi<sup>†</sup>  
Sharan Narang<sup>†</sup> Aakanksha Chowdhery<sup>†</sup> Denny Zhou<sup>†§</sup>  
<sup>†</sup>Google Research, Brain Team  
<sup>‡</sup>xuezhiw@google.com, <sup>§</sup>dennyzhou@google.com

Forbes

FORBES > INNOVATION > ENTERPRISE TECH

# The Hot New Job That Pays Six Figures: AI Prompt Engineering

Bloomberg

## AI's Hottest Job: Prompt Engineer

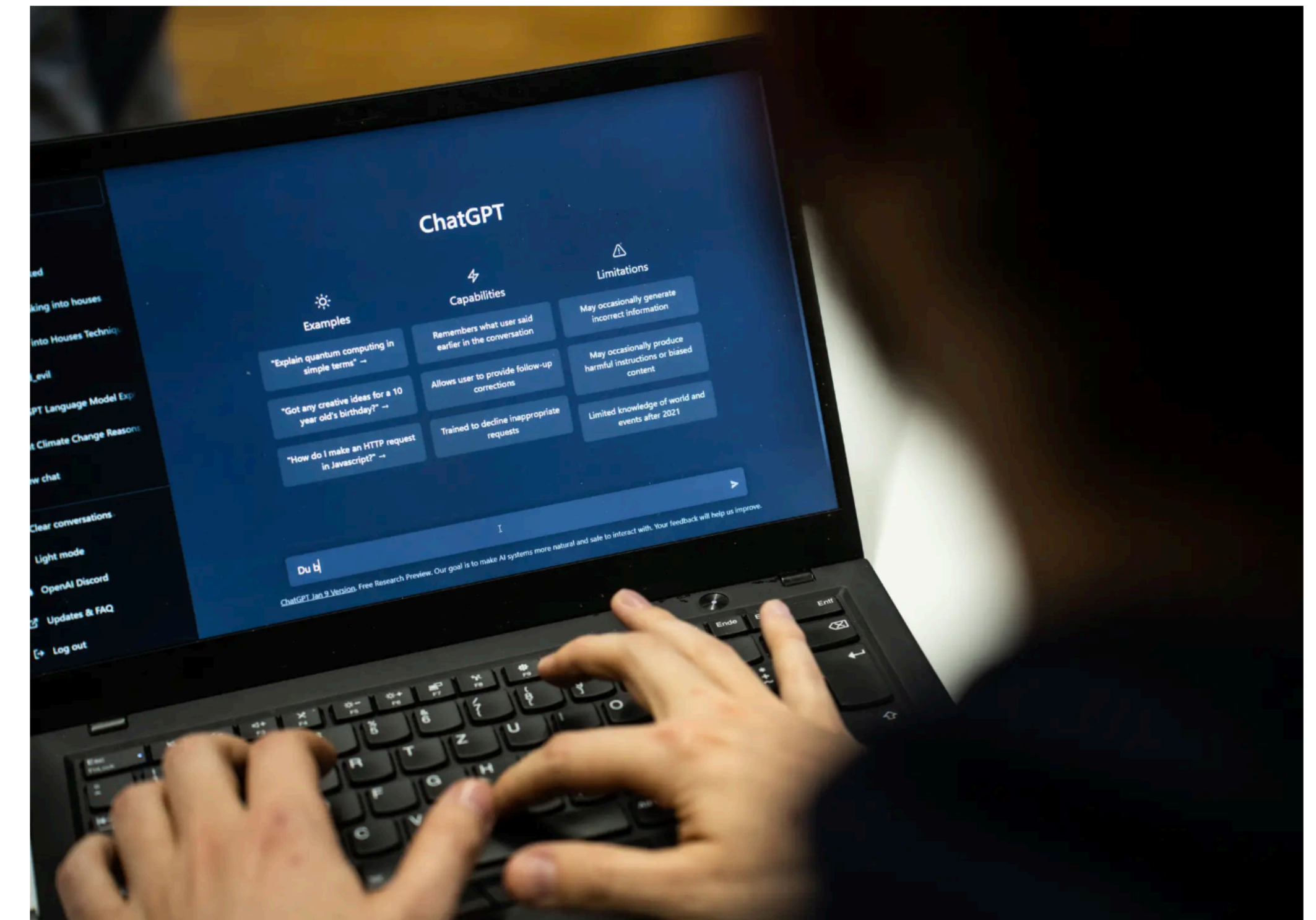
So-called AI whisperers can earn six-figure salaries, no programming experience necessary. Conrad Quilty-Harper met one of these prompt engineers to find out how to coax the best out of a large-language model. (Source: Bloomberg)

July 5th, 2023, 11:27 PM GMT+0800

## AI 'prompt engineer' jobs can pay up to \$375,000 a year and don't always require a background in tech

Britney Nguyen May 1, 2023, 11:34 PM GMT+8

Share Save



The rise of generative AI tools like ChatGPT is creating a hot market for "prompt engineers" who test and improve chatbot answers.

Getty Images

BUSINESS INSIDER



# Convergence of Vision and Language Models



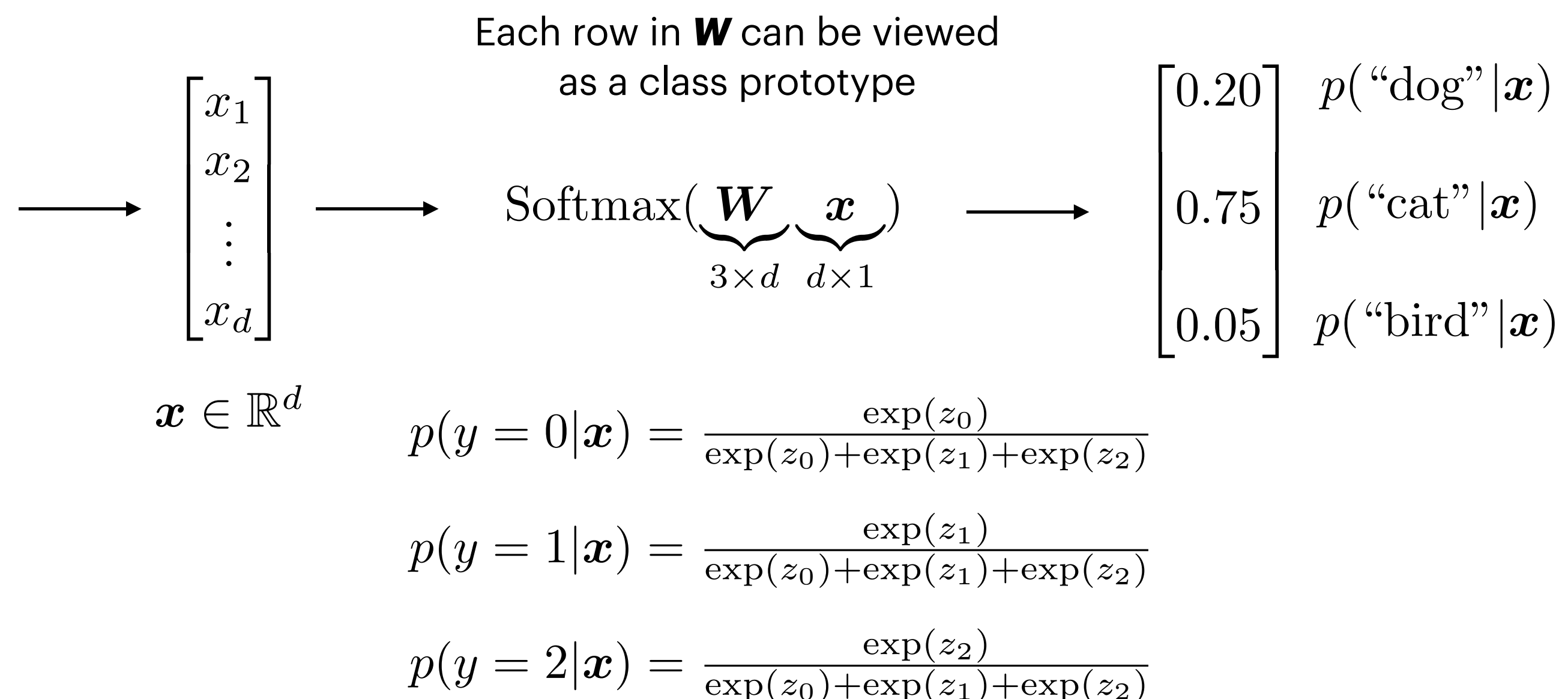
# Traditional vision models struggle to generalize



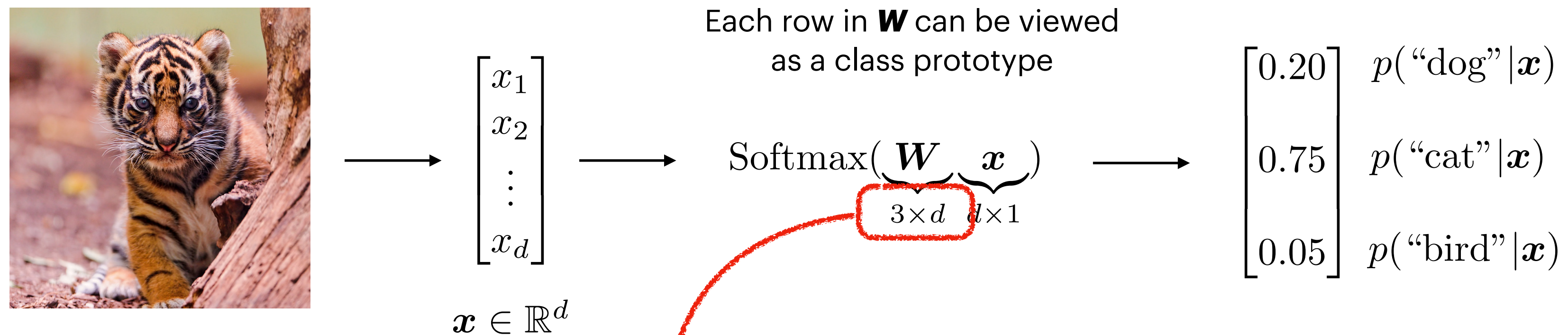
A classifier trained to recognise horse images would not be able to recognise zebra, though the latter is just like horse but with black-and-white stripes



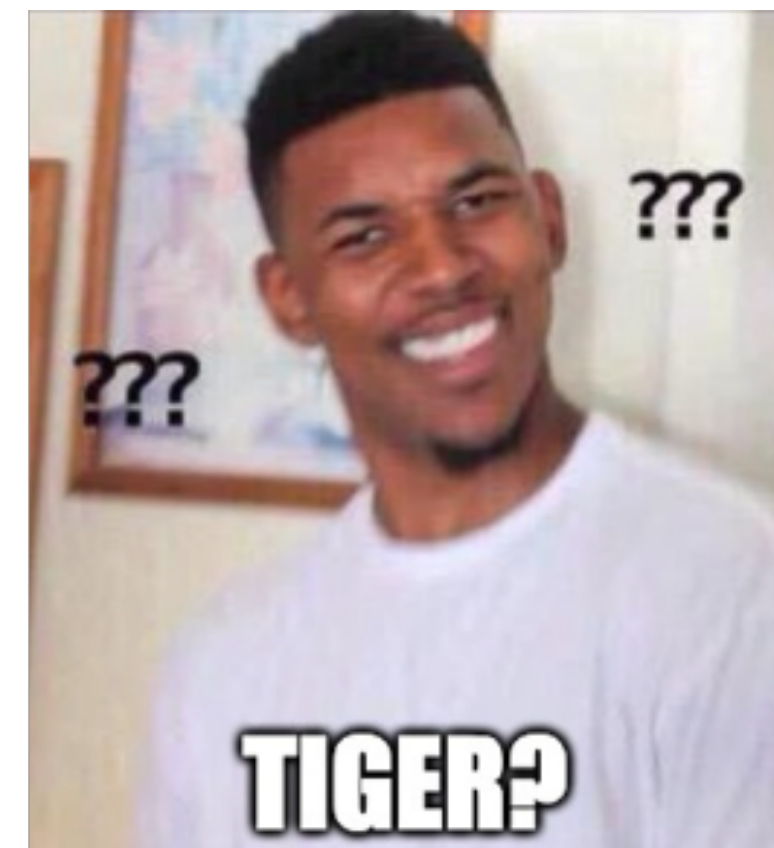
# Why traditional vision models struggle to generalize?



# Why traditional models struggle to generalize?



There is no class prototype for *tiger*  
Need to re-train the model!



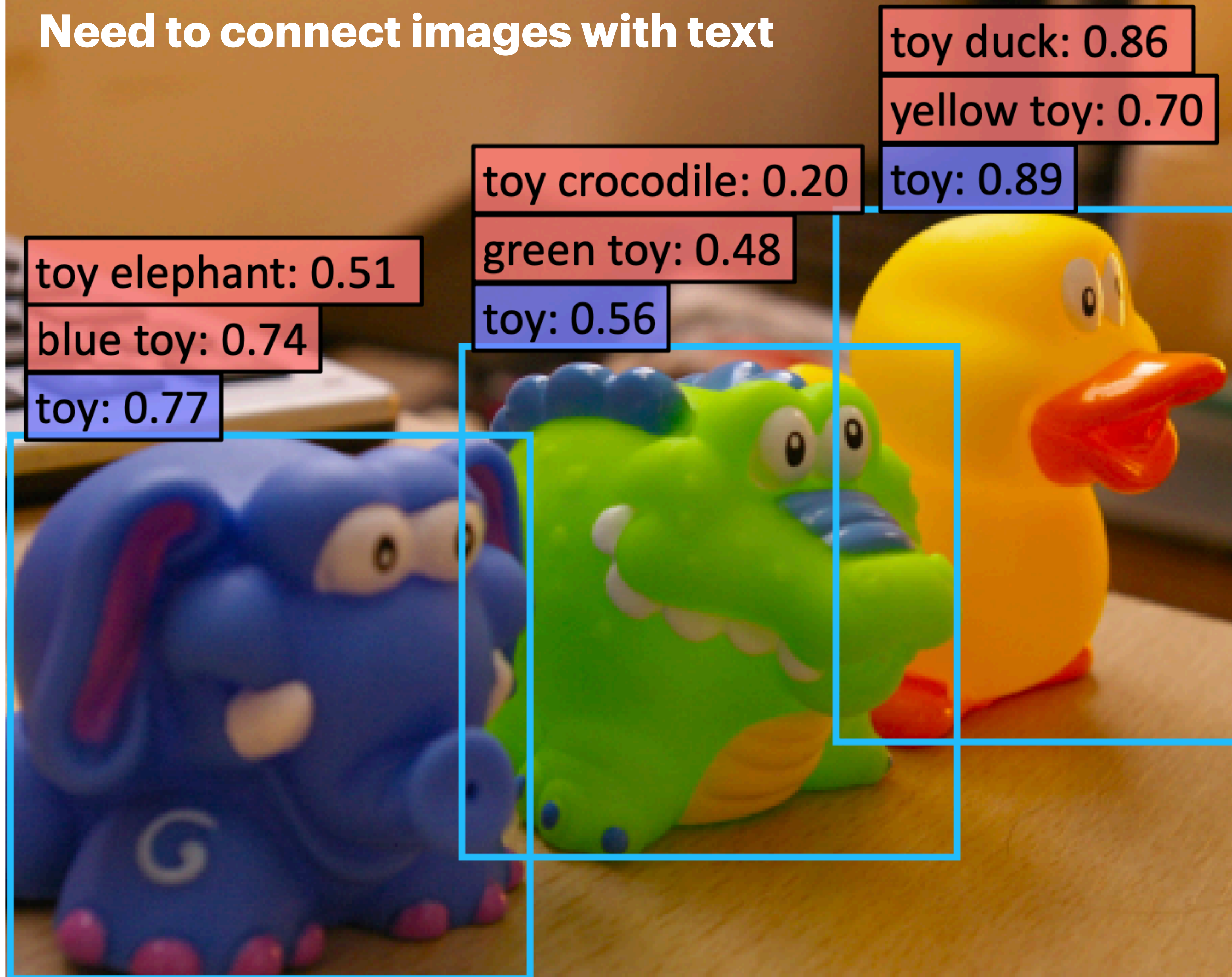


**Traditional models cannot  
handle open-vocabulary tasks**



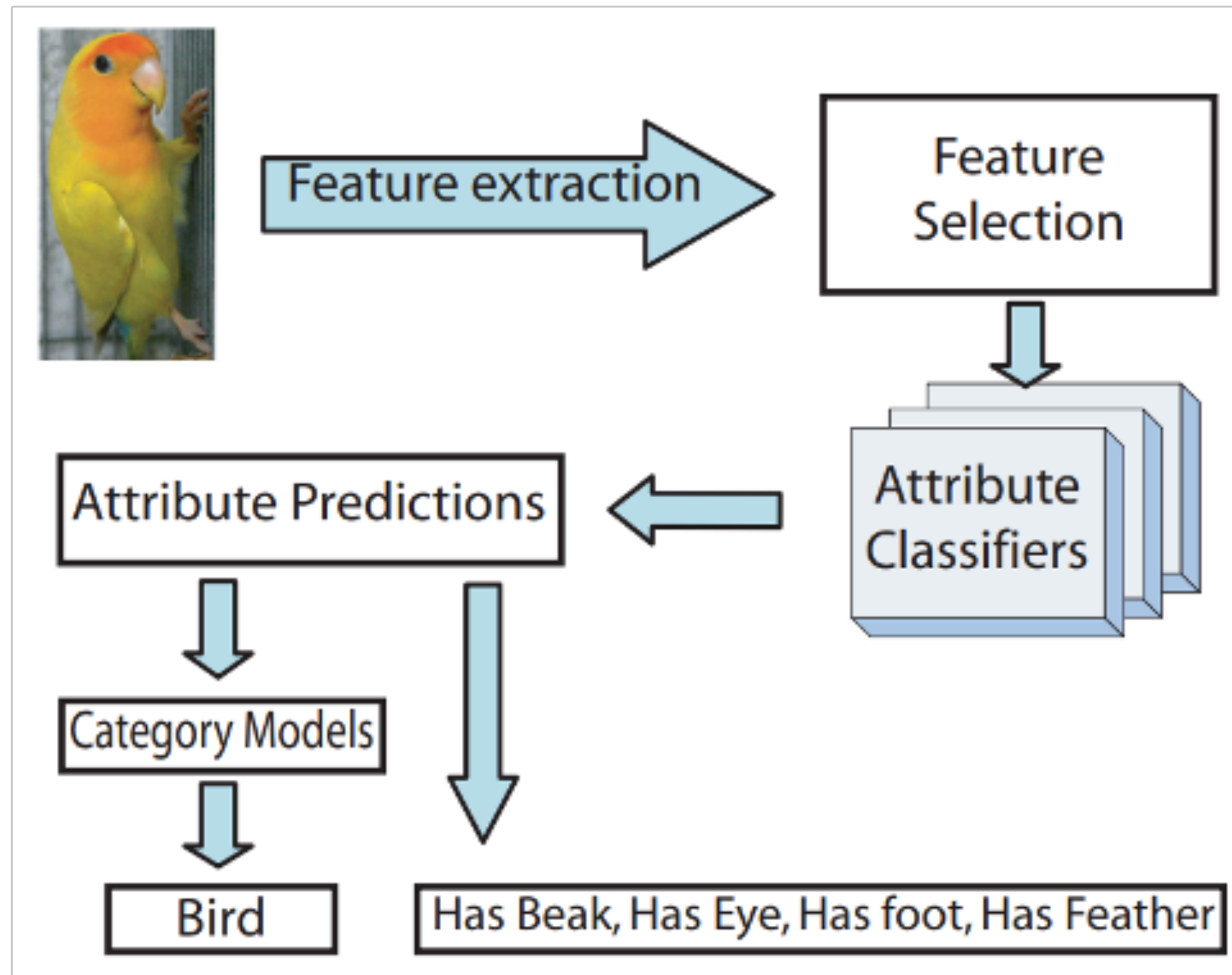


# Need to connect images with text



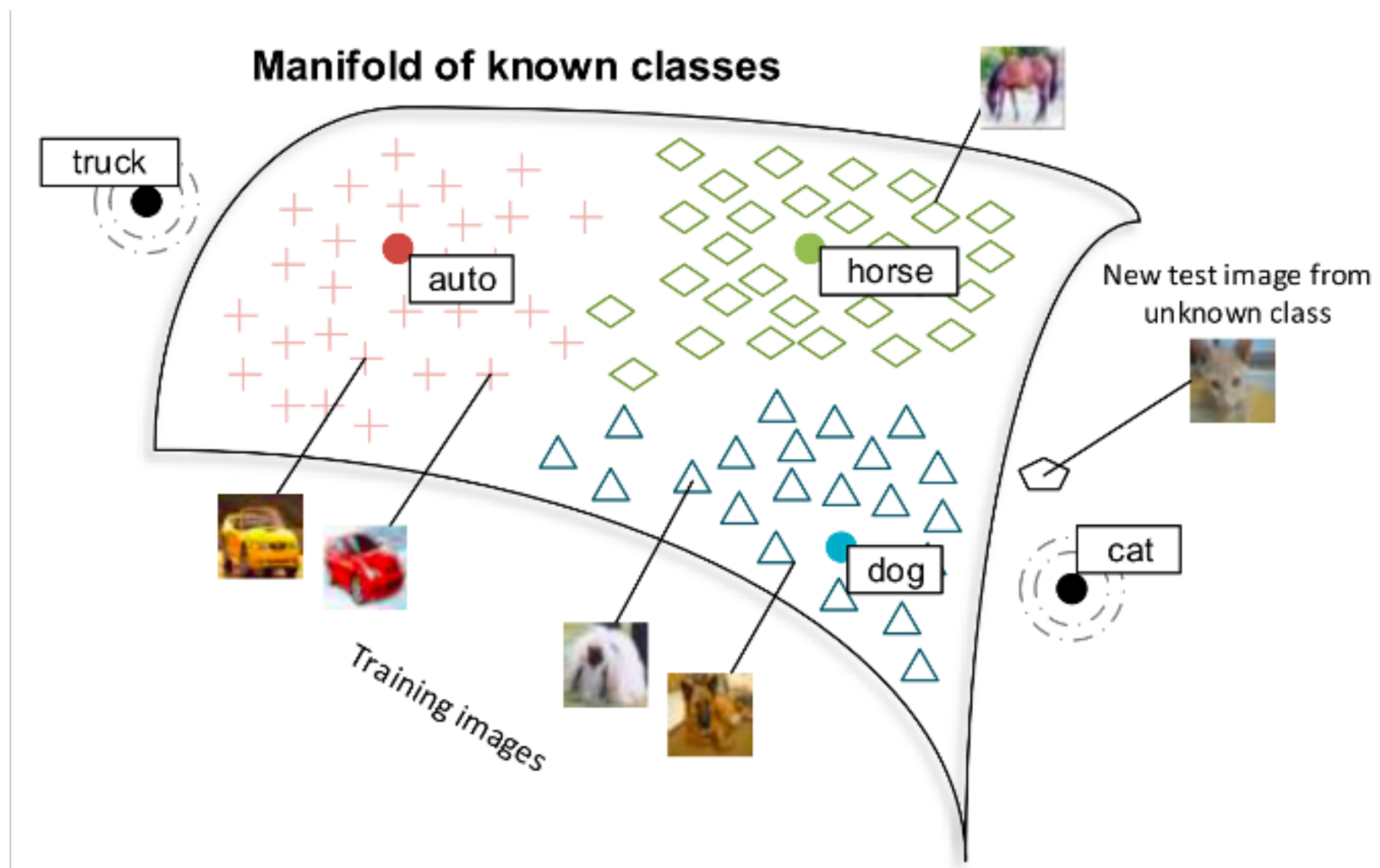


# Early methods



Associate classes with auxiliary information like **attributes**, which encode distinguishing properties of objects

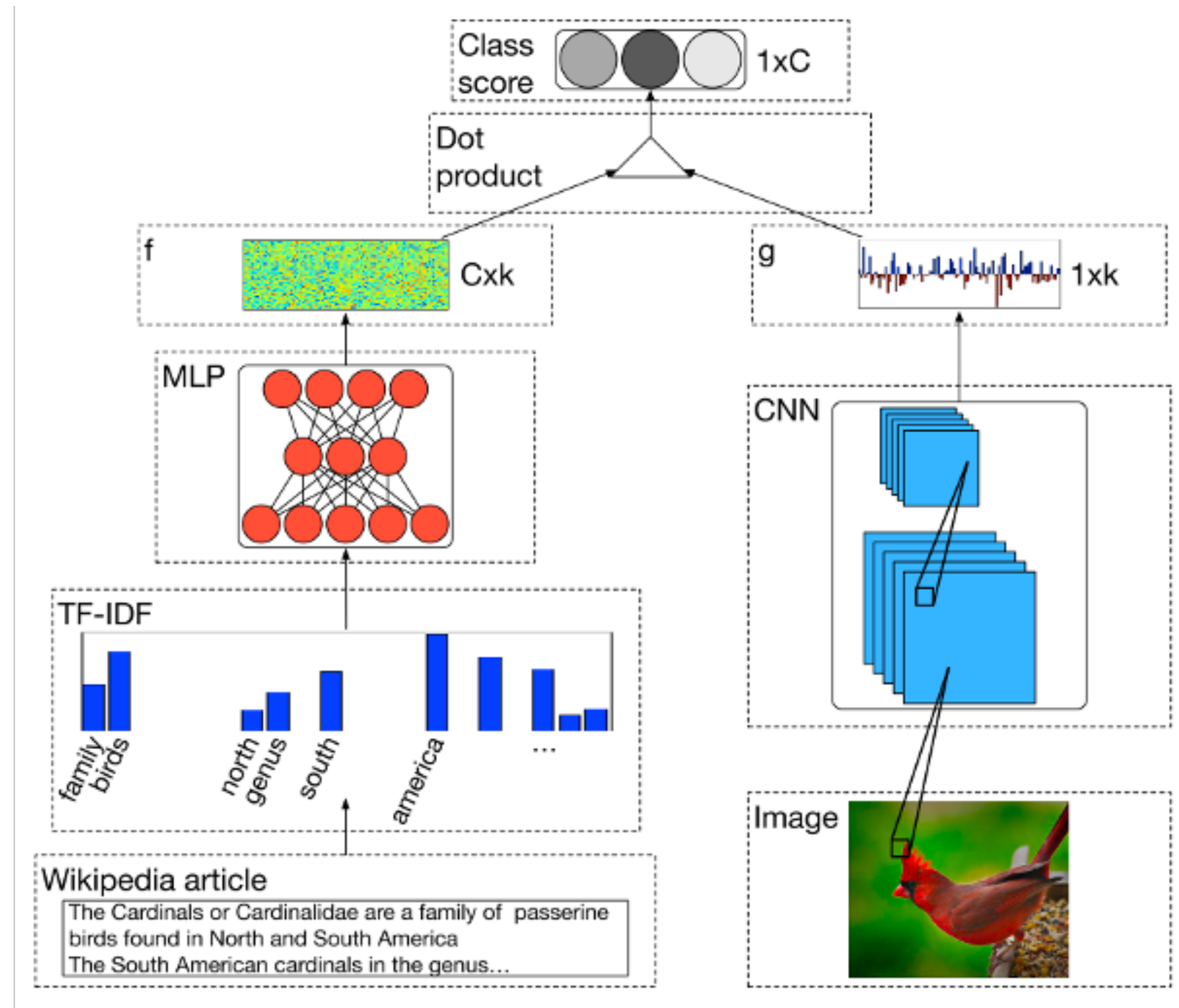
# Early methods



Associate images with semantic **word vectors** (i.e., word2vec)

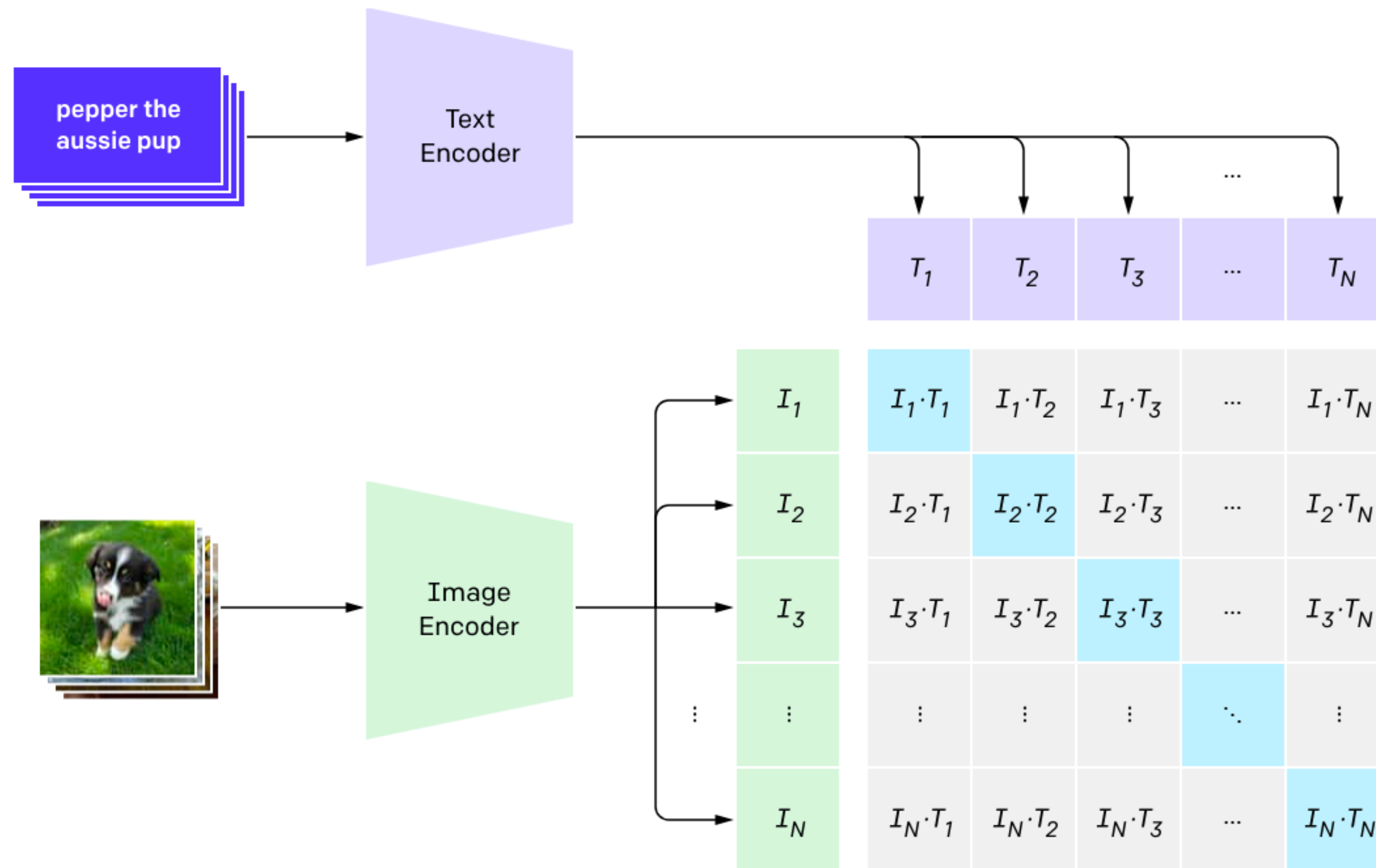


# Early methods



Learn a **joint embedding space** for images and text

# Today's methods



Learn a **joint embedding space** for images and text, **using large-capacity models and web-scale datasets**



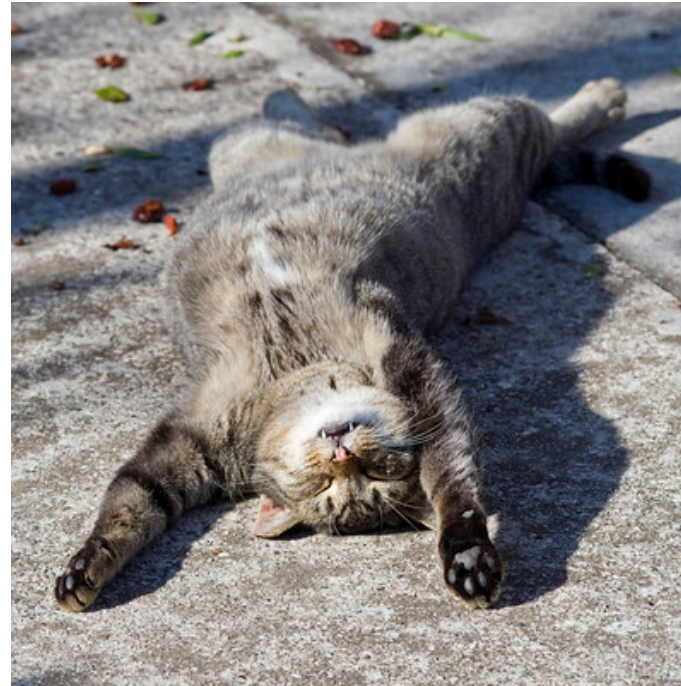
# Outline

- History
- Pre-training
- Prompting
- Applications

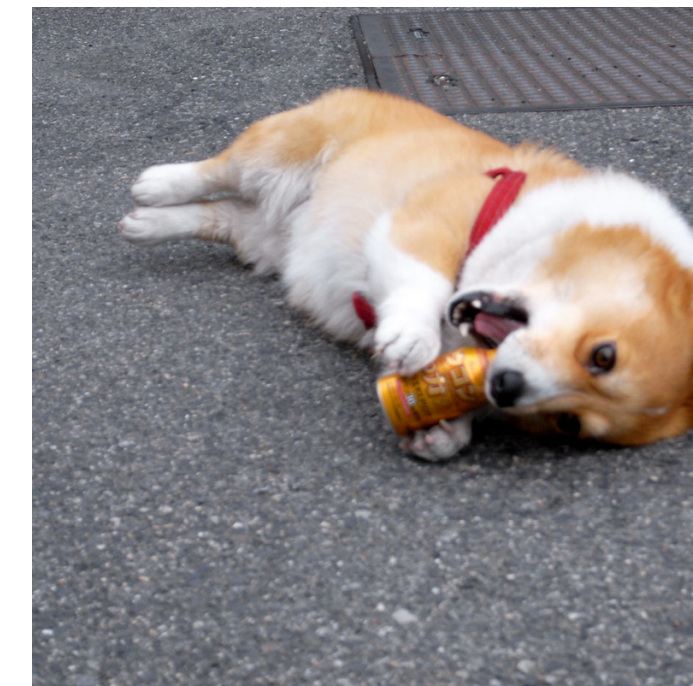
# Key idea: joint embedding space learning



cat sitting on the floor



cat laying down on the ground



a dog laying down with a  
bottle in mouth



# Contrastive learning

The goal is to associate each image with the correct label



Pig



Tiger



Panda



Hippo



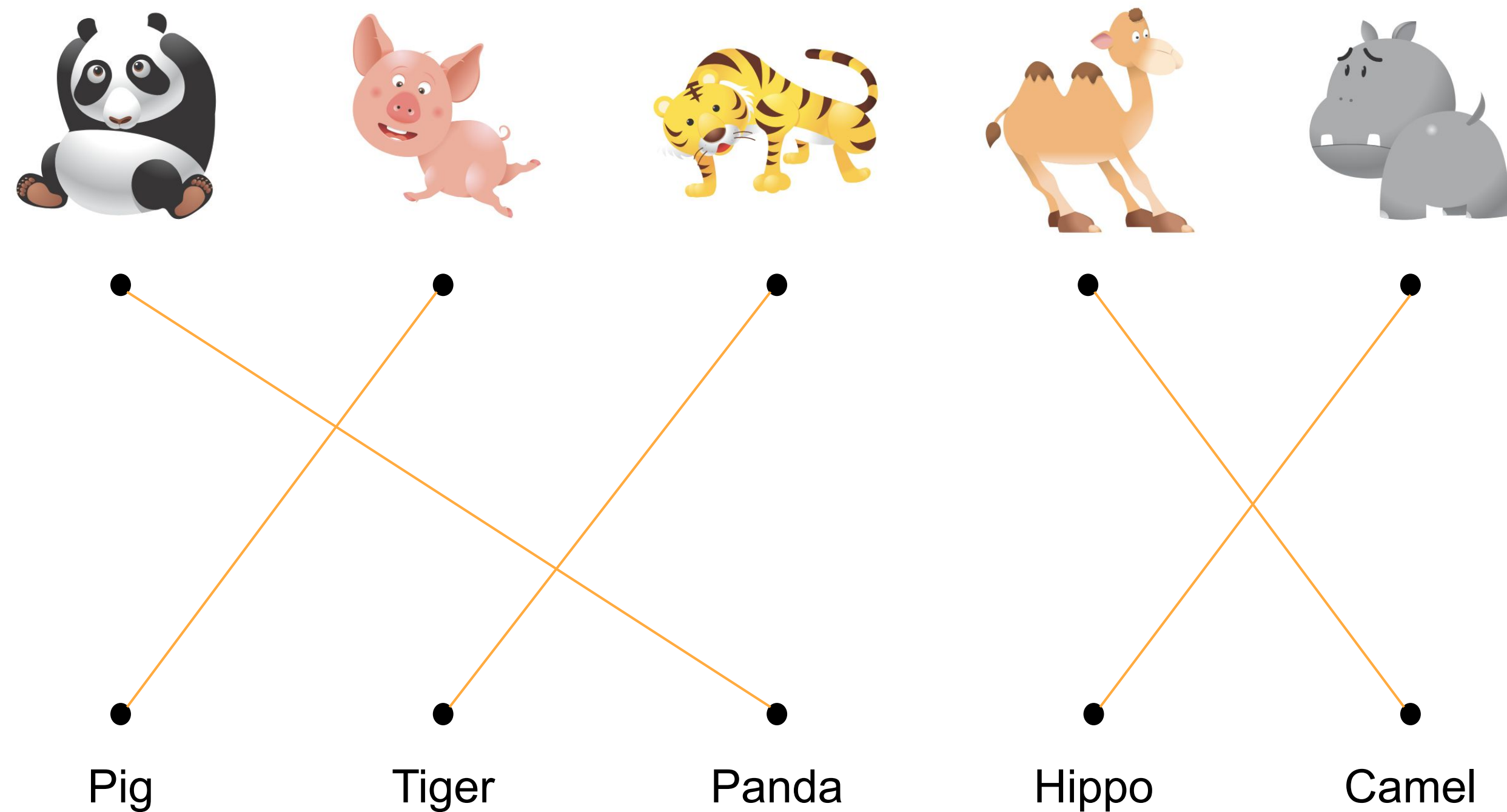
Camel

# Contrastive learning

Pull together matched pairs while push away unmatched pairs

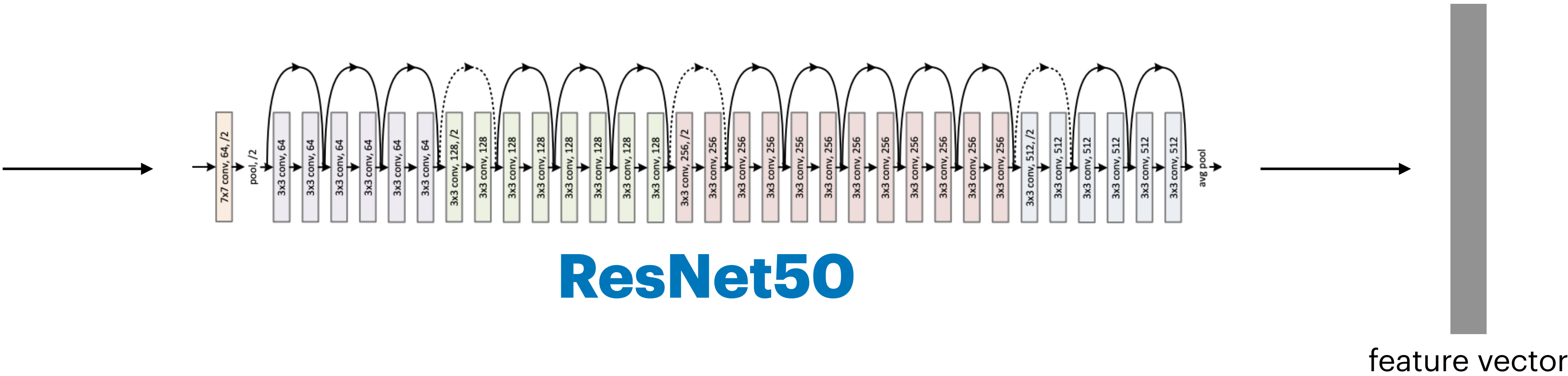
Reduce pair-wise feature distance  
(equivalent to increasing feature similarity)

Increase pair-wise feature distance  
(equivalent to decreasing feature similarity)

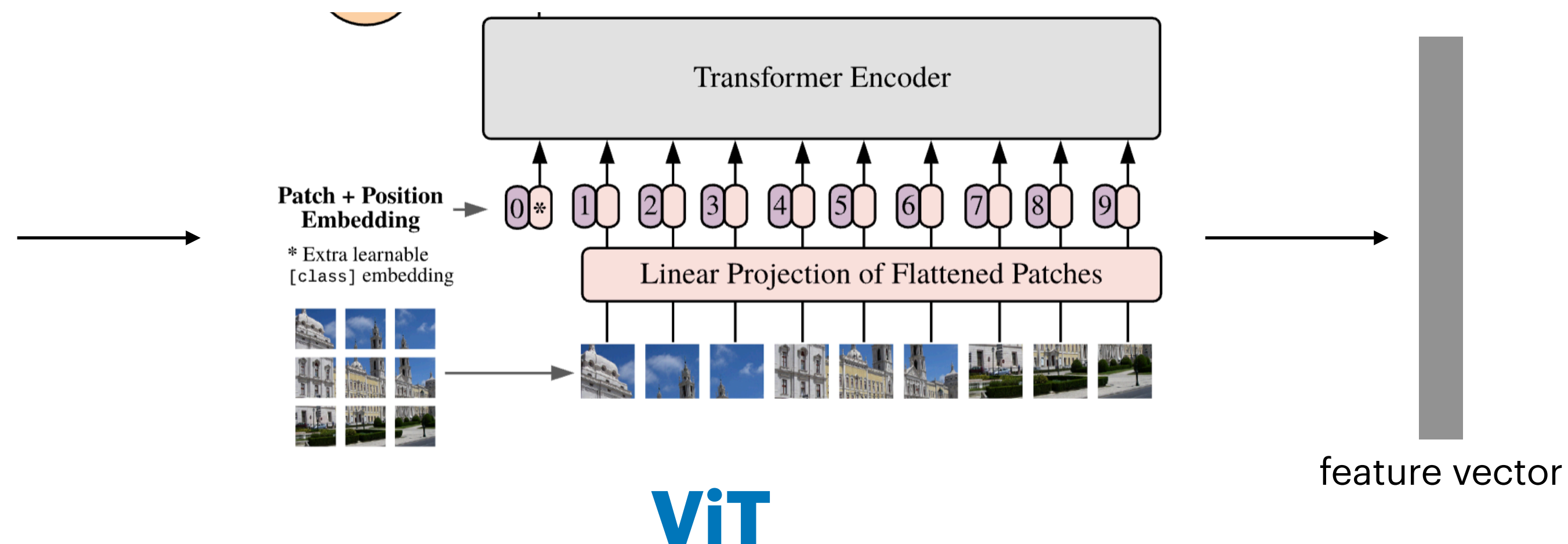




# Architecture: image encoder

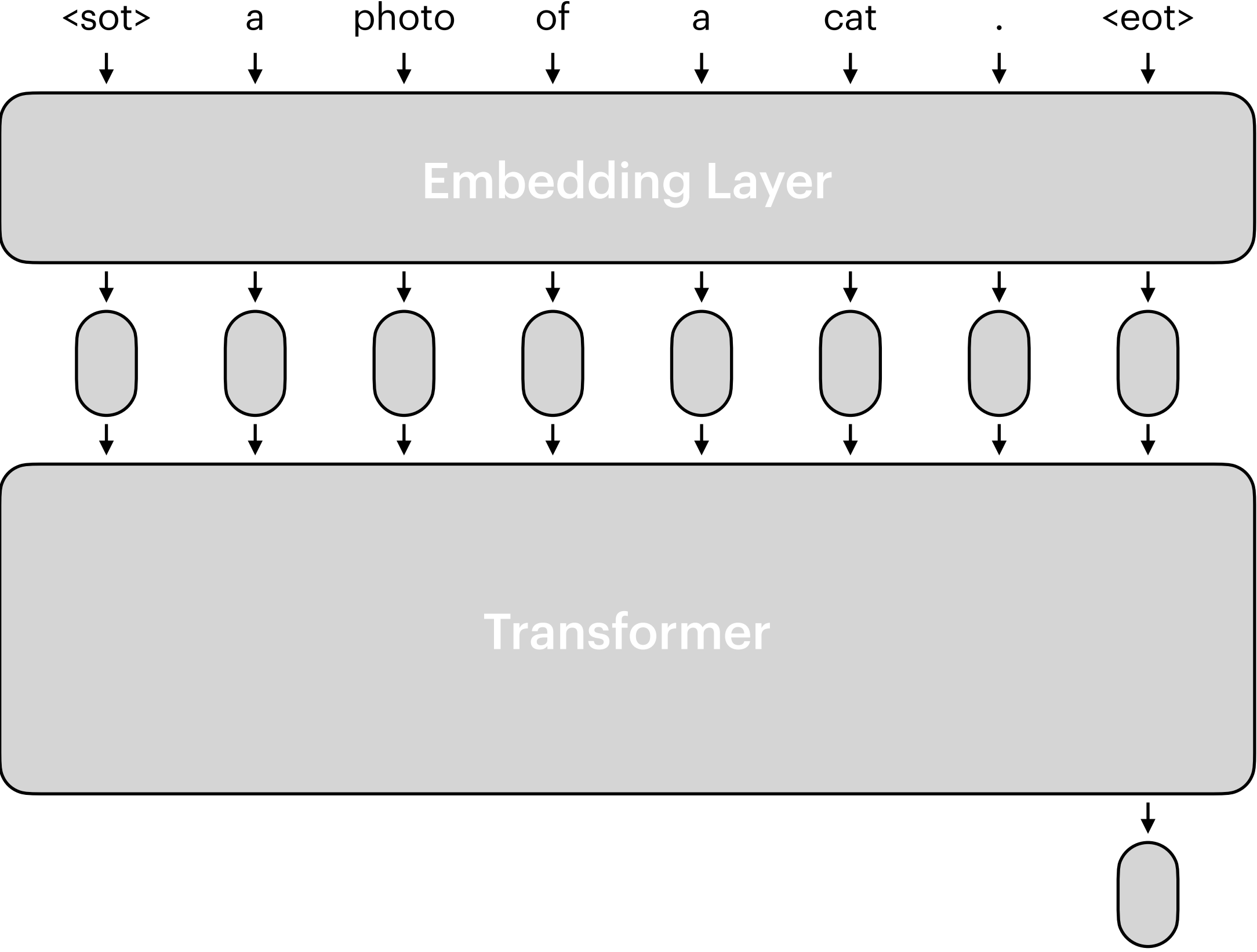


# Architecture: image encoder





# Architecture: text encoder



**Vocabulary**

word	index	word embedding
a	0	
and	1	
of	2	
⋮	⋮	⋮
<eot>	49,151	

# Data: LAION-5B

Backend url:

<https://knn5.laion>

Index:

laion\_5B

french cat



[Clip retrieval](#) works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embedddings

Display captions ☒

Display full

captions ☐

Display similarities

☐

Safe mode ☒

Hide duplicate urls

☒

Hide (near)

duplicate images ☒

Search over

image

Search with multilingual clip

☐



french cat



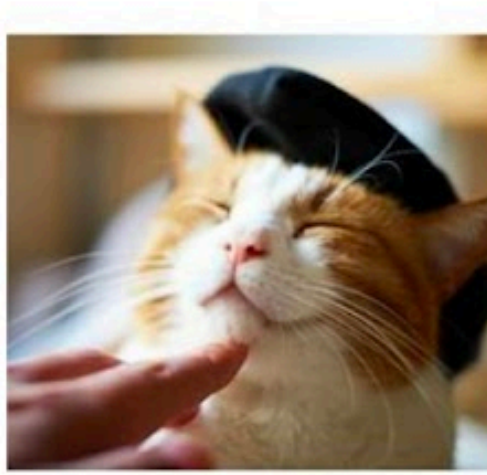
french cat



How to tell if your feline is french. He wears a b...



イケメン猫モデル「トキ・ナンタケット」がカッコいい - NAVER まとめ



Hilarious pics of funny cats! funnycatsgif.com



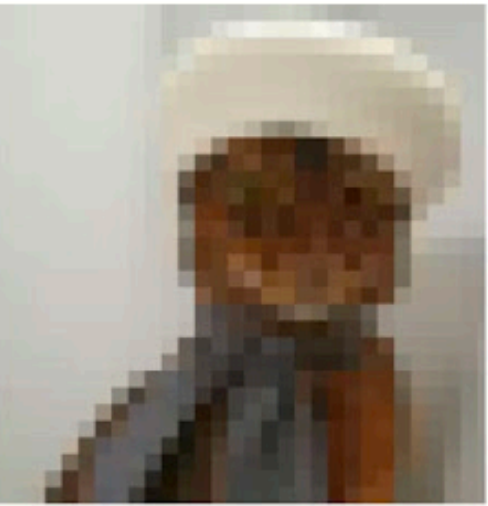
Hipster cat



網友挑戰「加幾筆畫出最創意貓咪圖片」，笑到岔氣之後我也手...



cat in a suit Georgian sells tomatoes

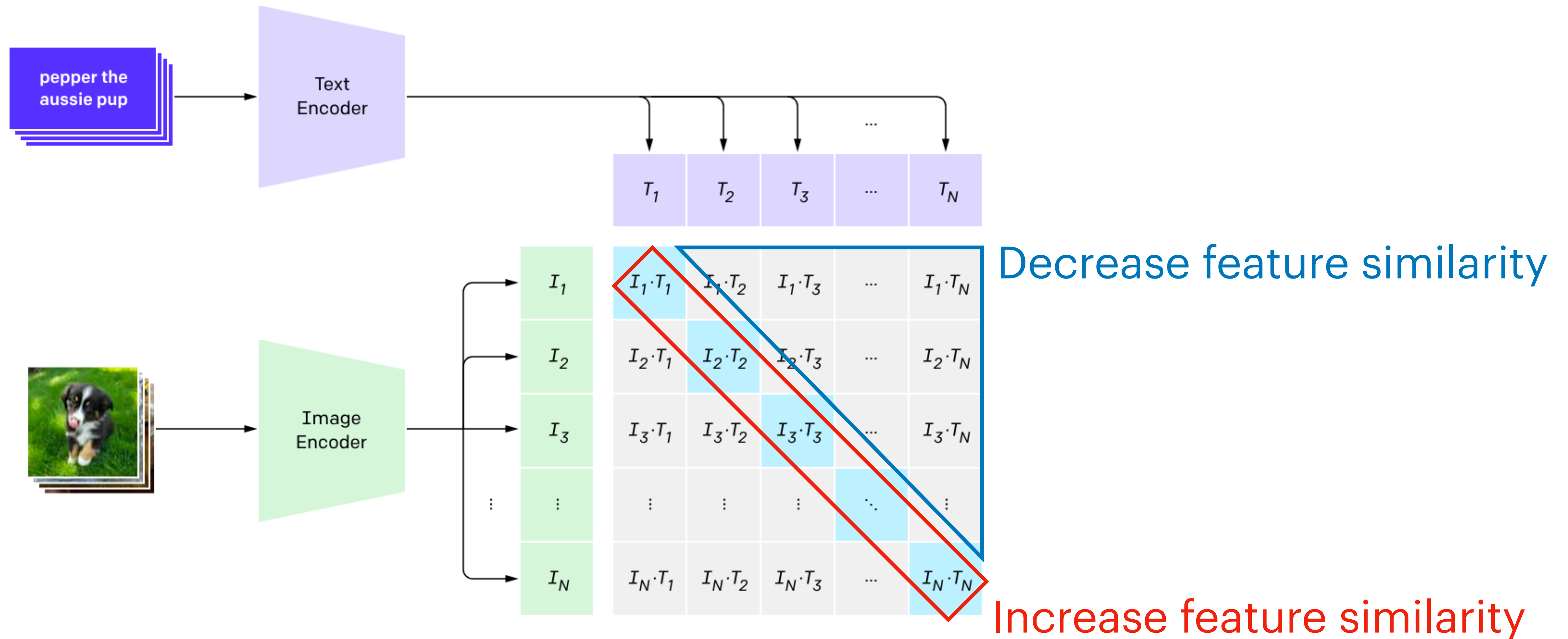


French Bread Cat Loaf Metal Print



# Contrastive Language-Image Pre-training (CLIP)

## 1. Contrastive pre-training



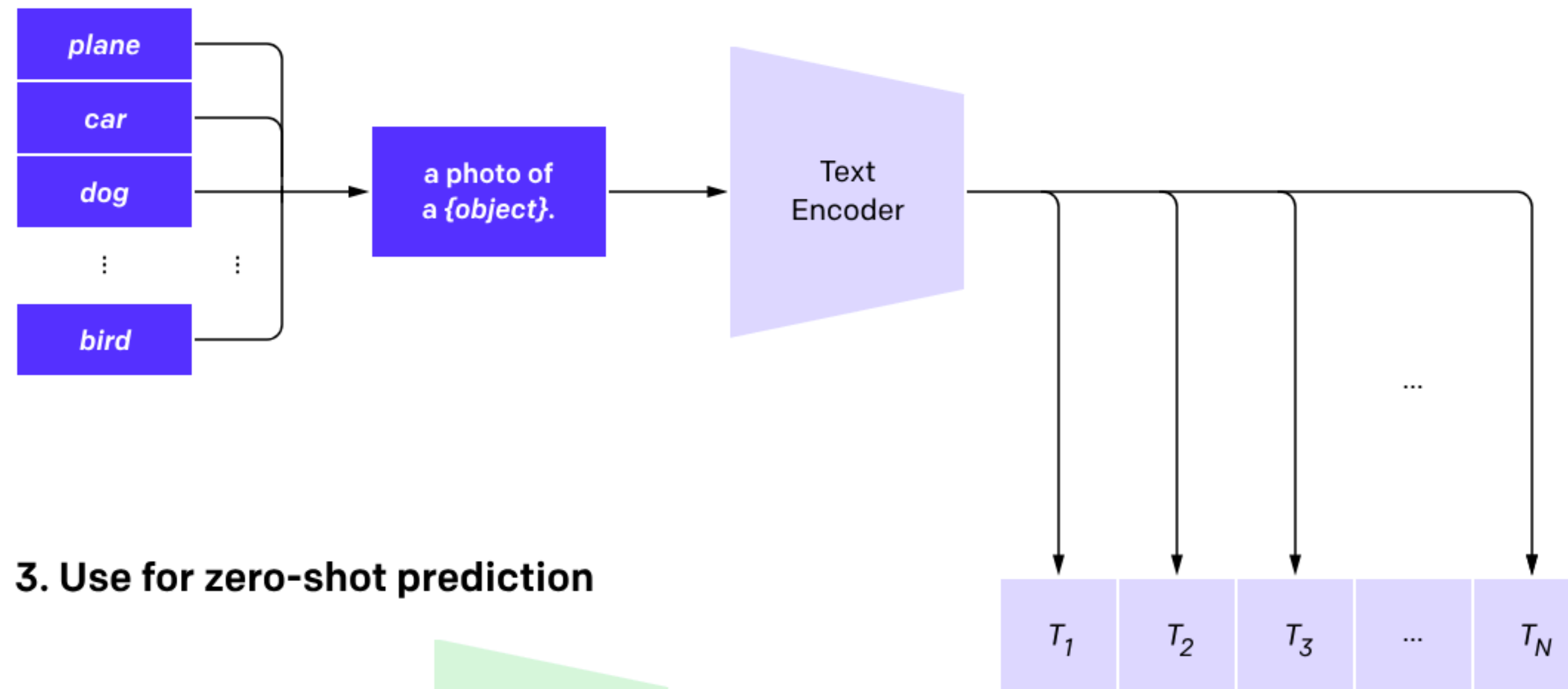
# Outline

- History
- Pre-training
- Prompting
- Applications

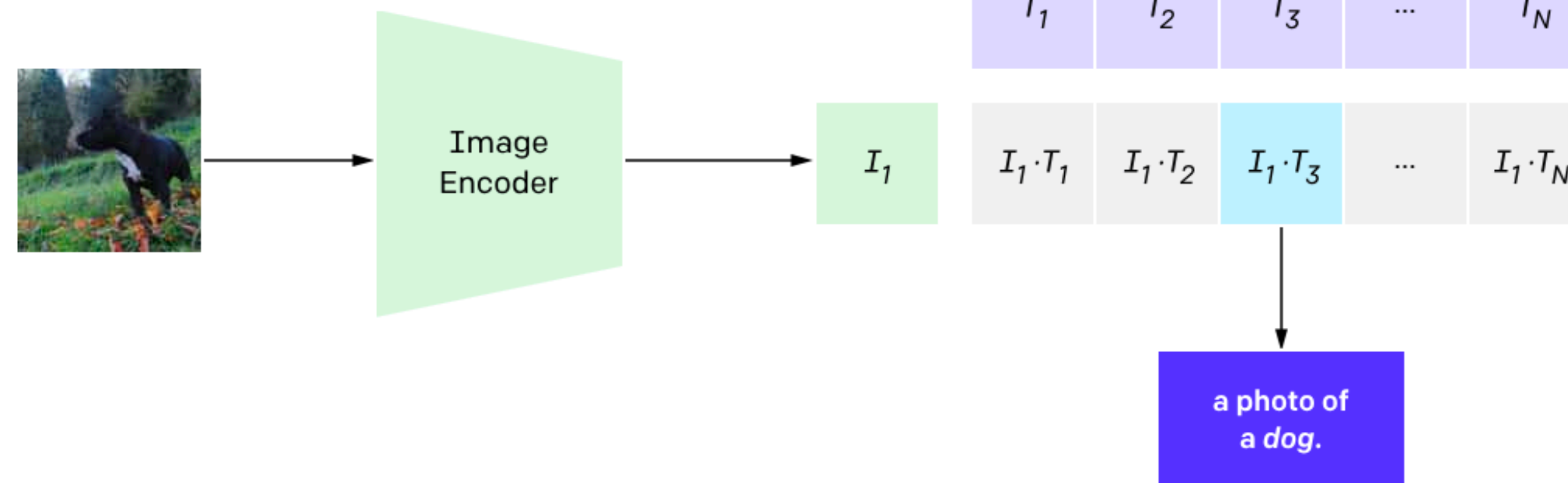


# Zero-shot prompting

## 2. Create dataset classifier from label text



## 3. Use for zero-shot prediction



FOOD101

guacamole (90.1%) Ranked 1 out of 101 labels



✓ a photo of **guacamole**, a type of food.

✗ a photo of **ceviche**, a type of food.

✗ a photo of **edamame**, a type of food.

✗ a photo of **tuna tartare**, a type of food.

✗ a photo of **hummus**, a type of food.

YOUTUBE-BB

airplane, person (89.0%) Ranked 1 out of 23



✓ a photo of a **airplane**.

✗ a photo of a **bird**.

✗ a photo of a **bear**.

✗ a photo of a **giraffe**.

✗ a photo of a **car**.

SUN397

television studio (90.2%) Ranked 1 out of 397



✓ a photo of a **television studio**.

✗ a photo of a **podium indoor**.

✗ a photo of a **conference room**.

✗ a photo of a **lecture room**.

✗ a photo of a **control room**.

EUROSAT

annual crop land (12.9%) Ranked 4 out of 10



✗ a centered satellite photo of **permanent crop land**.

✗ a centered satellite photo of **pasture land**.

✗ a centered satellite photo of **highway or road**.

✓ a centered satellite photo of **annual crop land**.

✗ a centered satellite photo of **brushland or shrubland**.



CIFAR-10

bird (40.9%) Ranked 1 out of 10 labels



- ✓ a photo of a **bird**.
- ✗ a photo of a **cat**.
- ✗ a photo of a **deer**.
- ✗ a photo of a **frog**.
- ✗ a photo of a **dog**.

FACIAL EMOTION RECOGNITION 2013 (FER2013)

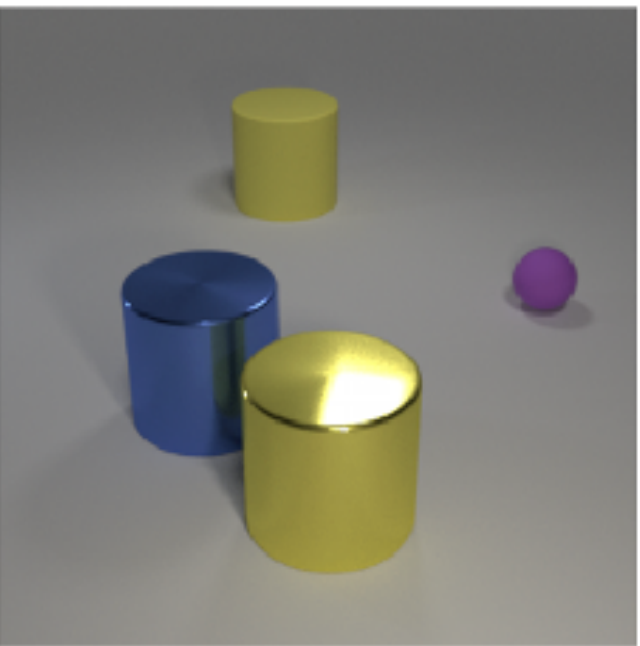
angry (8.2%) Ranked 5 out of 7



- ✗ a photo of a **happy** looking face.
- ✗ a photo of a **neutral** looking face.
- ✗ a photo of a **surprised** looking face.
- ✗ a photo of a **fearful** looking face.
- ✓ a photo of a **angry** looking face.

CLEVR COUNT

4 (17.1%) Ranked 2 out of 8



- ✗ a photo of **3** objects.
- ✓ a photo of **4** objects.
- ✗ a photo of **5** objects.
- ✗ a photo of **6** objects.
- ✗ a photo of **10** objects.

UCF101

Volleyball Spiking (99.3%) Ranked 1 out of 101



- ✓ a photo of a person **volleyball spiking**.
- ✗ a photo of a person **jump rope**.
- ✗ a photo of a person **long jump**.
- ✗ a photo of a person **soccer penalty**.
- ✗ a photo of a person **table tennis shot**.

STANFORD CARS

2012 Honda Accord Coupe (63.3%)    Ranked 1 out of 196



✓ a photo of a **2012 honda accord coupe**.

✗ a photo of a **2012 honda accord sedan**.

✗ a photo of a **2012 acura tl sedan**.

✗ a photo of a **2012 acura tsx sedan**.

✗ a photo of a **2008 acura tl type-s**.

KINETICS-700

country line dancing (99.0%)    Ranked 1 out of 700



✓ a photo of **country line dancing**.

✗ a photo of **square dancing**.

✗ a photo of **swing dancing**.

✗ a photo of **dancing charleston**.

✗ a photo of **salsa dancing**.

SUN

kennel indoor (98.6%)    Ranked 1 out of 723



✓ a photo of a **kennel indoor**.

✗ a photo of a **kennel outdoor**.

✗ a photo of a **jail cell**.

✗ a photo of a **jail indoor**.

✗ a photo of a **veterinarians office**.

FLOWERS-102

great masterwort (74.3%)    Ranked 1 out of 102



✓ a photo of a **great masterwort**, a type of flower.








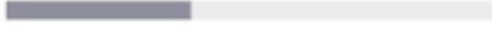


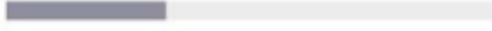


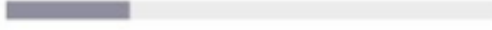


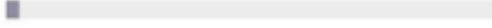

✗ a photo of a **bishop of llandaff**, a type of flower.

✗ a photo of a **pincushion flower**, a type of flower.

✗ a photo of a **globe flower**, a type of flower.

✗ a photo of a **prince of wales feathers**, a type of flower.

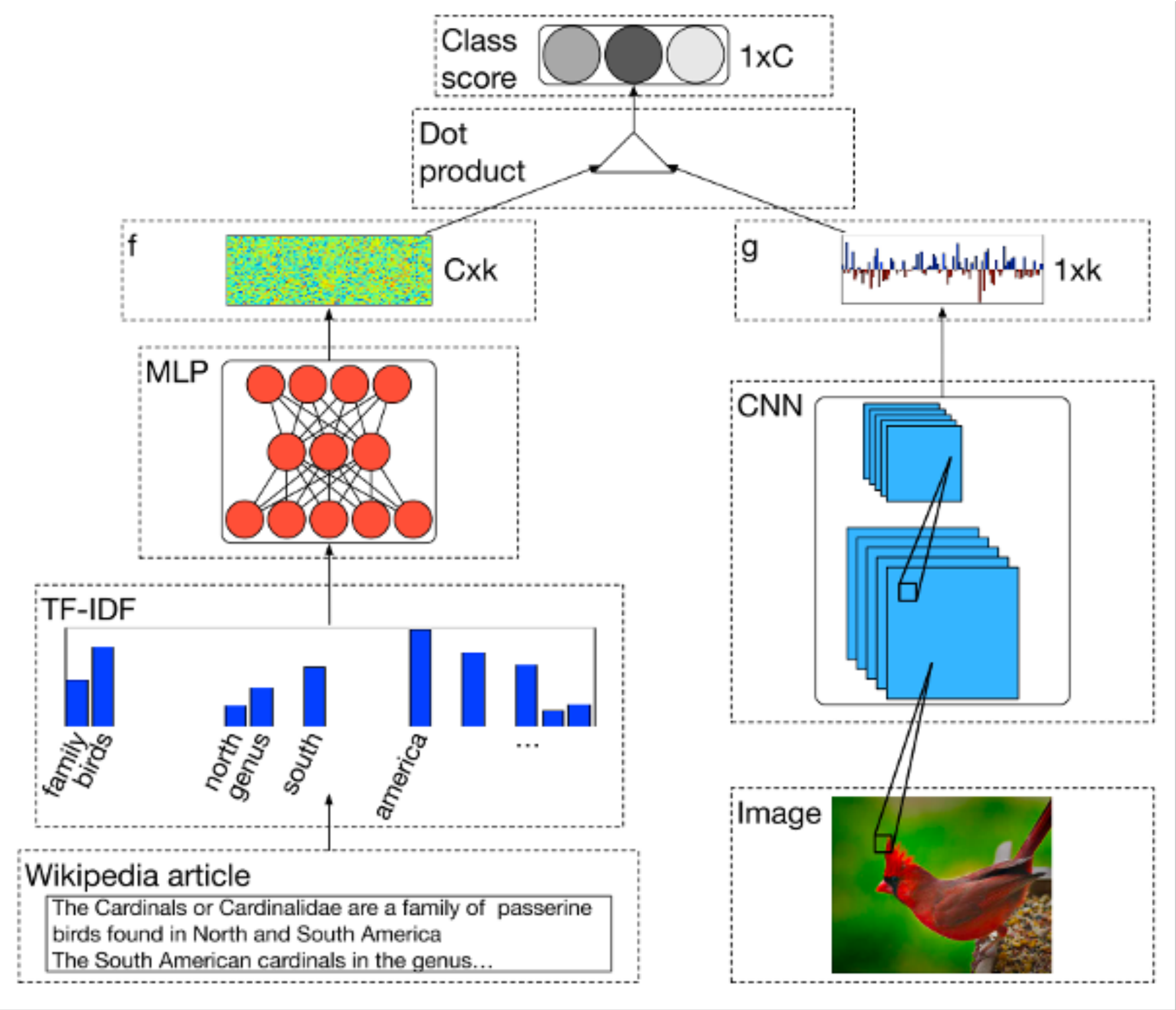


DATASET	IMAGENET RESNET101	CLIP VIT-L
 <p>ImageNet</p>	 <p>76.2%</p>	 <p>76.2%</p>
 <p>ImageNet V2</p>	 <p>64.3%</p>	 <p>70.1%</p>
 <p>ImageNet Rendition</p>	 <p>37.7%</p>	 <p>88.9%</p>
 <p>ObjectNet</p>	 <p>32.6%</p>	 <p>72.3%</p>
 <p>ImageNet Sketch</p>	 <p>25.2%</p>	 <p>60.2%</p>
 <p>ImageNet Adversarial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



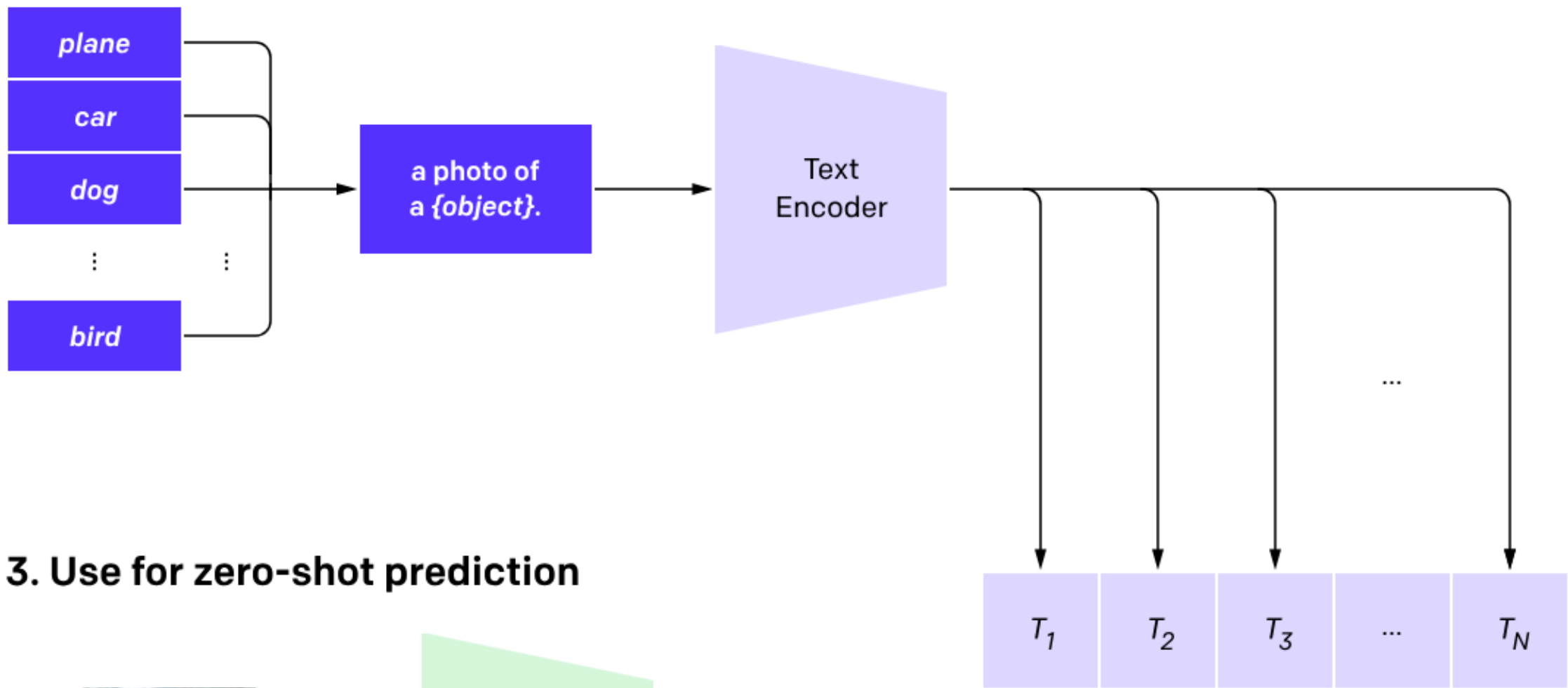
# Why CLIP works?

- Scaling (model & data)
- Transformer
- Contrastive learning

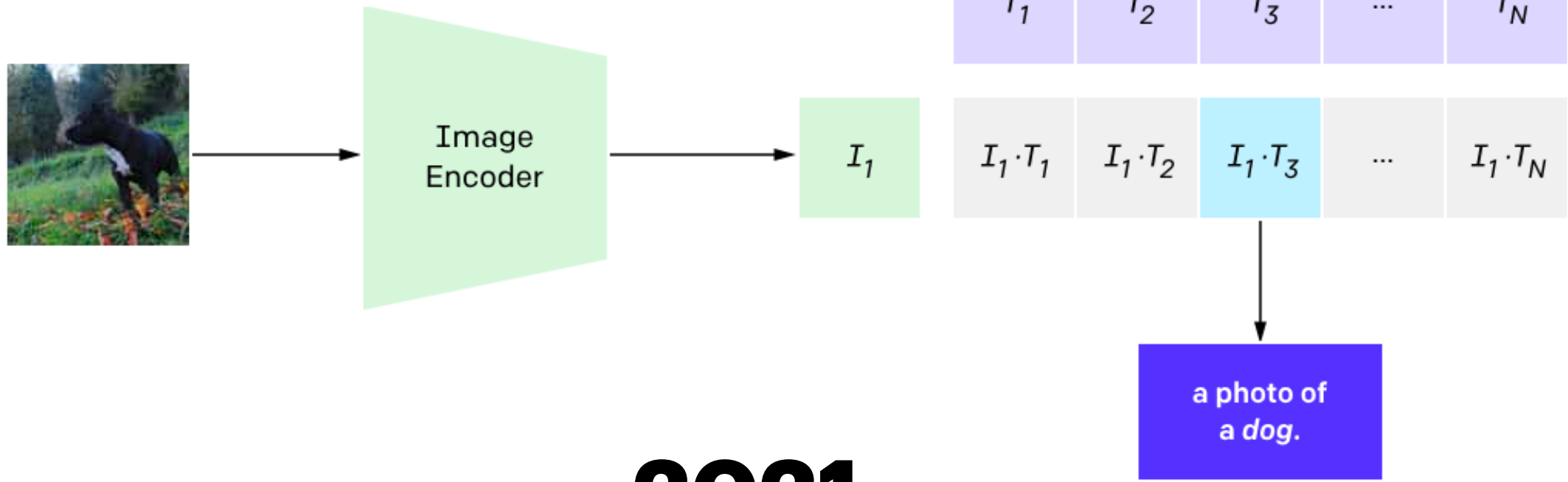


2015

## 2. Create dataset classifier from label text



## 3. Use for zero-shot prediction



2021



# Prompt engineering



Text prompts
A photo of a {dog}
A photo of a {cat}
A photo of a {bird}
...
A photo of a {tiger}


a bad photo of a {}.  
a photo of many {}.  
a sculpture of a {}.  
a photo of the hard to see {}.  
a low resolution photo of the {}.  
a rendering of a {}.  
graffiti of a {}.  
a bad photo of the {}.  
a cropped photo of the {}.  
a tattoo of a {}.  
the embroidered {}.  
a photo of a hard to see {}.  
a bright photo of a {}.  
a photo of a clean {}.  
a photo of a dirty {}.  
a dark photo of the {}.  
a drawing of a {}.  
a photo of my {}.  
the plastic {}.  
a photo of the cool {}.  
a close-up photo of a {}.  
a black and white photo of the {}.  
a painting of the {}.  
a painting of a {}.


a pixelated photo of the {}.  
a sculpture of the {}.  
a bright photo of the {}.  
a cropped photo of a {}.  
a plastic {}.  
a photo of the dirty {}.  
a jpeg corrupted photo of a {}.  
a blurry photo of the {}.  
a photo of the {}.  
a good photo of the {}.  
a rendering of the {}.  
a {} in a video game.  
a photo of one {}.  
a doodle of a {}.  
a close-up photo of the {}.  
a photo of a {}.  
the origami {}.  
the {} in a video game.  
a sketch of a {}.  
a doodle of the {}.  
a origami {}.  
a low resolution photo of a {}.  
the toy {}.  
a rendition of the {}.

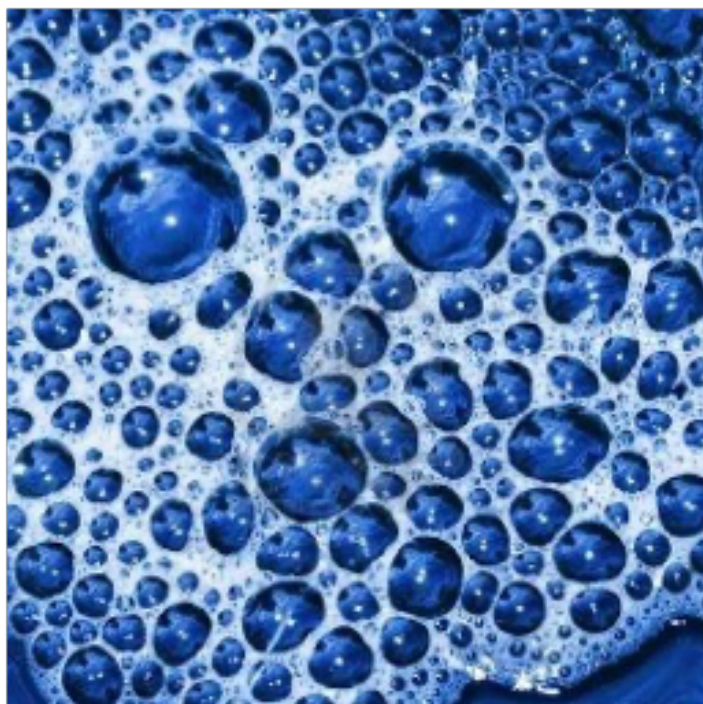
a photo of the clean {}.  
a photo of a large {}.  
a rendition of a {}.  
a photo of a nice {}.  
a photo of a weird {}.  
a blurry photo of a {}.  
a cartoon {}.  
art of a {}.  
a sketch of the {}.  
a embroidered {}.  
a pixelated photo of a {}.  
itap of the {}.  
a jpeg corrupted photo of the {}.  
a good photo of a {}.  
a plushie {}.  
a photo of the nice {}.  
a photo of the small {}.  
a photo of the weird {}.  
the cartoon {}.  
art of the {}.  
a drawing of the {}.  
a photo of the large {}.  
a black and white photo of a {}.  
the plushie {}.




# Prompt engineering is hard

Caltech101	Prompt	Accuracy
	a [CLASS].	82.68
	a photo of [CLASS].	80.81
	a photo of <b>a</b> [CLASS].	86.29
	[V] <sub>1</sub> [V] <sub>2</sub> ... [V] <sub>M</sub> [CLASS].	<b>91.83</b>

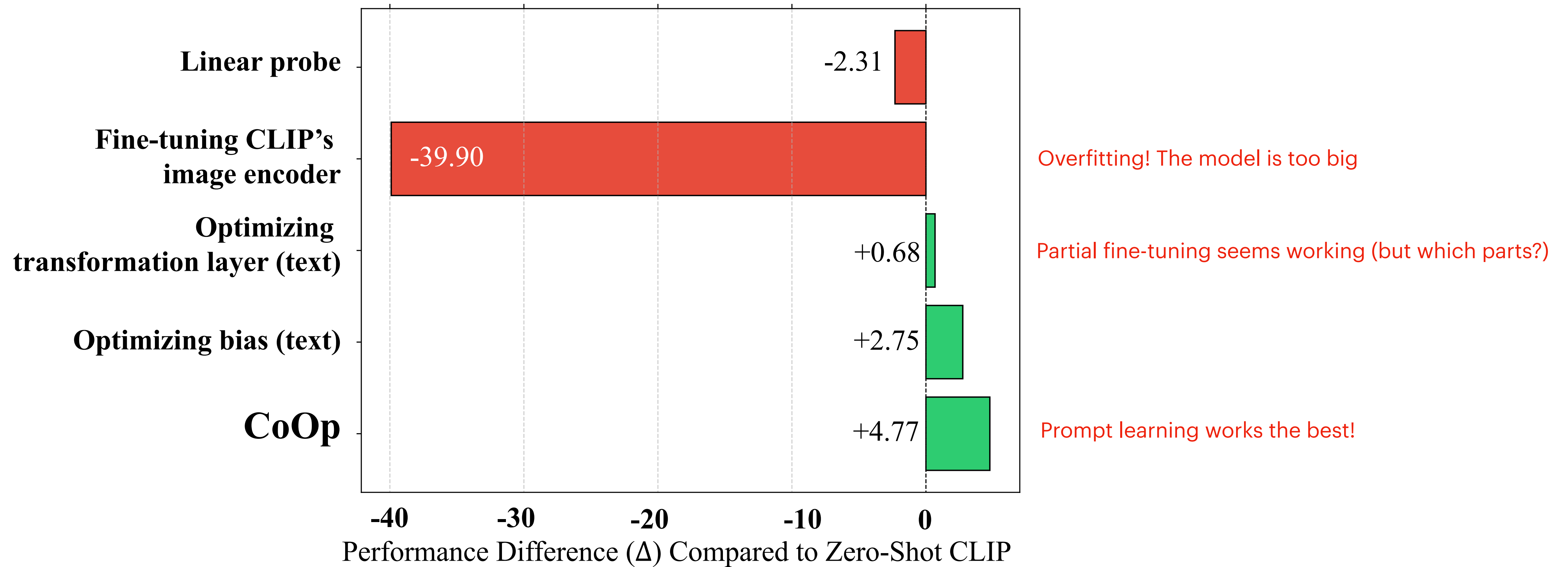
Flowers102	Prompt	Accuracy
	a photo of a [CLASS].	60.86
	a <b>flower</b> photo of a [CLASS].	65.81
	a photo of a [CLASS], <b>a type of flower</b> .	66.14
	[V] <sub>1</sub> [V] <sub>2</sub> ... [V] <sub>M</sub> [CLASS].	<b>94.51</b>

Describable Textures (DTD)	Prompt	Accuracy
	a photo of a [CLASS].	39.83
	a photo of a [CLASS] <b>texture</b> .	40.25
	[CLASS] texture.	42.32
	[V] <sub>1</sub> [V] <sub>2</sub> ... [V] <sub>M</sub> [CLASS].	<b>63.58</b>

EuroSAT	Prompt	Accuracy
	a photo of a [CLASS].	24.17
	a <b>satellite</b> photo of [CLASS].	37.46
	a centered satellite photo of [CLASS].	37.56
	[V] <sub>1</sub> [V] <sub>2</sub> ... [V] <sub>M</sub> [CLASS].	<b>83.53</b>

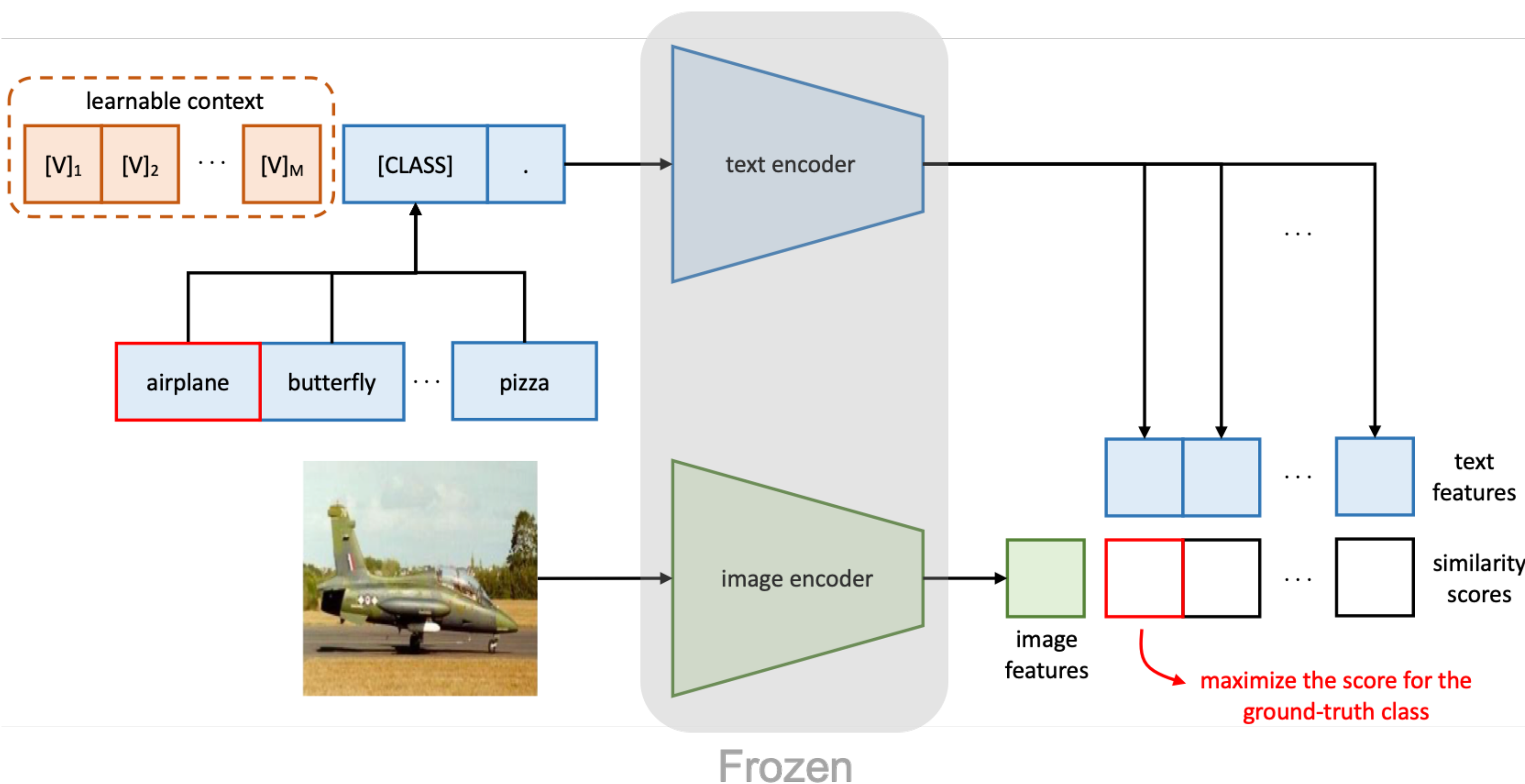


# Fine-tuning is also hard



# Context Optimization (CoOp)

/ku:p/



By forwarding a prompt  $\mathbf{t}$  to the text encoder  $g(\cdot)$ , we can obtain a classification weight vector representing a visual concept (still from the [EOS] token position). The prediction probability is computed as

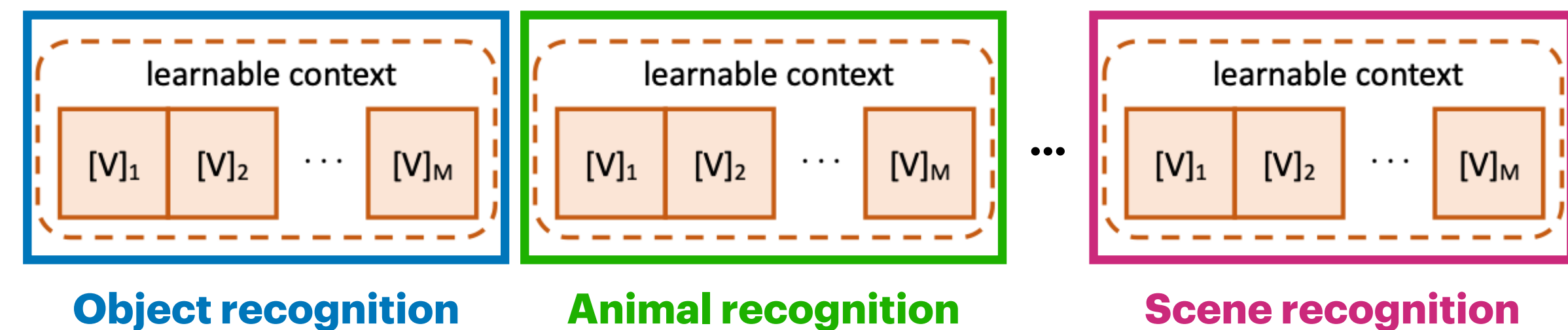
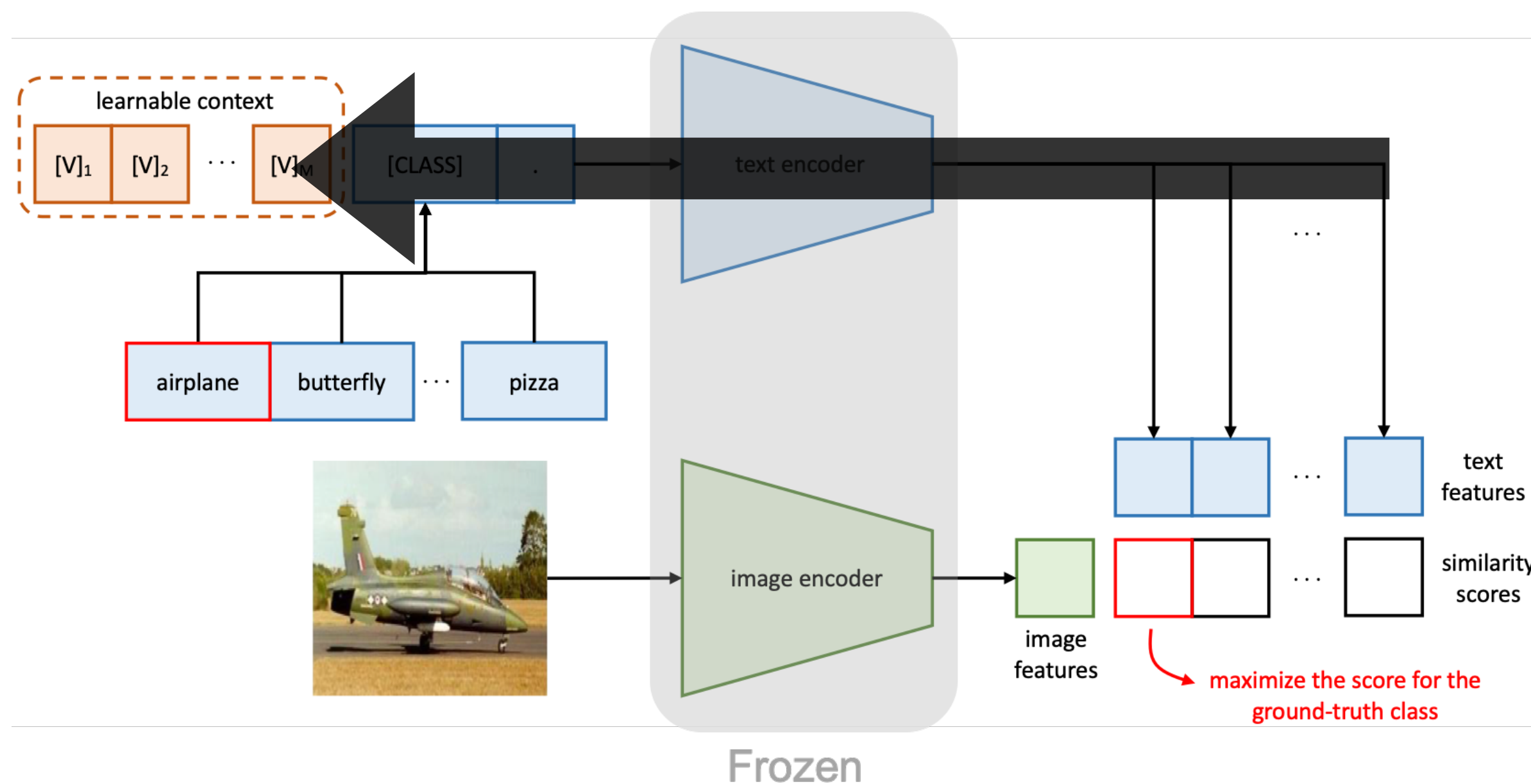
$$p(y = i | \mathbf{x}) = \frac{\exp(\cos(g(\mathbf{t}_i), \mathbf{f})/\tau)}{\sum_{j=1}^K \exp(\cos(g(\mathbf{t}_j), \mathbf{f})/\tau)}, \quad (3)$$

where the class token within each prompt  $\mathbf{t}_i$  is replaced by the corresponding word embedding vector(s) of the  $i$ -th class name.



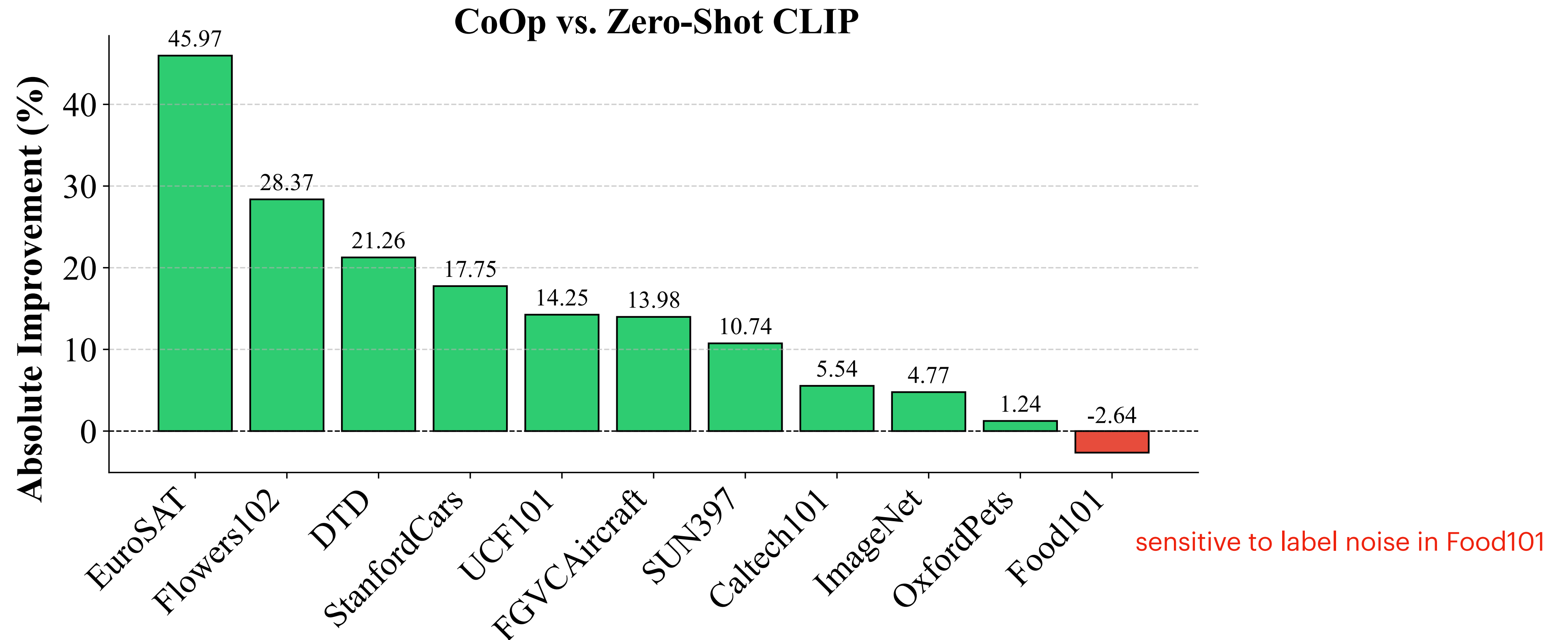
# Why do prompt learning?

- Enjoys rich gradient information
- Mitigates overfitting (few parameters)
- Reduces storage cost (one per task or user)



# Few-shot learning

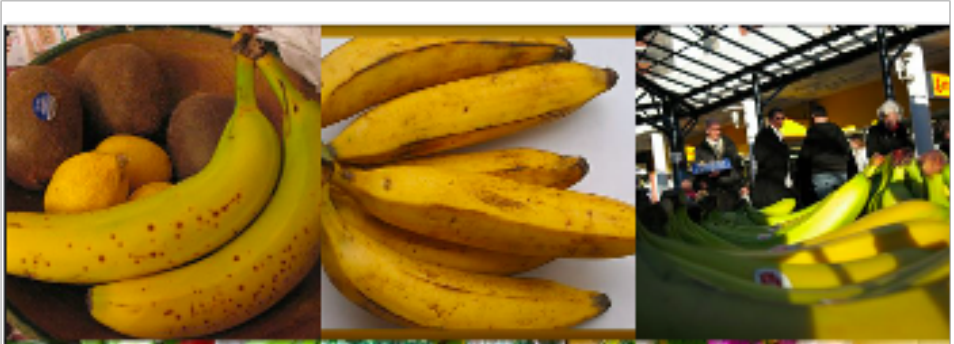


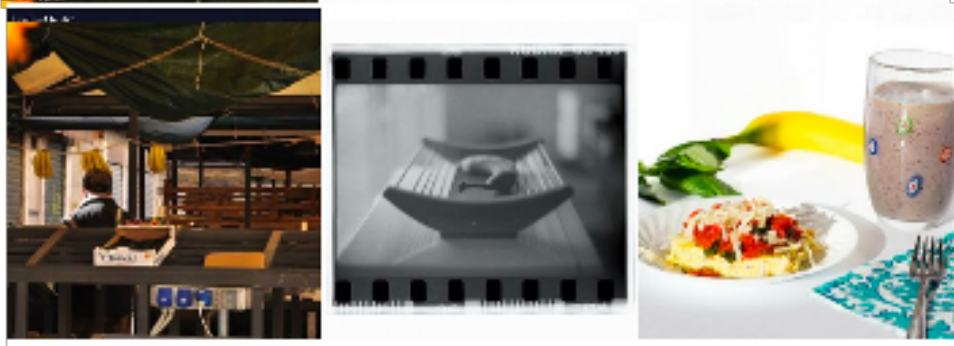

- Works on diverse tasks (objects, animals, scenes, actions, etc.)
- Significantly beats hand-crafted prompts (also needs labels for tuning)





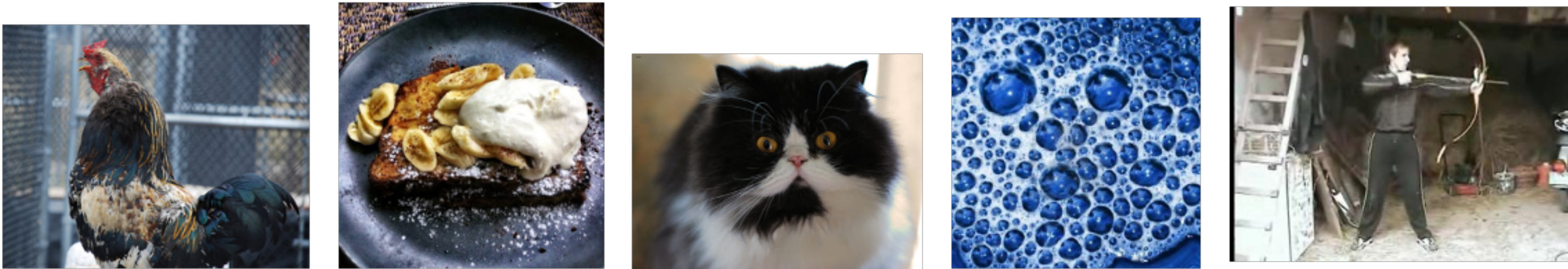
# Domain generalization

- Train on one dataset but test on a different one with domain shifts
- Still beats hand-crafted prompts despite being a learning-based approach

		CLIP	Ours
ImageNet (source)		58.18	<b>63.33</b>
V2 (target)		51.34	<b>55.40</b>
Sketch (target)		33.32	<b>34.67</b>
Adversarial (target)		21.65	<b>23.06</b>
Rendition (target)		56.00	<b>56.60</b>



# Interpretable? Not really



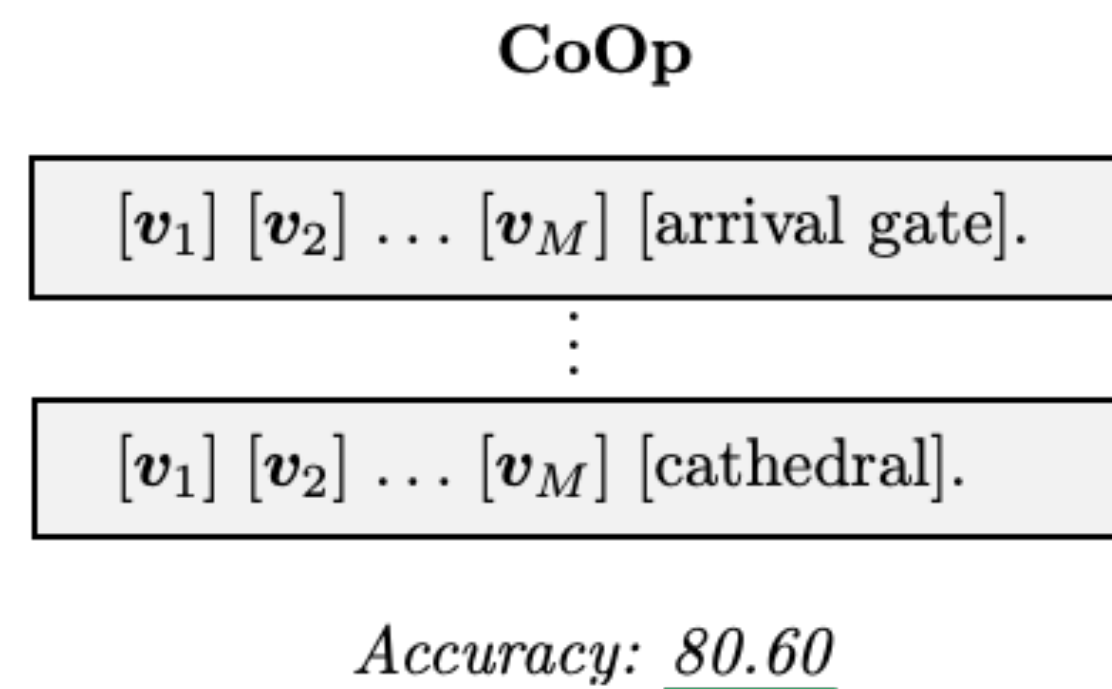
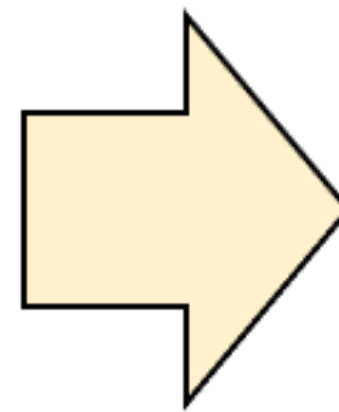
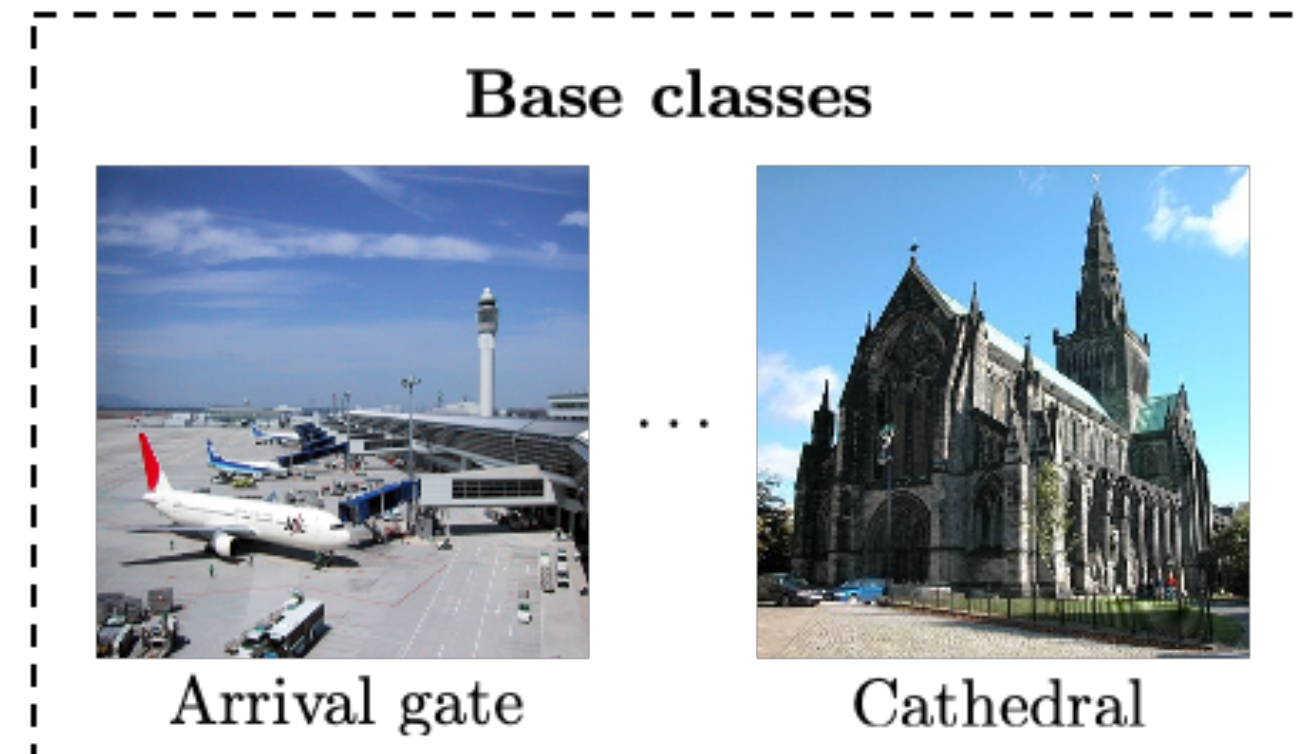
Finding 1: few are somewhat relevant

Finding 2: the whole prompt does not make much sense

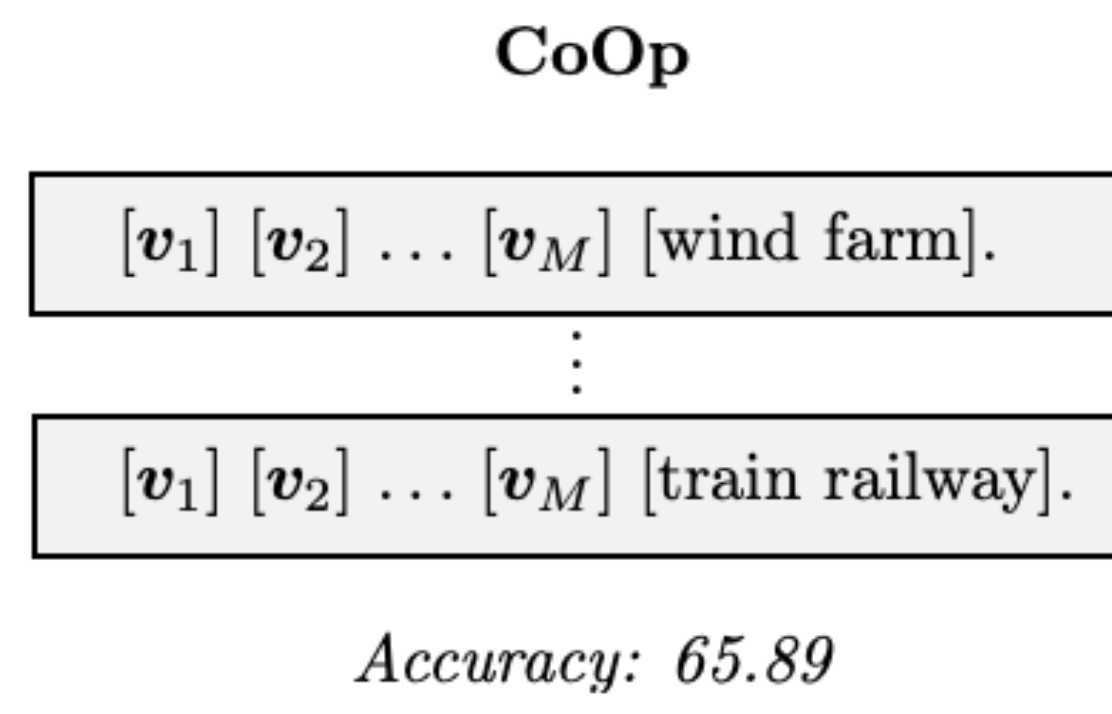
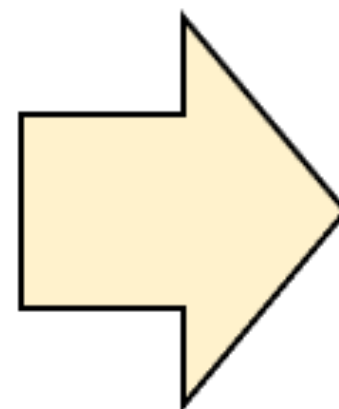
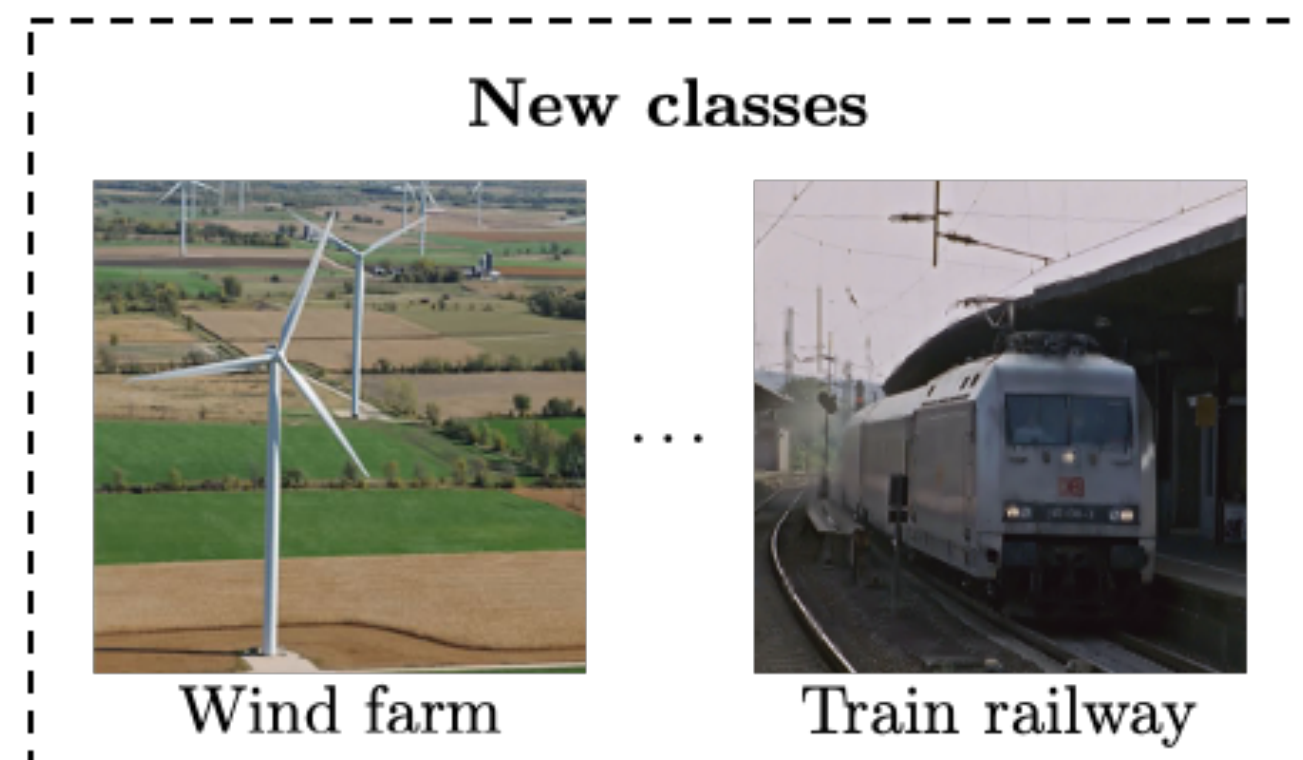
#	ImageNet	Food101	OxfordPets	DTD	UCF101
1	Potd (1.7136)	Lc (0.6752)	Tosc (2.5952)	Boxed (0.9433)	Meteorologist (1.5377)
2	That (1.4015)	Enjoyed (0.5305)	Judge (1.2635)	Seed (1.0498)	Exe (0.9807)
3	Filmed (1.2275)	Beh (0.5390)	Fluffy (1.6099)	Anna (0.8127)	Parents (1.0654)
4	Fruit (1.4864)	Matches (0.5646)	Cart (1.3958)	Mountain (0.9509)	Masterful (0.9528)
5	,... (1.5863)	Nytimes (0.6993)	Harlan (2.2948)	Eldest (0.7111)	Fe (1.3574)
6	°(1.7502)	Prou (0.5905)	Paw (1.3055)	Pretty (0.8762)	Thof (1.2841)
7	Excluded (1.2355)	Lower (0.5390)	Incise (1.2215)	Faces (0.7872)	Where (0.9705)
8	Cold (1.4654)	N/A	Bie (1.5454)	Honey (1.8414)	Kristen (1.1921)
9	Stery (1.6085)	Minute (0.5672)	Snuggle (1.1578)	Series (1.6680)	Imam (1.1297)
10	Warri (1.3055)	~ (0.5529)	Along (1.8298)	Coca (1.5571)	Near (0.8942)
11	Marvelcomics (1.5638)	Well (0.5659)	Enjoyment (2.3495)	Moon (1.2775)	Tummy (1.4303)
12	.: (1.7387)	Ends (0.6113)	Jt (1.3726)	Ih (1.0382)	Hel (0.7644)
13	N/A	Mis (0.5826)	Improving (1.3198)	Won (0.9314)	Boop (1.0491)
14	Lation (1.5015)	Somethin (0.6041)	Srsly (1.6759)	Replied (1.1429)	N/A
15	Muh (1.4985)	Seminar (0.5274)	Asteroid (1.3395)	Sent (1.3173)	Facial (1.4452)
16	.# (1.9340)	N/A	N/A	Piedmont (1.5198)	During (1.1755)



# Generalize beyond the training labels?



... only works for a subset of classes (overfitting)







ImageNet  
↓8.86%



Caltech101  
↓8.19%



Flowers102  
↓37.93%



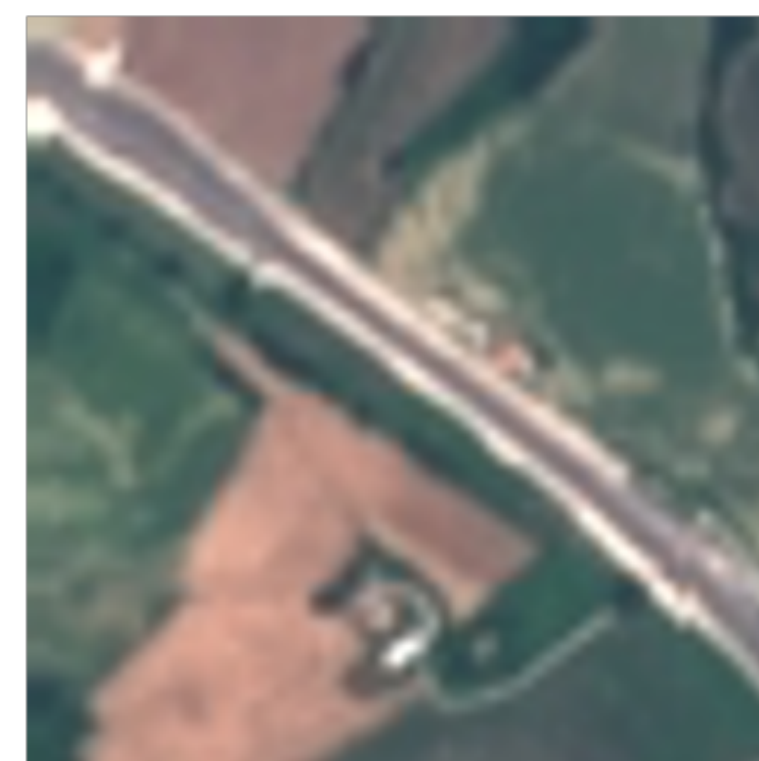
StanfordCars  
↓17.72%



FGVCAircraft  
↓18.14%



DTD  
↓38.26%



EuroSAT  
↓37.45%

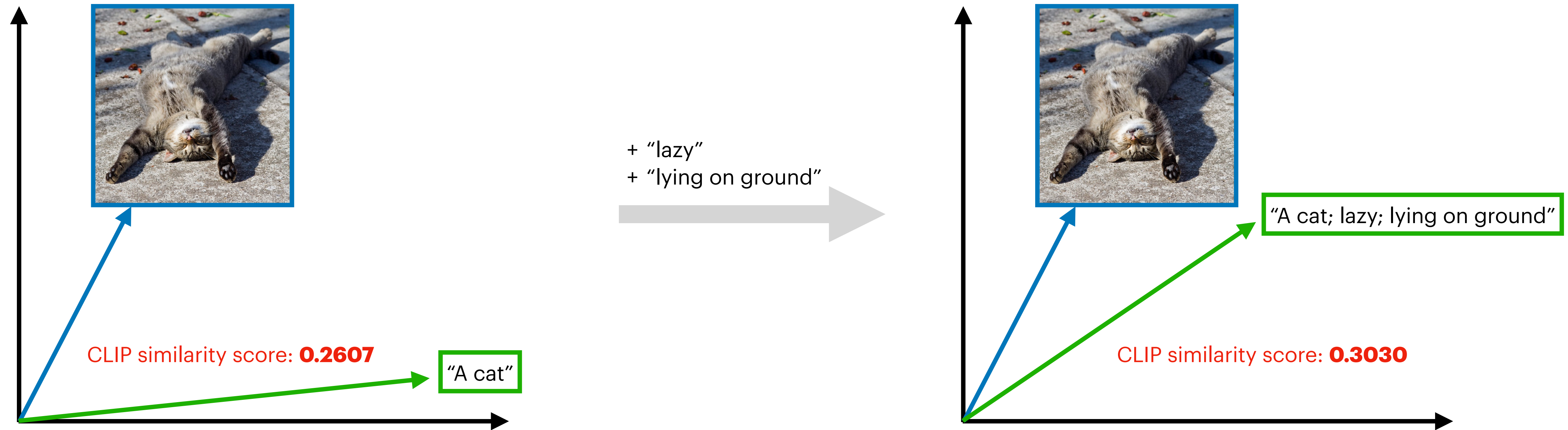


UCF101  
↓28.64%



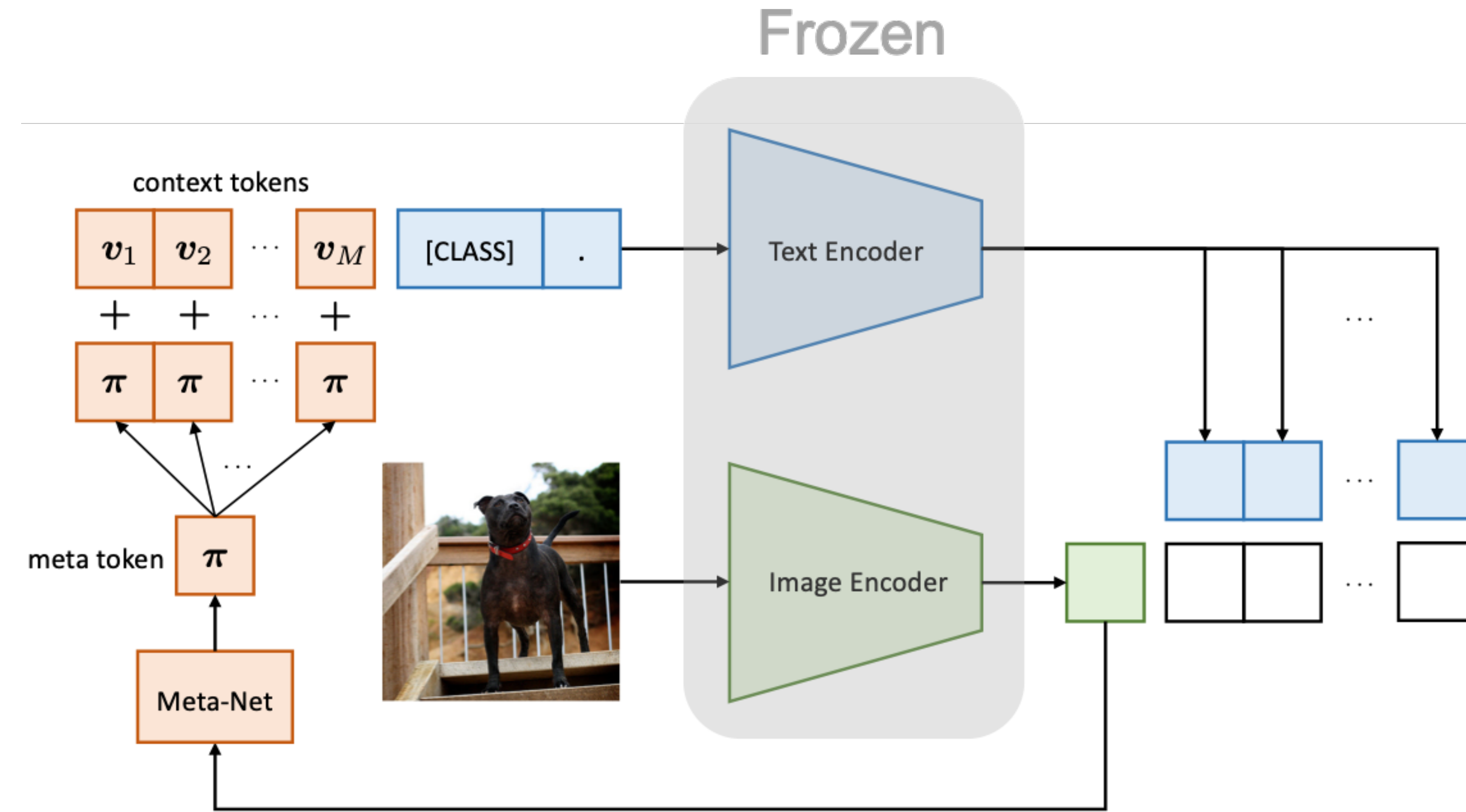
# What is a good prompt?

- Contains instance-specific information
- Pushes text features closer to image features



# Conditional Context Optimization (CoCoOp)

/kəʊˌku:p/



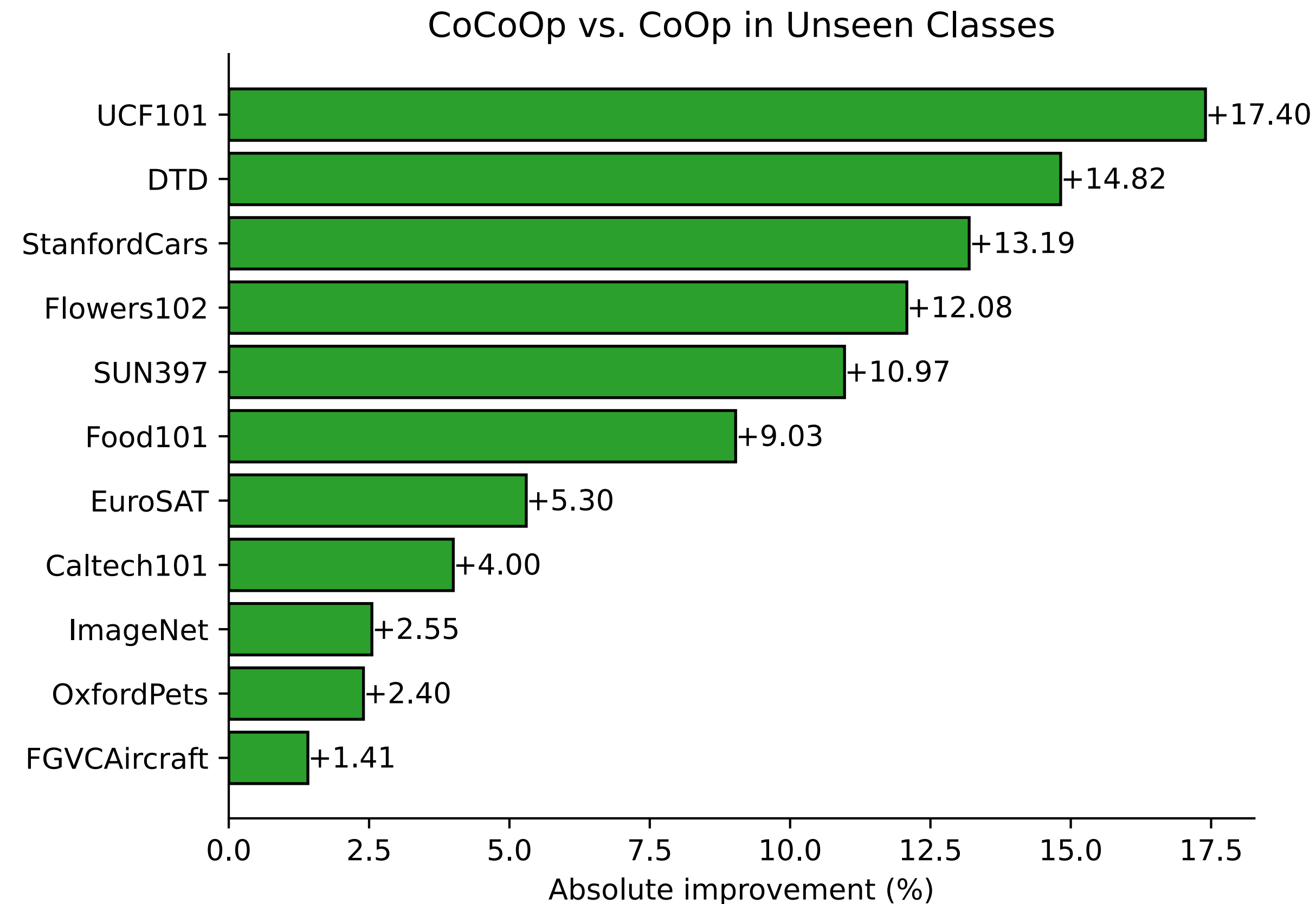
Let  $h_{\theta}(\cdot)$  denote the Meta-Net parameterized by  $\theta$ , each context token is now obtained by  $v_m(\mathbf{x}) = v_m + \pi$  where  $\pi = h_{\theta}(\mathbf{x})$  and  $m \in \{1, 2, \dots, M\}$ . The prompt for the  $i$ -th class is thus conditioned on the input, i.e.,  $t_i(\mathbf{x}) = \{v_1(\mathbf{x}), v_2(\mathbf{x}), \dots, v_M(\mathbf{x}), c_i\}$ . The prediction probability is computed as

$$p(y|\mathbf{x}) = \frac{\exp(\text{sim}(\mathbf{x}, g(t_y(\mathbf{x}))/\tau)}{\sum_{i=1}^K \exp(\text{sim}(\mathbf{x}, g(t_i(\mathbf{x}))/\tau)}. \quad (3)$$



# Key messages

## 1. Conditional prompt learning is more generalizable



# Key messages

- 1. Conditional prompt learning is more generalizable
- 2. Conditional prompt learning is more transferable

Table 2. **Comparison of prompt learning methods in the cross-dataset transfer setting.** Prompts applied to the 10 target datasets are learned from ImageNet (16 images per class). Clearly, CoCoOp demonstrates better transferability than CoOp.  $\Delta$  denotes CoCoOp’s gain over CoOp.

	Source					Target						
	ImageNet	Caltech101	OxfordPets	StanfordCars	Flowers102	Food101	FGVCAircraft	SUN397	DTD	EuroSAT	UCF101	Average
CoOp [62]	71.51	93.70	89.14	64.51	68.71	85.30	18.47	64.15	41.92	46.39	66.55	63.88
CoCoOp	71.02	94.43	90.14	65.32	71.88	86.06	22.94	67.36	45.73	45.37	68.21	65.74
$\Delta$	-0.49	+0.73	+1.00	+0.81	+3.17	+0.76	+4.47	+3.21	+3.81	-1.02	+1.66	+1.86



# Key messages

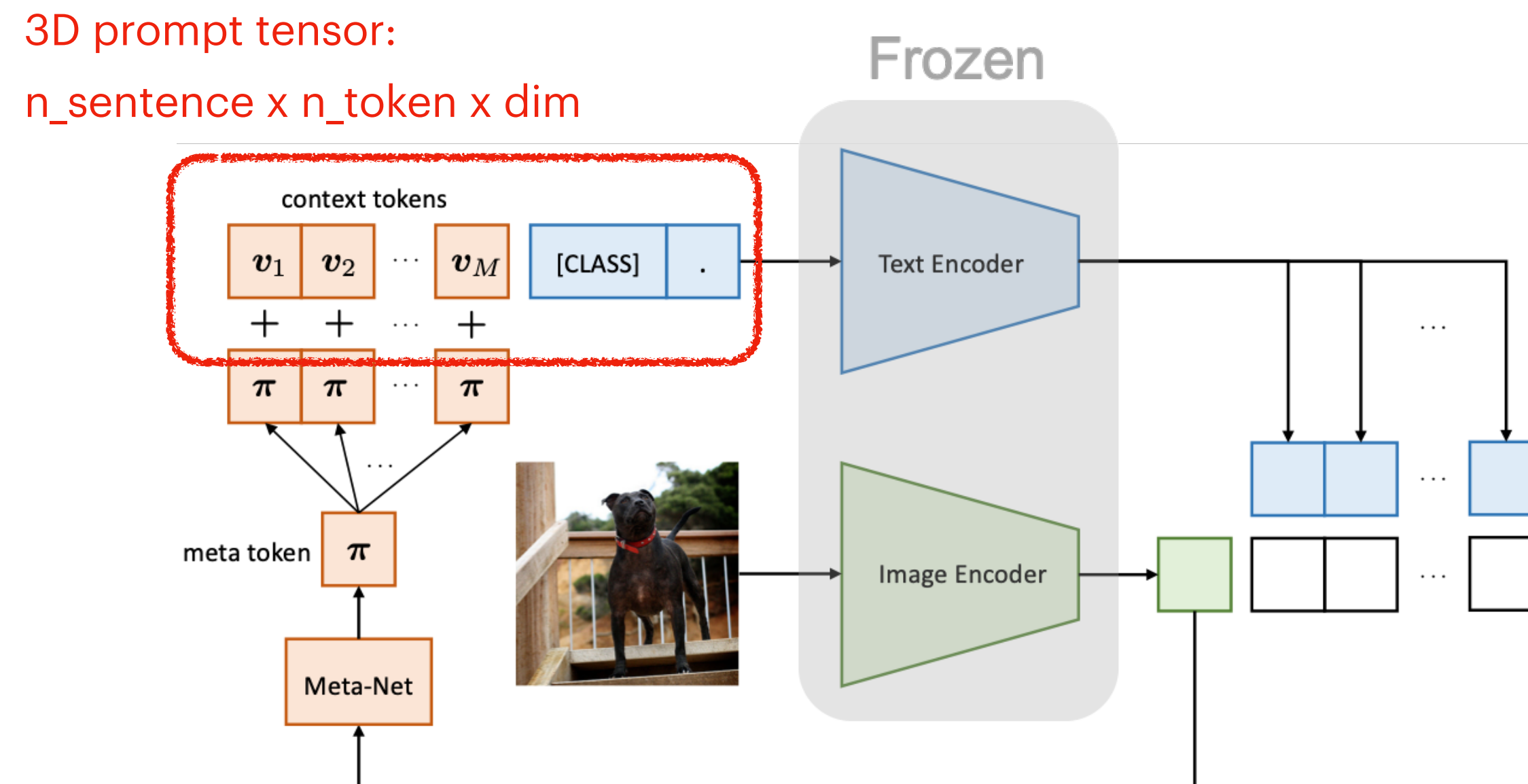
- 1. Conditional prompt learning is more generalizable
- 2. Conditional prompt learning is more transferable
- 3. Conditional prompt learning is more robust

Table 3. **Comparison of manual and learning-based prompts in domain generalization.** CoOp and CoCoOp use as training data 16 images from each of the 1,000 classes on ImageNet. In general, CoCoOp is more domain-generalizable than CoOp.

	Learnable?	Source	Target			
		ImageNet	ImageNetV2	ImageNet-Sketch	ImageNet-A	ImageNet-R
CLIP [40]		66.73	60.83	46.15	47.77	73.96
CoOp [62]	✓	<b>71.51</b>	<b>64.20</b>	47.99	49.71	75.21
CoCoOp	✓	71.02	64.07	<b>48.75</b>	<b>50.63</b>	<b>76.18</b>

# Key messages

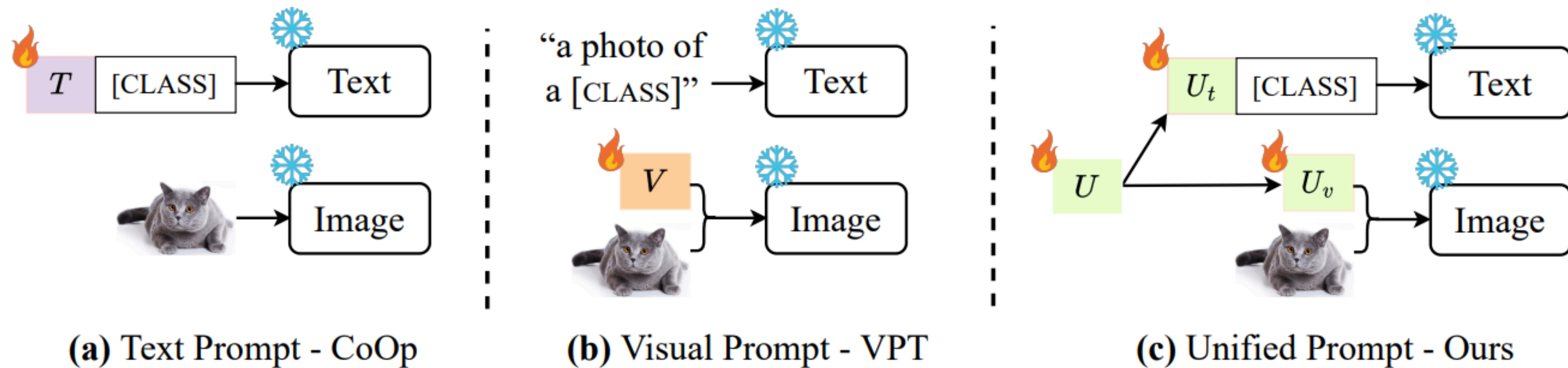
1. Conditional prompt learning is more generalizable
2. Conditional prompt learning is more transferable
3. Conditional prompt learning is more robust
4. Conditional prompt learning is **very slow to train (batch\_size=1)**





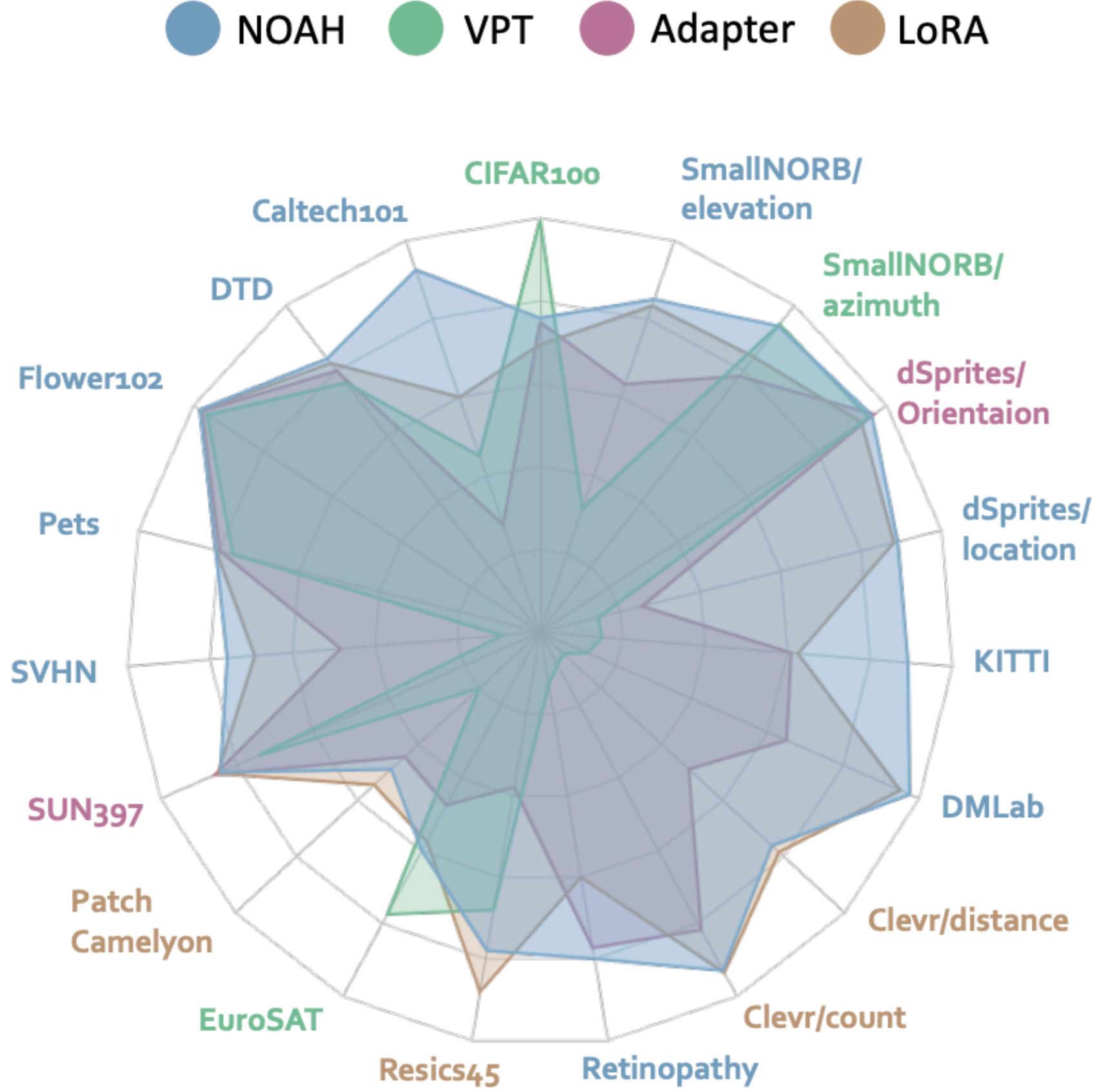
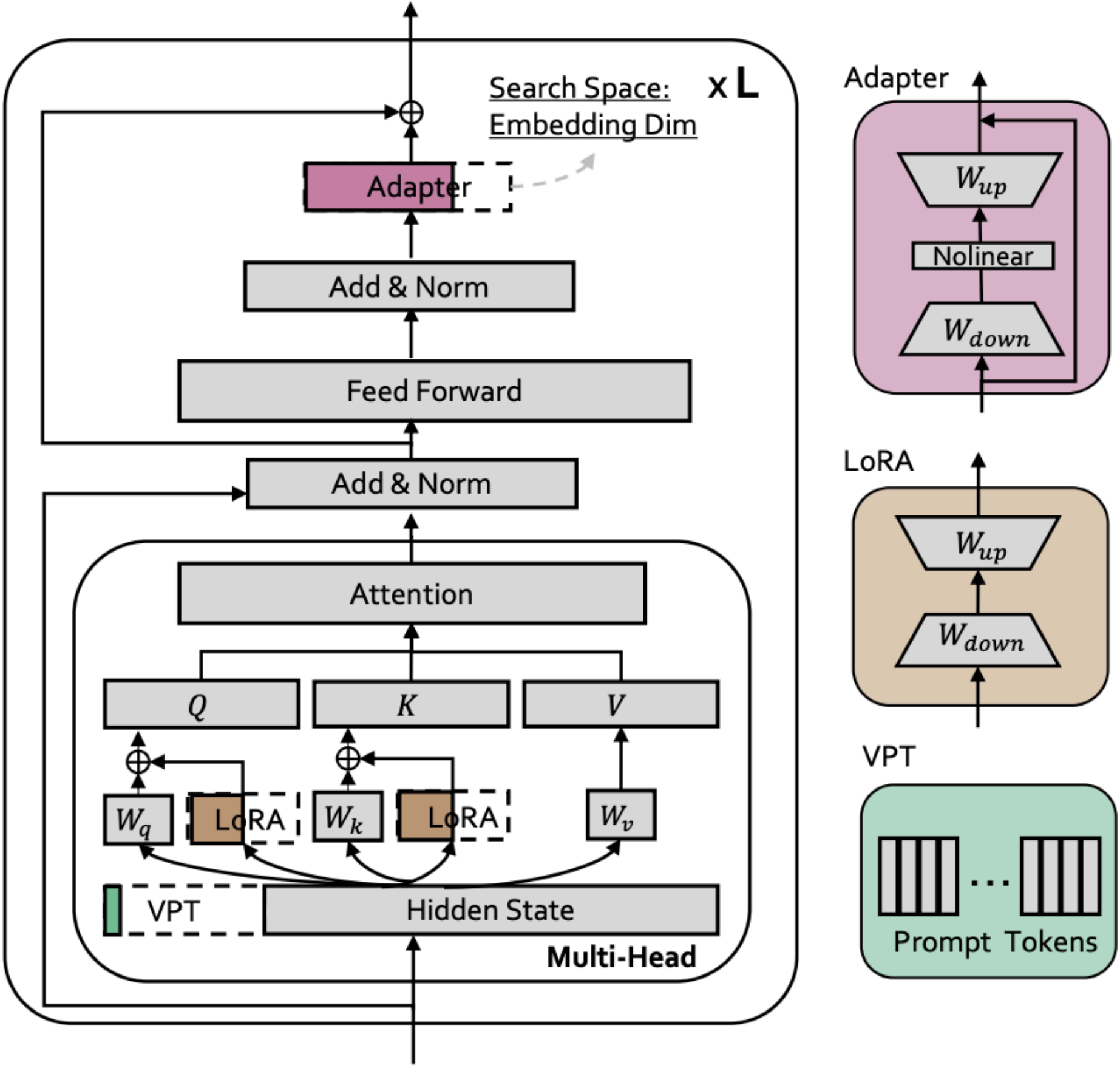
# Multimodal prompt learning

- Idea: simultaneously adjust text and image features
- Same performance but much faster training



#	Method	Source	Target				Average	OOD Average
		ImageNet	-V2	-S	-A	-R		
1	CoOp	71.51	64.20	47.99	49.71	75.21	61.72	59.28
2	CoCoOp	71.02	64.07	<b>48.75</b>	50.63	76.18	62.13	59.91
3	VPT-shallow	68.98	62.10	47.68	47.19	76.10	60.38	58.27
4	VPT-deep	70.57	63.67	47.66	43.85	74.42	60.04	57.40
5	UPT	<b>72.63</b>	<b>64.35</b>	48.66	<b>50.66</b>	<b>76.24</b>	<b>62.51</b>	<b>59.98</b>

# Have more compute? Do prompt search





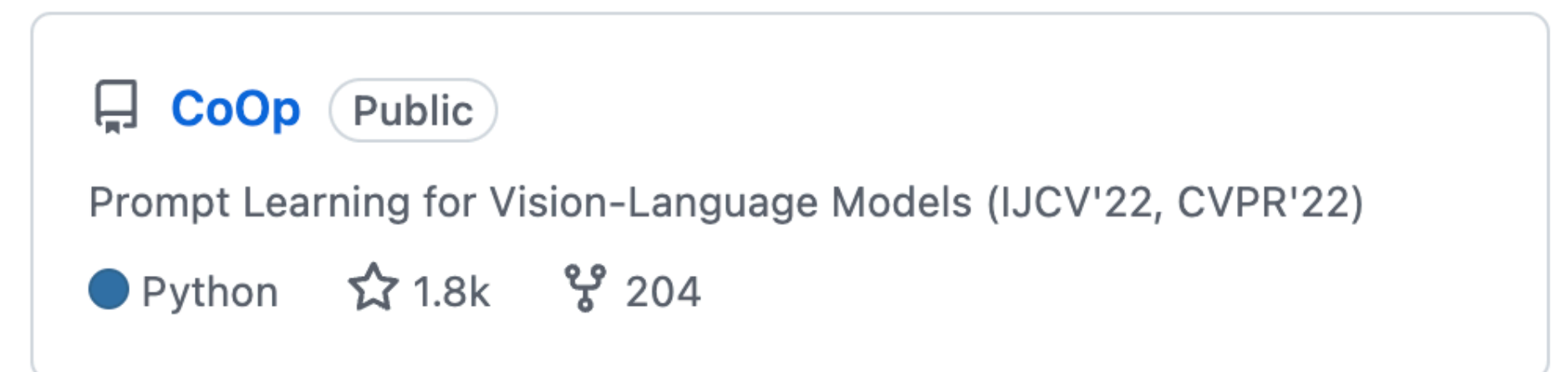
# Take home messages

- VLMs largely reshaped the landscape of visual recognition
- Deploying VLMs in the real world is a non-trivial problem
- Prompt learning is a data-efficient adaptation method
- Conditional prompt learning works better but is too slow
- Multimodal prompt learning strikes a good balance between performance and speed
- Do NAS to search for the best adaptation modules if more compute is available

## Relevant prompting papers

- Learning to Prompt for Vision-Language Models
- Conditional Prompt Learning for Vision-Language Models
- Unified Vision and Language Prompt Learning
- Neural Prompt Search

Open-source code: <https://github.com/KaiyangZhou/CoOp>



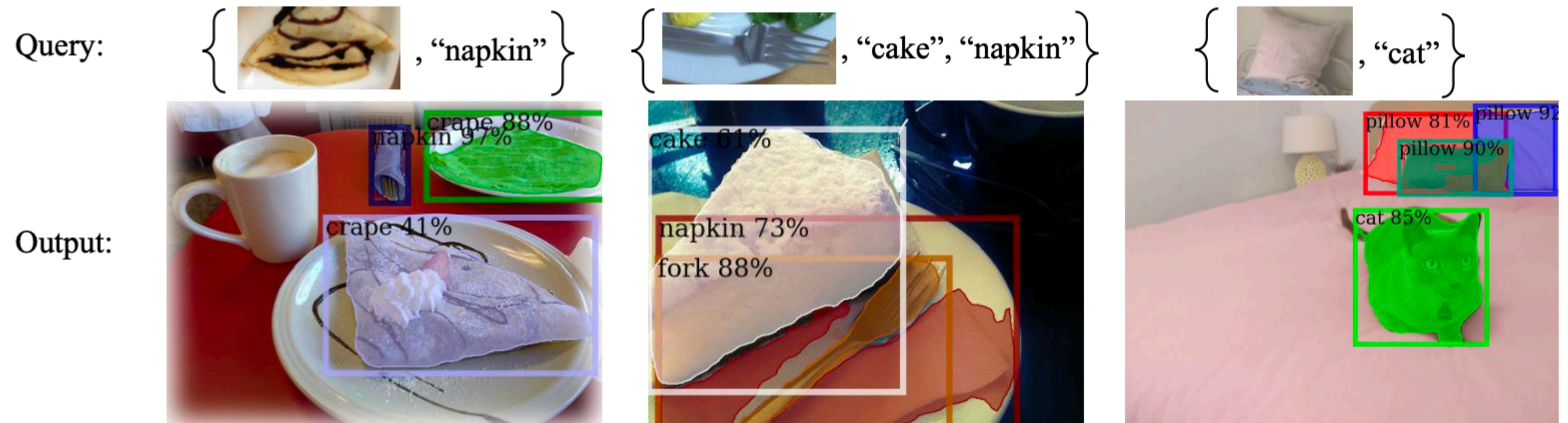
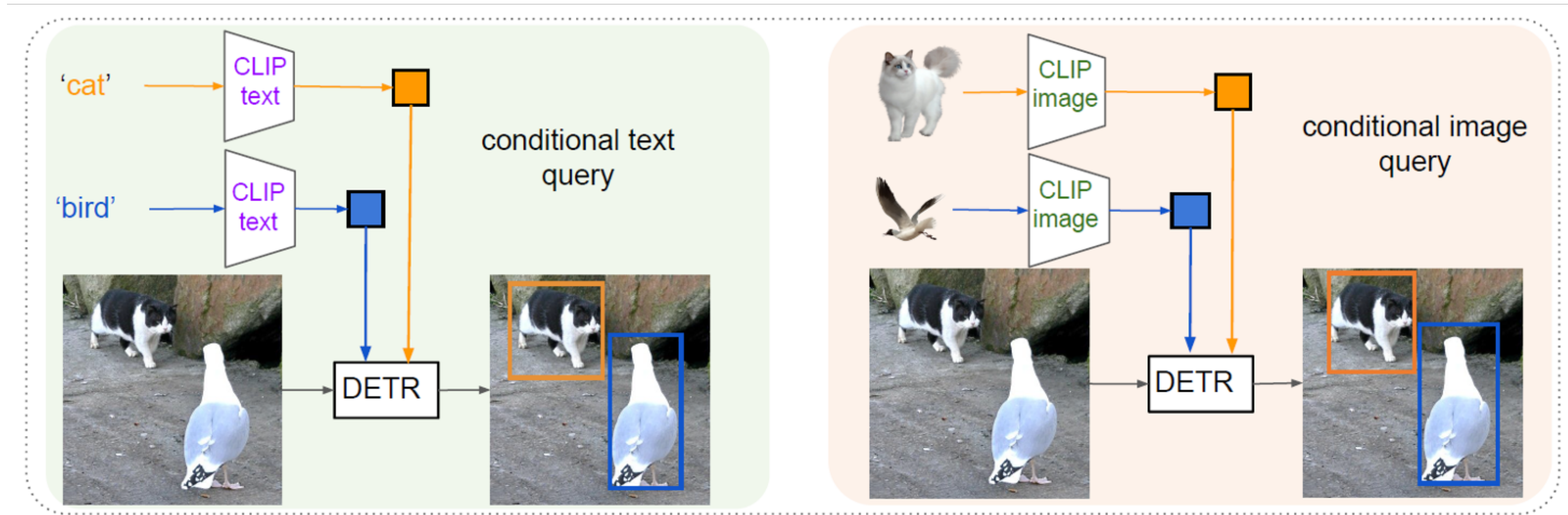
# Outline

- History
- Pre-training
- Prompting
- Applications

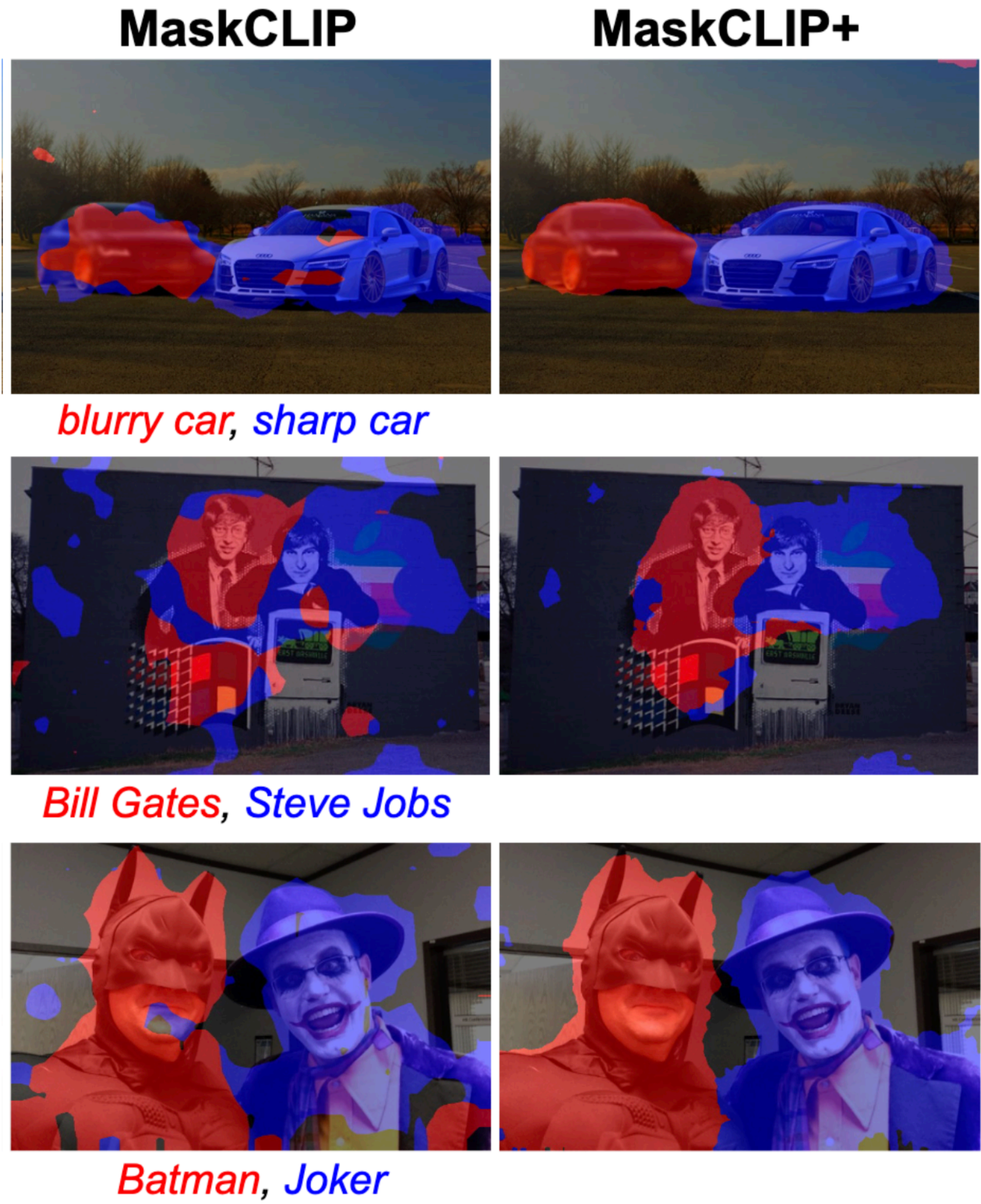
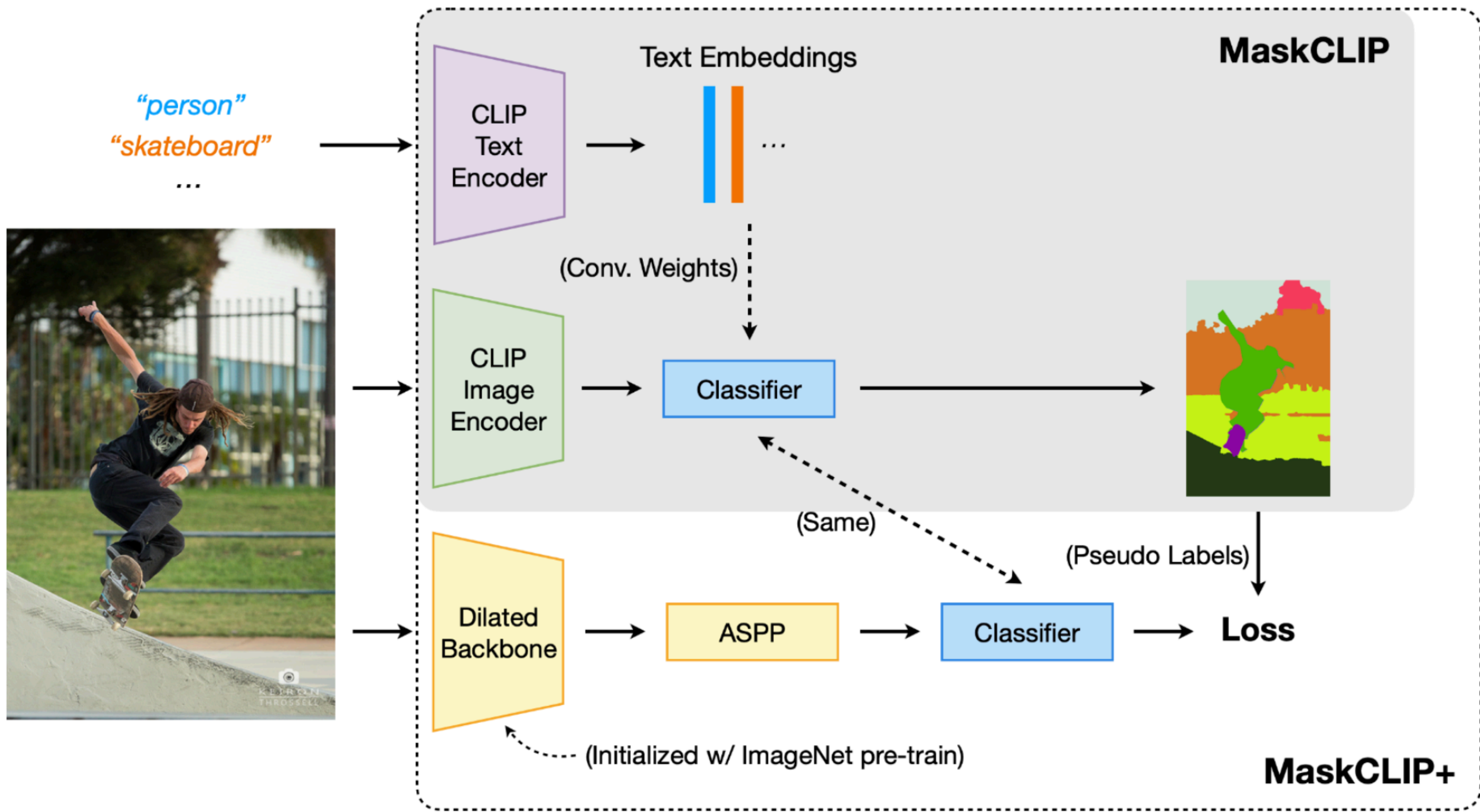


# Open-Vocabulary Perception





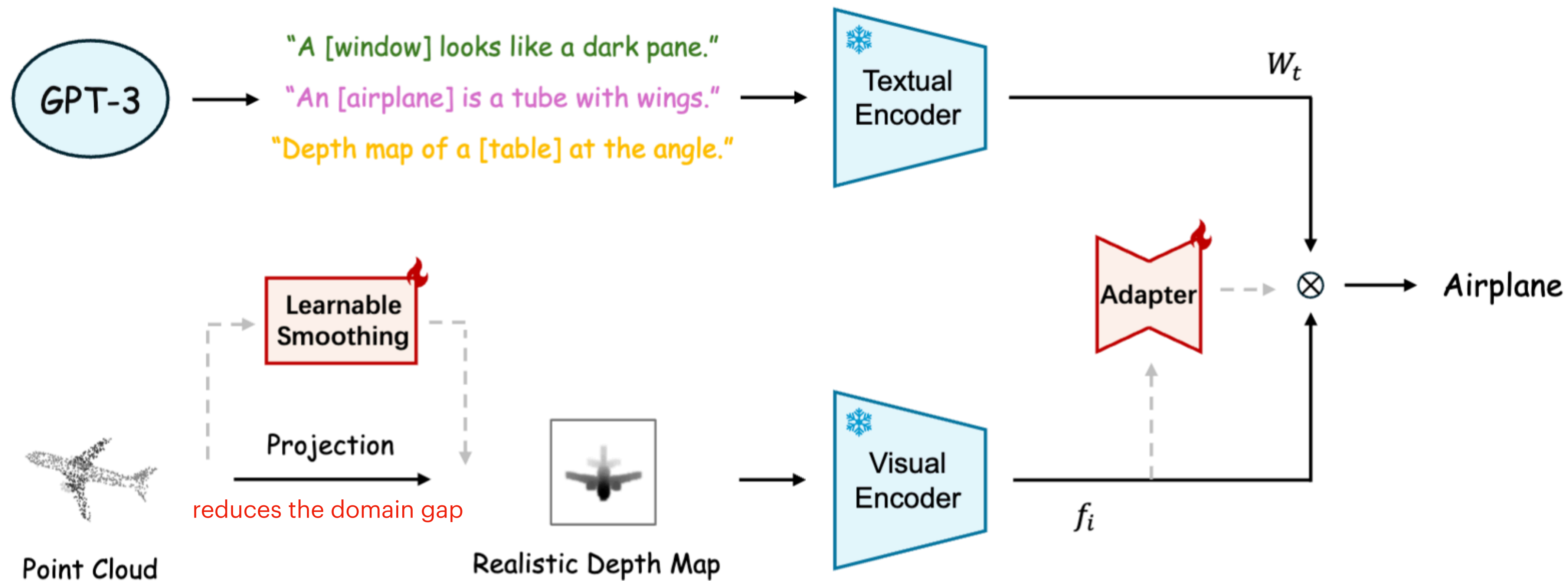


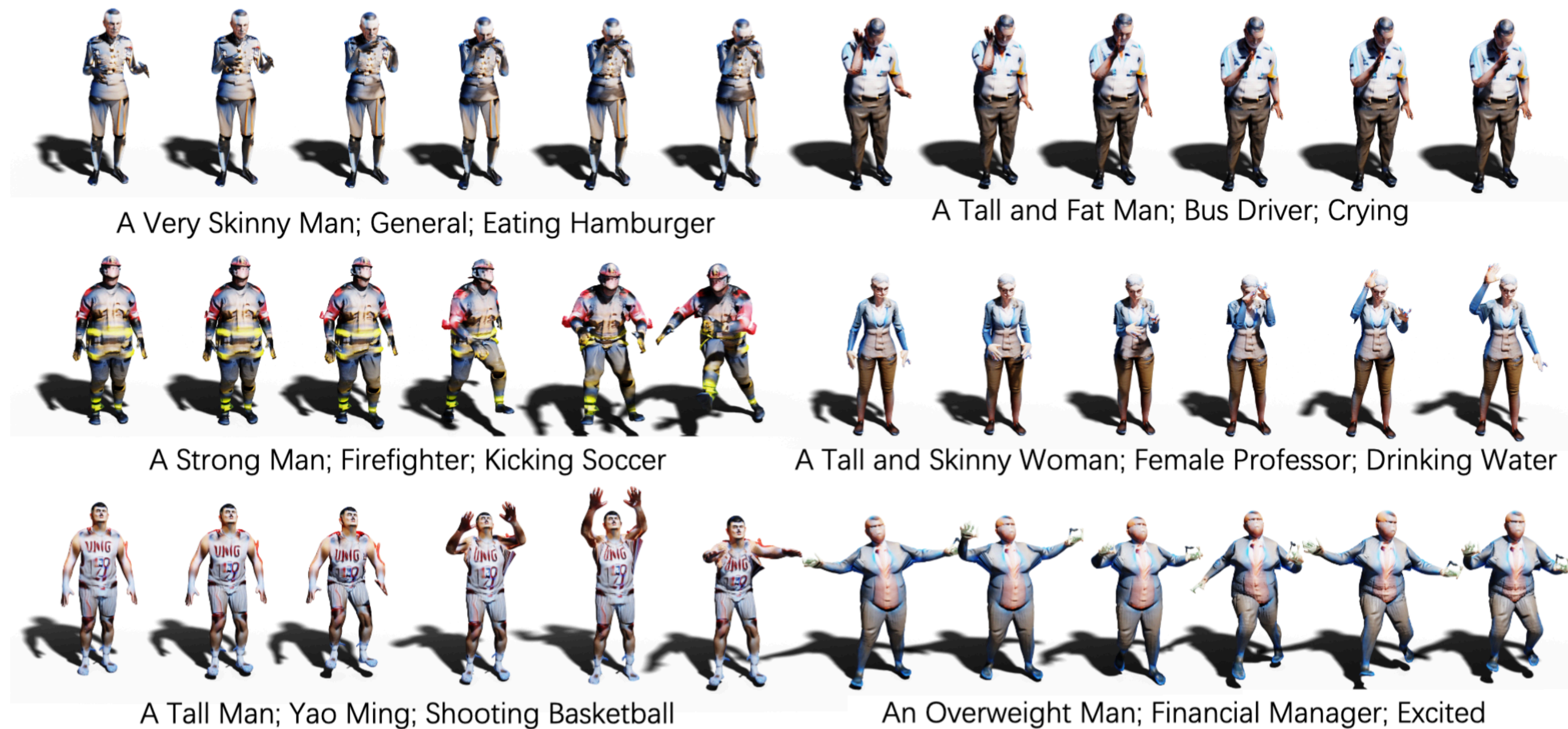
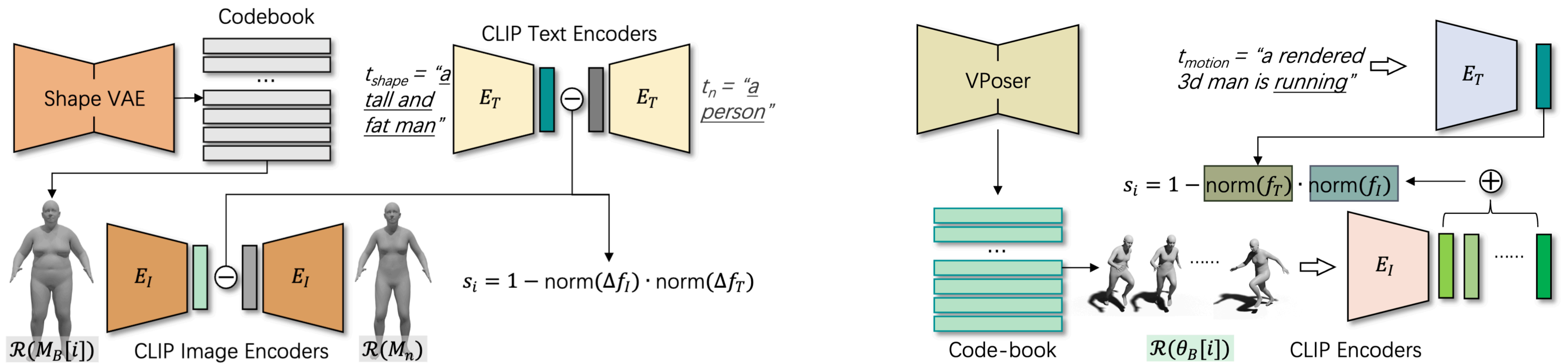




# 3D Understanding and Generation

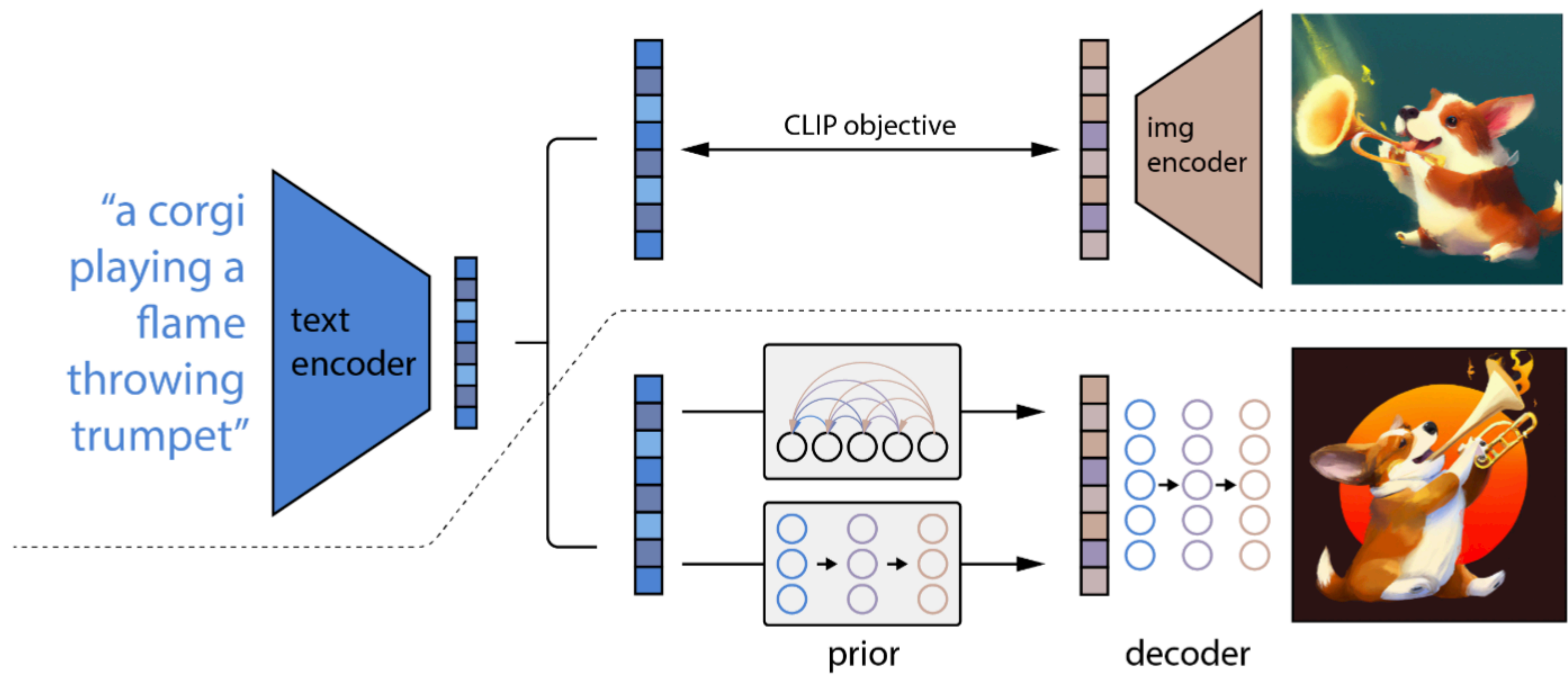






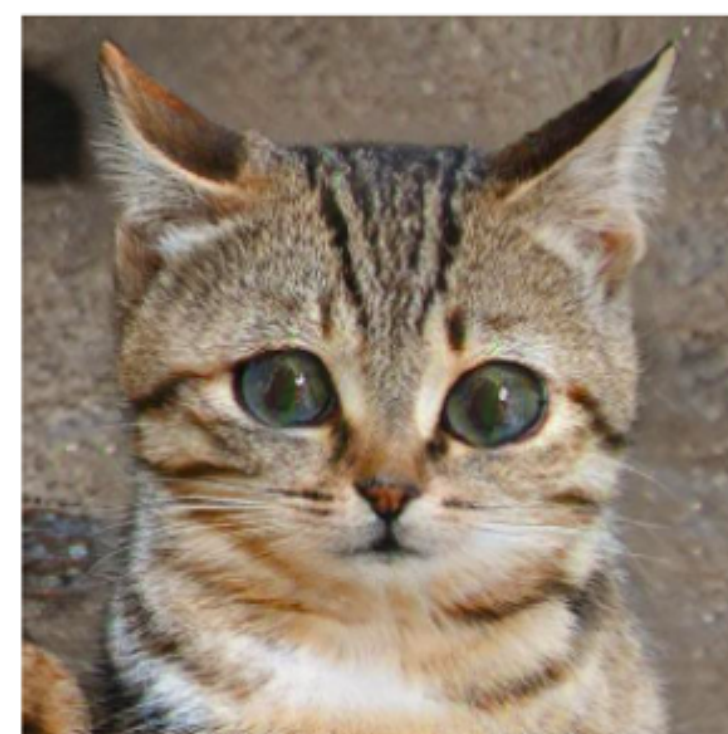
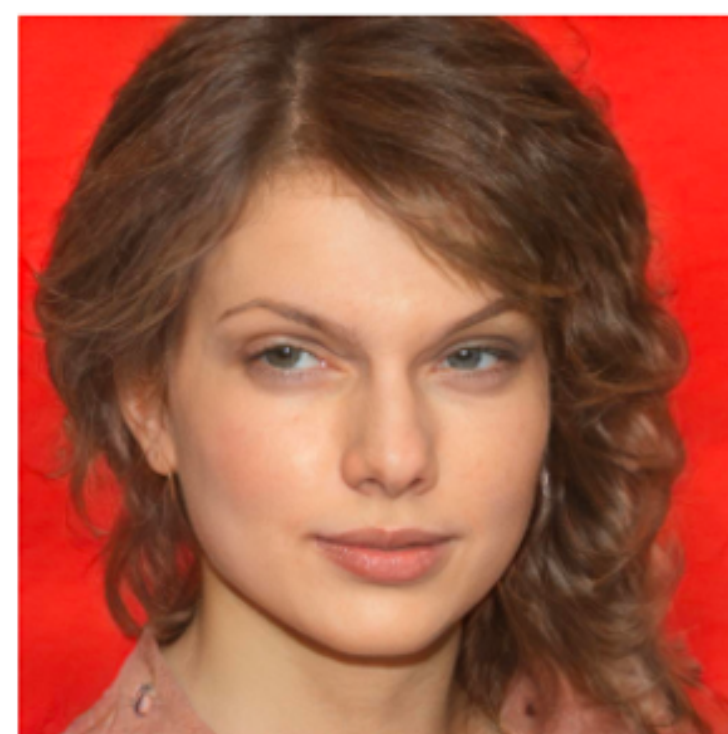
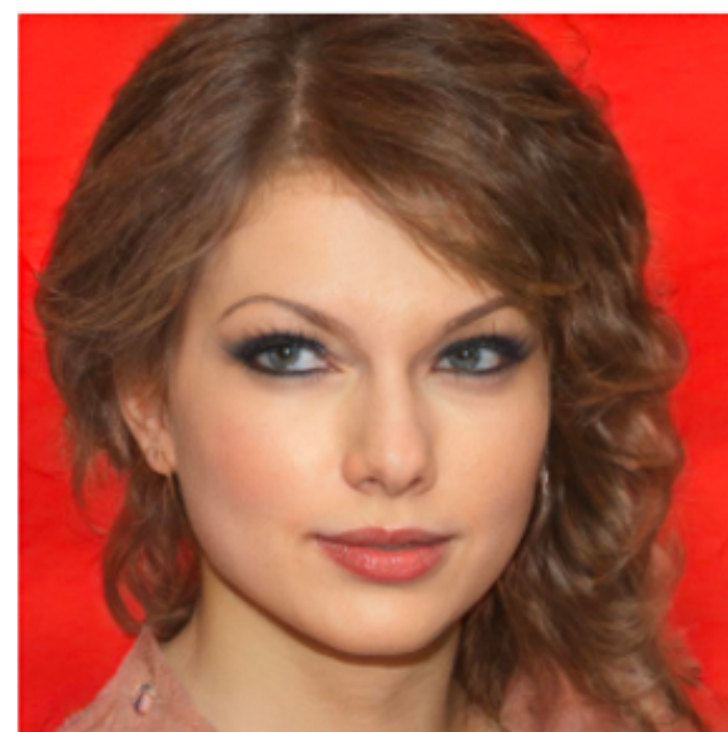


# Generative Models and Creativity



panda mad scientist mixing sparkling chemicals, artstation





“Emma Stone”

“Mohawk hairstyle”

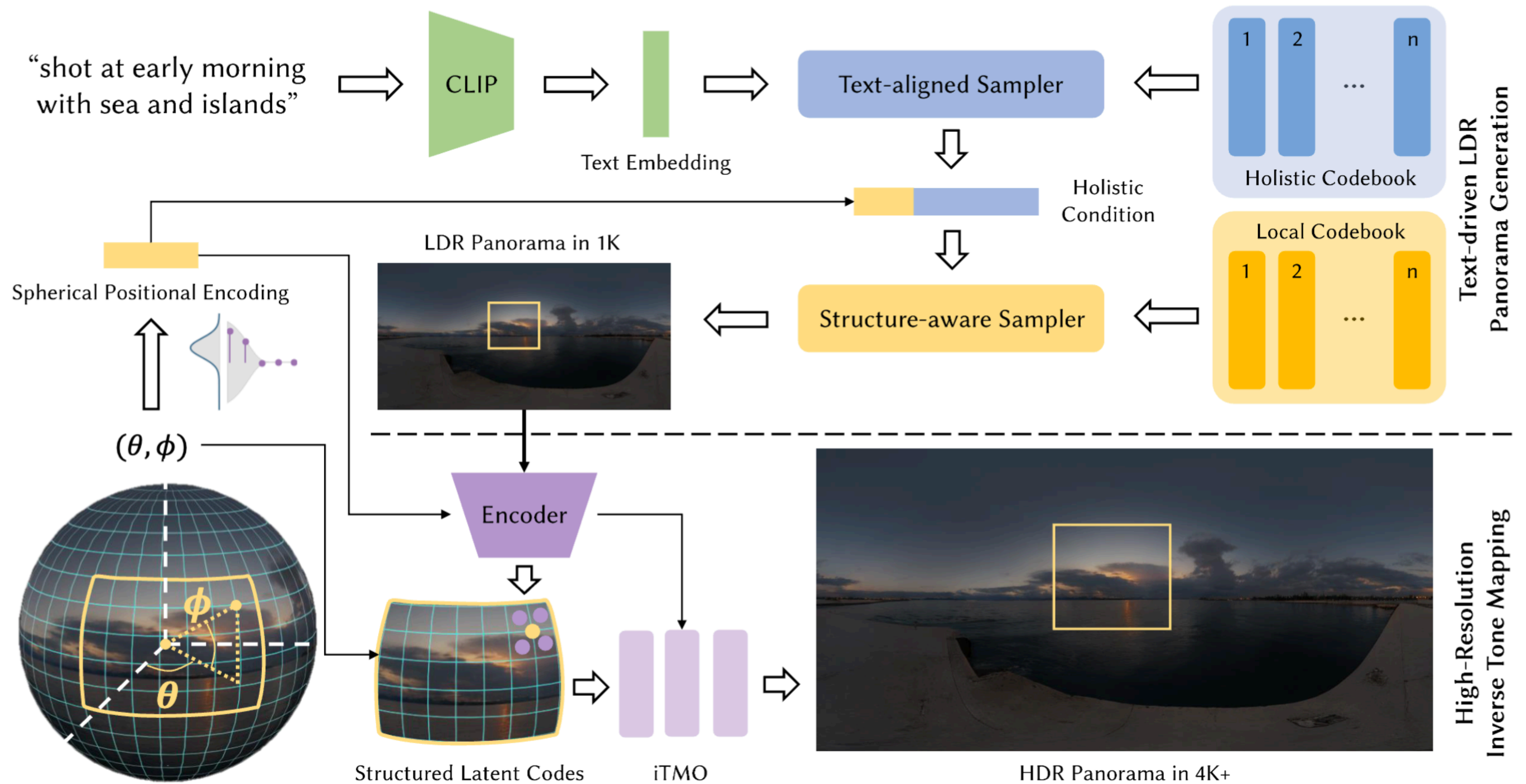
“Without makeup”

“Cute cat”

“Lion”

“Gothic church”







# Recap

- **History:**
  - evolution of vision and language models, convergence to VLMs
- **Pre-training**
  - contrastive learning, dual encoders, image-text pairs
- **Prompting**
  - prompt engineering, prompt learning
- **Applications**
  - open-vocabulary perception, 3D, GenAI