# Vision and Language

Berkeley
N L P

## Eric Wallace

with thanks to Rudy Corona & Daniel Fried
CS 288, 4/12/2022

# What is Language Grounding?

‣ Language often refers *to the world*

# What is Language Grounding?

▸ Language often refers *to the world*

▸ Grounding is tying language to non-linguistic things (e.g., databases, vision, sound)

# What is Language Grounding?

▸ Language often refers *to the world*

▸ Grounding is tying language to non-linguistic things (e.g., databases, vision, sound)

▸ Today we will talk about grounding into *visual* environments:



*"Add the tomatoes and mix"*



*"Take me to the shop on the corner"*

# Grounding

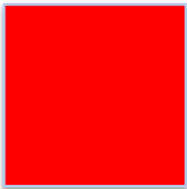‣ (Some) possible things to map language to:

# Grounding

‣ (Some) possible things to map language to:

- **Low-level percepts**: *red* means this set of RGB values, *loud* means lots of decibels on our microphone, *soft* means these properties on our haptic sensor…
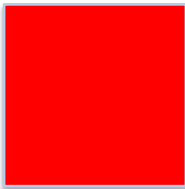
# Grounding

▸ (Some) possible things to map language to:

- **Low-level percepts**: *red* means this set of RGB values, *loud* means lots of decibels on our microphone, *soft* means these properties on our haptic sensor…

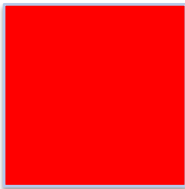- **High-level percepts**: *cat* means this type of pattern

# Grounding

‣ (Some) possible things to map language to:

- **Low-level percepts**: *red* means this set of RGB values, *loud* means lots of decibels on our microphone, *soft* means these properties on our haptic sensor…

- **High-level percepts**: *cat* means this type of pattern

- **Embodiment (effects on the world)**: *go left* means the robot turns left, *speed up* means increasing actuation
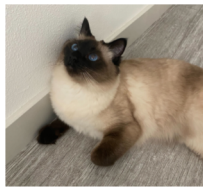
# Grounding

▸ (Some) possible things to map language to:

- **Low-level percepts**: *red* means this set of RGB values, *loud* means lots of decibels on our microphone, *soft* means these properties on our haptic sensor…

- **High-level percepts**: *cat* means this type of pattern

- **Embodiment (effects on the world)**: *go left* means the robot turns left, *speed up* means increasing actuation

- **Social (effects on others)**: polite language is correlated with longer forum discussions

# Grounding

‣ (Some) possible things to map language to:

- **Low-level percepts**: *red* means this set of RGB values, *loud* means lots of decibels on our microphone, *soft* means these properties on our haptic sensor…

- **High-level percepts**: *cat* means this type of pattern

- **Embodiment (effects on the world)**: *go left* means the robot turns left, *speed up* means increasing actuation

- **Social (effects on others)**: polite language is correlated with longer forum discussions

For a nice taxonomy, related work, and examples, see *Experience Grounds Language* [Bisk et al. 2020]

A Gallery of Tasks

# Image Captioning



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.

**Microsoft COCO Captions**: Chen et al. 2015

# Conditional Generation (2D)



vibrant portrait painting of Salvador Dalí with a robotic half face

a shiba inu wearing a beret and black turtleneck

a close up of a handpalm with leaves growing from it

an espresso machine that makes coffee from human souls, artstation

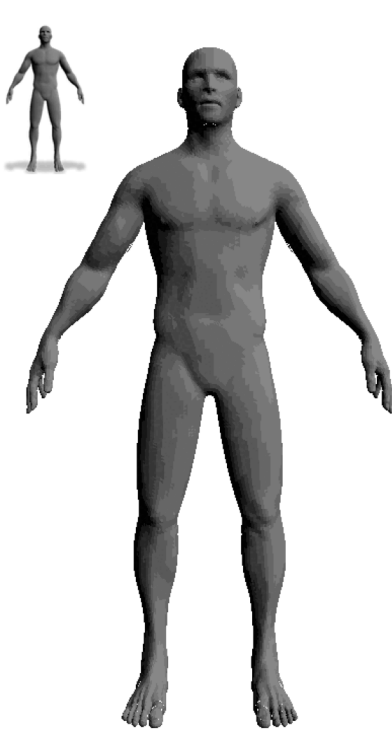panda mad scientist mixing sparkling chemicals, artstation

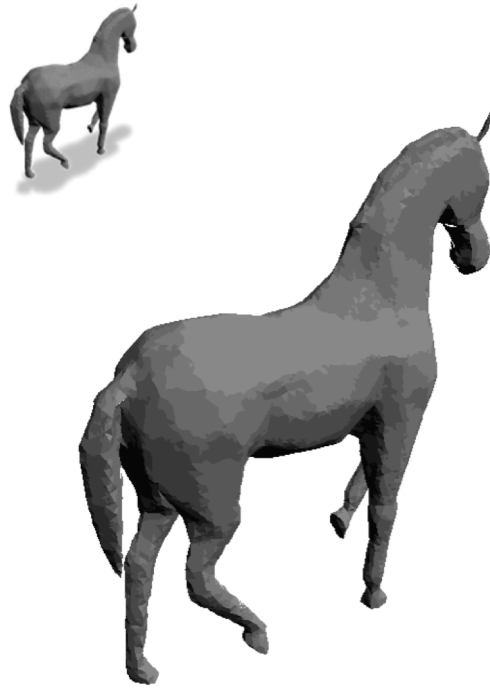a corgi's head depicted as an explosion of a nebula

**DALL-E 2**: Ramesh et al. 2022
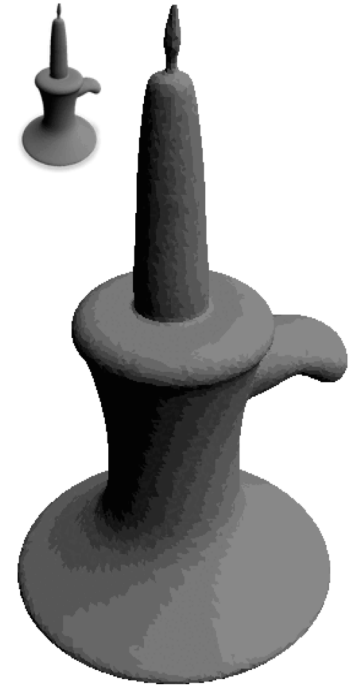
# Conditional Generation (3D)



"Iron Man"

"Astronaut Horse"

"Colorful Crochet Candle"

**Text2Mesh**: Michel et al. 2021

# Visual Question Answering



What is the dog wearing?
life jacket          collar
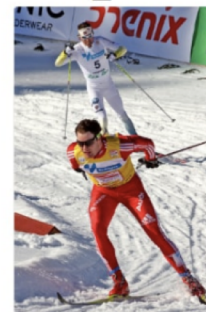
How many skiers are there?
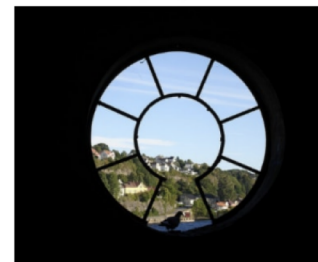2          1

What number is on the train?
7907          8551
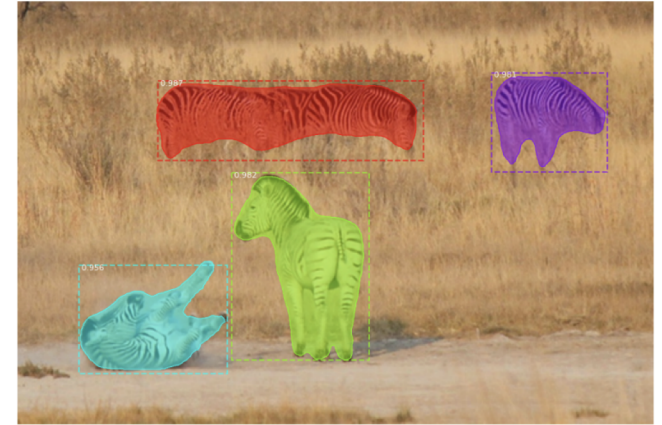
What is sitting in the window?
bird          clock

**VQA 2.0**: Goyal et al. 2017

# Object Detection (2D)



(a) Query: "street lamp"

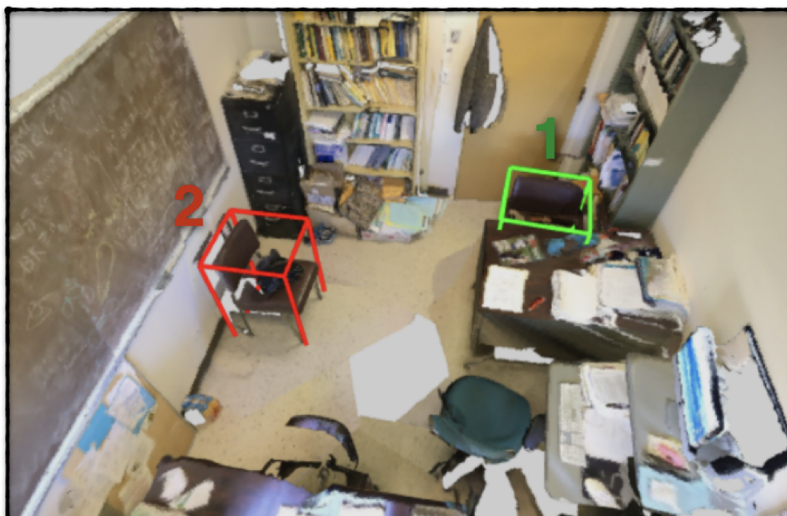(b) Query: "major league logo"

(c) Query: "zebras on savanna"

**MDETR**: Kamath et al. 2021

# Object Detection (3D)



1. "The chair closest to the door."
2. "The chair under the chalkboard."

1. "The office chair that is green."
2. "Choose the brown office chair pushed under the desk."

**ReferIt3D**: Achlioptas et al. 2020

# Vision and Language Navigation



"Place a clean ladle on a counter"

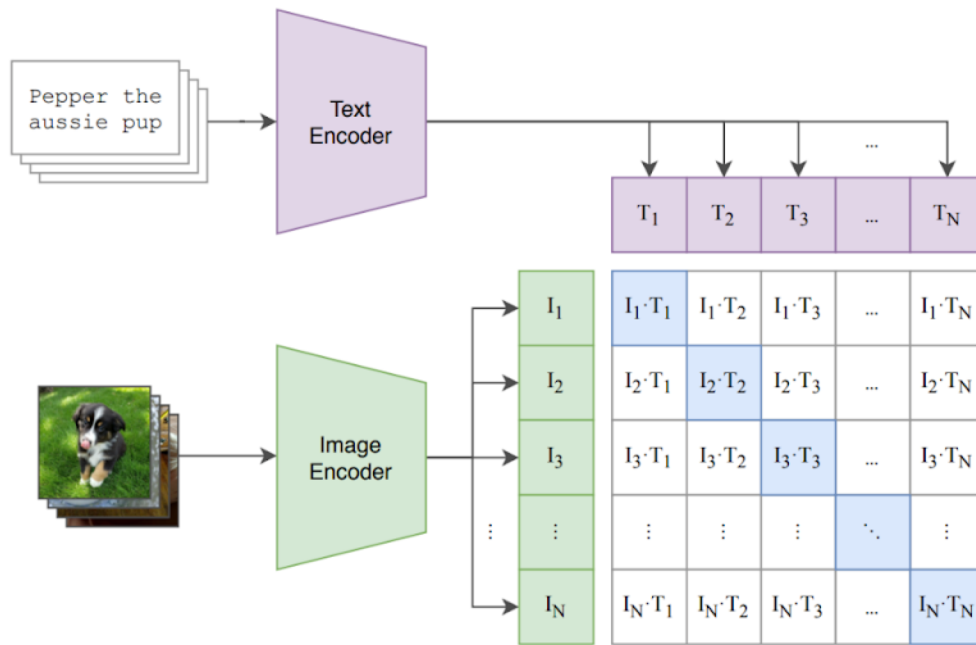**ALFRED**: Shridhar et al. 2020

# CLIP

# CLIP

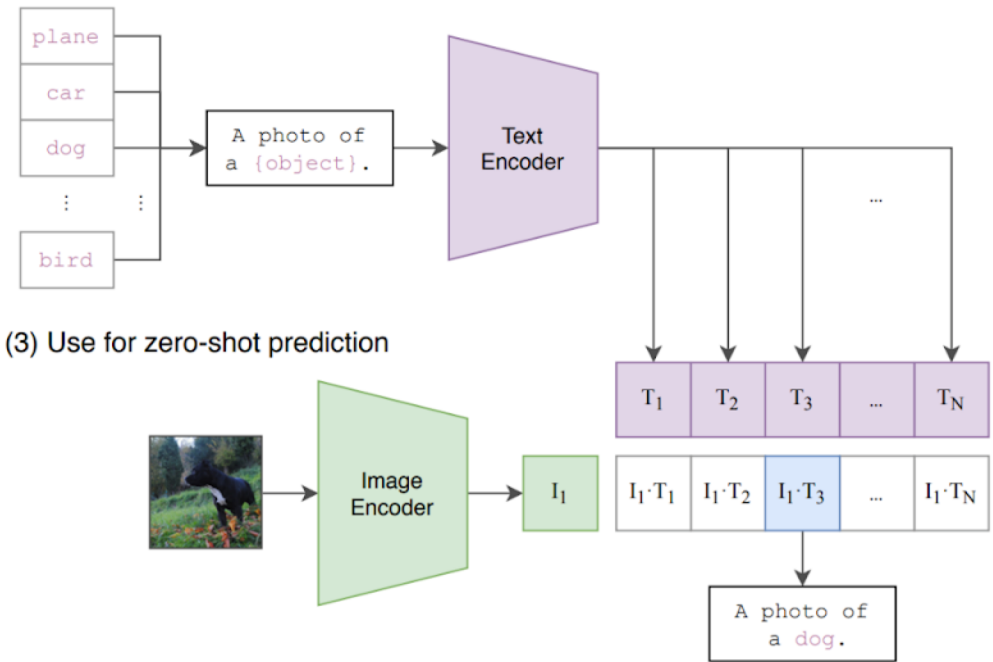(an encoder for putting images and text into the same embedding space)

# Embedding Images and Language



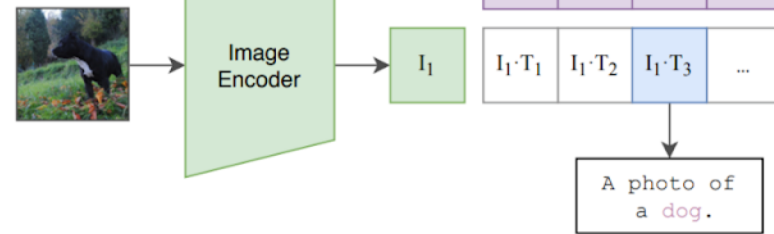(1) Contrastive pre-training

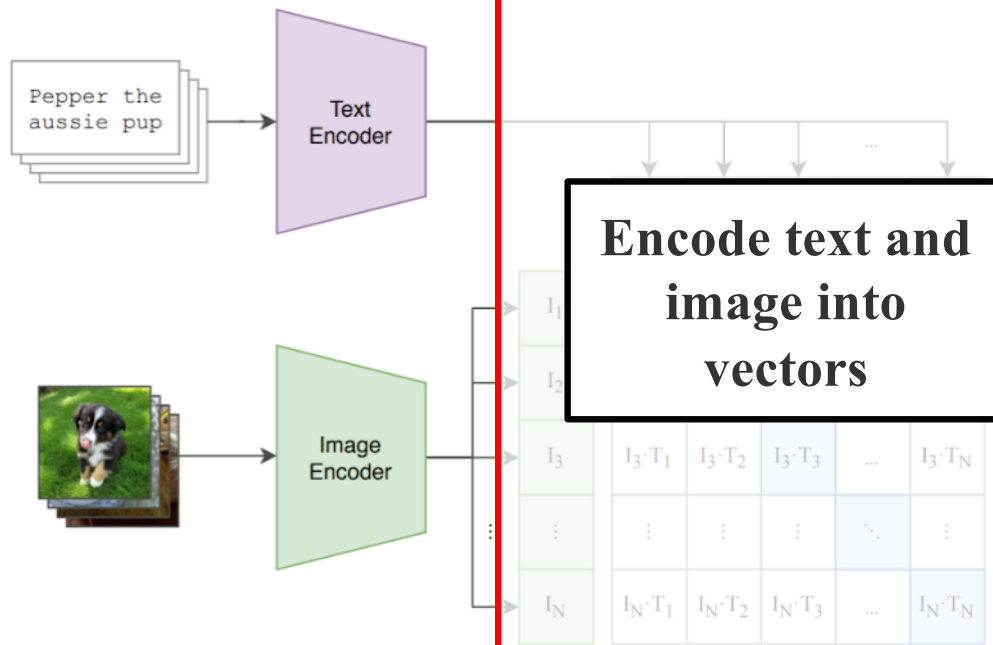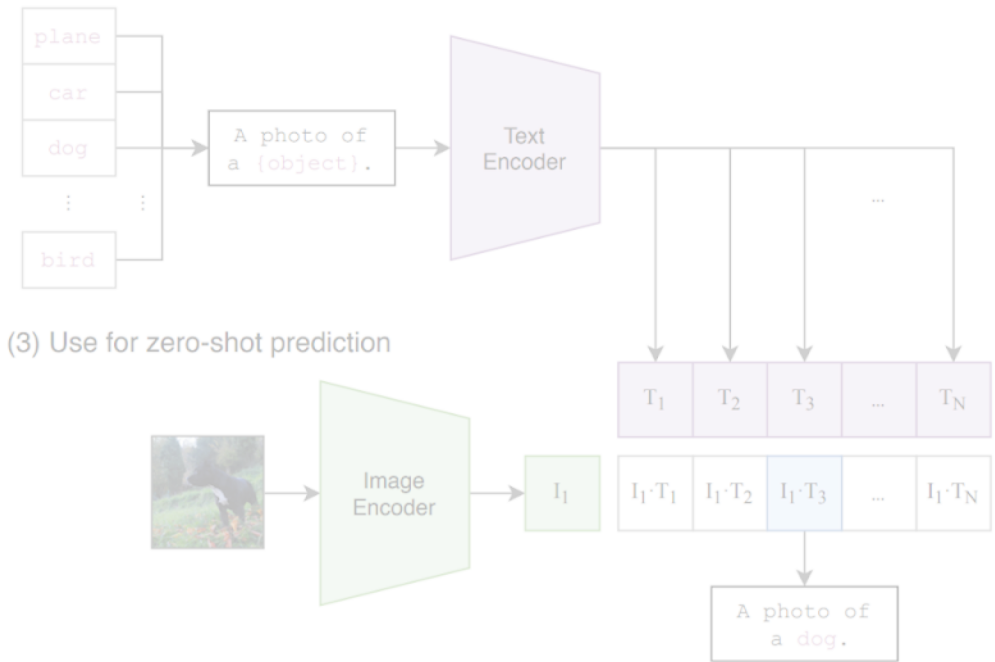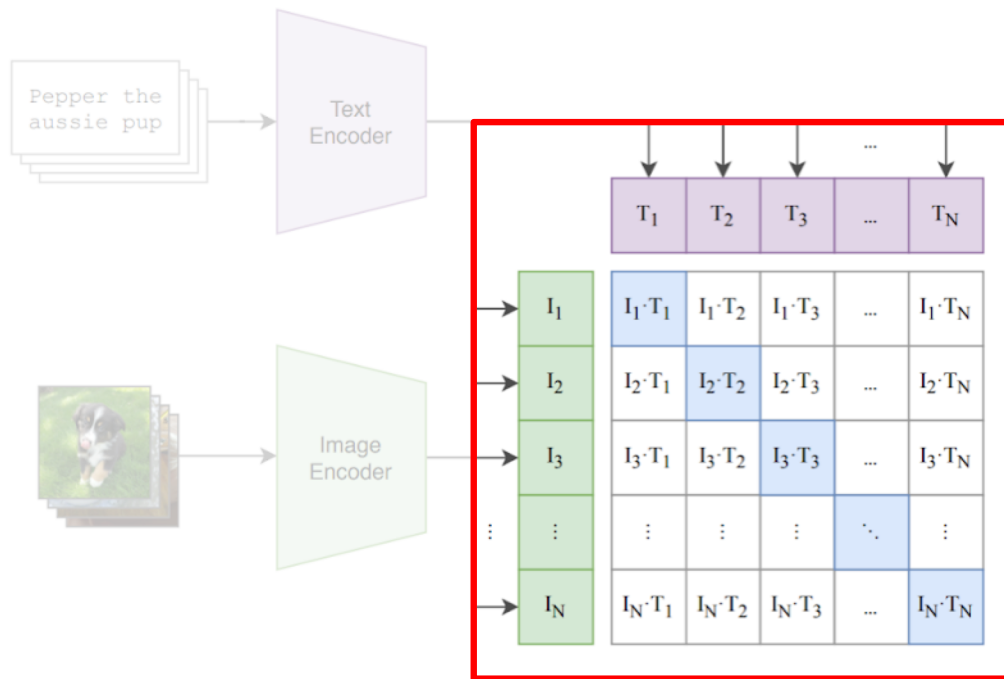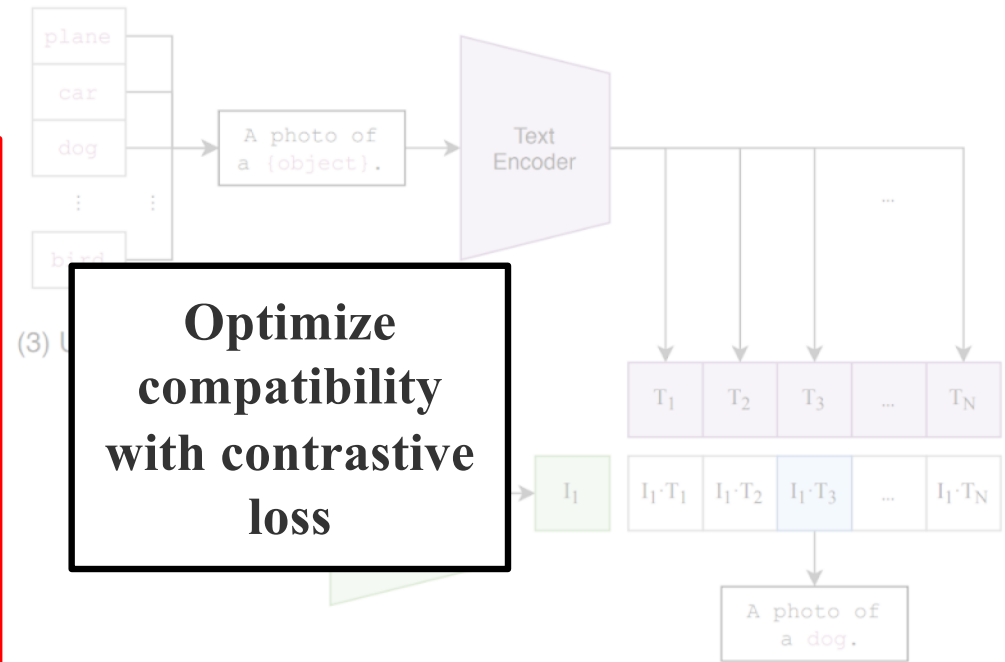(2) Create dataset classifier from label text

(3) Use for zero-shot prediction

CLIP: Radford et al. 2021

# Embedding Images and Language



(1) Contrastive pre-training

Pepper the aussie pup → Text Encoder

Image Encoder

$I_3 \cdot T_1$, $I_3 \cdot T_2$, $I_3 \cdot T_3$, ..., $I_3 \cdot T_N$

$I_N \cdot T_1$, $I_N \cdot T_2$, $I_N \cdot T_3$, ..., $I_N \cdot T_N$

**Encode text and image into vectors**

(2) Create dataset classifier from label text

plane, car, dog, bird → A photo of a {object}. → Text Encoder

(3) Use for zero-shot prediction

Image Encoder → $I_1$

$I_1 \cdot T_1$, $I_1 \cdot T_2$, $I_1 \cdot T_3$, ..., $I_1 \cdot T_N$

A photo of a dog.

**CLIP**: Radford et al. 2021

# Embedding Images and Language



(1) Contrastive pre-training

Pepper the aussie pup → Text Encoder

Image Encoder

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

(2) Create dataset classifier from label text

plane
car
dog
⋮
bird

A photo of a {object}. → Text Encoder

**Optimize compatibility with contrastive loss**

(3) U...

|       | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|-------|-------|-------|-------|-----|-------|
| $I_1$ | $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |

A photo of a dog.

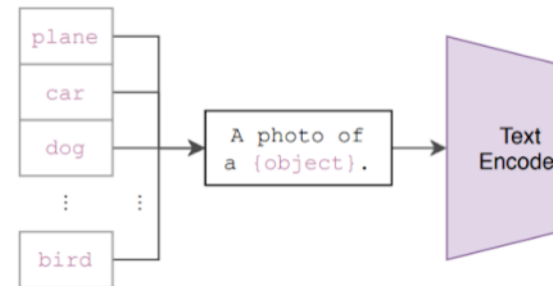**CLIP**: Radford et al. 2021

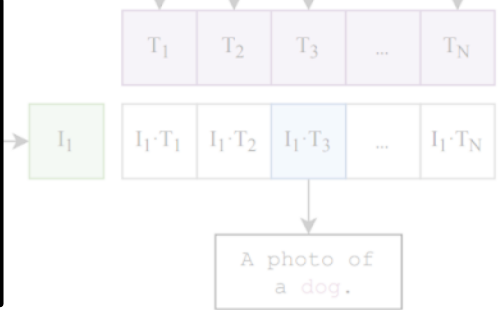# Embedding Images and Language



(1) Contrastive pre-training

(2) Create dataset classifier from label text

(3) Use for zero-shot prediction

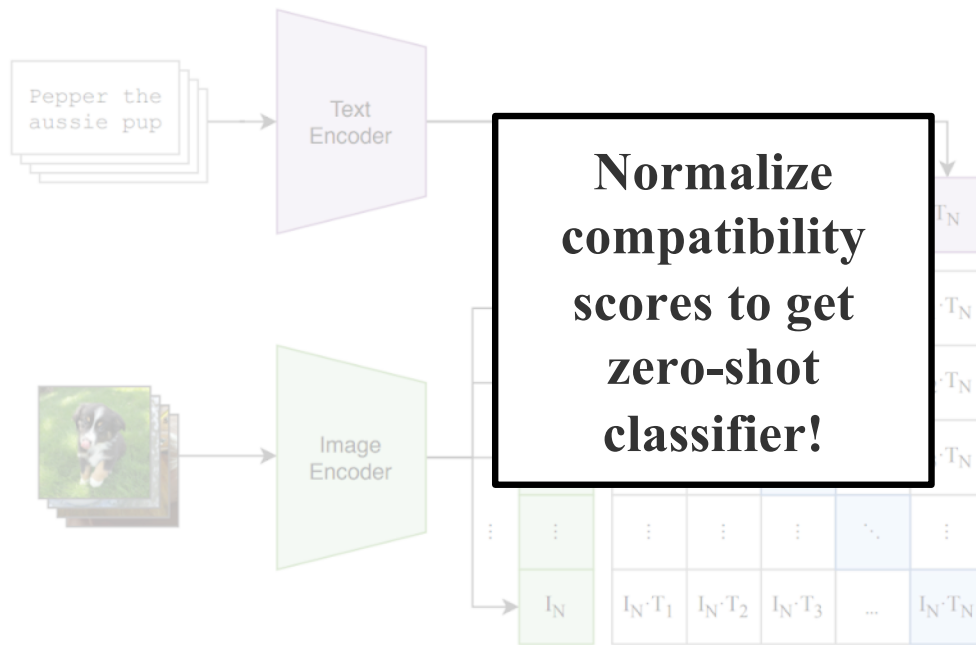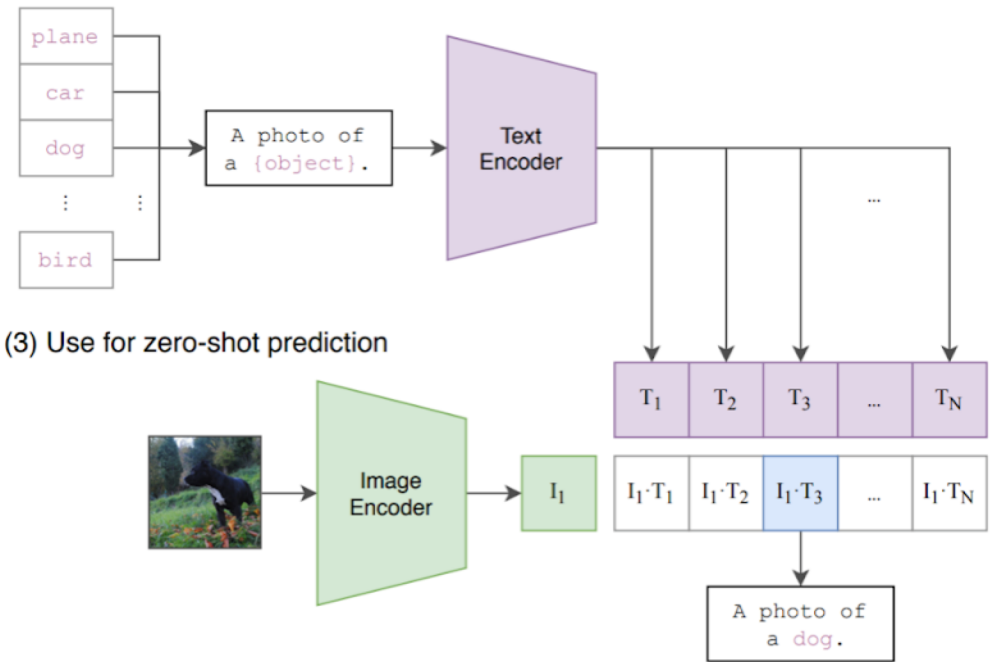**Classification dataset created with templated prompts**
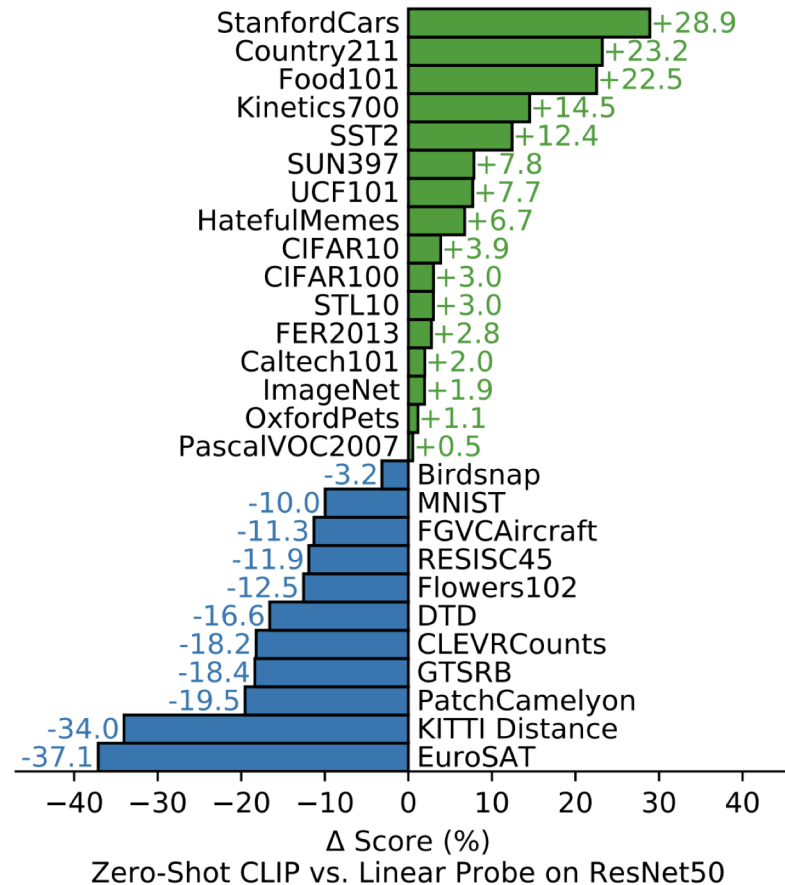
# Embedding Images and Language



**Normalize compatibility scores to get zero-shot classifier!**

**CLIP**: Radford et al. 2021

# Embedding Images and Language



Δ Score (%)
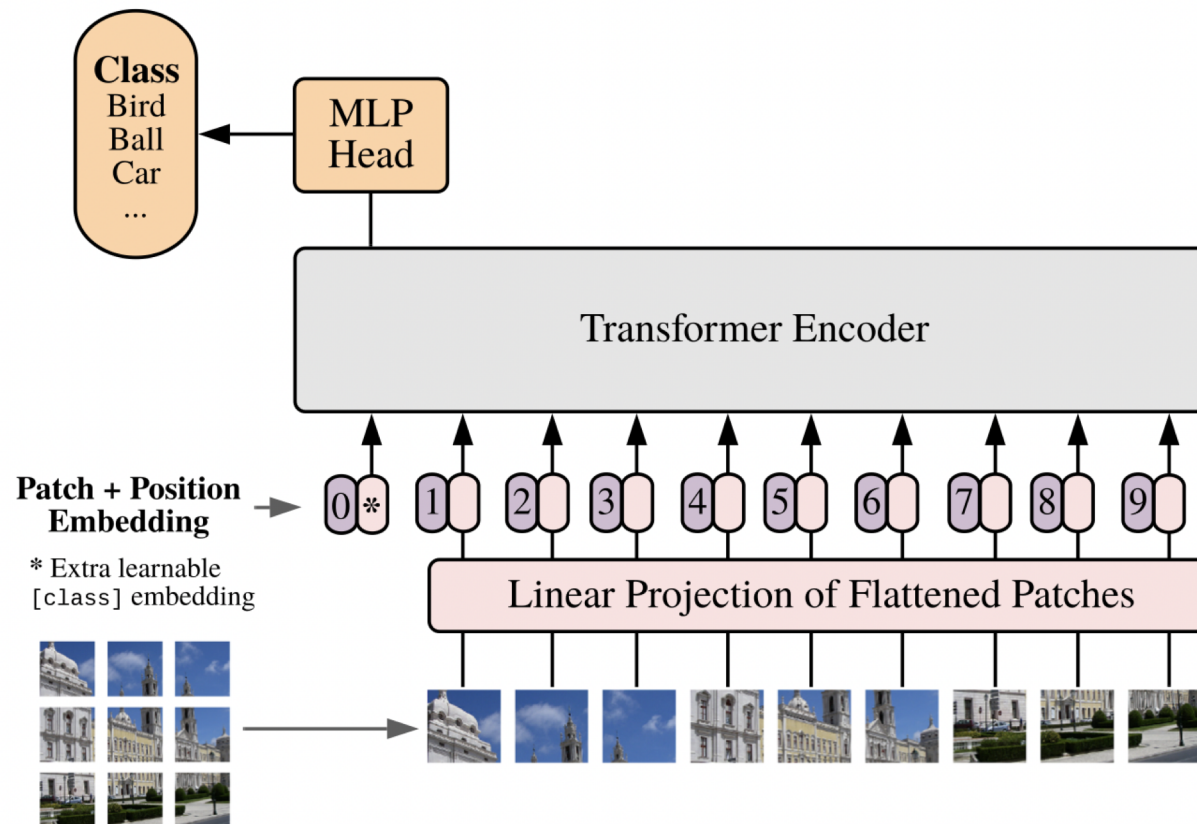Zero-Shot CLIP vs. Linear Probe on ResNet50
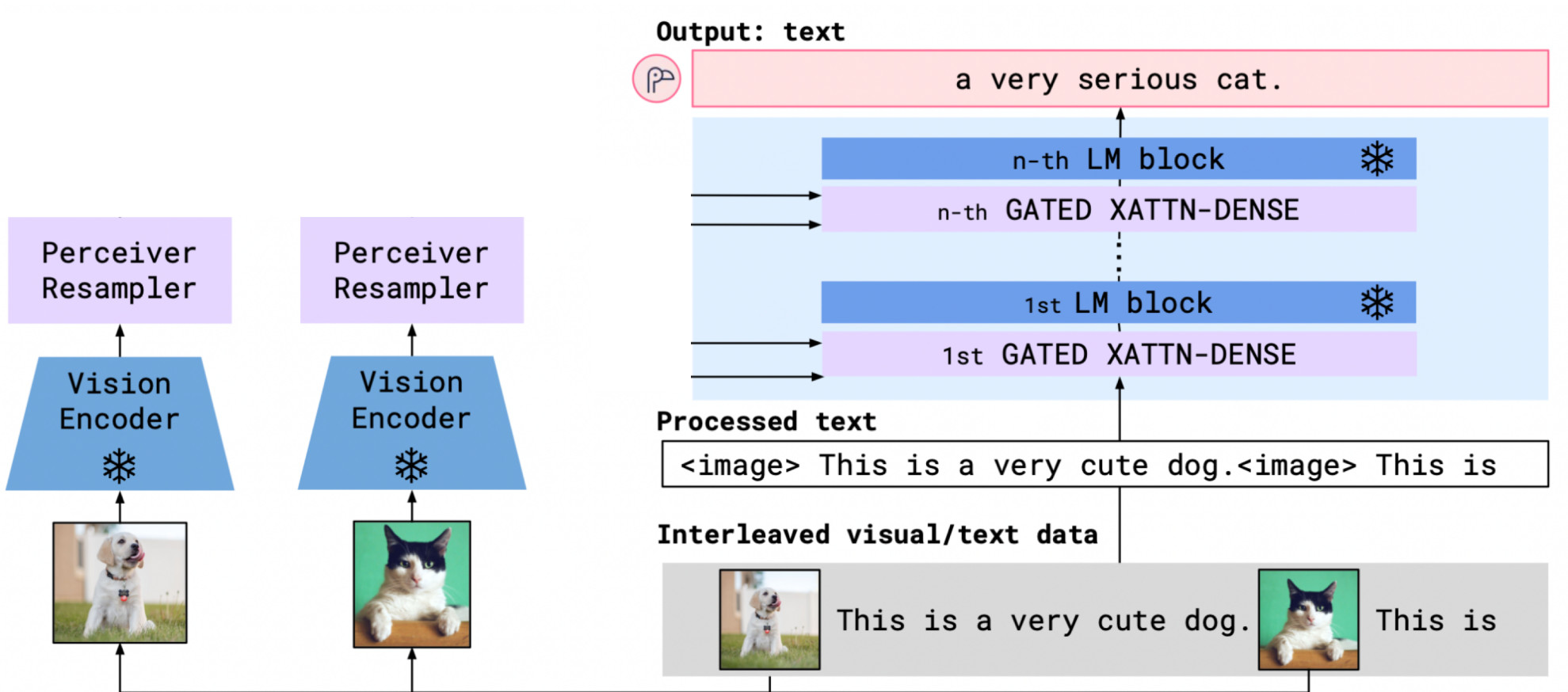
**CLIP**: Radford et al. 2021

# Joint Vision-Language Models

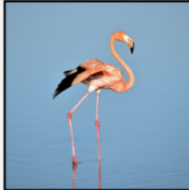# Vision Transformers with Patches



Vision Transformer (ViT)

# Flamingo 🦩

**Output: text**

a very serious cat.

n-th LM block ❄

n-th GATED XATTN-DENSE

⋮

1st LM block ❄

1st GATED XATTN-DENSE

**Processed text**

`<image>` This is a very cute dog.`<image>` This is

**Interleaved visual/text data**

This is a very cute dog. This is

Perceiver Resampler

Perceiver Resampler

Vision Encoder ❄

Vision Encoder ❄

# Flamingo 🦩

| Input Prompt | Completion |
|---|---|

**Input Prompt row 1:**

 This is a chinchilla. They are mainly found in Chile.  This is a shiba. They are very popular in Japan.  This is

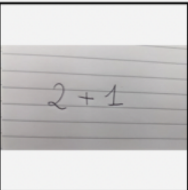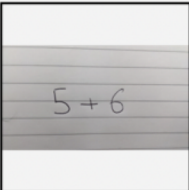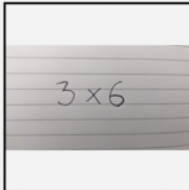→ **a flamingo. They are found in the Caribbean and South America.**

**Input Prompt row 2:**

 What is the title of this painting? Answer: The Hallucinogenic Toreador.  Where is this painting displayed? Answer: Louvres Museum, Paris.  What is the name of the city where this was painted? Answer:

→ **Arles.**

**Input Prompt row 3:**

 Output: "Underground"  Output: "Congress"  Output:

→ **"Soulomes"**

**Input Prompt row 4:**

 2+1=3  5+6=11  

→ **3x6=18**

Output: ▲

# Flamingo 🦩



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?
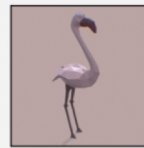
It looks like a computer.

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.

---



What is the common thing about these three images?

They are all flamingos.

What is the difference between these three images?

The first one is a cartoon, the second one is a real flamingo, and the third one is a 3D model of a flamingo.

---



This is an apple with a sticker on it.

What does the sticker say?

The sticker says "iPod".

Where is the photo taken?

It looks like it's taken in a backyard.

Do you think it is printed or handwritten?

It looks like it's handwritten.

What color is the sticker?

It's white.

---



This is a cityscape. It looks like Chicago.

What makes you think this is Chicago?

I think it's Chicago because of the Shedd Aquarium in the background.



What about this one? Which city is this and what famous landmark helped you recognise the city?

This is Tokyo. I think it's Tokyo because of the Tokyo Tower.

# Vector-Quantized Vision-Language

| 1 | 5 | 2 | 6 |
|---|---|---|---|
| 9 | 13 | 10 | 14 |
| 3 | 7 | 4 | 8 |
| 11 | 15 | 12 | 16 |

**DALL-E 1**: Ramesh et al. 2021

# Vector-Quantized Vision-Language



Neural Discrete Representation Learning: van Oord et al. 2017

# Vector-Quantized Vision-Language

## Step 2

Learn Joint
Language and Code Distribution

"A kitten
with a pink
background"

| 1 | 5 | 2 | 6 |
|----|----|----|----|
| 9 | 13 | 10 | 14 |
| 3 | 7 | 4 | 8 |
| 11 | 15 | 12 | 16 |

# Vector-Quantized Vision-Language

## Step 2

Learn Joint
Language and Code Distribution

"A kitten
with a pink
background"

| 1 | 5 | 2 | 6 |
| 9 | 13 | 10 | 14 |
| 3 | 7 | 4 | 8 |
| 11 | 15 | 12 | 16 |



**Generating Long Sequences with Sparse Transformers**: Child et al. 2019

Reduced to language modeling
problem!

**DALL-E 1**: Ramesh et al. 2021