# vision & language

CS 685, Spring 2022
Advanced Natural Language Processing
http://people.cs.umass.edu/~miyyer/cs685/

## Mohit Iyyer

College of Information and Computer Sciences
University of Massachusetts Amherst

# image captioning



a red truck is parked on a street lined with trees

# visual question answering



- Is this truck considered "vintage"?
- Does the road look new?
- What kind of tree is behind the truck?

we've seen how to compute representations of words and sentences. what about images?
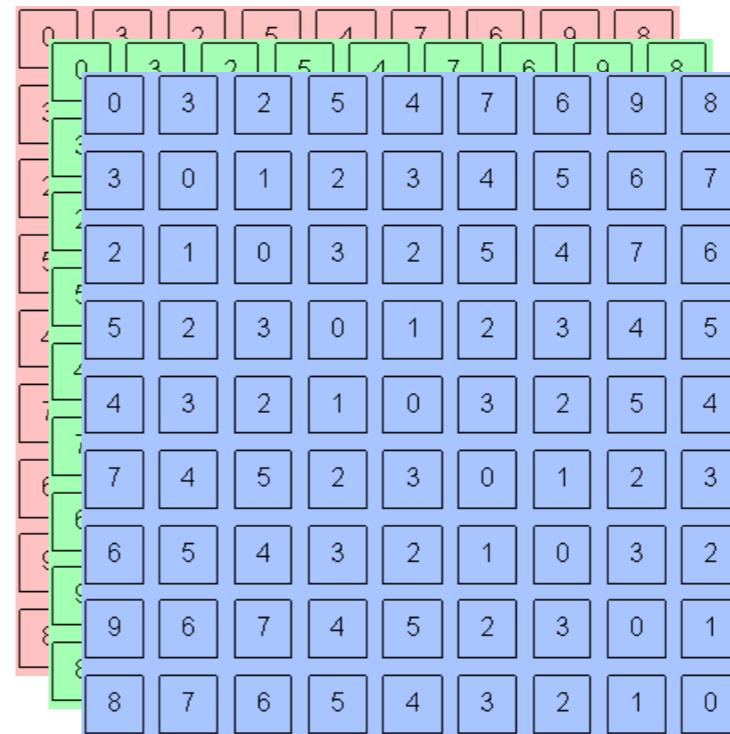
# grayscale images are matrices



La Gare Montparnasse, 1895

| 0 | 3 | 2 | 5 | 4 | 7 | 6 | 9 | 8 |
|---|---|---|---|---|---|---|---|---|
| 3 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 2 | 1 | 0 | 3 | 2 | 5 | 4 | 7 | 6 |
| 5 | 2 | 3 | 0 | 1 | 2 | 3 | 4 | 5 |
| 4 | 3 | 2 | 1 | 0 | 3 | 2 | 5 | 4 |
| 7 | 4 | 5 | 2 | 3 | 0 | 1 | 2 | 3 |
| 6 | 5 | 4 | 3 | 2 | 1 | 0 | 3 | 2 |
| 9 | 6 | 7 | 4 | 5 | 2 | 3 | 0 | 1 |
| 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |

what range of values can each pixel take?
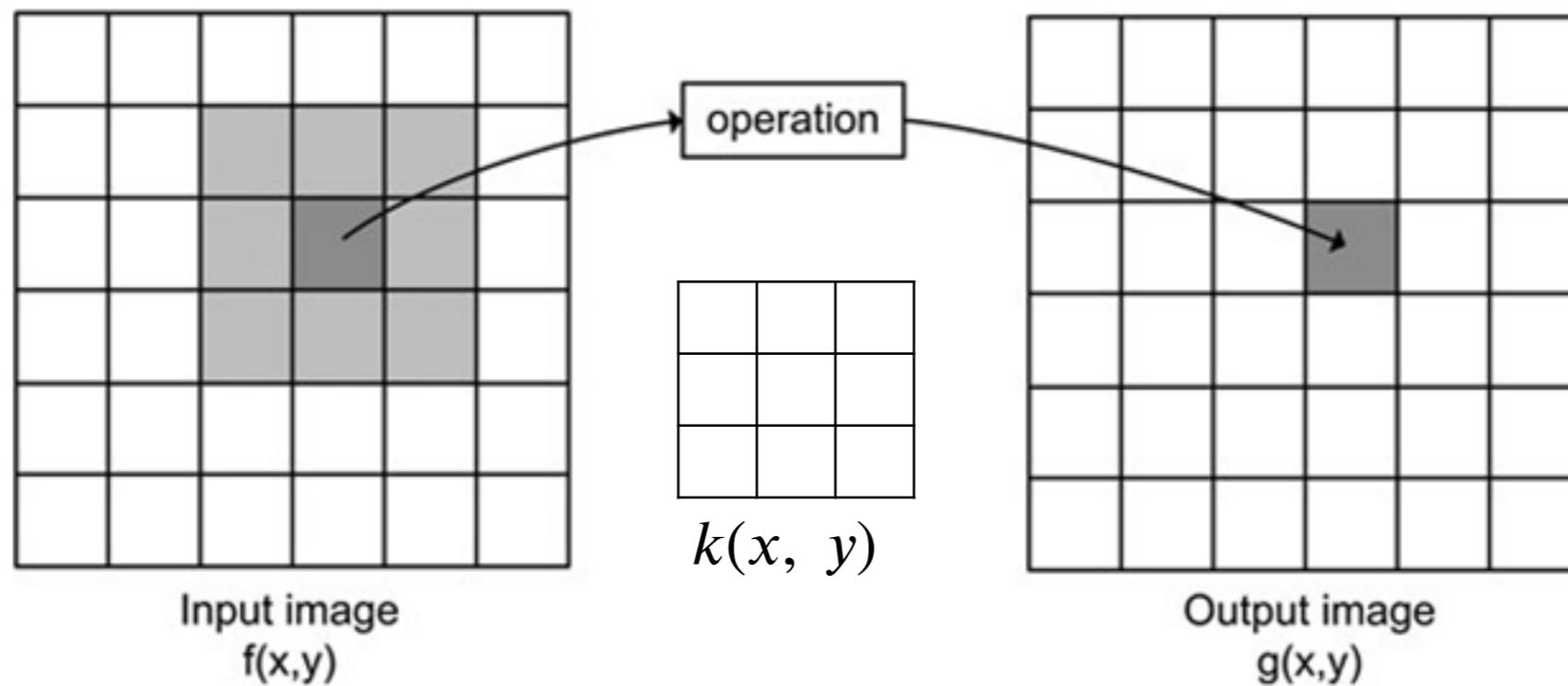
# color images are tensors



*channel x height x width*

Channels are usually RGB: Red, Green, and Blue
Other color spaces: HSV, HSL, LUV, XYZ, Lab, CMYK, etc

# Convolution operator



Input image
f(x,y)

$k(x, \ y)$

Output image
g(x,y)

$$g(x, \ y) = \sum_{v} \sum_{u} k(u, v) f(x \ - u, \ y - v)$$

(filter, kernel)

Input image  *  Weights  →  Output image

| 4 | 5 | 7 | 6 | 6 |
| 3 | 2 | 8 | 0 | 7 |
| 6 | 7 | 7 | 1 | 5 |
| 3 | 0 | 1 | 1 | 1 |
| 4 | 3 | 2 | 1 | 7 |

*

| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 0 | 0 |

| | 11 | 2 | 15 | |
| | 13 | 8 | 12 | |
| | ? | | | |

# demo:
http://setosa.io/ev/image-kernels/

# Convolutional Layer (with 4 filters)

weights:
4x1x9x9

Input: 1x224x224

Output: 4x224x224

if zero padding,
and stride = 1

Convolution

# Convolutional Layer (with 4 filters)

weights:
4x1x9x9

Input: 1x224x224
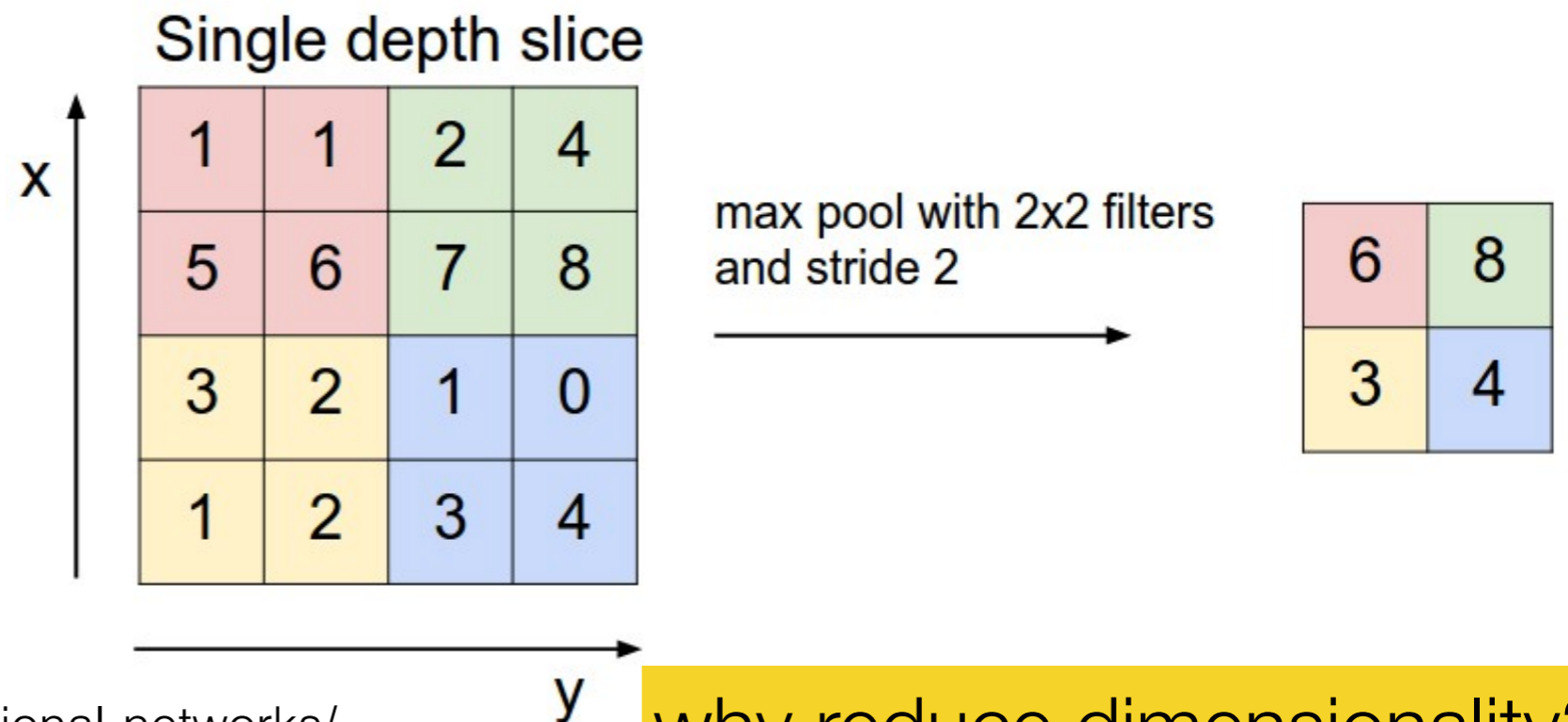
Output: 4x112x112

if zero padding,
but stride = 2



Convolution

# pooling layers also used to reduce dimensionality

*Convolutional Layers:*
slide a set of small filters over the image



32
32
3

*Pooling Layers:*
reduce dimensionality of representation

Single depth slice

x

| 1 | 1 | 2 | 4 |
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

max pool with 2x2 filters and stride 2

| 6 | 8 |
| 3 | 4 |

y

why reduce dimensionality?

# Alexnet

## ImageNet Classification with Deep Convolutional Neural Networks

**Alex Krizhevsky**
University of Toronto
kriz@cs.utoronto.ca

**Ilya Sutskever**
University of Toronto
ilya@cs.utoronto.ca

**Geoffrey E. Hinton**
University of Toronto
hinton@cs.utoronto.ca

the paper that started the deep learning revolution!

# image classification

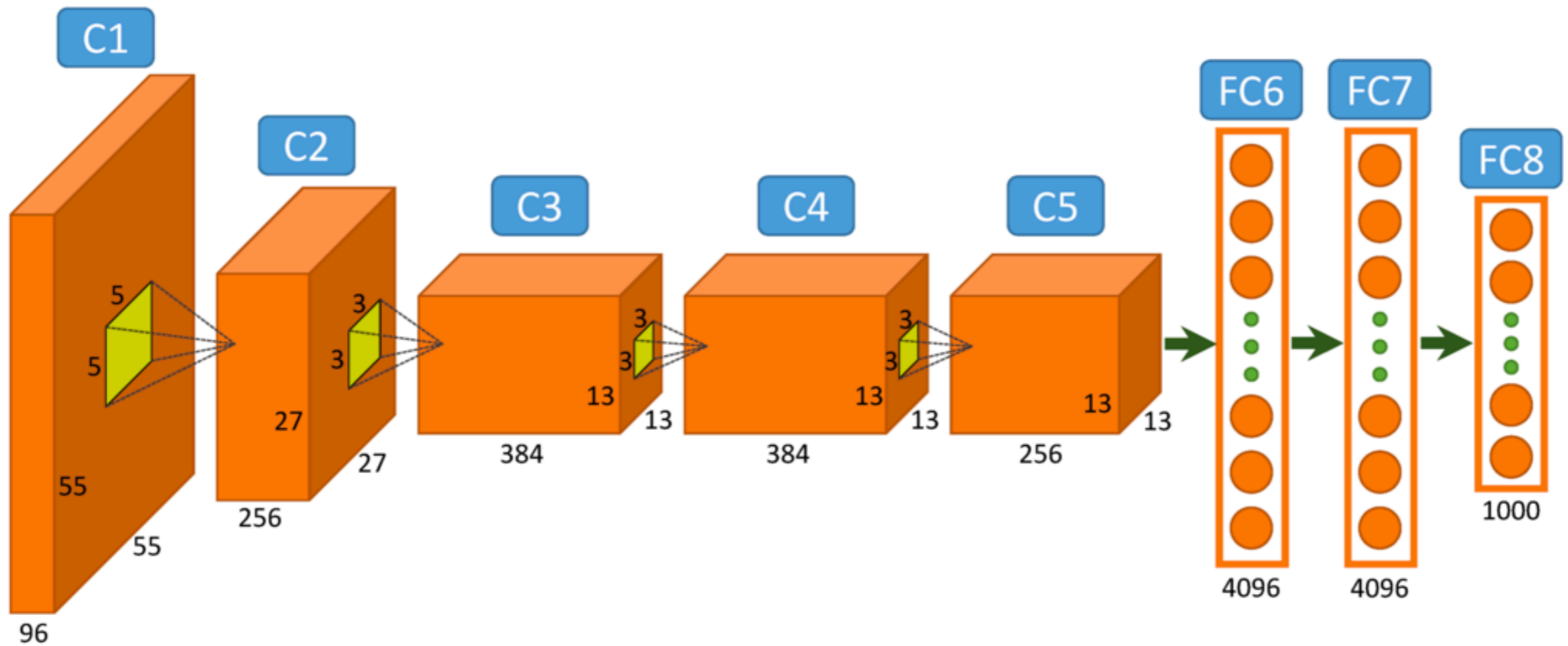Classify an image into 1000 possible classes:
e.g. Abyssinian cat, Bulldog, French Terrier, Cormorant, Chickadee,
red fox, banjo, barbell, hourglass, knot, maze, viaduct, etc.



cat, tabby cat  (0.71)
Egyptian cat (0.22)
red fox (0.11)
…..

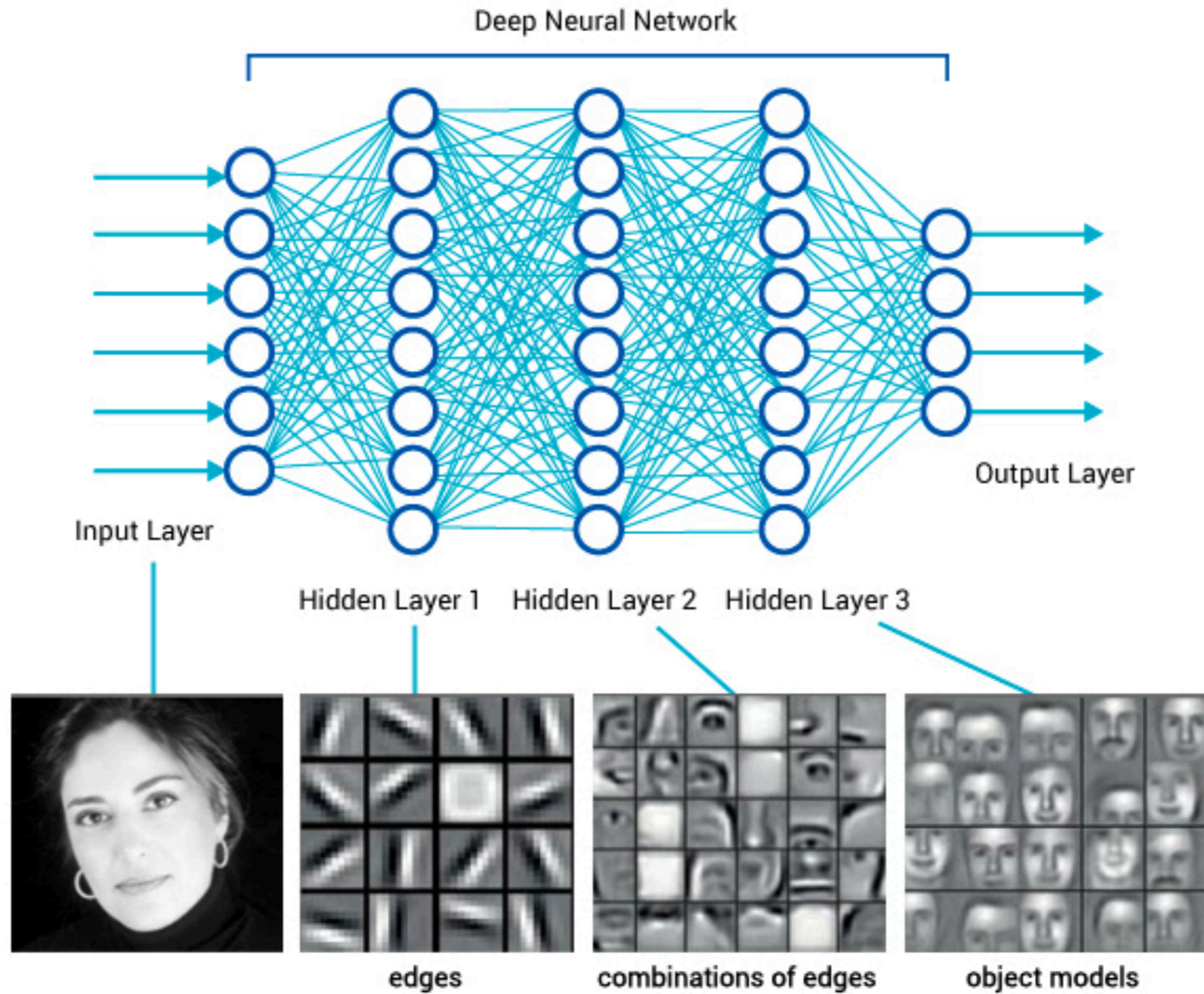train on the ImageNet challenge dataset, ~1.2 million images

# Alexnet

# Alexnet



conv+pool

conv+pool

conv   conv   conv

linear   linear

C1   C2   C3   C4   C5   FC6   FC7   FC8

linear+ softmax

https://www.saagie.com/fr/blog/object-detection-part1

# What is happening?



Deep Neural Network

Input Layer

Hidden Layer 1   Hidden Layer 2   Hidden Layer 3

Output Layer

edges    combinations of edges    object models

https://www.saagie.com/fr/blog/object-detection-part1

# Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)

VGG, 19 layers
(ILSVRC 2014)

ResNet, 152 layers
(ILSVRC 2015)

152 layers

22 layers

19 layers

11.7

16.4

25.8

28.2

8 layers

8 layers

shallow

3.57

6.7

7.3

ILSVRC'15
ResNet

ILSVRC'14
GoogleNet

ILSVRC'14
VGG

ILSVRC'13

ILSVRC'12
AlexNet

ILSVRC'11

ILSVRC'10

Slide by Mohammad Rastegari

# ImageNet pretraining -> Instagram pretraining

Bigger models are saturated
on ImageNet, but with more
data bigger models do better



Biggest network was pretrained on
3.5B Instagram images

Trained on 336 GPUs for 22 days

Mahajan et al, "Exploring the Limits of Weakly Supervised Pretraining", arXiv 2018

at the end of the day, we generate a fixed size vector from an image and run a classifier over it

*CNN* $\Bigg($  $\Bigg)$ = 

softmax: predict 'truck'

key insight: this vector is useful for many more tasks than just image classification! we can use it for *transfer learning*

*CNN*  =

# simple visual QA

- i = *CNN*(image) > use an existing network trained for image classification and freeze weights

- q = *RNN*(question) > learn weights

- answer = softmax(linear([i;q]))

why isn't this a good way of doing visual QA?

How many benches are shown?

# visual attention

- Use the question representation $q$ to determine where in the image to look



How many benches are shown?

softmax:
predict answer

attention over final convolutional
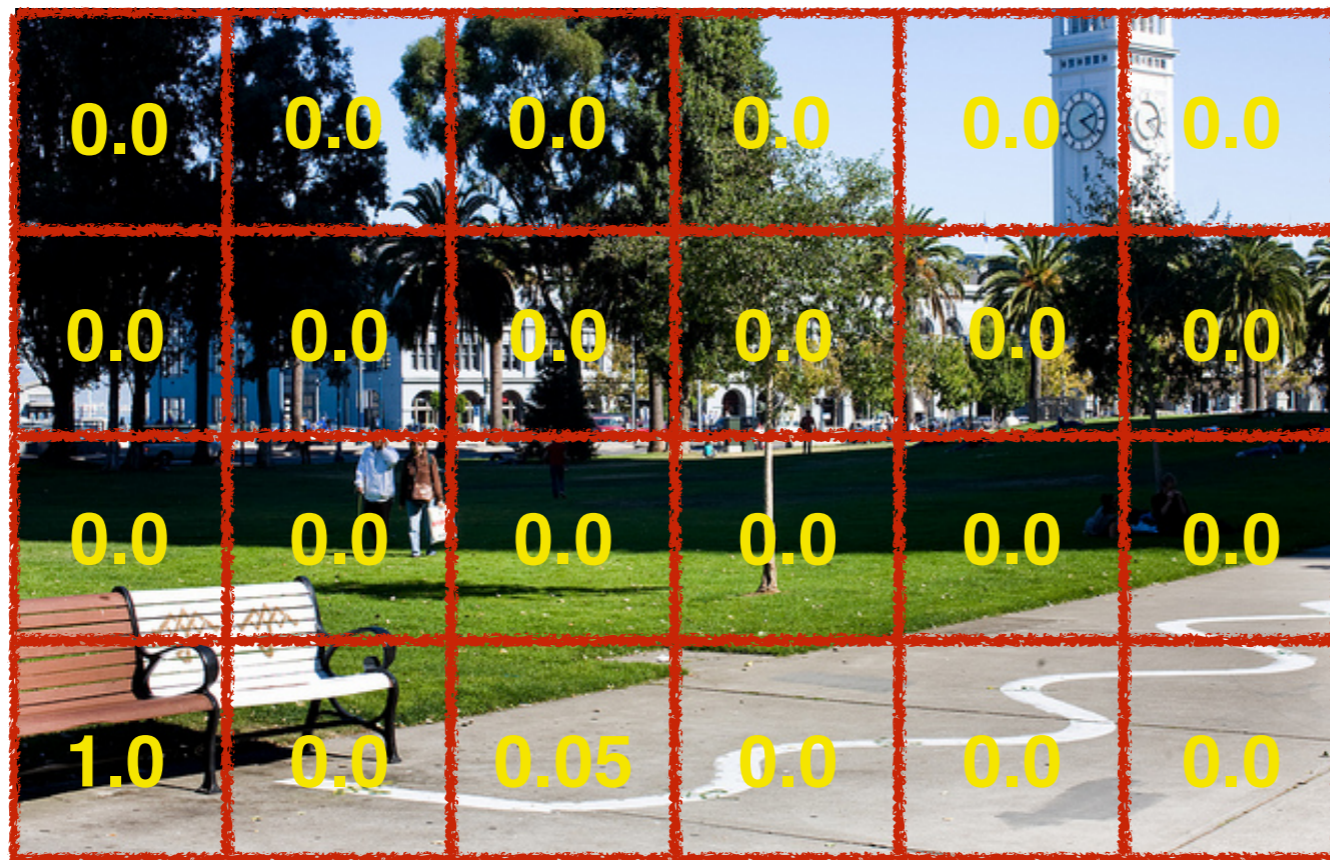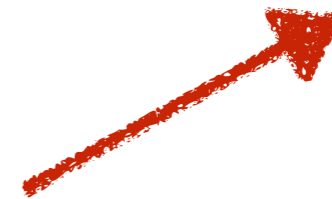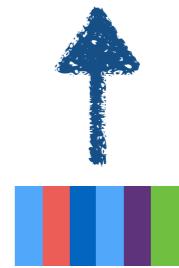layer in network: 196 boxes, captures
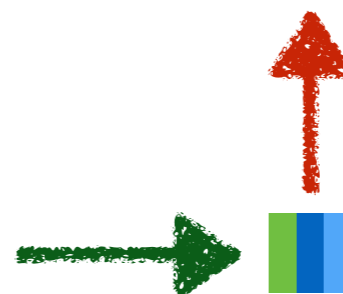color and positional information

| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 0.05 | 0.05 | 0.0 | 0.0 | 0.0 | 0.0 |
| 0.2 | 0.1 | 0.05 | 0.0 | 0.0 | 0.0 |
| 0.3 | 0.2 | 0.05 | 0.0 | 0.0 | 0.0 |

How many benches are shown?

softmax:
predict answer

attention over final convolutional
layer in network: 196 boxes, captures
color and positional information

| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 0.05 | 0.05 | 0.0 | 0.0 | 0.0 | 0.0 |
| 0.2 | 0.1 | 0.05 | 0.0 | 0.0 | 0.0 |
| 0.3 | 0.2 | 0.05 | 0.0 | 0.0 | 0.0 |

how can we
compute these
attention
scores?

How many benches are shown?

# hard attention

attention over final convolutional layer in network: 196 boxes, captures color and positional information

softmax: predict answer

we can use *reinforcement learning* to focus on just one box

| | | | | | |
|---|---|---|---|---|---|
| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 1.0 | 0.0 | 0.05 | 0.0 | 0.0 | 0.0 |

How many benches are shown?

# Grounded question answering

# Neural nets learn lexical groundings



Is there a red shape above a **circle**?

yes

[Iyyer et al. 2014, Bordes et al. 2014, Yang et al. 2015, Malinowski et al., 2015]

Slide credit: Jacob Andreas

# Semantic parsers learn composition

*Is there a red shape above a circle?*

yes

[Wong & Mooney 2007, Kwiatkowski et al. 2010, Liang et al. 2011, A et al. 2013]

Slide credit: Jacob Andreas

# Neural module networks learn both!

Is there a red shape above a circle?

yes

# Neural module networks
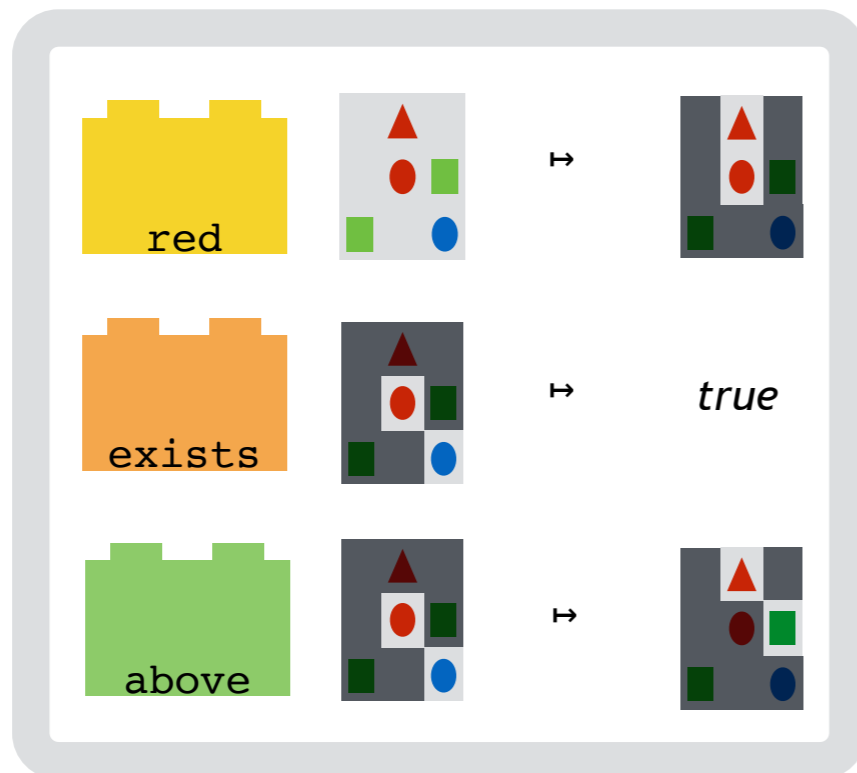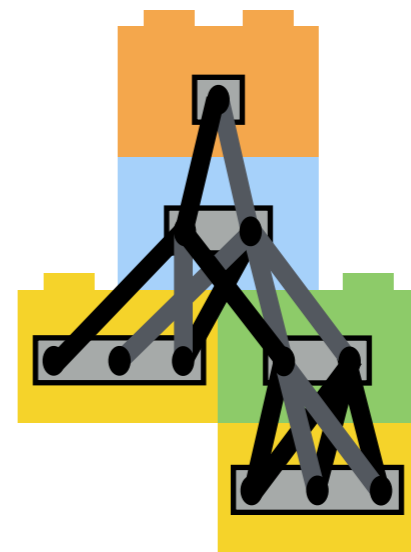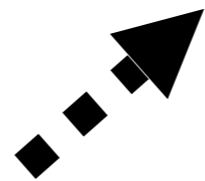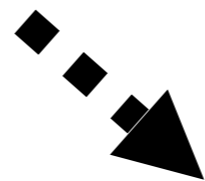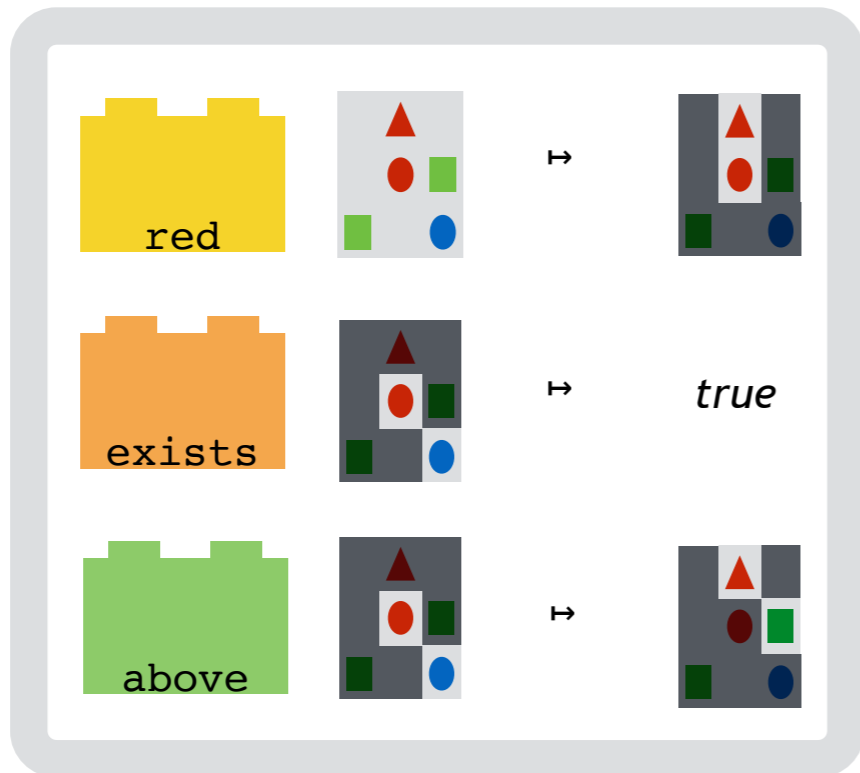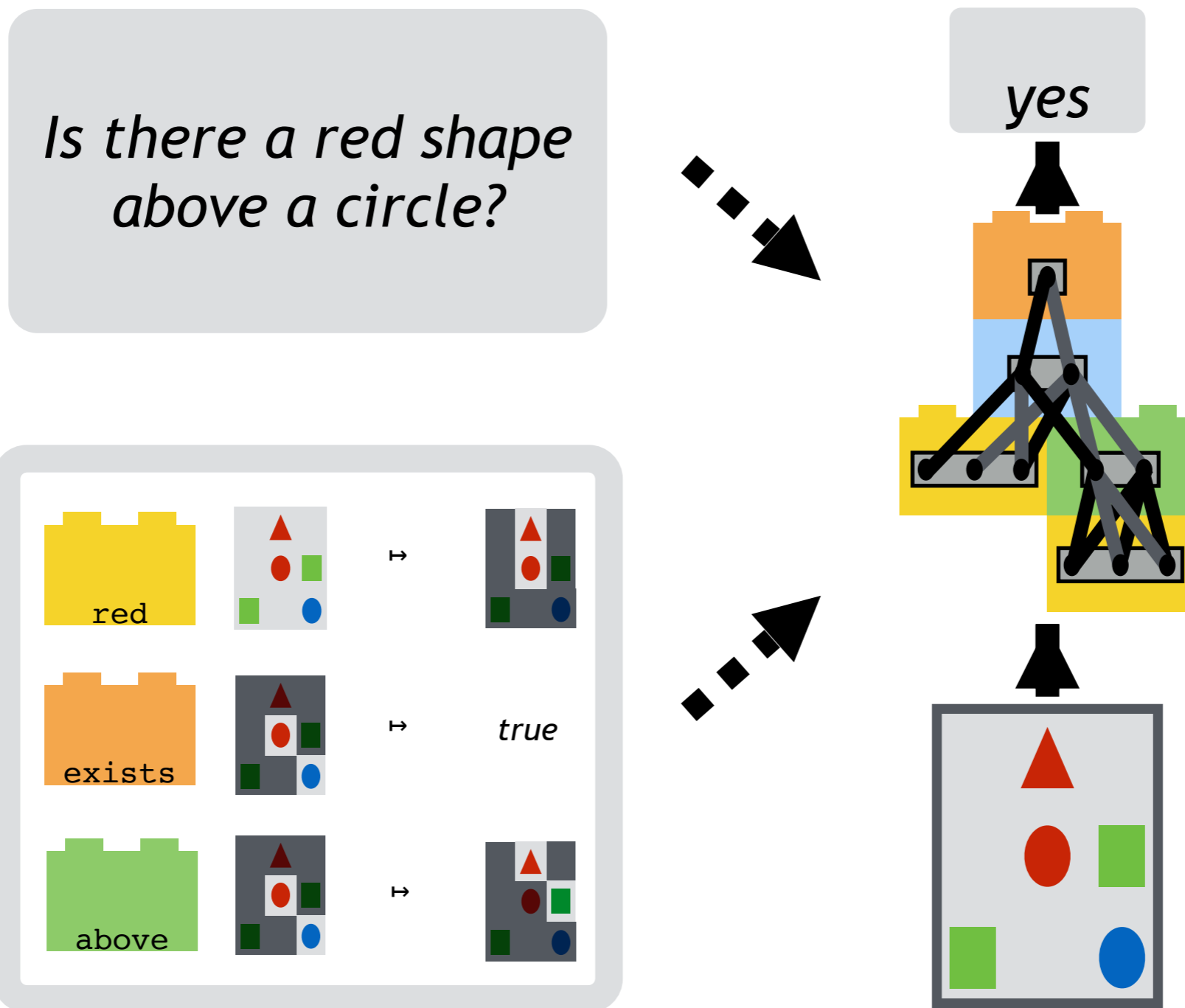
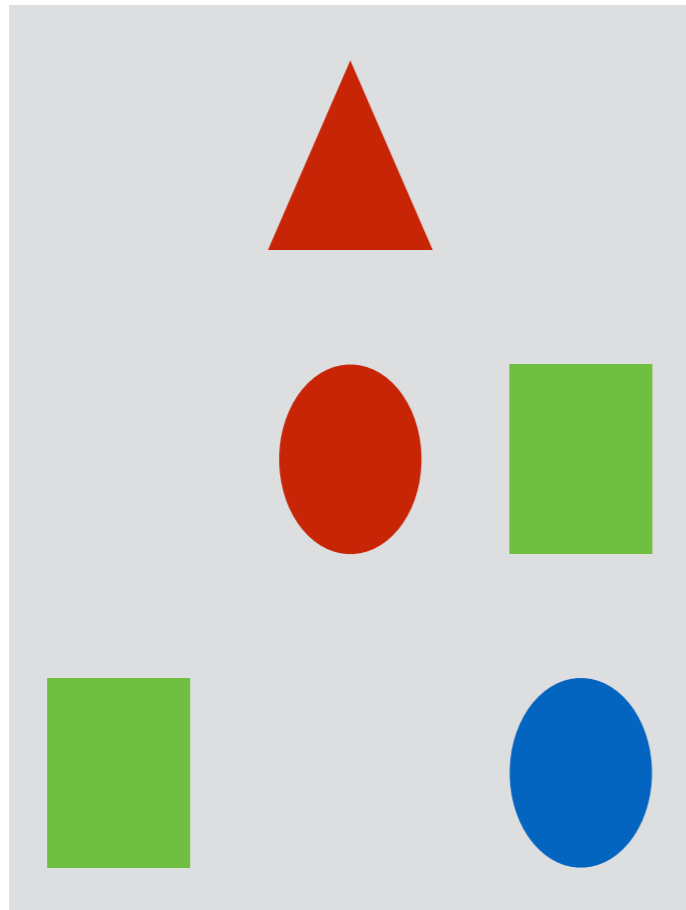*Is there a red shape above a circle?*

# Neural module networks



*Is there a red shape above a circle?*

Slide credit: Jacob Andreas

# Neural module networks

*Is there a red shape above a circle?*

yes

red

exists ↦ *true*

above

Slide credit: Jacob Andreas

# Sentence meanings are computations

# NLVR²: natural language for visual reasoning! (Suhr et al., 2018)



**TRUE OR FALSE:** the left image contains twice the number of dogs as the right image, and at least two dogs in total are standing.
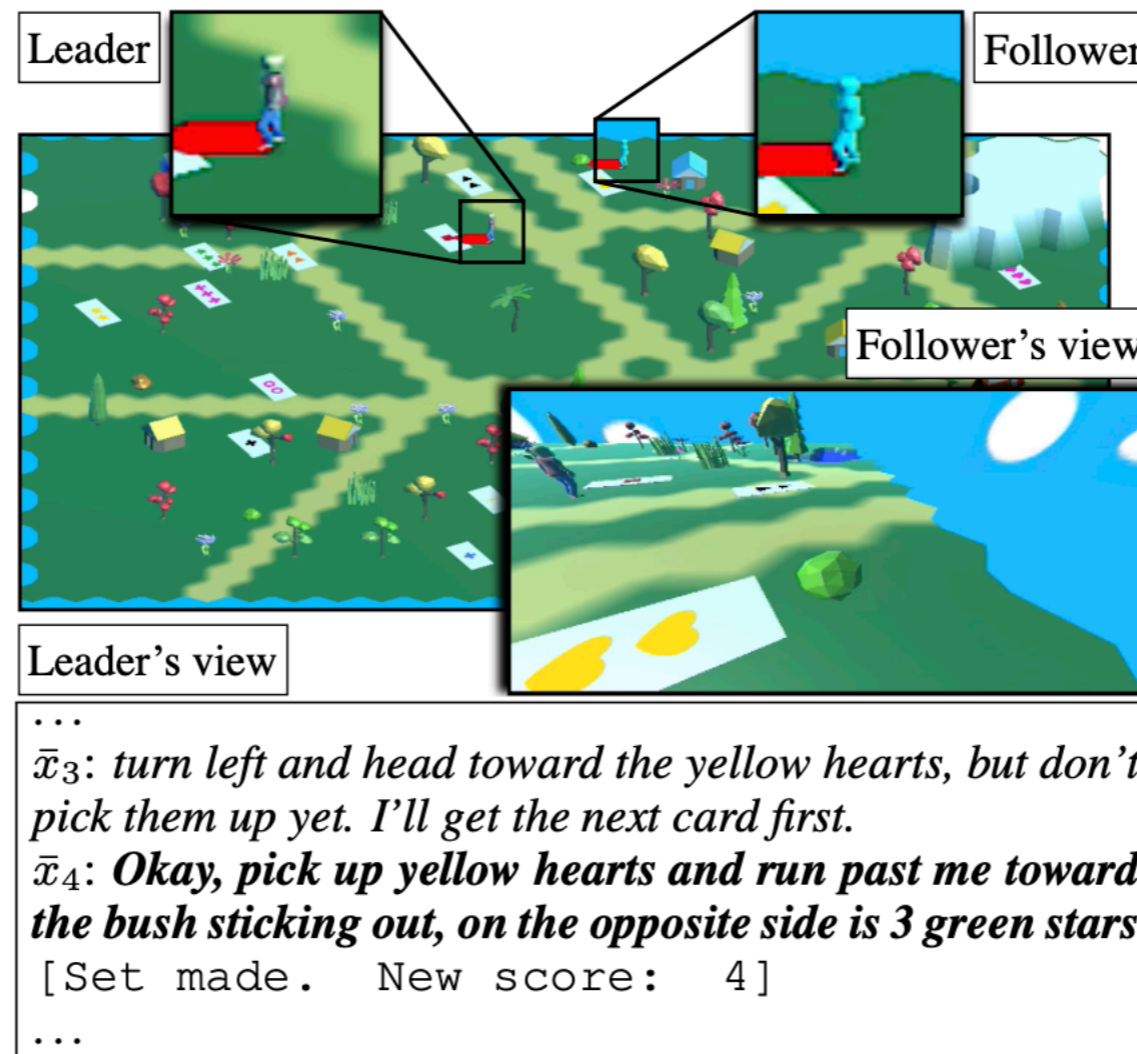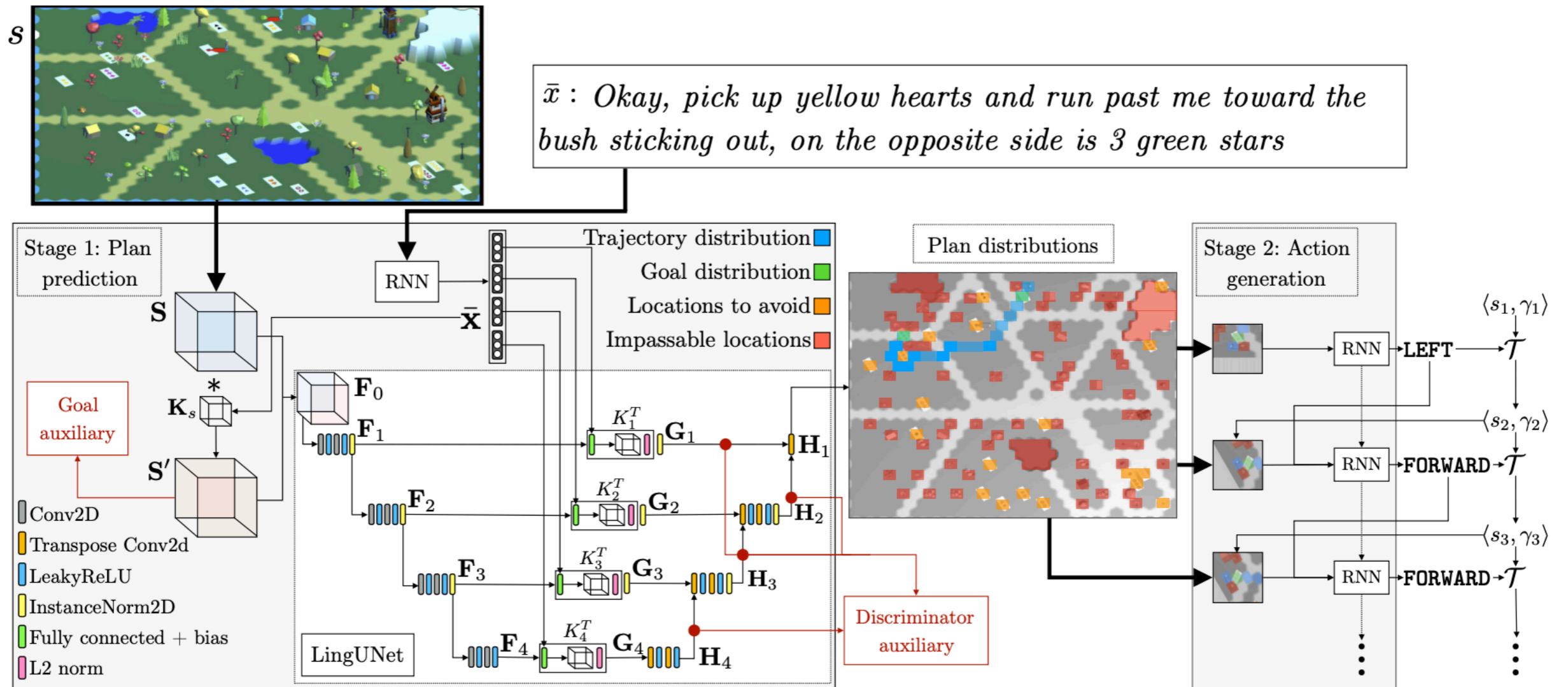
Figure 1: A snapshot from an interaction in CEREAL-BAR. The current instruction is in bold. The large image shows the entire environment. This overhead view is available only to the leader. The follower sees a first-person view only (bottom right). The zoom boxes (top) show the leader and follower.

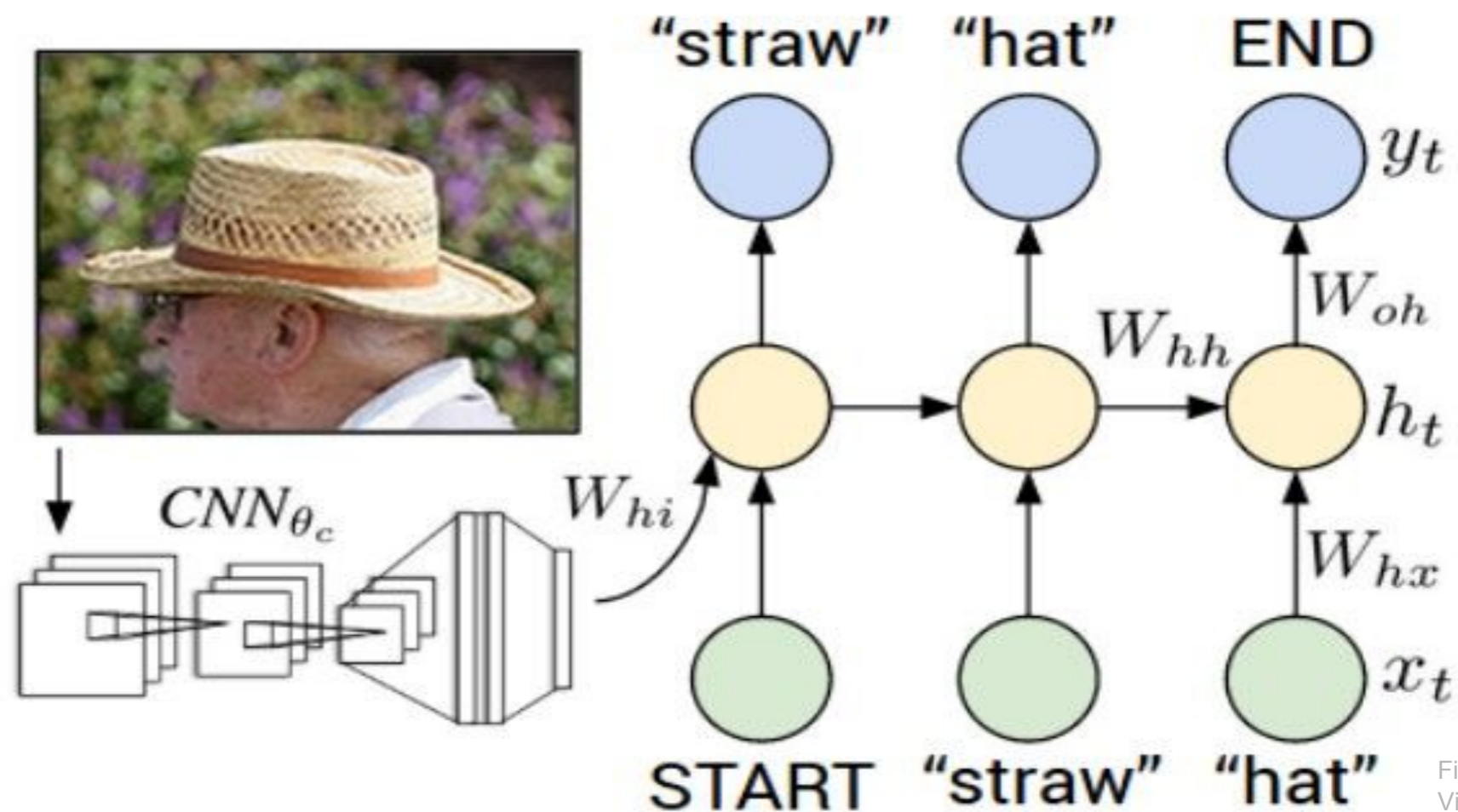$\bar{x}$ : *Okay, pick up yellow hearts and run past me toward the bush sticking out, on the opposite side is 3 green stars*

Stage 1: Plan prediction

RNN

Trajectory distribution
Goal distribution
Locations to avoid
Impassable locations

$\mathbf{S}$

$\bar{\mathbf{x}}$

Goal auxiliary

$\mathbf{K}_s$ *

$\mathbf{S}'$

$\mathbf{F}_0$
$\mathbf{F}_1$
$\mathbf{F}_2$
$\mathbf{F}_3$
$\mathbf{F}_4$

$K_1^T$ $\mathbf{G}_1$
$K_2^T$ $\mathbf{G}_2$
$K_3^T$ $\mathbf{G}_3$
$K_4^T$ $\mathbf{G}_4$

$\mathbf{H}_1$
$\mathbf{H}_2$
$\mathbf{H}_3$
$\mathbf{H}_4$

Conv2D
Transpose Conv2d
LeakyReLU
InstanceNorm2D
Fully connected + bias
L2 norm

LingUNet

Discriminator auxiliary

Plan distributions

Stage 2: Action generation

RNN **LEFT** $\rightarrow \mathcal{T}$
$\langle s_1, \gamma_1 \rangle$

RNN **FORWARD** $\rightarrow \mathcal{T}$
$\langle s_2, \gamma_2 \rangle$

RNN **FORWARD** $\rightarrow \mathcal{T}$
$\langle s_3, \gamma_3 \rangle$

Suhr et al., 2019 ("CEREALBAR")

# Image Captioning

Explain Images with Multimodal Recurrent Neural Networks, Mao et al.
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei
Show and Tell: A Neural Image Caption Generator, Vinyals et al.
Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.
Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

test image

This image is CC0 public domain

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

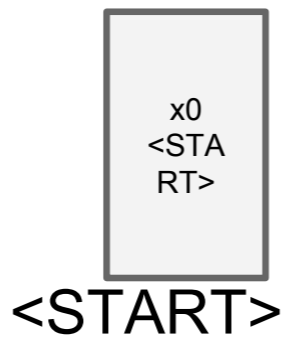conv-512
conv-512
maxpool

FC-4096
FC-4096
FC-1000
softmax

test image

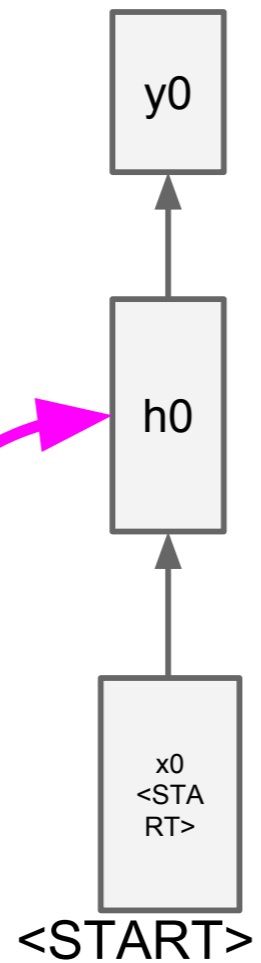this is our ImageNet CNN, now used as a feature extractor

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096
FC-1000
softmax
X

test image

this is our ImageNet CNN, now used as a feature extractor

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

test image

x0
<START>

<START>

test image

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

v

**Wih**

y0

h0

x0 <START>

<START>

**before:**

$h = \tanh(Wxh * x + Whh * h)$

**now:**

$h = \tanh(Wxh * x + Whh * h + \mathbf{Wih * v})$

let's use the image features to create a conditional LM

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

test image

y0    y1

h0 → h1

x0
<START>    straw

<START>

# Image Captioning: Failure Cases

*A woman is holding a cat in her hand*



*A person holding a computer mouse on a desk*
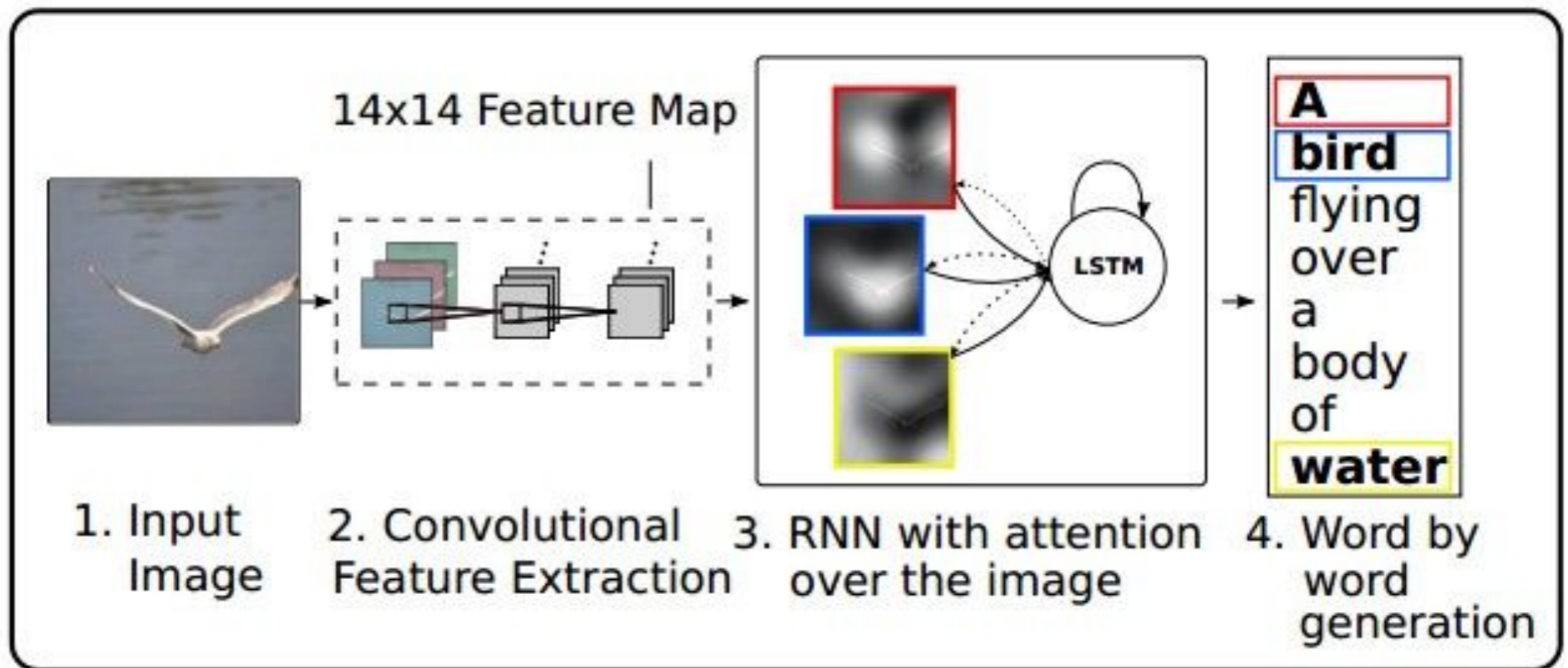


*A woman standing on a beach holding a surfboard*
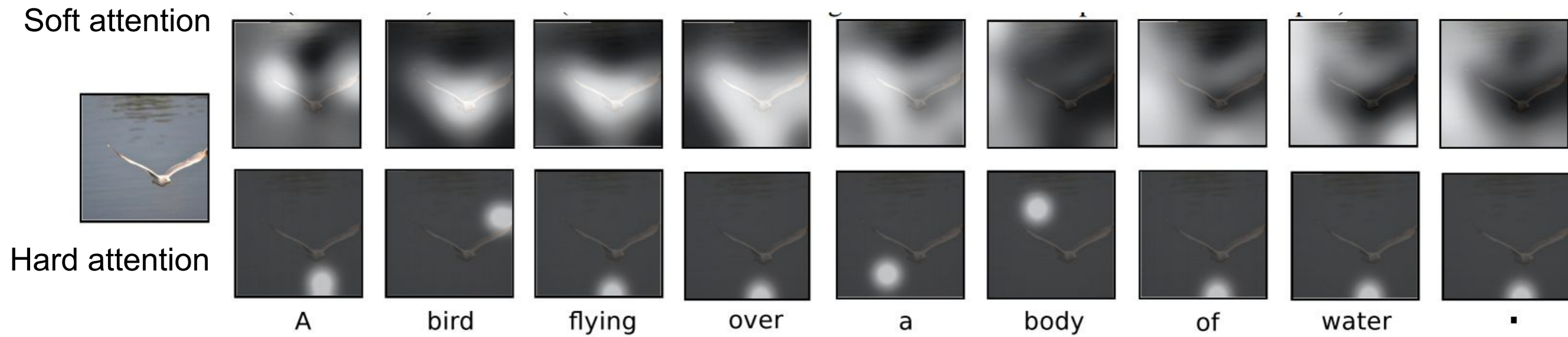


*A bird is perched on a tree branch*



*A man in a baseball uniform throwing a ball*
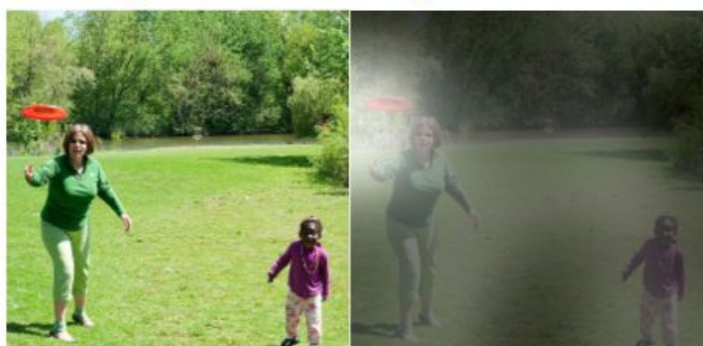
# Image Captioning with Attention

RNN focuses its attention at a different spatial location when generating each word

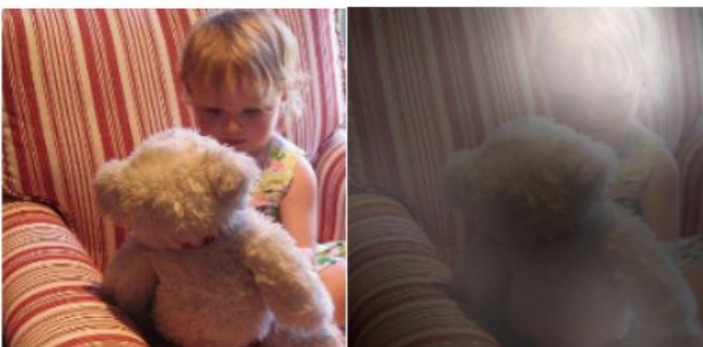# Image Captioning with Attention

Soft attention

Hard attention

A    bird    flying    over    a    body    of    water    .

Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with Attention



A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

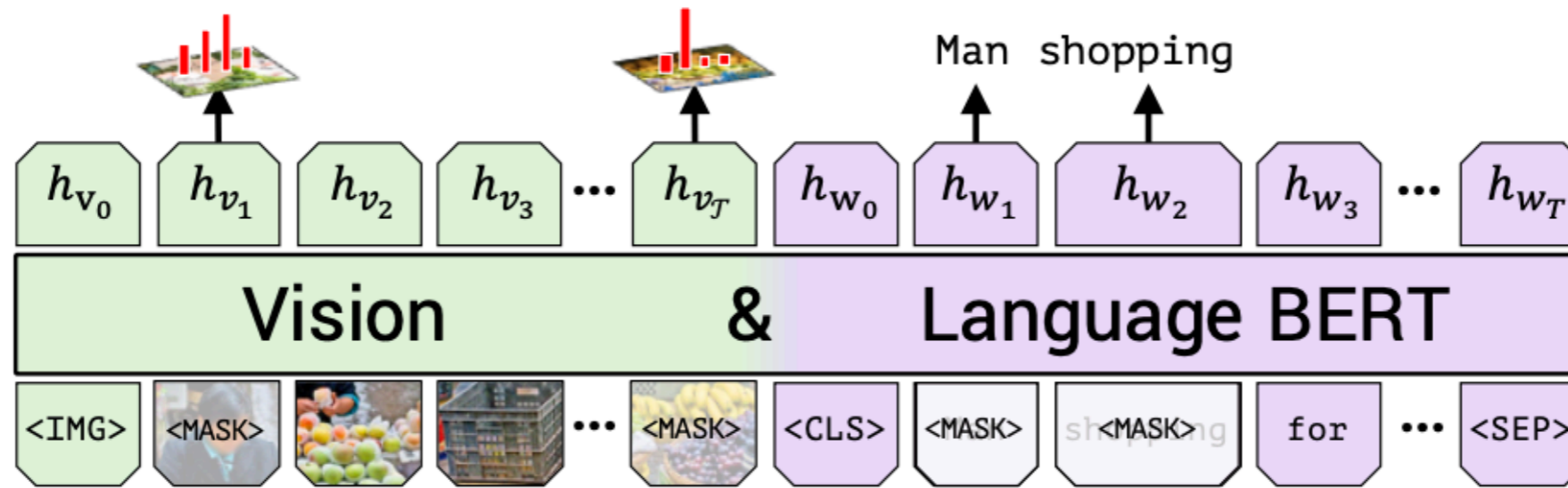A little <u>girl</u> sitting on a bed with a teddy bear.

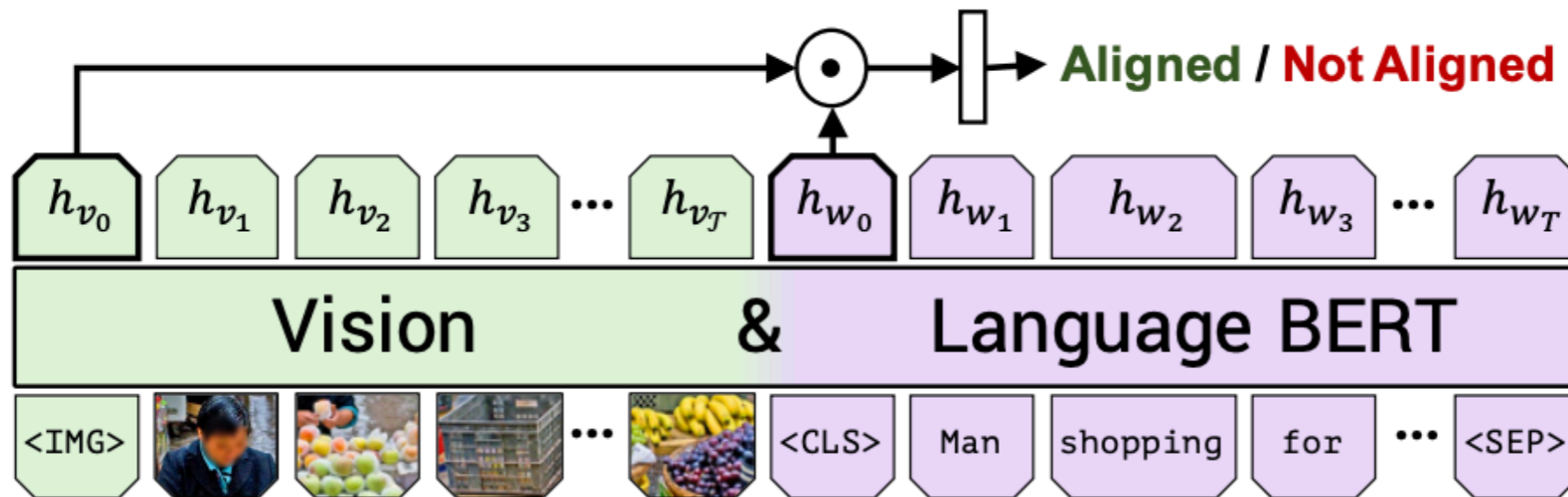A group of <u>people</u> sitting on a boat in the water.

A giraffe standing in a forest with <u>trees</u> in the background.

Xu et al, "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# VilBERT (vision and language BERT)



(a) Masked multi-modal learning

(b) Multi-modal alignment prediction

Lu et al., 2019 ("VilBERT")
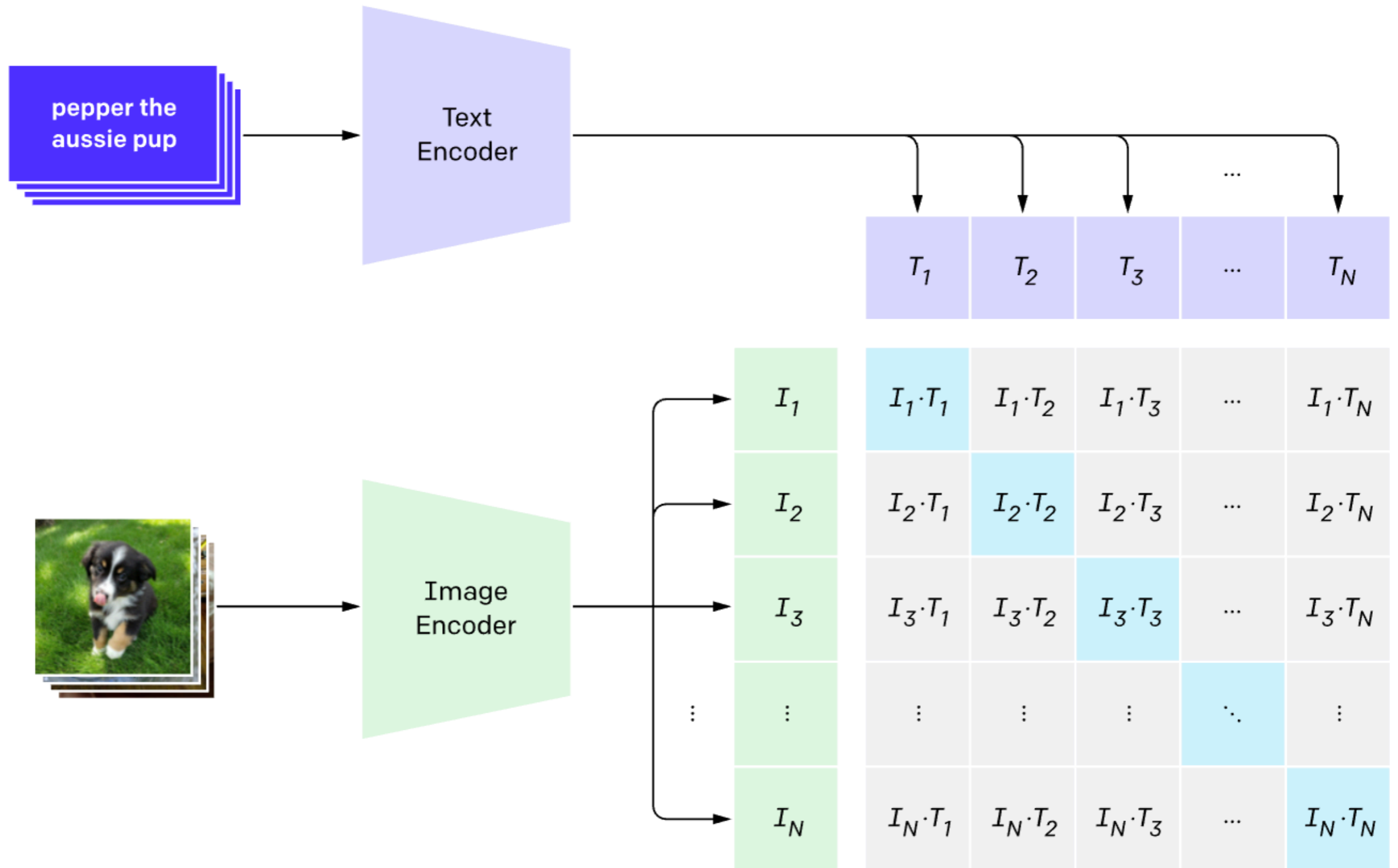
# OpenAI's CLIP: Contrastive language-image pretraining

- VilBERT and similar methods (e.g. LXMERT) rely on small labeled datasets like MSCOCO and Visual Genome (~100K images each)

- OpenAI collect 400 million (image, text) pairs from the web

- Then, they train an image encoder and a text encoder with a simple contrastive loss: given a collection of images and text, predict which (image, text) pairs actually occurred in the dataset
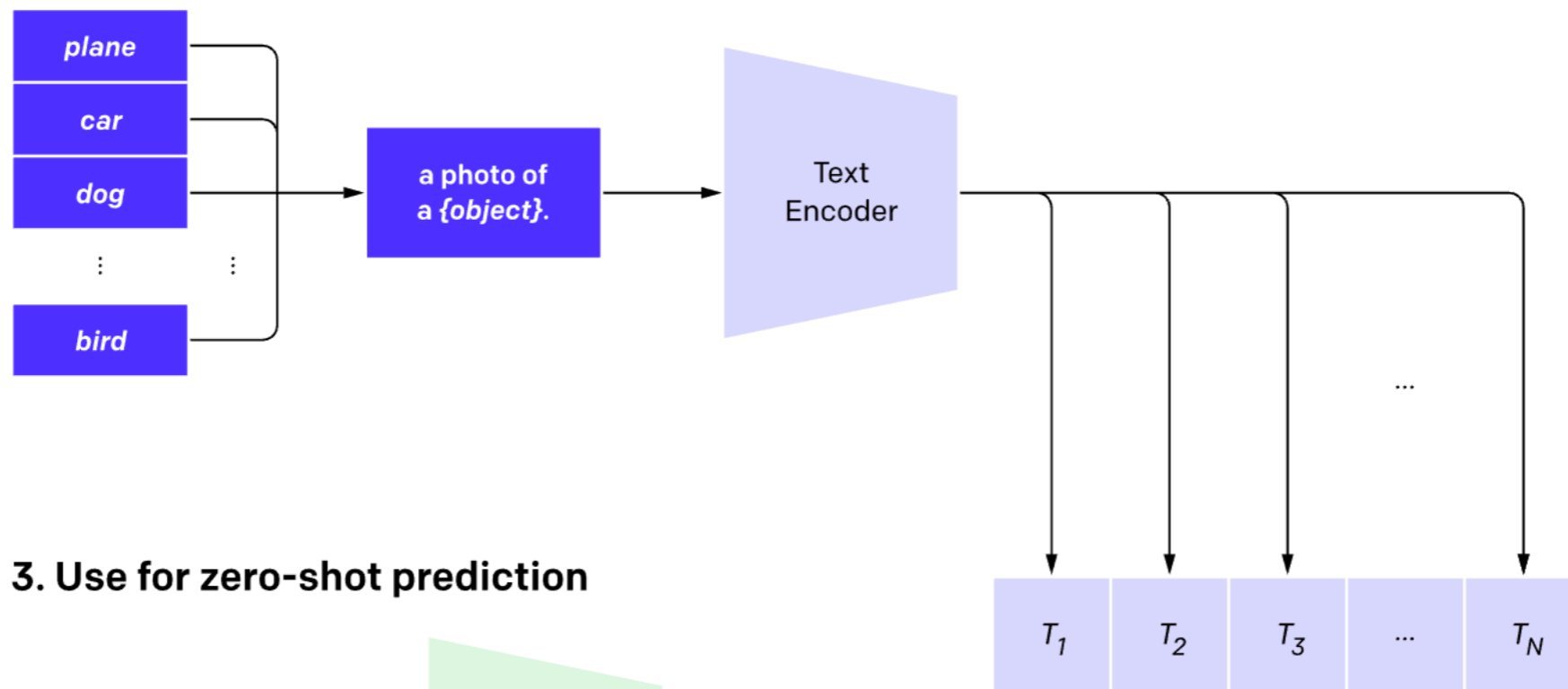
Radford et al., 2021 ("CLIP")
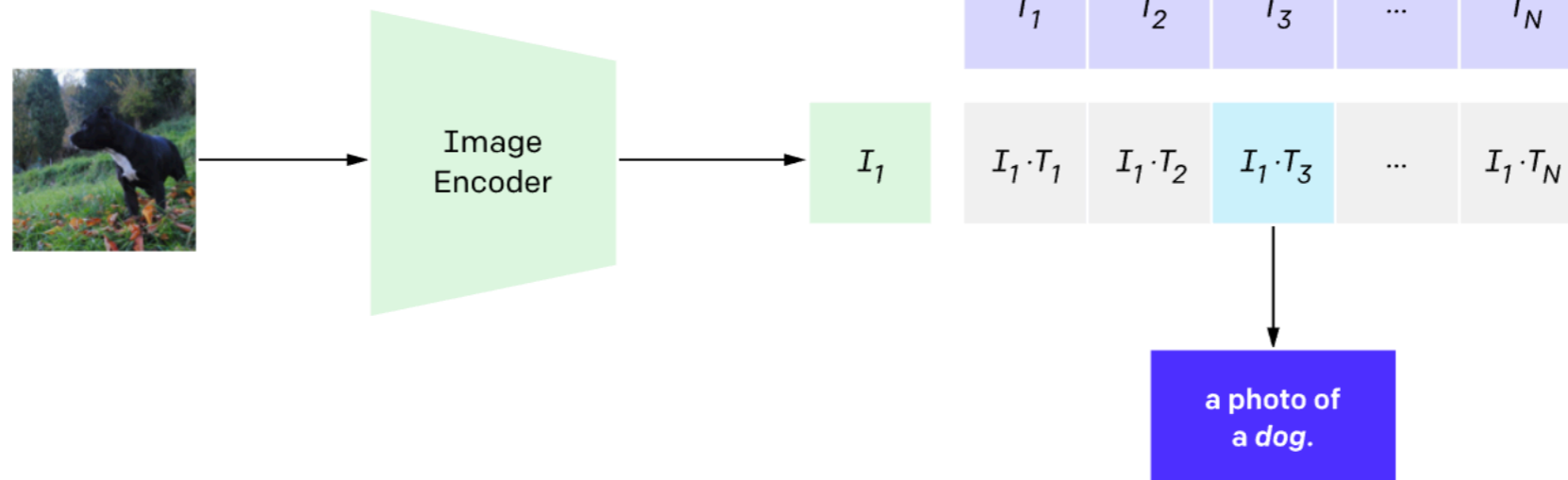
# 1. Contrastive pre-training

# Similar to GPT-3, you can use CLIP for zero-shot learning



**2. Create dataset classifier from label text**

plane / car / dog / ... / bird → a photo of a {object}. → Text Encoder → $T_1$ $T_2$ $T_3$ ... $T_N$

**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*.

| 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% |