

Representation Learning

Semester 2, 2025

Kris Ehinger

Outline

- Contrastive learning
- CLIP
- Visualisation and interpretation
- Invariance, generalisation

Learning outcomes

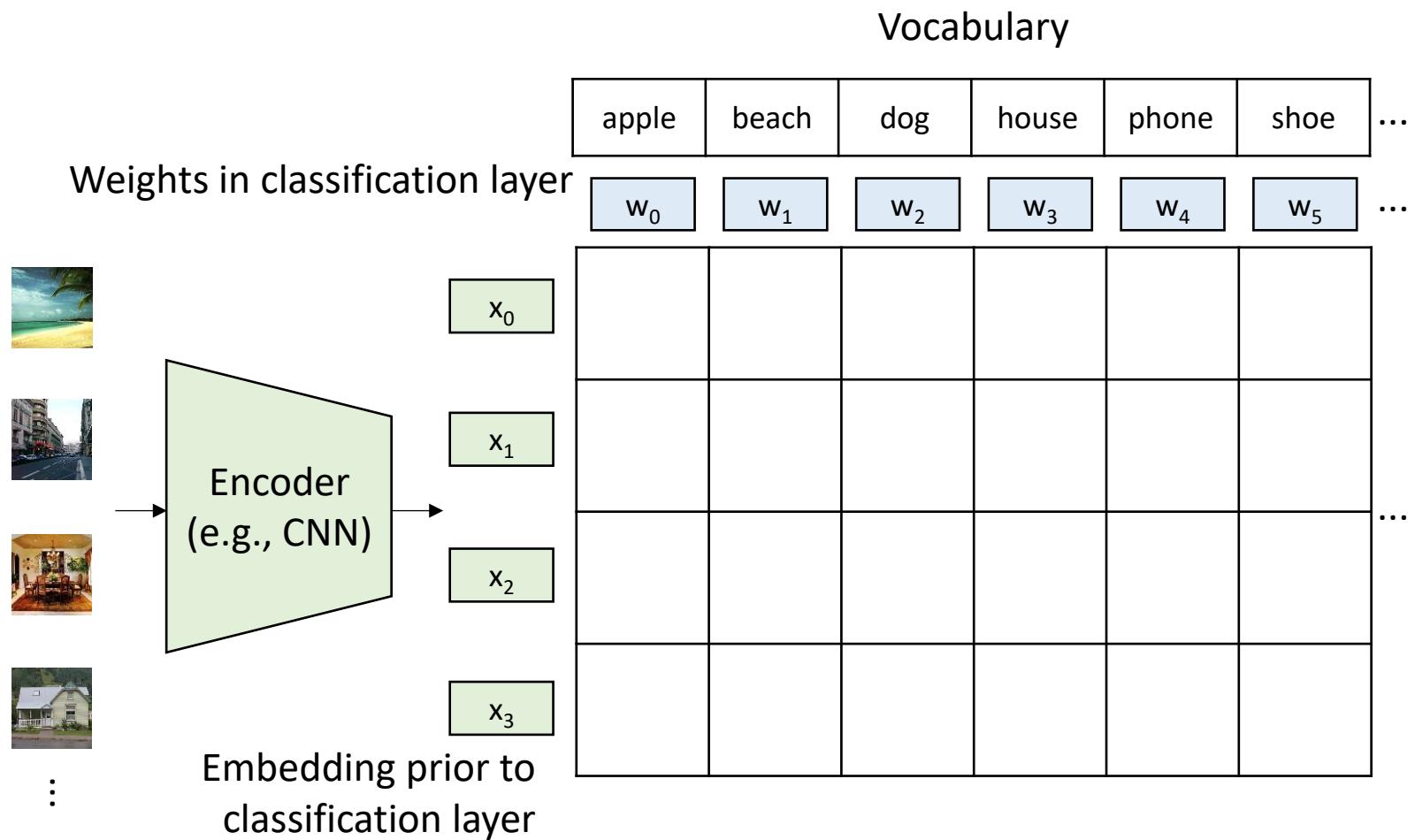
- Explain the concepts behind contrastive learning and CLIP
- Explain methods to visualise or explain the decisions of an image recognition algorithm
- Evaluate computer vision algorithms in terms of their invariance (or tolerance) to image variation

Contrastive learning

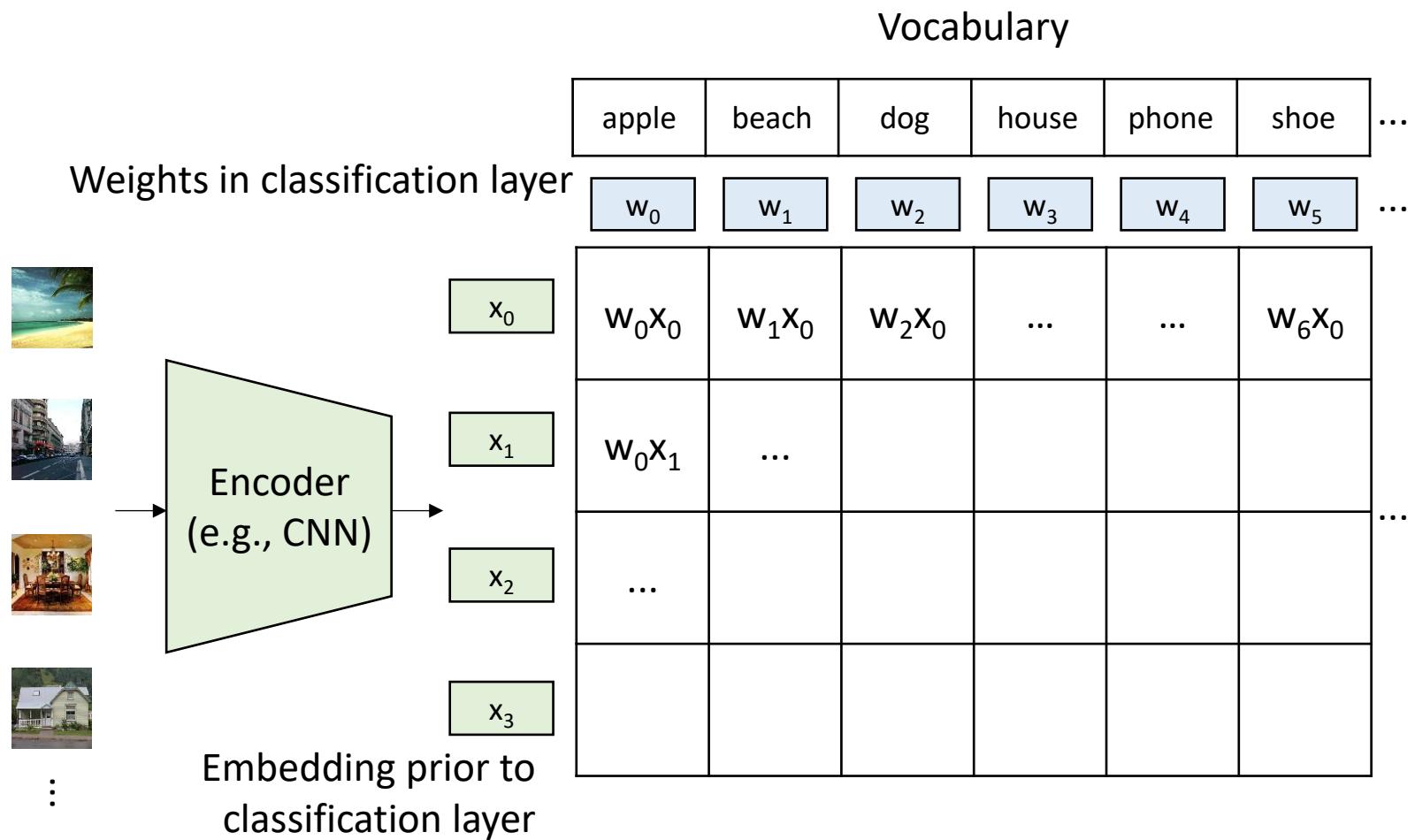
Contrastive learning

- Self-supervised training method
- Recall:
 - Supervised = train from ground truth labels
 - Unsupervised = train without labels (e.g., clustering)
 - Self-supervised = generate "labels" from data itself
- Example: mask part of the image and train a model to predict the missing part
- Advantage: no need for manual annotation

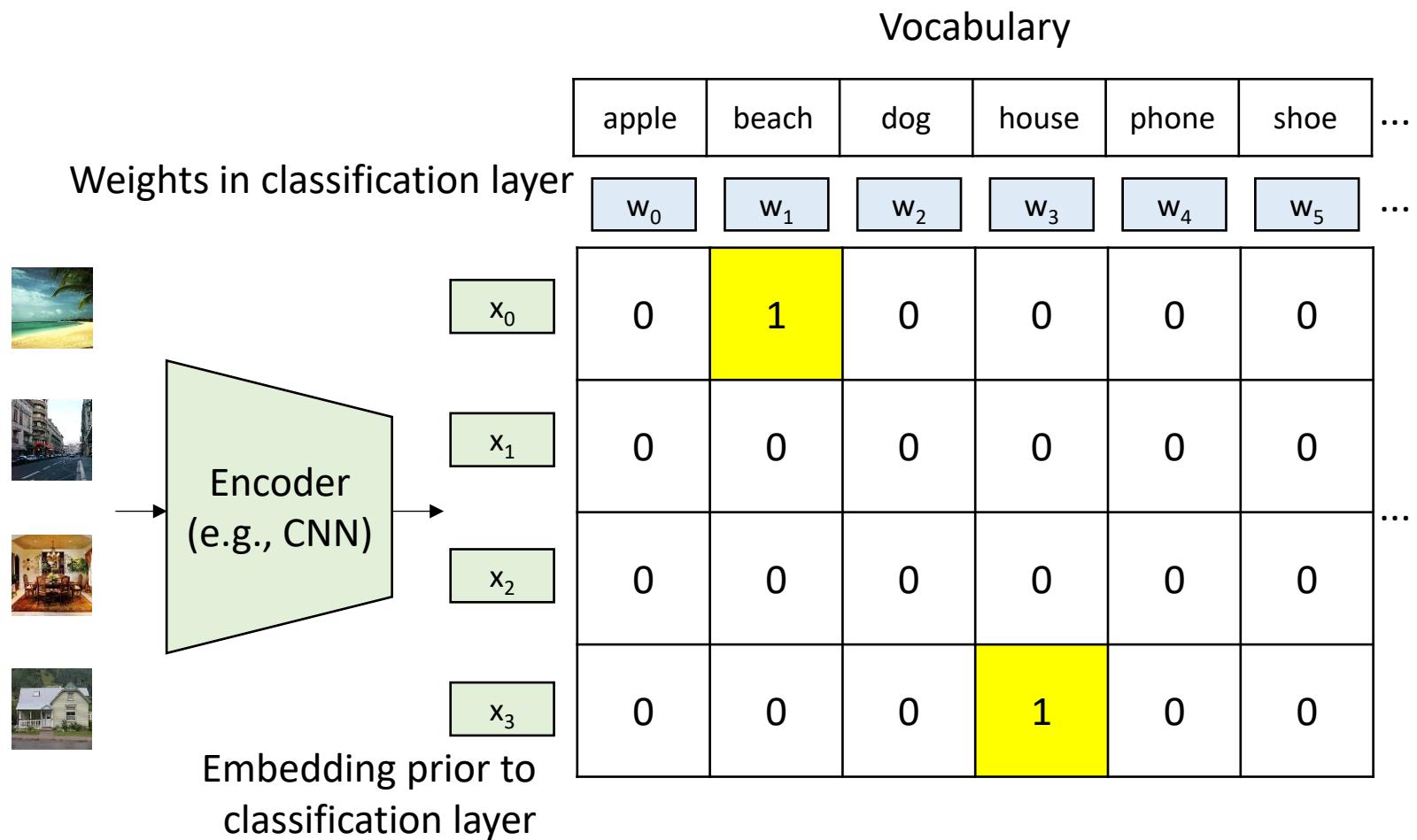
Contrastive learning (w/ labels)



Contrastive learning (w/ labels)



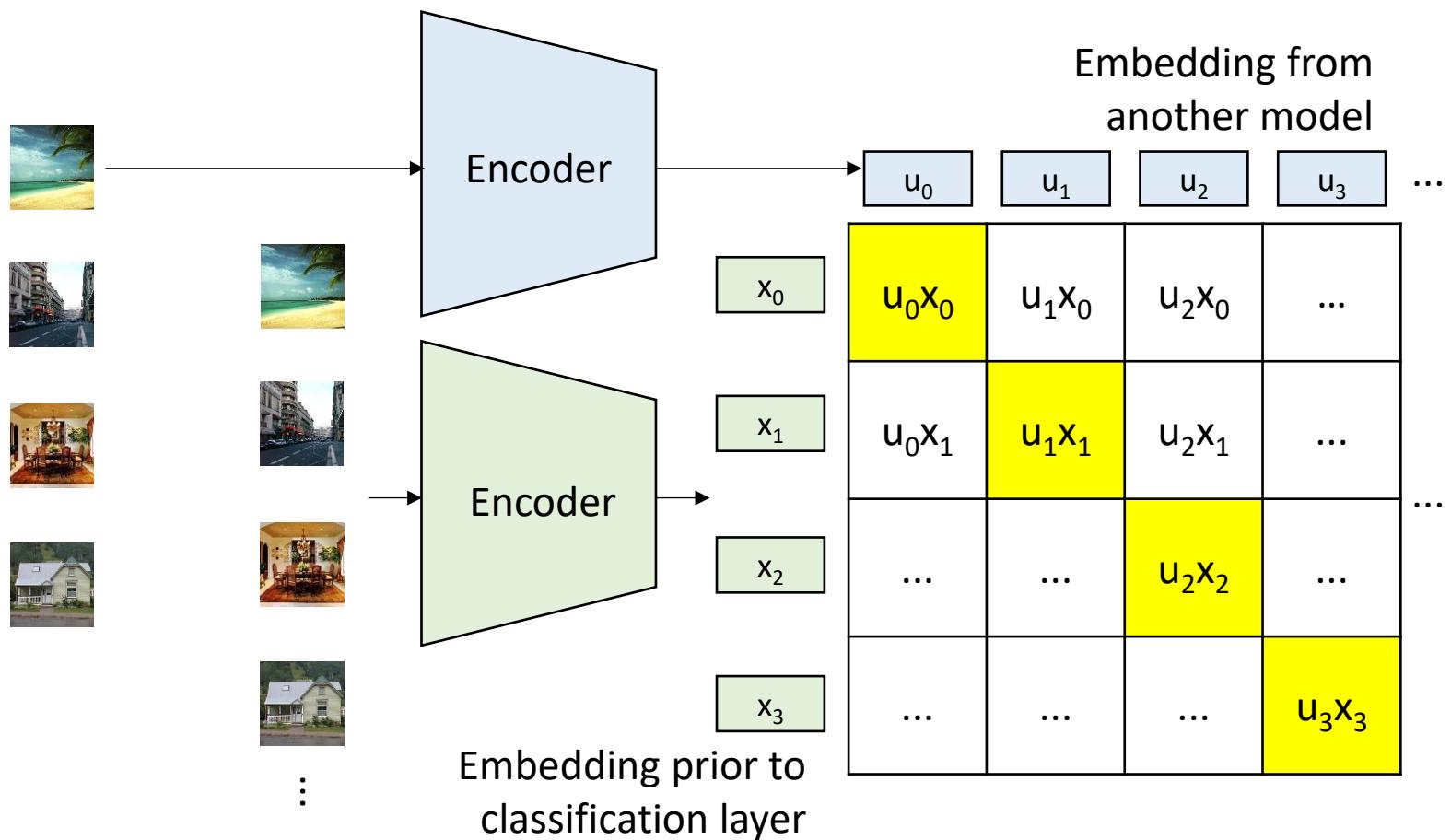
Contrastive learning (w/ labels)



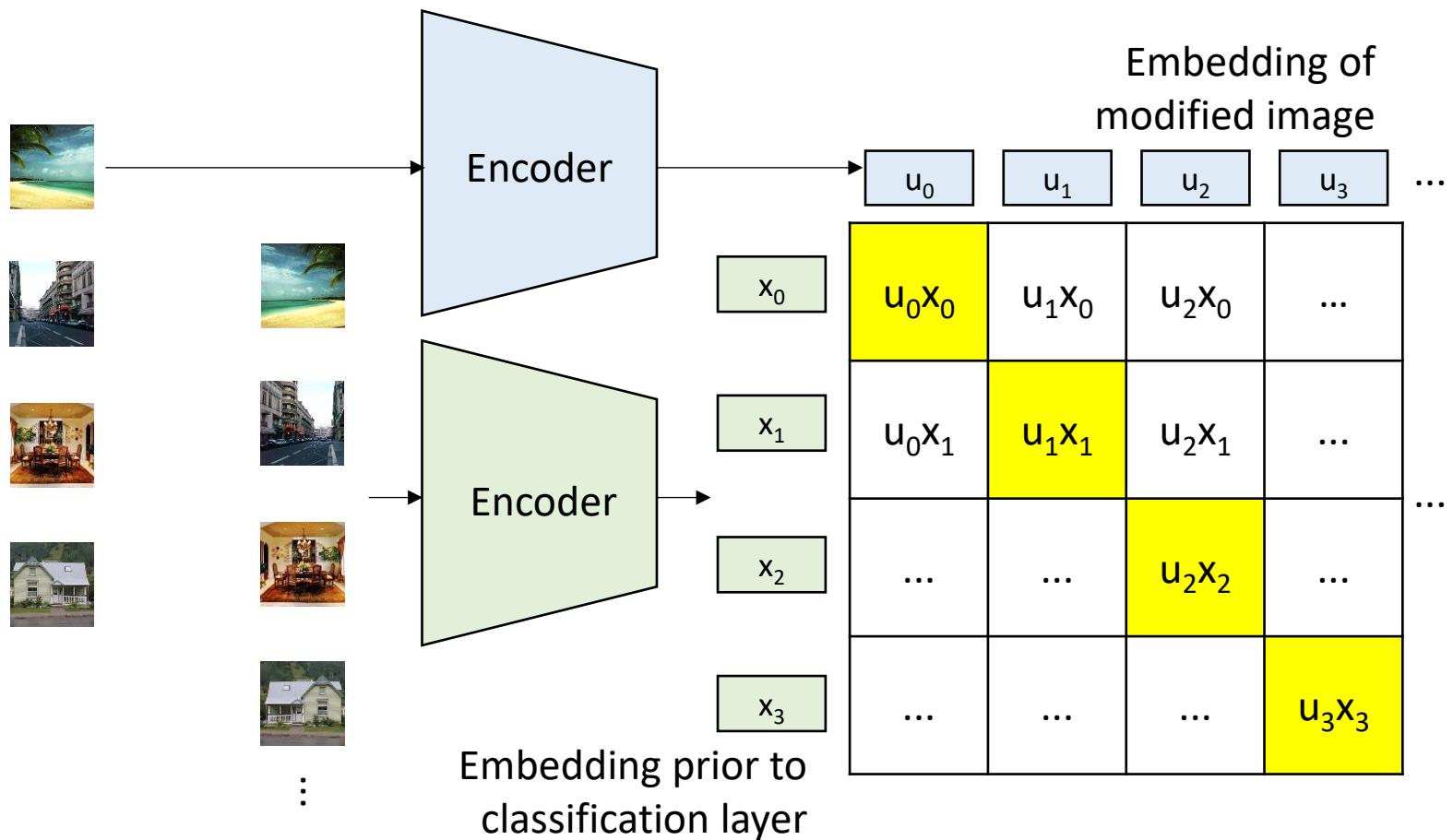
Contrastive learning (w/ labels)

- Goal is to:
 - Maximize similarity of positive pairs
 - Minimize similarity of negative pairs
- What if you have no labels? Make the vocabulary from the data.

Contrastive learning (no labels)



Contrastive learning (no labels)



Contrastive loss

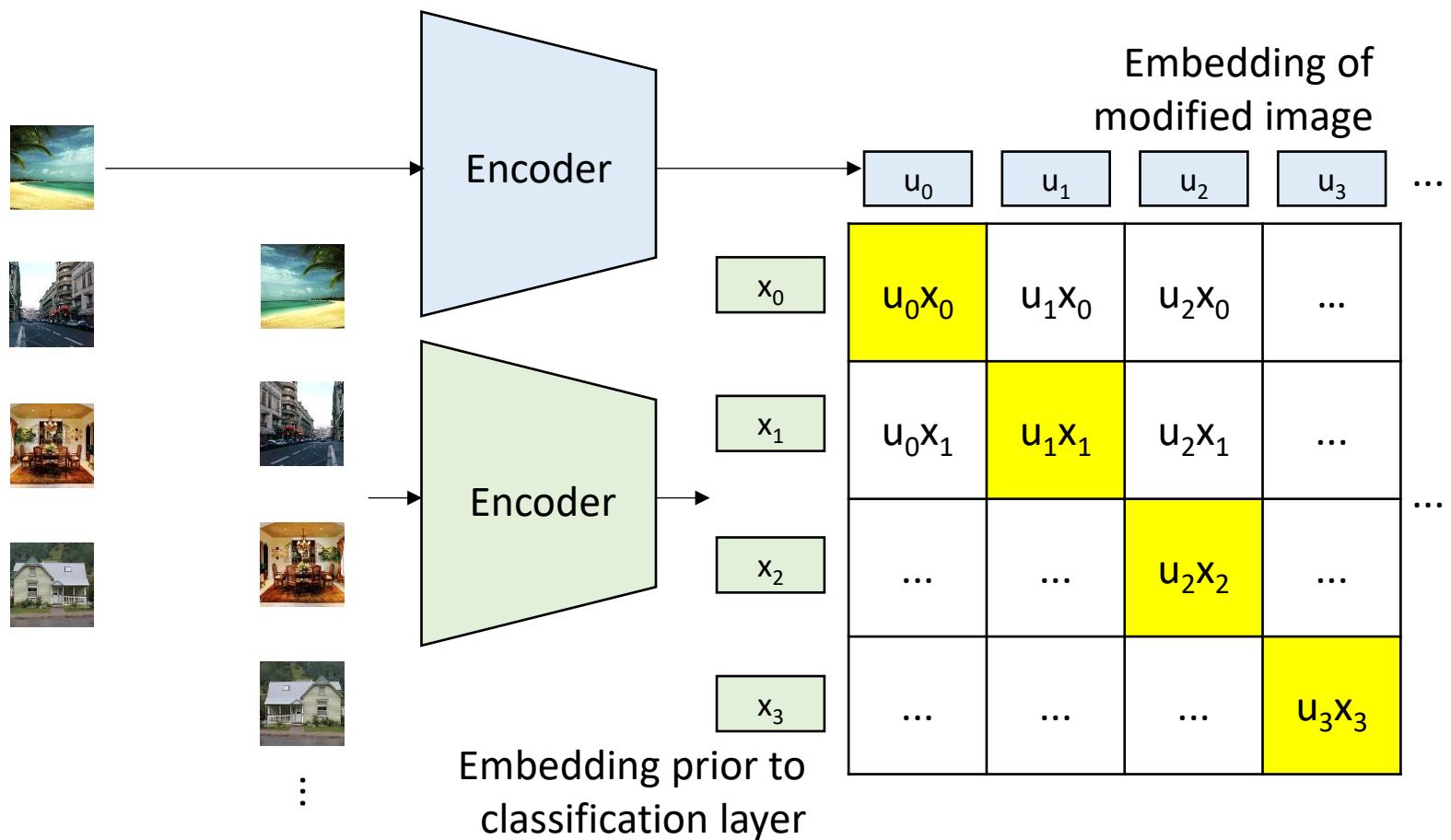
- Contrastive loss function

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} [\mathbf{k} \neq i] \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)},$$

Similarity to match Temperature
 Similarity to all non-matches

- To reduce loss, increase similarity of matched pairs and decrease similarity of non-matched pairs
 - Temperature = smoothness parameter (**low value means similarity drops sharply**)

Contrastive learning (no labels)



Contrastive learning

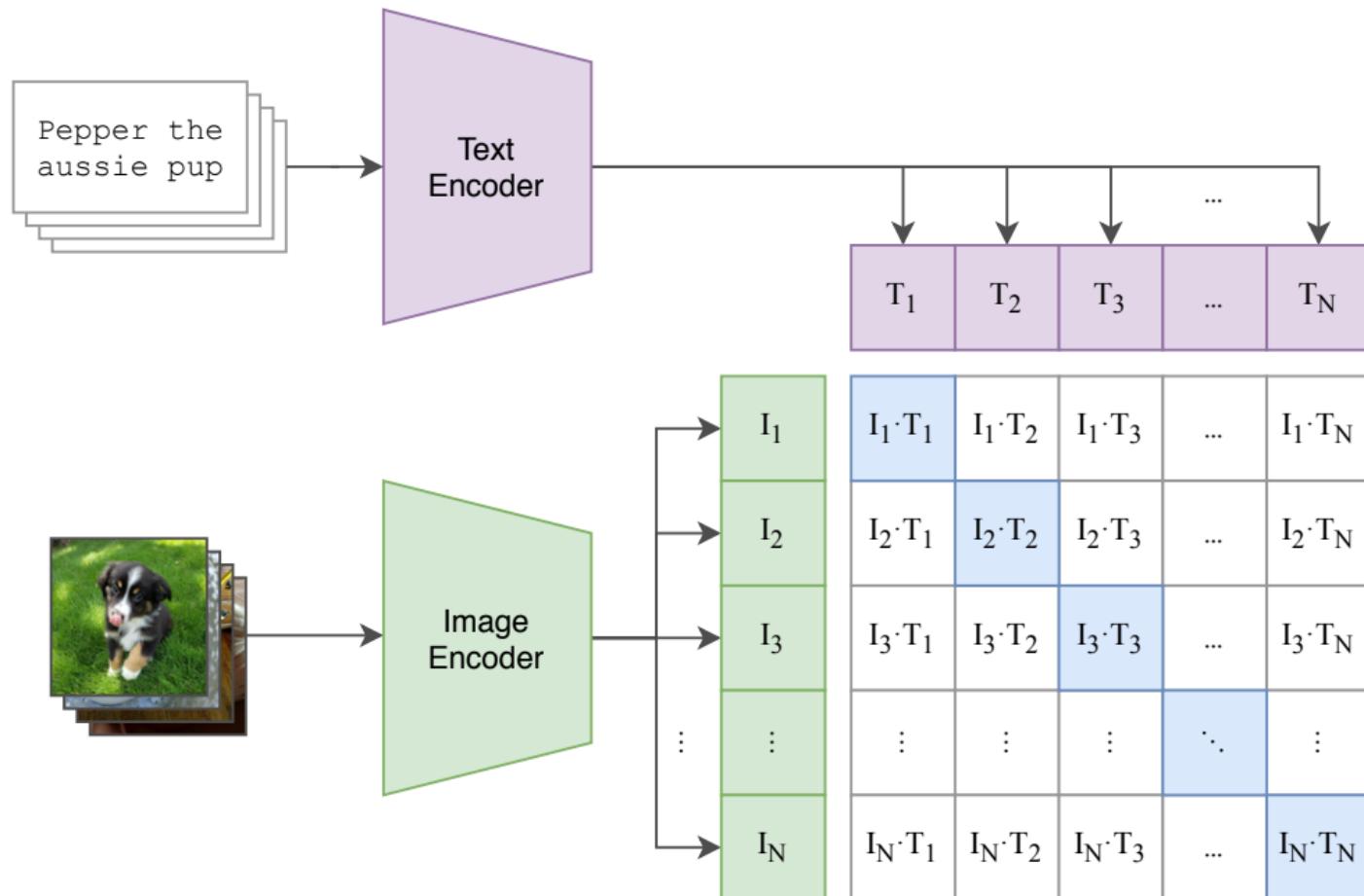
- Self-supervised learning: there is a training signal, but generated from the data itself (e.g., two variants of the same image)
- Many different implementations with different approaches to sampling examples, backpropagating loss

CLIP

CLIP

- CLIP = Contrastive Language-Image Pretraining
- Self-supervised learning on a large dataset of images + text captions
 - Reduces requirement for manual labelling of images
 - Wider variety of concepts that could be learned (not just 1000 object classes)

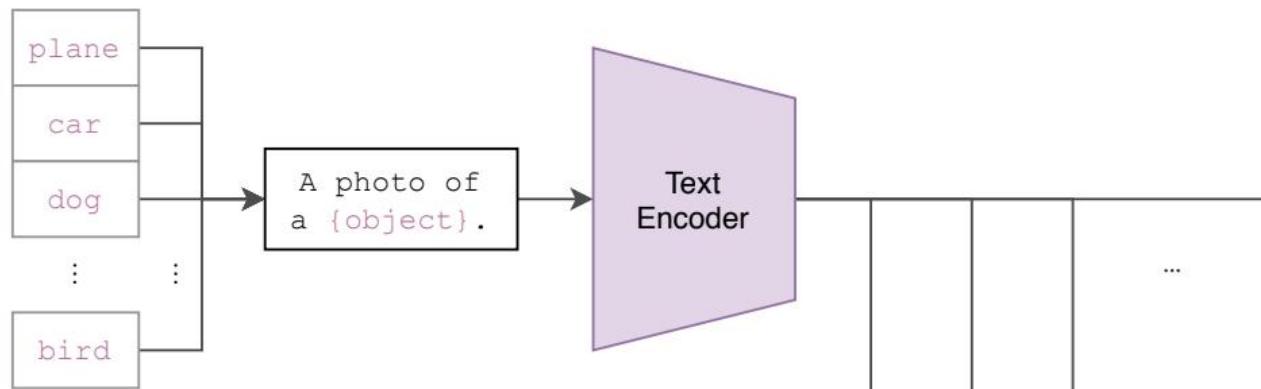
CLIP



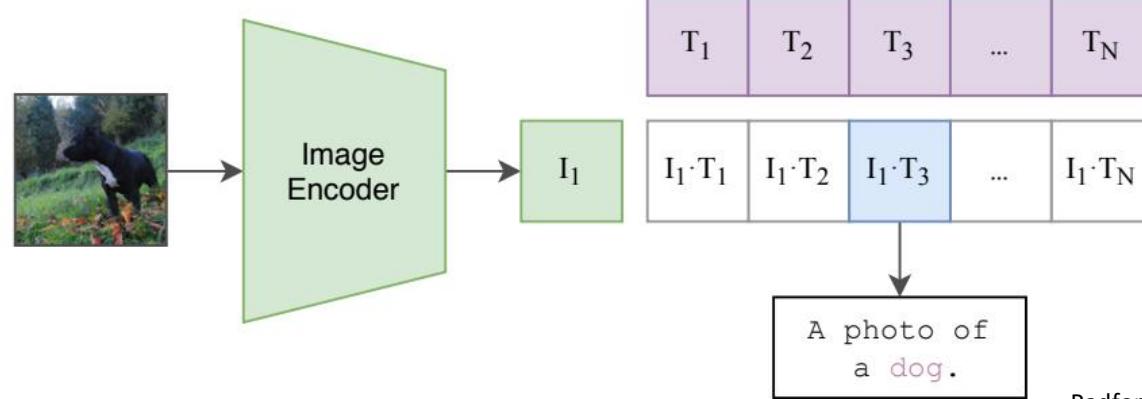
Radford, et al. (2021)

Using CLIP as a classifier

(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

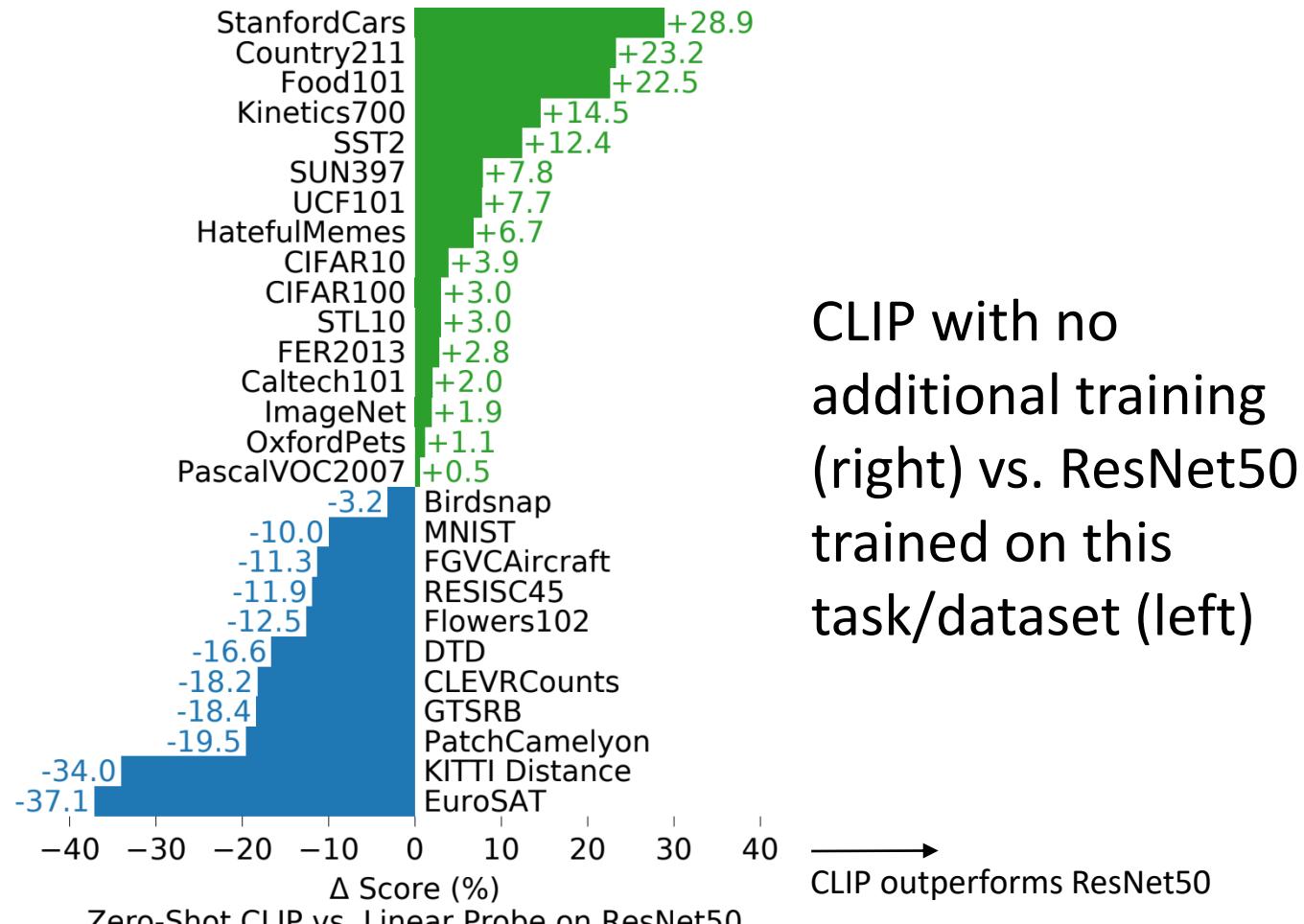


Radford, et al. (2021)

CLIP

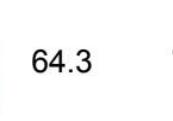
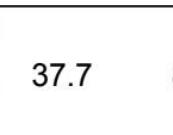
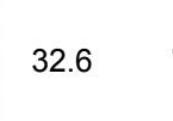
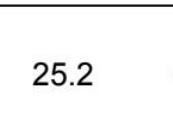
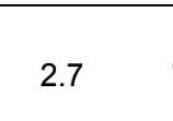
- Trained on WIT, proprietary dataset of 400 million image-text pairs (mostly from internet?)
- Train on a batch of N image-text pairs:
 - N^2 possible pairs
 - N are true matches
 - $N^2 - N$ are false matches
- Measure cosine similarity between image and text embedding
- Contrastive learning with multi-class N-pair loss

Zero-shot classification



Radford, et al. (2021)

Robustness of CLIP

	Dataset Examples						ImageNet	Zero-Shot ResNet101	CLIP	Δ Score
ImageNet							76.2	76.2		0%
ImageNetV2							64.3	70.1		+5.8%
ImageNet-R							37.7	88.9		+51.2%
ObjectNet							32.6	72.3		+39.7%
ImageNet Sketch							25.2	60.2		+35.0%
ImageNet-A							2.7	77.1		+74.4%

Radford, et al. (2021)

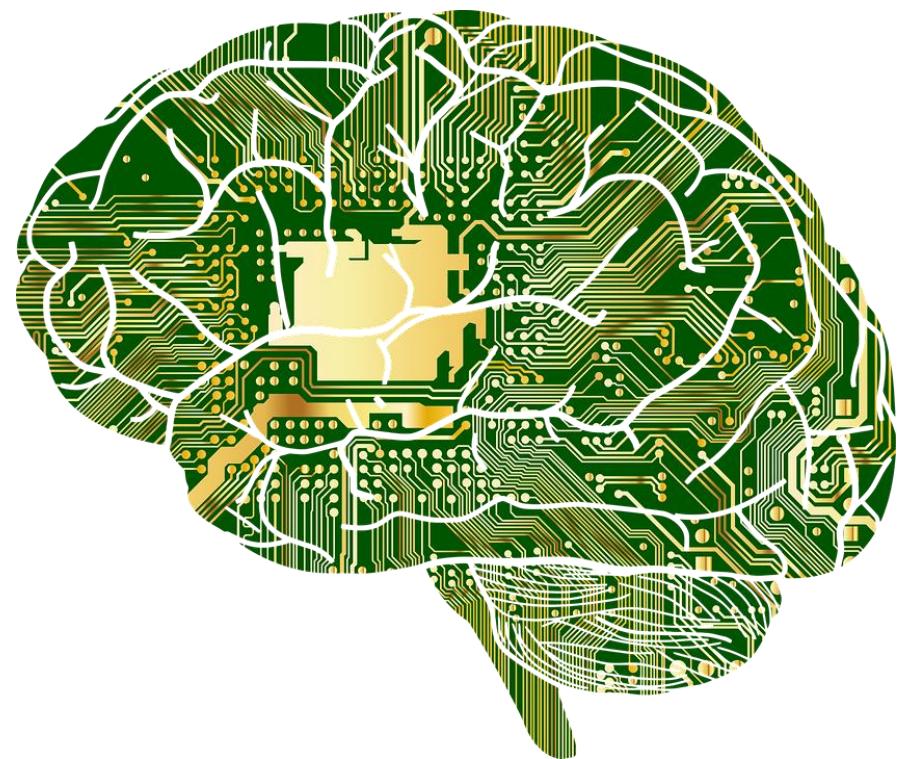
Summary

- CLIP is a particularly useful model for generating embeddings:
 - Trained on a very large image dataset (400M)
 - Can generate text pseudo-labels due to shared embedding space
- Common uses:
 - Auto-labelling large image datasets (e.g., LAION)
 - Embeddings for LLMs

Model visualisation

Understanding vision models

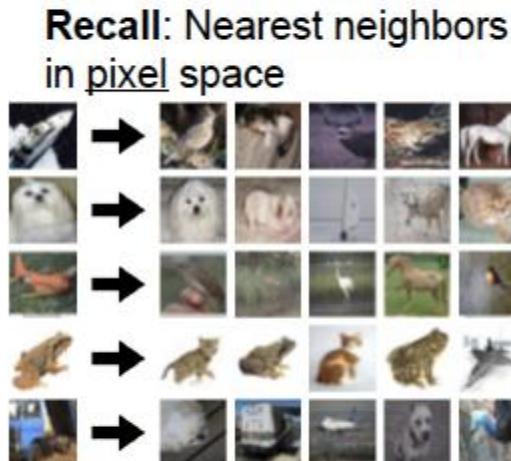
How do we know what a computer vision model is thinking?



Visualising embeddings

- How are images organised in this feature space?
- It's a high-dimensional space, so we can't just plot all images in this space
 - Can use dimensionality reduction
 - Or look at local regions (what images are near neighbours in this space?)

Nearest neighbours



Test image L2 Nearest neighbors in feature space



Krizhevsky, Sutskever, & Hinton (2012)

Feature space visualisation

- Options for dimensionality reduction
- PCA (principal component analysis)
 - Show the dimensions with the most variance
 - Simple but often hard to interpret since only a few dimensions can be visualised simultaneously
- t-SNE (t-distributed stochastic neighbor embedding)
 - Flatten high-dimensional data into 2D or 3D so that near neighbours stay nearby



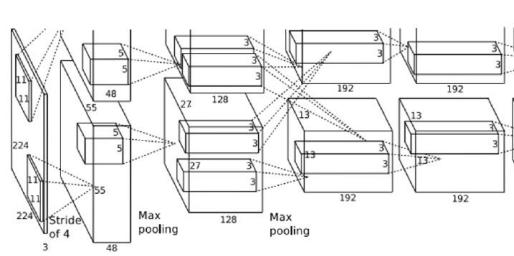
<https://cs.stanford.edu/people/karpathy/cnnembed/>

"Where" is the model looking?

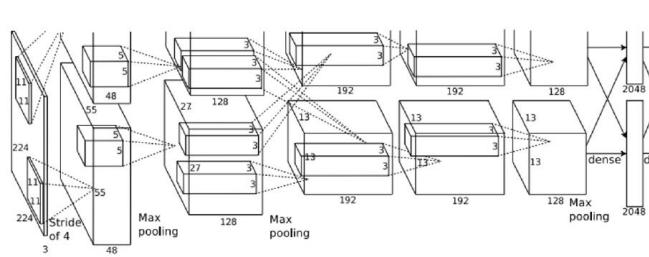
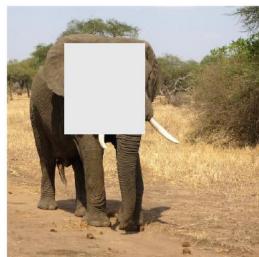
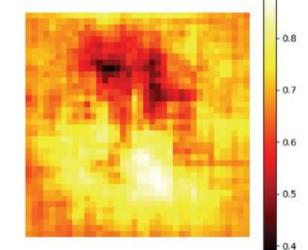
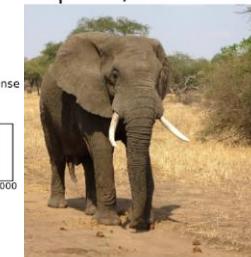
- What parts of an image are most important for determining the class label?
 - Can help show what features the model uses to determine class
 - Can help debug problems (e.g., using background to classify object, label confusion when there are multiple objects)

Visualising image regions

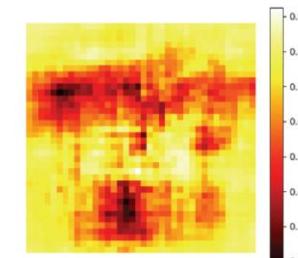
- Occlusion method: mask image and see how much class probability changes



African elephant, *Loxodonta africana*

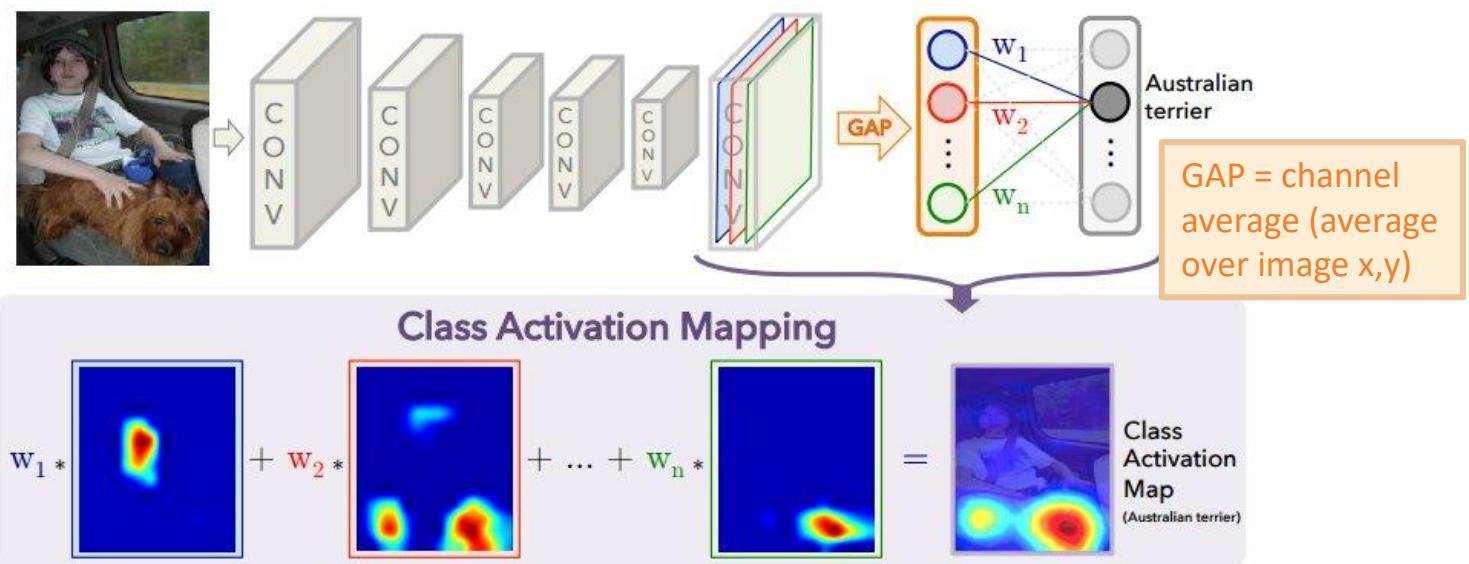


go-kart

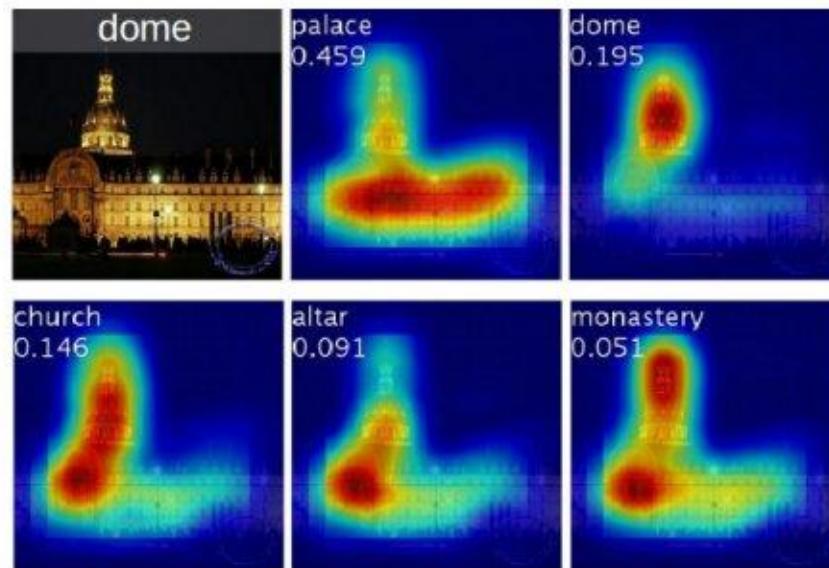


CAM (Class Activation Mapping)

- Add a Global Average Pooling (GAP) layer before classification layer, use weights of this layer to determine *where* the class-relevant features are



CAM (Class Activation Mapping)



Class activation maps of top 5 predictions



Class activation maps for one object class

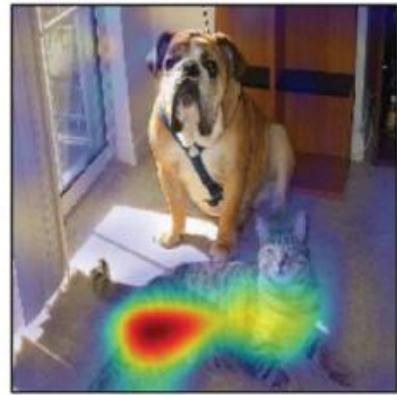
Grad-CAM



(a) Original Image



(b) Guided Backprop ‘Cat’



(c) Grad-CAM ‘Cat’



(g) Original Image



(h) Guided Backprop ‘Dog’



(i) Grad-CAM ‘Dog’

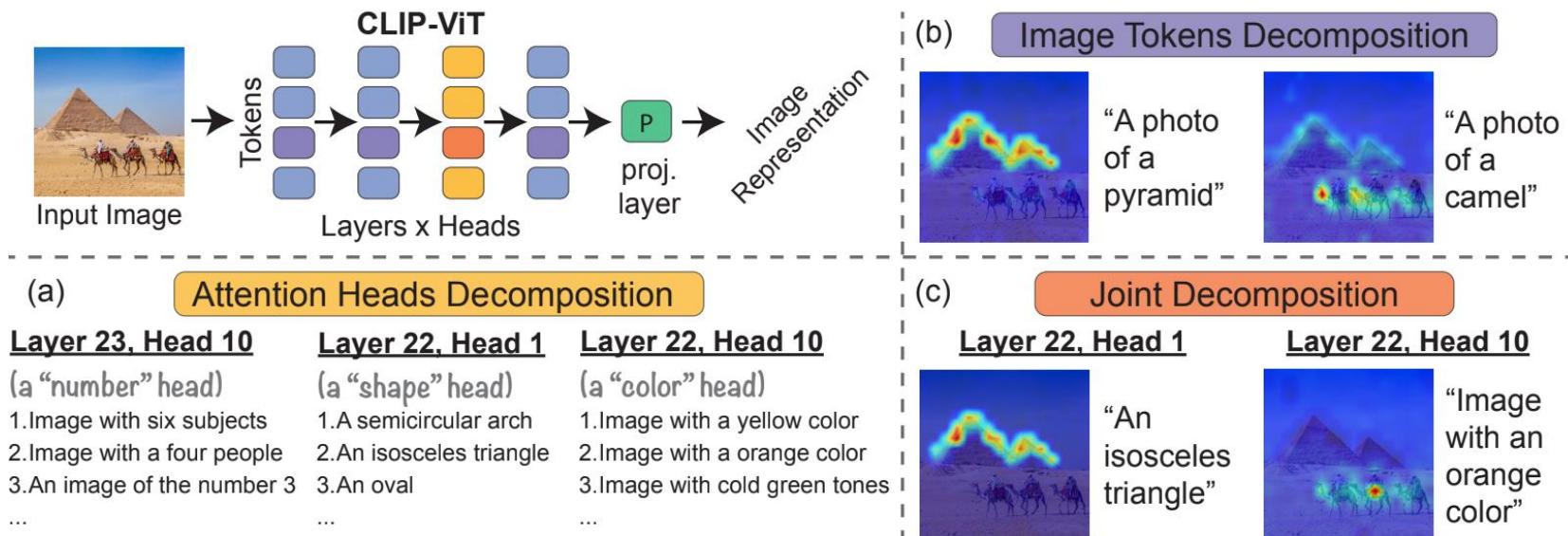
Selvaraju, et al. (2016)

"What" does the model encode?

- What features does a particular (layer/neuron/attention head/...) encode?
- Visualise the activation maps, filters, or show images that produce a high response
 - Or, for models like CLIP, output text that produces a high response

"What" does the model encode?

Break down response into contributions of each attention head ("what") and patch location ("where"):



"What" does the model encode?

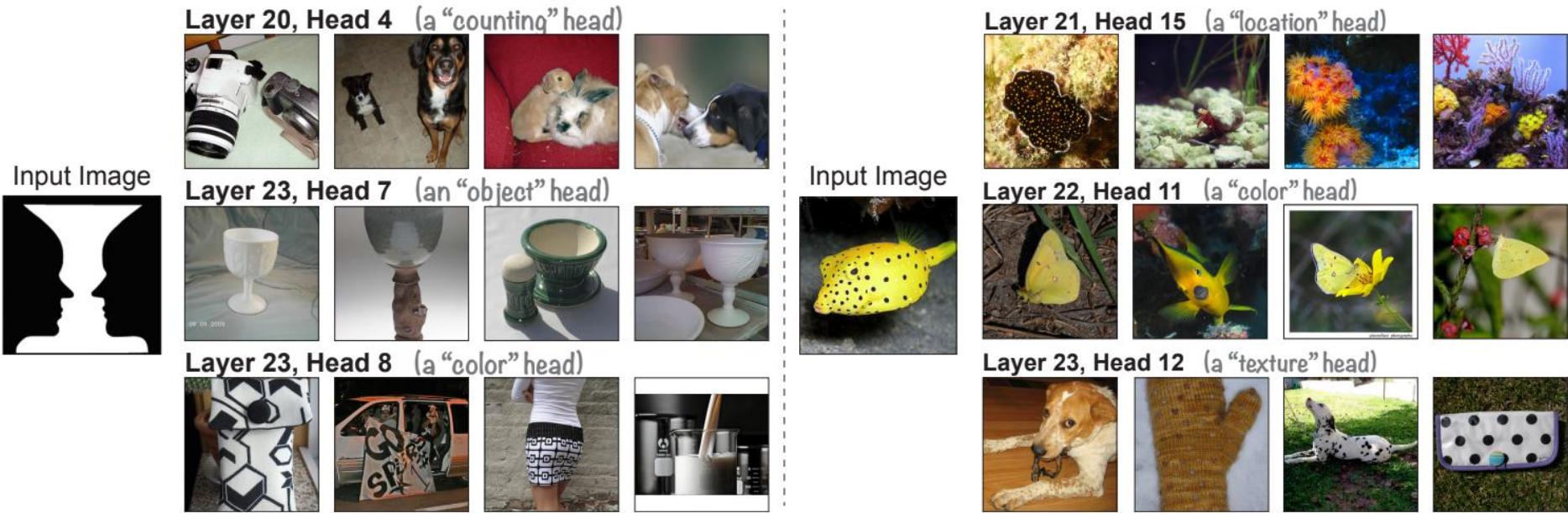
L21.H11 (“Geo-locations”)	L23.H10 (“Counting”)	L22.H8 (“Letters”)
Photo captured in the Arizona desert	Image with six subjects	A photo with the letter V
Picture taken in Alberta, Canada	Image with four people	A photo with the letter F
Photo taken in Rio de Janeiro, Brazil	An image of the number 3	A photo with the letter D
Picture taken in Cyprus	An image of the number 10	A photo with the letter T
Photo taken in Seoul, South Korea	The number fifteen	A photo with the letter X
L22.H11 (“Colors”)	L22.H6 (“Animals”)	L22.H3 (“Objects”)
A charcoal gray color	Curious wildlife	An image of legs
Sepia-toned photograph	Majestic soaring birds	A jacket
Minimalist white backdrop	An image with dogs	A helmet
High-contrast black and white	Image with a dragonfly	A scarf
Image with a red color	An image with cats	A table
L23.H12 (“Textures”)	L22.H1 (“Shapes”)	L22.H2 (“Locations”)
Artwork with pointillism technique	A semicircular arch	Urban park greenery
Artwork with woven basket design	An isosceles triangle	Cozy home interior
Artwork featuring barcode arrangement	An oval	Urban subway station
Image with houndstooth patterns	Rectangular object	Energetic street scene
Image with quilted fabric patterns	A sphere	Tranquil boating on a lake

"What" does the model encode?

Layer 23, Head 0	Layer 23, Head 1
Intricate wood carving Nighttime illumination Image with woven fabric design Image with shattered glass reflections A photo of food	Photograph taken in a retro diner Intense athlete Detailed illustration of a futuristic bioreactor Image with holographic retro gaming aesthetics Antique historical artifact
Layer 23, Head 2	Layer 23, Head 3
Image showing prairie grouse Image with a penguin A magnolia An image with dogs An image with cats	Bustling city square Serene park setting Warm and cozy indoor scene Modern airport terminal Remote hilltop hut
Layer 23, Head 4	Layer 23, Head 5
Playful siblings A photo of a young person Image with three people A photo of a woman A photo of a man	Intertwined tree branches Flowing water bodies A meadow A smoky plume Blossoming springtime blooms

"What" does the model encode?

Nearest neighbors of an input image, according to different attention heads:

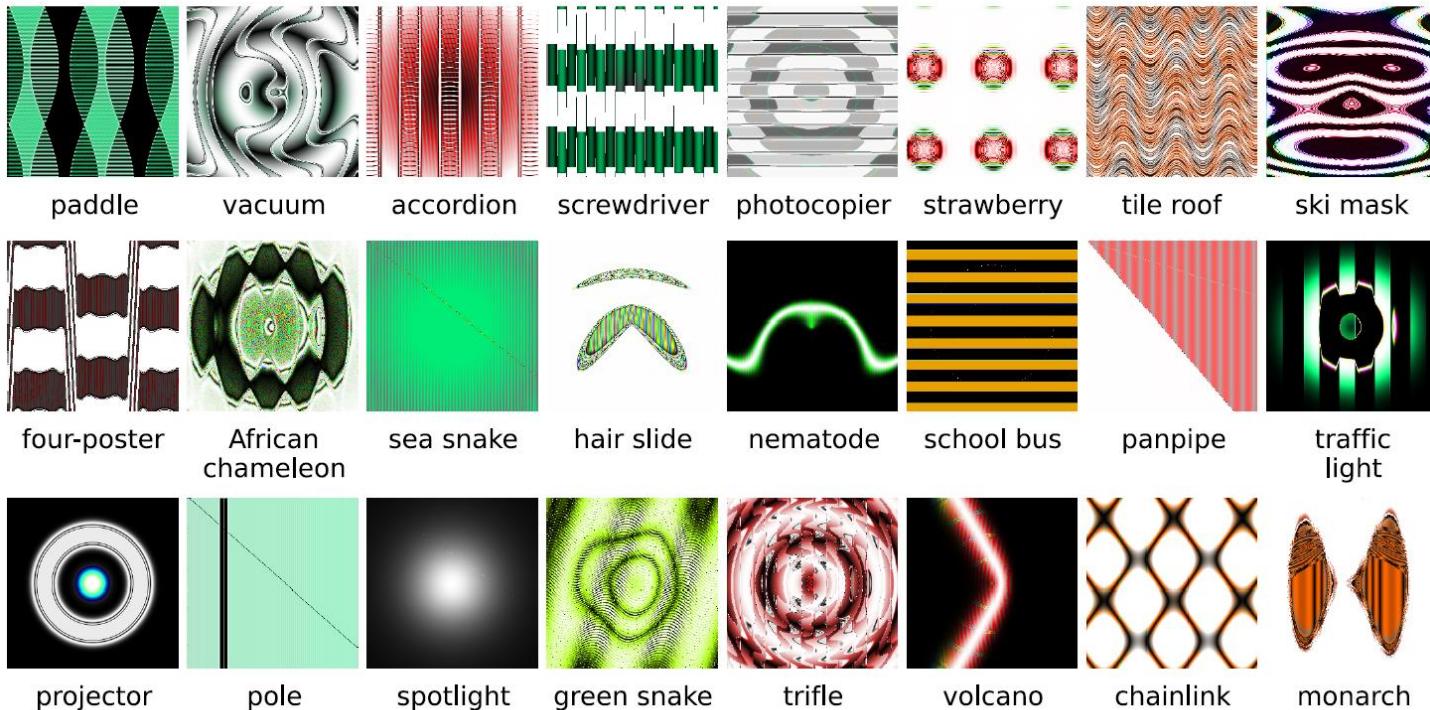


Visualising classes

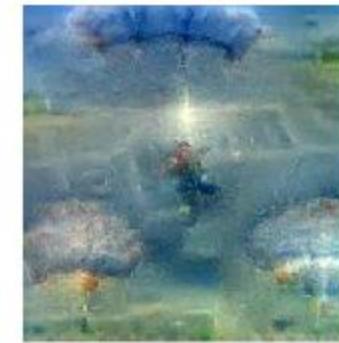
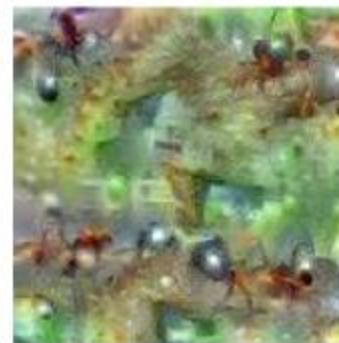
- Usually based on gradient ascent – synthesize image that maximises class label response
 - Initialise an image with zeros or small random noise
 - Run image through network, compute gradient
 - Update *image pixels* in a direction that minimises loss
- Problem: there are many possible arrays of pixels that can generate very high model response; not all of these will look like realistic images

Visualising classes

High-confidence classifications



Visualising classes



Google AI Blog: Mordvintsev, Olah, & Tyka (2015)

<https://ai.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>

Visualising classes



Nguyen, et al. (2016)

Summary

- Various ways to visualise what a model is doing
 - Visualise the feature space
 - Visualise the classes
 - Check where in the image a model is looking
 - Check what each part of a model (e.g., neuron or attention head) represents

Invariance & generalisation

Invariance / tolerance

- Are the features learned by CNNs invariant to
 - Lighting?
 - Translation?
 - Image plane rotation?
 - Scale?
 - 3D rotation / pose?

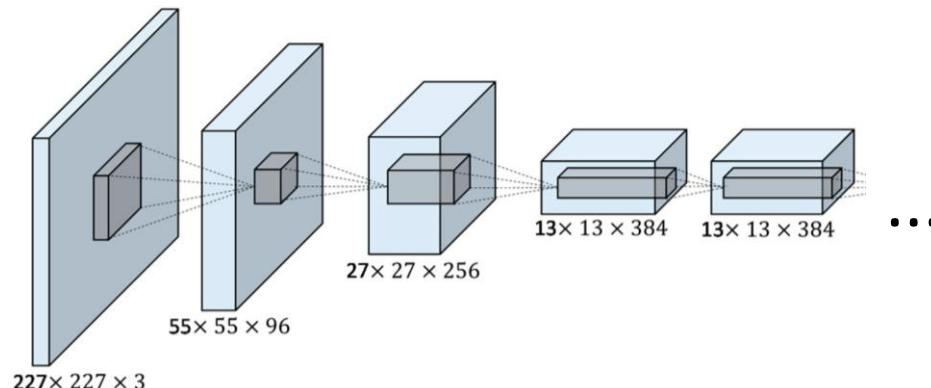
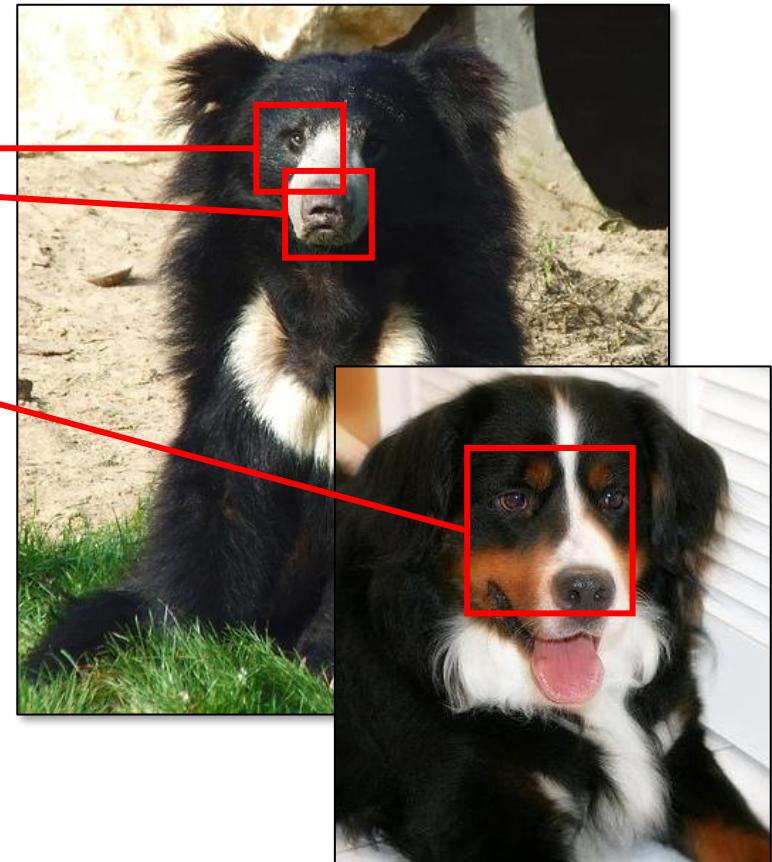
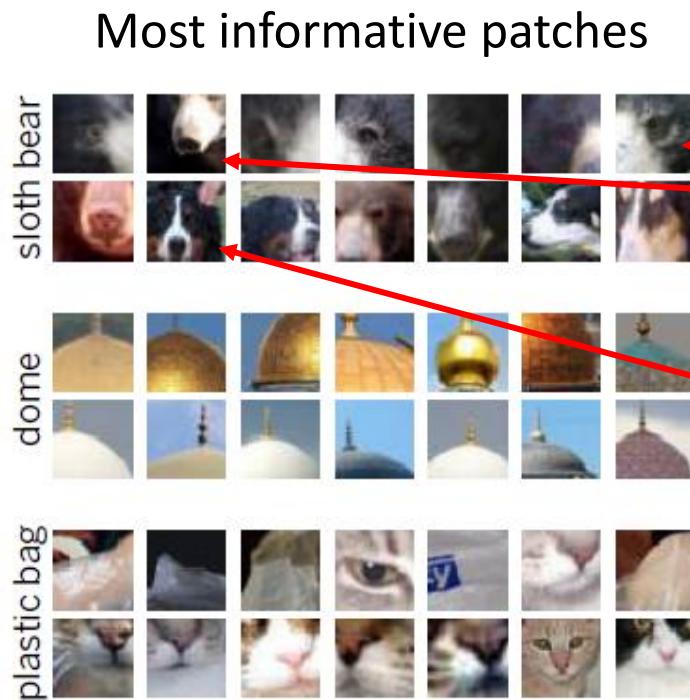


Image: Han, Zhong, Cao, & Zhang (2017)

Visualising classes



BagNet, trained on ImageNet

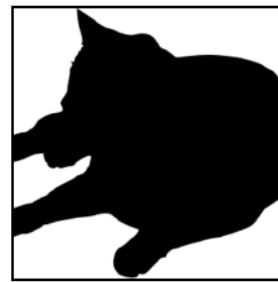
Shape and texture



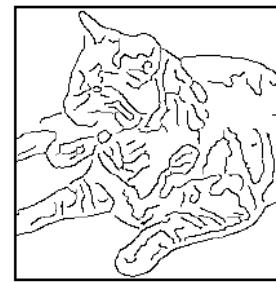
original



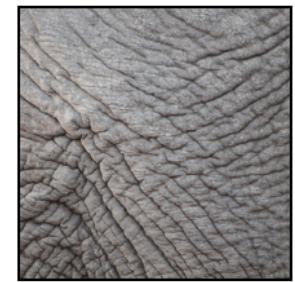
greyscale



silhouette

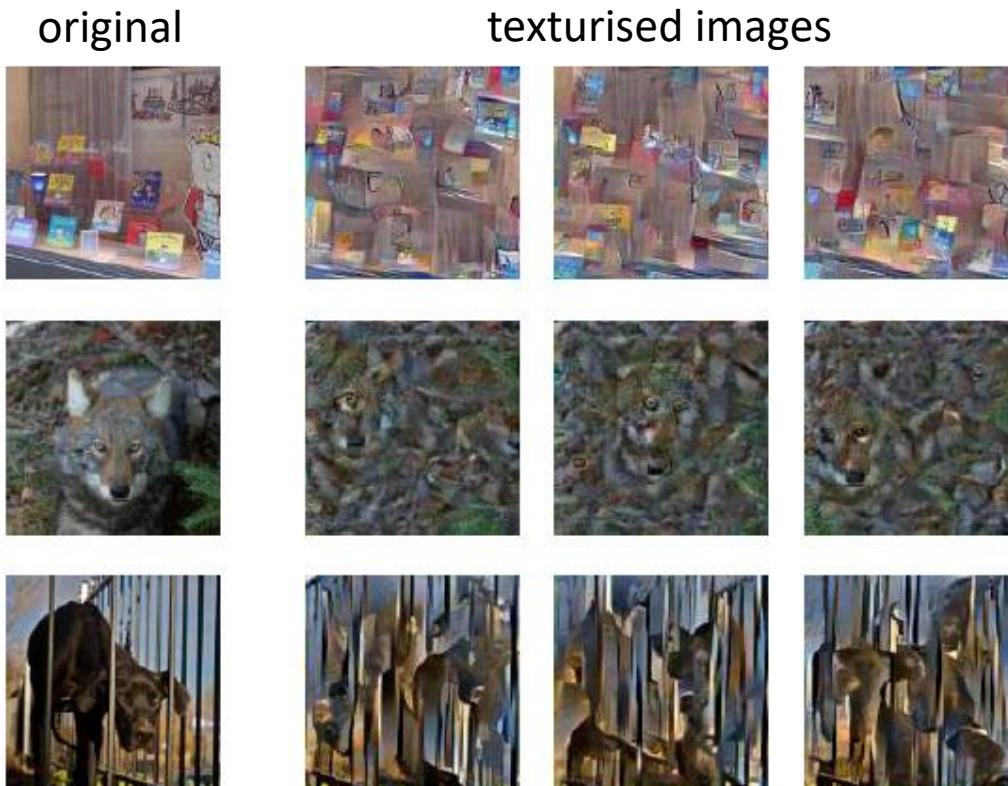


edges



texture

Shape and texture



VGG-16, trained on
ImageNet

Performance drop:
 $90\% \rightarrow 79\%$

Generalisation

- Models are very sensitive to some types of noise

Train on regular images:



Can recognize:



Can't recognize:

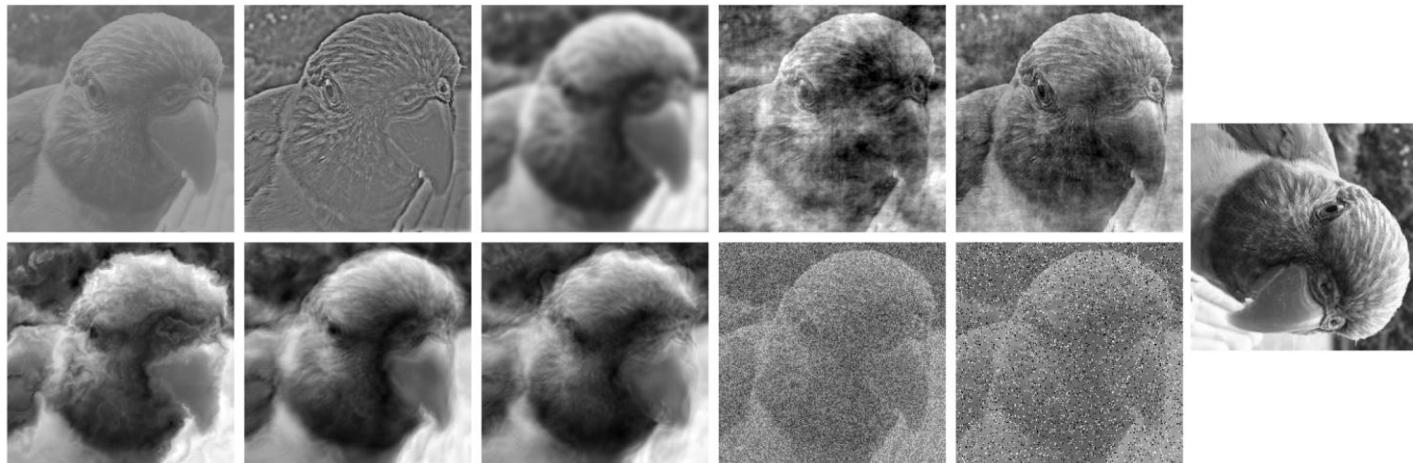


Image: Geirhos, Temme, Rauber, Schütt, Bethge, & Wichmann (2018)

Generalisation

Evaluation condition

	colour	88.5	96.7	90.6	50.0	83.1	86.1	84.2	90.8	10.4	8.1	97.9	95.4	72.3	93.0	91.1	92.4	94.9	10.2	11.2	95.5	95.9
	greyscale	86.6	87.8	95.6	94.1	86.2	93.2	87.8	90.5	10.3	9.8	94.0	96.8	96.2	93.3	95.7	94.3	90.9	11.4	12.8	94.8	95.1
	contrast (5%)	47.6	13.1	14.2	89.4	19.6	39.8	17.1	10.2	28.6	29.0	46.3	51.7	95.1	50.5	79.1	59.4	45.2	34.6	37.9	90.9	88.2
	low-pass (std=7)	48.5	18.9	16.1	16.4	78.4	11.9	16.0	9.8	6.9	6.6	16.0	18.6	14.4	87.2	20.5	13.8	13.5	7.1	9.3	74.7	74.9
	high-pass (std=0.7)	49.8	21.1	24.7	29.9	11.7	92.6	27.7	8.3	10.4	20.6	25.1	22.8	29.2	25.0	94.3	27.5	28.3	18.9	19.8	91.4	90.7
	phase noise (90°)	57.4	23.3	28.3	31.2	27.0	46.6	81.4	24.4	7.4	8.9	30.8	31.4	30.6	31.4	43.4	87.4	24.1	7.8	7.6	82.9	82.6
	rotation (90°)	78.5	36.5	43.3	39.9	31.8	40.4	37.7	89.0	8.5	8.0	38.5	41.9	40.3	35.2	40.1	40.5	89.0	8.3	8.8	80.1	80.5
	salt-and-pepper noise (0.2)	NA	6.1	6.4	5.8	7.9	6.2	6.2	6.4	79.4	6.2	6.2	6.1	6.3	5.4	5.8	5.7	6.2	89.6	6.2	78.6	13.6
	uniform noise (0.35)	45.6	6.2	7.3	6.9	9.0	7.3	6.2	6.0	10.2	80.3	84.6	83.3	85.0	84.6	83.7	82.5	83.8	85.4	89.8	11.0	71.5
human observers	A ¹	A ²	A ³	A ⁴	A ⁵	A ⁶	A ⁷	A ⁸	A ⁹	B ¹	B ²	B ³	B ⁴	B ⁵	B ⁶	B ⁷	B ⁸	B ⁹	C ¹	C ²		

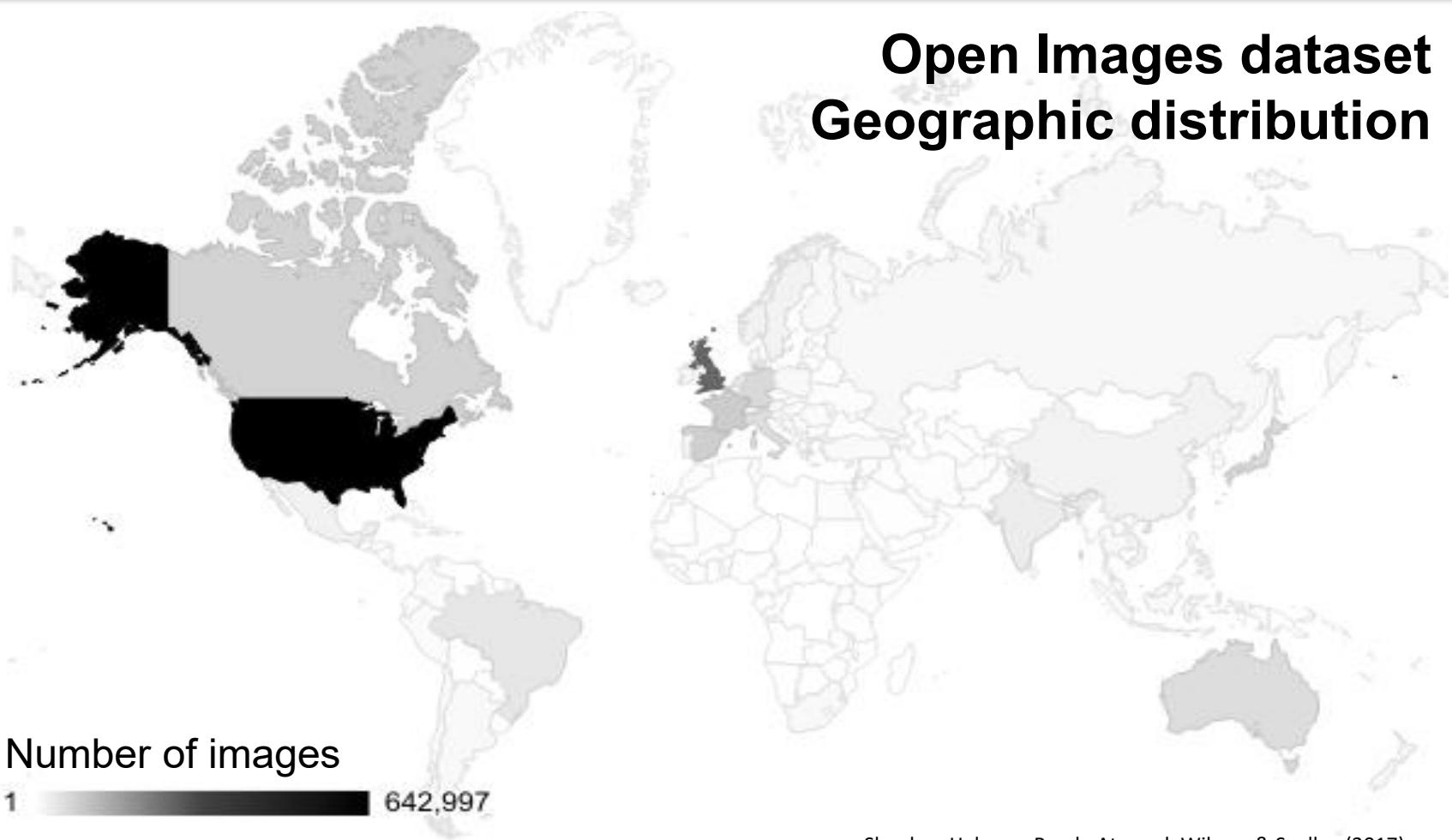


= manipulation included in training data

Model

ResNet-50 retrained on ImageNet subset with various manipulations

Cultural bias?





Ground truth: Soap

Azure: food, cheese, bread, cake, sandwich
Clarifai: food, wood, cooking, delicious, healthy
Google: food, dish, cuisine, comfort food, spam
Amazon: food, confectionary, sweets, burger
Watson: food, food product, turmeric, seasoning
Tencent: food, dish, matter, fast food, nutrient



Ground truth: Soap

UK, 1890 \$/month

Azure: toilet, design, art, sink
Clarifai: people, faucet, healthcare, lavatory, wash closet
Google: product, liquid, water, fluid, bathroom accessory
Amazon: sink, indoors, bottle, sink faucet
Watson: gas tank, storage tank, toiletry, dispenser, soap dispenser
Tencent: lotion, toiletry, soap dispenser, dispenser, after shave

Invariance / tolerance

- CNNs are tolerant to variation included in the training data
- But may not generalise well to variation outside their training distribution
- SOTA image training datasets have biases (geographic, wealth) which impact generalisation