

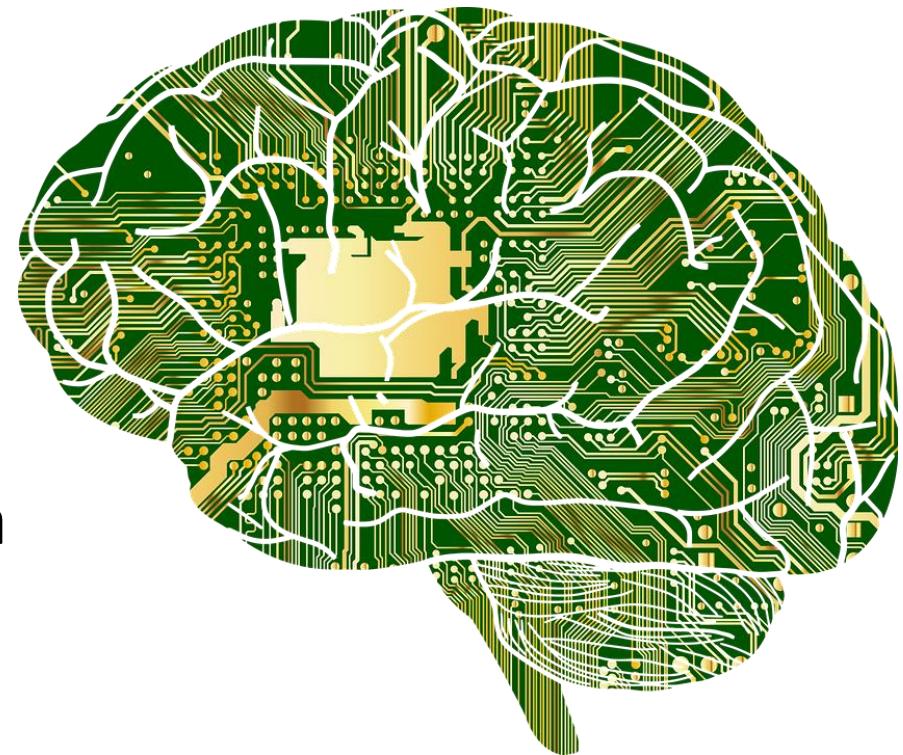
Deep Learning II

Semester 2, 2022

Kris Ehinger

Understanding CNNs

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



Outline

- ImageNet results
- CNN visualisation
- Invariance, generalisation

Learning outcomes

- Explain methods that can be used to evaluate computer vision algorithms
- Explain strengths and weaknesses of CNNs for image recognition
- Evaluate computer vision algorithms in terms of their invariance (or tolerance) to image variation

ImageNet results

What affects performance?

Easiest
classes:



Hardest
classes:



Russakovsky et al. (2014)

Model visualisation

Feature representation

“AlexNet”: Krizhevsky, Sutskever, & Hinton (2012)

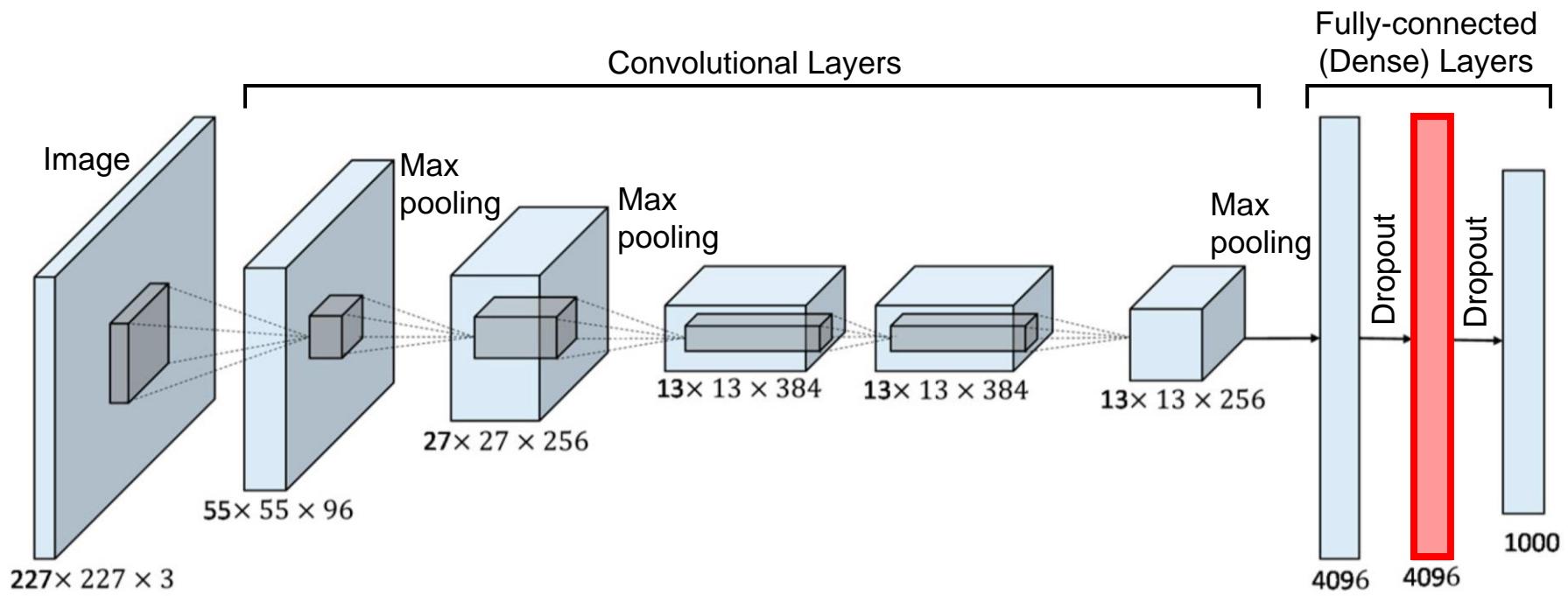
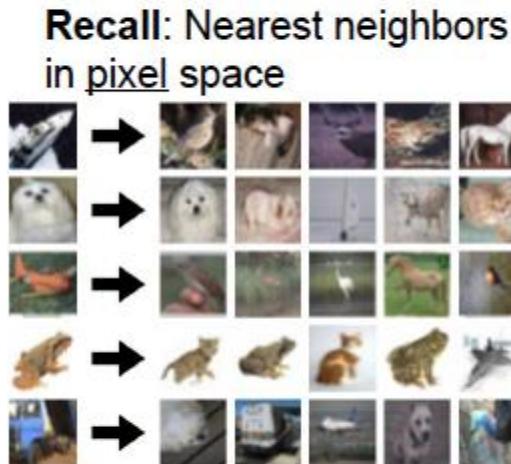


Image: Han, Zhong, Cao, & Zhang (2017)

Visualising features

- 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?)
- What do individual neuron responses represent?

Nearest neighbours



Test image L2 Nearest neighbors in feature space



Feature space visualisation

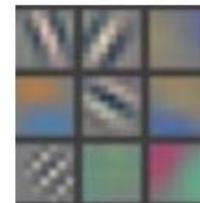
- Options for dimensionality reduction
- PCA (principal component analysis)
 - Show the dimensions with the most variance
 - Simple but often hard to interpret since only a few dimensions can be visualised simultaneously
- t-SNE (t-distributed stochastic neighbor embedding)
 - Flatten high-dimensional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



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

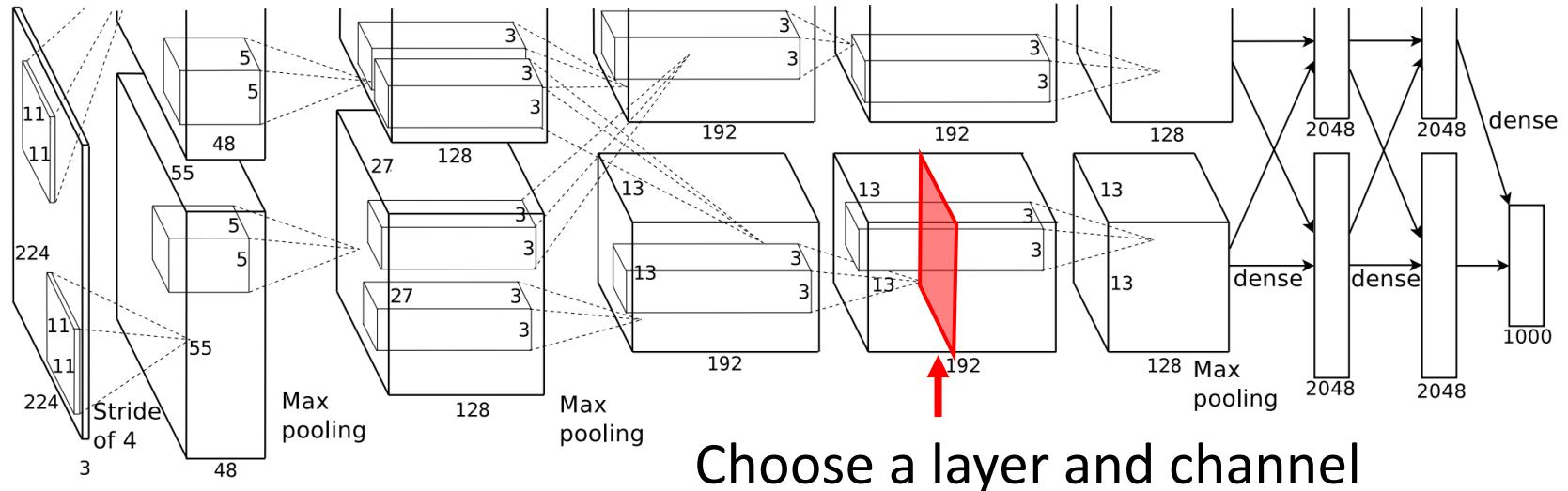
Visualising convolutional kernels

- Visualising the first convolutional layer kernels is easy because the input channels are RGB
- Visualising layer conv kernels is harder because the channels are high-dimensional and represent complex features



Layer 1

Maximally activating patches

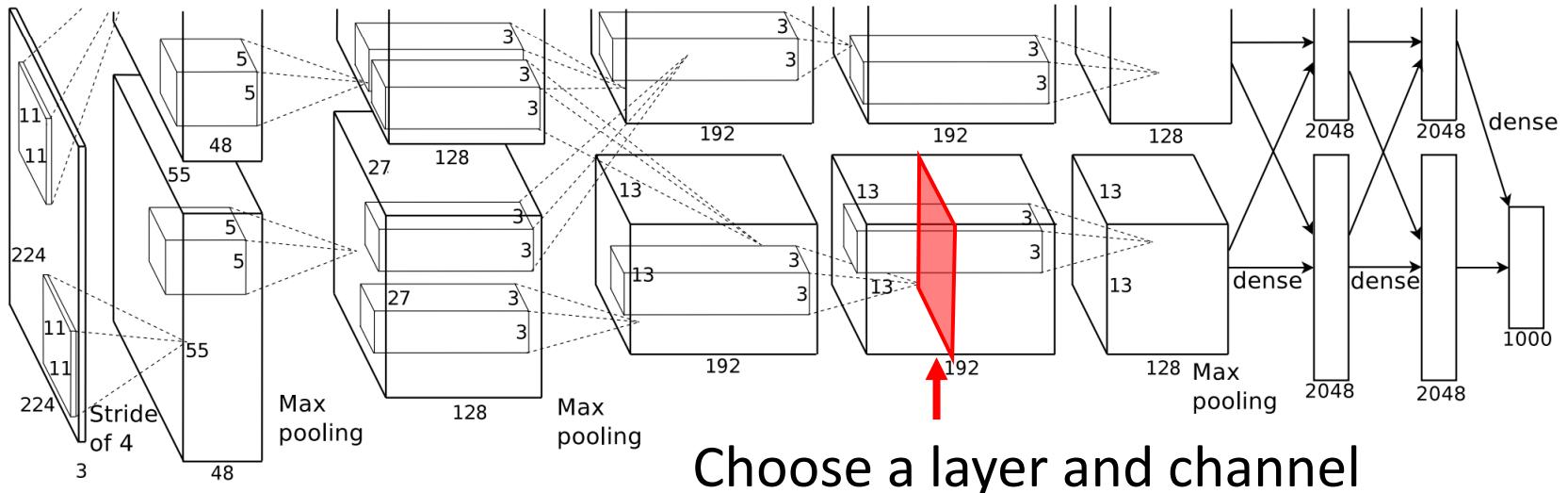


Run many images through the network and find patches that give the highest response in this channel:



Zeiler & Fergus (2014)

Guided backprop



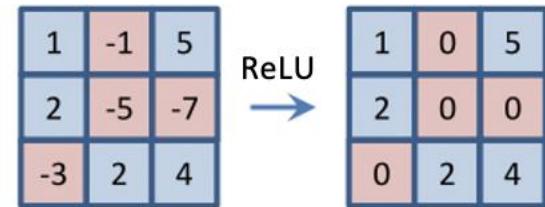
- Compute gradient of neuron value with respect to image pixels
- Which pixels matter most for correct classification?

Zeiler & Fergus (2014)

Guided backprop

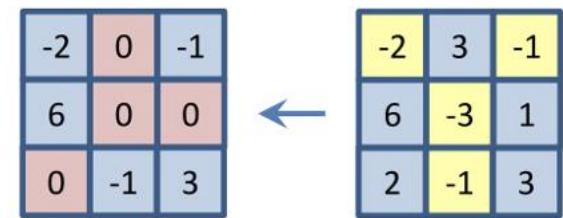
- ReLU activation function means neurons with response < 0 are set to 0

Forward pass



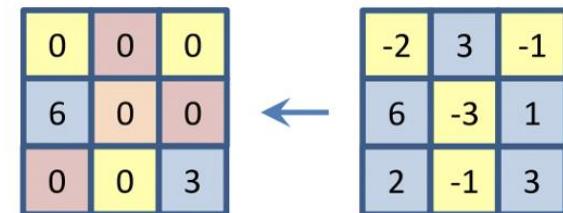
- Traditional backprop does not pass back gradient when neuron response is 0

Backward pass:
backpropagation

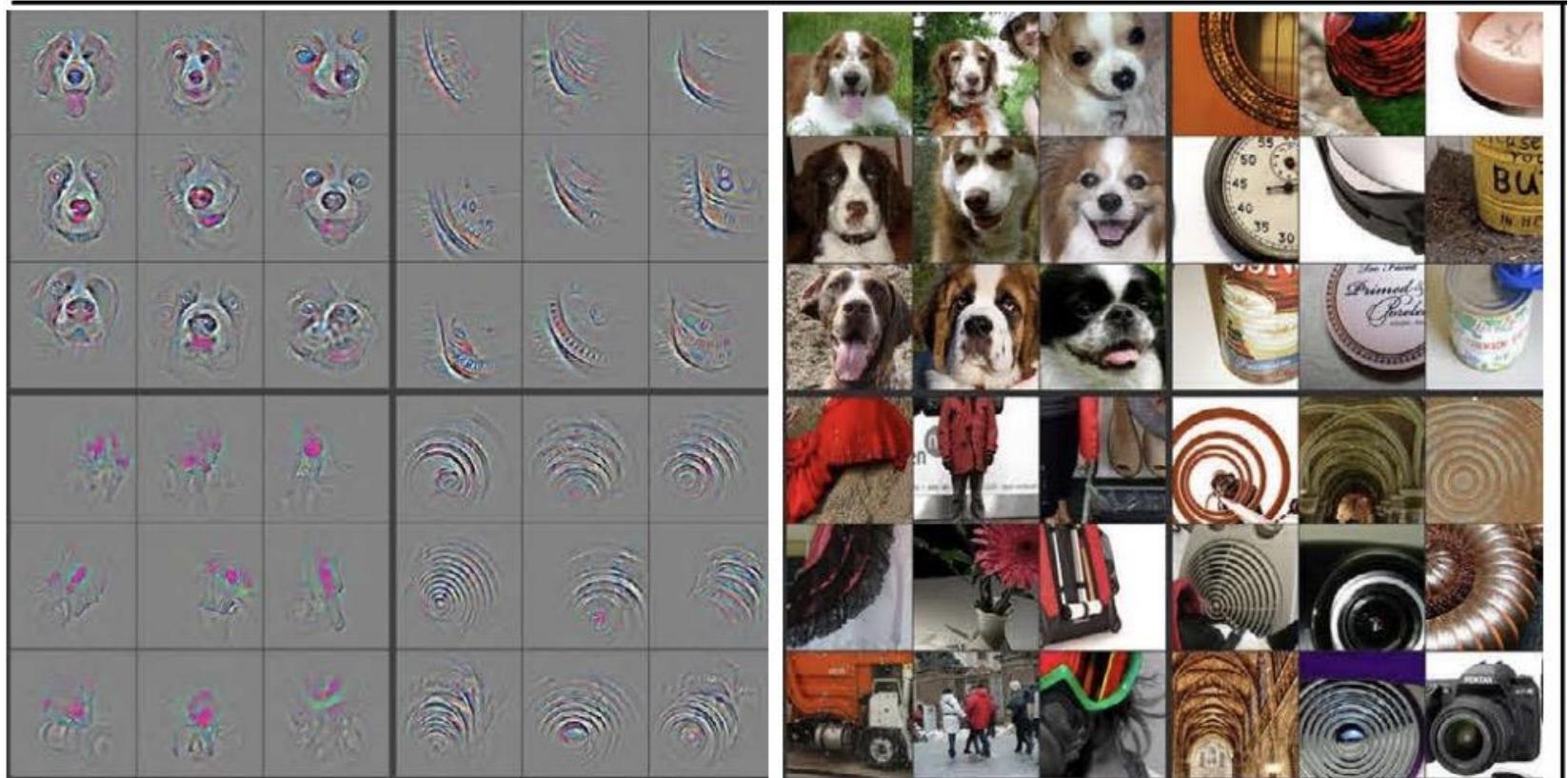


- Guided backprop also does not pass back negative gradient

Backward pass:
*guided
backpropagation*



Visualising convolutional kernels



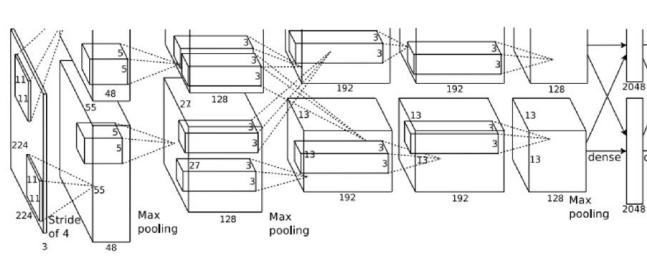
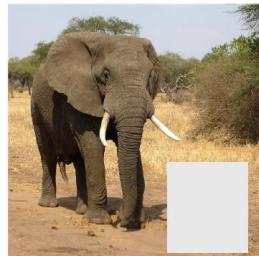
Zeiler & Fergus (2014)

Visualising image regions

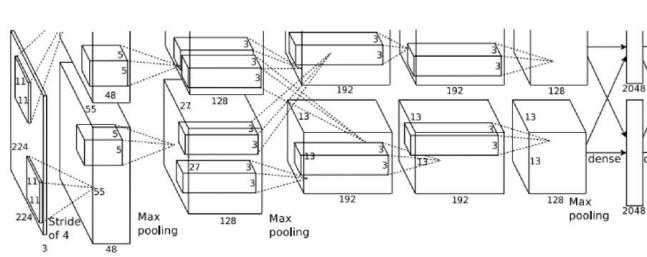
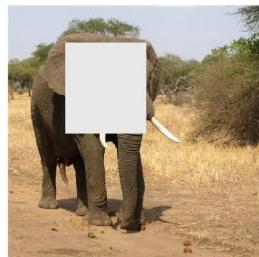
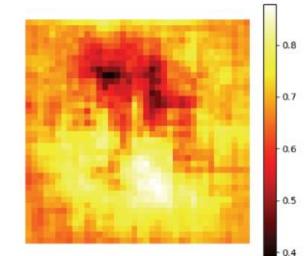
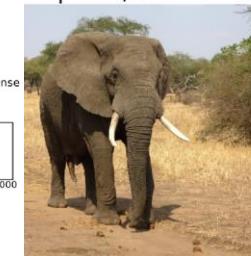
- What parts of an image are most important for determining the class label?
 - Can help show what features the model uses to determine class
 - Can help debug problems (e.g., using background to classify object, label confusion when there are multiple objects)

Visualising image regions

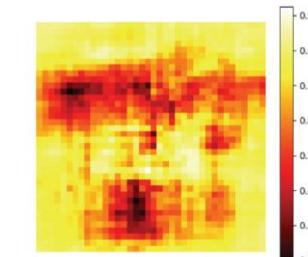
- Occlusion method: mask image and see how much class probability changes



African elephant, *Loxodonta africana*



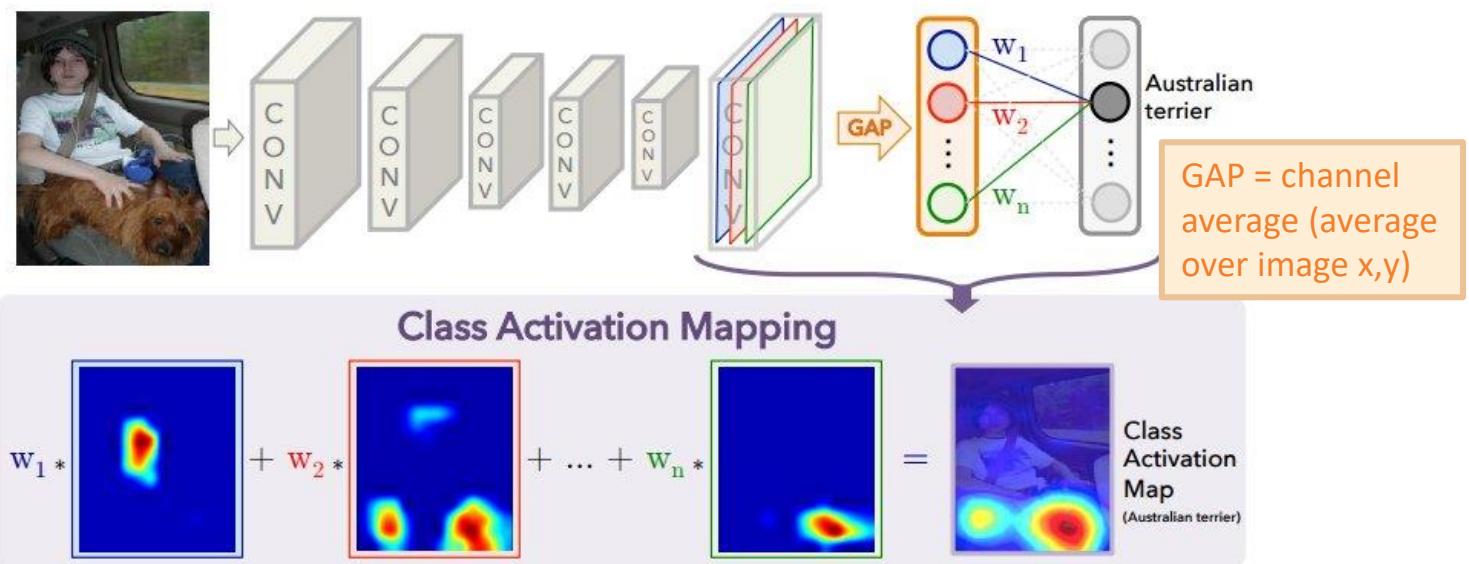
go-kart



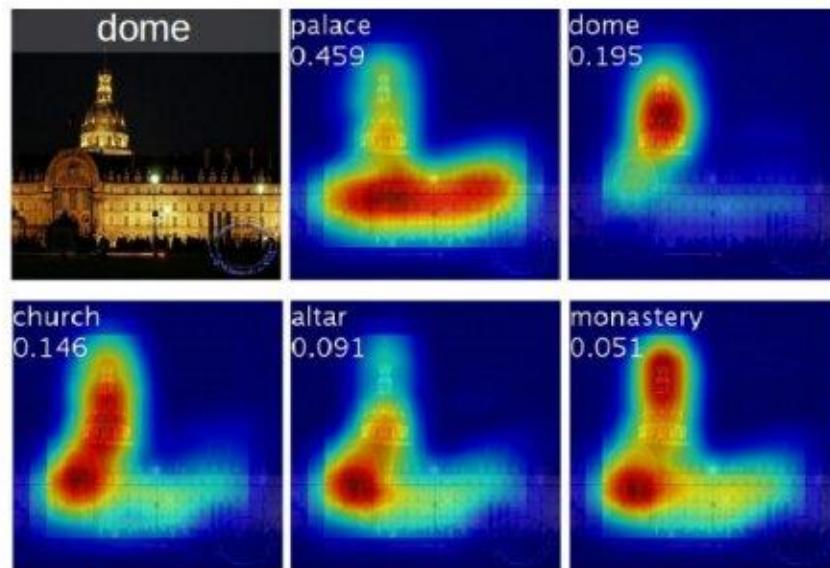
Zeiler & Fergus (2014)

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

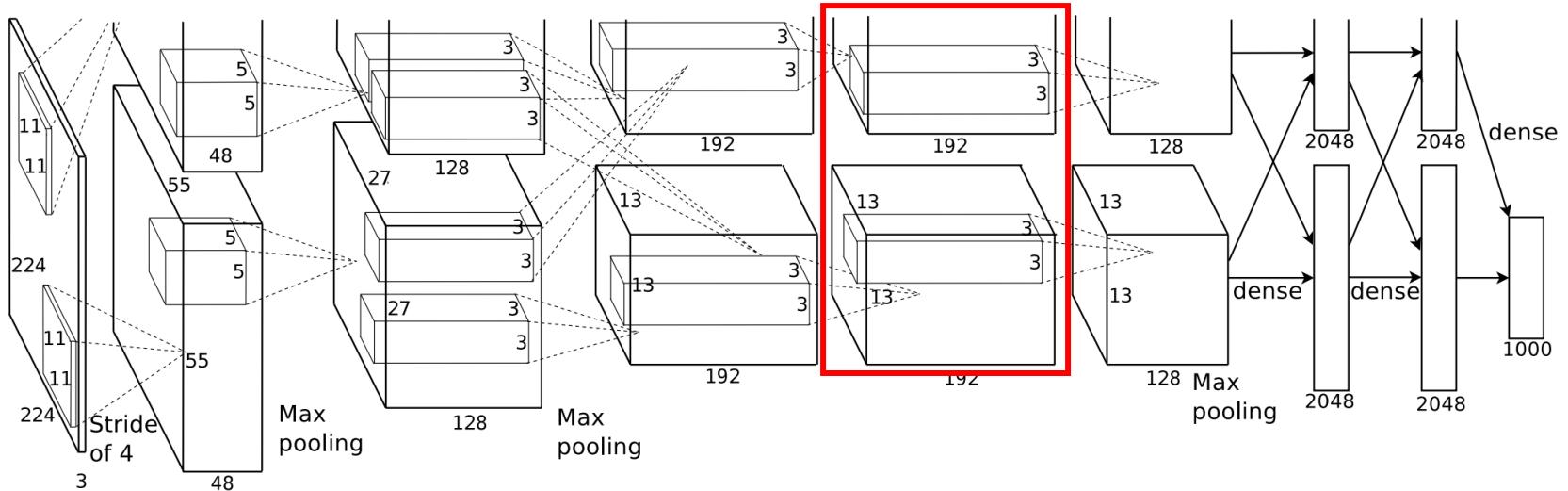


Class activation maps for one object class

CAM (Class Activation Mapping)

- Disadvantages of CAM:
 - Most models don't use GAP, so the GAP layer must be added to a pretrained network and then finetuned
 - Only allows visualisation of the last layer
- More flexible alternative: Grad-CAM (Gradient-Weighted Class Activation Mapping)

Grad-CAM

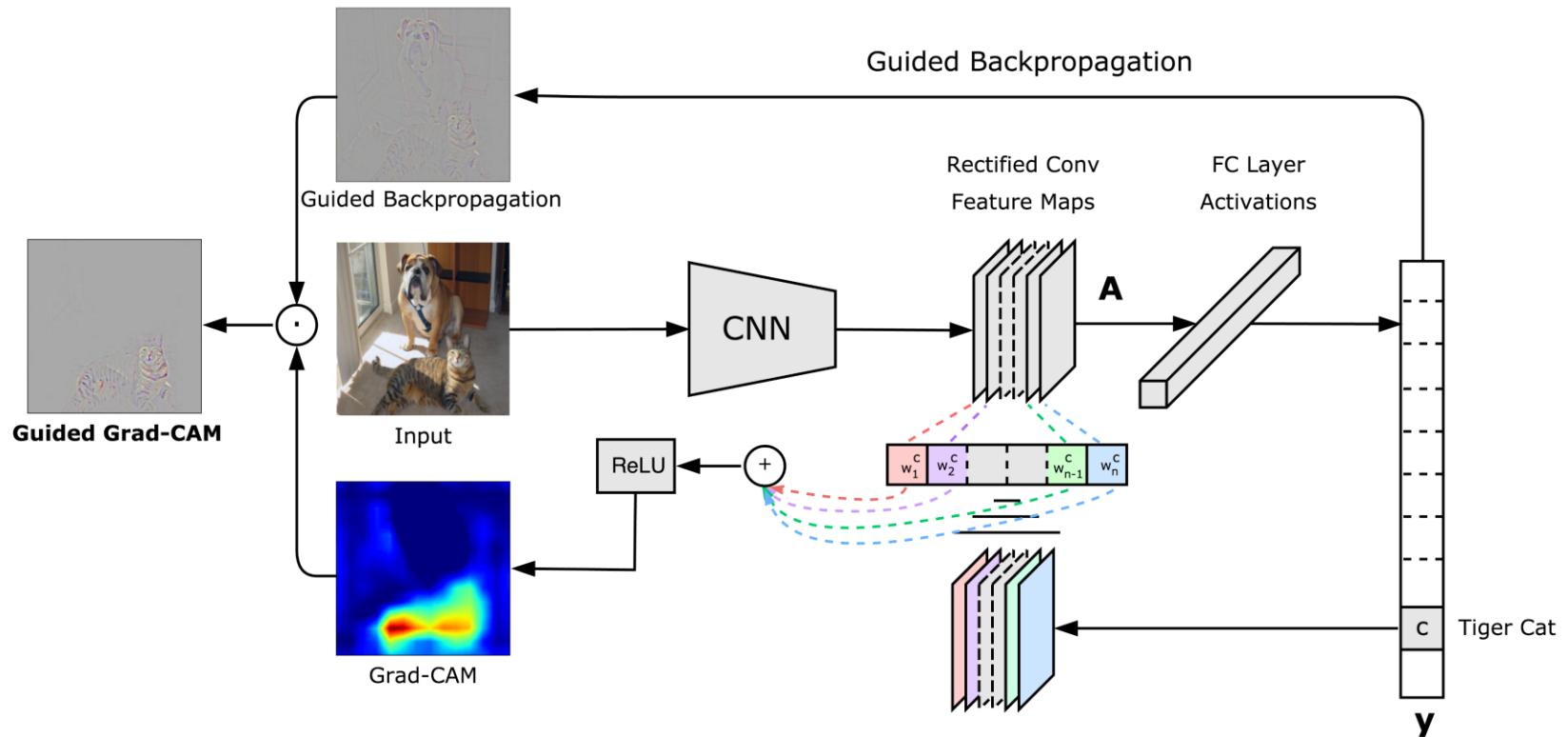


- Take response from some layer $A \in \mathbb{R}^{H \times W \times K}$
- Compute gradient of class score with respect to layer response
- Global Average Pool (average over image x,y) the gradients to get a vector of weights α_k (1 weight per channel)
- Compute activation map $ReLU(\sum_k \alpha_k A_{h,w,k})$

Selvaraju, et al. (2016)

Grad-CAM

Gradient-Weighted Class Activation Mapping (Grad-CAM)



Selvaraju, et al. (2016)

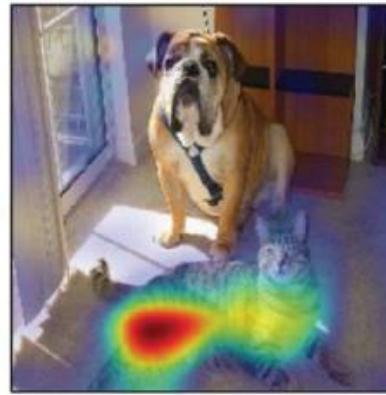
Grad-CAM



(a) Original Image



(b) Guided Backprop ‘Cat’



(c) Grad-CAM ‘Cat’



(g) Original Image



(h) Guided Backprop ‘Dog’



(i) Grad-CAM ‘Dog’

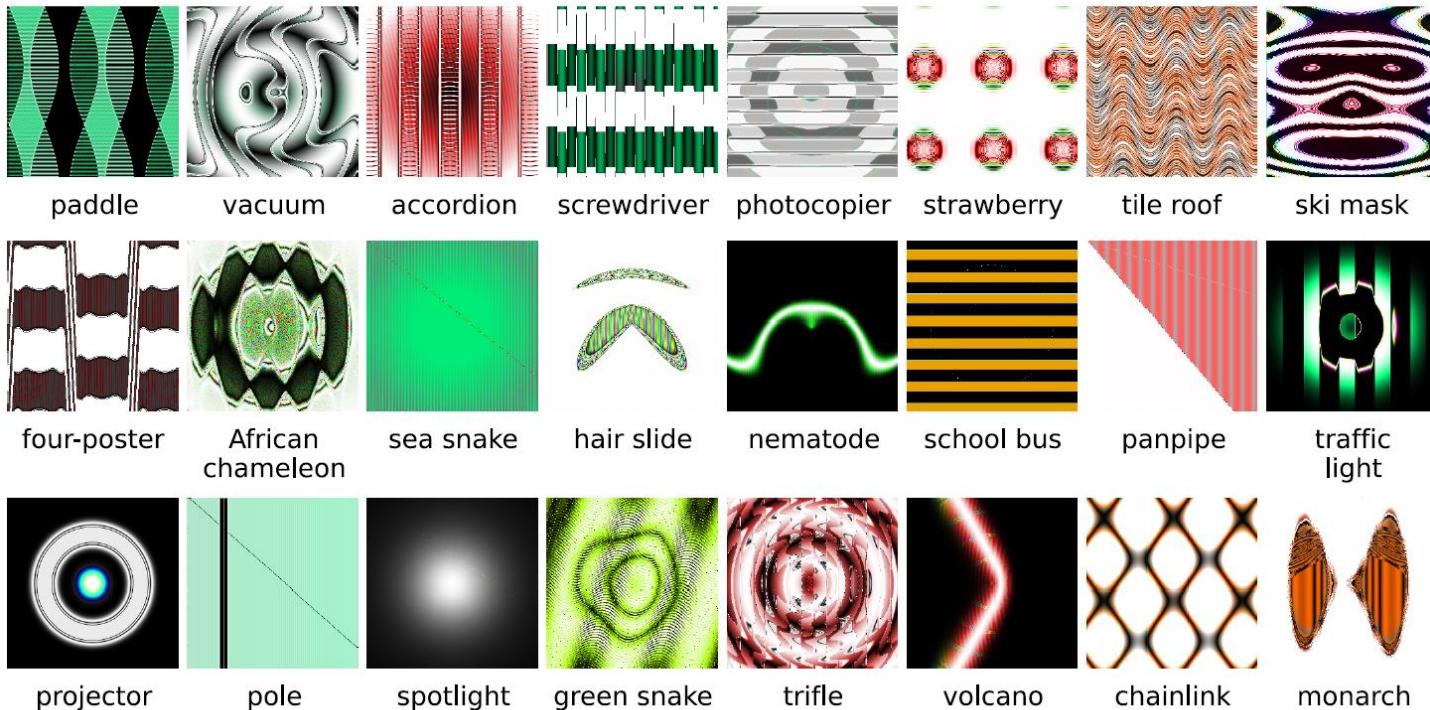
Selvaraju, et al. (2016)

Visualising classes

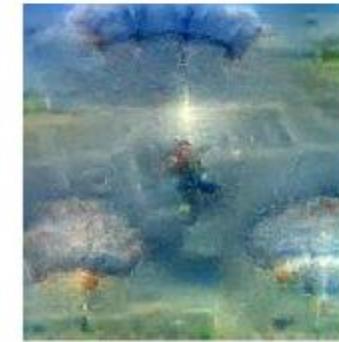
- Usually based on gradient ascent – synthesize image that maximises class label response
 - Initialise an image with zeros or small random noise
 - Run image through network, compute gradient
 - Update *image pixels* in a direction that minimises loss
- Problem: there are many possible arrays of pixels 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
 - Feature space visualisations
 - Visualising image regions that support a class decision
 - Class visualisations
- Each method provides different information, so it's usually best to try multiple approaches

Invariance & generalisation

Invariance / tolerance

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

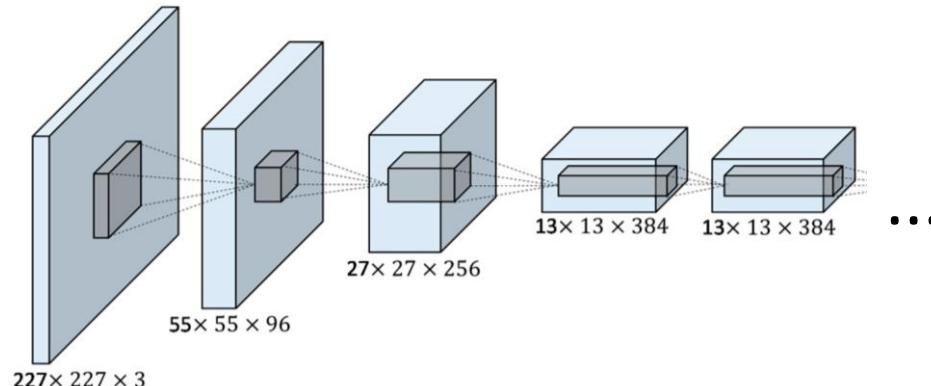
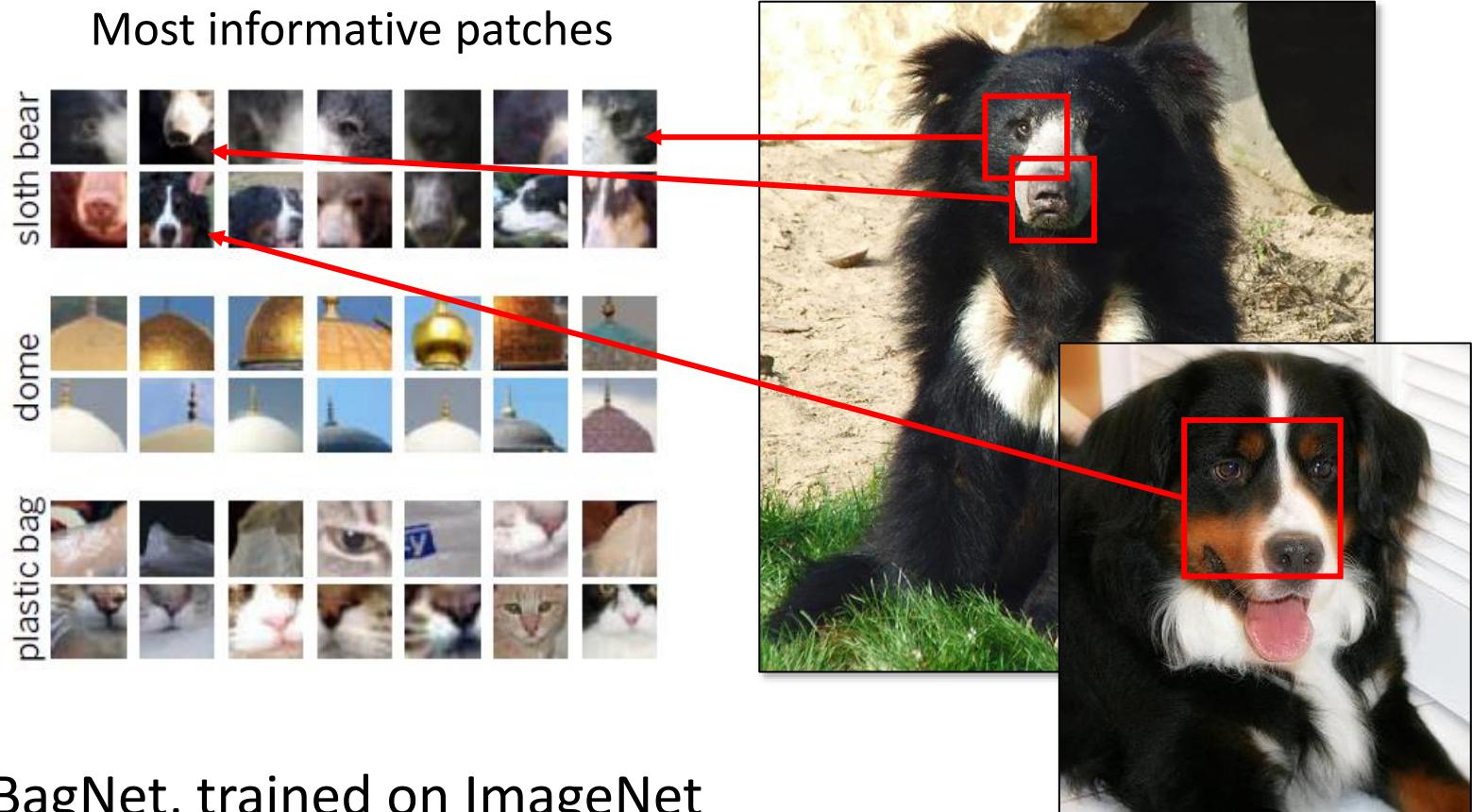


Image: Han, Zhong, Cao, & Zhang (2017)

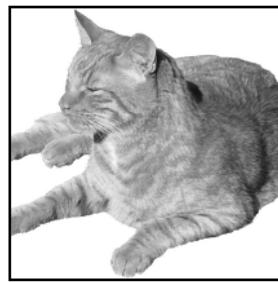
Visualising classes



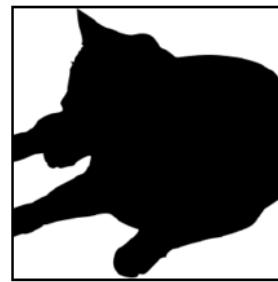
Shape and texture



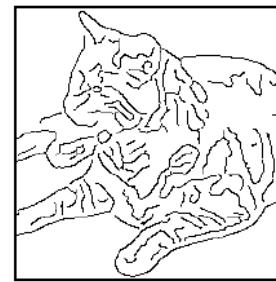
original



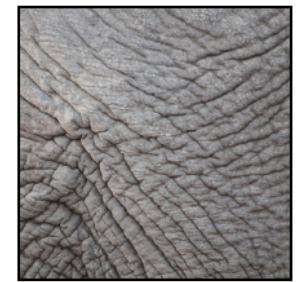
greyscale



silhouette

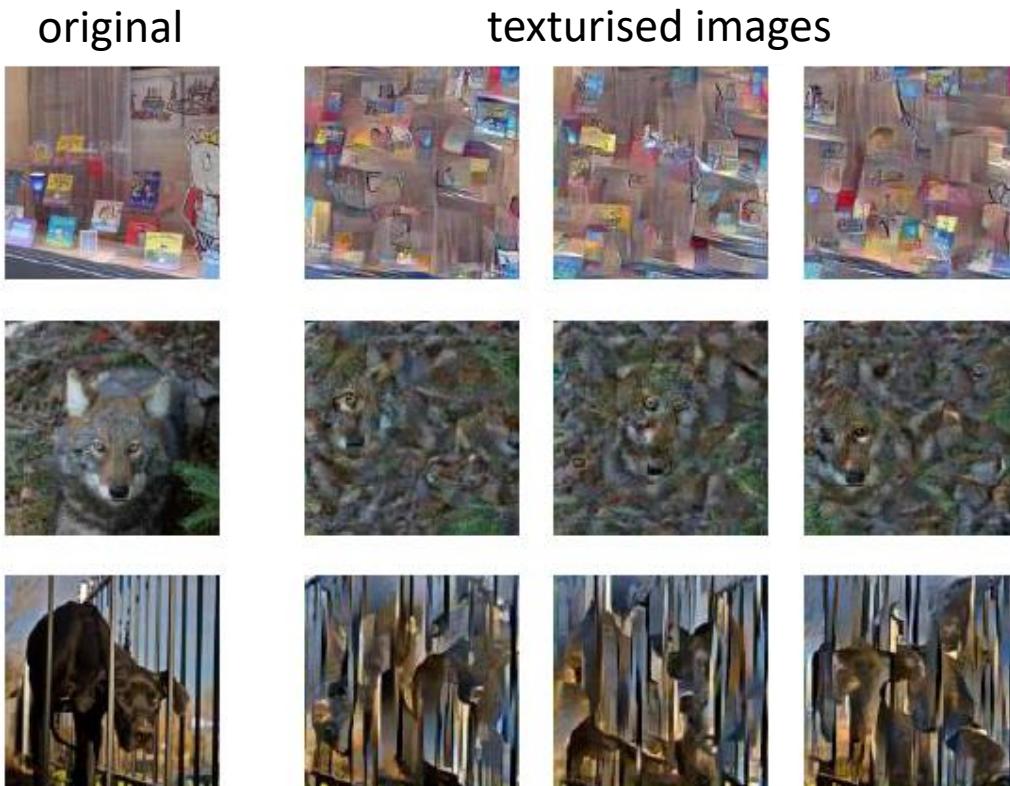


edges



texture

Shape and texture



VGG-16, trained on
ImageNet

Performance drop:
90% → 79%

Generalisation

- Models are very sensitive to some types of noise

Train on regular images:



Can recognize:



Can't recognize:

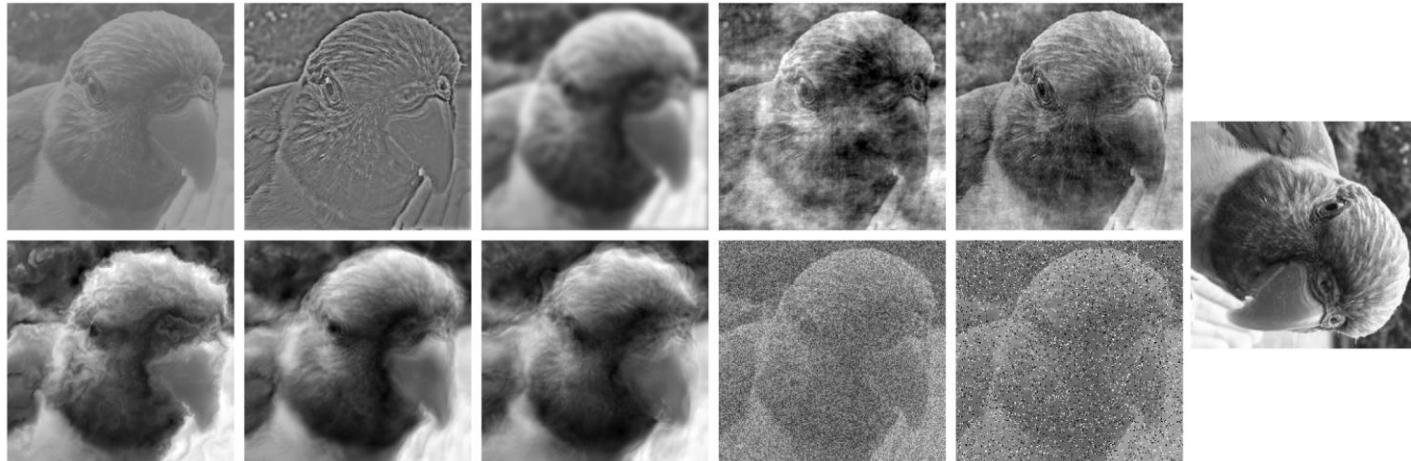


Image: Geirhos, Temme, Rauber, Schütt, Bethge, & Wichmann (2018)

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

Invariance / tolerance

- CNNs are tolerant to variation included in the training data
- But often not tolerant to variation that didn't appear in the training data
- Classification tends to rely on recognizing a few key features / local texture elements

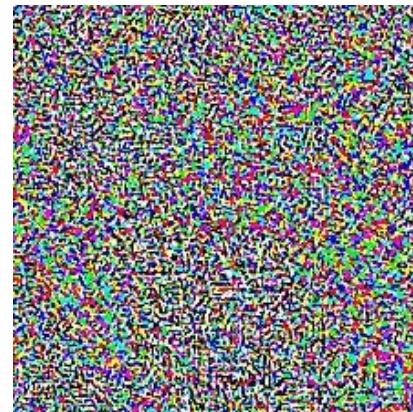
Adversarial images

- Adding small amounts of noise to an image can completely change the model's perception



Original image: “panda”
(57.7% confidence)

+



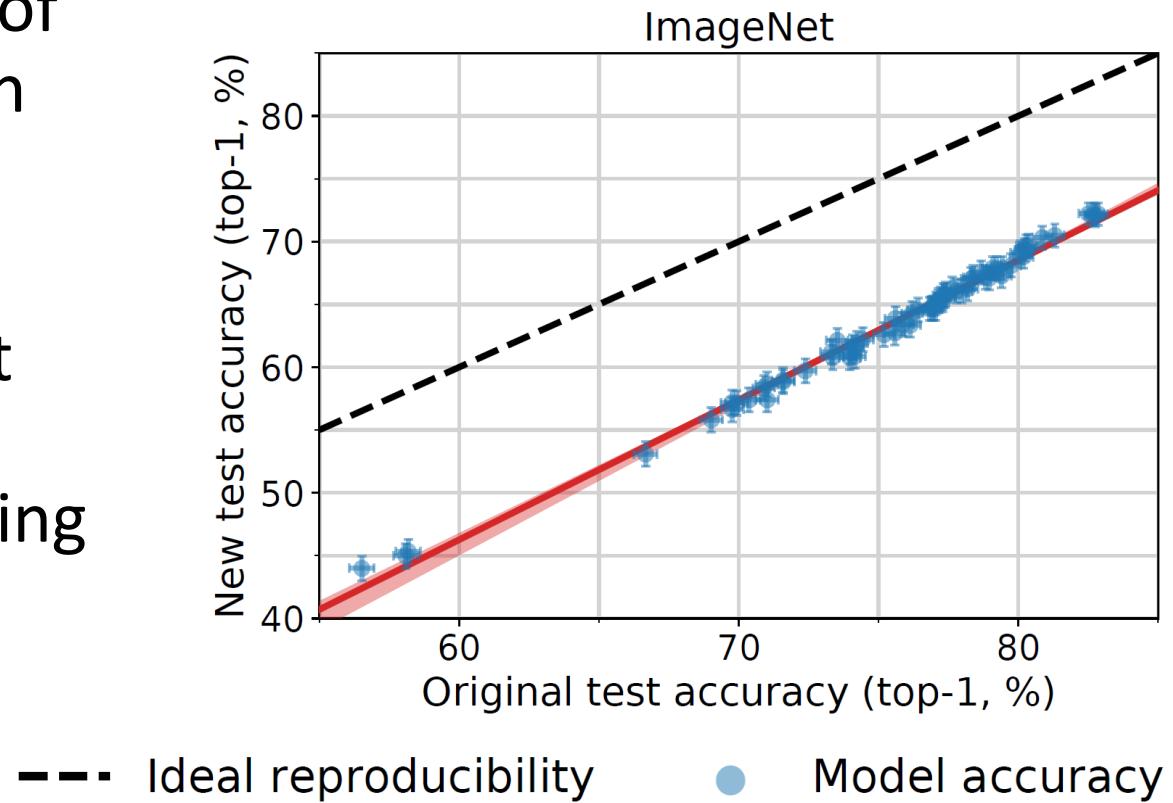
=



New image: “gibbon”
(99.3% confidence)

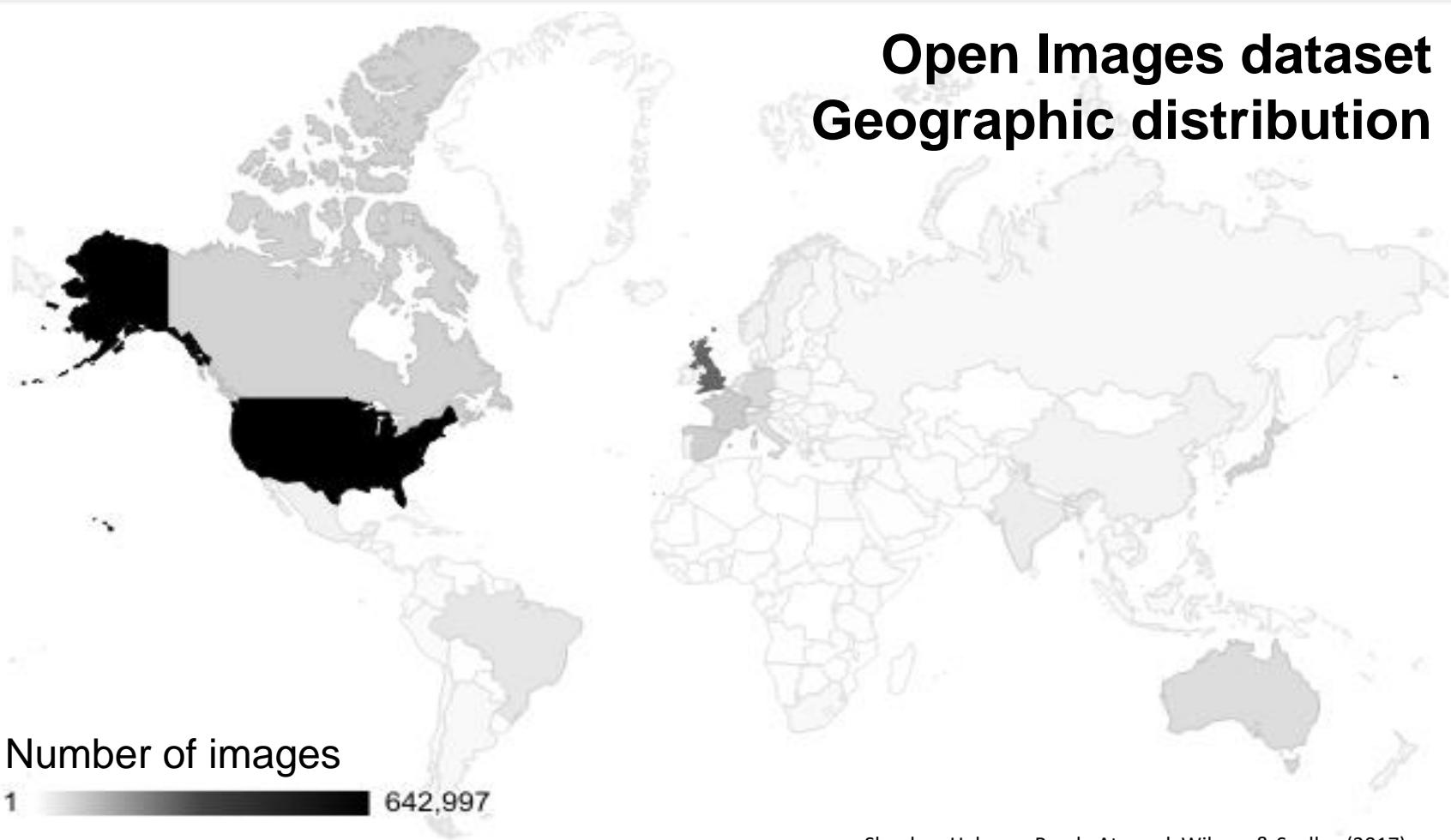
Generalisation

- Performance of top models on ImageNet vs. ImageNetV2
- Drop of about 10% suggests some overfitting to quirks of ImageNet



Geographic bias?

**Open Images dataset
Geographic distribution**



Summary

- ImageNet-trained features are used for a variety of visual tasks
- But ImageNet-trained CNNs have some issues:
 - Very sensitive to noise
 - Recognize local texture features but not global shape
 - Generalisation errors on very similar datasets
 - Biases due to dataset construction