
Machine Learning For Design

Lecture 4 - Machine Learning for Images / Part 2

Alessandro Bozzon

23/02/2022

mlfd-io@tudelft.nl
www.ml4design.com

Admin

-
- Very few questions for Week 2 :(
 - We will publish few quizzes for Week 2 today

 - First group assignment next week!
 - Deadline next Tuesday

**How do
humans see?**

Hubel and Wiesel, 1959

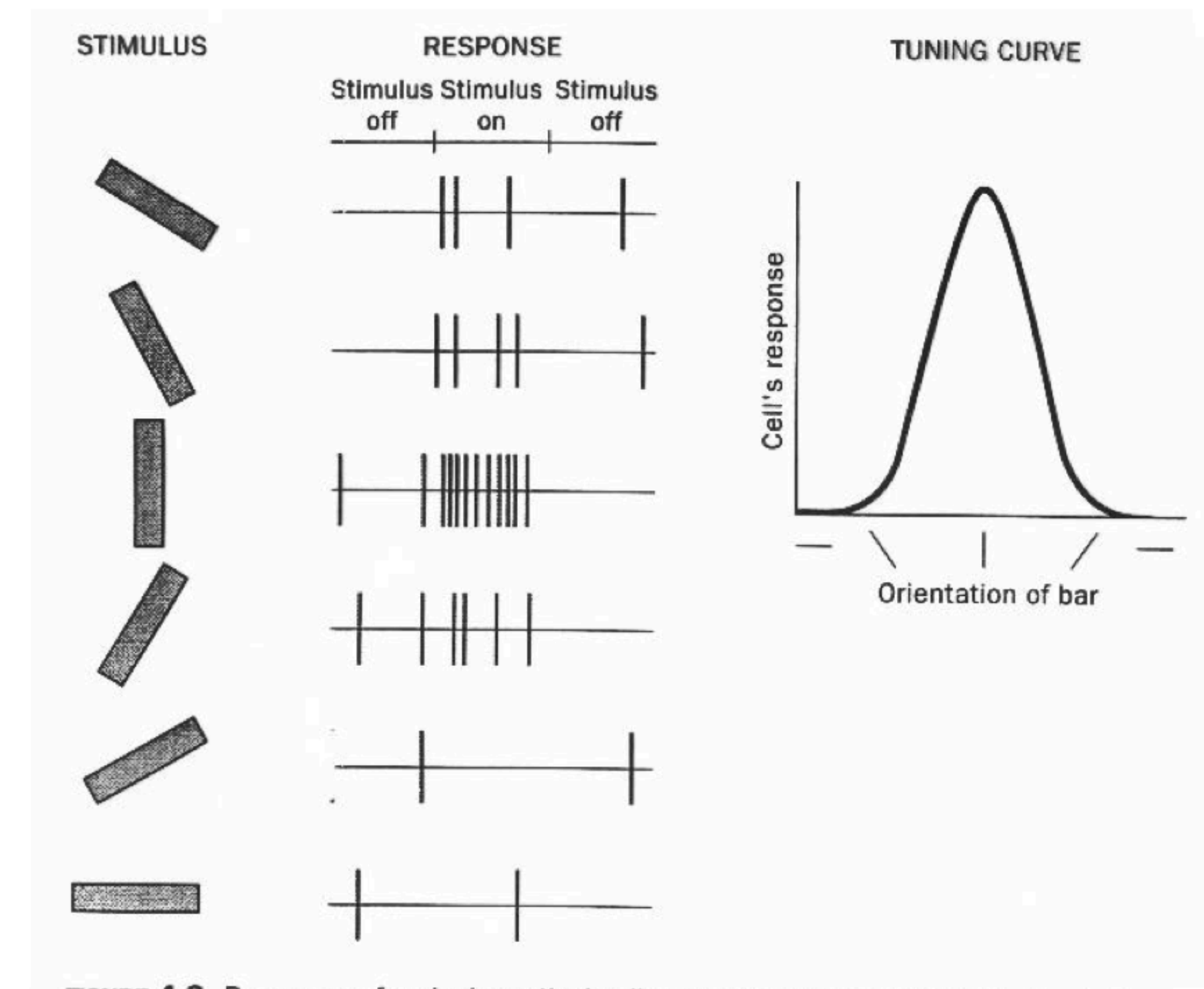
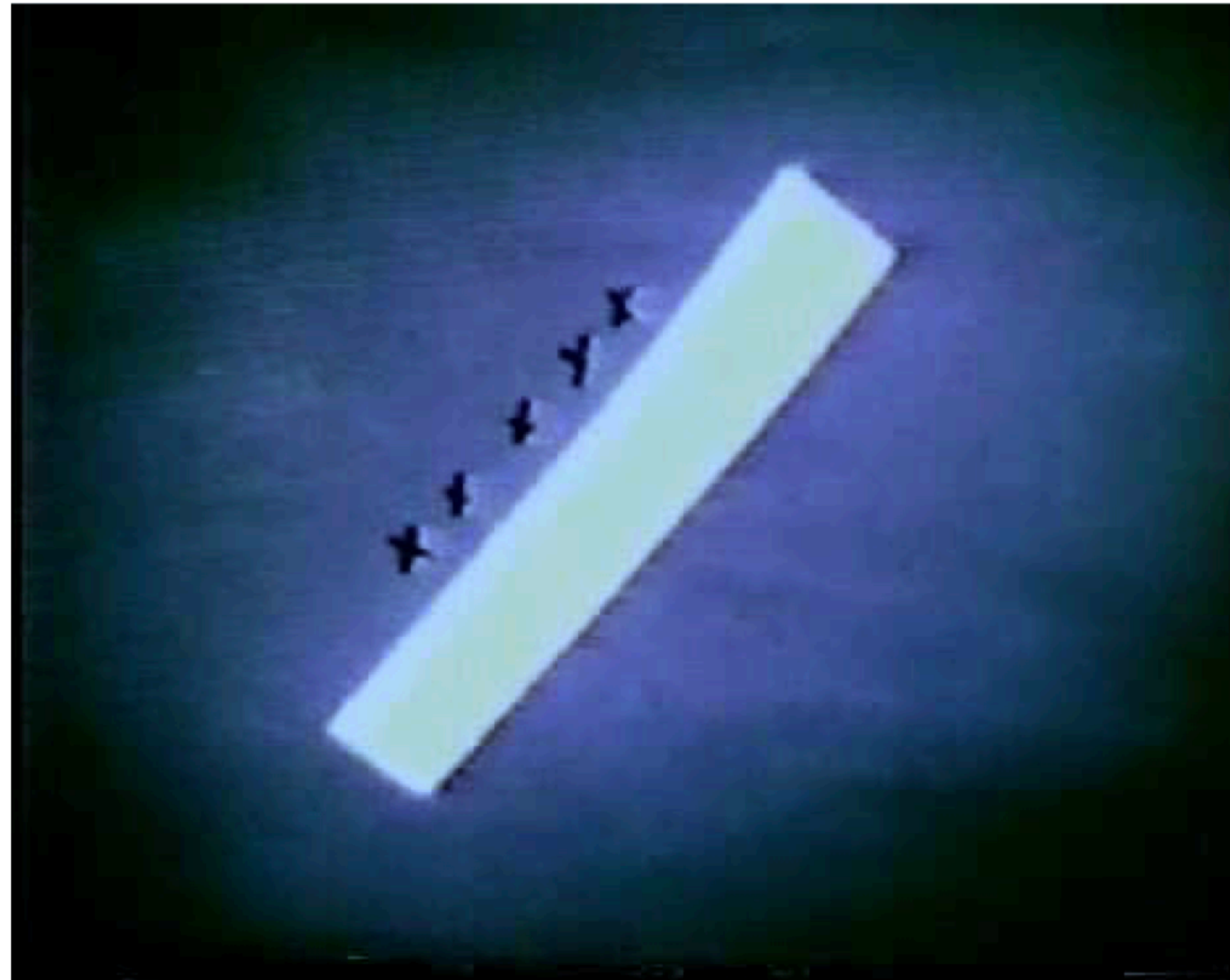
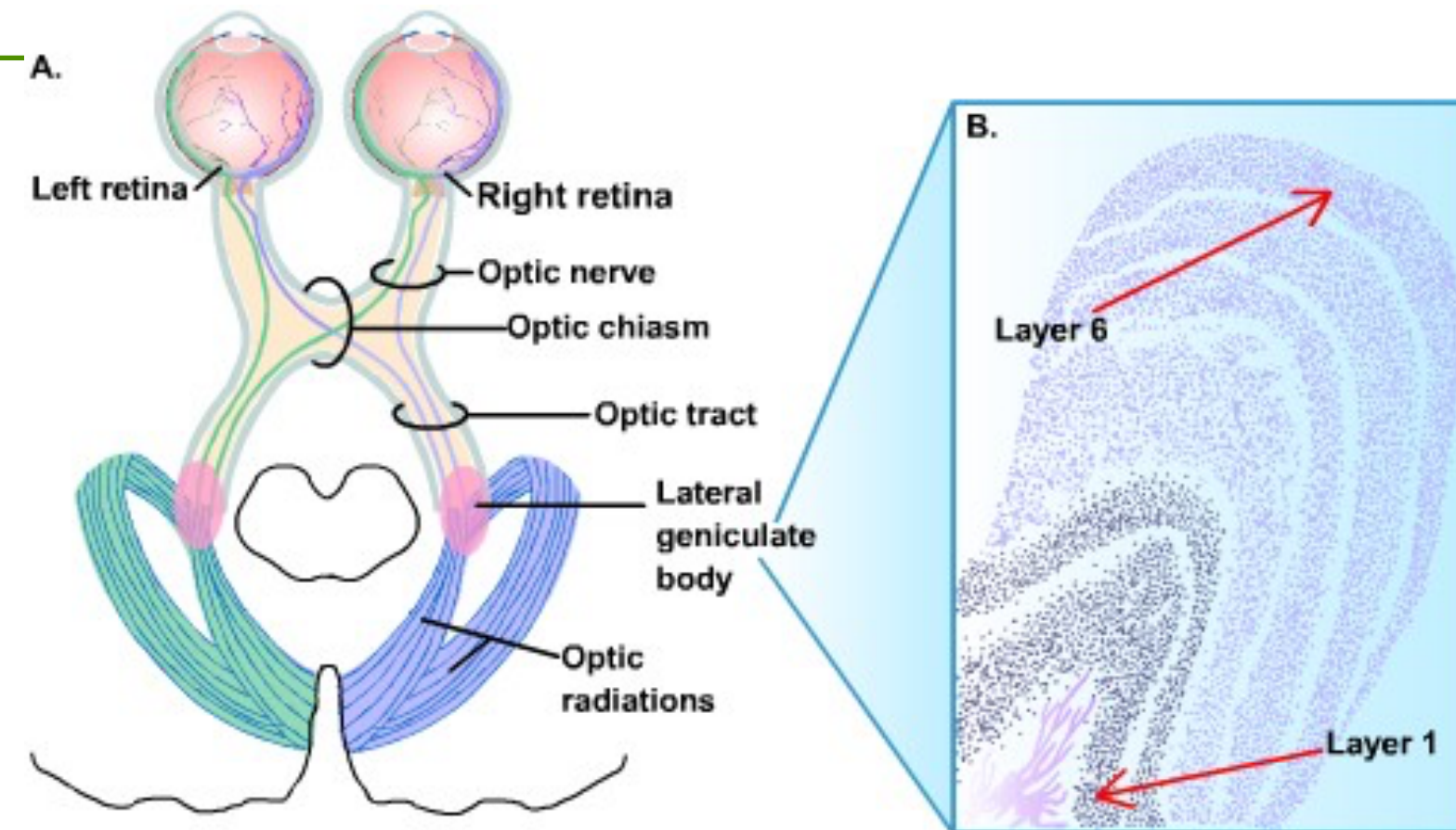


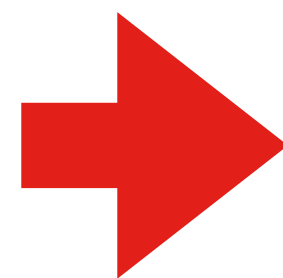
FIGURE 4.8 Response of a single cortical cell to bars presented at various orientations.

<https://www.youtube.com/watch?v=IOHayh06LJ4>

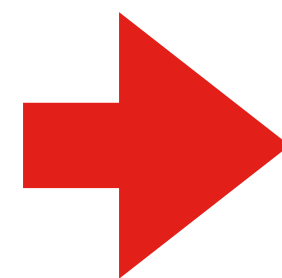
Neural Pathways



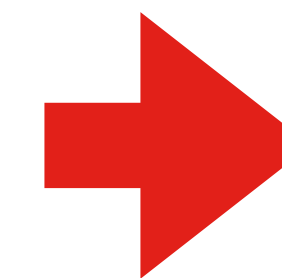
Edges



Simple Shapes



Complex Shapes

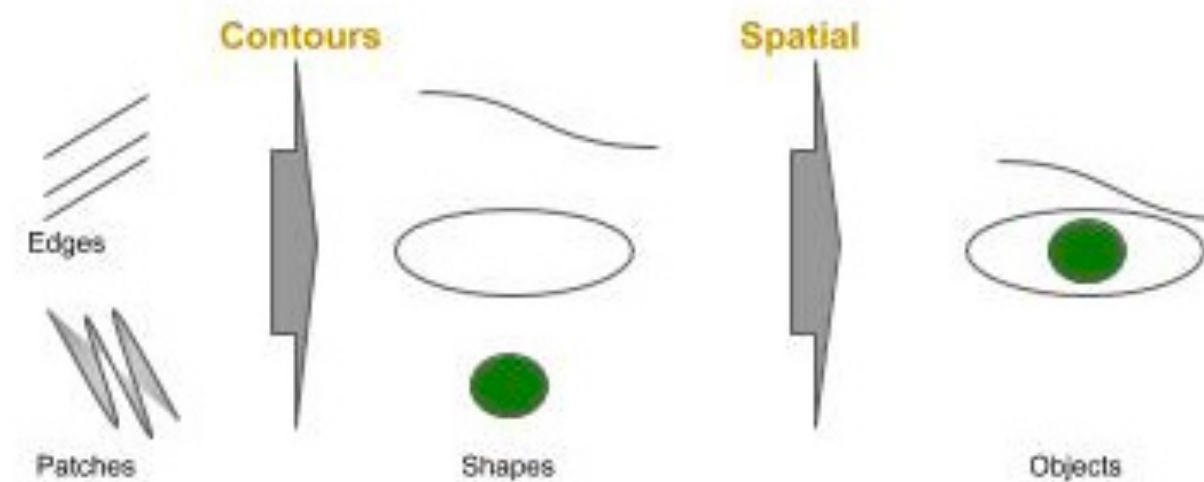


Faces and Objects



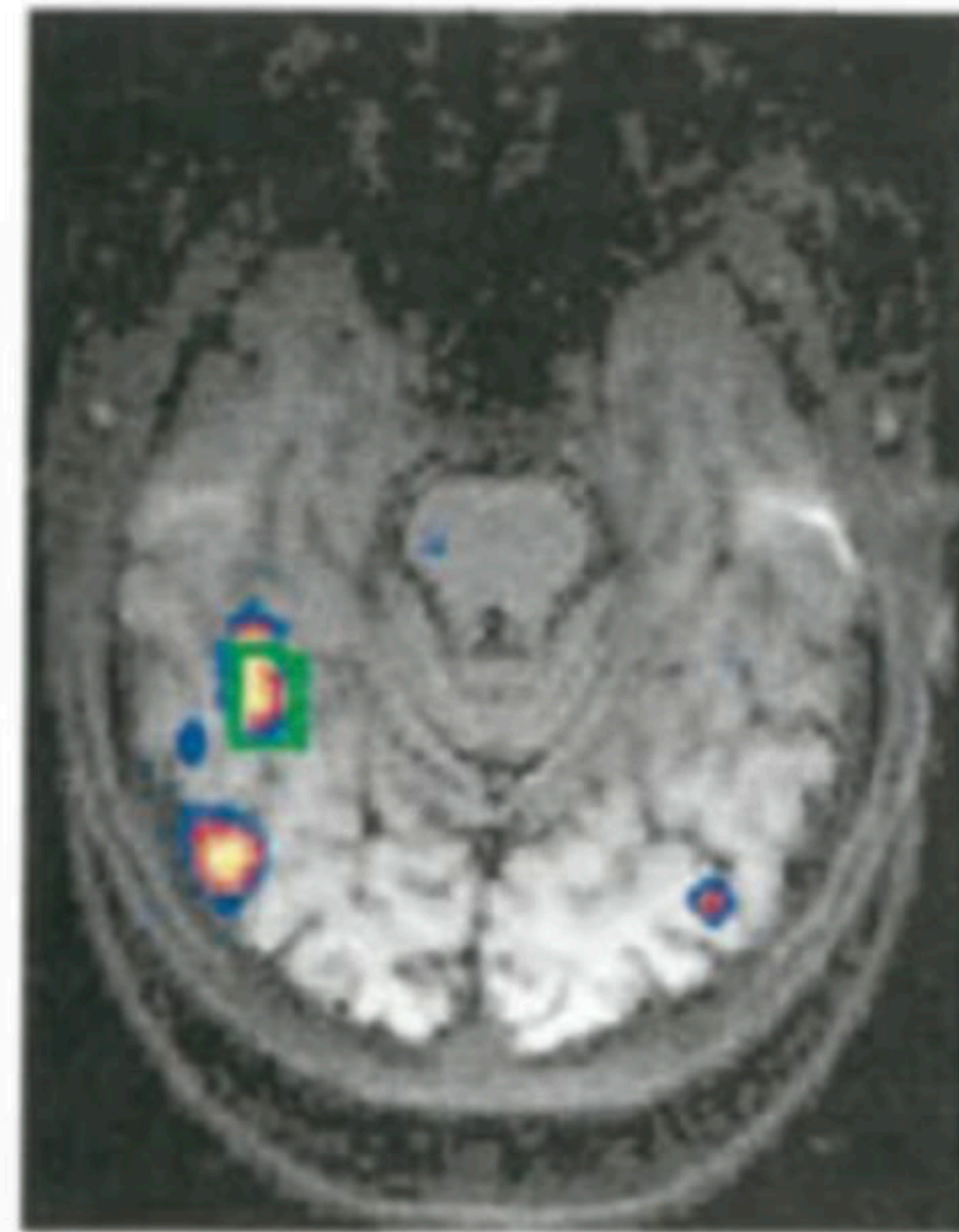
Lower layers

Upper layers



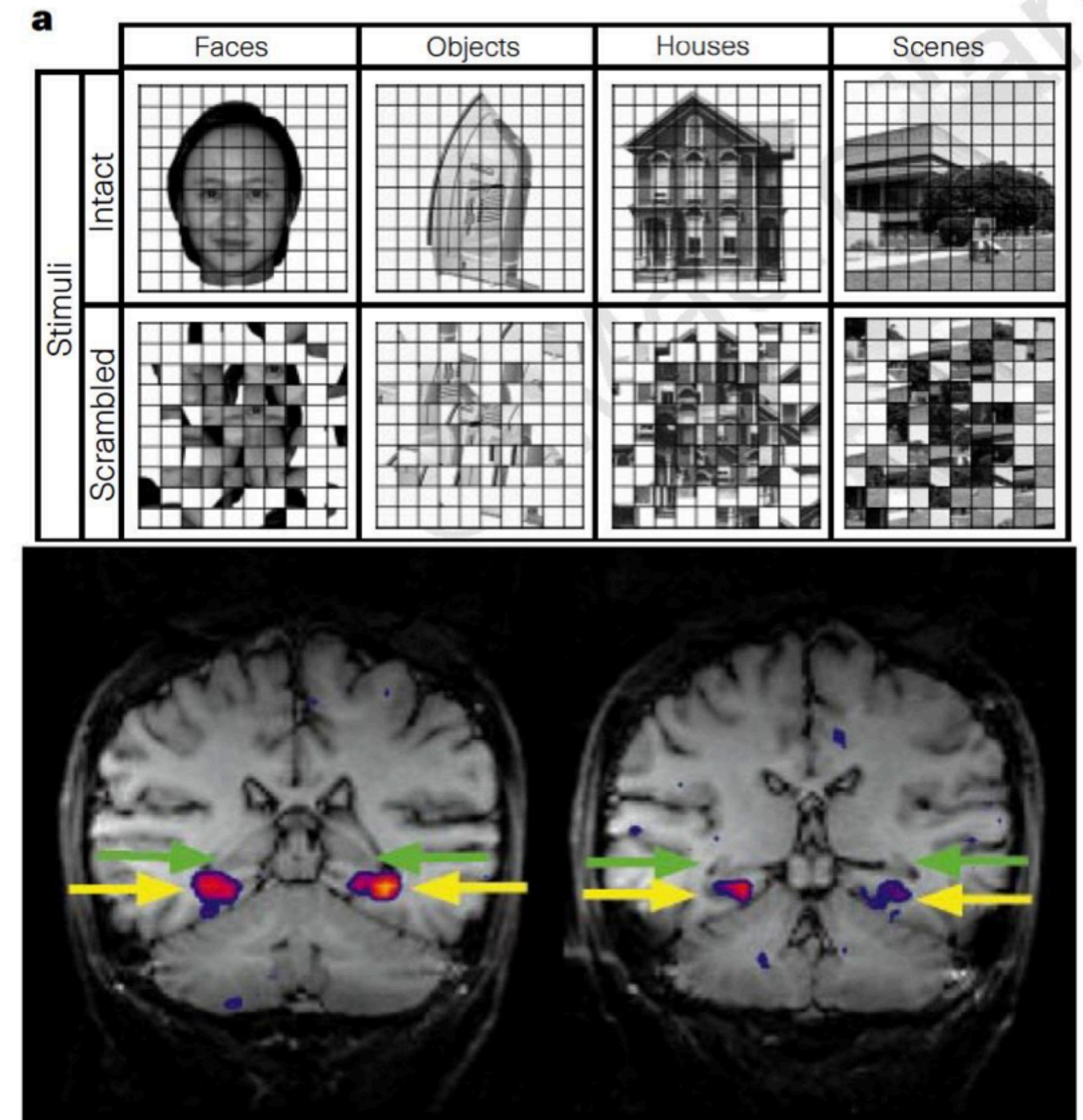
Neural Correlation of Objects & Scene Recognition

Faces > Houses



% signal change

Kanwisher et al. J. Neuro. 1997



Epstein & Kanwisher, Nature, 1998

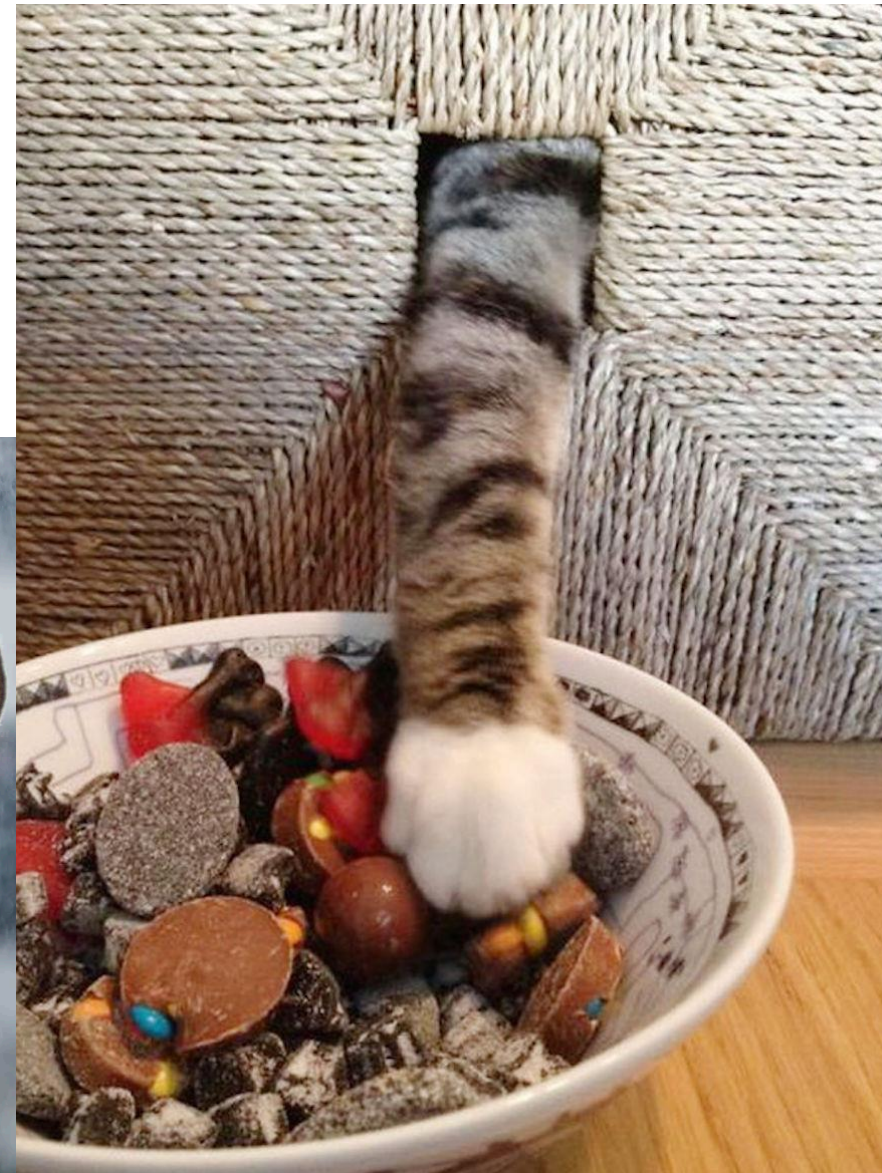
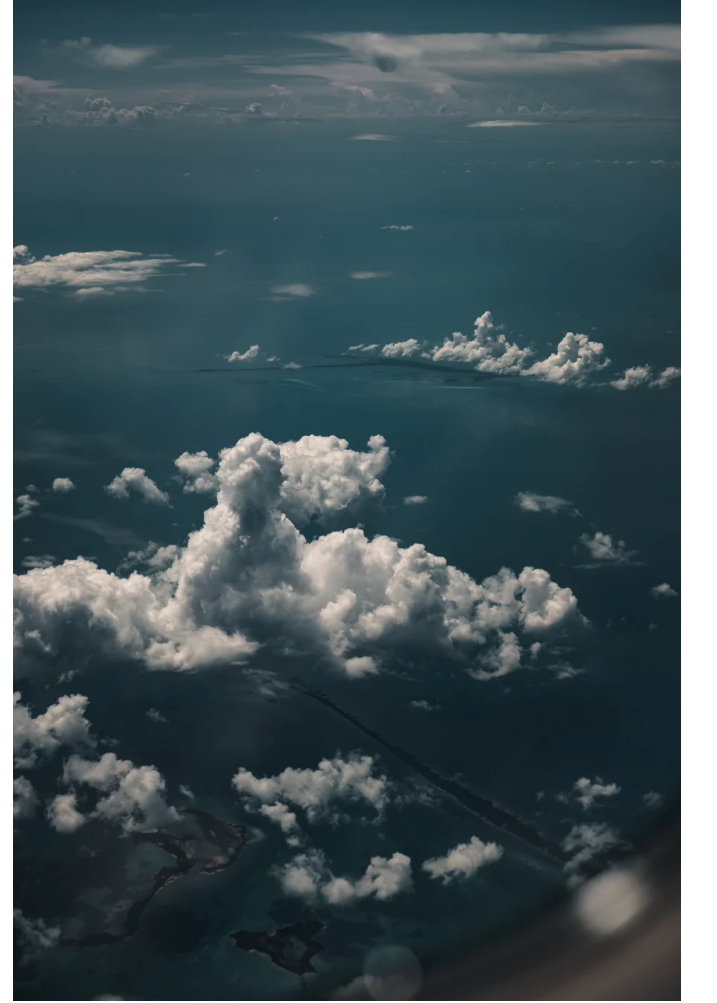
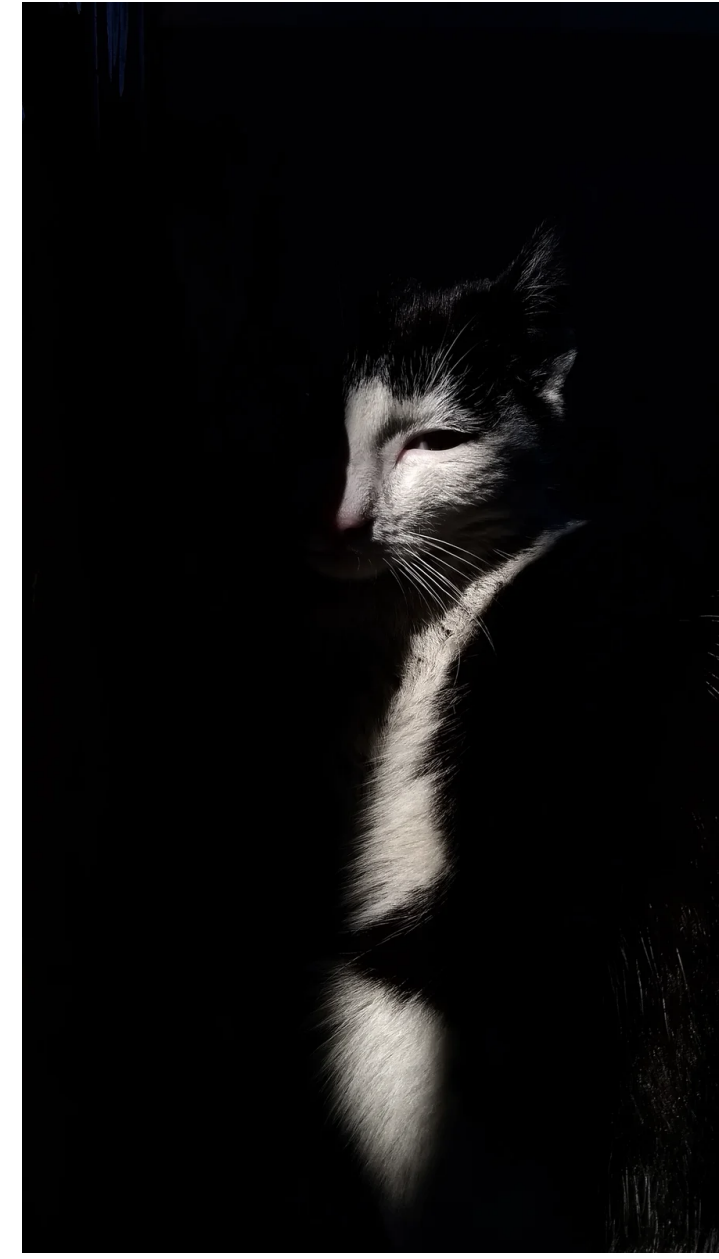
**Why is
machine
vision hard?**

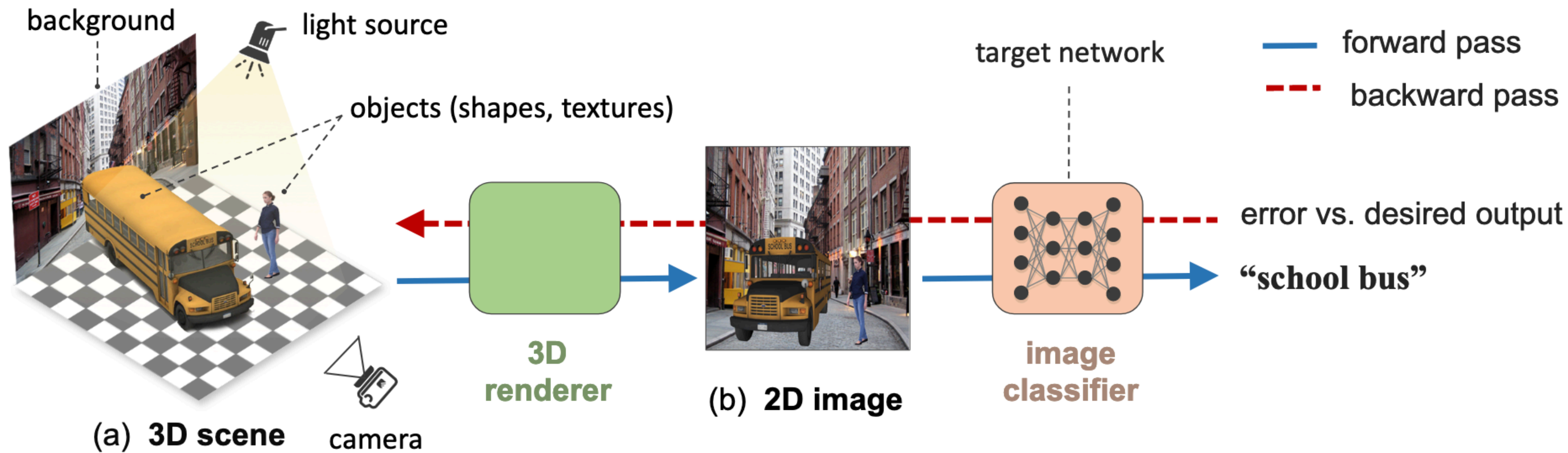
The deformable and truncated cat



Figure 1. **The deformable and truncated cat.** Cats exhibit (almost) unconstrained variations in shape and layout.

Parkhi et al. *The truth about cats and dogs*. 2011





Strike (with) a Pose: Neural Networks Are Easily Fooled by Strange Poses of Familiar Objects. Alcorn et al. 2019

Computer Vision Challenges

■ Viewpoint Variation

- A single instance of an object can be oriented in many ways with respect to the camera

■ Scale variation

- Visual classes often exhibit variation in their size (size in the real world, not only in terms of their extent in the image)

■ Deformation

- Many objects of interest are not rigid bodies and can be deformed in extreme ways

■ Occlusion

- The objects of interest can be occluded. Sometimes only a small portion of an object (as little as few pixels) could be visible

■ Illumination Condition

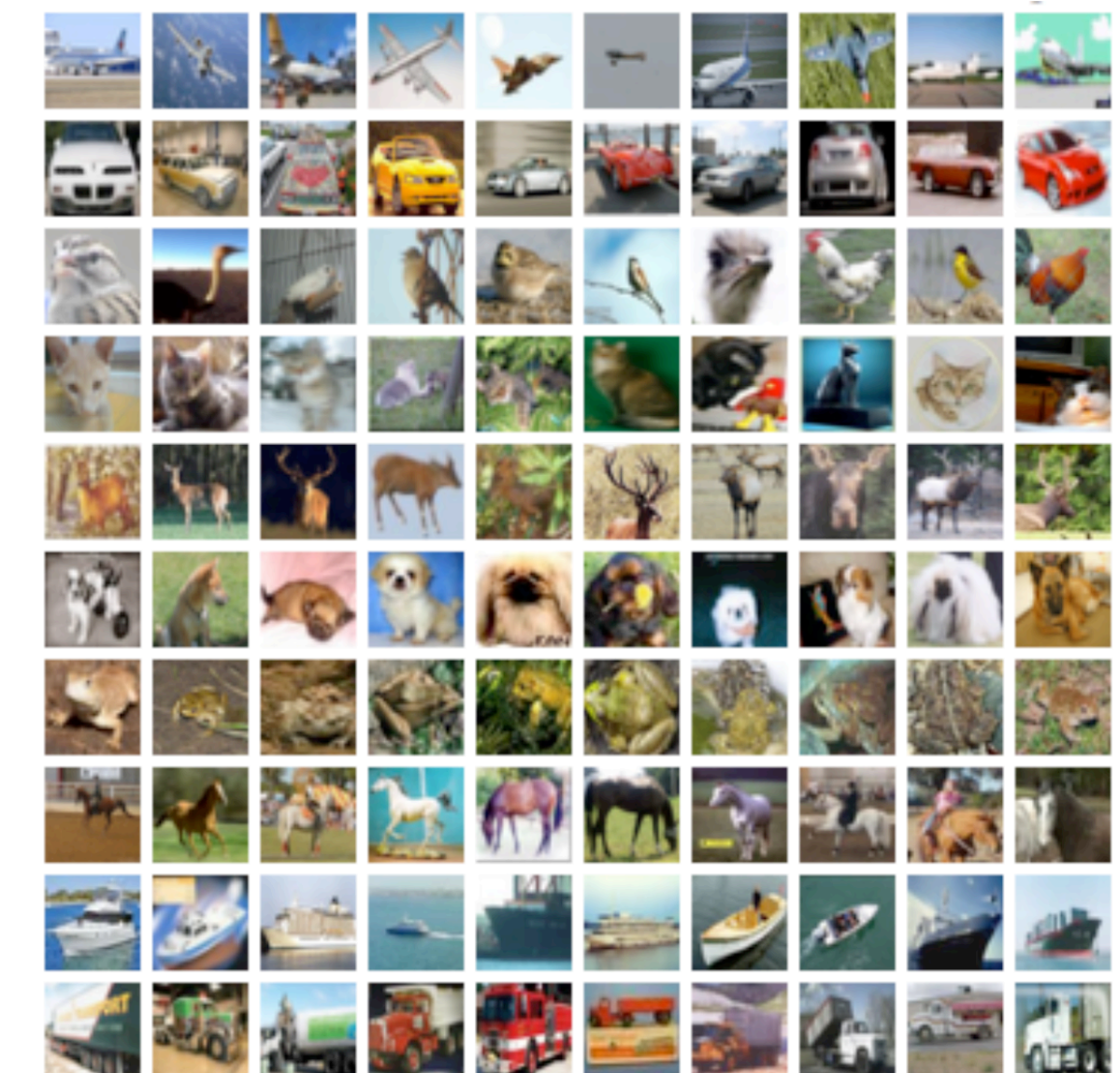
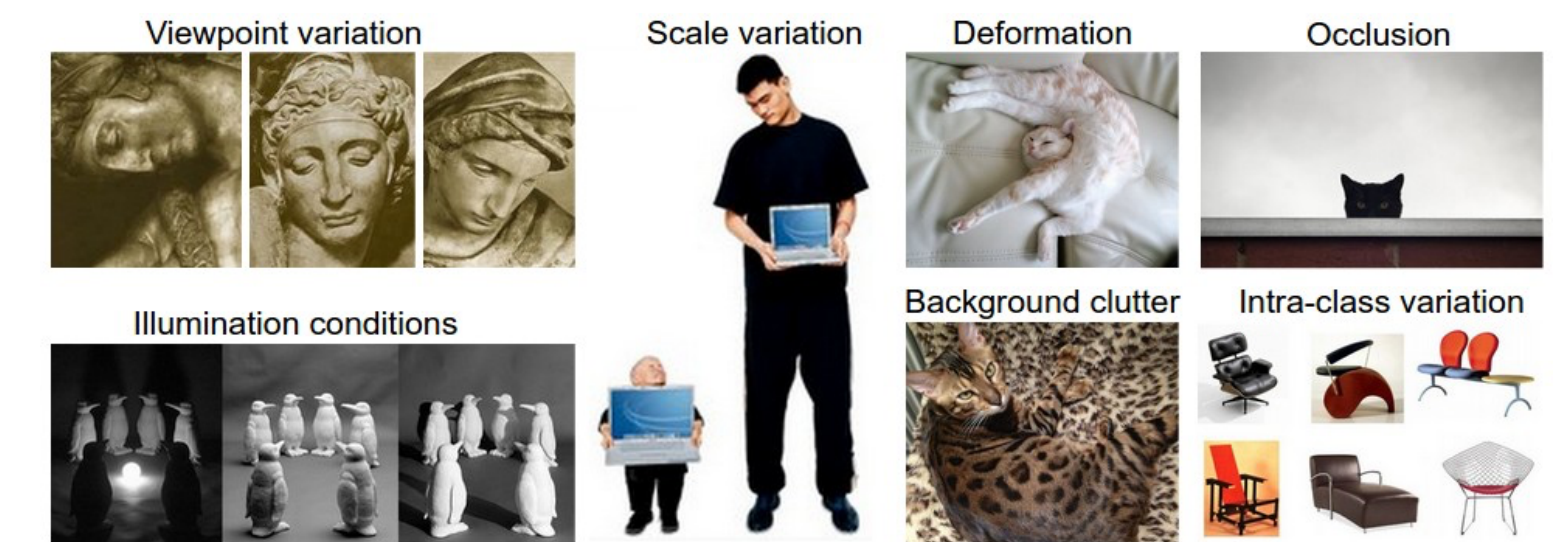
- The effects of illumination are drastic on the pixel level

■ Background clutter

- The objects of interest may blend into their environment, making them hard to identify

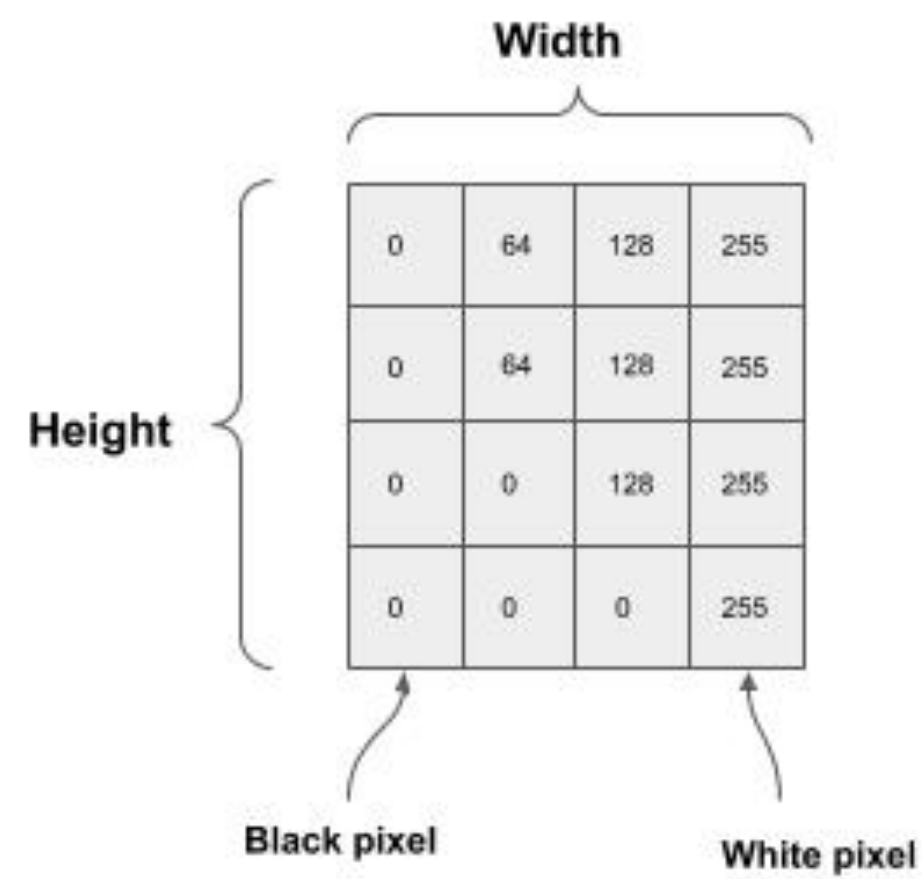
■ Intra-class variation

- The classes of interest can often be relatively broad, such as chair. There are many different types of these objects, each with their own appearance



**How CV
models
work?**

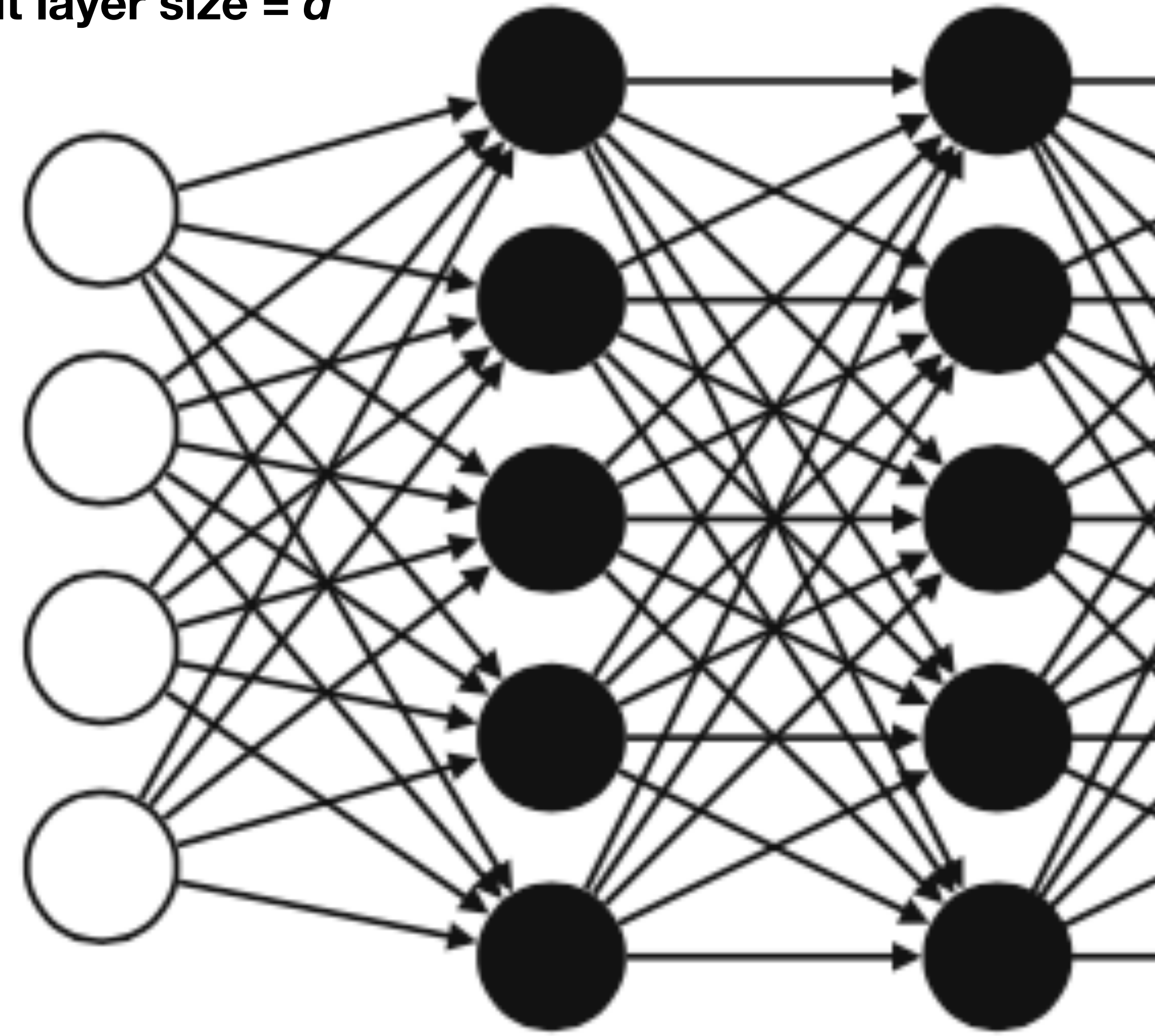
Flattening



$$d = \text{width} \times \text{height}$$

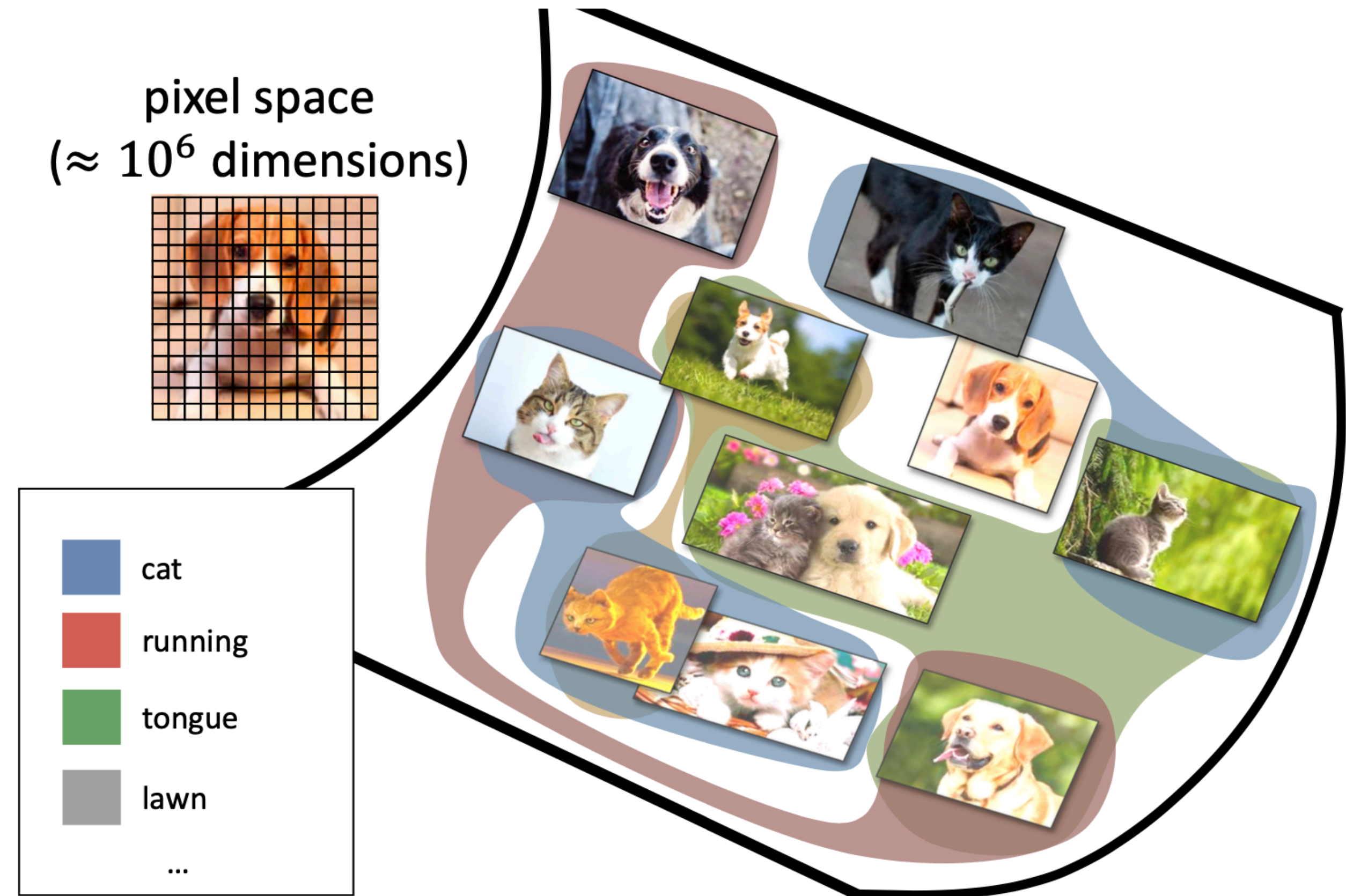


$$\text{Input layer size} = d$$

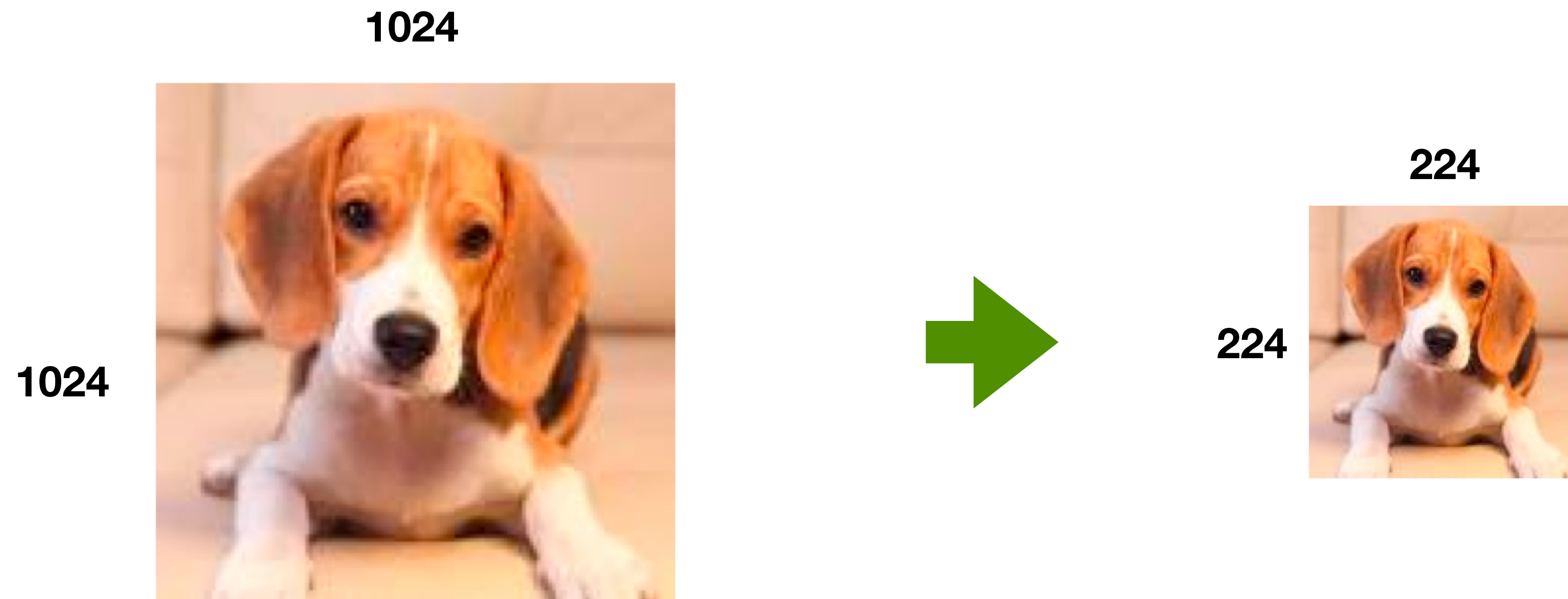


Course of dimensionality

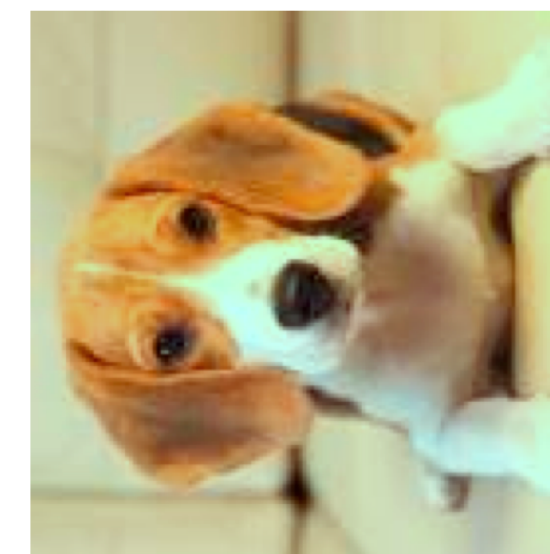
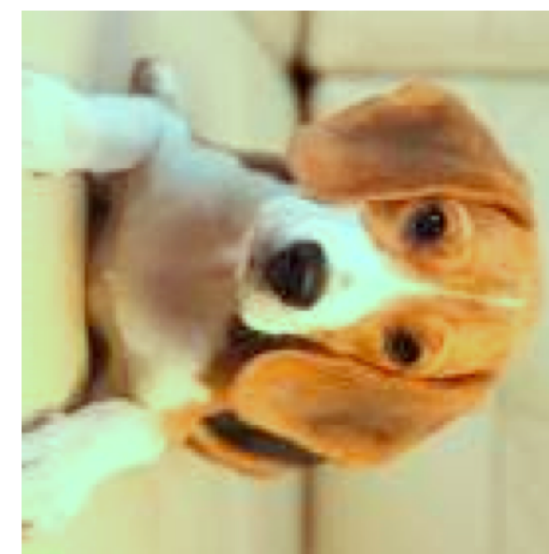
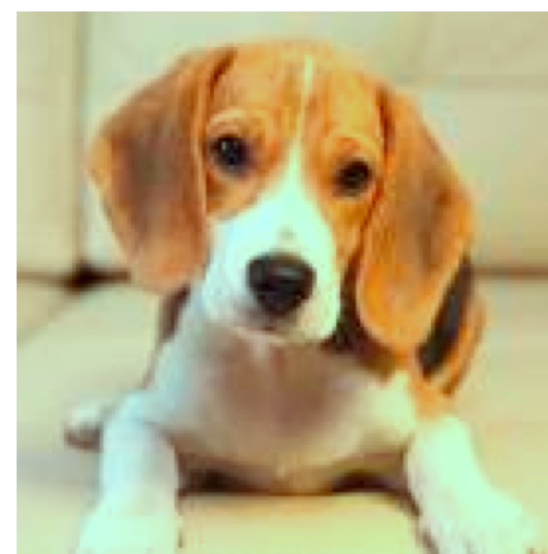
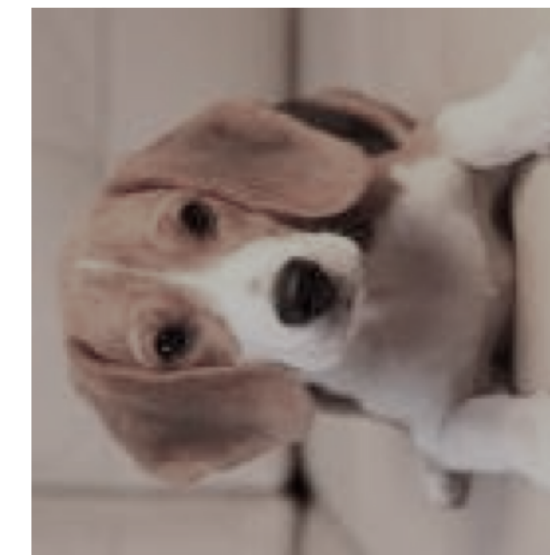
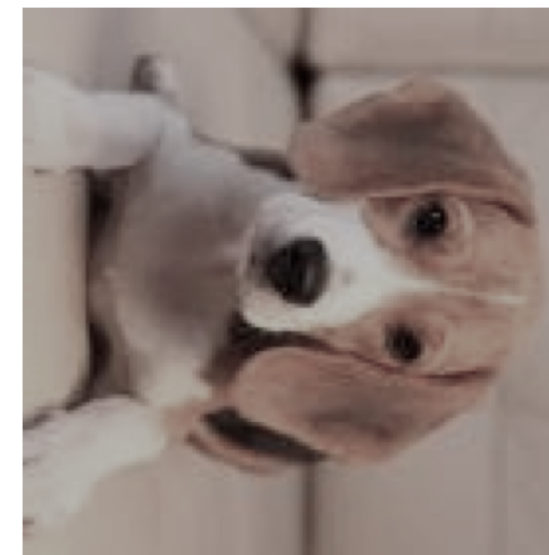
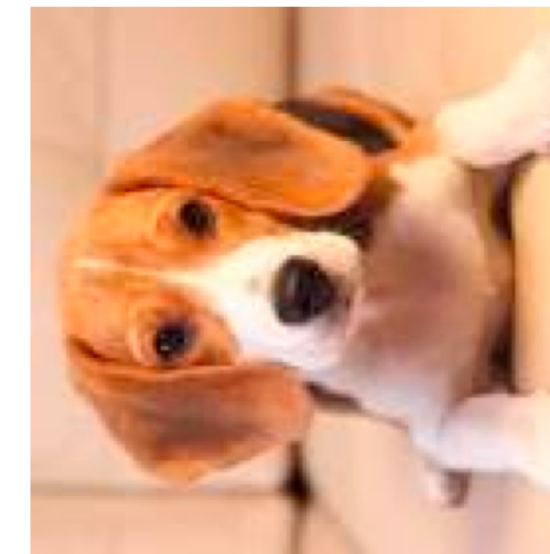
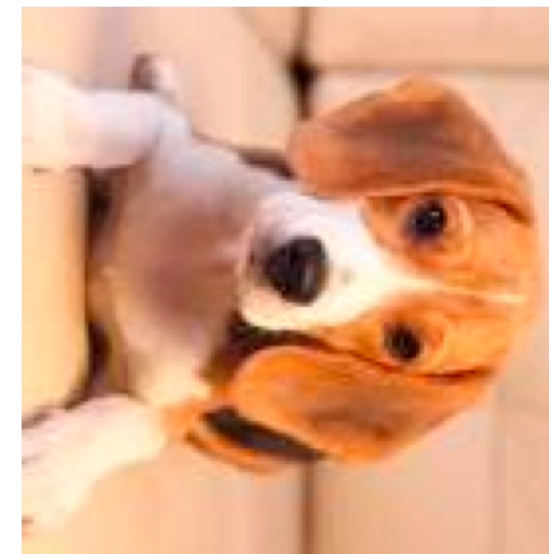
- High dimensionality
 - A 1024×768 image has $d = 786432!$
 - A tiny 32×32 image has $d = 1024$
- Decision boundaries in pixel space are extremely complex
- We will need “big” ML models with lots of parameters
 - For example, linear regressors need d parameters



Downsampling



What about generalisation?



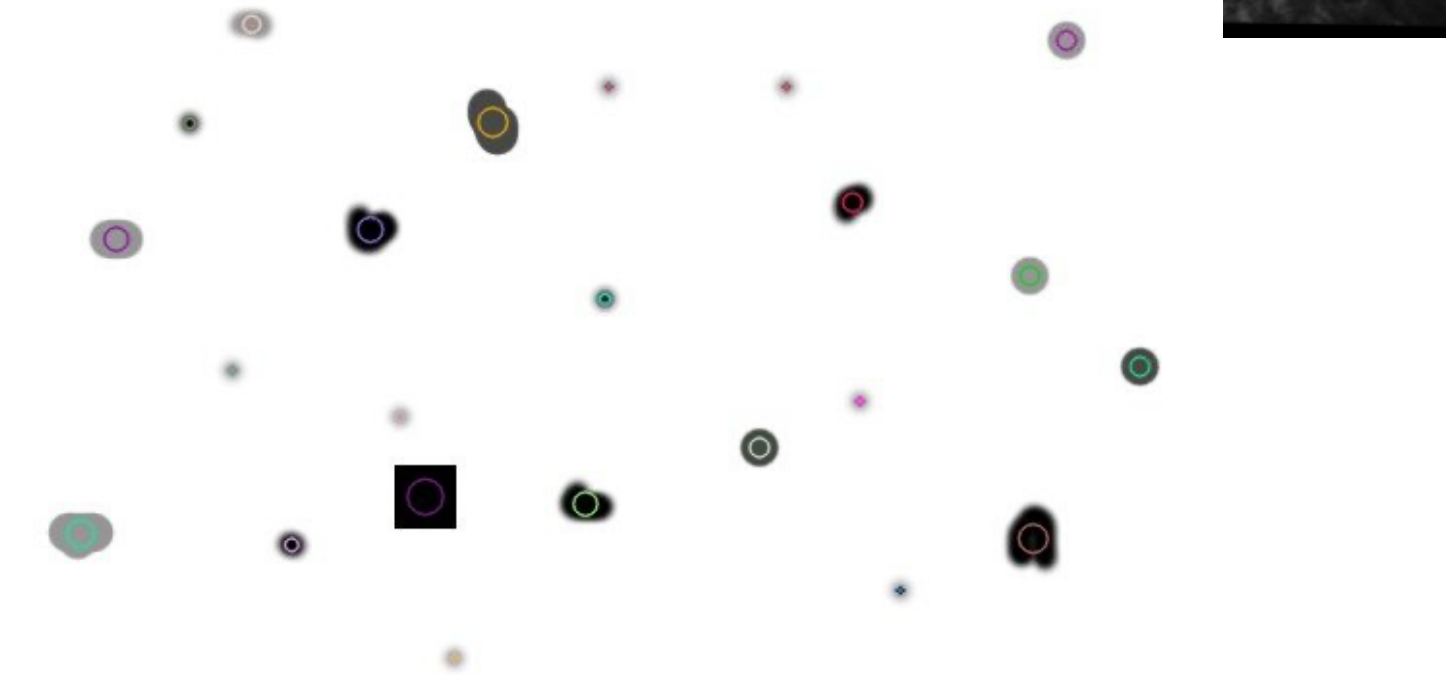
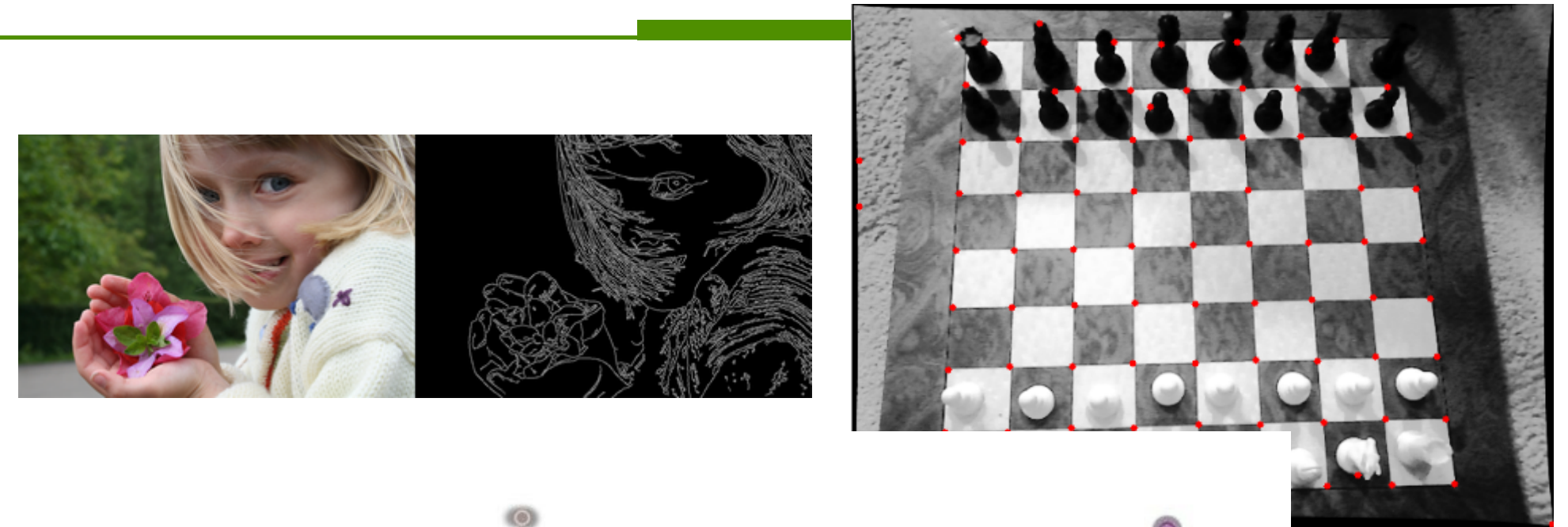
The “old days”: Feature Extraction

■ Feature

- A relevant piece of information about the content of an image
 - e.g. edges, corners, blobs (regions), ridges

■ A good feature

- Is repeatable
- Identifiable
- can be easily tracked and compared
- Is consistent across different scales, lighting conditions, and viewing angles
- Is still visible in noisy images or when only part of an object is visible
- can distinguish objects from one another



Feature after looking at one image



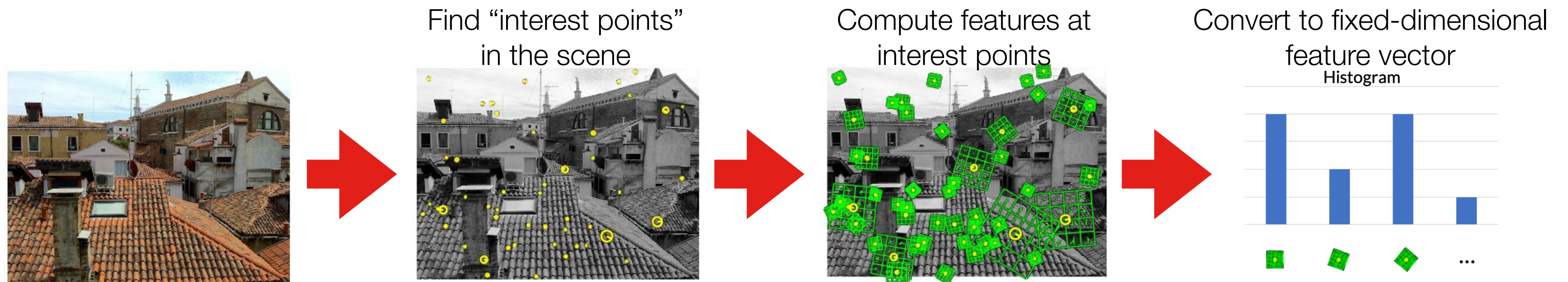
Feature after looking at thousands of images



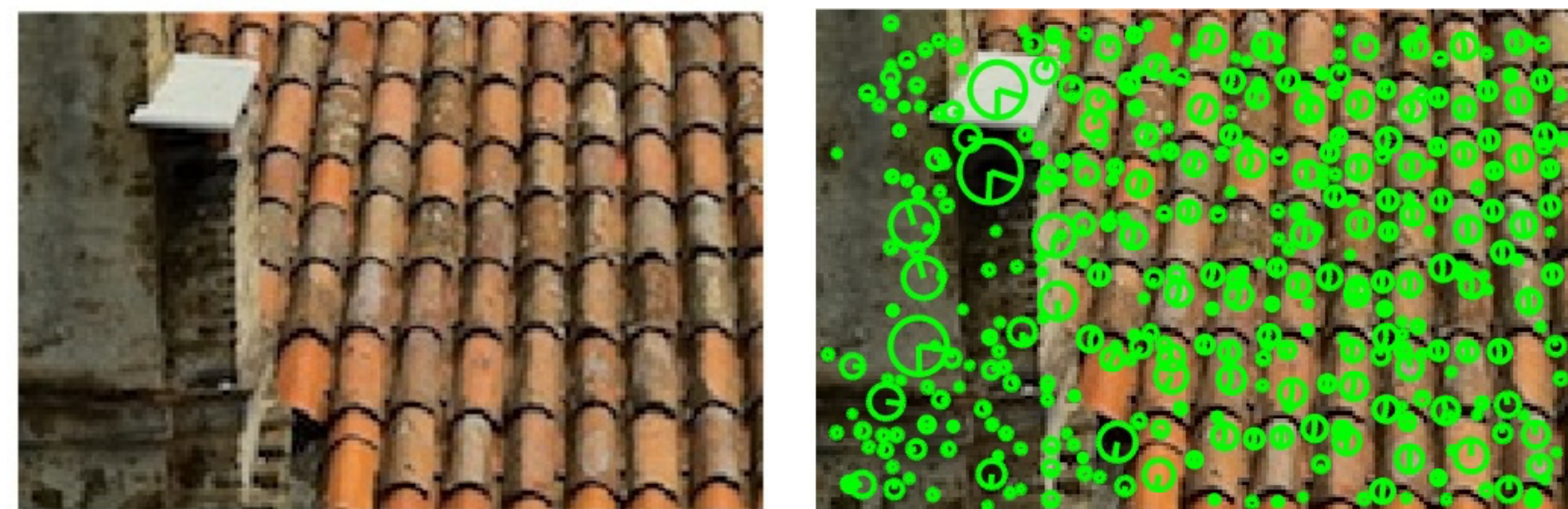
Feature Extraction Techniques

<https://www.vlfeat.org>

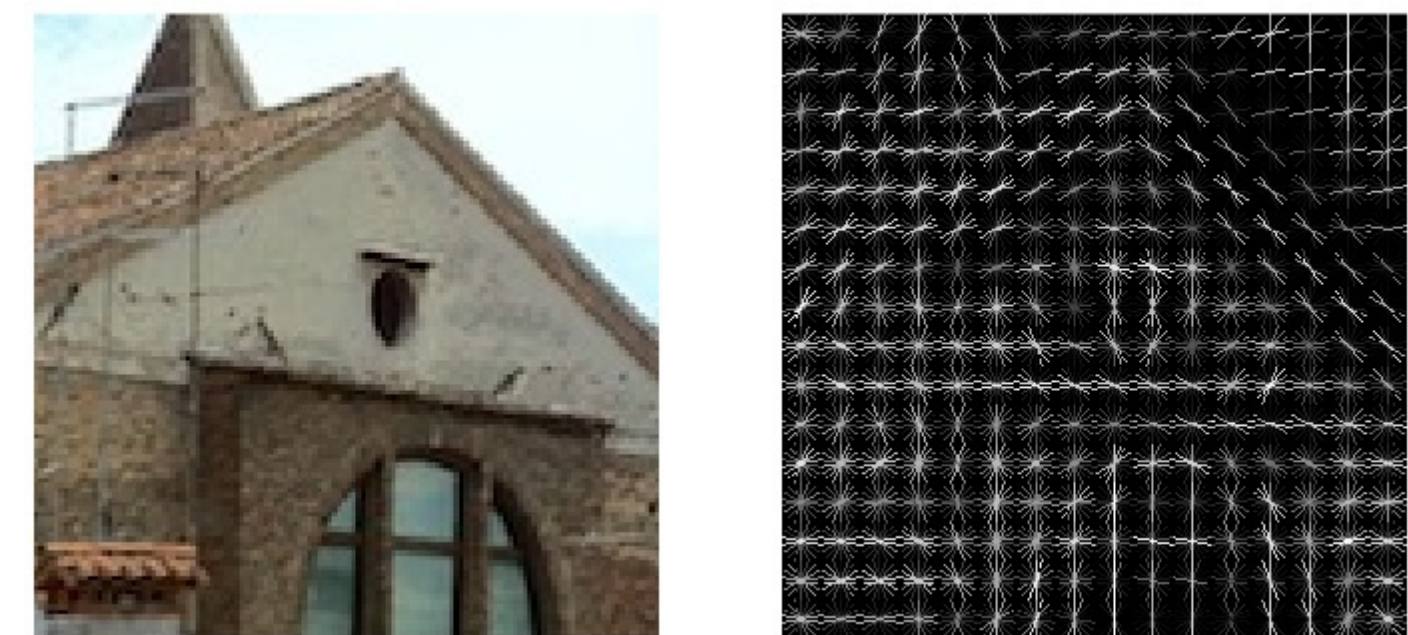
Scale-Invariant Feature Transform (SIFT)



Co-variant feature detector

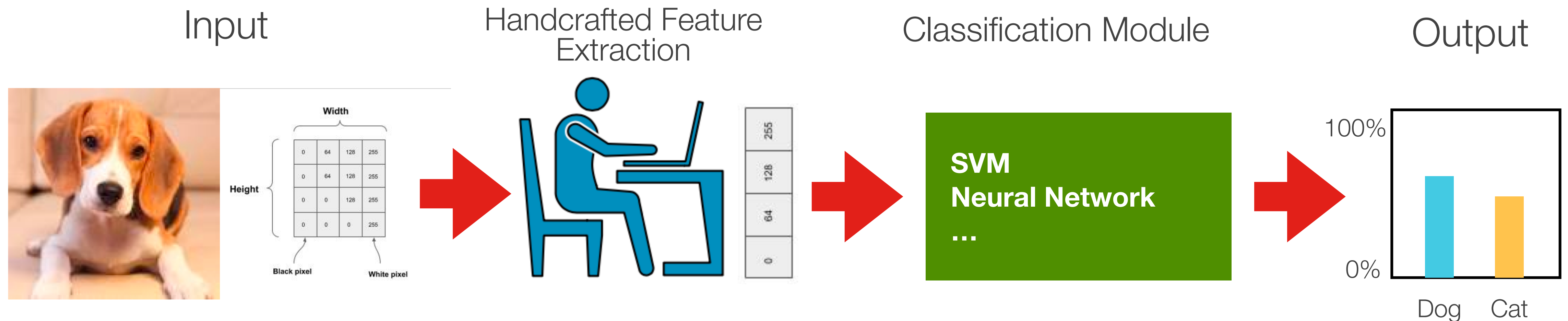


Histogram and oriented gradients



The “old days”: Feature Engineering

- Machine learning models are only as good as the features you provide
- To figure out which features you should use for a specific problem
 - rely on domain knowledge (or partner with domain experts)
 - experiment to create features that make machine learning algorithms work better



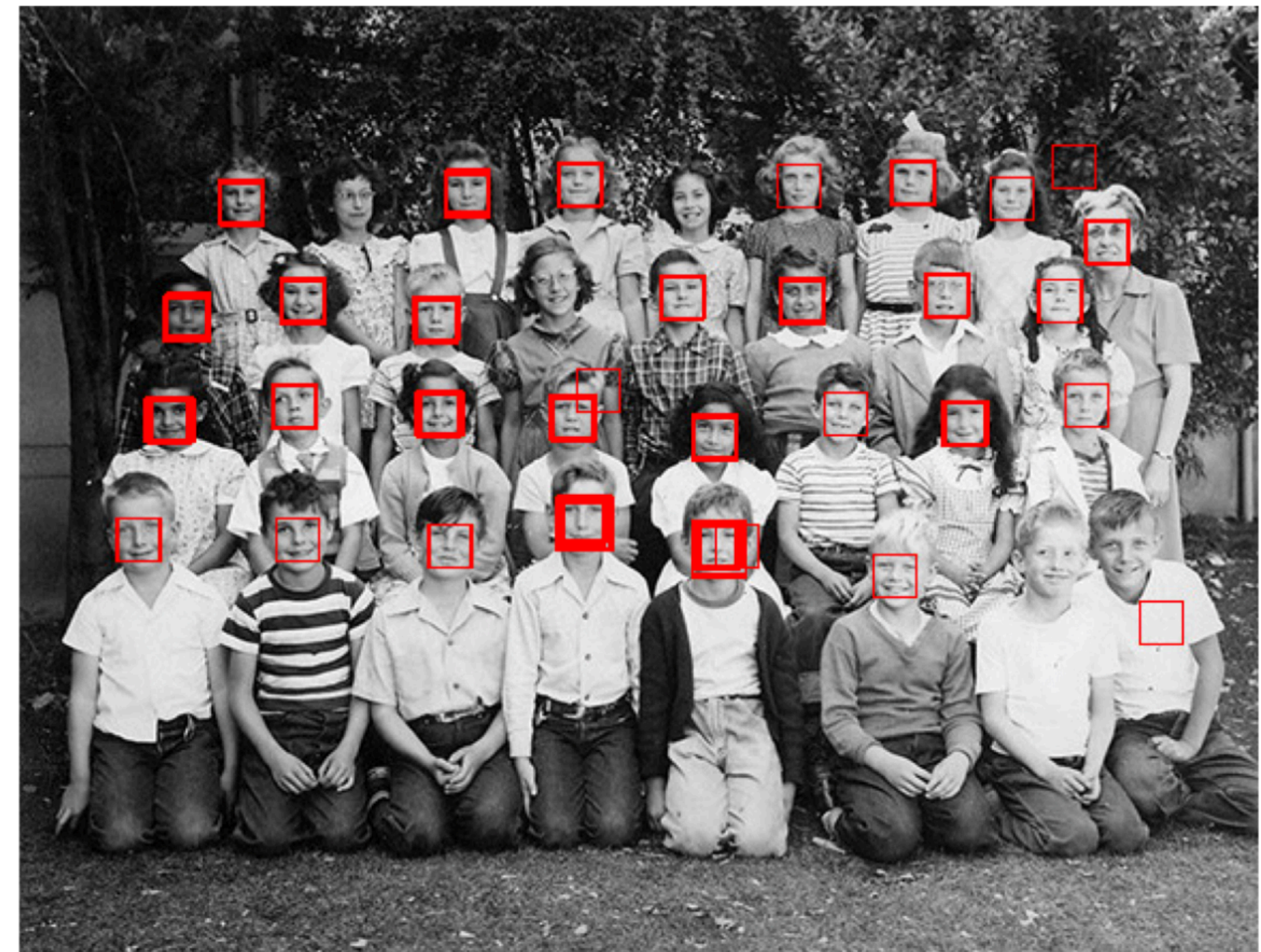
Performance

Object Detection (~2007)



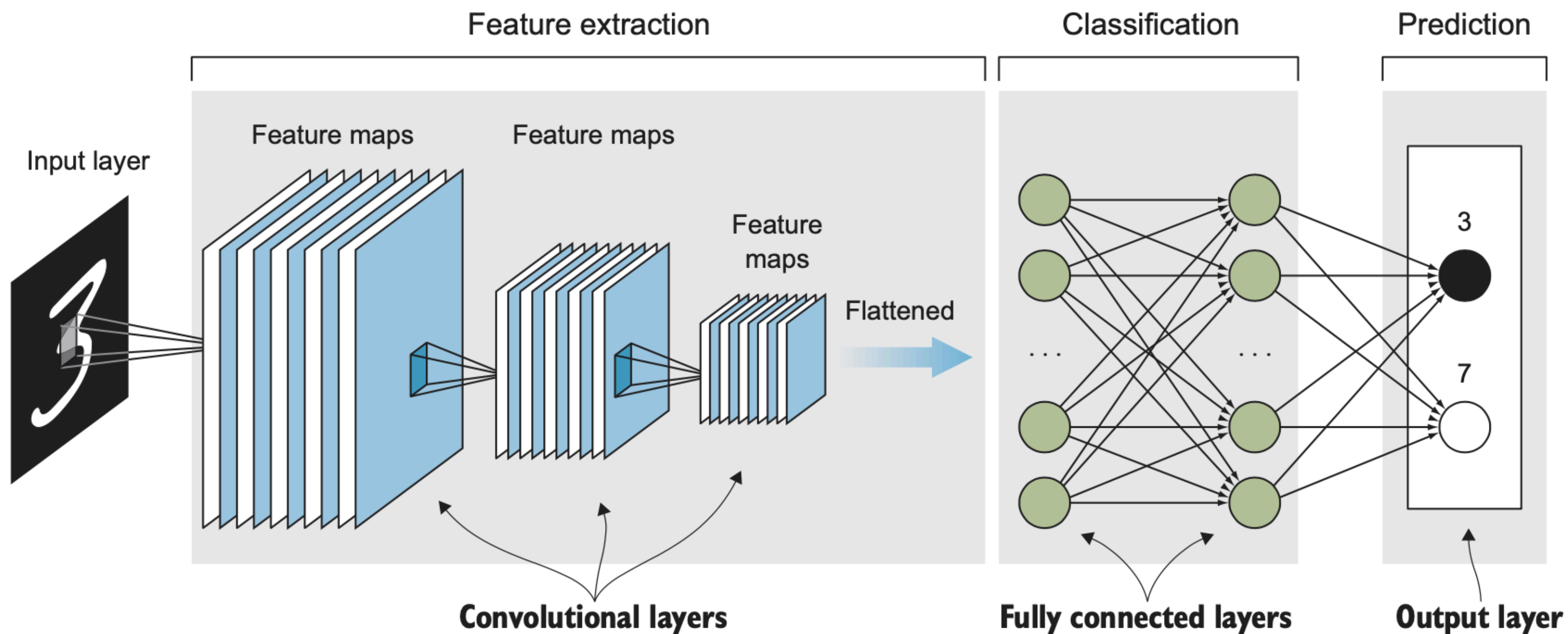
Felzenszwalb, Ramanan, McAllester. A Discriminatively Trained, Multiscale, Deformable Part Model. CVPR 2008 (DPM v1)

Face Detection (~2013)



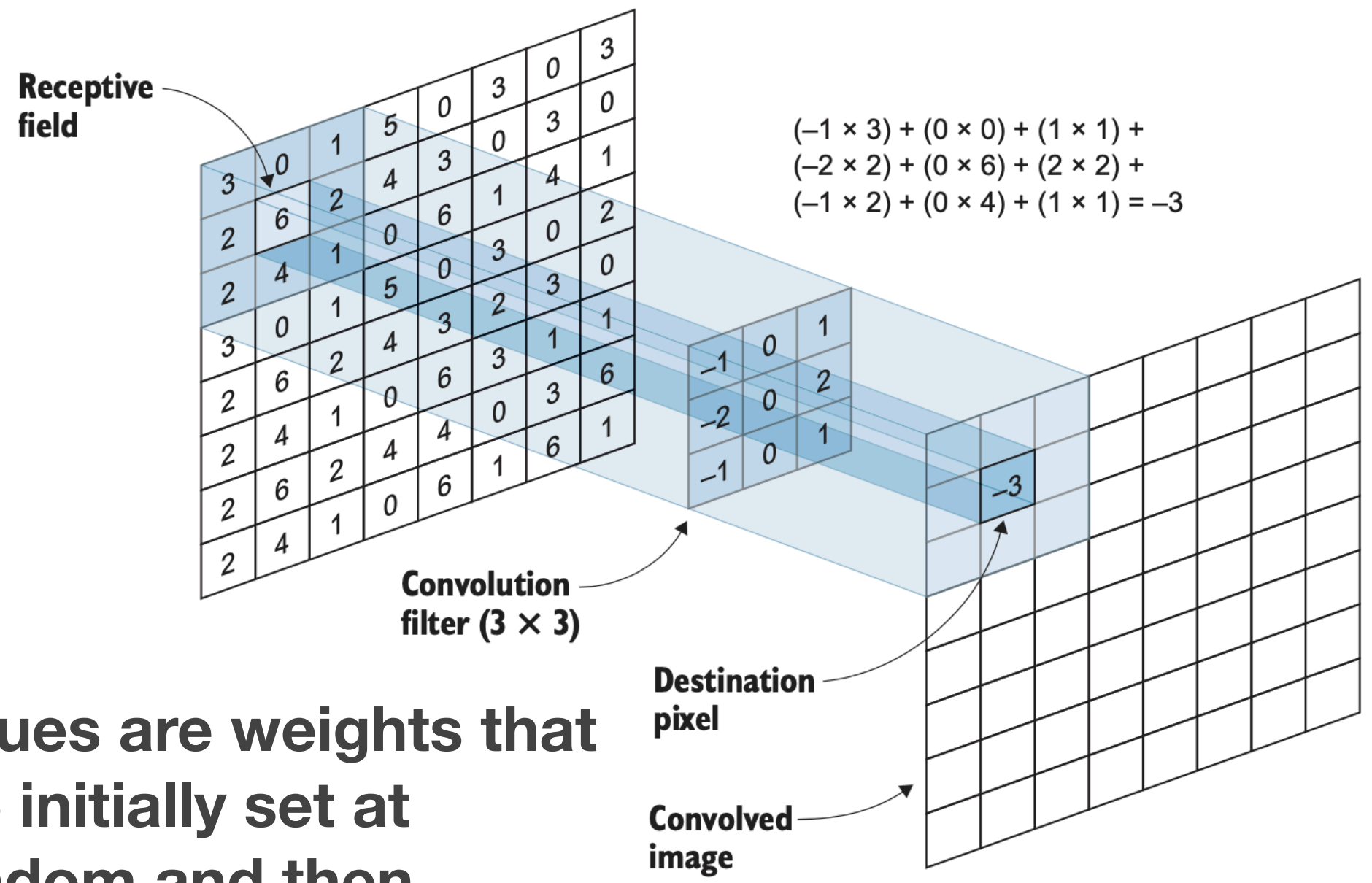
<https://github.com/alexdemartos/ViolaAndJones>

Convolutional Neural Networks




- CNNs exploit image properties to drastically reduce the number of model parameters
- Feature maps
 - **Automatically** extracted hierarchical
 - Retain spatial association between pixels
- **Translation invariance**
 - a dog is a dog even if its image is shifted by a few pixels
- Local interactions
 - all processing happens within very small image windows
 - within each layer, far-away pixels cannot influence nearby pixels

Convolution & Feature Maps



Values are weights that are initially set at random and then learned

Input image




Convolution kernel with optimized weights

0	-1	0
-1	4	-1
0	-1	0

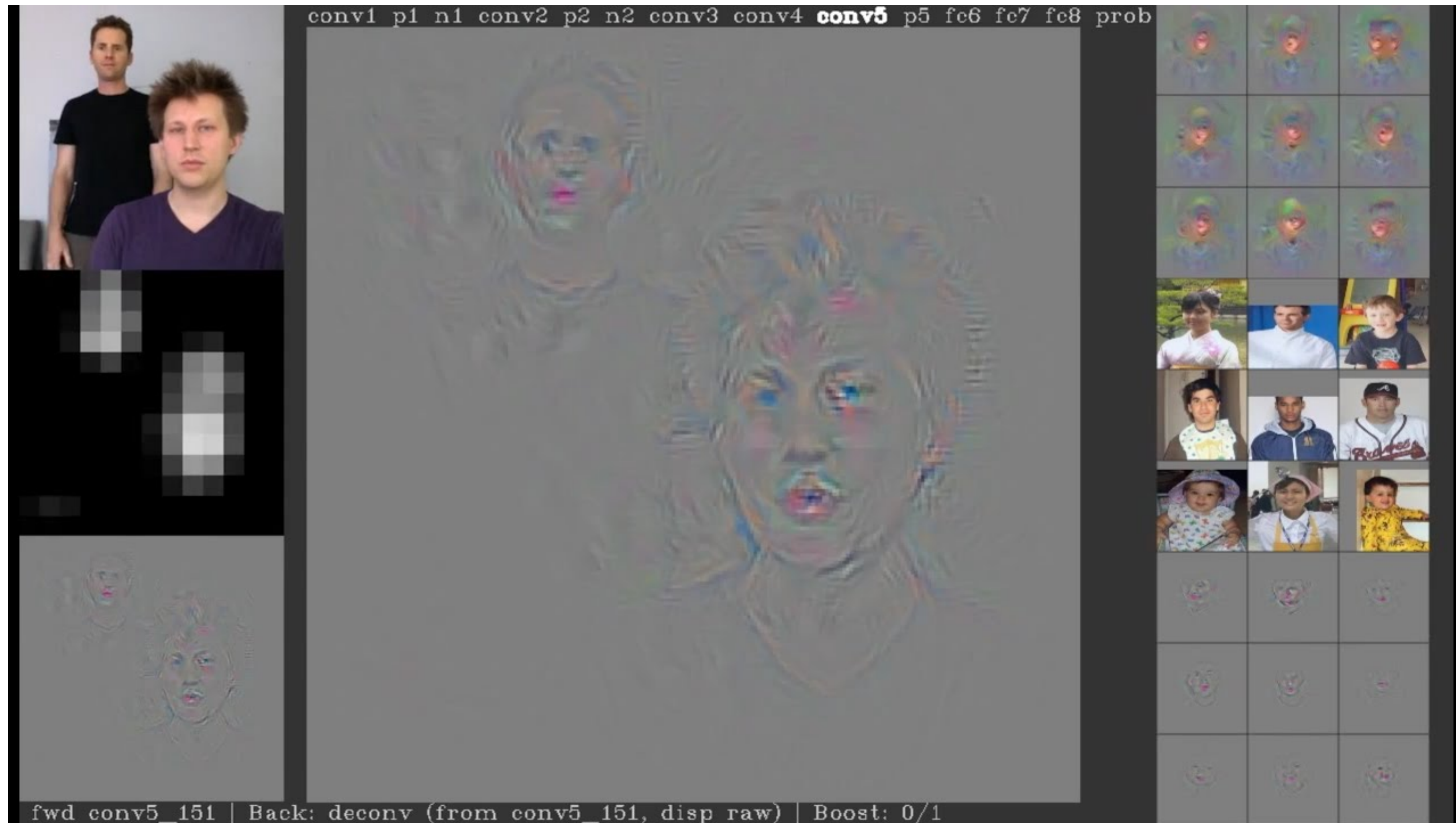
* =

Convolved image (feature map)

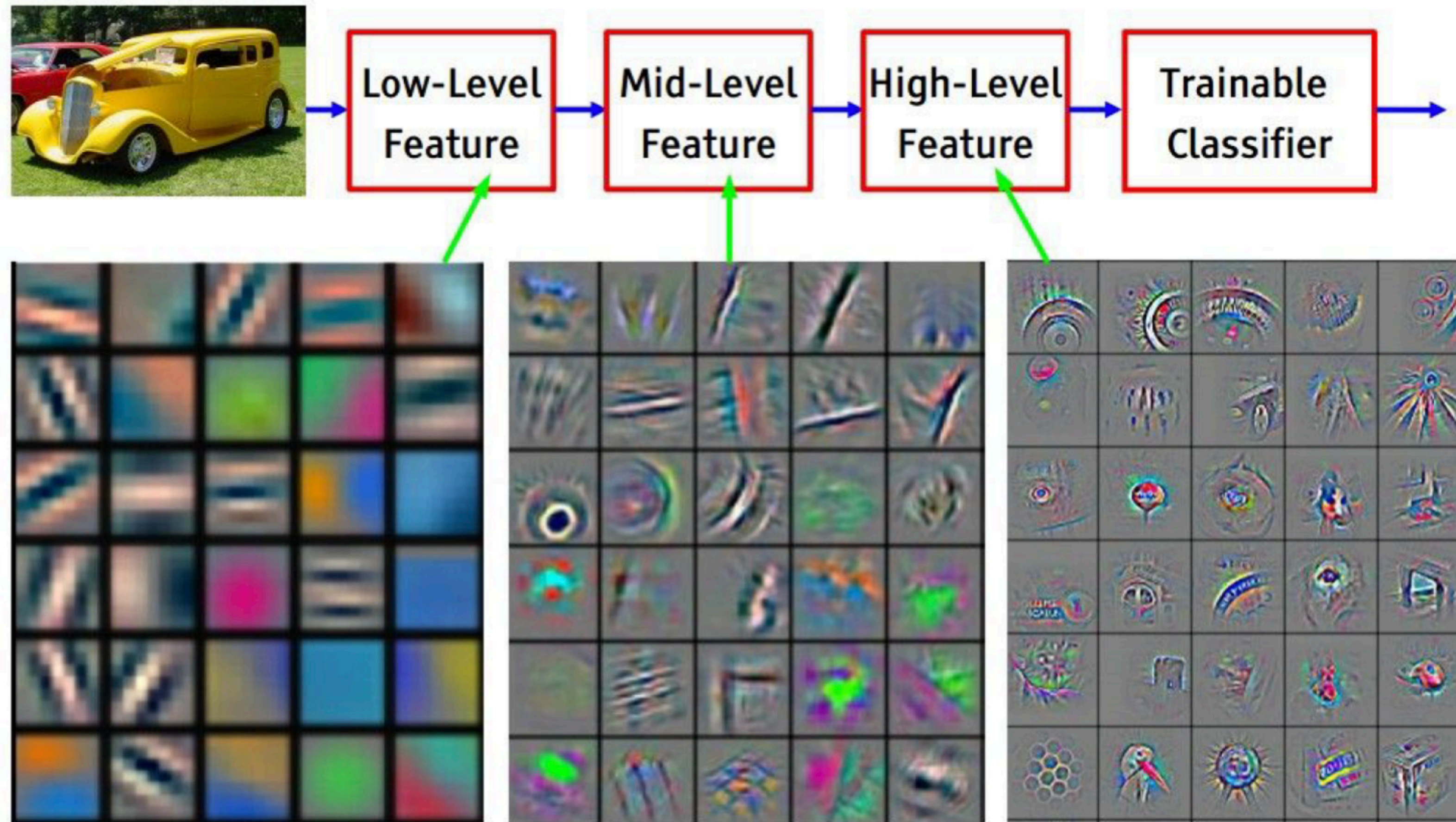


■ Try this
<https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html>

What do CNN learn?



Feature Visualisation



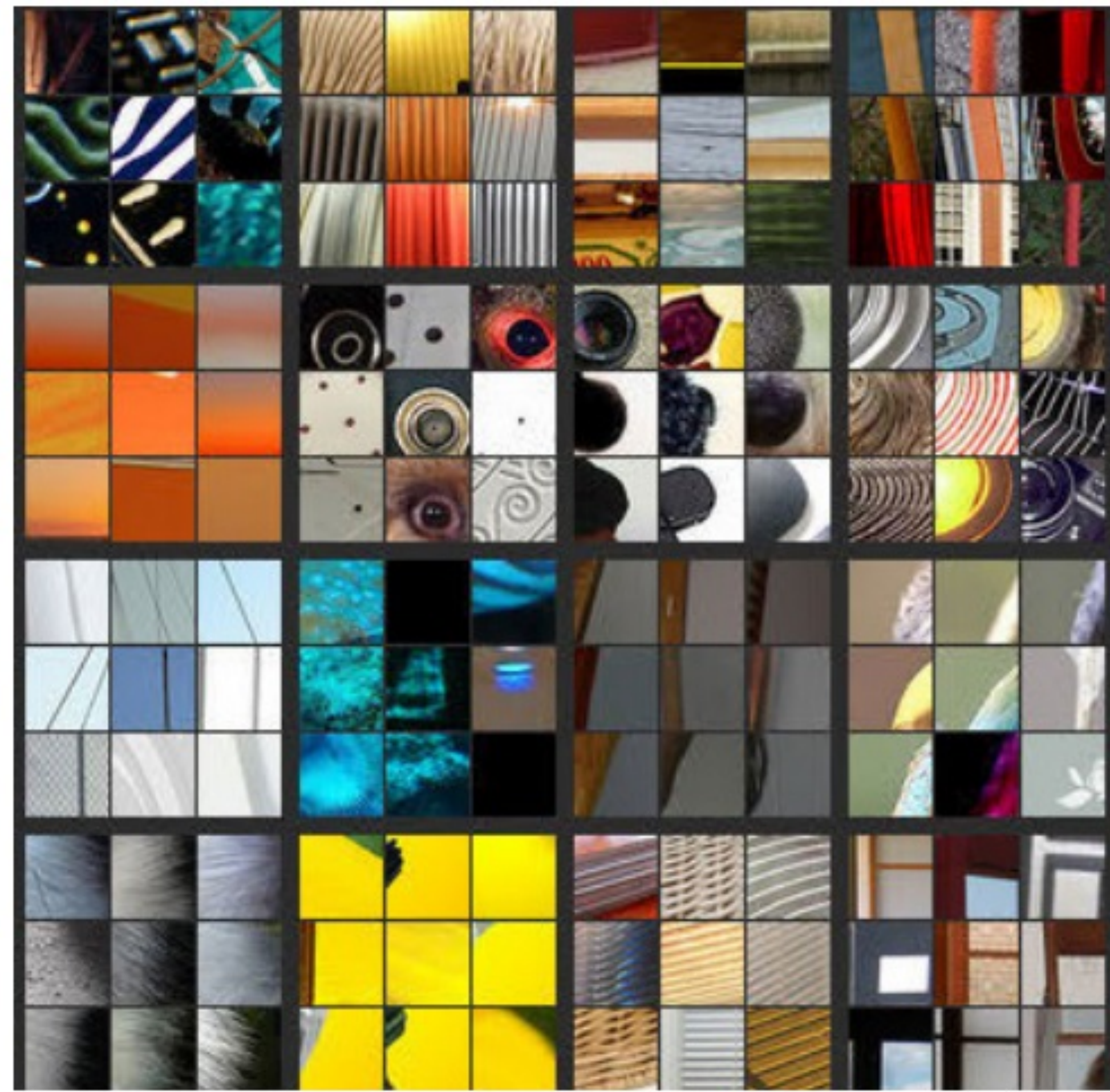
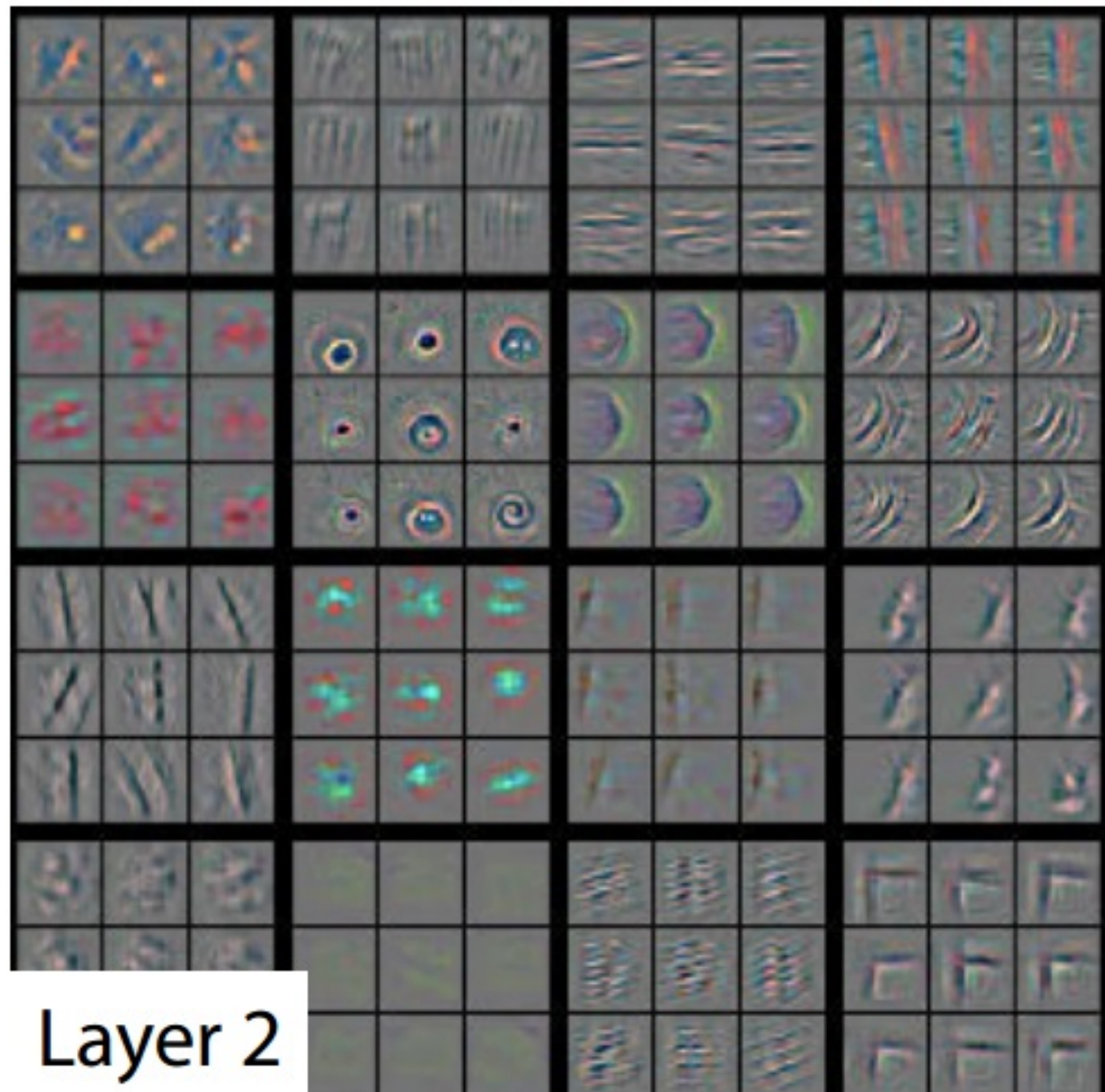
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



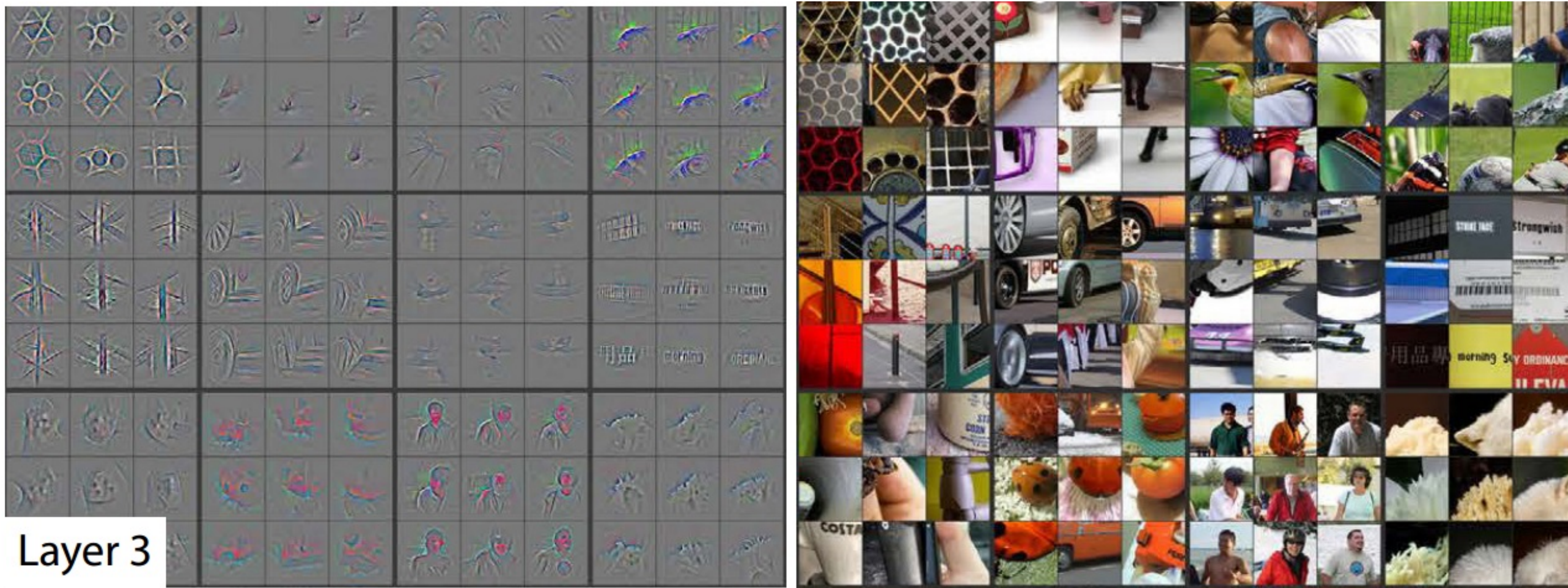
Layer 1



Visualizing and Understanding Convolutional Network. Zeiler and Fergus, ECCV 2014



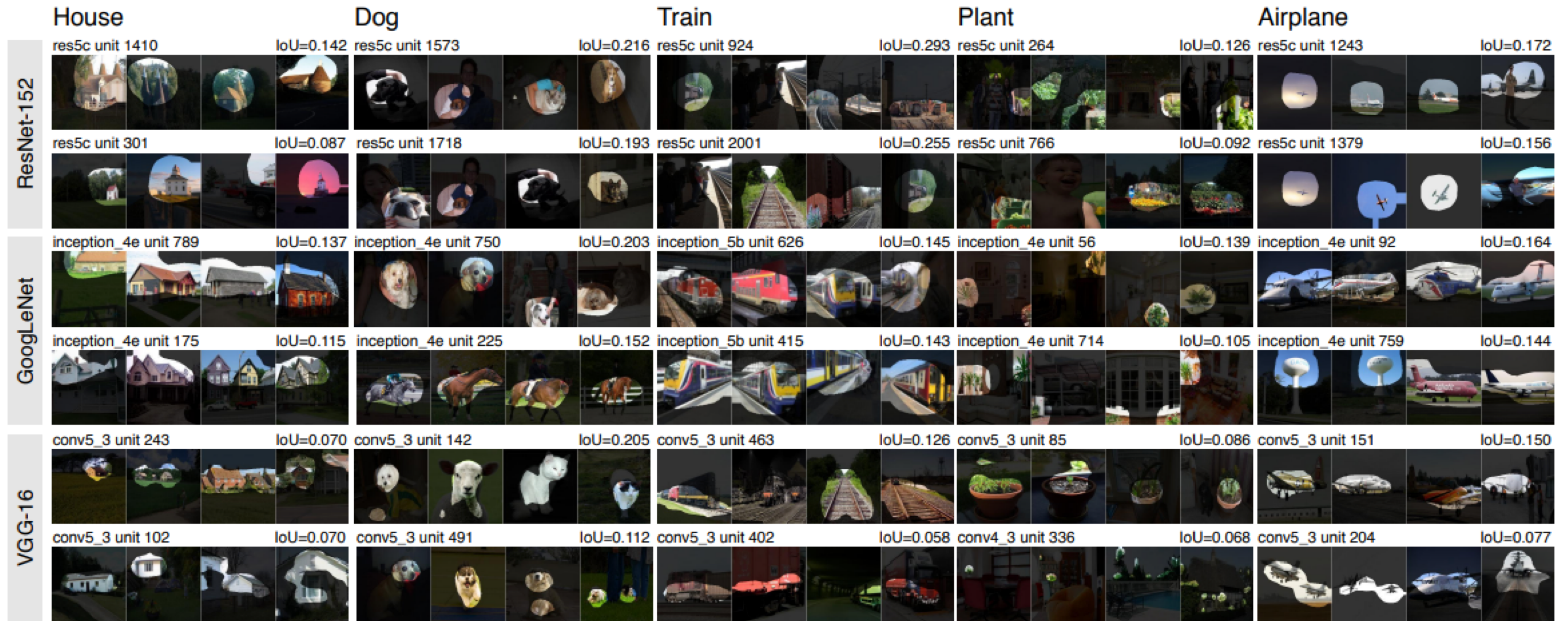
Visualizing and Understanding
Convolutional Network.
Zeiler and Fergus, ECCV 2014



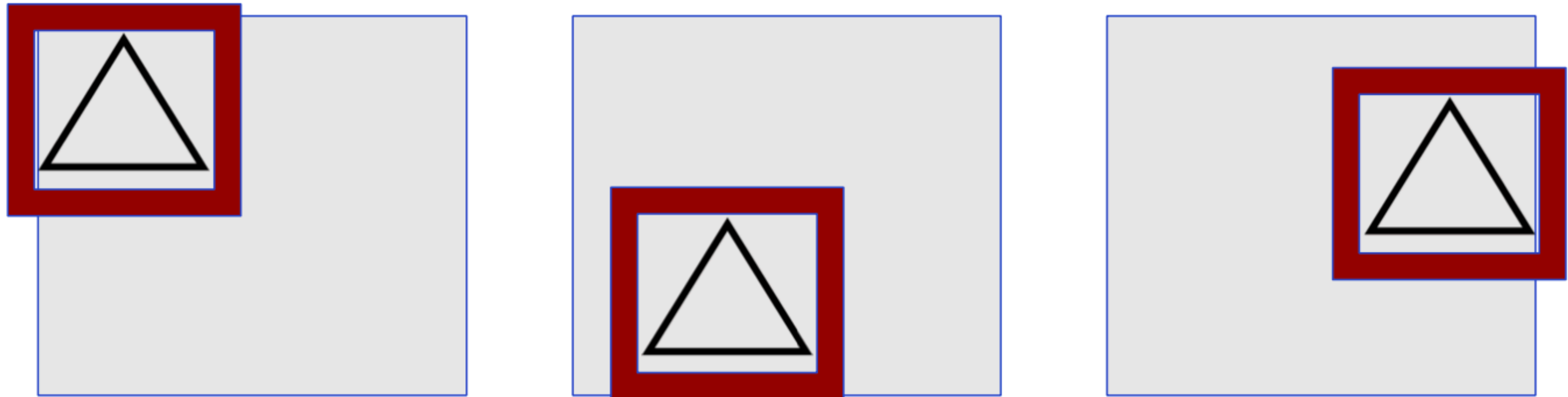
Layer 3

Visualizing and Understanding Convolutional Network.
Zeiler and Fergus, ECCV 2014

Network dissection



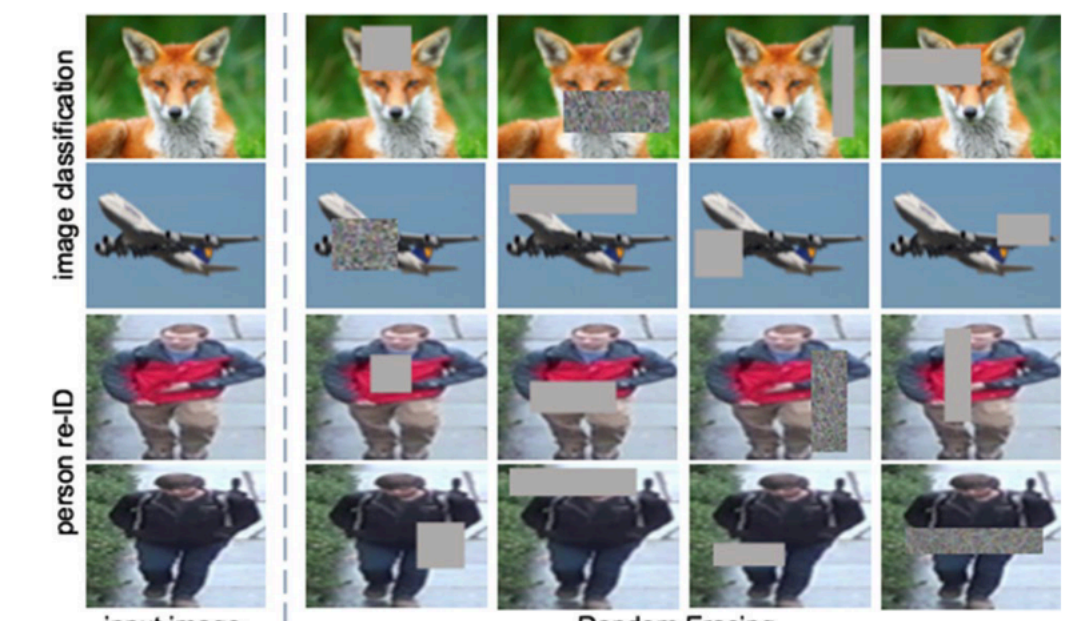
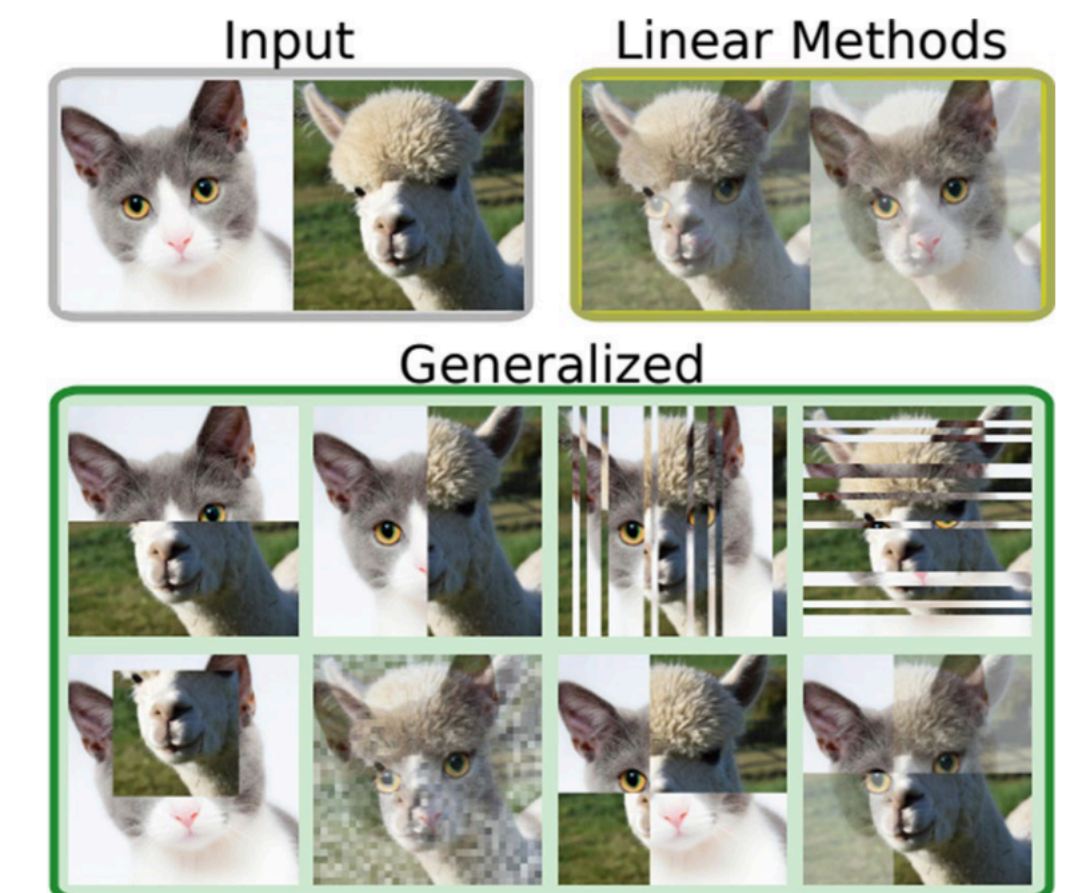
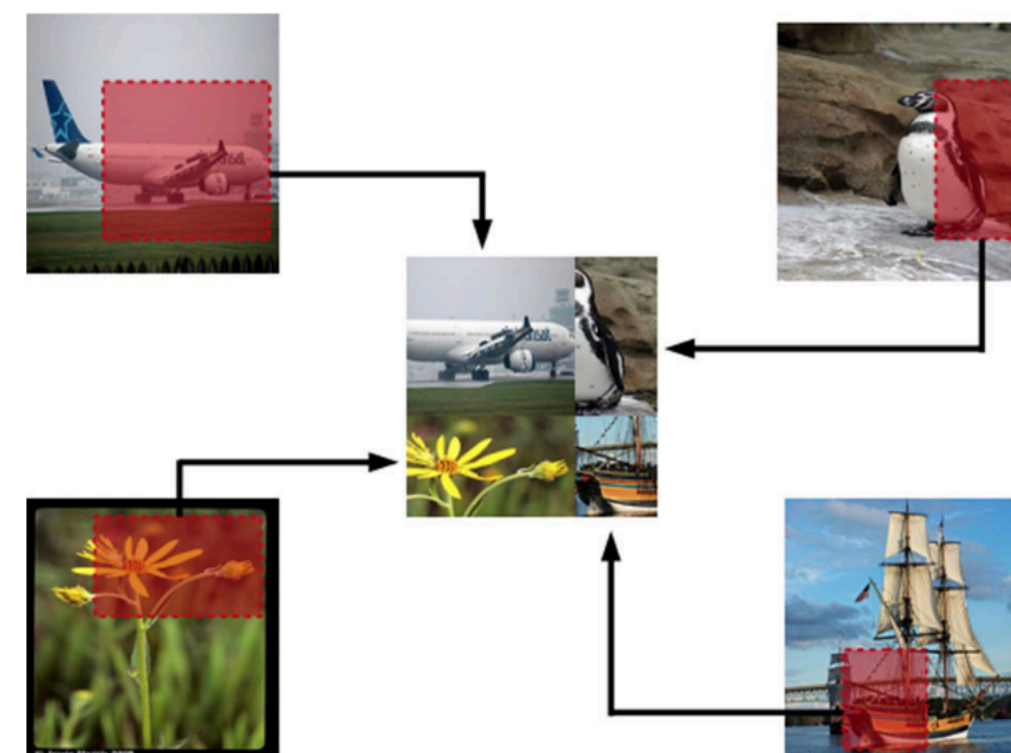
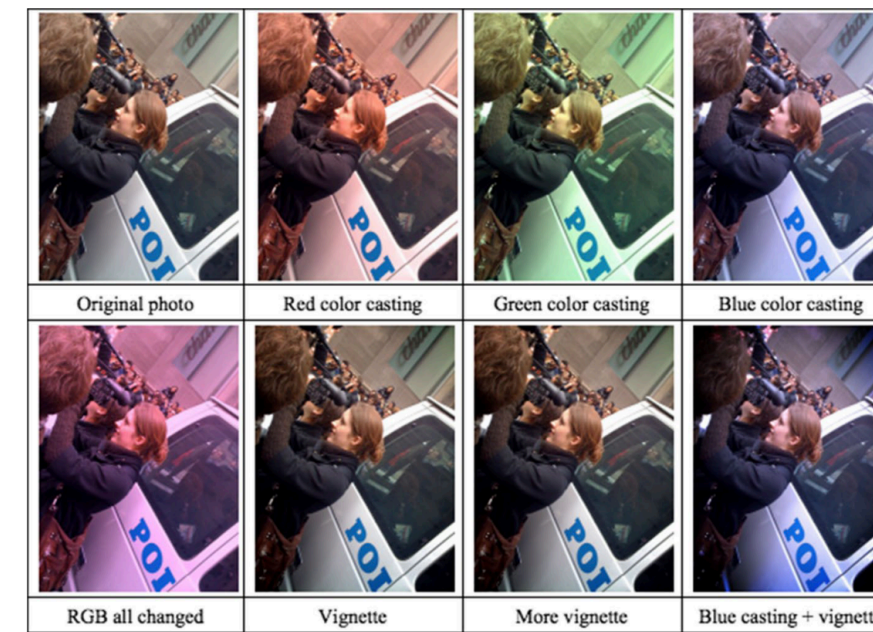
Translation Invariance



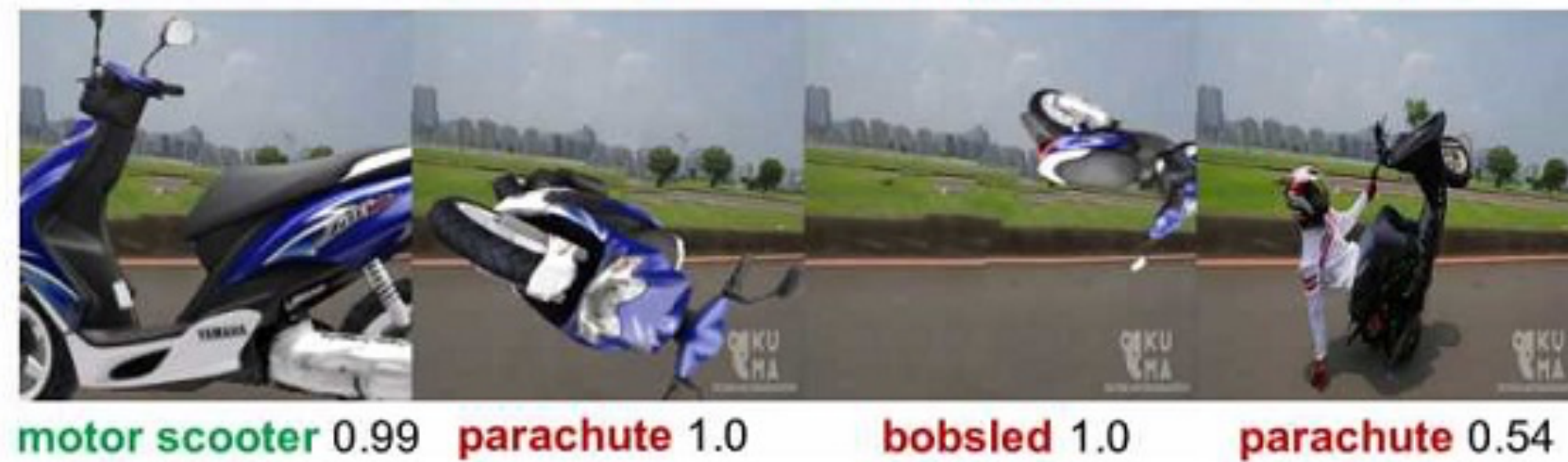
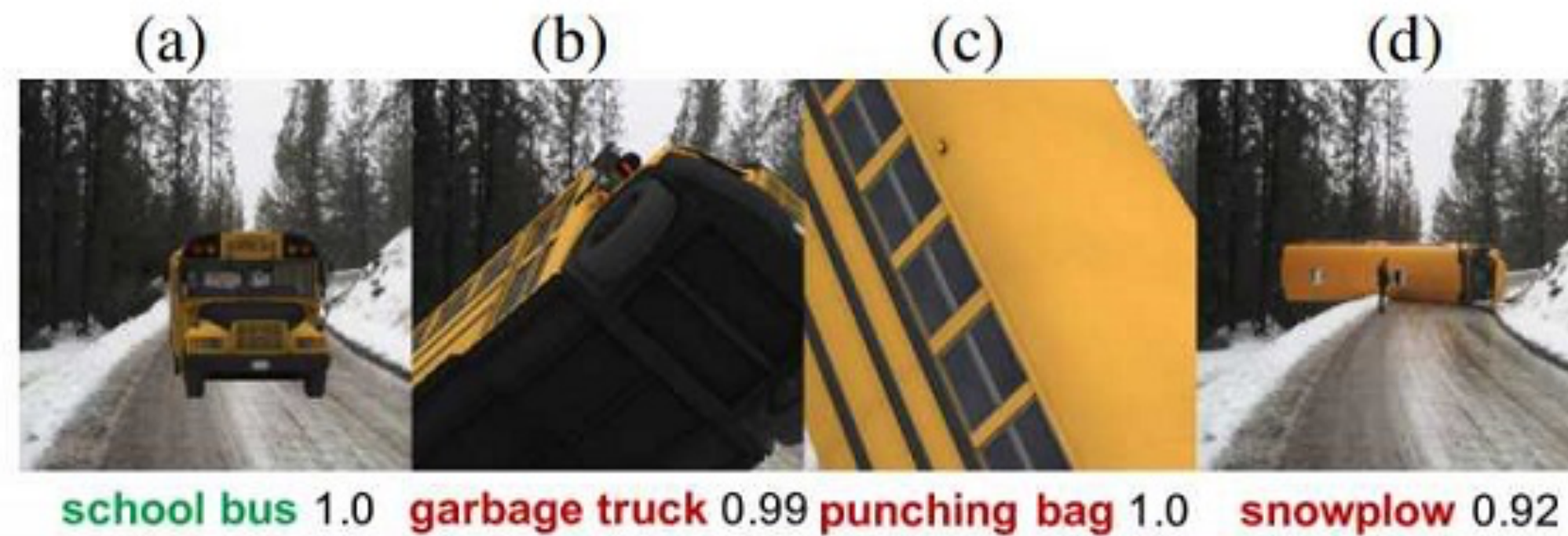
- But not rotation and scaling invariance!

Data Augmentation

- Generate variations of the input data, to improve generalisability (out of distribution inputs)
 - Improve invariance (rotation, scaling, distortion)
- Geometric
 - Flipping
 - Color space
 - Cropping
 - Rotation
 - Translation
 - Noise Injection
- Color space transformation
- Mixing Images
- Random erasing
- Adversarial training
- GAN-based image generation



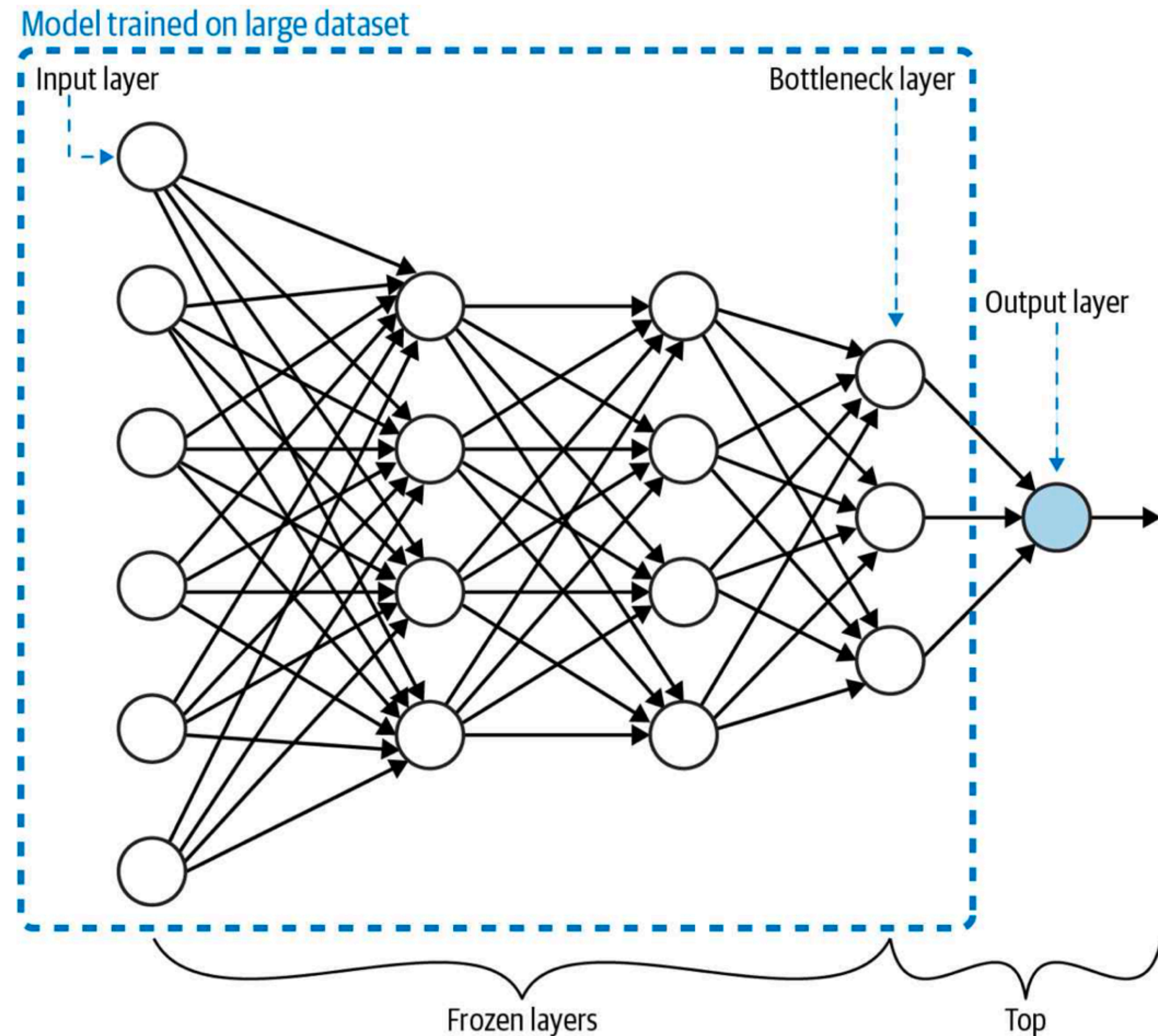
Robustness to input variation



Strike (with) a Pose: Neural Networks Are Easily Fooled by Strange Poses of Familiar Objects. Alcorn et al. 2019

<https://arxiv.org/pdf/1811.11553.pdf>

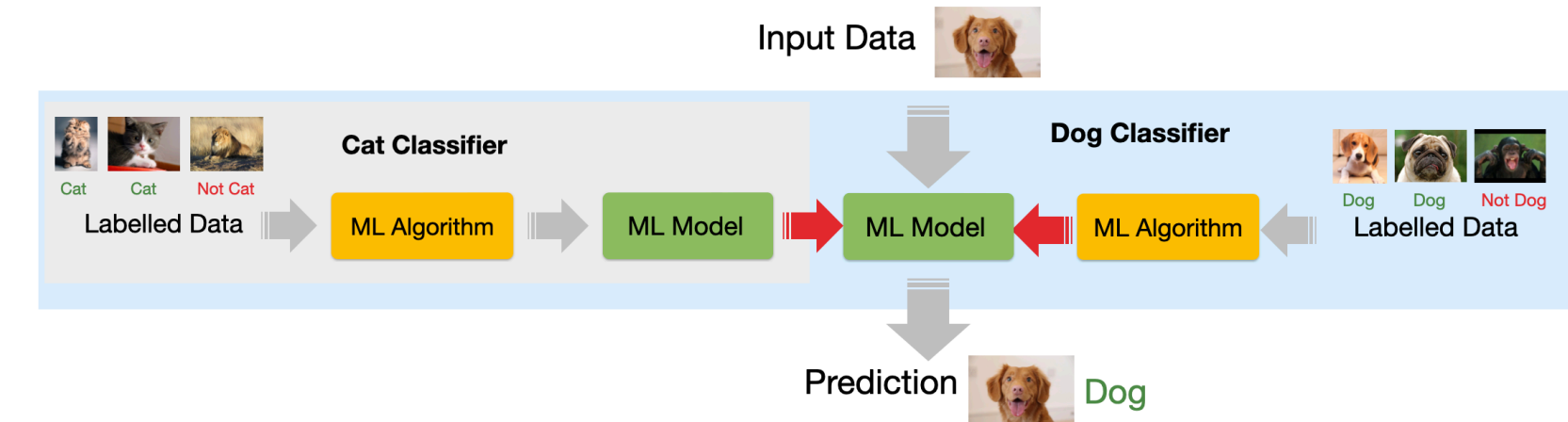
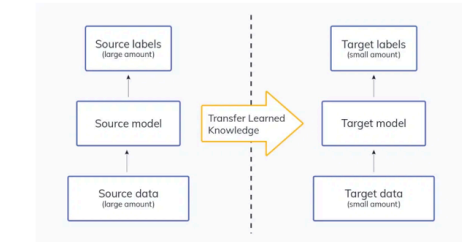
Transfer Learning



Transfer Learning

Reuse a model trained for one task is re-purposed (tuned) on a different but related task

Useful in tasks lacking abundant data



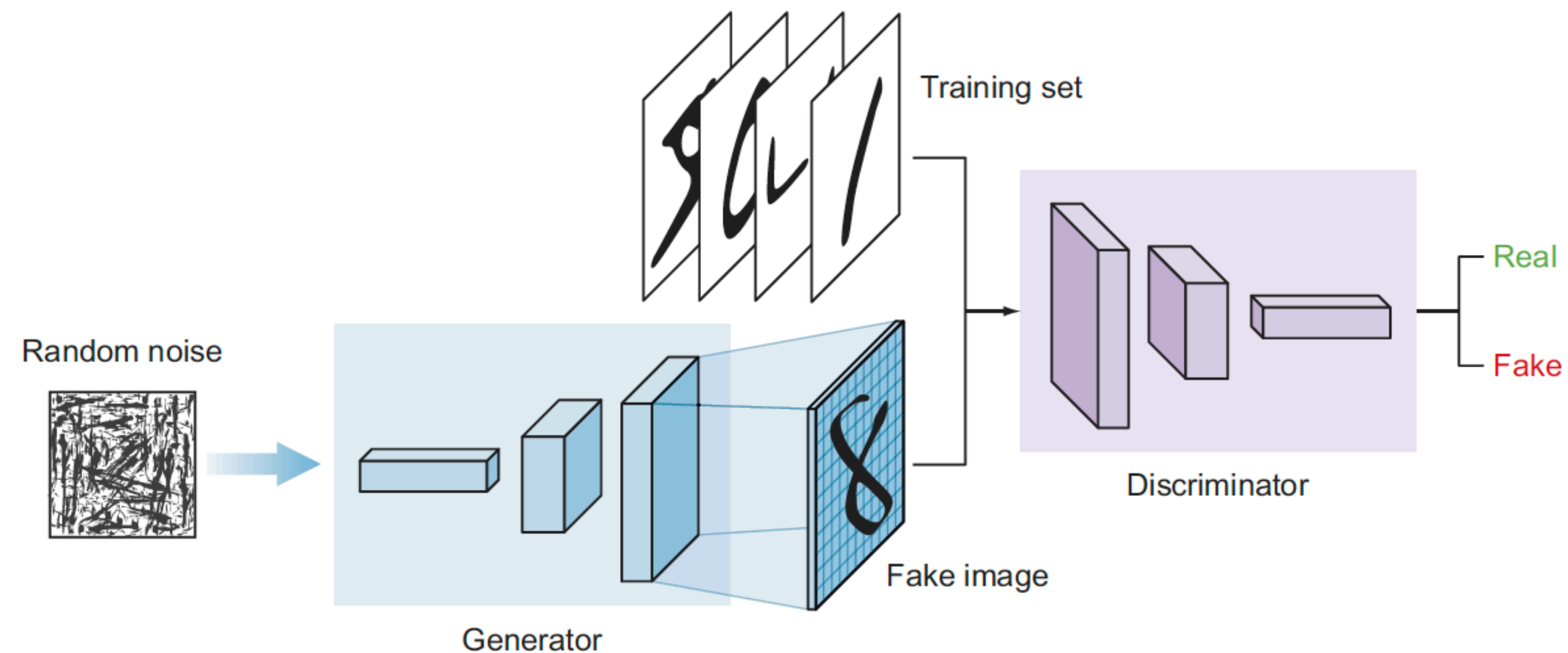
35

- **Problem:** training custom ML models requires extremely large datasets
- **Transfer learning:** take a model that has been trained on the **same type of data for a similar task** and apply it to a specialised task using our own custom data.
 - **Same data:** same data modality. same types of images (e.g. professional pictures vs. Social media pictures)
 - **Similar tasks:** if you need a new object classification model, use a model pre-trained for object classification

Advanced Computer Vision Techniques

Generative Adversarial Networks

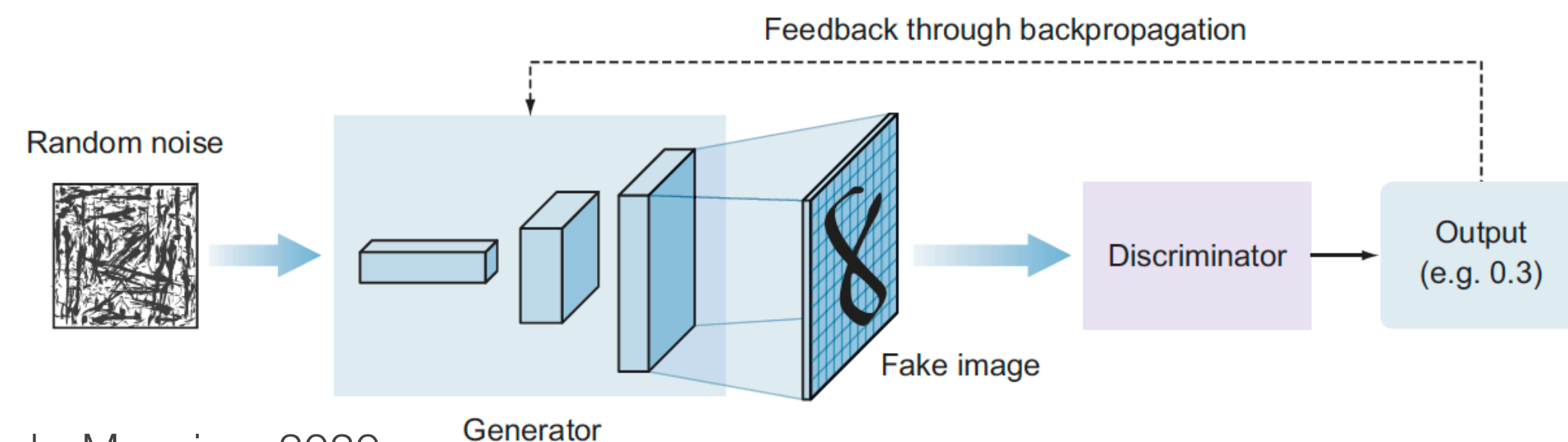
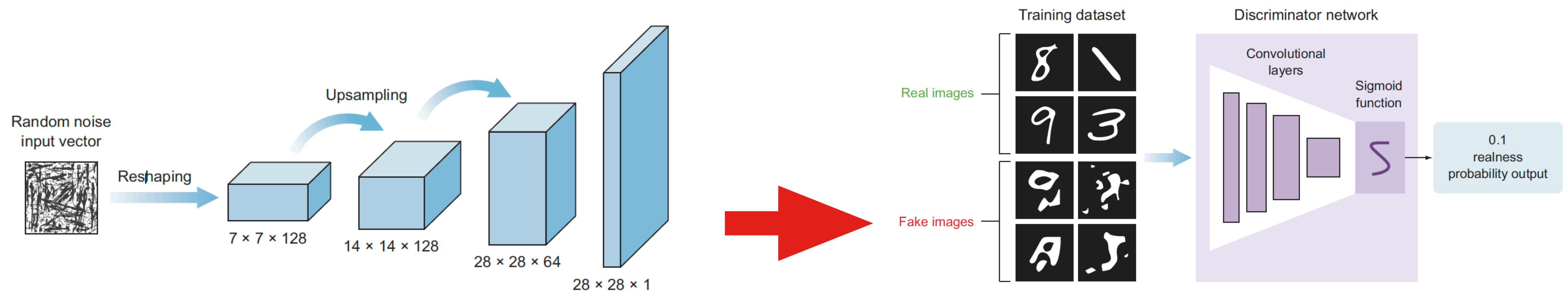
- Learn patterns from the training dataset and create new images that have a similar distribution of the training set
- Two deep neural networks that compete with each other
 - The generator tries to convert random noise into observations that look as if they have been sampled from the original dataset
 - The discriminator tries to predict whether an observation comes from the original dataset or is one of the generator's forgeries



Generative Adversarial Networks

- The generator's architecture looks like an inverted CNN that starts with a narrow input and is upsampled a few times until it reaches the desired size

- The discriminator's model is a typical classification neural network that aims to classify images generated by the generator as real or fake



Which face is real? - <https://www.whichfaceisreal.com/>

PLAY

ABOUT

METHODS

LEARN

PRESS

CONTACT

BOOK

CALLING BS

Click on the person who is real.



■ Try this

<https://thispersondoesnotexist.com/> 38

Image super-resolution GAN

- A good technical summary

<https://blog.paperspace.com/image-super-resolution/>



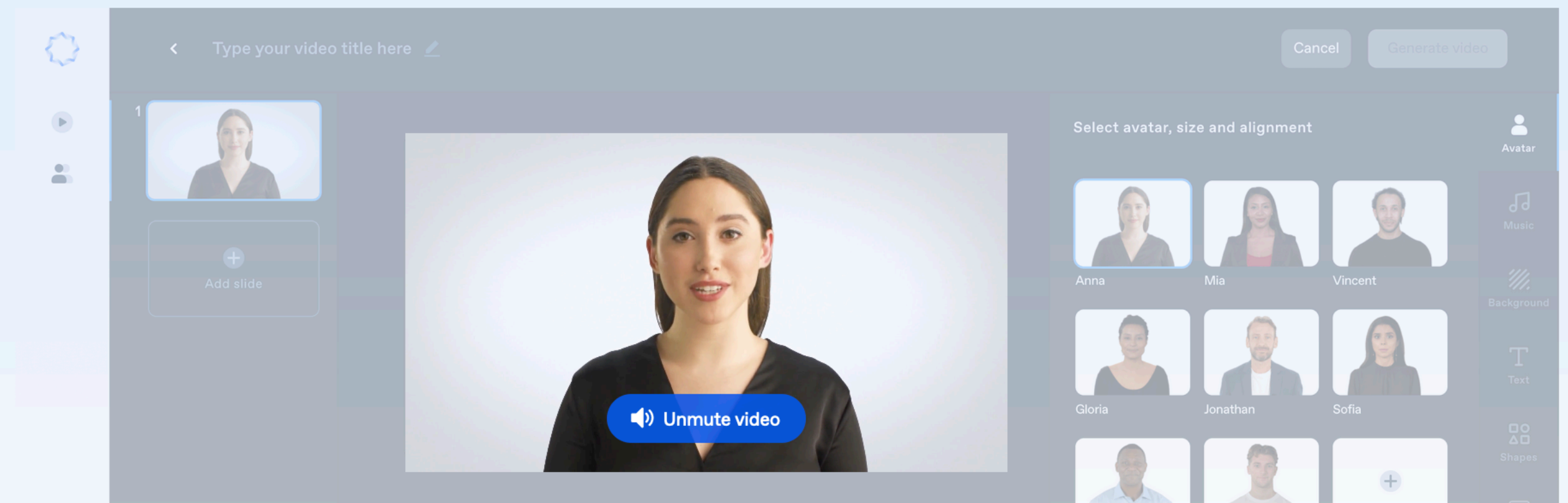
<https://newatlas.com/super-resolution-weizmann-institute/23486/>

Synthetic Video Generation



**Say goodbye to cameras,
microphones and actors!**

Create professional AI videos from text in 60+ languages.



Text-to-image Generation

TEXT PROMPT an illustration of a baby daikon radish in a tutu walking a dog

AI-GENERATED IMAGES



Edit prompt or view more images ↓

TEXT PROMPT an armchair in the shape of an avocado. . . .

AI-GENERATED IMAGES



Edit prompt or view more images ↓

- <https://openai.com/blog/dall-e/>



- ML-generated painting sold for \$432,500
- The network trained on a dataset of 15,000 portraits painted between the fourteenth and twentieth centuries
- Network “learned” the style, and generated a new painting

Neural Style Transfer



Content Image

+



Style Image

=



Stylized Result

<https://replicate.com/rinongal/stylegan-nada>



Deepfakes



Admin

Week 3 Tasks

- Have fun with the first assignment !
- Please contribute Week 3 questions - we will share the link later
- See you on Friday!

Machine Learning For Design

Lecture 3 - Machine Learning for Images

Alessandro Bozzon

16/02/2022

mfd-io@tudelft.nl
www.ml4design.com

Credits

- CMU Computer Vision course - Matthew O'Toole. <http://16385.courses.cs.cmu.edu/spring2022/>
- CIS 419/519 Applied Machine Learning. Eric Eaton, Dinesh Jayaraman. <https://www.seas.upenn.edu/~cis519/spring2020/>
- Deep Learning Patterns and Practices - Andrew Ferlitsch, Manning, 2021
- Machine Learning Design Patterns - Lakshmanan, Robinson, Munn, 2020
- Grokking Machine Learning. Luis G. Serrano. Manning, 2021
- Deep Learning for Vision Systems. Mohamed Elgendy. Manning, 2020