

Lecture 13

Introduction to Convolutional Neural Networks Part 1

STAT 479: Deep Learning, Spring 2019

Sebastian Raschka

<http://stat.wisc.edu/~sraschka/teaching/stat479-ss2019/>

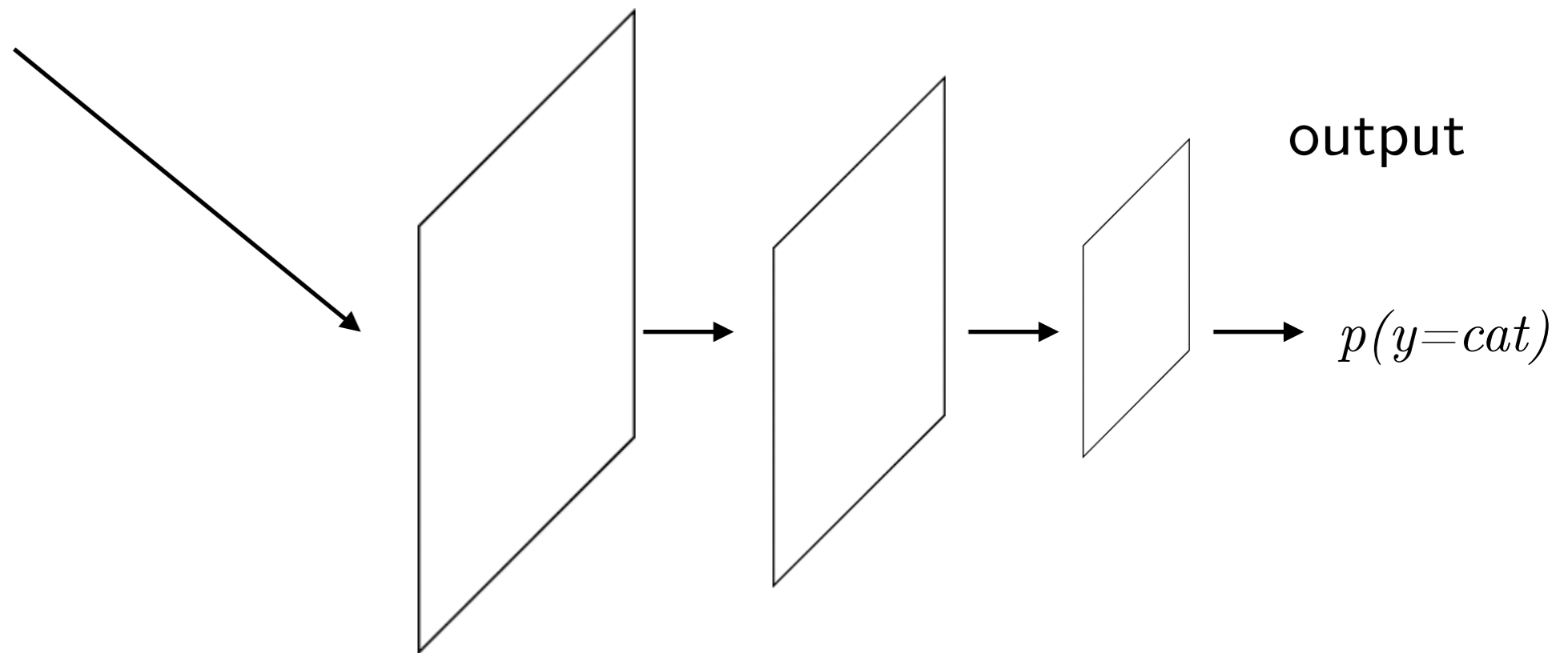
CNNs for Image Classification



Image Source: twitter.com/%2Fcats&psig=AOvVaw30_o-PCM-K21DiMAJQimQ4&ust=1553887775741551



Image Source: <https://www.pinterest.com/pin/244742560974520446>



Some Applications of CNNs

Predicting people's age

(160,000 images from the RenRen Social Network)

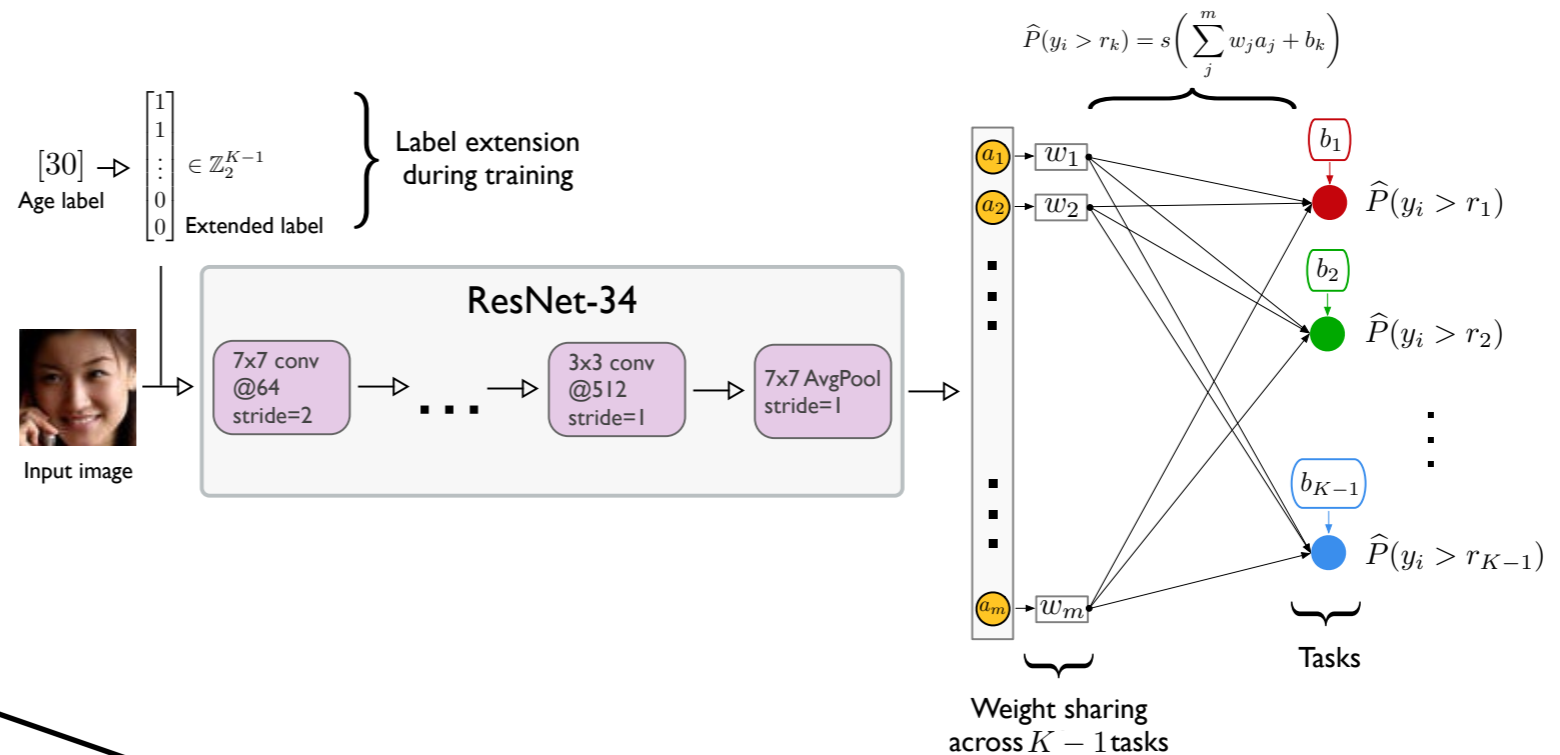


Figure 2. Illustration of the Consistent Rank Logits CNN (CORAL-CNN) used for age prediction. From the estimated probability values, the binary labels are obtained via Eq. (5) and converted to the age label via Eq. (1).

Table 1. Age prediction errors on the test sets *without* task importance weighting.

Method	Random Seed	MORPH-2		AFAD		UTKFace		CACD	
		MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE
CE-CNN	0	3.40	4.88	3.98	5.55	6.57	9.16	6.18	8.86
	1	3.39	4.87	4.00	5.57	6.24	8.69	6.10	8.79
	2	3.37	4.87	3.96	5.50	6.29	8.78	6.13	8.87
	AVG \pm SD	3.39 \pm 0.02	4.89 \pm 0.01	3.98 \pm 0.02	5.54 \pm 0.04	6.37 \pm 0.18	8.88 \pm 0.25	6.14 \pm 0.04	8.84 \pm 0.04
OR-CNN (Niu et al., 2016)	0	2.98	4.26	3.66	5.10	5.71	8.11	5.53	7.91
	1	2.98	4.26	3.69	5.13	5.80	8.12	5.53	7.98
	2	2.96	4.20	3.68	5.14	5.71	8.11	5.49	7.89
	AVG \pm SD	2.97 \pm 0.01	4.24 \pm 0.03	3.68 \pm 0.02	5.13 \pm 0.02	5.74 \pm 0.05	8.08 \pm 0.06	5.52 \pm 0.02	7.93 \pm 0.05
CORAL-CNN (ours)	0	2.68	3.75	3.49	4.82	5.46	7.61	5.56	7.80
	1	2.63	3.66	3.46	4.83	5.46	7.63	5.37	7.64
	2	2.61	3.64	3.52	4.91	5.48	7.63	5.25	7.53
	AVG \pm SD	2.64 \pm 0.04	3.68 \pm 0.06	3.49 \pm 0.03	4.85 \pm 0.05	5.47 \pm 0.01	7.62 \pm 0.01	5.39 \pm 0.16	7.66 \pm 0.14

Cao, Wenzhi, Vahid Mirjalili, and Sebastian Raschka. "Consistent Rank Logits for Ordinal Regression with Convolutional Neural Networks." *arXiv preprint arXiv:1901.07884* (2019).

Some Applications of CNNs

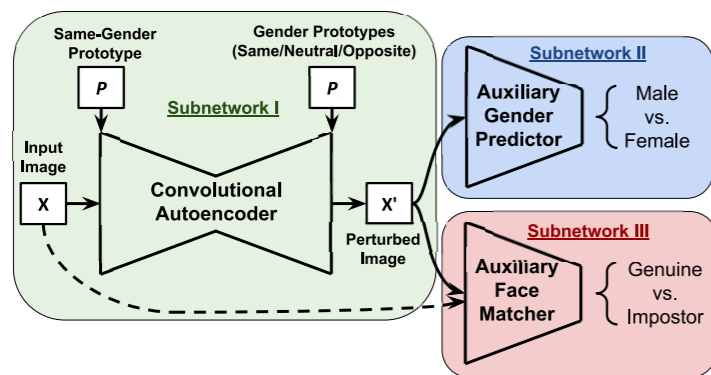


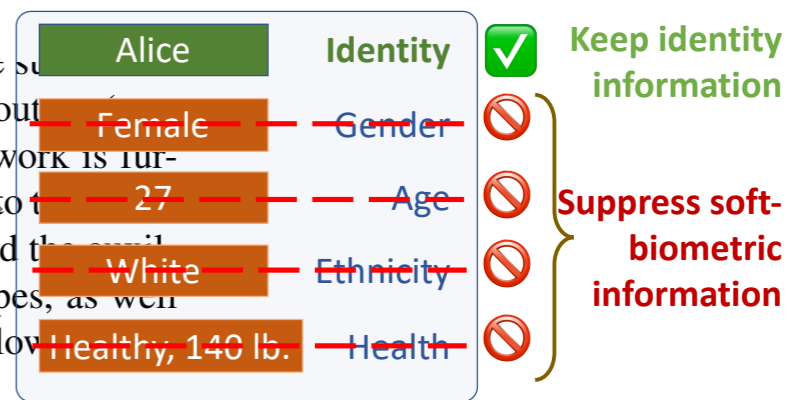
Figure 1. Schematic representation of the semi-adversarial neural network architecture designed to derive perturbations that are able to confound gender classifiers while still allowing biometric matchers to perform well. The overall network consists of three sub-components: a convolutional autoencoder (subnetwork I), an auxiliary gender classifier (subnetwork II), and an auxiliary matcher (subnetwork III).

2.2.1 Convolutional autoencoder

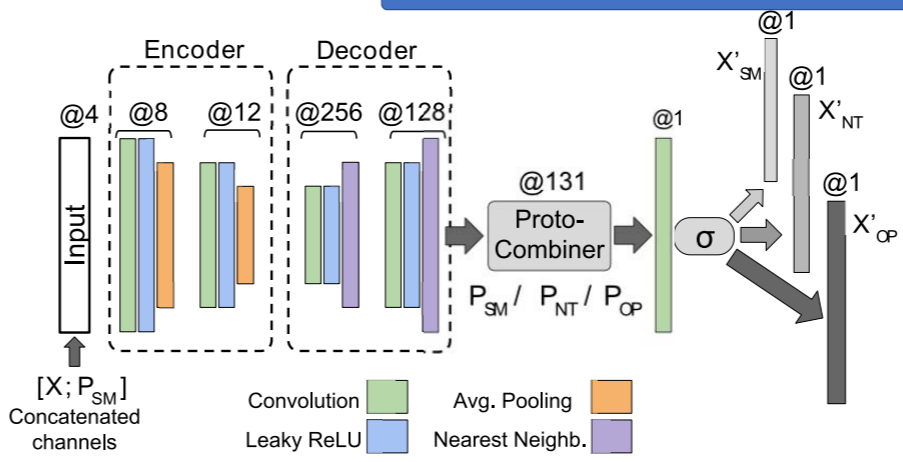
The architecture of the convolutional autoencoder sub-

ages and the perturbed images by minimizing the su-
 element-wise cross entropy loss between the input and output (perturbed) images. After pre-
 trained by passing it through the autoencoder sub-network is fur-
 ther trained by passing it through the auxiliary gender predictor and the auxiliary face matcher (Fig. 1). The autoencoder, as well as the two sub-networks, are described in the following sections.

GDPR Compliance



Specific Purpose: identification and/or recognition
Data Minimization



Vahid Mirjalili, Sebastian Raschka, Anoop Namboodiri, and Arun Ross (2018) *Semi-adversarial networks: Convolutional autoencoders for imparting privacy to face images*. Proc. of 11th IAPR International Conference on Biometrics (ICB 2018), Gold Coast, Australia.

Vahid Mirjalili, Sebastian Raschka, and Arun Ross (2018) *Gender Privacy: An Ensemble of Semi Adversarial Networks for Confounding Arbitrary Gender Classifiers*. 9th IEEE International Conference on Biometrics: Theory, Applications, and Systems

Some Applications of CNNs



Dermatologist-level classification of skin cancer

An artificial intelligence trained to classify images of skin lesions as benign lesions or malignant skin cancers achieves the accuracy of board-certified dermatologists.

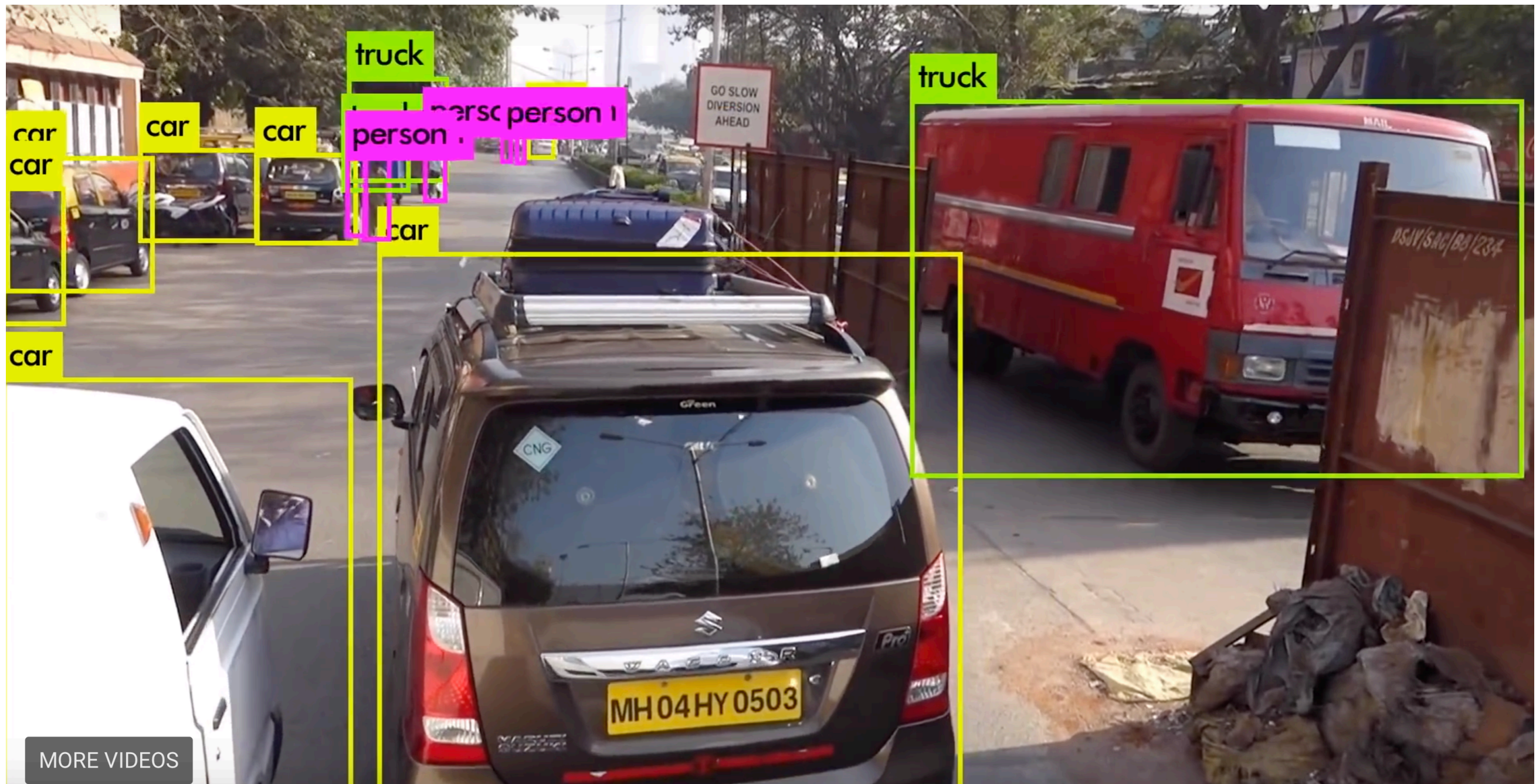
In this work, we pretrain a deep neural network at general object recognition, then fine-tune it on a dataset of ~130,000 skin lesion images comprised of over 2000 diseases.

[FULL NATURE ARTICLE >](#)

[OPEN-ACCESS PDF >](#)

Esteva, Andre, Brett Kuprel, Roberto A. Novoa, Justin Ko, Susan M. Swetter, Helen M. Blau, and Sebastian Thrun. "Dermatologist-level classification of skin cancer with deep neural networks." *Nature* 542, no. 7639 (2017): 115.

Object Detection



Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 779-788).

Object Segmentation

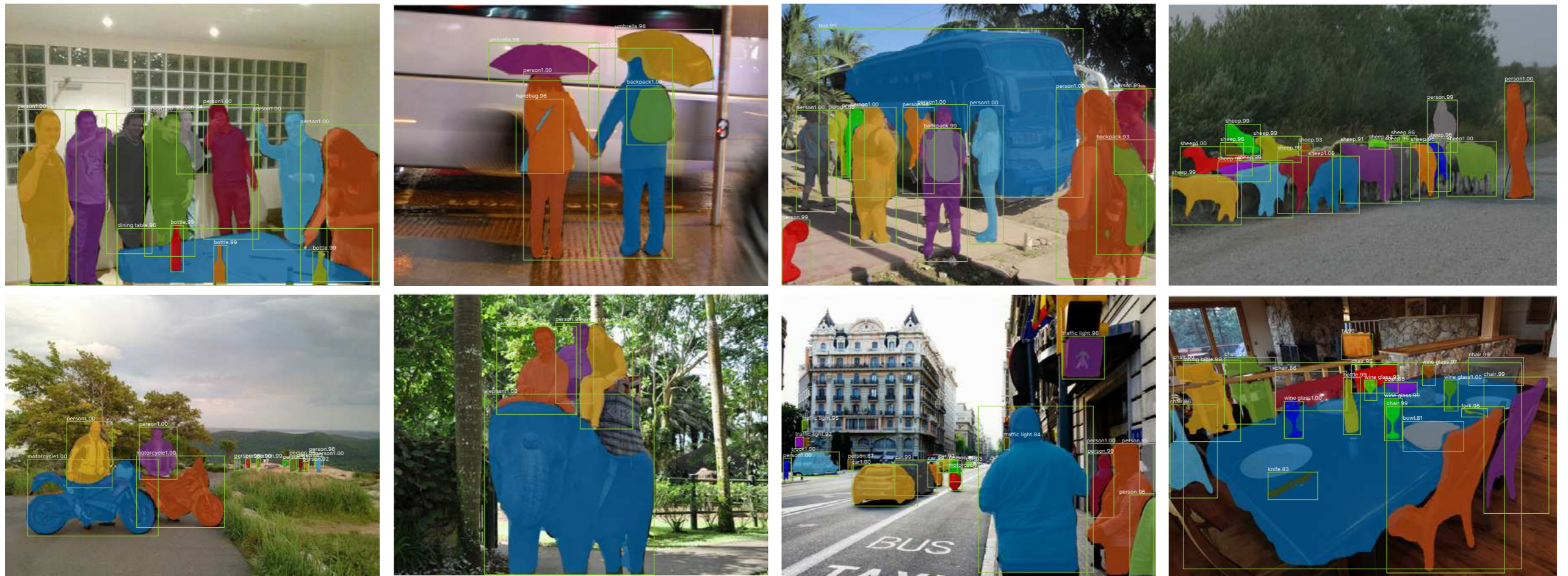


Figure 2. **Mask R-CNN** results on the COCO test set. These results are based on ResNet-101 [15], achieving a *mask AP* of 35.7 and running at 5 fps. Masks are shown in color, and bounding box, category, and confidences are also shown.

He, Kaiming, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. "Mask R-CNN." In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2961-2969. 2017.

Why Image Classification is Hard

Different lighting, contrast, viewpoints, etc.

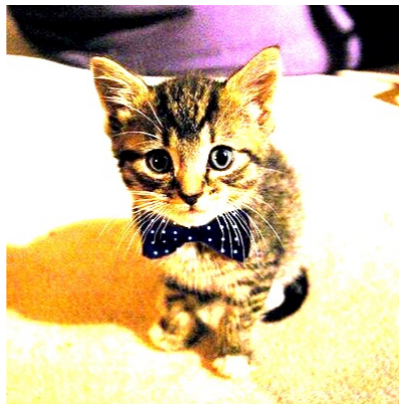
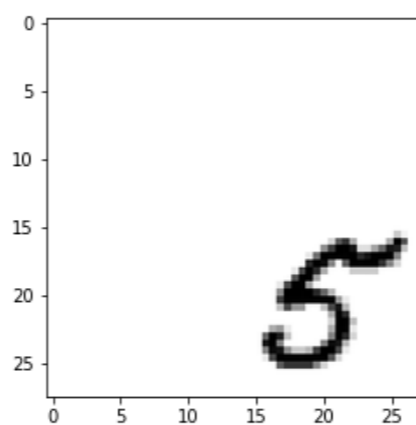
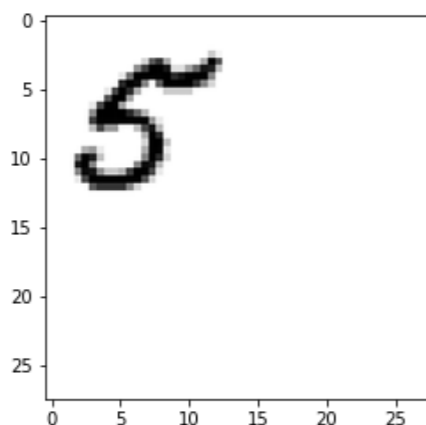


Image Source:
twitter.com/%2Fcats&psig=AOvVaw30_o-PCM-K21DiMAJQimQ4&ust=1553887775741551

Image Source: https://www.123rf.com/photo_76714328_side-view-of-tabby-cat-face-over-white.html

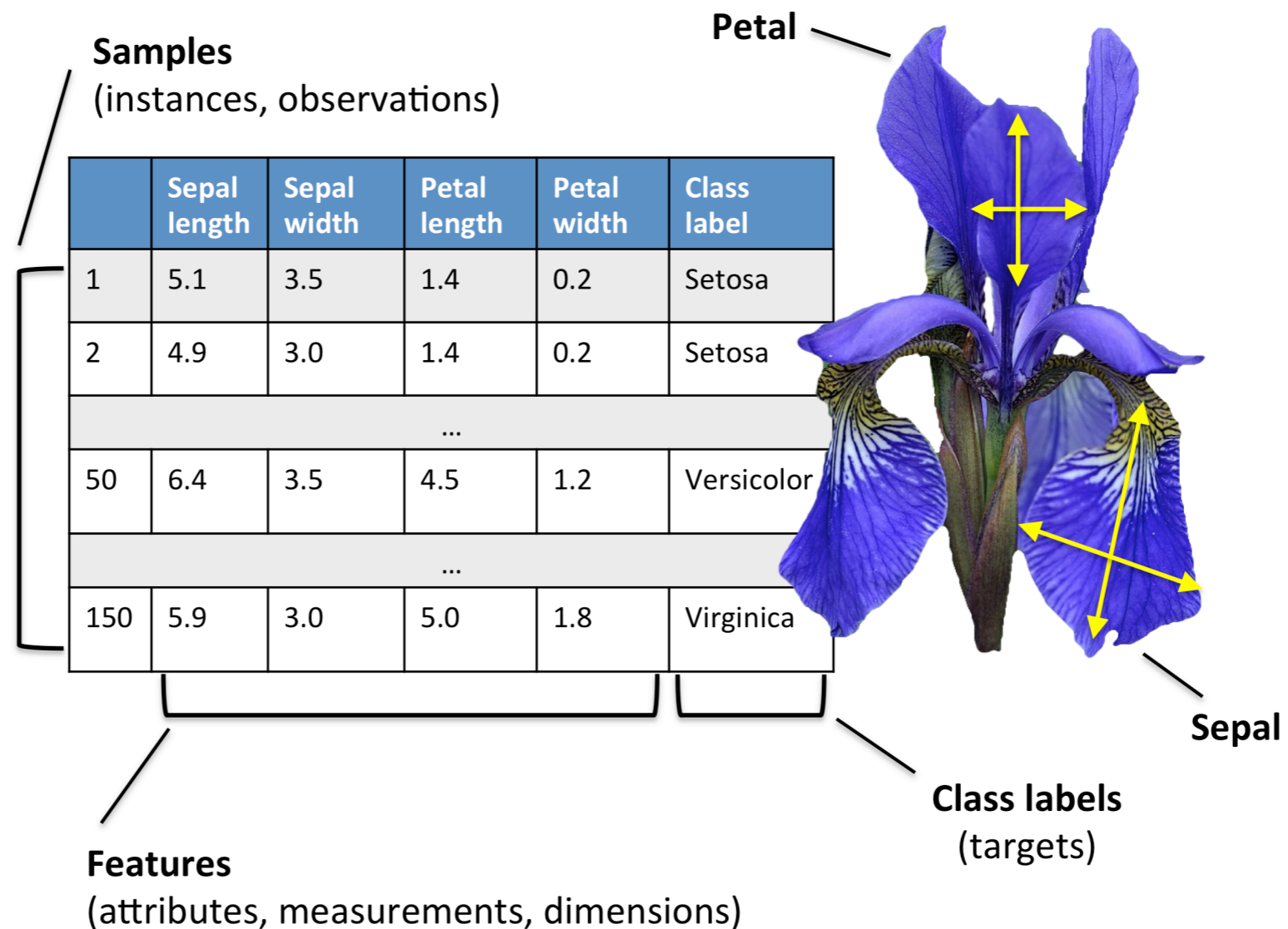
Or even simple translation



This is hard for traditional methods like multi-layer perceptrons, because the prediction is basically based on a sum of pixel intensities

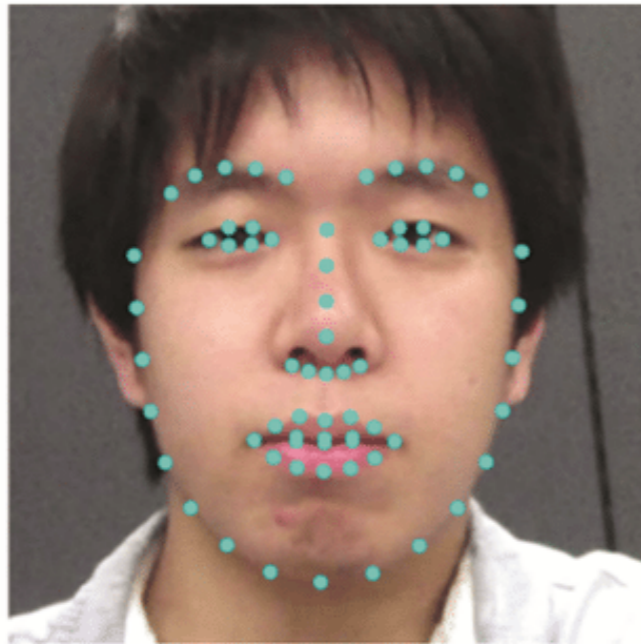
Traditional Approaches

a) Use hand-engineered features



Traditional Approaches

a) Use hand-engineered features



(a) Detected facial keypoints



(b) Facial organ keypoints

Sasaki, K., Hashimoto, M., & Nagata, N. (2016). Person Invariant Classification of Subtle Facial Expressions Using Coded Movement Direction of Keypoints. In *Video Analytics. Face and Facial Expression Recognition and Audience Measurement* (pp. 61-72). Springer, Cham.

Traditional Approaches

b) Preprocess images (centering, cropping, etc.)



Image Source: <https://www.tokkoro.com/2827328-cat-animals-nature-feline-park-green-trees-grass.html>

Main Concepts Behind Convolutional Neural Networks

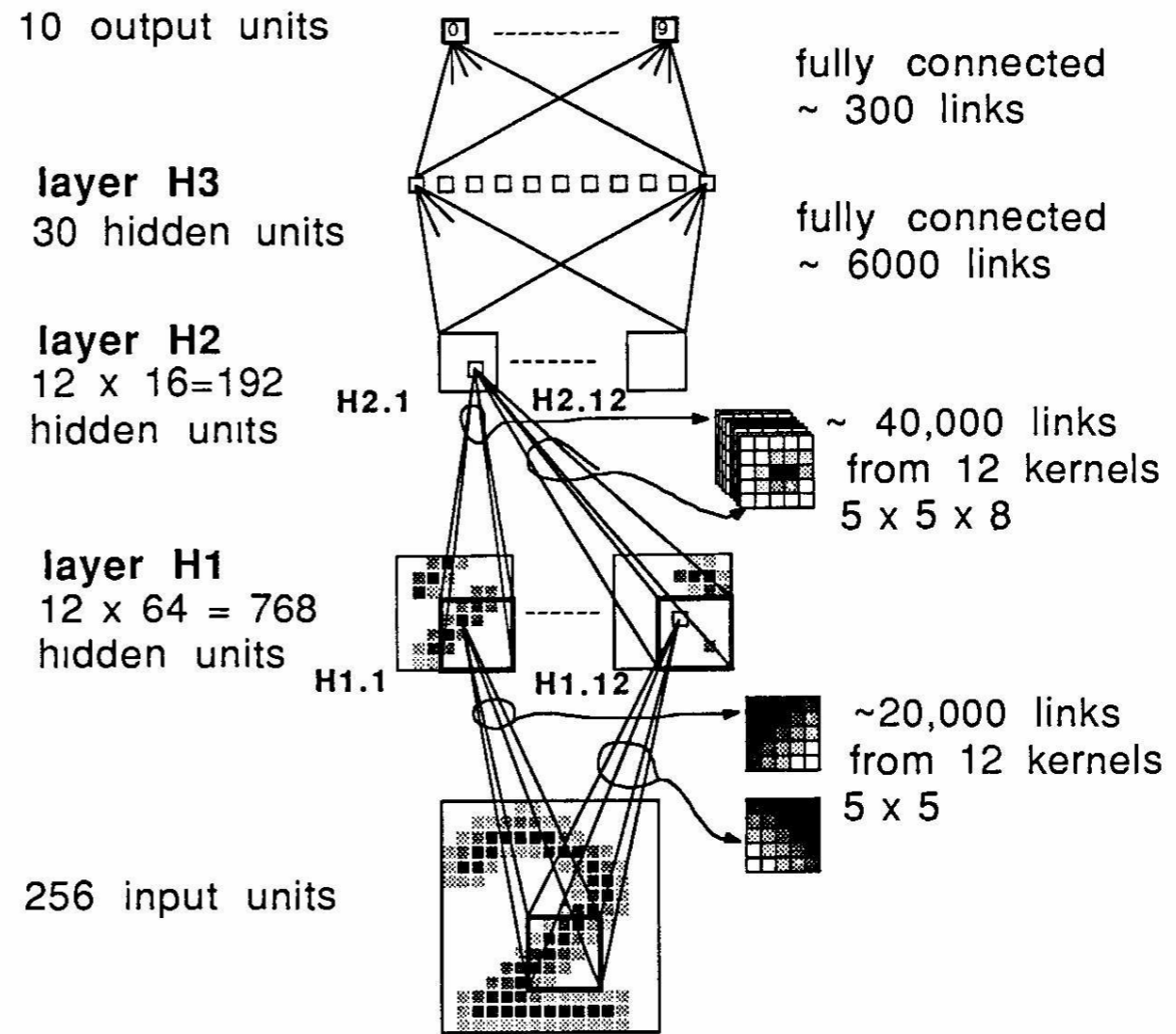
- **Sparse-connectivity:** A single element in the feature map is connected to only a small patch of pixels. (This is very different from connecting to the whole input image, in the case of multi-layer perceptrons.)
- **Parameter-sharing:** The same weights are used for different patches of the input image.

Image Source: <https://www.tokkoro.com/2827328-cat-animals-nature-feline-park-green-trees-grass.html>

Convolutional Neural Networks

80322-4129 80206
 40004 14310
 3787R 05453
 5502 75216
 35460 44209

1011913485726803226414186
 6359720299299722510046701
 3084111591010615406103631
 1064111030475262009979966
 8912056728557131427955460
 2018750189112993089970984
 0109707597331972015519055
 1075318255182814358090943
 1787521655460554603546055
 18255108503047520439401



Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard and L. D. Jackel: [Backpropagation Applied to Handwritten Zip Code Recognition](#), Neural Computation, 1(4):541-551, Winter 1989.

Convolutional Neural Networks

PROC. OF THE IEEE, NOVEMBER 1998

7

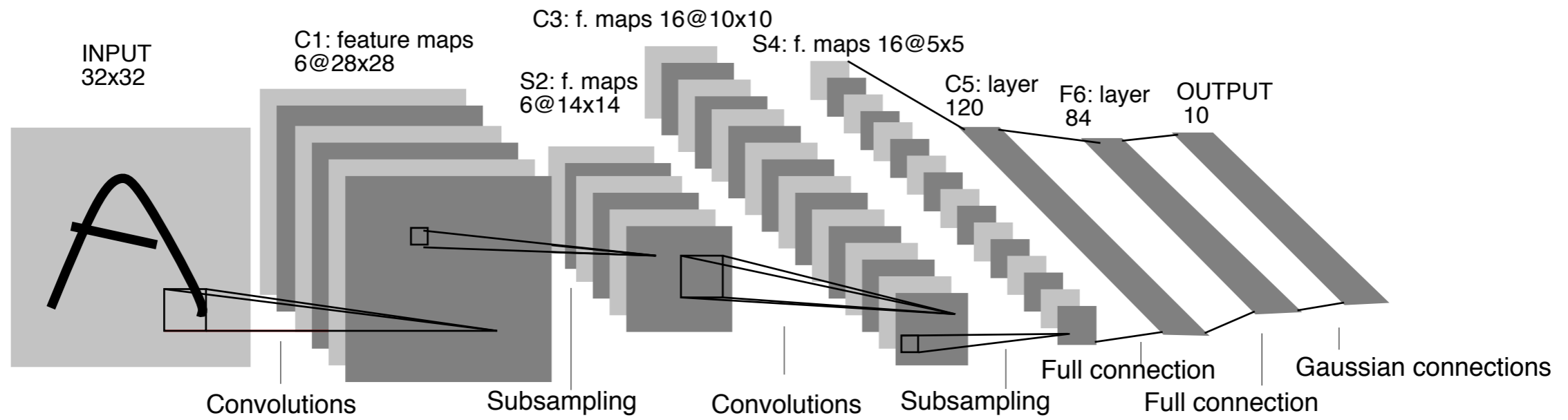


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Yann LeCun, Léon Bottou, Yoshua Bengio and Patrick Haffner: [Gradient Based Learning Applied to Document Recognition](#), Proceedings of IEEE, 86(11):2278–2324, 1998.

Convolutional Neural Networks

PROC. OF THE IEEE, NOVEMBER 1998

7

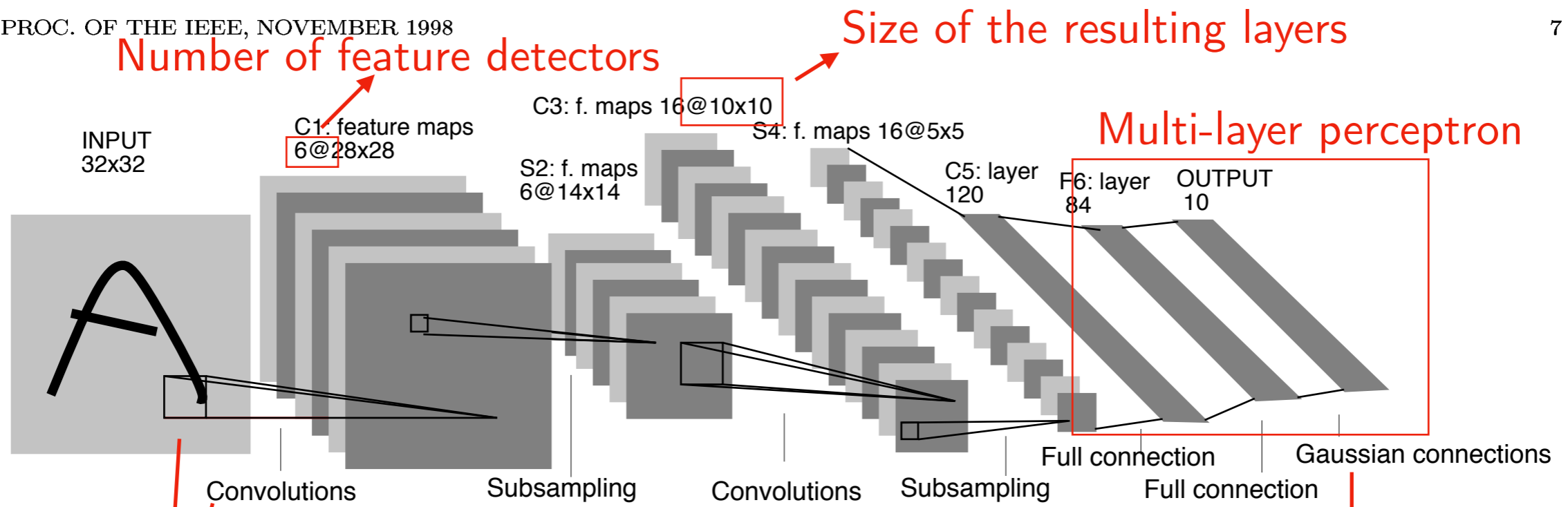


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

nowadays called "pooling"

"Feature detectors" (weight matrices)
that are being reused ("weight sharing")
=> also called "kernel" or "filter"

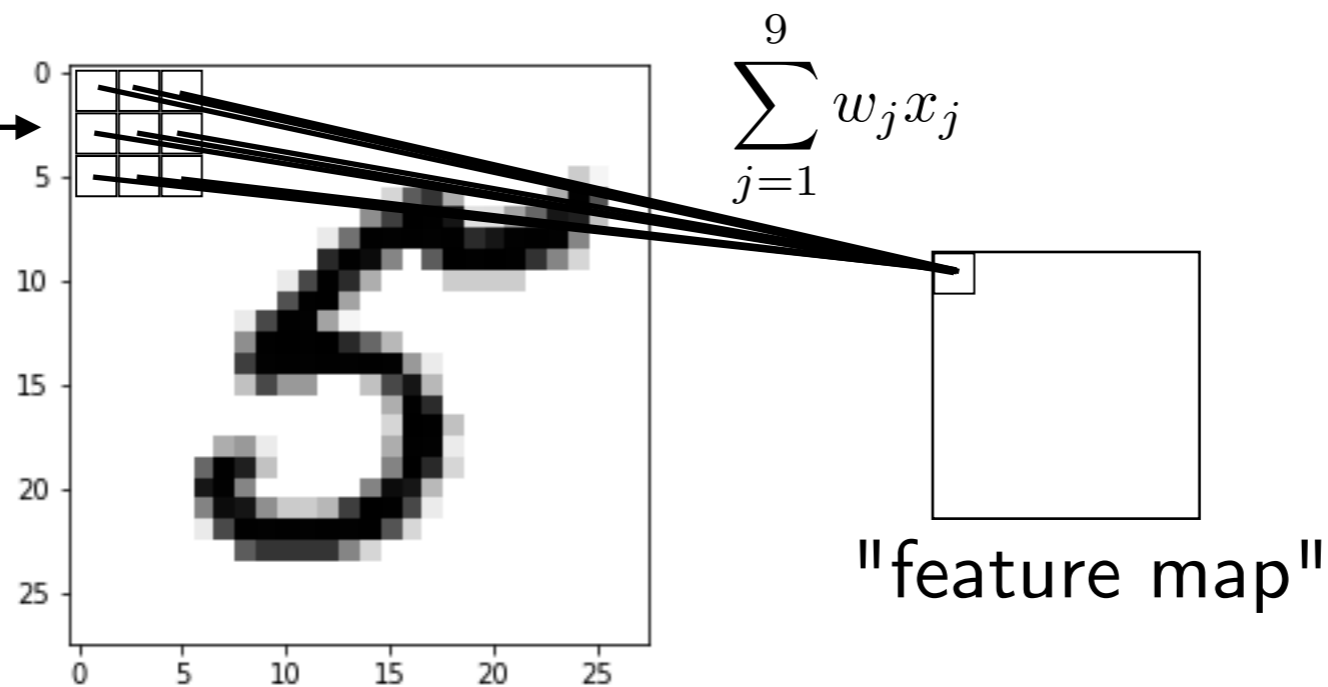
basically a fully-connected
layer + MSE loss
(nowadays better to use
fc-layer + softmax
+ cross entropy)

Yann LeCun, Léon Bottou, Yoshua Bengio and Patrick Haffner: [Gradient Based Learning Applied to Document Recognition](#), Proceedings of IEEE, 86(11):2278–2324, 1998.

Weight Sharing

A "feature detector" (kernel) slides over the inputs to generate a feature map

The pixels are referred to as "receptive field"

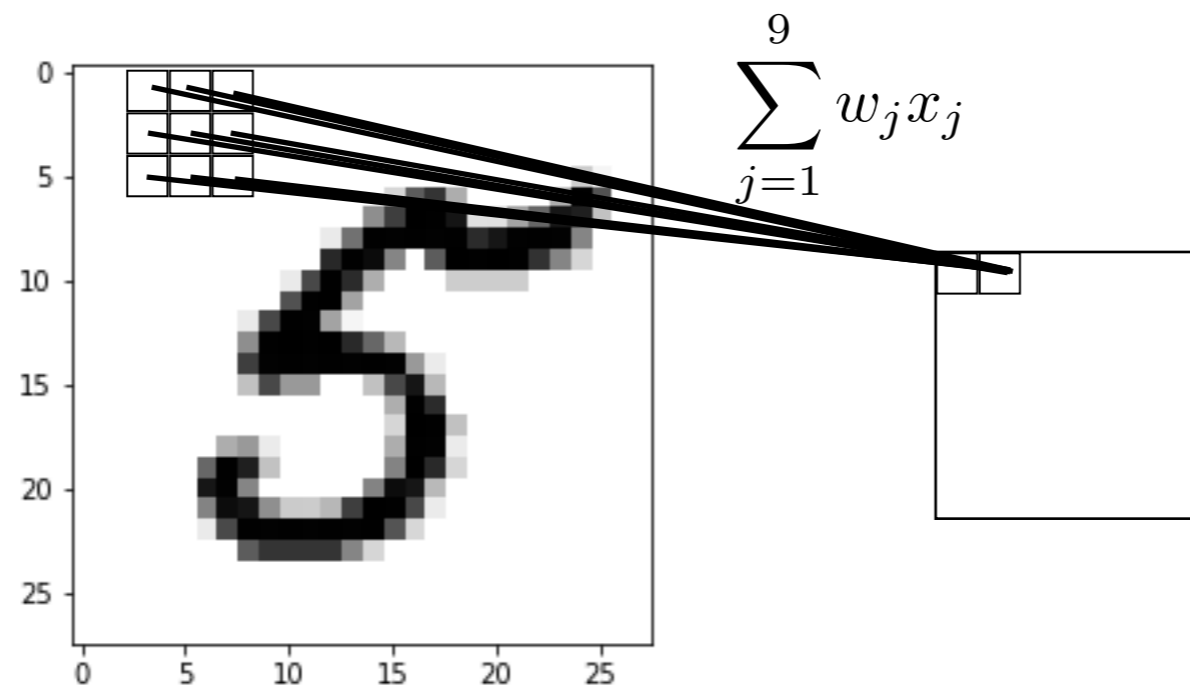


Rationale: A feature detector that works well in one region may also work well in another region

Plus, it is a nice reduction in parameters to fit

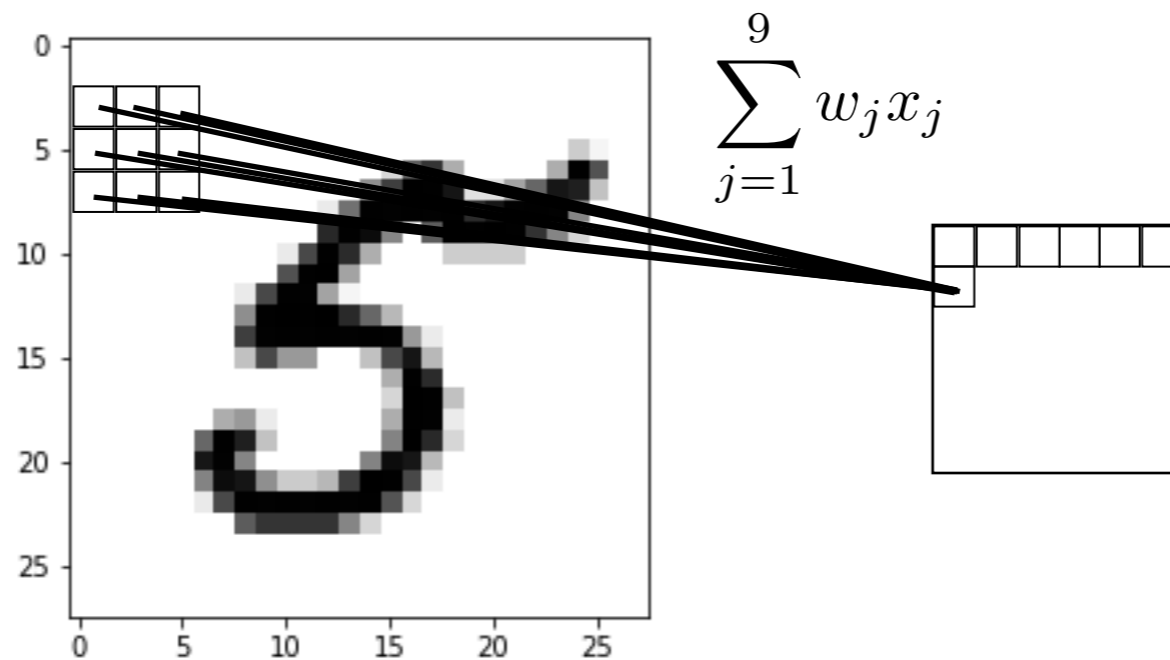
Weight Sharing

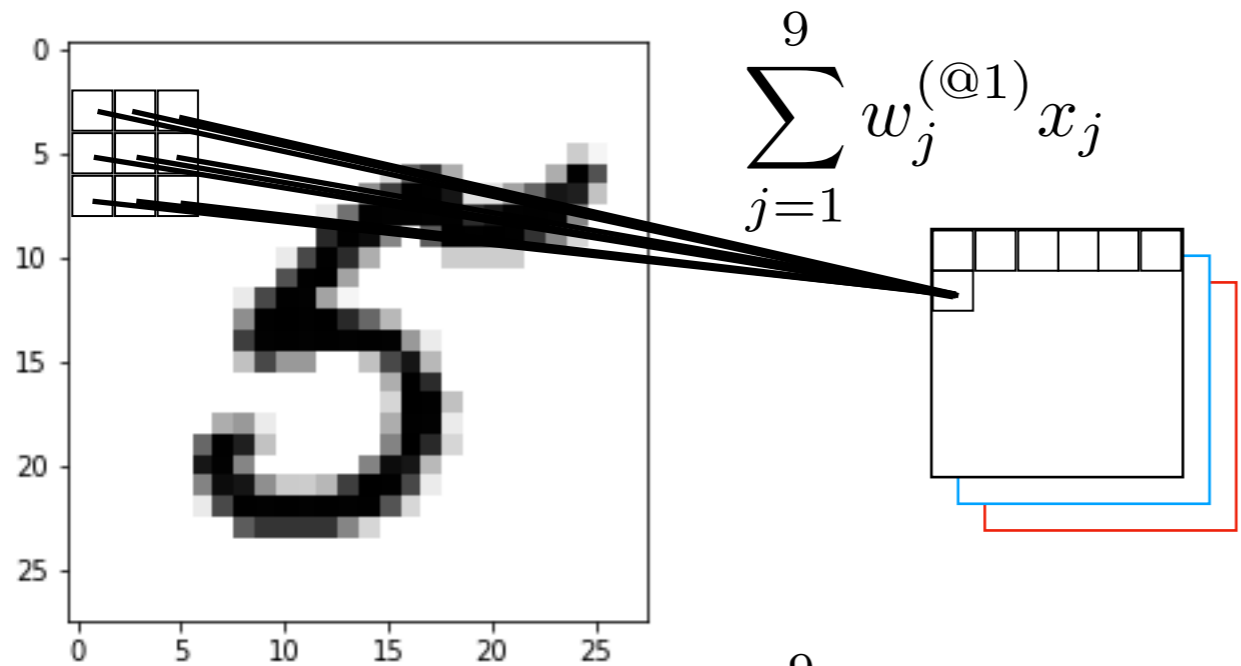
A "feature detector" (kernel) slides over the inputs to generate a feature map



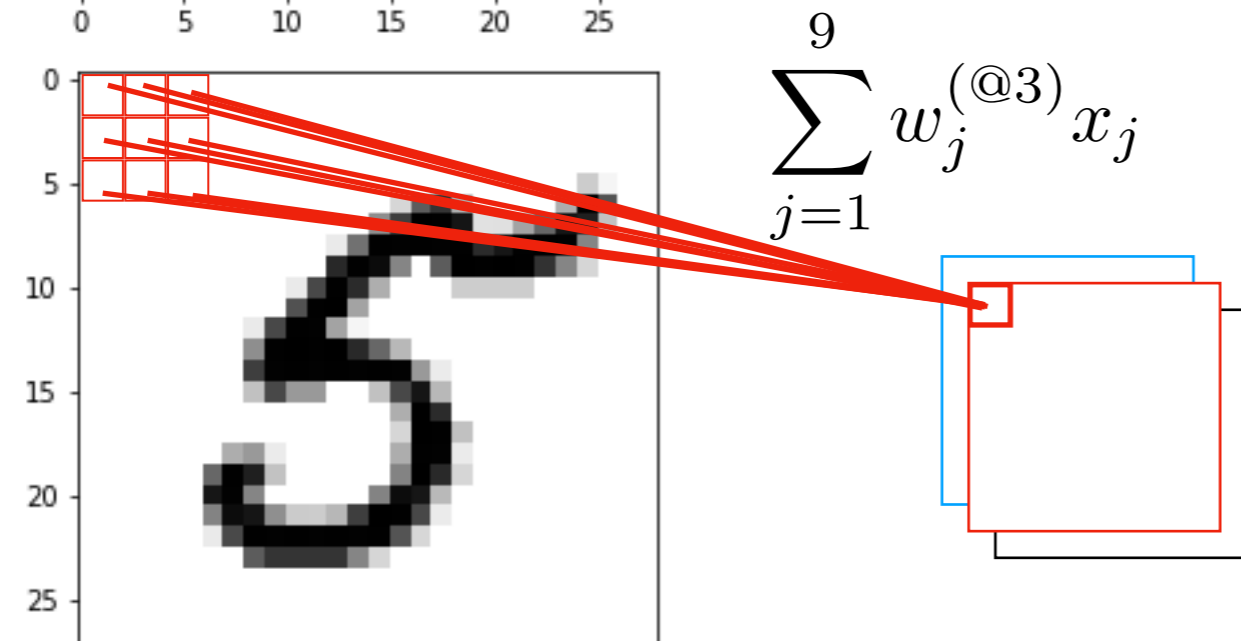
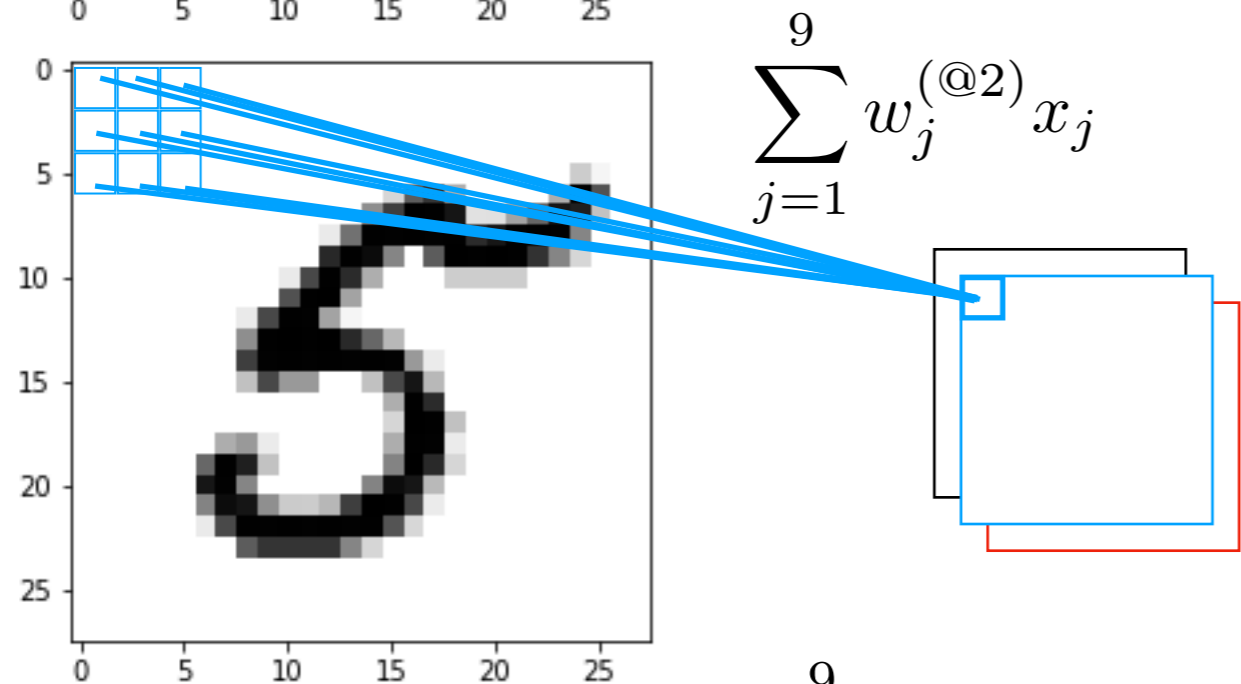
Weight Sharing

A "feature detector" (kernel) slides over the inputs to generate a feature map



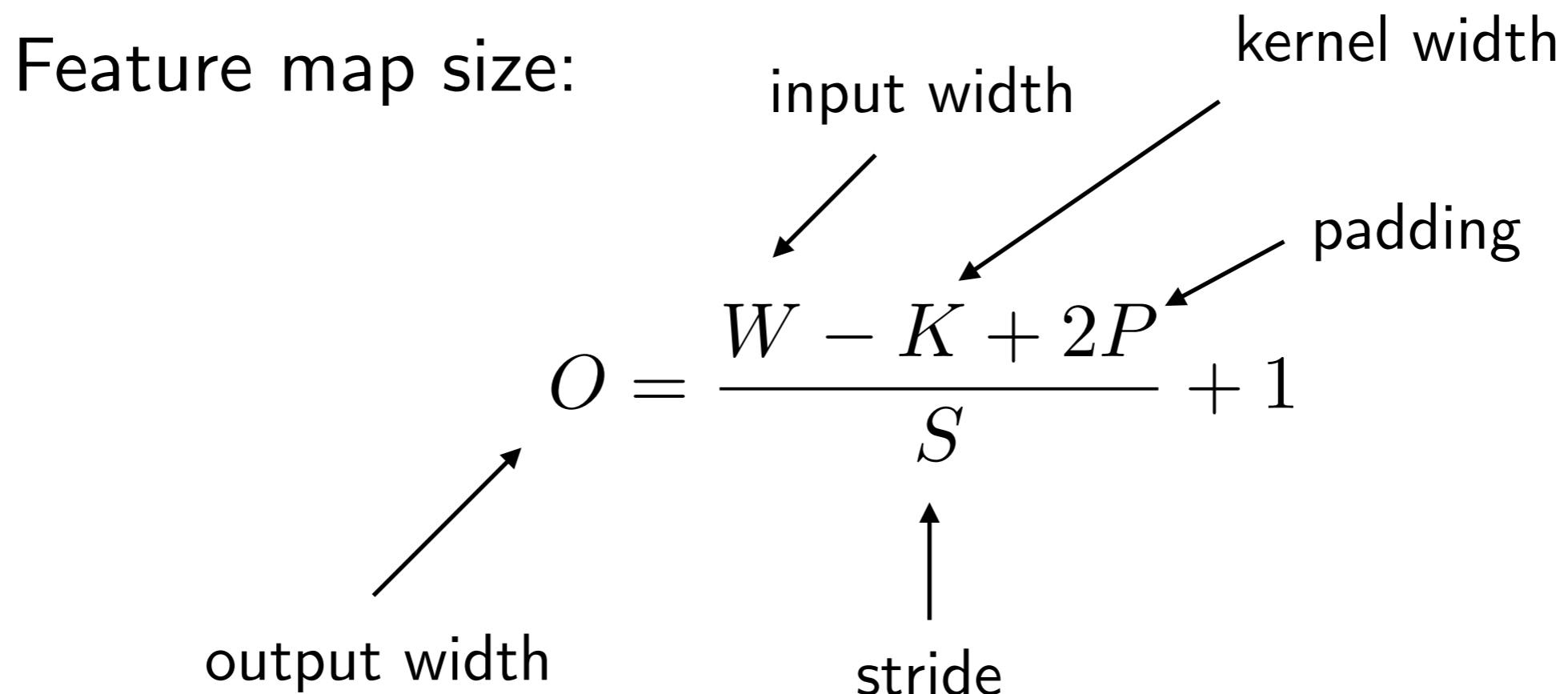


Multiple "feature detectors" (kernels) are used to create multiple feature maps



Size Before and After Convolutions

Feature map size:

$$O = \frac{W - K + 2P}{S} + 1$$


The diagram illustrates the formula for the output width O of a convolution. The variables are defined as follows:

- W : input width
- K : kernel width
- P : padding
- S : stride

The output width O is calculated as the integer division of $W - K + 2P$ by S , plus 1.

Hidden Layers

PROC. OF THE IEEE, NOVEMBER 1998

7

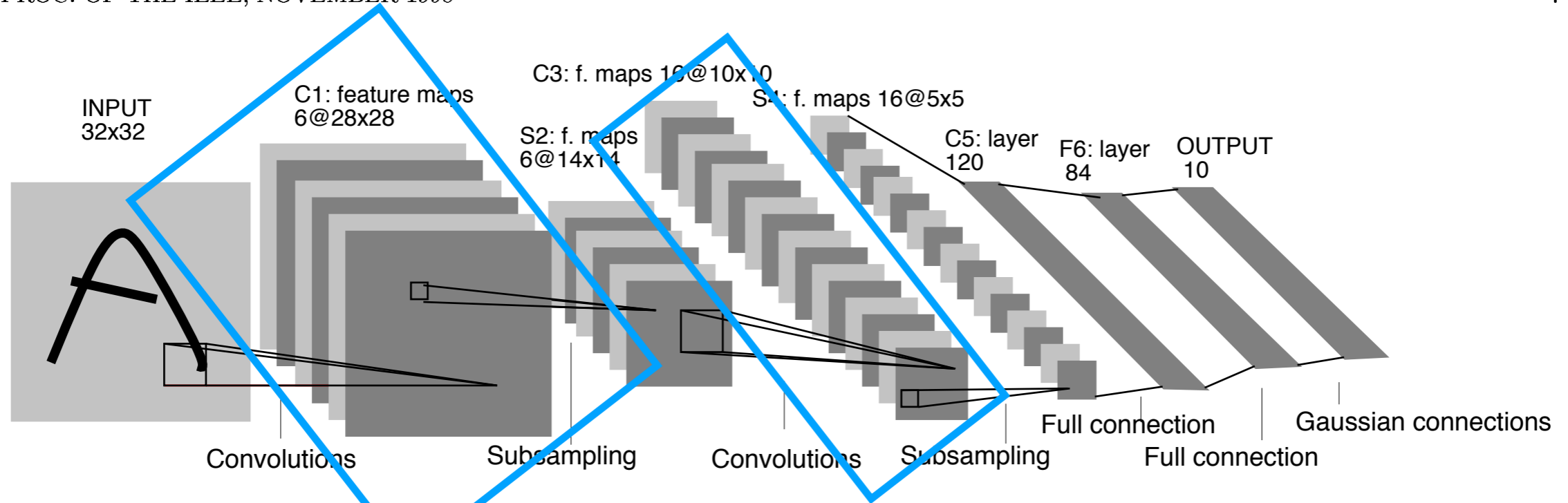


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

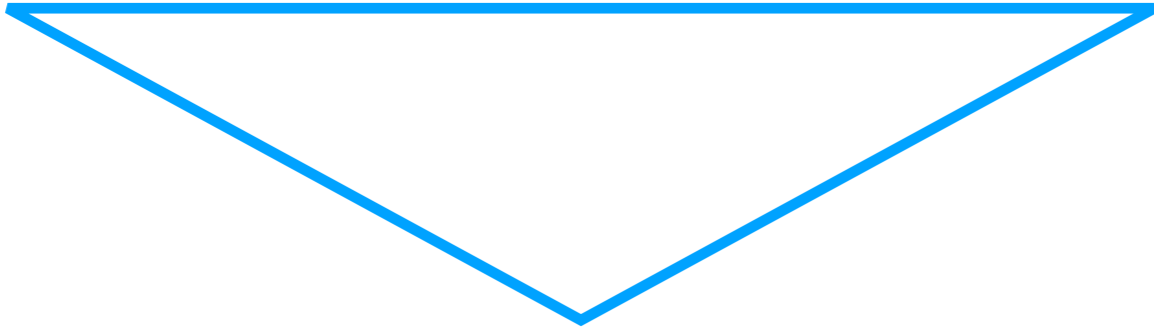
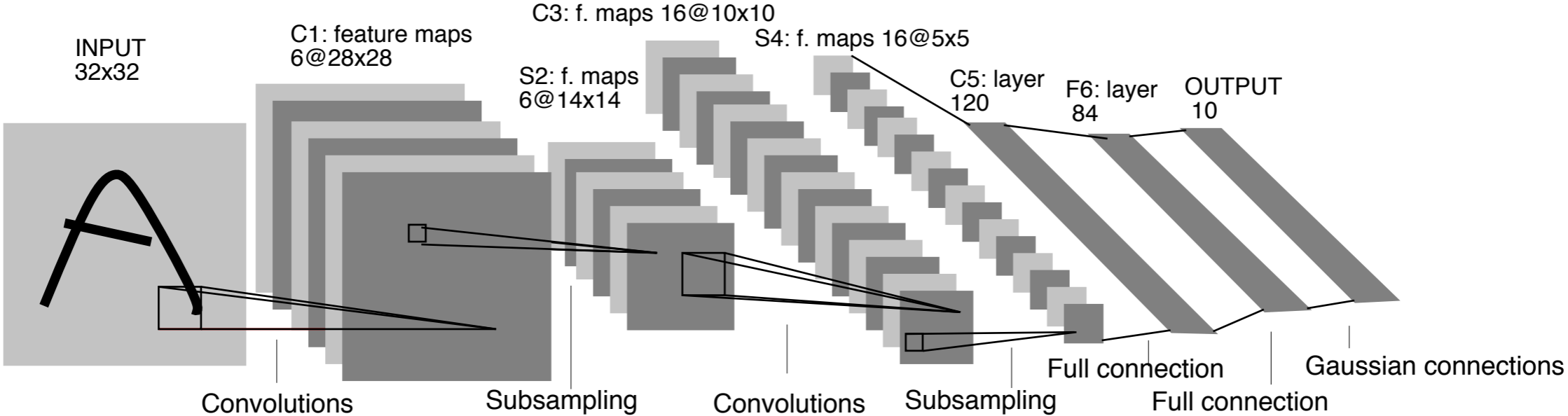
Each "bunch" of feature maps represents one hidden layer in the neural network.

Counting the FC layers, this network has 5 layers

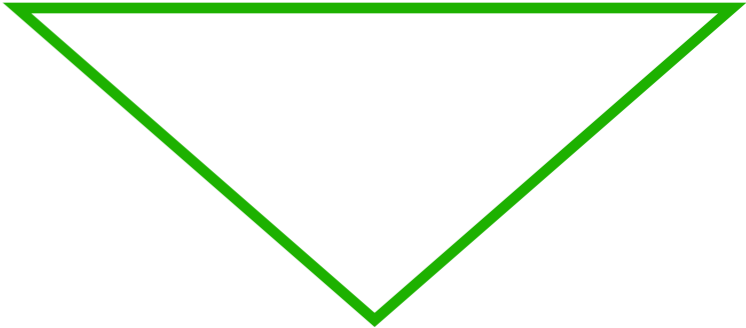
Hidden Layers

PROC. OF THE IEEE, NOVEMBER 1998

7

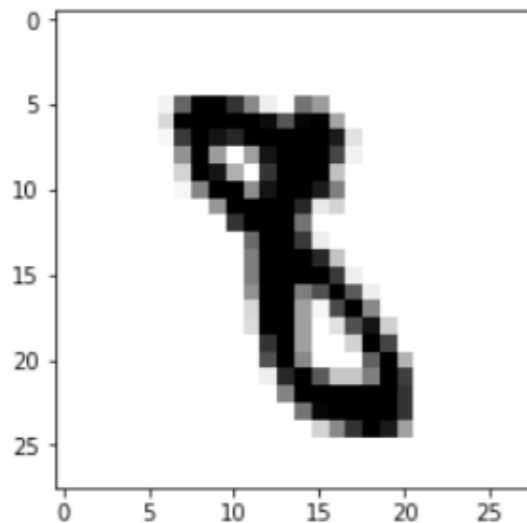


"Automatic feature extractor"



"Regular classifier"

Kernel Dimensions



For a grayscale image with a 5x5 feature detector (kernel), we have the following dimensions (number of parameters to learn)

```
a.shape
```

```
(1, 28, 28)
```

```
import torch
```

```
conv = torch.nn.Conv2d(in_channels=1,  
                       out_channels=8,  
                       kernel_size=(5, 5),  
                       stride=(1, 1))
```

```
conv.weight.size()
```

```
torch.Size([8, 1, 5, 5])
```

```
conv.bias.size()
```

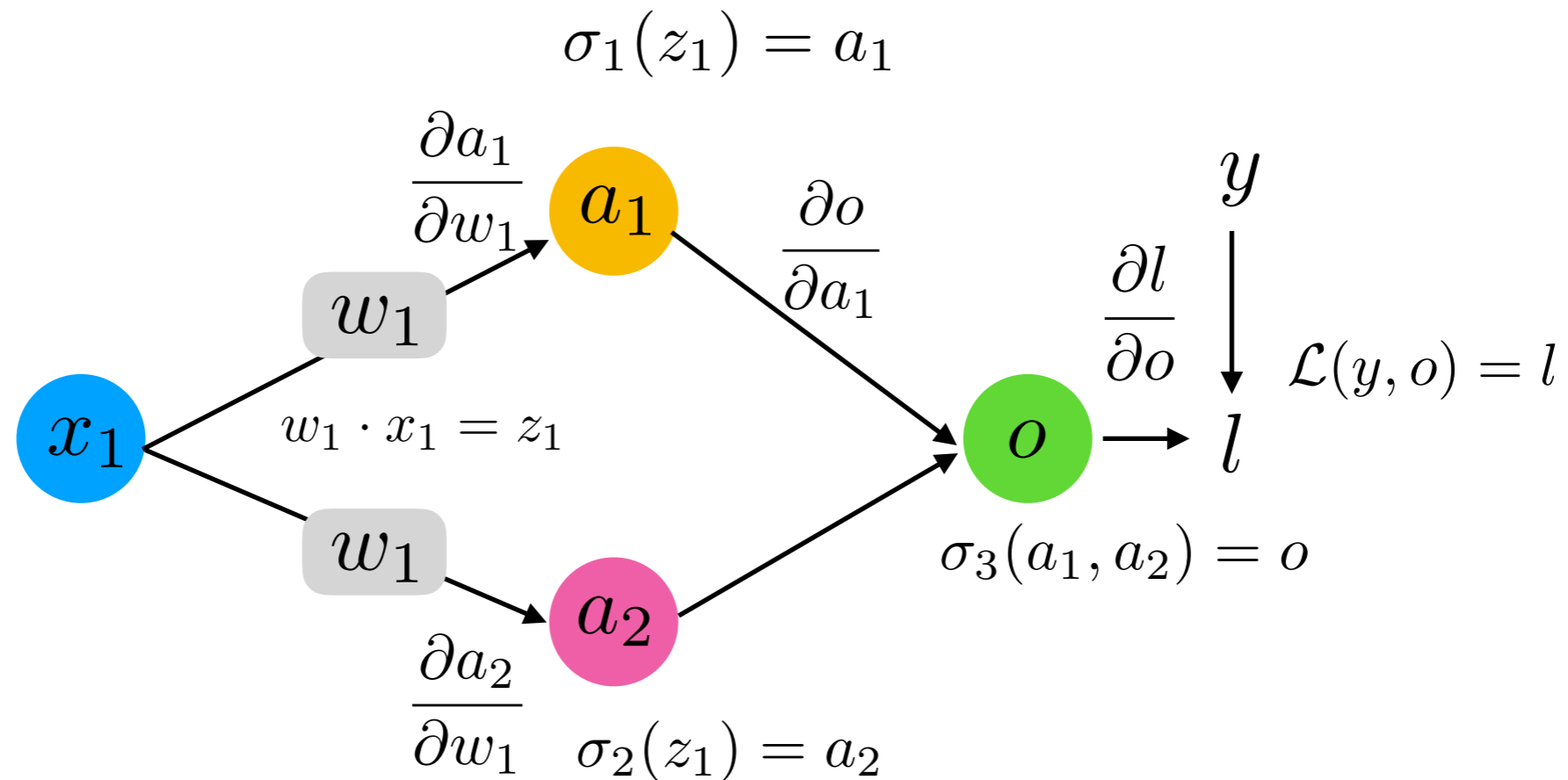
```
torch.Size([8])
```

What do you think is the output size for this 28x28 image?

Backpropagation in CNNs

Same overall concept as before: Multivariable chain rule, but now with an additional weight sharing constraint

Remember Lecture 6? Graph with Weight Sharing



Upper path

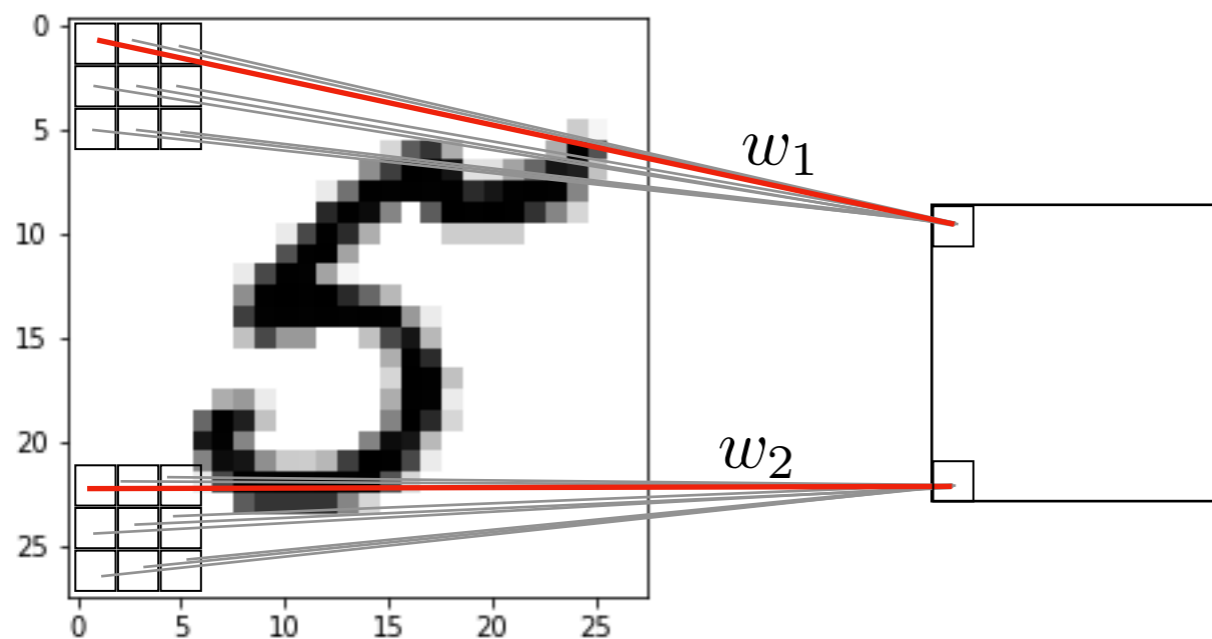
$$\frac{\partial l}{\partial w_1} = \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_1} \cdot \frac{\partial a_1}{\partial w_1} + \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_2} \cdot \frac{\partial a_2}{\partial w_1} \quad (\text{multivariable chain rule})$$

Lower path

Backpropagation in CNNs

Same overall concept as before: Multivariable chain rule, but now with an additional weight sharing constraint

Due to weight sharing: $w_1 = w_2$



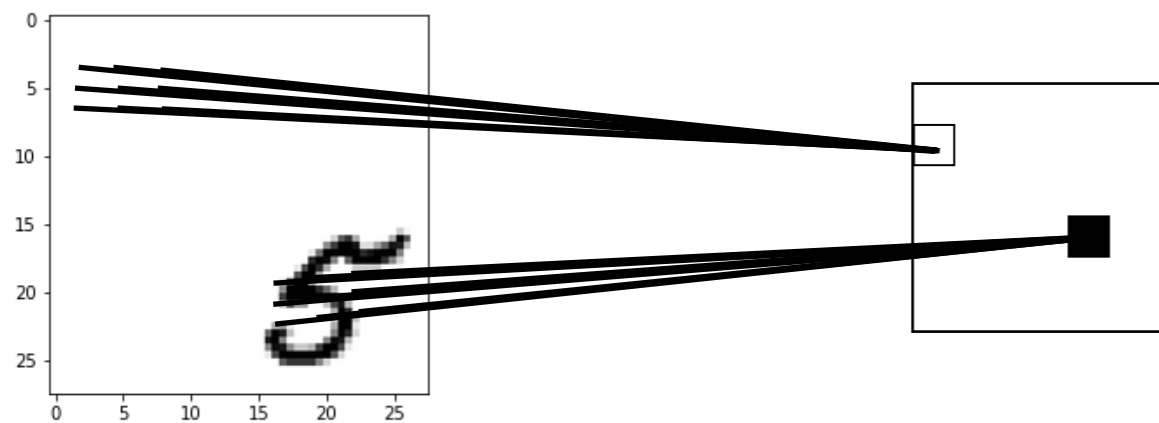
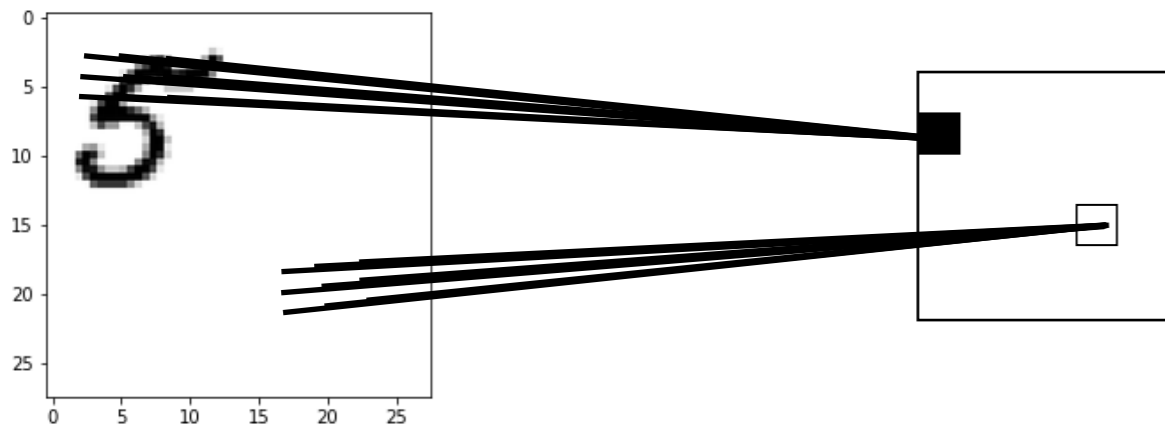
weight update:

$$w_1 := w_2 := w_1 - \eta \cdot \frac{1}{2} \left(\frac{\partial \mathcal{L}}{\partial w_1} + \frac{\partial \mathcal{L}}{\partial w_2} \right)$$

Optional averaging

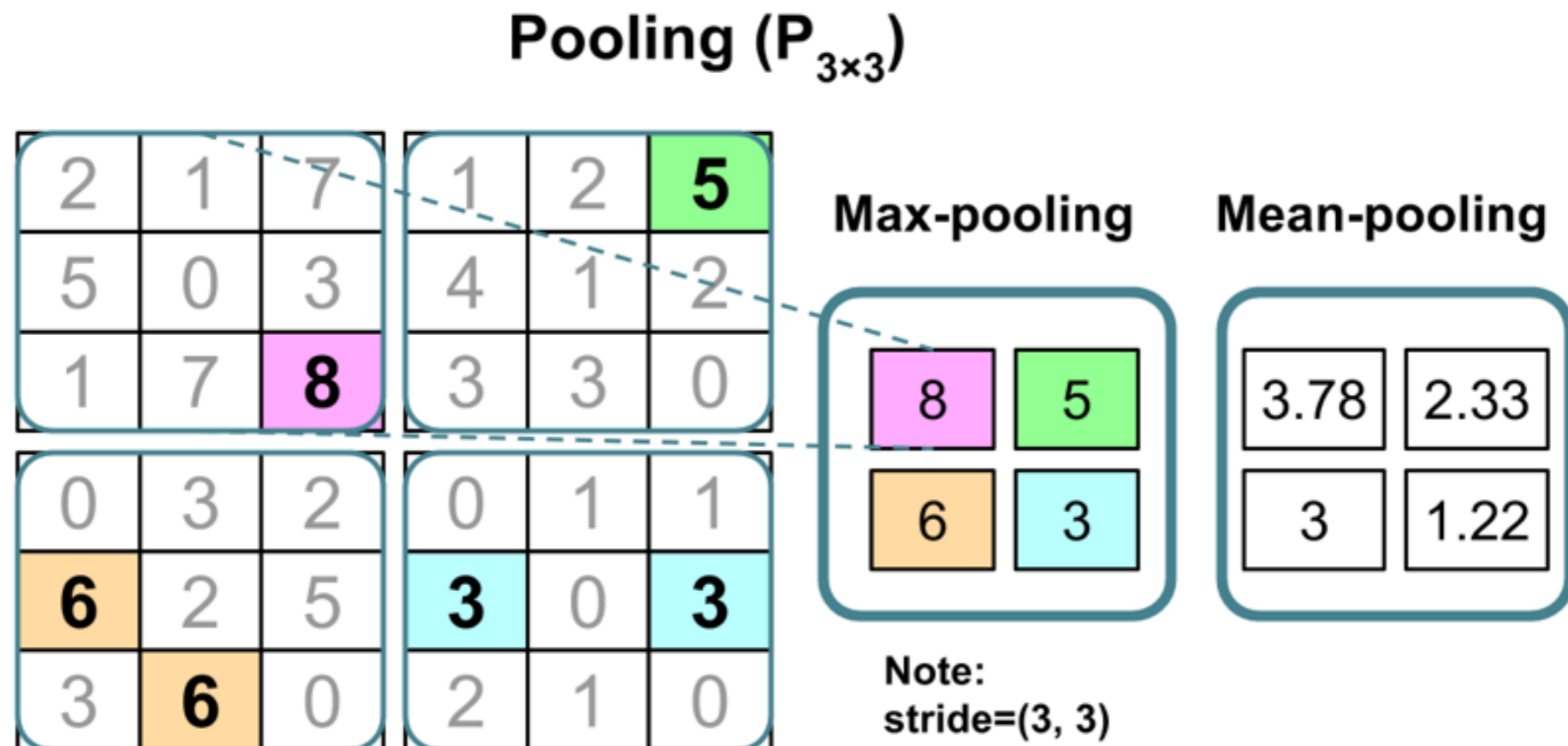
CNNs and Translation/Rotation/Scale Invariance

Note that CNNs are not really invariant to scale, rotation, translation, etc.



The activations are still dependent on the location, etc.

Pooling Layers Can Help With Local Invariance

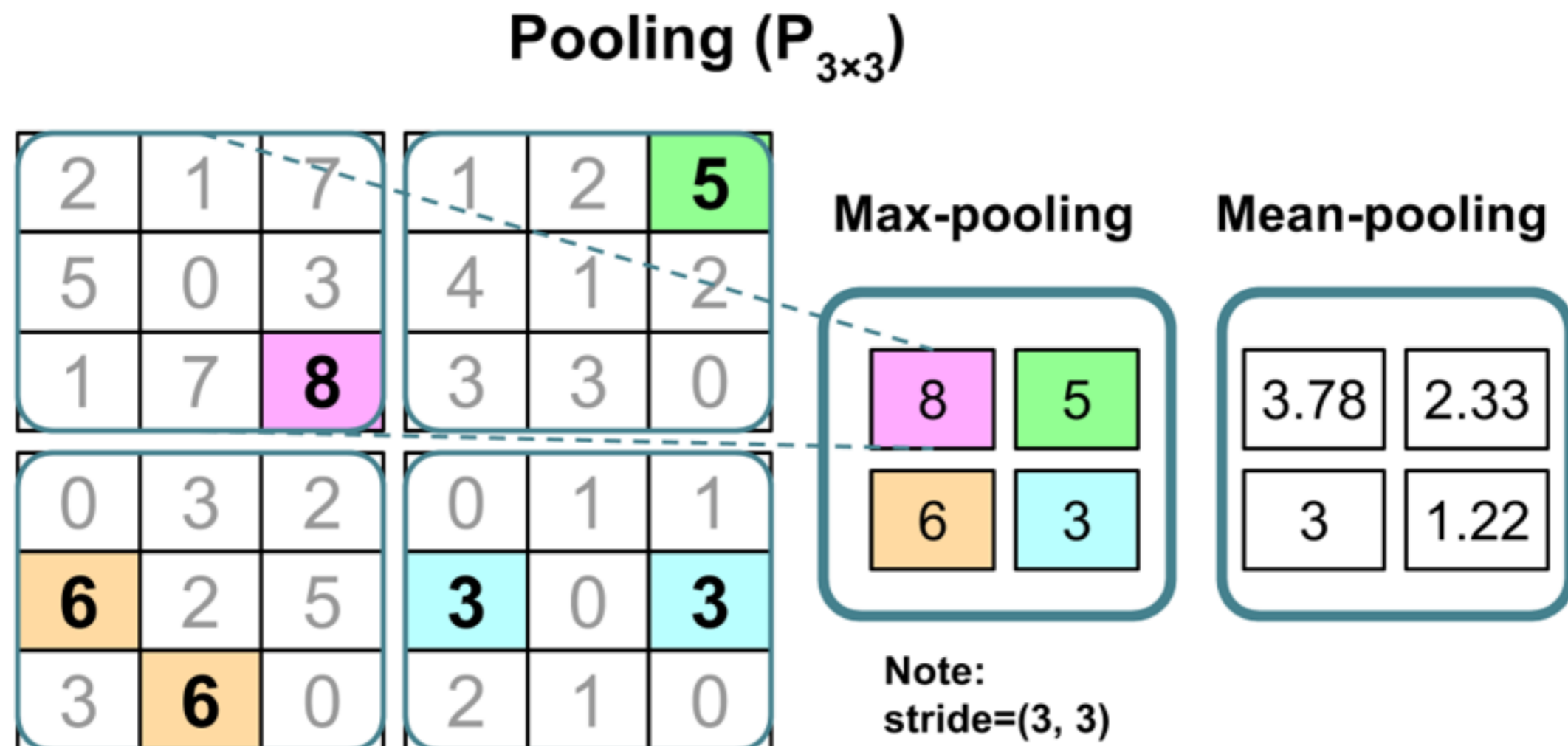


Downside: Information is lost.

May not matter for classification, but applications where relative position is important (like face recognition)

In practice for CNNs: some image preprocessing still recommended

Pooling Layers Can Help With Local Invariance



Note that typical pooling layers do not have any learnable parameters

Downside: Information is lost.

May not matter for classification, but applications where relative position is important (like face recognition)

In practice for CNNs: some image preprocessing still recommended

Main Breakthrough for CNNs: AlexNet & ImageNet

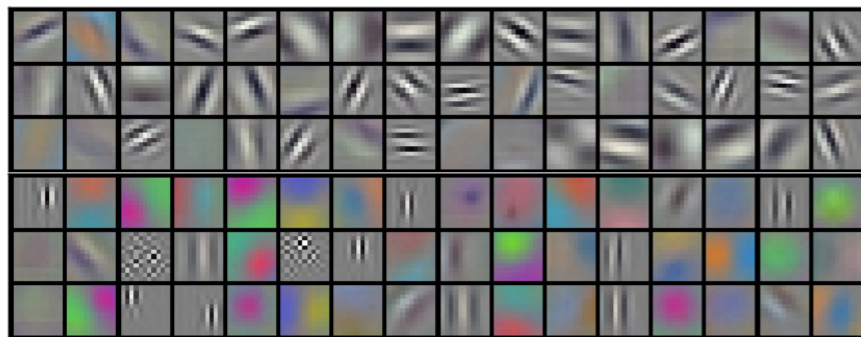


Figure 3: 96 convolutional kernels of size $11 \times 11 \times 3$ learned by the first convolutional layer on the $224 \times 224 \times 3$ input images. The top 48 kernels were learned on GPU 1 while the bottom 48 kernels were learned on GPU 2. See Section 6.1 for details.

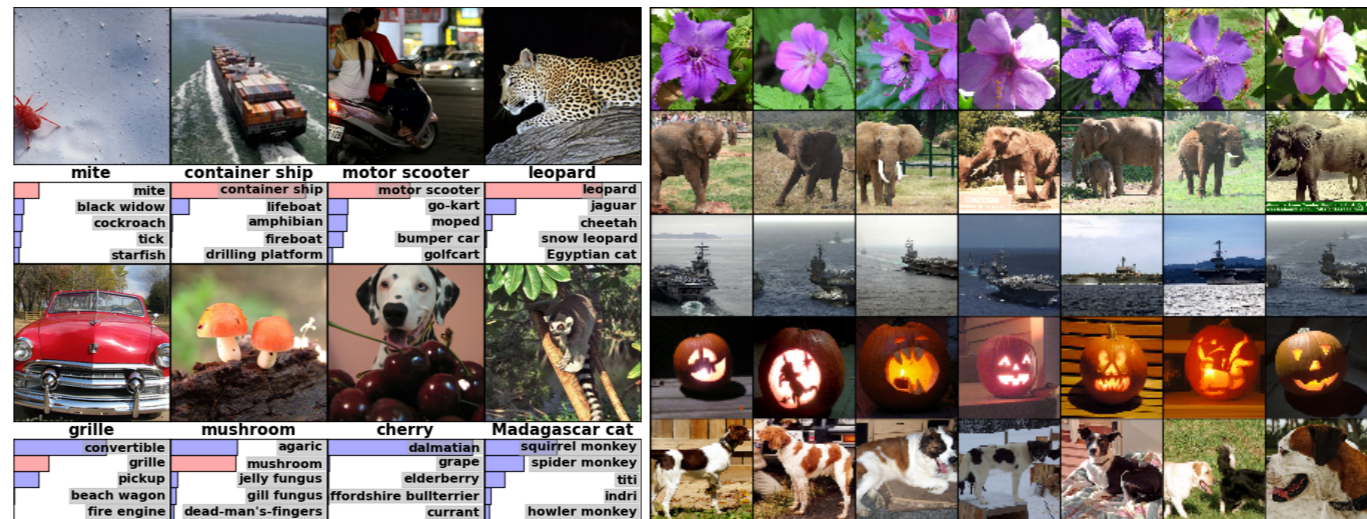
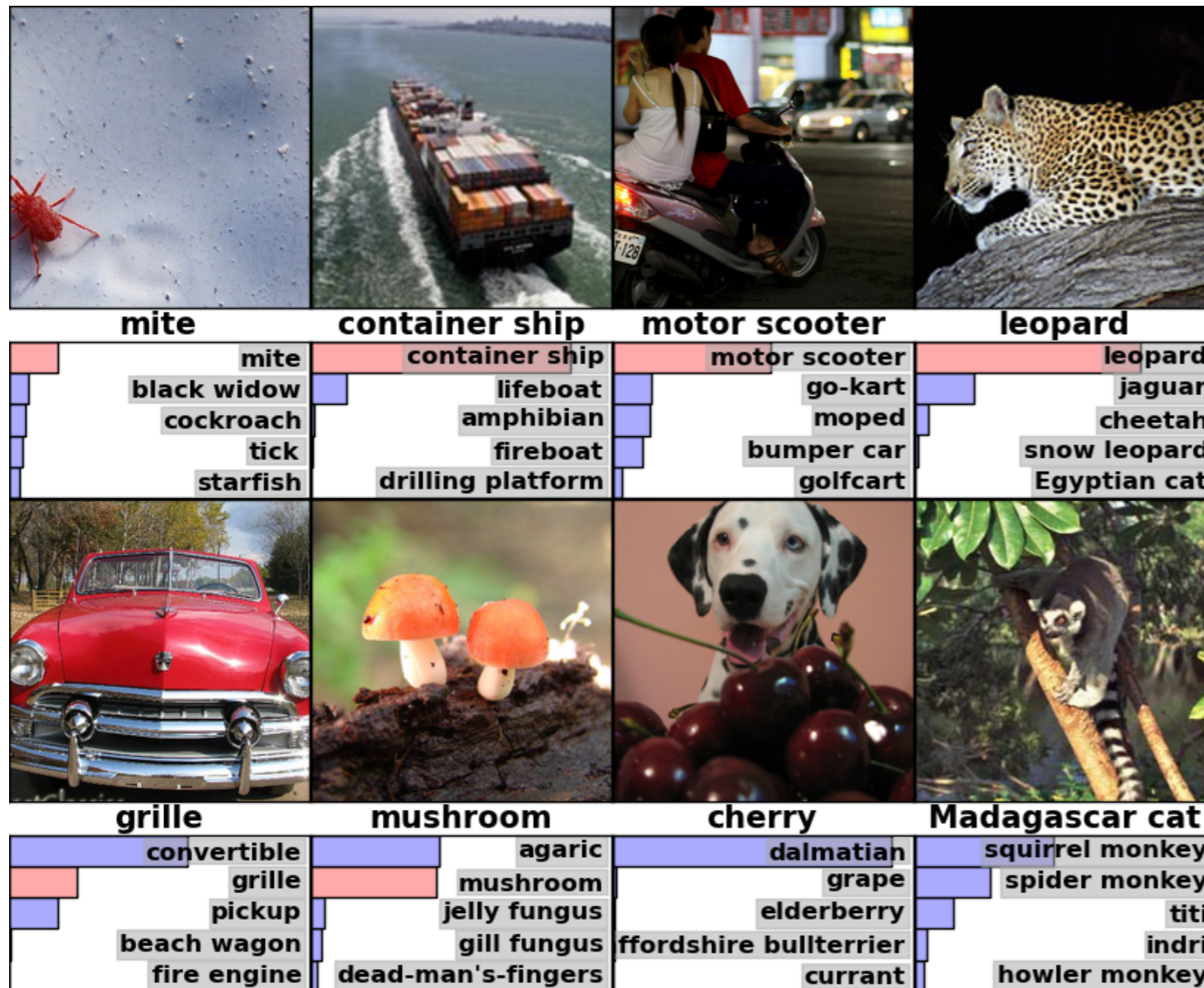


Figure 4: **(Left)** Eight ILSVRC-2010 test images and the five labels considered most probable by our model. The correct label is written under each image, and the probability assigned to the correct label is also shown with a red bar (if it happens to be in the top 5). **(Right)** Five ILSVRC-2010 test images in the first column. The remaining columns show the six training images that produce feature vectors in the last hidden layer with the smallest Euclidean distance from the feature vector for the test image.

Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems* (pp. 1097-1105).

Main Breakthrough for CNNs: AlexNet & ImageNet

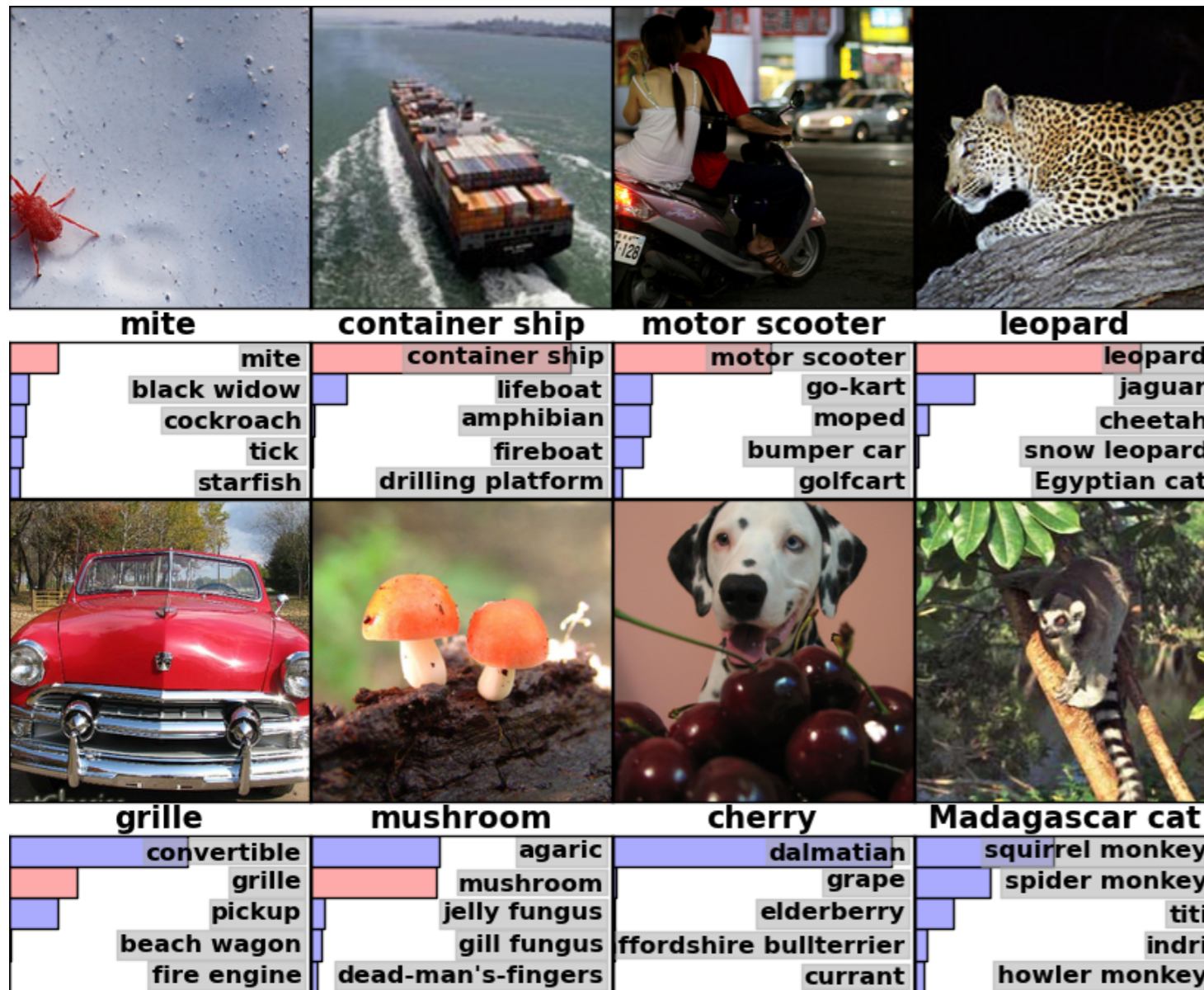


The ImageNet set that was used has ~1.2 million images and 1000 classes

Accuracy is measured as top-5 performance:
Correct prediction if the true label matches one of the top 5 predictions of the model

Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems* (pp. 1097-1105).

Main Breakthrough for CNNs: AlexNet & ImageNet



Note that the actual network inputs were still 224x224 images (random crops from downsampled 256x256 images)

224x224 is still a good/ reasonable size today (224*224*3 = 150,528 features)

Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems* (pp. 1097-1105).

Common CNN Architectures

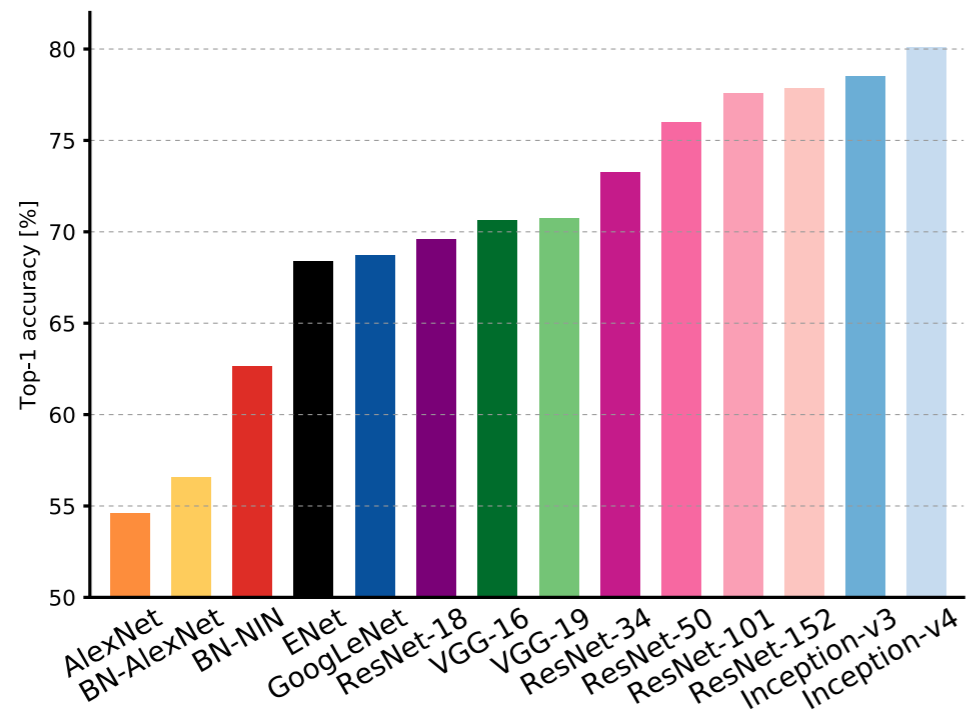


Figure 1: **Top1 vs. network.** Single-crop top-1 validation accuracies for top scoring single-model architectures. We introduce with this chart our choice of colour scheme, which will be used throughout this publication to distinguish effectively different architectures and their correspondent authors. Notice that networks of the same group share the same hue, for example ResNet are all variations of pink.

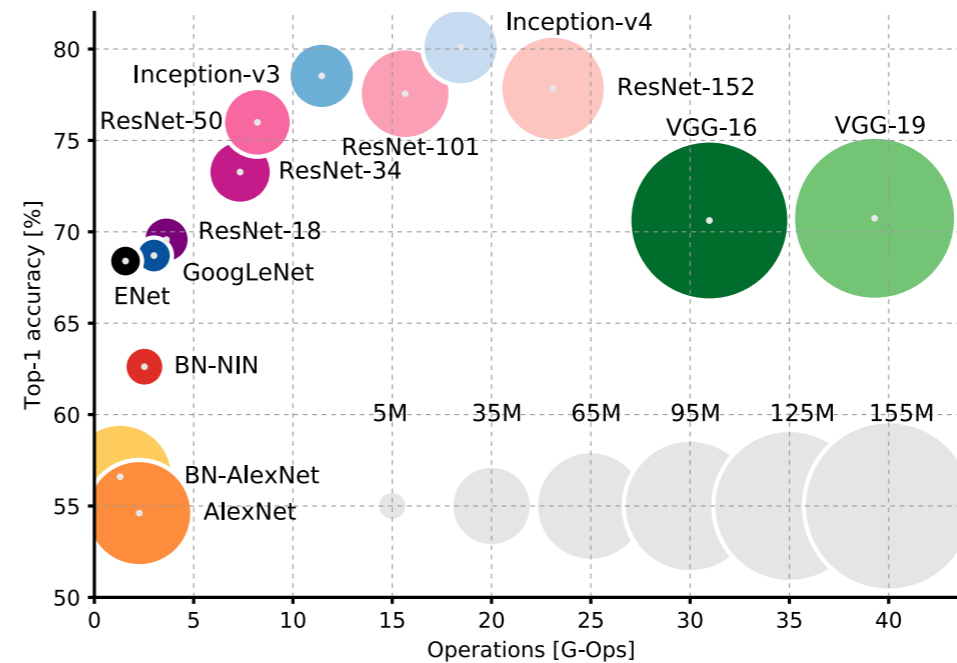


Figure 2: **Top1 vs. operations, size \propto parameters.** Top-1 one-crop accuracy versus amount of operations required for a single forward pass. The size of the blobs is proportional to the number of network parameters; a legend is reported in the bottom right corner, spanning from 5×10^6 to 155×10^6 params. Both these figures share the same y-axis, and the grey dots highlight the centre of the blobs.

Canziani, A., Paszke, A., & Culurciello, E. (2016). An analysis of deep neural network models for practical applications. *arXiv preprint arXiv:1605.07678*.

Convolutions with Color Channels

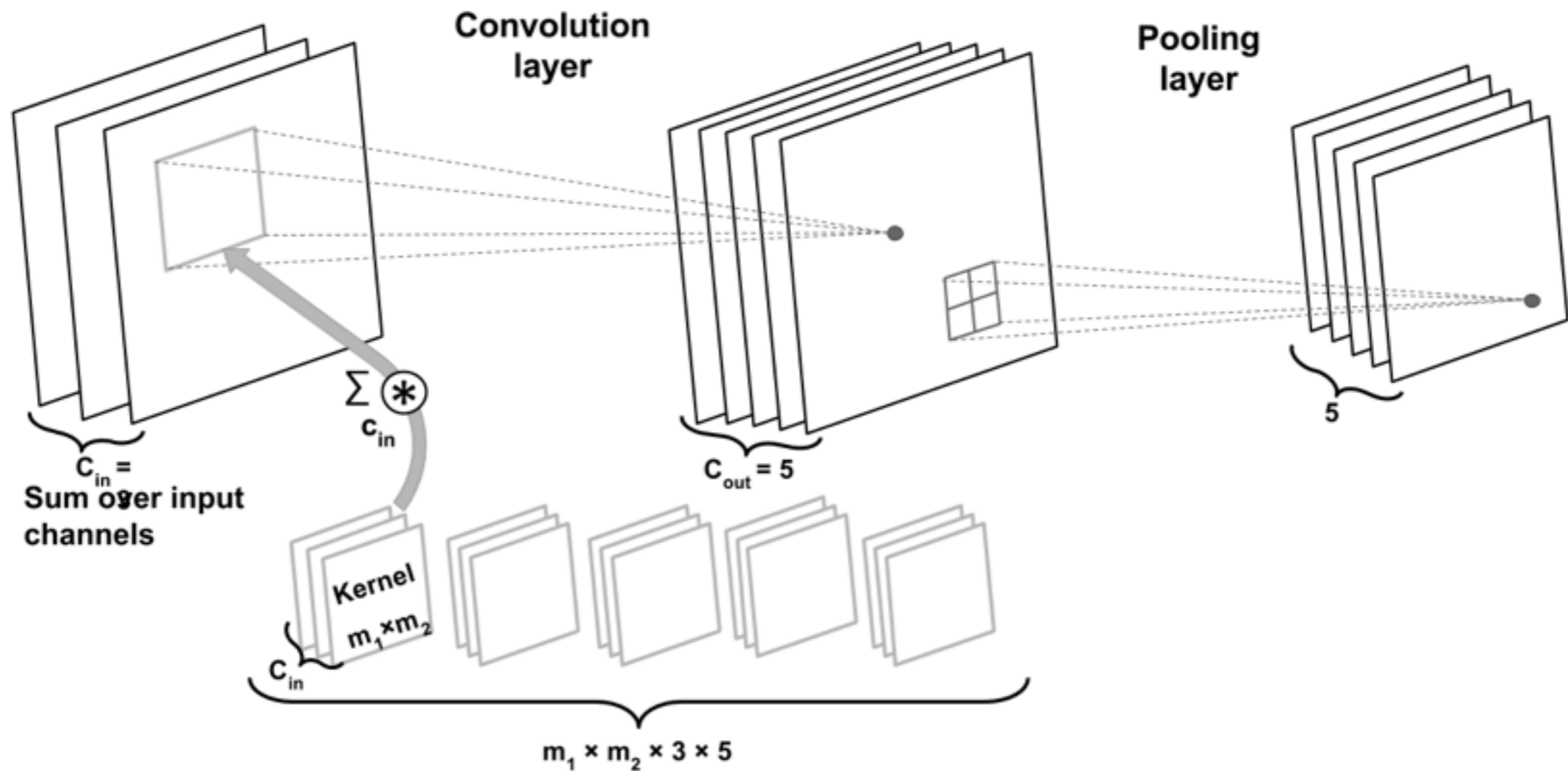


Image dimension: $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2 \times C_{in}}$ in NWHC format,
 CUDA & PyTorch use NCWH

Convolutions with Color Channels

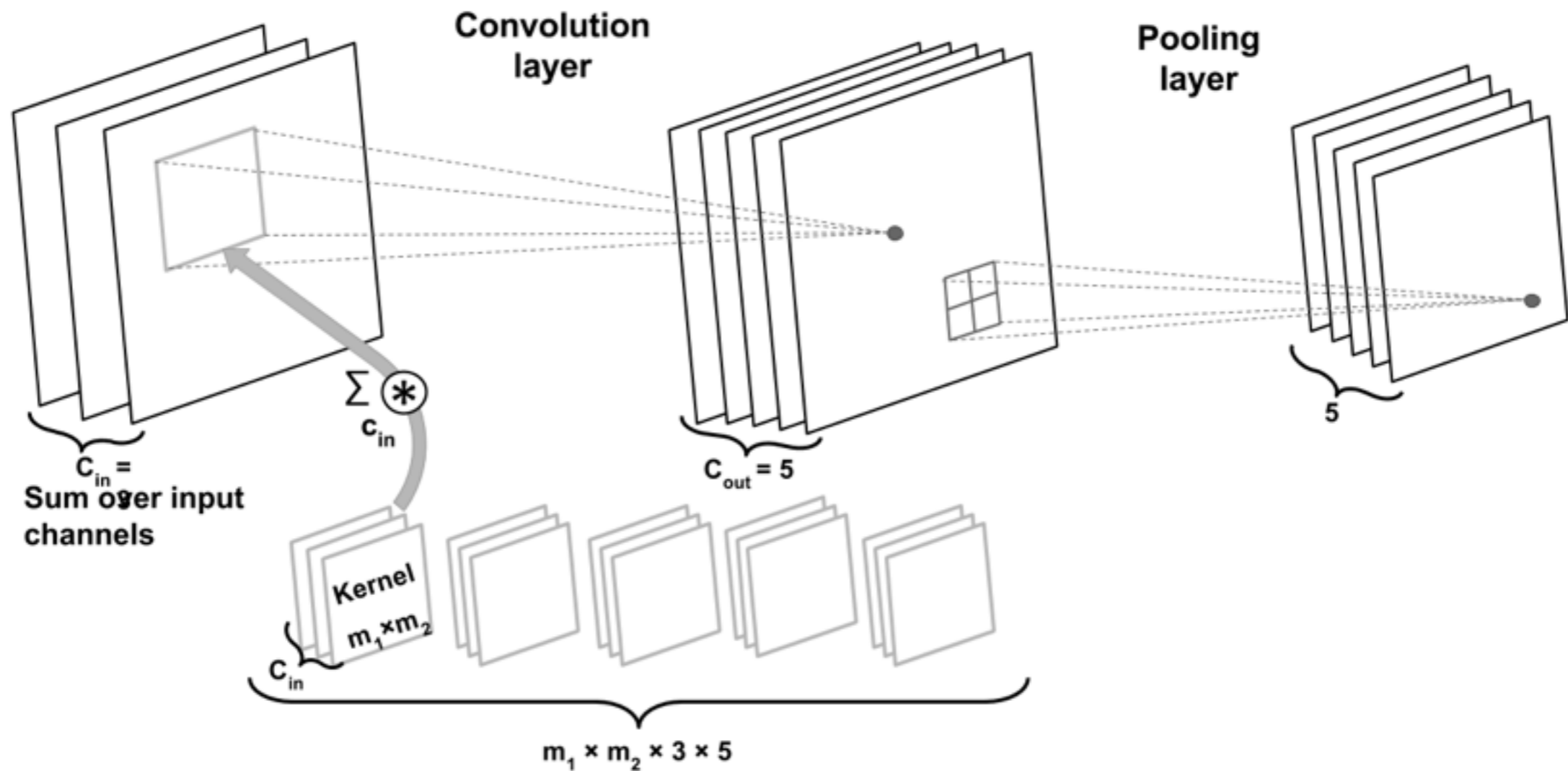
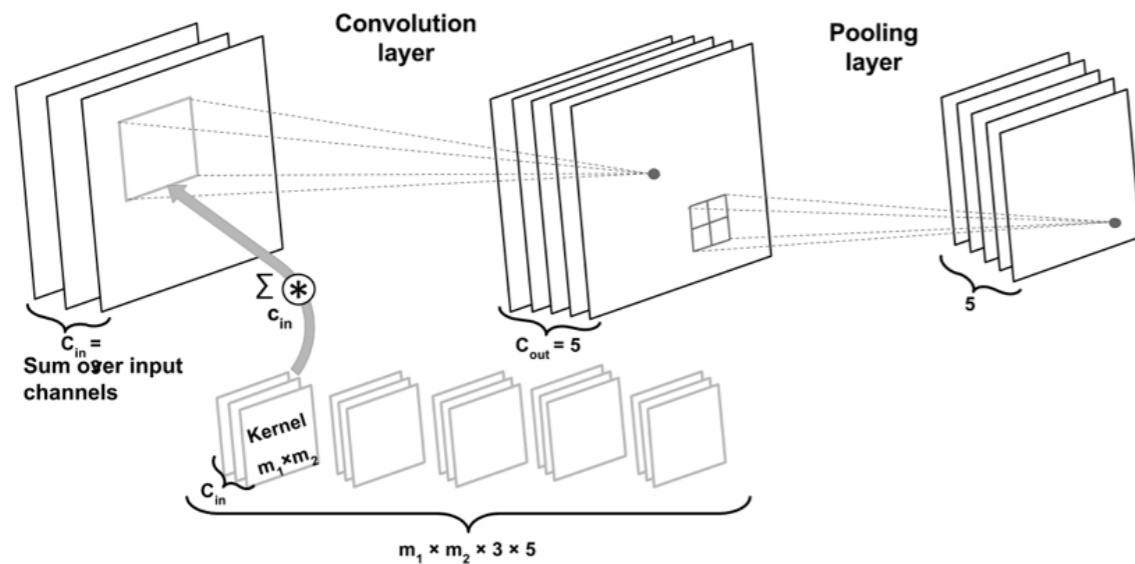


Image dimensions: $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2 \times C_{in}}$

Kernel dimensions: $\mathbf{W} \in \mathbb{R}^{m_1 \times m_2 \times C_{in} \times C_{out}}$ $\mathbf{b} \in \mathbb{R}^{C_{out}}$

Convolutions with Color Channels



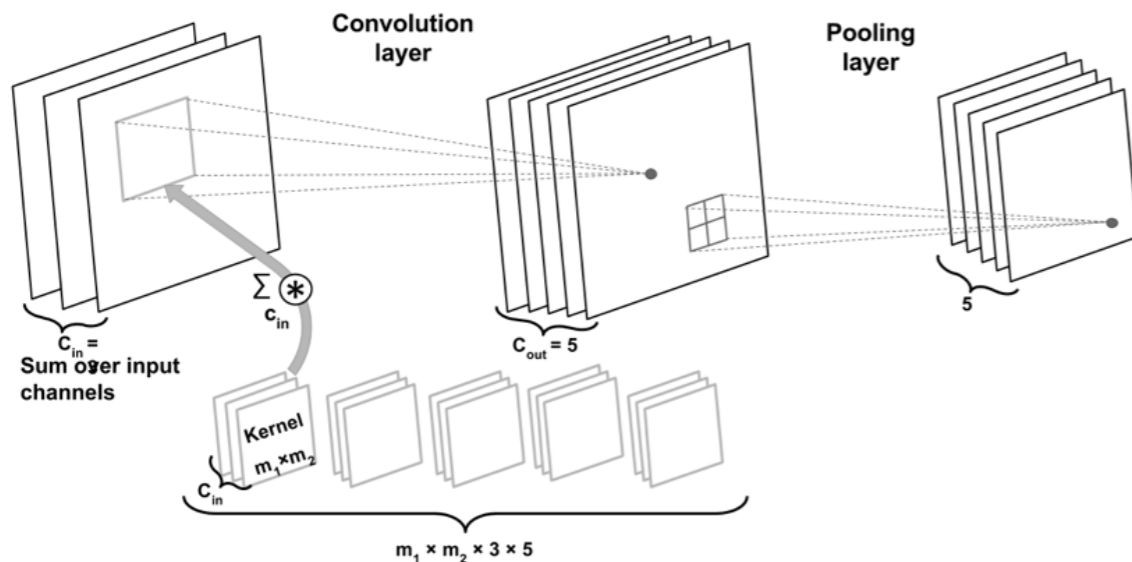
Number of parameters:

$$m_1 \times m_2 \times 3 \times 5 + 5$$

Assume 5x5 kernel:

$$5 \times 5 \times 3 \times 5 + 5 = 380$$

Convolutions with Color Channels



Assume "same" padding (more on padding later), such that the hidden layer has the same height and width as the original images

If we use a CNN:

Number of parameters:

$$m_1 \times m_2 \times 3 \times 5 + 5$$

Assume 128x128 images and

5x5 kernel:

$$5 \times 5 \times 3 \times 5 + 5 = 380$$

If we use a fully connected layer:

Number of parameters:

$$(n_1 \times n_2 \times 3) \times (n_1 \times n_2 \times 5) + (n_1 \times n_2 \times 5)$$

$$(n_1 \times n_2)^2 \times 3 \times 5 + n_1 \times n_2 \times 5$$

$$(128 \times 128)^2 \times 3 \times 5 + 128 \times 128 \times 5$$

$$= 4,026,613,760$$

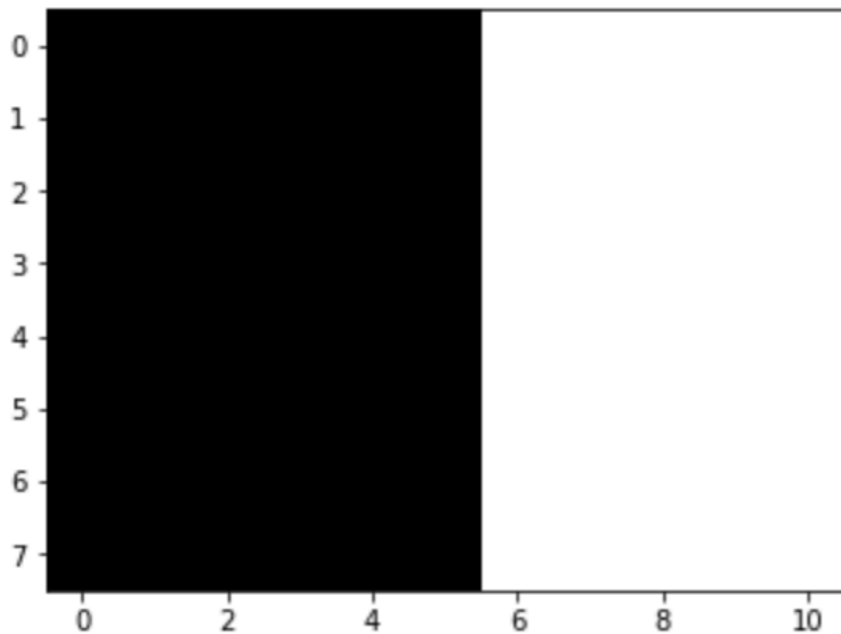
What a CNN Can See

Simple example: vertical edge detector

```
conv.weight[0, 0, :, :] = torch.tensor([[1, 0, -1],  
                                         [1, 0, -1],  
                                         [1, 0, -1]]).float()
```

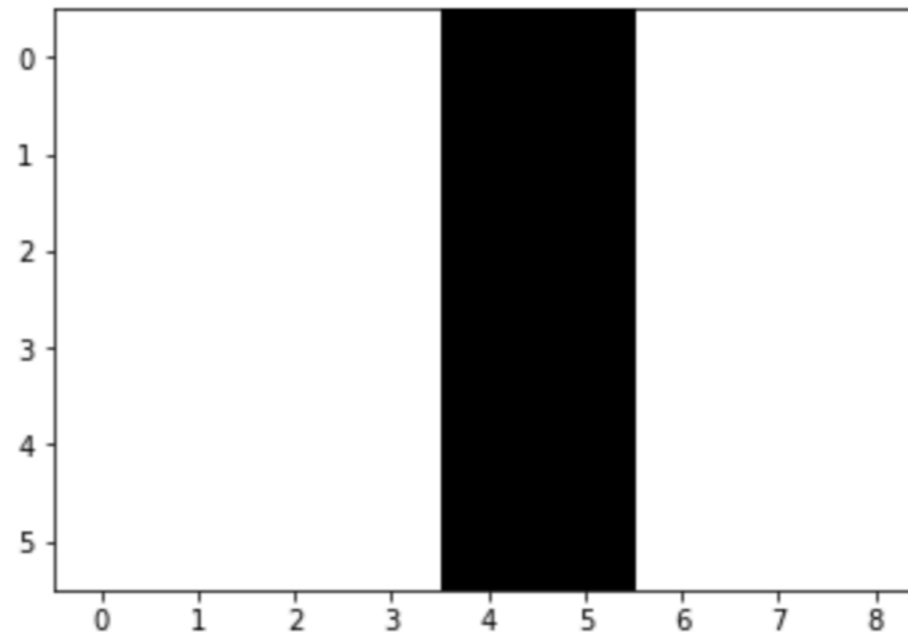
```
t = torch.tensor([  
  [0., 0., 0., 0., 0., 0., 1., 1., 1., 1., 1.],  
  [0., 0., 0., 0., 0., 0., 1., 1., 1., 1., 1.],  
  [0., 0., 0., 0., 0., 0., 1., 1., 1., 1., 1.],  
  [0., 0., 0., 0., 0., 0., 1., 1., 1., 1., 1.],  
  [0., 0., 0., 0., 0., 0., 1., 1., 1., 1., 1.],  
  [0., 0., 0., 0., 0., 0., 1., 1., 1., 1., 1.],  
  [0., 0., 0., 0., 0., 0., 1., 1., 1., 1., 1.],  
  [0., 0., 0., 0., 0., 0., 1., 1., 1., 1., 1.],  
])
```

```
plt.imshow(t, cmap='gray');
```



(From classical computer vision research)

```
tt = torch.zeros([1, 1] + list(t.size()))  
tt[0, 0, :, :] = t  
after = conv(tt)  
plt.imshow(after[0, 0, :, :].detach().numpy(), cmap='gray');
```



What a CNN Can See

Simple example: vertical edge detector

```
conv = torch.nn.Conv2d(in_channels=1,  
                       out_channels=1,  
                       kernel_size=(3, 3))
```

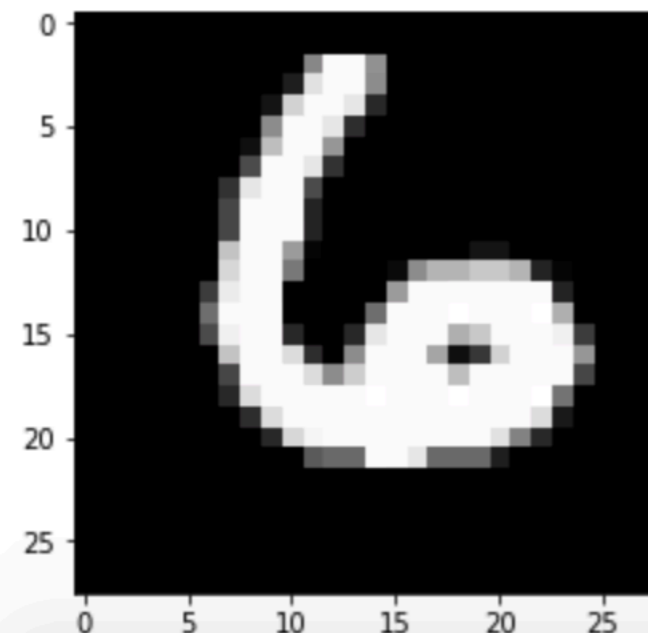
```
conv.weight.size()
```

```
torch.Size([1, 1, 3, 3])
```

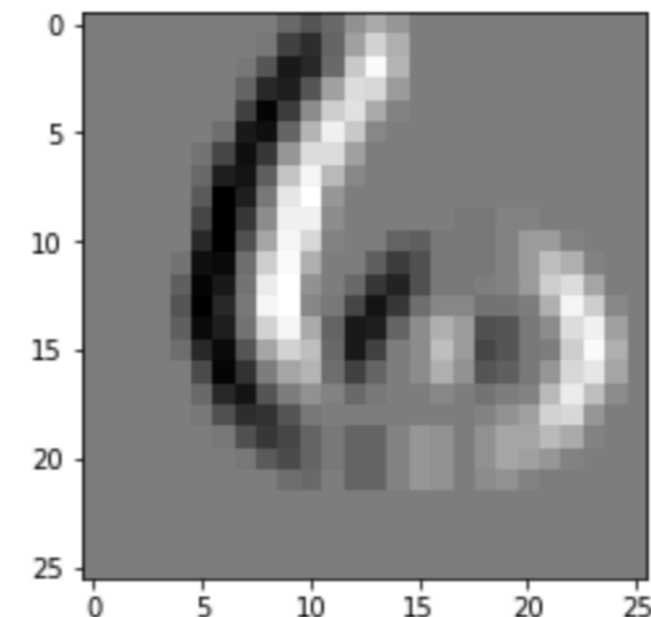
```
conv.weight[0, 0, :, :] = torch.tensor([[1, 0, -1],  
                                         [1, 0, -1],  
                                         [1, 0, -1]]).float()  
conv.bias[0] = torch.tensor([0.]).float()
```

```
images_after = conv(images)
```

```
plt.imshow(images[5, 0], cmap='gray');
```



```
plt.imshow(images_after[5, 0].detach().numpy(), cmap='gray');
```



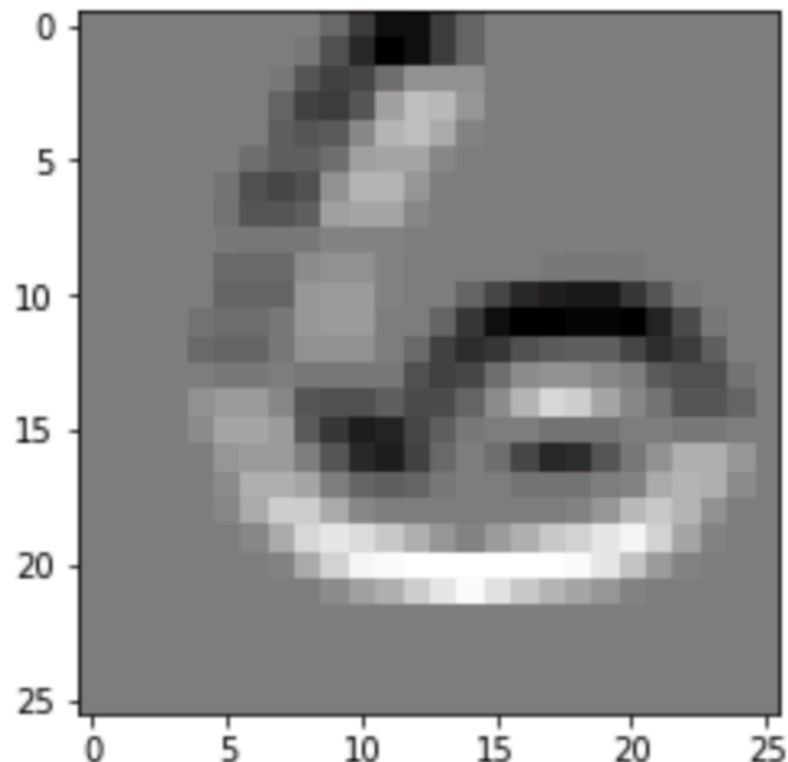
What a CNN Can See

Simple example: horizontal edge detector

```
: conv.weight[0, 0, :, :] = torch.tensor([[1, 1, 1],  
                                           [0, 0, 0],  
                                           [-1, -1, -1]]).float()  
conv.bias[0] = torch.tensor([0.]).float()
```

```
: images_after2 = conv(images)
```

```
: plt.imshow(images_after2[5, 0].detach().numpy() , cmap='gray');
```



A CNN can learn whatever it finds best based on optimizing the objective (e.g., minimizing a particular loss to achieve good classification accuracy)

What a CNN Can See

Which patterns from the training set activate the feature map?

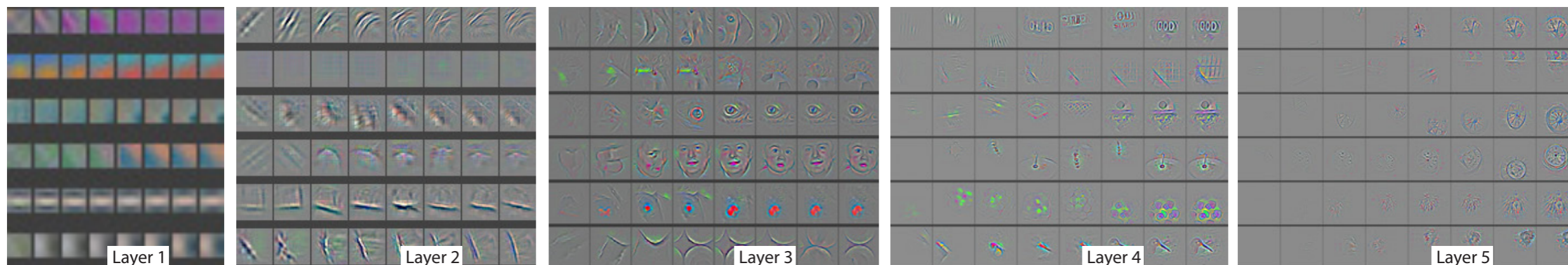


Fig. 4. Evolution of a randomly chosen subset of model features through training. Each layer's features are displayed in a different block. Within each block, we show a randomly chosen subset of features at epochs [1,2,5,10,20,30,40,64]. The visualization shows the strongest activation (across all training examples) for a given feature map, projected down to pixel space using our deconvnet approach. Color contrast is artificially enhanced and the figure is best viewed in electronic form.

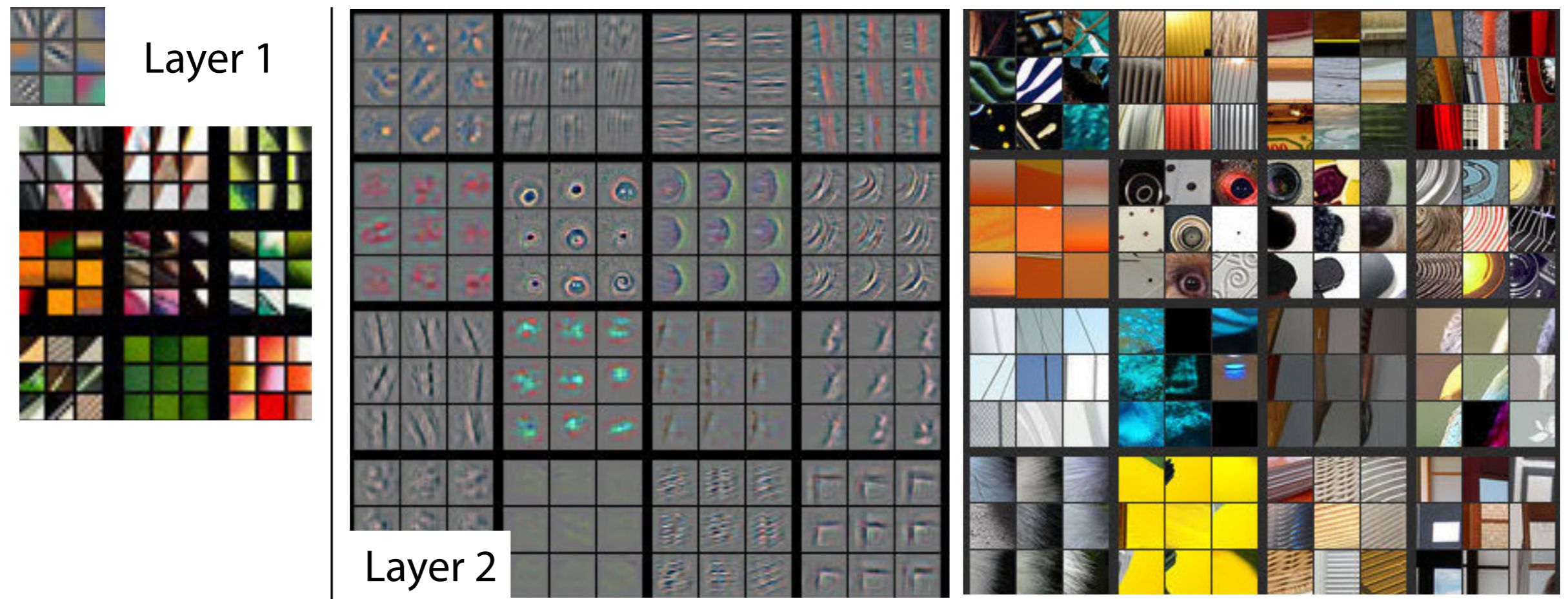
[Zeiler, M. D., & Fergus, R. \(2014, September\). Visualizing and understanding convolutional networks. In *European conference on computer vision* \(pp. 818-833\). Springer, Cham.](#)

Method: backpropagate strong activation signals in hidden layers to the input images, then apply "unpooling" to map the values to the original pixel space for visualization

What a CNN Can See

Which patterns from the training set activate the feature map?

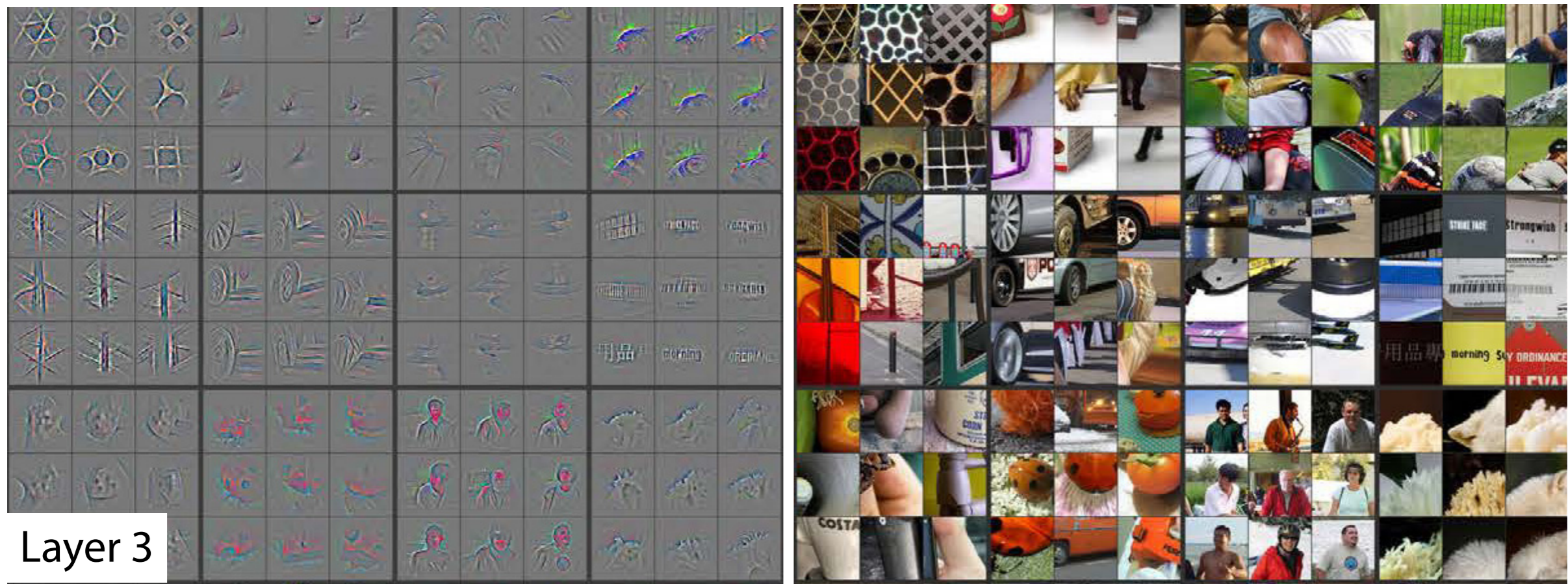
Zeiler, M. D., & Fergus, R. (2014, September). Visualizing and understanding convolutional networks. In *European conference on computer vision* (pp. 818-833). Springer, Cham.



What a CNN Can See

Which patterns from the training set activate the feature map?

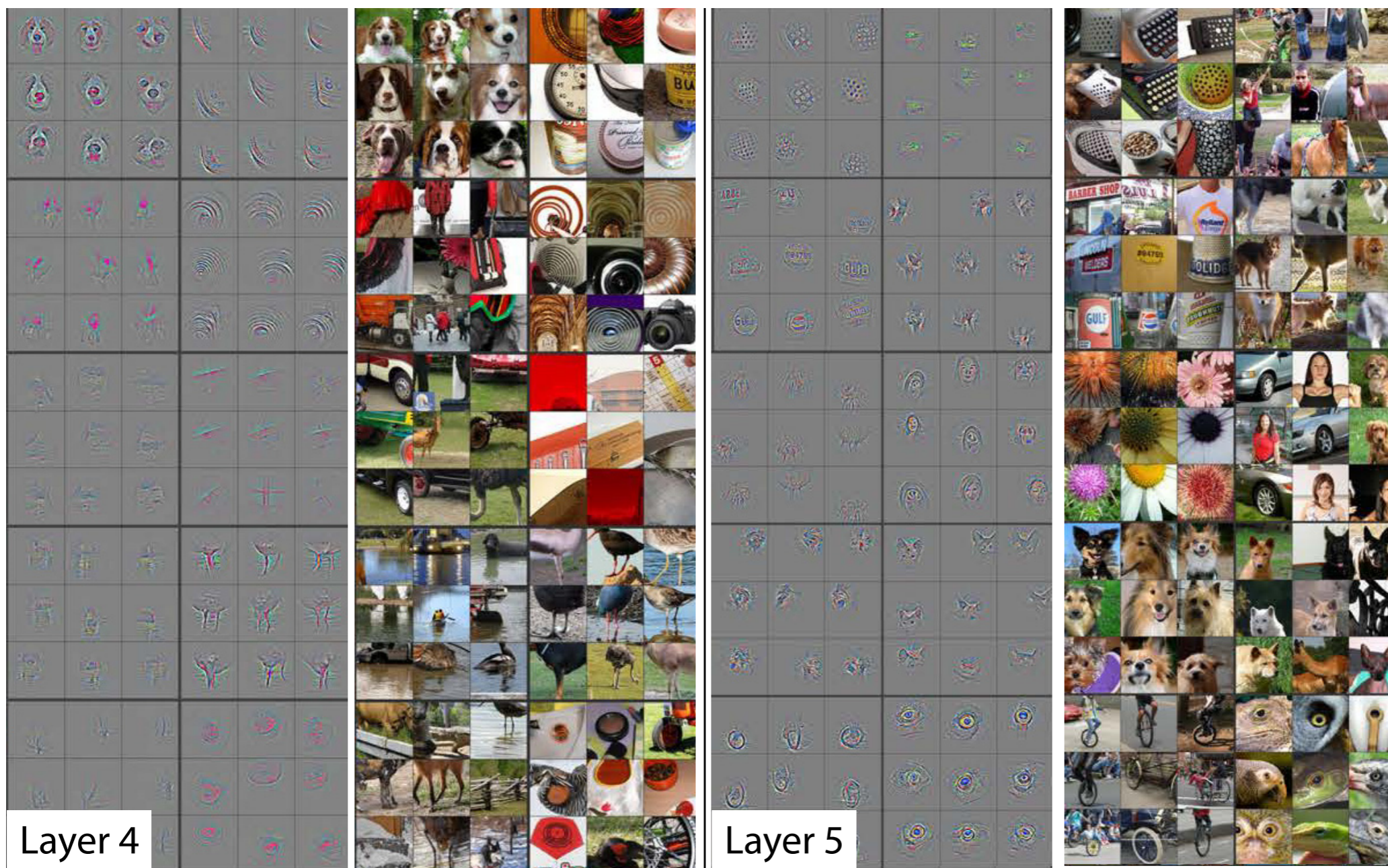
Zeiler, M. D., & Fergus, R. (2014, September). Visualizing and understanding convolutional networks. In *European conference on computer vision* (pp. 818-833). Springer, Cham.



What a CNN Can See

Which patterns from the training set activate the feature map?

Zeiler, M. D., & Fergus, R. (2014, September). Visualizing and understanding convolutional networks. In *European conference on computer vision* (pp. 818-833). Springer, Cham.



End of part 1