# Agenda

- **Convolutional & pooling layers**

- **Convolutional neural networks**

- **Feature visualization**

- **Applications**

# Images as 2D Arrays

- Grayscale image is a 2D array of pixel values

- Color images are 3D array
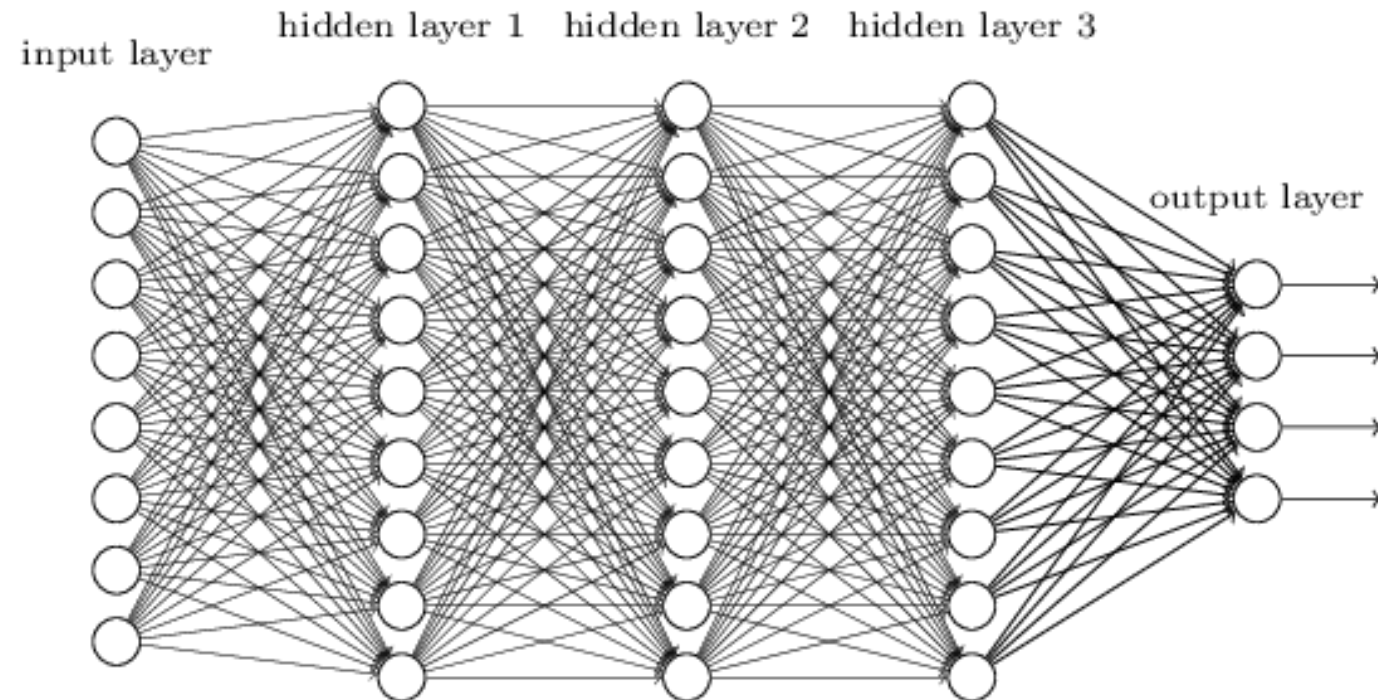  - 3rd dimension is color (e.g., RGB)
  - Called "channels"



| 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 |

# Convolution Neural Network: CNN

We know it is good to learn a small model.
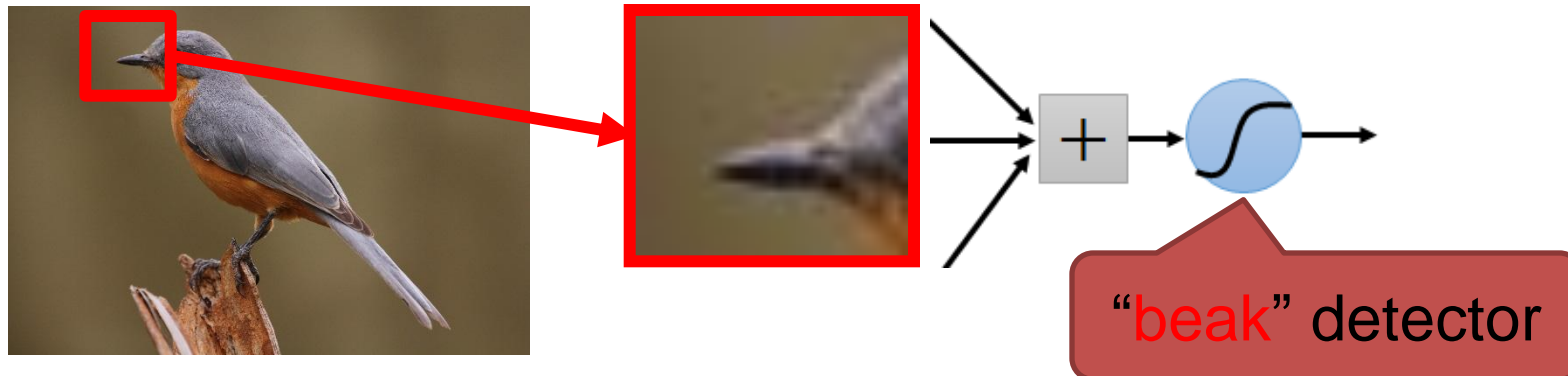From this fully connected model, do we really need all the edges?
Can some of these be shared?

# Consider learning an image:
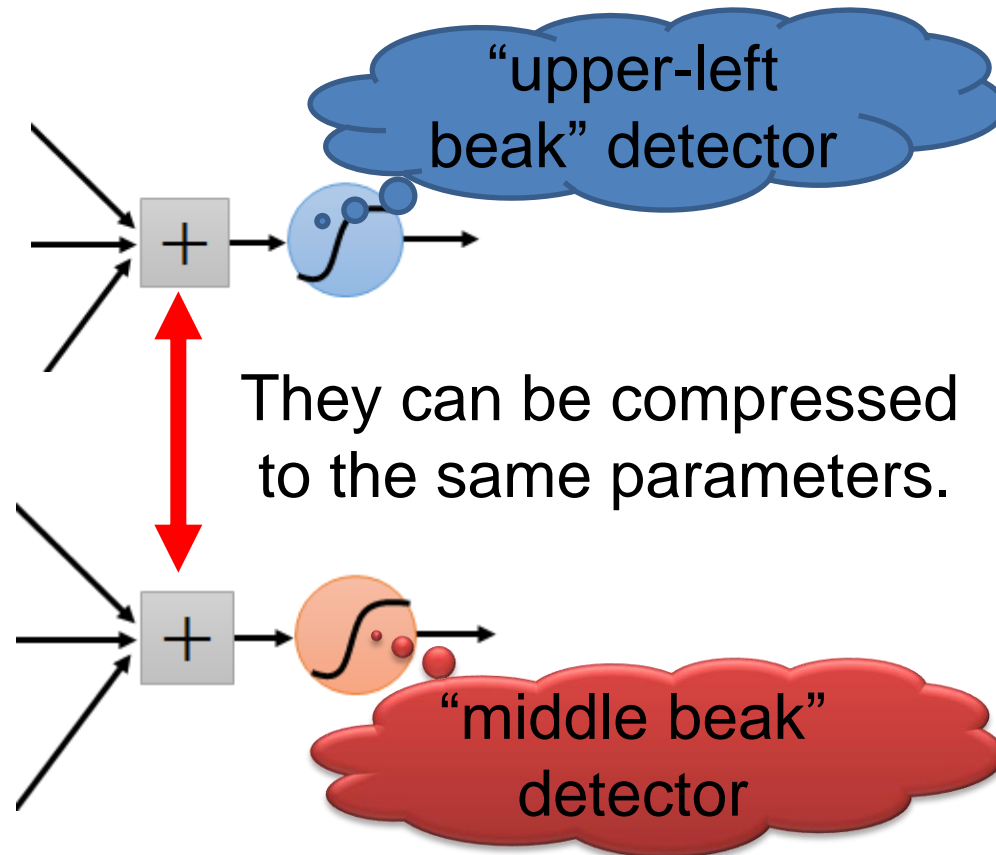
Some patterns are much smaller than the whole image

Can represent a small region with fewer parameters



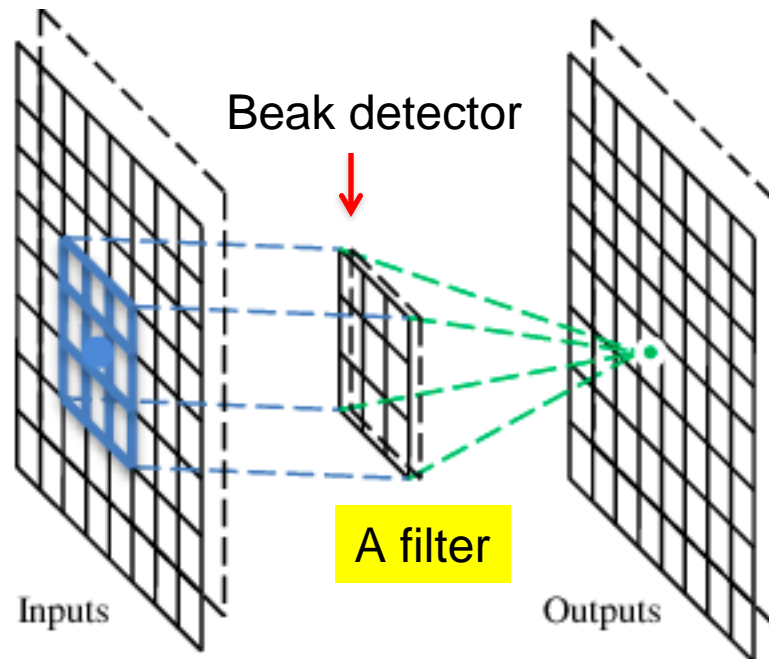"beak" detector

Same pattern appears in different places:
They can be compressed!
What about training a lot of such "small" detectors and each detector must "move around".



"upper-left beak" detector

They can be compressed to the same parameters.

"middle beak" detector

# A convolutional layer

A CNN is a neural network with some convolutional layers (and some other layers).  A convolutional layer has a number of filters that does convolutional operation.

Beak detector
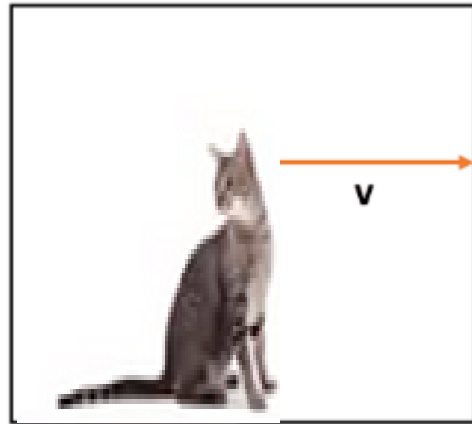
A filter

Inputs

Outputs

# What Is Translation?

A **translation** is a geometric transformation that shifts all points in a given direction and by the same distance. Alternatively, it can be interpreted as sliding the origin of the coordinate system by the same amount but in the opposite direction.

The translation can be expressed mathematically as the vector sum of a constant vector $\mathbf{v}$ to each point $\mathbf{x}$:

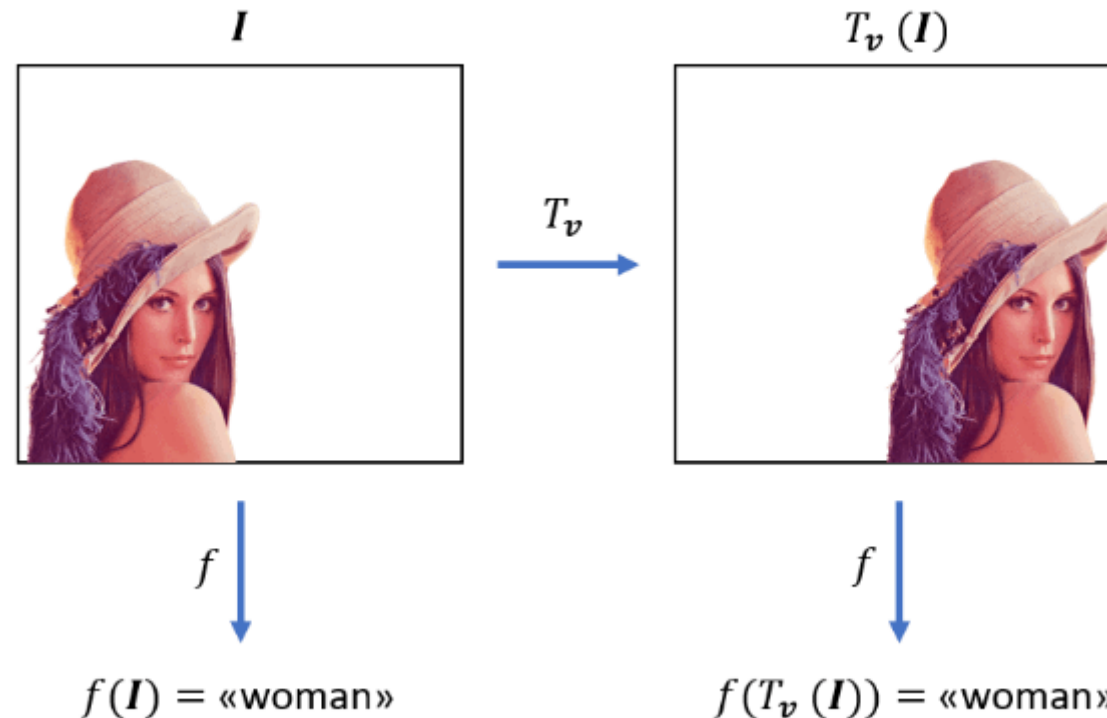$$T_{\mathbf{v}}(\mathbf{x}) = \mathbf{x} + \mathbf{v}$$
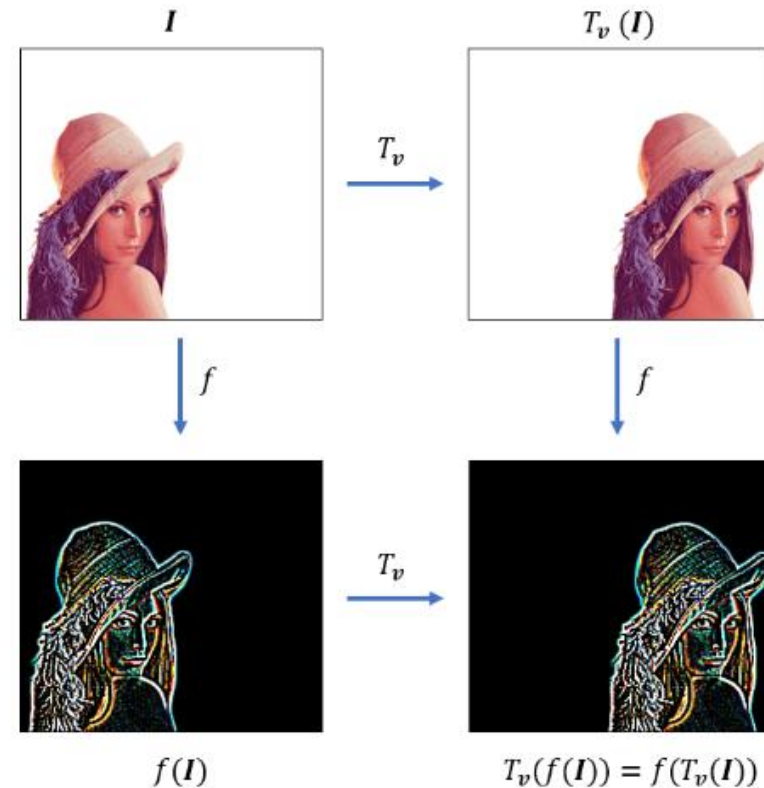


(a)

(b)

# Translation Invariance

- **Translation invariance -** is a property of a system or model where its performance or output does not change when the input data is translated or shifted.

$$I \qquad\qquad T_v(I)$$

$$T_v$$

$$f \qquad\qquad\qquad f$$

$$f(I) = \text{«woman»} \qquad\qquad f(T_v(I)) = \text{«woman»}$$

# Translation Equivariance

**Translation equivariance -** is a property of a function or an operation where the output's structure does not change even if the input is translated.
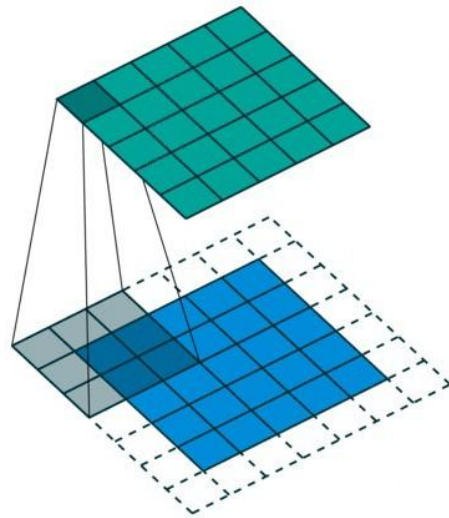


$T_v(f(I)) = f(T_v(I))$

A function $f$ is said equivariant to a function $g$ if and only if:
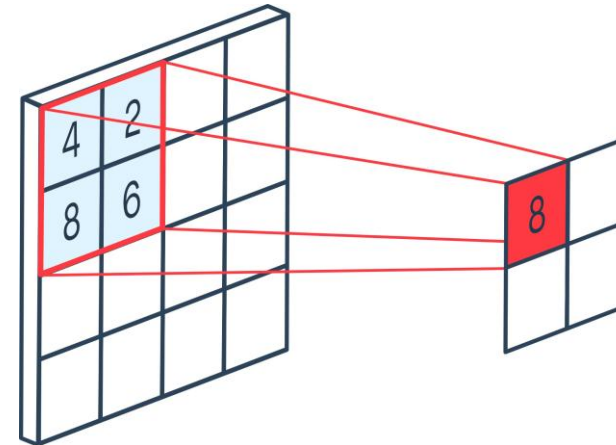
$$f(g(x)) = g(f(x))$$

In other words, $f$ is equivariant to $g$ if the order of application does not change the result of the composite function.

# Structure in Images

- Use layers that capture structure



**Convolution layers**
(Capture equivariance)

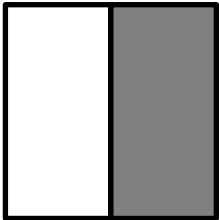**Pooling layers**
(Capture invariance)

# Vertical edge detection

| 10 | 10 | 10 | 0 | 0 | 0 |
|----|----|----|---|---|---|
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |

\*

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

=

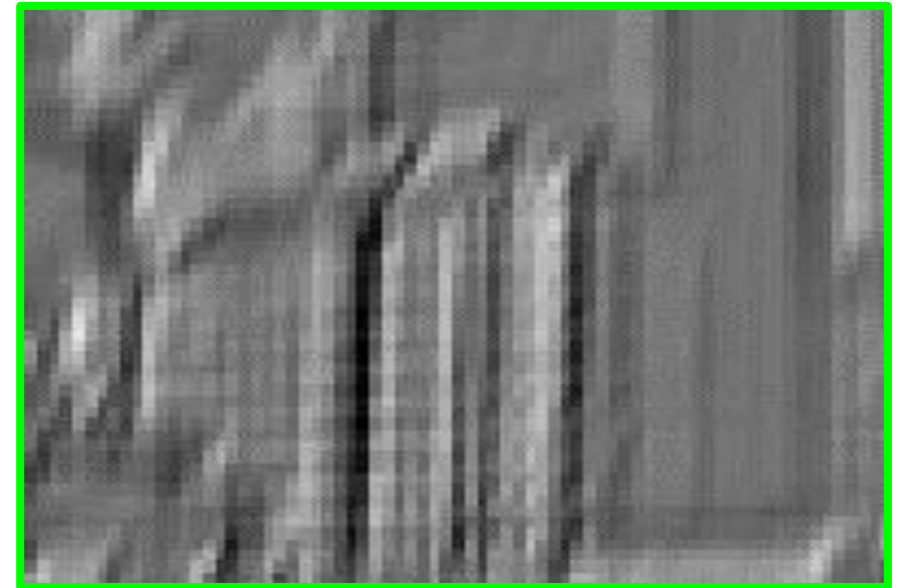| 0 | 30 | 30 | 0 |
|---|----|----|---|
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |

# Convolution Filters



$$\text{output}[i, j] = \sum_{\tau=0}^{k-1} \sum_{\gamma=0}^{k-1} \text{filter}[\tau, \gamma] \cdot \text{image}[i + \tau, j + \gamma]$$

# 2D Convolution Filters

- **Given:**
  - A 2D input $x$
  - A 2D $h \times w$ kernel $k$

- The 2D convolution is:

$$y[s,t] = \sum_{\tau=0}^{h-1} \sum_{\gamma=0}^{w-1} k[\tau, \gamma] \cdot x[s+\tau, t+\gamma]$$

# 2D Convolution Filters



Example Edge Detection Kernels



Result of Convolution with Horizontal Kernel
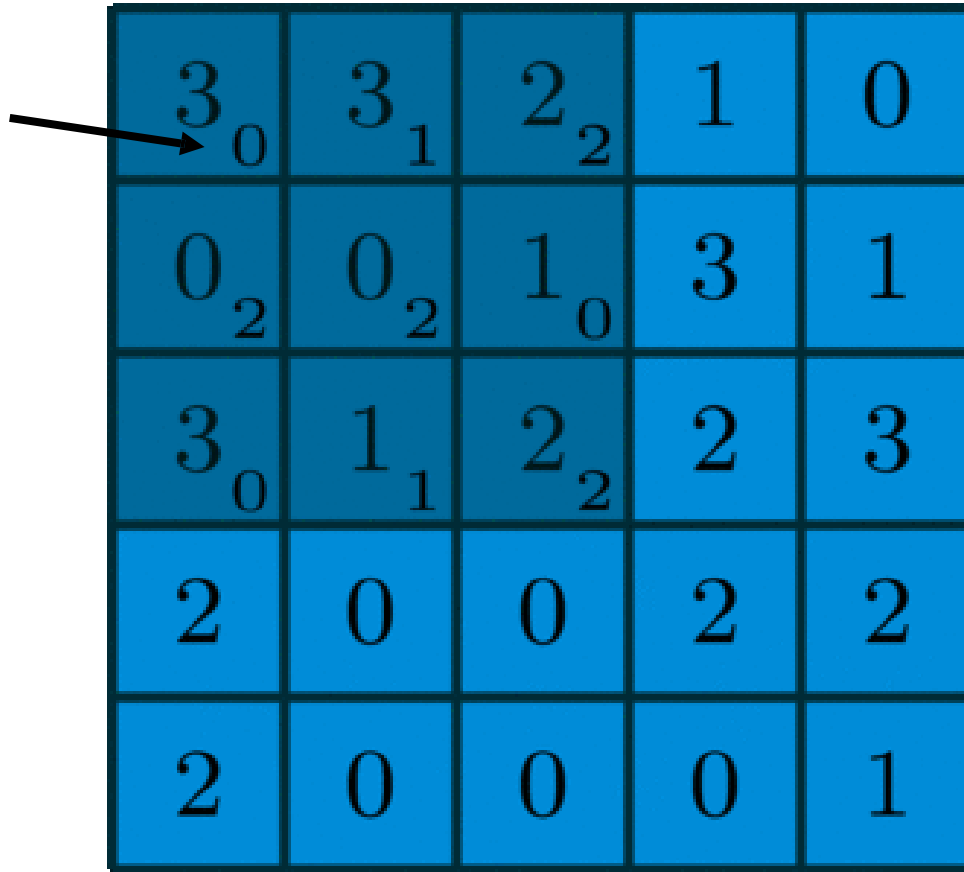
# 2D Convolution Filters

- Historically (until late 1980s), kernel parameters were handcrafted
    - E.g., "edge detectors"

- In convolutional neural networks, they are learned
    - Essentially a linear layer with fewer "connections"
    - Backpropagate as usual!
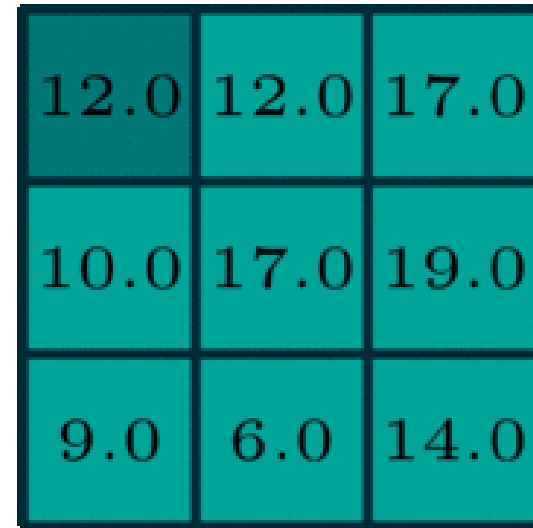
# Convolution Layers

Learnable parameters

# Convolution Layer Parameters

- **Stride:** How many pixels to skip (if any)
  - **Default:** Stride of 1 (no skipping)

Filter



$$OutDimension = \frac{InputDimension}{StrideDimension}$$

# Convolution Layer Parameters

- **Padding:** Add zeros to edges of image to capture ends
  - **Default:** No padding



stride = 1, zero-padding = 1

stride = 2, zero-padding = 1

# Convolution Layer Parameters

- **Summary:** Hyperparameters
    - Kernel size
    - Stride
    - Amount of zero-padding
    - Output channels

- Together, these determine the relationship between the input tensor shape and the output tensor shape

- Typically, also use a single bias term for each convolution filter

# Convolution Layers

Weight sharing

Local connectivity

filter size, stride

# input channels

# filters = #output (activation) maps

# Example

- Kernel size 3, stride 2, padding 1

- 3 input channels
  - Hence kernel size 3×3×3

- 2 output channels
  - Hence 2 kernels

- Total # of parameters:
  - $(3×3×3 + 1)×2 = 56$

# Pooling Layers

# Pooling Layers

$$\text{output}[0,0] = \max_{0 \le \tau < k} \max_{0 \le \gamma < k} \text{image}[0 + \tau, 0 + \gamma]$$

# Pooling Layers



$$\text{output}[0,1] = \max_{0 \le \tau < k} \max_{0 \le \gamma < k} \text{image}[0 + \tau, 1 + \gamma]$$

# Pooling Layers



$$\text{output}[0,2] = \max_{0 \le \tau < k} \max_{0 \le \gamma < k} \text{image}[0 + \tau, 2 + \gamma]$$

# Pooling Layers

$$\text{output}[i, j] = \max_{0 \le \tau < k} \max_{0 \le \gamma < k} \text{image}[i + \tau, j + \gamma]$$

# Pooling Layers

- **Summary:** Hyperparameters
    - Kernel size
    - Stride (usually >1)
    - Amount of zero-padding
    - Pooling function (almost always "max")

- Together, these determine the relationship between the input tensor shape and the output tensor shape

- **Note:** Unlike convolution, pooling operates on channels separately
    - Thus, $n$ input channels → $n$ output channels

# Summary: Convolution vs. Pooling

- **Convolution layers:** Translation equivariant
  - If object is translated, convolution output is translated by same amount
  - Produce "image-shaped" features that retain associations with input pixels

- **Pooling layers:** Translation invariant
  - Binning to make outputs insensitive to translation
  - Also reduces dimensionality

- Combined in modern architectures
  - Convolution to construct equivariant features
  - Pooling to enable invariance

# The whole CNN



cat dog ......

Fully Connected
Feedforward network

Flattened

Convolution

Max Pooling

Convolution

Max Pooling

Can repeat many times

# Example

- Suppose we want to predict whether an image depicts Cartesian axes

$$\begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & -2 & 0 & -2 \end{bmatrix}$$

input image

?

target (binary) label

# Example

- **Step 1:** Convolve the image with two filters
  - No padding, stride 1

- **Step 2:** Run max pooling

$$\begin{bmatrix} -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \end{bmatrix}, \begin{bmatrix} -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ 1 & 1 & 1 \\ -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}$$

convolution filters

# Example

$$
\begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & -2 & 0 & -2 \end{bmatrix}
\qquad
\begin{bmatrix} -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \end{bmatrix}
$$

# Example

$$\begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix}_2$$

$$\begin{bmatrix} -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \end{bmatrix}$$

$$\begin{bmatrix} 2 & \cdot \\ \cdot & \cdot \end{bmatrix}$$

$$\left( 0 \times \frac{-1}{2} \right) + (1 \times 1) + \left( 0 \times \frac{-1}{2} \right)$$

$$\left( 0 \times \frac{-1}{2} \right) + (1 \times 1) + \left( 0 \times \frac{-1}{2} \right)$$

$$\left( 1 \times \frac{-1}{2} \right) + (1 \times 1) + \left( 1 \times \frac{-1}{2} \right) = 2$$

# Example

$$\begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & . & 0 & . \end{bmatrix} \quad \begin{bmatrix} -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \end{bmatrix} \quad \Longrightarrow \quad \begin{bmatrix} 2 & -2 \\ . & . \end{bmatrix}$$

# Example

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & -2 & 0 \end{bmatrix}$$

$$\begin{bmatrix} -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \end{bmatrix}$$

$\longrightarrow$

$$\begin{bmatrix} 2 & -2 \\ -1 & \cdot \end{bmatrix}$$

# Example

$$\begin{bmatrix} 0 & & & \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & -2 & 0 & -2 \end{bmatrix} \begin{bmatrix} -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \end{bmatrix} \longrightarrow \begin{bmatrix} 2 & -2 \\ -1 & 1 \end{bmatrix}$$

# Example



Input image     Convolution     Features

# Example



$$\begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & -2 & 0 & -2 \end{bmatrix}$$

$$\begin{bmatrix} -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \end{bmatrix}$$

$$\begin{bmatrix} -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ 1 & 1 & 1 \\ -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}$$

$$\begin{bmatrix} 2 & -2 \\ -1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} -1 & -\frac{1}{2} \\ \frac{7}{2} & 4 \end{bmatrix}$$

$$\begin{bmatrix} 2 \\ 4 \end{bmatrix}$$
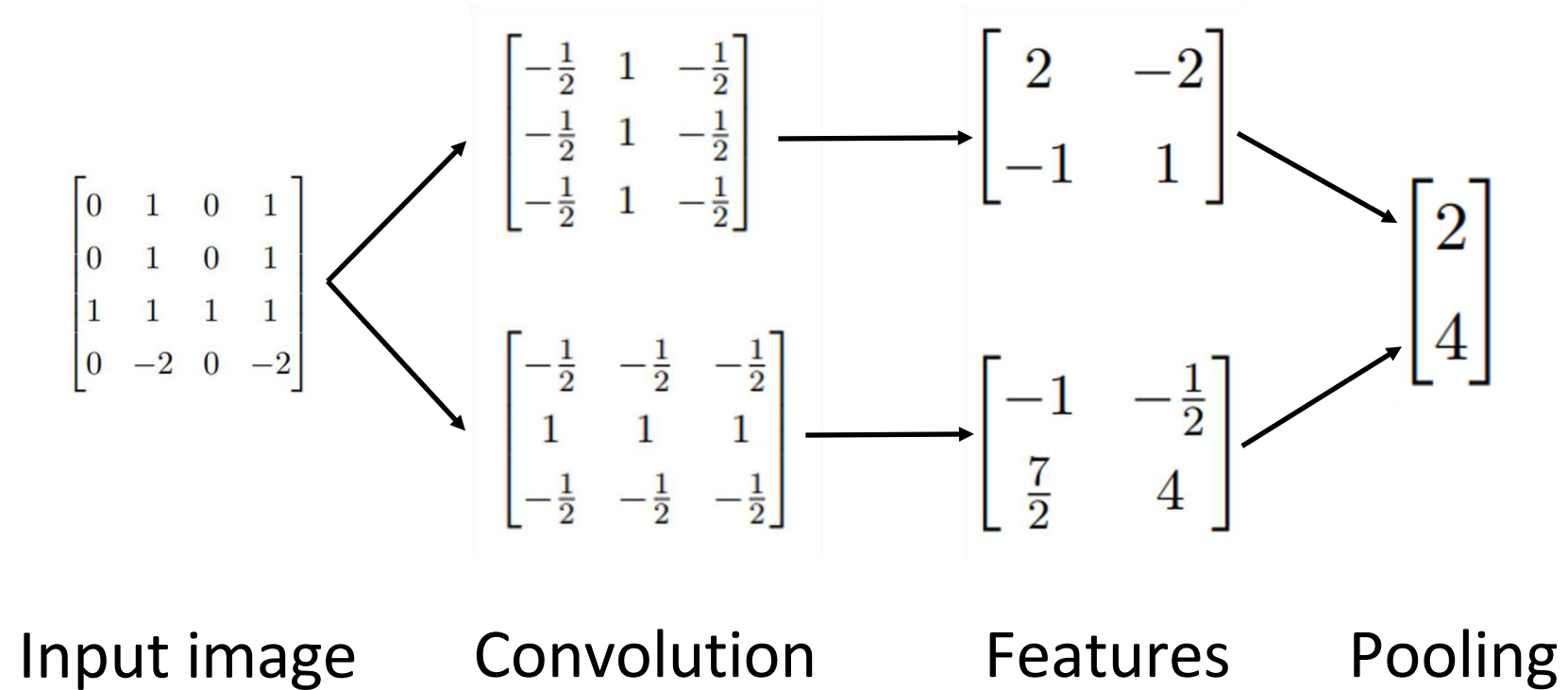
Input image     Convolution     Features     Pooling

# Example



$$\begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & -2 & 0 & -2 \end{bmatrix}$$

$$\begin{bmatrix} -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \end{bmatrix}$$

$$\begin{bmatrix} -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ 1 & 1 & 1 \\ -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}$$

$$\begin{bmatrix} 2 & -2 \\ -1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} -1 & -\frac{1}{2} \\ \frac{7}{2} & 4 \end{bmatrix}$$

$$\begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

Clear vertical line

Clear horizontal line

Input image        Convolution        Features        Pooling

$n \times n$ image $\quad f \times f$ filter

padding $p$ $\quad$ stride $s$

$$\left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \quad \times \quad \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$$

# Agenda

- **Convolutional & pooling layers**

- **Convolutional neural networks**

- **Feature visualization**

- **Applications**

# Example Architecture: AlexNet

- **ImageNet dataset**
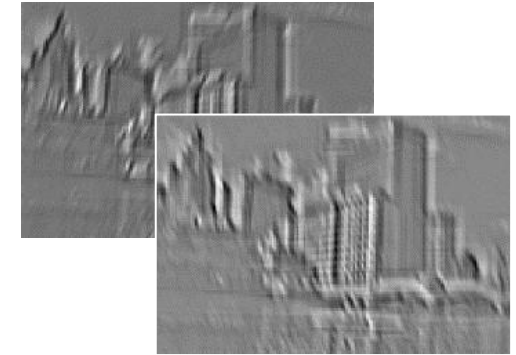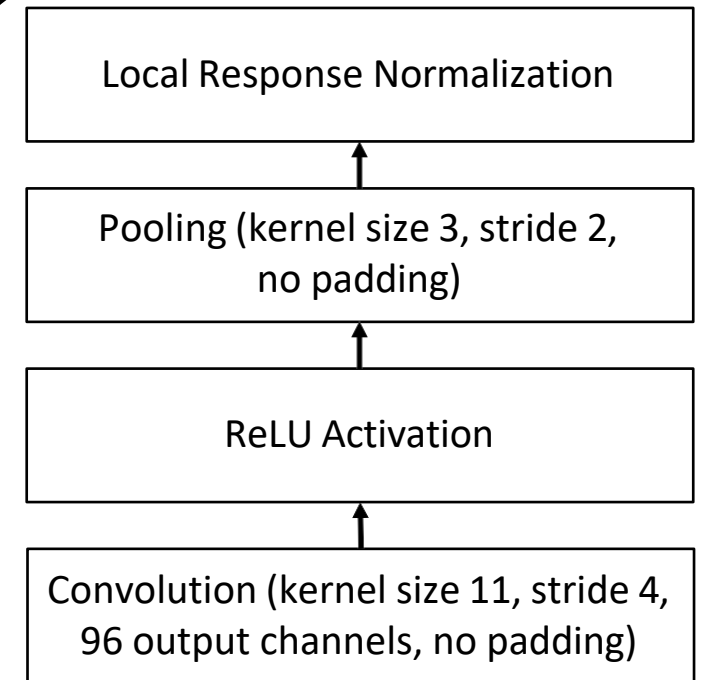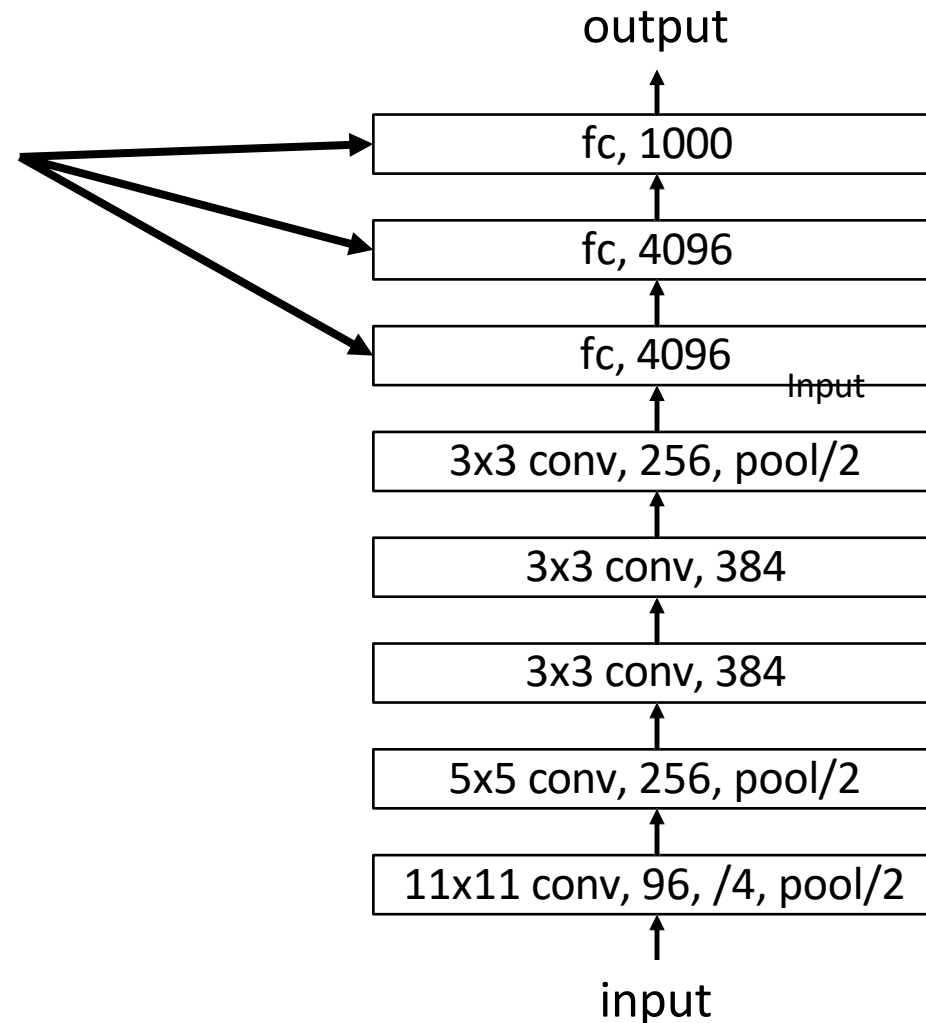  - 1000 class image classification problem (e.g., grey fox, tabby cat, barber chair)
  - >1M image-label pairs gathered from internet and crowdsourced labels

- **AlexNet Architecture (Krizhevsky 2012)**
  - Historically important architecture
  - Image classification network (~60M parameters)
  - Trained using GPUs on ImageNet dataset
  - Huge improvement in performance compared to prior state-of-the-art

# Example Architecture: AlexNet



**output**

**Fully connected (i.e., linear) layers**

| fc, 1000 |
| fc, 4096 |
| fc, 4096 | Input |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 5x5 conv, 256, pool/2 |
| 11x11 conv, 96, /4, pool/2 |

**input**

Local Response Normalization

Pooling (kernel size 3, stride 2, no padding)

ReLU Activation

Convolution (kernel size 11, stride 4, 96 output channels, no padding)

# Example Architecture: AlexNet

# Aside: Local Response Normalization

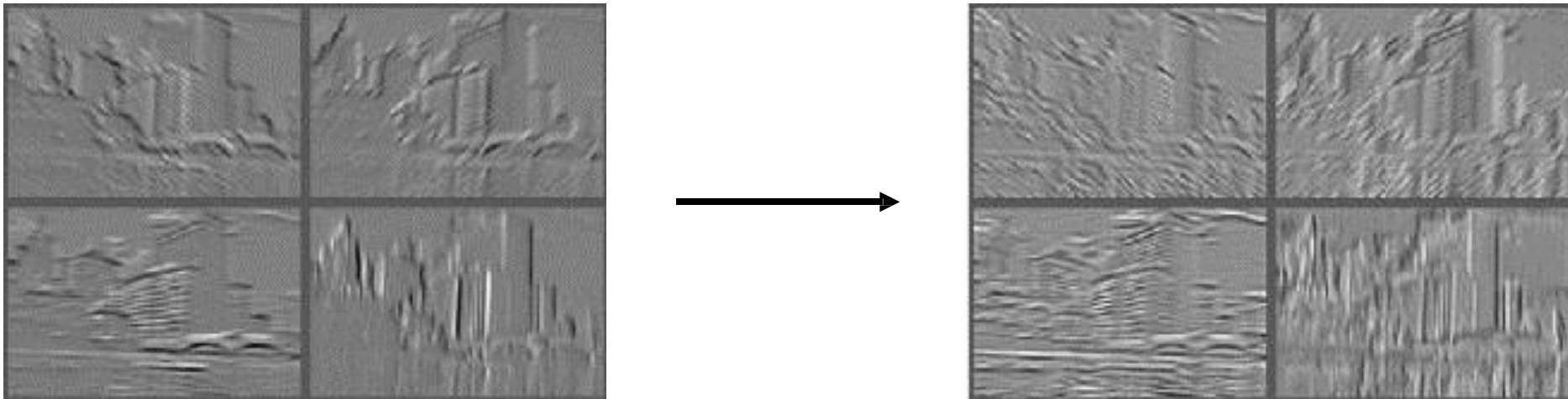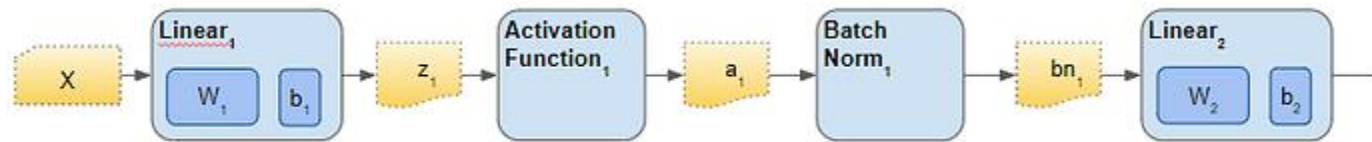- Highlights areas where the feature maps change

- Historically a standard layer, but no longer used

- Also called "contrastive normalization"

# Convolutional Neural Networks

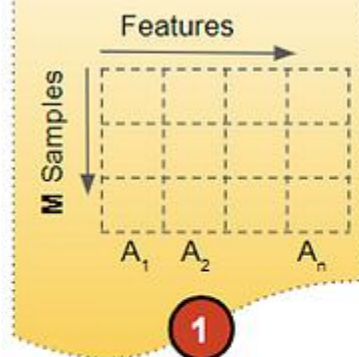- "Convolutional layer" often refers to sequence of layers

- **Modern sequence of layers**
  - Convolution -> Batch Normalization -> Pooling -> ReLU
  - Convolution -> Batch Normalization -> ReLU -> Pooling

- Can also omit pooling (especially for very deep neural networks)



https://towardsdatascience.com/batch-norm-explained-visually-how-it-works-and-why-neural-networks-need-it-b18919692739

**Batch Norm**

**Mini-batch: Activations**

Features →

M Samples ↓

$A_1$ $A_2$ $A_n$

**(1)**

**Mean and Std Dev**

$$\mu_i = \frac{1}{M} \sum A_i$$

$$\sigma_i = \sqrt{\frac{1}{M} \sum (A_i - \mu)^2}$$

**(2)**

**Normalize**

$$\hat{A}_i = \frac{A_i - \mu_i}{\sigma_i}$$

**(3)**

**Scale and Shift**

$$B\tilde{N}_i = \gamma \odot \hat{A}_i + \beta$$

Beta (β)   Gamma (γ)

**(4)**

Output (bn)

**Moving Average**

$$\mu_{mov_i} = \alpha \mu_{mov_i} + (1 - \alpha)\mu_i$$

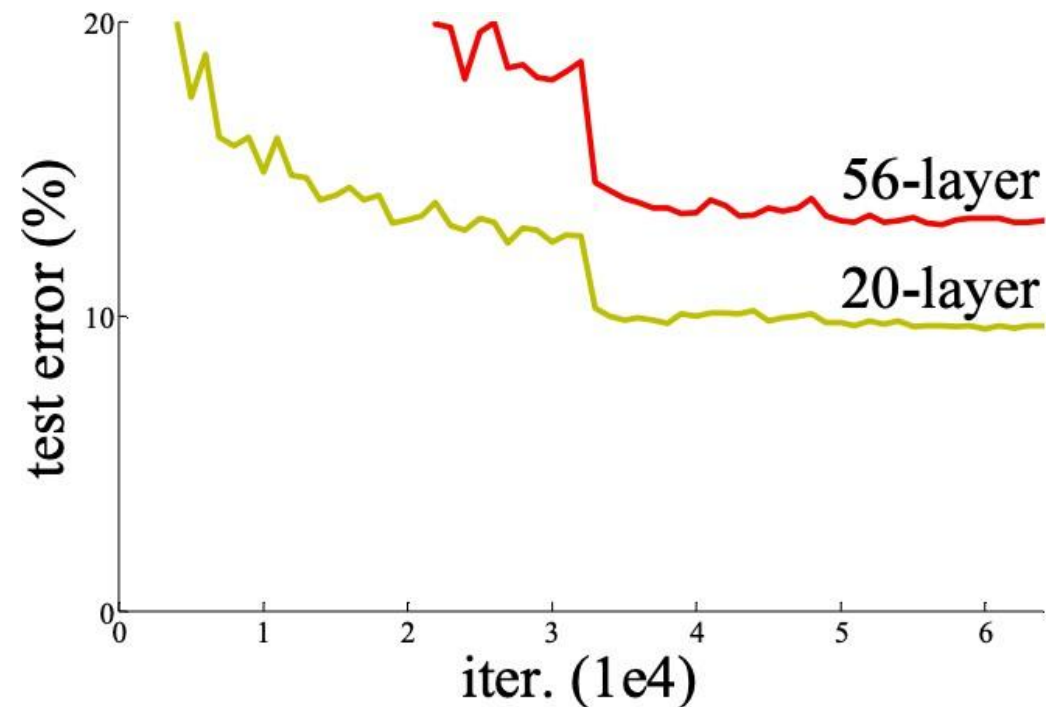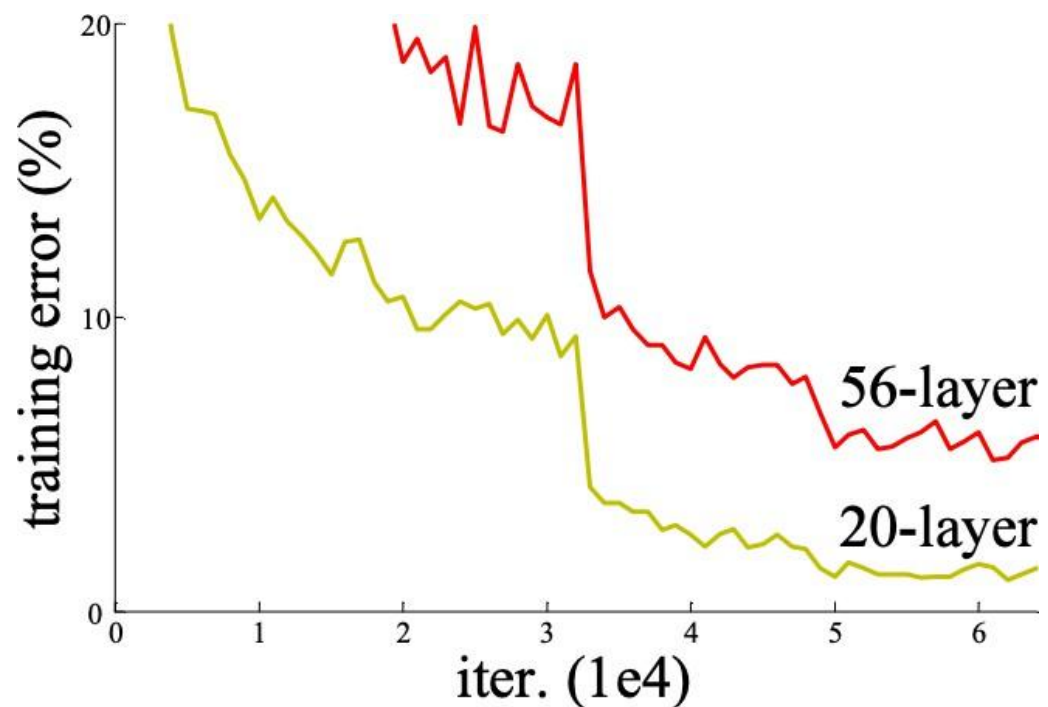$$\sigma_{mov_i} = \alpha \sigma_{mov_i} + (1 - \alpha)\sigma_i$$

Moving Avg (mean)   Moving Avg (Var)

**(5)**

# Residual Connections

- **Challenges with deeper networks**
  - Overfitting?
  - No, 56 layer network underfits!

# Residual Connections

- **Challenges with deep networks**
  - Overfitting?
  - No, 56 layer network underfits!

- **Optimization/representation**
  - Difficulty representing the identity function!

- **Solution:** "Skip" connections
  - Facilitate direct feedback from loss
  - Easy to represent identity function

$H(x) = F(x) + x$

F(x) + x

relu

conv

F(x)

relu

conv

X
identity

X

Residual block

# Residual Connections

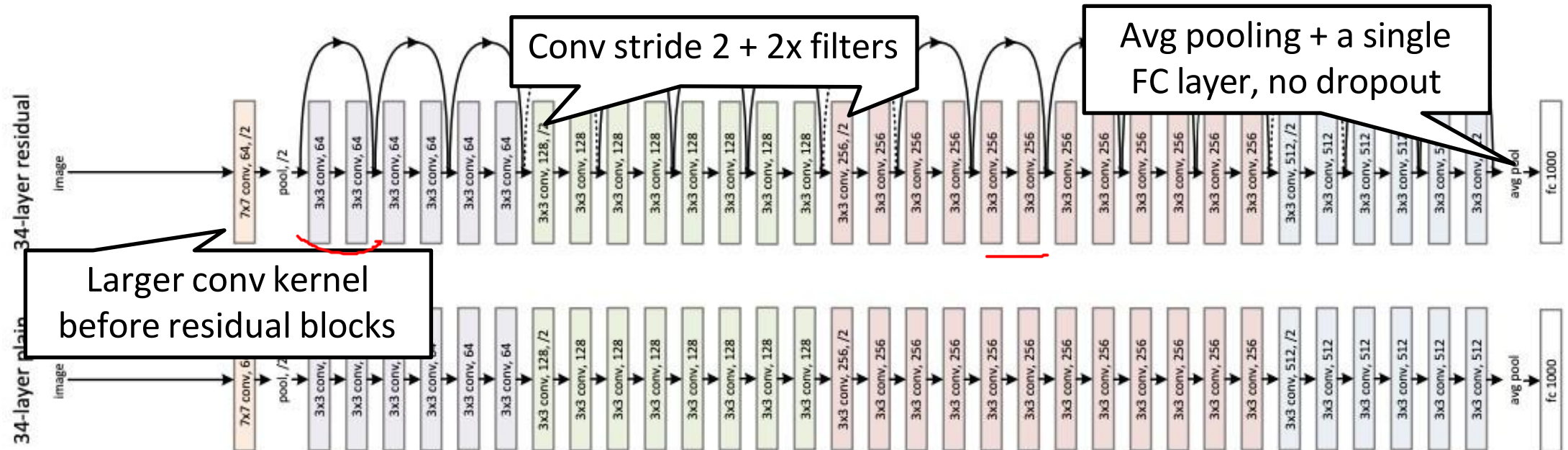- **Residual layers:** Given any convolutional layer $F(x)$, use

$$H(x) = F(x) + x$$

- Two views of residual connections:
  - **View 1:** Providing shortcuts to gradients on the backward pass
  - **View 2:** Allow each "residual block" to fit the residual error (boosting!)

$$F(x) = H(x) - x$$

# Residual Networks

- Stack lots of residual blocks!
  - Kernel size 3, no padding, stride 1, no pooling
  - Reduce feature dimensions by using stride 2 once every $K$ blocks
  - Maintains feature size to build very deep nets

Conv stride 2 + 2x filters

Avg pooling + a single FC layer, no dropout

Larger conv kernel before residual blocks

Image credit: He et al, Residual Nets, 2015

# Residual Networks

- For deeper networks, improve efficiency through 1x1 convolutions
- Many other improvements since 2015!
  - E.g., "ResNeXt", "Identity Mappings", "ConvNeXt" etc.
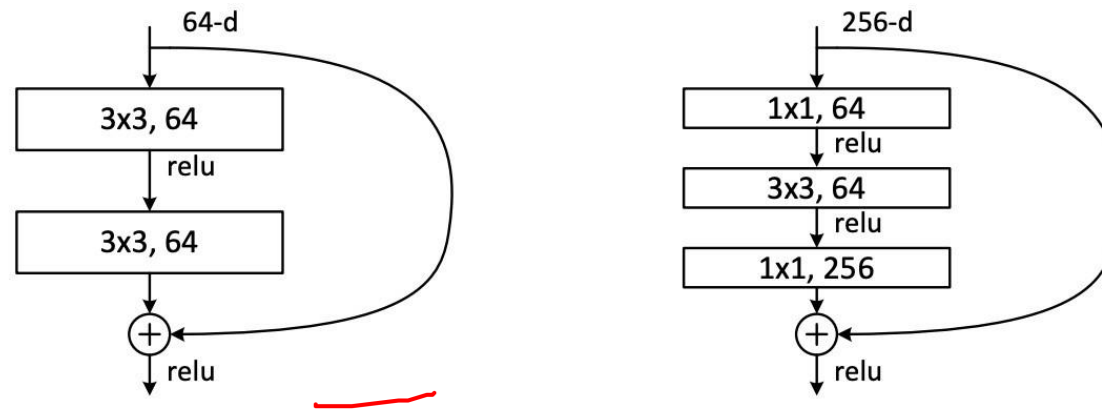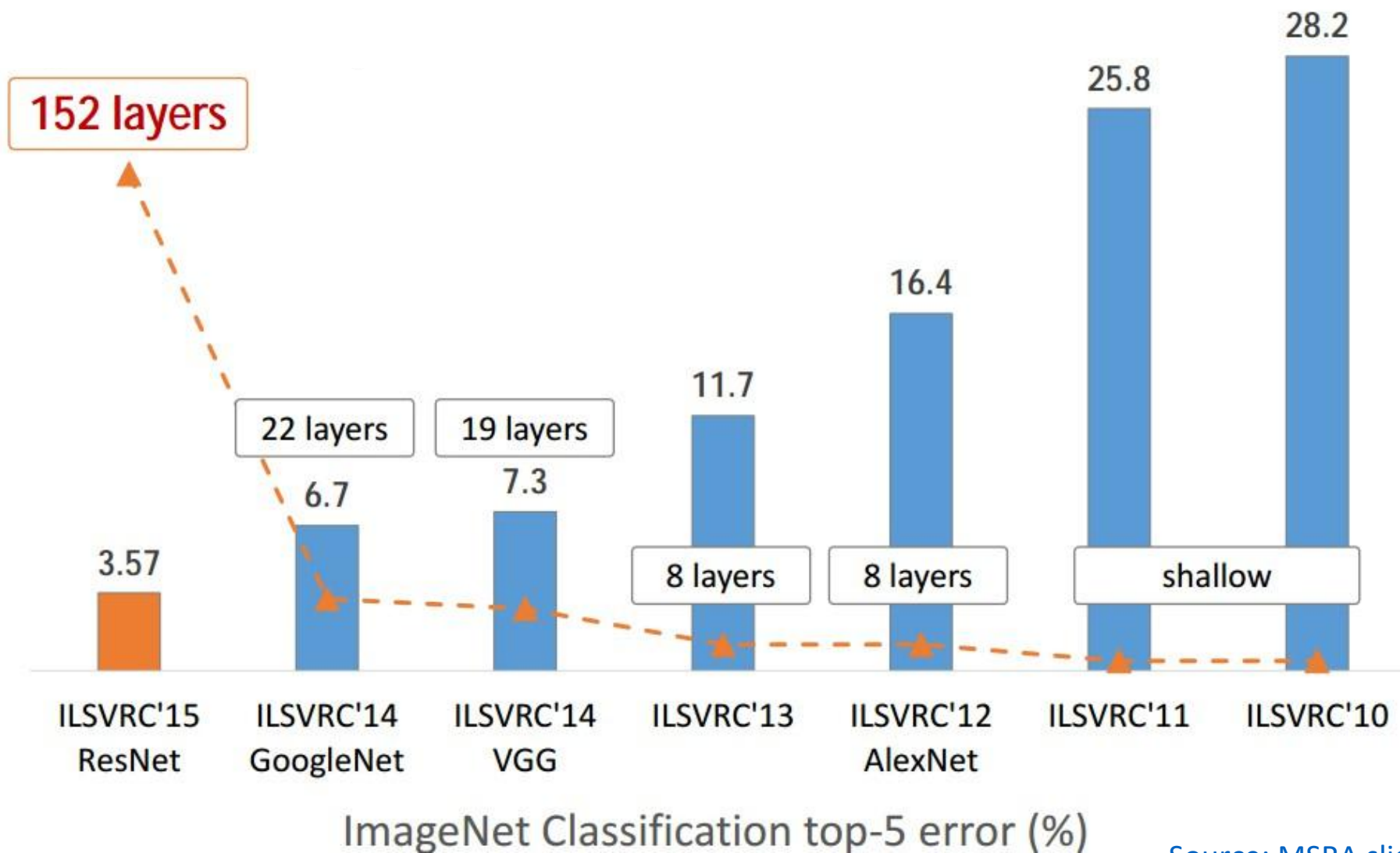


Figure 5. A deeper residual function $\mathcal{F}$ for ImageNet. Left: a building block (on 56×56 feature maps) as in Fig. 3 for ResNet-34. Right: a "bottleneck" building block for ResNet-50/101/152.

Image credit: He et al, Residual Nets, 2015

# Evolution of Neural Networks



**152 layers**

3.57

ILSVRC'15
ResNet

22 layers
6.7

ILSVRC'14
GoogleNet

19 layers
7.3

ILSVRC'14
VGG

11.7
8 layers

ILSVRC'13

16.4
8 layers

ILSVRC'12
AlexNet

25.8
shallow

ILSVRC'11

28.2

ILSVRC'10

ImageNet Classification top-5 error (%)

# Evolution of Neural Networks

AlexNet, 8 layers
(ILSVRC 2012)
~60M params

VGG, 19 layers
(ILSVRC 2014)
~140M params

ResNet, 152 layers
(ILSVRC 2015)

Less computation
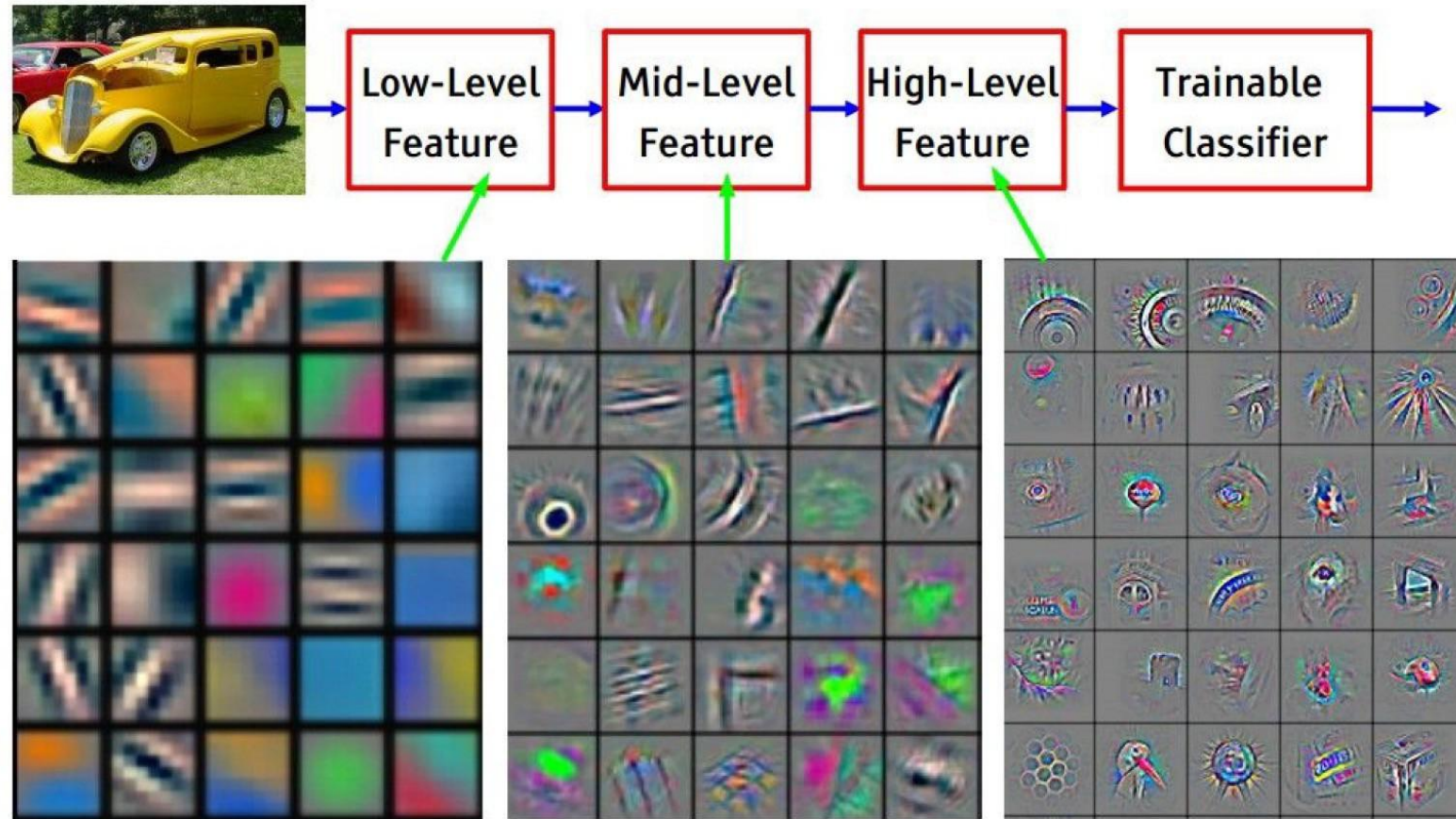in forward pass
than VGGNet!

Back to 60M params

GoogleNet, 22 layers
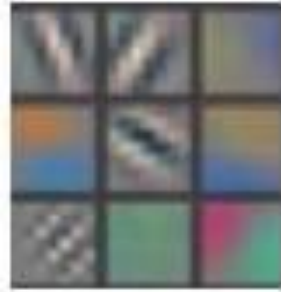(ILSVRC 2014)
~5M params

# Agenda

- **Convolutional & pooling layers**

- **Convolutional neural networks**

- **Feature visualization**

- **Applications**
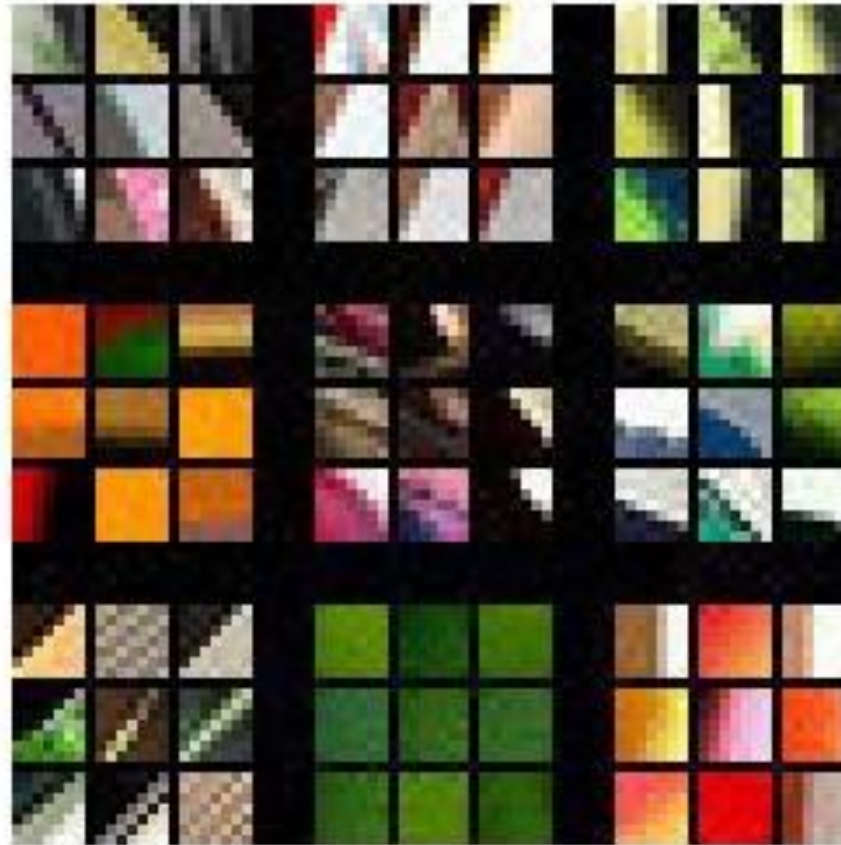
# Feature Visualization



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

Slide credit: Yann LeCun

# Layer 1



Visualizing and Understanding Convolutional Networks [Zeiler and Fergus, ECCV 2014]

# Layer 2



Visualizing and Understanding Convolutional Networks [Zeiler and Fergus, ECCV 2014]

# Layer 3



Visualizing and Understanding Convolutional Networks [Zeiler and Fergus, ECCV 2014]

# Layer 3



Layer 4

Layer 5

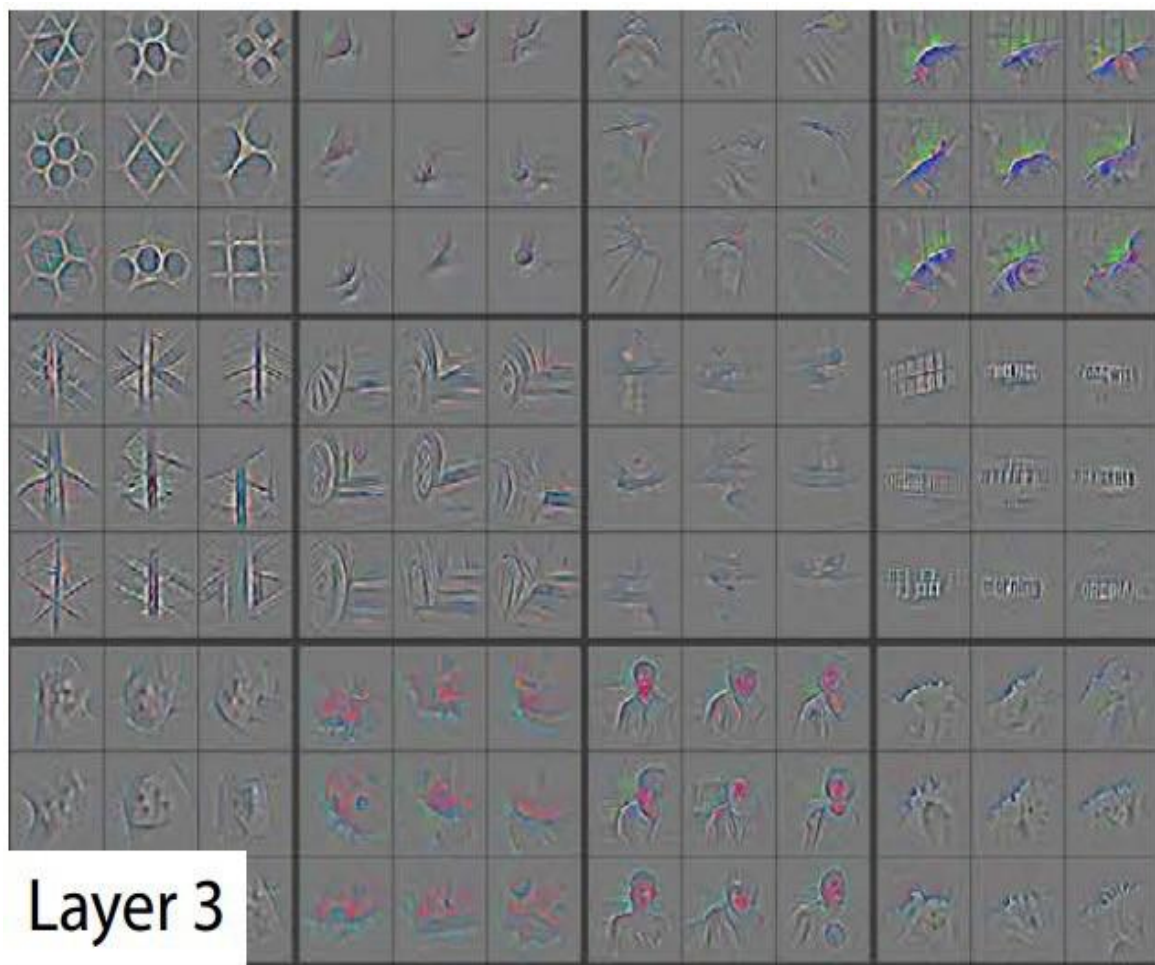Visualizing and Understanding Convolutional Networks [Zeiler and Fergus, ECCV 2014]
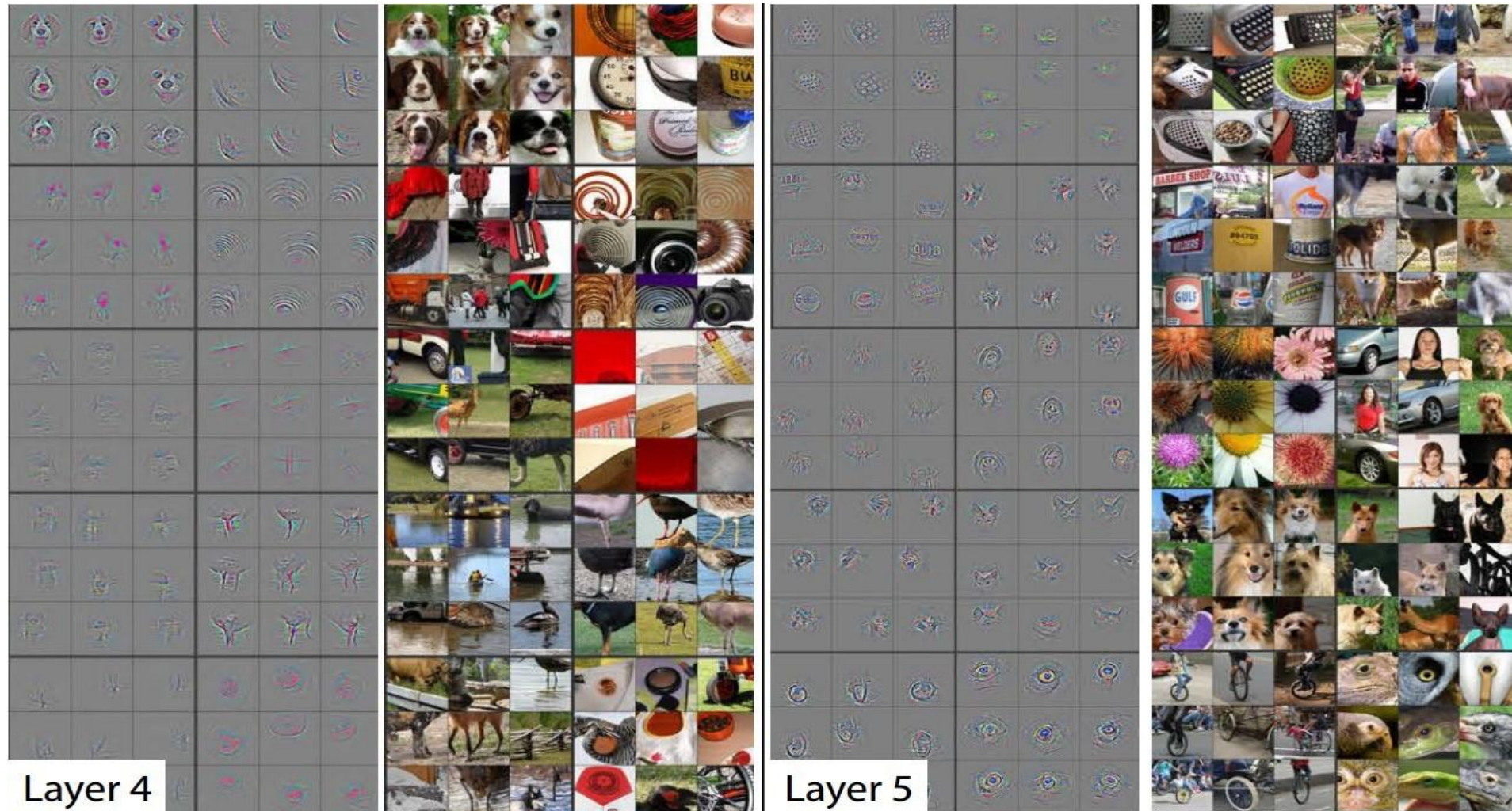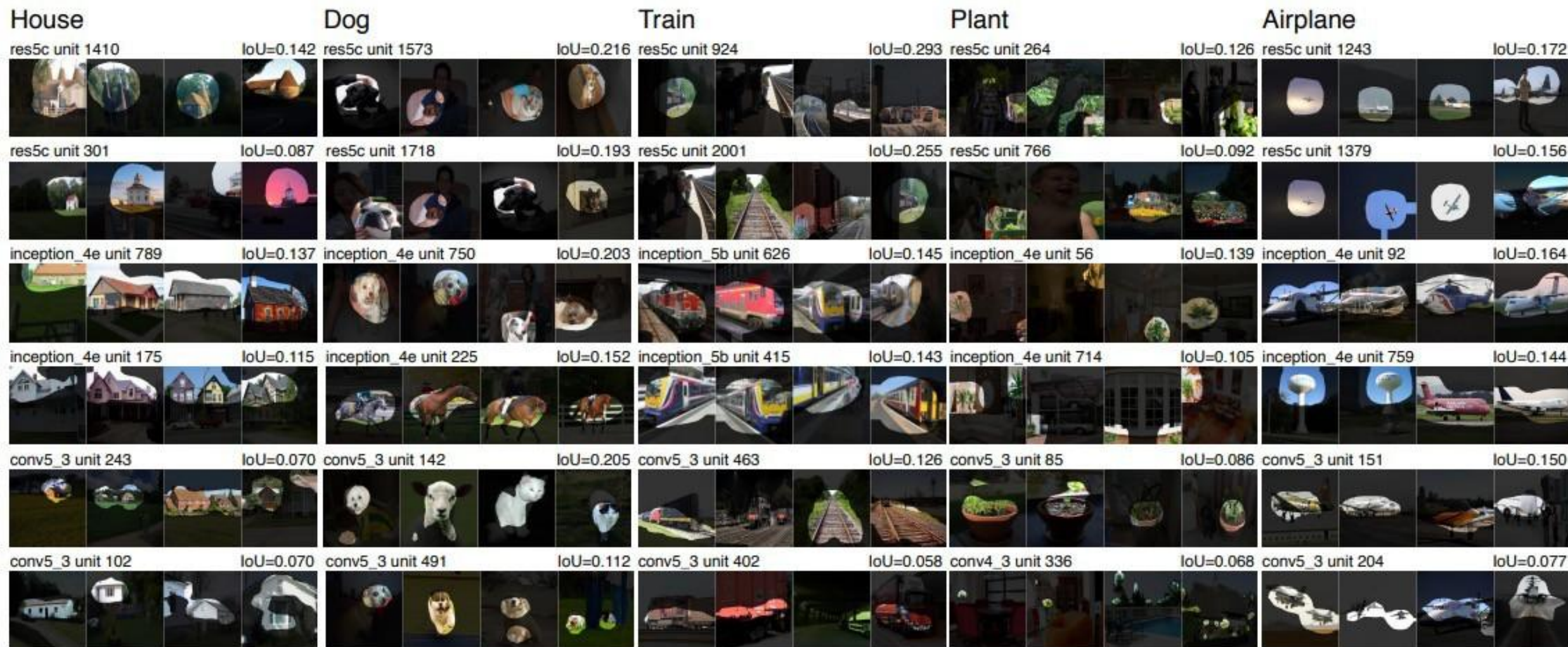
# Neural Network Dissection

# What About Small Datasets?

- **Transfer learning:** We can reuse trained concepts!
  - Since CNNs trained on ImageNet appear to learn general features
  - We can reuse these models in some way to perform new tasks

- **Strategy 1:** Feature extraction
  - Remove final (softmax) layer and replace with a new one
  - Train only the new layer

- **Strategy 2:** Finetuning
  - Do the same thing but train end-to-end

# What About Small Datasets?

- **New dataset is similar to the original dataset**
  - Can use very small datasets
  - Both strategies work


- **New dataset is different from original dataset**
  - Transfer learning still works!
  - Moderate-sized datasets
  - Finetune end-to-end
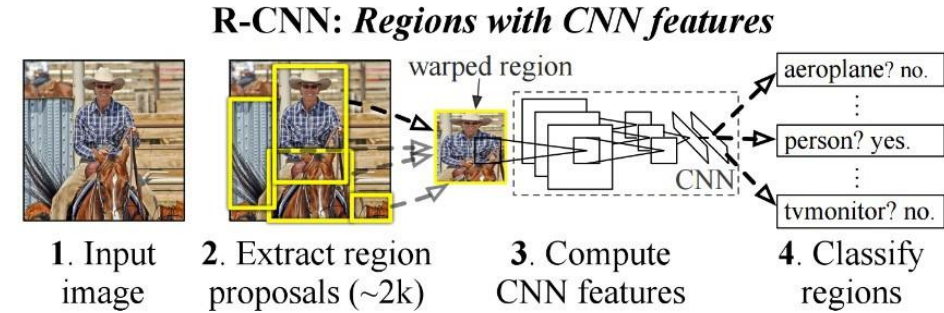  - **Examples:** Medical images, audio spectrograms, etc.

# Agenda

- **Convolutional & pooling layers**

- **Convolutional neural networks**

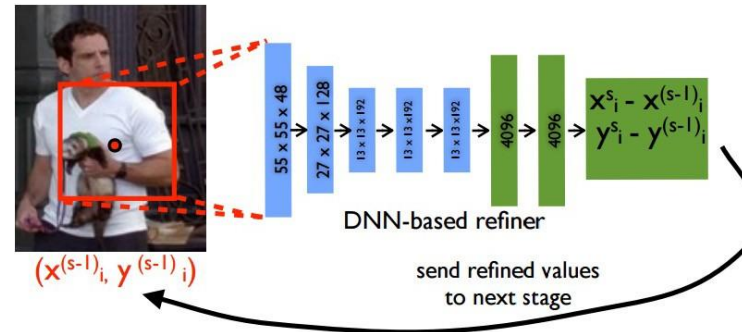- **Feature visualization**

- **Applications**

# Applications

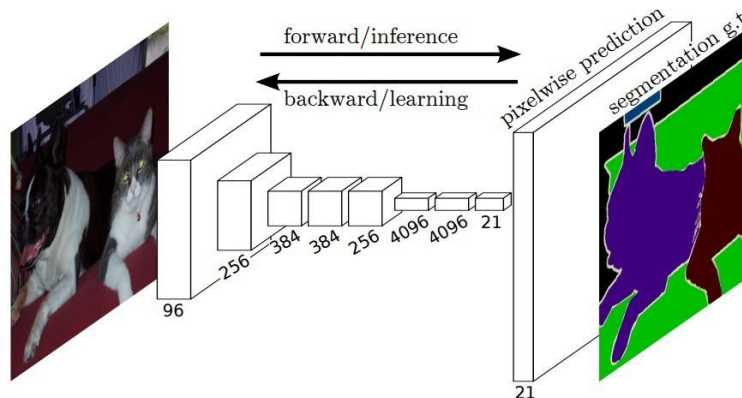**Object detection**

[Girshick et al. CVPR14]

**R-CNN:** *Regions with CNN features*

warped region

aeroplane? no.
.
person? yes.
.
tvmonitor? no.

CNN

1. Input image
2. Extract region proposals (~2k)
3. Compute CNN features
4. Classify regions

**Pose detection (regression)**

[Toshev et al. CVPR14]

$55 \times 55 \times 48$
$27 \times 27 \times 128$
$13 \times 13 \times 192$
$13 \times 13 \times 192$
$13 \times 13 \times 192$
4096
4096

$x^s_i - x^{(s-1)}_i$
$y^s_i - y^{(s-1)}_i$

DNN-based refiner

$(x^{(s-1)}_i, y^{(s-1)}_i)$

send refined values to next stage

**Semantic segmentation**

[Long et al. CVPR15]

forward/inference

backward/learning

pixelwise prediction

segmentation g.t.

256  384  384  256  4096  4096  21

96

21

Examples courtesy Jia-Bin Huang

# Applications

**Similarity metric learning**



[Chopra et al. CVPR05]

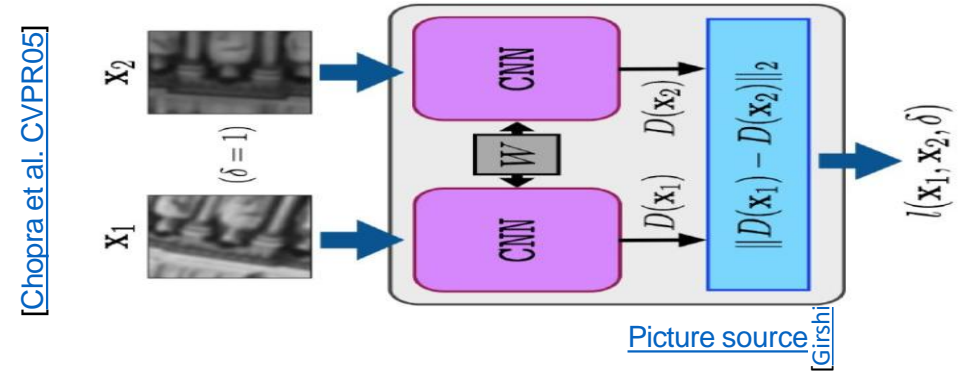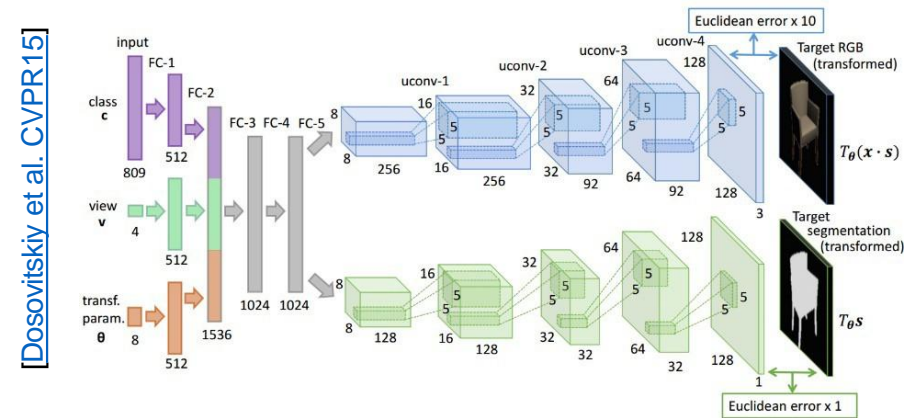Picture source [Girshi

**Image generation**



[Dosovitskiy et al. CVPR15]

**Low-level image processing: (superresolution, deblurring, image quality etc.)**



[Dong et al. ECCV 2014]

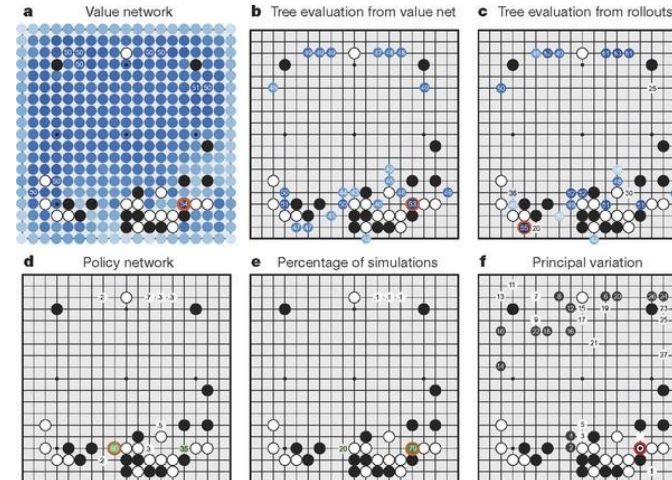Examples courtesy Jia-Bin Huang

# Applications: Game Playing

CNN + Reinforcement learning

[Mnih et al, Nature' 15]

[Silver et al, Nature '16]

Policy network

$p_{\sigma/\rho}$ (a|s)
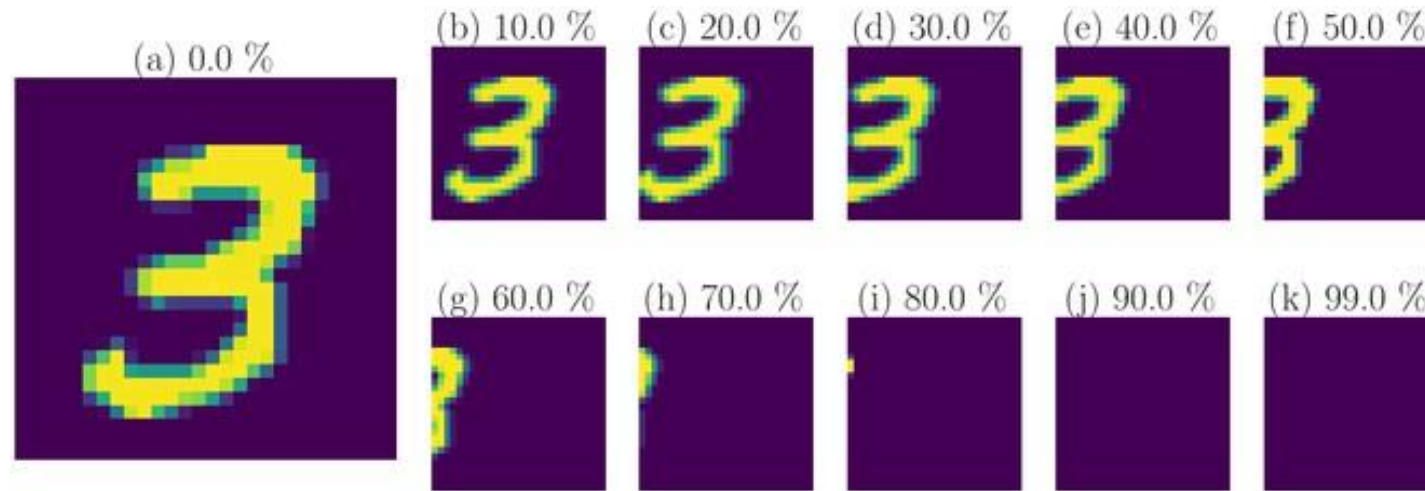
s

**a** Value network    **b** Tree evaluation from value net    **c** Tree evaluation from rollouts

**d** Policy network    **e** Percentage of simulations    **f** Principal variation

# Applications: Art Generation



See if you can tell artist originals from machine style imitations at:
http://turing.deepart.io/

Paper: Gatys et al, "Neural ... Style", arXiv '15
Code (torch): https://github.com/jcjohnson/neural-style

# Structure in Images

- **Translation invariance -** is a property of a system or model where its performance or output does not change when the input data is translated or shifted.
  - Consider image classification (e.g., labels are cat, dog, etc.)
  - **Invariance:** If we translate an image, it does not change the category label



**Source:** Ott et al., Learning in the machine: To share or not to share?

# Structure in Images

- **Translation equivariance -** is a property of a function or an operation where the output's structure does not change even if the input is translated.
  - Consider object detection (e.g., find the position of the cat in an image)
  - **Equivariance:** If we translate an image, the the object is translated similarly