



# Lecture 16: Convolutional Neural Networks

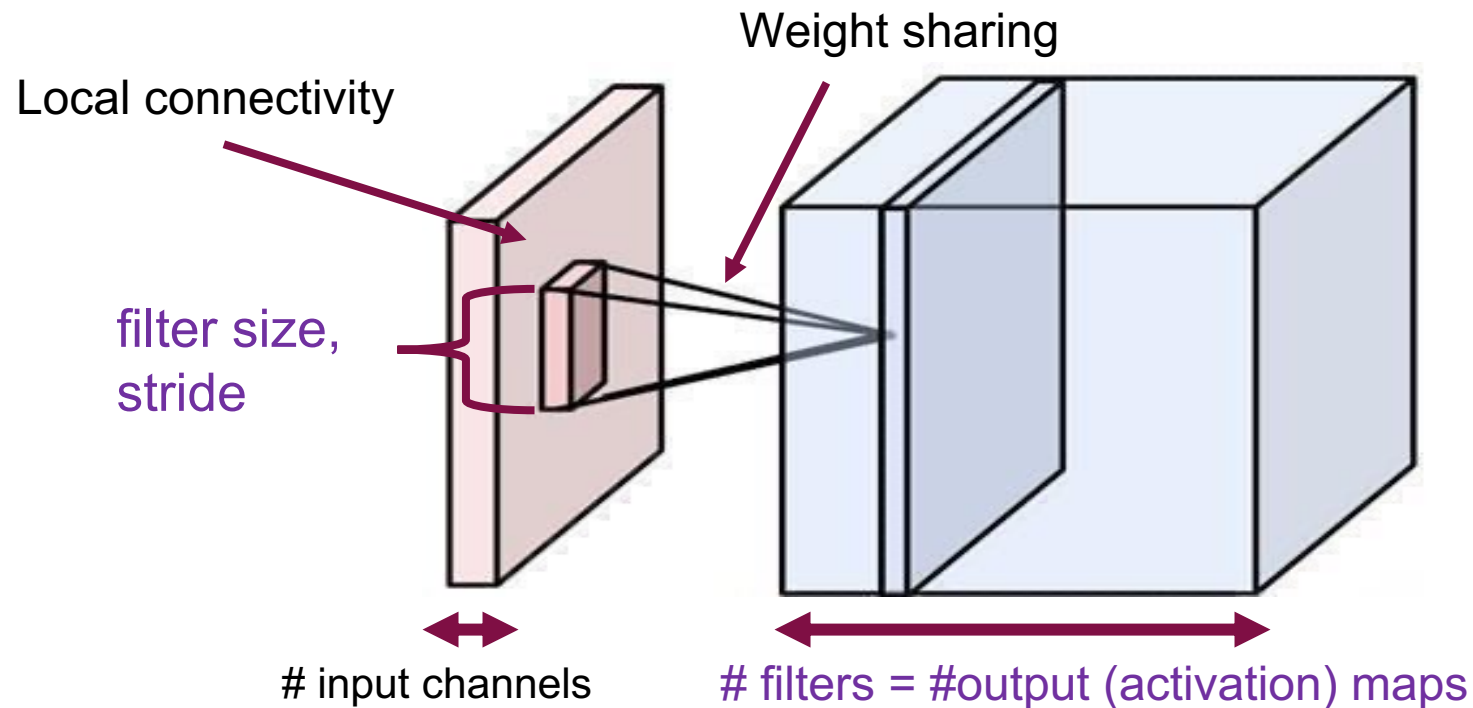
Mar 15, 2023

CIS 4190/5190

Spring 2023

# Convolutional Layer Summary

- Local connectivity
- Weight sharing
- Handling multiple input/output channels
- Retains location associations



# Stride

Filter

1	0
0	0.5

Input	Stride X →					
	0	0	0	0	0	0
	0	1	0	0.5	0.5	0
	0	0	0.5	1	0	0
	0	0	1	0.5	1	0
	0	1	0.5	0.5	1	0
	0	0	0	0	0	0
Stride Y ↓						

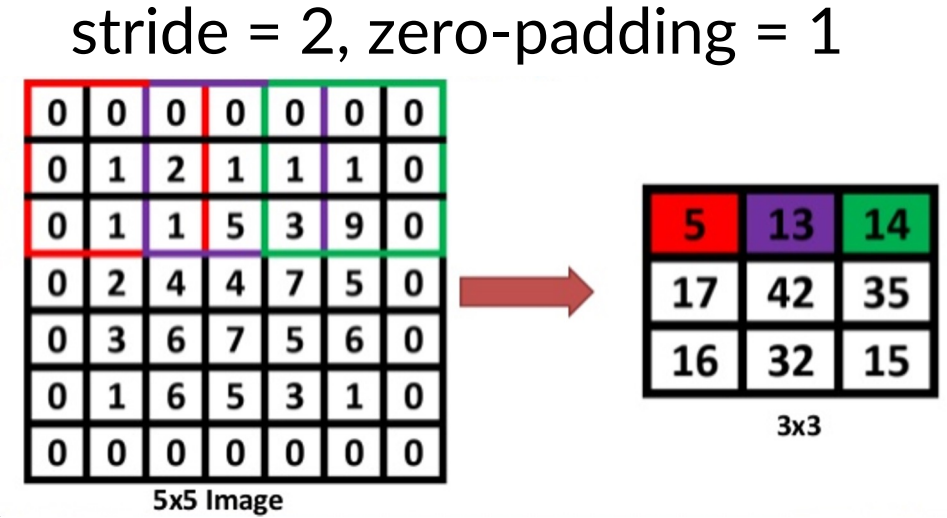
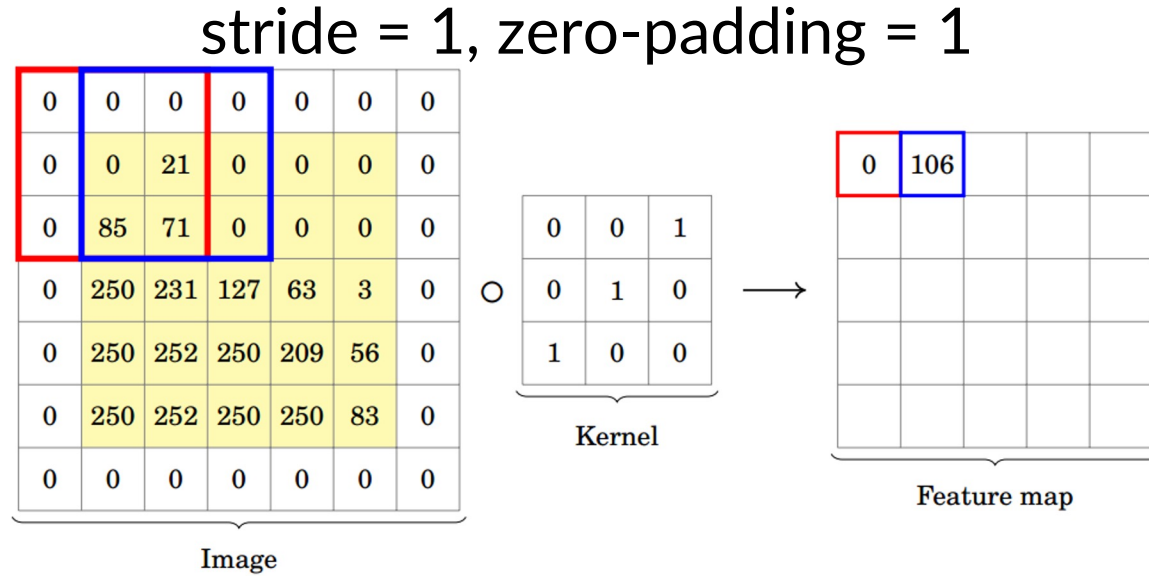
Output with stride 1

0.5	0	0.25	0.25	0
0	1.25	0.5	0.5	0.5
0	0.5	0.75	1.5	0
0.5	0.25	1.25	1	1
0.	1	0.5	0.5	1

Output with stride 2

0.50	0.25	0.00
0.00	0.75	0.00
0.00	0.50	1.00

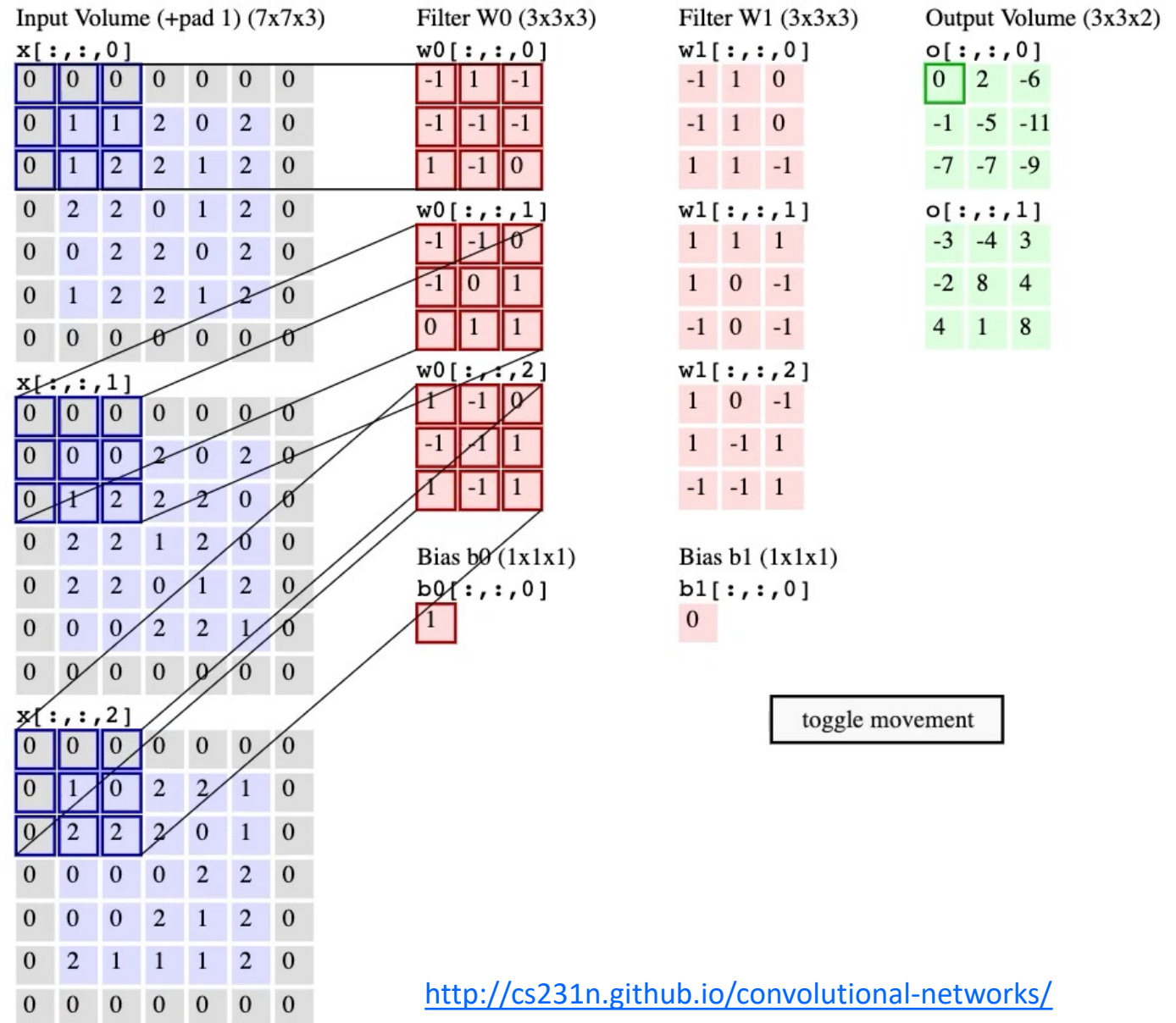
# Zero-Padding



The kernel size, amount of zero-padding, and stride, together determine the output spatial dimensions

# Convolution Filter Bank Demo

- Notes:
  - Multiple (3) inputs
    - Hence kernels of size 3x3x3
  - Multiple (2) outputs (hence 2 kernels)
  - And one bias parameter for each kernel
  - Stride 2, zero-padding 1
- Net #parameters in the bank:
  - $(3 \times 3 \times 3 + 1) \times 2 = 56$



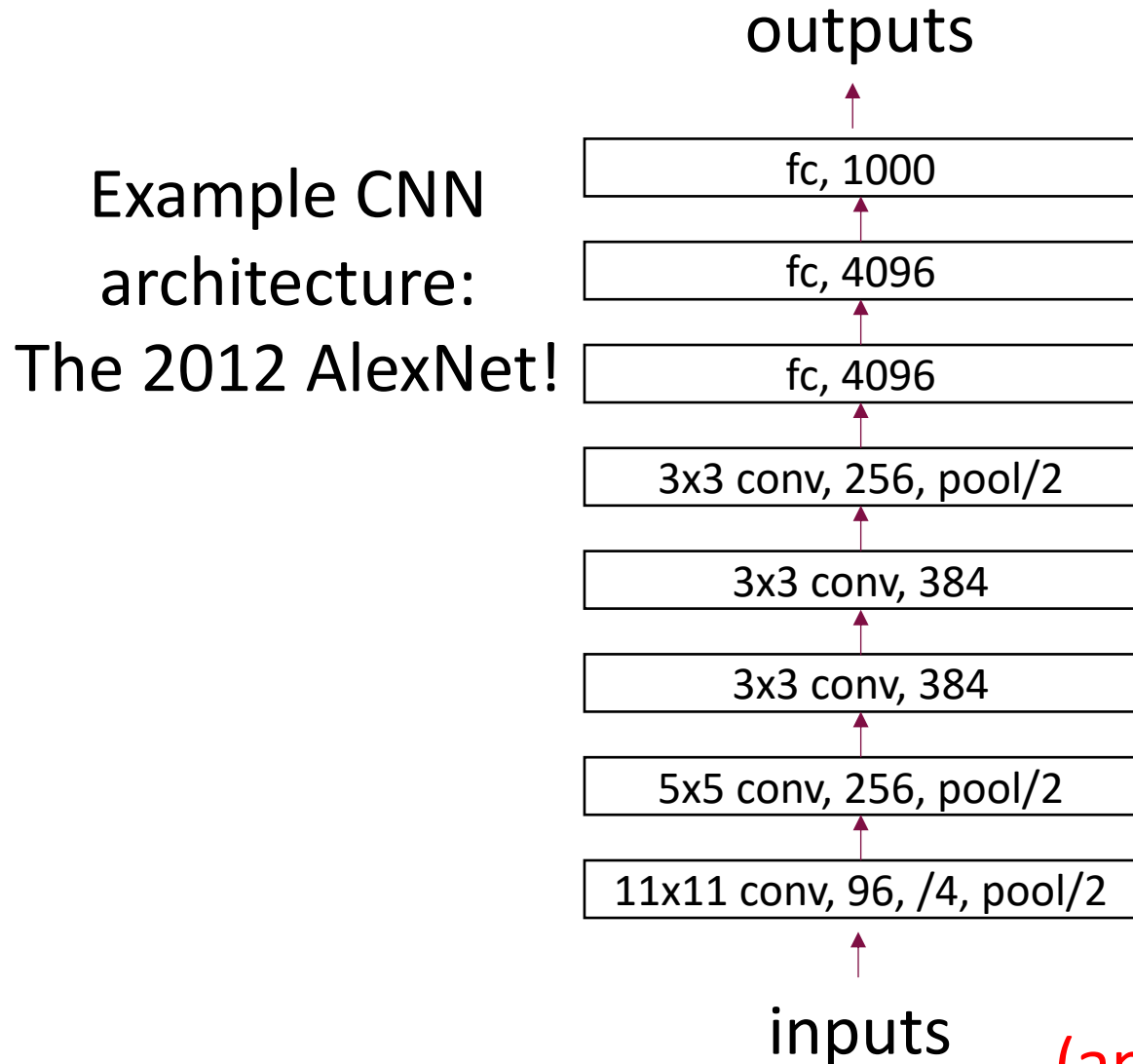
# Can we back-propagate through a convolution?

- Yes!
- A convolution is after all a special case of a linear operation  $Y = WX$ , with local connections and shared weights.
- Differentiable w.r.t. its inputs, as well as w.r.t. its weights.

# The paraphernalia around convolutions inside CNNs

Pooling, Normalization, Activation Functions ...

# Convolutions inside a neural network



“8 layers”, really “8 layer blocks”  
“5 convolution blocks” followed  
by 3 fully connected layers

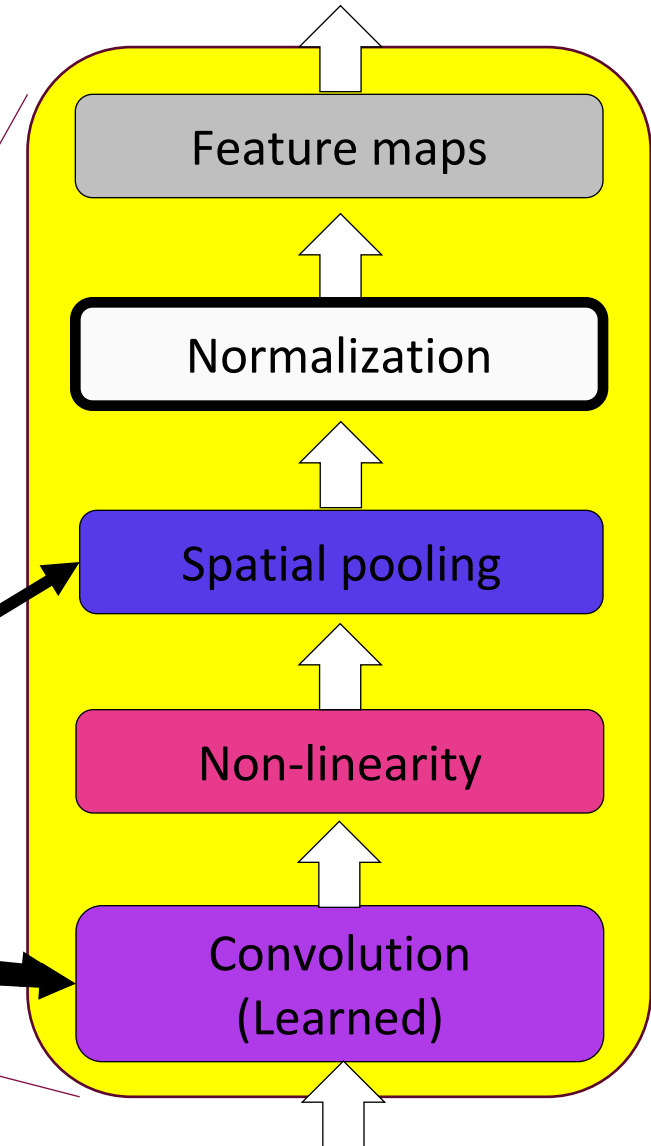
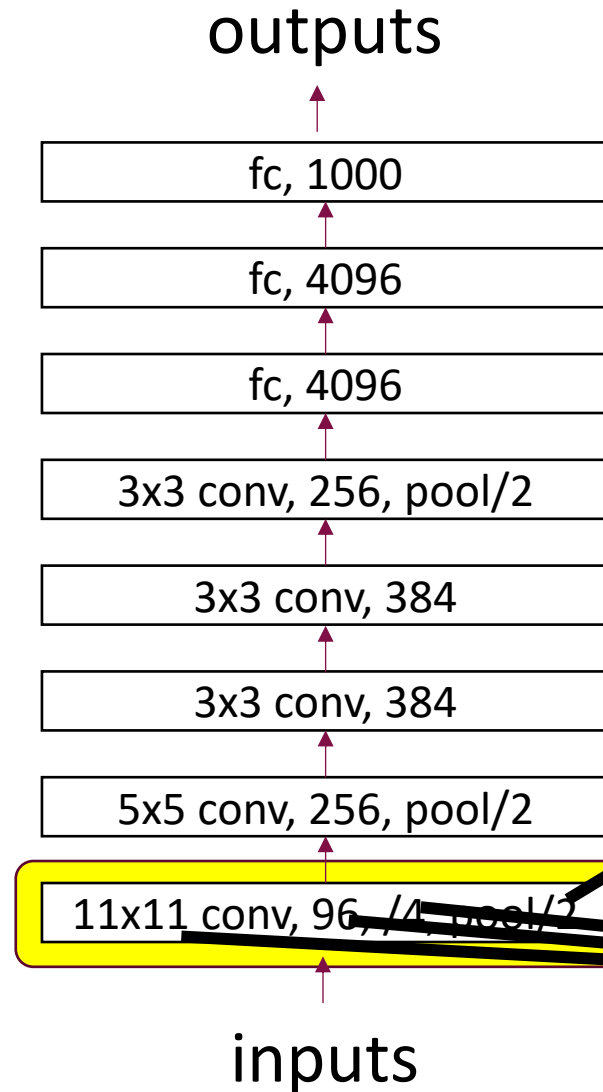
More on AlexNet soon!

But first, what is a “convolution block”?  
(and what are all the numbers in each layer?)

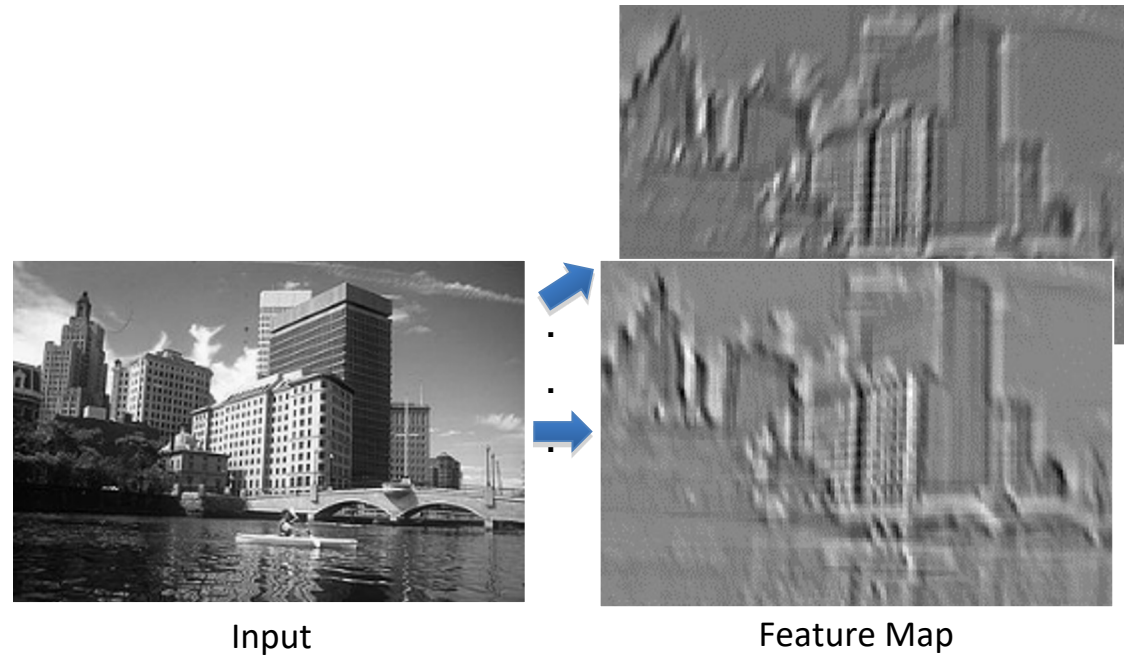
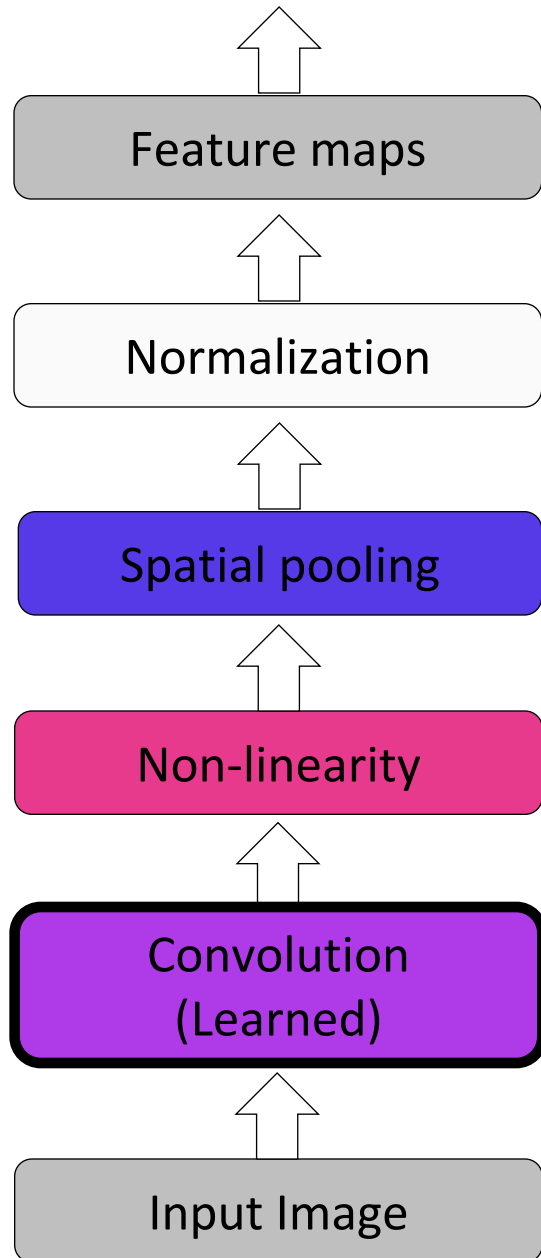


# Typical accompaniments to “convolution layers”

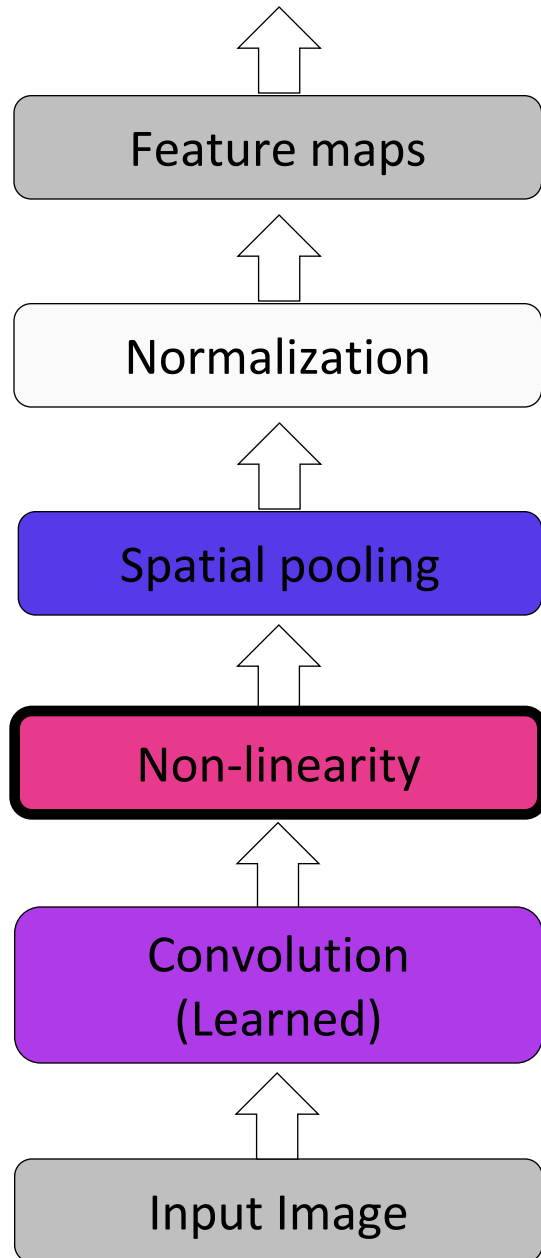
Example CNN  
architecture



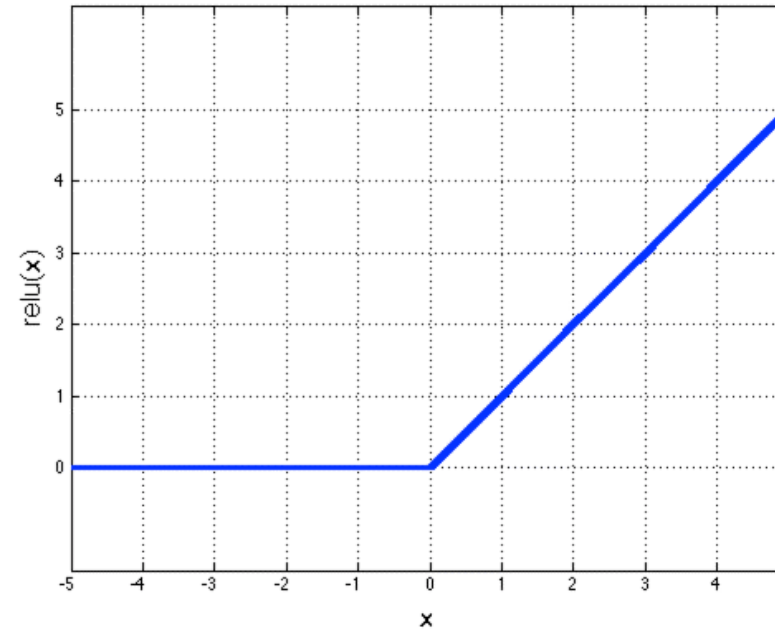
# Convolve → activation function → pool → normalize



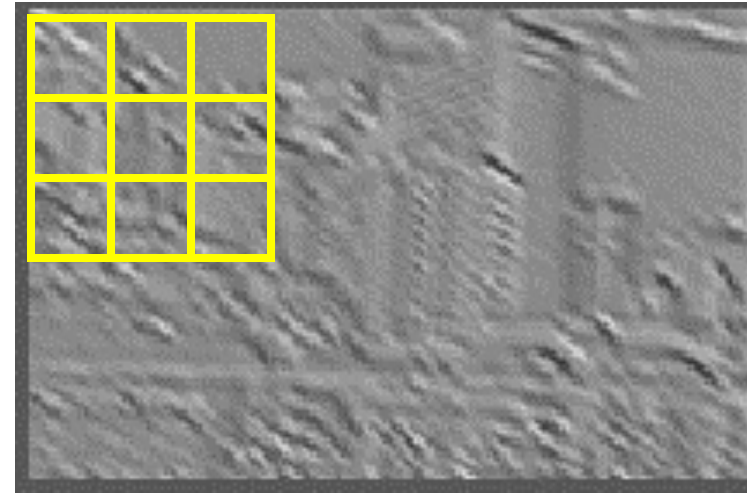
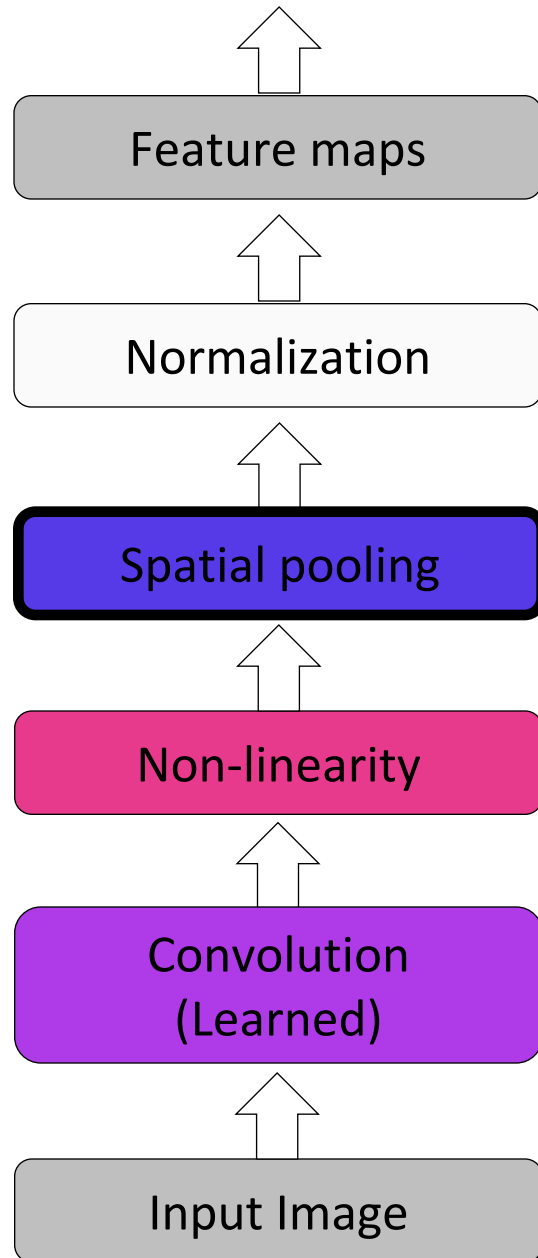
# Convolve → activation function → pool → normalize



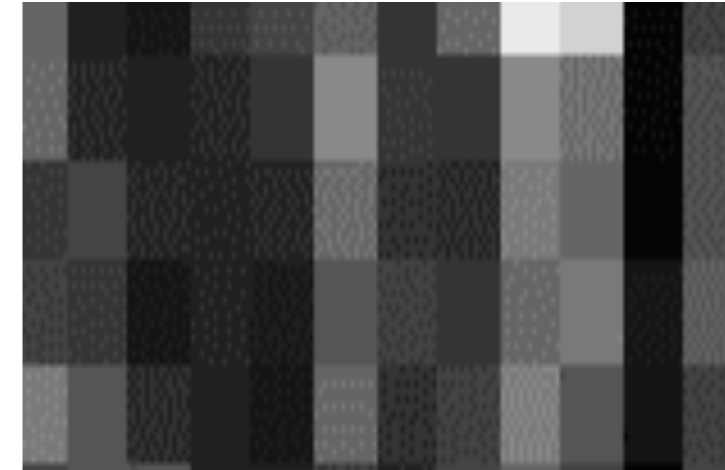
## Rectified Linear Unit (ReLU)



# Convolve → activation function → pool → normalize



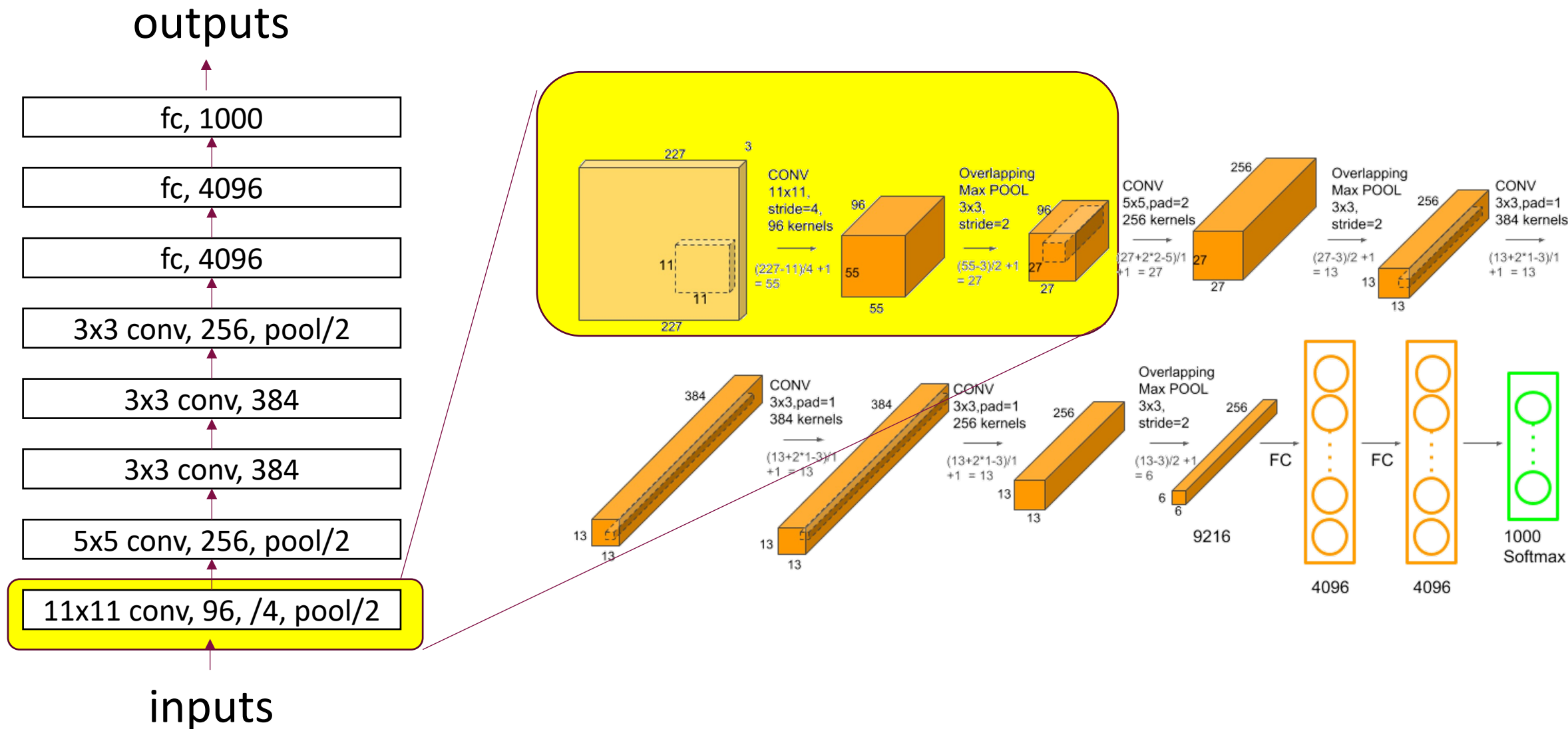
Max pooling



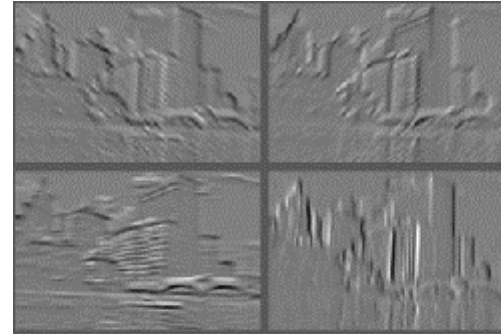
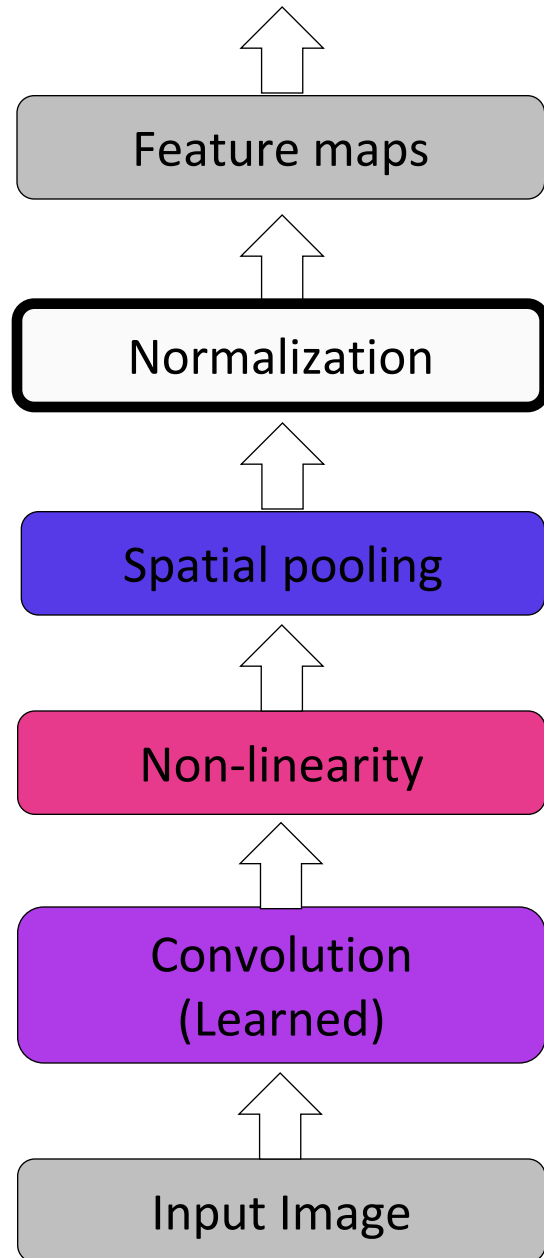
Max-pooling: *translation invariance*.  
Often applied with a stride.  
No learnable parameters.

Convolution provides equivariance to shift  
Pooling provides invariance to shift

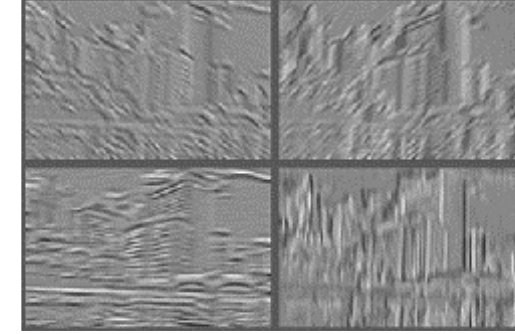
# Back to AlexNet



# Convolve → activation function → pool → normalize



Feature Maps



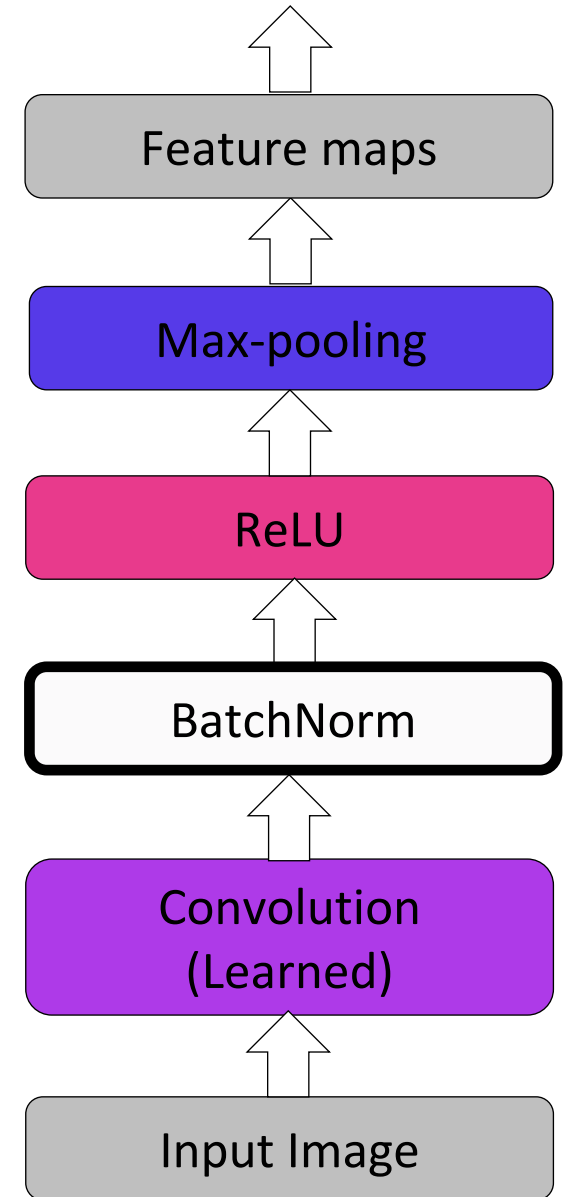
Feature Maps  
After Contrast  
Normalization

“Contrast normalization” highlights areas where the feature maps change. Used to be a standard component in neural networks. Not used in modern architectures.

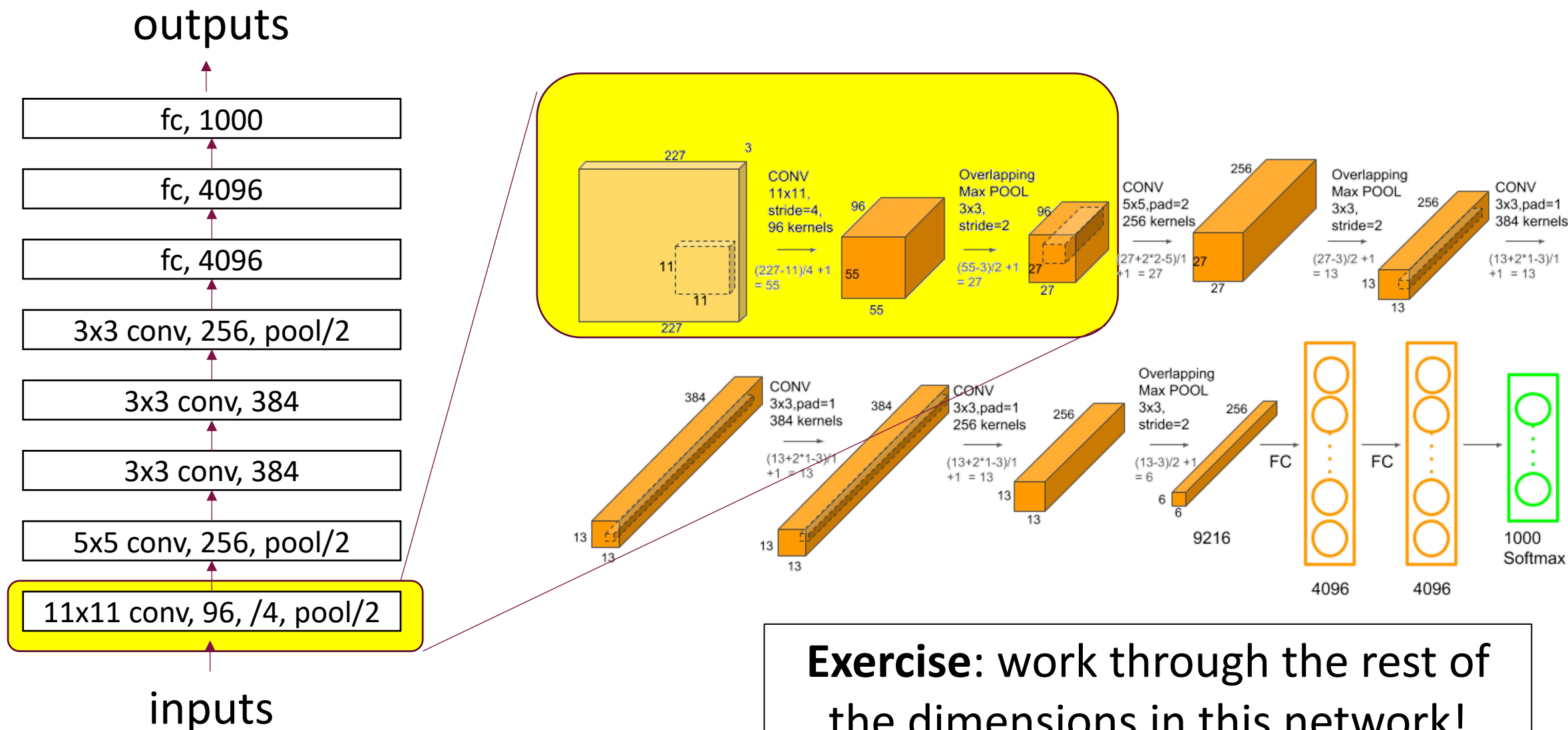
# Modern variants

- BatchNorm is very commonly used.
- Most common variants of a convolutional block:
  - Conv-BatchNorm-Maxpool-ReLU, or
  - Conv-BatchNorm-ReLU-Maxpool
- Sometimes even no Maxpool, to keep feature map spatial dimensions large. Often in very deep networks.

Often, when people say “convolution layer”, it is implicit that they mean a full convolutional block with various layers following the actual convolutional layer



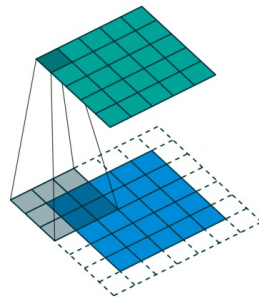
# Back to AlexNet



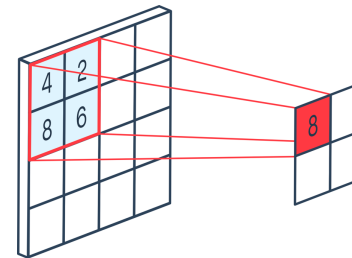


# Summary: Image-specific operations in neural nets

- Machinery to convert image matrices into vectors of reasonable dimensions, retaining useful location associations. Two main workhorses:
  - **Convolution layers** – Location-independent processing. Shift equivariance.
    - Convolutions produce “image”-like feature maps, which retain associations with input pixels.
  - **Pooling layers** – Binning to make outputs insensitive to translation and reduce dimensionality. Shift invariance.
    - A dog is a dog even if its image is shifted by a few pixels.



Convolution layers

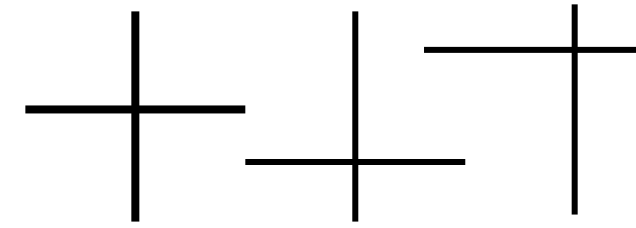


Pooling layers



# A Convolution Exercise

Suppose we want to find out whether the following image depicts Cartesian axes.



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

As a step towards this, we convolve the image with two filters (no padding, stride of 1).

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

Compute the output by hand.

# A Convolution Exercise

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

# A Convolution Exercise

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

$$\begin{aligned} &\left(0 \times \frac{-1}{2}\right) + (1 \times 1) + \left(0 \times \frac{-1}{2}\right) \\ &\left(0 \times \frac{-1}{2}\right) + (1 \times 1) + \left(0 \times \frac{-1}{2}\right) \\ &\left(0 \times \frac{-1}{2}\right) + (1 \times 1) + \left(0 \times \frac{-1}{2}\right) = 2 \end{aligned}$$

# A Convolution Exercise

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

# A Convolution Exercise

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

# A Convolution Exercise

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

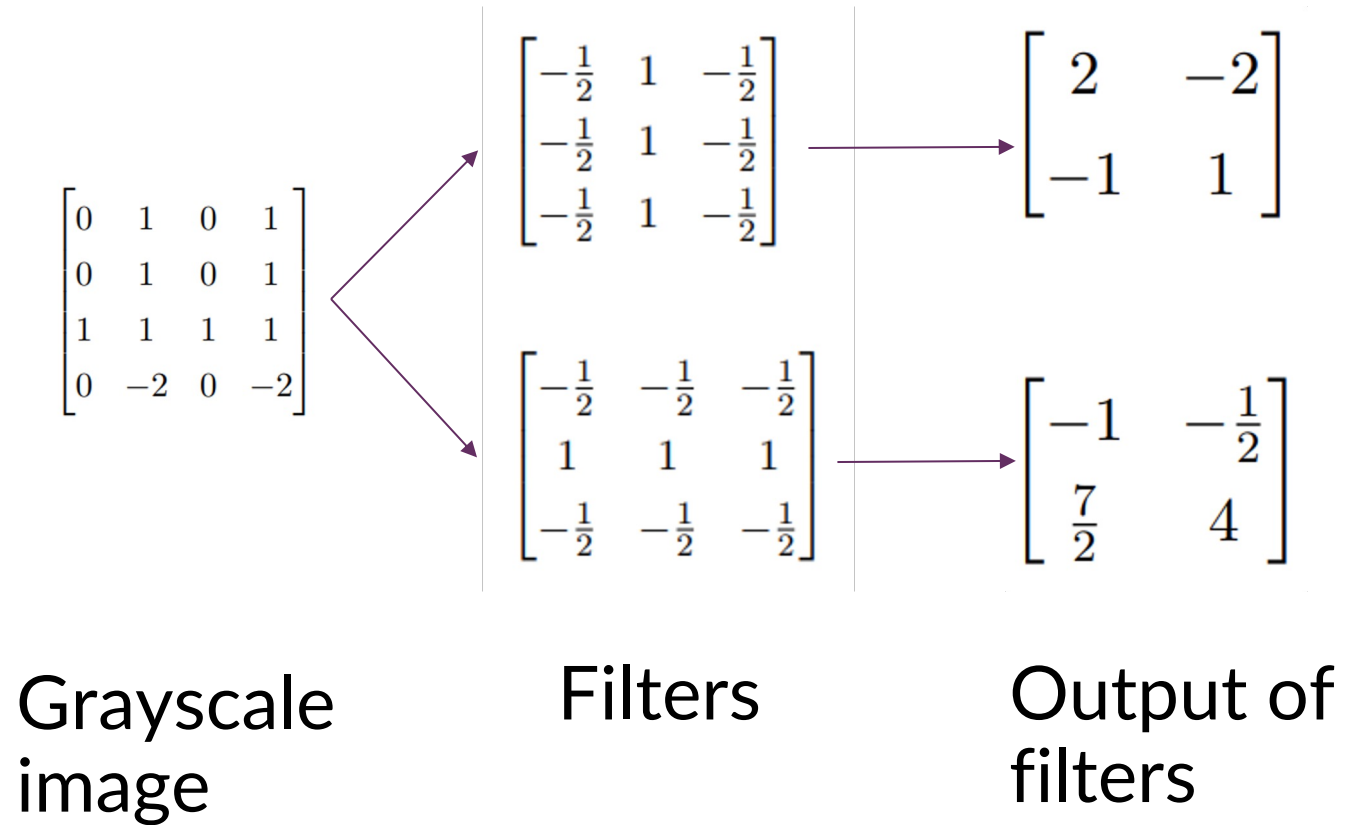
$$\begin{bmatrix} -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 & -\frac{1}{2} \end{bmatrix}$$



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

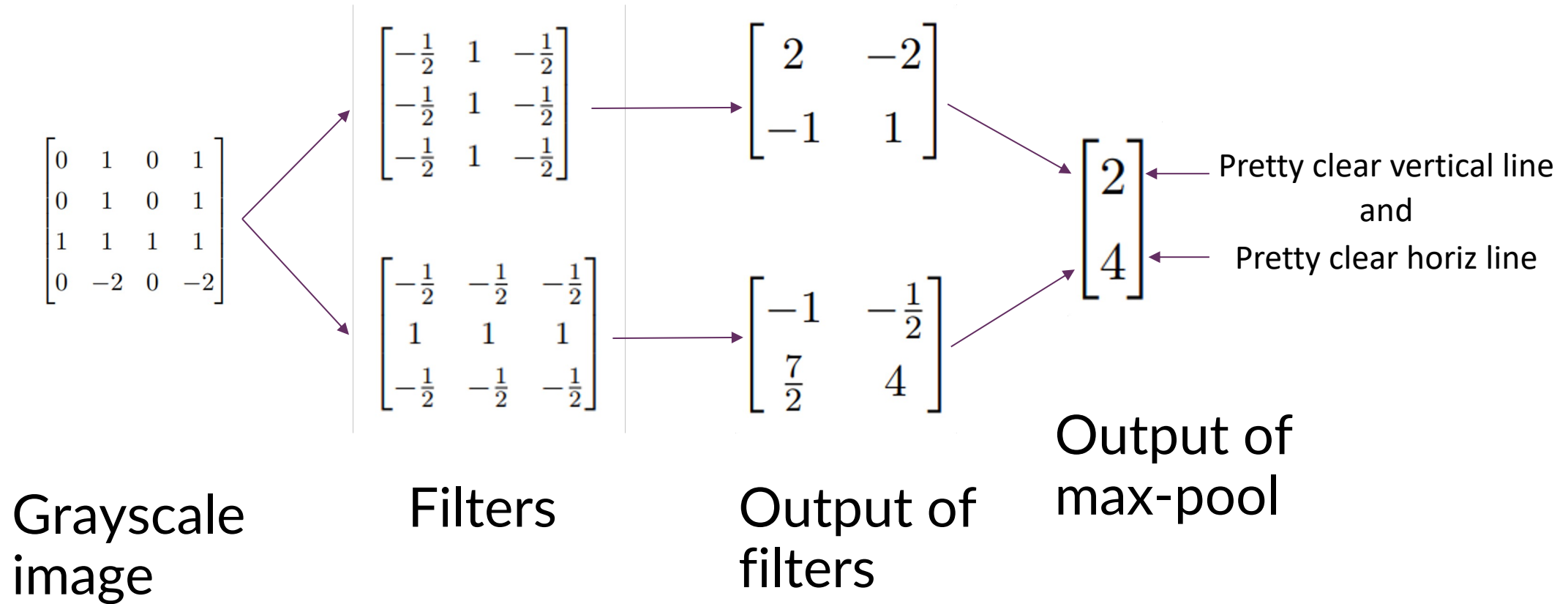


# Convolution Exercise Solution



# Convolutional Exercise Solution

Next, what happens if we run max-pooling on the filter outputs?

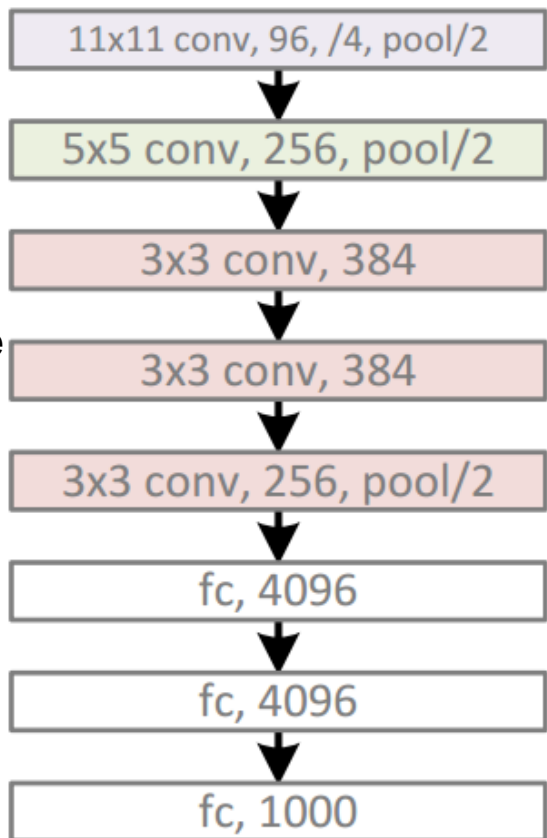


# Example architectures

# Example architectures

AlexNet, 8 layers  
(ILSVRC 2012)

~60M params



“Standard” scheme

[Conv-ReLU-pool?]

[Conv-ReLU-pool?]

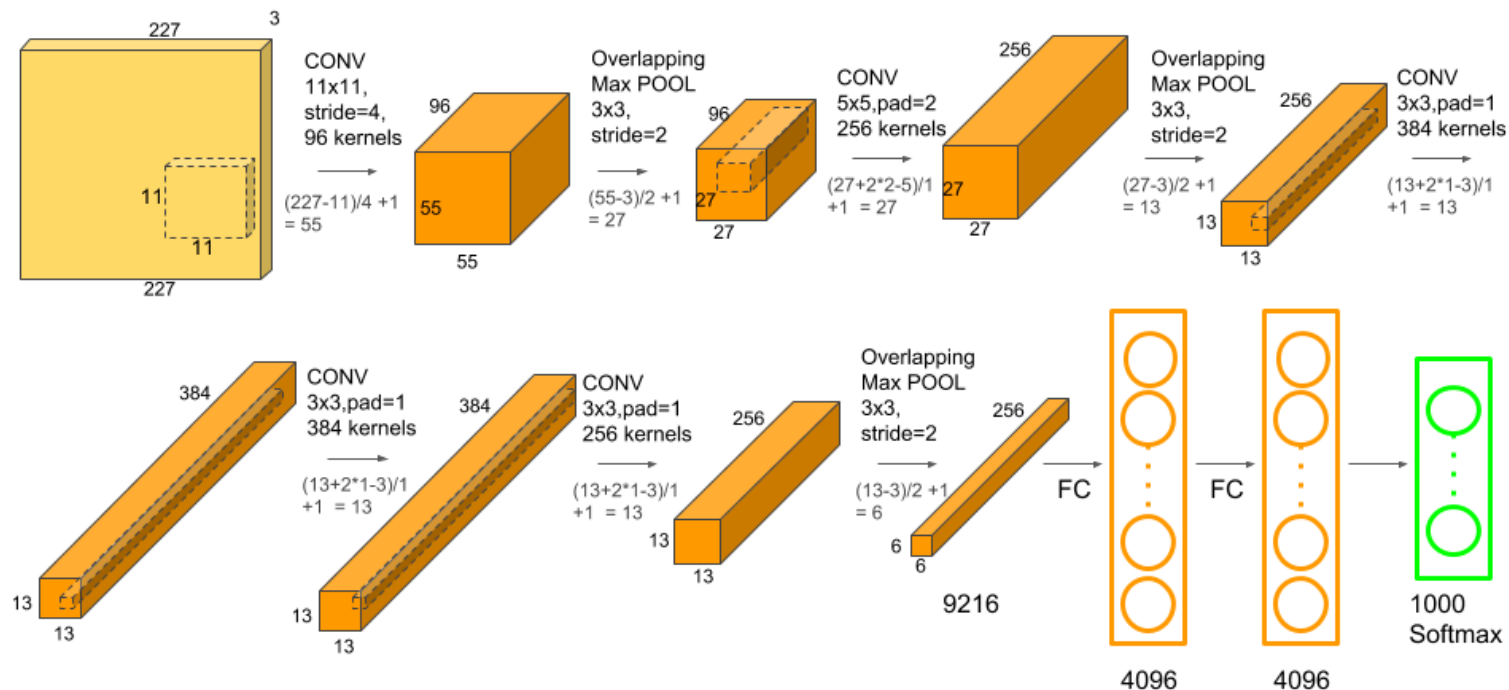
[Conv-ReLU-pool?]

...

Fully connected

...

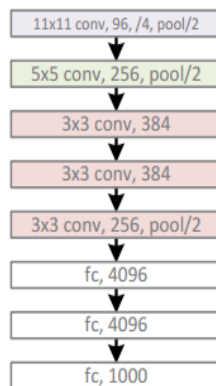
Fully connected



# Example architectures

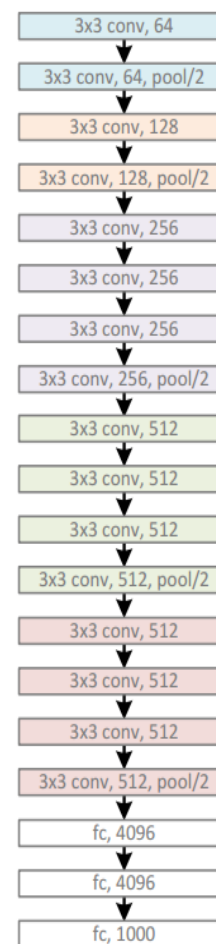
AlexNet, 8 layers  
(ILSVRC 2012)

~60M params



VGG, 19 layers  
(ILSVRC 2014)

~140M params



**“Standard” scheme**

[Conv-ReLU-pool?]  
[Conv-ReLU-pool?]  
[Conv-ReLU-pool?]

...

Fully connected

...

Fully connected



# Example architectures

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



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



ResNet, 152 layers  
(ILSVRC 2015)

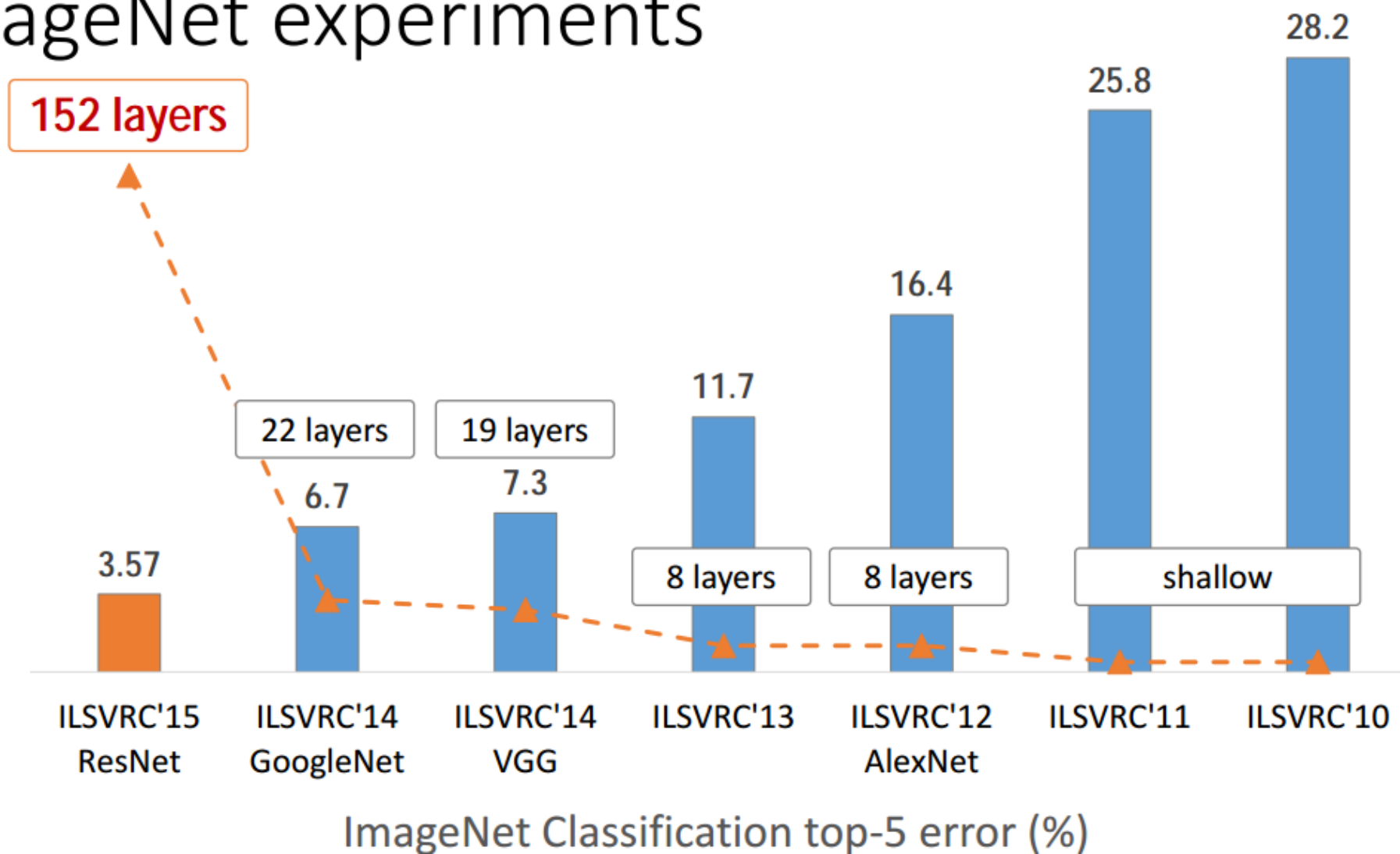
Less computation in  
forward pass than  
VGGNet!

Back to 60M params

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



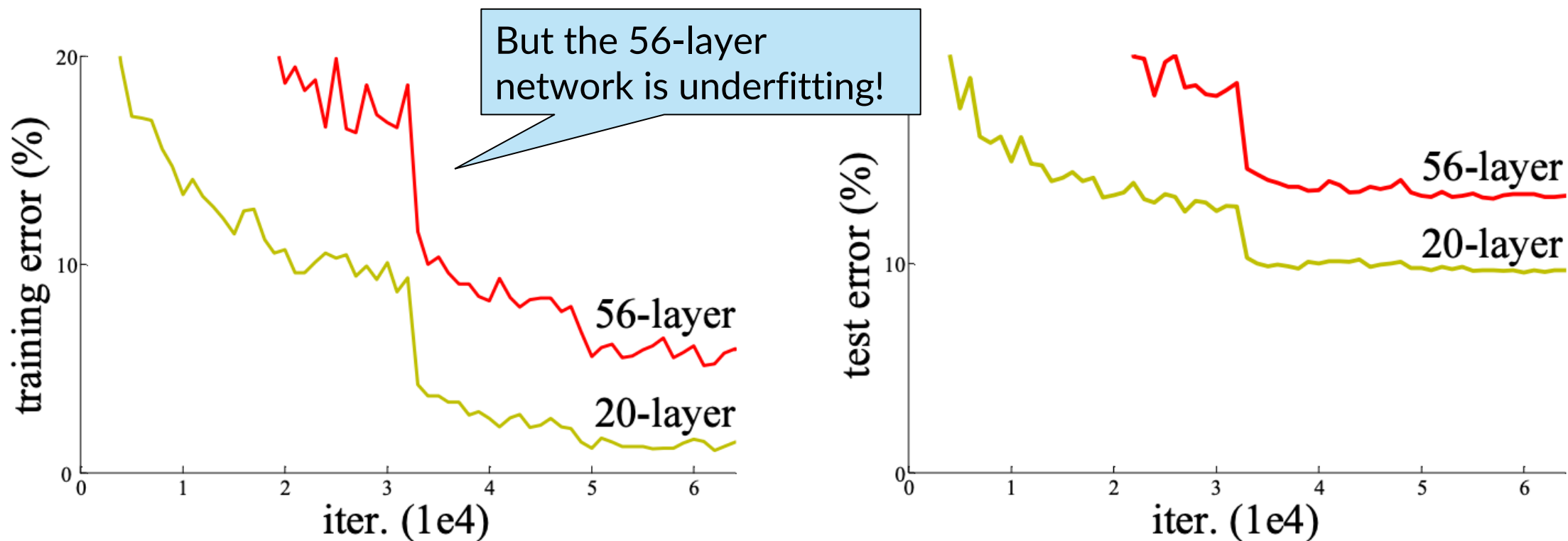
# ImageNet experiments





# Residual Network

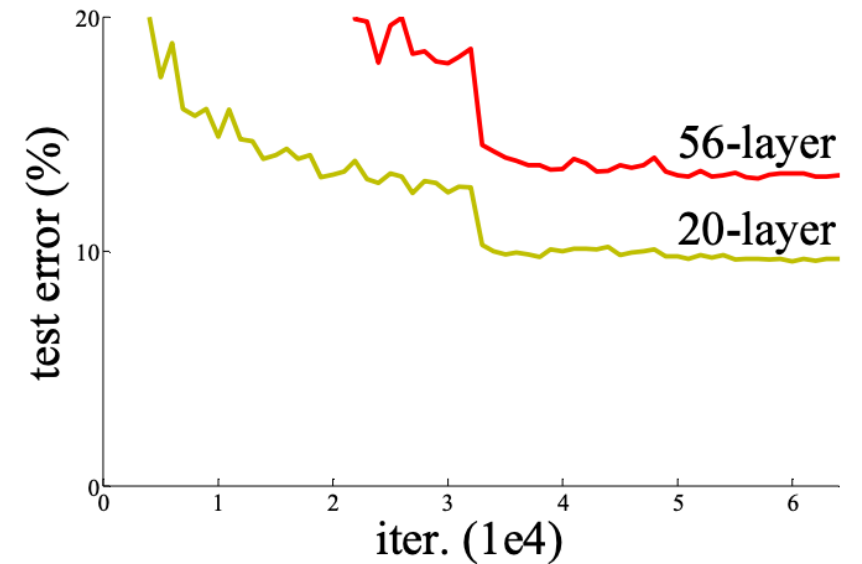
- Q: Why are deeper networks not always better?
- Hypothesis 1: Because of overfitting.



# Residual Network

- Q: Why are deeper networks not always better?
- Hypothesis 2: Because of optimization issues with deeper networks.

Idea: *Skip connections* that facilitate more direct feedback from the loss to the weights.

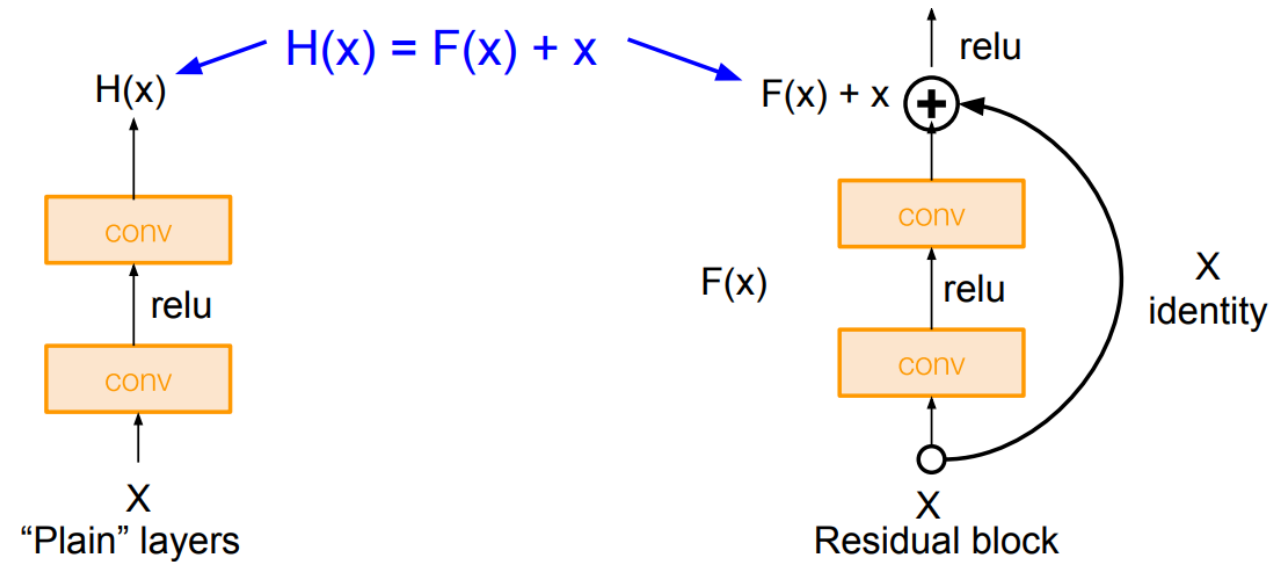


# Residual Network

Two views of residual connections:

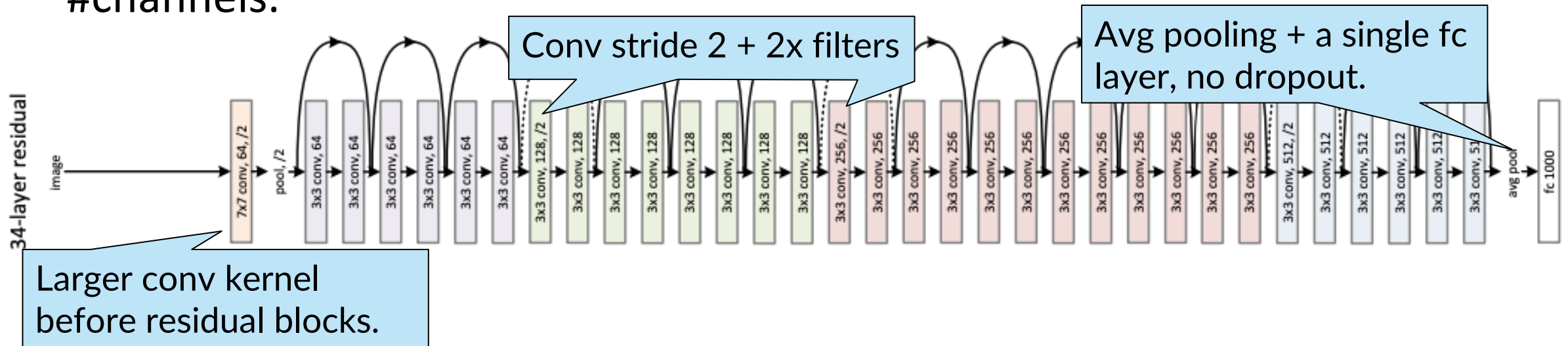
1. Providing shortcuts to gradients on the backward pass.
2. Allowing each “residual block” to fit the residual error function (recall gradient boosting!)

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



# Residual Network

- Stack lots of residual blocks.
- Zero-padded stride-1 3x3 convolutions + no max-pooling  $\Rightarrow$  maintains feature map size to build very deep nets.
- Reduce dimensions through stride 2 once every  $K$  blocks, increase #channels.



# Residual block designs

- For deeper networks, improve efficiency through 1x1 convolutions.

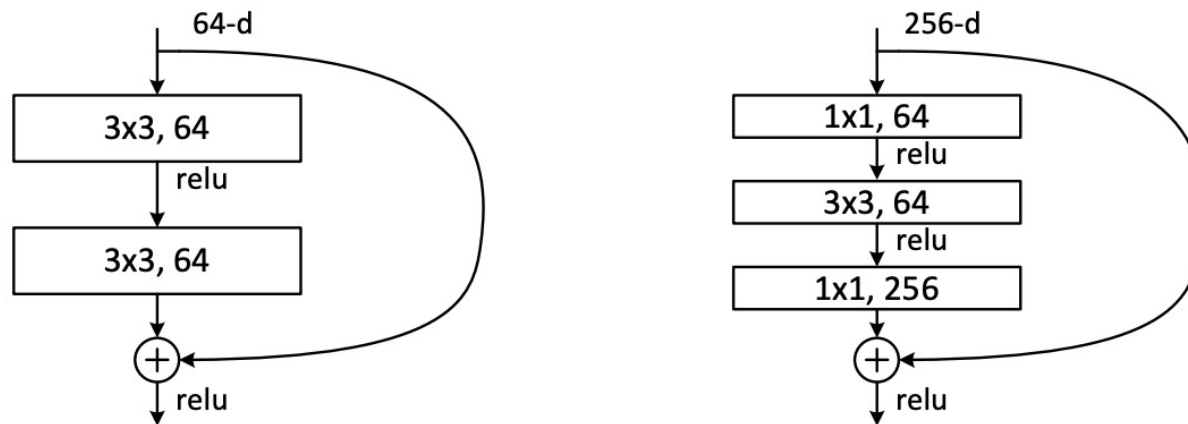


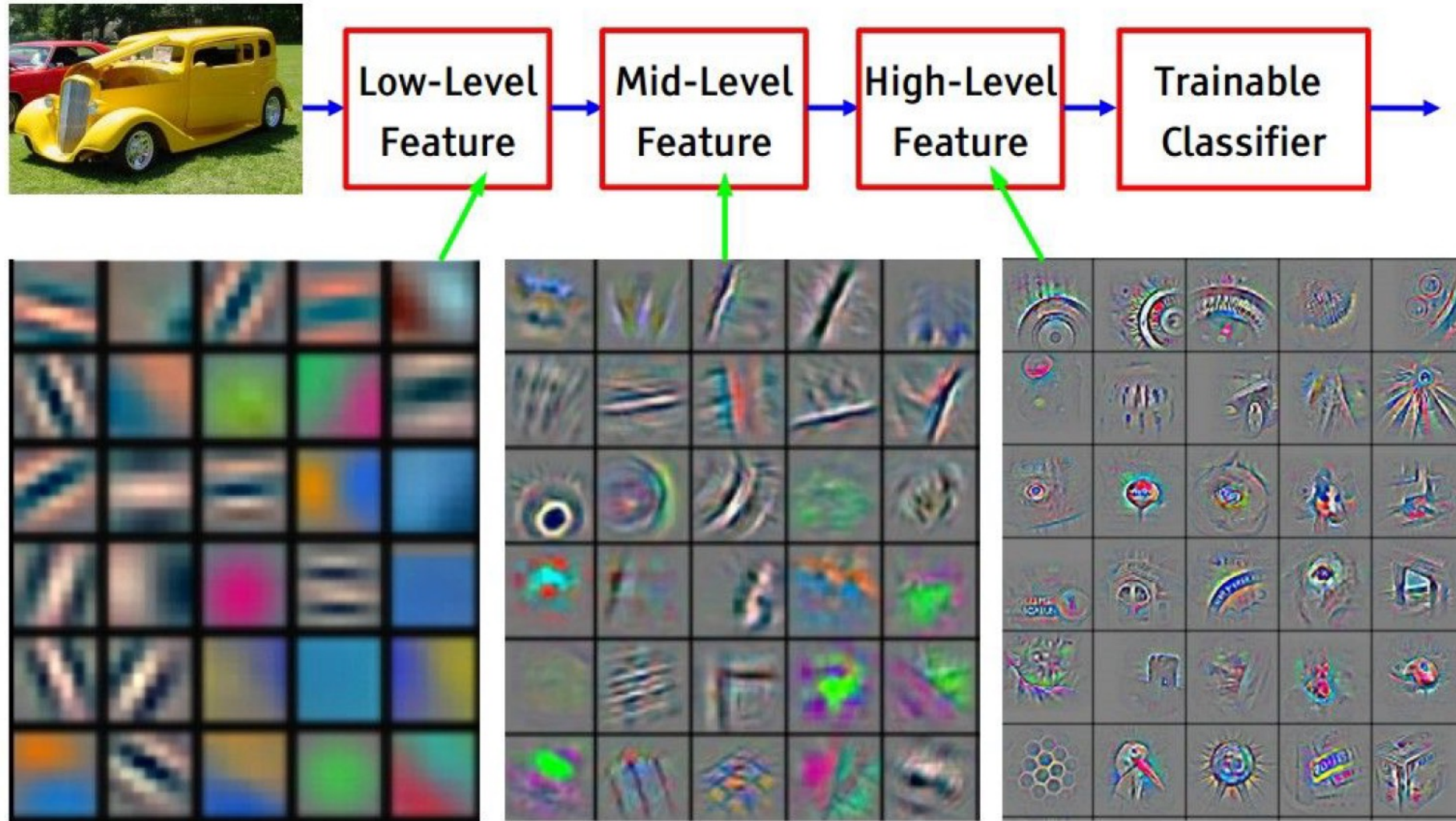
Figure 5. A deeper residual function  $\mathcal{F}$  for ImageNet. Left: a building block (on  $56 \times 56$  feature maps) as in Fig. 3 for ResNet-34. Right: a “bottleneck” building block for ResNet-50/101/152.

Many other improvements since 2015! E.g. “ResNeXt”, “Identity Mappings”, “ConvNeXt” etc.

# What do CNNs learn?

Visualizing and Understanding CNNs

# Feature visualization



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



# Layer 1



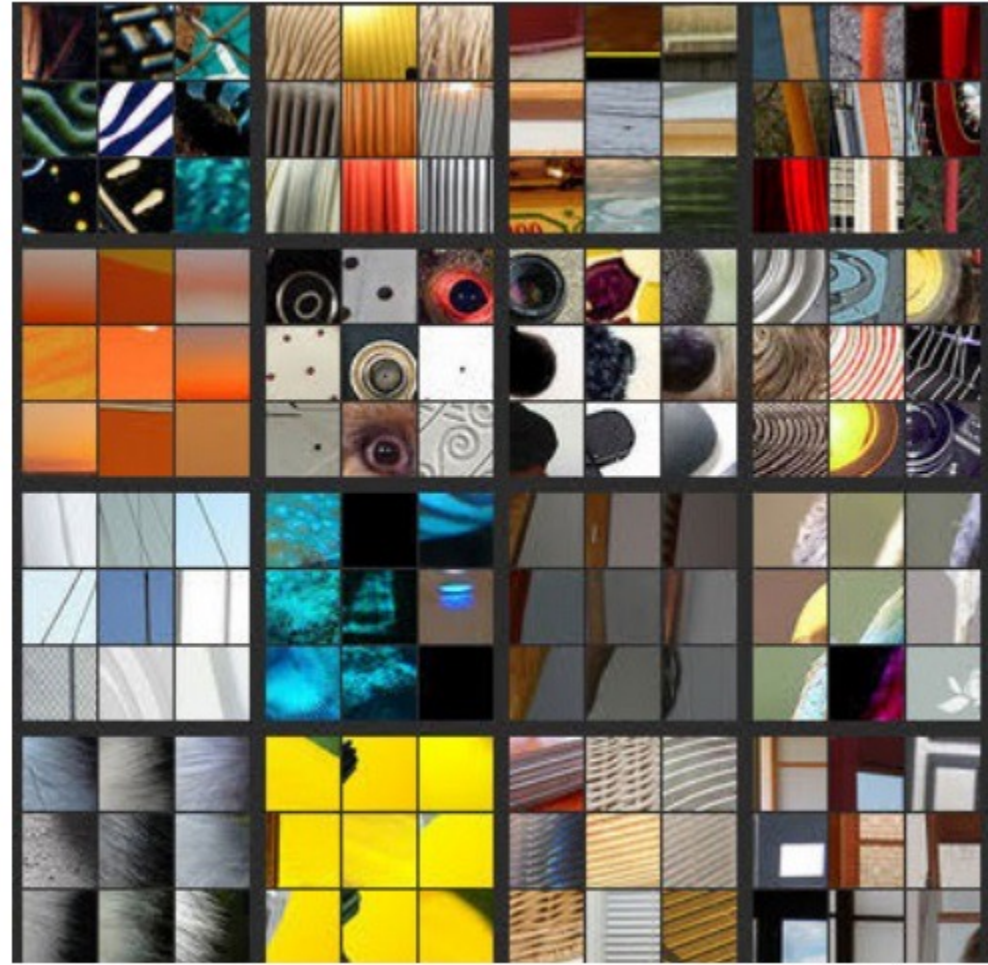
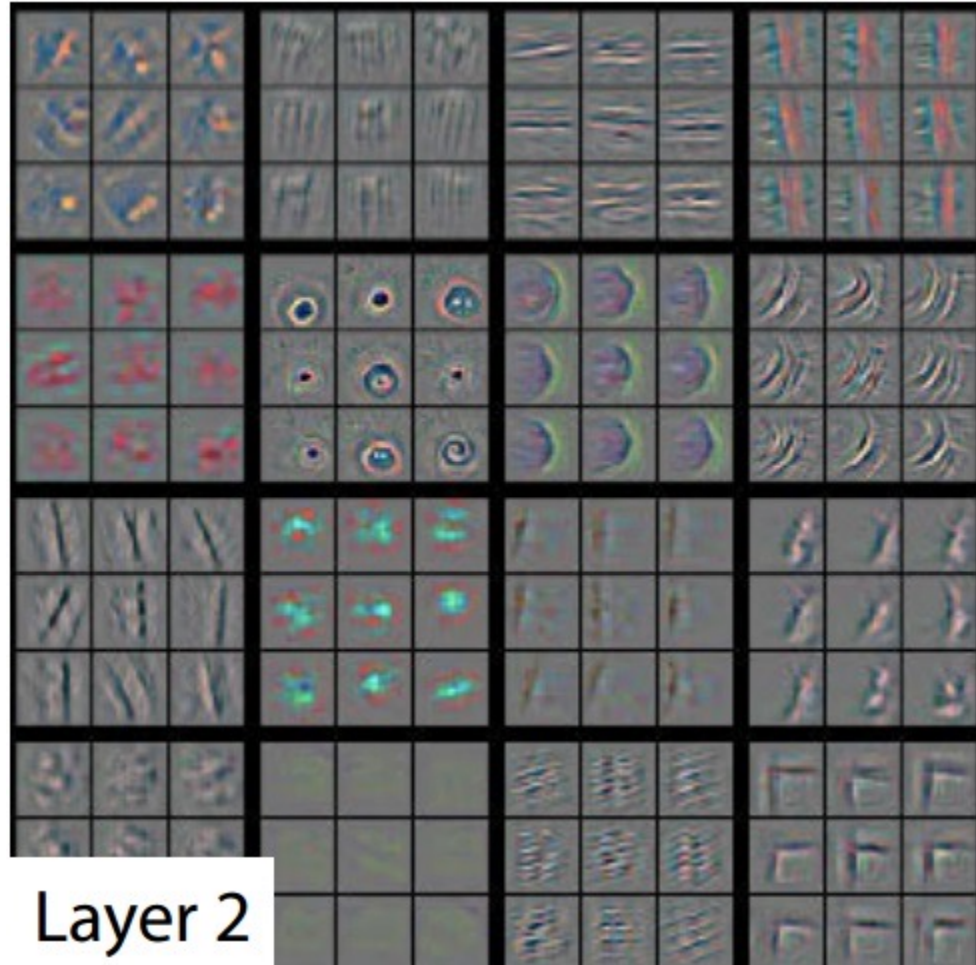
## Layer 1



Slide credit: Jia-Bin Huang



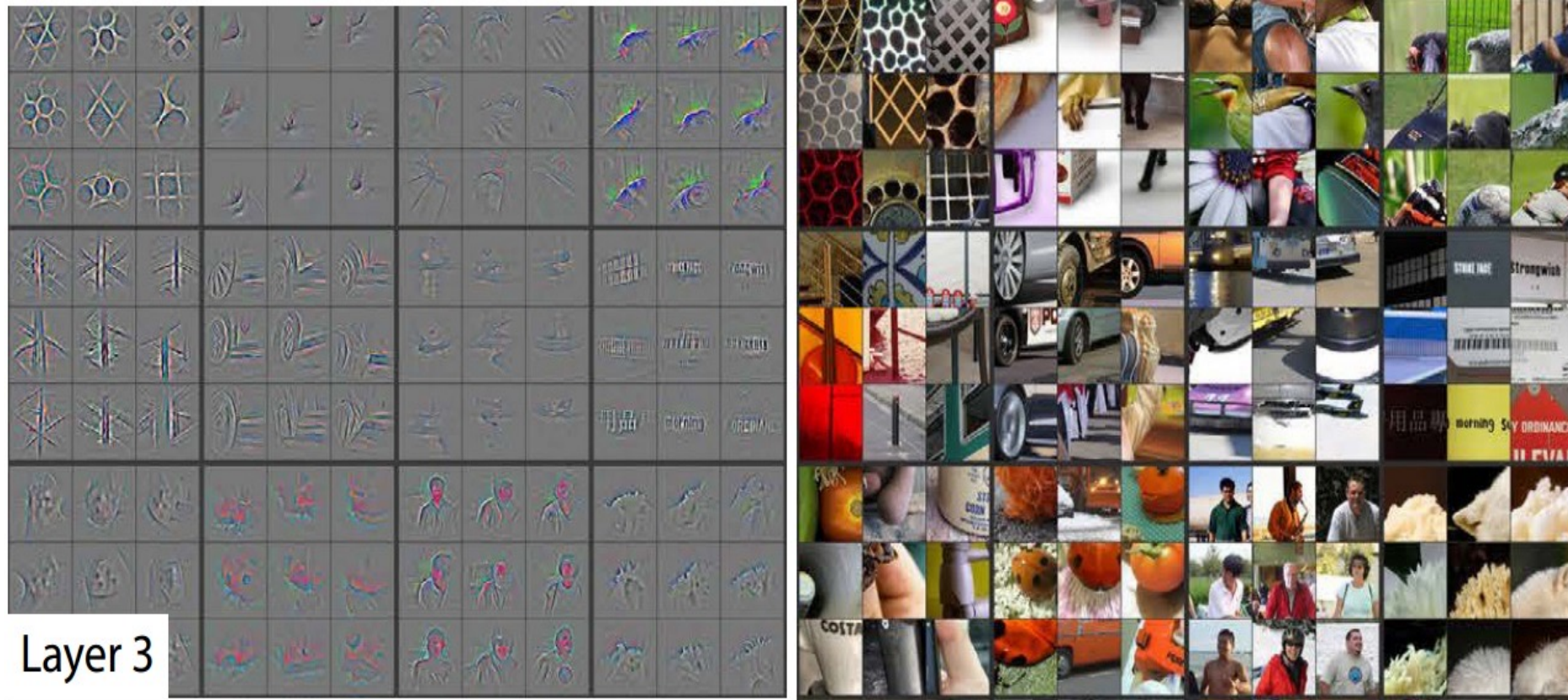
# Layer 2



Slide credit: Jia-Bin Huang

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

# Layer 3

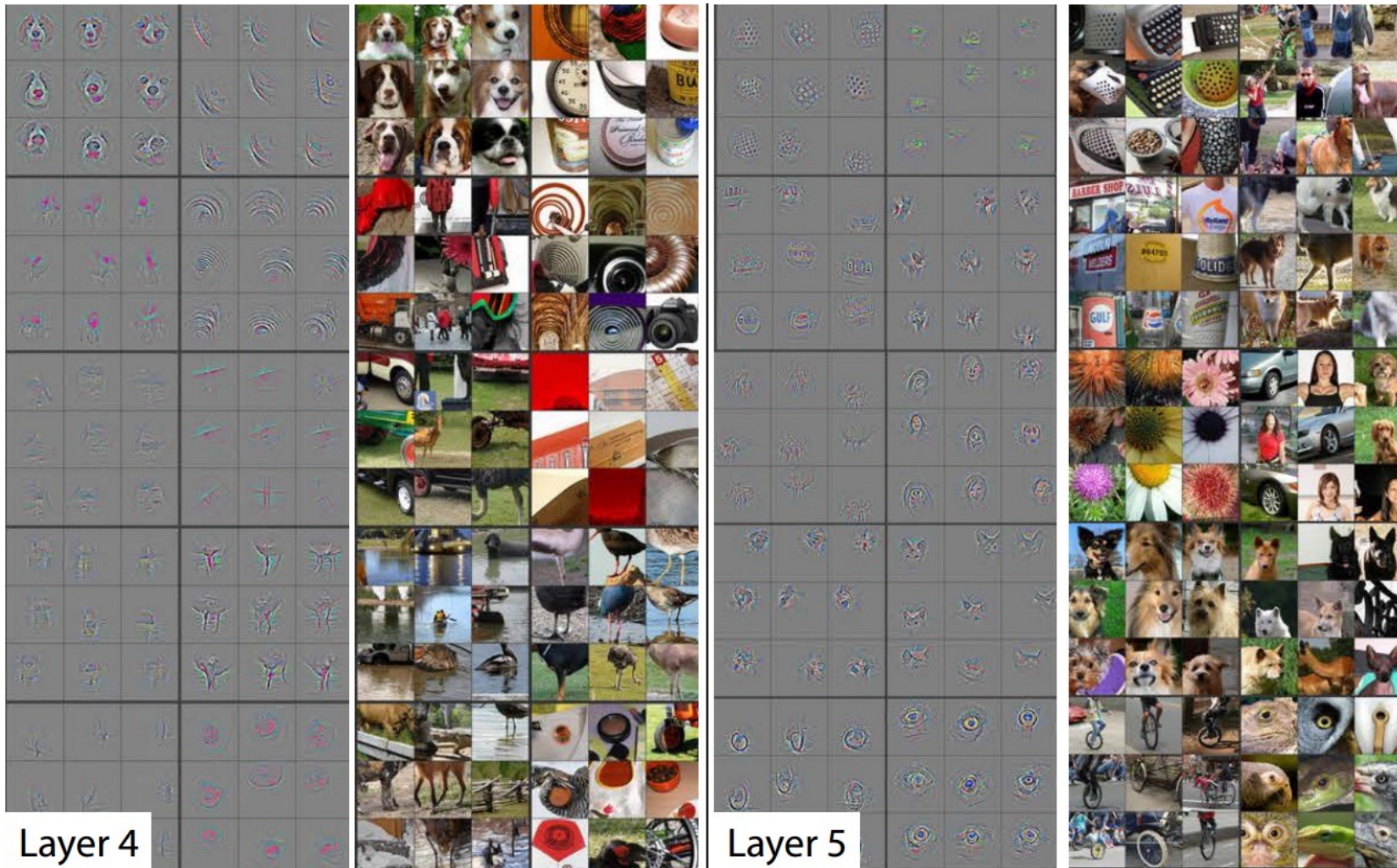


Slide credit: Jia-Bin Huang

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



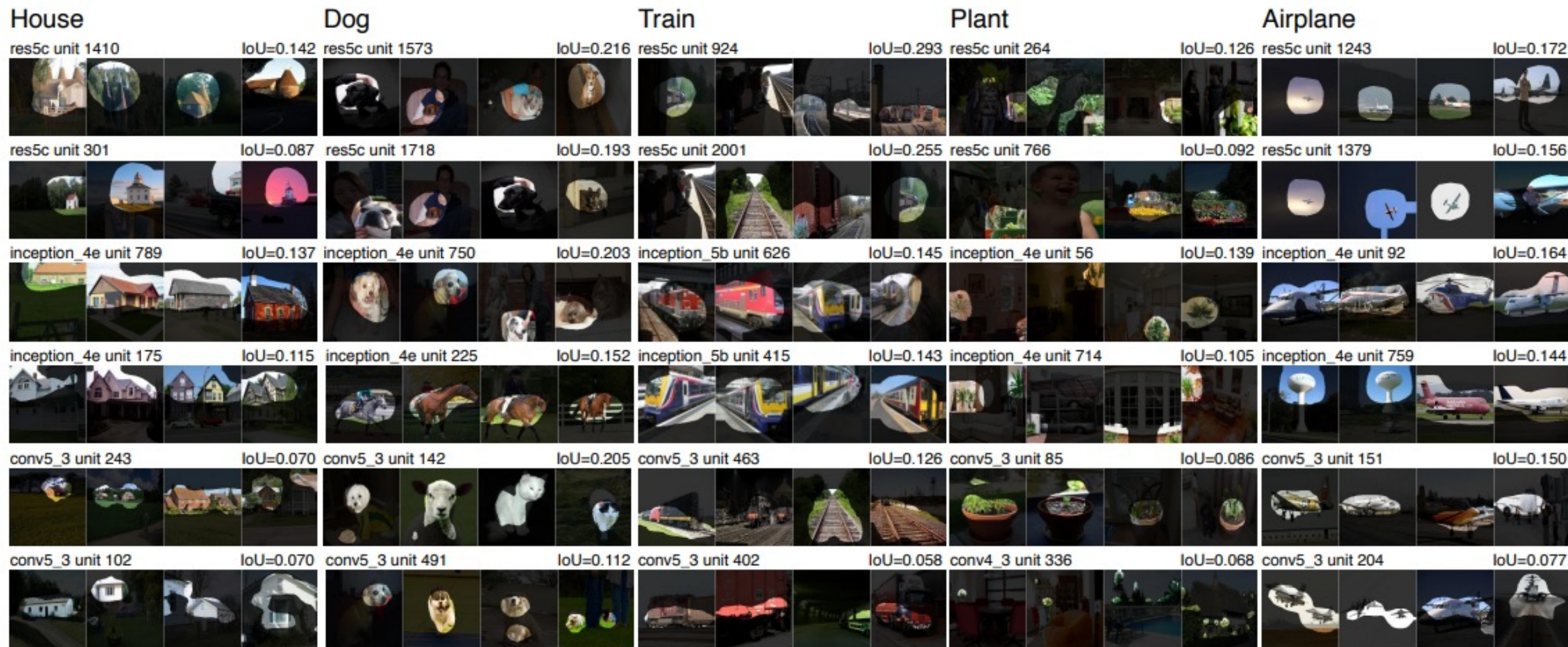
# Layer 4 and 5



Slide credit: Jia-Bin Huang



# Network dissection



CNNs with small datasets

# Can we reuse trained concepts?

Since CNN's trained for ImageNet object category classification appear to learn many apparently general features, why not reuse these models in some way to perform new tasks?

# Transfer learning with CNNs

What if your task doesn't have Imagenet-sized data?

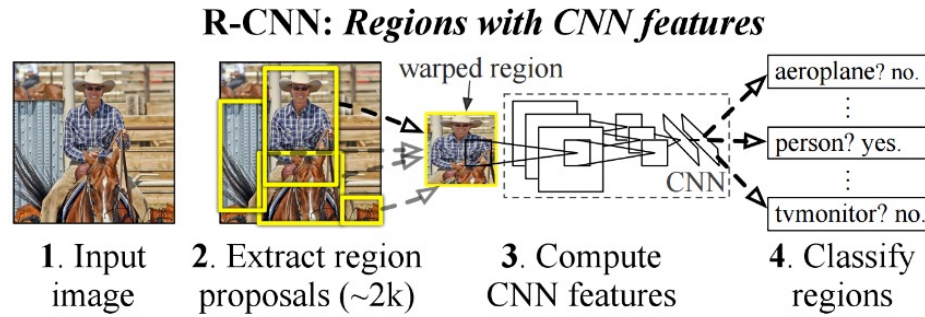
For tasks close to original task, can make do with small datasets + feature extraction or shallow finetuning.

For tasks far from original task, you will need to use moderate-sized datasets + deeper finetuning

Slide credit: Fei-Fei Li and Andrej Karpathy

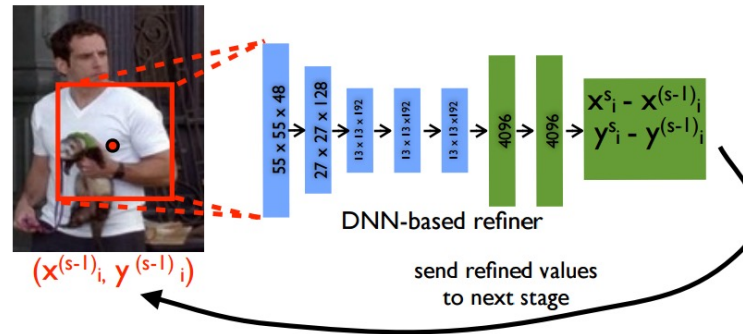
# Some sample applications

[Girshick et al. CVPR14]



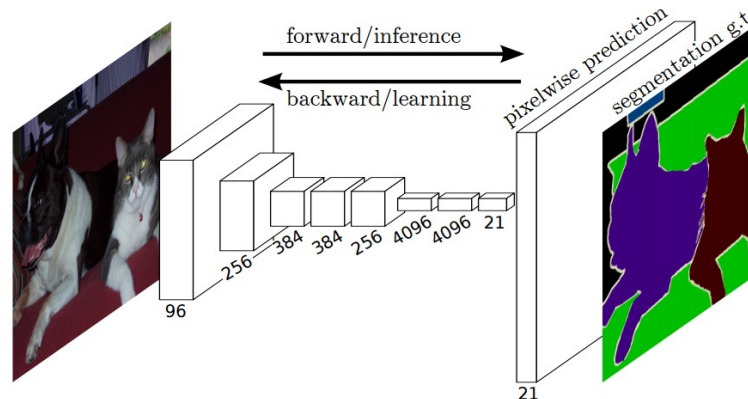
Object detection

[Toshev et al. CVPR14]



Pose detection (regression)

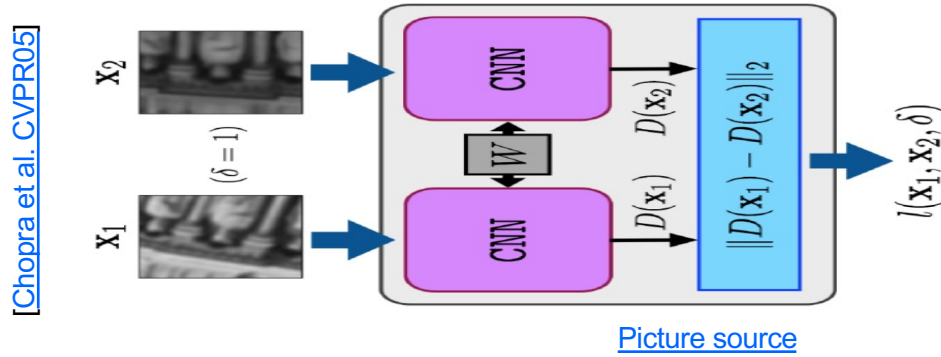
[Long et al. CVPR15]



Semantic segmentation



# Some sample applications



Similarity metric learning

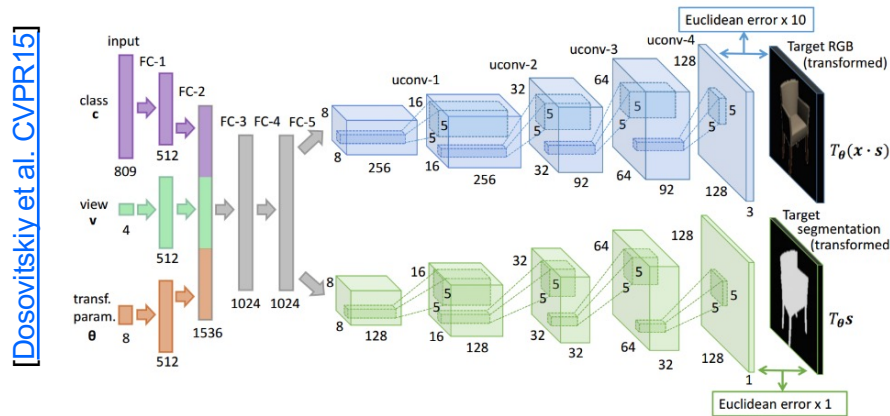
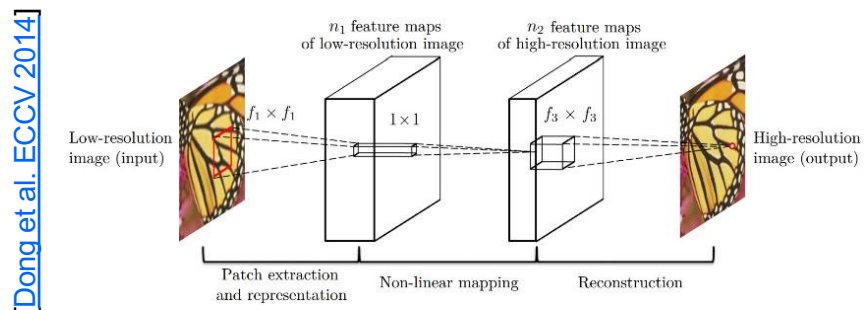


Image generation

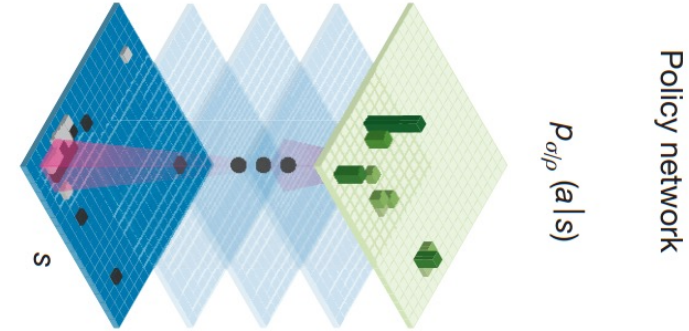


Low-level image processing:  
(superresolution, deblurring,  
image quality etc.)

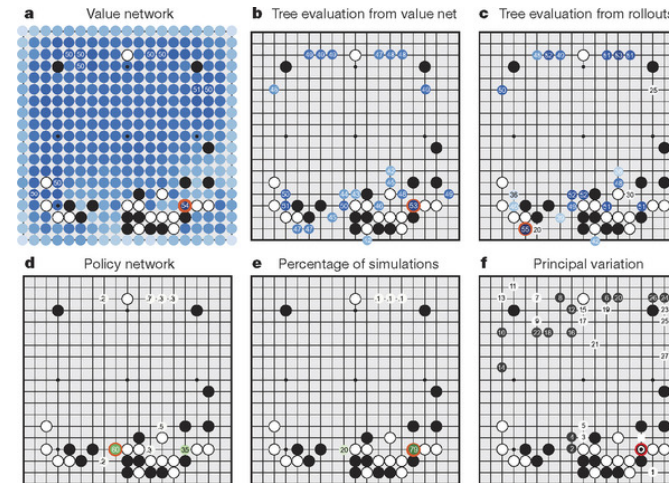
# Game playing!

CNN + Reinforcement learning

[Mnih et al., Nature '15]



[Silver et al., Nature '16]



# ConvNet Art!



See if you can tell  
artists' 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>

# Pytorch Training Loop



# Pytorch Training Loop

```
22 def train(args, model, device, train_loader, optimizer, epoch):
23     model.train()
24     for batch_idx, (data, target) in enumerate(train_loader):
25         data, target = data.to(device), target.to(device)
26         optimizer.zero_grad()
27         output = model(data)
28         loss = F.nll_loss(output, target)
29         loss.backward()
30         optimizer.step()
31         if batch_idx % args.log_interval == 0:
32             print('Train Epoch: {} [{}/{} ({:.0f}%)]\tLoss: {:.6f}'.format(
33                 epoch, batch_idx * len(data), len(train_loader.dataset),
34                 100. * batch_idx / len(train_loader), loss.item()))
```

Looping over mini-batches

Zero out all old gradients

Runs forward pass `model.forward(data)`

Loss computation

Backpropagation

Gradient step

# Pytorch Training Loop

```
83 def main():
84     torch.manual_seed(1)
85     device = torch.device("cuda")
86     train_loader = torch.utils.data.DataLoader( Load dataset
87         datasets.MNIST('../data', train=True, download=True,
88             transform=transforms.Compose([
89                 transforms.ToTensor(),
90                 transforms.Normalize((0.1307,), (0.3081,))
91             ])),
92         batch_size=64, shuffle=True)
93
94     model = Net().to(device)
95     optimizer = optim.Adam(model.parameters(), lr=1e-4)
96     scheduler = optim.lr_scheduler.StepLR(optimizer, step_size=1, gamma=0.9)
97     Loop over epochs (full passes over data)
98     for epoch in range(1, 15):
99         Minibatch SGD for one epoch
100         train(model, device, train_loader, optimizer, epoch)
101         scheduler.step()
102         Update base learning rate
```

# Pytorch Model

- To use your model (once it has been trained):

```
model.eval()      # puts model in evaluation mode  
label = model(input)  # forward pass to compute  
outputs
```