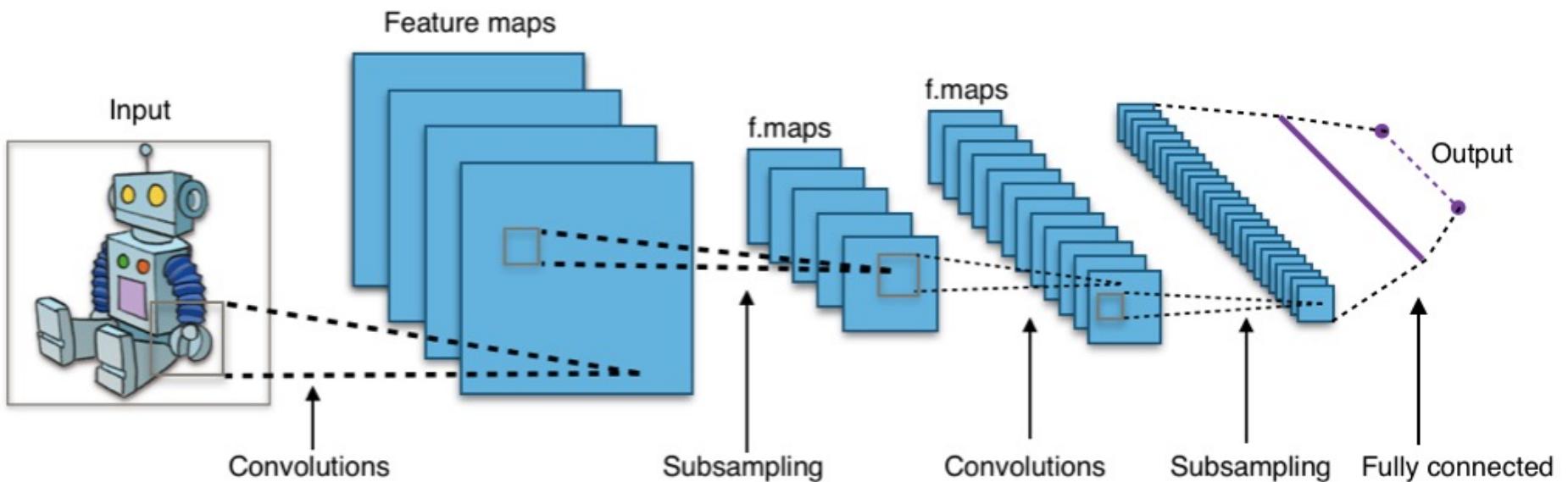


Convolutional neural networks

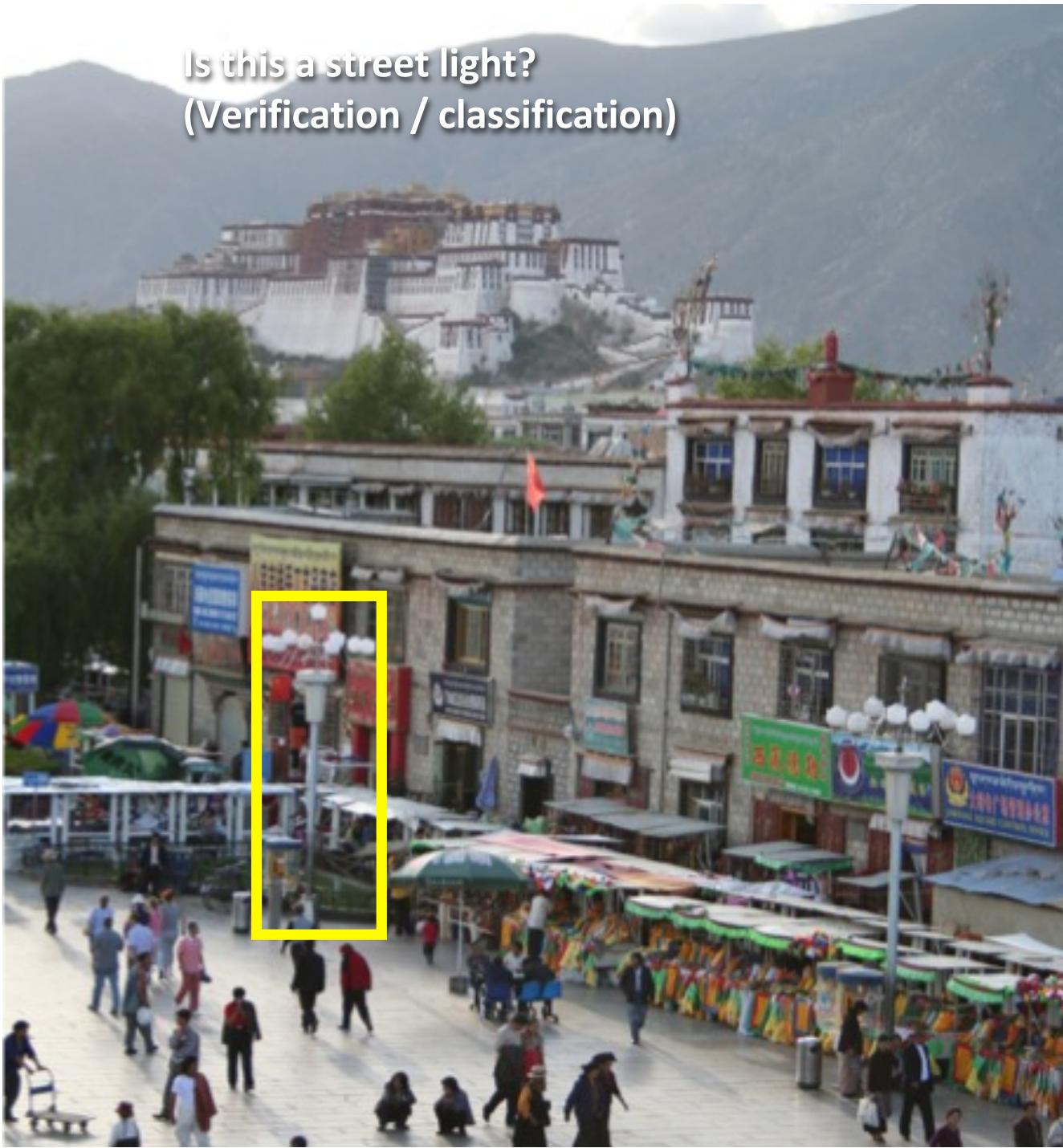


Overview of today's lecture

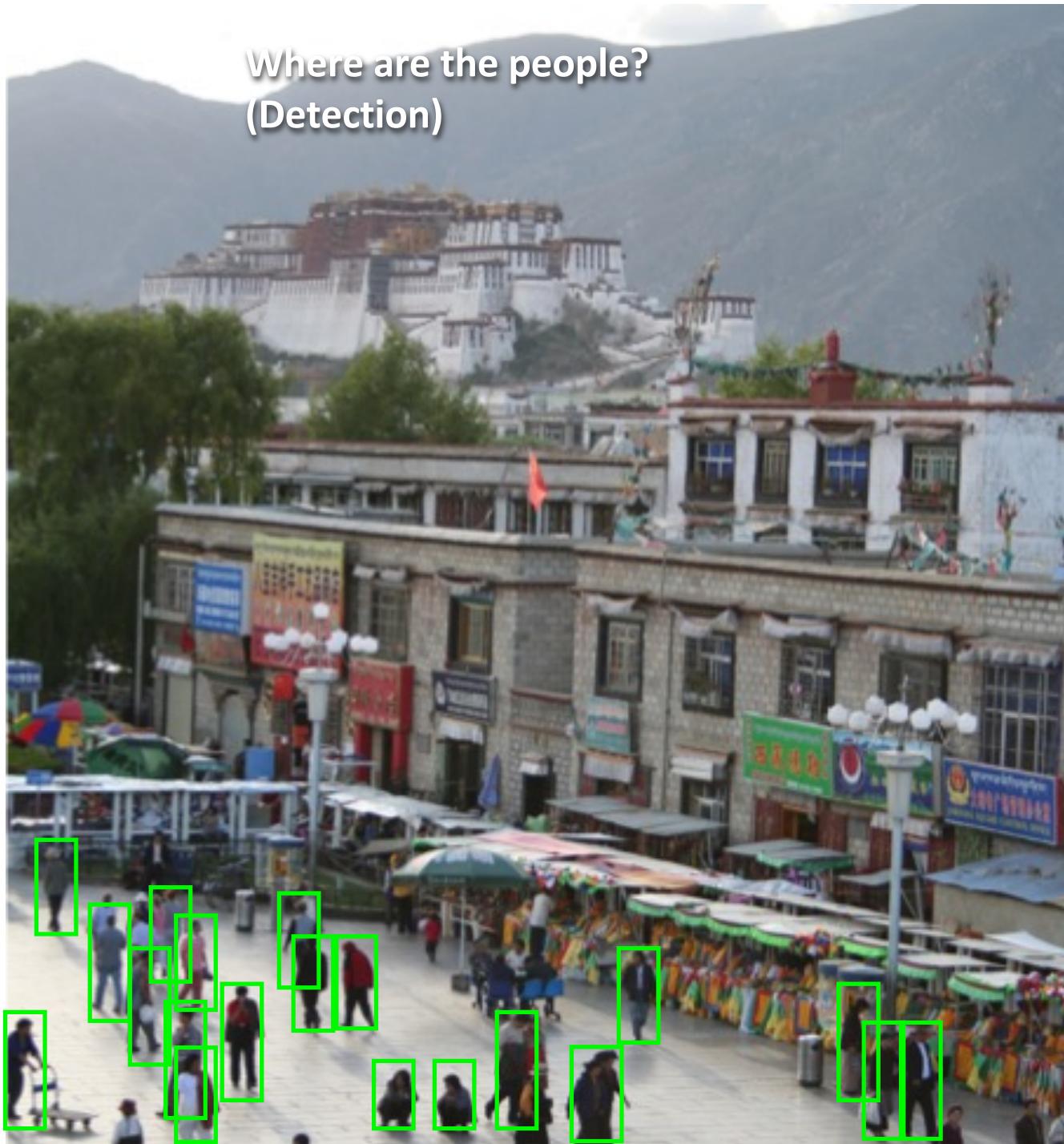
- What is recognition?
- Convolutional Neural Networks (CNNs / ConvNets)
- Training ConvNets.

Recognition

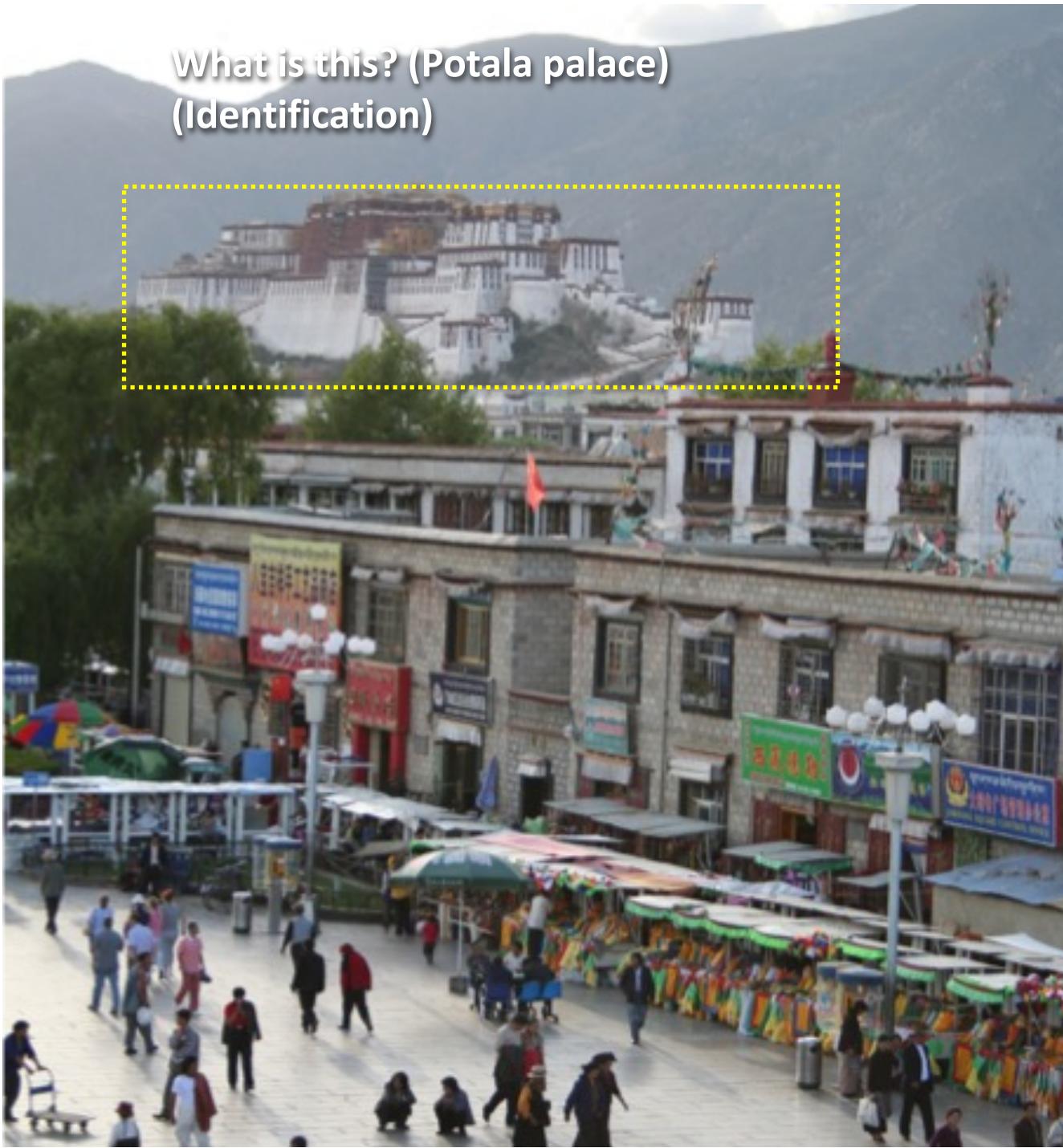
Is this a street light?
(Verification / classification)



Where are the people? (Detection)



What is this? (Potala palace)
(Identification)



Sky

What's in the scene? (semantic segmentation)

Mountain

Trees

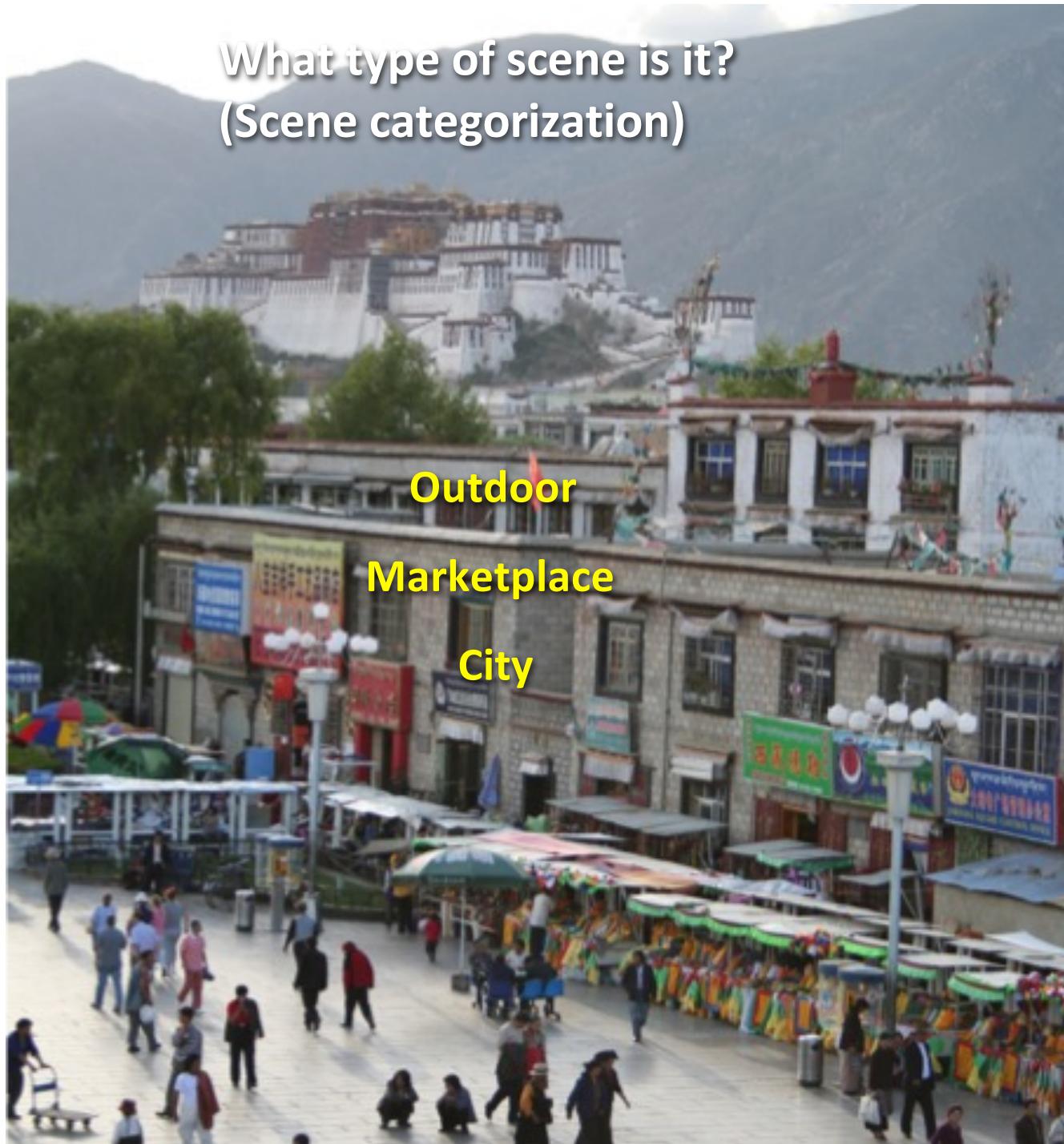
Building

Vendors

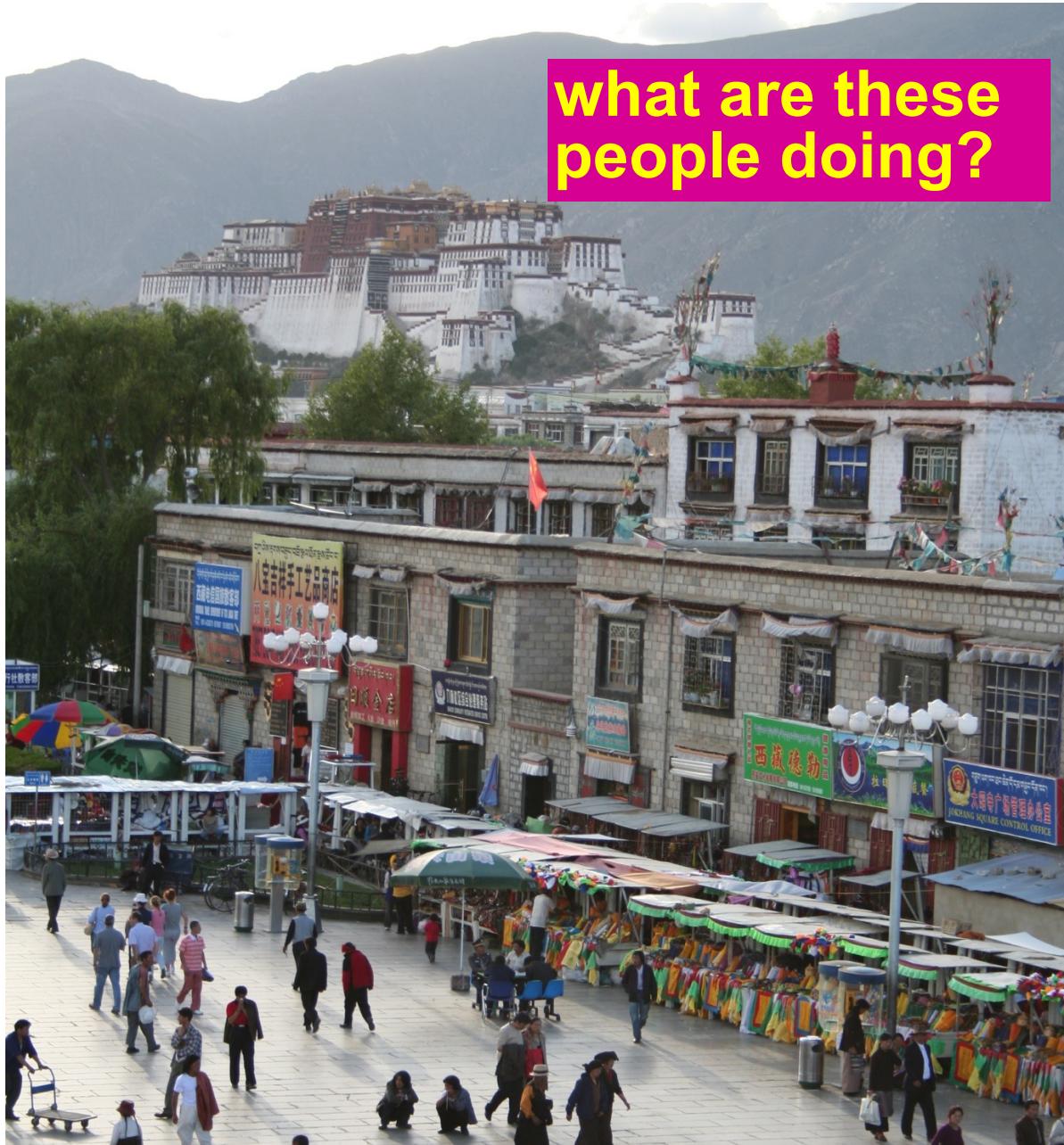
People

Ground

What type of scene is it? (Scene categorization)



Activity / Event Recognition



what are these
people doing?

Convolutional Neural Networks

Motivation

The screenshot shows the MIT Technology Review website's homepage for the year 2013. At the top, there is a navigation bar with links for HOME, MENU, CONNECT, THE LATEST, POPULAR, MOST SHARED, and a user icon. The main title "10 BREAKTHROUGH TECHNOLOGIES 2013" is prominently displayed, with "MIT Technology Review" to its left. Below the title, there are ten cards, each representing a breakthrough technology. The first card, "Deep Learning", is circled in red. The other cards are: "Temporary Social Media", "Prenatal DNA Sequencing", "Additive Manufacturing", "Baxter: The Blue-Collar Robot", "Memory Implants", "Smart Watches", "Ultra-Efficient Solar", "Big Data from", and "Supergrids". Each card has a brief description and a right-pointing arrow.

MIT Technology Review

10 BREAKTHROUGH TECHNOLOGIES 2013

Introduction The 10 Technologies Past Years

Deep Learning

With massive amounts of computational power, machines can now recognize objects and translate speech in real time. Artificial intelligence is finally getting smart.

Temporary Social Media

Messages that quickly self-destruct could enhance the privacy of online communications and make people freer to be spontaneous.

Prenatal DNA Sequencing

Reading the DNA of fetuses will be the next frontier of the genomic revolution. But do you really want to know about the genetic problems or musical aptitude of your unborn child?

Additive Manufacturing

Skeptical about 3-D printing? GE, the world's largest manufacturer, is on the verge of using the technology to make jet parts.

Baxter: The Blue-Collar Robot

Rodney Brooks's newest creation is easy to interact with, but the complex innovations behind the robot show just how hard it is to get along with people.

Memory Implants

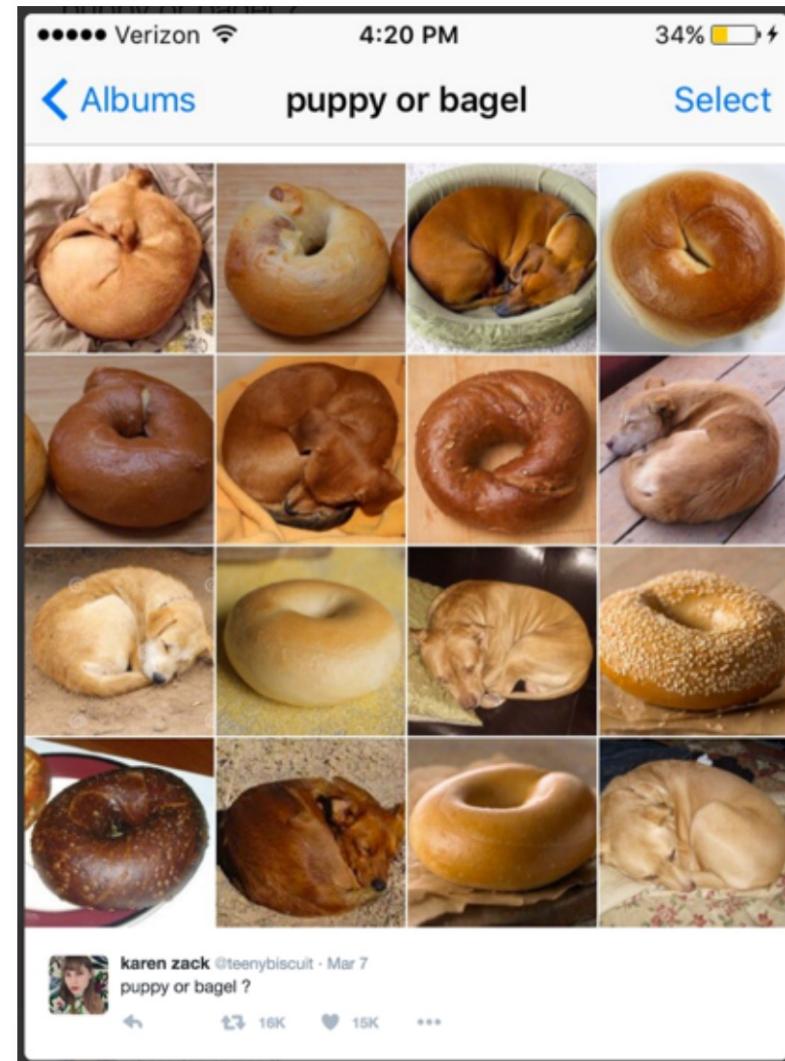
Smart Watches

Ultra-Efficient Solar

Big Data from

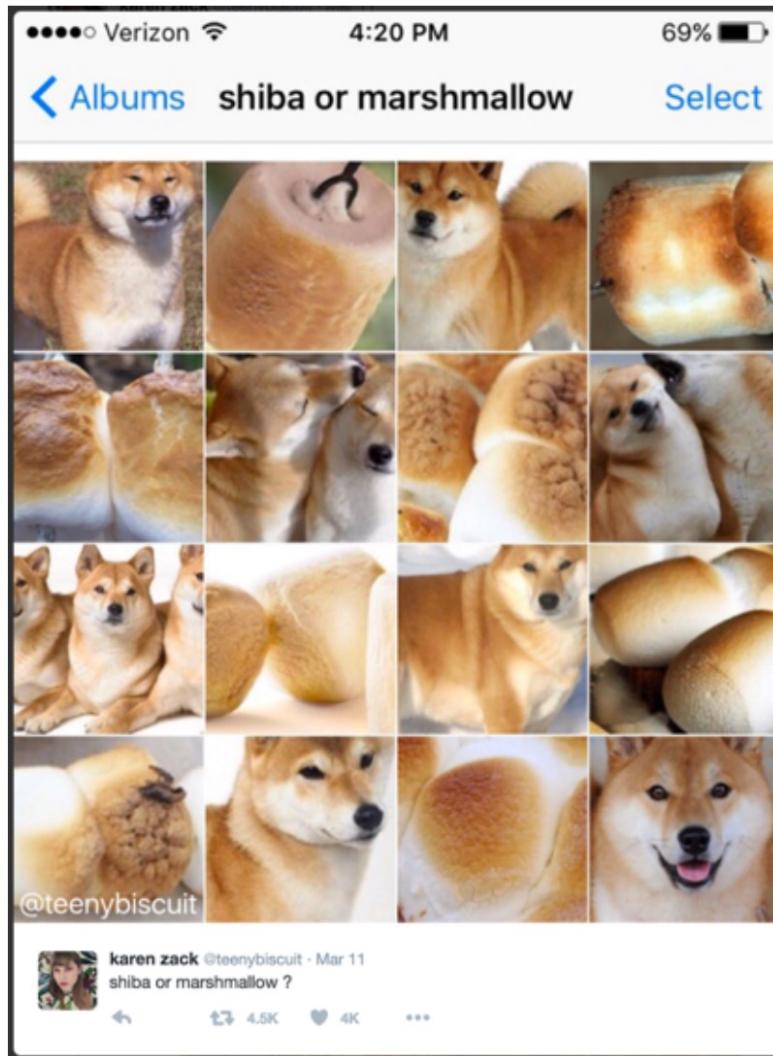
Supergrids

(Unrelated) Dog vs Food



[Karen Zack, @teenybiscuit]

(Unrelated) Dog vs Food



[Karen Zack, @teenybiscuit]

CNNs in 2012: “SuperVision” (aka “AlexNet”)

“AlexNet” — Won the ILSVRC2012 Challenge

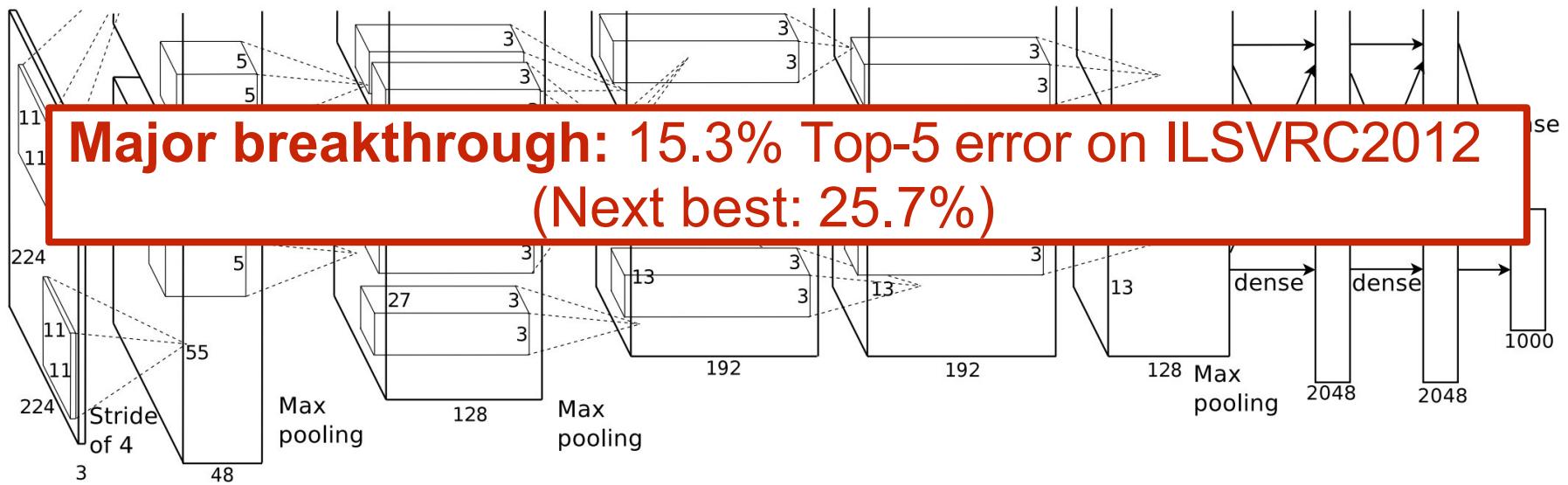
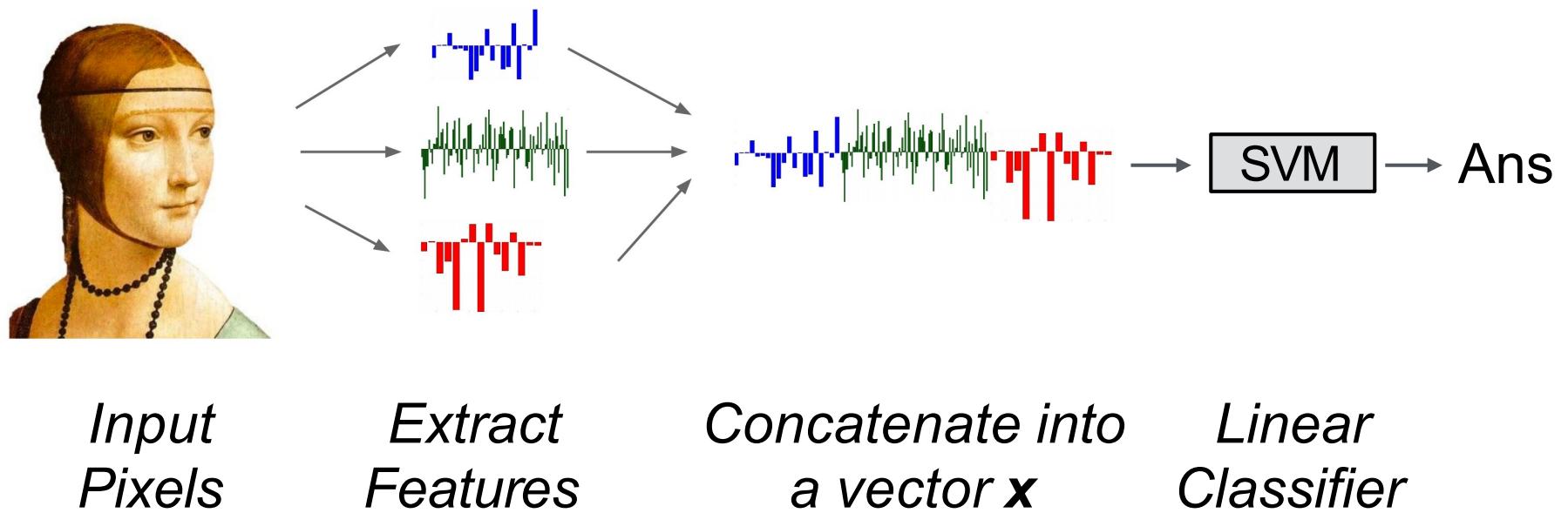


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network’s input is 150,528-dimensional, and the number of neurons in the network’s remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

[Krizhevsky, Sutskever, Hinton. NIPS 2012]

Recap: Before Deep Learning



The last layer of (most) CNNs are linear classifiers



*Input
Pixels*

*Perform everything with a big neural
network, trained end-to-end*

Key: perform enough processing so that by the time you get to the end of the network, the classes are linearly separable

ConvNets

They're just neural networks with
3D activations and weight sharing

What is an image?

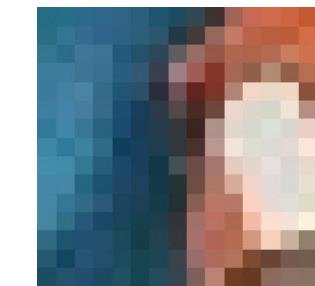


What is an image?



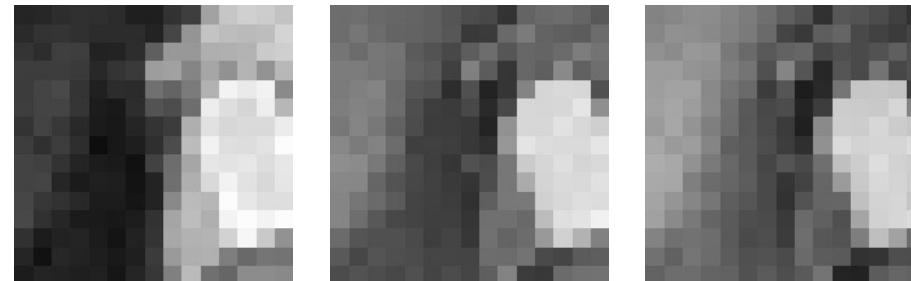
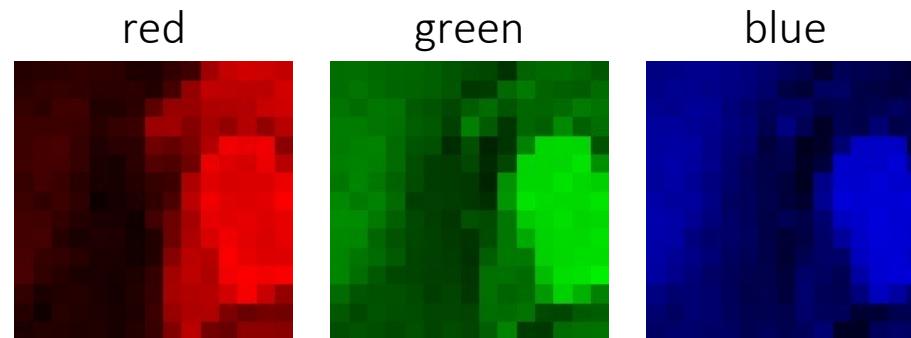
A (color) image
is a 3D tensor of
numbers.

What is an image?



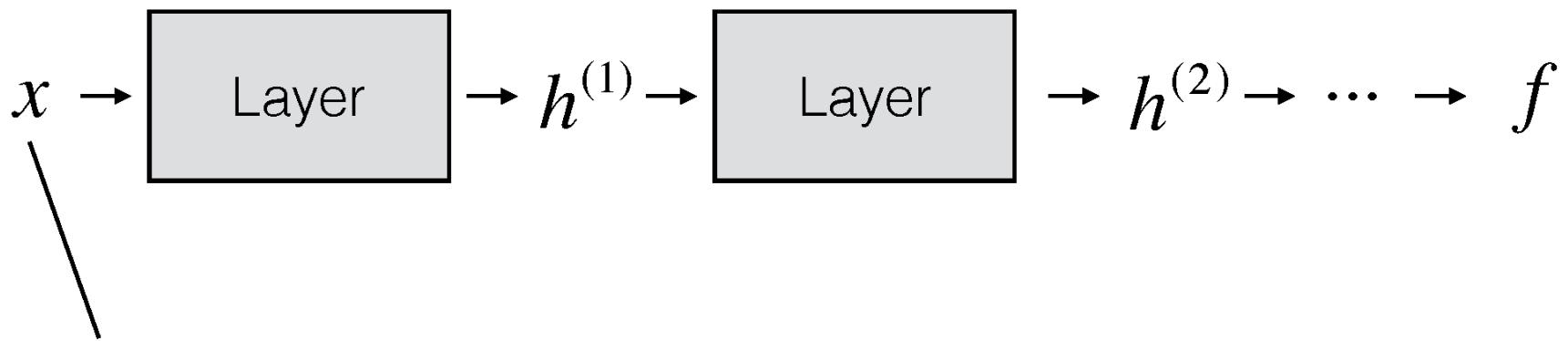
color image patch

How many bits are
the intensity values?



Each
channel
is a 2D array
of numbers.

What shape should the activations have?



- The input is an image, which is 3D (RGB channel, height, width)
- We could flatten it to a 1D vector, but then we lose structure
- What about keeping everything in 3D?

ConvNets

They're just neural networks with
3D activations and weight sharing

3D Activations

before:

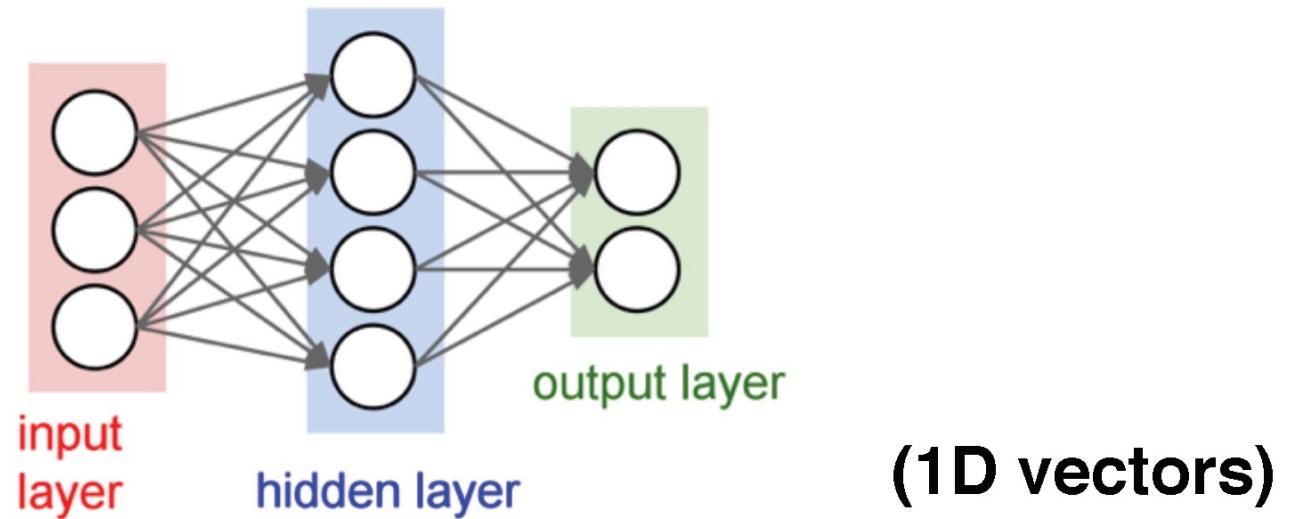
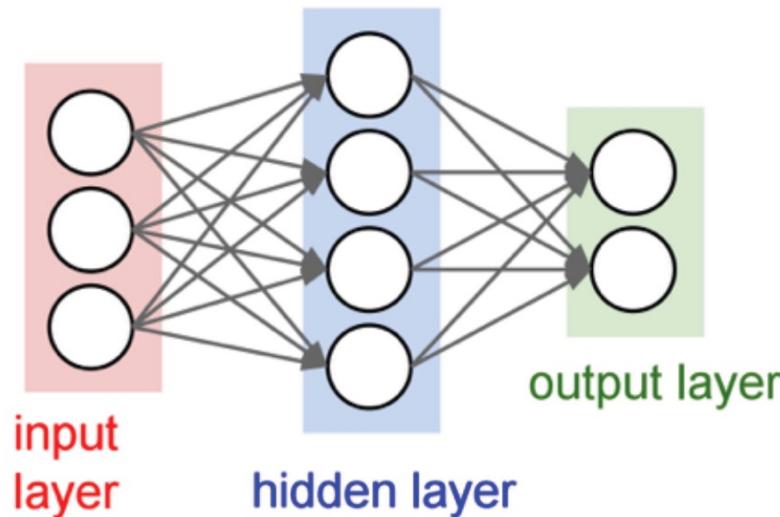


Figure: Andrej Karpathy

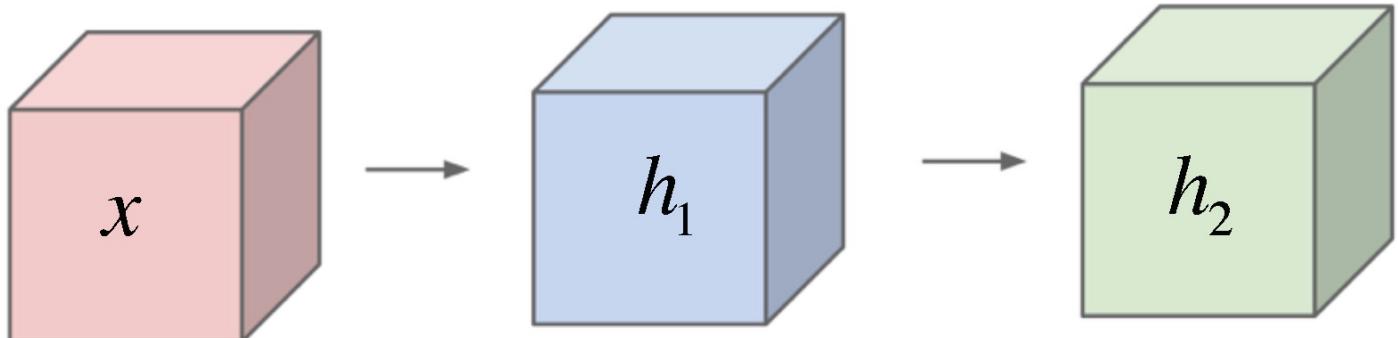
3D Activations

before:



(1D vectors)

now:



(3D arrays)

Figure: Andrej Karpathy

3D Activations

All Neural Net
activations
arranged in **3**
dimensions:

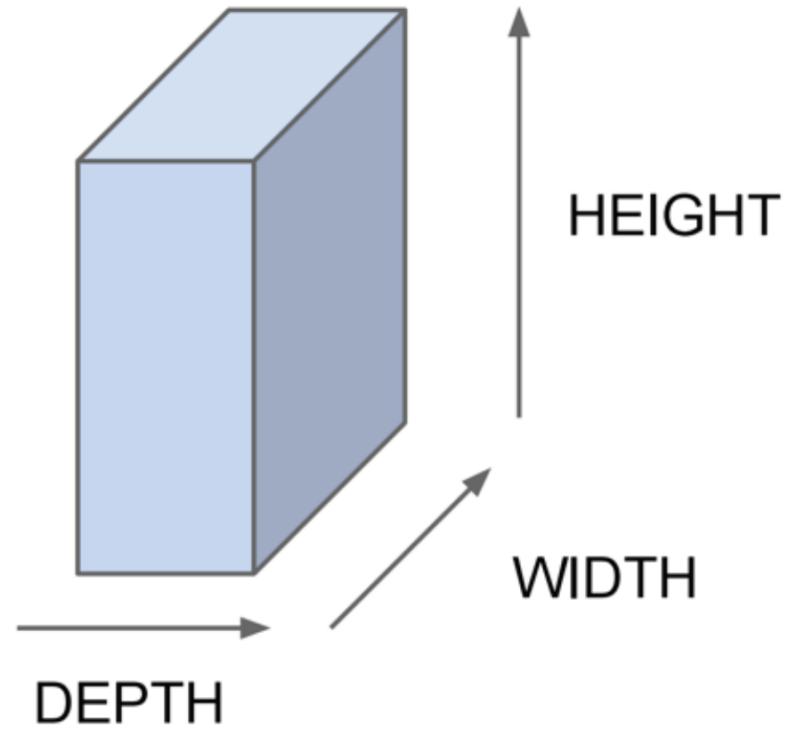
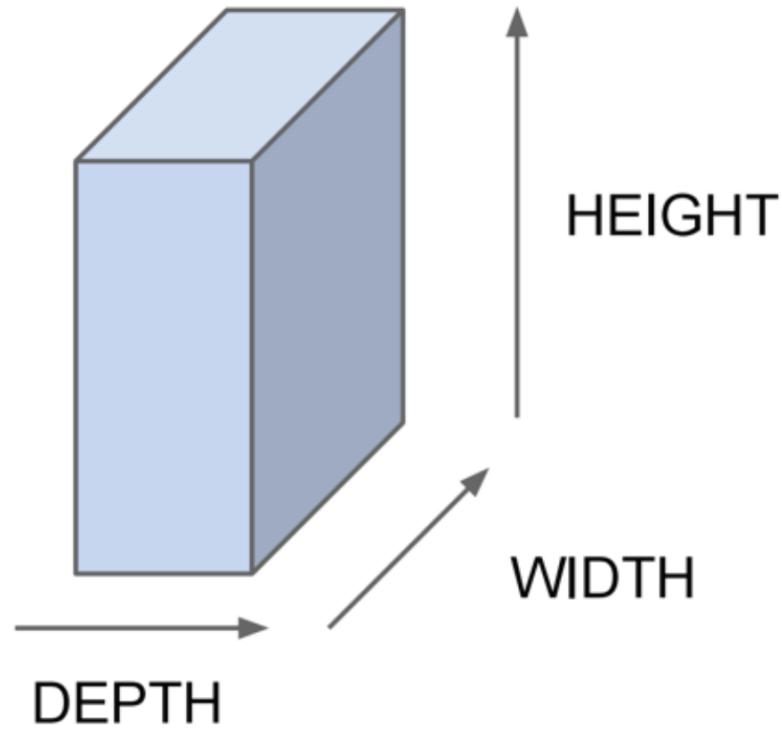


Figure: Andrey Karpathy

3D Activations

All Neural Net activations arranged in **3 dimensions**:

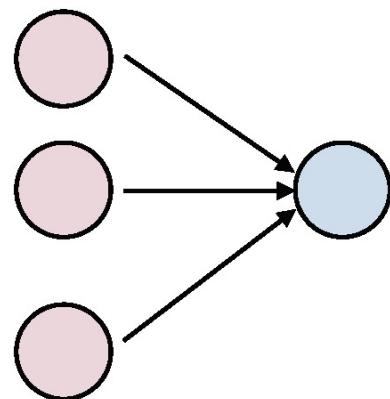


For example, a CIFAR-10 image is a $3 \times 32 \times 32$ volume (3 depth — RGB channels, 32 height, 32 width)

Figure: Andrey Karpathy

3D Activations

1D Activations:



3D Activations:

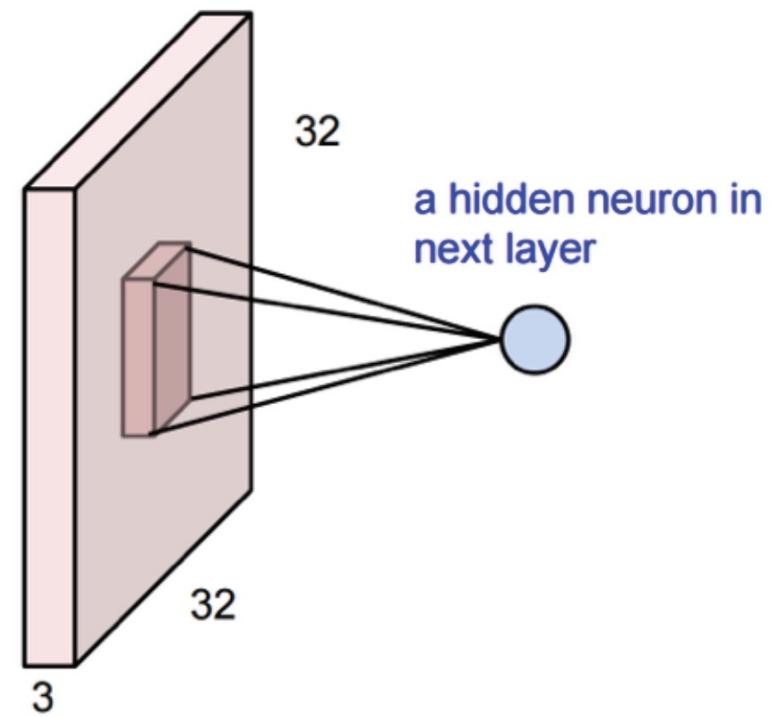
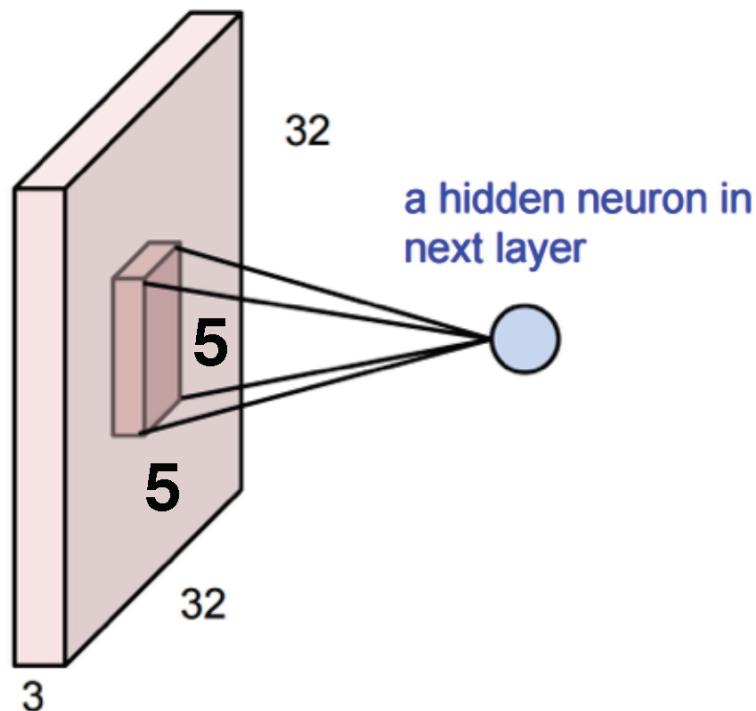


Figure: Andrey Karpathy

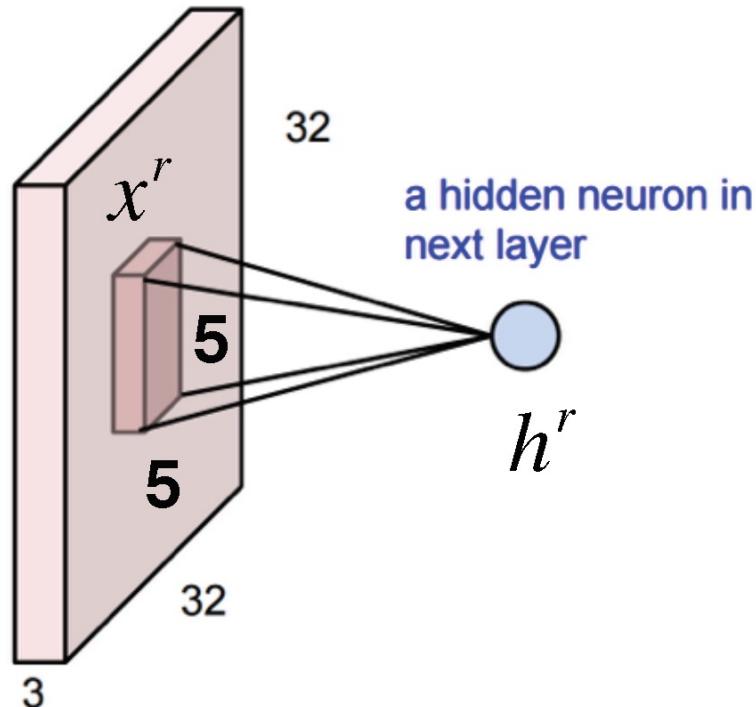
3D Activations



- The input is $3 \times 32 \times 32$
- This neuron depends on a $3 \times 5 \times 5$ chunk of the input
- The neuron also has a $3 \times 5 \times 5$ set of weights and a bias (scalar)

Figure: Andrej Karpathy

3D Activations



Example: consider the region of the input “ x^r ”

With output neuron h^r

Then the output is:

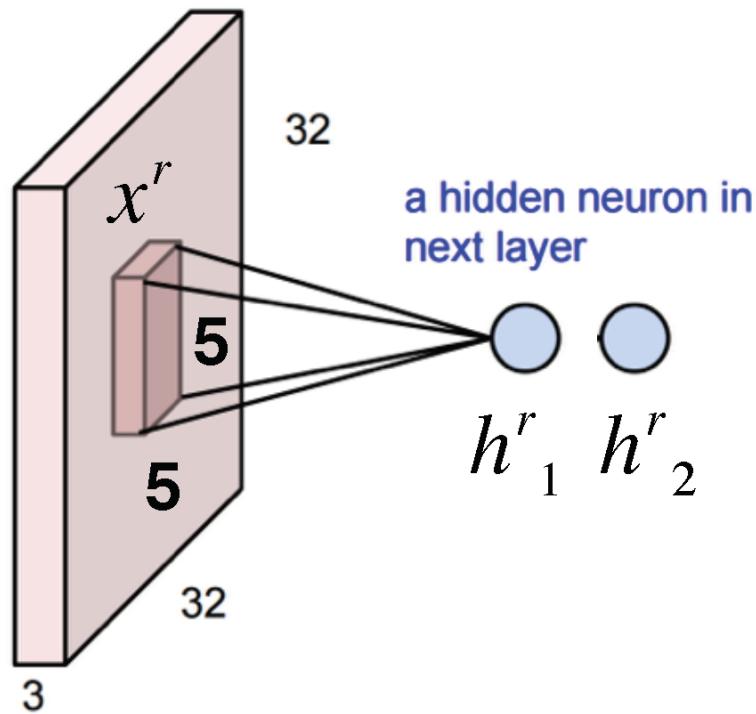
$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$



Sum over 3 axes

Figure: Andrej Karpathy

3D Activations



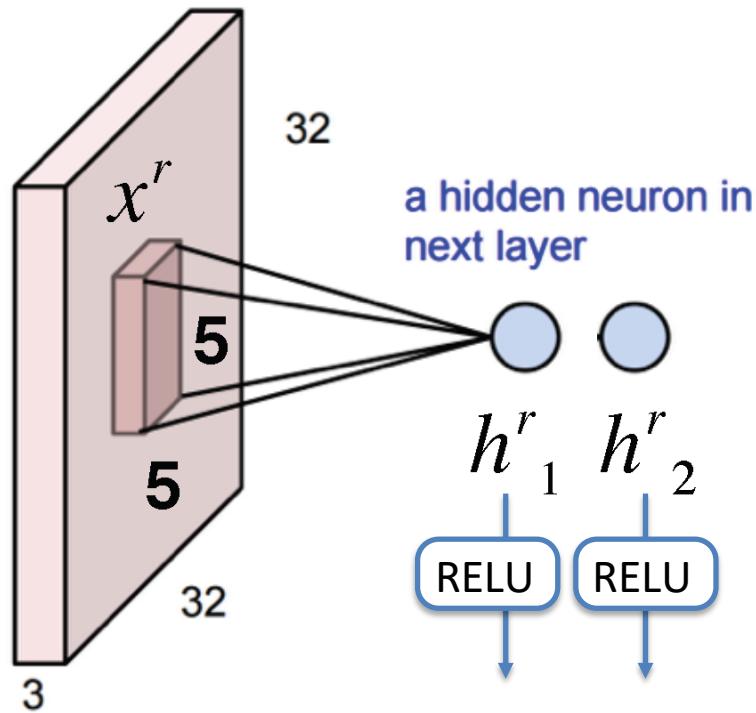
With **2** output neurons

$$h^r_1 = \sum_{ijk} x^r_{ijk} W_{1ijk} + b_1$$

$$h^r_2 = \sum_{ijk} x^r_{ijk} W_{2ijk} + b_2$$

Figure: Andrej Karpathy

3D Activations



Apply non-linear function unit,
For instance, rectified linear unit:

RELU(h^r_1)

With **2** output neurons

$$h^r_1 = \sum_{ijk} x^r_{ijk} W_{1ijk} + b_1$$

$$h^r_2 = \sum_{ijk} x^r_{ijk} W_{2ijk} + b_2$$

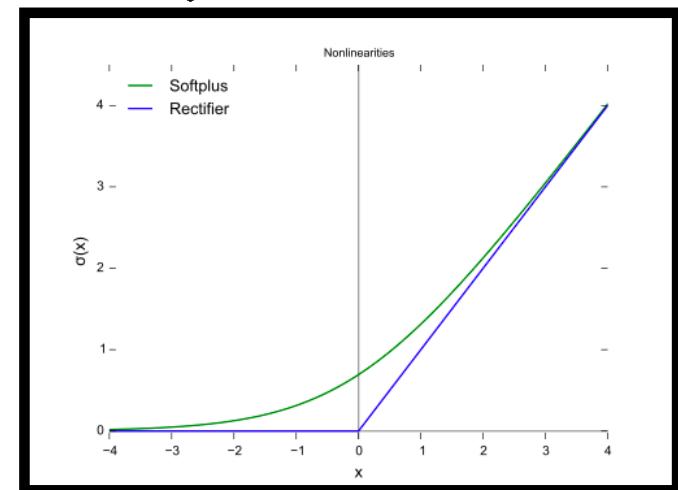
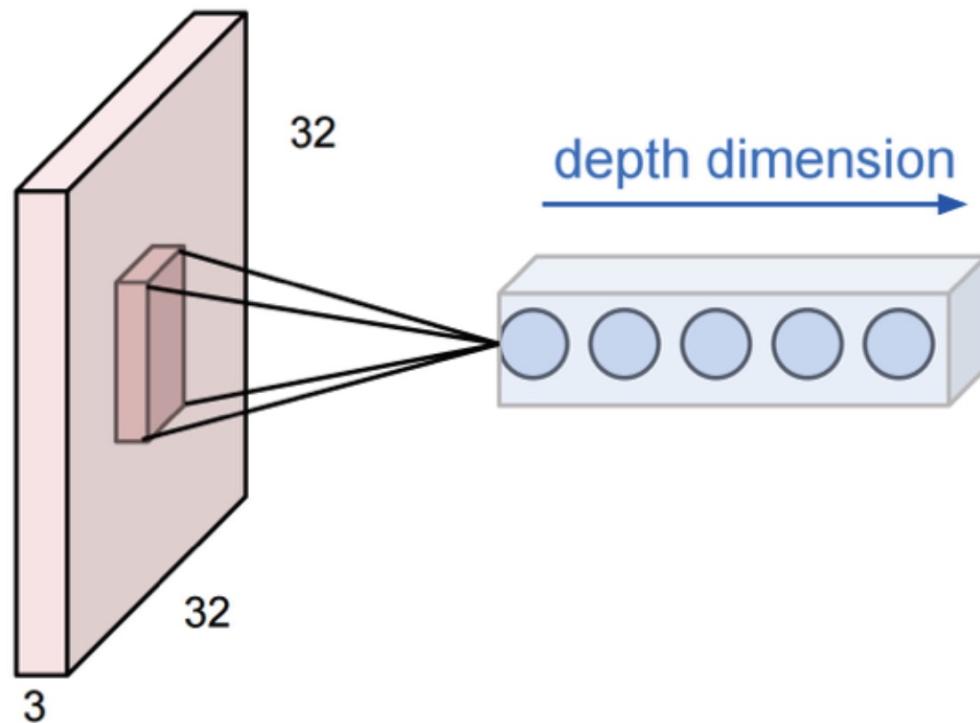


Figure: Andrej Karpathy

3D Activations

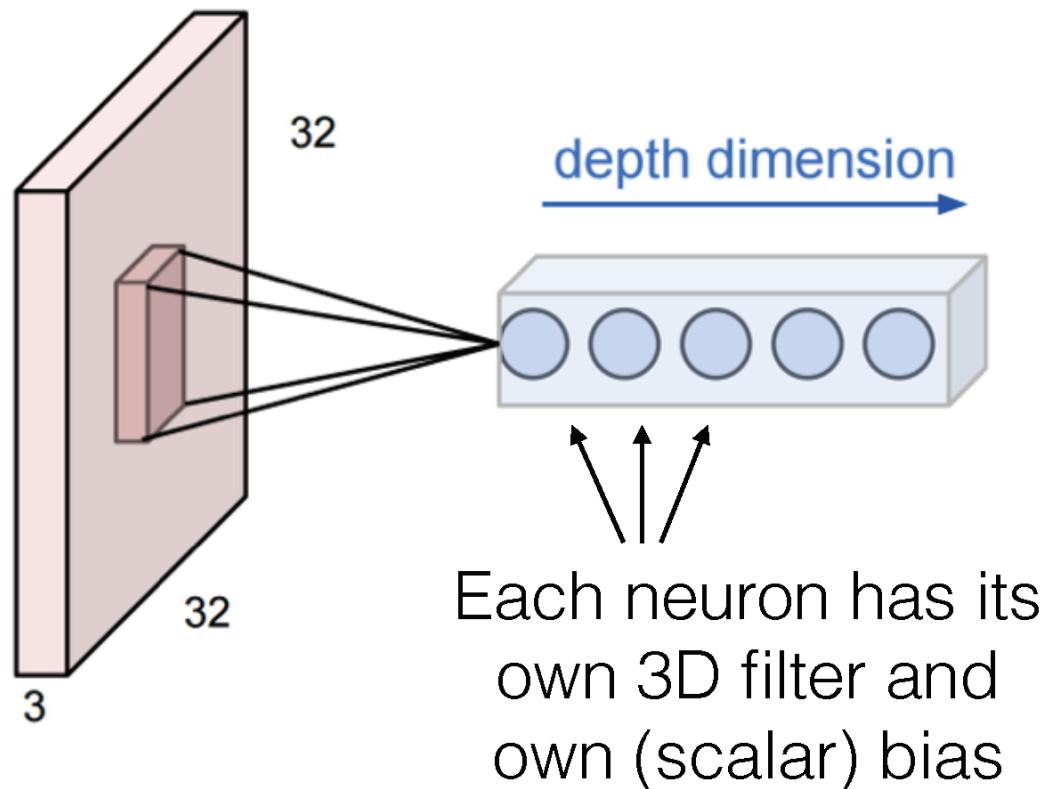


We can keep adding more outputs

These form a column in the output volume:
[depth x 1 x 1]

Figure: Andrej Karpathy

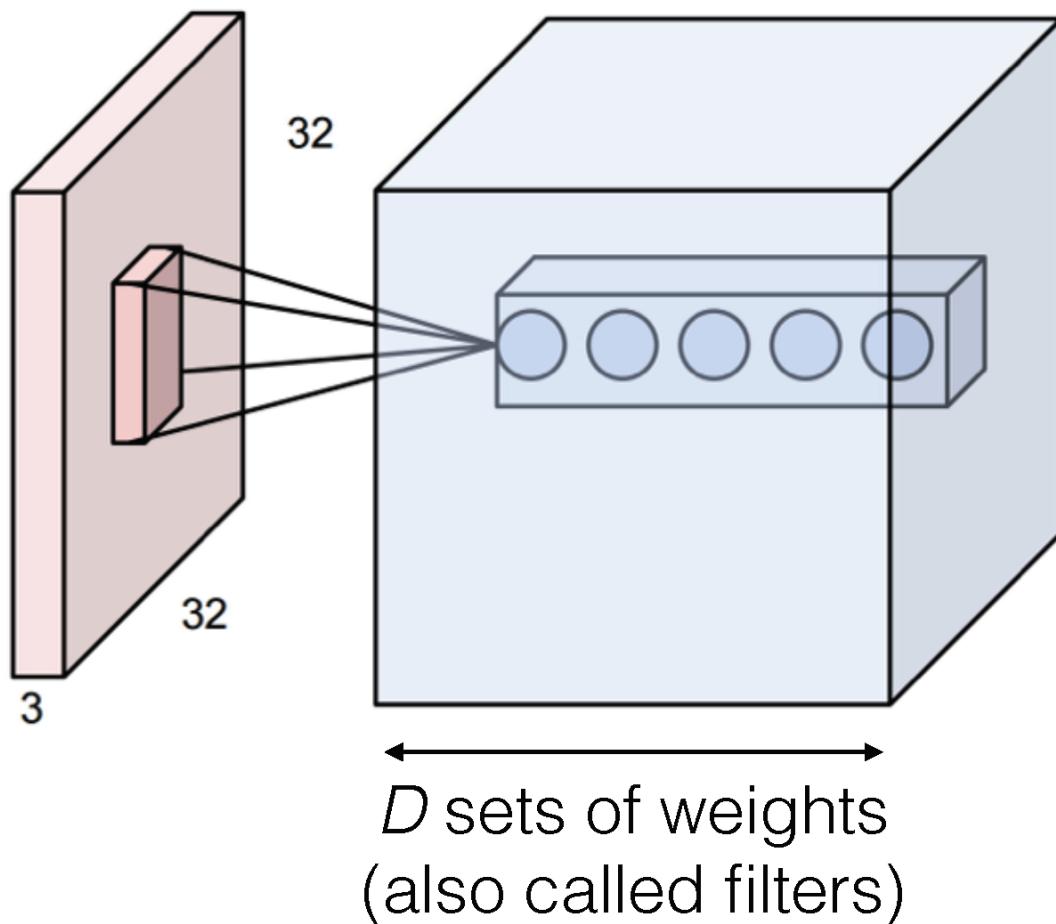
3D Activations



We can keep adding more outputs

These form a column in the output volume:
[depth x 1 x 1]

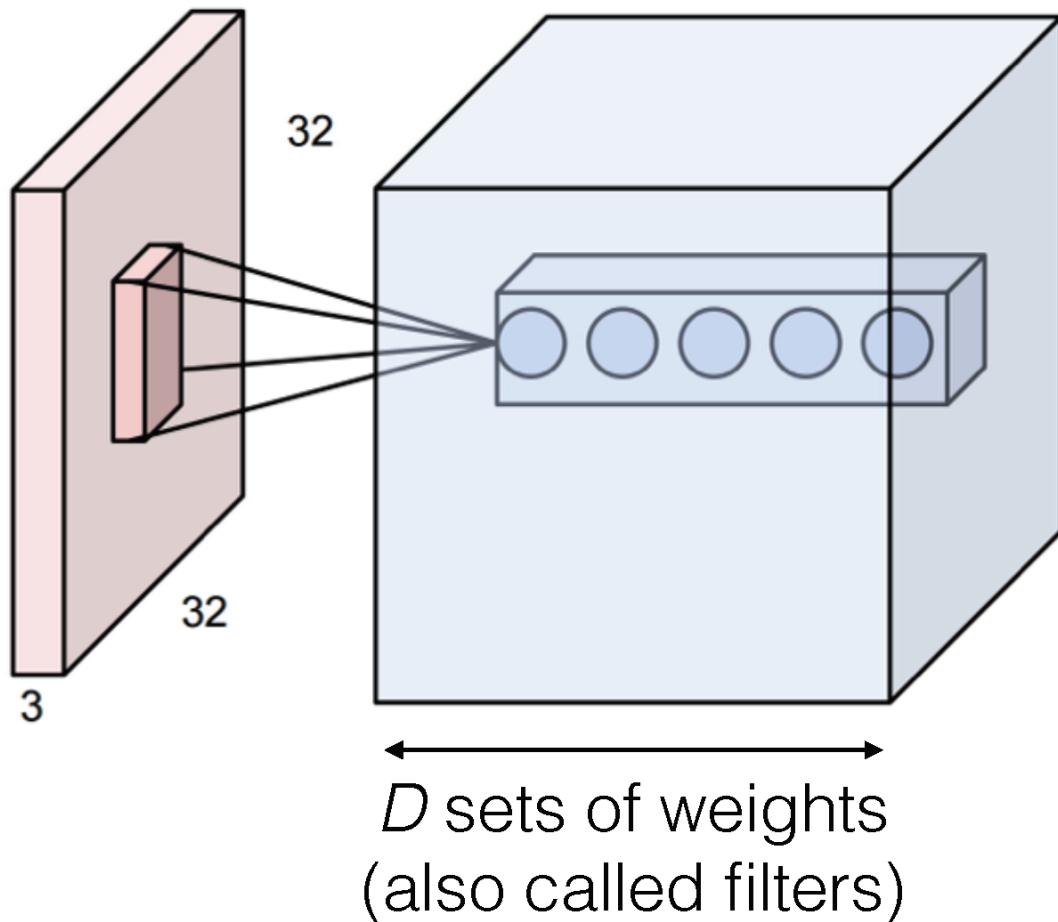
3D Activations



Now repeat this
across the input

Figure: Andrej Karpathy

3D Activations

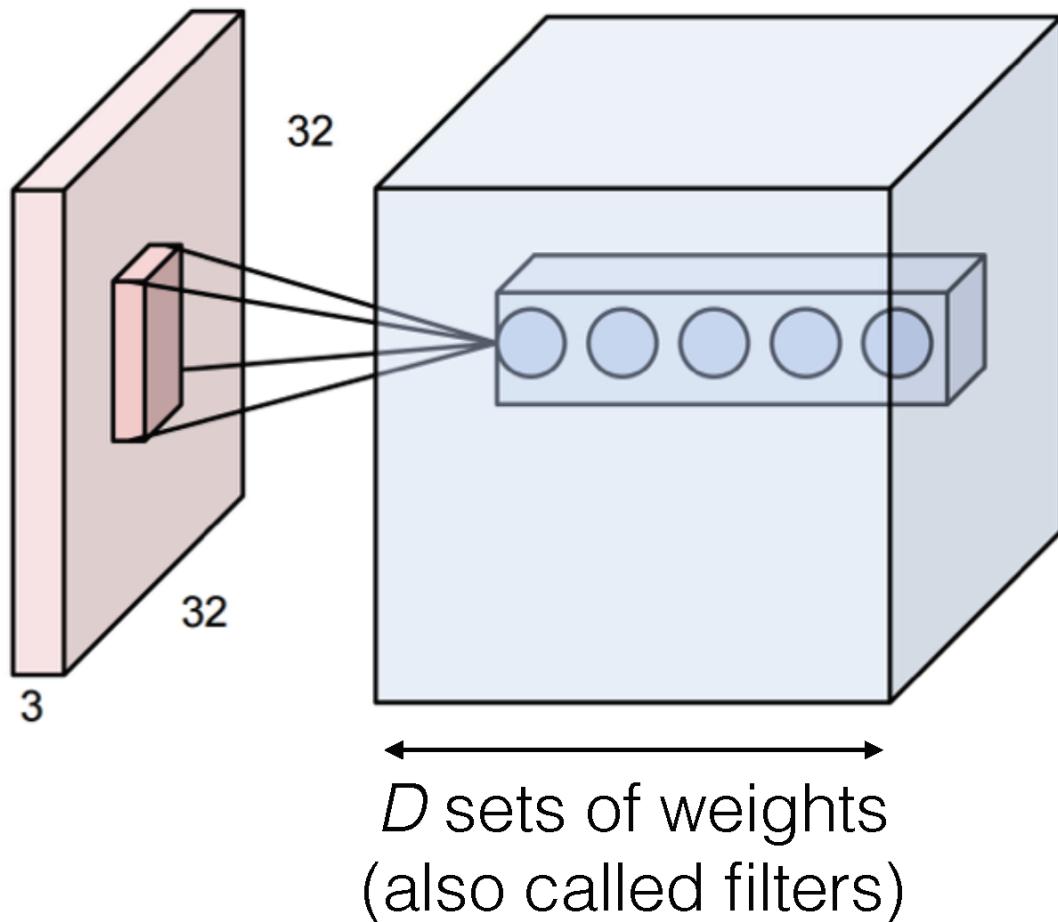


Now repeat this across the input

Weight sharing:
Each filter shares
the same weights
(but each depth
index has its own
set of weights)

Figure: Andrey Karpathy

3D Activations

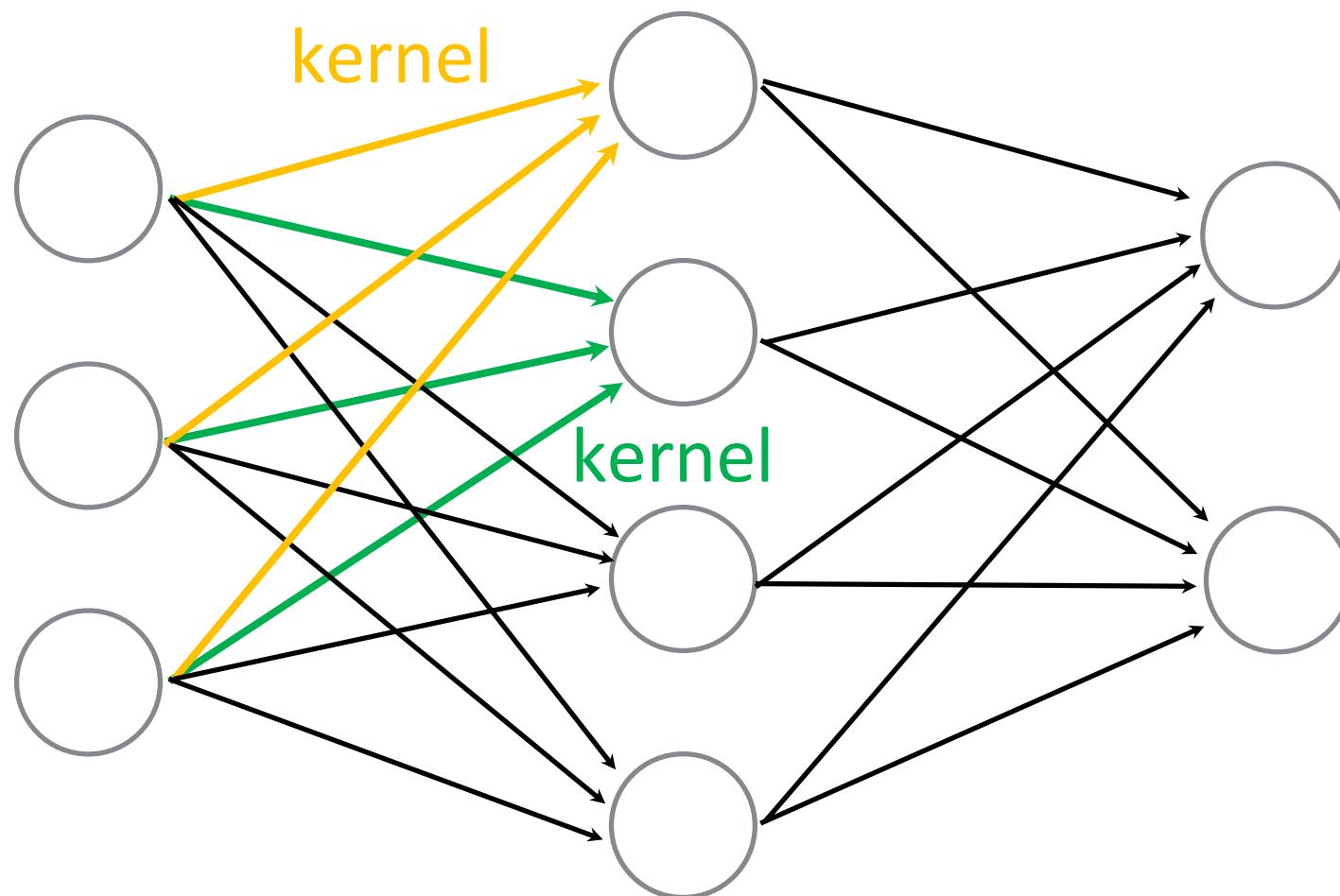


With weight sharing,
this is called
convolution

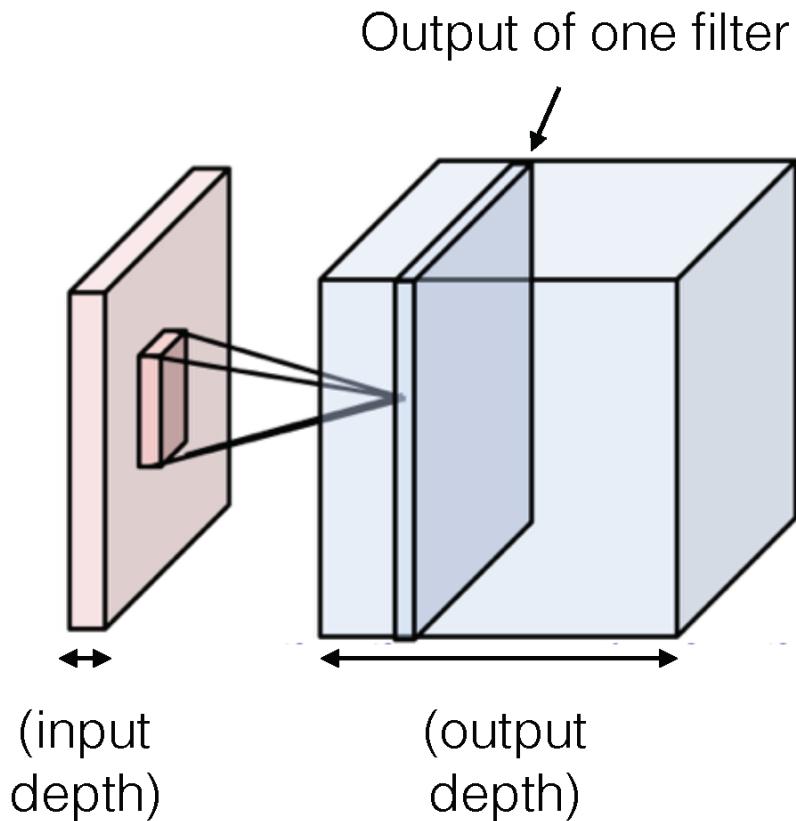
Without weight sharing,
this is called a
locally connected layer

Figure: Andrej Karpathy

Fully connected and Convolution



3D Activations

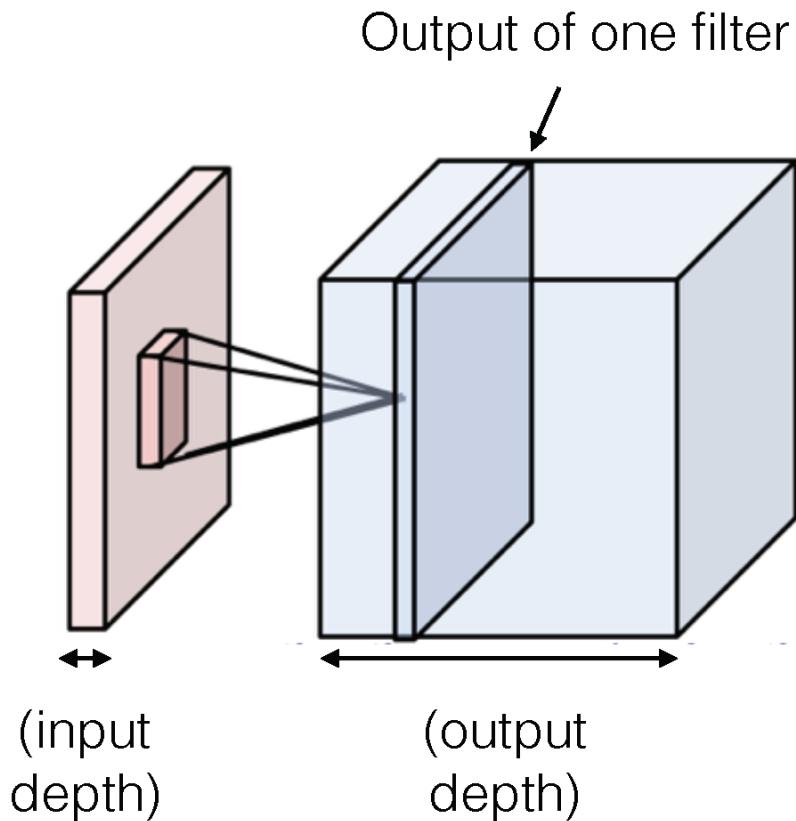


One set of weights gives
one slice in the output

To get a 3D output of depth D ,
use D different filters

In practice, ConvNets use
many filters (~ 64 to 1024)

3D Activations



One set of weights gives
one slice in the output

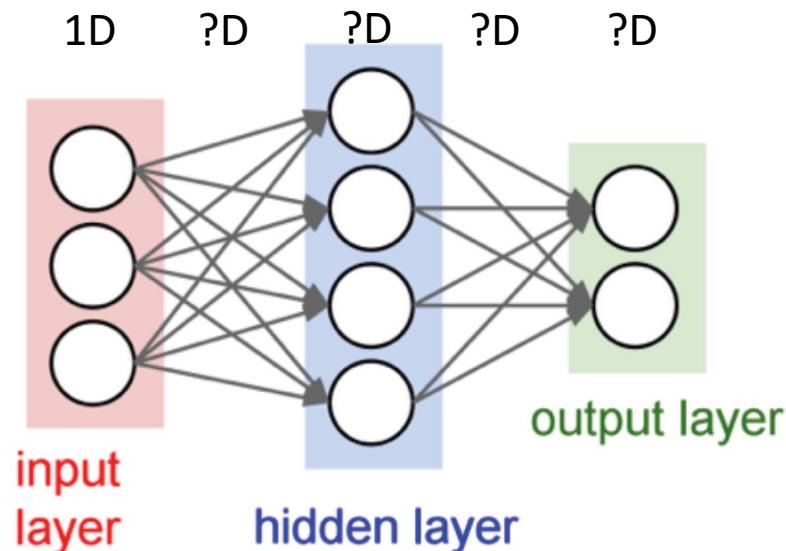
To get a 3D output of depth D ,
use D different filters

In practice, ConvNets use
many filters (~ 64 to 1024)

All together, the weights are **4** dimensional:
(output depth, input depth, kernel height, kernel width)

3D Activations

before:



now:

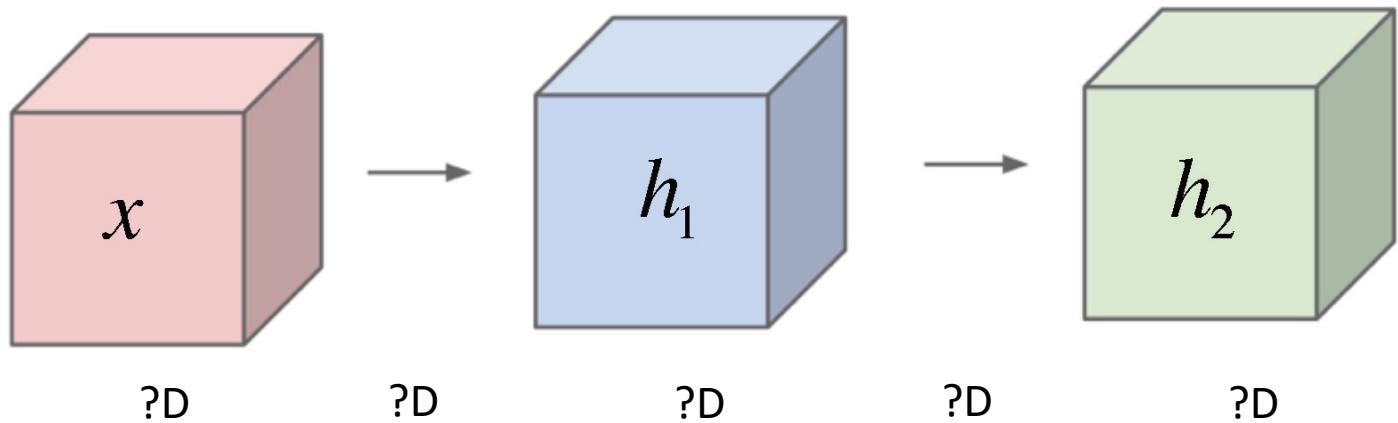
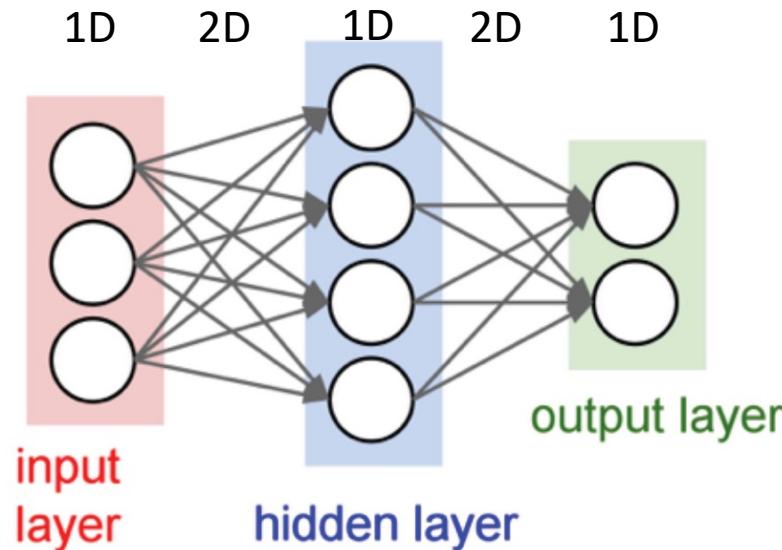


Figure: Andrej Karpathy

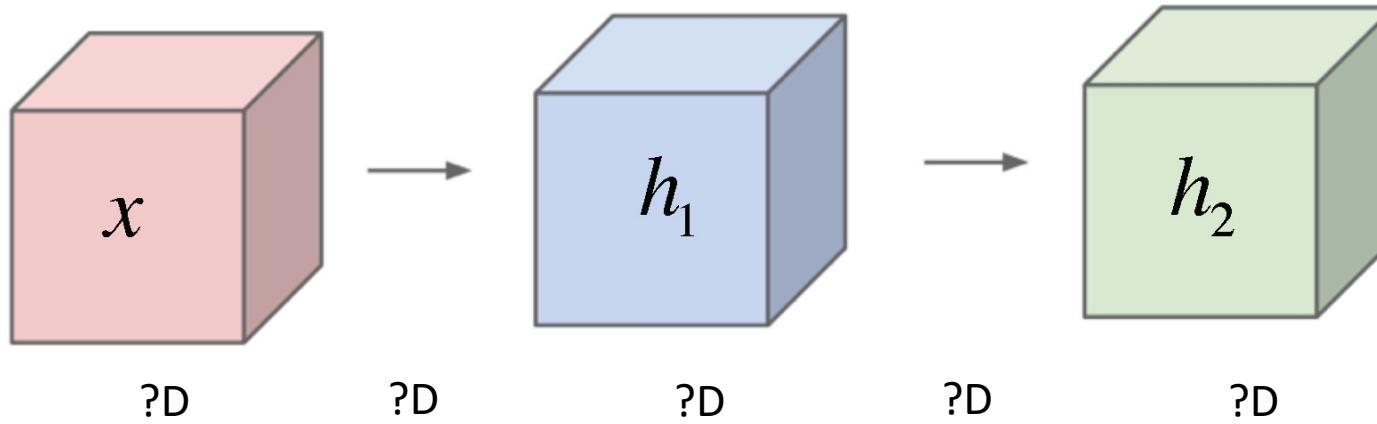
3D Activations

before:



(1D vectors)

now:

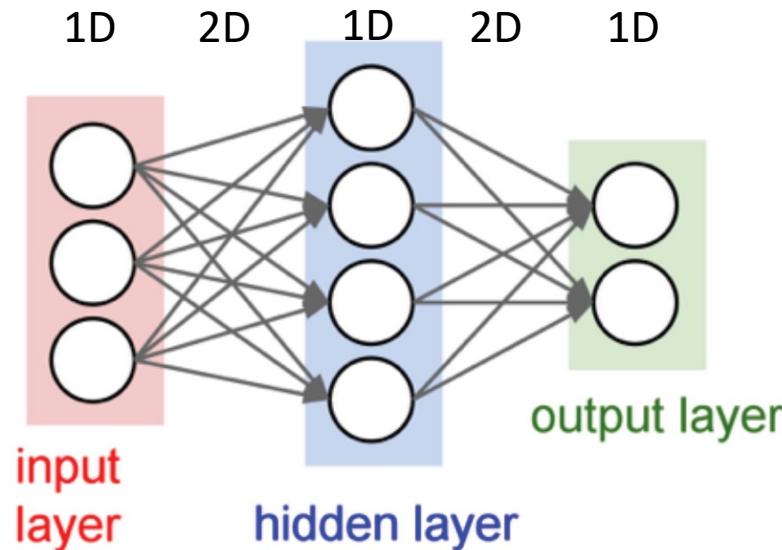


(3D arrays)

Figure: Andrej Karpathy

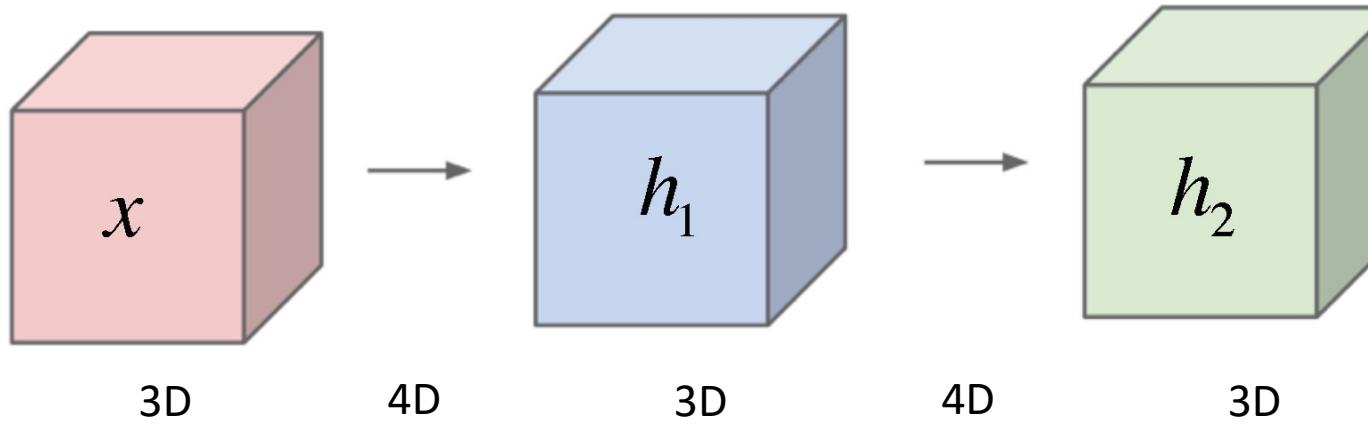
3D Activations

before:



(1D vectors)

now:



(3D arrays)

Figure: Andrej Karpathy

3D Activations

We can unravel the 3D cube and show each layer separately:

(Input)

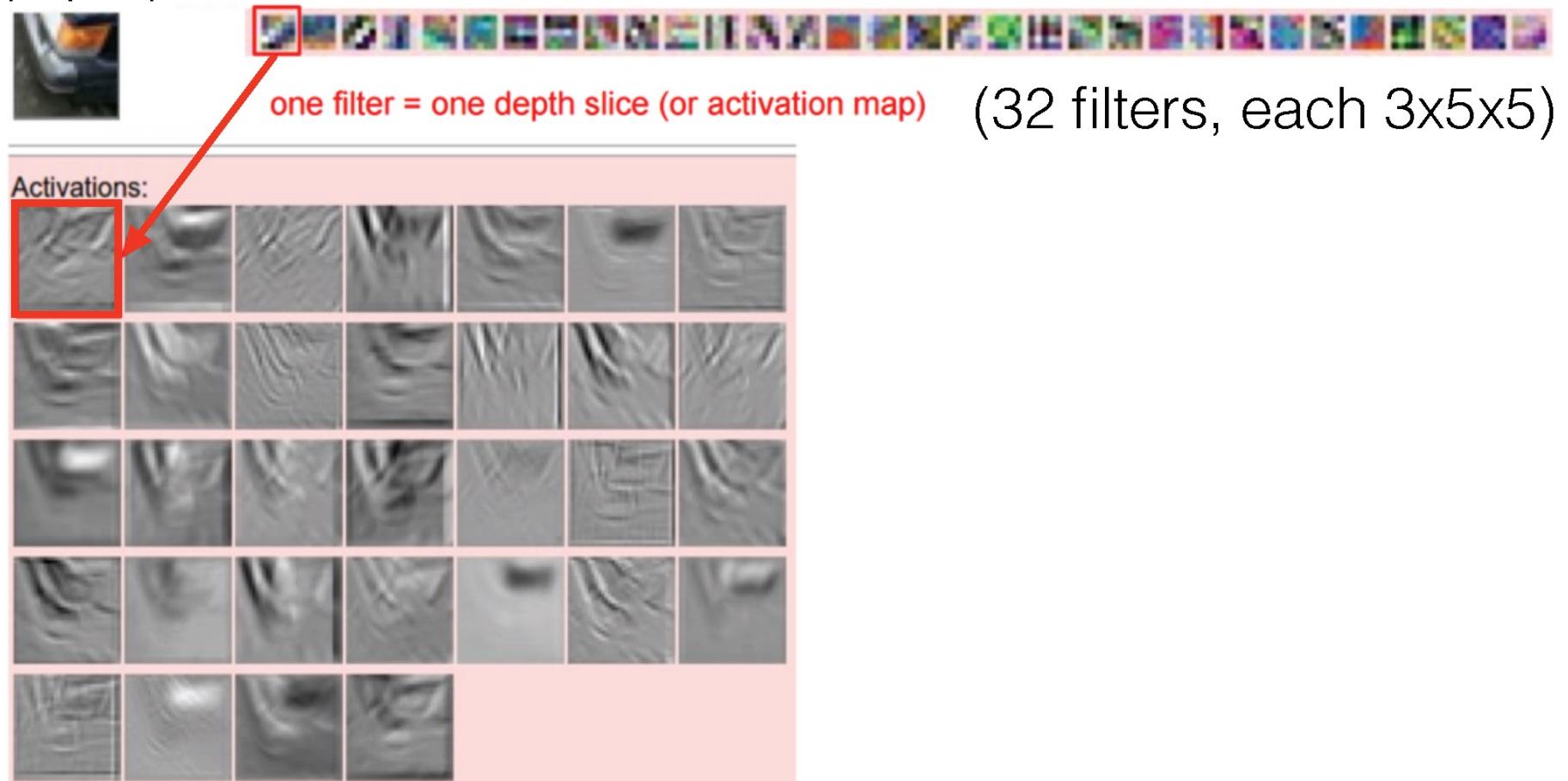


Figure: Andrej Karpathy

3D Activations

We can unravel the 3D cube and show each layer separately:

(Input)



Figure: Andrej Karpathy

3D Activations

We can unravel the 3D cube and show each layer separately:

(Input)

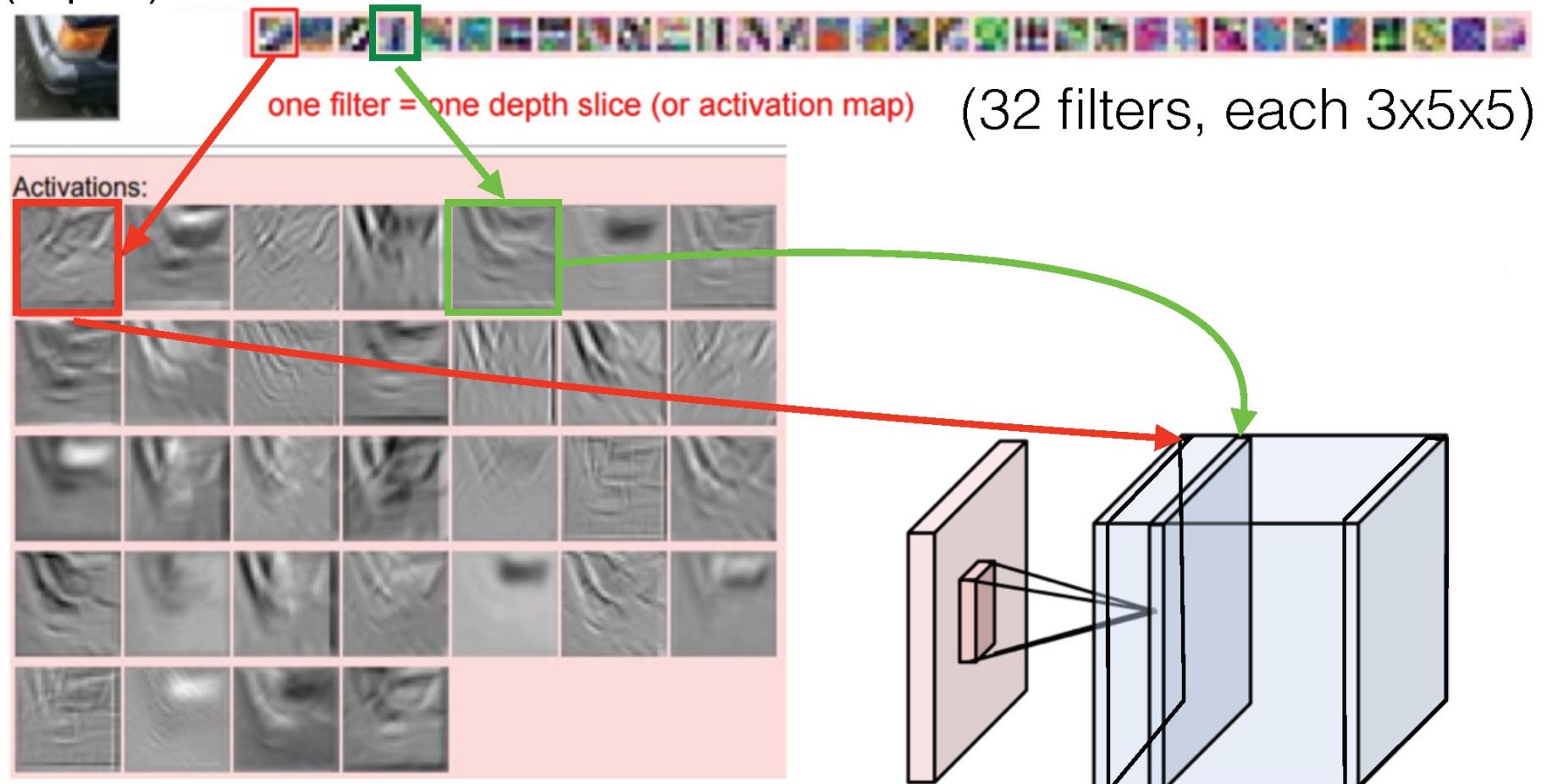


Figure: Andrej Karpathy

3D Activations

We can unravel the 3D cube and show each layer separately:

(Input)

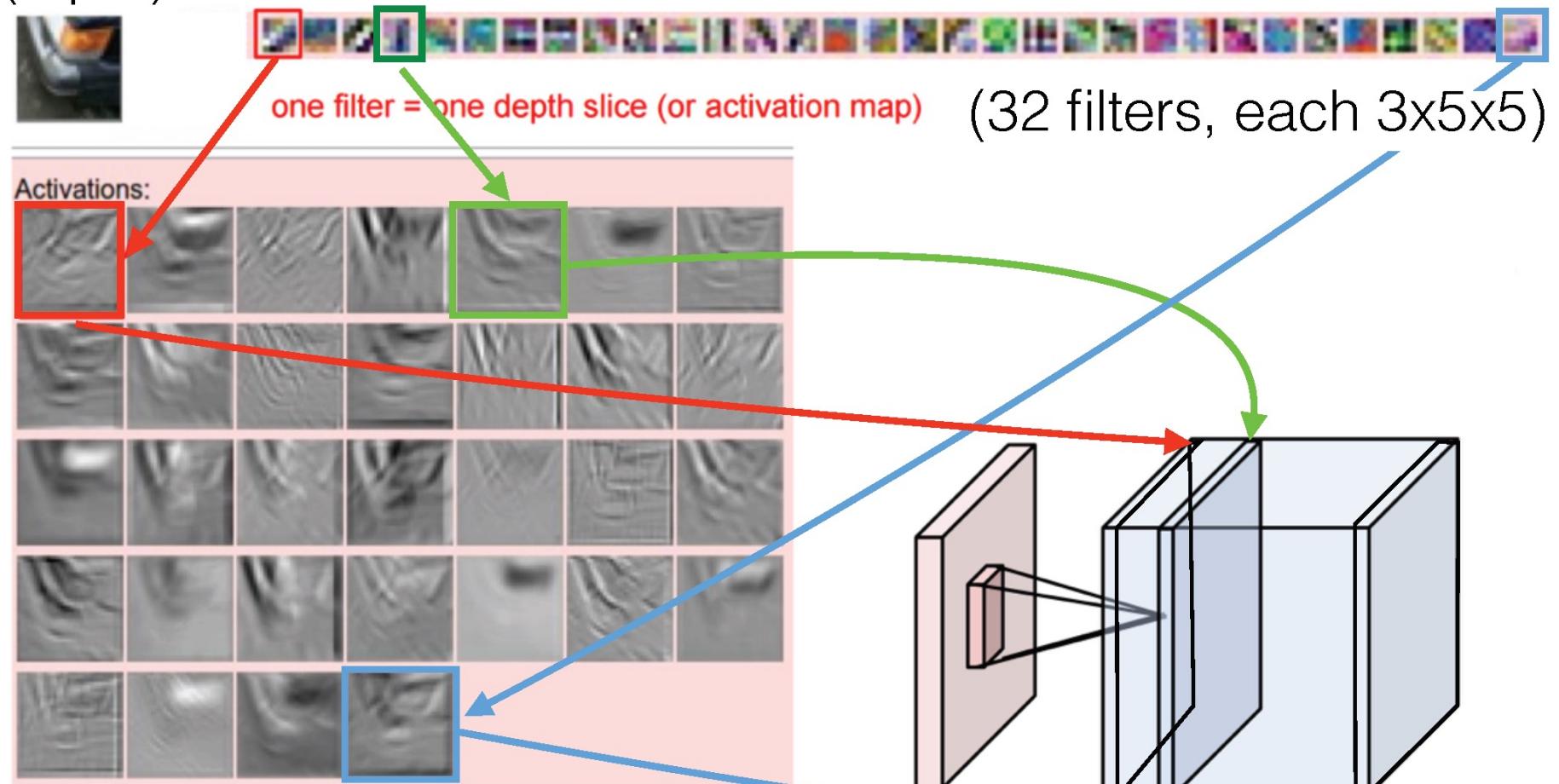
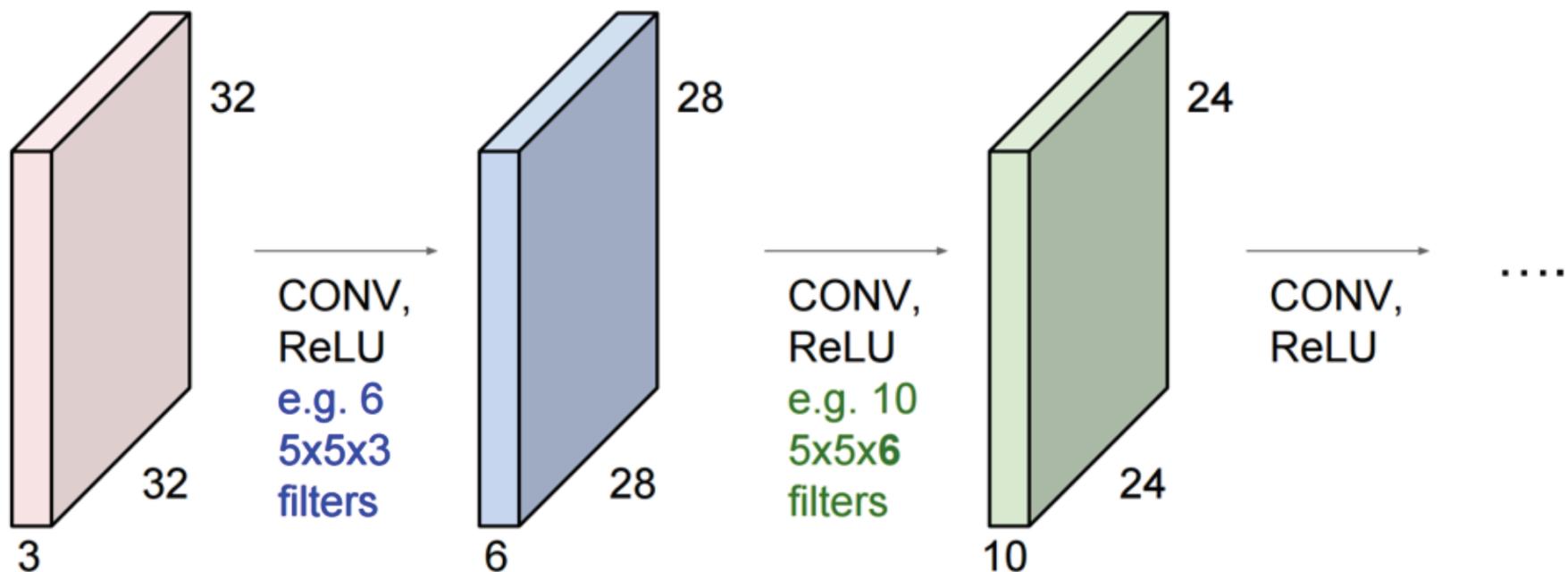


Figure: Andrej Karpathy

(Recap)

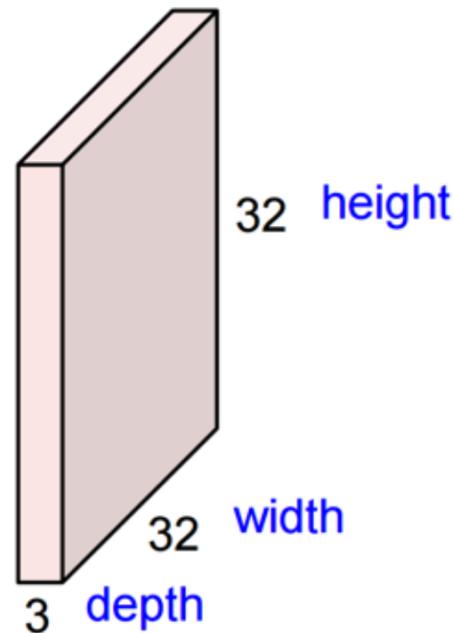
A **ConvNet** is a sequence of convolutional layers, interspersed with activation functions (and possibly other layer types)



(Recap)

Convolution Layer

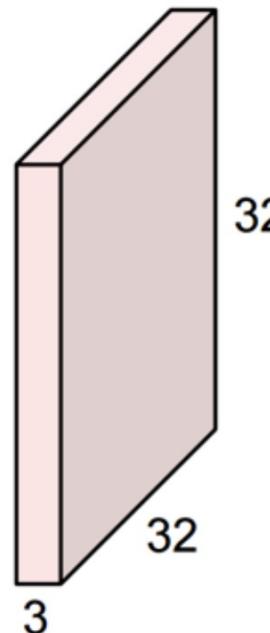
32x32x3 image



(Recap)

Convolution Layer

32x32x3 image



5x5x3 filter

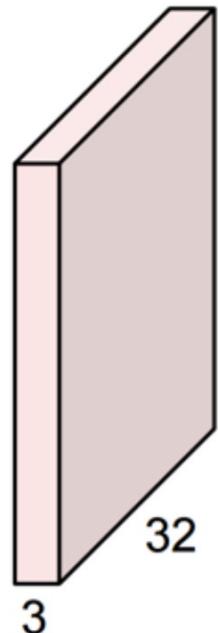


Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

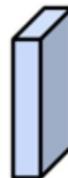
(Recap)

Convolution Layer

$32 \times 32 \times 3$ image



$5 \times 5 \times 3$ filter

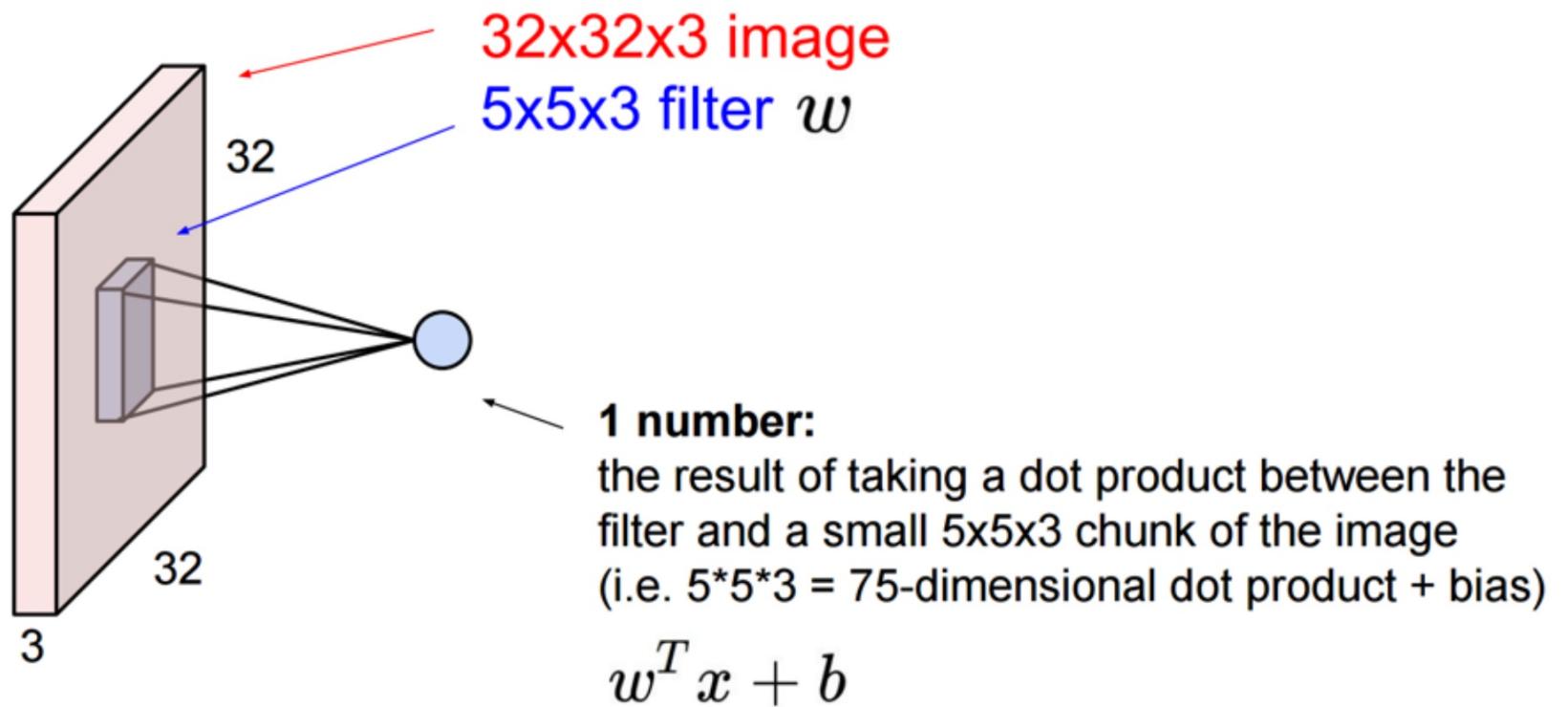


Filters always extend the full depth of the input volume

Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

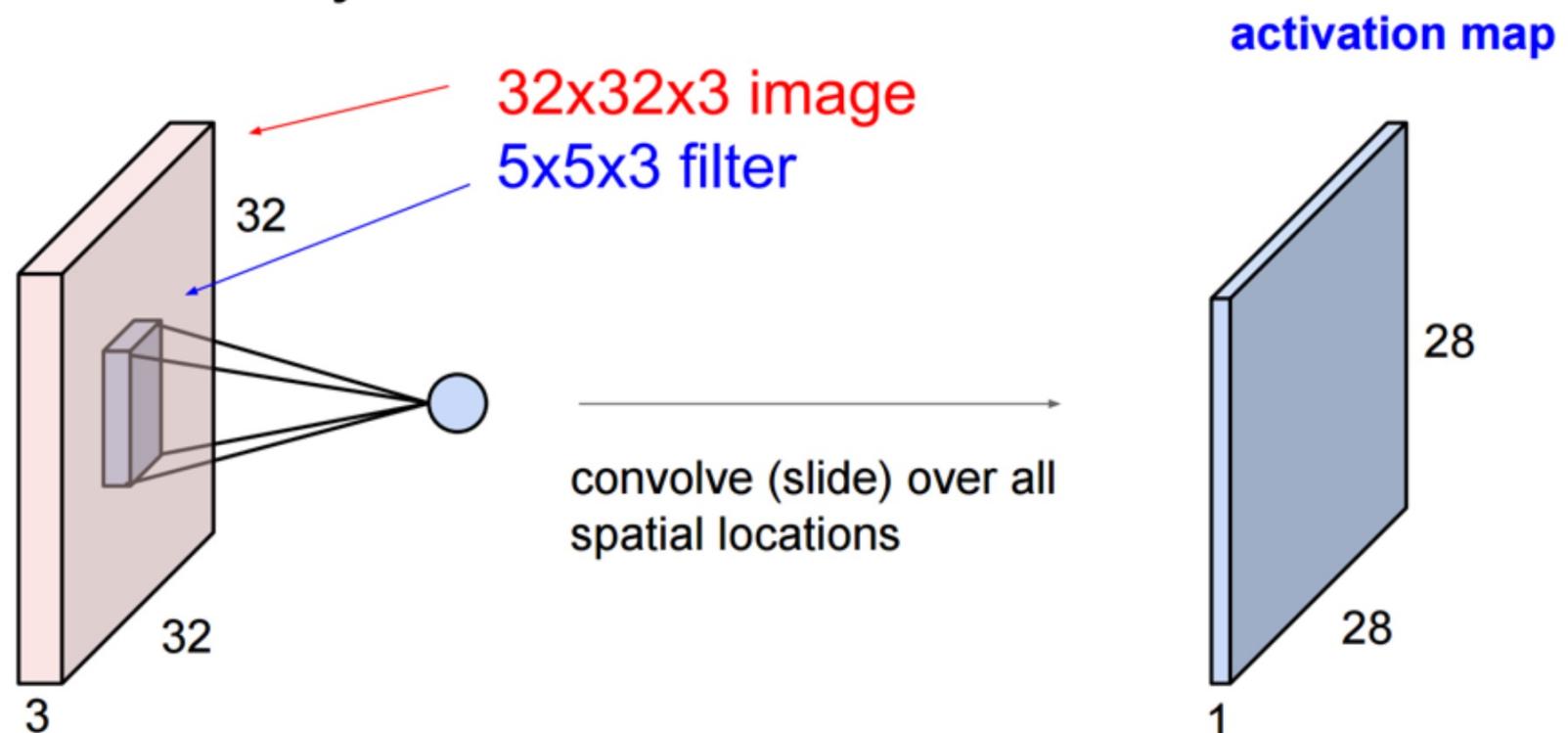
(Recap)

Convolution Layer



(Recap)

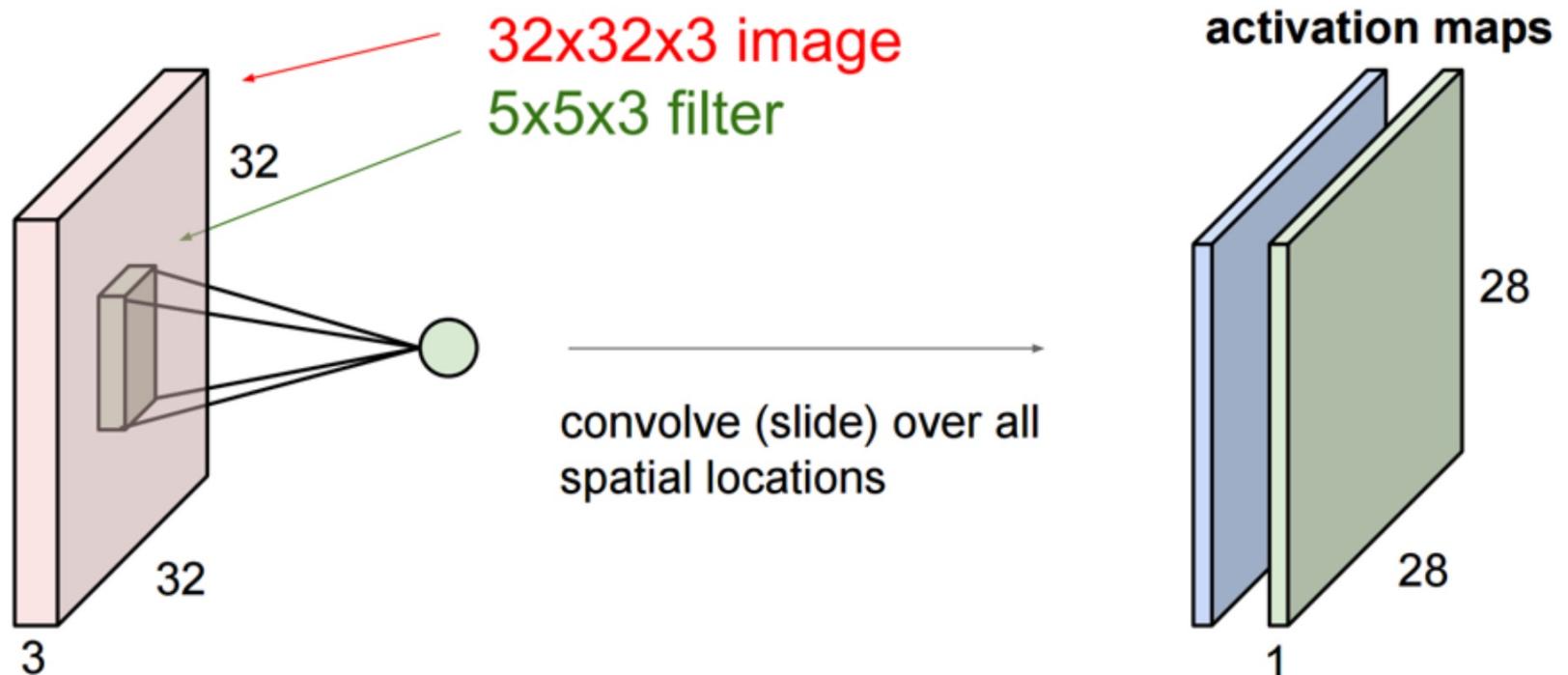
Convolution Layer



(Recap)

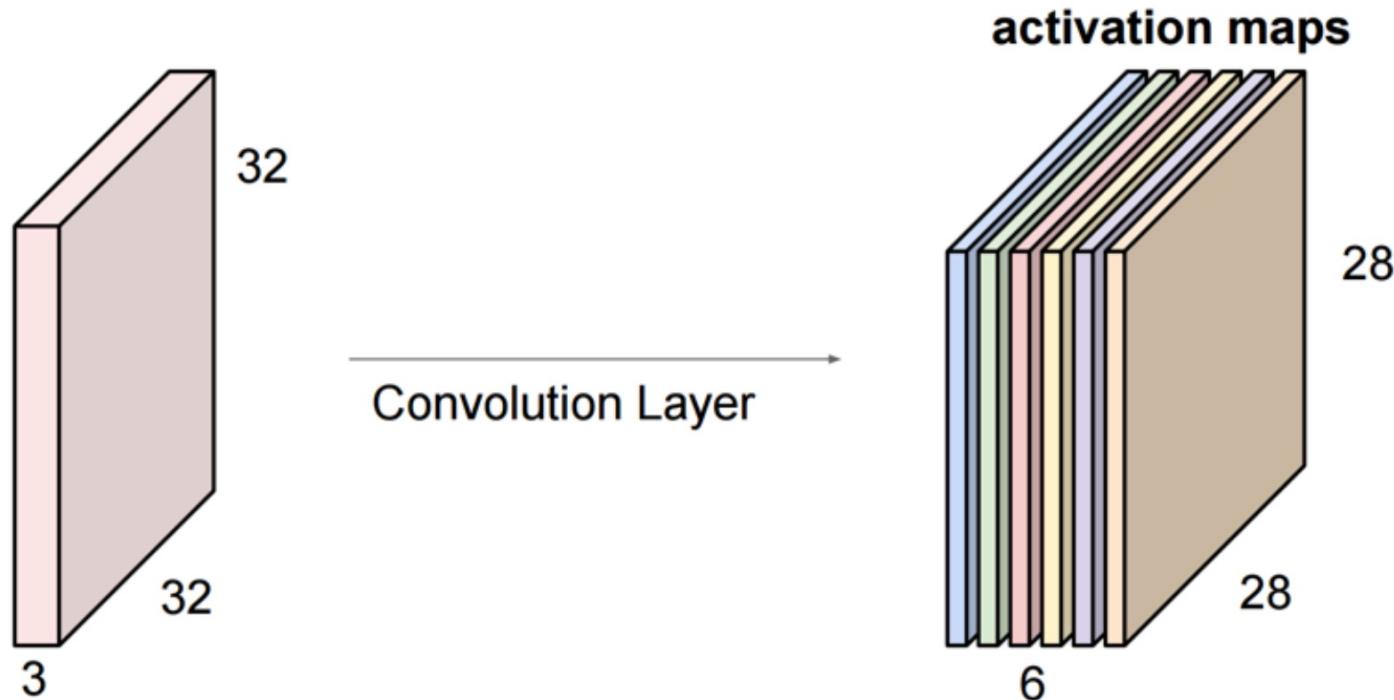
Convolution Layer

consider a second, green filter



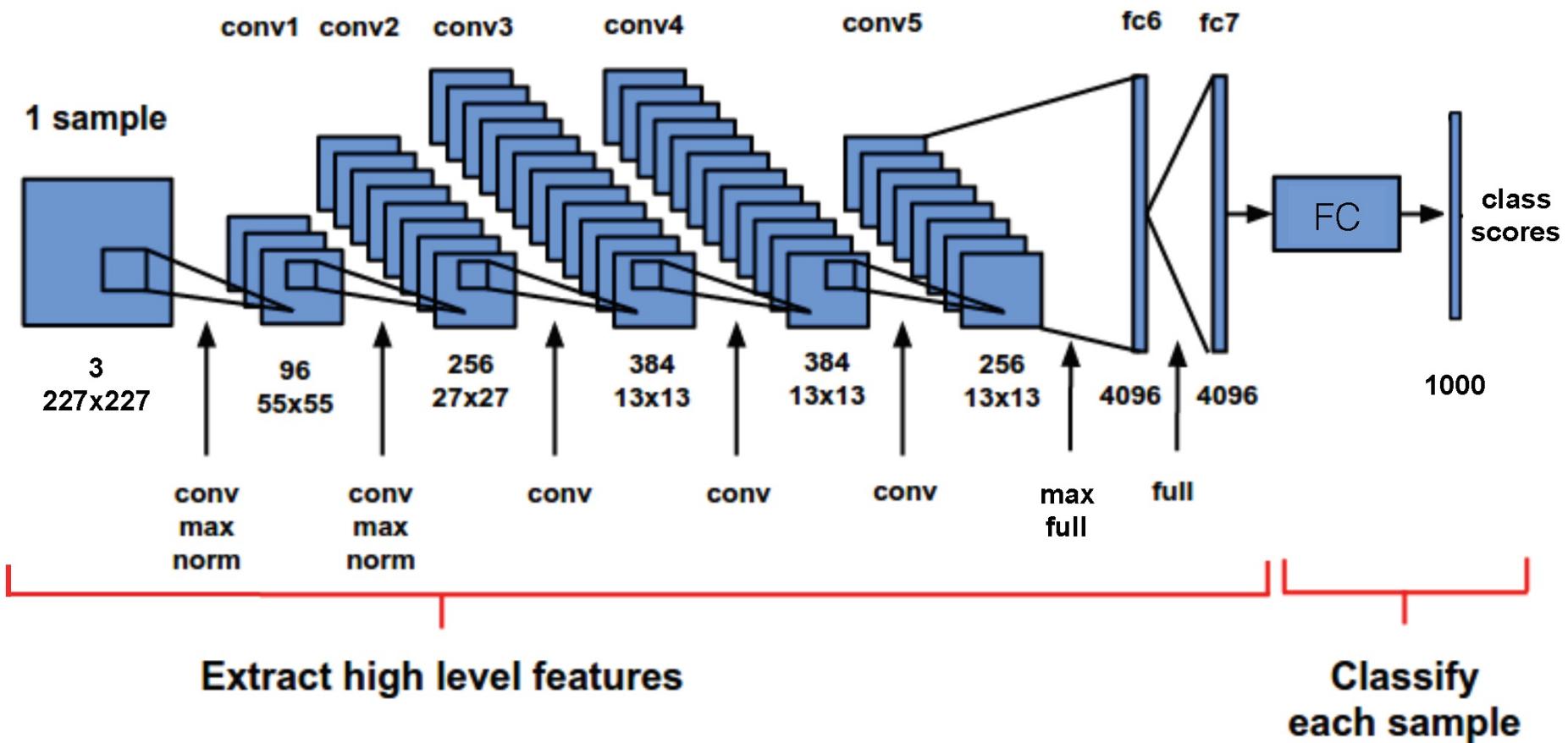
(Recap)

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

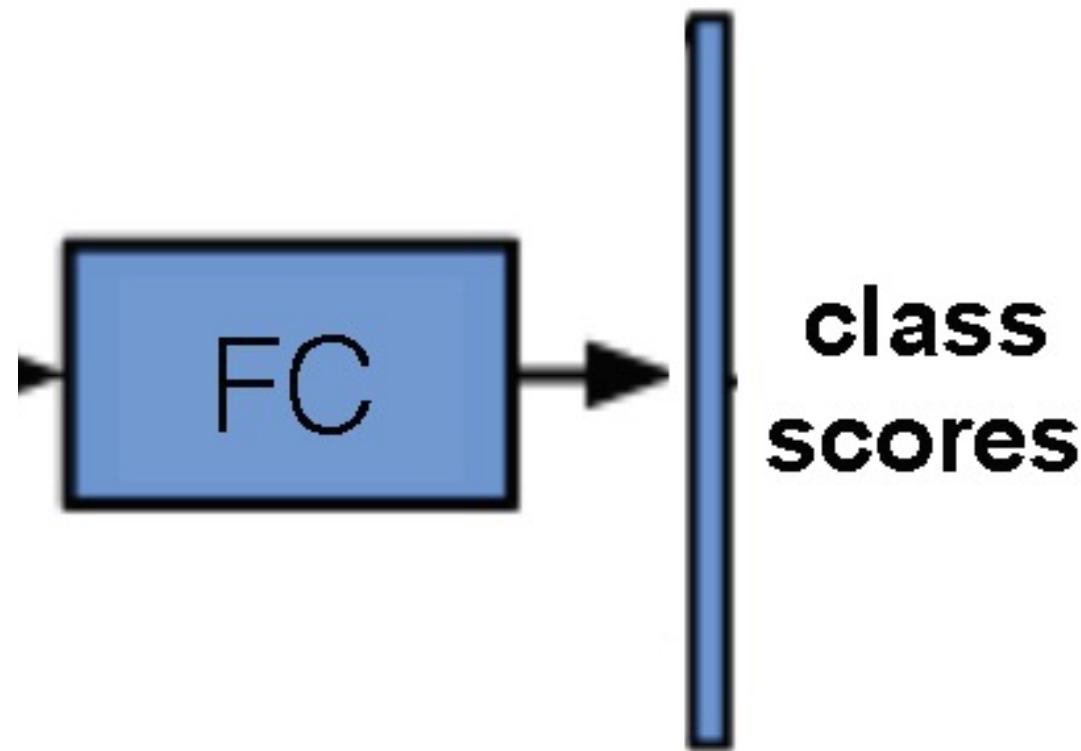


We stack these up to get a “new image” of size $28 \times 28 \times 6$!

Typical entire architecture

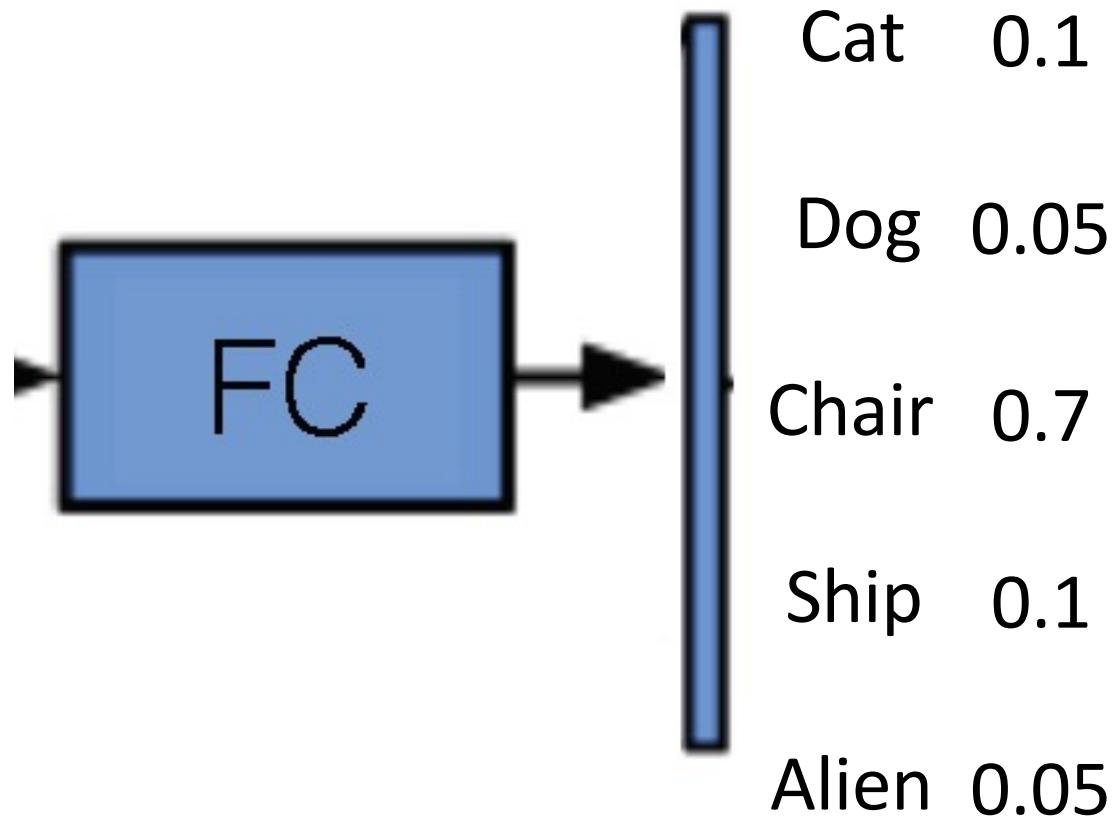


Typical entire architecture

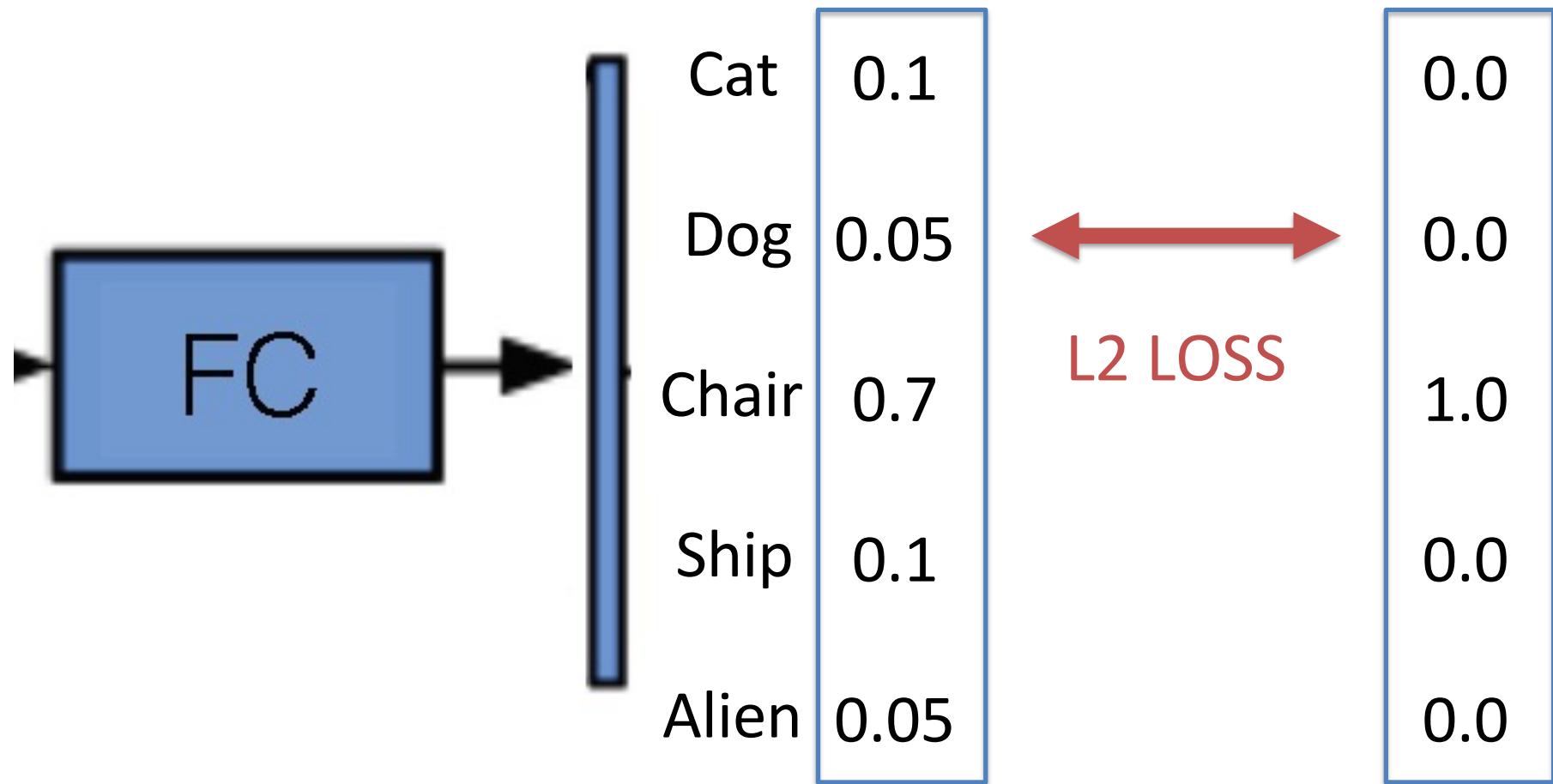


1000

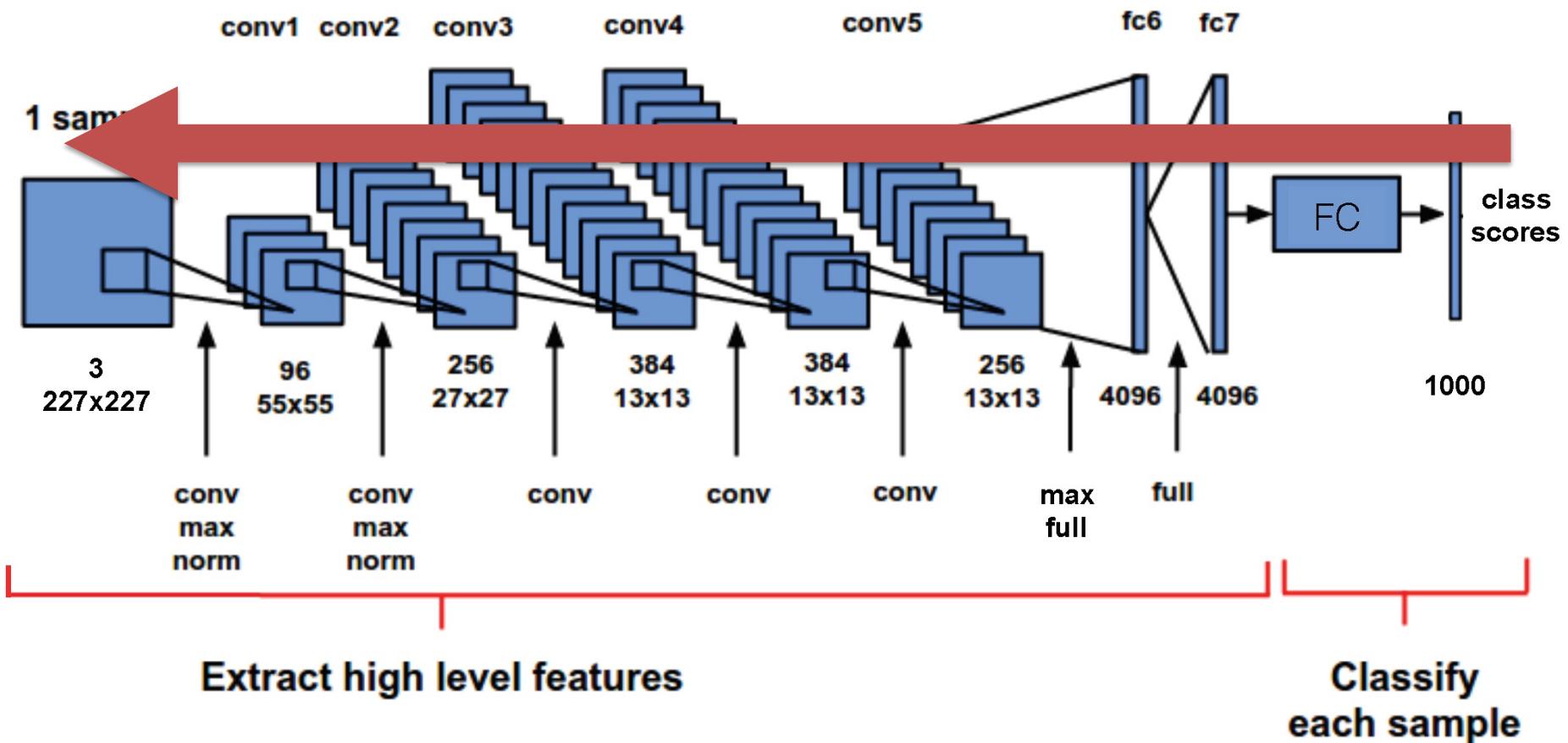
Typical entire architecture



Typical entire architecture



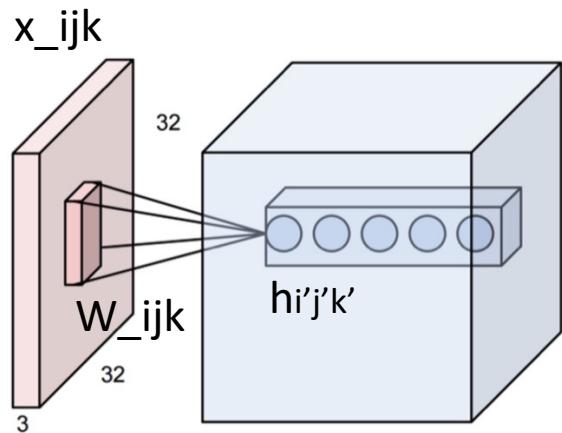
Typical entire architecture



Demos

- <http://cs231n.stanford.edu/>
- <http://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html>

Back propagation



Given gradients of the right layer

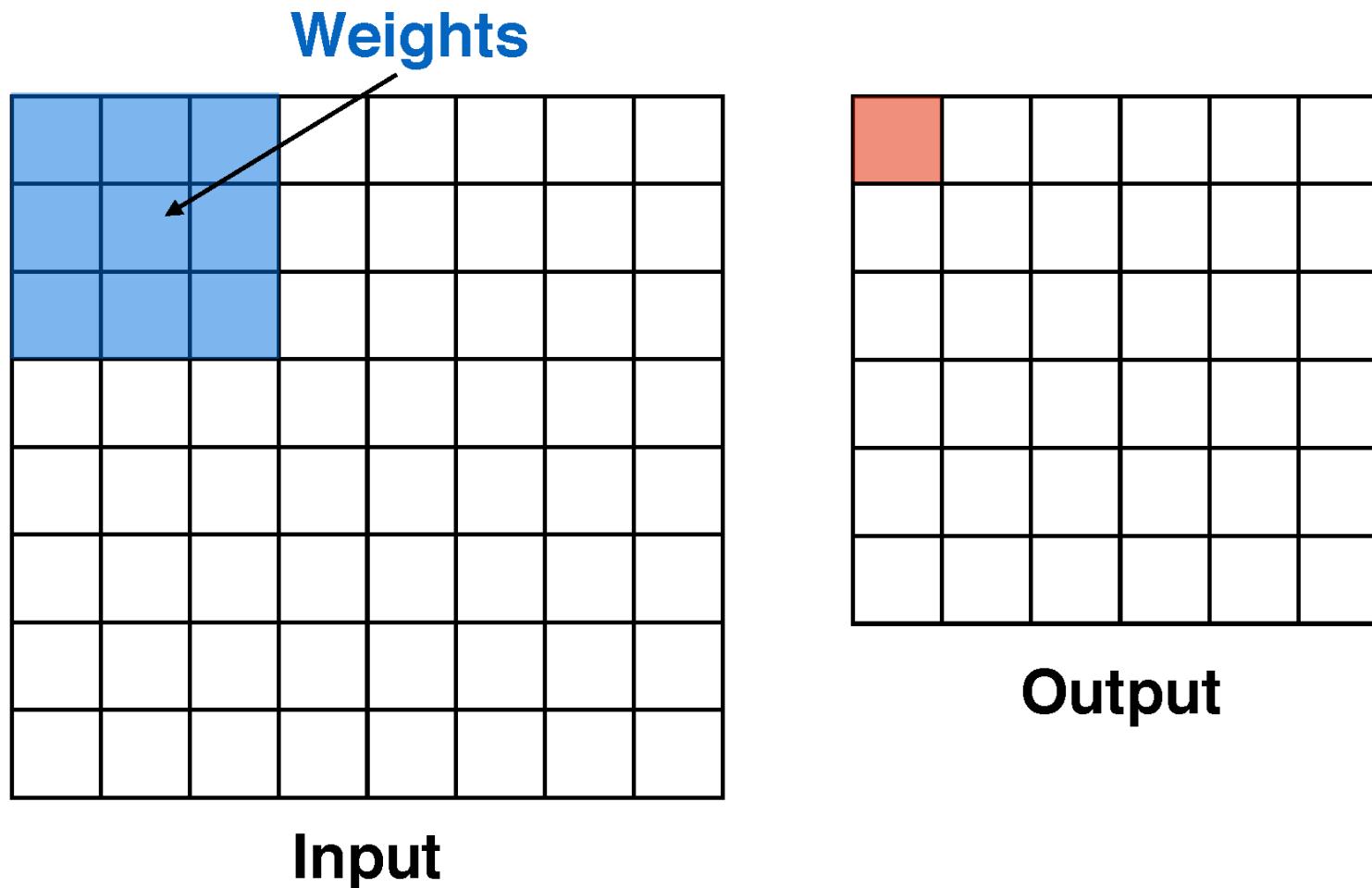
$$\frac{\partial \mathcal{L}}{\partial h_{i'j'k'}}$$

$$h_{i'j'k'} = \sum_{ijk} x_{ijk} W_{ijk} + b$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial x_{i'j'k'}} &= \frac{\partial \mathcal{L}}{\partial h_{i1j1k1}} \frac{\partial h_{i1j1k1}}{\partial x_{i'j'k'}} + \\ &\quad \frac{\partial \mathcal{L}}{\partial h_{i2j2k2}} \frac{\partial h_{i2j2k2}}{\partial x_{i'j'k'}} + \\ &\quad \frac{\partial \mathcal{L}}{\partial h_{i3j3k3}} \frac{\partial h_{i3j3k3}}{\partial x_{i'j'k'}} + \dots \end{aligned}$$

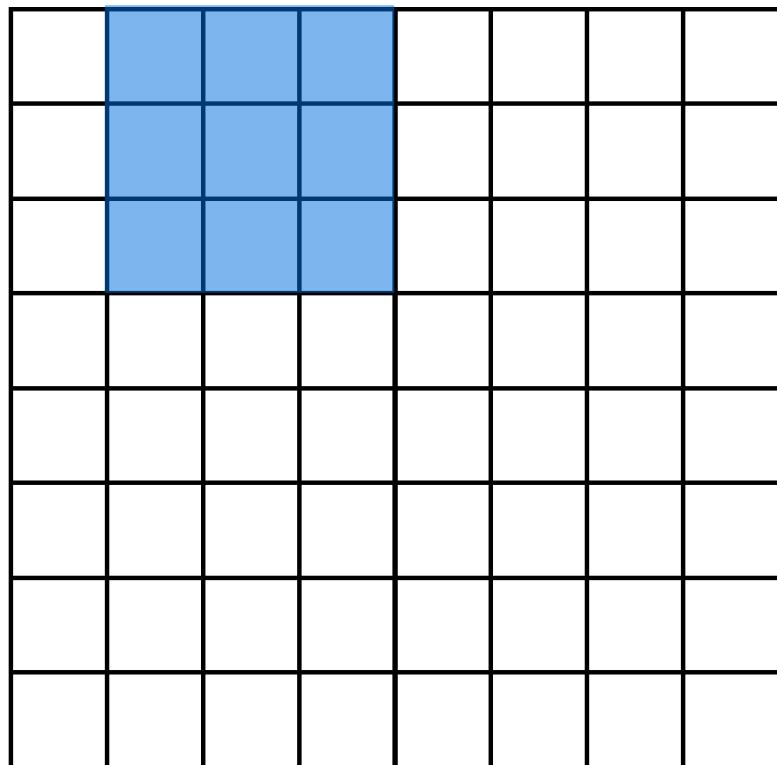
Convolution: Stride

During convolution, the weights “slide” along the input to generate each output

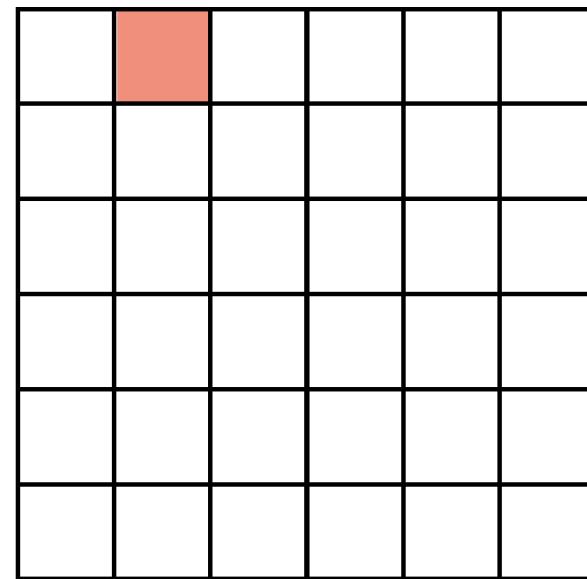


Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



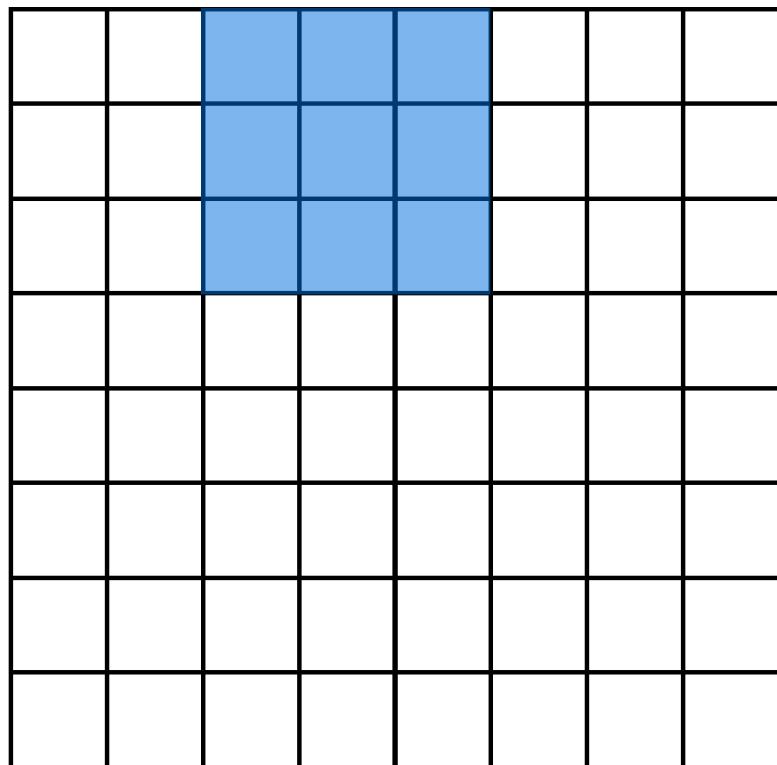
Input



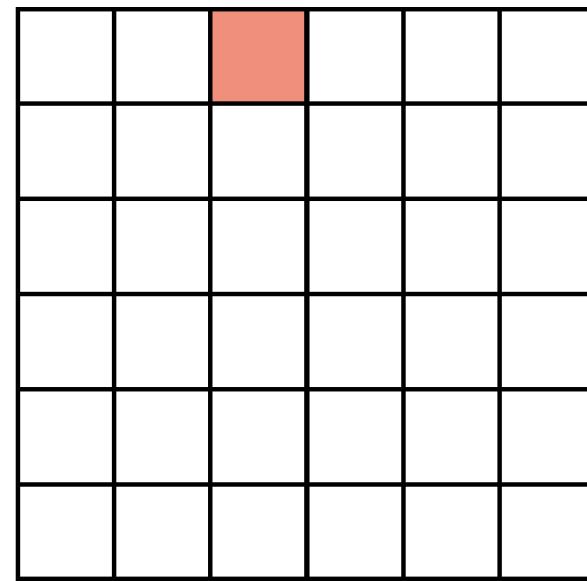
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



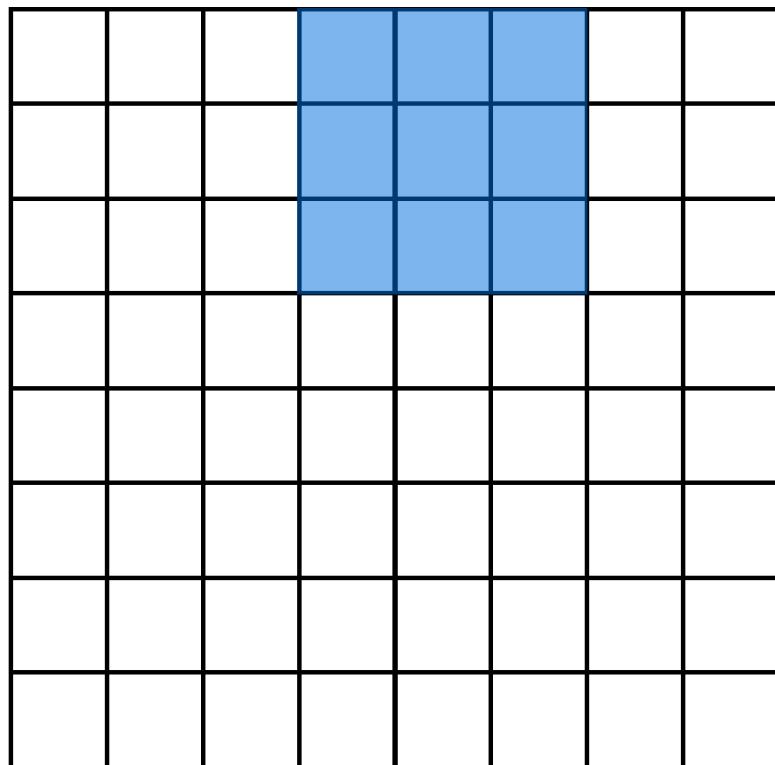
Input



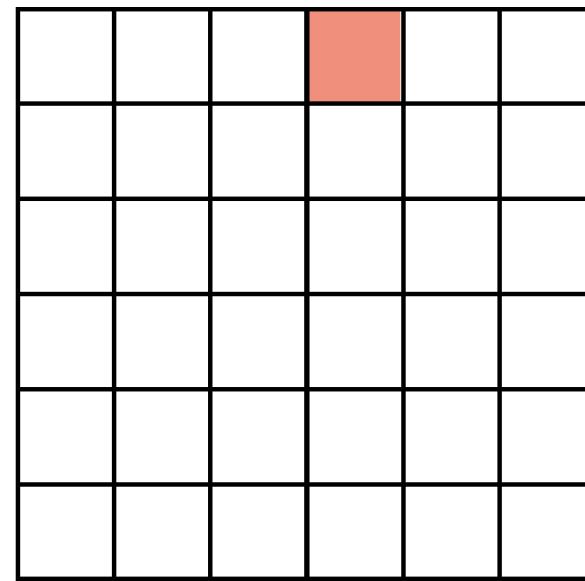
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



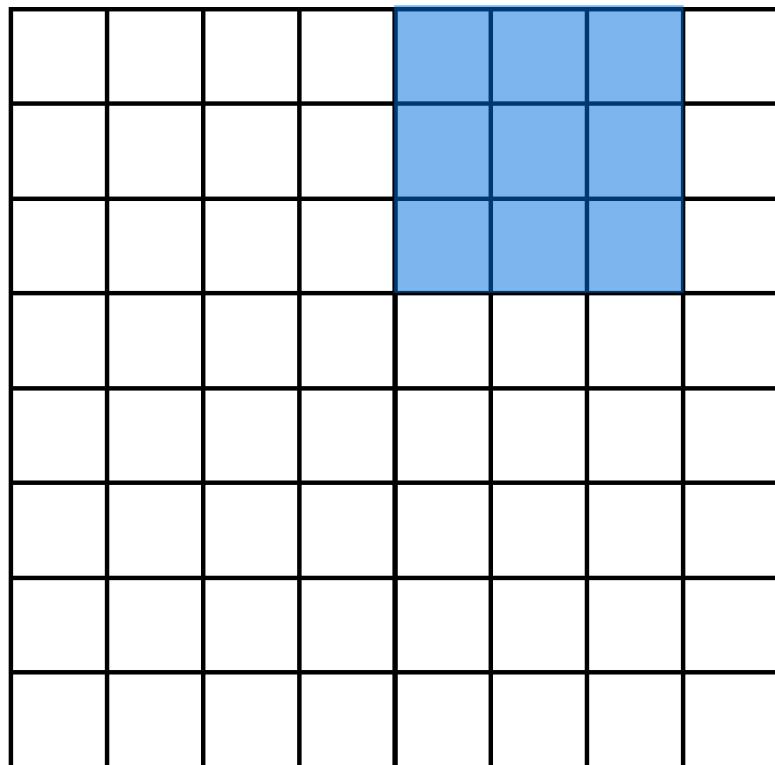
Input



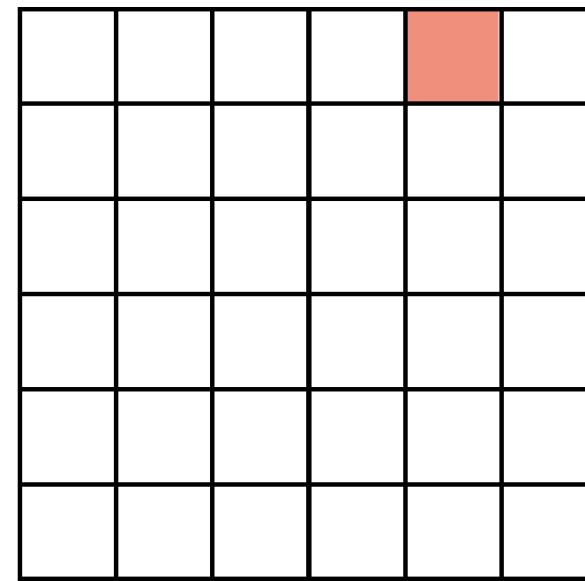
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



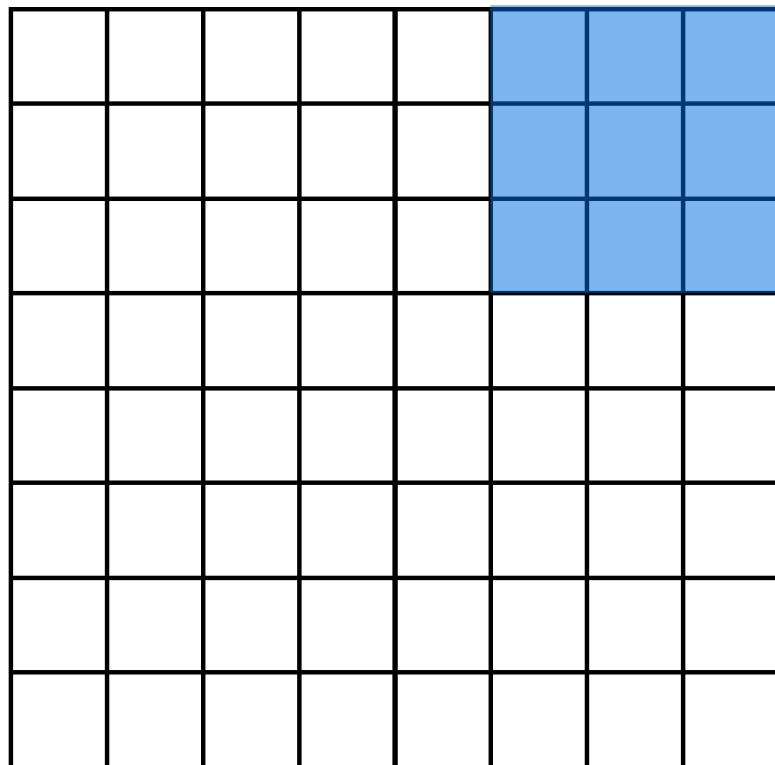
Input



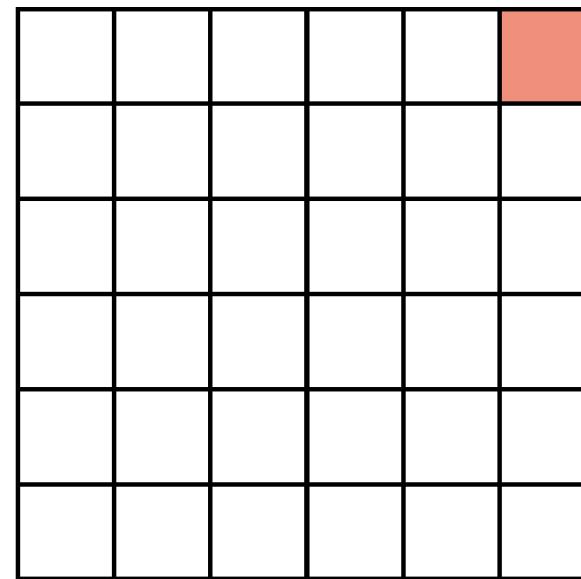
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



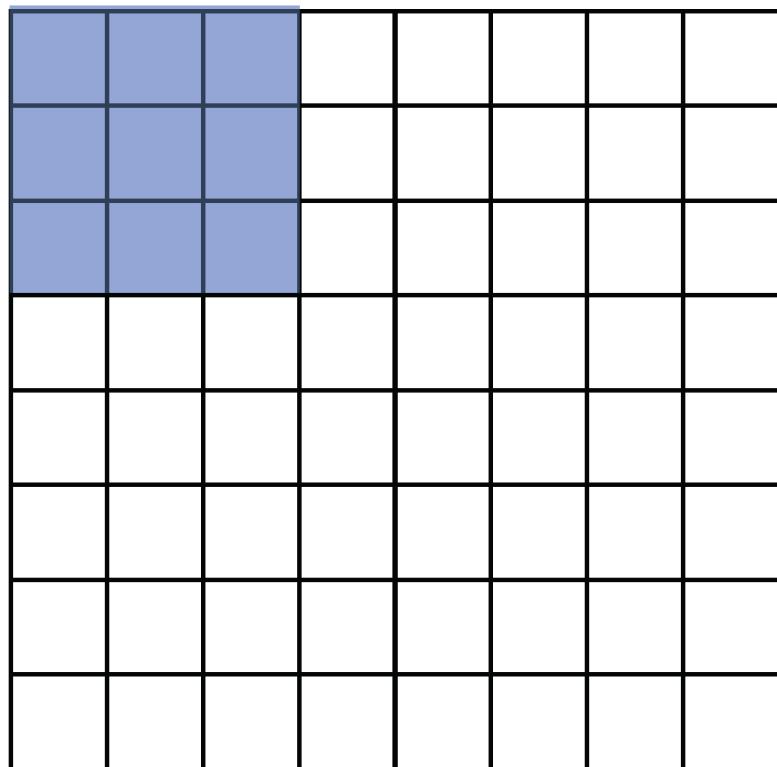
Input



Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



Input

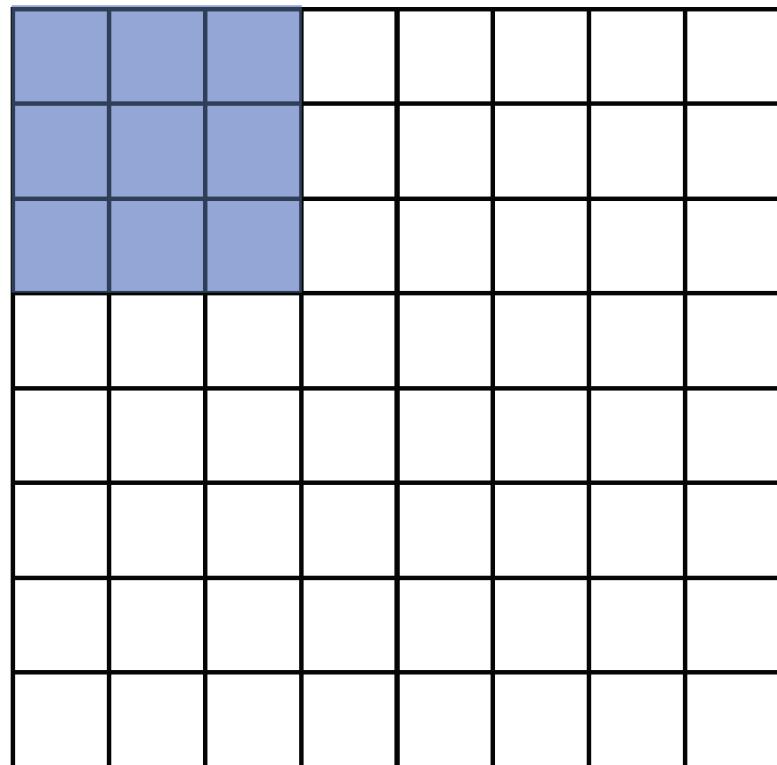
Recall that at each position,
we are doing a **3D** sum:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

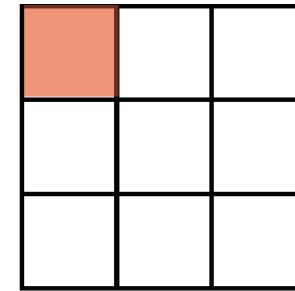
(channel, row, column)

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



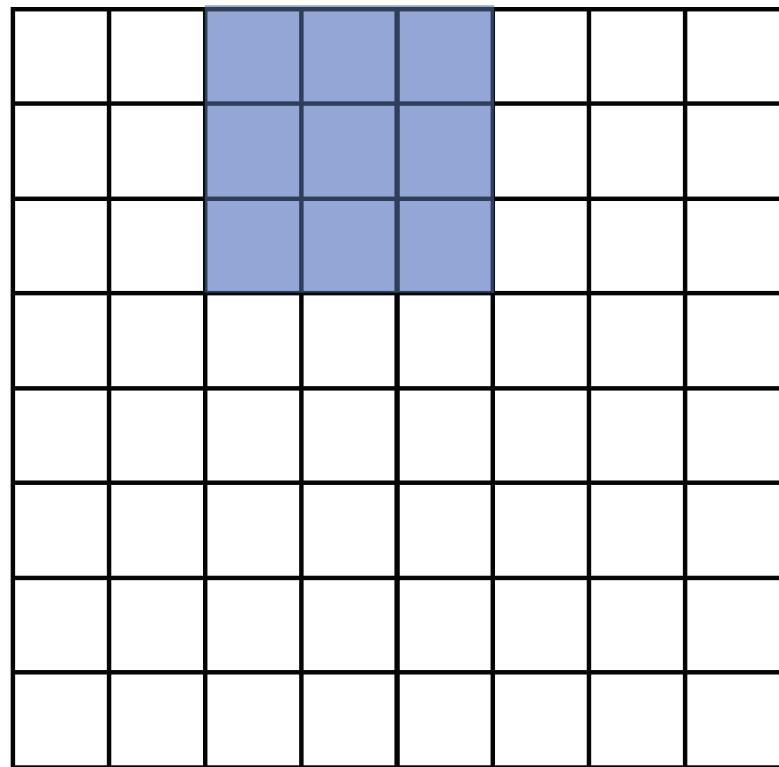
Input



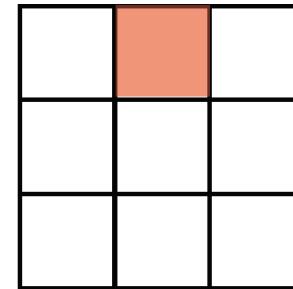
Output

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



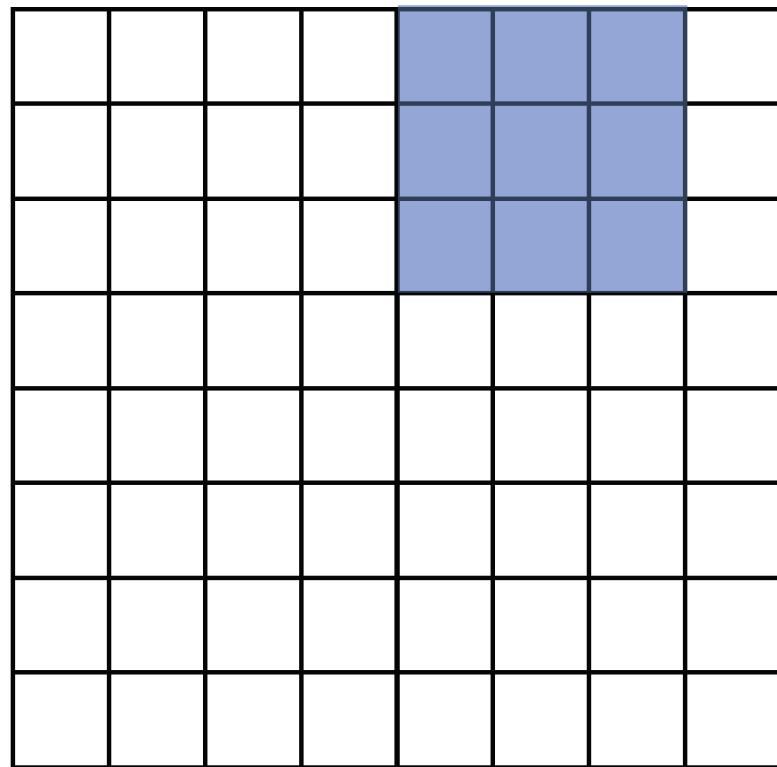
Input



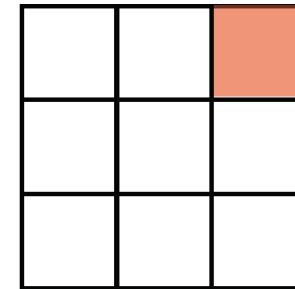
Output

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



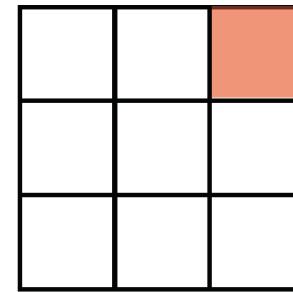
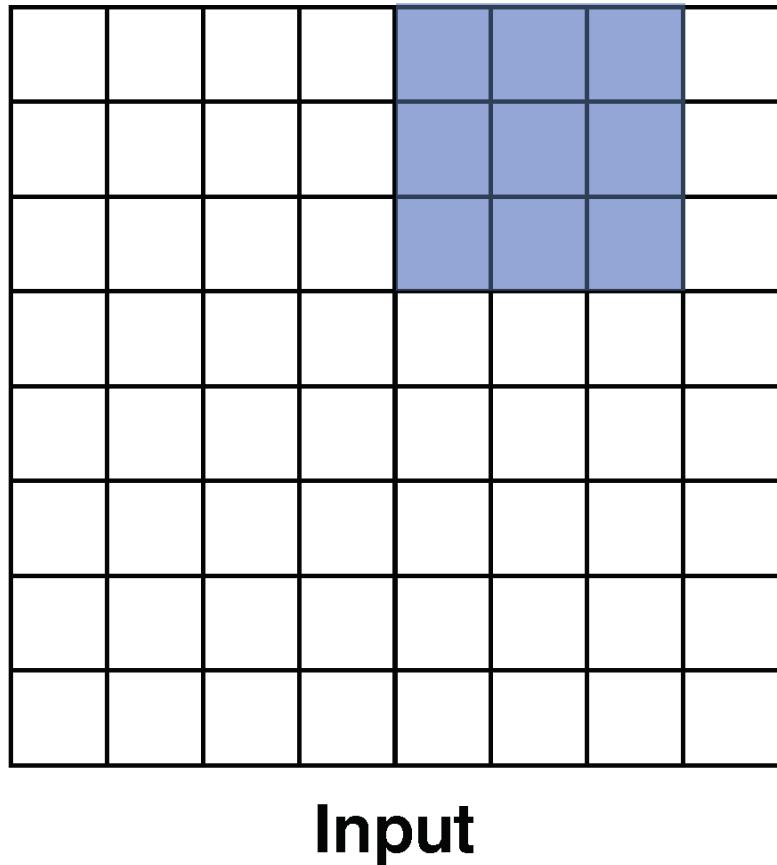
Input



Output

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



- *Notice that with certain strides, we may not be able to cover all of the input*
- *The output is also half the size of the input*

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1, stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1, stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1, stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1, stride = 2**

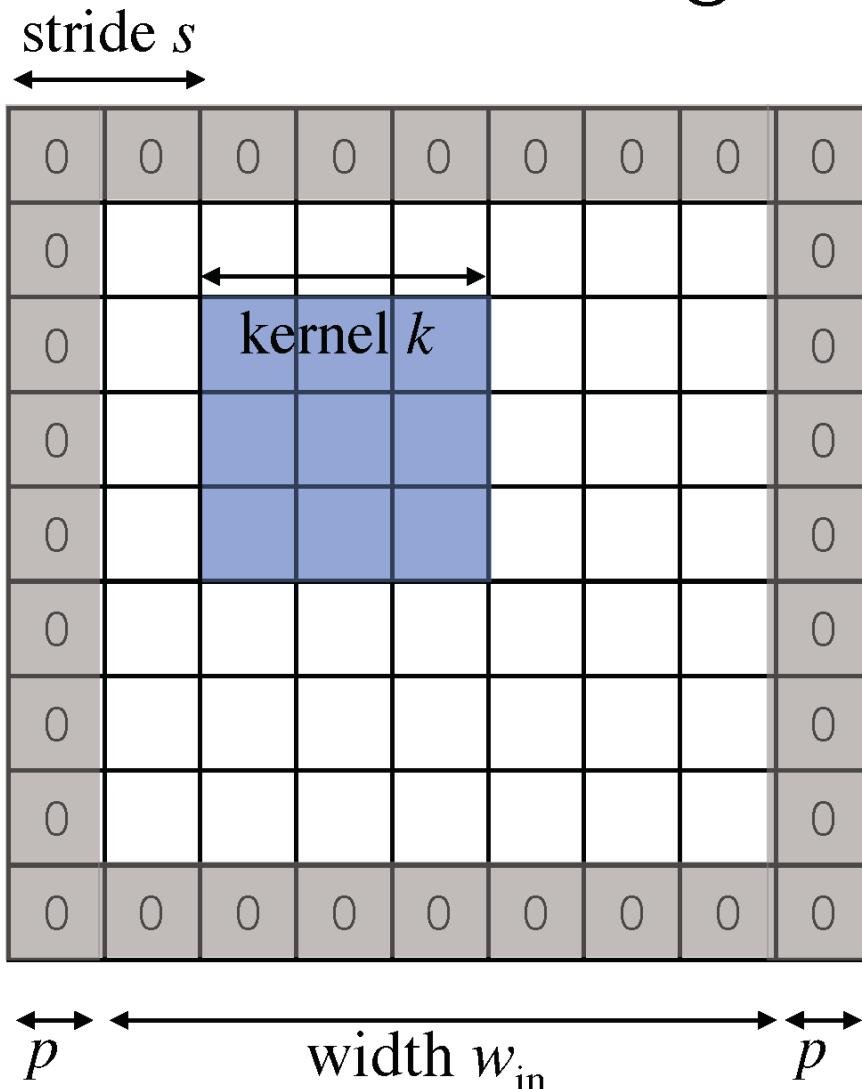
0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution:

How big is the output?

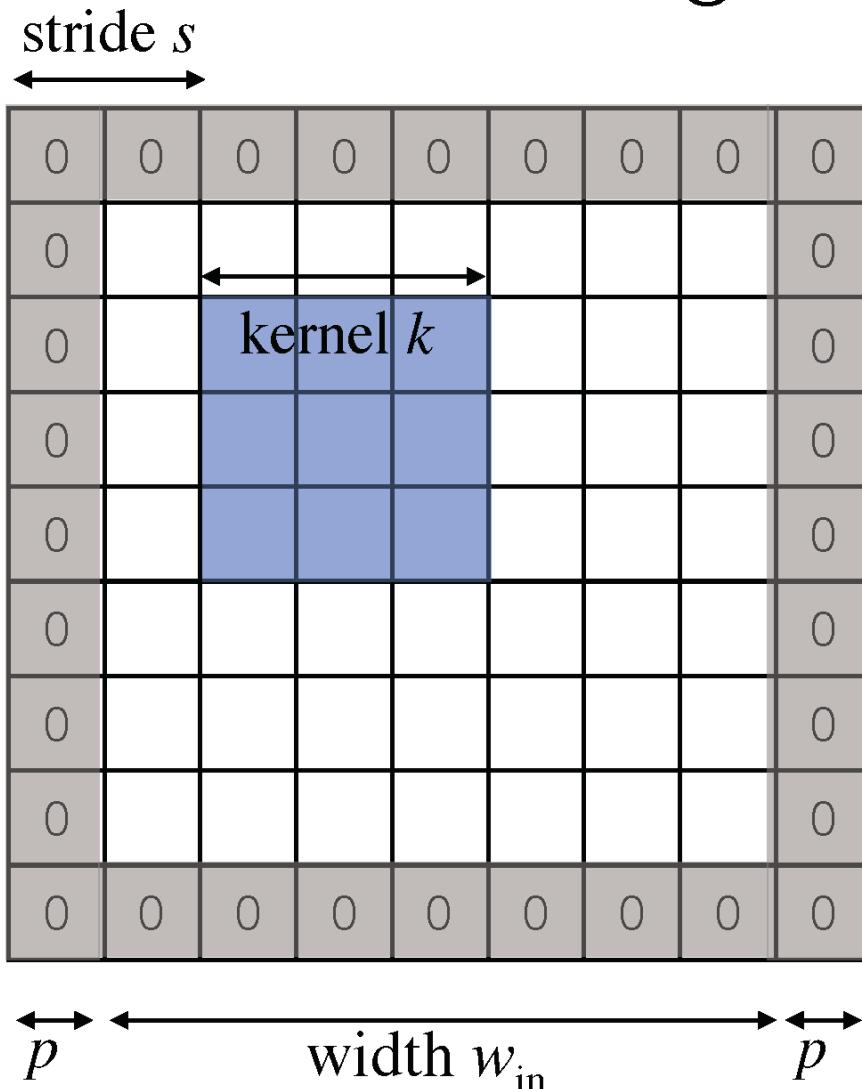


In general, the output has size:

$$w_{\text{out}} = \left\lfloor \frac{w_{\text{in}} + 2p - k}{s} \right\rfloor + 1$$

Convolution:

How big is the output?



Example: $k=3$, $s=1$, $p=1$

$$\begin{aligned}w_{out} &= \left\lfloor \frac{w_{in} + 2p - k}{s} \right\rfloor + 1 \\&= \left\lfloor \frac{w_{in} + 2 - 3}{1} \right\rfloor + 1 \\&= w_{in}\end{aligned}$$

VGGNet [Simonyan 2014]
uses filters of this shape

Pooling

For most ConvNets, **convolution** is often followed by **pooling**:

- Creates a smaller representation while retaining the most important information
- The “max” operation is the most common
- Why might “avg” be a poor choice?

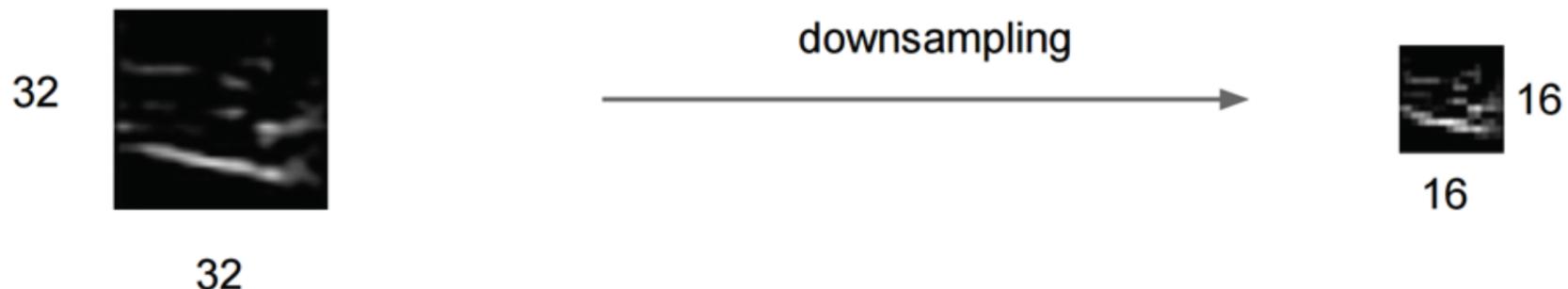
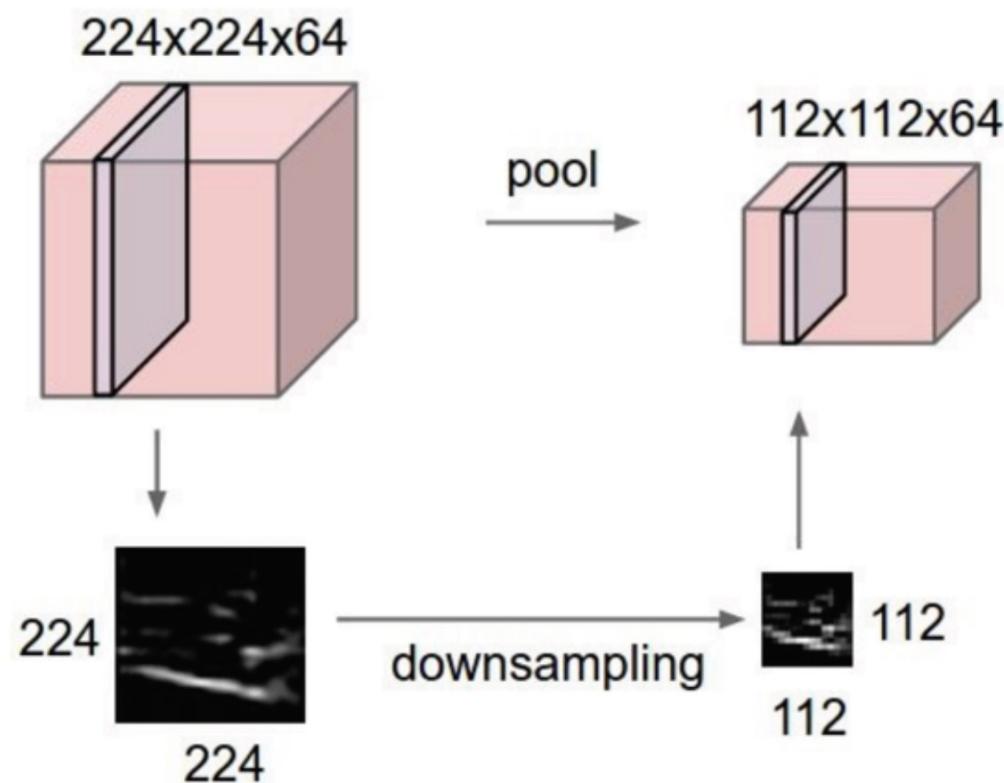


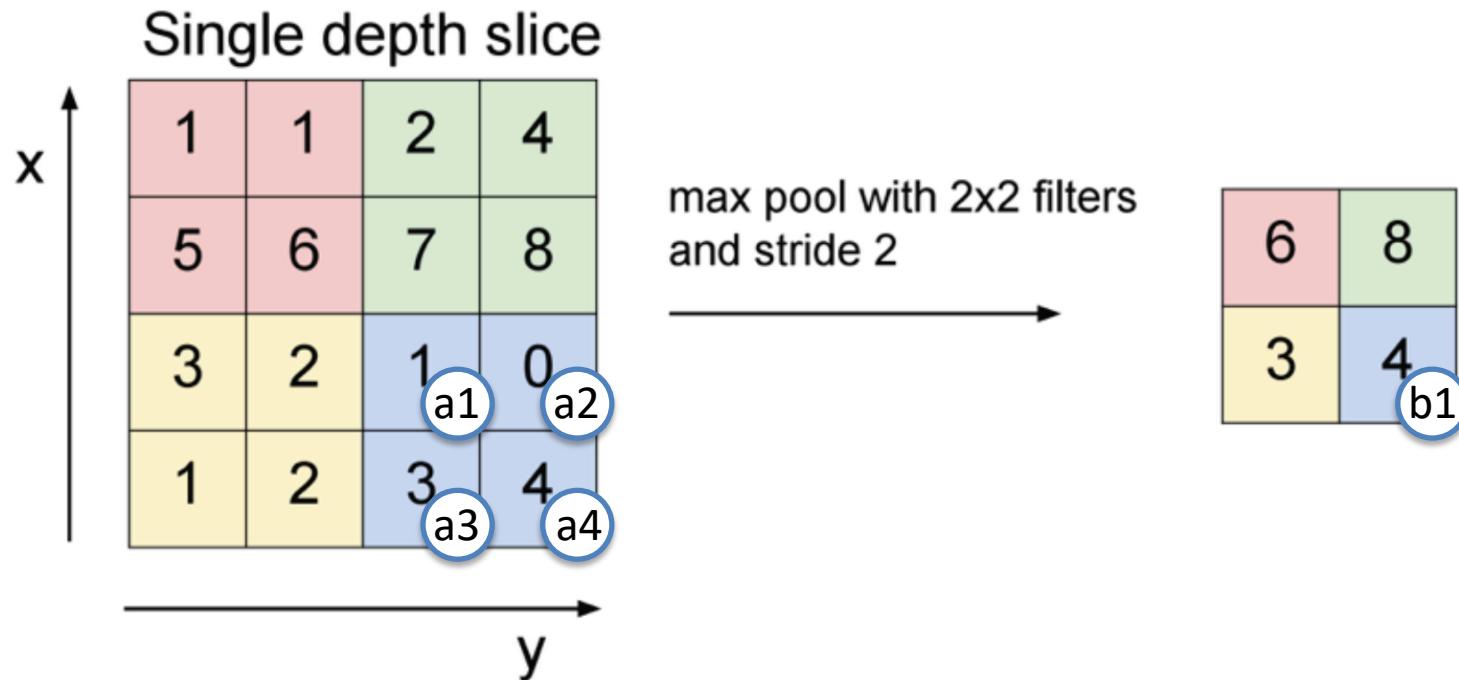
Figure: Andrej Karpathy

Pooling

- makes the representations smaller and more manageable
- operates over each activation map independently:



Max Pooling

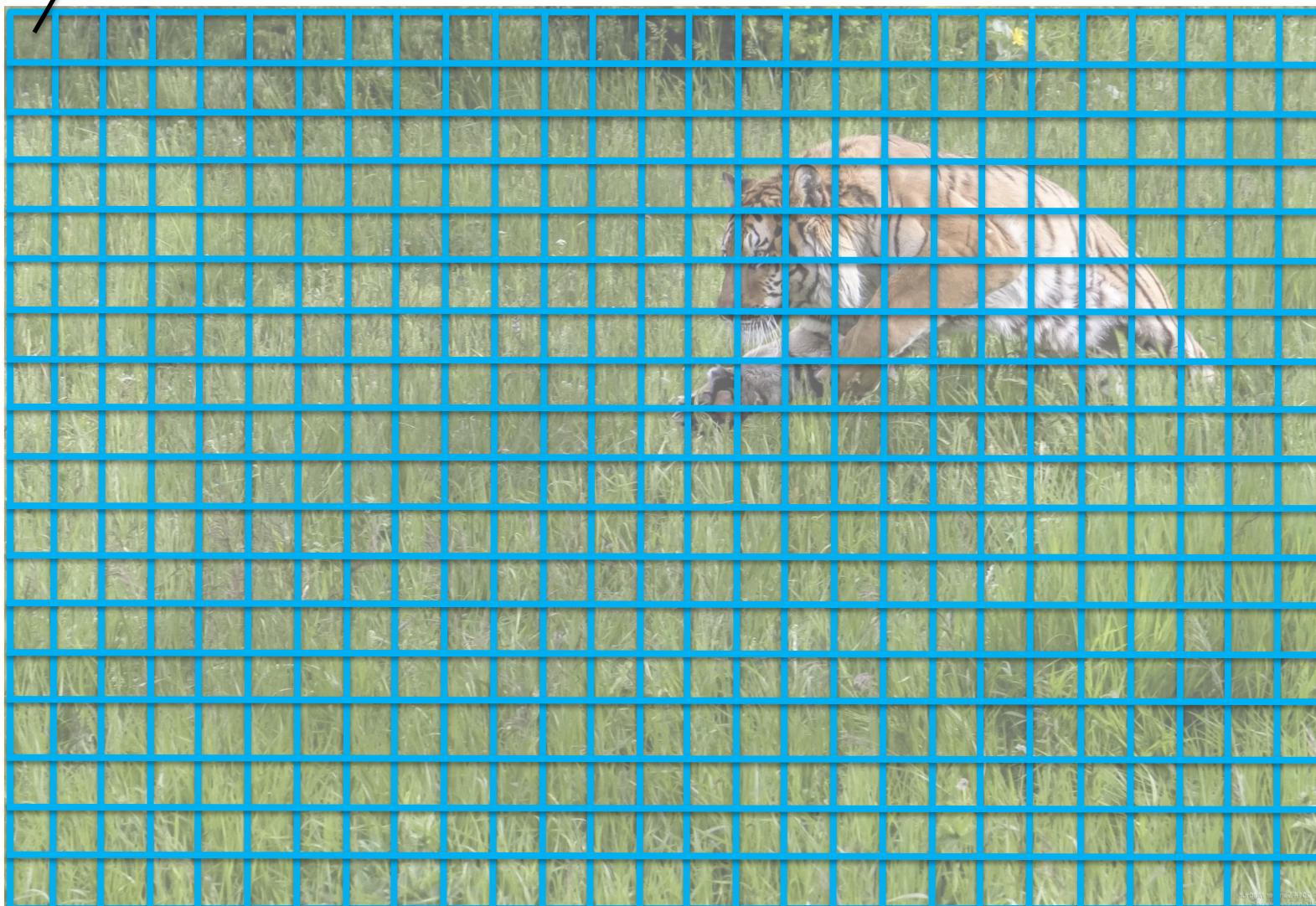


What's the backprop rule for max pooling?

- In the forward pass, store the index that took the max
- The backprop gradient is the input gradient at that index

Figure: Andrej Karpathy

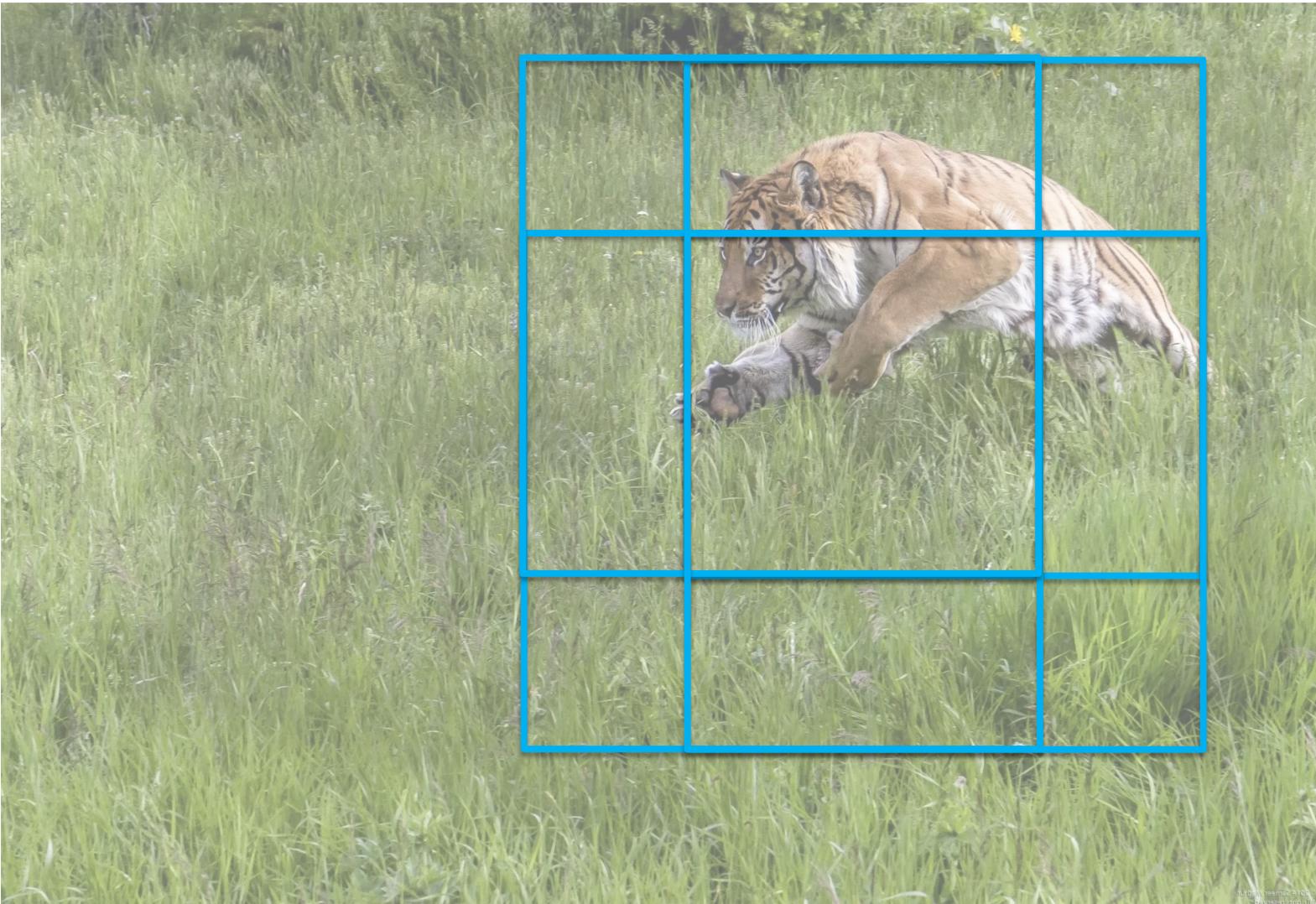
Max Pooling



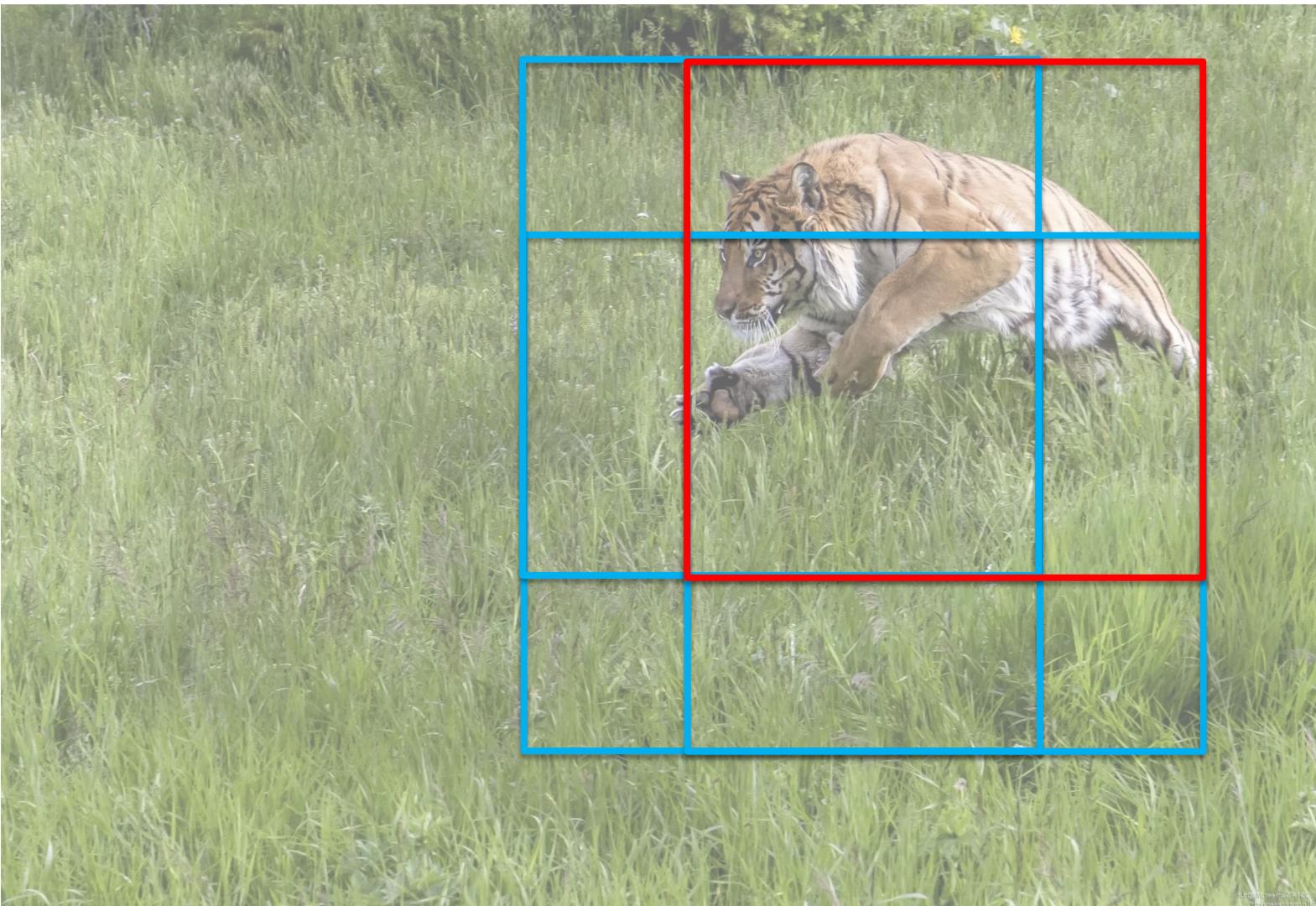
Max Pooling



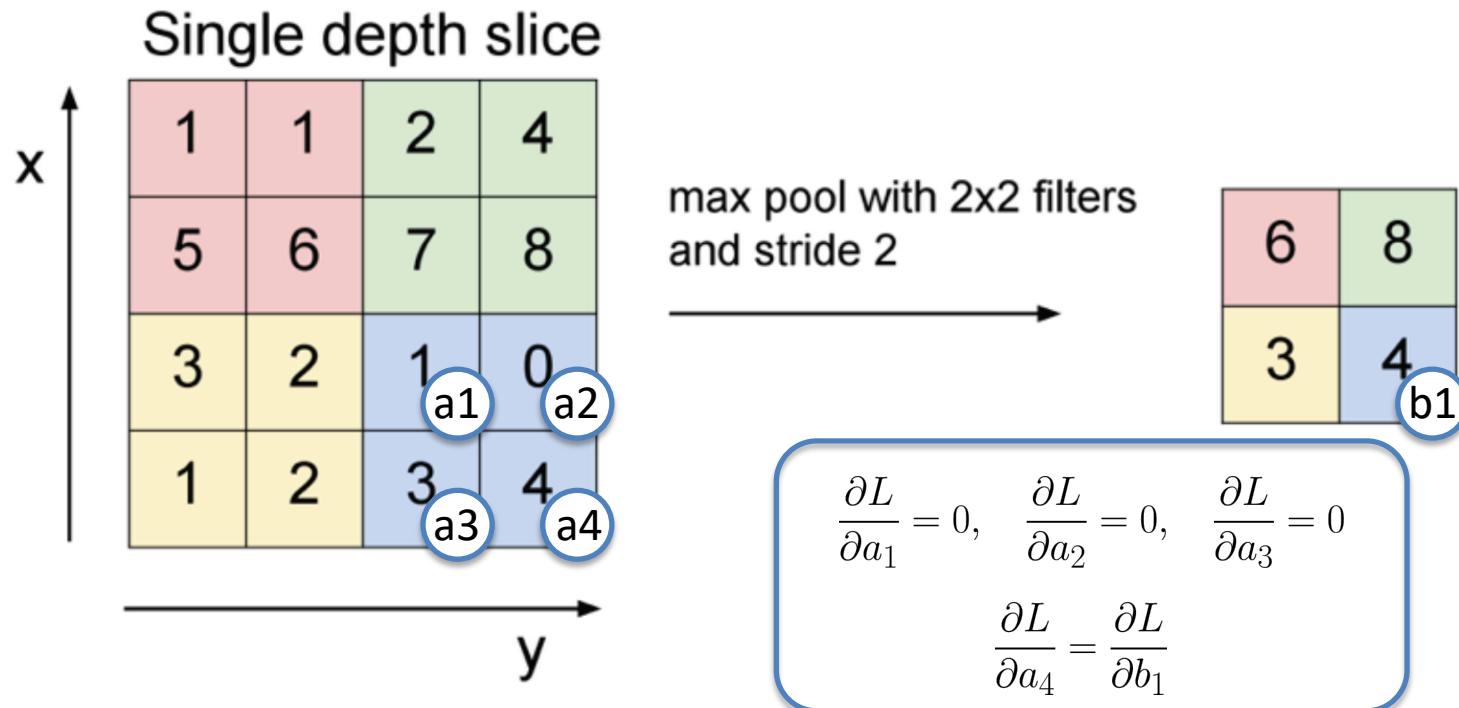
Max Pooling



Max Pooling



Max Pooling

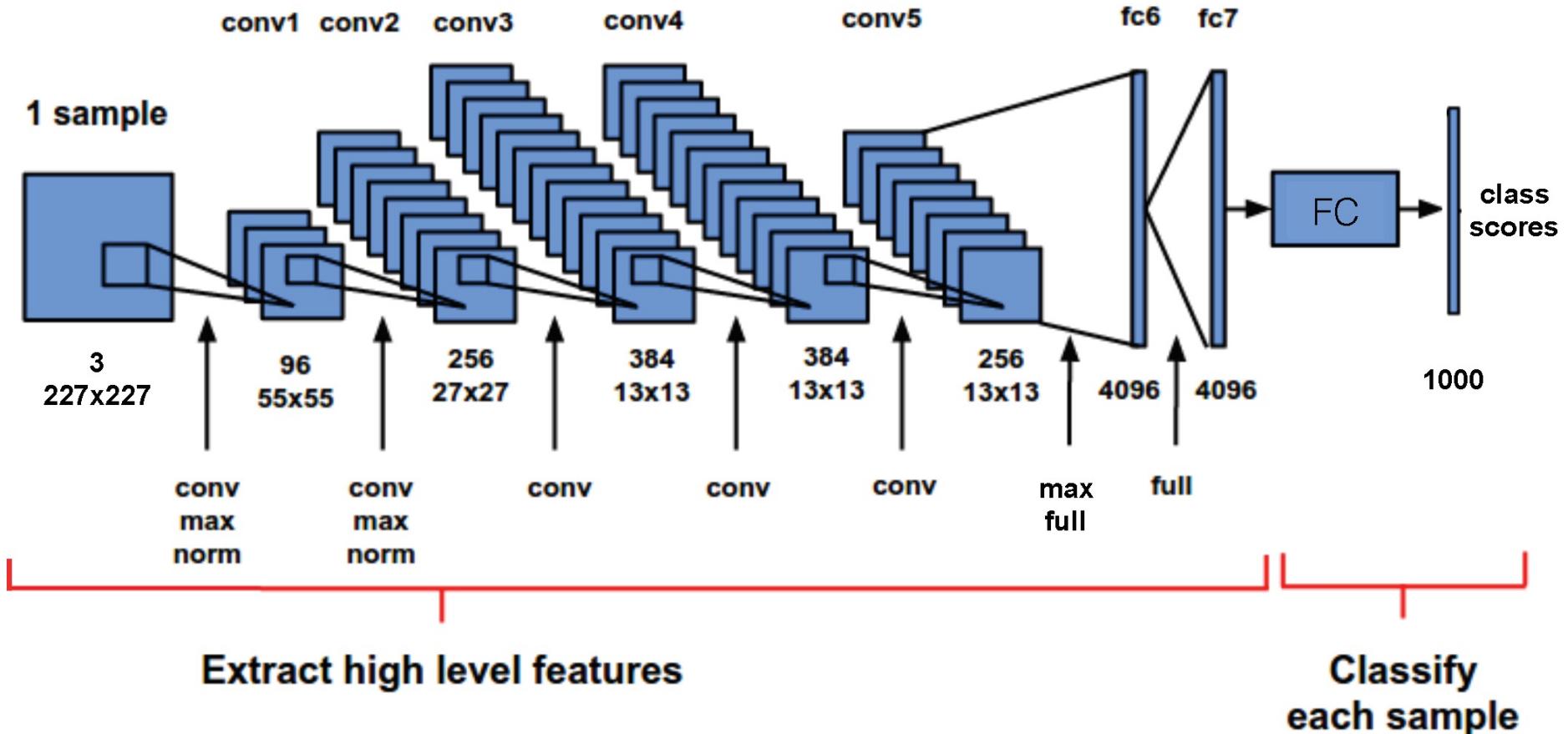


What's the backprop rule for max pooling?

- In the forward pass, store the index that took the max
- The backprop gradient is the input gradient at that index

Figure: Andrej Karpathy

Example: AlexNet [Krizhevsky 2012]



“max”: max pooling

“norm”: local response normalization

“full”: fully connected

Figure: [Karnowski 2015] (with corrections)

Going deeper with convolutions

Christian Szegedy

Google Inc.

Wei Liu

University of North Carolina, Chapel Hill

Yangqing Jia

Google Inc.

Pierre Sermanet

Google Inc.

Scott Reed

University of Michigan

Dragomir Anguelov

Google Inc.

Dumitru Erhan

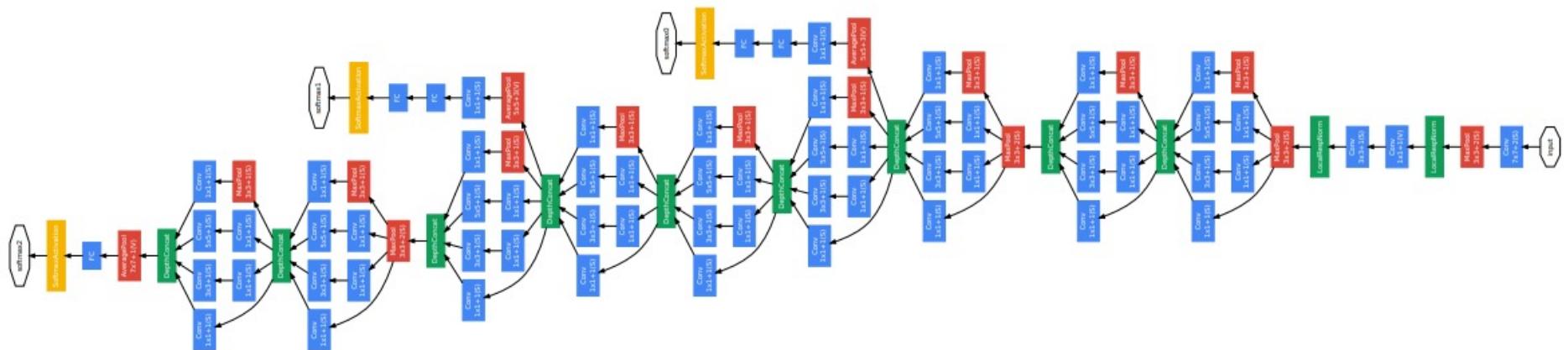
Google Inc.

Vincent Vanhoucke

Google Inc.

Andrew Rabinovich

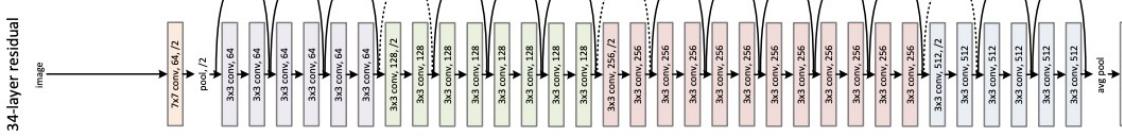
Google Inc.



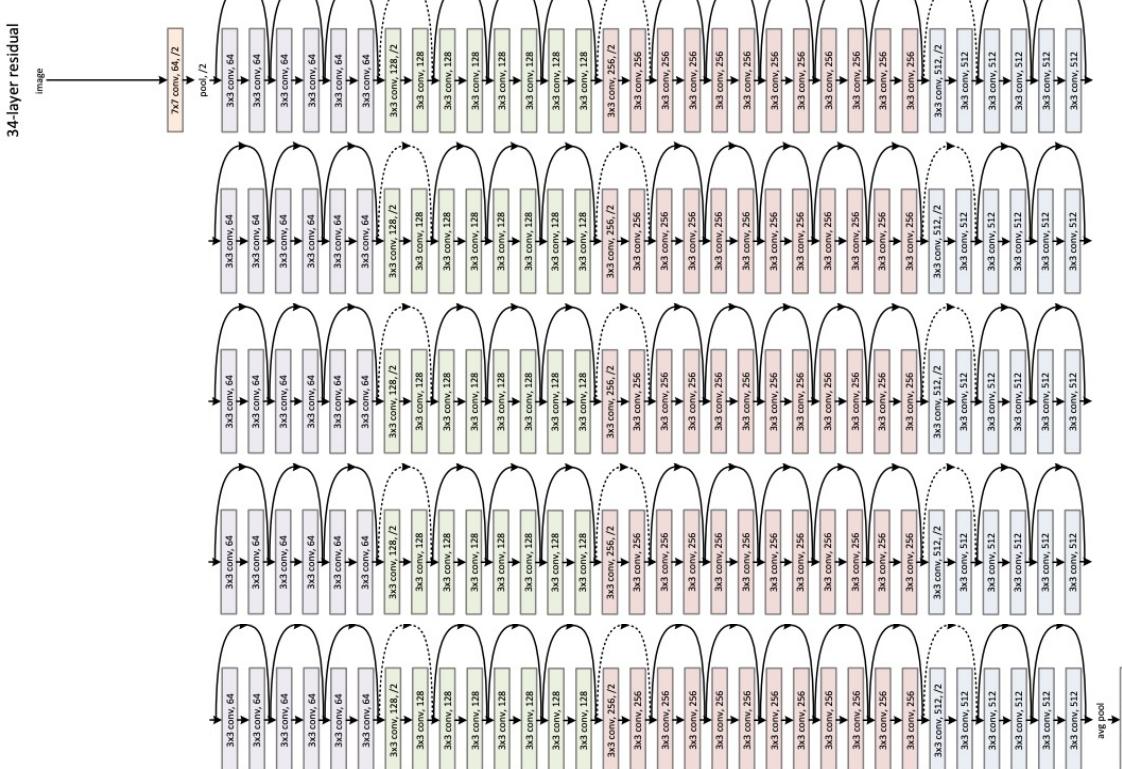
22 layers, but 12 times fewer parameters than AlexNet

2015

34 layers ResNet



152 layers ResNet



Example ConvNet

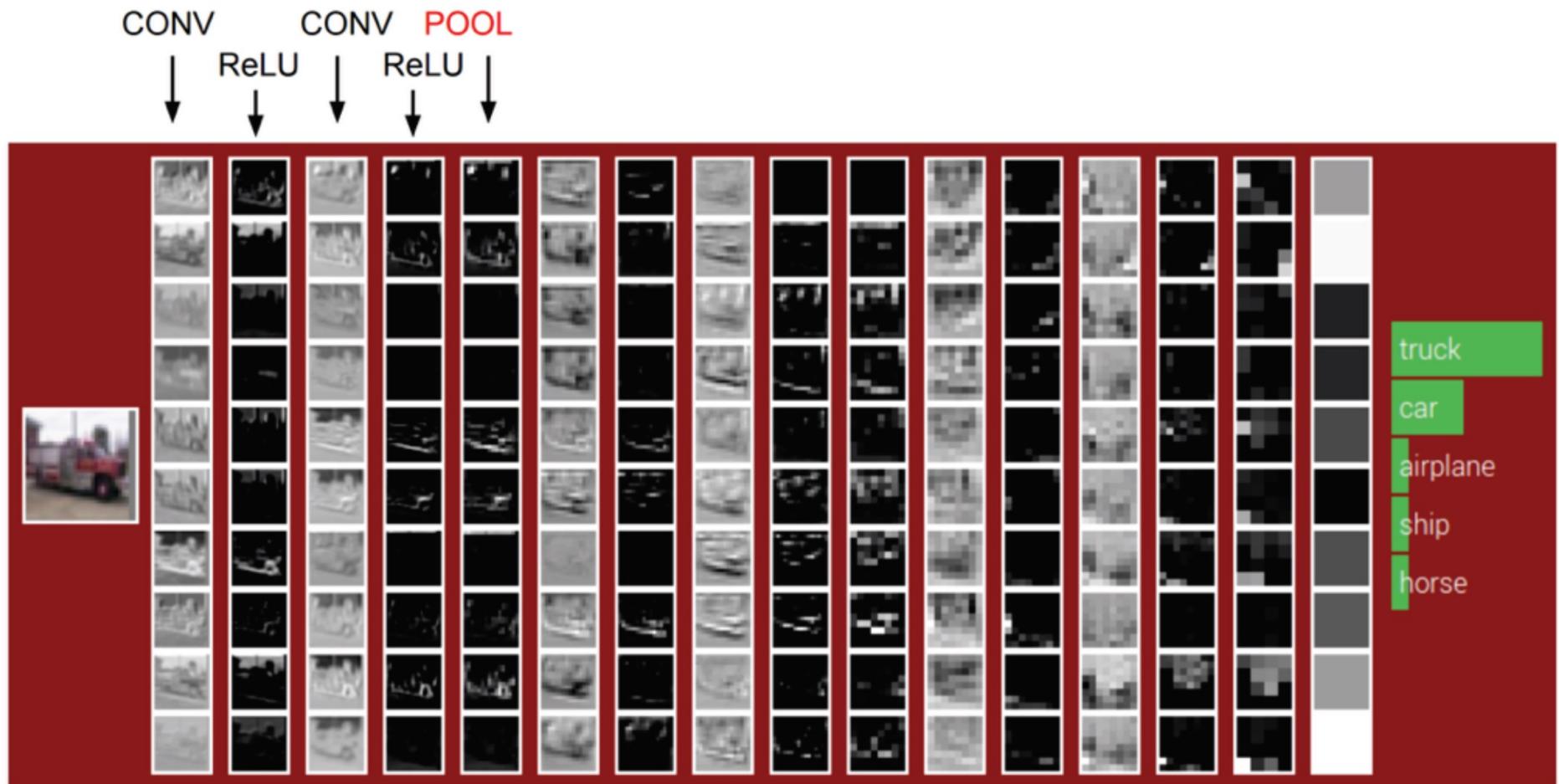


Figure: Andrej Karpathy

Example ConvNet

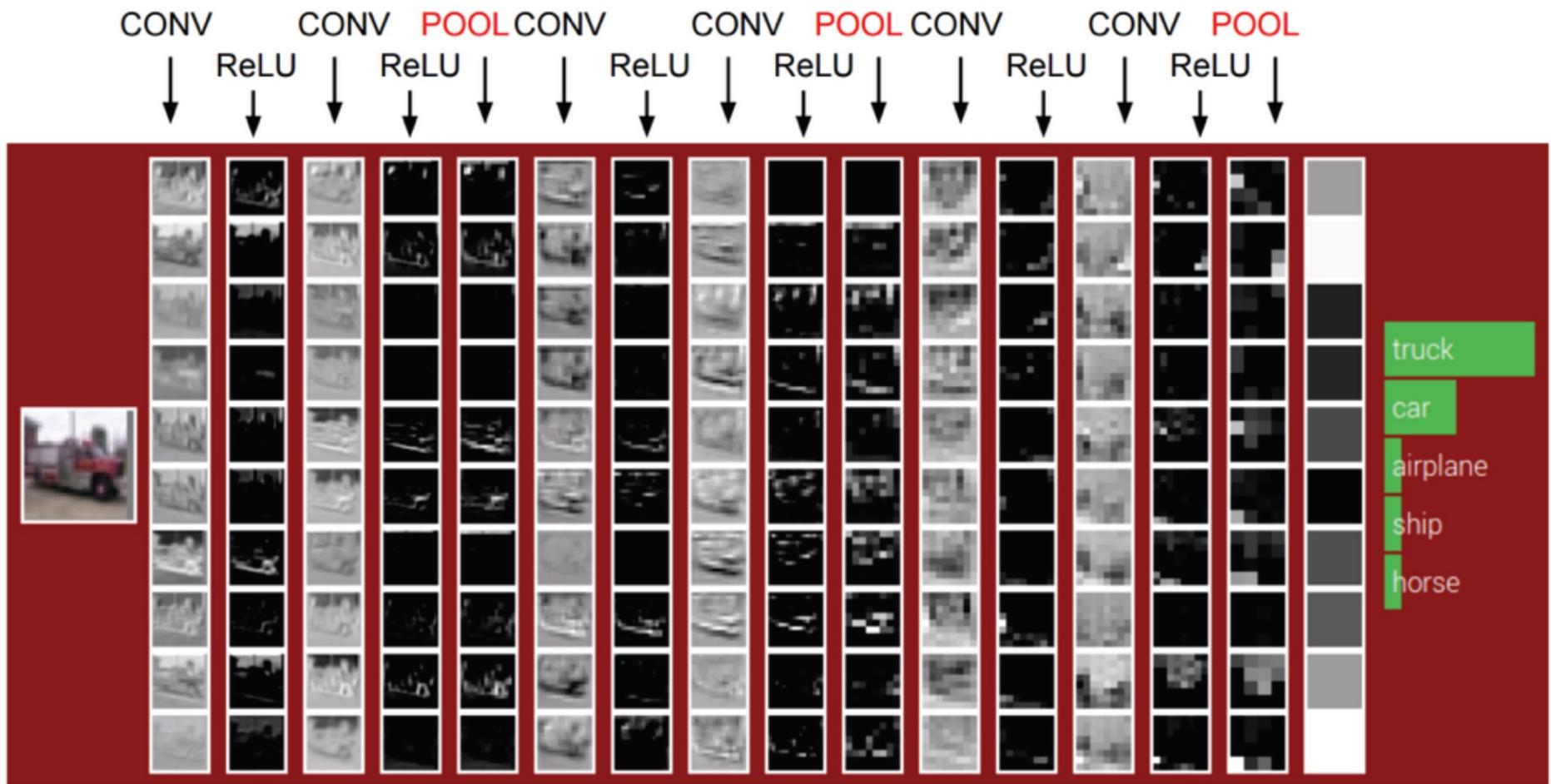


Figure: Andrej Karpathy

Example ConvNet

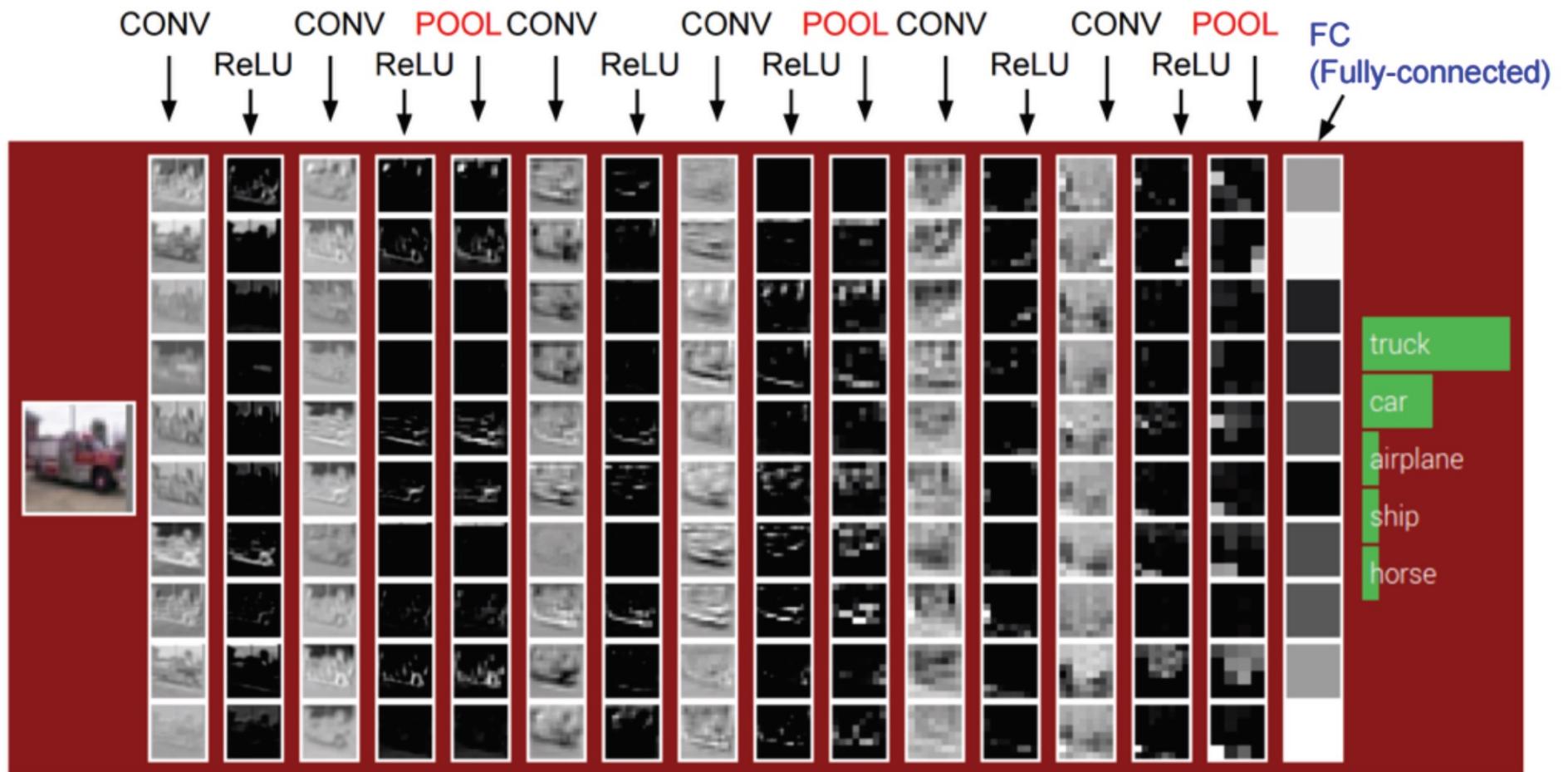
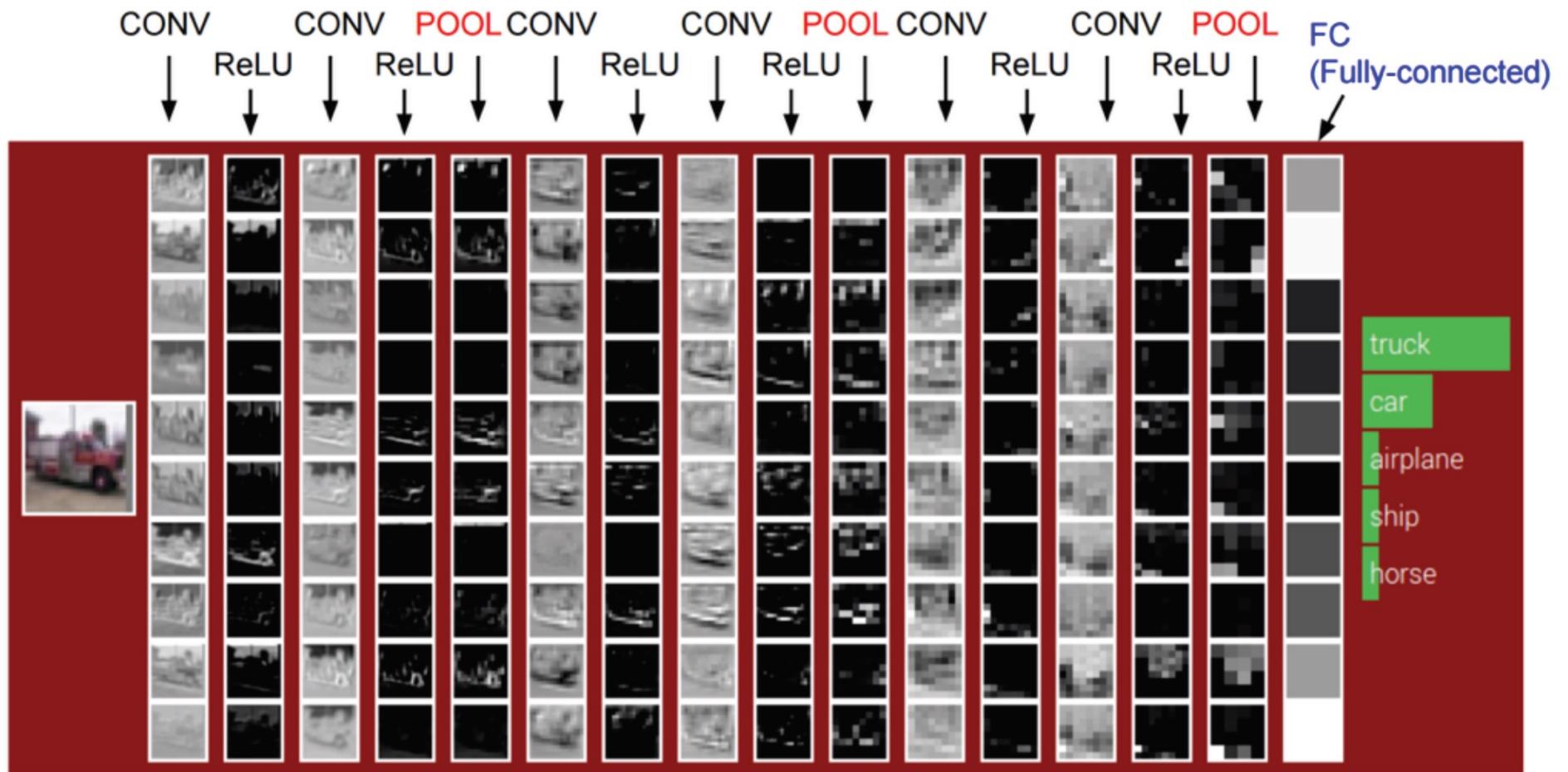


Figure: Andrej Karpathy

Example ConvNet



10x3x3 conv filters, stride 1, pad 1

2x2 pool filters, stride 2

Figure: Andrej Karpathy

Clarifai: ILSVRC 2013 winner

- Refinement of AlexNet

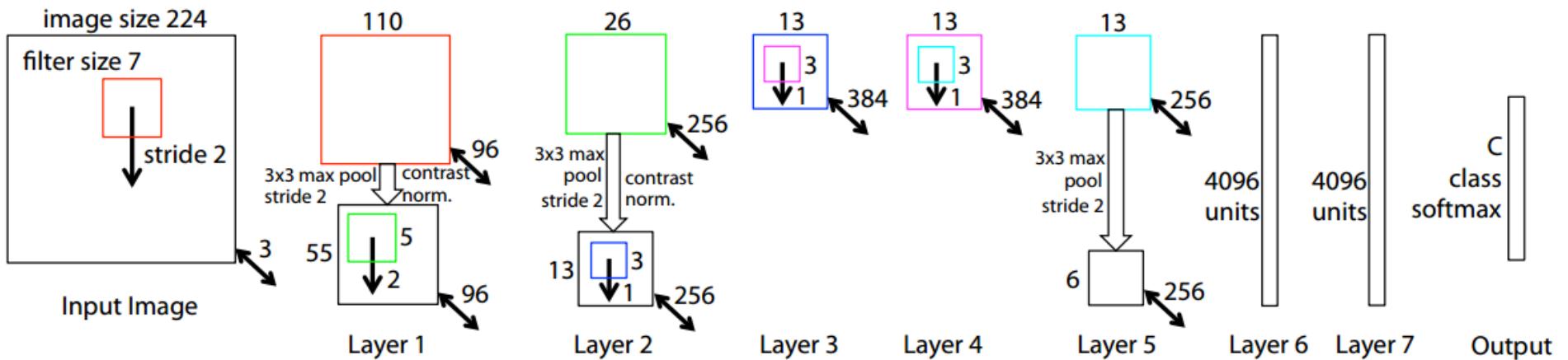
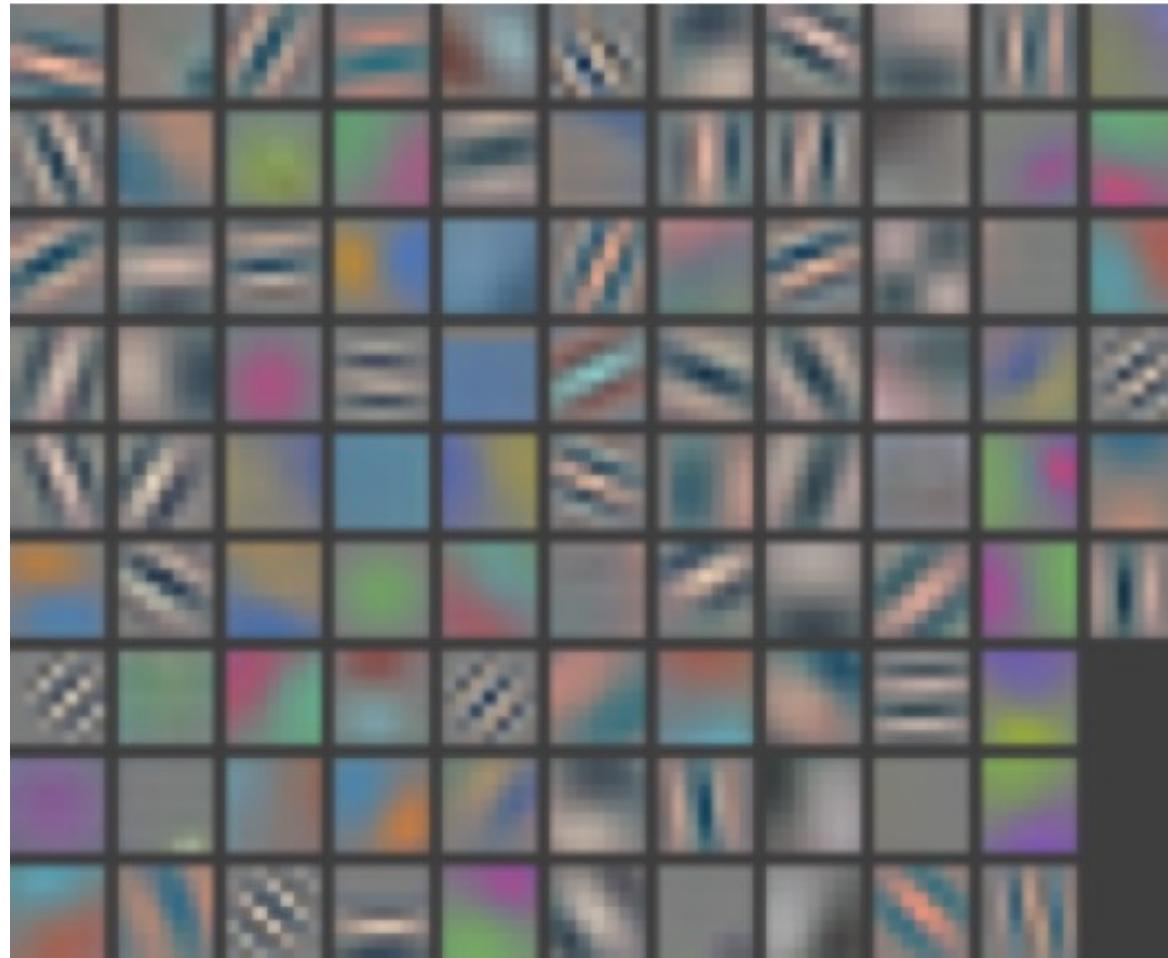


Figure 3. Architecture of our 8 layer convnet model. A 224 by 224 crop of an image (with 3 color planes) is presented as the input. This is convolved with 96 different 1st layer filters (red), each of size 7 by 7, using a stride of 2 in both x and y. The resulting feature maps are then: (i) passed through a rectified linear function (not shown), (ii) pooled (max within 3x3 regions, using stride 2) and (iii) contrast normalized across feature maps to give 96 different 55 by 55 element feature maps. Similar operations are repeated in layers 2,3,4,5. The last two layers are fully connected, taking features from the top convolutional layer as input in vector form ($6 \cdot 6 \cdot 256 = 9216$ dimensions). The final layer is a C -way softmax function, C being the number of classes. All filters and feature maps are square in shape.

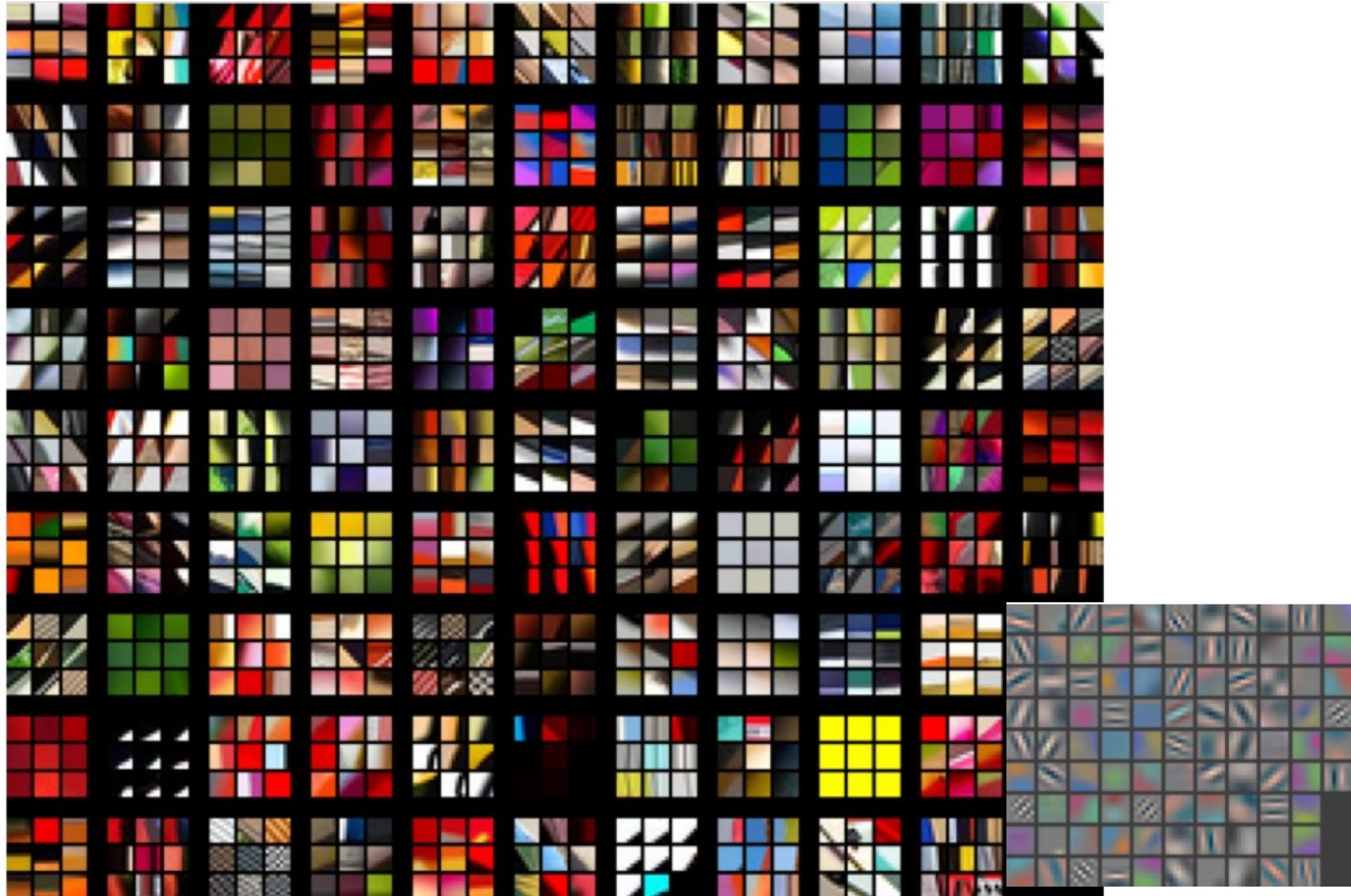
M. Zeiler and R. Fergus, [Visualizing and Understanding Convolutional Networks](#),
ECCV 2014 (Best Paper Award winner)

Layer 1 Filters



M. Zeiler and R. Fergus, [Visualizing and Understanding Convolutional Networks](#),
ECCV 2014 (Best Paper Award winner)

Layer 1: Top-9 Patches

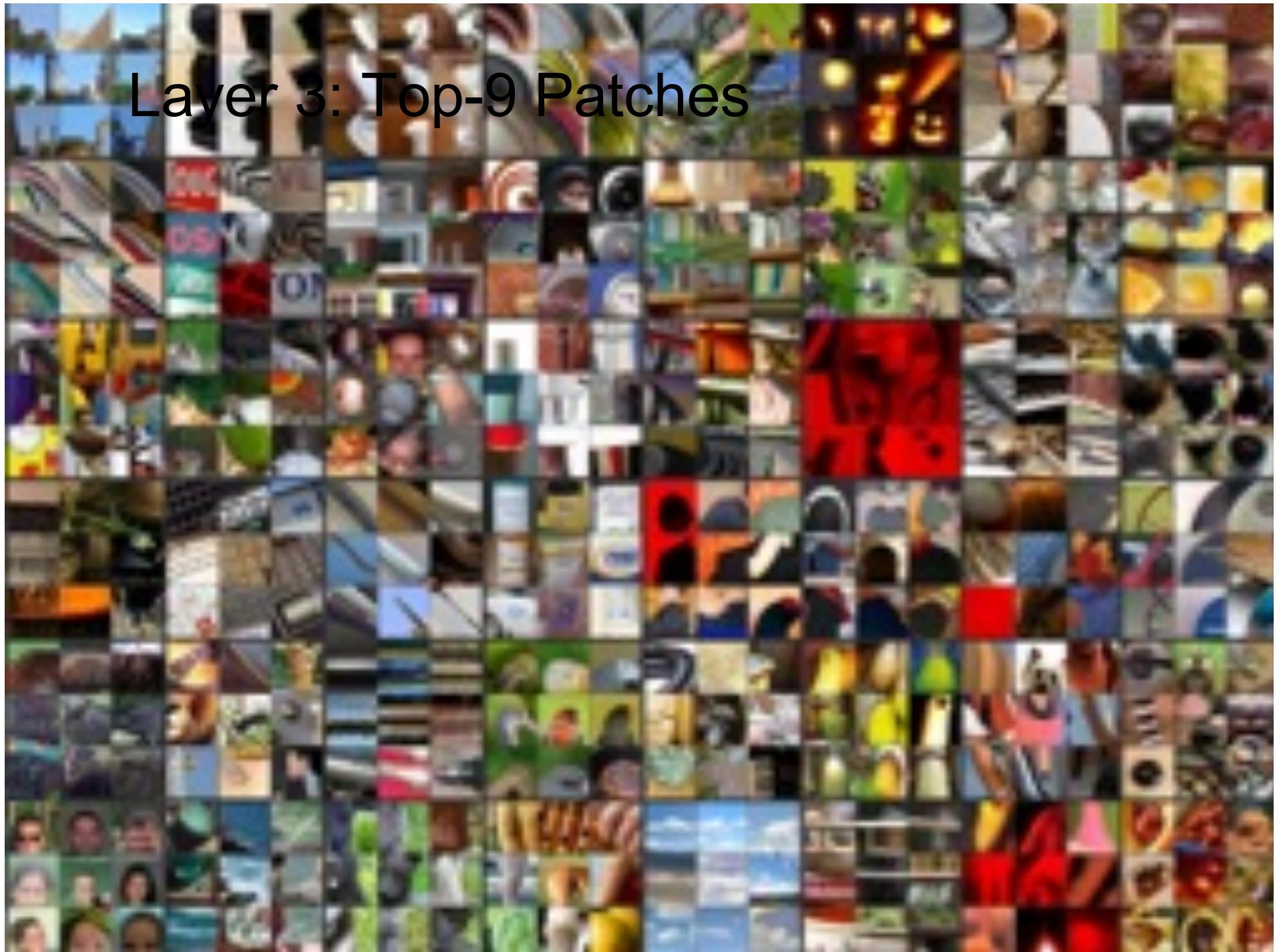


Layer 2: Top-9 Patches



- Patches from validation images that give maximal activation of a given feature map

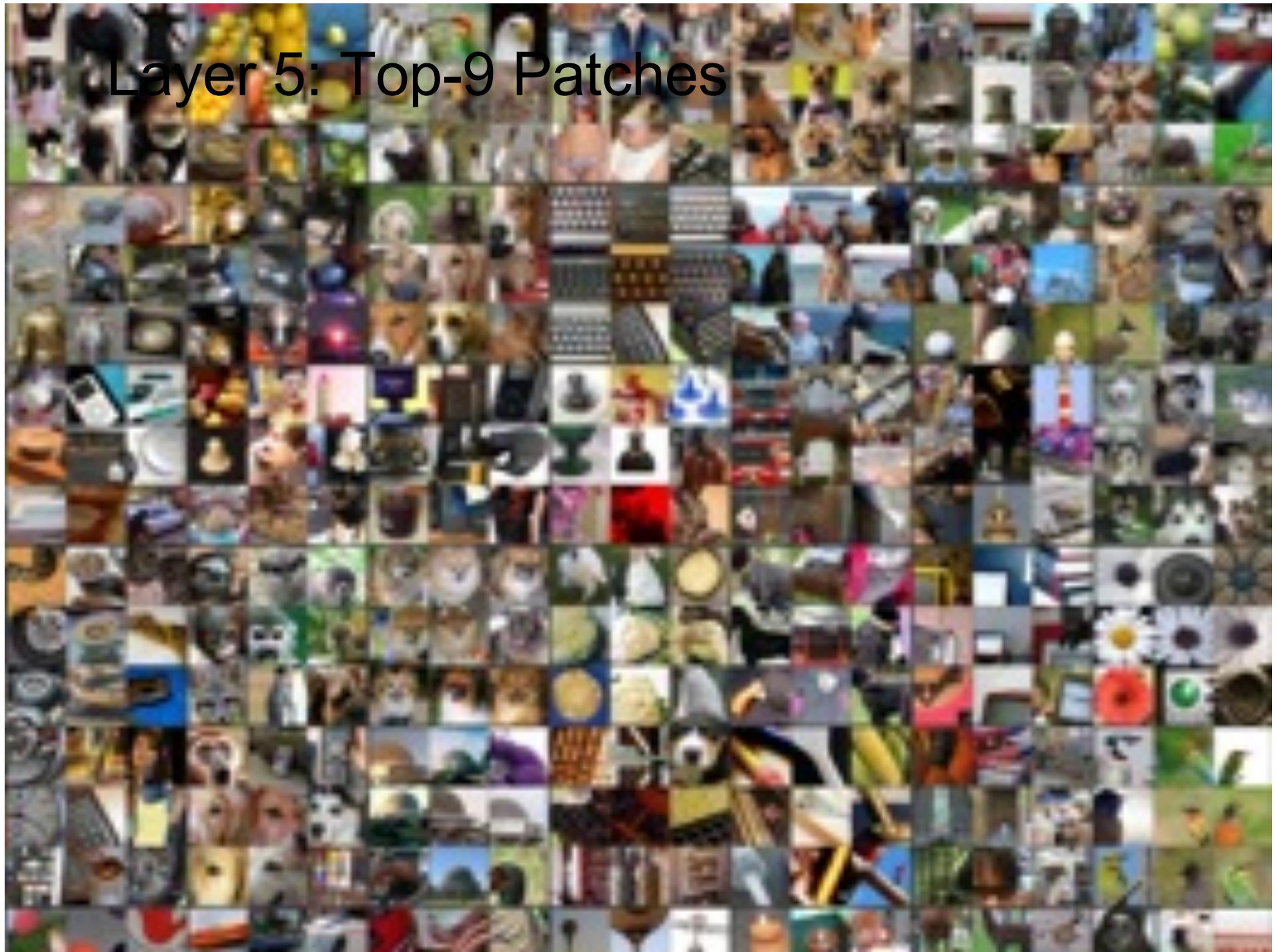
Layer 3: Top-9 Patches



Layer 4: Top-9 Patches



Layer 5: Top-9 Patches



Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification

Kaiming He

Xiangyu Zhang

Shaoqing Ren

Jian Sun

Microsoft Research

{kahe, v-xiangz, v-shren, jiansun}@microsoft.com



GT: horse cart
1: horse cart
2: minibus
3: oxcart
4: stretcher
5: half track



GT: birdhouse
1: birdhouse
2: sliding door
3: window screen
4: mailbox
5: pot



GT: forklift
1: forklift
2: garbage truck
3: tow truck
4: trailer truck
5: go-kart



GT: coucal
1: coucal
2: indigo bunting
3: lorikeet
4: walking stick
5: custard apple



GT: komondor
1: komondor
2: patio
3: llama
4: mobile home
5: Old English sheepdog



GT: yellow lady's slipper
1: yellow lady's slipper
2: slug
3: hen-of-the-woods
4: stinkhorn
5: coral fungus



GT: torch
1: stage
2: spotlight
3: torch
4: microphone
5: feather boa



GT: banjo
1: acoustic guitar
2: shoji
3: bow tie
4: cowboy hat
5: banjo



GT: go-kart
1: go-kart
2: crash helmet
3: racer
4: sports car
5: motor scooter



Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification

Kaiming He

Xiangyu Zhang

Shaoqing Ren

Jian Sun

Microsoft Research

{kahe, v-xiangz, v-shren, jiansun}@microsoft.com



GT: horse cart
1: horse cart
2: minibus
3: oxcart
4: stretcher
5: half track



GT: birdhouse
1: birdhouse
2: sliding door
3: window screen
4: mailbox
5: pot



GT: forklift
1: forklift
2: garbage truck
3: tow truck
4: trailer truck
5: go-kart



GT: mountain tent
1: sleeping bag
2: mountain tent
3: parachute
4: ski
5: flagpole



GT: geyser
1: geyser
2: volcano
3: sandbar
4: breakwater
5: leatherback turtle



GT: microwave
1: microwave
2: washer
3: toaster
4: stove
5: dishwasher



GT: coucal
1: coucal
2: indigo bunting
3: lorikeet
4: walking stick
5: custard apple



GT: komondor
1: komondor
2: patio
3: llama
4: mobile home
5: Old English sheepdog



GT: yellow lady's slipper
1: yellow lady's slipper
2: slug
3: hen-of-the-woods
4: stinkhorn
5: coral fungus



GT: sunscreen
1: hair spray
2: ice lolly
3: sunscreen
4: water bottle
5: lotion



GT: flute
1: flute
2: oboe
3: panpipe
4: trombone
5: bassoon



GT: wooden spoon
1: wok
2: frying pan
3: spatula
4: wooden spoon
5: hot pot



GT: torch
1: stage
2: spotlight
3: torch
4: microphone
5: feather boa



GT: banjo
1: acoustic guitar
2: shoji
3: bow tie
4: cowboy hat
5: banjo



GT: go-kart
1: go-kart
2: crash helmet
3: racer
4: sports car
5: motor scooter

Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification

Kaiming He

Xiangyu Zhang

Shaoqing Ren

Jian Sun

Microsoft Research

{kahe, v-xiangz, v-shren, jiansun}@microsoft.com



(Error %)

	team	top-5 (test)
in competition ILSVRC 14	MSRA, SPP-nets [11]	8.06
	VGG [25]	7.32
	GoogLeNet [29]	6.66
post-competition	VGG [25] (arXiv v5)	6.8
	Baidu [32]	5.98
	MSRA, PReLU-nets	4.94

Table 7. The **multi-model** results for the ImageNet 2012 test set.

Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification

Kaiming He

Xiangyu Zhang

Shaoqing Ren

Jian Sun

Microsoft Research

{kahe, v-xiangz, v-shren, jiansun}@microsoft.com



GT: letter opener
1: drumstick
2: candle
3: wooden spoon
4: spatula
5: ladle



GT: letter opener
1: Band Aid
2: ruler
3: rubber eraser
4: pencil box
5: wallet



GT: letter opener
1: fountain pen
2: ballpoint
3: hammer
4: can opener
5: ruler



GT: spotlight
1: grand piano
2: folding chair
3: rocking chair
4: dining table
5: upright piano



GT: spotlight
1: acoustic guitar
2: stage
3: microphone
4: electric guitar
5: banjo



GT: spotlight
1: altar
2: candle
3: perfume
4: restaurant
5: confectionery



GT: restaurant
1: wine bottle
2: candle
3: red wine
4: French loaf
5: wooden spoon



GT: restaurant
1: goblet
2: plate
3: candle
4: red wine
5: dining table



GT: restaurant
1: plate
2: meat loaf
3: ice cream
4: chocolate sauce
5: potpie

(Error %)

	team	top-5 (test)
in competition ILSVRC 14	MSRA, SPP-nets [11]	8.06
	VGG [25]	7.32
	GoogLeNet [29]	6.66
post-competition	VGG [25] (arXiv v5)	6.8
	Baidu [32]	5.98
	MSRA, PReLU-nets	4.94

Table 7. The multi-model results for the ImageNet 2012 test set.