

Lecture 15

Convolutional Neural Networks for Object Recognition



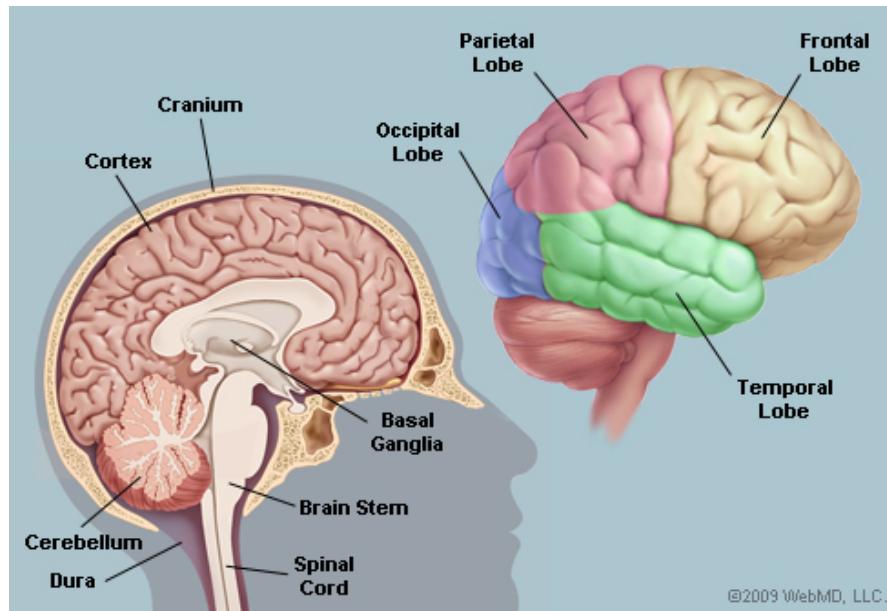
- Neural Networks
- Convolutional Neural Networks (CNNs) and image classification
- Object Detection with CNNs
- Segmentation with CNNs

Slides of this lecture are courtesy Andrej Karpathy & Justin Johnson

Part 1. Neural Networks

Understand the Human Brain

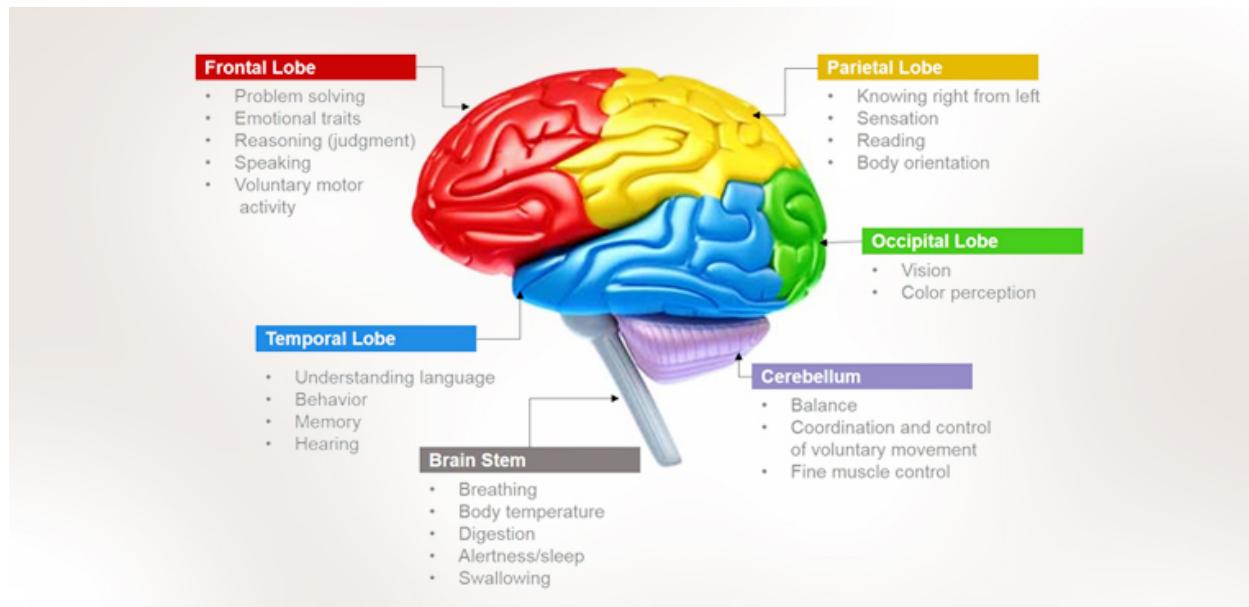
- Neuroscientist
 - Try to understand the structure and mechanism of human brain biologically



©2009 WebMD, LLC.

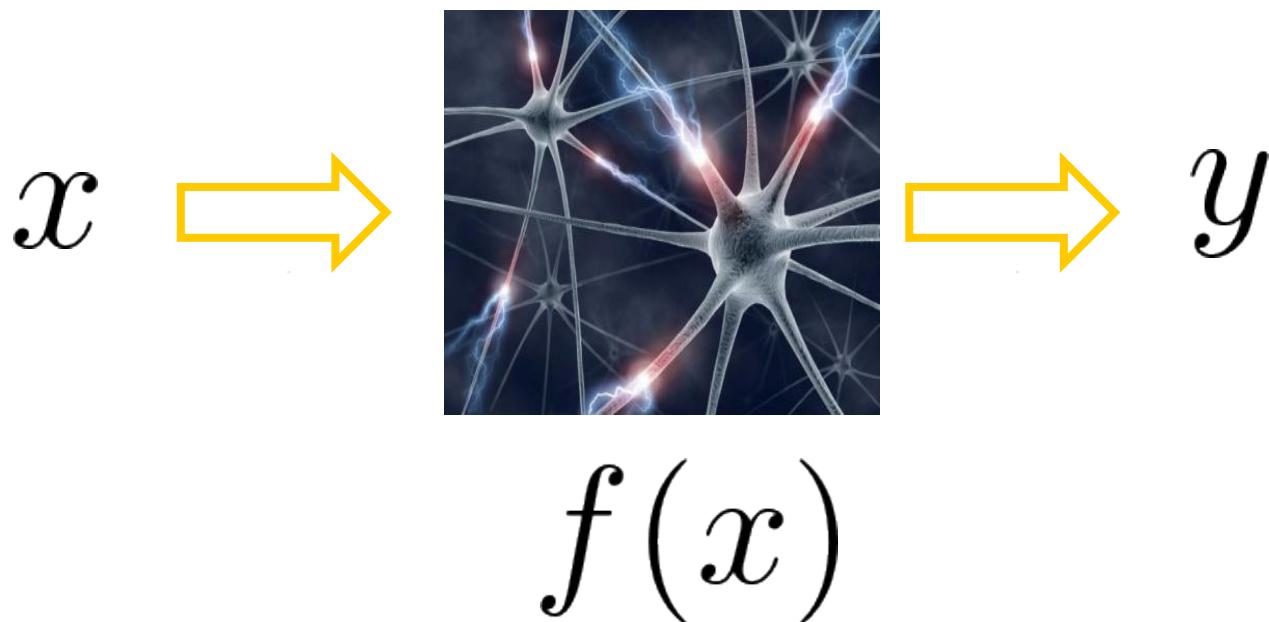
Understand the Human Brain

- Cognitive scientist
 - Try to understand how the human brain works by experiments



Understand the Human Brain

- Computer Scientist
 - Try to model the human brain with computational models



Understand the Human Brain

- What is the form of the function $f(x)$?
 - No idea!
 - Concatenate simple functions (neurons)



x



$$f(x)$$



$$y \in \{+1, -1\}$$

Dog

Neural Network: concatenation of functions

Linear score function: $f = Wx$

Neural Network: without the brain stuff

Linear score function: $f = Wx$

2-layer Neural Network

$$f = f_2(f_1(x)) = W_2 \max(0, W_1 x)$$

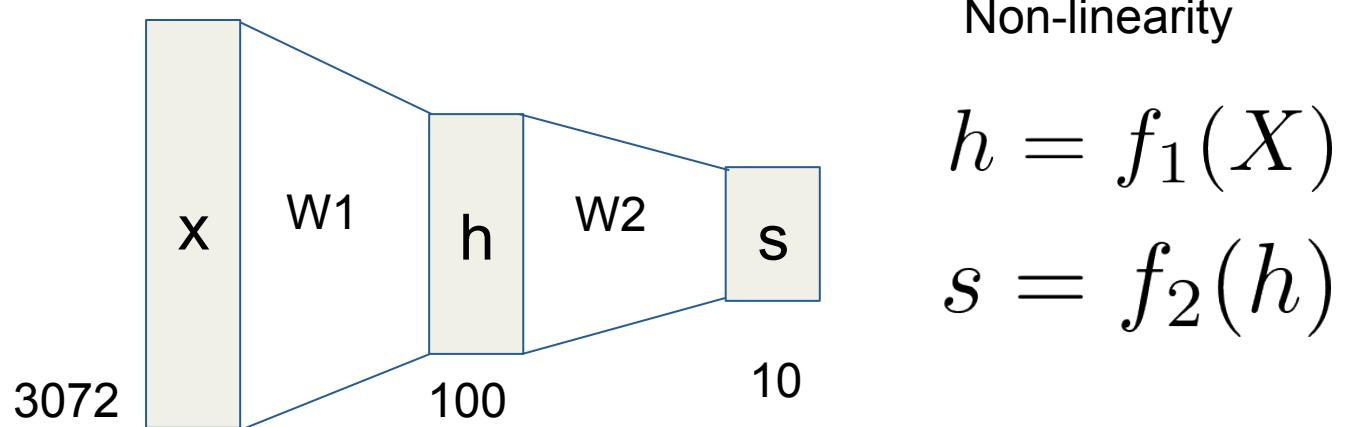
Non-linearity

Neural Network: without the brain stuff

Linear score function: $f = Wx$

2-layer Neural Network

$$f = f_2(f_1(x)) = W_2 \max(0, W_1 x)$$



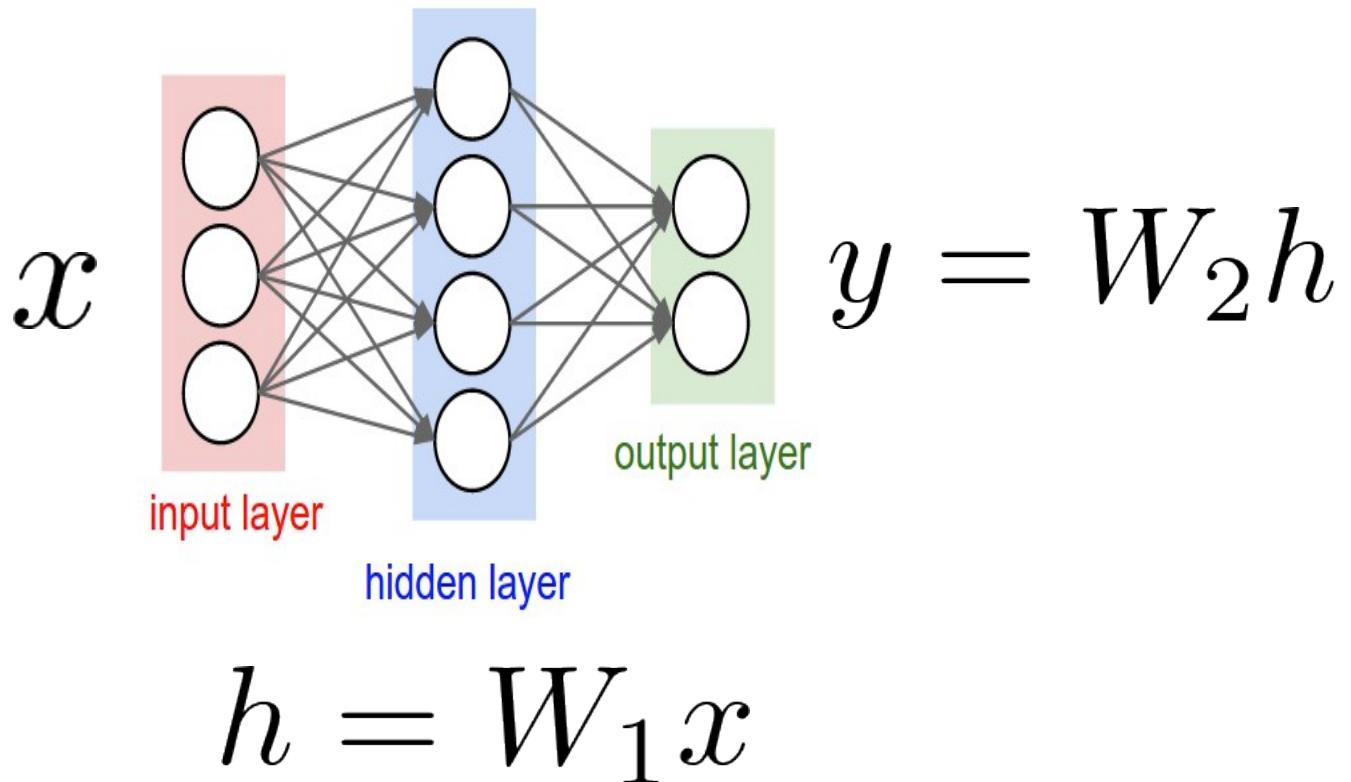
Neural Network: without the brain stuff

Linear score function: $f = Wx$

2-layer Neural Network $f = W_2 \max(0, W_1 x)$
or 3-layer Neural Network

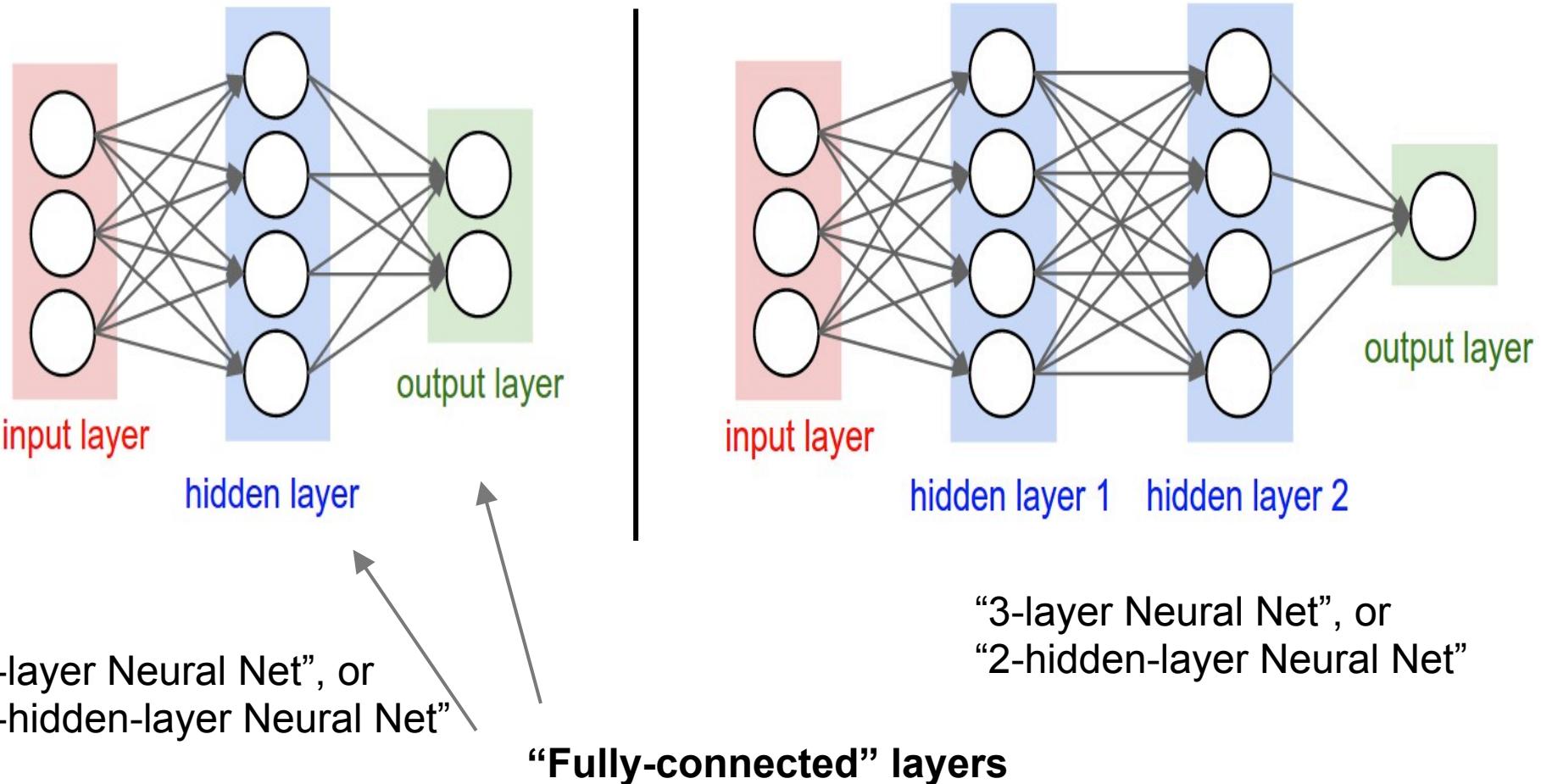
$f = W_3 \max(0, W_2 \max(0, W_1 x))$

Neural Networks: Architectures



**“Fully-connected” layers
(matrix multiplication)**

Neural Networks: Architectures

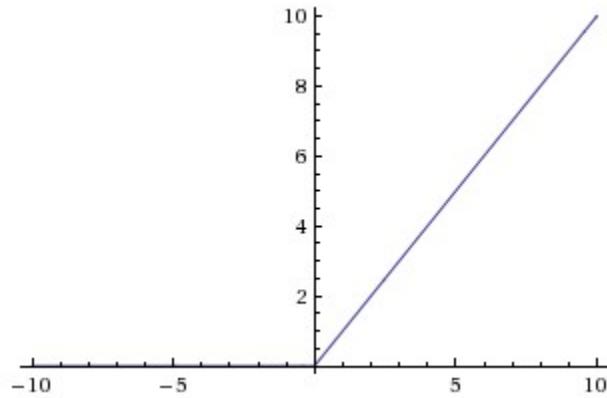


Activation Functions

2-layer Neural Network

$$f = f_2(f_1(x)) = W_2 \max(0, W_1 x)$$

ReLU $\max(0, x)$

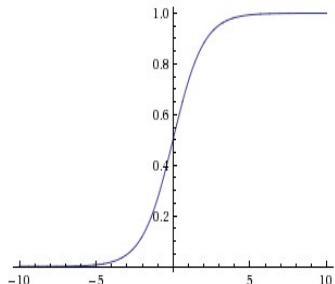


Introduce non-linearity to the network

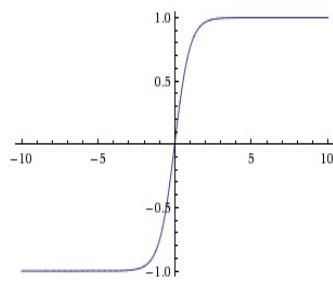
Activation Functions

Sigmoid

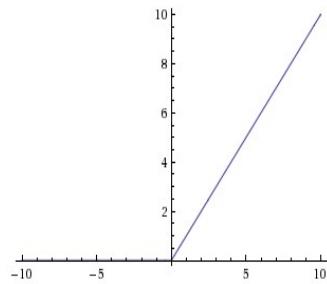
$$\sigma(x) = 1/(1 + e^{-x})$$



$$\tanh \quad \tanh(x)$$

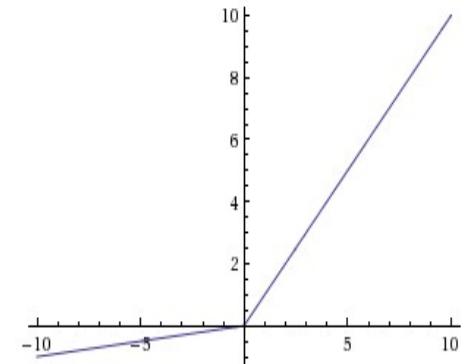


$$\text{ReLU} \quad \max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

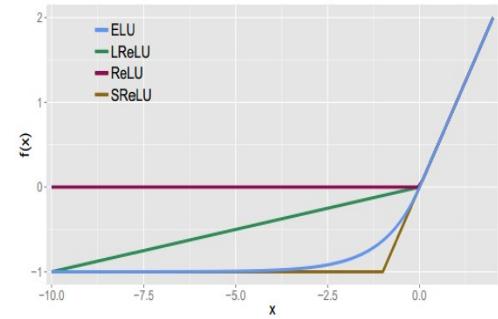


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

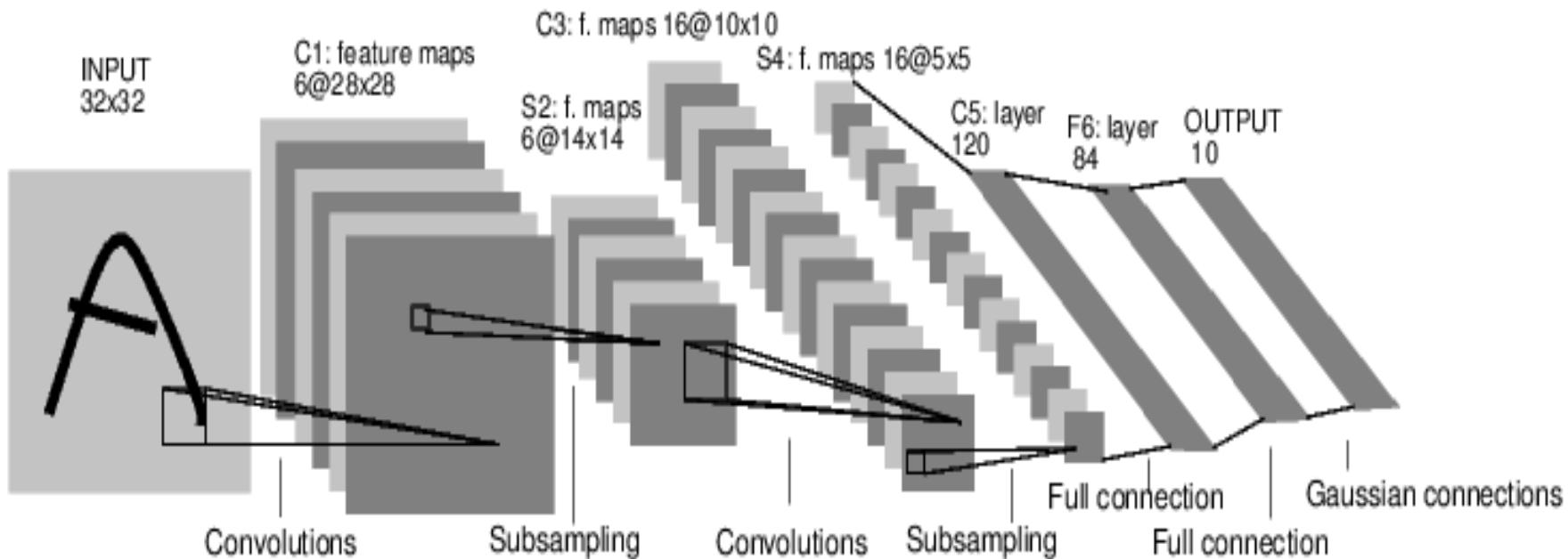
ELU

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$



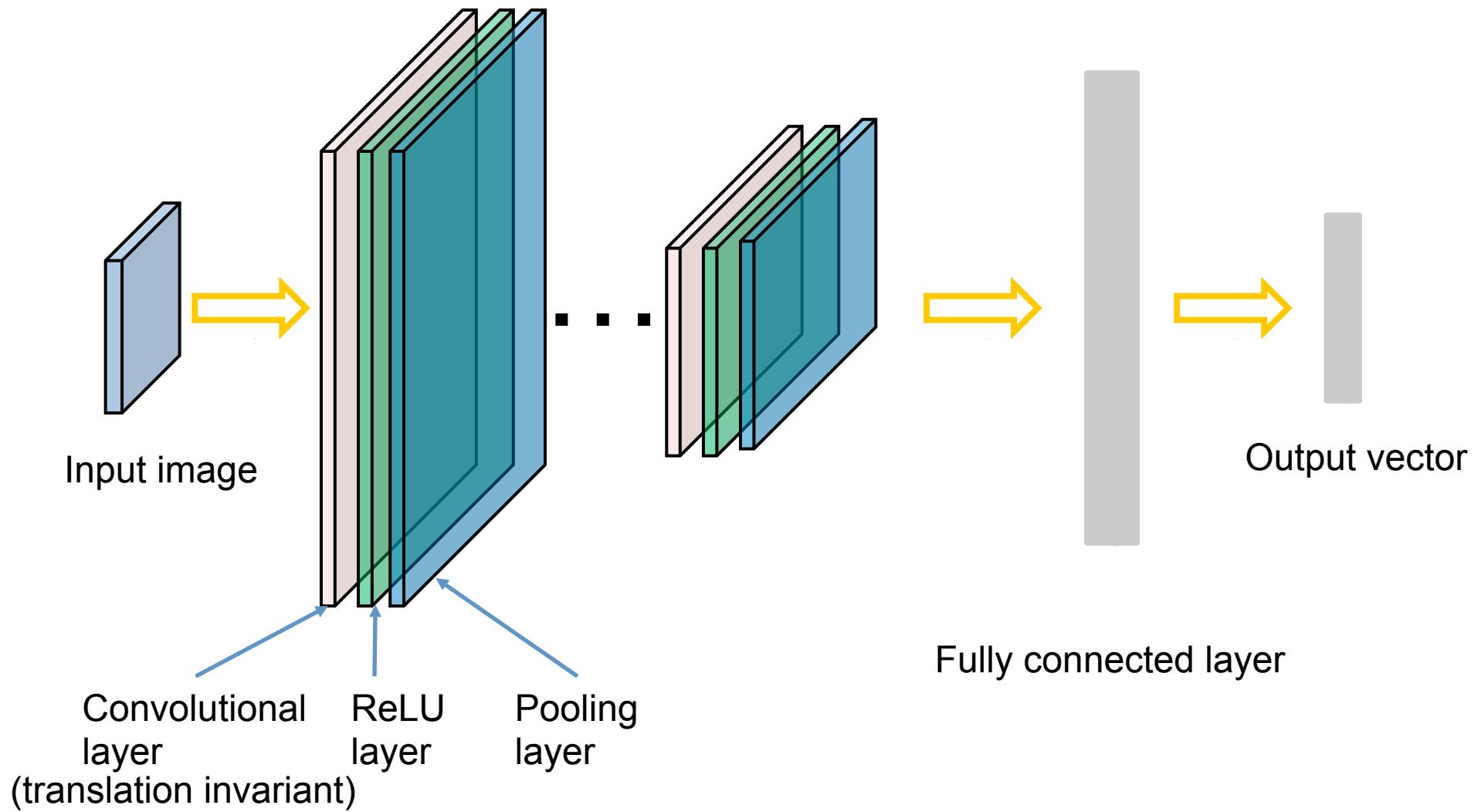
Part 2. Convolutional Neural Networks

Convolutional Neural Networks



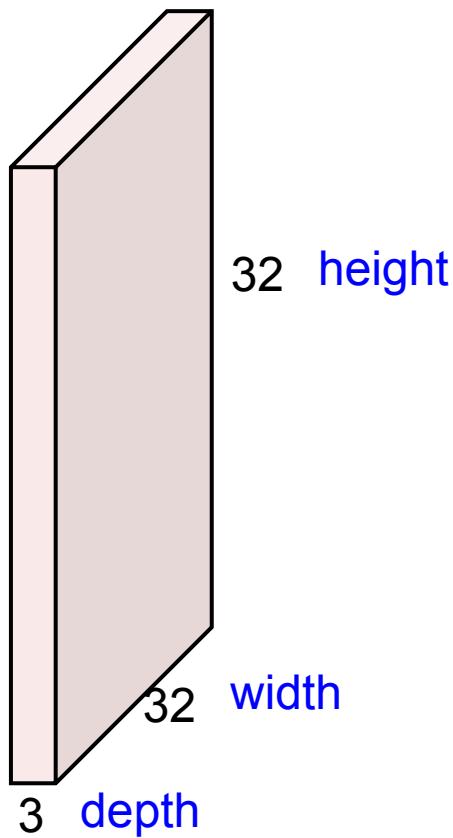
[LeNet-5, LeCun 1980]

Convolutional Neural Networks



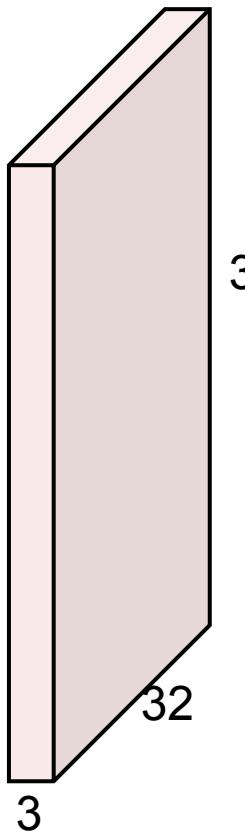
Convolutional Layer

32x32x3 image



Convolutional Layer

32x32x3 image

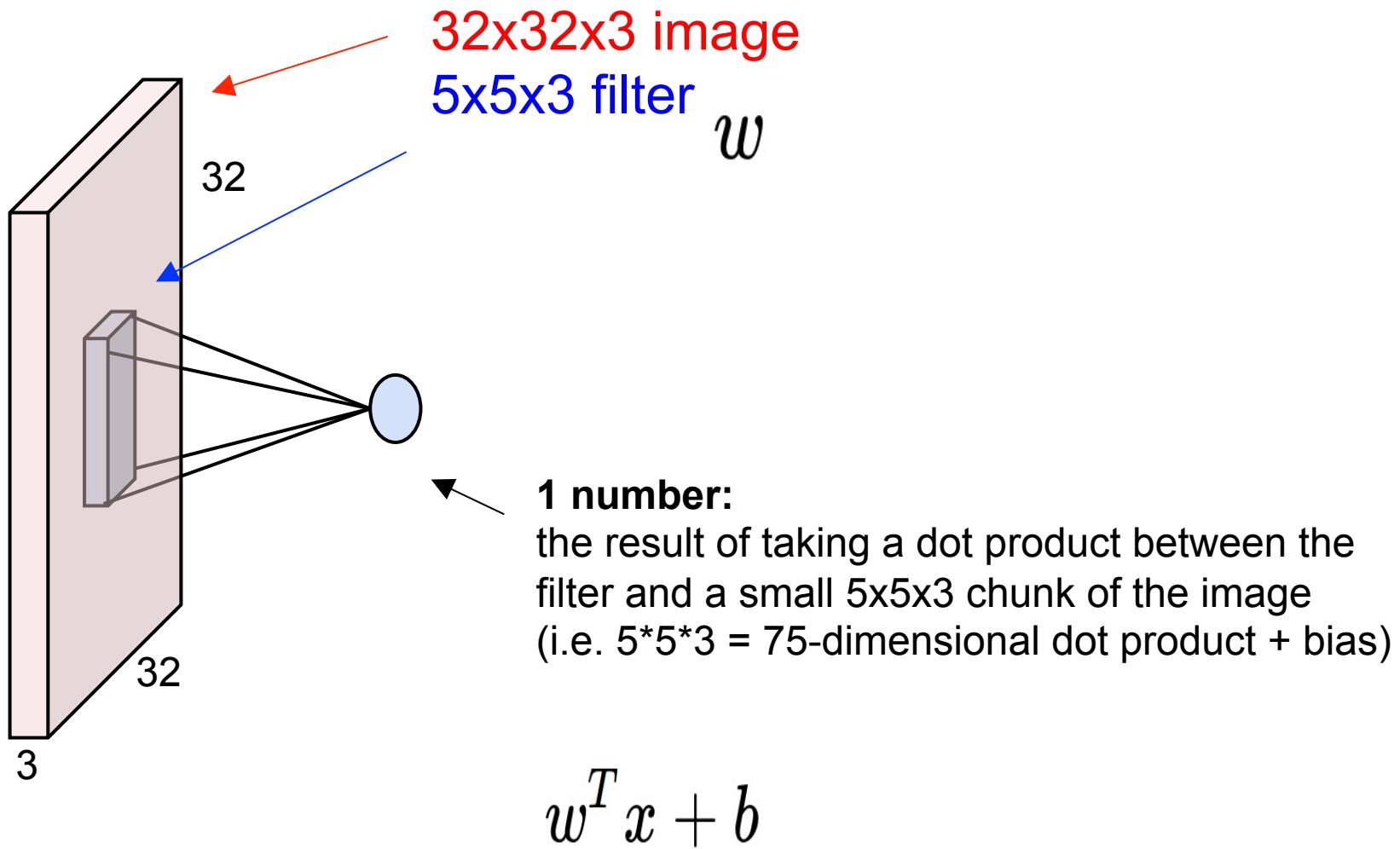


5x5x3 filter

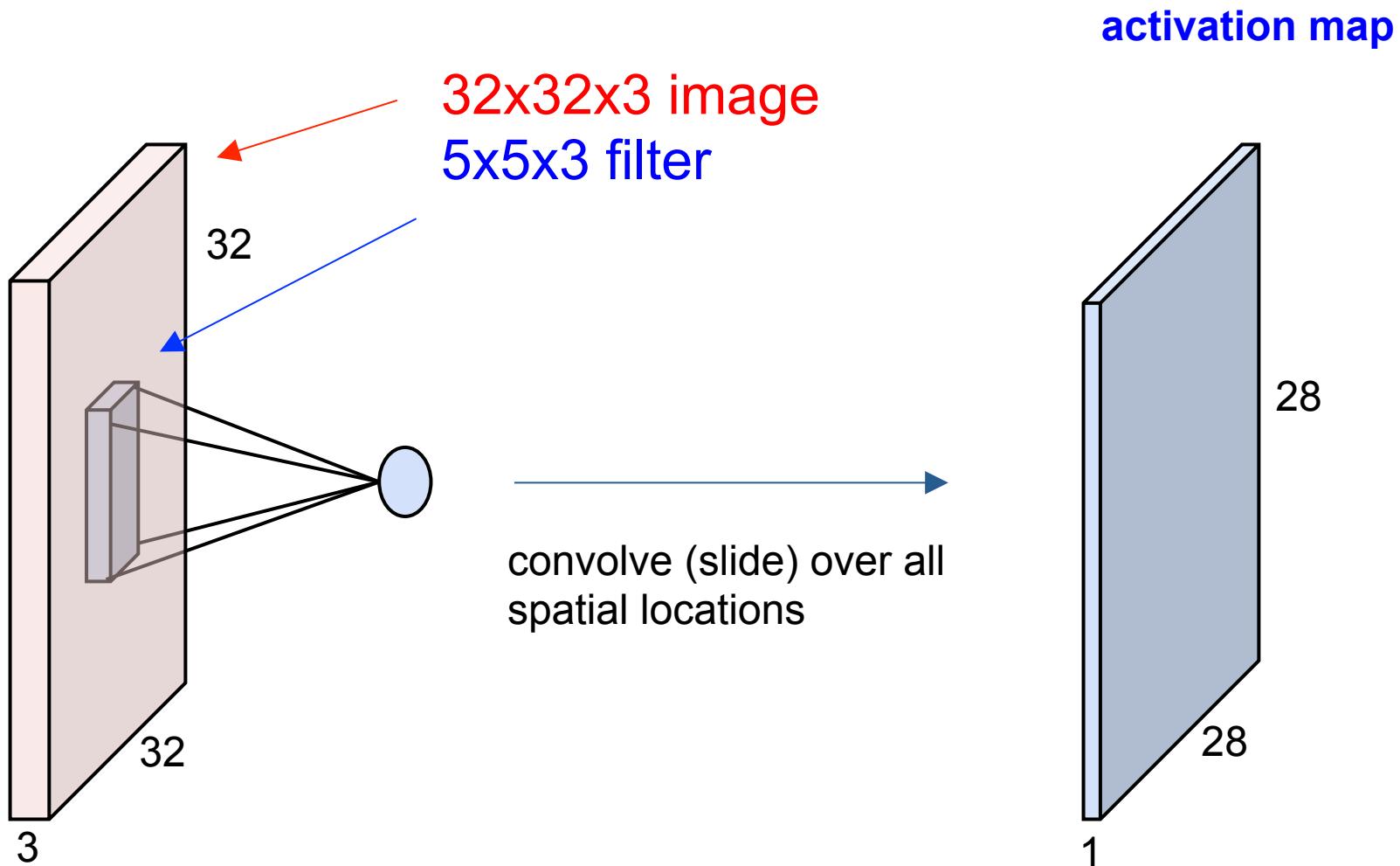


Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Convolutional Layer

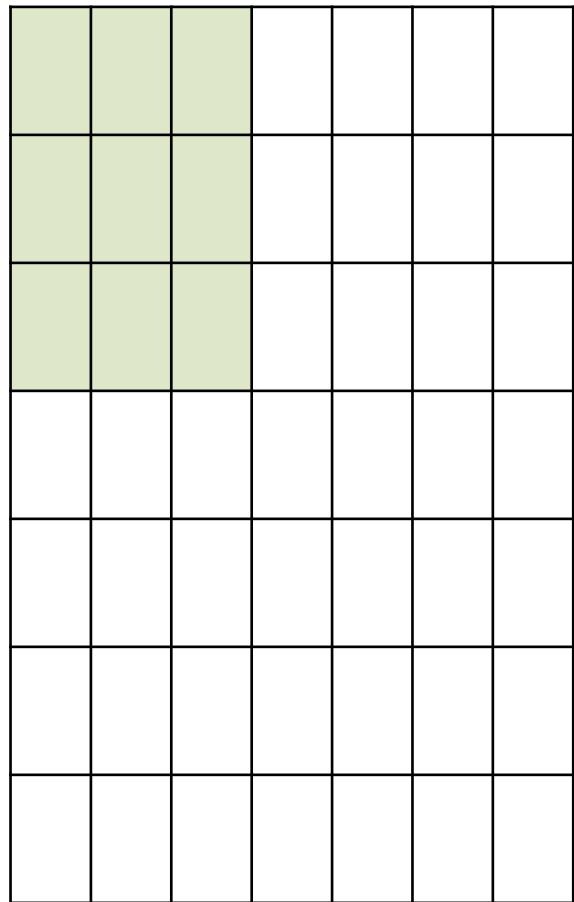


Convolutional Layer



A closer look at spatial dimensions:

7

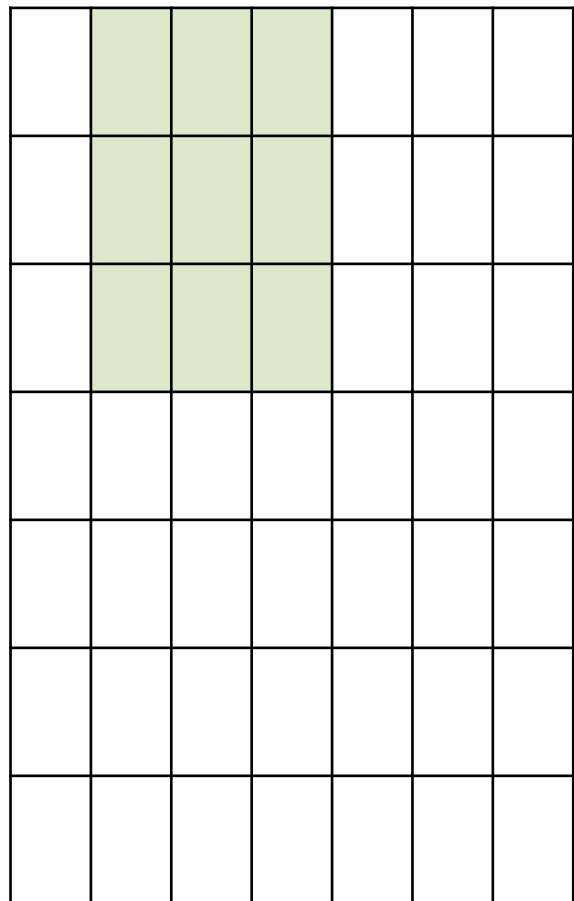


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

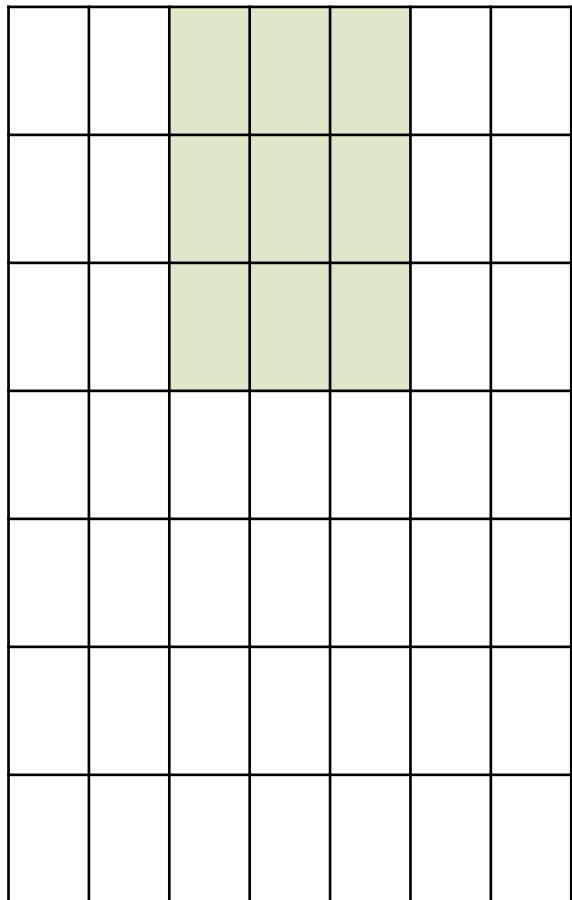


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

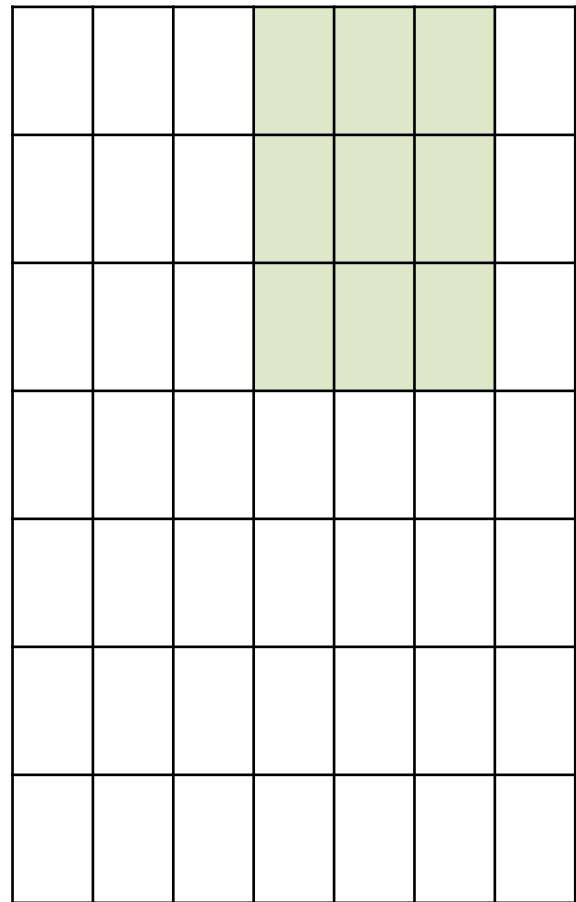


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

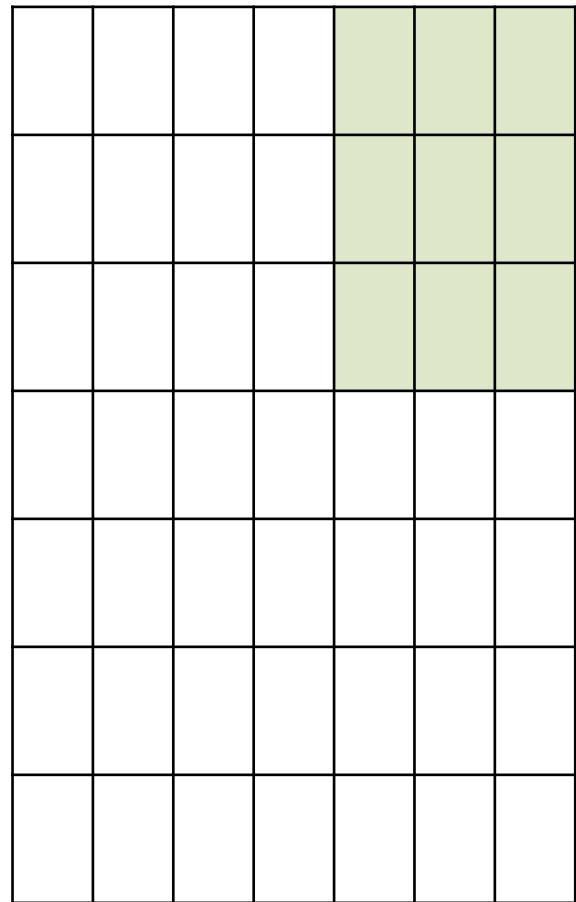


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

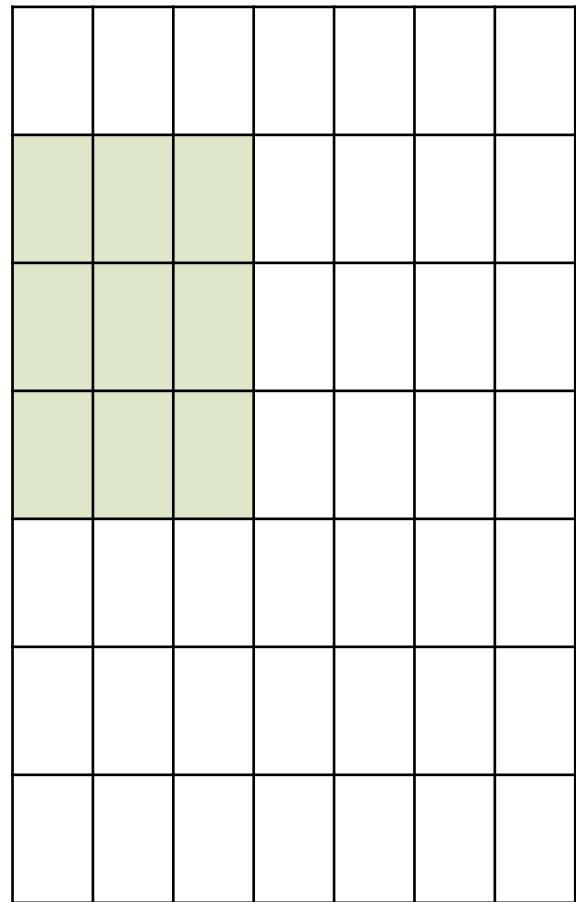


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

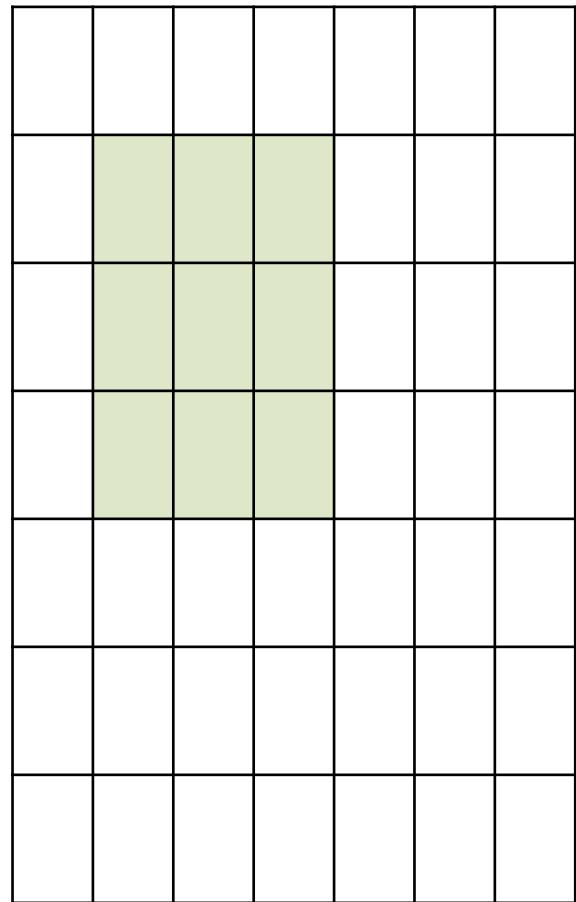


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

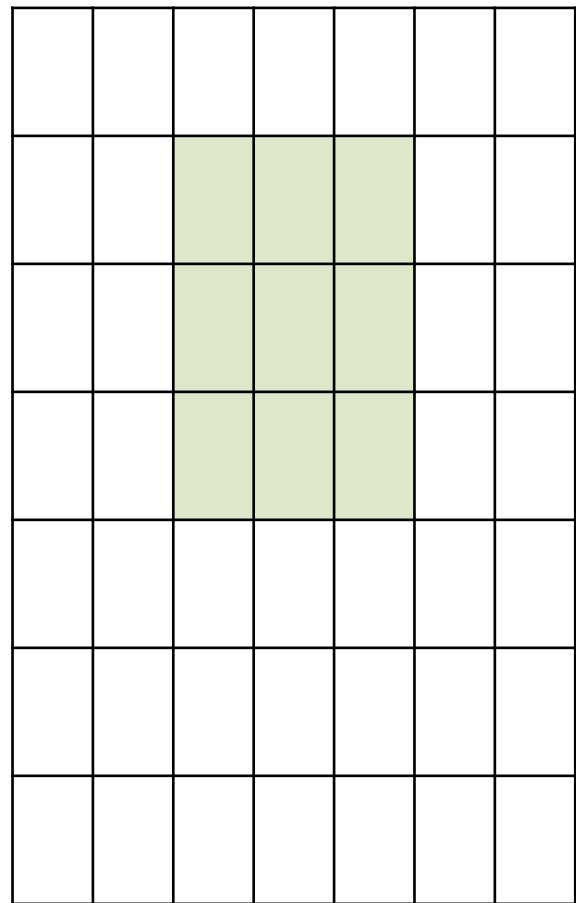


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

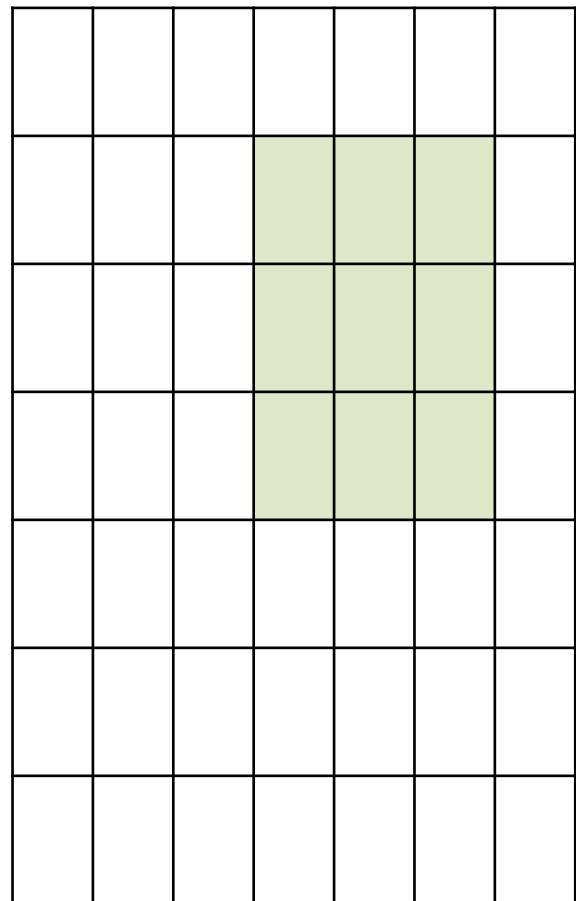


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

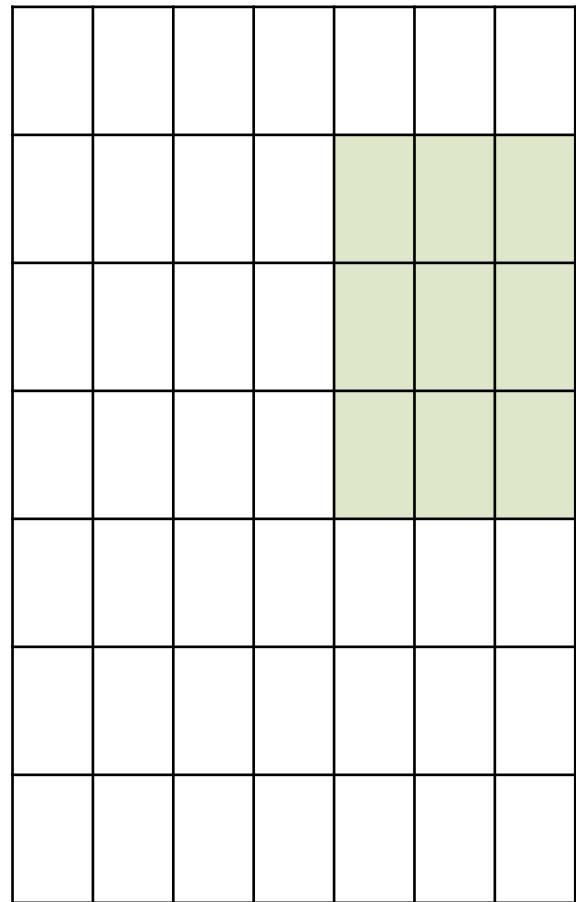


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7



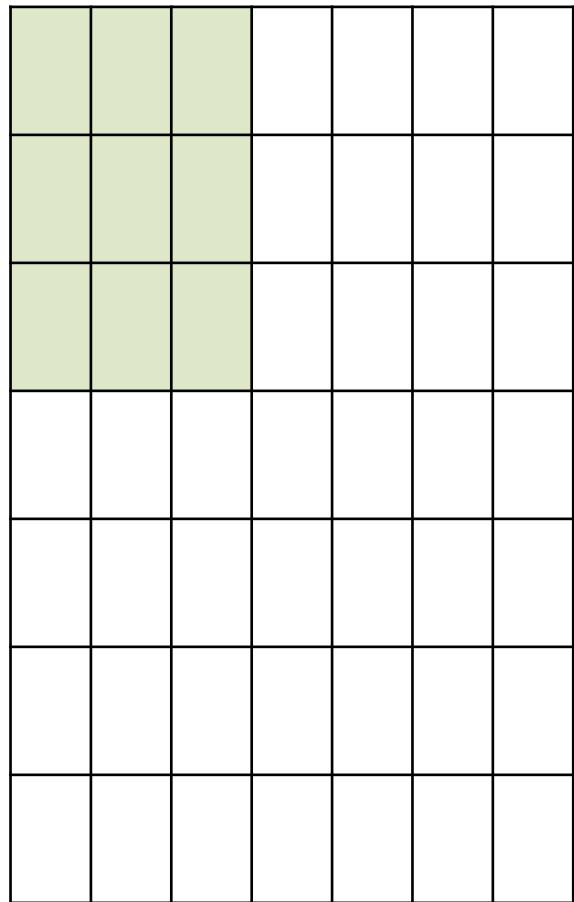
7x7 input (spatially)
assume 3x3 filter

=> 5x5 output

7

A closer look at spatial dimensions:

7

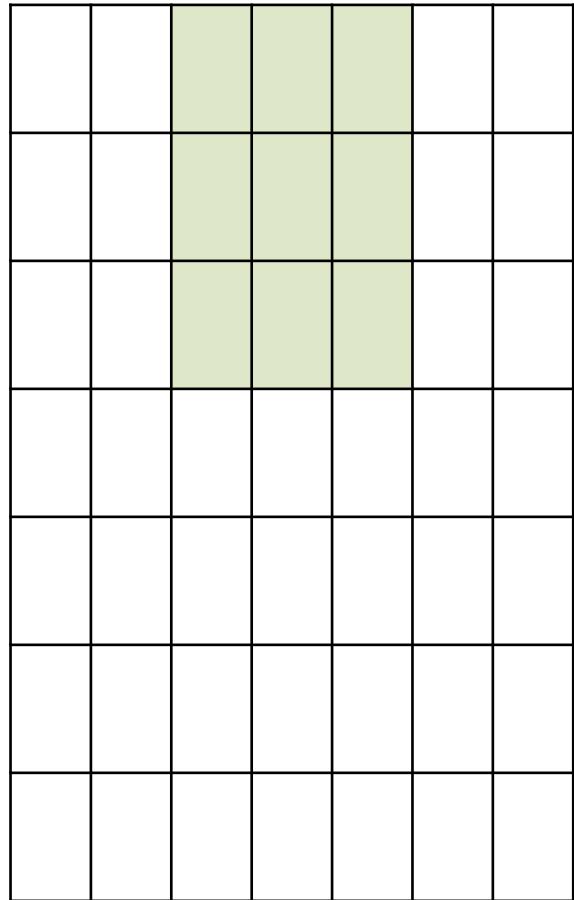


7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

7

A closer look at spatial dimensions:

7

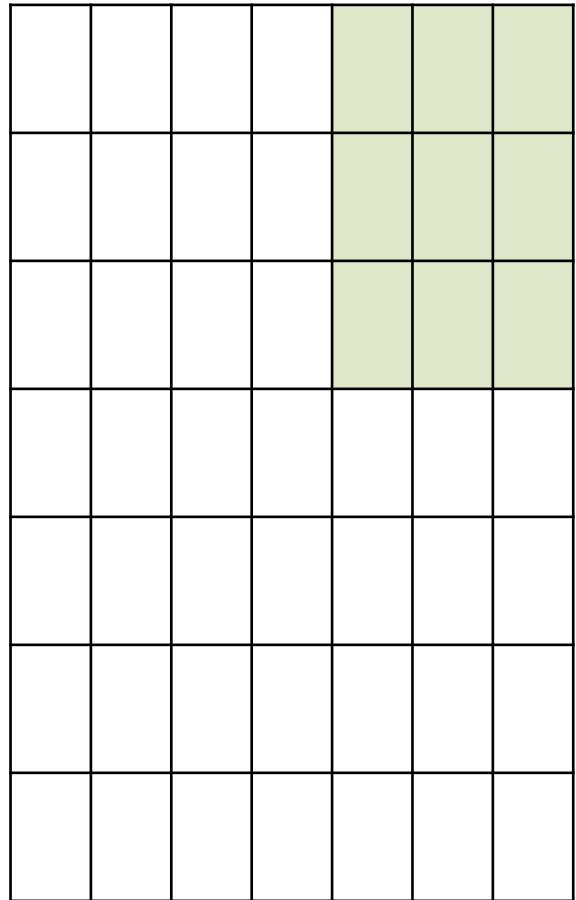


7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

7

A closer look at spatial dimensions:

7

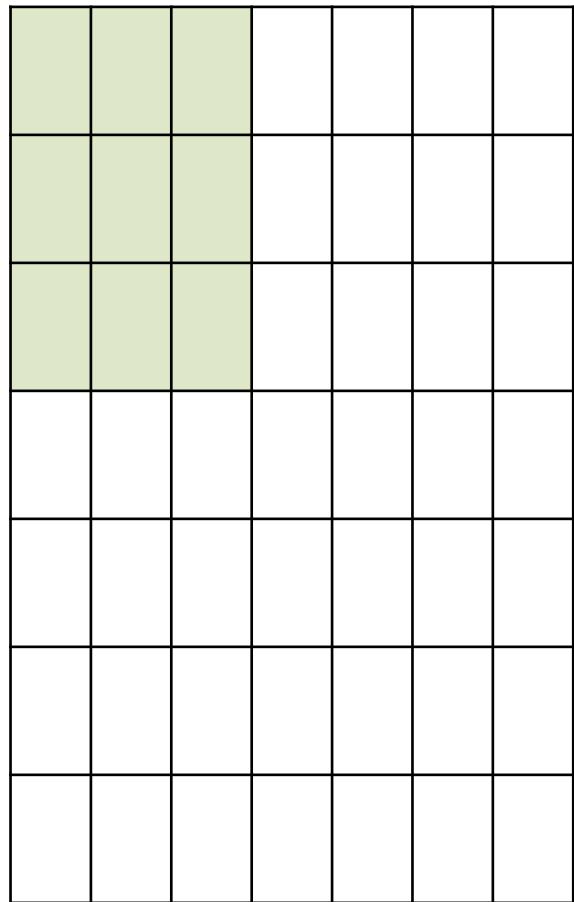


7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
=> 3x3 output!

A closer look at spatial dimensions:

7

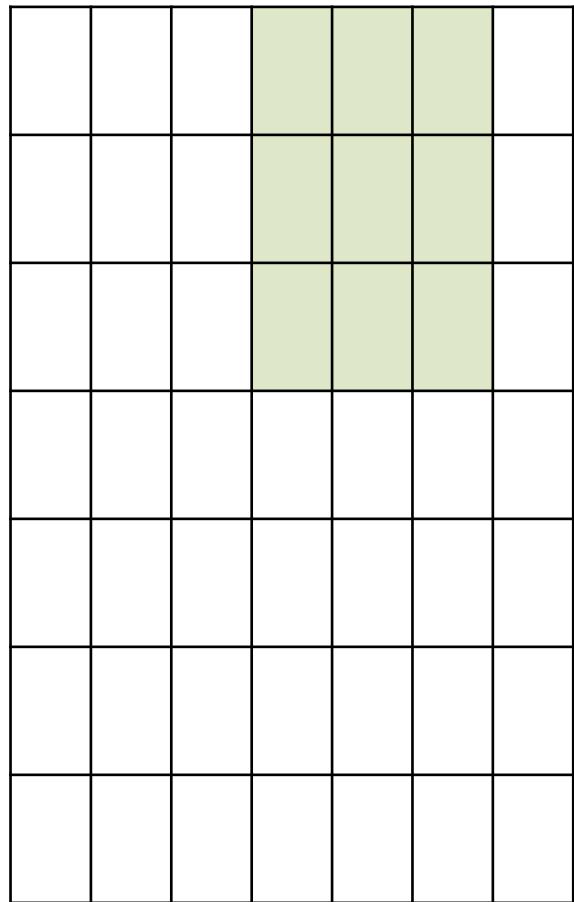


7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

A closer look at spatial dimensions:

7

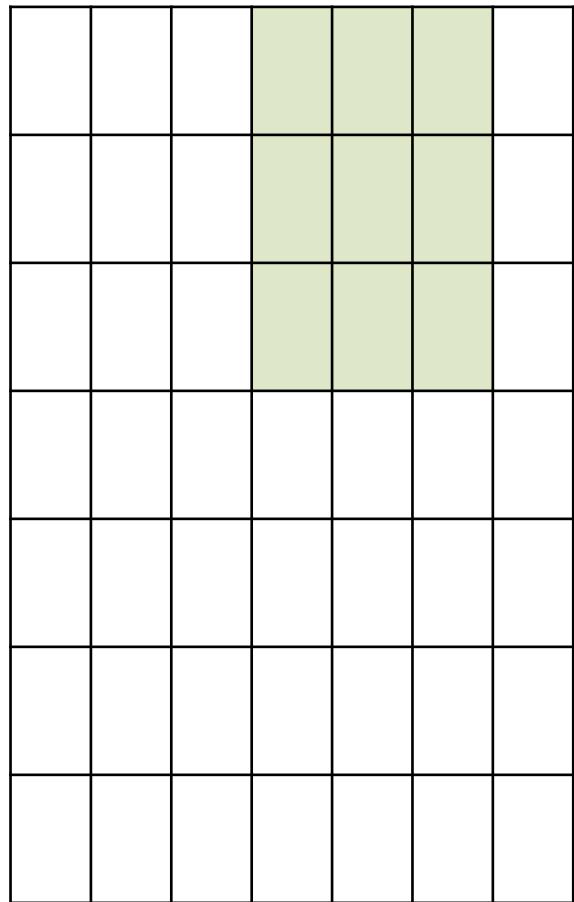


7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

7

A closer look at spatial dimensions:

7

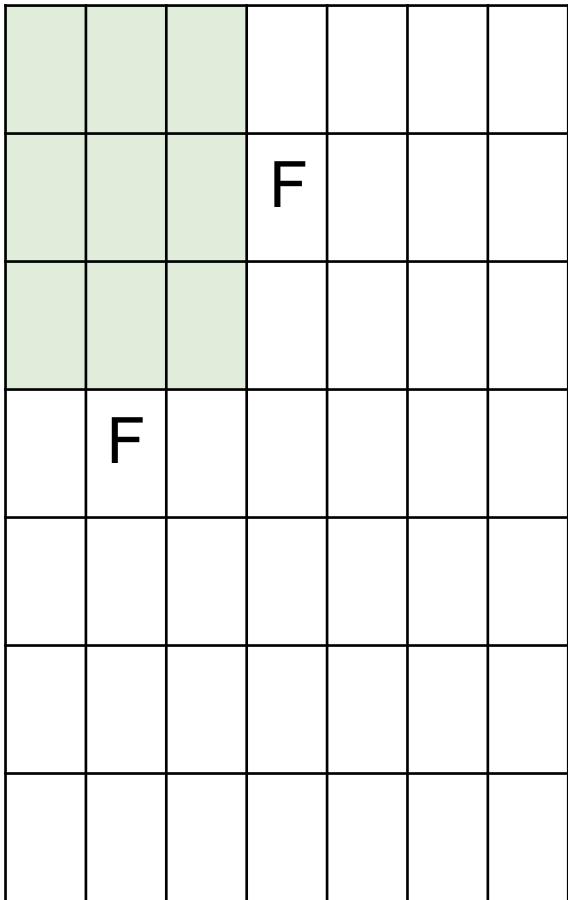


7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

7

doesn't fit!
cannot apply 3x3 filter on
7x7 input with stride 3.

N



N

Output size:

$$(N - F) / \text{stride} + 1$$

e.g. $N = 7, F = 3$:

$$\text{stride } 1 \Rightarrow (7 - 3)/1 + 1 = 5$$

$$\text{stride } 2 \Rightarrow (7 - 3)/2 + 1 = 3$$

$$\text{stride } 3 \Rightarrow (7 - 3)/3 + 1 = 2.33$$

In practice: Common to zero pad the border

e.g. input 7x7

3x3 filter, applied with stride 1

pad with 1 pixel border => what is the output?

(recall:)
 $(N - F) / \text{stride} + 1$

In practice: Common to zero pad the border

| | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|--|--|--|--|--|--|--|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | |
| 0 | | | | | | | | | | | | | |
| 0 | | | | | | | | | | | | | |
| 0 | | | | | | | | | | | | | |
| 0 | | | | | | | | | | | | | |
| | | | | | | | | | | | | | |
| | | | | | | | | | | | | | |
| | | | | | | | | | | | | | |

e.g. input 7x7

3x3 filter, applied with stride 1

pad with 1 pixel border => what is the output?

7x7 output!

In practice: Common to zero pad the border

| | | | | | | | | |
|---|---|---|---|---|---|--|--|--|
| 0 | 0 | 0 | 0 | 0 | 0 | | | |
| 0 | | | | | | | | |
| 0 | | | | | | | | |
| 0 | | | | | | | | |
| 0 | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |

e.g. input 7x7

3x3 filter, applied with stride 1

pad with 1 pixel border => what is the output?

7x7 output!

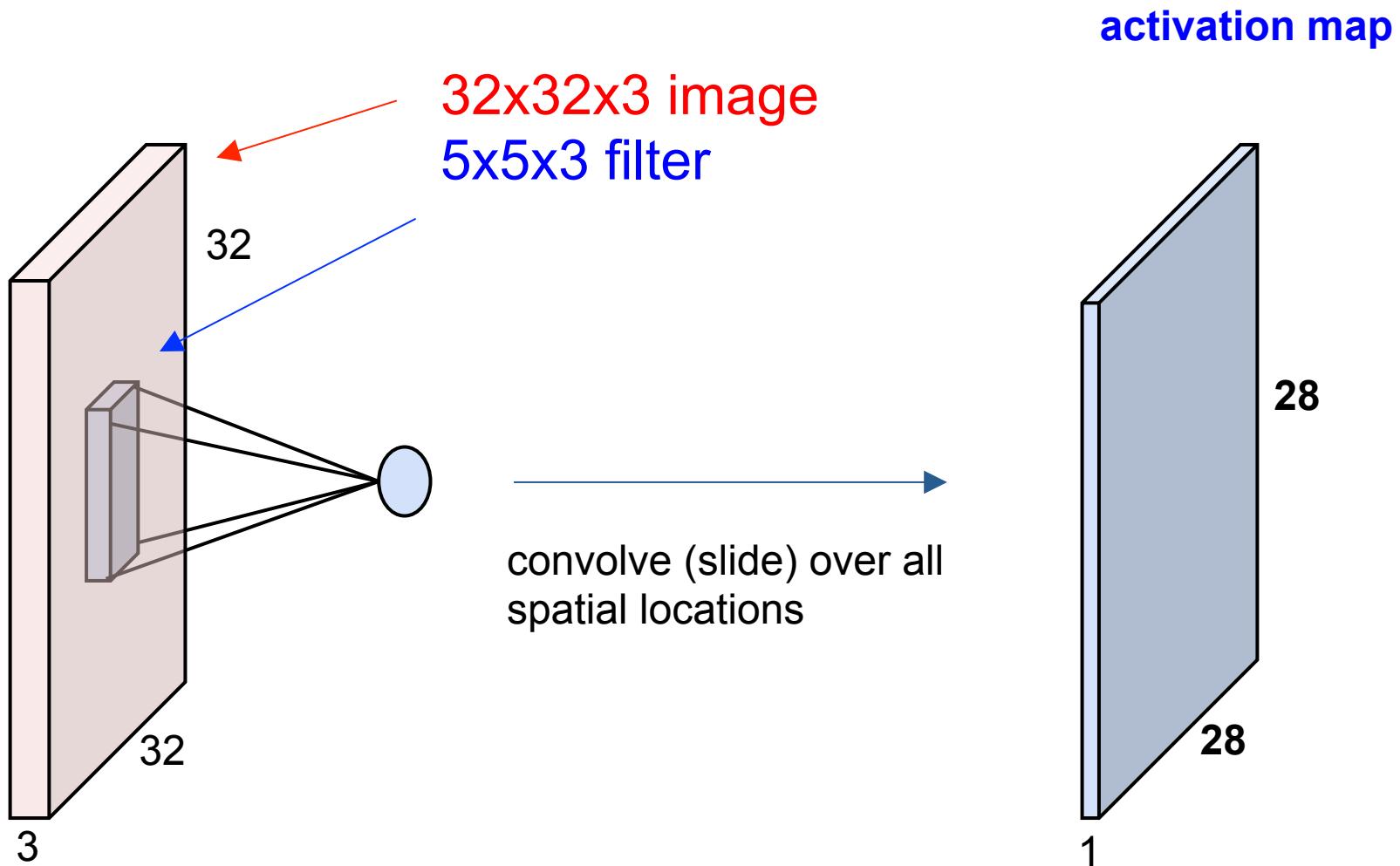
in general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with $(F-1)/2$. (will preserve size spatially)

e.g. $F = 3 \Rightarrow$ zero pad with 1

$F = 5 \Rightarrow$ zero pad with 2

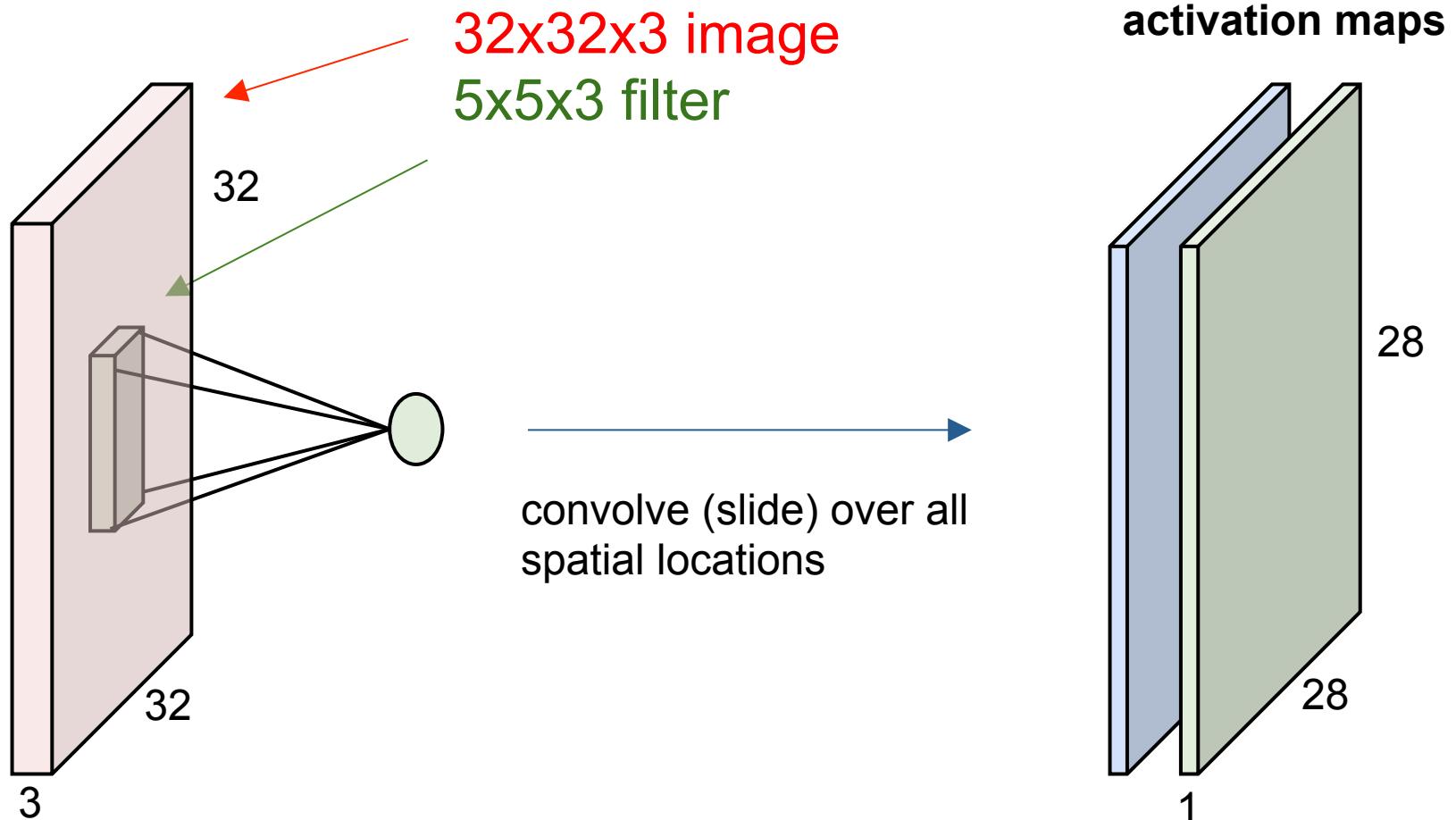
$F = 7 \Rightarrow$ zero pad with 3

A closer look at spatial dimensions:

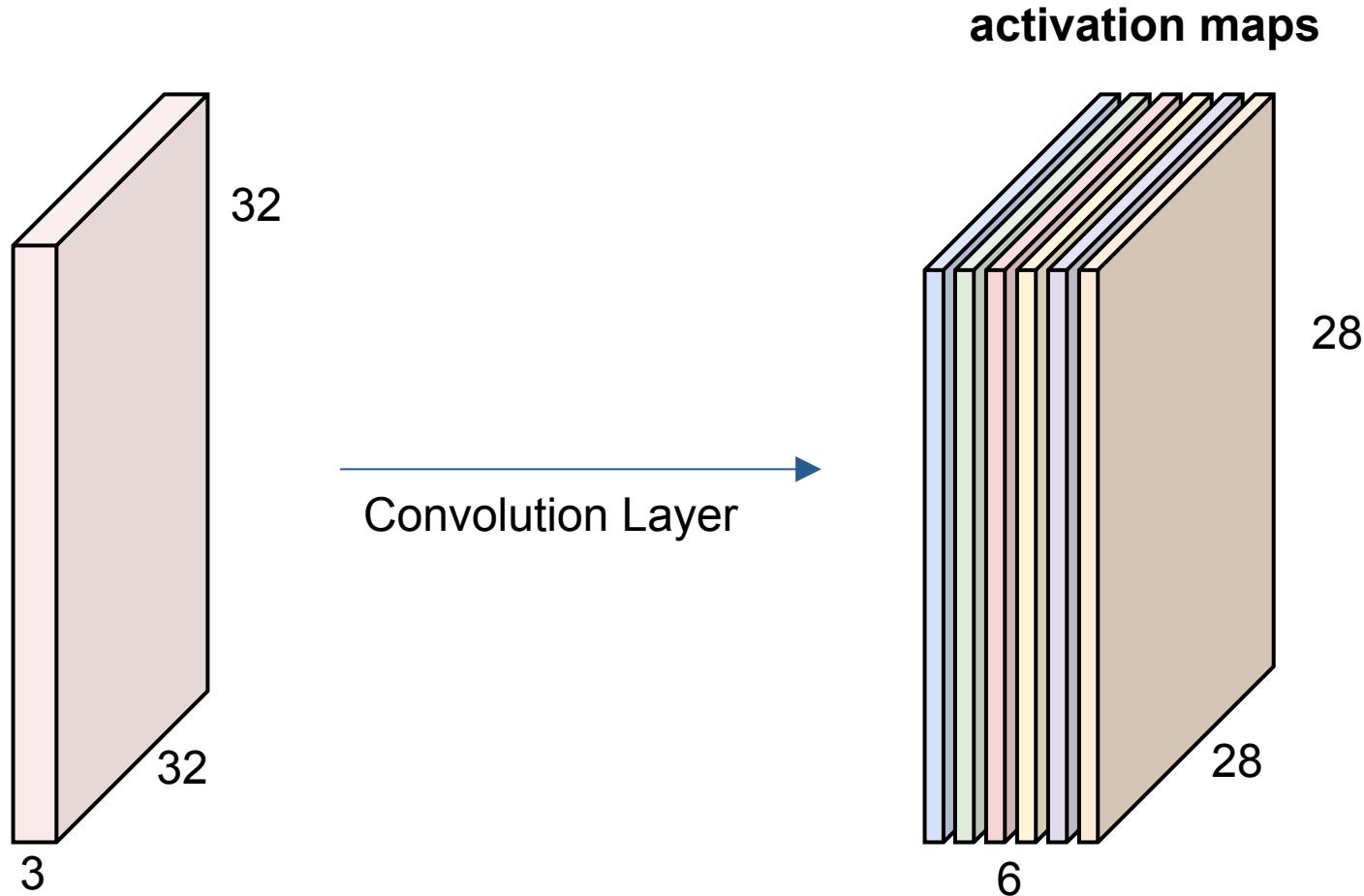


Convolutional Layer

consider a second, green filter

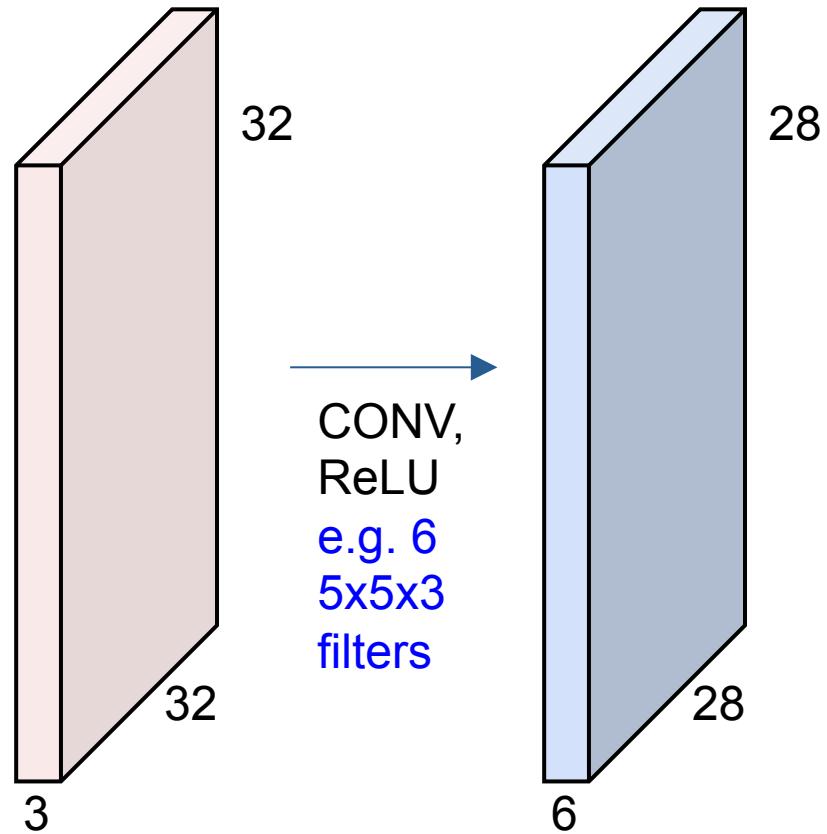


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

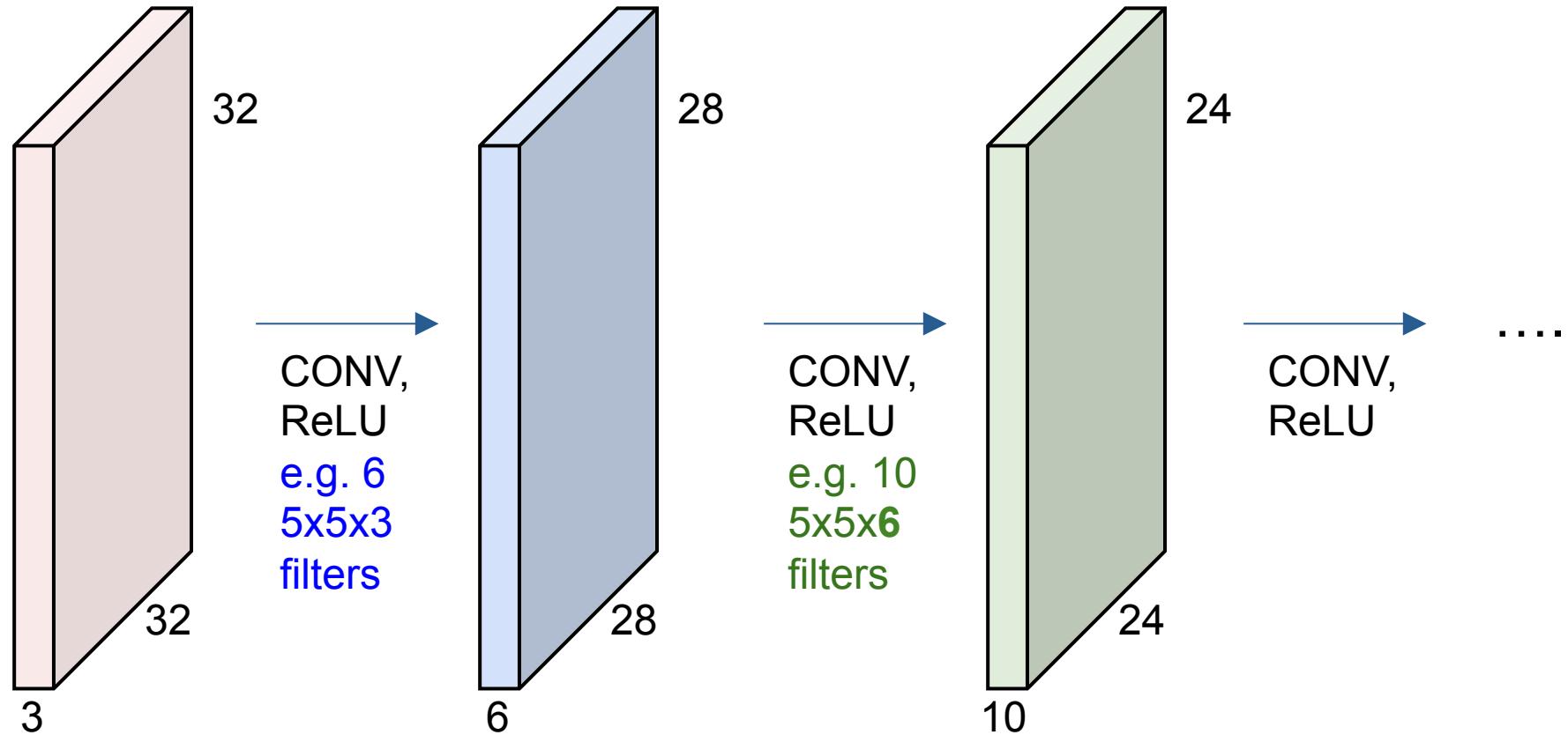


We stack these up to get a “new image” of size $28 \times 28 \times 6$!

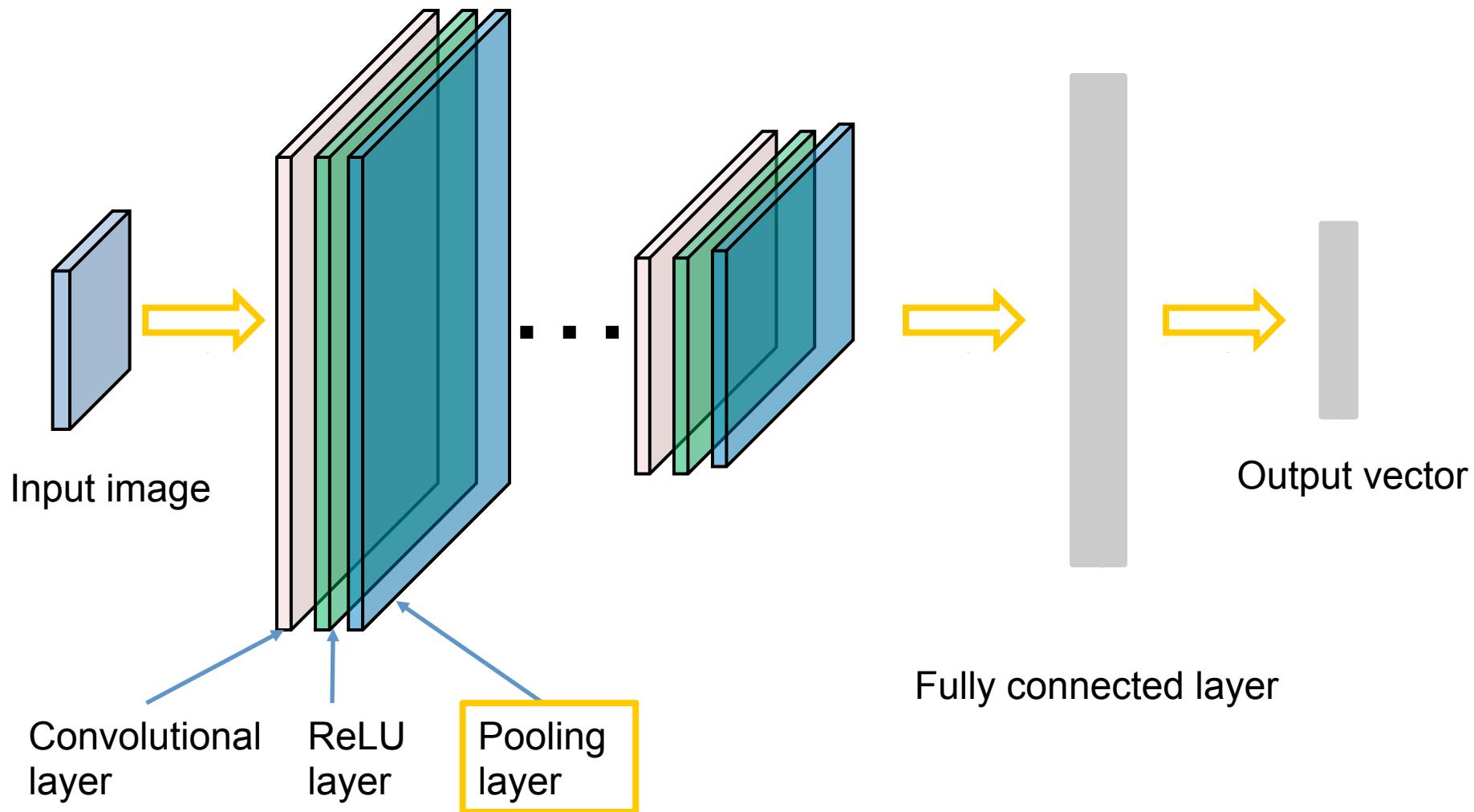
Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions



Preview: ConvNet is a sequence of Convolutional Layers, interspersed with activation functions

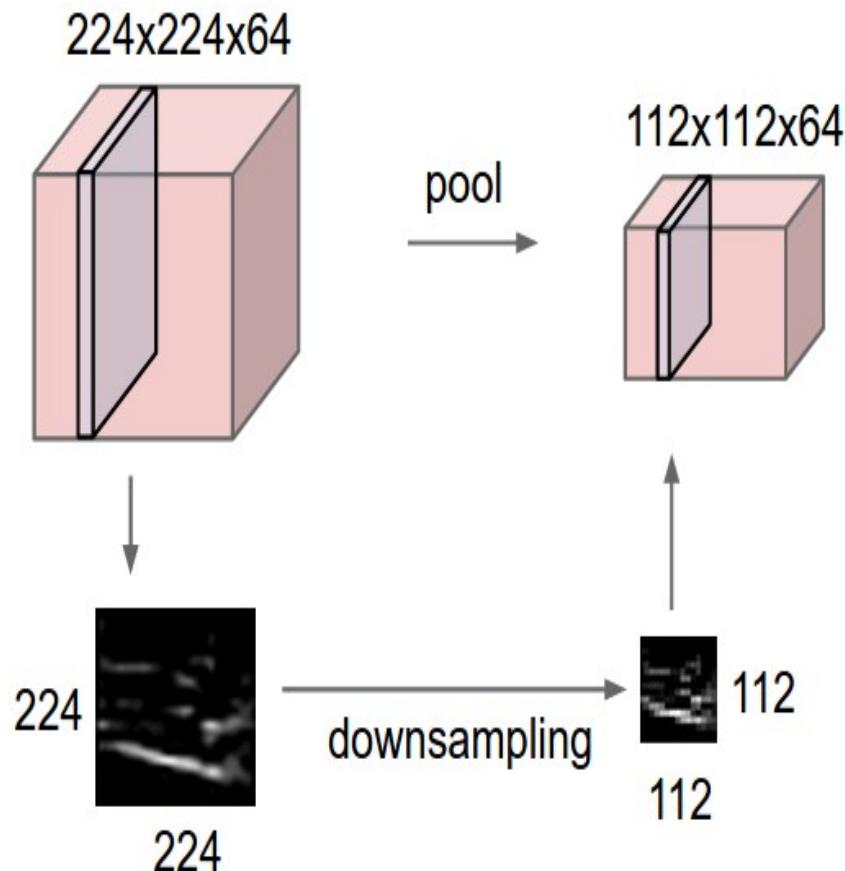


Convolutional Neural Networks



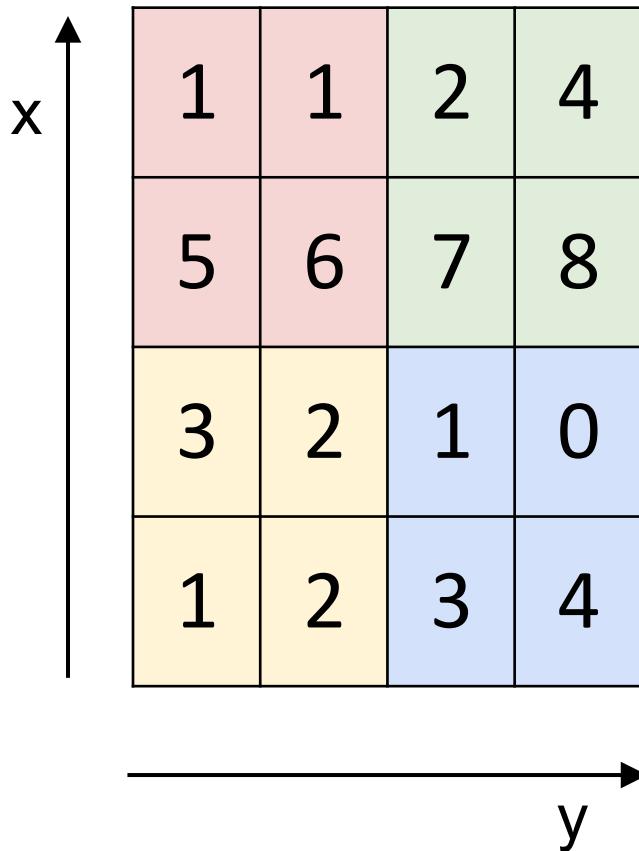
Pooling layer

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

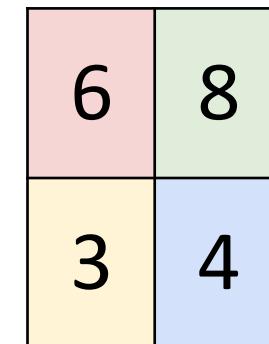


MAX POOLING

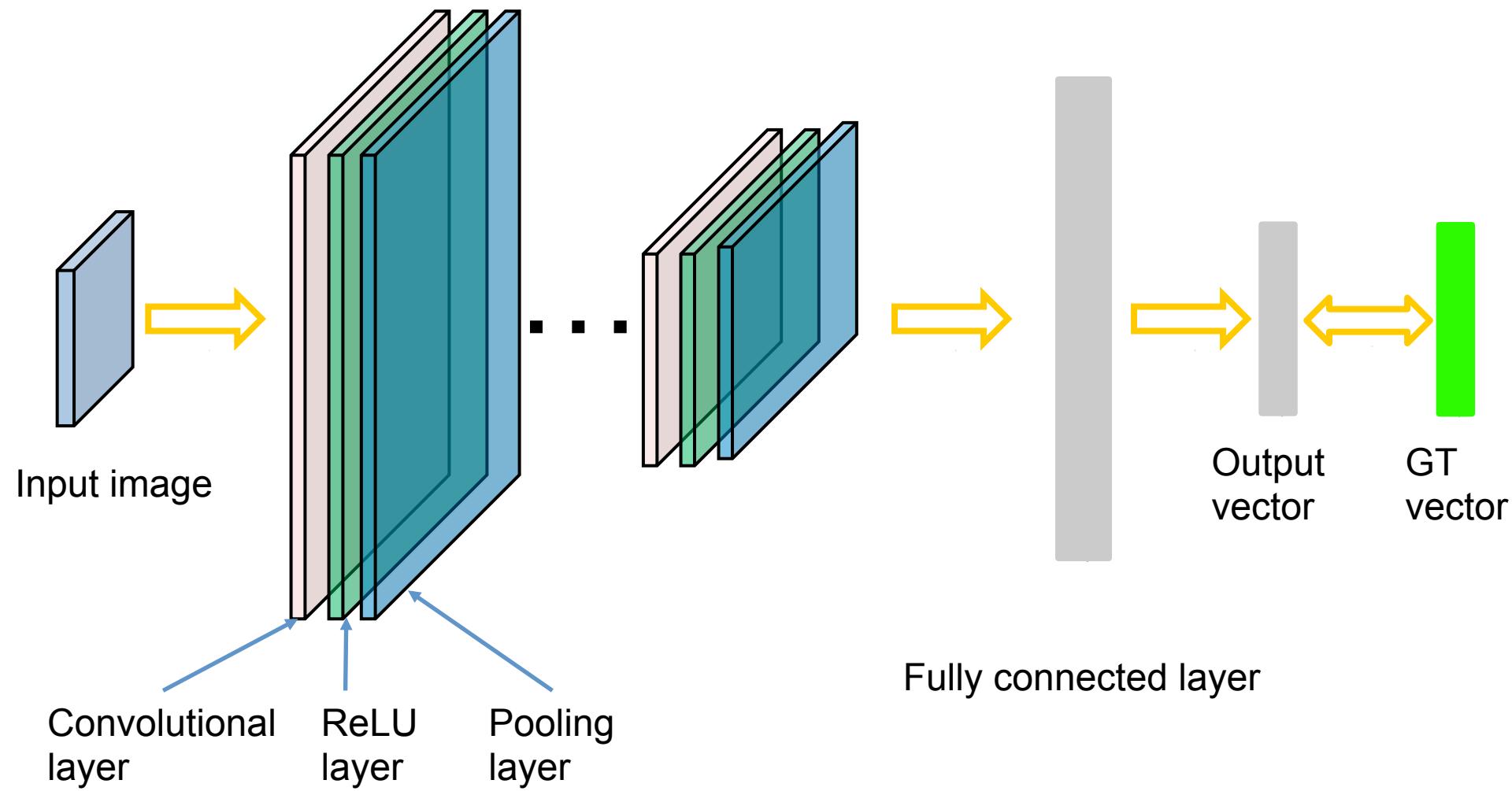
Single depth slice



max pool with 2x2 filters
and stride 2

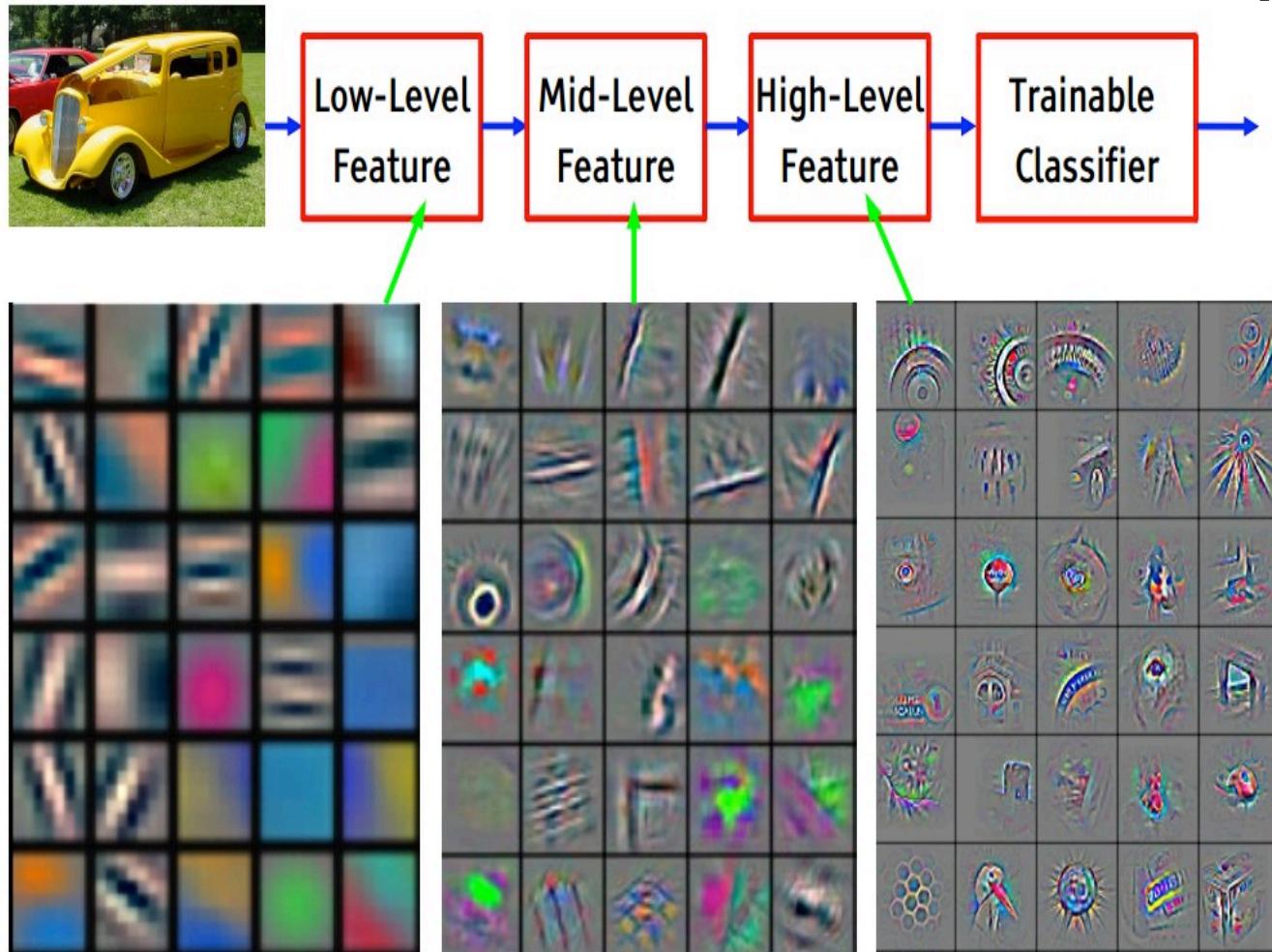


Training: back-propogate errors

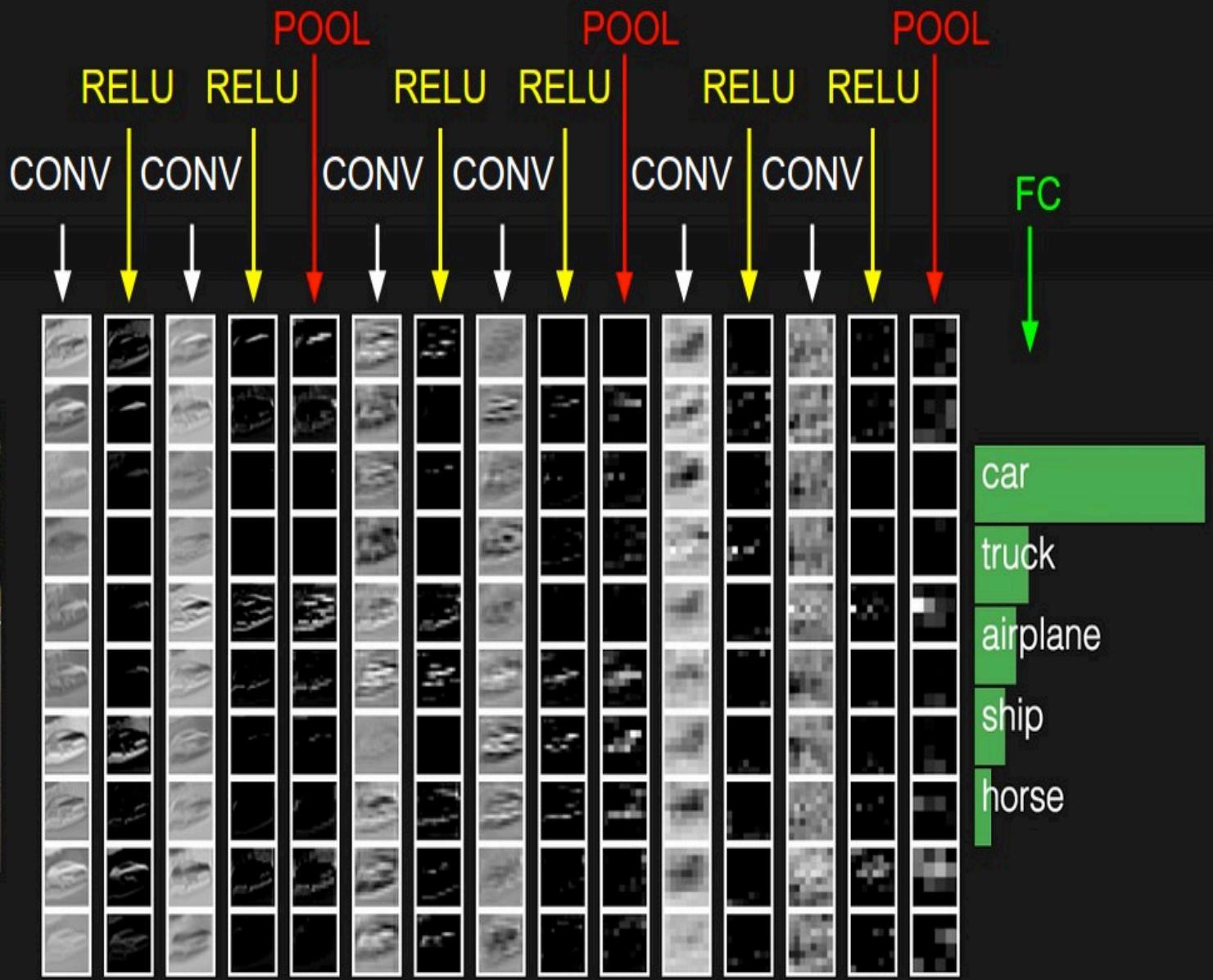


Preview

[From recent
Yann LeCun
slides]

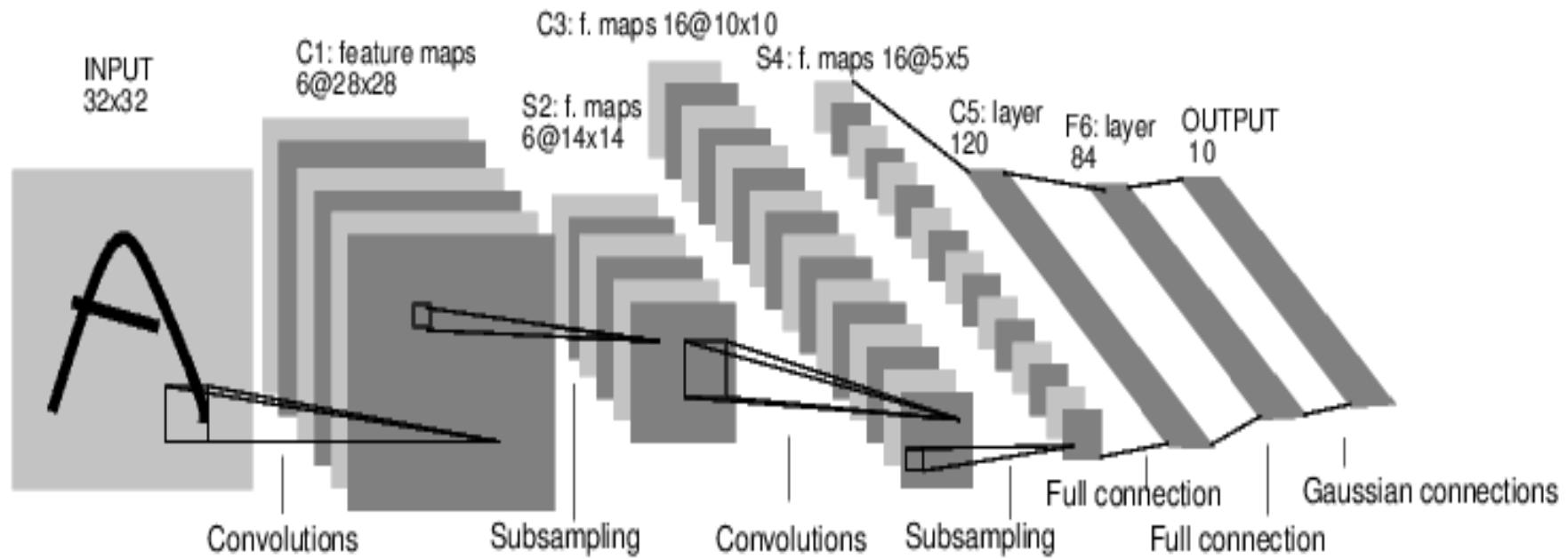


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



Case Study: LeNet-5

[LeCun et al., 1998]

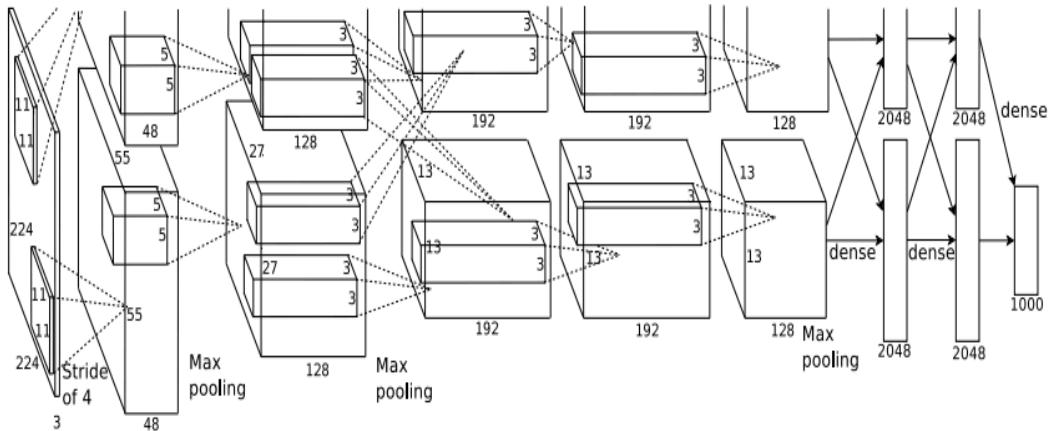


Conv filters were 5x5, applied at stride 1

Subsampling (Pooling) layers were 2x2 applied at stride 2
i.e. architecture is [CONV-POOL-CONV-POOL-CONV-FC]

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

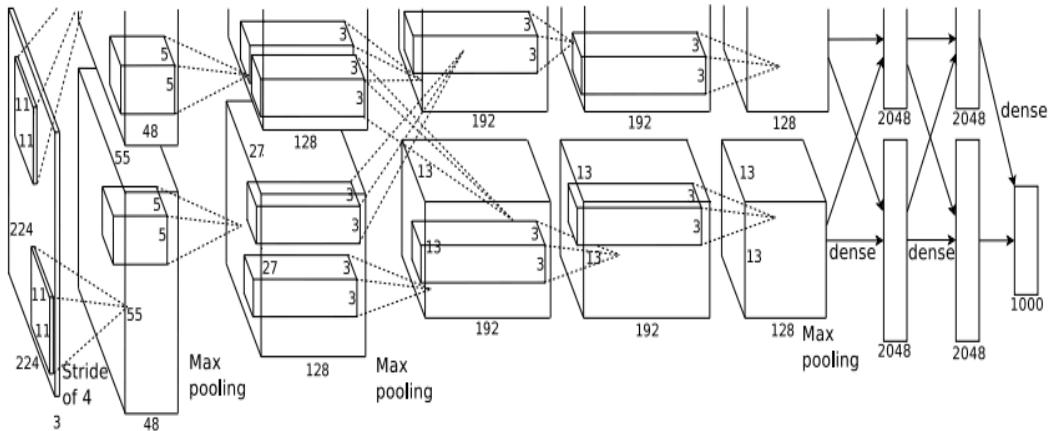
First layer (CONV1): 96 11x11 filters applied at stride 4

=>

Q: what is the output volume size? Hint: $(227-11)/4+1 = 55$

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4

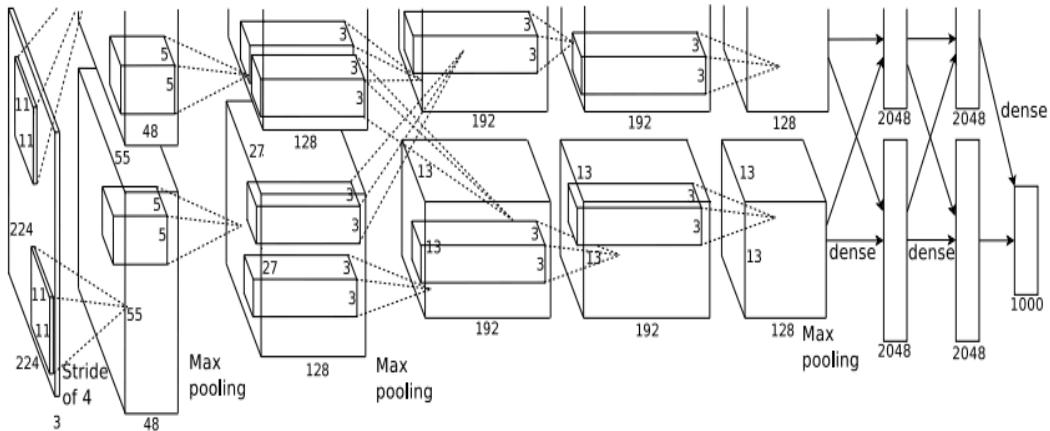
=>

Output volume **[55x55x96]**

Q: What is the total number of parameters in this layer?

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4

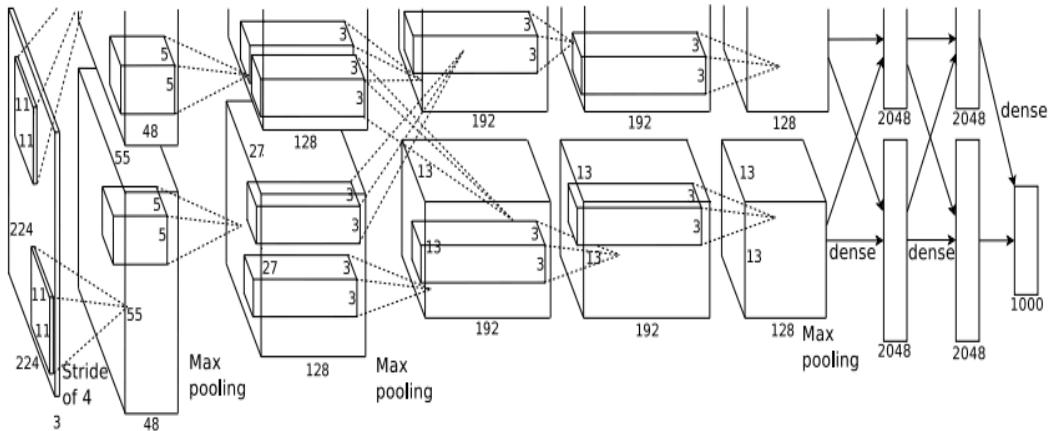
=>

Output volume **[55x55x96]**

Parameters: $(11 \times 11 \times 3) \times 96 = 35K$

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

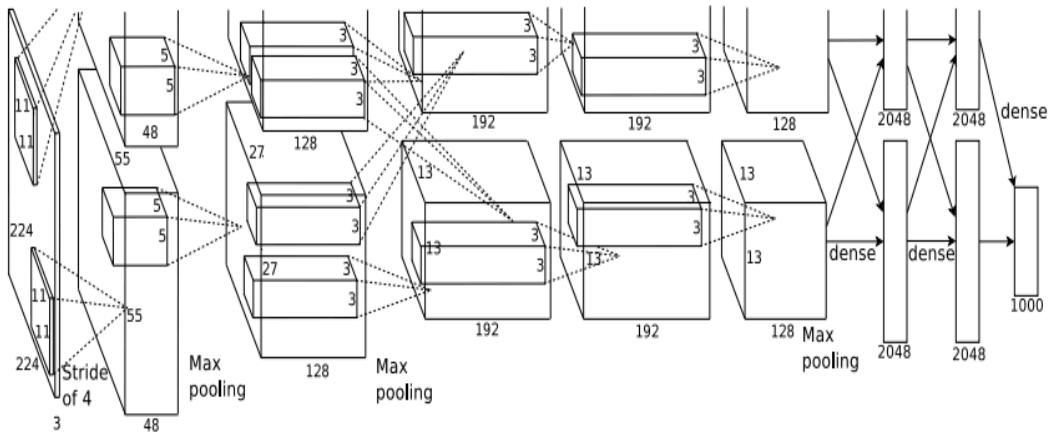
After CONV1: 55x55x96

Second layer (POOL1): 3x3 filters applied at stride 2

Q: what is the output volume size? Hint: $(55-3)/2+1 = 27$

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

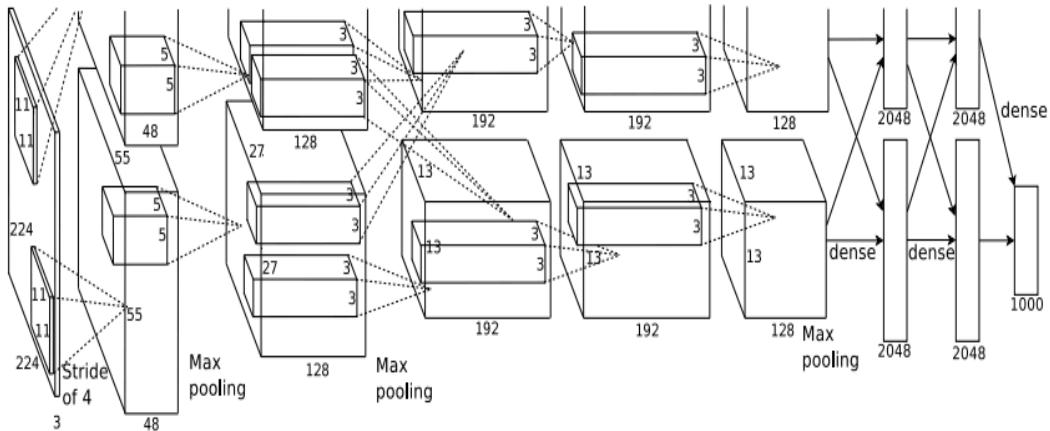
Second layer (POOL1): 3x3 filters applied at stride 2

Output volume: 27x27x96

Q: what is the number of parameters in this layer?

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

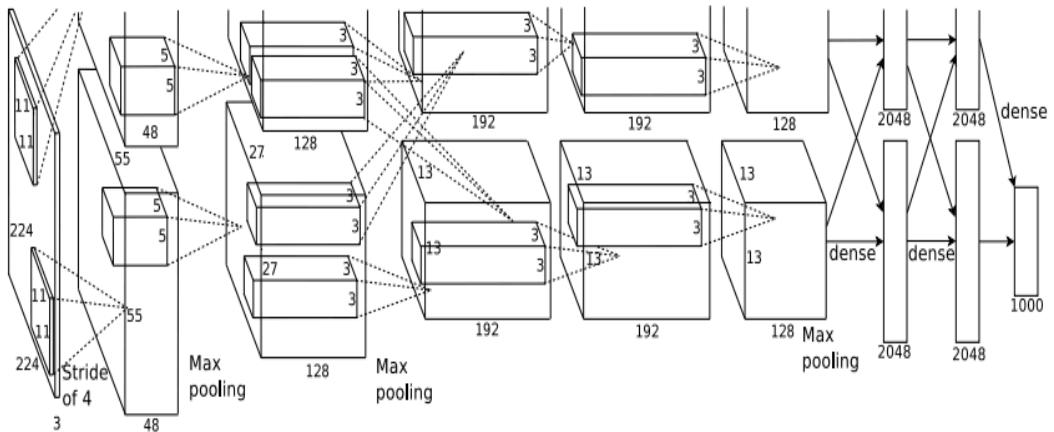
Second layer (POOL1): 3x3 filters applied at stride 2

Output volume: 27x27x96

Parameters: 0!

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

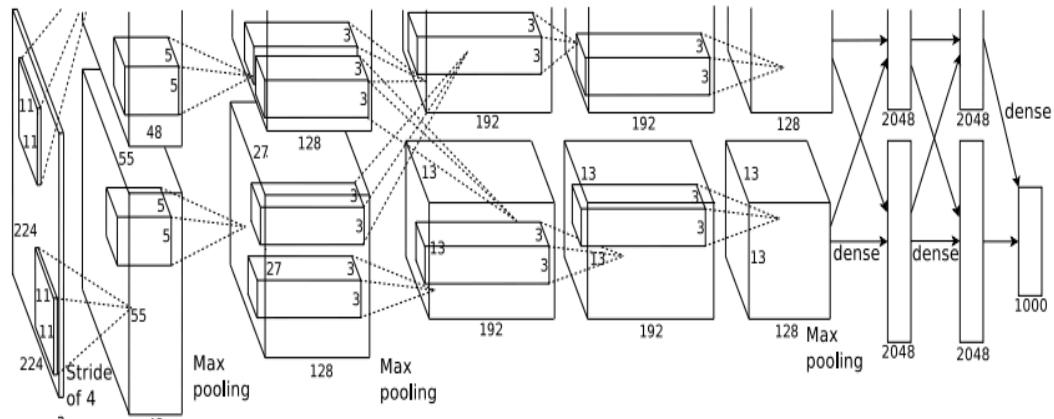
After POOL1: 27x27x96

...

Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:
[227x227x3] INPUT



[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

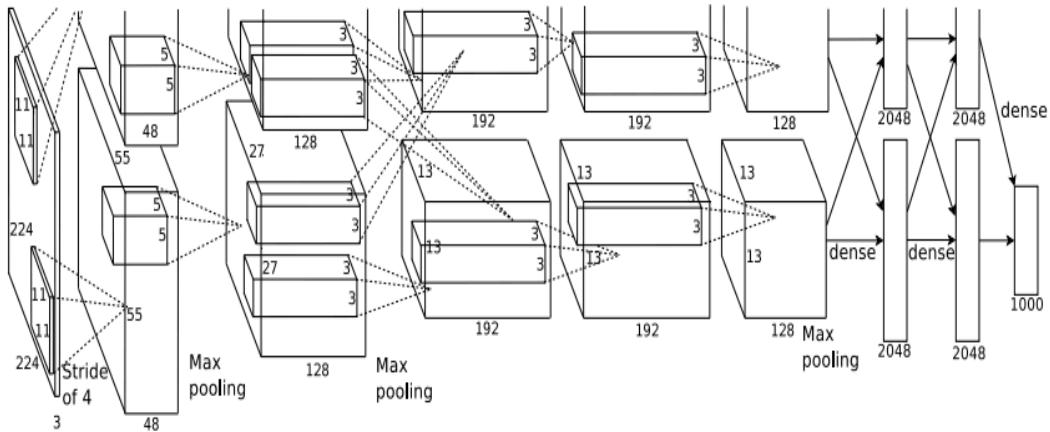
[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)

Case Study: AlexNet

[Krizhevsky et al. 2012]



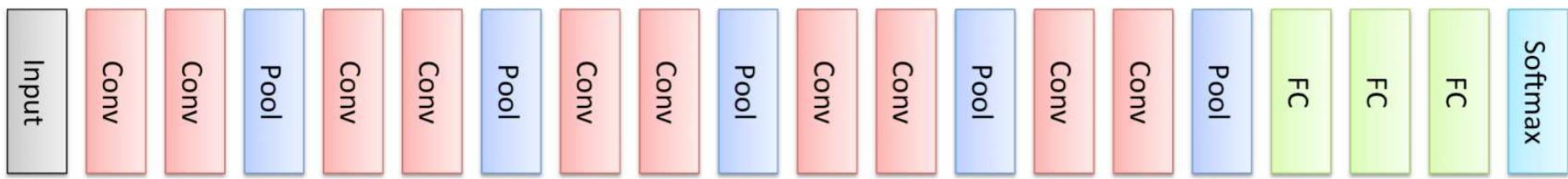
Details/Retrospectives:

- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

VGGNet



Only 3x3 CONV stride 1, pad 1
and 2x2 MAX POOL stride 2

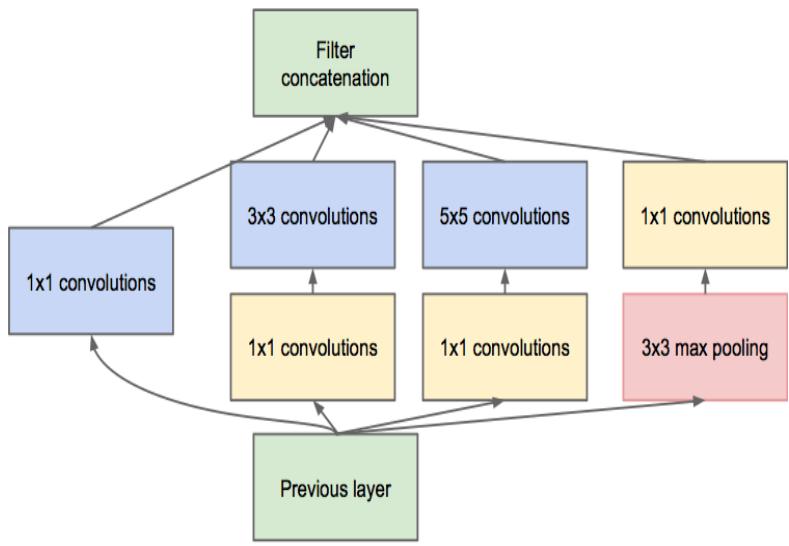
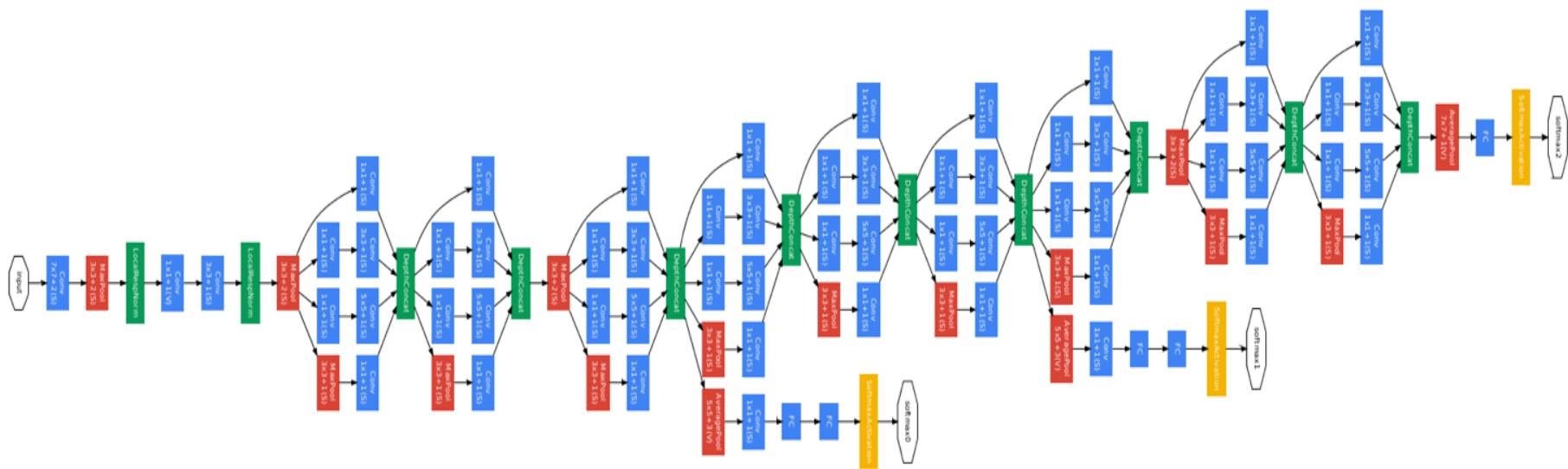
11.2% top 5 error in ILSVRC 2013
→
7.3% top 5 error

Case Study: VGGNet [Simonyan and Zisserman, 2014]

INPUT: [224x224x3] memory: $224 \times 224 \times 3 = 150K$ params: 0
CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2M$ params: $(3 \times 3 \times 3) \times 64 = 1,728$
CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2M$ params: $(3 \times 3 \times 64) \times 64 = 36,864$
POOL2: [112x112x64] memory: $112 \times 112 \times 64 = 800K$ params: 0
CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6M$ params: $(3 \times 3 \times 64) \times 128 = 73,728$
CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6M$ params: $(3 \times 3 \times 128) \times 128 = 147,456$
POOL2: [56x56x128] memory: $56 \times 56 \times 128 = 400K$ params: 0
CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 128) \times 256 = 294,912$
CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 256) \times 256 = 589,824$
CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 256) \times 256 = 589,824$
POOL2: [28x28x256] memory: $28 \times 28 \times 256 = 200K$ params: 0
CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 256) \times 512 = 1,179,648$
CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
POOL2: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: 0
CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
POOL2: [7x7x512] memory: $7 \times 7 \times 512 = 25K$ params: 0
FC: [1x1x4096] memory: 4096 params: $7 \times 7 \times 512 \times 4096 = 102,760,448$
FC: [1x1x4096] memory: 4096 params: $4096 \times 4096 = 16,777,216$
FC: [1x1x1000] memory: 1000 params: $4096 \times 1000 = 4,096,000$ (not counting biases)

Case Study: GoogLeNet

[Szegedy et al., 2014]



Inception module

ILSVRC 2014 winner (6.7% top 5 error)

Case Study: GoogLeNet

| type | patch size/ stride | output size | depth | #1×1 | #3×3 reduce | #3×3 | #5×5 reduce | #5×5 | pool proj | params | ops |
|----------------|-----------------------|----------------|-------|------|----------------|------|----------------|------|--------------|--------|------|
| convolution | 7×7/2 | 112×112×64 | 1 | | | | | | | 2.7K | 34M |
| max pool | 3×3/2 | 56×56×64 | 0 | | | | | | | | |
| convolution | 3×3/1 | 56×56×192 | 2 | | 64 | 192 | | | | 112K | 360M |
| max pool | 3×3/2 | 28×28×192 | 0 | | | | | | | | |
| inception (3a) | | 28×28×256 | 2 | 64 | 96 | 128 | 16 | 32 | 32 | 159K | 128M |
| inception (3b) | | 28×28×480 | 2 | 128 | 128 | 192 | 32 | 96 | 64 | 380K | 304M |
| max pool | 3×3/2 | 14×14×480 | 0 | | | | | | | | |
| inception (4a) | | 14×14×512 | 2 | 192 | 96 | 208 | 16 | 48 | 64 | 364K | 73M |
| inception (4b) | | 14×14×512 | 2 | 160 | 112 | 224 | 24 | 64 | 64 | 437K | 88M |
| inception (4c) | | 14×14×512 | 2 | 128 | 128 | 256 | 24 | 64 | 64 | 463K | 100M |
| inception (4d) | | 14×14×528 | 2 | 112 | 144 | 288 | 32 | 64 | 64 | 580K | 119M |
| inception (4e) | | 14×14×832 | 2 | 256 | 160 | 320 | 32 | 128 | 128 | 840K | 170M |
| max pool | 3×3/2 | 7×7×832 | 0 | | | | | | | | |
| inception (5a) | | 7×7×832 | 2 | 256 | 160 | 320 | 32 | 128 | 128 | 1072K | 54M |
| inception (5b) | | 7×7×1024 | 2 | 384 | 192 | 384 | 48 | 128 | 128 | 1388K | 71M |
| avg pool | 7×7/1 | 1×1×1024 | 0 | | | | | | | | |
| dropout (40%) | | 1×1×1024 | 0 | | | | | | | | |
| linear | | 1×1×1000 | 1 | | | | | | | 1000K | 1M |
| softmax | | 1×1×1000 | 0 | | | | | | | | |

Fun features:

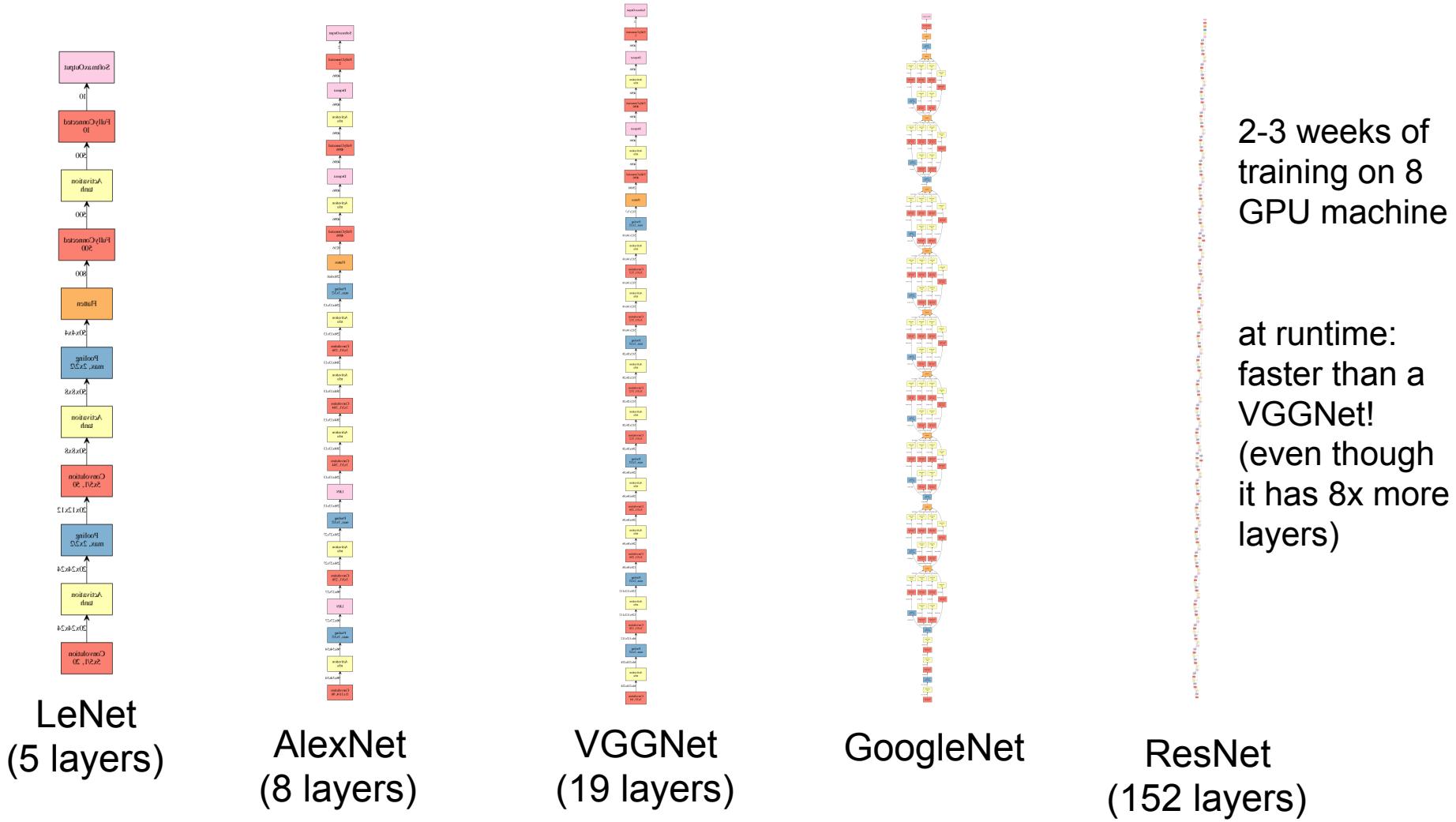
- Only 5 million params!
(Removes FC layers completely)

Compared to AlexNet:

- 12X less params
- 2x more compute
- 6.67% (vs. 16.4%)

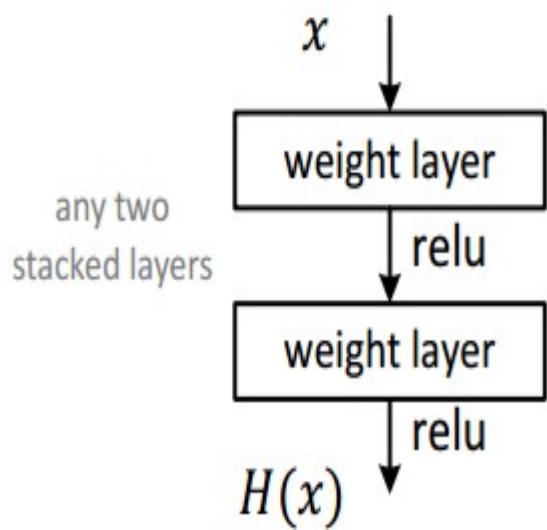
Case Study: ResNet [He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)

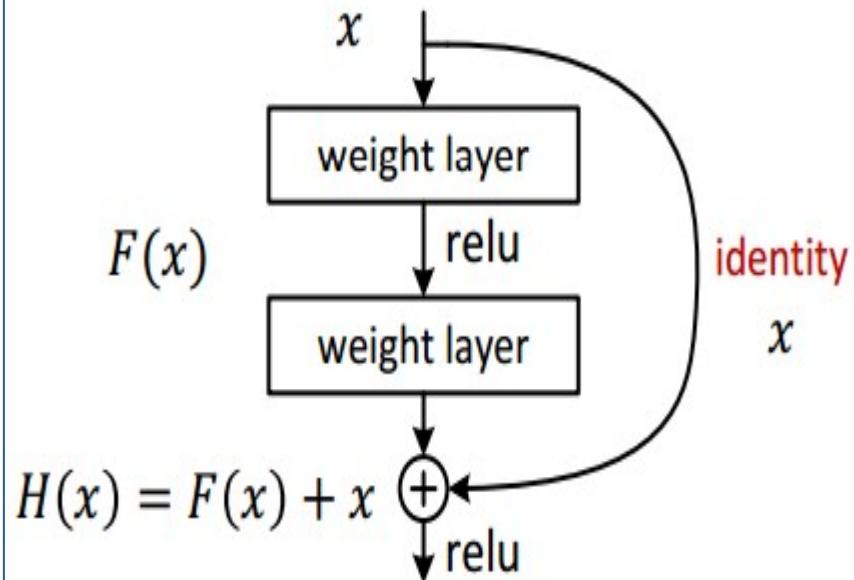


Case Study: ResNet [He et al., 2015]

- Plain net



- Residual net



ILSVRC image classification (CLS) task

Steel drum



1000 object classes

1,431,167 images

CLS-LOC

ILSVRC image classification (CLS) task

Steel drum



Output:

- Scale
- T-shirt
- Steel drum
- Drumstick
- Mud turtle



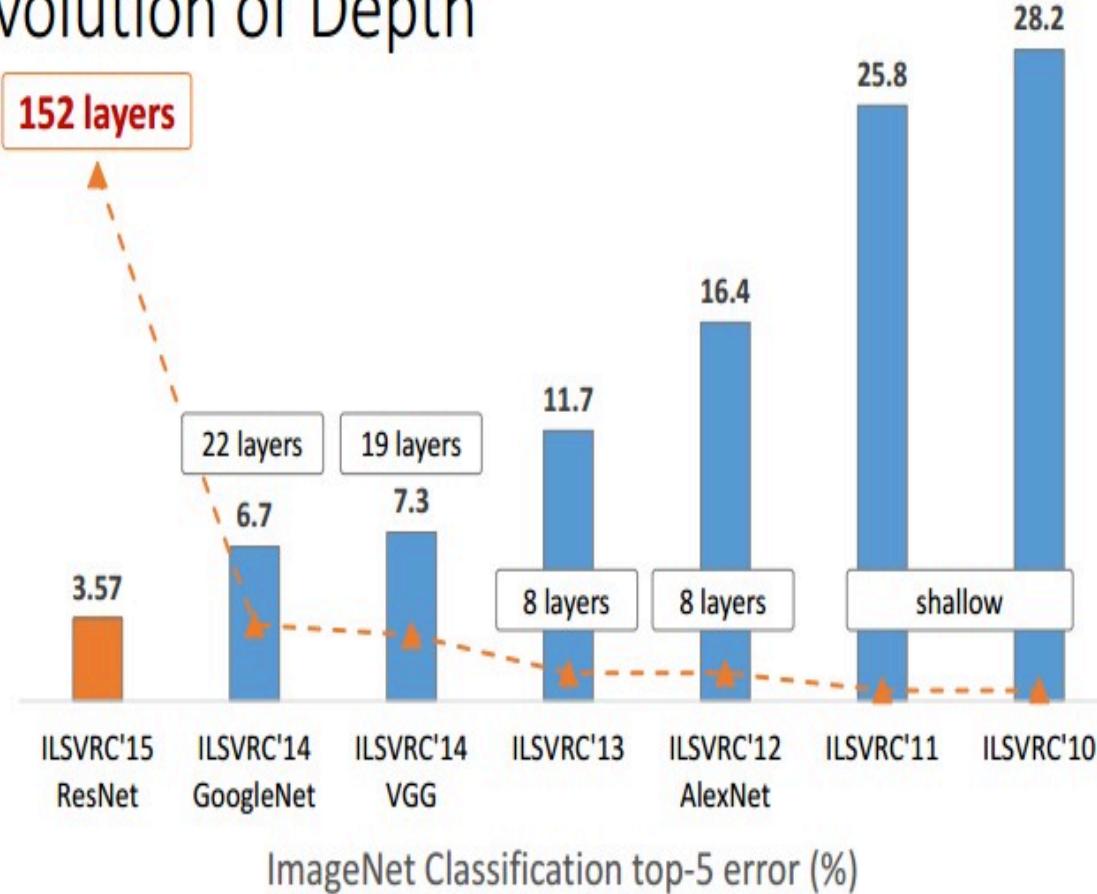
Output:

- Scale
- T-shirt
- Giant panda
- Drumstick
- Mud turtle



$$\text{Error} = \frac{1}{100,000} \sum_{\substack{100,000 \\ \text{images}}} 1[\text{incorrect on image i}]$$

Revolution of Depth

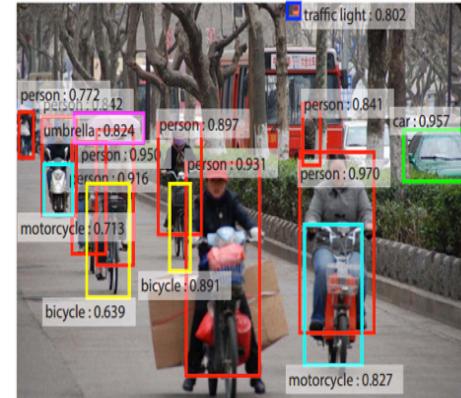
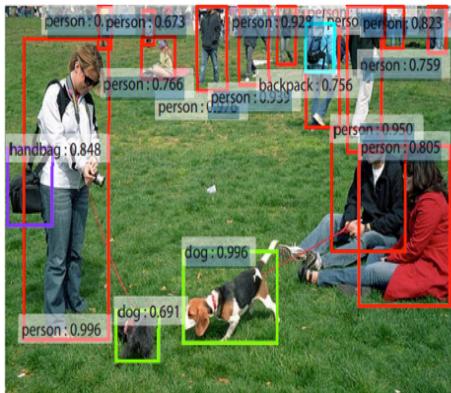


Part 2 Summary

- ConvNets stack CONV,POOL,FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like
 $[(\text{CONV-RELU})^* \text{N-POOL?}]^* \text{M-}(\text{FC-RELU})^* \text{K,SOFTMAX}$
where N is usually up to ~5, M is large, $0 \leq K \leq 2$.
 - but recent advances such as ResNet/GoogLeNet challenge this paradigm

Part 3. Object Detection with CNNs

Object Detection



Results from Faster R-CNN, Ren et al 2015

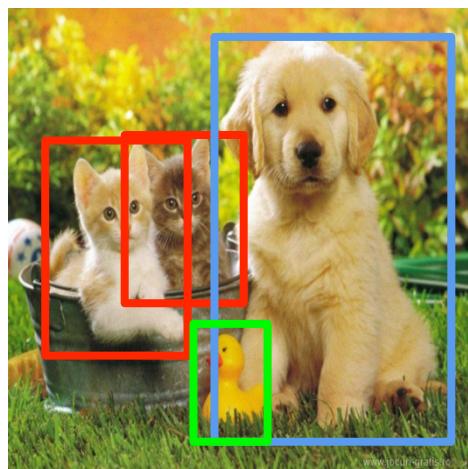
Computer Vision Tasks

Classification



CAT

Object Detection



CAT, DOG, DUCK

Instance Segmentation



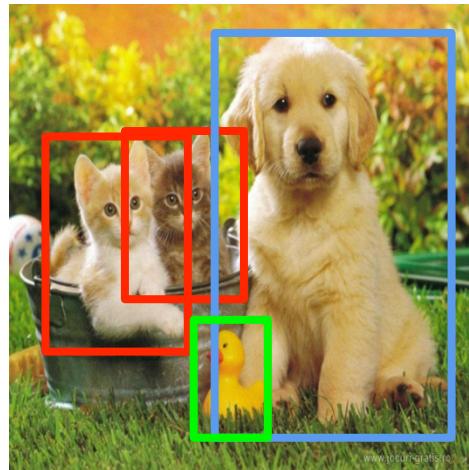
CAT, DOG, DUCK

Computer Vision Tasks

Classification



Object Detection



Instance
Segmentation



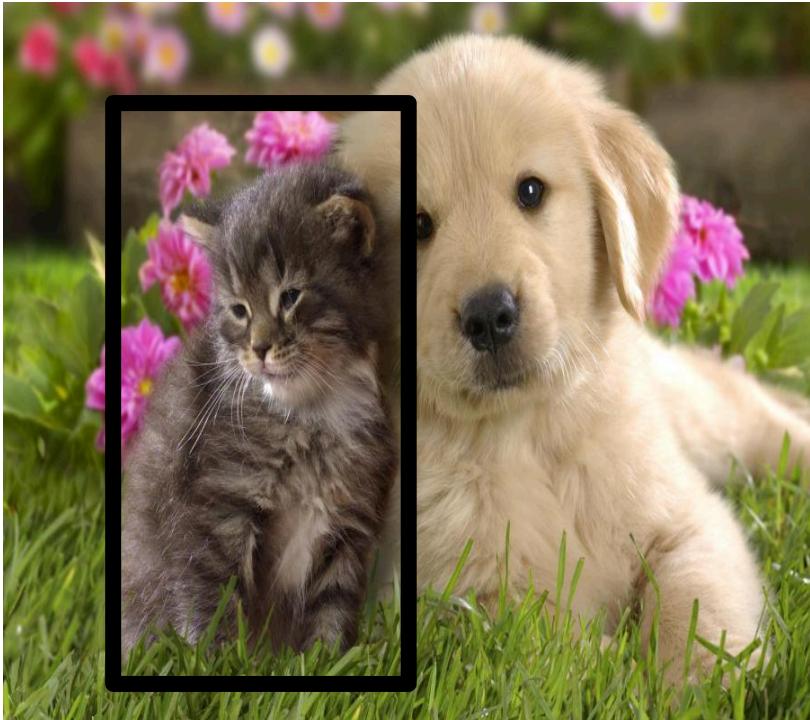
Detection as Classification



CAT? NO

DOG? NO

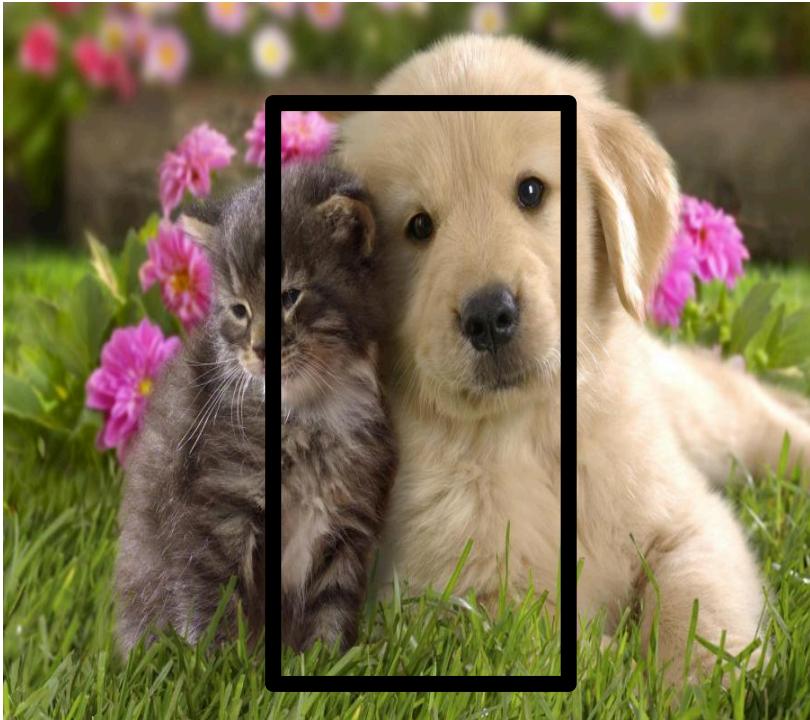
Detection as Classification



CAT? YES!

DOG? NO

Detection as Classification



CAT? NO

DOG? NO

Detection as Classification

Problem: Need to test many positions and scales

Solution: If your classifier is fast enough, just do it

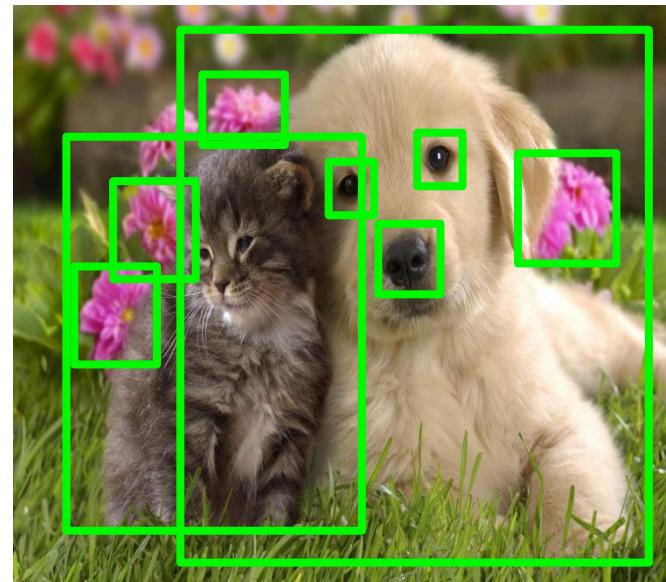
Detection as Classification

Problem: Need to test many positions and scales,
and use a computationally demanding classifier (CNN)

Solution: Only look at a tiny subset of possible positions

Region Proposals

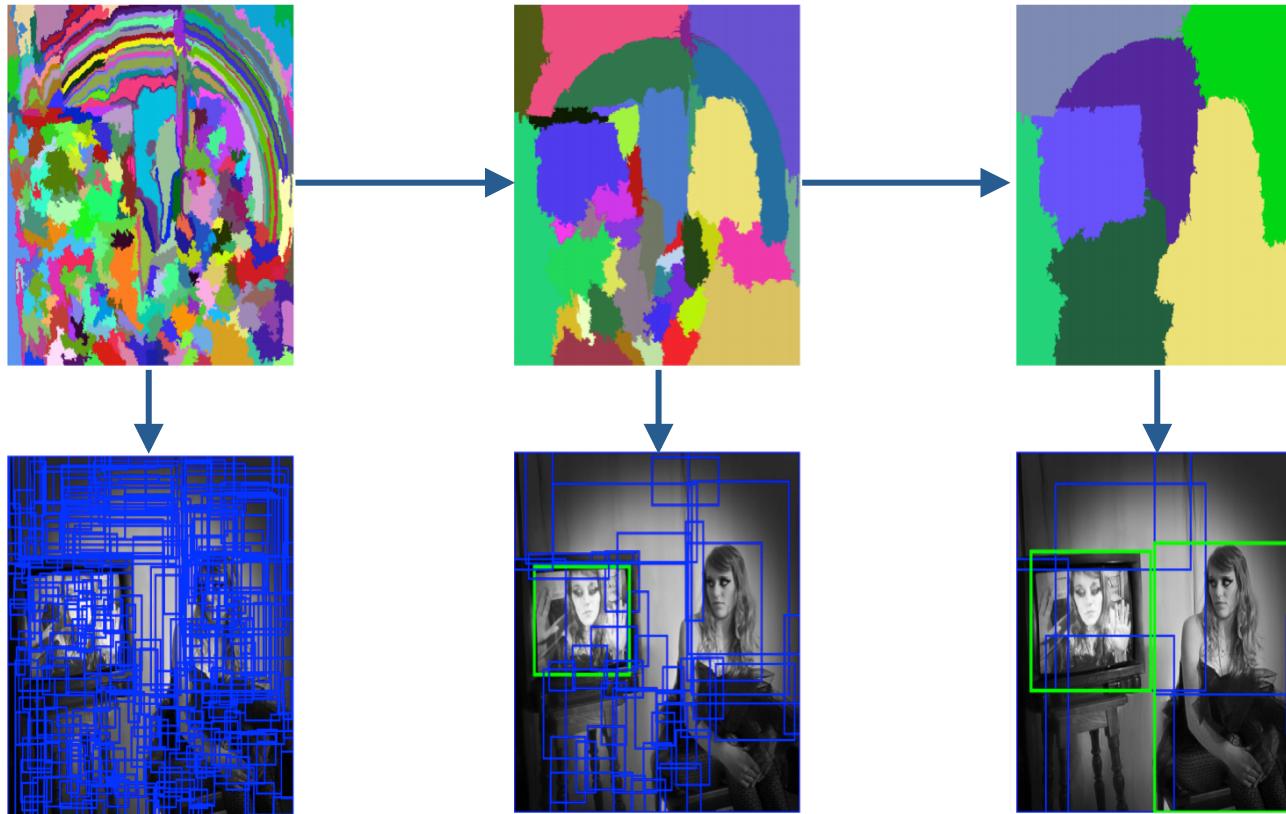
- Find “blobby” image regions that are likely to contain objects
- “Class-agnostic” object detector
- Look for “blob-like” regions



Region Proposals: Selective Search

Bottom-up segmentation, merging regions at multiple scales

Convert
regions
to boxes



Uijlings et al, "Selective Search for Object Recognition", IJCV 2013

Region Proposals: Many other choices

| Method | Approach | Outputs Segments | Outputs Score | Control #proposals | Time (sec.) | Repeatability | Recall Results | Detection Results |
|-----------------------|----------------|------------------|---------------|--------------------|-------------|---------------|----------------|-------------------|
| Bing [18] | Window scoring | | ✓ | ✓ | 0.2 | *** | * | . |
| CPMC [19] | Grouping | ✓ | ✓ | ✓ | 250 | - | ** | * |
| EdgeBoxes [20] | Window scoring | | ✓ | ✓ | 0.3 | ** | *** | *** |
| Endres [21] | Grouping | ✓ | ✓ | ✓ | 100 | - | *** | ** |
| Geodesic [22] | Grouping | ✓ | | ✓ | 1 | * | *** | ** |
| MCG [23] | Grouping | ✓ | ✓ | ✓ | 30 | * | *** | *** |
| Objectness [24] | Window scoring | | ✓ | ✓ | 3 | . | * | . |
| Rahtu [25] | Window scoring | | ✓ | ✓ | 3 | . | . | * |
| RandomizedPrim's [26] | Grouping | ✓ | | ✓ | 1 | * | * | ** |
| Rantalankila [27] | Grouping | ✓ | | ✓ | 10 | ** | . | ** |
| Rigor [28] | Grouping | ✓ | | ✓ | 10 | * | ** | ** |
| SelectiveSearch [29] | Grouping | ✓ | ✓ | ✓ | 10 | ** | *** | *** |
| Gaussian | | | | ✓ | 0 | . | . | * |
| SlidingWindow | | | | ✓ | 0 | *** | . | . |
| Superpixels | | ✓ | | | 1 | * | . | . |
| Uniform | | | | ✓ | 0 | . | . | . |

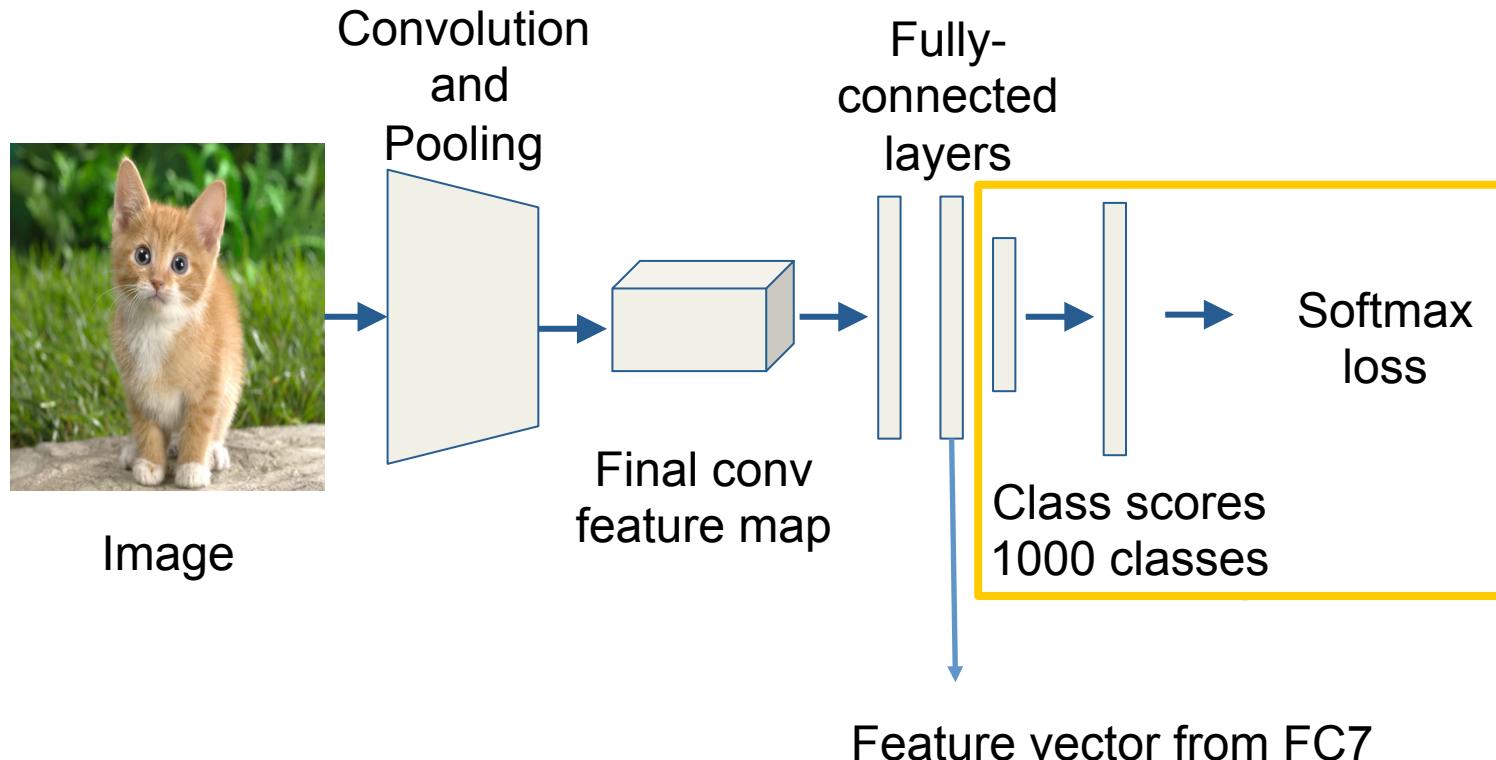
Region Proposals: Many other choices

| Method | Approach | Outputs Segments | Outputs Score | Control #proposals | Time (sec.) | Repeatability | Recall Results | Detection Results |
|-----------------------|----------------|------------------|---------------|--------------------|-------------|---------------|----------------|-------------------|
| Bing [18] | Window scoring | | ✓ | ✓ | 0.2 | *** | * | . |
| CPMC [19] | Grouping | ✓ | ✓ | ✓ | 250 | - | ** | * |
| EdgeBoxes [20] | Window scoring | | ✓ | ✓ | 0.3 | ** | *** | *** |
| Endres [21] | Grouping | ✓ | ✓ | ✓ | 100 | - | *** | ** |
| Geodesic [22] | Grouping | ✓ | | ✓ | 1 | * | *** | ** |
| MCG [23] | Grouping | ✓ | ✓ | ✓ | 30 | * | *** | *** |
| Objectness [24] | Window scoring | | ✓ | ✓ | 3 | . | * | . |
| Rahtu [25] | Window scoring | | ✓ | ✓ | 3 | . | . | * |
| RandomizedPrim's [26] | Grouping | ✓ | | ✓ | 1 | * | * | ** |
| Rantalankila [27] | Grouping | ✓ | | ✓ | 10 | ** | . | ** |
| Rigor [28] | Grouping | ✓ | | ✓ | 10 | * | ** | ** |
| SelectiveSearch [29] | Grouping | ✓ | ✓ | ✓ | 10 | ** | *** | *** |
| Gaussian | | | | ✓ | 0 | . | . | * |
| SlidingWindow | | | | ✓ | 0 | *** | . | . |
| Superpixels | | ✓ | | | 1 | * | . | . |
| Uniform | | | | ✓ | 0 | . | . | . |

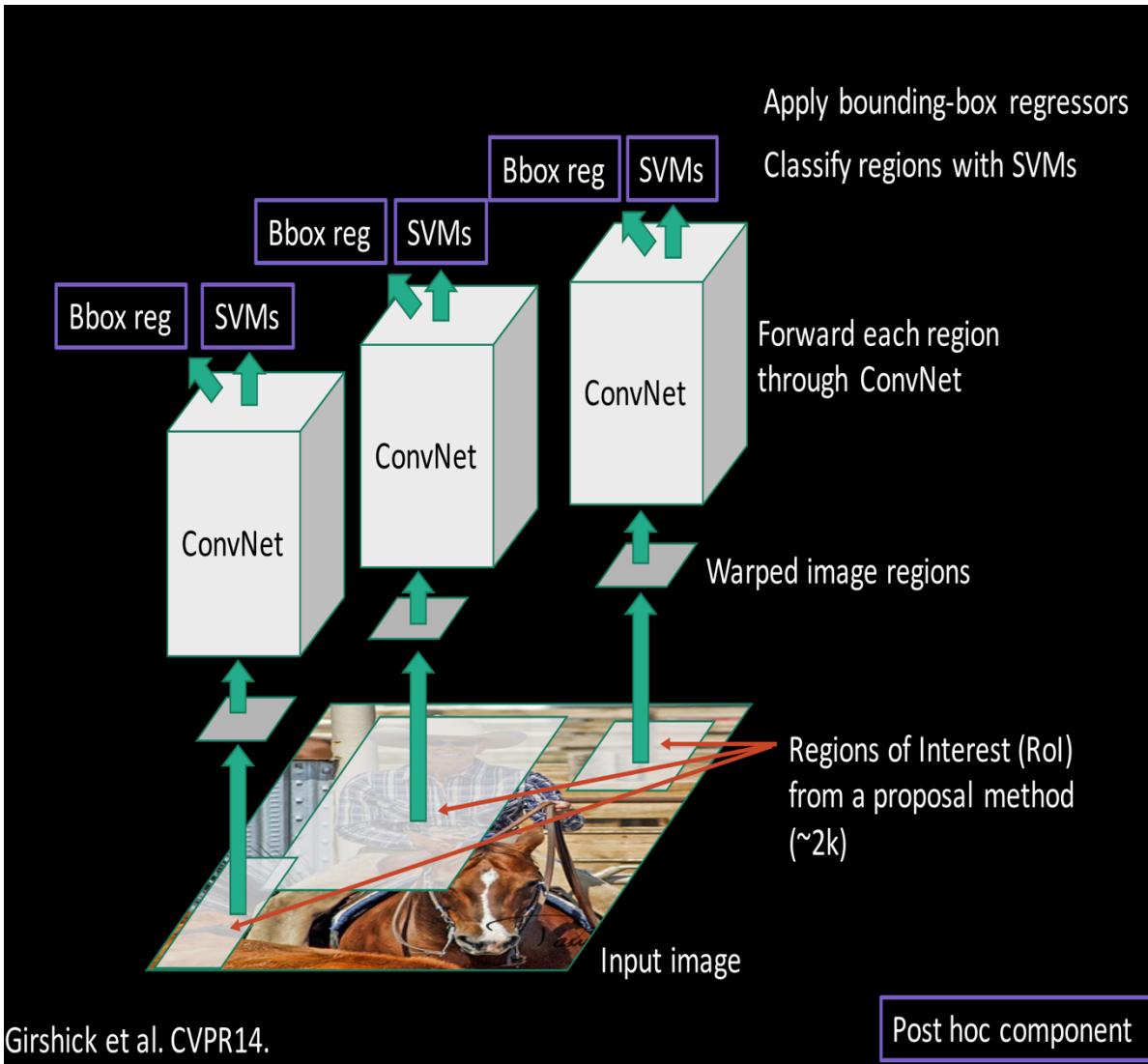
Object Detection: Datasets

| | PASCAL VOC (2010) | ImageNet Detection (ILSVRC 2014) | MS-COCO (2014) |
|--------------------------------------|-------------------------|---|-------------------|
| Number of classes | 20 | 200 | 80 |
| Number of images (train + val) | ~20k | ~470k | ~120k |
| Mean objects per image | 2.4 | 1.1 | 7.2 |

Use CNN to extract features



Putting it together: R-CNN

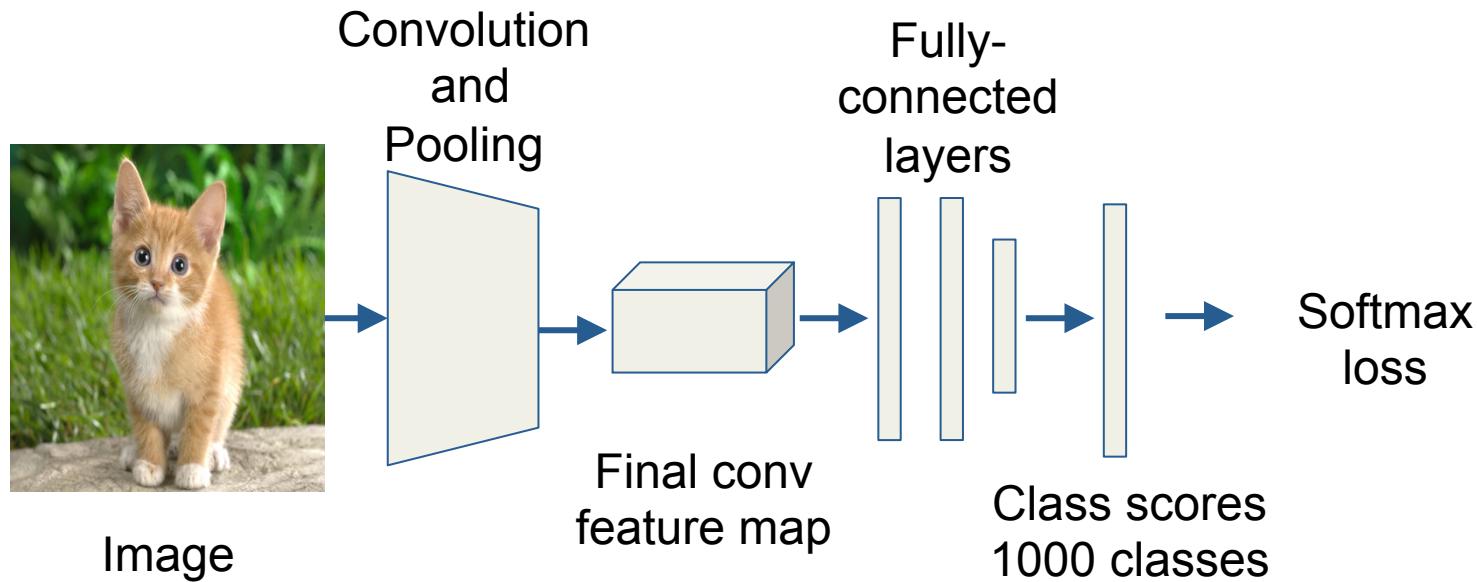


Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014

Slide credit: Ross Girshick

R-CNN Training

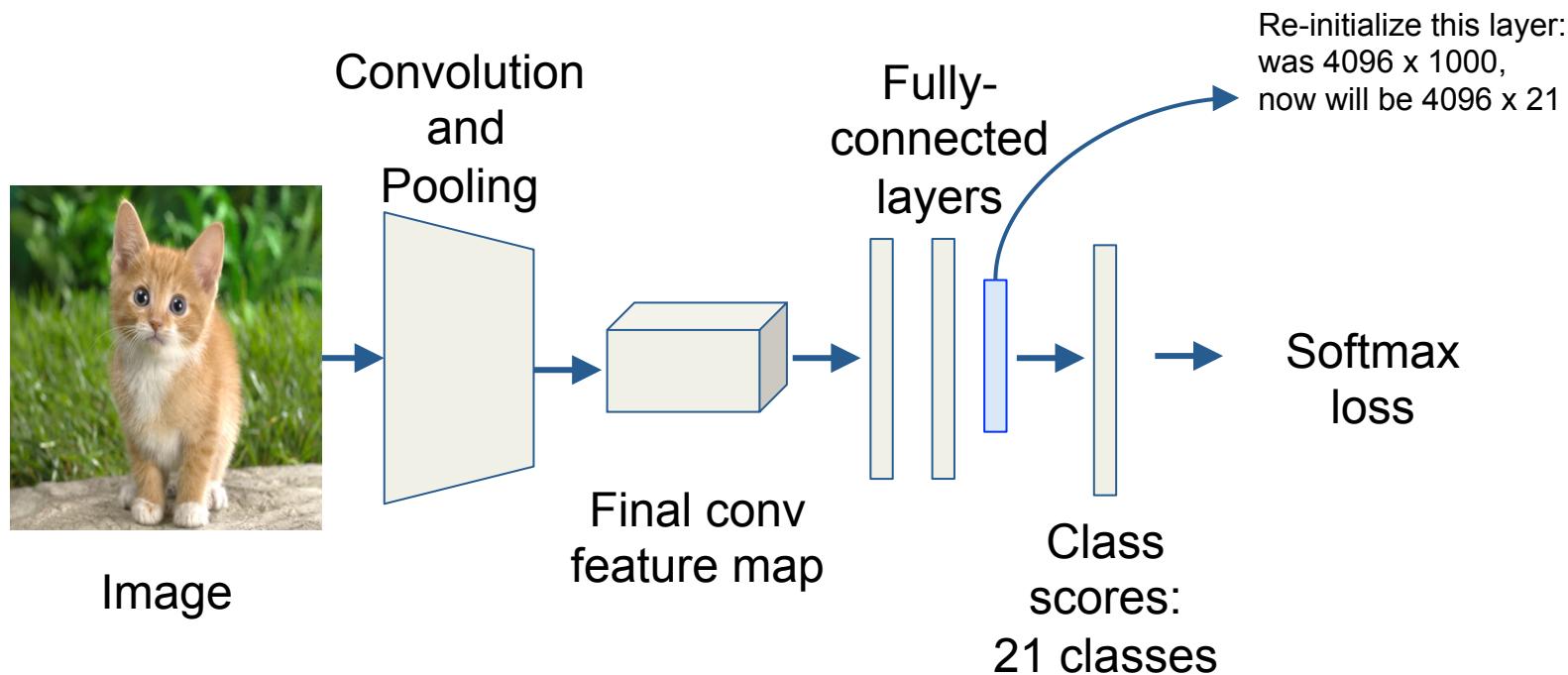
Step 1: Train (or download) a classification model for ImageNet (AlexNet)



R-CNN Training

Step 2: Fine-tune model for detection

- Instead of 1000 ImageNet classes, want 20 object classes + background
- Throw away final fully-connected layer, reinitialize from scratch
- Keep training model using positive / negative regions from detection images



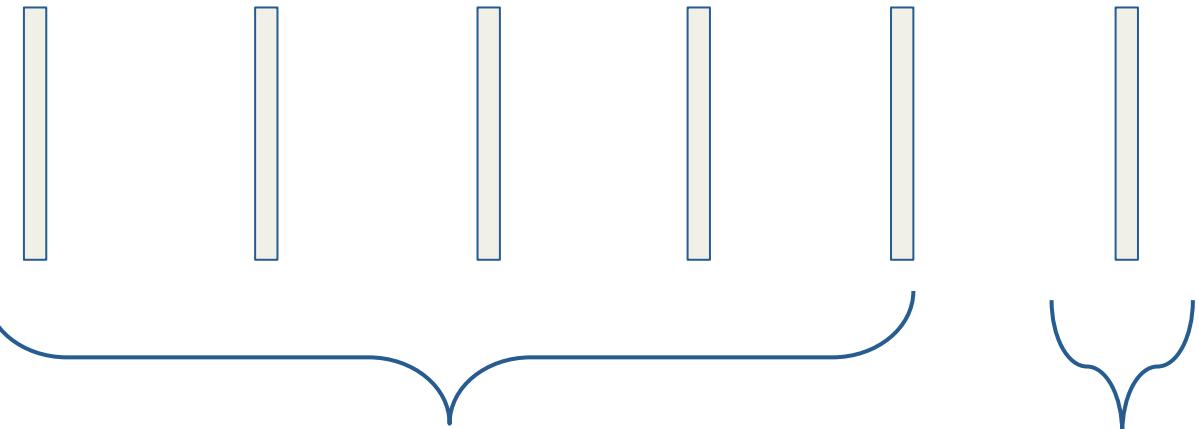
R-CNN Training

Step 3: Train one binary SVM per class to classify region features

Training image regions



Cached region features



Positive samples for cat SVM

Negative samples for cat SVM

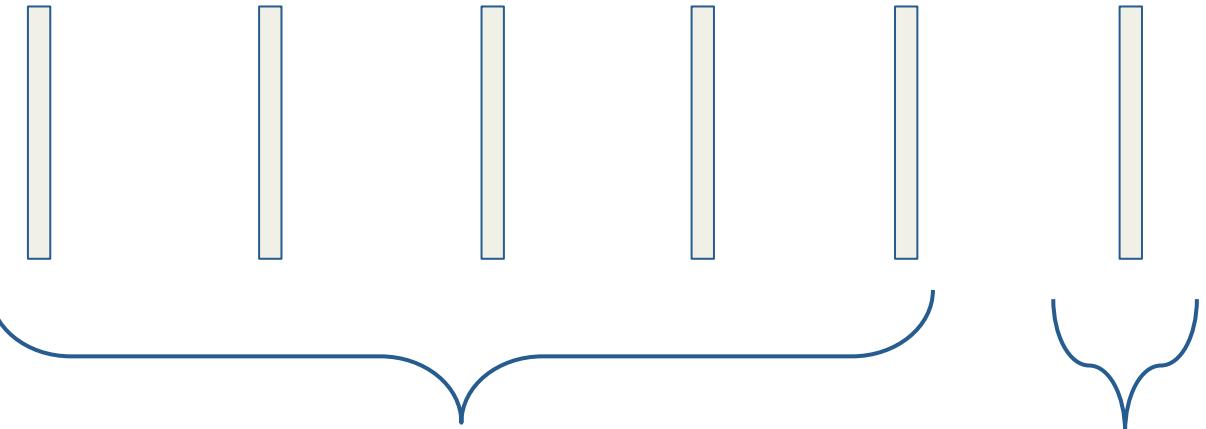
R-CNN Training

Step 3: Train one binary SVM per class to classify region features

Training image regions



Cached region features



Negative samples for
dog SVM

Positive samples for dog
SVM

R-CNN Training

Step 4 (optional: bbox regression): For each class, train a linear regression model to map from cached features to offsets to GT boxes to make up for “slightly wrong” proposals

Training image regions



Cached region features



Regression targets
(dx , dy , dw , dh)
Normalized coordinates

(0, 0, 0, 0)
Proposal is good

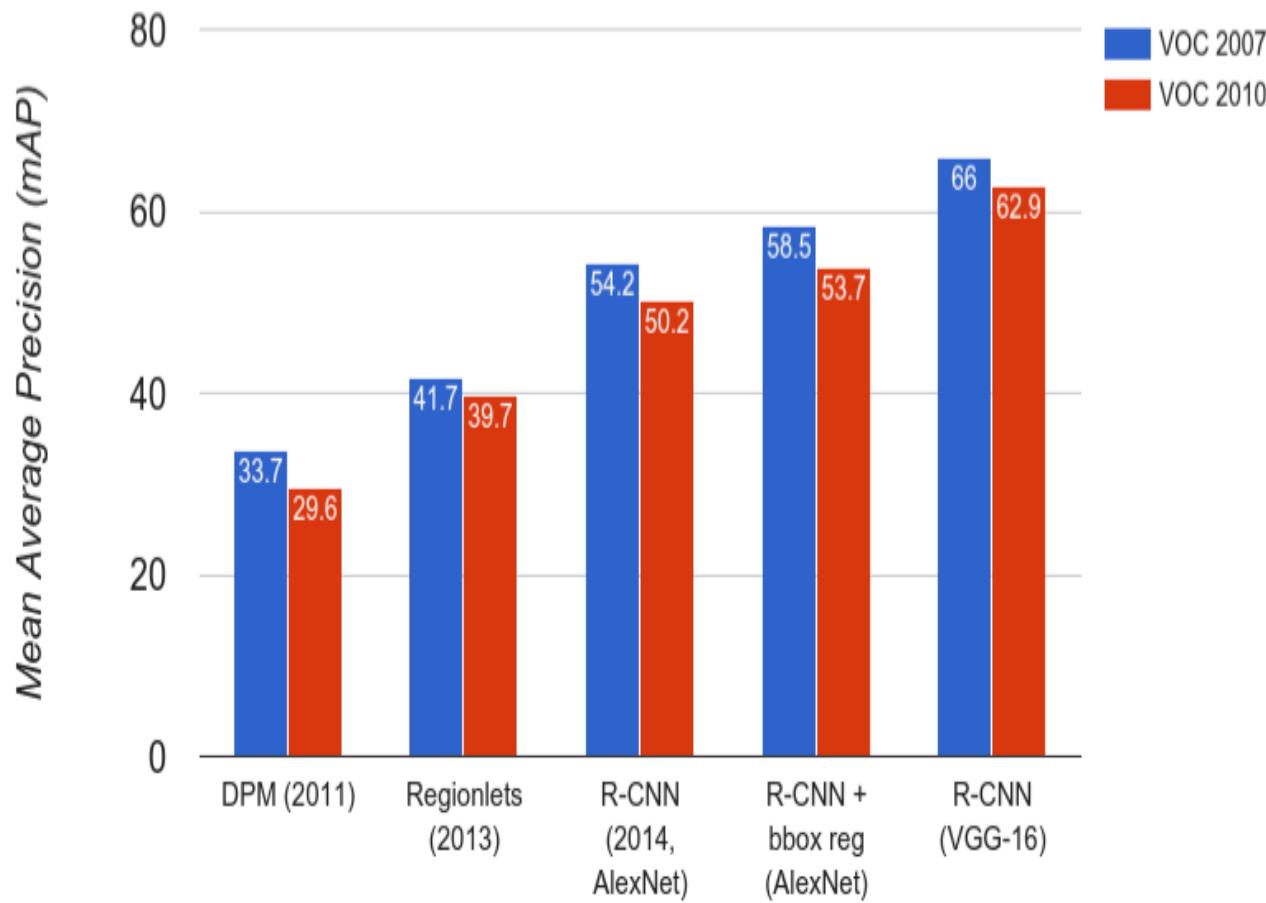
(.25, 0, 0, 0)
Proposal too far to left

(0, 0, -0.125, 0)
Proposal too wide

Object Detection: Evaluation

- We use a metric called “mean average precision” (mAP)
- Compute average precision (AP) separately for each class, then average over classes
- A detection is a true positive if it has IoU with a ground-truth box greater than some threshold (usually 0.5) (mAP@0.5)
- Combine all detections from all test images to draw a precision / recall curve for each class; AP is area under the curve
- mAP is a number from 0 to 100; high is good

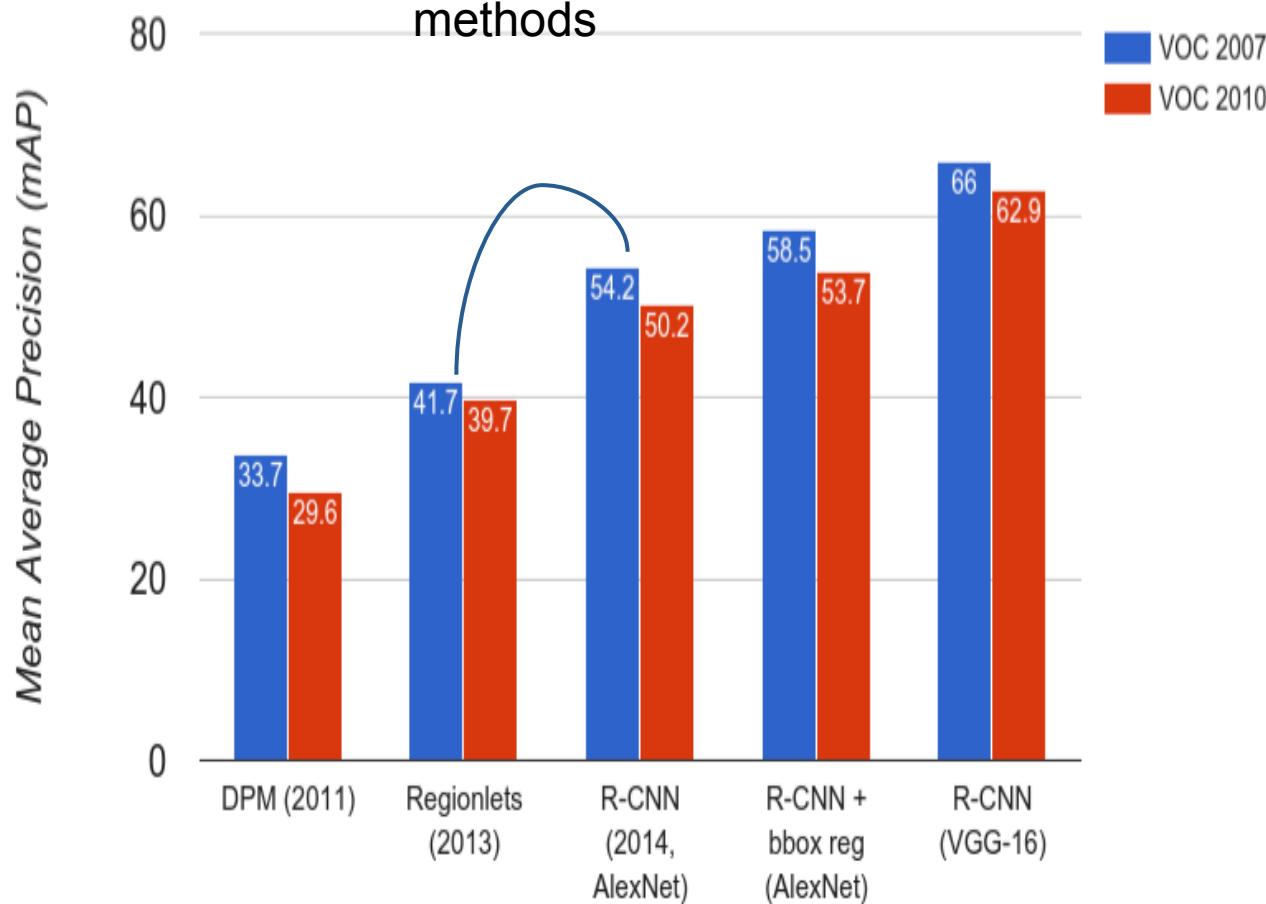
R-CNN Results



Wang et al, "Regionlets for Generic Object Detection", ICCV 2013

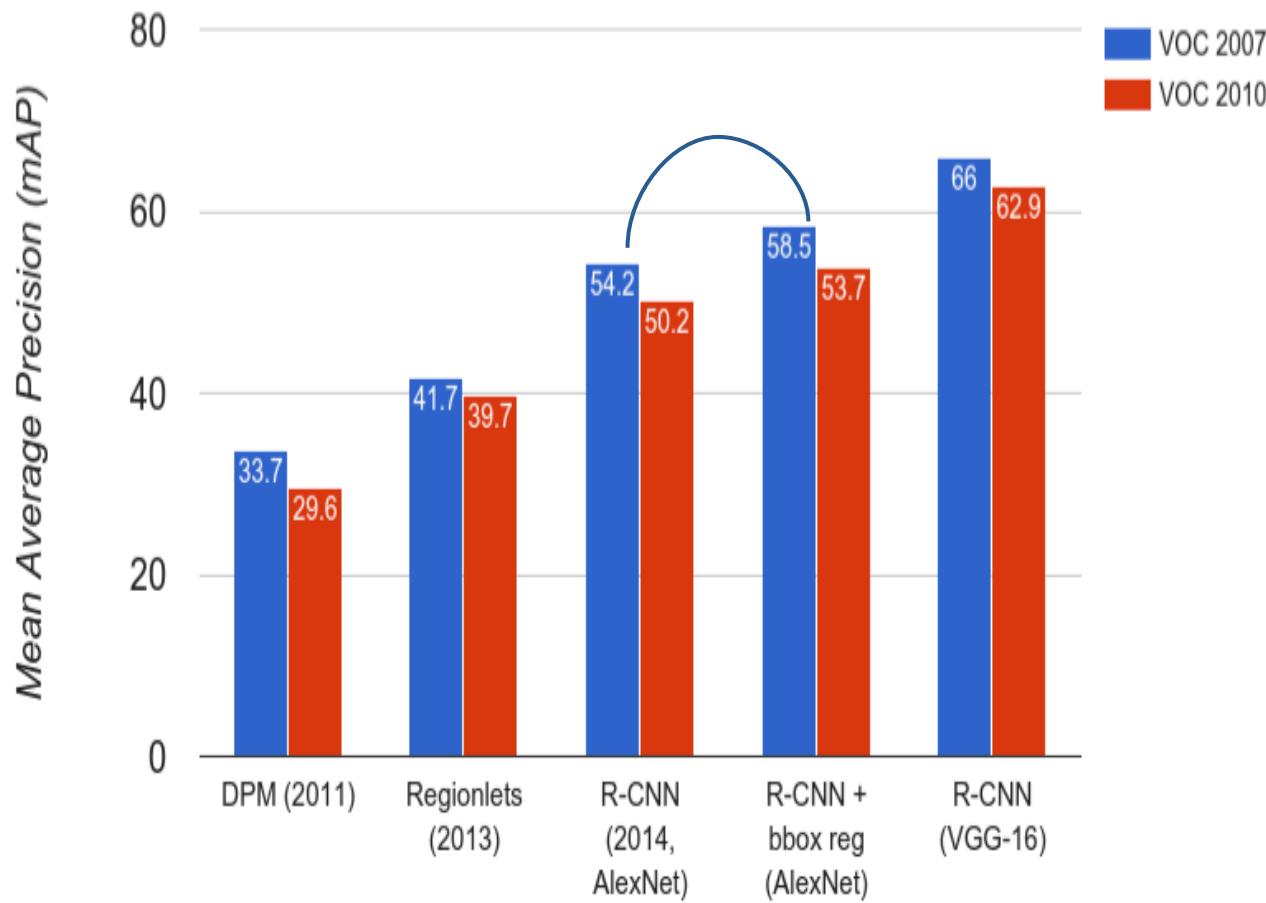
R-CNN Results

Big improvement
compared to pre-CNN
methods



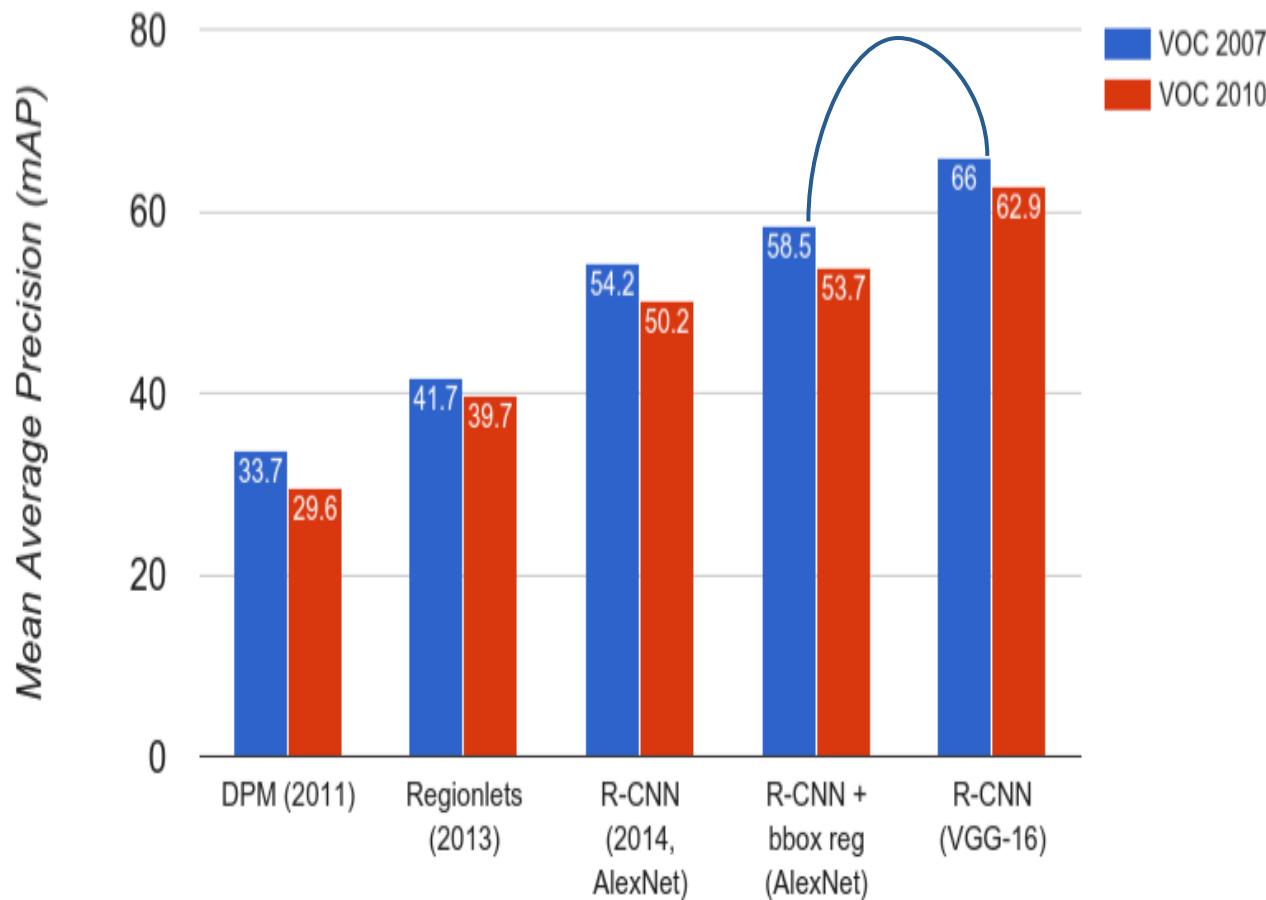
R-CNN Results

Bounding box
regression helps a bit



R-CNN Results

Features from a deeper
network help a lot



Part 3. Summary

Object Detection:

- Find a variable number of objects by classifying image regions
- Before CNNs: dense multiscale sliding window (HoG, DPM)
- Avoid dense sliding window with region proposals
- R-CNN: Selective Search + CNN classification / regression
- Deeper networks do better

Part 4. Segmentation with CNNs

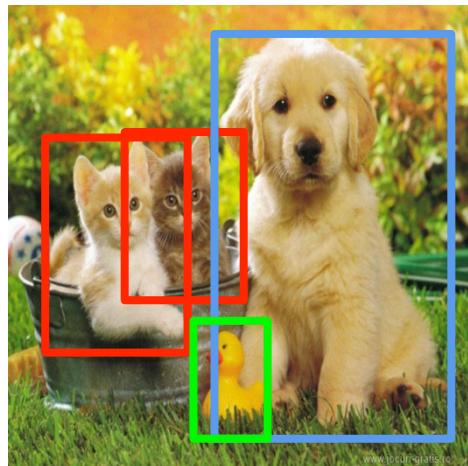
Computer Vision Tasks

Classification



CAT

Object Detection



CAT, DOG, DUCK

Instance Segmentation



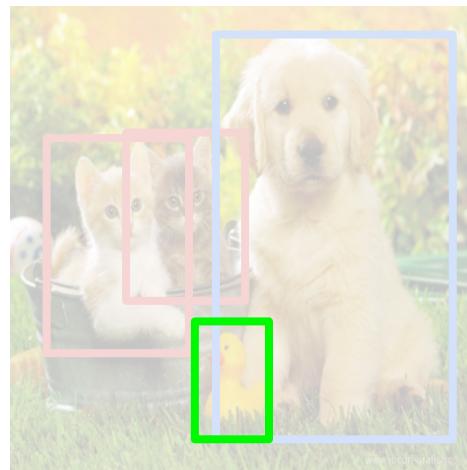
CAT, DOG, DUCK

Computer Vision Tasks

Classification



Object Detection



Segmentation

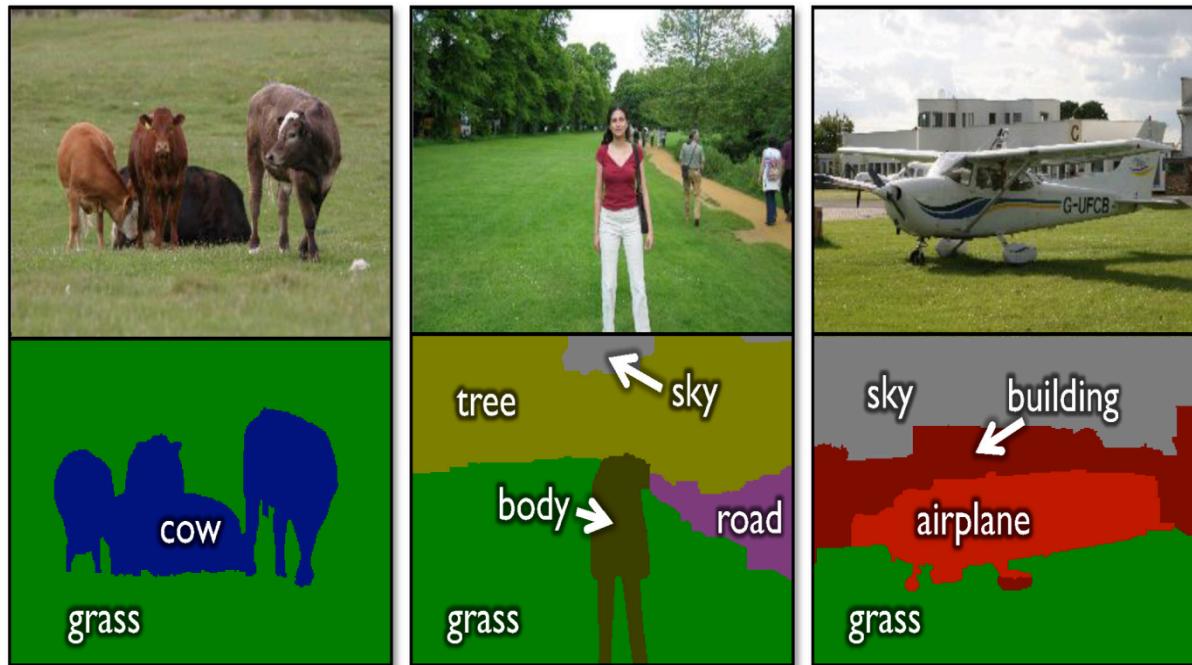


Semantic Segmentation

Label every pixel!

Don't differentiate instances (cows)

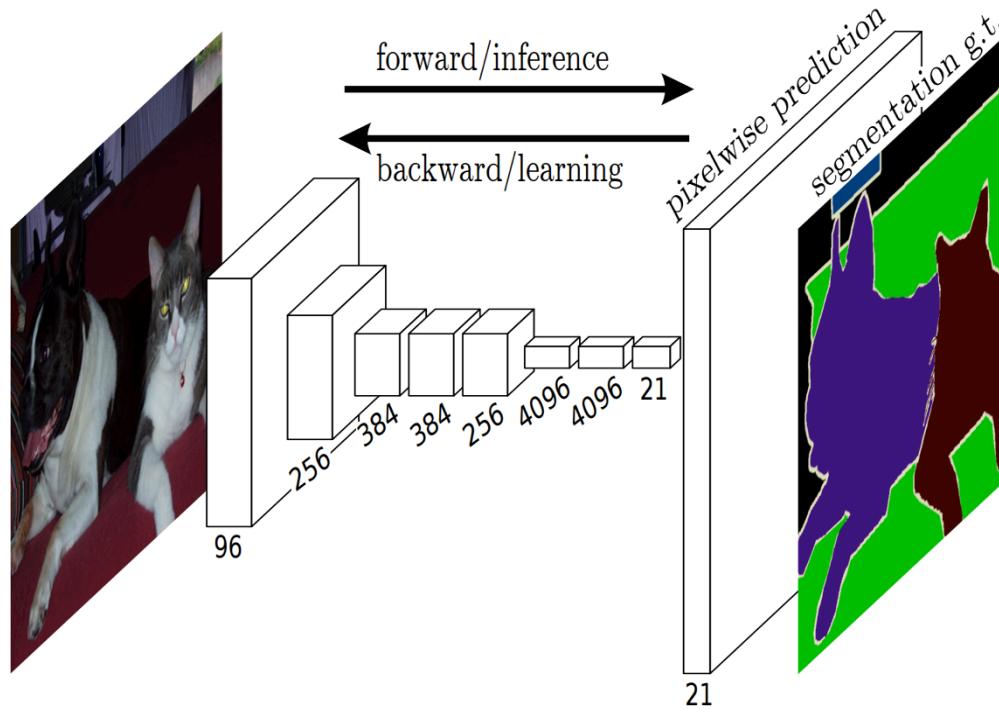
Classic computer vision problem



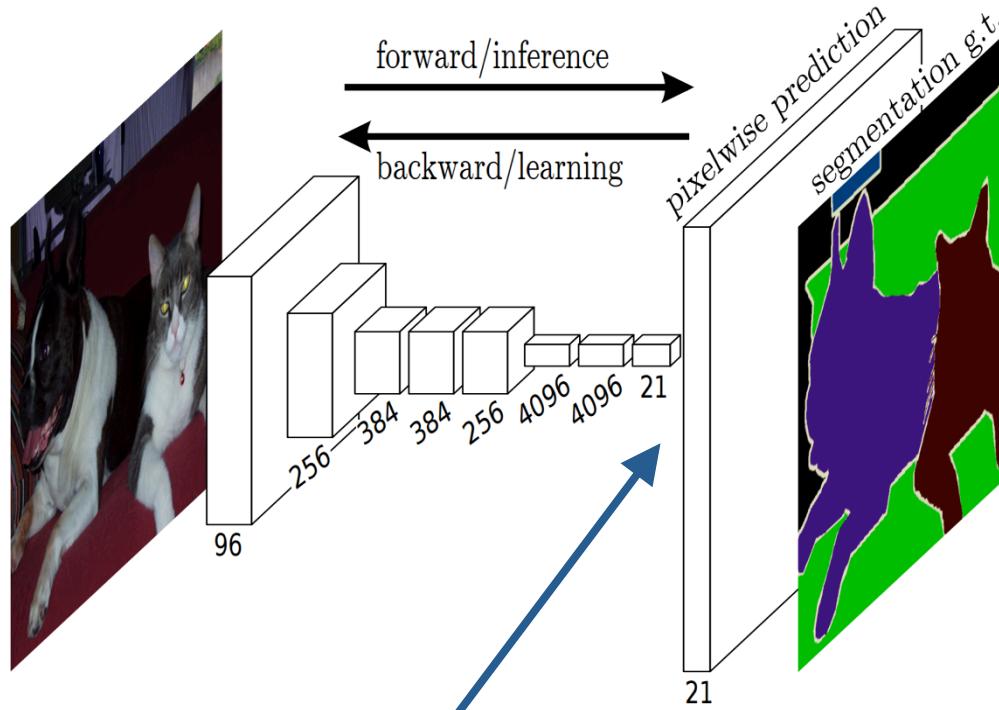
| object classes | building | grass | tree | cow | sheep | sky | airplane | water | face | car |
|----------------|----------|-------|------|------|-------|------|----------|-------|------|------|
| bicycle | flower | sign | bird | book | chair | road | cat | dog | body | boat |

Figure credit: Shotton et al, "TextonBoost for Image Understanding: Multi-Class Object Recognition and Segmentation by Jointly Modeling Texture, Layout, and Context", IJCV 2007

Semantic Segmentation: Upsampling



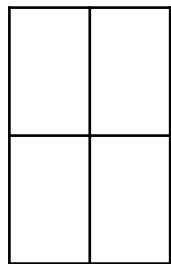
Semantic Segmentation: Upsampling



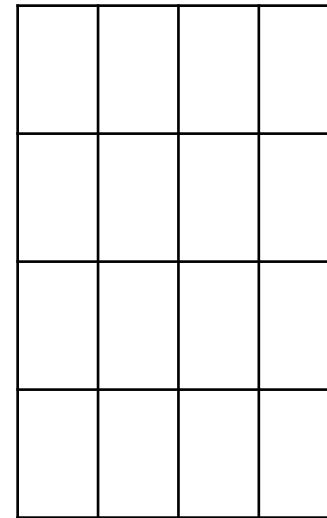
Learnable
upsampling!

Learnable Upsampling: “Deconvolution”

3 x 3 “deconvolution”, stride 2 pad 1



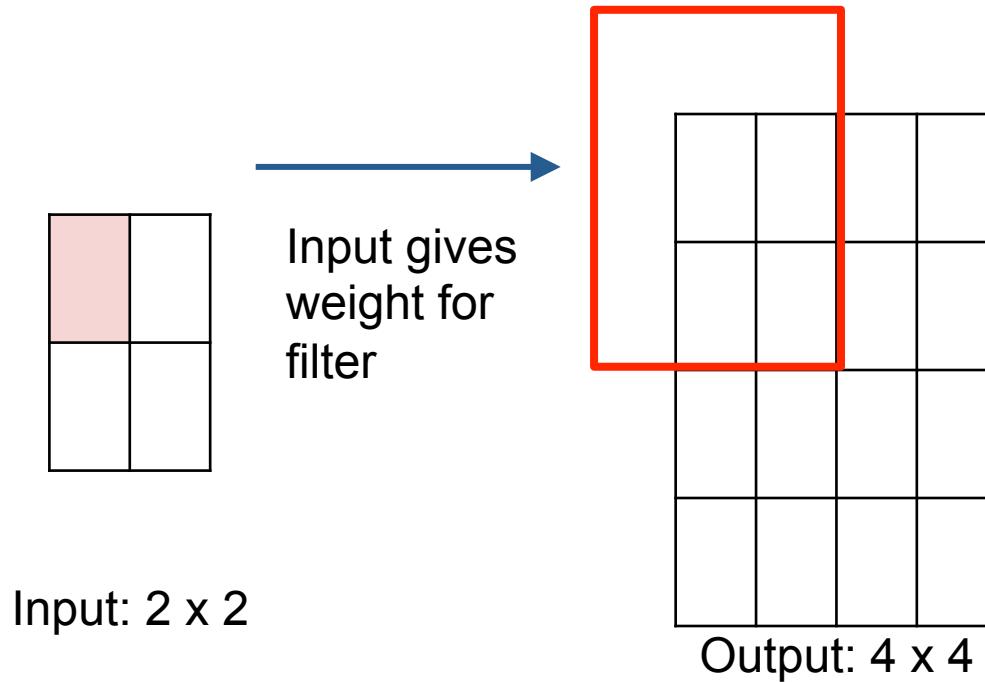
Input: 2 x 2



Output: 4 x 4

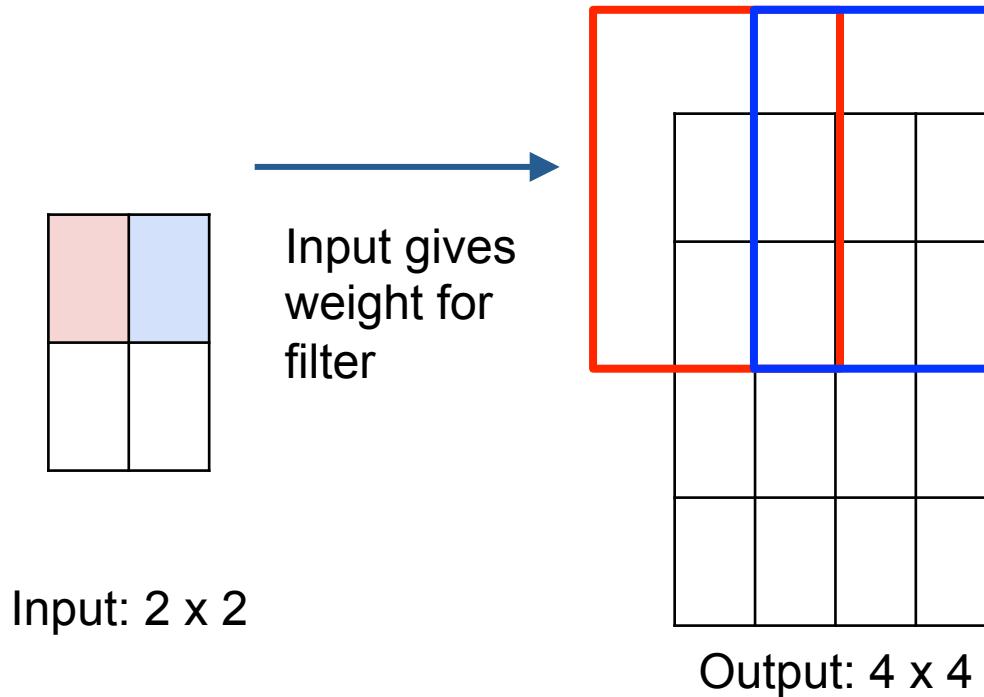
Learnable Upsampling: “Deconvolution”

3 x 3 “deconvolution”, stride 2 pad 1



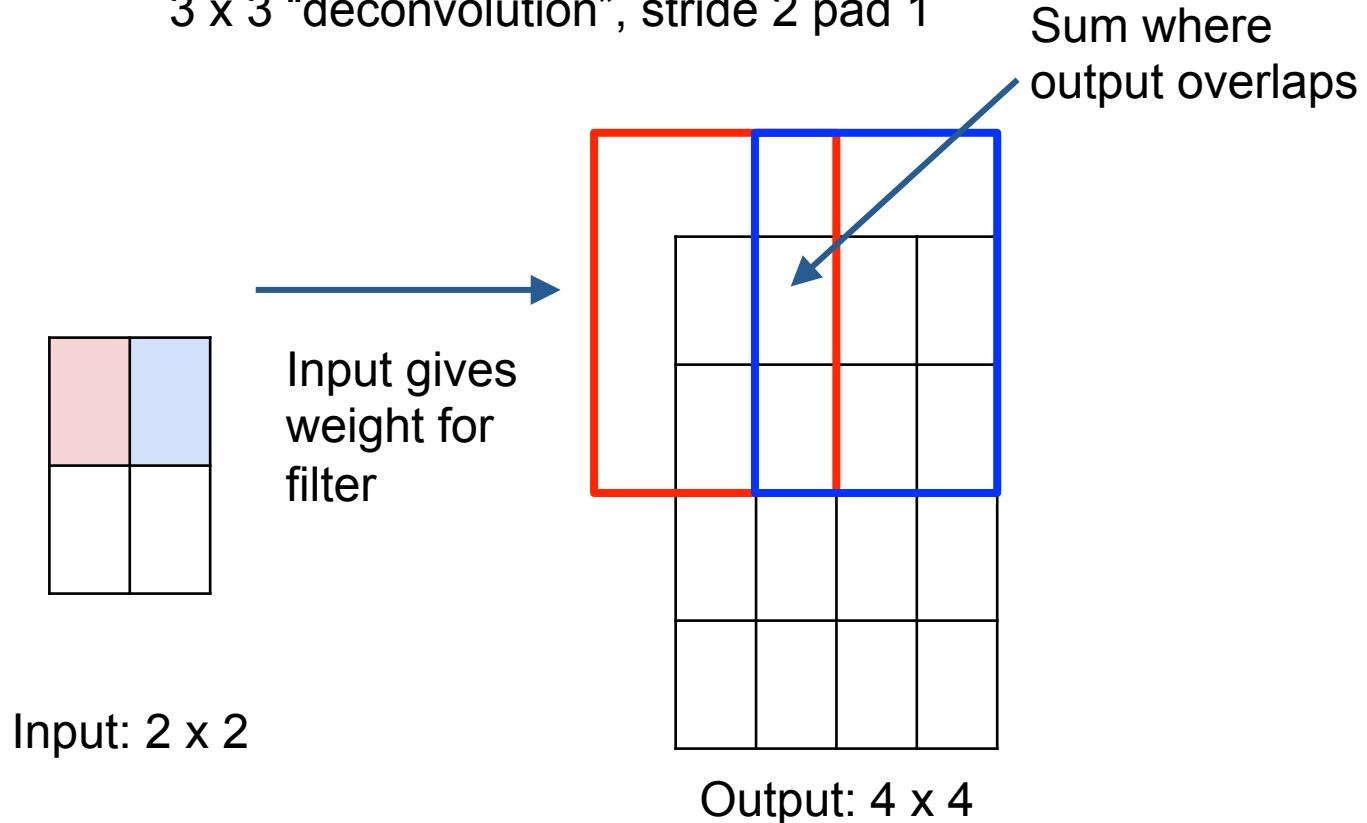
Learnable Upsampling: “Deconvolution”

3 x 3 “deconvolution”, stride 2 pad 1

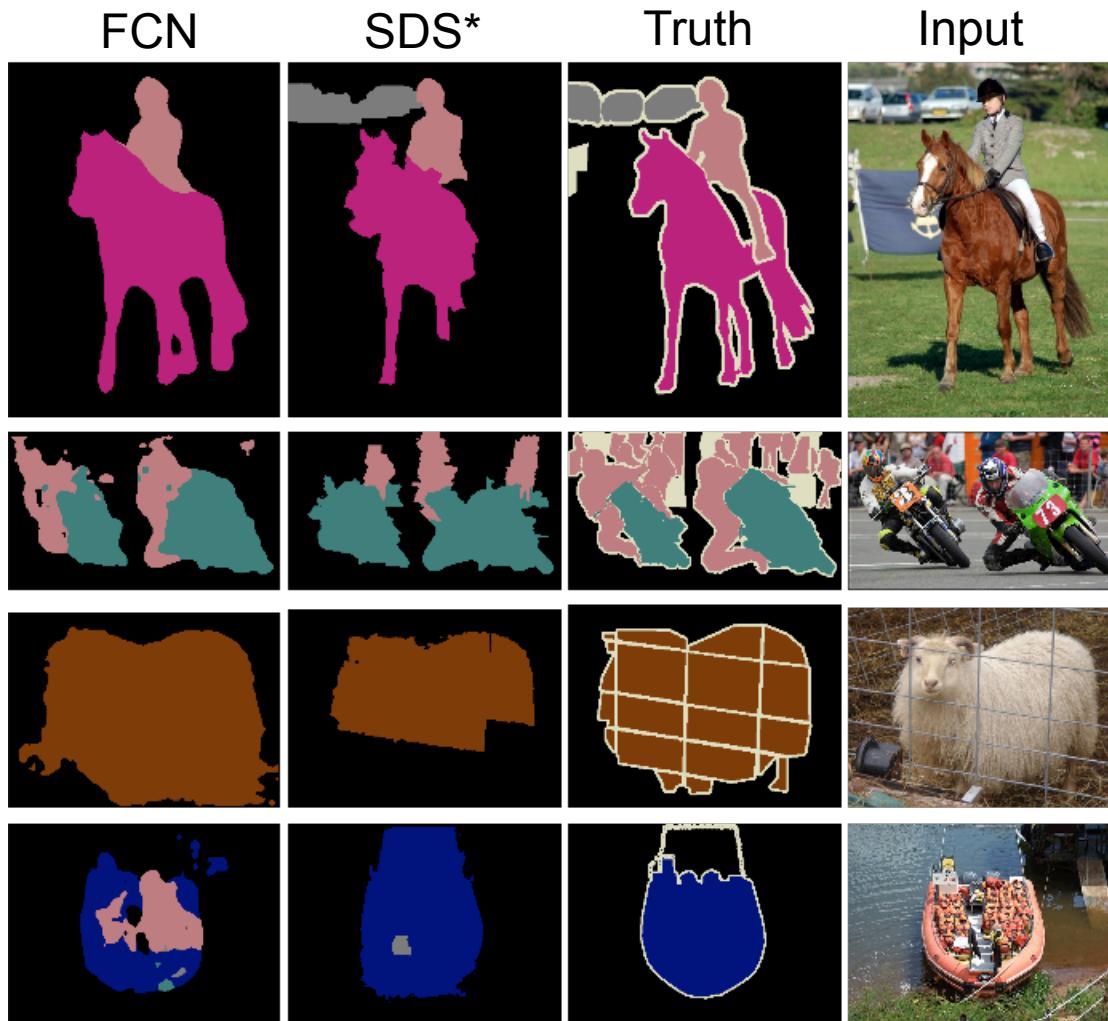


Learnable Upsampling: “Deconvolution”

3 x 3 “deconvolution”, stride 2 pad 1



Fully convolutional network



Relative to prior state-of-the-art SDS:

- 20% relative improvement for mean IoU
- 286× faster

*Simultaneous Detection and Segmentation
Hariharan et al. ECCV14

Part 4 Summary

- Semantic segmentation
 - Classify all pixels
 - Fully convolutional models, downsample then upsample
 - Learnable upsampling: fractionally strided convolution
 - Skip connections can help