

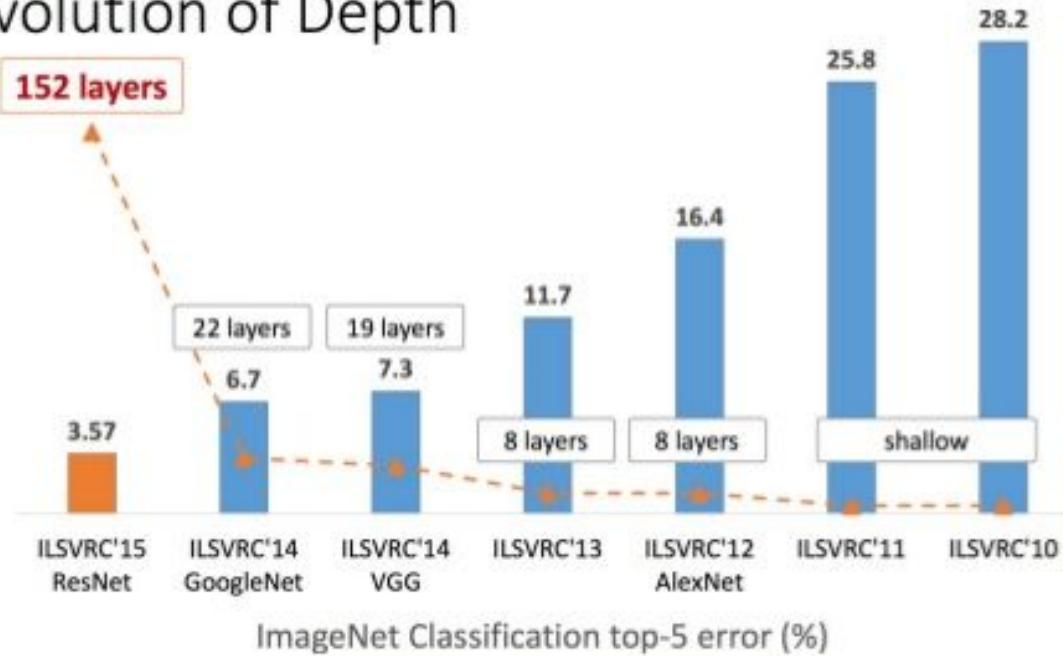
Artificial Intelligence and Machine learning course

Deep Learning lesson n.2

Prof. Barbara Caputo
DAUIN Politecnico di Torino

Quick Recap:

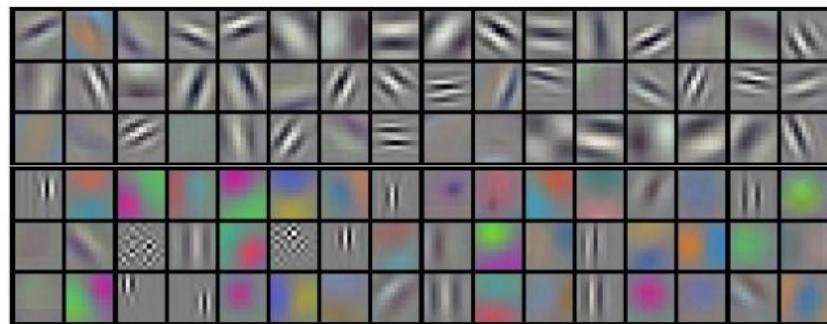
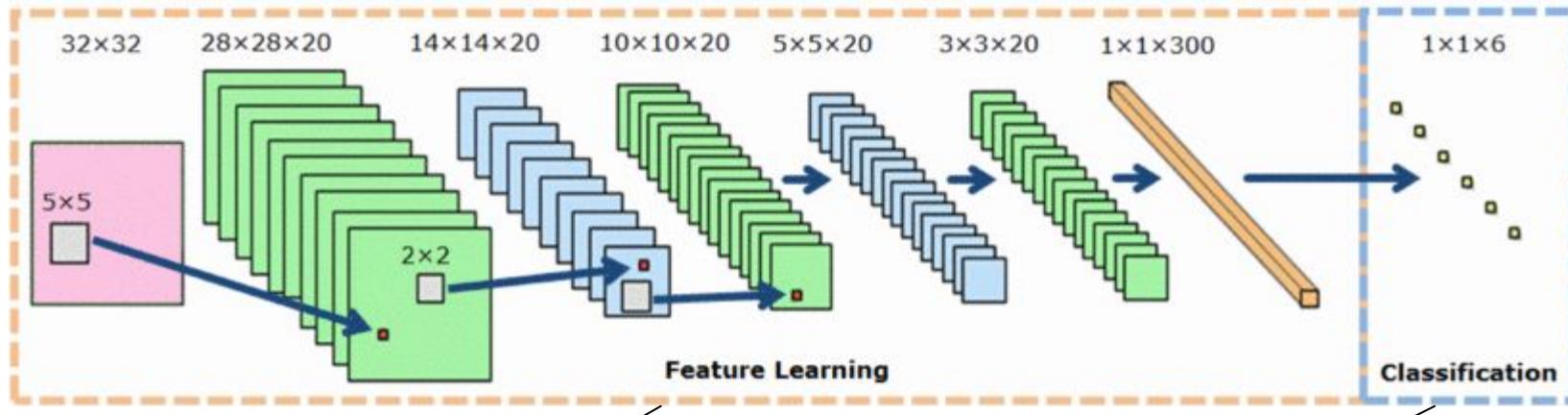
Revolution of Depth



He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. ["Deep Residual Learning for Image Recognition."](#) arXiv preprint arXiv:1512.03385 (2015). [\[slides\]](#)

80

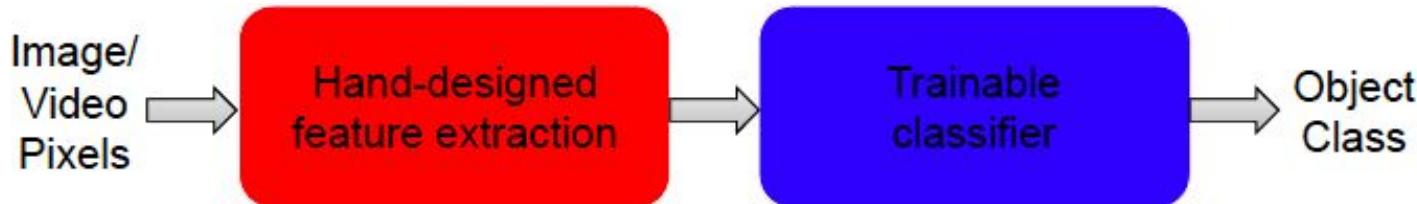
Deep Convolutional Neural Networks for Image Classification



Quick Recap:

“Shallow” vs. “deep” architectures

Traditional recognition: “Shallow” architecture

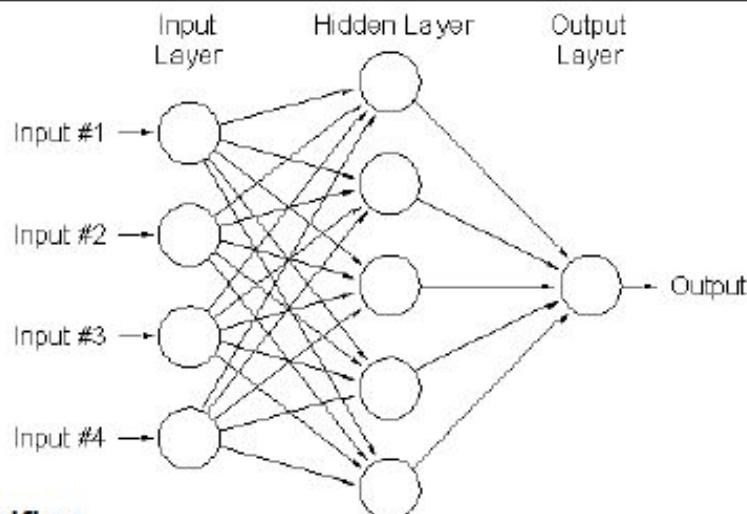


Deep learning: “Deep” architecture



Quick Recap:

Background: Multi-Layer Neural Networks



- Nonlinear classifier
- **Training:** find network weights w to minimize the error between true training labels y_i and estimated labels $f_w(\mathbf{x}_i)$:

$$E(\mathbf{w}) = \sum_{i=1}^N (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2$$

- Minimization can be done by gradient descent provided f is differentiable
 - This training method is called **back-propagation**

Quick Recap:

Forward propagation is the process of computing the output of the network given its input.

Backpropagation is the algorithm used for Neural Networks training. It works by updating the network weights via **Gradient Descent** and **derivatives chain rule** applied on a **loss**

Quick Recap:

Prob. of k class
given the input
(Softmax):

$$p(c_k=1|x) = \frac{e^{o_k}}{\sum_j e^{o_j}}$$

Quick Recap:

Prob. of k class
given the input
(Softmax):

$$p(c_k=1|x) = \frac{e^{o_k}}{\sum_j e^{o_j}}$$

Loss (negative
log-likelihood):

$$L(x, y; \theta) = -\sum_j y_j \log p(c_j|x)$$

$$y = [0^1 0^0 \dots 0^k 1^0 \dots 0^c]$$

Quick Recap:

Prob. of k class
given the input
(Softmax):

$$p(c_k=1|\mathbf{x}) = \frac{e^{o_k}}{\sum_j e^{o_j}}$$

Loss (negative
log-likelihood):

$$L(\mathbf{x}, y; \theta) = -\sum_j y_j \log p(c_j|\mathbf{x})$$

$$\mathbf{y} = [0^1 0^0 \dots 0^k 1^0 \dots 0^c]$$

Derivative of Loss
w.r.t output:

$$\frac{\partial L}{\partial o} = p(c|\mathbf{x}) - \mathbf{y}$$

Quick Recap:

Prob. of k class
given the input
(Softmax):

$$p(c_k=1|\mathbf{x}) = \frac{e^{o_k}}{\sum_j e^{o_j}}$$

Loss (negative
log-likelihood):

$$L(\mathbf{x}, y; \theta) = -\sum_j y_j \log p(c_j|\mathbf{x})$$

$$\mathbf{y} = [0^1 0^0 \dots 0^k 1^0 \dots 0^c]$$

Derivative of Loss
w.r.t output:

$$\frac{\partial L}{\partial o} = p(c|\mathbf{x}) - \mathbf{y}$$

Chain rule:

$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial W^3}$$

Quick Recap:

Prob. of k class
given the input
(Softmax):

$$p(c_k=1|\mathbf{x}) = \frac{e^{o_k}}{\sum_j e^{o_j}}$$

Loss (negative
log-likelihood):

$$L(\mathbf{x}, y; \boldsymbol{\theta}) = -\sum_j y_j \log p(c_j|\mathbf{x})$$

$$\mathbf{y} = [0^1 0^0 \dots 0^k 1^0 \dots 0^c]$$

Derivative of Loss
w.r.t output:

$$\frac{\partial L}{\partial o} = p(c|\mathbf{x}) - \mathbf{y}$$

Chain rule:

$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial W^3}$$

Update rule:

$$w_i \leftarrow w_i - \eta \frac{\partial L}{\partial w_i}$$

Quick Recap:

Curse of dimensionality problem...

Quick Recap:

Curse of dimensionality problem..

Overfitting... :-(

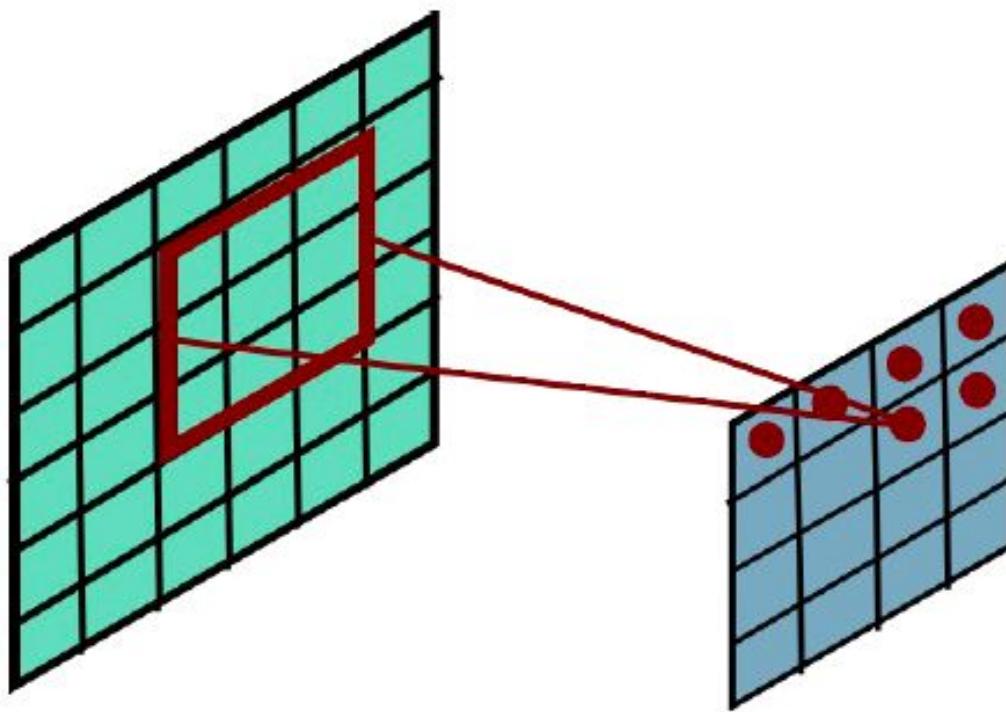
Quick Recap:

Curse of dimensionality problem..

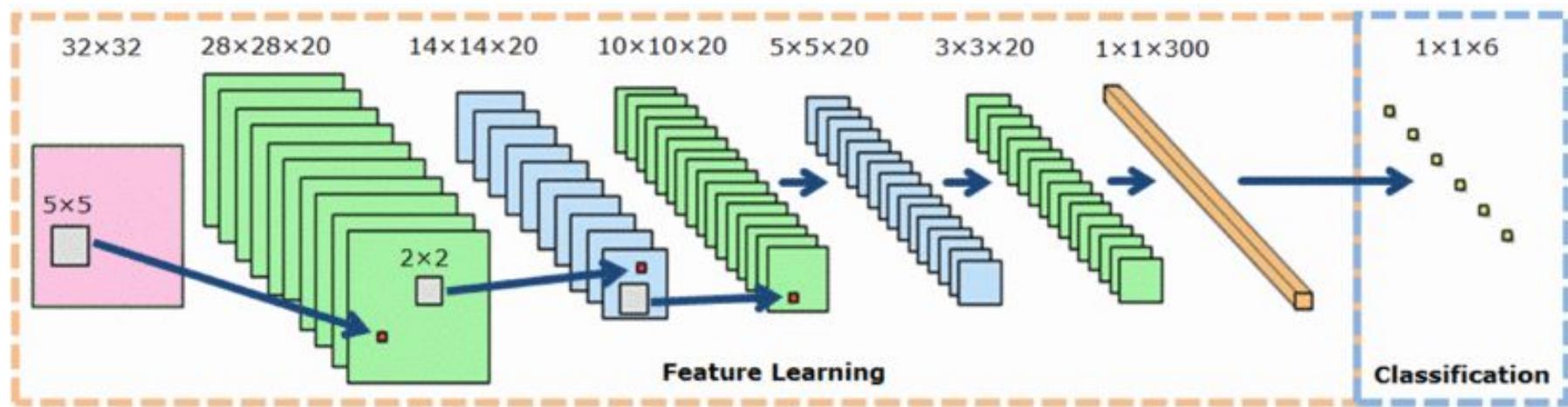
Overfitting...:(

Solution? :-)

Convolutional Layer

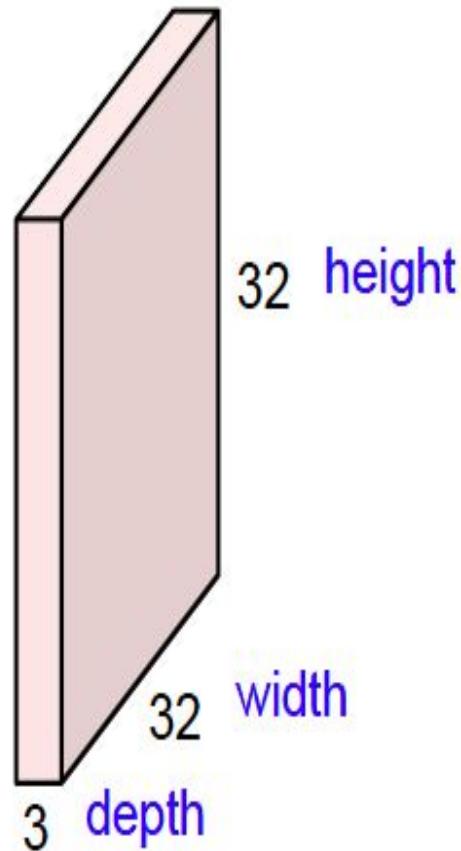


Convolutional Neural Networks



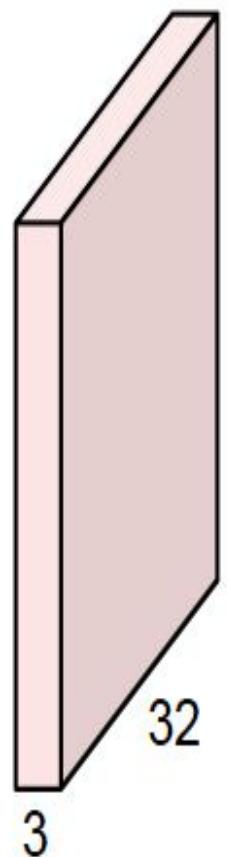
Convolution Layer

32x32x3 image

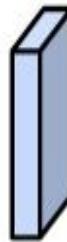


Convolution Layer

32x32x3 image



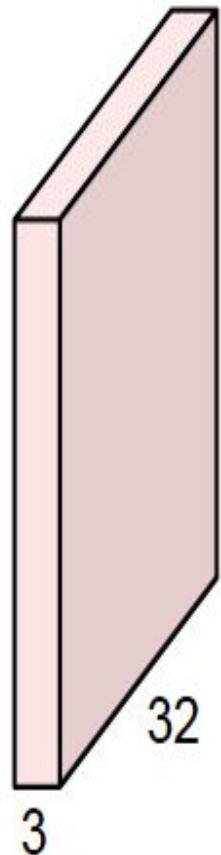
5x5x3 filter



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

Convolution Layer

32x32x3 image



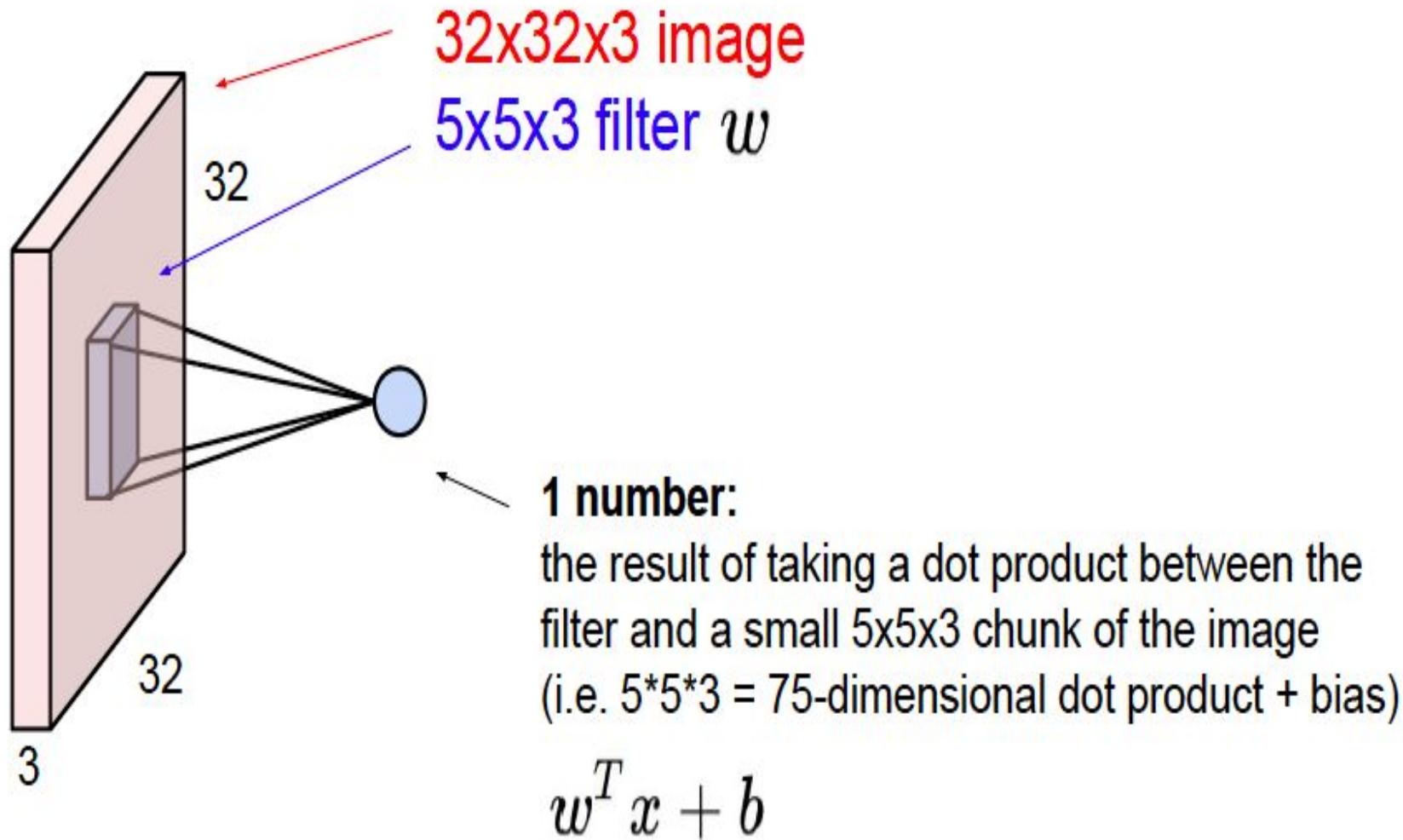
5x5x3 filter



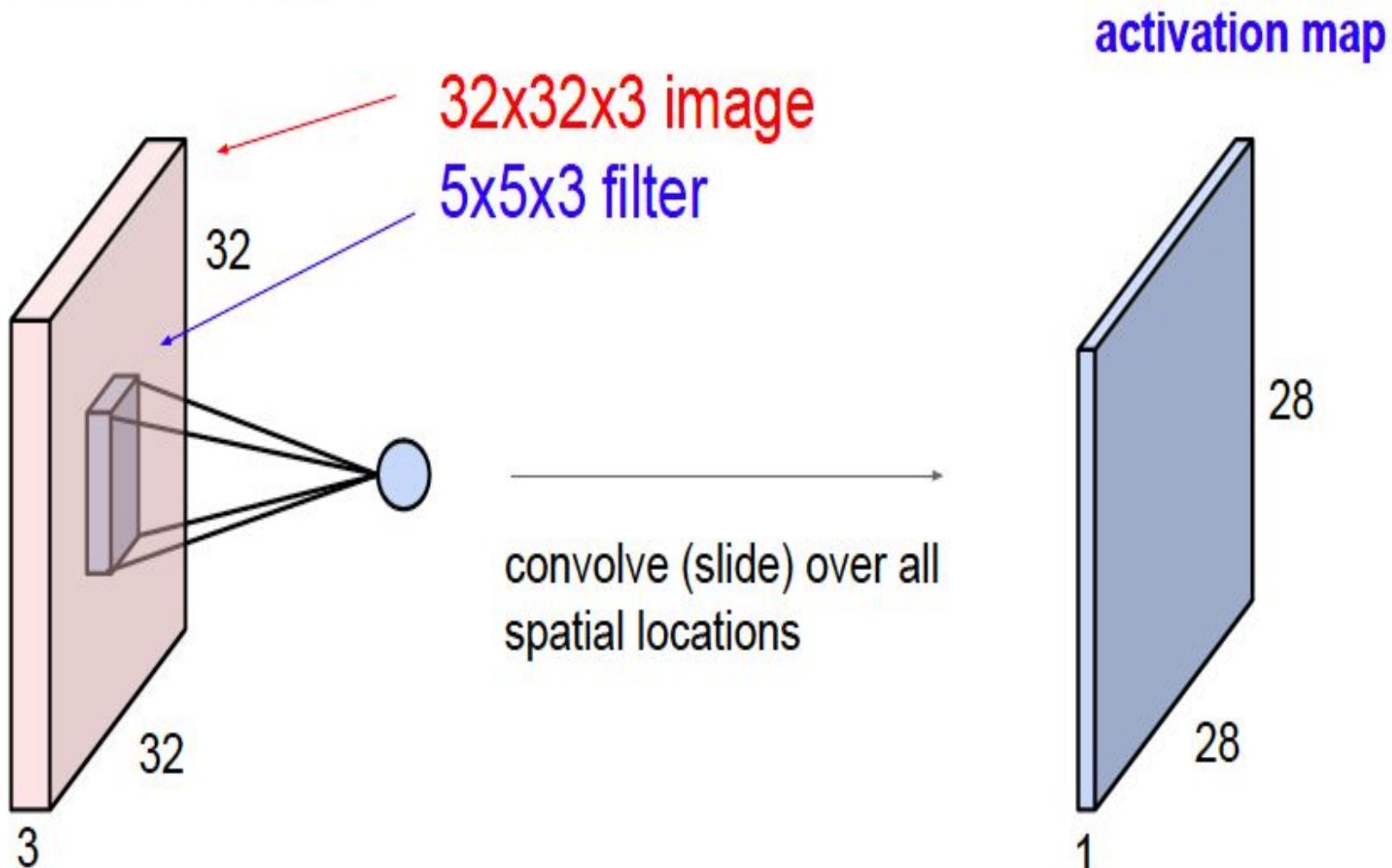
Filters always extend the full depth of the input volume

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

Convolution Layer

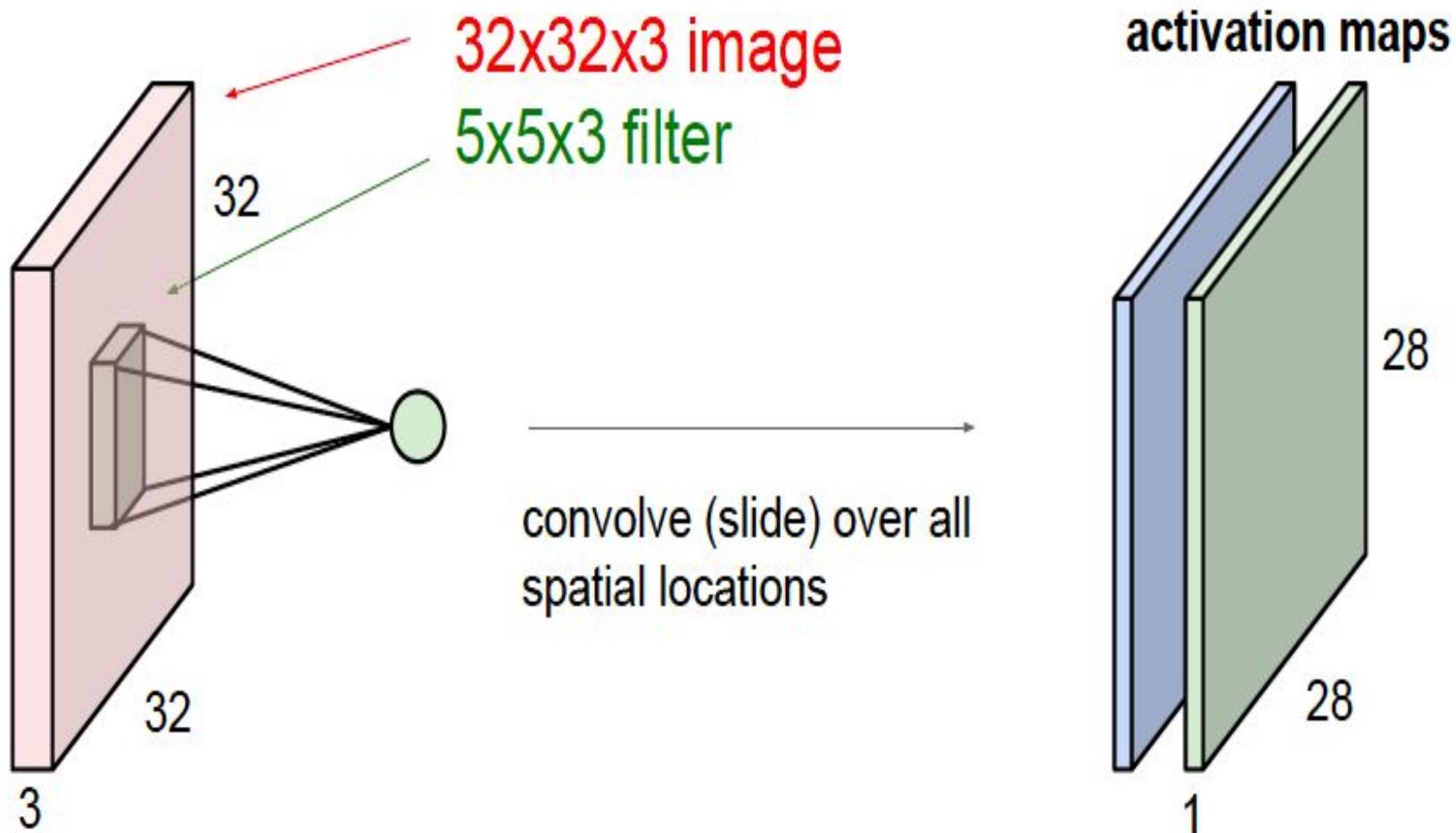


Convolution Layer

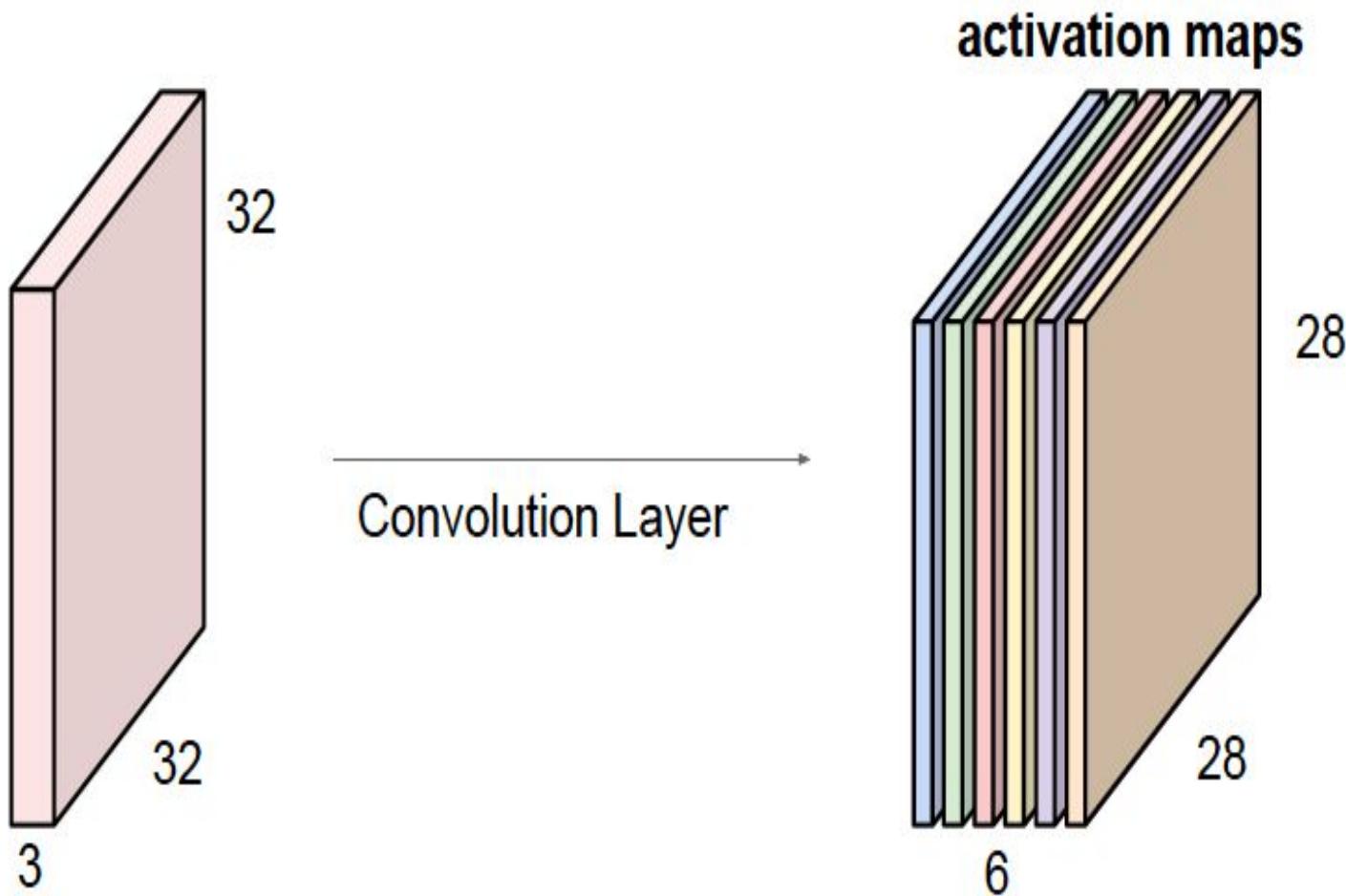


Convolution Layer

consider a second, green filter

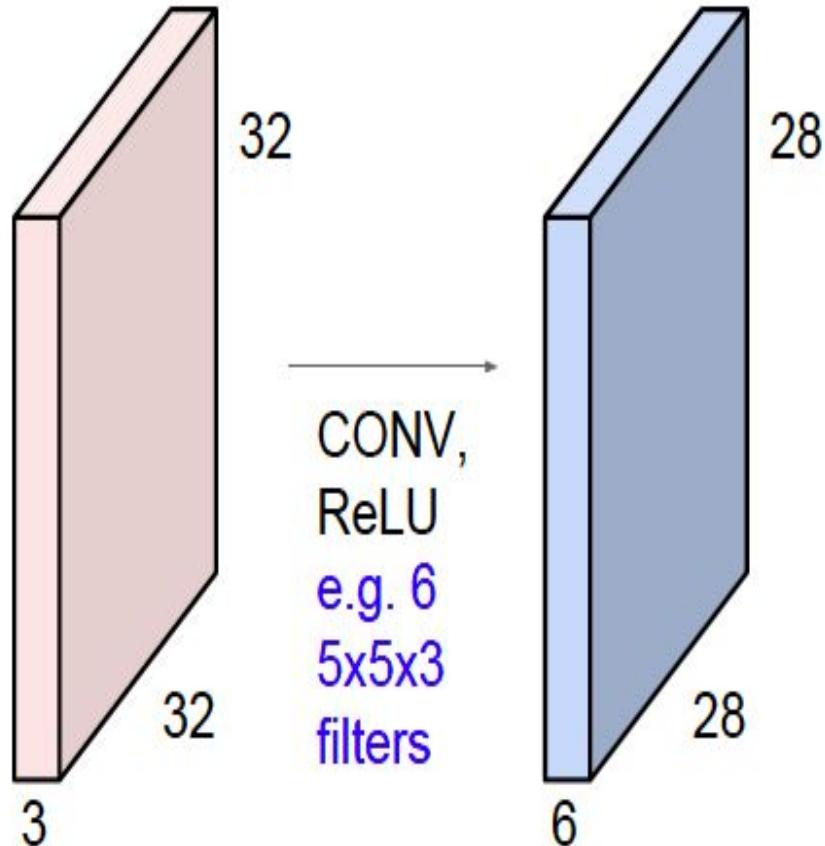


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

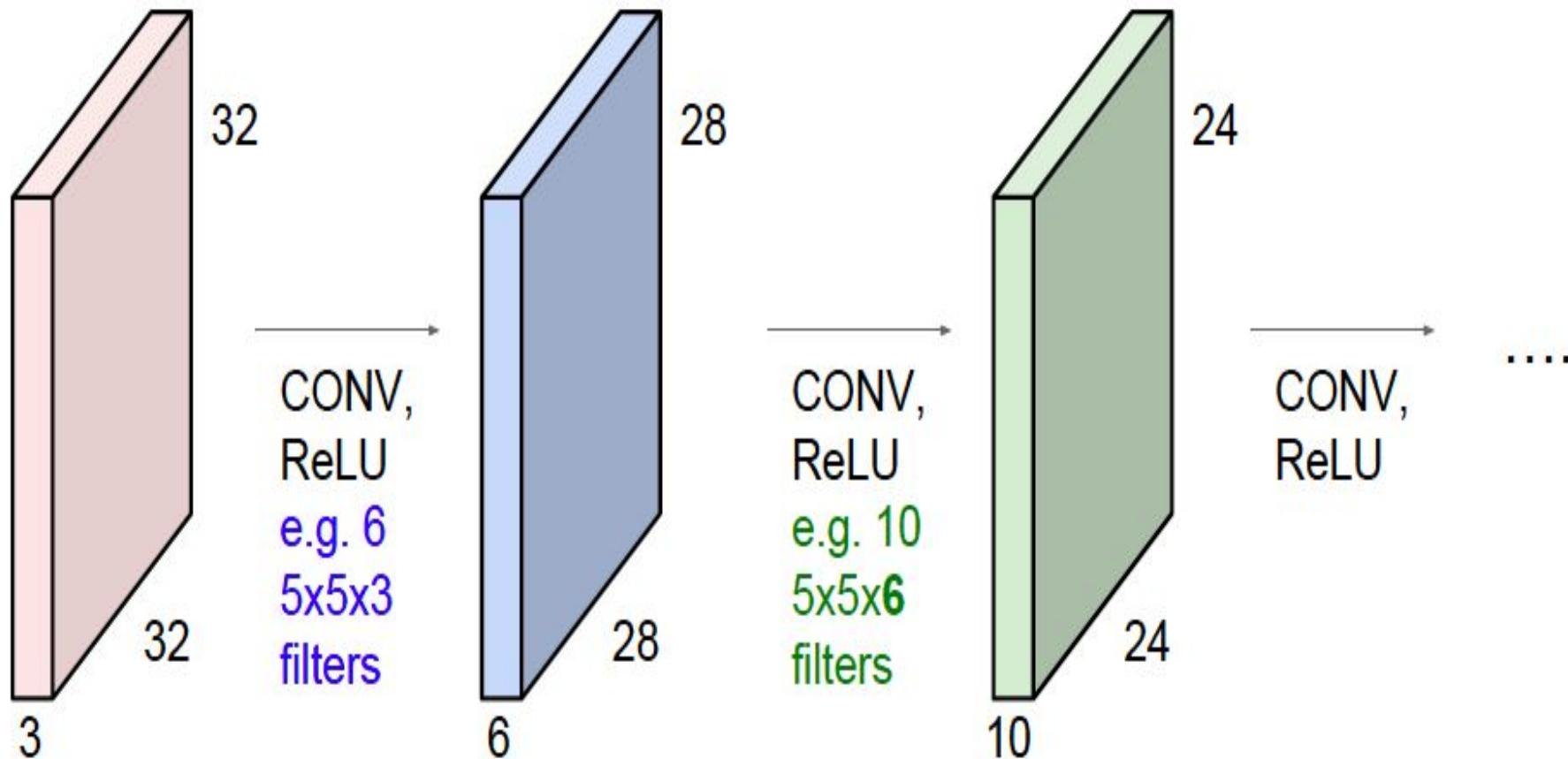


We stack these up to get a “new image” of size 28x28x6!

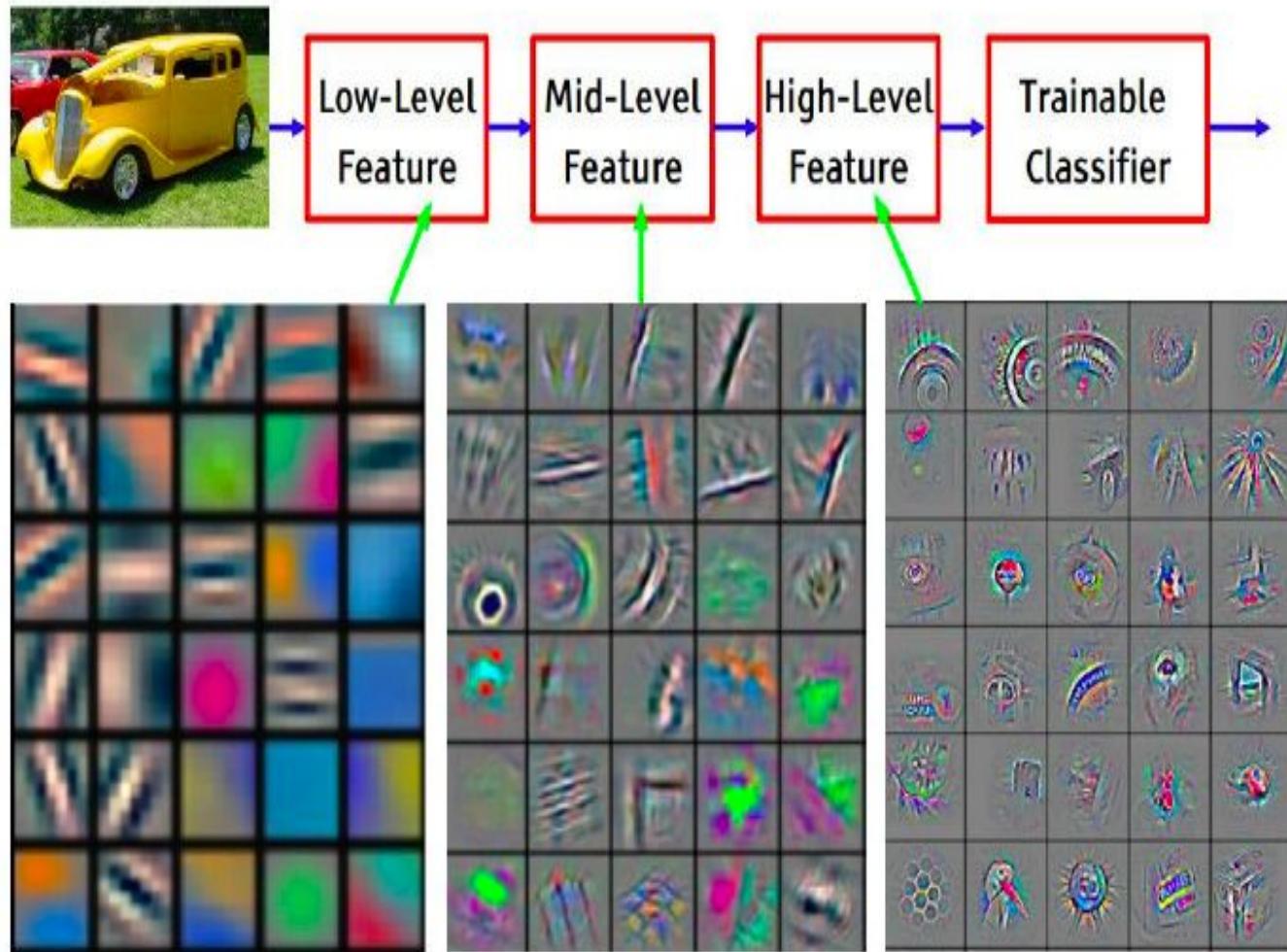
Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions



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



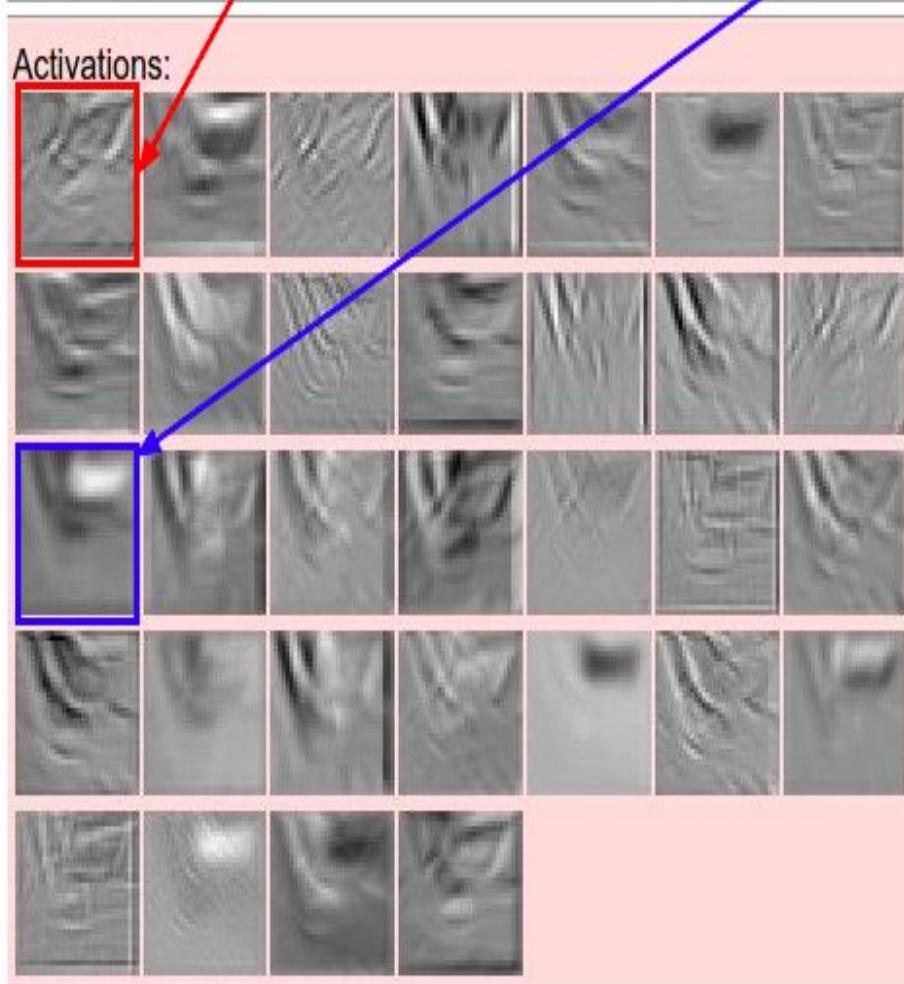
Preview



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



one filter =>
one activation map



example 5×5 filters
(32 total)

We call the layer convolutional
because it is related to convolution
of two signals:

$$f[x,y] * g[x,y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1, n_2] \cdot g[x - n_1, y - n_2]$$

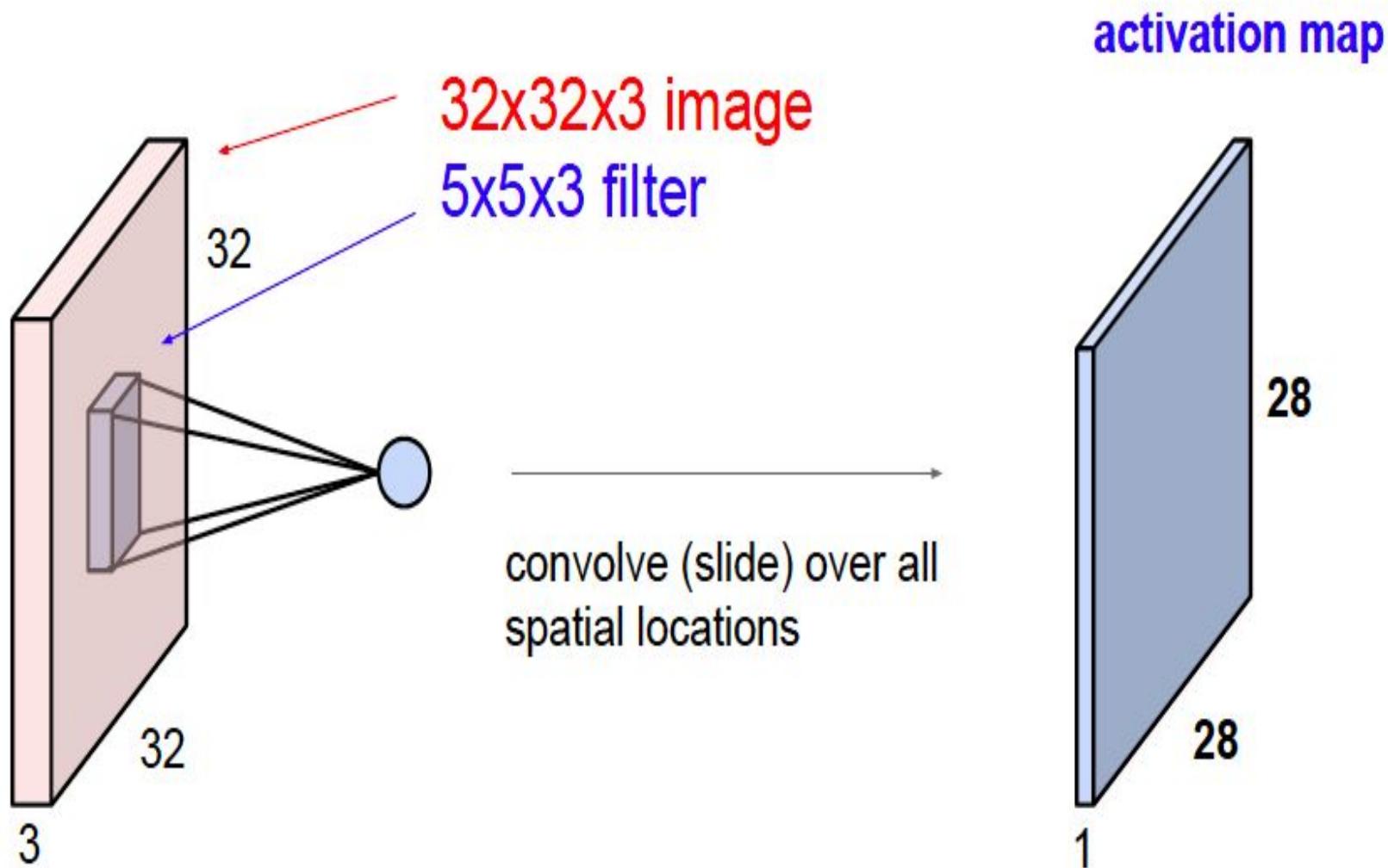


elementwise multiplication and sum of
a filter and the signal (image)

preview:

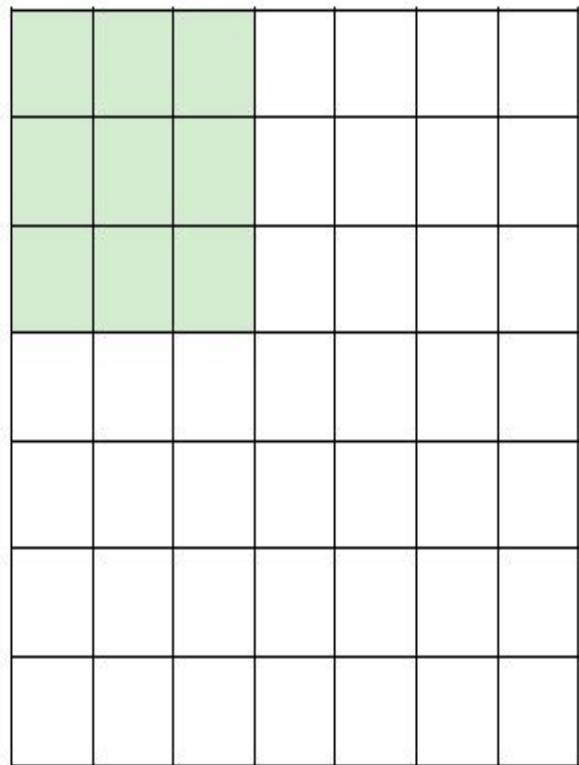


A closer look at spatial dimensions:



A closer look at spatial dimensions:

7

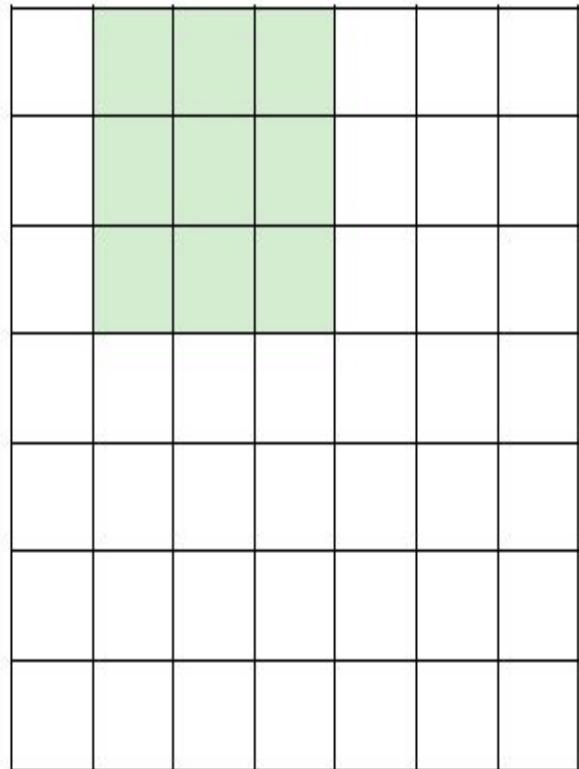


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

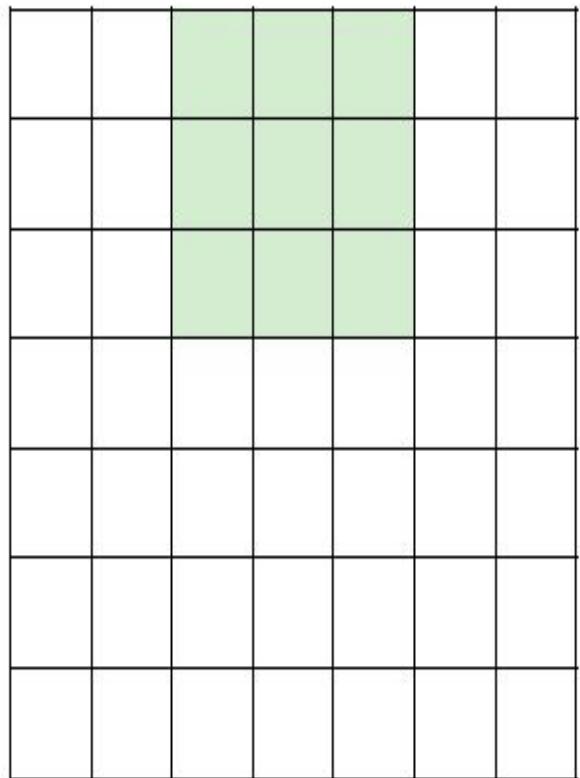


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

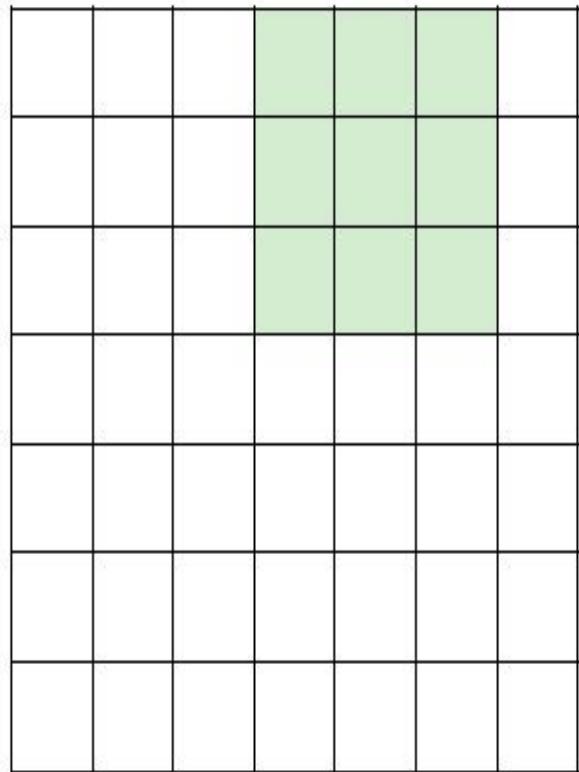


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

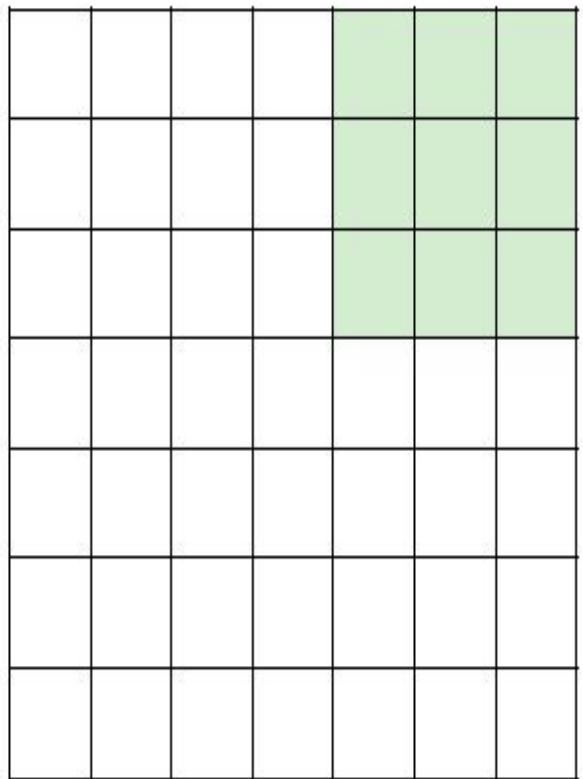


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

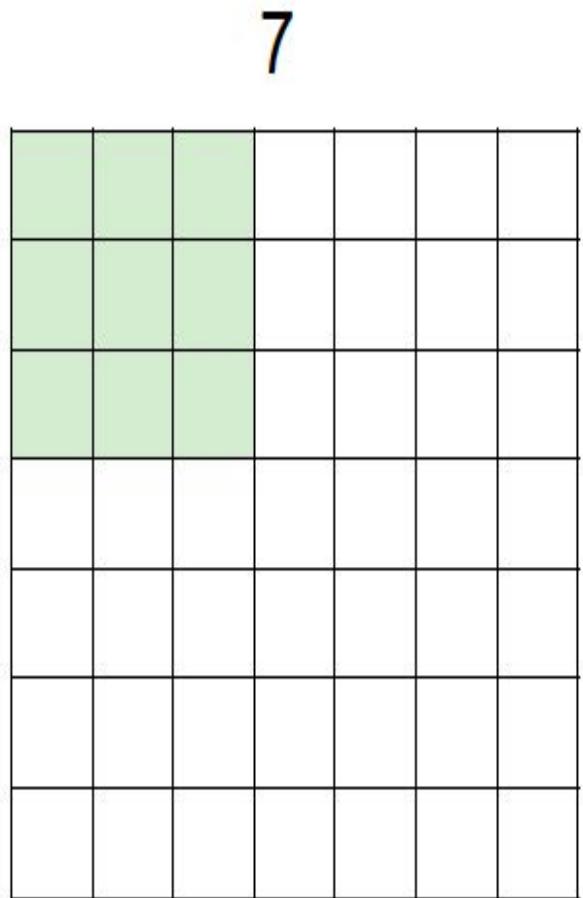
7



7x7 input (spatially)
assume 3x3 filter

=> 5x5 output

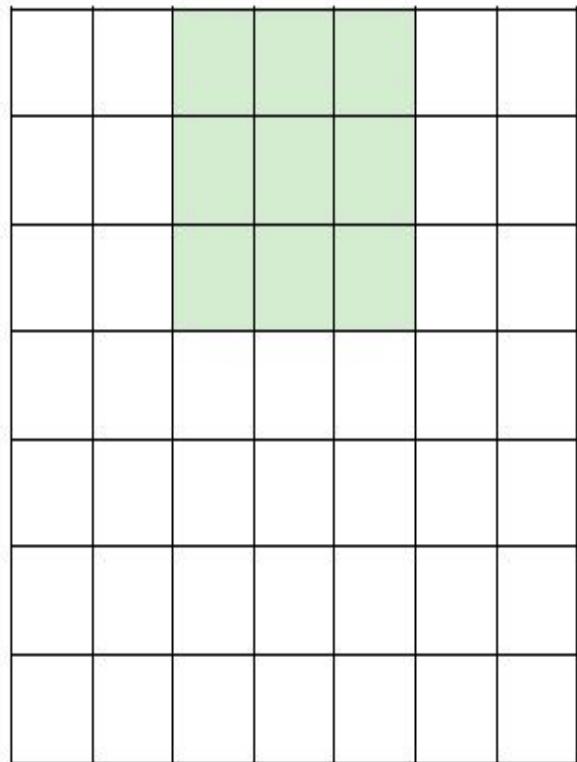
A closer look at spatial dimensions:



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

A closer look at spatial dimensions:

7

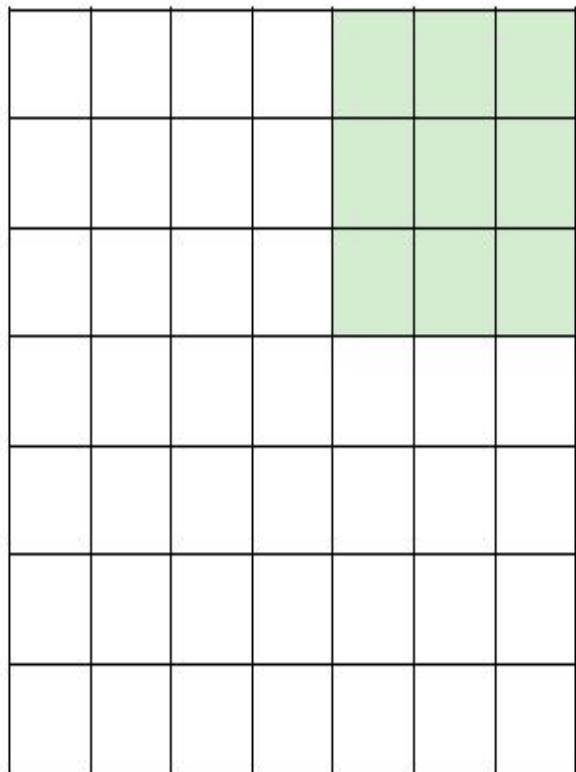


7

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

A closer look at spatial dimensions:

7

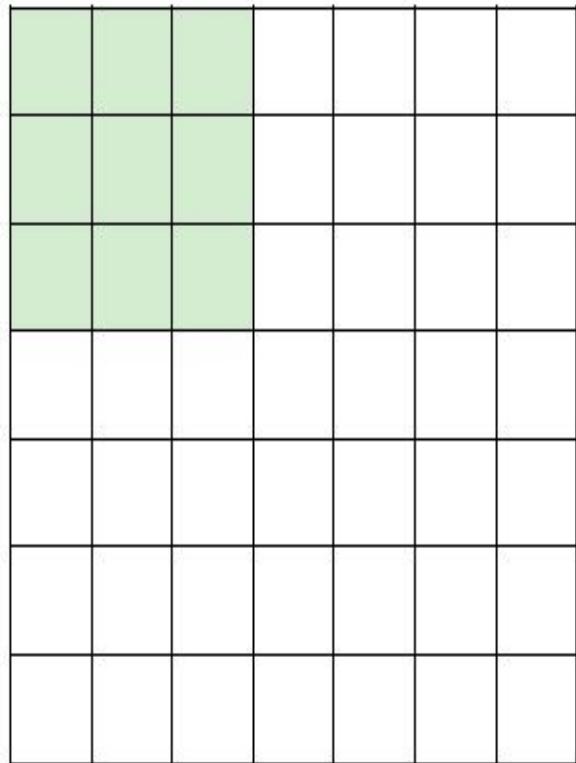


7

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

A closer look at spatial dimensions:

7

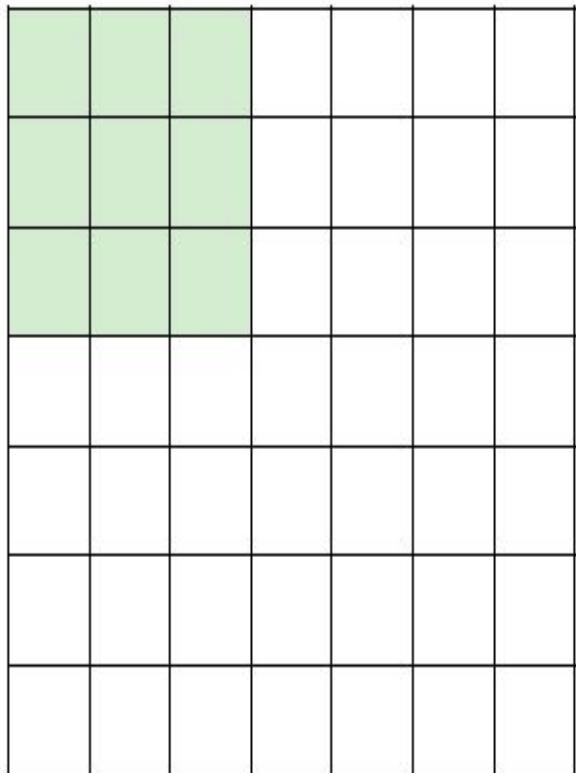


7

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

A closer look at spatial dimensions:

7

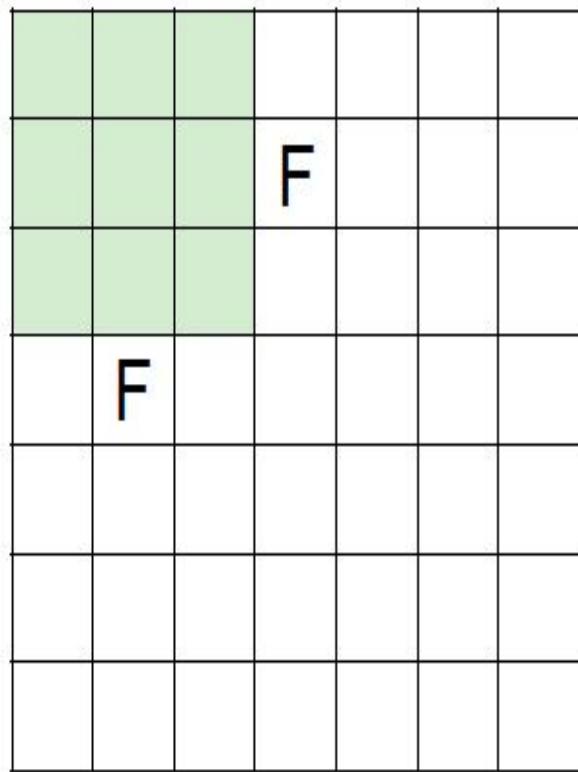


7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

7

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

N



N

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

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

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

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

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

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?

(recall:)

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

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 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)

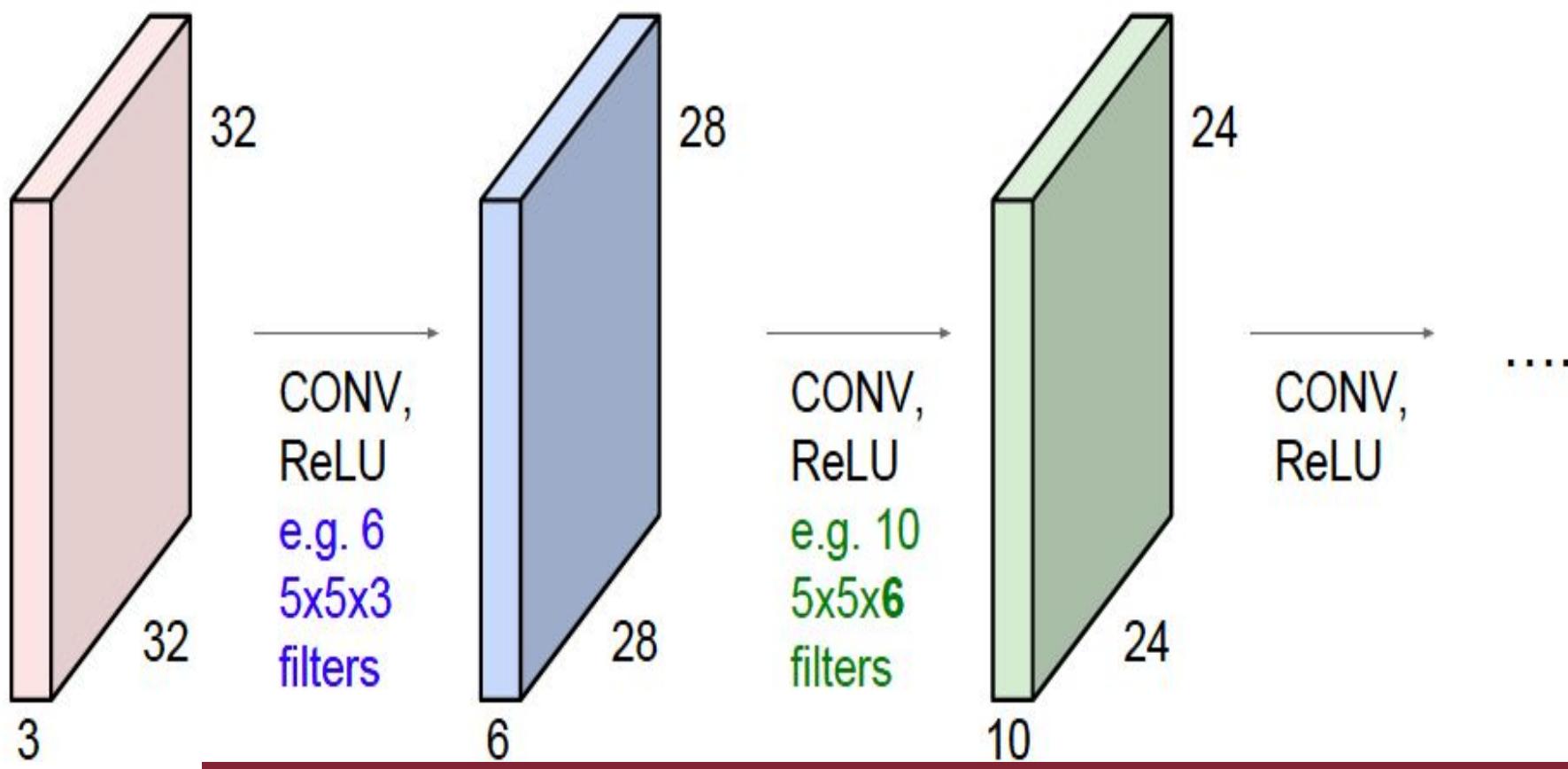
e.g. $F = 3 \Rightarrow$ zero pad with 1

$F = 5 \Rightarrow$ zero pad with 2

$F = 7 \Rightarrow$ zero pad with 3

Remember back to...

E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially!
(32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.

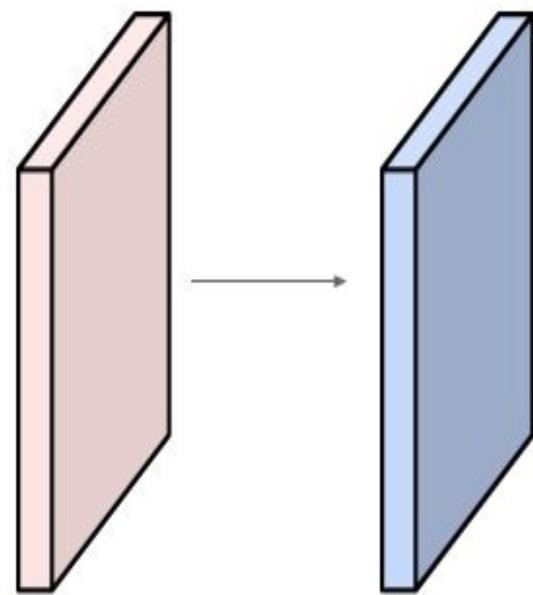


Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

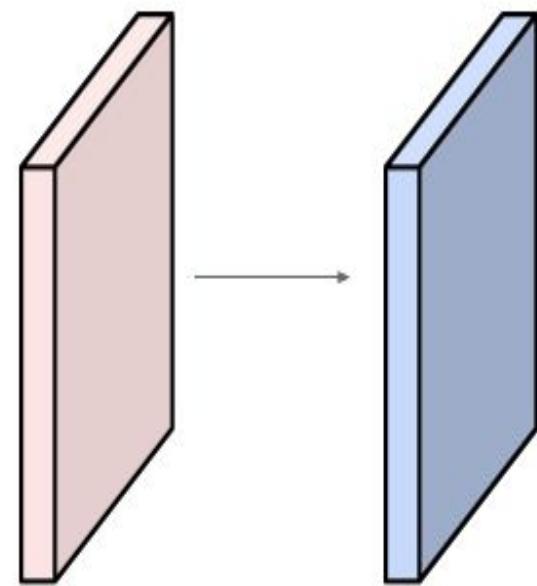
Output volume size: ?



Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2



Output volume size:

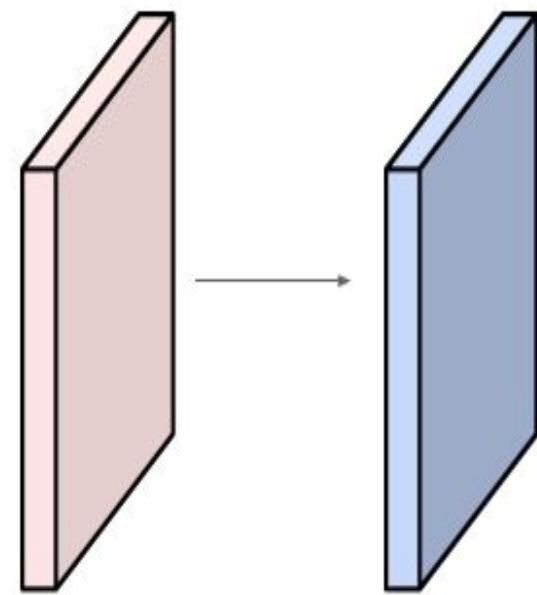
$(32+2*2-5)/1+1 = 32$ spatially, so

32x32x10

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

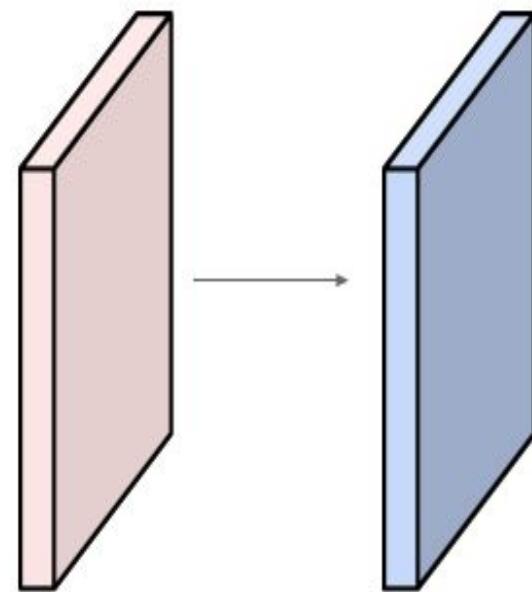


Number of parameters in this layer?

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

each filter has $5*5*3 + 1 = 76$ params (+1 for bias)
 $\Rightarrow 76*10 = 760$

Summary. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

Common settings:

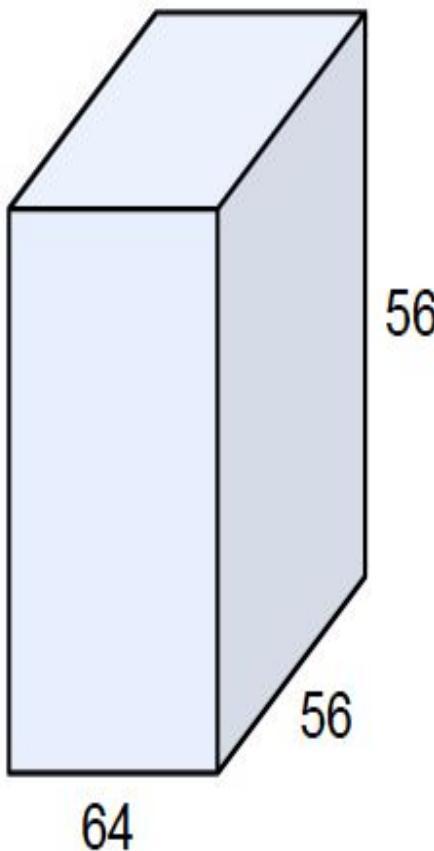
Summary. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

$K = (\text{powers of } 2, \text{ e.g. } 32, 64, 128, 512)$

- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ?$ (whatever fits)
- $F = 1, S = 1, P = 0$

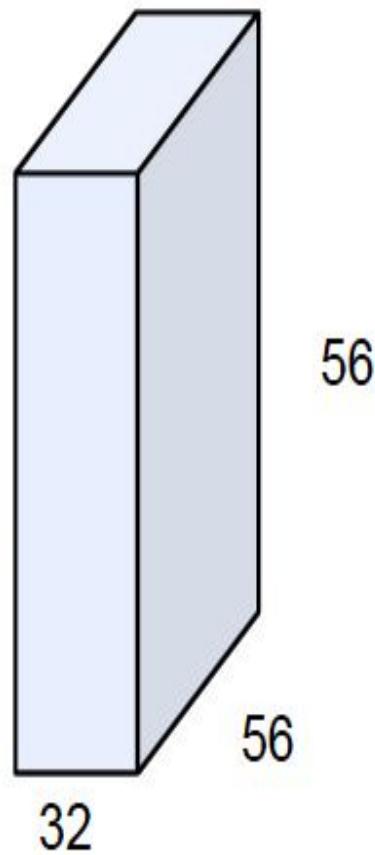
(btw, 1x1 convolution layers make perfect sense)



1x1 CONV
with 32 filters

→

(each filter has size
 $1 \times 1 \times 64$, and performs a
64-dimensional dot
product)



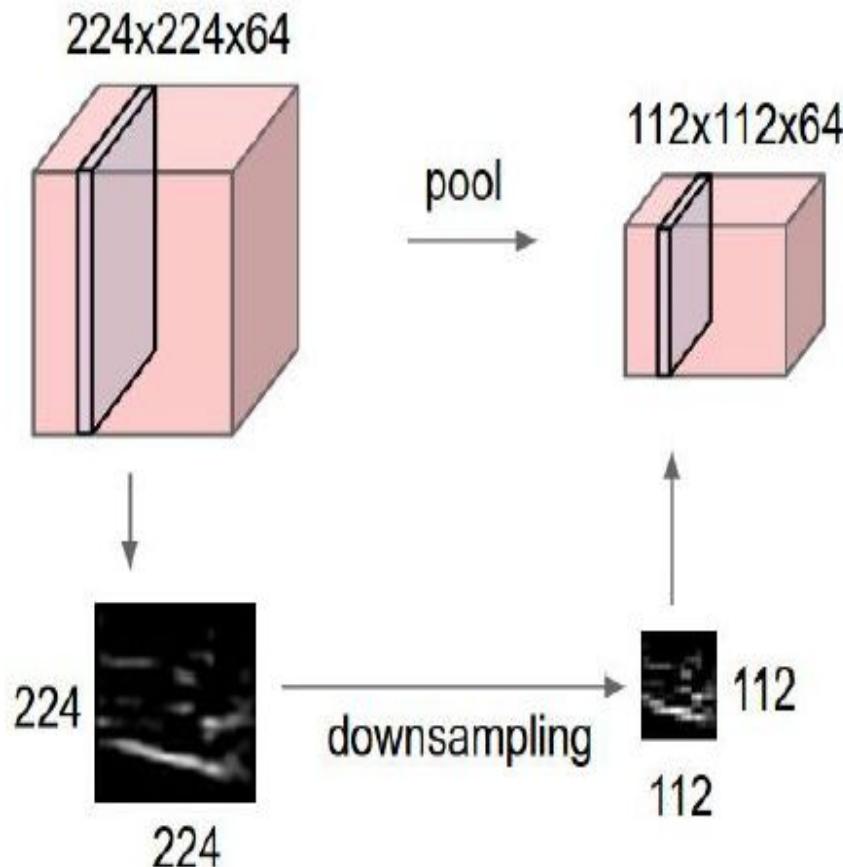
two more layers to go: POOL/FC



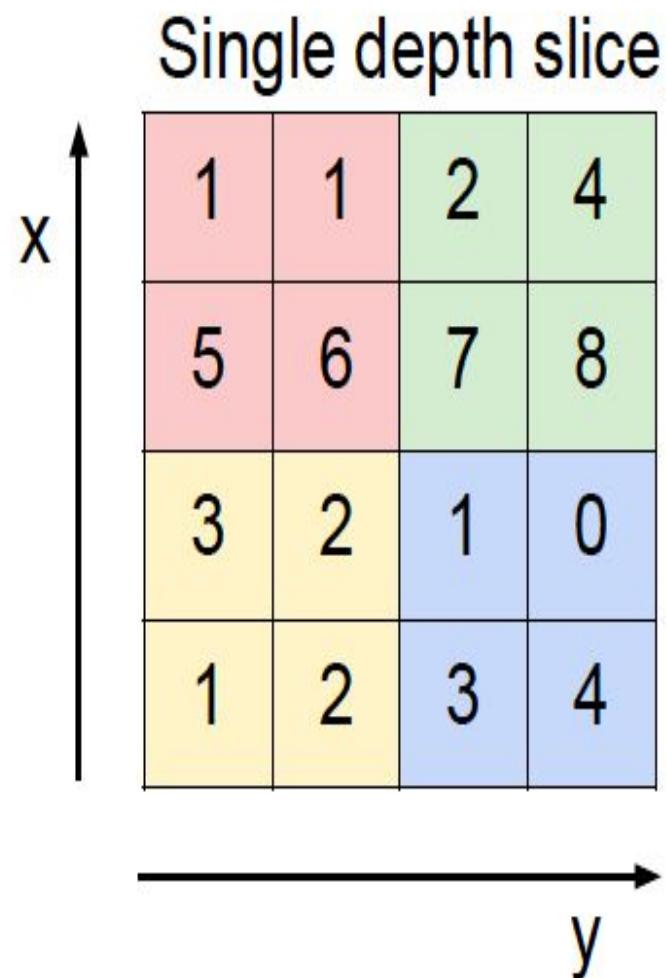
car
truck
airplane
ship
horse

Pooling layer

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



MAX POOLING



max pool with 2x2 filters
and stride 2

6	8
3	4

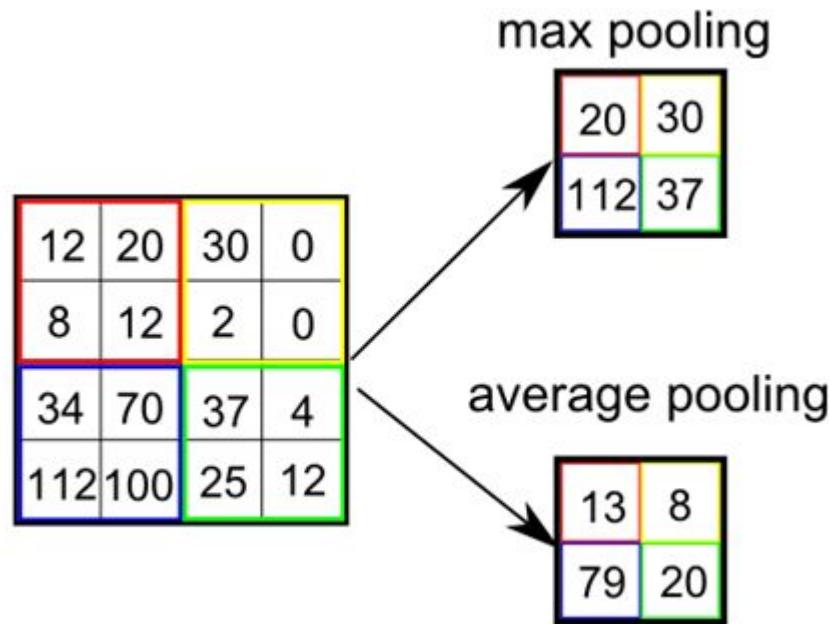
- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
 - their spatial extent F ,
 - the stride S ,
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F)/S + 1$
 - $H_2 = (H_1 - F)/S + 1$
 - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

Common settings:

- Accepts a volume of size $W_1 \times H_1 \times D_1$ $F = 2, S = 2$
- Requires three hyperparameters:
 - their spatial extent F ,
 - the stride S , $F = 3, S = 2$
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F)/S + 1$
 - $H_2 = (H_1 - F)/S + 1$
 - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

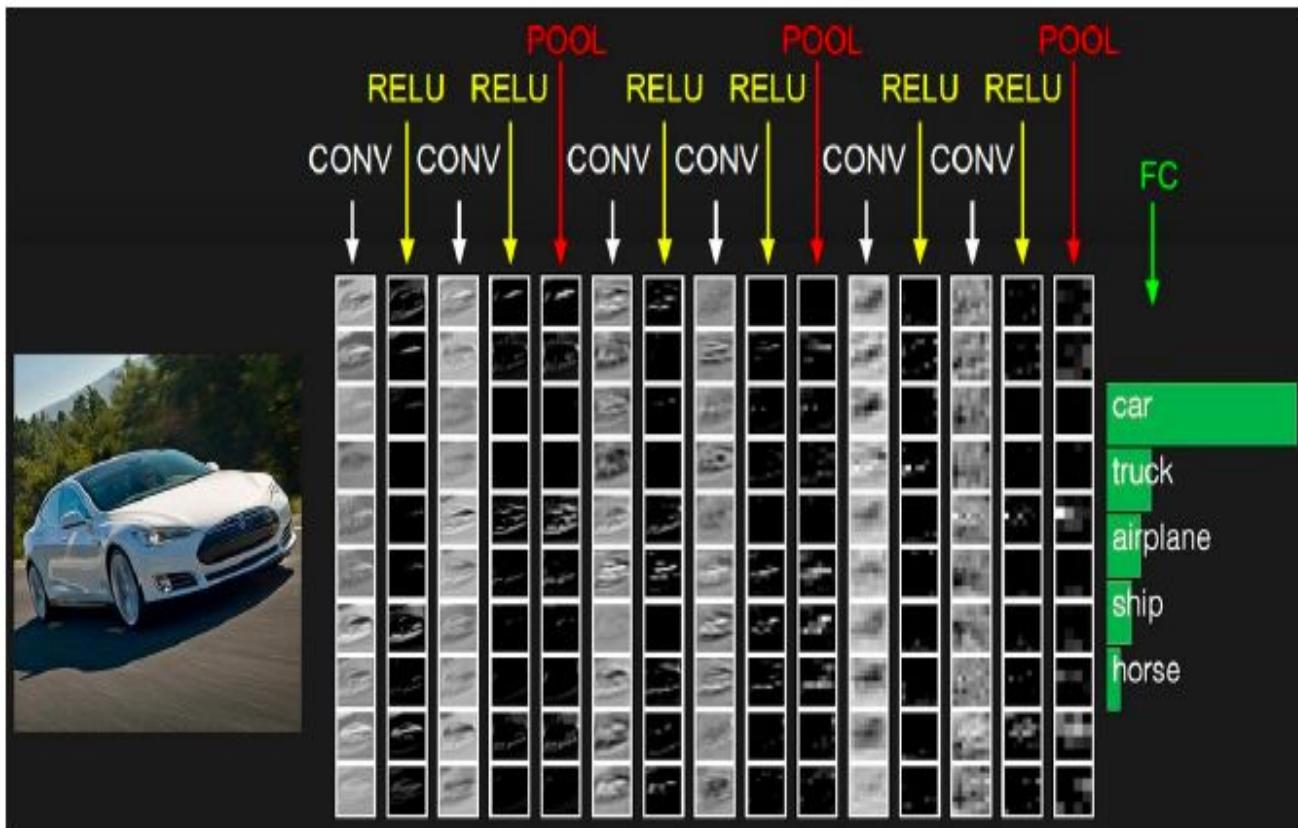
Average Pooling

- Sometimes is convenient to *average* pool instead of max
- The output of avg pooling depends on each input
- Used on recent architectures in place of FC layers

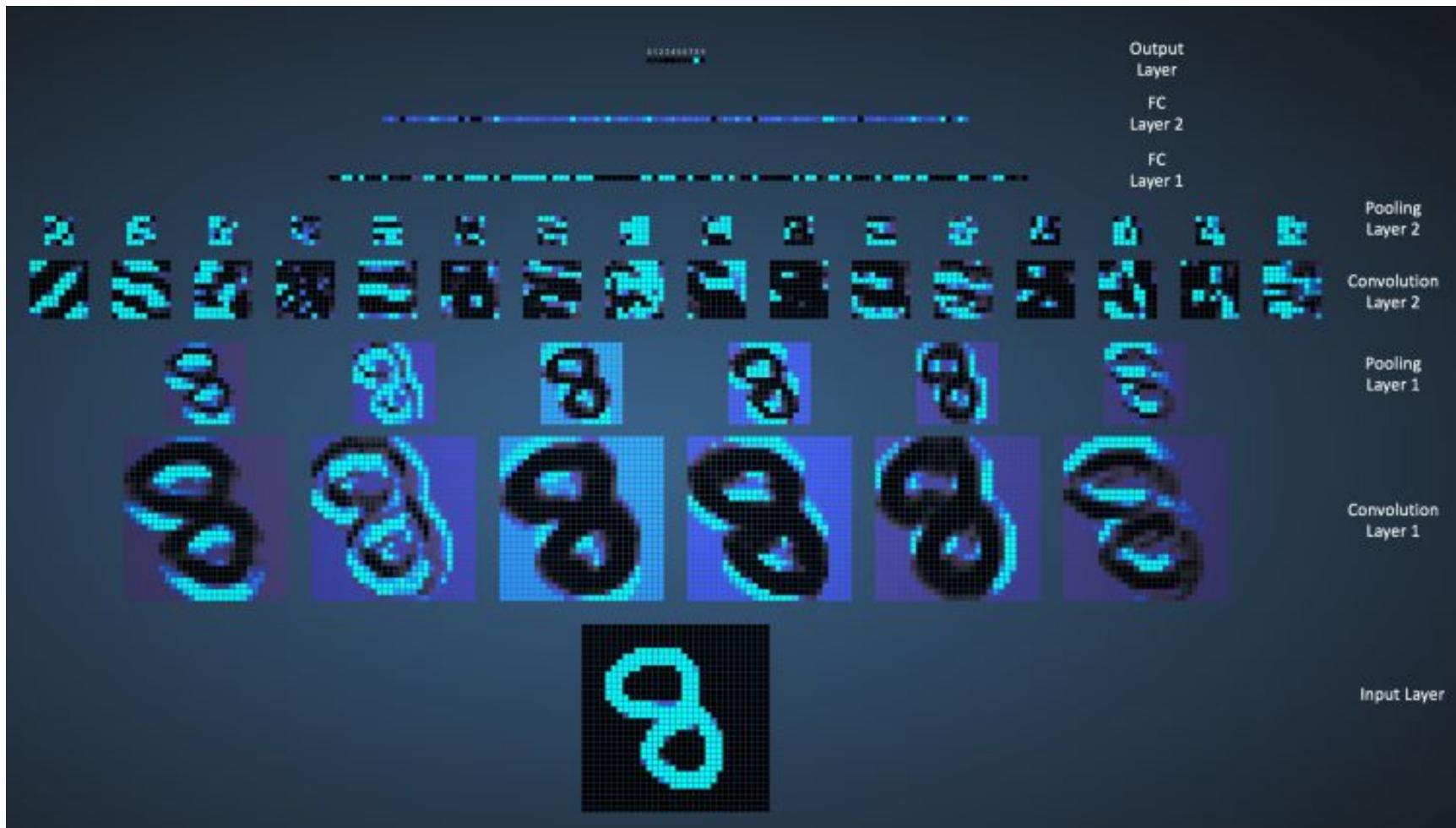


Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



Activation maps: example



5 minutes break

Regularization

Do we need it?

Why do we need it?

Regularization

Why do we need it?

The training data contains information about the regularities in the mapping from input to output. But it also contains sampling error.

There will be accidental regularities just because of the particular training cases that were chosen.

When we fit the model, it cannot tell which regularities are real and which are caused by sampling error...

...So it fits both kinds of regularity. If the model is very flexible it can model the sampling error really well.

Regularization

Why do we need it?

The training data contains information about the regularities in the mapping from input to output. But it also contains sampling error.

There will be accidental regularities just because of the particular training cases that were chosen.

When we fit the model, it cannot tell which regularities are real and which are caused by sampling error...

...So it fits both kinds of regularity. If the model is very flexible it can model the sampling error really well.

Overfitting again! :-(

Regularization

Solution:

Penalize large weights using penalties or constraints on their squared values (*L2 penalty*) or absolute values (*L1 penalty*).

Regularization

L2 penalty Solution:

$$\widetilde{L}(\mathbf{w}) = L(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^2$$

Regularization

L2 penalty Solution:

$$\widetilde{L}(\mathbf{w}) = L(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^2$$

Applying gradient descent to this new cost function we obtain:

Regularization

L2 penalty Solution:

$$\widetilde{L}(\mathbf{w}) = L(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^2$$

Applying gradient descent to this new cost function we obtain:

$$w_i \leftarrow w_i - \eta \frac{\partial L}{\partial w_i} - \eta \lambda w_i.$$

Regularization

L2 penalty Solution:

$$\widetilde{L}(\mathbf{w}) = L(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^2$$

Applying gradient descent to this new cost function we obtain:

$$w_i \leftarrow w_i - \eta \frac{\partial L}{\partial w_i} - \eta \lambda w_i.$$

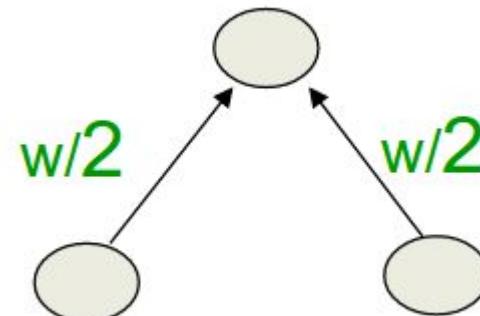
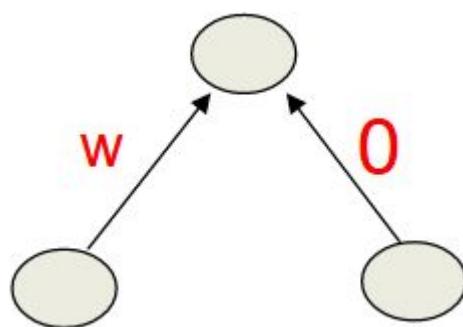


..Also called *Weight Decay!*

Regularization

Effect of L2 regularization:

- It prevents the network from using weights that it does not need.
- This often improves generalization a lot because it helps to stop the network from fitting the sampling error.
- It makes a smoother model in which the output changes more slowly as the input changes.
- If the network has two very similar inputs it prefers to put half the weight on each rather than all the weight on one.



Regularization

- On big&wide layers, overfitting could happen in the form of co-adaptation
- It happens especially on dense layers → FC
- We could reduce the number of neurons in a layer, but it would reduce capacity too! :-(
- How could we keep high capacity and low overfitting at the same time?

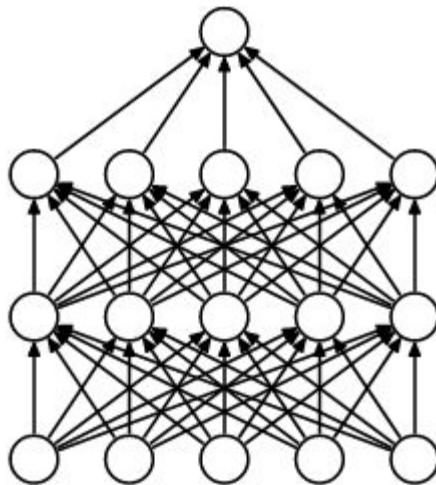
Regularization

- On big&wide layers, overfitting could happen in the form of co-adaptation
- It happens especially on dense layers → FC
- We could reduce the number of neurons in a layer, but it would reduce capacity too! :-(
- How could we keep high capacity and low overfitting at the same time?

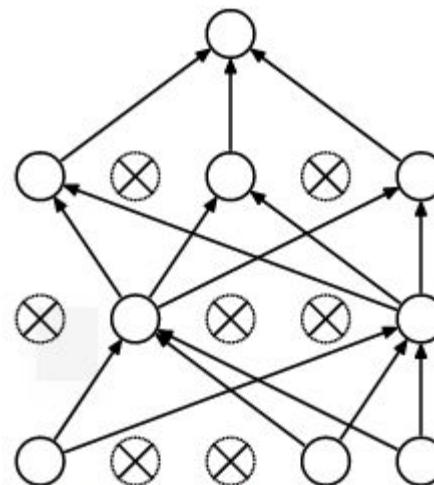
Dropout method!

Regularization

Dropout:



(a) Standard Neural Net



(b) After applying dropout.

Source: "Dropout: A Simple Way to Prevent Neural Networks from Overfitting". Nitish Srivastava et. al.

Regularization

Dropout:

- We drop-out some neurons at *training time* and we keep all neurons at test time
- It helps avoiding co-adaptation of neurons and thus reduce overfitting
- It has an “ensemble networks” effect
- Usual value is 0.5: half of neurons removed at training time
- More training time
- Used on several powerful architectures
- Suggested every time you have a FC layer
- Not clear effect on conv layers - it could be used as a source of noise, and it helps sometimes for fighting overfitting

Overfitting & regularization: The best advice

- As Hinton said, if your network is not overfitting you are not learning enough!
- His suggestion: overfit the network model on data, then regularize the hell out of it!
- This means: increase the depth and the wideness of the network until we see overfitting, then we fight overfit with *Weight Decay & Dropout*

Normalization problem

- Any training algorithm usually works better on normalized data
- We could normalize input data, but..
- After each conv / FC layer, data is not normalized anymore!
- This slows down the training process, brings convergence problems, lead to sub-optimal results
- Vanishing gradient nightmare!
- ...Let normalize data after each conv / FC layer!

Normalization problem

- Any training algorithm usually works better on normalized data
- We could normalize input data, but..
- After each conv / FC layer, data is not normalized anymore!
- This slows down the training process, brings convergence problems, lead to sub-optimal results
- Vanishing gradient nightmare!
- ...Let normalize data after each conv / FC layer!

Batch normalization layer!

Normalization problem

Batch normalization layer:

- In order to normalize X data:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \text{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Normalization problem

Batch normalization layer:

- In order to normalize X data:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

- We don't know the mean and the variance!
- ...but we can calculate them on each batch at training time:
learning means and var
- At test time we use learned statistics to approximate true dataset mean and variance

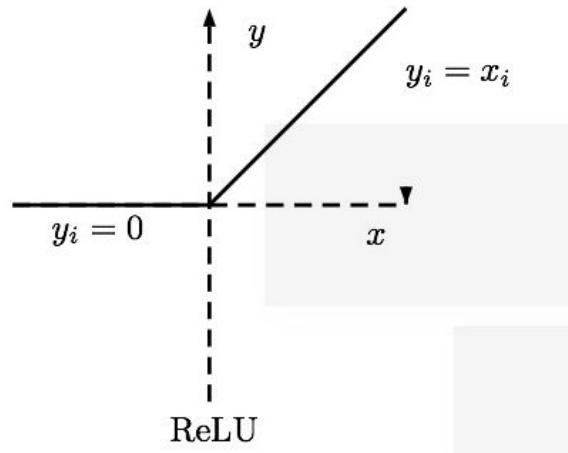
Normalization problem

Batch normalization layer:

- Widely used on all recent CNN architectures after each conv and FC layer
- It normalize the flowing data in the whole network
- Very good remedy to vanishing gradient problem up to ~20 layers
- It increase memory consumption, but...
- ...it speeds up network convergence A LOT!
- Always use it if you don't have a good reason to not use it ;-)

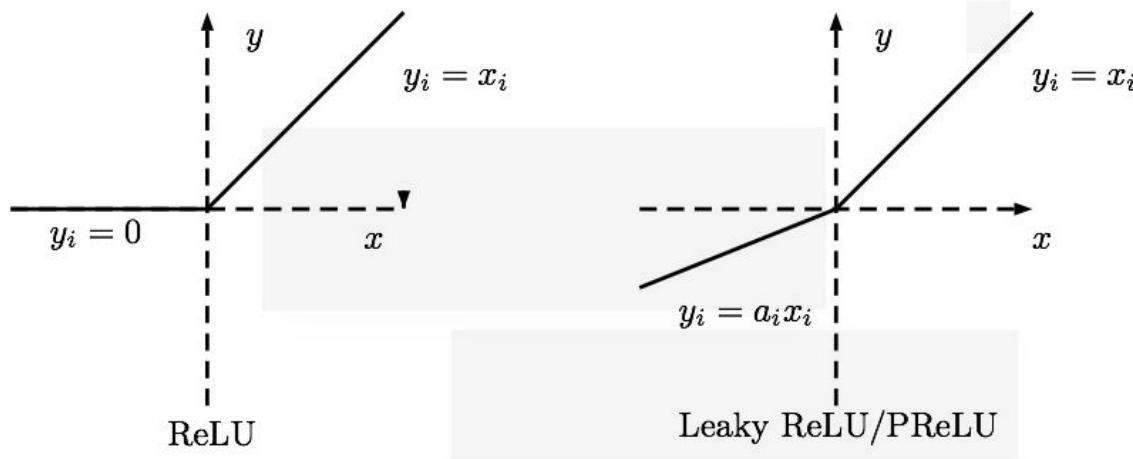
Leaky ReLU

- ReLU is the standard non linear activation unit
- ...but getting zero output for a whole range of input is not always nice!
- ReLU units could “die” if they get always negative input



Leaky ReLU

- ReLU is the standard non linear activation unit
- ...but getting zero output for a whole range of input is not always nice!
- ReLU units could “die” if they get always negative input
- LeakyReLU solution



- It avoids ReLU death and it helps reducing vanishing gradient problem

Recap

Standard Convolutional Neural Networks layers:

- Convolutional, Pooling, ReLU, Fully Connected

Advanced Convolutional Neural Networks layers:

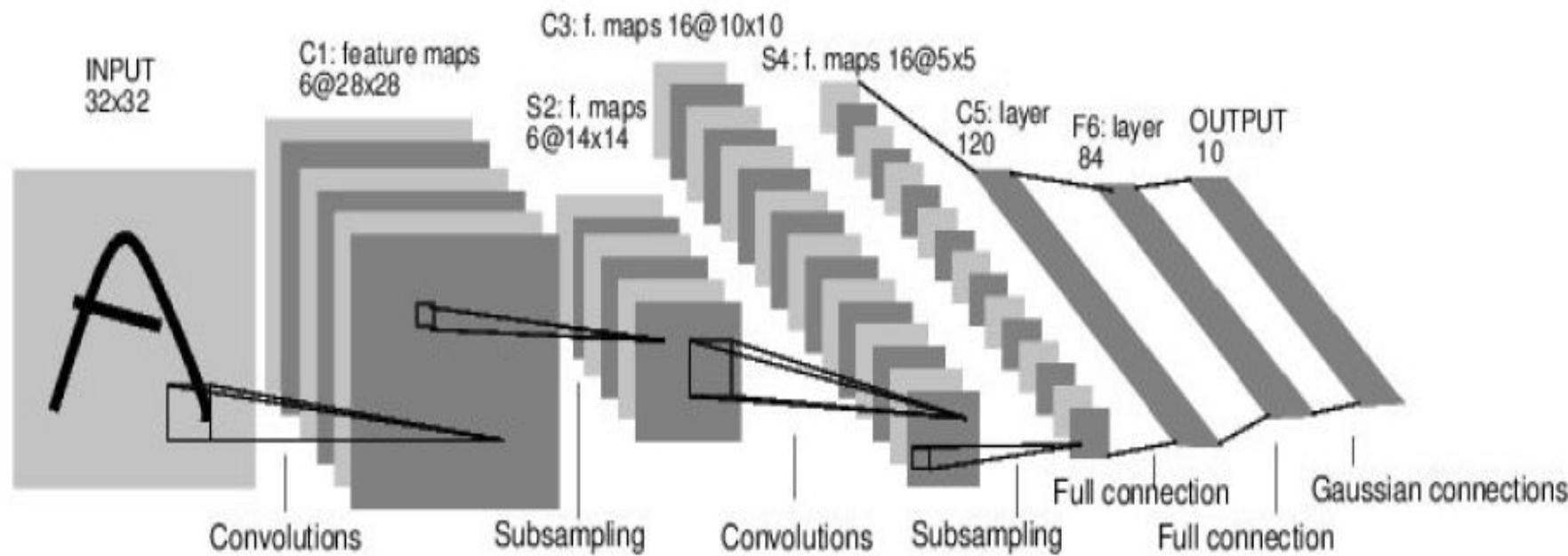
- Batch Normalization, Average Pooling, LeakyReLU

Regularization & improving techniques:

- Weight Decay, Dropout, batchnorm

Case Study: LeNet-5

[LeCun et al., 1998]

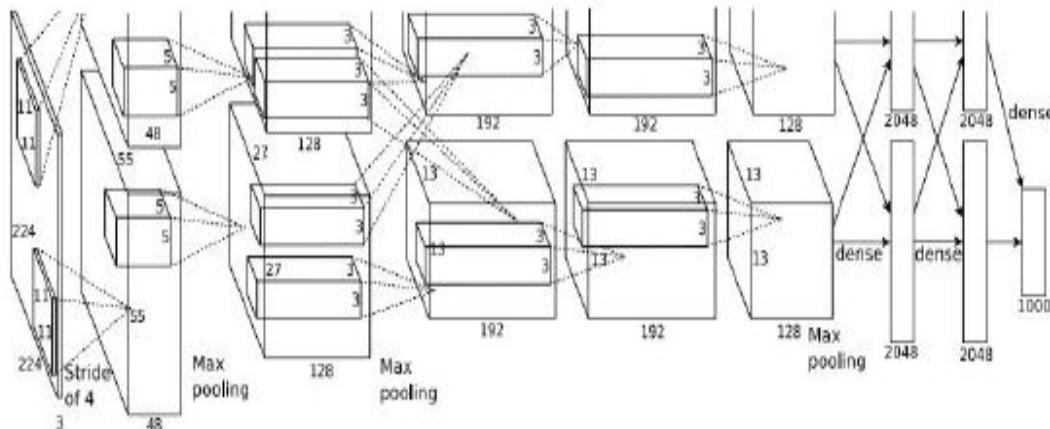


Conv filters were 5×5 , applied at stride 1

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

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

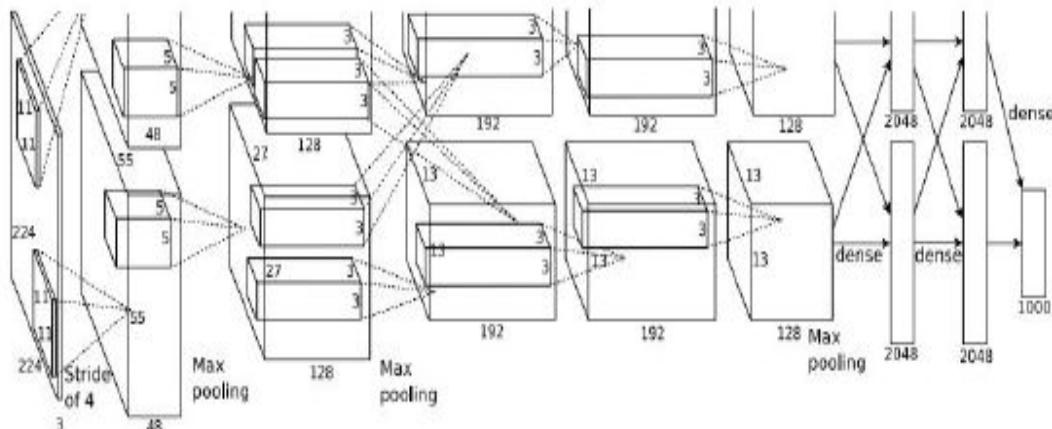
First layer (CONV1): 96 11x11 filters applied at stride 4

=>

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

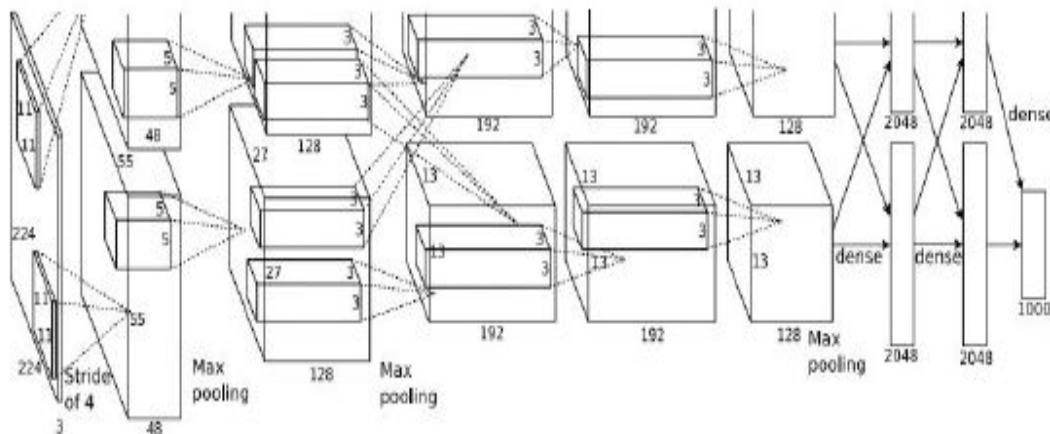
=>

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

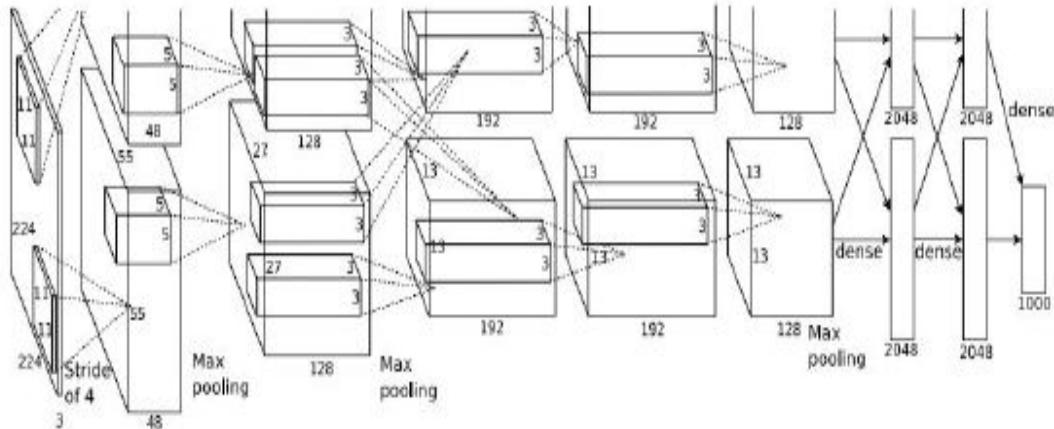
=>

Output volume **[55x55x96]**

Parameters: **(11*11*3)*96 = 35K**

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

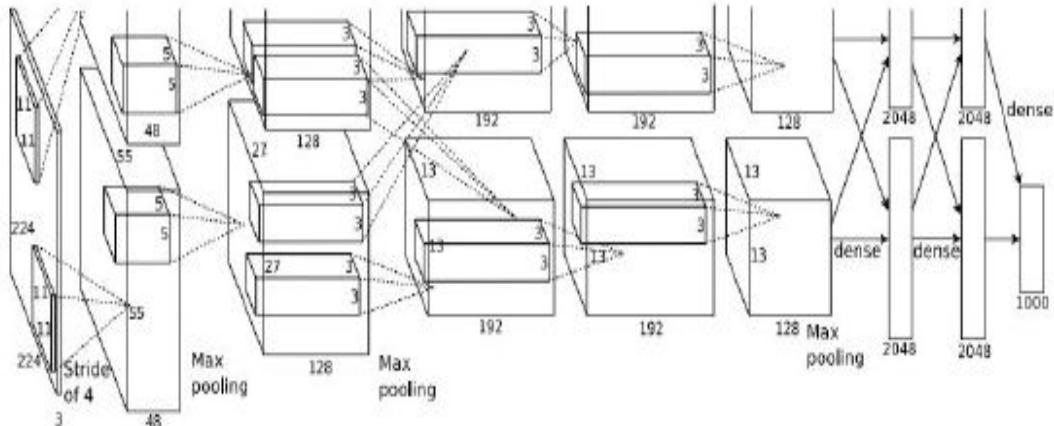
After CONV1: 55x55x96

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

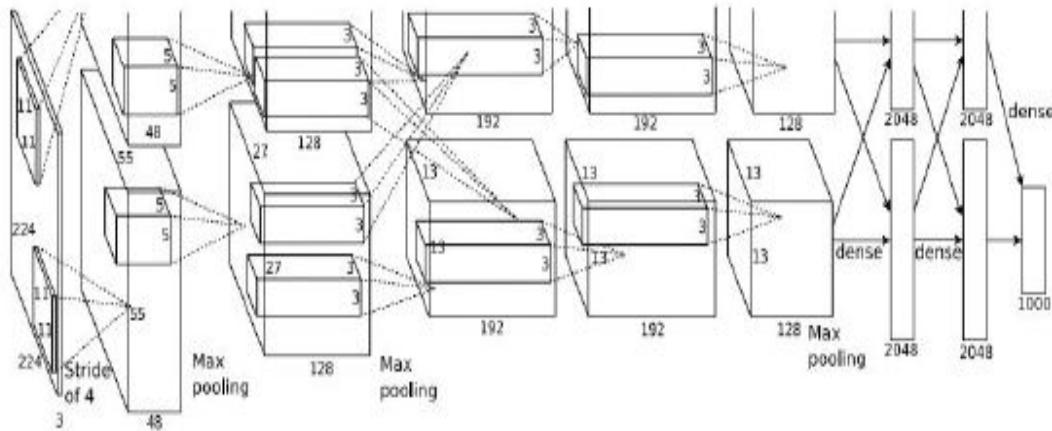
Second layer (POOL1): 3x3 filters applied at stride 2

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

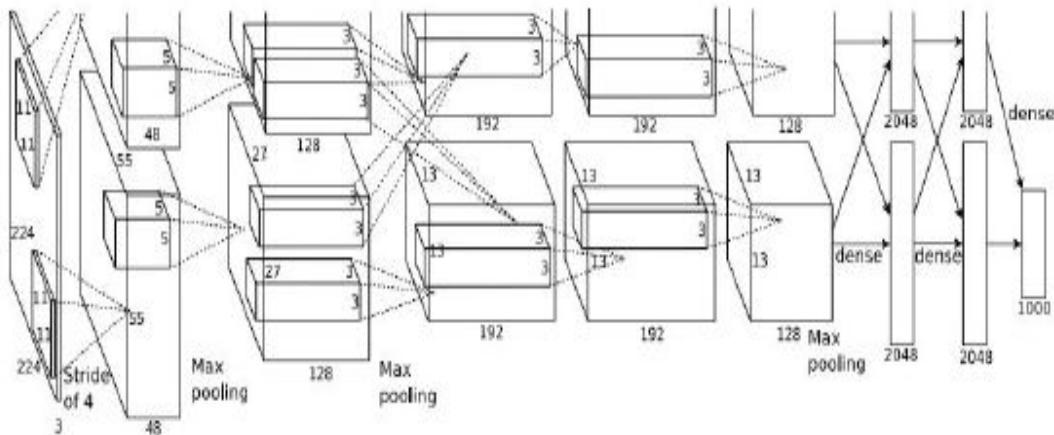
Second layer (POOL1): 3x3 filters applied at stride 2

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

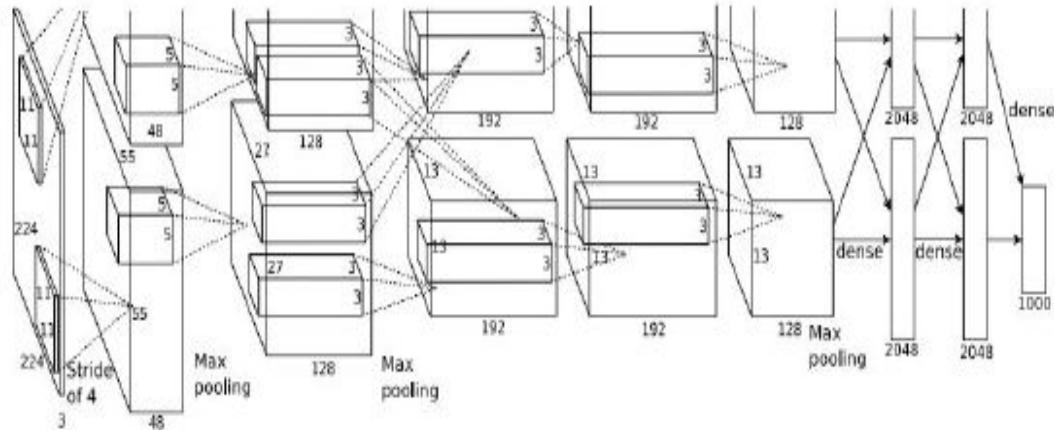
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[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]

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

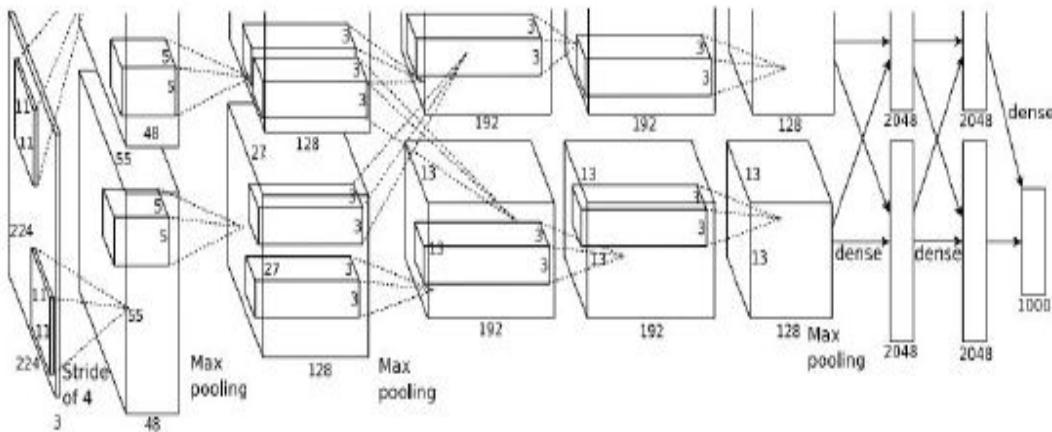
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

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

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



Details/Retrospectives:

- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 15.4%

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

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

best model

11.2% top 5 error in ILSVRC 2013

->

7.3% top 5 error

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 2: Number of parameters (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

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

ConvNet Configuration			
B	C	D	
13 weight layers	16 weight layers	16 weight layers	19
put (224 x 224 RGB image)			
conv3-64	conv3-64	conv3-64	cc
conv3-64	conv3-64	conv3-64	cc
	maxpool		
conv3-128	conv3-128	conv3-128	co1
conv3-128	conv3-128	conv3-128	co1
	maxpool		
conv3-256	conv3-256	conv3-256	co1
conv3-256	conv3-256	conv3-256	co1
	conv1-256	conv3-256	co1
		conv3-256	co1
	maxpool		
conv3-512	conv3-512	conv3-512	co1
conv3-512	conv3-512	conv3-512	co1
	conv1-512	conv3-512	co1
		conv3-512	co1
	maxpool		
conv3-512	conv3-512	conv3-512	co1
conv3-512	conv3-512	conv3-512	co1
	conv1-512	conv3-512	co1
		conv3-512	co1
	maxpool		
FC-4096			
FC-4096			
FC-1000			
soft-max			

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

TOTAL memory: 24M * 4 bytes \approx 93MB / image (only forward! ~ 2 for bwd)
TOTAL params: 138M parameters

ConvNet Configuration			
B	C	D	E
13 weight layers	16 weight layers	16 weight layers	19 weight layers
put (224 x 224 RGB image)			
conv3-64	conv3-64	conv3-64	cc
conv3-64	conv3-64	conv3-64	cc
	maxpool		
conv3-128	conv3-128	conv3-128	co1
conv3-128	conv3-128	conv3-128	co1
	maxpool		
conv3-256	conv3-256	conv3-256	co1
conv3-256	conv3-256	conv3-256	co1
	conv1-256	conv3-256	co1
		conv3-256	co1
	maxpool		
conv3-512	conv3-512	conv3-512	co1
conv3-512	conv3-512	conv3-512	co1
	conv1-512	conv3-512	co1
		conv3-512	co1
	maxpool		
conv3-512	conv3-512	conv3-512	co1
conv3-512	conv3-512	conv3-512	co1
	conv1-512	conv3-512	co1
		conv3-512	co1
	maxpool		
FC-4096			
FC-4096			
FC-1000			
soft-max			

INPUT: [224x224x3] memory: 224*224*3=150K params: 0 (not counting biases)

CONV3-64: [224x224x64] memory: **224*224*64=3.2M** params: $(3*3*3)*64 = 1,728$

CONV3-64: [224x224x64] memory: **224*224*64=3.2M** params: $(3*3*64)*64 = 36,864$

POOL2: [112x112x64] memory: 112*112*64=800K params: 0

CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: $(3*3*64)*128 = 73,728$

CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: $(3*3*128)*128 = 147,456$

POOL2: [56x56x128] memory: 56*56*128=400K params: 0

CONV3-256: [56x56x256] memory: 56*56*256=800K params: $(3*3*128)*256 = 294,912$

CONV3-256: [56x56x256] memory: 56*56*256=800K params: $(3*3*256)*256 = 589,824$

CONV3-256: [56x56x256] memory: 56*56*256=800K params: $(3*3*256)*256 = 589,824$

POOL2: [28x28x256] memory: 28*28*256=200K params: 0

CONV3-512: [28x28x512] memory: 28*28*512=400K params: $(3*3*256)*512 = 1,179,648$

CONV3-512: [28x28x512] memory: 28*28*512=400K params: $(3*3*512)*512 = 2,359,296$

CONV3-512: [28x28x512] memory: 28*28*512=400K params: $(3*3*512)*512 = 2,359,296$

POOL2: [14x14x512] memory: 14*14*512=100K params: 0

CONV3-512: [14x14x512] memory: 14*14*512=100K params: $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory: 14*14*512=100K params: $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory: 14*14*512=100K params: $(3*3*512)*512 = 2,359,296$

POOL2: [7x7x512] memory: 7*7*512=25K params: 0

FC: [1x1x4096] memory: 4096 params: $7*7*512*4096 = \textcolor{blue}{102,760,448}$

FC: [1x1x4096] memory: 4096 params: $4096*4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params: $4096*1000 = 4,096,000$

Note:

Most memory is in early CONV

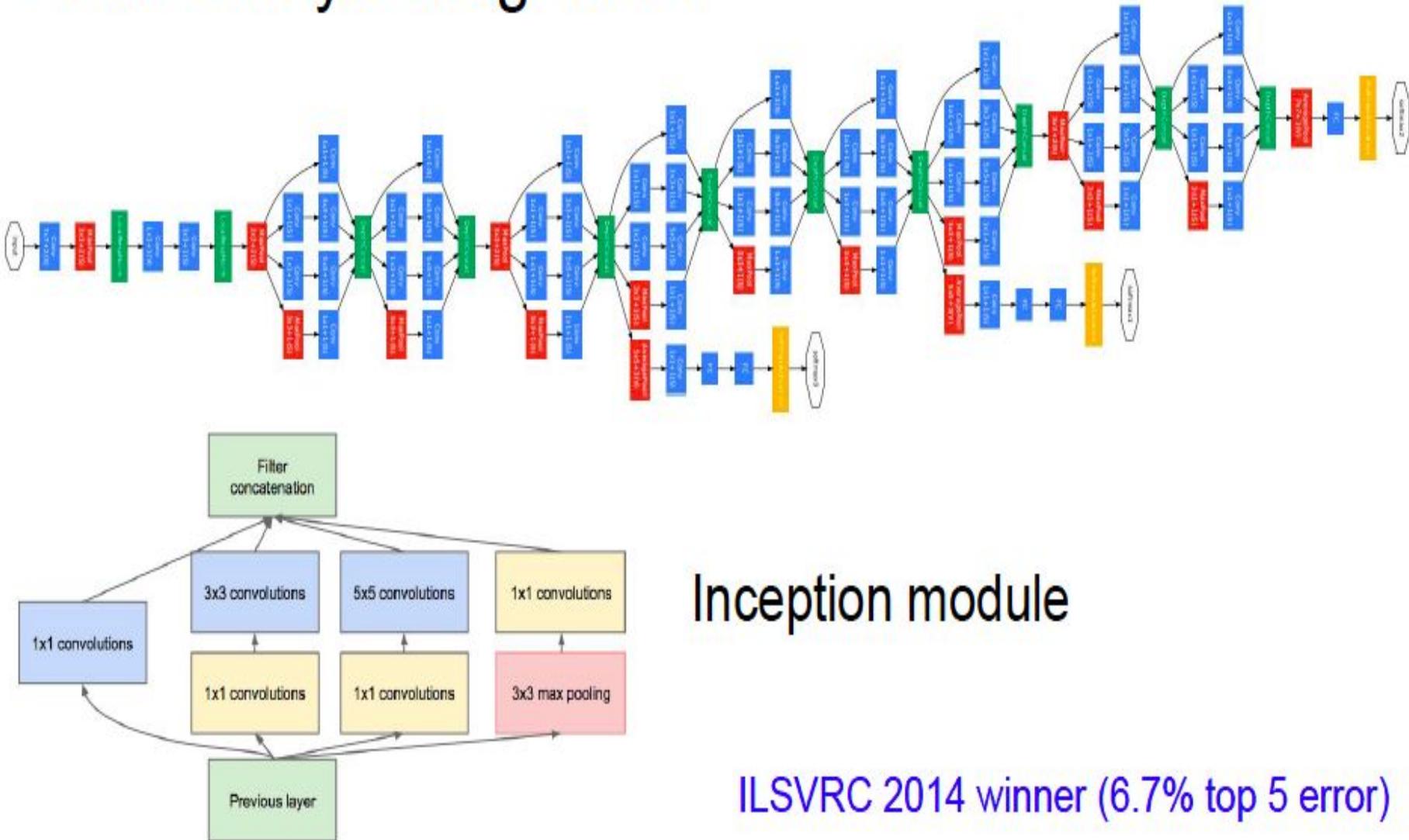
Most params are in late FC

TOTAL memory: 24M * 4 bytes ~= 93MB / image (only forward! ~*2 for bwd)

TOTAL params: 138M parameters

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)

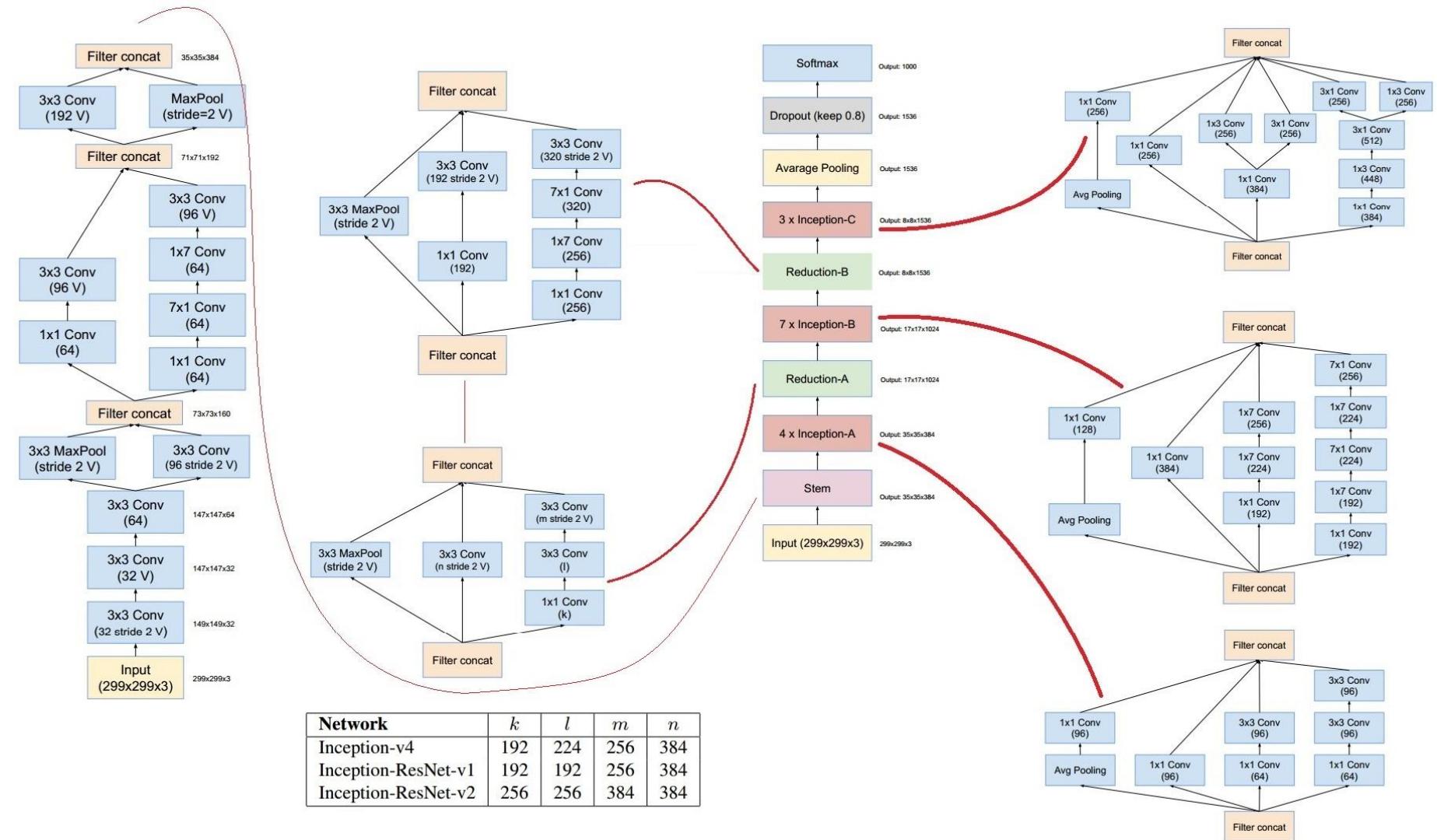
Compared to AlexNet:

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

Case study: Inception networks



Case study: Inception V4



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

- Increasing the depth of an architecture leads to improved accuracy...
- Until a point!
- 34-layer plain network exhibit strong vanishing gradient problem
- Batchnorm helps, but...
- How could we increase network depth definitely avoiding this issues?

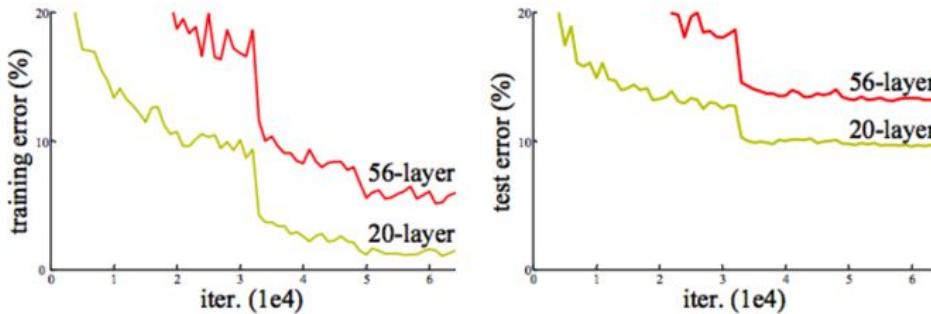
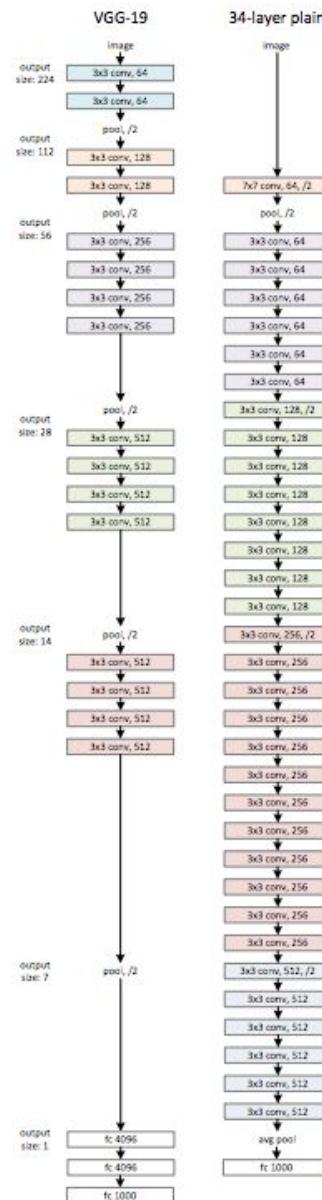


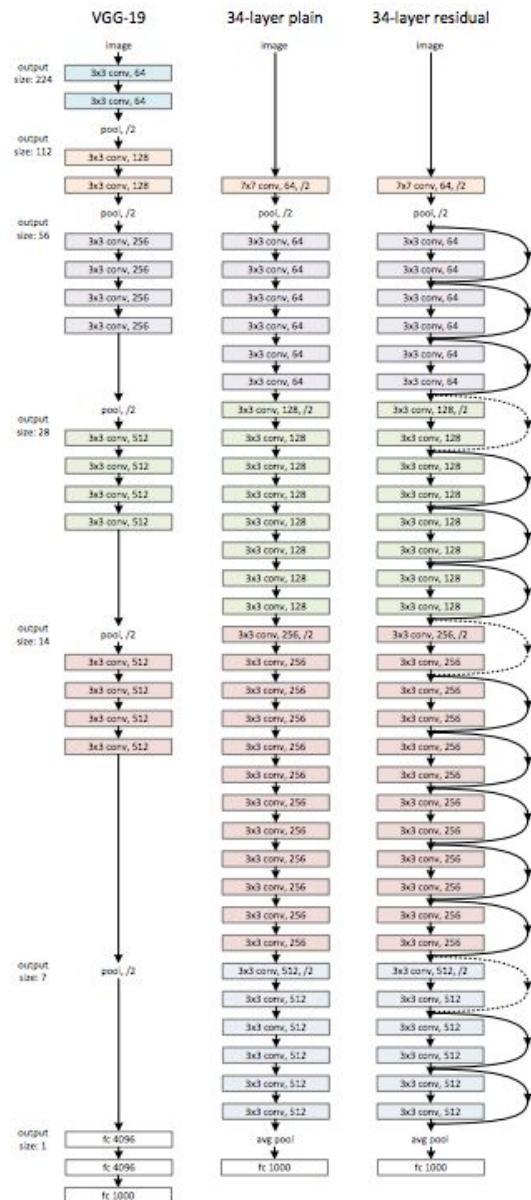
Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.



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

- Increasing the depth of an architecture leads to improved accuracy...
- Until a point!
- 34-layer plain network exhibit strong vanishing gradient problem
- Batchnorm helps, but...
- How could we increase network depth definitely avoiding this issues?

- Skip Connections!
- Residual networks!



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

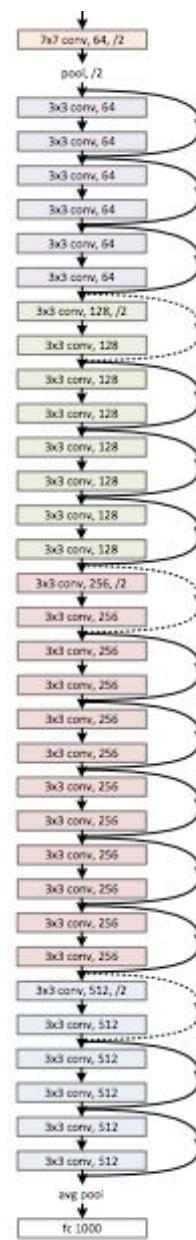
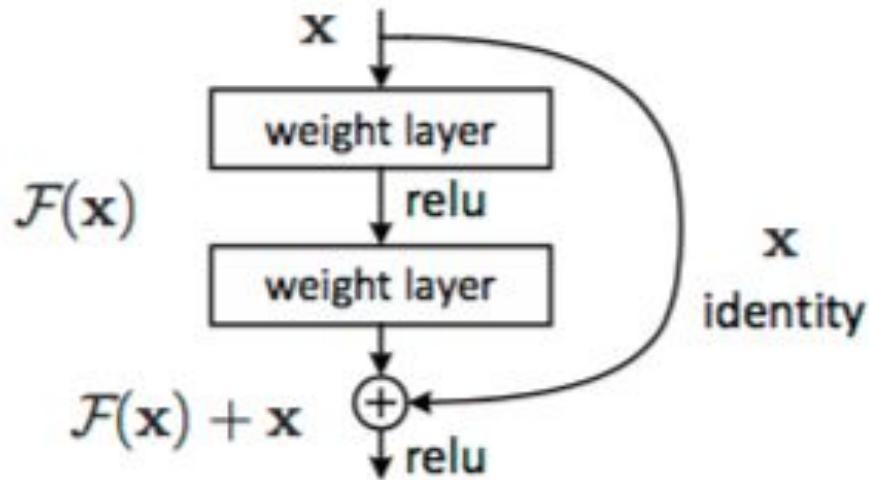
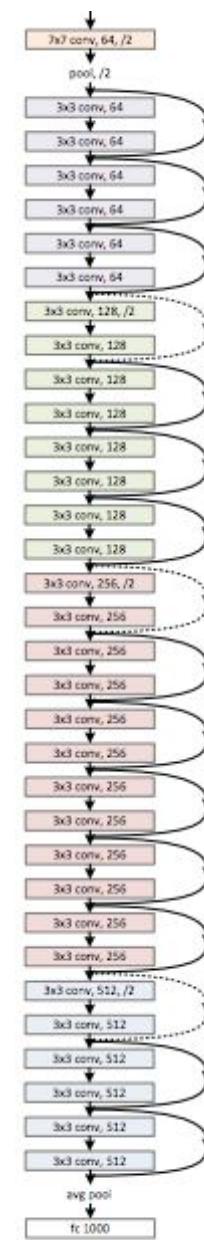
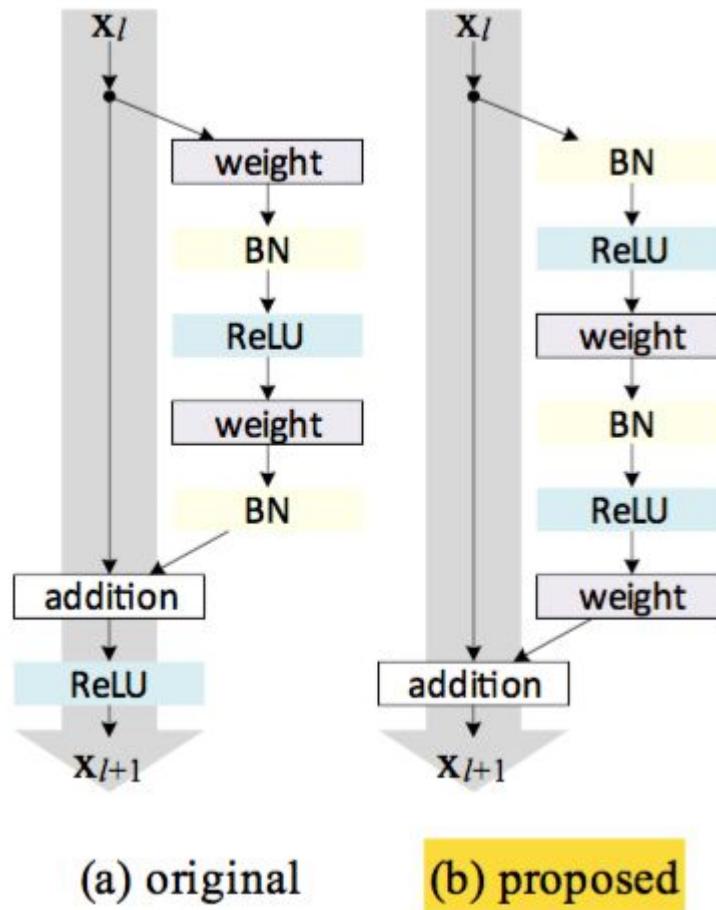


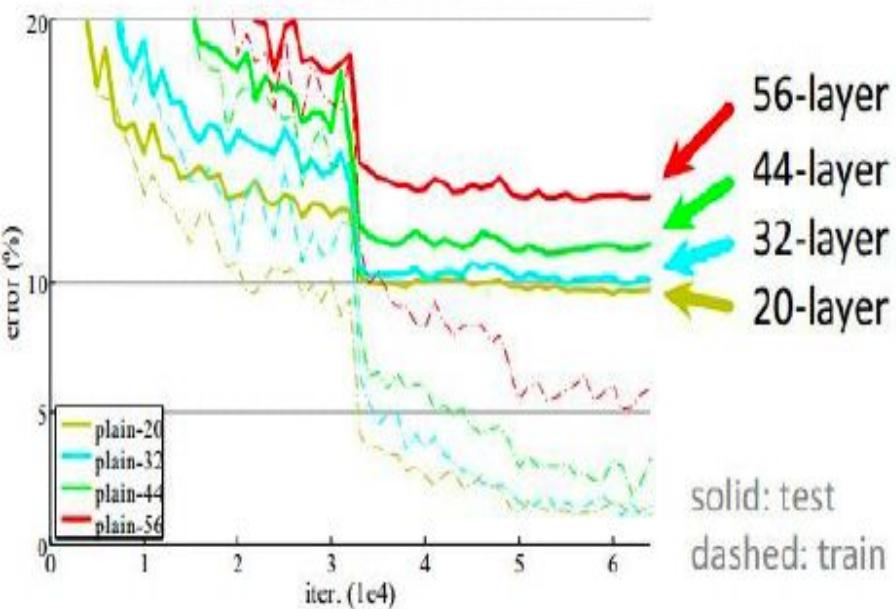
Figure 2. Residual learning: a building block.

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

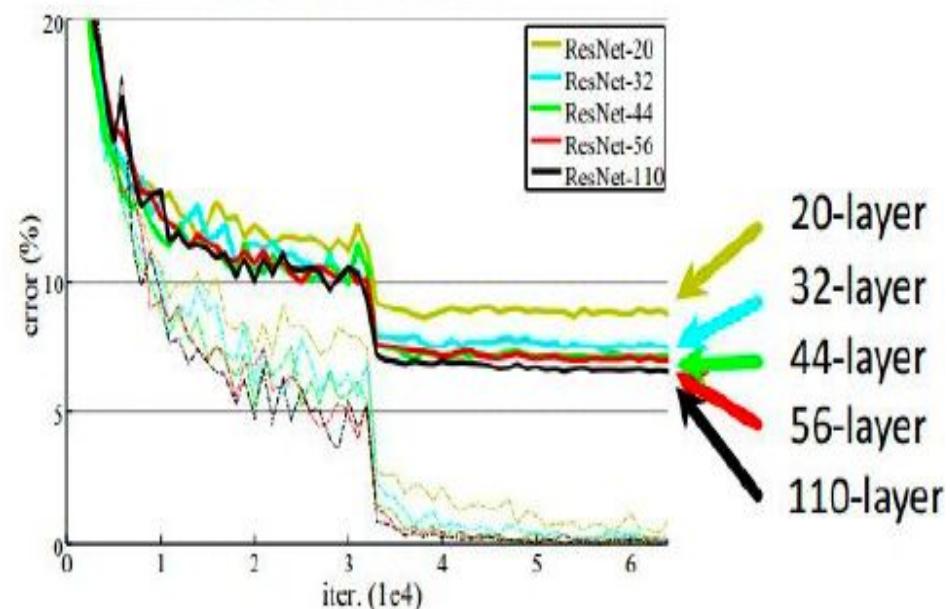


CIFAR-10 experiments

CIFAR-10 plain nets



CIFAR-10 ResNets



Case Study: ResNet

[He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)



MSRA @ ILSVRC & COCO 2015 Competitions

- **1st places** in all five main tracks
 - ImageNet Classification: “Ultra-deep” (quote Yann) **152-layer** nets
 - ImageNet Detection: **16%** better than 2nd
 - ImageNet Localization: **27%** better than 2nd
 - COCO Detection: **11%** better than 2nd
 - COCO Segmentation: **12%** better than 2nd

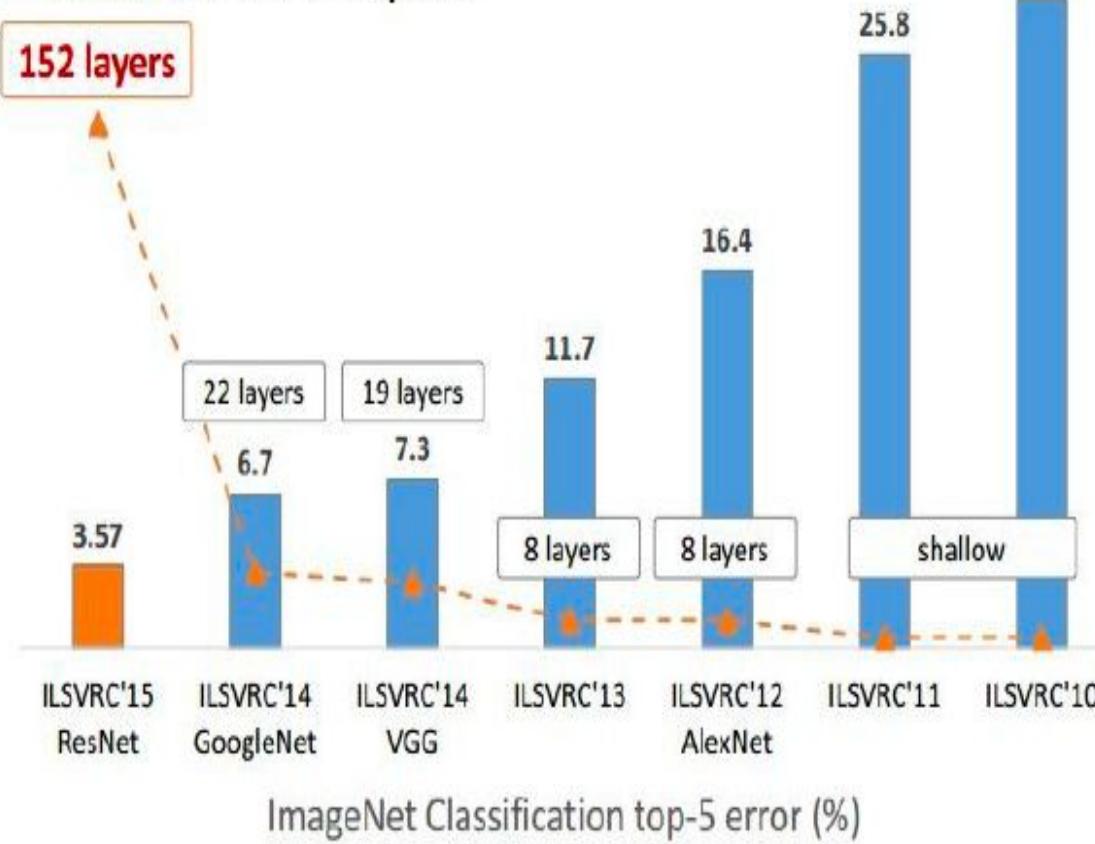
*improvements are relative numbers



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. “Deep Residual Learning for Image Recognition”, arXiv 2015.

Slide from Kaiming He's recent presentation <https://www.youtube.com/watch?v=1PGLj-uKT1w>

Revolution of Depth



Case Study: ResNet

[He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)

Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)



ResNet, **152 layers**
(ILSVRC 2015)

Microsoft
Research



2-3 weeks of training
on 8 GPU machine

at runtime: faster
than a VGGNet!
(even though it has
8x more layers)



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

(slide from Kaiming He's recent presentation)

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

$[(CONV-RELU)^*N-POOL?]^*M-(FC-RELU)^*K,SOFTMAX$

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

That's all!