

DEEP LEARNING

FFR135, Artificial Neural Networks

Olof Mogren

Chalmers University of Technology

October 2016

DEEP LEARNING

- Artificial neural networks
- Many layers of abstractions
- Outperforms traditional methods in:
 - Image classification
 - Natural language processing
 - Machine translation
 - Sentiment analysis
 - Speech recognition
 - Reinforcement learning



SEMI-RECENT PROGRESS

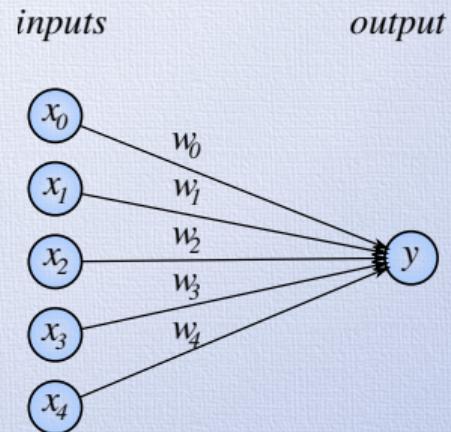
- 2006: Depth breakthrough:
layerwise pretrained Restricted
Boltzmann Machines
- GPUs
- Practical use
*Real applications from Google,
Facebook, Tesla, Microsoft, Apple,
and others!*



A fast learning algorithm for deep belief nets; Hinton, Osindero, Tehi; Neural Computation; 2006

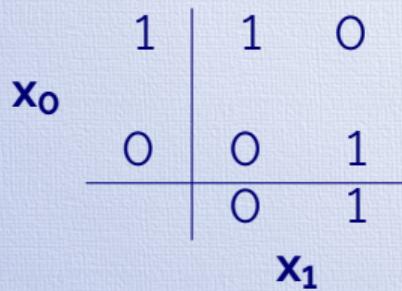
PERCEPTRON

- 1943, McCulloch & Pitts (neuron model)
- 1958, Rosenblatt (perceptron)
- Linear (binary) classification of inputs
- Can not learn any non-linear function (e.g. XOR)

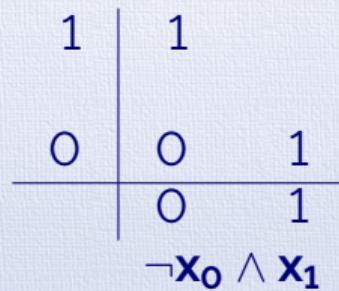


MODELLING XOR

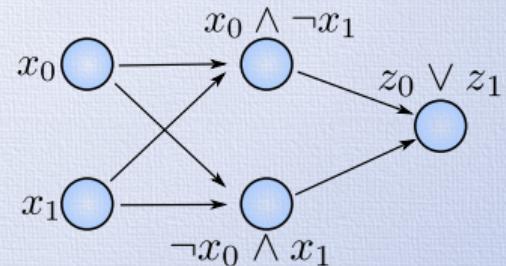
MODELLING XOR



$$x_0 \wedge \neg x_1$$

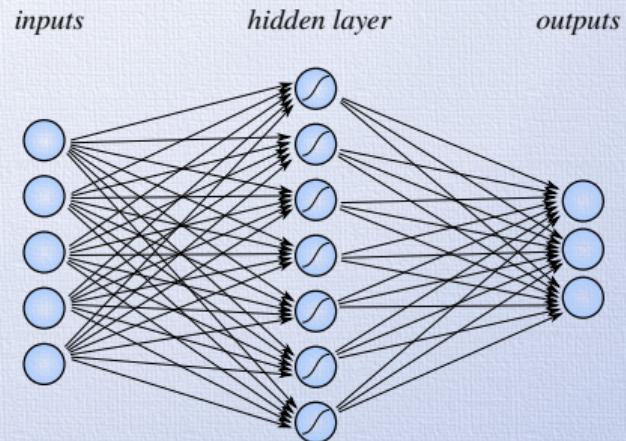


$$\neg x_0 \wedge x_1$$



MULTI-LAYER PERCEPTRON

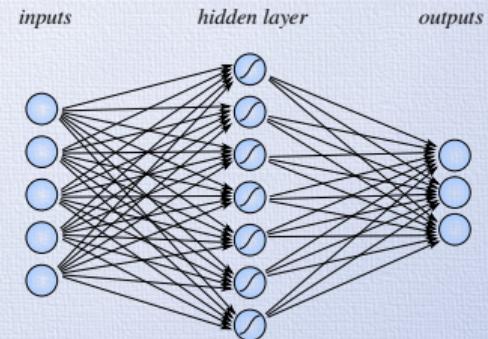
- Combining layers lets us represent non-linear functions
- Each layer:
 - Linear transformation:
 $\mathbf{a} = W\mathbf{x} + \mathbf{b}$
 - Non-linear (element-wise) activation: $\mathbf{h} = g(\mathbf{a})$



MODELLING FUNCTIONS

- Universal function approximation
- Stacking layers: function composition
- Apply error/loss function to output
- Continuously differentiable; chain rule
- Propagating errors (backpropagation)
- (Mini-batch) Stochastic gradient descent (SGD)

details



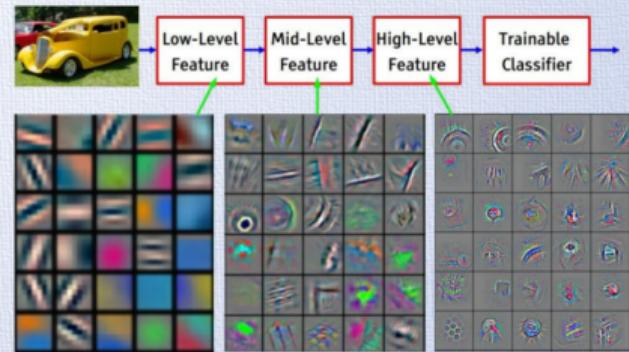
MOTIVATION OF DEPTH

- More compact representation (exponentially)
- There are boolean functions that require
 - **Polynomial** number of units (**deep** architecture)
 - **Exponential** number of units (**shallow** architecture)
 - E.g., parity function (for n input bits):
 - efficiently represented with depth **$O(\log n)$**
 - but **$O(2^n)$** gates if represented by a depth two circuit (Yao, 1985)

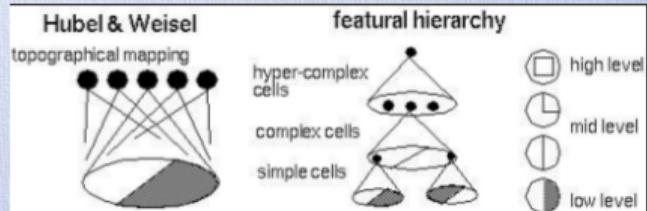
Exploring Strategies for Training Deep Neural Networks; Larochelle, Bengio, Louradour, Lamblin; JMLR 2009

LEARNING LEVELS OF REPRESENTATION

- Each layer:
non-linear transformation of inputs:
 $\mathbf{h} = \text{sigmoid}(\mathbf{W}\mathbf{x} + \mathbf{b})$
- Learning representations; abstractions
- No feature engineering!

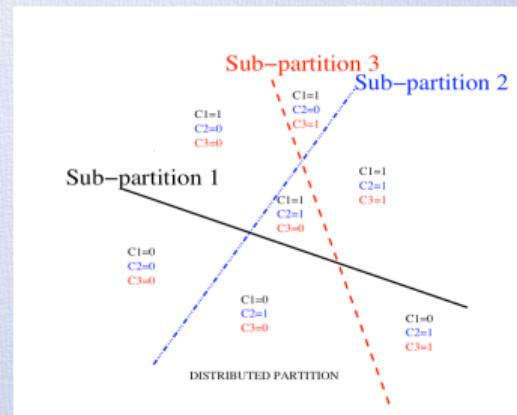


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



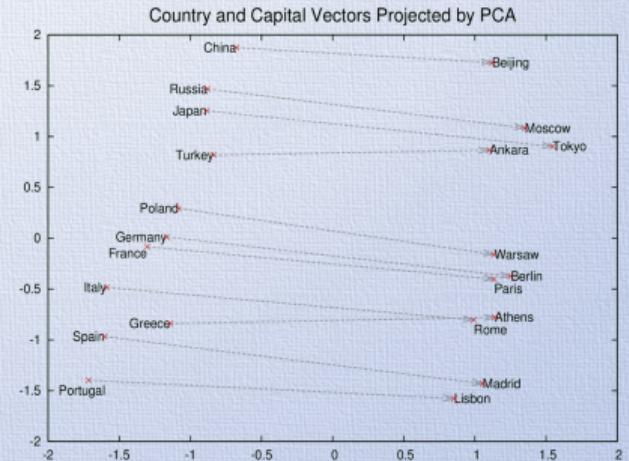
DISTRIBUTED REPRESENTATIONS

- E.g.: big, yellow, Volkswagen
- Non-distributed representations:
 n binary parameters $\rightarrow n$ values
- E.g.: Clustering, n-grams, decision trees, etc.
- NNs learn distributed representations
- Distributed representations:
 n binary parameters $\rightarrow 2^n$ possible values



EXAMPLE: WORD EMBEDDINGS

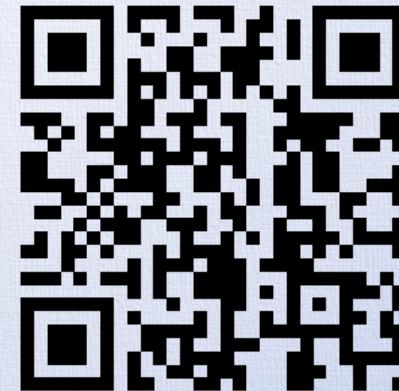
- Distributed representations for words
- word2vec, glove, etc.



DEEP LEARNING IN JAVASCRIPT

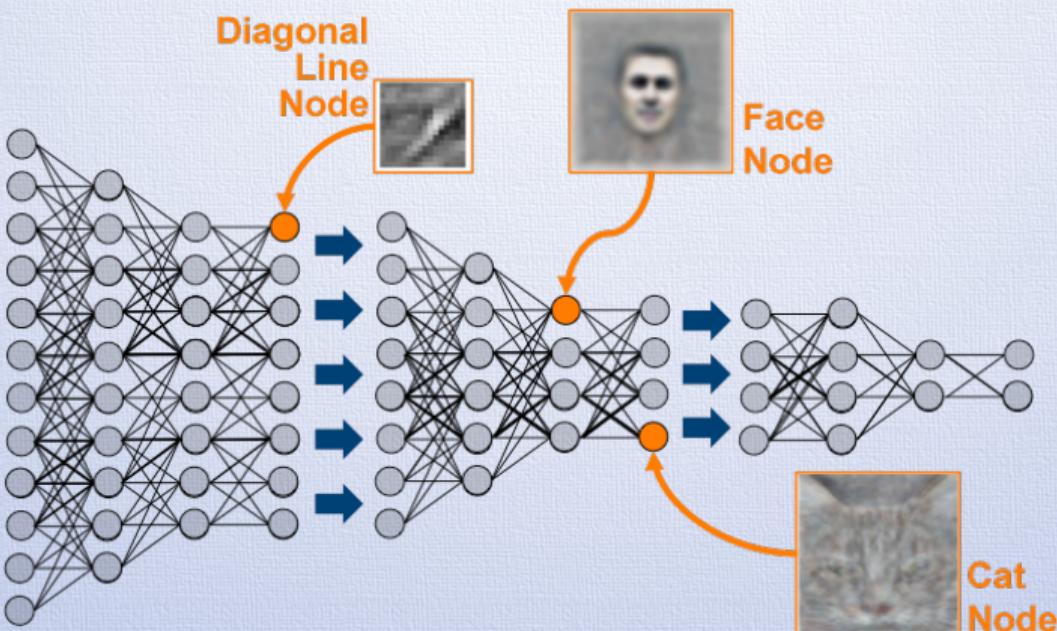


cs231n.stanford.edu



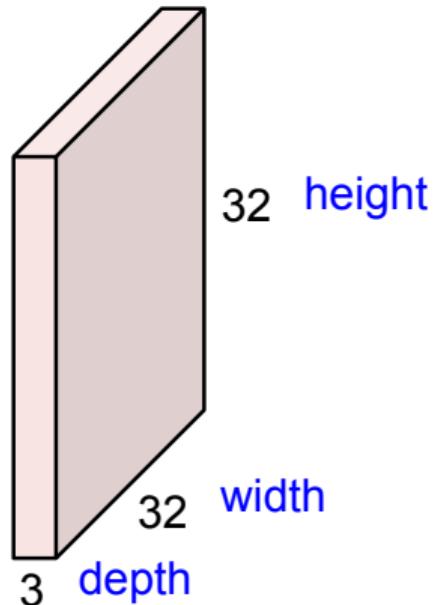
playground.tensorflow.org

LEVELS OF ABSTRACTIONS



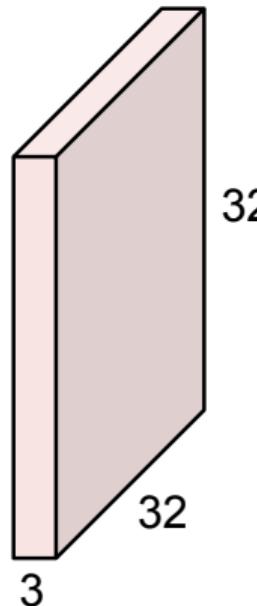
Convolution Layer

32x32x3 image



Convolution Layer

32x32x3 image



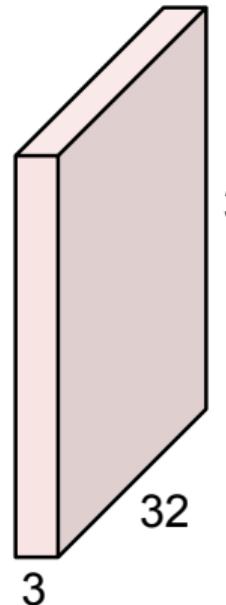
5x5x3 filter



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

Convolution Layer

32x32x3 image



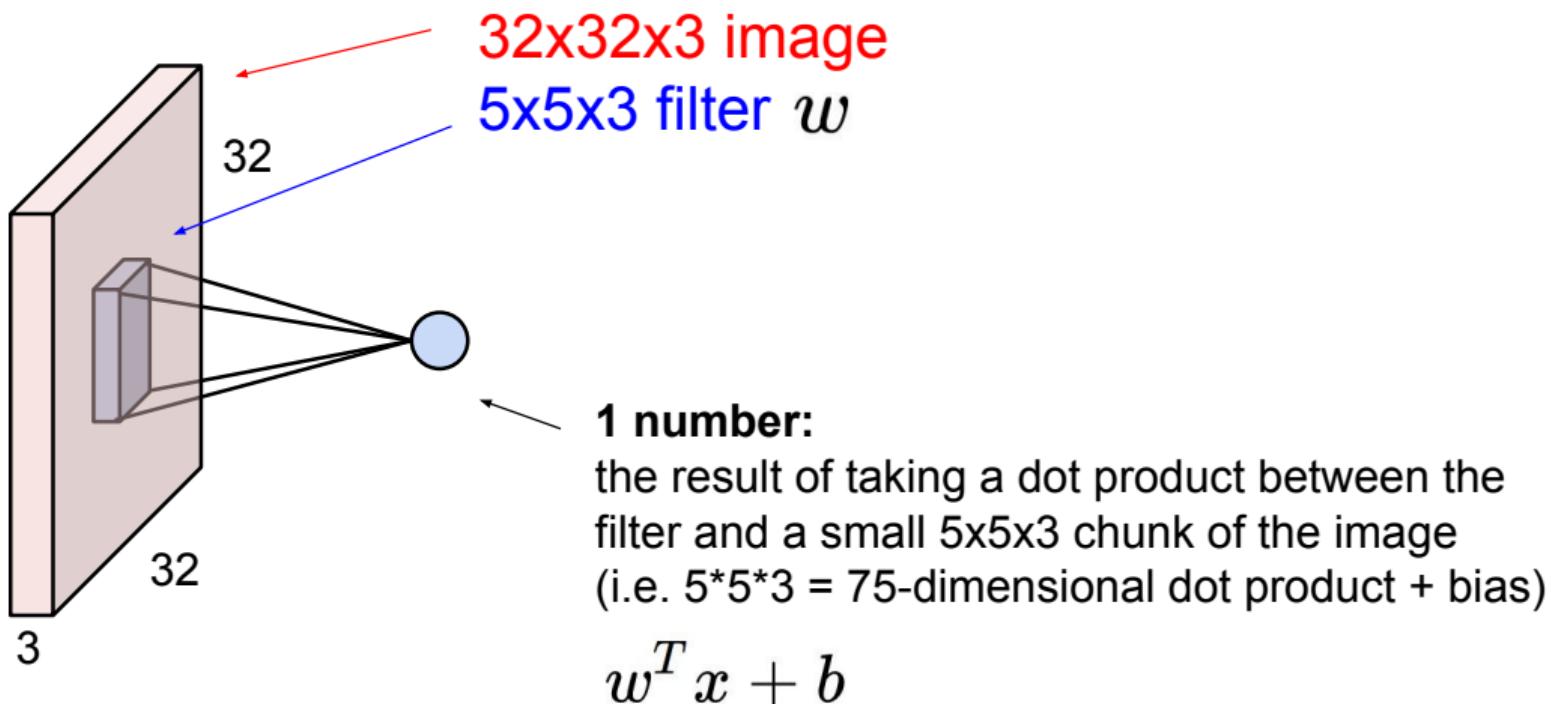
5x5x3 filter



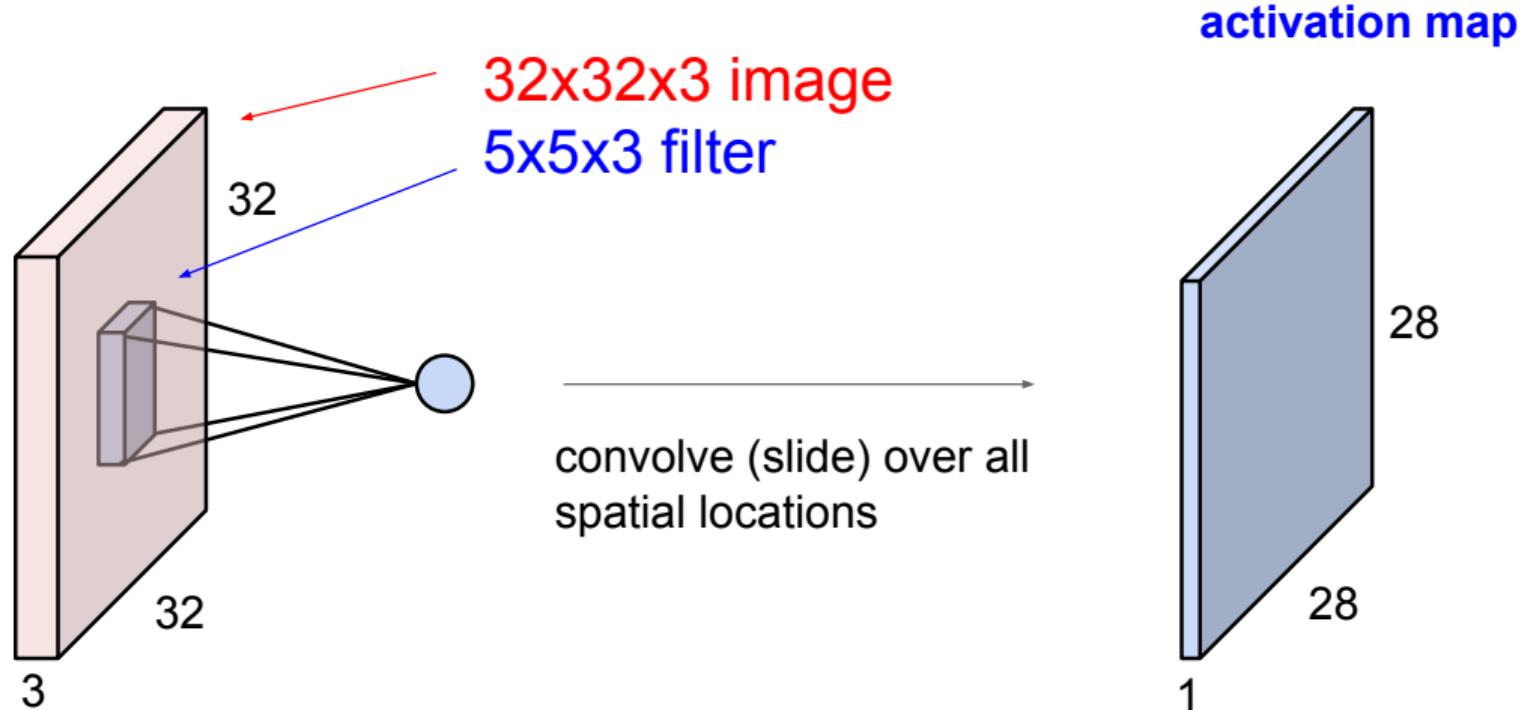
Filters always extend the full depth of the input volume

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

Convolution Layer

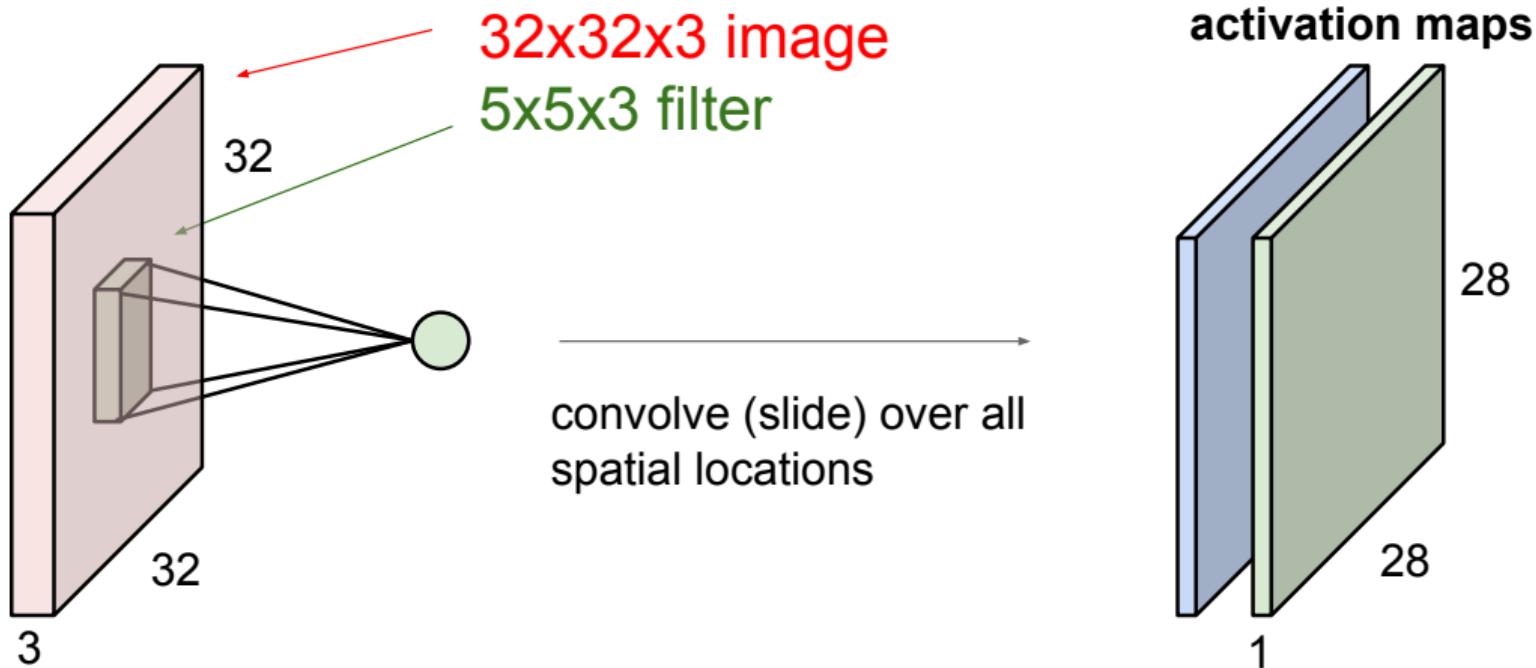


Convolution Layer

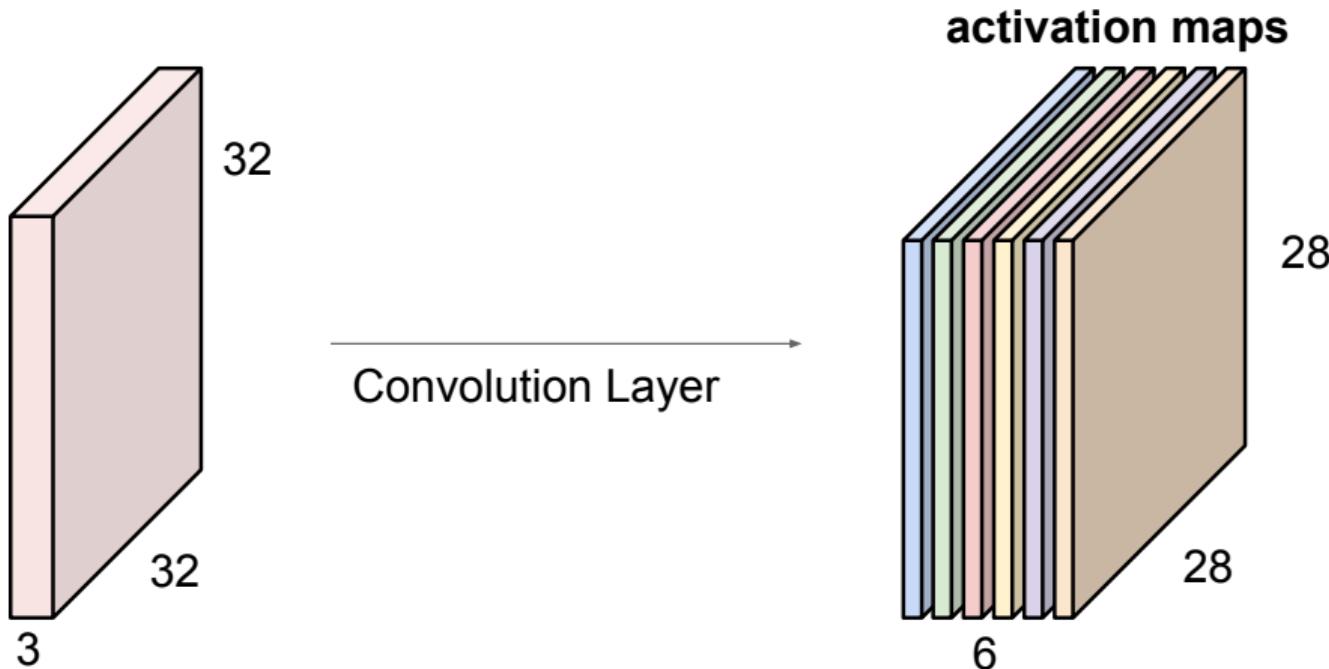


Convolution Layer

consider a second, green filter

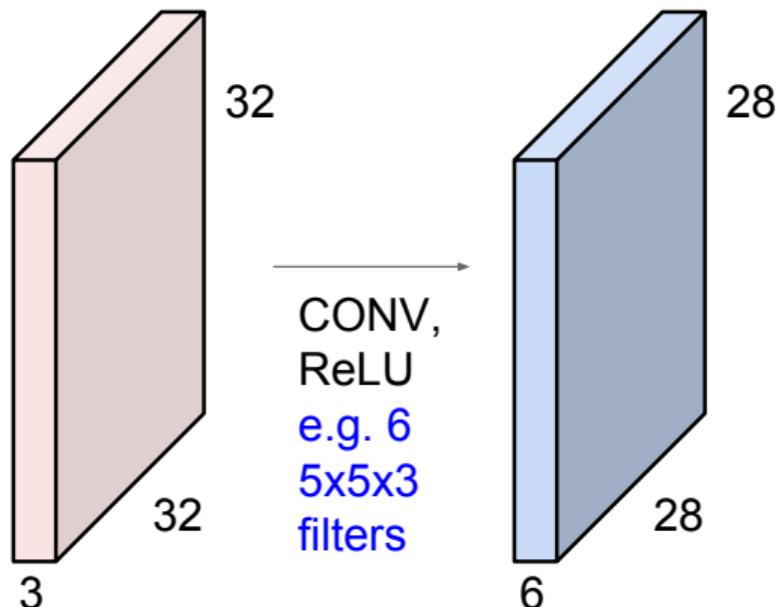


For example, if we had 6 5×5 filters, we'll get 6 separate activation maps:

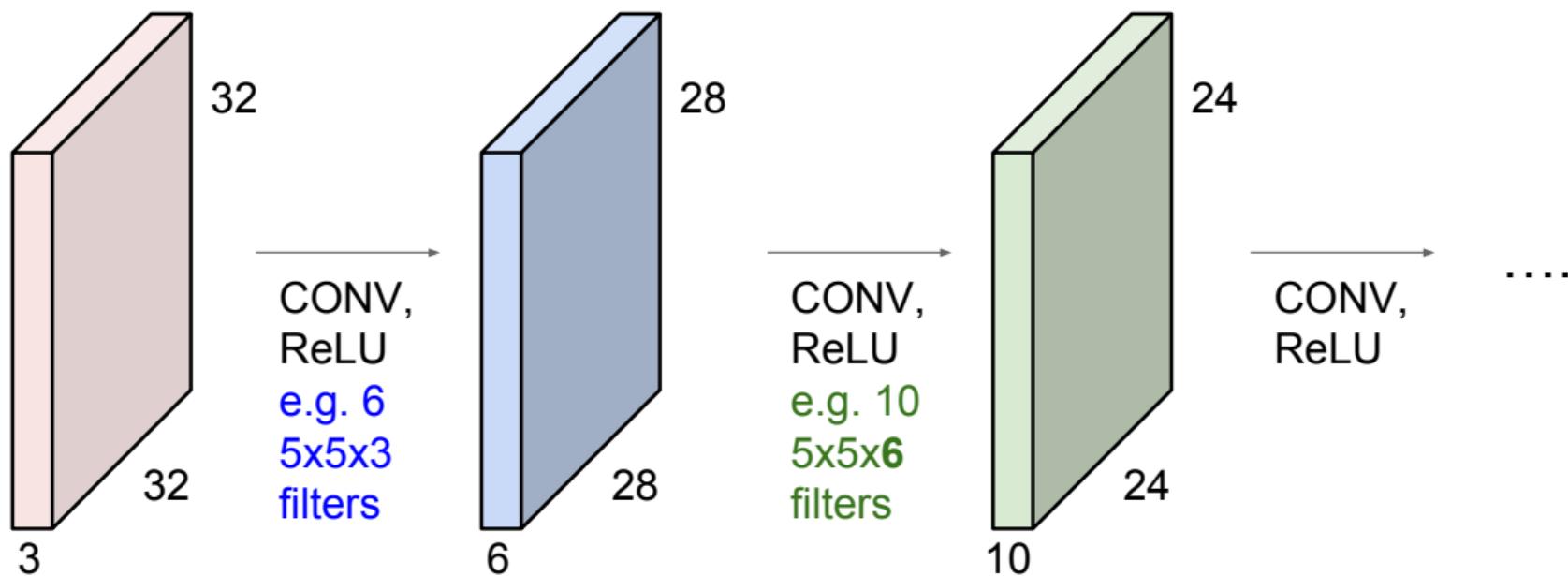


We stack these up to get a “new image” of size $28 \times 28 \times 6$!

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

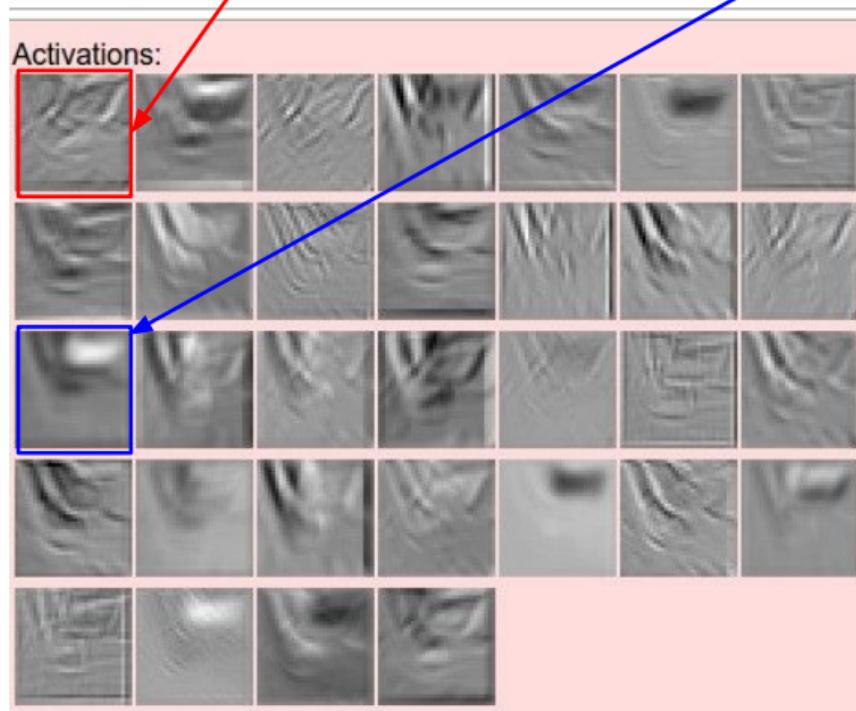


Preview: ConvNet is a sequence of Convolutional Layers, interspersed with activation functions





one filter =>
one activation map



example 5x5 filters
(32 total)

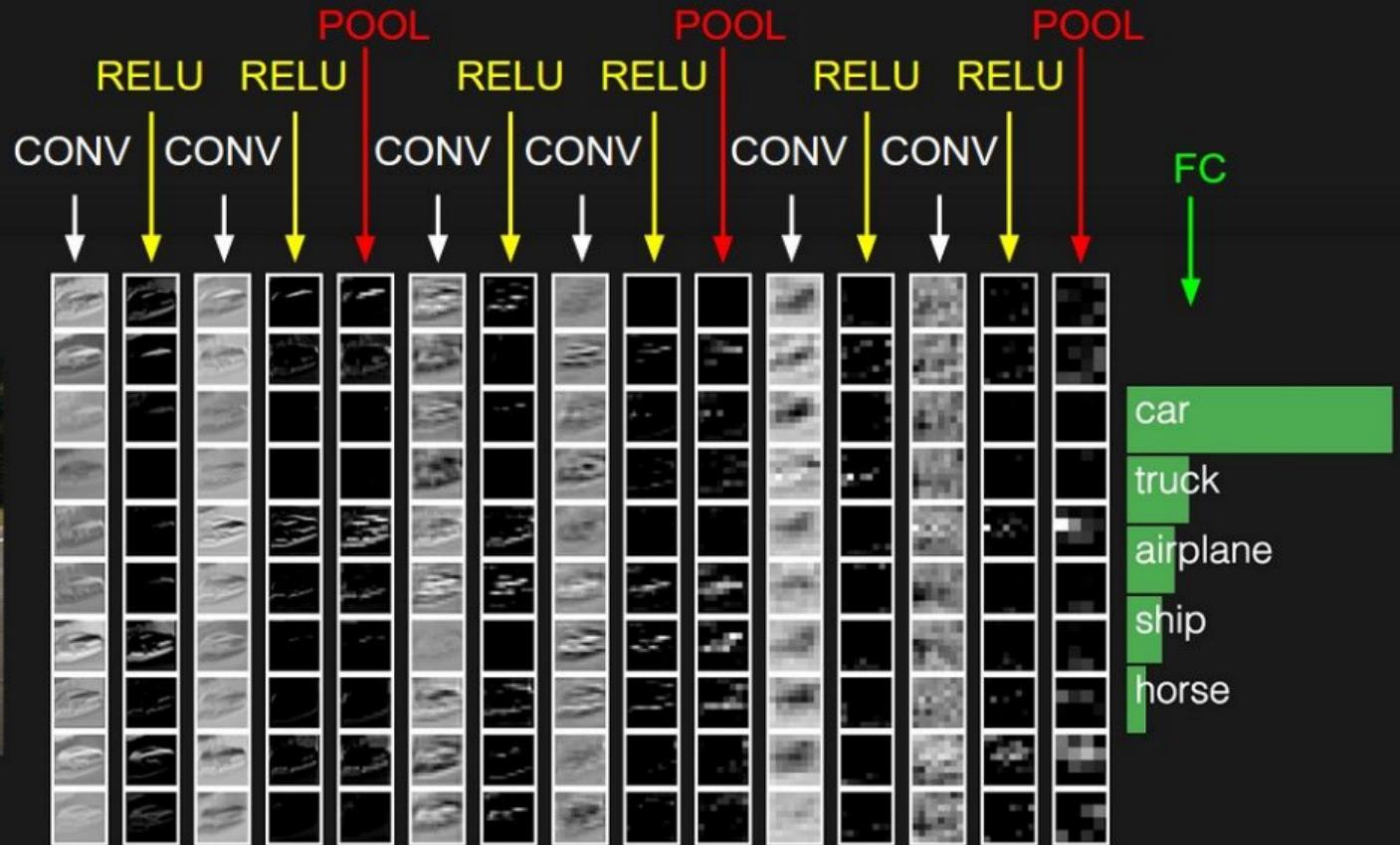
We call the layer convolutional
because it is related to convolution
of two signals:

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

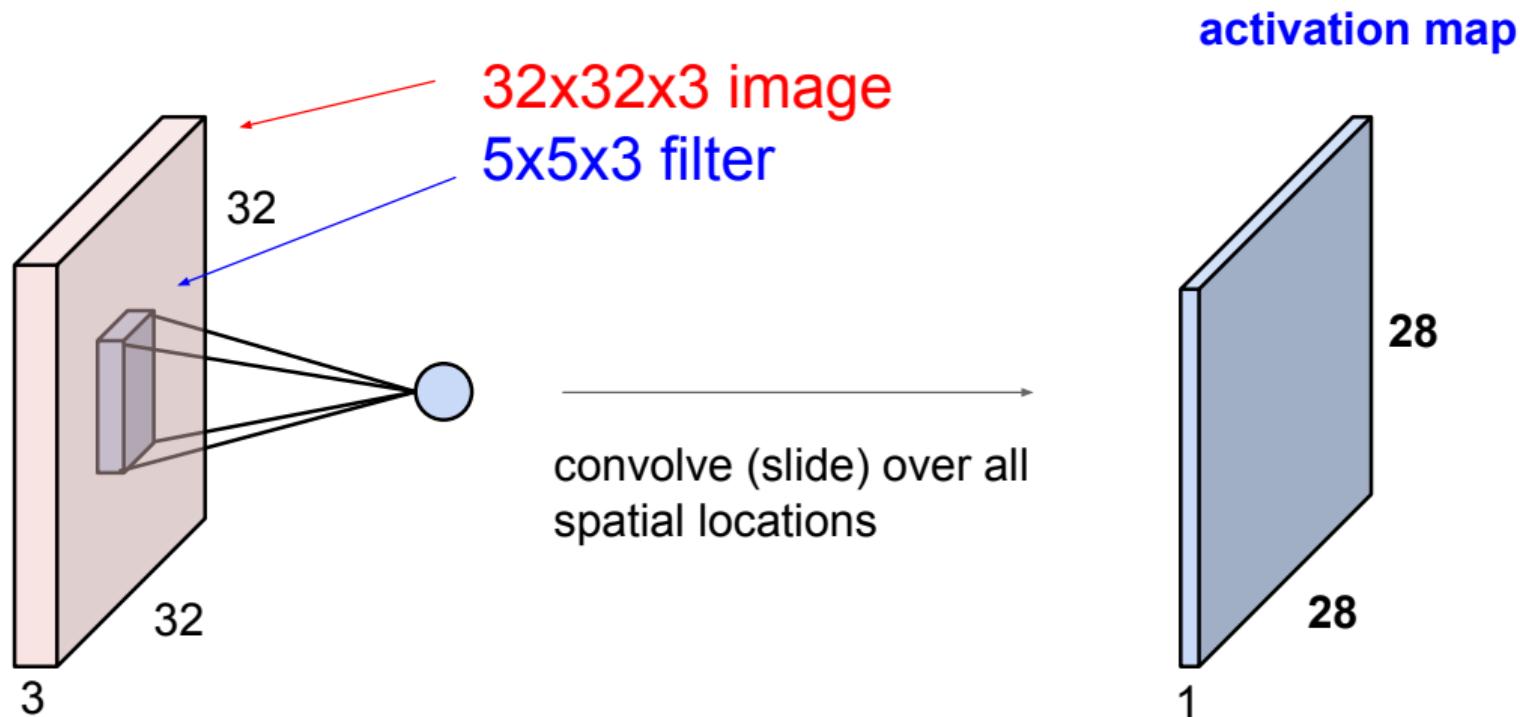


elementwise multiplication and sum of
a filter and the signal (image)

preview:

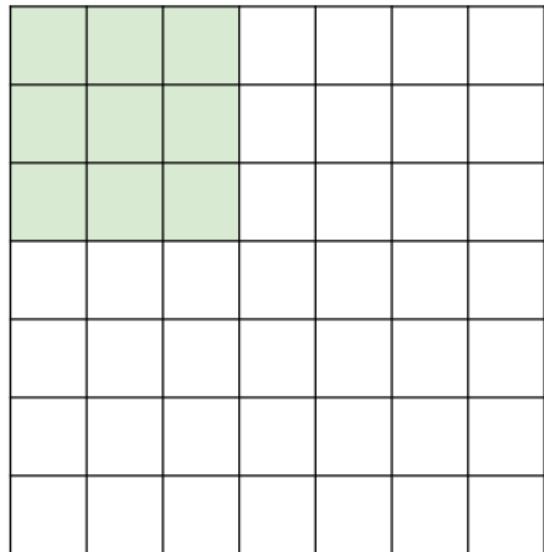


A closer look at spatial dimensions:



A closer look at spatial dimensions:

7

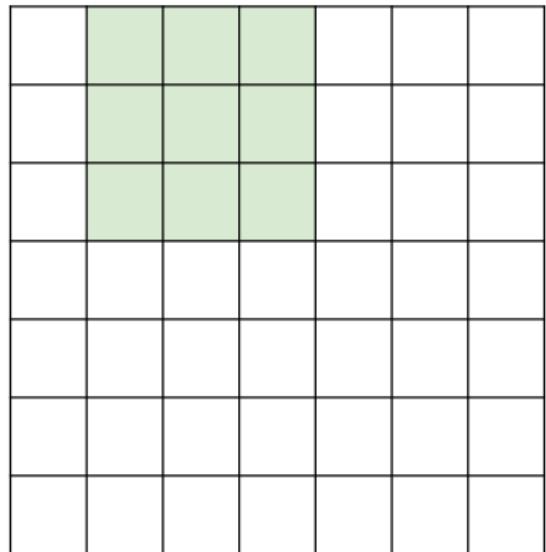


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

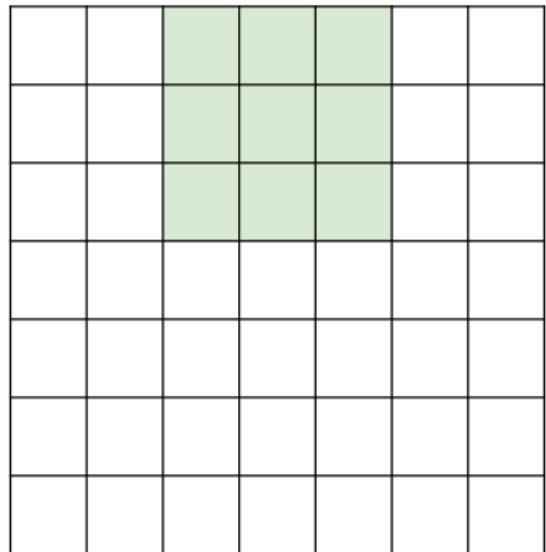


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

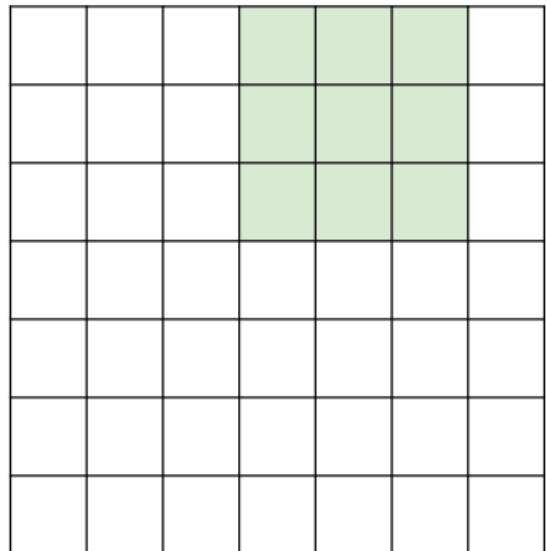


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

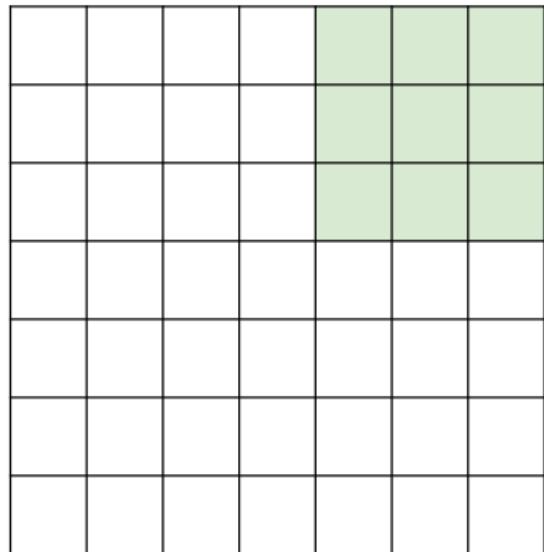


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7



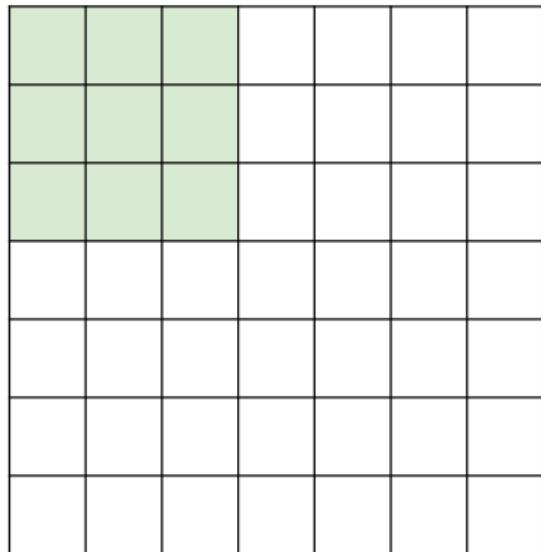
7x7 input (spatially)
assume 3x3 filter

=> **5x5 output**

7

A closer look at spatial dimensions:

7

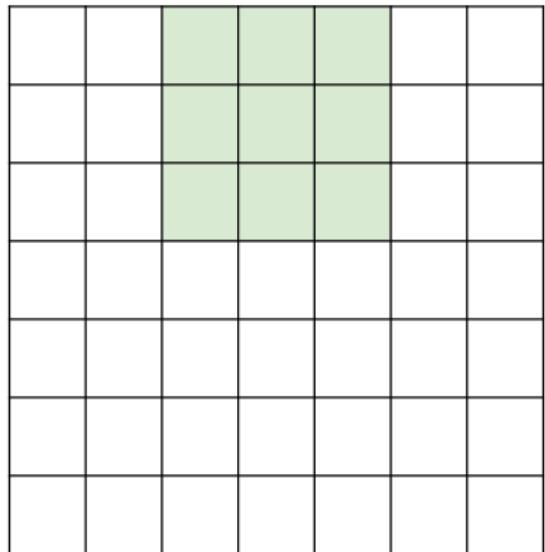


7

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

A closer look at spatial dimensions:

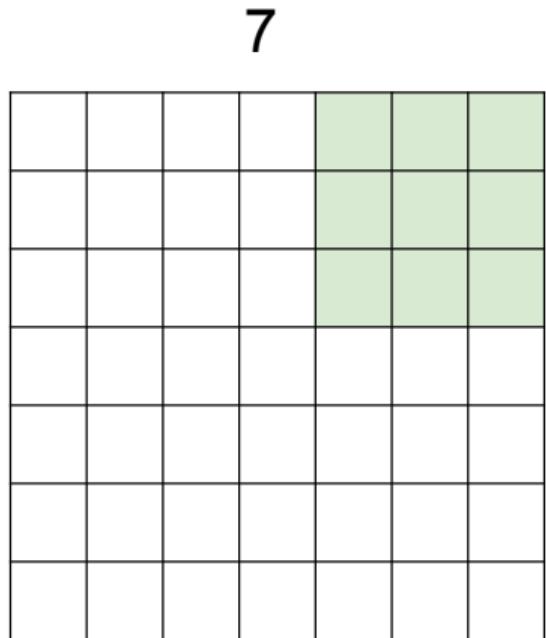
7



7

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

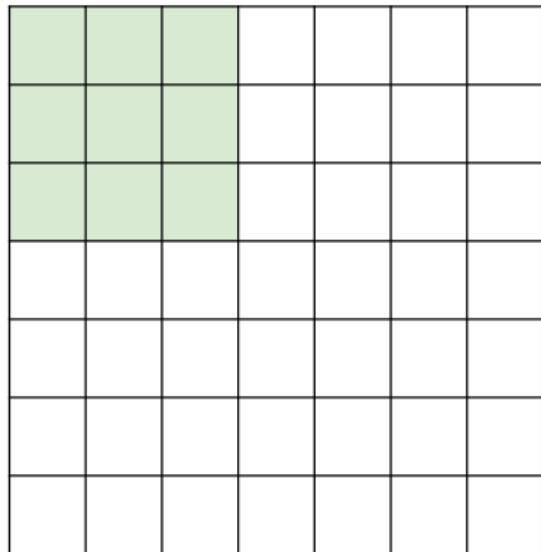
A closer look at spatial dimensions:



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

A closer look at spatial dimensions:

7

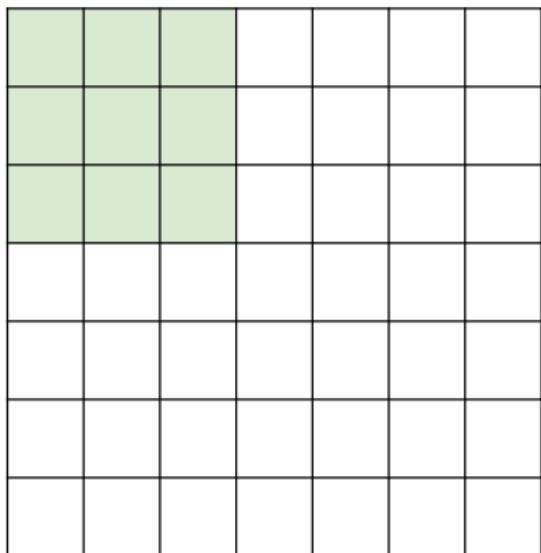


7

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

A closer look at spatial dimensions:

7

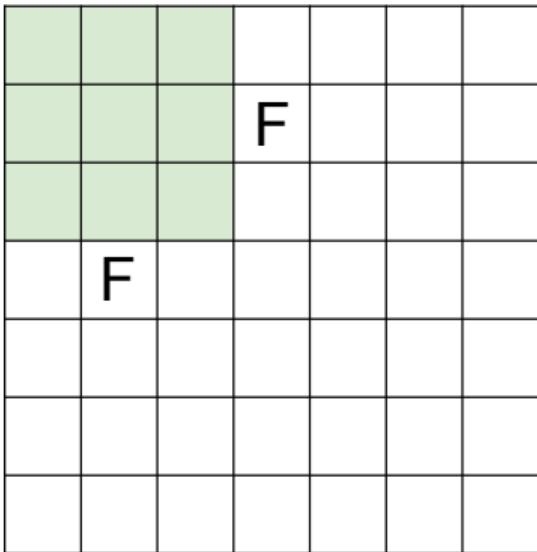


7

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

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

N



N

Output size:
(N - F) / 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)

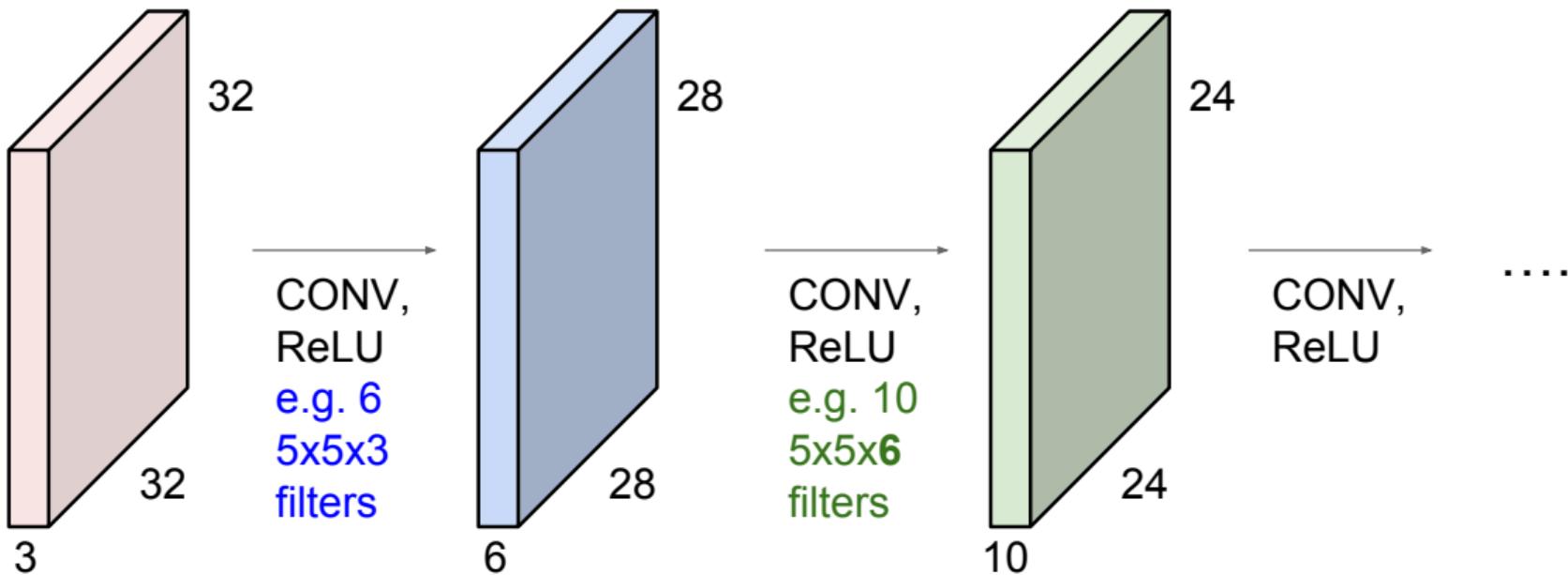
e.g. $F = 3 \Rightarrow$ zero pad with 1

$F = 5 \Rightarrow$ 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 \rightarrow 28 \rightarrow 24 ...). Shrinking too fast is not good, doesn't work well.

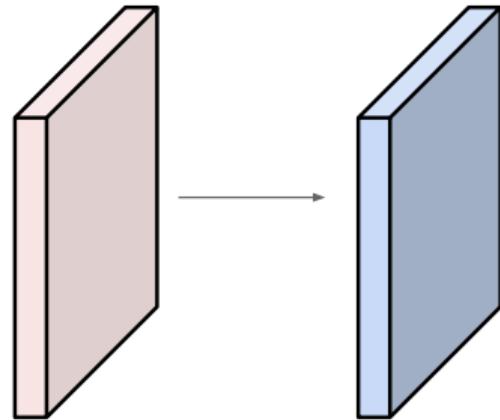


Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

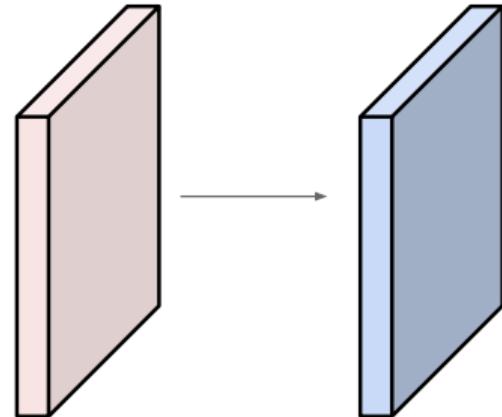
Output volume size: ?



Examples time:

Input volume: **32x32x3**

10 **5x5** filters with stride **1**, pad **2**



Output volume size:

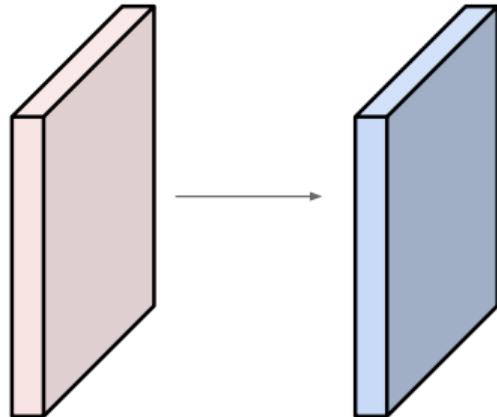
$(32+2*2-5)/1+1 = 32$ spatially, so

32x32x10

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

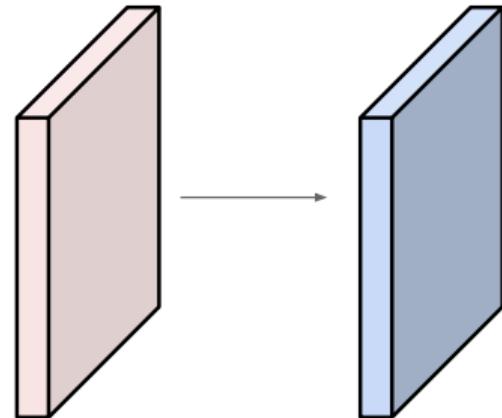


Number of parameters in this layer?

Examples time:

Input volume: **32x32x3**

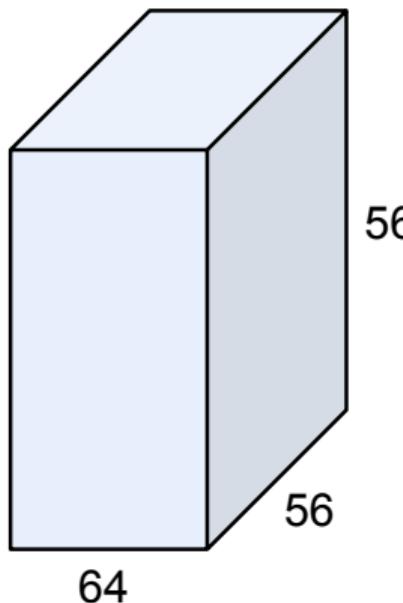
10 **5x5** filters with stride 1, pad 2



Number of parameters in this layer?

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

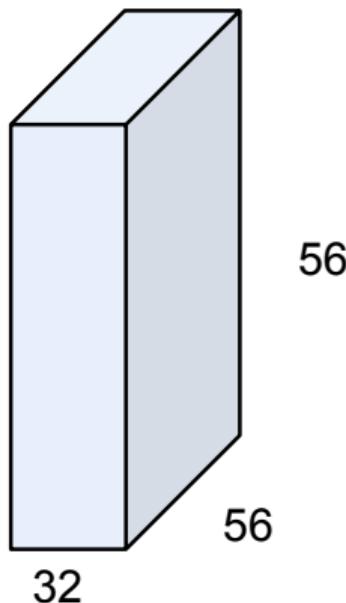
(btw, 1x1 convolution layers make perfect sense)



1x1 CONV
with 32 filters

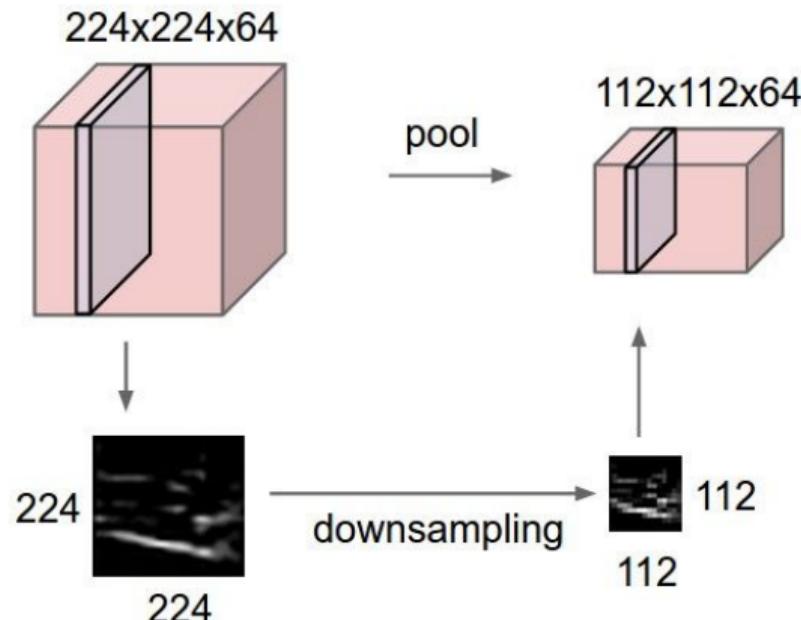
→

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

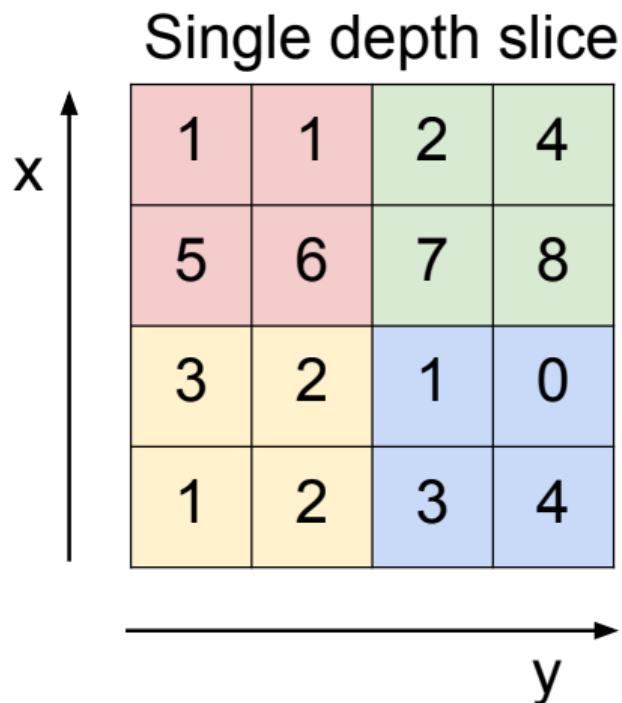


Pooling layer

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



MAX POOLING

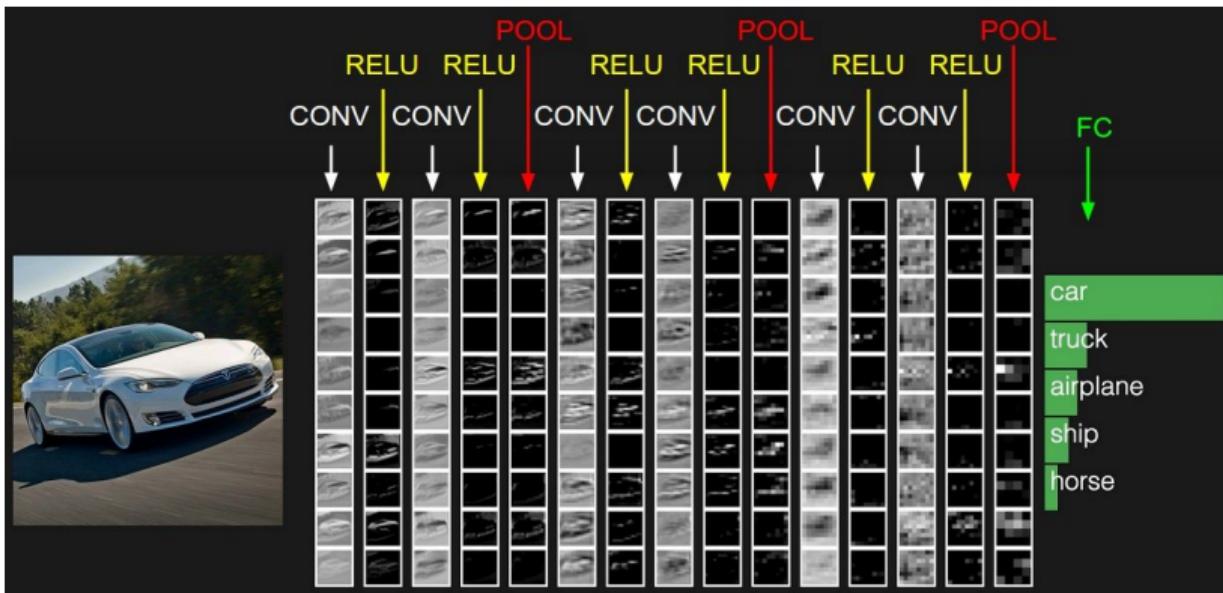


max pool with 2x2 filters
and stride 2

6	8
3	4

Fully Connected Layer (FC layer)

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



DROPOUT

- During training:
- For each postactivation h_i , with probability p let $h_i = 0$
- Redundancy
- Equivalent to learning an ensemble of networks

Improving neural networks by preventing co-adaptation of feature detectors;
Hinton, Srivastava, Krizhevsky, Sutskever, Salakhutdinov; (2012); arXiv:1207.0580

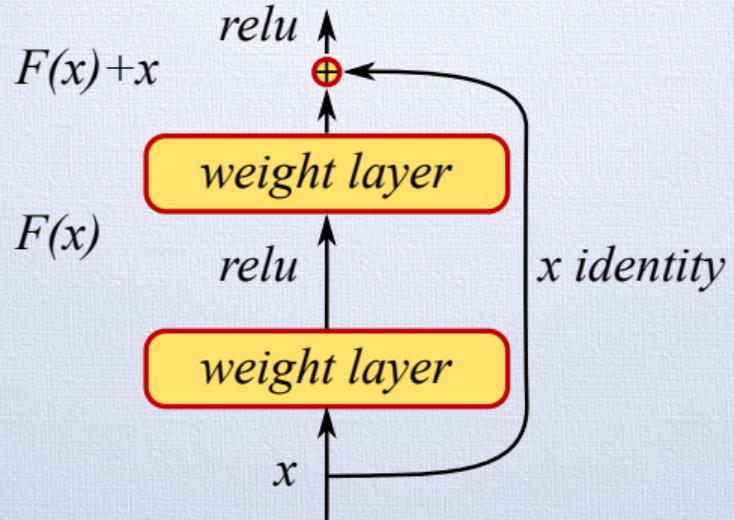
more on regularization

BATCH NORMALIZATION

- For each batch
- Normalize inputs to every layer to zero mean, unit variance.
- Helps with covariance shift

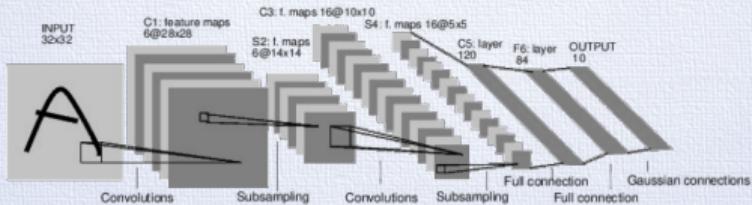
Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift; Ioffe, Szegedy; arXiv:1502.03167

RESIDUAL CONNECTIONS



Deep Residual Learning for Image Recognition; He, Zhang, Ren, Sun;
arXiv:1512.03385

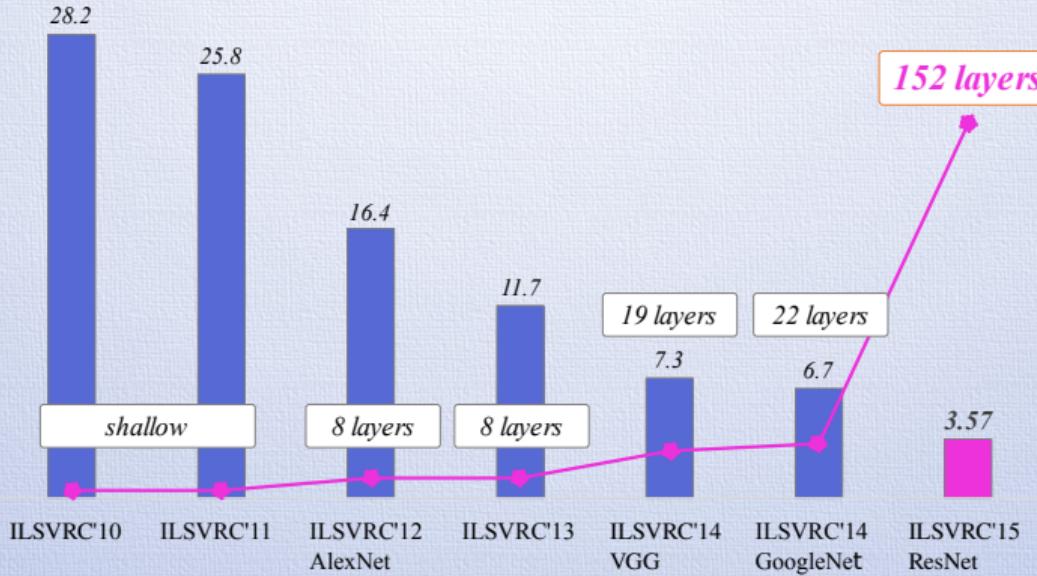
DEEPER AND DEEPER



[LeNet-5, LeCun 1980]

- 1998: LeNet-5; 3 layers
- 2012: AlexNet; 8 layers
- 2014: GoogLeNet; 22 layers (illustration)
- 2015: Residual Nets; 152 layers
- “Surpassed” human performance in 2015

DEPTH DEVELOPMENT



ImageNet Classification top-5 error (%)

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

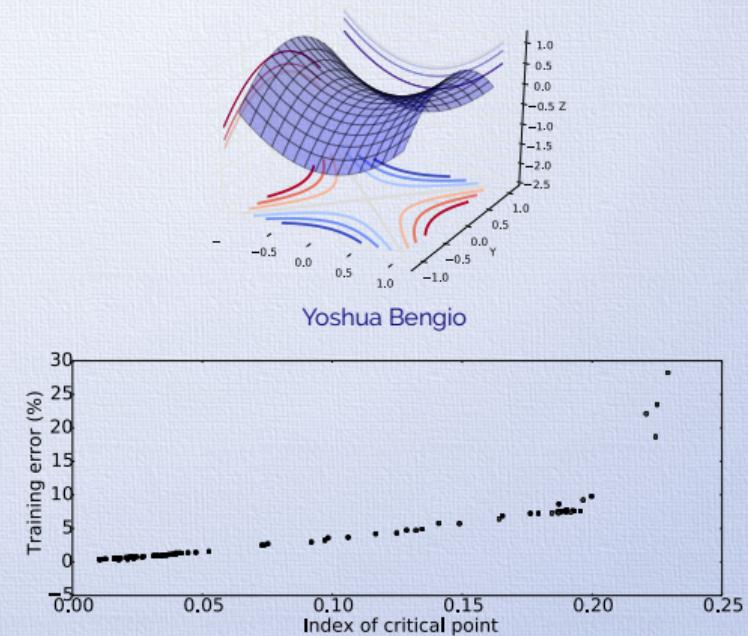
<http://mogren.one/>

NON-CONVEX OPTIMIZATION

- Loss function non-convex
- Low-D: **local minima** dominate
- High-D: **saddle points** dominate
- Local minima are close to global minimum
- Convexity not needed

The loss surfaces of multilayer networks;
Choromanska, et.al.; AISTATS 2015

Identifying and attacking the saddle point problem in high-dimensional non-convex optimization; Dauphin, et.al.; NIPS 2014

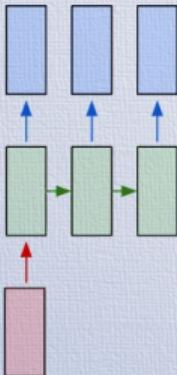


SEQUENCE MODELLING

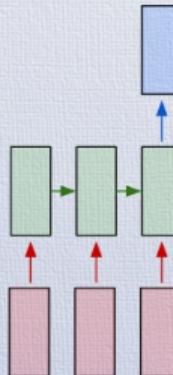
one to one



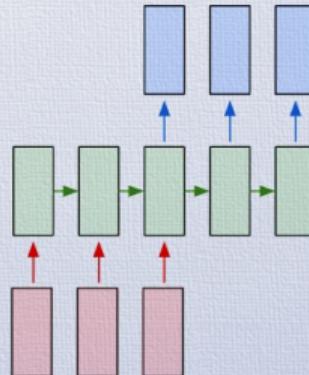
one to many



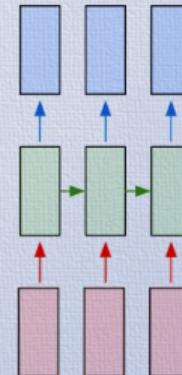
many to one



many to many

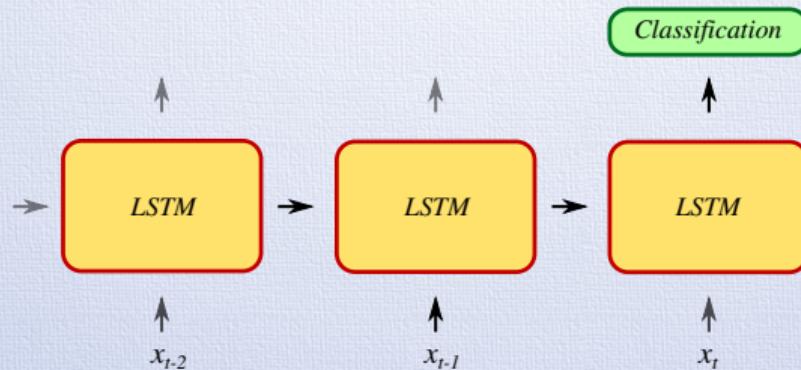


many to many



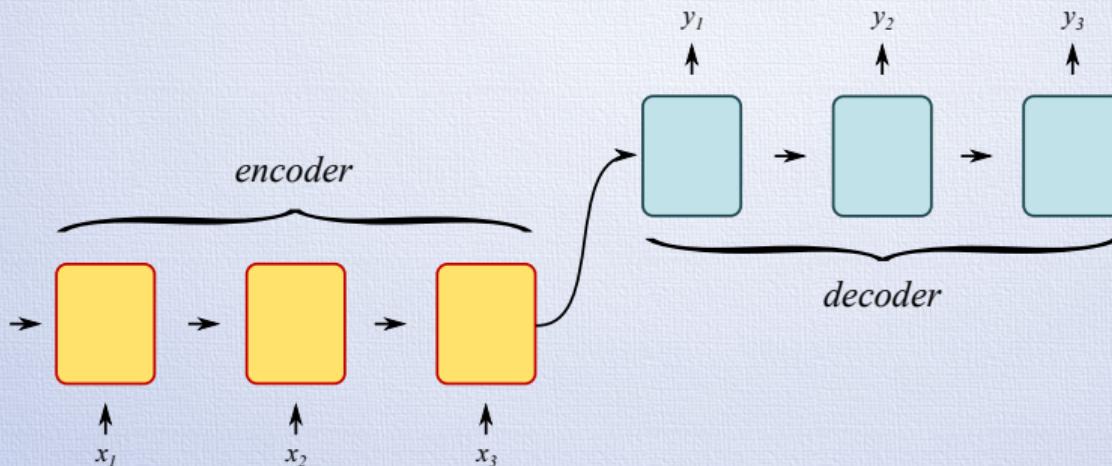
Andrej Karpathy
[details](#)

SENTIMENT ANALYSIS



- Binary sequence classification

NEURAL MACHINE TRANSLATION, NMT



Sequence to sequence learning with neural networks; Sutskever, Vinyals, Le; NIPS 2014

Neural machine translation by jointly learning to align and translate; Bahdanau, Cho, Bengio; ICLR 2015

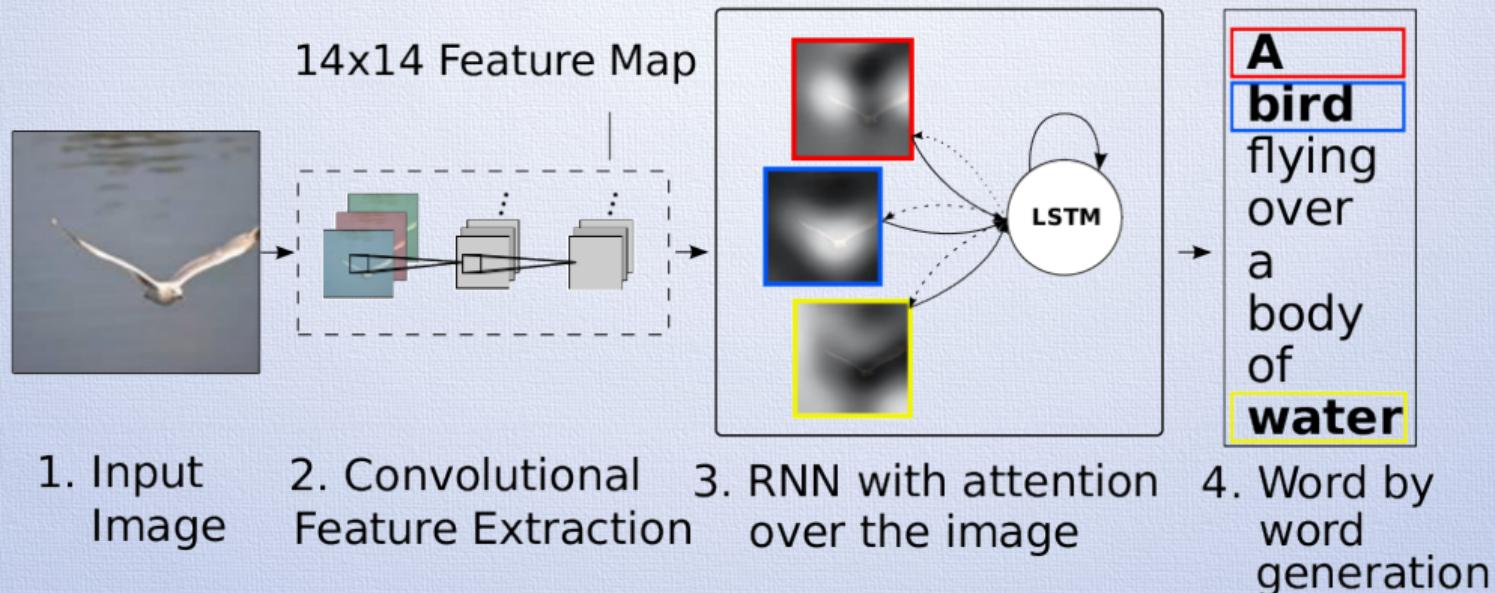
RECENT ADVANCES IN NMT

- Subwords (BPE) (Sennrich et.al., ACL 2016)
- 8 layers deep LSTM model.
- Quantized weights $\in \{-1, 0, +1\}$
- Downpour SGD: parallel training
- 8GPUs, one host.
- Human evaluation:
results comparable to human translators!



Google's neural machine translation system: Bridging the gap between human and machine translation; Yonghui Wu, et.al.; arXiv 1609.08144

CAPTION GENERATION



[more](#)

<http://mogren.one/>



<http://mogren.one/>

APPENDIX

WEIGHT DECAY (L2)

L2-Norm penalty term on weights when training.

$$\Omega(\theta) = \sum_k \sum_i \sum_j (W_{i,j}^{(k)})^2 = \sum_k \|W^{(k)}\|_F^2$$

- Applied on weights
- Can be interpreted as a Gaussian prior on the weights

back to dropout

WEIGHT DECAY (L1)

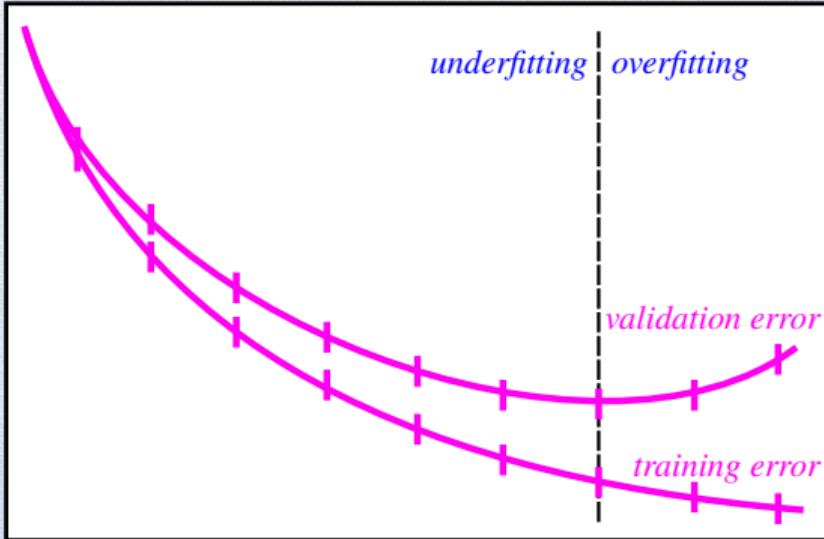
L1-Norm penalty term on weights when training.

$$\Omega(\theta) = \sum_k \sum_i \sum_j |W_{i,j}^{(k)}|$$

- Will force that some weights be exactly zero (sparsity).

back to dropout

EARLY STOPPING



- Keep track of validation error during training.
- Save the model at intervals.
- When validation error goes up, roll back network.

[back to dropout](#)

<http://mogren.one/>



<http://mogren.one/>

LEARNING TO PLAY ATARI BREAK-OUT



[online](#) [offline](#) [back to rl](#)

<http://mogren.one/>

DEEP REINFORCEMENT LEARNING

- Learning a policy
- *Infrequent* reward signal
- Deep Q-Learning:
Approximating the action-value function
- Atari games.
- Alpha Go
- Autonomous driving





<http://mogren.one/>

ATTENTION VISUALIZATION



A woman is throwing a frisbee in a park.

[back to caption introduction](#)



<http://mogren.one/>

MEMORY NETWORKS

- Attention refers back to internal memory; state of encoder
- Neural Turing Machines
- (End-To-End) Memory Networks:
explicit memory mechanisms
(out of scope today)

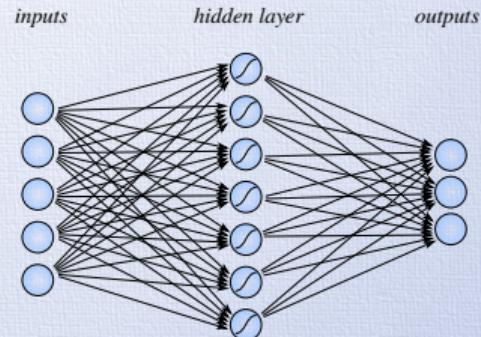
back



<http://mogren.one/>

LEARNING

- Forward pass (function application(s))
- Compute error for output
- Compute gradients (backpropagation)
derivative of stacked layers: chain rule
- Update weights (a small step)
(minibatch stochastic gradient descent)

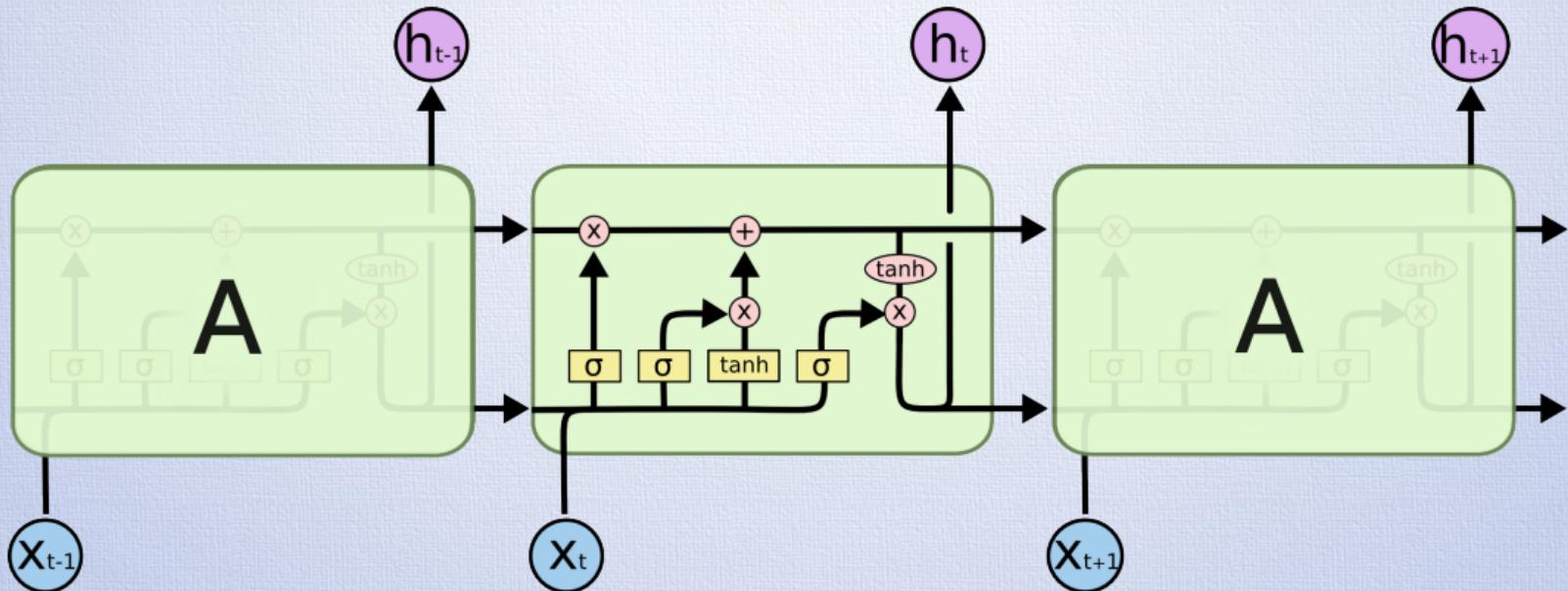


back to modelling



<http://mogren.one/>

LSTM

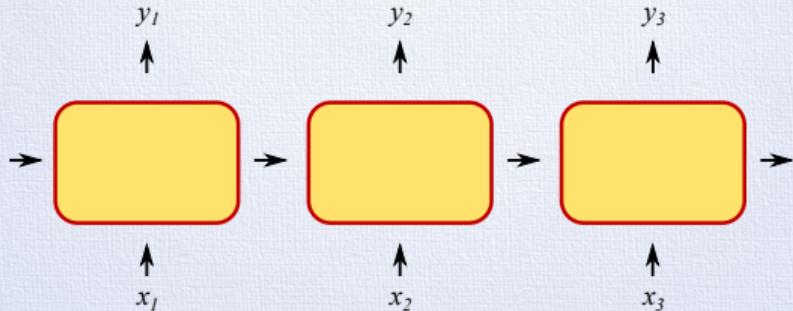


Christopher Olah

back to rnn

<http://mogren.one/>

MODELLING LANGUAGE USING RNNs



- Language models: $P(\text{word}_i | \text{word}_1, \dots, \text{word}_{i-1})$
- Recurrent Neural Networks
- “Long Short-Term Memory” (LSTM)
- Fixed vector representation for sequences
- Language generation (sampling; beam search)

[back to rnn](#) | [click for demo](#)

CHAR-RNNs

(Karpathy et.al. 2014)

- One LSTM cell per character
 - Word vocabulary independence (learn OOV terms)
 - Character vocabulary is small
-
- Learns to write like Shakespeare
 - Or like linux programmers
 - Mostly correct, (brackets opened/closed, indentation looks good)
 - Or like Wikipedia markup
 - **Or the text-files on my harddrive**

[back to rnn](#)

WITH A LITTLE TRAINING

X. Z

rbndtr Jltcub aoes re) sneecietioon ine gugu soiitid
as) iidoyan tsstaprm f tss
aBeAur:.e, 3a2o0annsdnss 3adr firet nte
tte tetaneg,.1 Ene)13.<

oides nor:

60

5. fipnapn::5, UdeirBseess-ed'iesr.:
ngsever, tciph:dsctum
se ralr msute 5emriipee vmetit Iooms m:eieeaaoeedes?
niei?nr rieanoitt 1 Dsehane
epd, ioe tee, onemseLsl

WITH SOME MORE TRAINING

u5ual ove repant org meucien pashe.. / ohe the man.
of namiof in loine
are proak ois unae if doubou u rew}

tr L is Yave tudtisat (theaon ir olcilar
it in, dil riRer nanes and is prome urdic unpu cunre
inh Is surgend the fcare celaoao ale Larcawanbusio
= 1 maxvand bis / o -pseire
counmt bhen-a }
etitwl ind tuvabToin maMatias otyraoss epronk

[inn os on:

The alulitat o/ pesede thav.

<http://mogren.one/>

ci Vill has 1 neighbors.

Graph has 2 vertard. Number on coykut: 2583Moher.

Number }411185n169V71}305911 sdmgap Kart:

CL Volvo Penta Volvo Penta in you hax on Volvo

TrucladionTenkoy i meance can 06:715g333139 >-C vexter:

unbe Bactizaces processed deviewroced an than quardy

processed before: 150

erreight rasunce qual. Number stam ald Boat.

Decided on. 55012418571607Meag.b>mlwyt arluto

aufortion!

Dlidlis Bont!

- C2 ((Rekroce: Prepors: Frapers. Number much: and
yimile have 2 vertices anbuct cont87a

b9 N1khwear has 1 neighbors.

Graph has 6 neighbors.

Graph has 9 vertices and 699 edges.

This on centroeds basbanam has no soluth. Number of centroids proclesc. Number of centroids processed before: 235

rearnt has 1 neighbors.

Graph has 4 vertices and 7 edges.

This instance has no solution!

Decided on centroid: aent. Nerticed least has no soam too strast. Number of too small (size: 5pgabl. Number of centroids processed resuase Tas instance cast. Number of centroids processed before: 342 daser ha



<http://mogren.one/>