



CS489/698: Intro to ML

Lecture 12: Training Neural Networks and
CNNs part 1

focal
Systems



UNIVERSITY OF
WATERLOO

Outline

- Training
- Convolutional Layers



Outline

- Training
- Convolutional Layers



SGD Formalized

Algorithm 8.1 Stochastic gradient descent (SGD) update at training iteration k

Require: Learning rate ϵ_k .

Require: Initial parameter θ

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

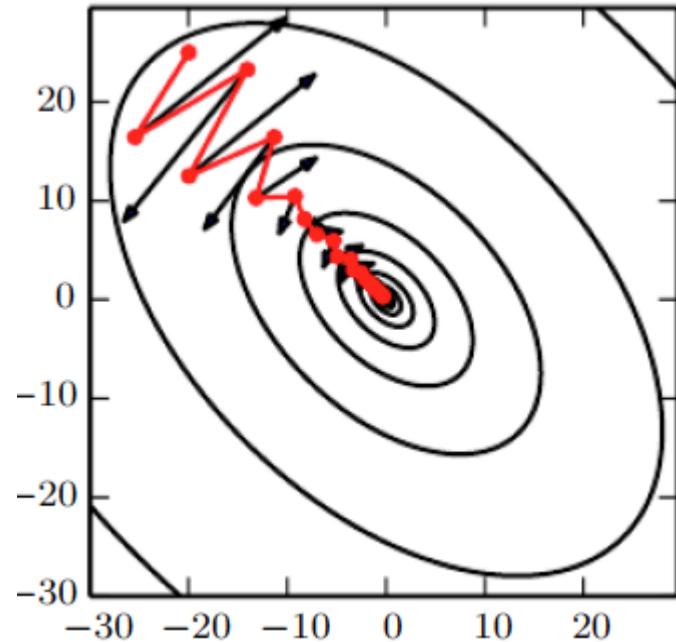
 Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$

end while

Momentum

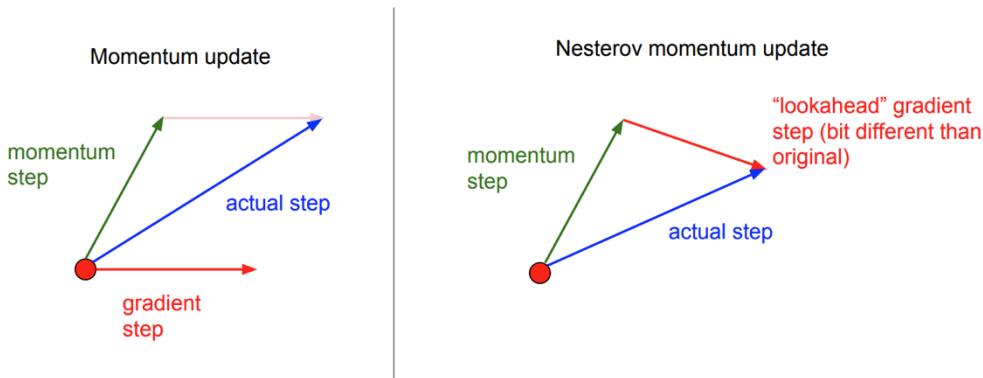
- Adds a velocity term
- $v \leftarrow \alpha v - \epsilon \nabla_{\theta} \left(\frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta), y^{(i)}) \right),$
 $\theta \leftarrow \theta + v.$
- Adds a speedup of at most $\frac{\epsilon \|g\|}{1 - \alpha}$.
where $0 < m < 1$ is a new global parameter which must be determined by trial and error.
Momentum simply adds a fraction m of the previous weight update to the current one.
When the gradient keeps pointing in the same direction, this will increase the size of the steps taken towards the minimum. It is therefore often necessary to reduce the global learning rate μ when using a lot of momentum (m close to 1). If you combine a high learning rate with a lot of momentum, you will rush past the minimum with huge steps
- Momentum constant usually 0.9, 0.5, 0.99 corresponding to a 10x, 2x, 100x increase in max speed



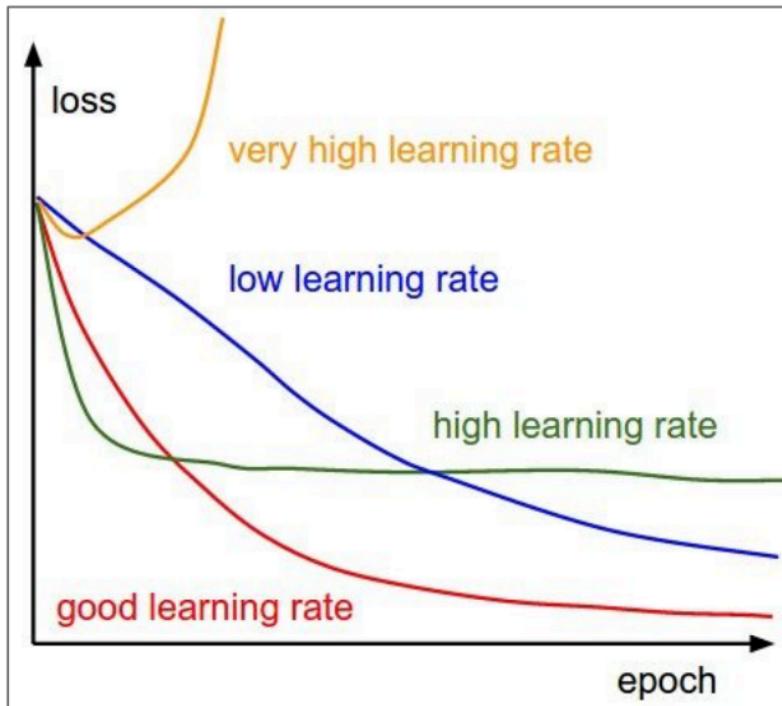
Nesterov Momentum

- Applies momentum first then gradient

$$\begin{aligned}\mathbf{v} &\leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\theta} \left(\frac{1}{m} \sum_{i=1}^m L(\mathbf{f}(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)}) \right), \\ \theta &\leftarrow \theta + \mathbf{v}.\end{aligned}$$



Setting the learning rate



=> Learning rate decay over time!

step decay:

e.g. decay learning rate by half every few epochs.

exponential decay:

$$\alpha = \alpha_0 e^{-kt}$$

1/t decay:

$$\alpha = \alpha_0 / (1 + kt)$$

- Also good to try a setting for 100 iterations and see which is best on validation set

focal
Systems



Adagrad

- Adapting the learning rate per parameter
- If gradient is always high, decreases learning rate
- If gradient is low, increases learning rate
- meant for fast convergence on convex problems
- suffers from premature decay

AdaGrad is an optimization method that allows different step sizes for different features. It increases the influence of rare but informative features

Adagrad

Algorithm 8.4 The AdaGrad algorithm

Require: Global learning rate ϵ

Require: Initial parameter θ

Require: Small constant δ , perhaps 10^{-7} , for numerical stability

Initialize gradient accumulation variable $r = \mathbf{0}$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Accumulate squared gradient: $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$

 Compute update: $\Delta\theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$. (Division and square root applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta\theta$

end while

RMSProp

- Replaces average with exponential decay
- Prevents premature decay

Algorithm 8.6 RMSProp algorithm with Nesterov momentum

Require: Global learning rate ϵ , decay rate ρ , momentum coefficient α .

Require: Initial parameter θ , initial velocity v .

Initialize accumulation variable $r = \mathbf{0}$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets $y^{(i)}$.

 Compute interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$

 Compute gradient: $g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$

 Accumulate gradient: $r \leftarrow \rho r + (1 - \rho)g \odot g$

 Compute velocity update: $v \leftarrow \alpha v - \frac{\epsilon}{\sqrt{r}} \odot g$. ($\frac{1}{\sqrt{r}}$ applied element-wise)

 Apply update: $\theta \leftarrow \theta + v$

end while

Adam - Adaptive Momentum

- Applies momentum with parameterized learning rates

Adaptive Moment Estimation (Adam) [15] is another method that computes adaptive learning rates for each parameter. In addition to storing an exponentially decaying average of past squared gradients v_t like Adadelta and RMSprop, Adam also keeps an exponentially decaying average of past gradients m_t , similar to momentum:

Algorithm 8.7 The Adam algorithm

Require: Step size ϵ (Suggested default: 0.001)
Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$.
 (Suggested defaults: 0.9 and 0.999 respectively)
Require: Small constant δ used for numerical stabilization. (Suggested default: 10^{-8})
Require: Initial parameters θ
Initialize 1st and 2nd moment variables $s = \mathbf{0}$, $r = \mathbf{0}$
Initialize time step $t = 0$
while stopping criterion not met **do**
 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.
 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
 $t \leftarrow t + 1$
 Update biased first moment estimate: $\hat{s} \leftarrow \rho_1 s + (1 - \rho_1) \mathbf{g}$
 Update biased second moment estimate: $\hat{r} \leftarrow \rho_2 r + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$
 Correct bias in first moment: $\hat{s} \leftarrow \frac{\hat{s}}{1 - \rho_1^t}$
 Correct bias in second moment: $\hat{r} \leftarrow \frac{\hat{r}}{1 - \rho_2^t}$
 Compute update: $\Delta\theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r}} + \delta}$ (operations applied element-wise)
 Apply update: $\theta \leftarrow \theta + \Delta\theta$
end while
Agastya Kalra

How to choose?

- SGD < SGD+Momentum < SGD+Nesterov Momentum
- Adam is a good default
- RMSProp is good for RNNs, but also good default
- SGD + Nesterov momentum is best if you have time/resources to optimize learning rate
- More of an Art

focal
Systems



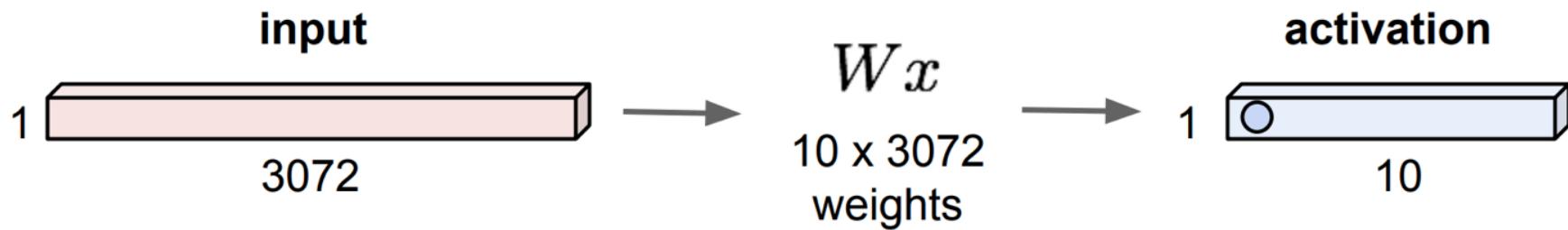
UNIVERSITY OF
WATERLOO

Outline

- Training
- Convolutional Layers

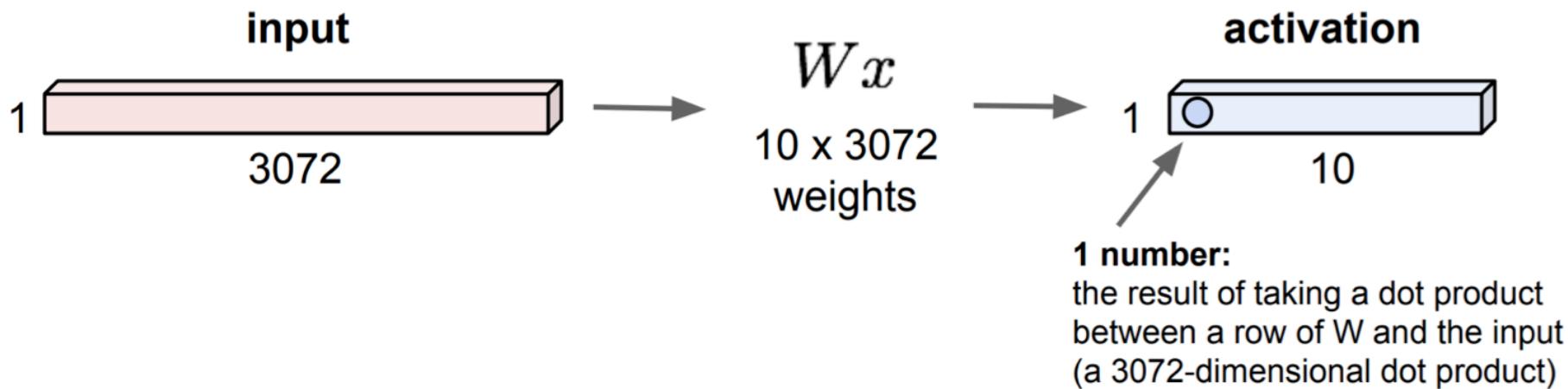
Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



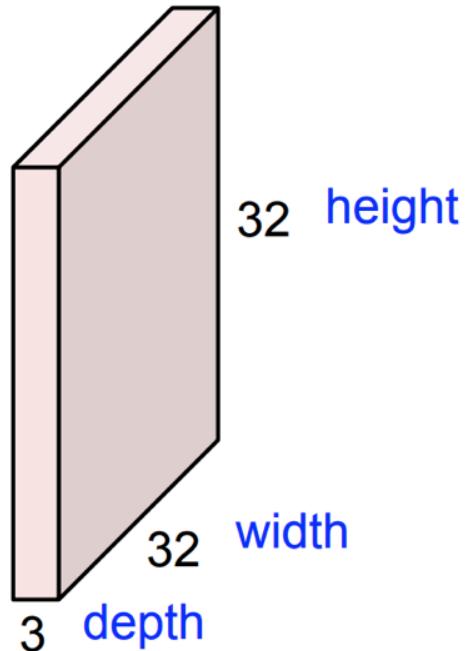
Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



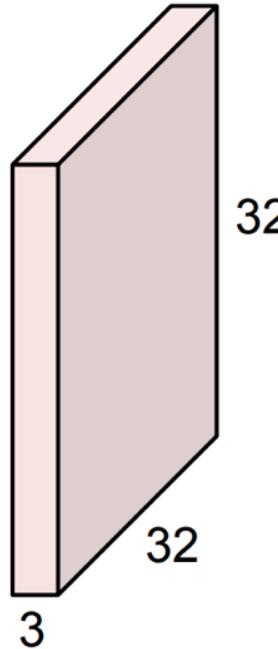
Convolutional Layer

32x32x3 image -> preserve spatial structure



Convolutional Layer

32x32x3 image



5x5x3 filter



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

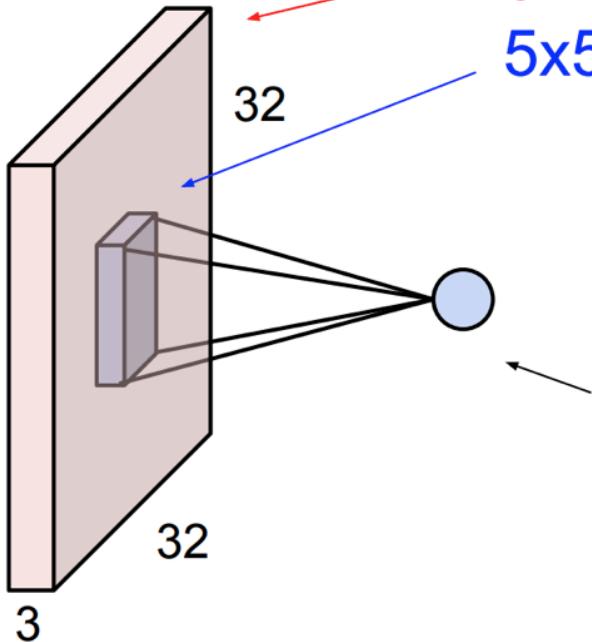
focal
Systems



UNIVERSITY OF
WATERLOO

Convolutional Layer

Convolution Layer



32x32x3 image
5x5x3 filter w

1 number:
the result of taking a dot product between the
filter and a small 5x5x3 chunk of the image
(i.e. $5 \times 5 \times 3 = 75$ -dimensional dot product + bias)

$$w^T x + b$$

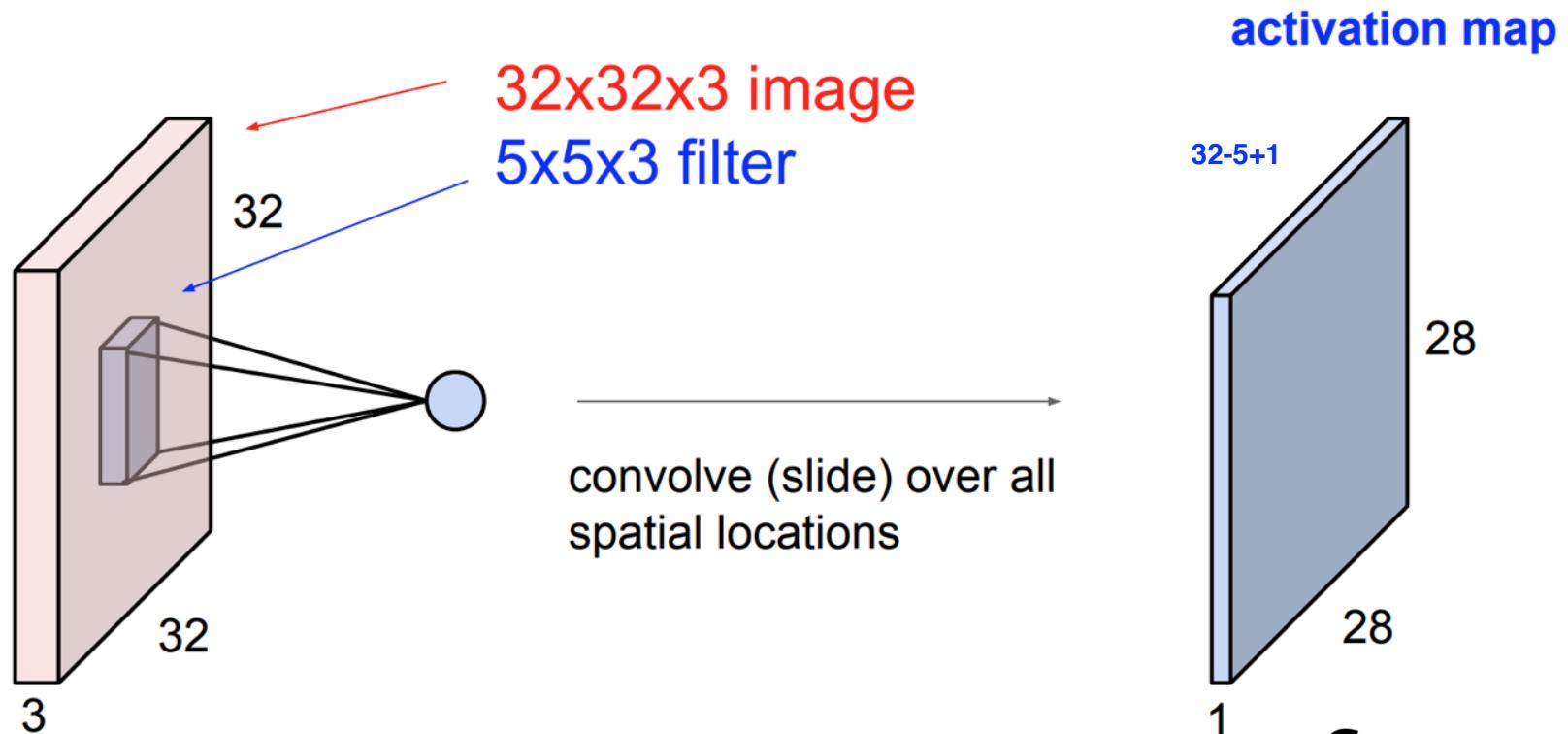
focal
Systems



UNIVERSITY OF
WATERLOO

Convolutional Layer

Convolution Layer



focal
Systems

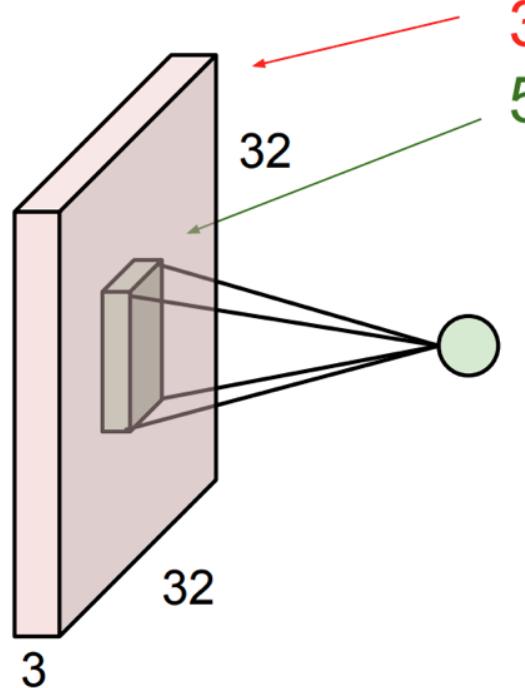


UNIVERSITY OF
WATERLOO

Convolutional Layer

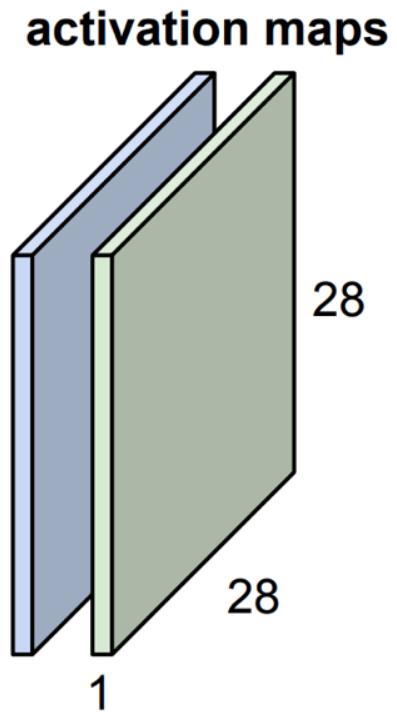
Convolution Layer

consider a second, green filter



32x32x3 image
5x5x3 filter

convolve (slide) over all
spatial locations



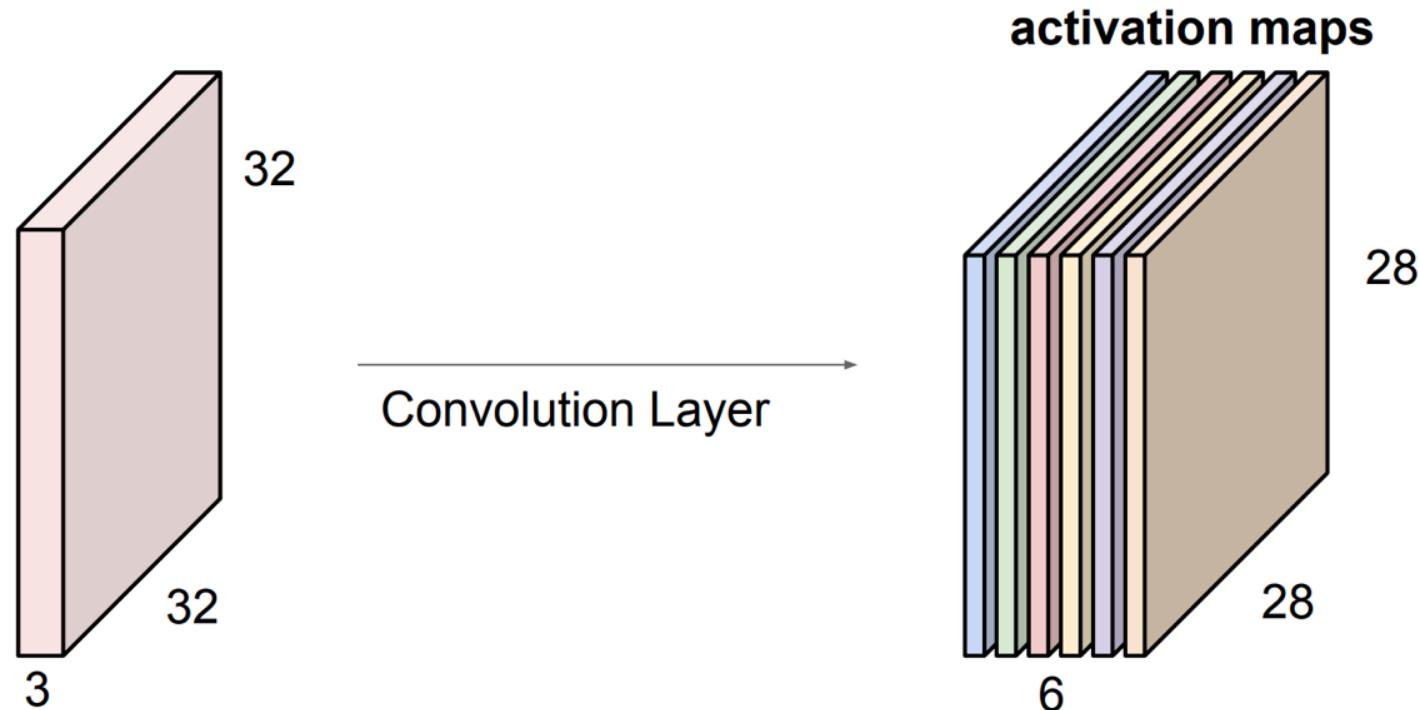
focal
Systems



UNIVERSITY OF
WATERLOO

Convolutional Layer

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!

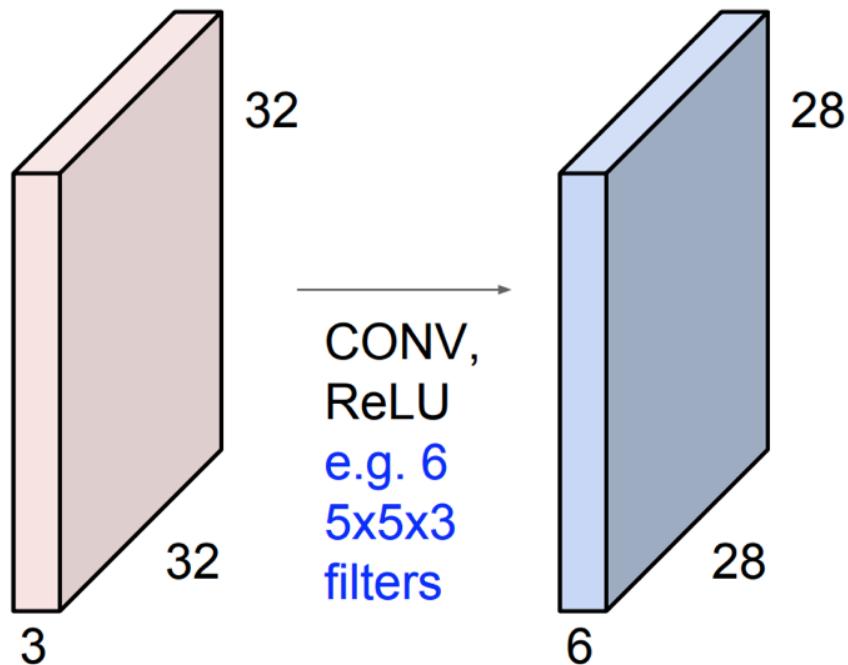
focal
Systems



UNIVERSITY OF
WATERLOO

Convolutional Layer

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



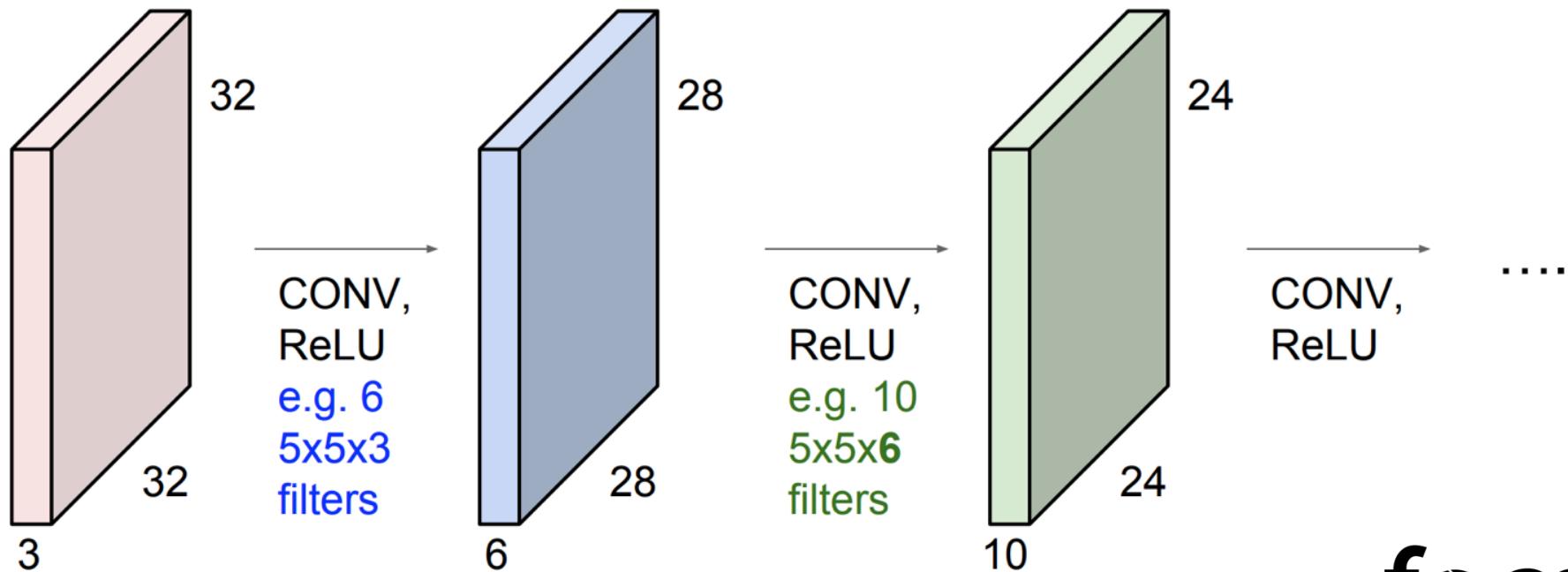
focal
Systems



UNIVERSITY OF
WATERLOO

Convolutional Layer

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



focal
Systems



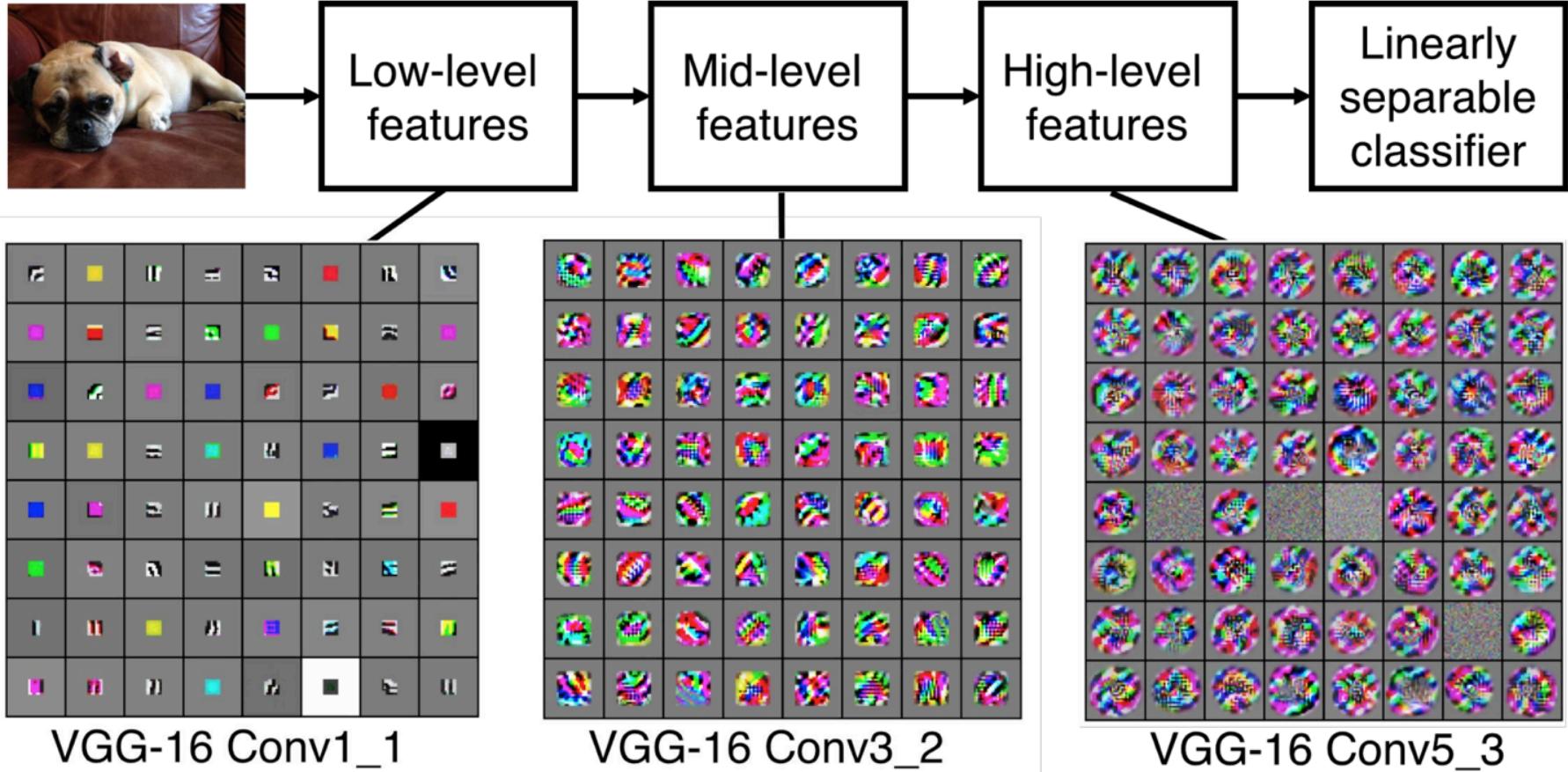
UNIVERSITY OF
WATERLOO

Convolutional Layer

Preview

[Zeiler and Fergus 2013]

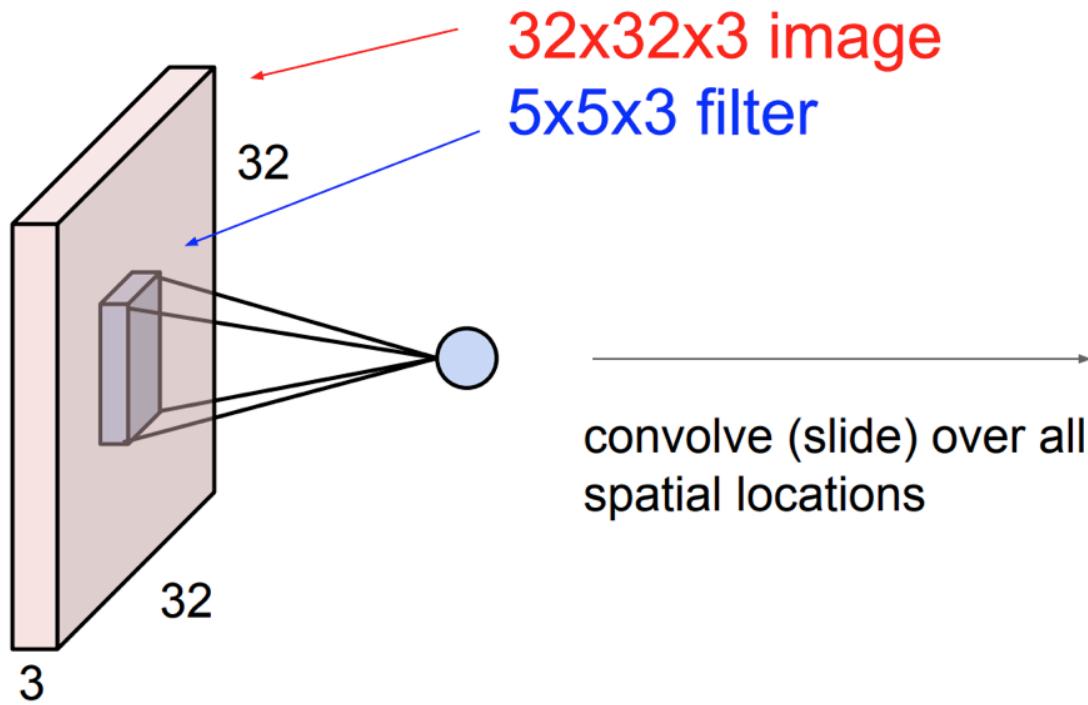
Visualization of VGG-16 by Lane McIntosh. VC architecture from [Simonyan and Zisserman 2014]



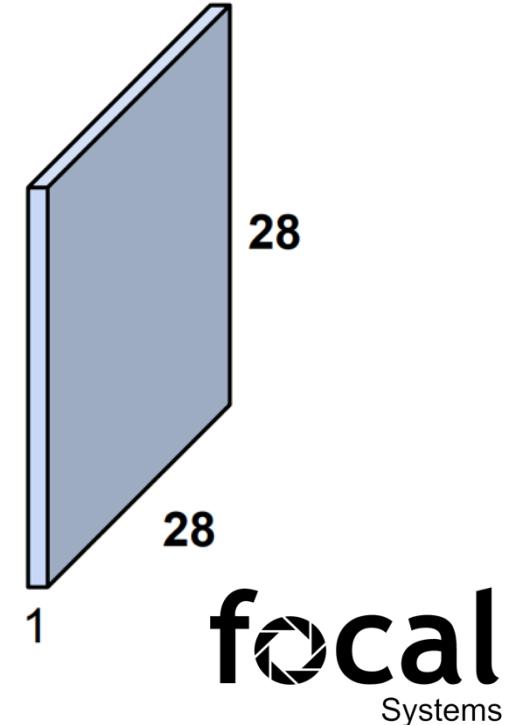
http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture5.pdf

Convolutional Layer

A closer look at spatial dimensions:



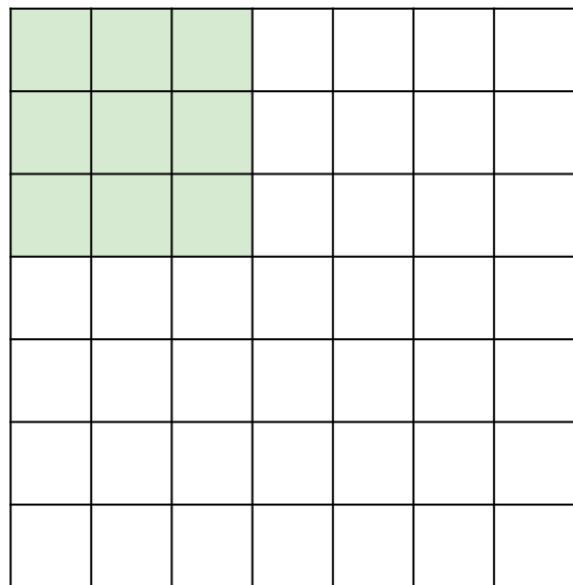
activation map



Convolutional Layer

A closer look at spatial dimensions:

7

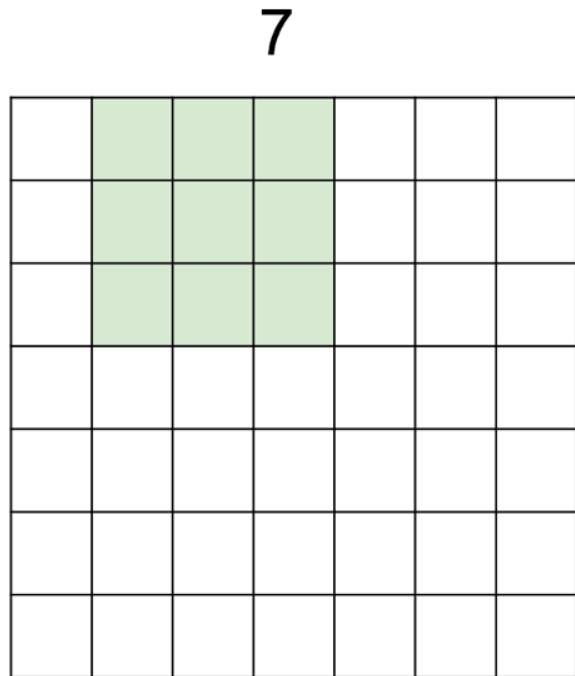


7x7 input (spatially)
assume 3x3 filter

7

Convolutional Layer

A closer look at spatial dimensions:

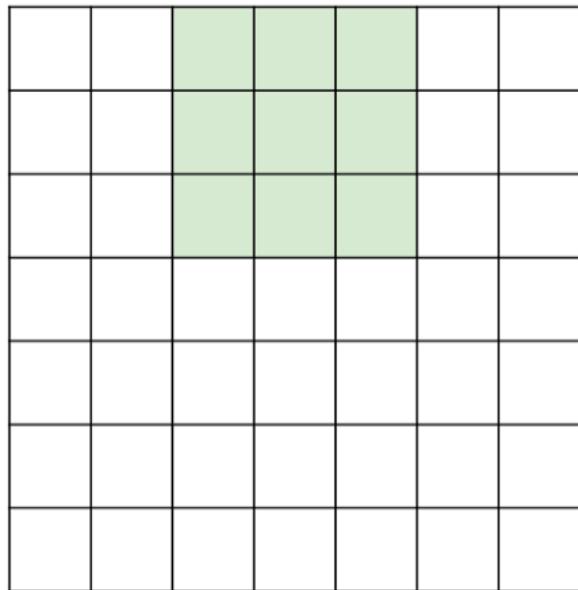


7x7 input (spatially)
assume 3x3 filter

Convolutional Layer

A closer look at spatial dimensions:

7

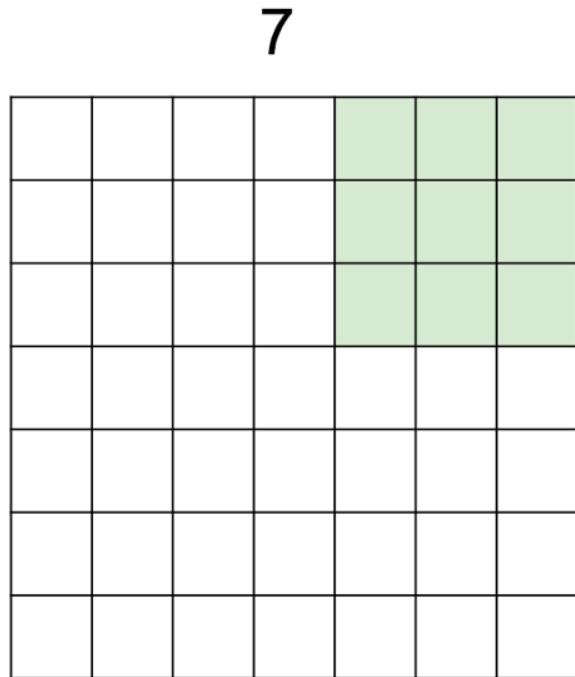


7x7 input (spatially)
assume 3x3 filter

7

Convolutional Layer

A closer look at spatial dimensions:

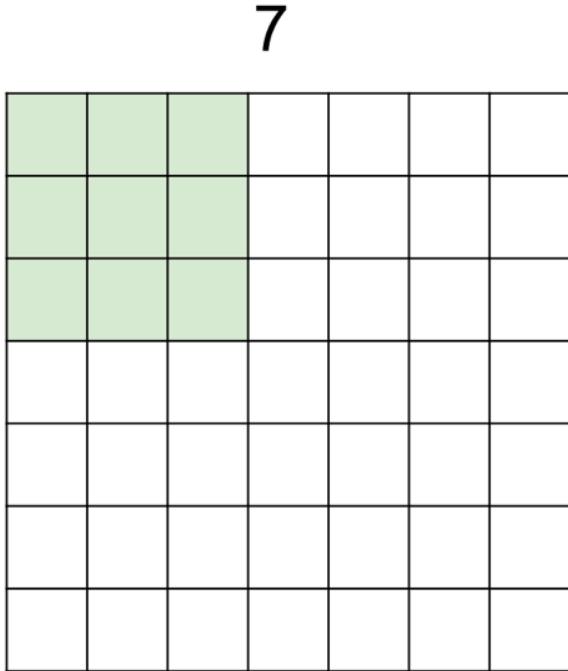


7x7 input (spatially)
assume 3x3 filter

=> **5x5 output**

Convolutional Layer

A closer look at spatial dimensions:



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

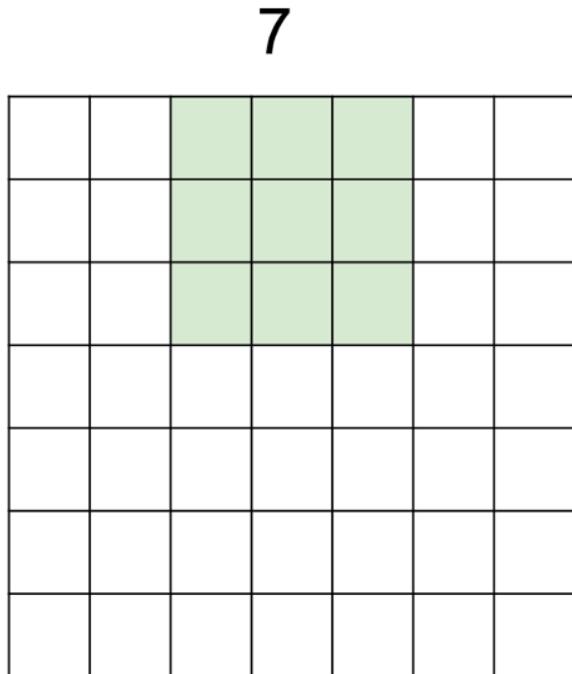
focal
Systems



UNIVERSITY OF
WATERLOO

Convolutional Layer

A closer look at spatial dimensions:

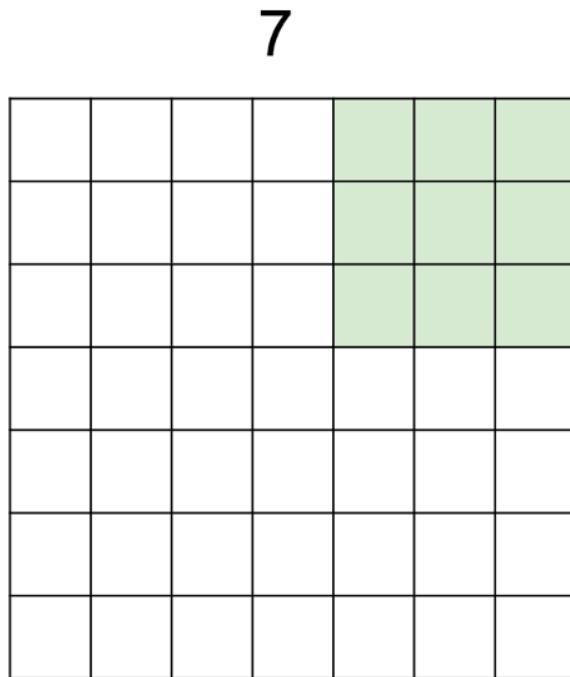


7

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

Convolutional Layer

A closer look at spatial dimensions:

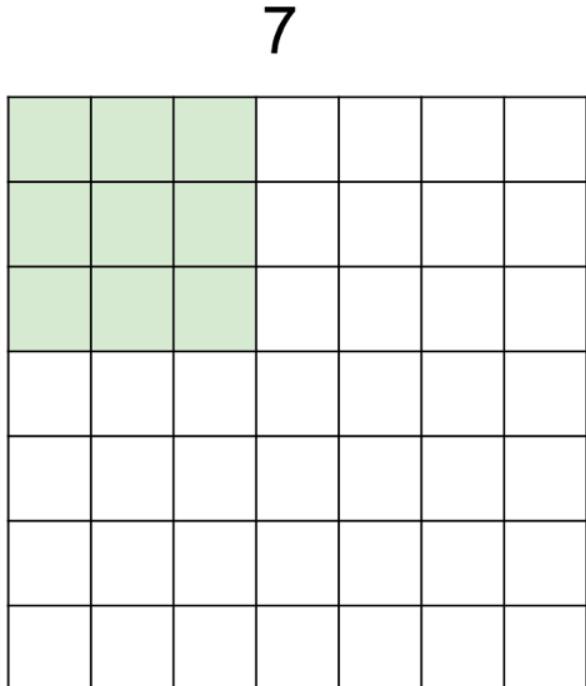


7

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

Convolutional Layer

A closer look at spatial dimensions:

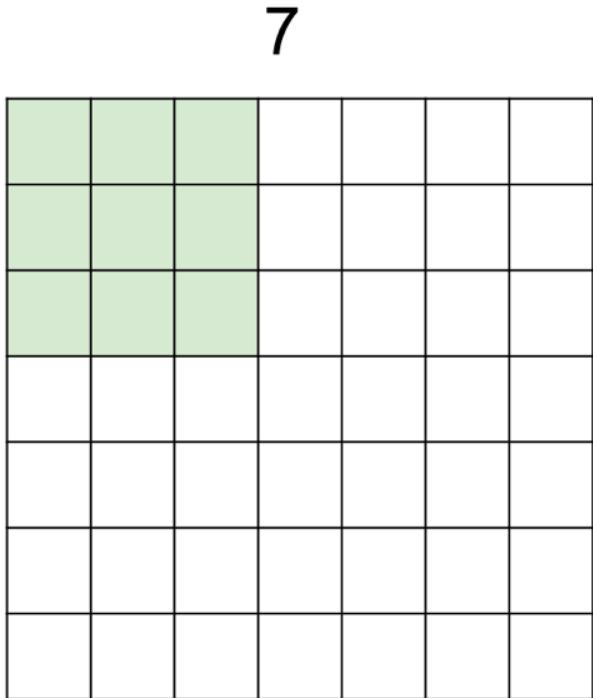


7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**



Convolutional Layer

A closer look at spatial dimensions:



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

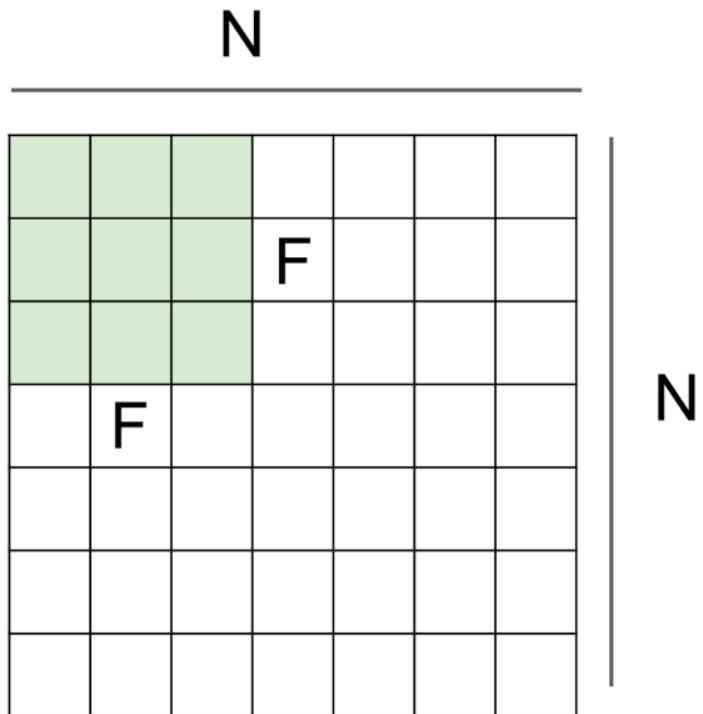
doesn't fit!
cannot apply 3x3 filter on
7x7 input with stride 3.

focal
Systems



UNIVERSITY OF
WATERLOO

Convolutional Layer



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

e.g. $N = 7$, $F = 3$:
stride 1 => $(7 - 3)/1 + 1 = 5$
stride 2 => $(7 - 3)/2 + 1 = 3$
stride 3 => $(7 - 3)/3 + 1 = 2.33$:\

focal
Systems



UNIVERSITY OF
WATERLOO

Convolutional Layer

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:)

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

focal
Systems



UNIVERSITY OF
WATERLOO

Convolutional Layer

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!

focal
Systems



UNIVERSITY OF
WATERLOO

Convolutional Layer

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

focal
Systems

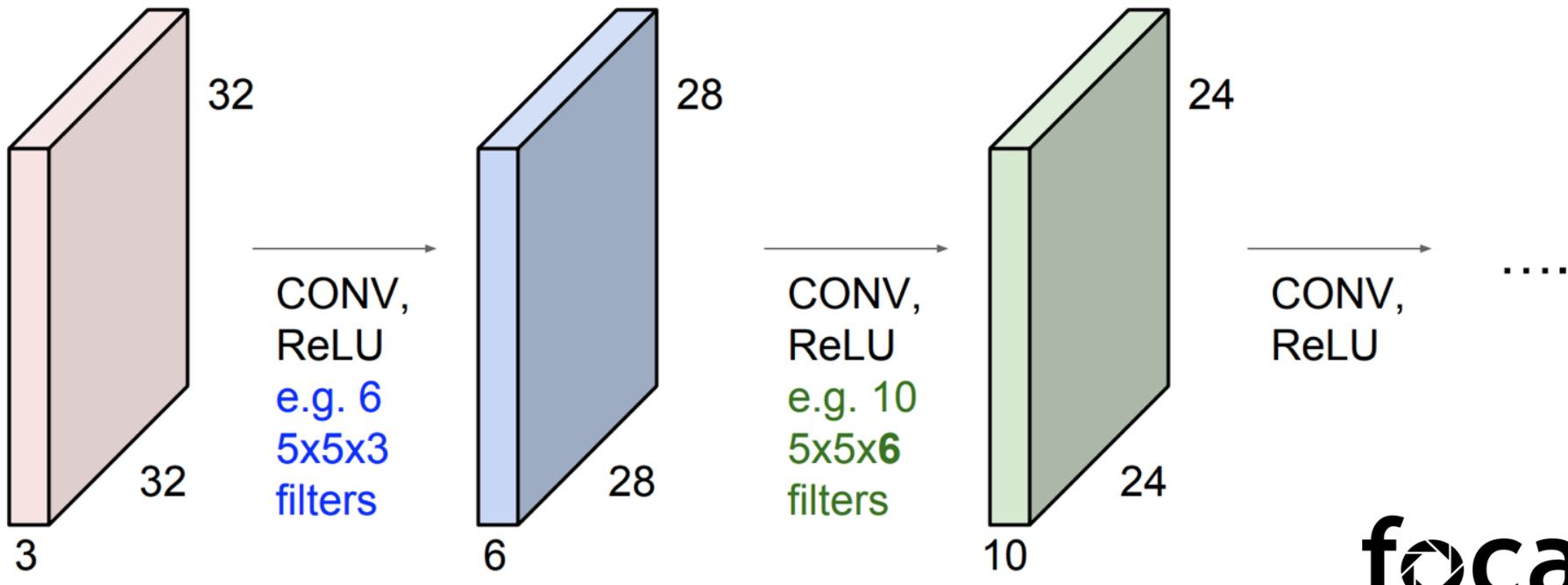


UNIVERSITY OF
WATERLOO

Convolutional Layer

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.



focal
Systems



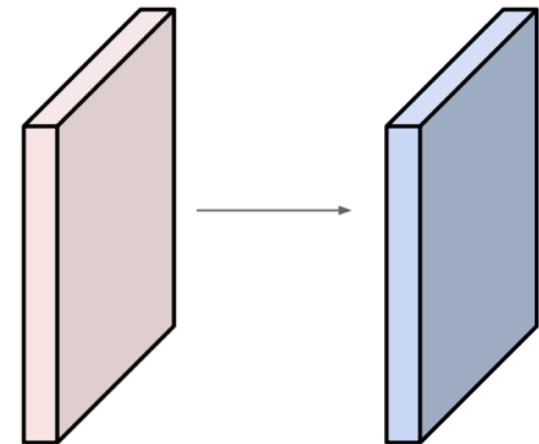
UNIVERSITY OF
WATERLOO

Convolutional Layer

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2



Output volume size: ?

focal
Systems



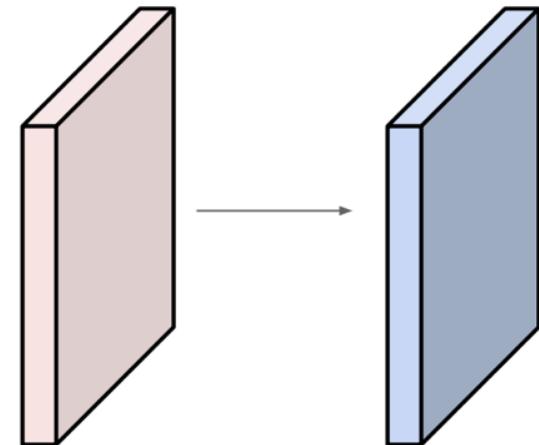
UNIVERSITY OF
WATERLOO

Convolutional Layer

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

focal
Systems



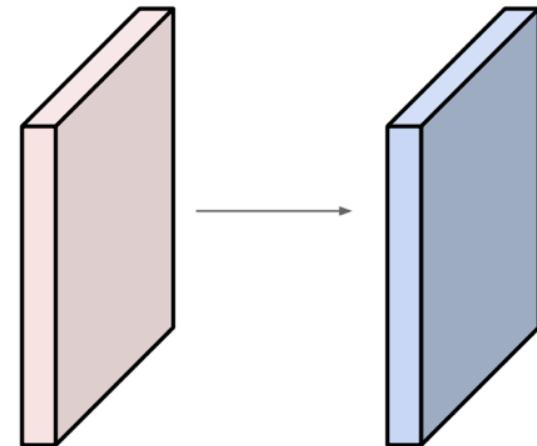
UNIVERSITY OF
WATERLOO

Convolutional Layer

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

focal
Systems



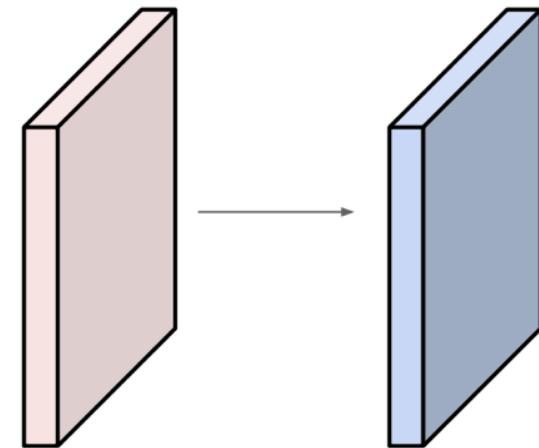
UNIVERSITY OF
WATERLOO

Convolutional 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)

$$\Rightarrow 76*10 = 760$$

focal
Systems



UNIVERSITY OF
WATERLOO

Convolutional Layer

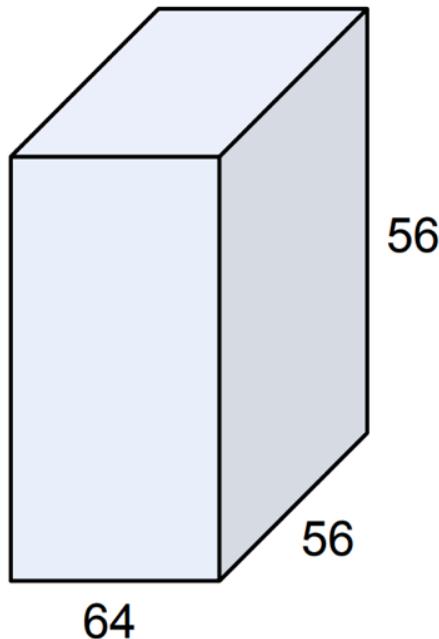
Summary. To summarize, the Conv Layer:

- 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.



Convolutional Layer

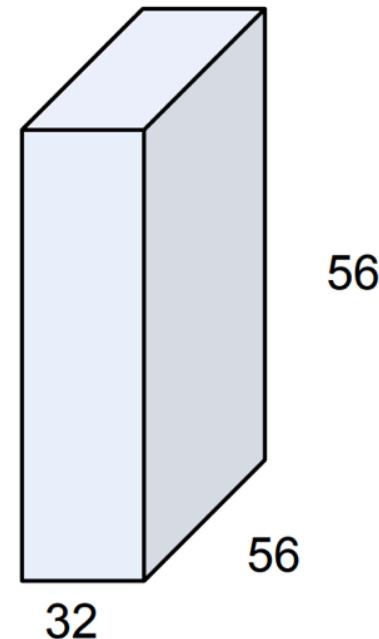
(btw, 1x1 convolution layers make perfect sense)



1x1 CONV
with 32 filters

→

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



focal
Systems



UNIVERSITY OF
WATERLOO



Deep Learning for Retail