

Convolutional Neural Networks

Kietikul Jearanaitanakij

Department of Computer Engineering, KMITL

(Slides are adapted from cs231n @Stanford University)

Convolutional Neural Networks

- Why multilayer neural networks do not satisfy our needs?
 - Fully-connected architecture is **too complex** for computations.
 - There are a lot of weights to update. Very **long training time**.
 - Having too many weights tends to encounter a **poor generalization**.
- The convolutional neural networks can vastly reduce the amount of parameters in the network.

History

A bit of history...

The **Mark I Perceptron** machine was the first implementation of the perceptron algorithm.

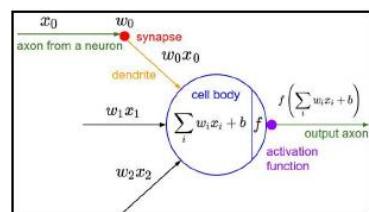
The machine was connected to a camera that used 20×20 cadmium sulfide photocells to produce a 400-pixel image.

recognized
letters of the alphabet

update rule:

$$w_i(t+1) = w_i(t) + \alpha(d_j - y_j(t))x_{j,i}$$

$$f(x) = \begin{cases} 1 & \text{if } w \cdot x + b > 0 \\ 0 & \text{otherwise} \end{cases}$$



Frank Rosenblatt, ~1957: Perceptron



This image by Rocky Acosta is licensed under [CC-BY 3.0](#)

Multilayer neural network

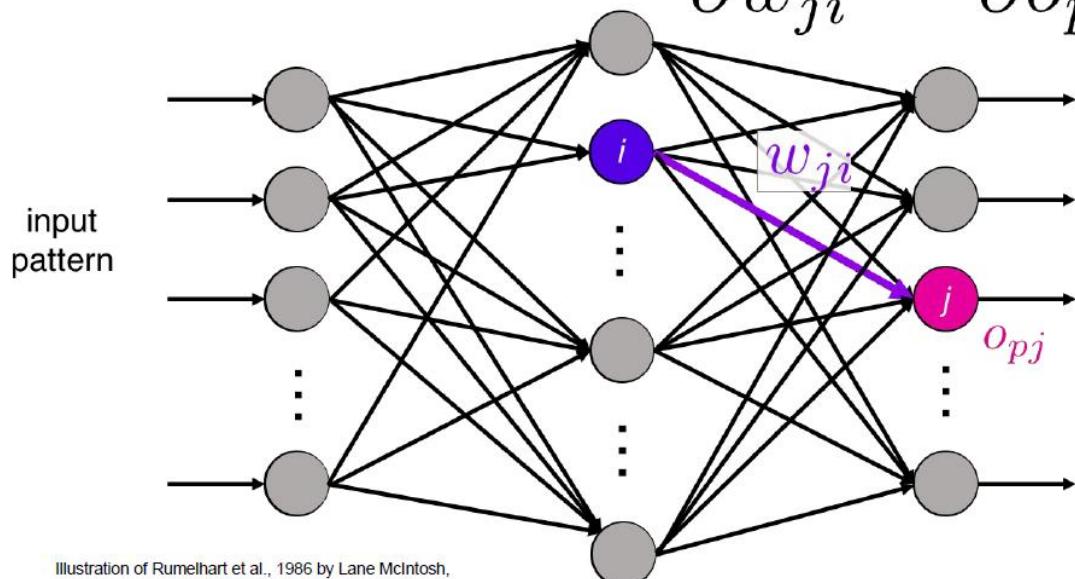


Illustration of Rumelhart et al., 1986 by Lane McIntosh,
copyright CS231n 2017

Error Gradient

Local Gradient

$$\frac{\partial E_p}{\partial w_{ji}} = \frac{\partial E_p}{\partial o_{pj}} \frac{\partial o_{pj}}{\partial w_{ji}}$$

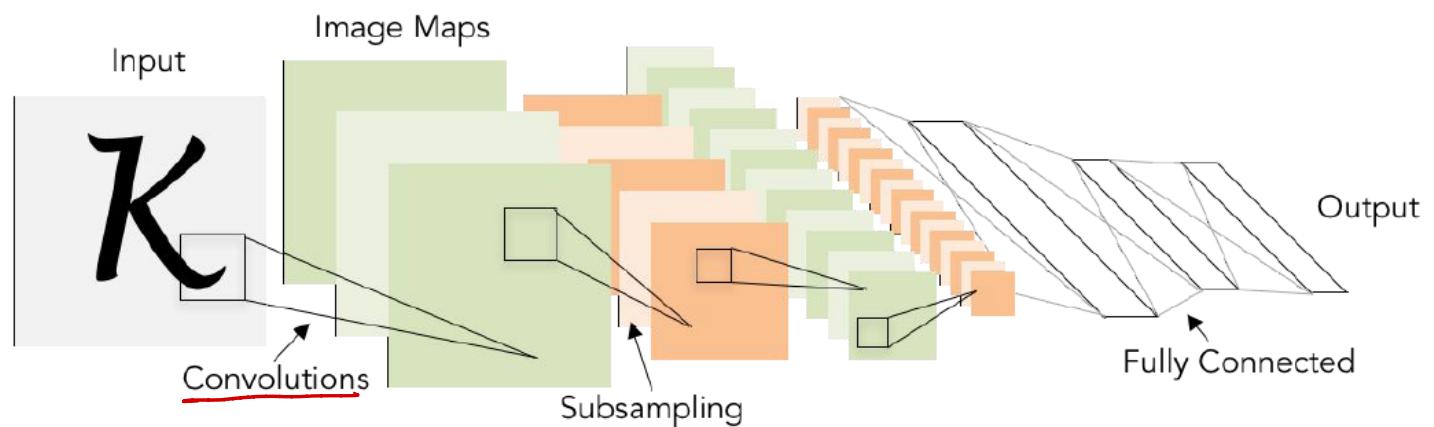
recognizable math

Rumelhart et al., 1986: First time back-propagation became popular

Gradient-based learning applied to document recognition

[LeCun, Bottou, Bengio, Haffner 1998]

LeNet



"การนิ่ง/การขุ่น"

LeNet-5

First strong results

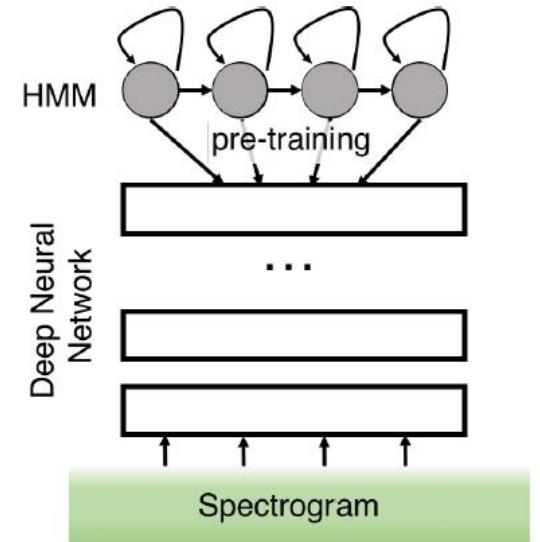
Acoustic Modeling using Deep Belief Networks

Abdel-rahman Mohamed, George Dahl, Geoffrey Hinton, 2010

Context-Dependent Pre-trained Deep Neural Networks

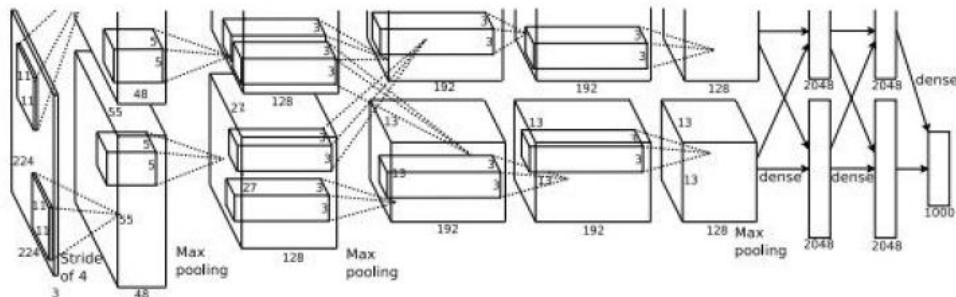
for Large Vocabulary Speech Recognition

George Dahl, Dong Yu, Li Deng, Alex Acero, 2012

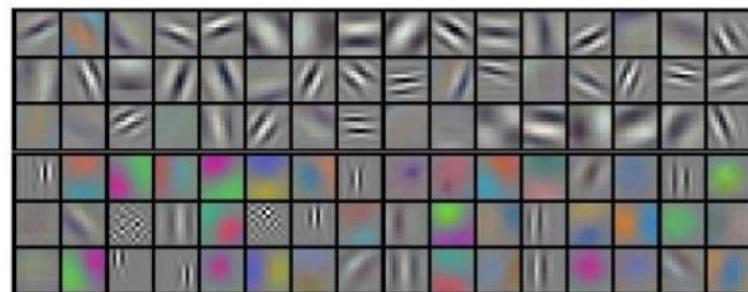


Imagenet classification with deep convolutional neural networks

Alex Krizhevsky, Ilya Sutskever, Geoffrey E Hinton, 2012

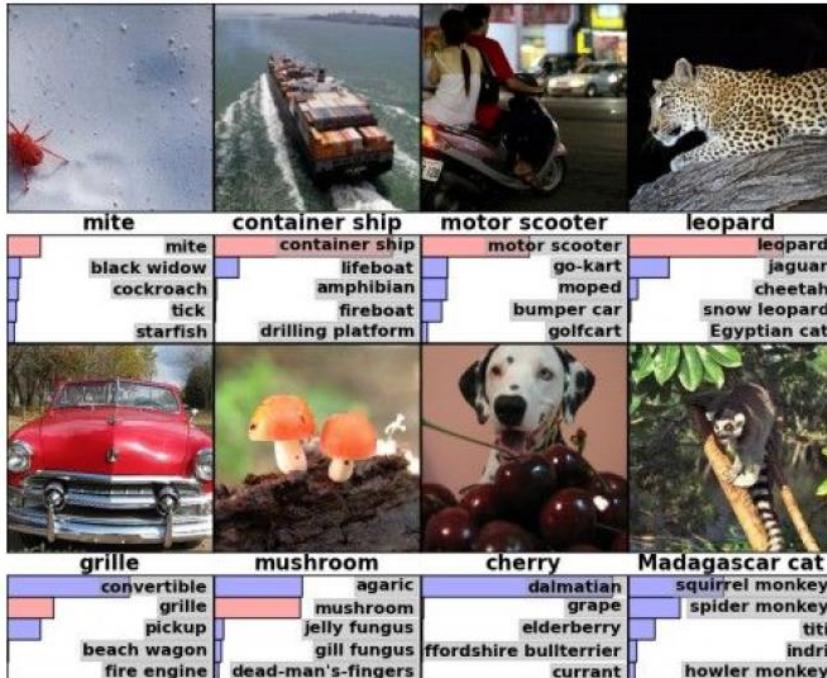


AlexNet



Fast-forward to today: ConvNets are everywhere

Classification



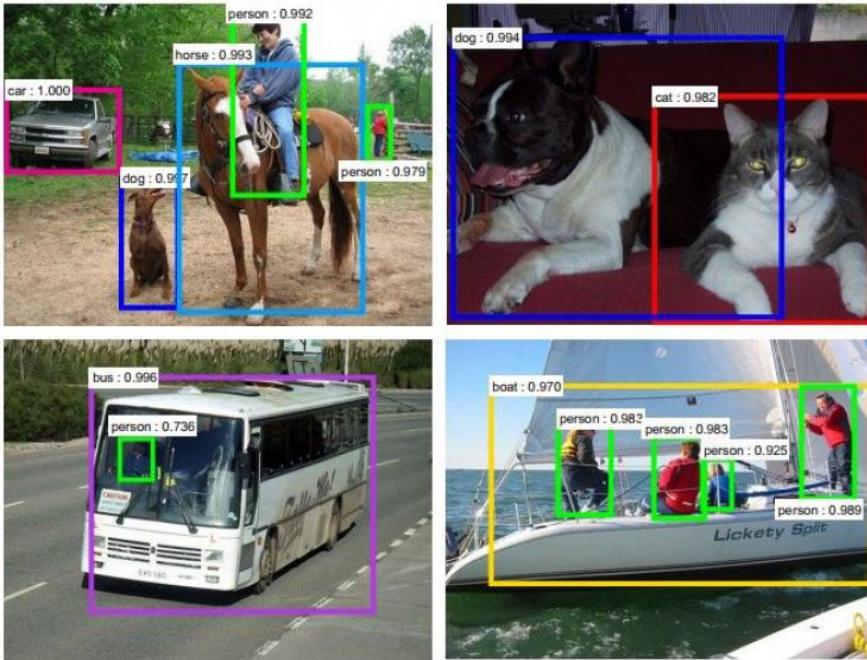
Retrieval



Figures copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

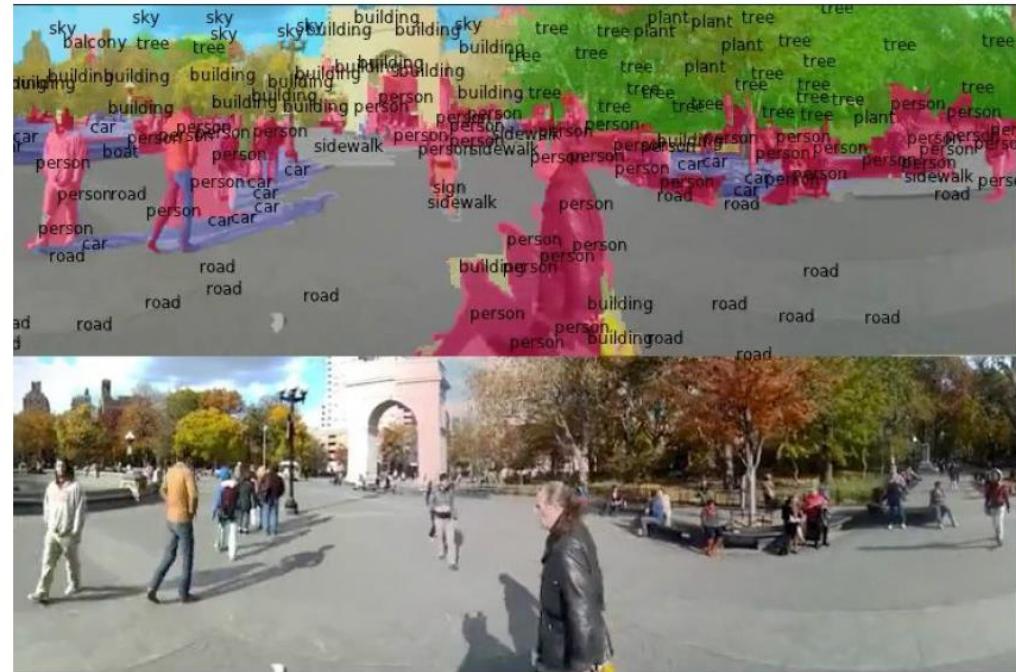
Fast-forward to today: ConvNets are everywhere

Detection



Figures copyright Shaoqing Ren, Kaiming He, Ross Girshick, Jian Sun, 2015. Reproduced with permission.

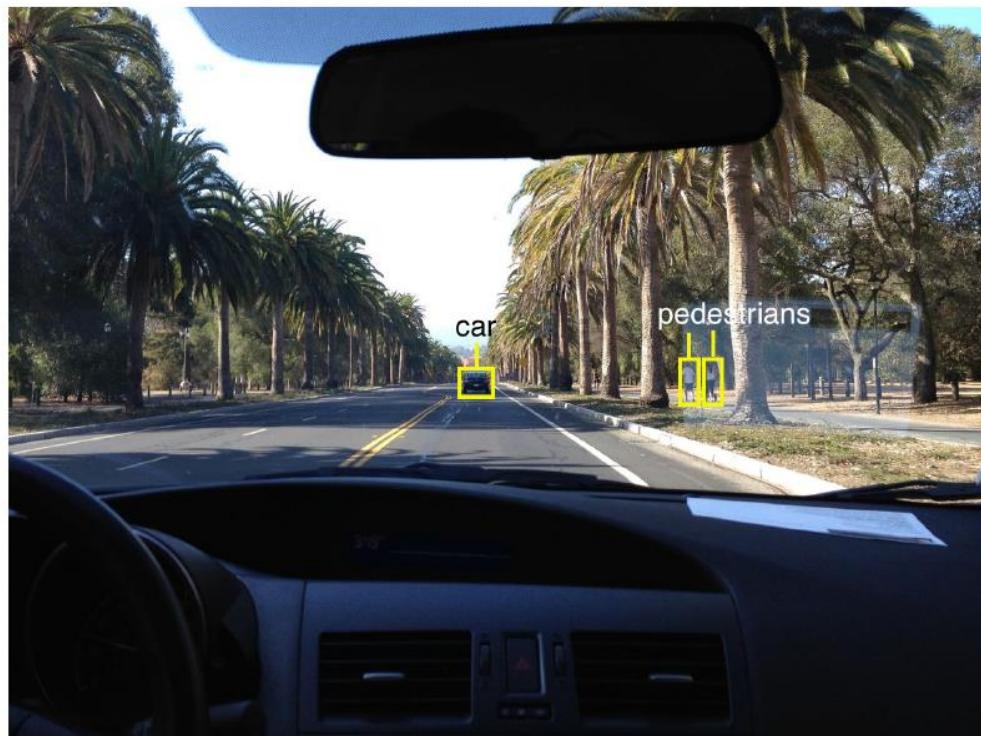
Segmentation



Figures copyright Clement Farabet, 2012.
Reproduced with permission.

[Farabet et al., 2012]

Fast-forward to today: ConvNets are everywhere



self-driving cars

Photo by Lane McIntosh. Copyright CS231n 2017.



[This image](#) by GBPublic_PR is licensed under [CC-BY 2.0](#)

NVIDIA Tesla line

(these are the GPUs on rye01.stanford.edu)

Note that for embedded systems a typical setup would involve NVIDIA Tegras, with integrated GPU and ARM-based CPU cores.

Image Captioning

[Vinyals et al., 2015]
[Karpathy and Fei-Fei, 2015]

No errors



A white teddy bear sitting in the grass

Minor errors



A man in a baseball uniform throwing a ball

Somewhat related



A woman is holding a cat in her hand



A man riding a wave on top of a surfboard



A cat sitting on a suitcase on the floor



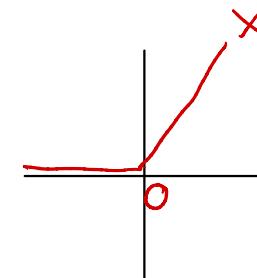
A woman standing on a beach holding a surfboard

All images are CC0 Public domain:
<https://pixabay.com/en/luggage-antique-cat-1643010/>
<https://pixabay.com/en/teddy-plush-bears-cute-teddy-bear-1623436/>
<https://pixabay.com/en/surf-wave-summer-sport-litoral-1668716/>
<https://pixabay.com/en/woman-female-model-portrait-adult-983967/>
<https://pixabay.com/en/handstand-lake-meditation-496008/>
<https://pixabay.com/en/baseball-player-shortstop-infield-1045263/>

Captions generated by Justin Johnson using [Neuraltalk2](#)

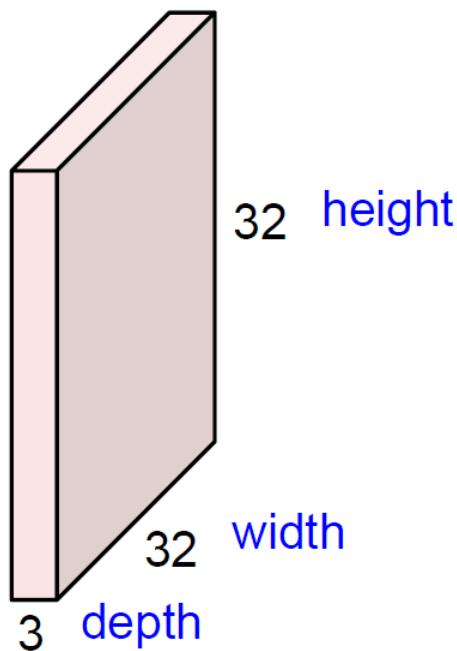
Layers in Convolutional Neural Networks

- CNN : Convolutional Neural Network
- There are usually 4 types of layers in CNN.
 1. Convolution layer
 2. Pooling layer
 3. ReLU (Rectified Linear Unit) layer
 4. Fully connected layer



Convolution Layer

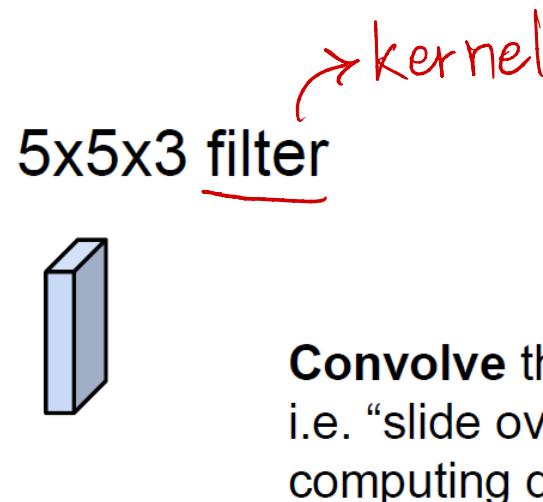
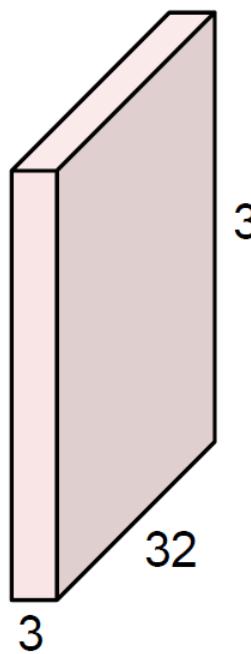
32x32x3 image -> preserve spatial structure



Input

Convolution Layer

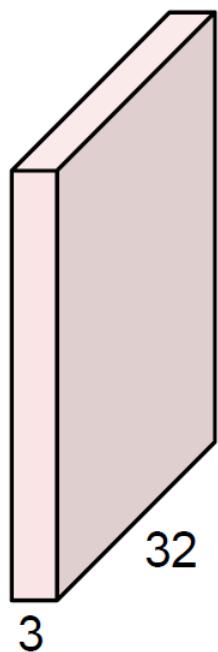
32x32x3 image



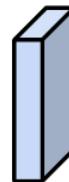
Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Convolution Layer

32x32x3 image



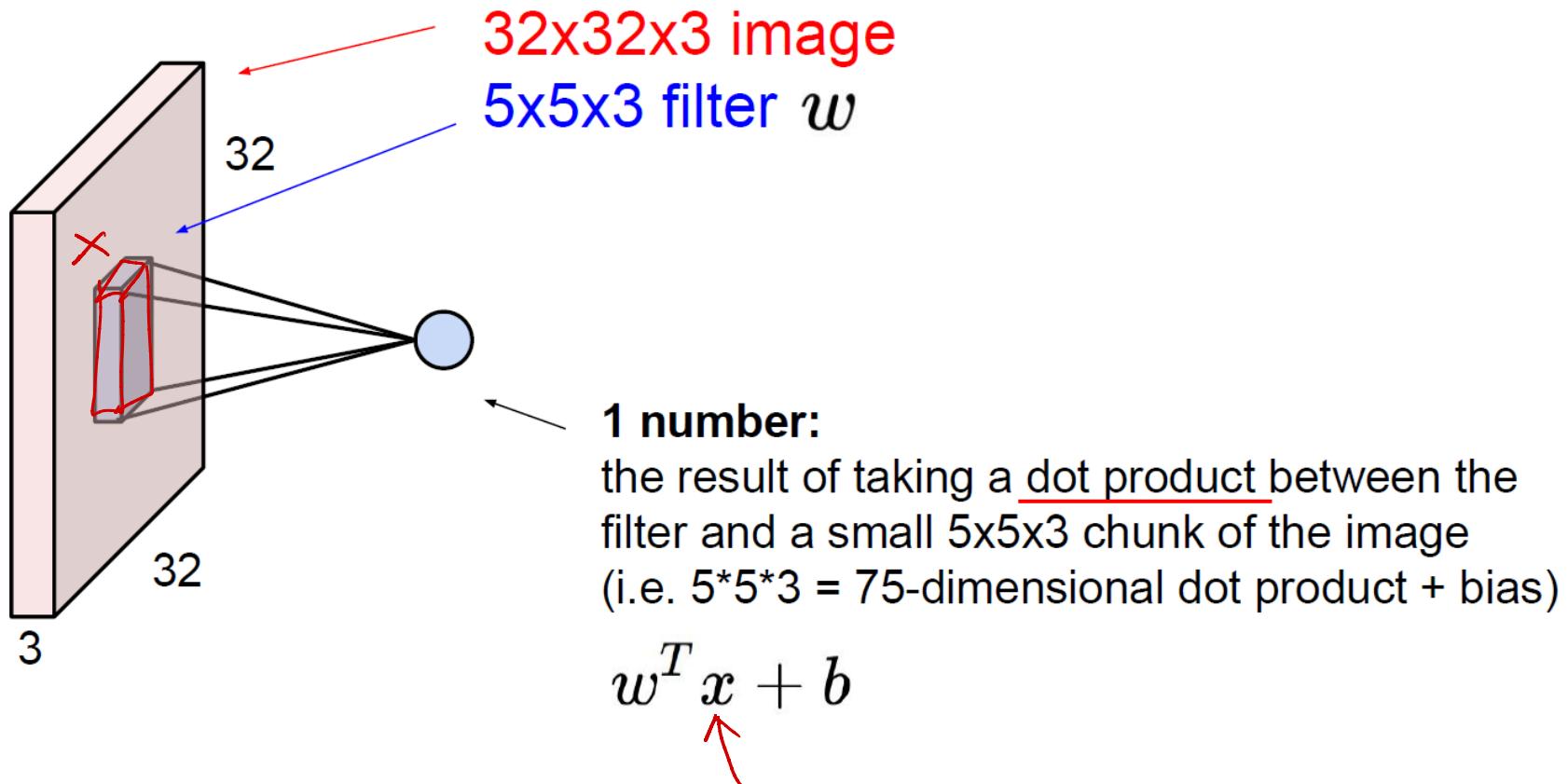
5x5x3 filter



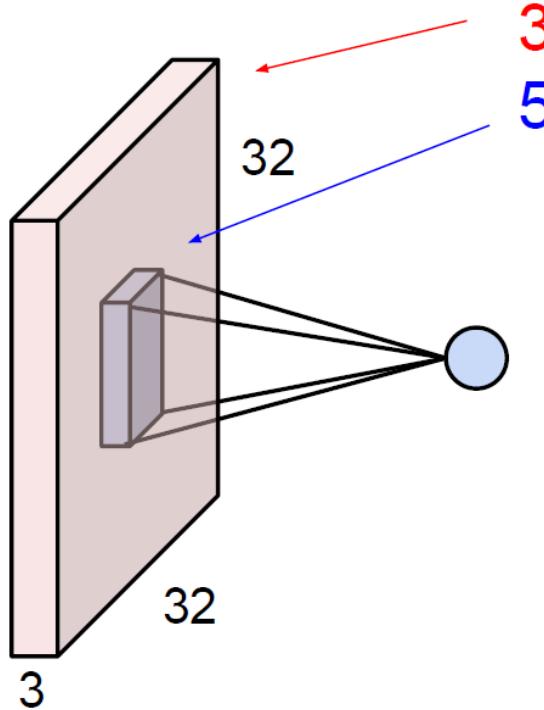
Filters always extend the full depth of the input volume

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

Convolution Layer



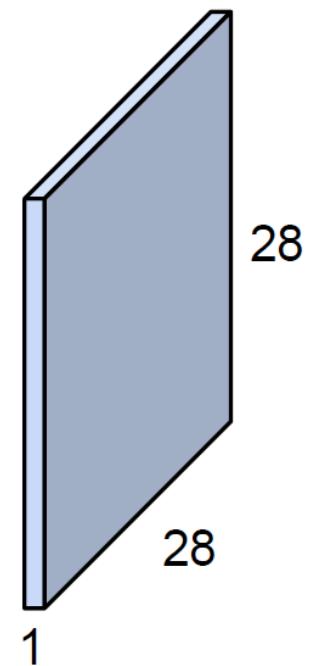
Convolution Layer



32x32x3 image
5x5x3 filter

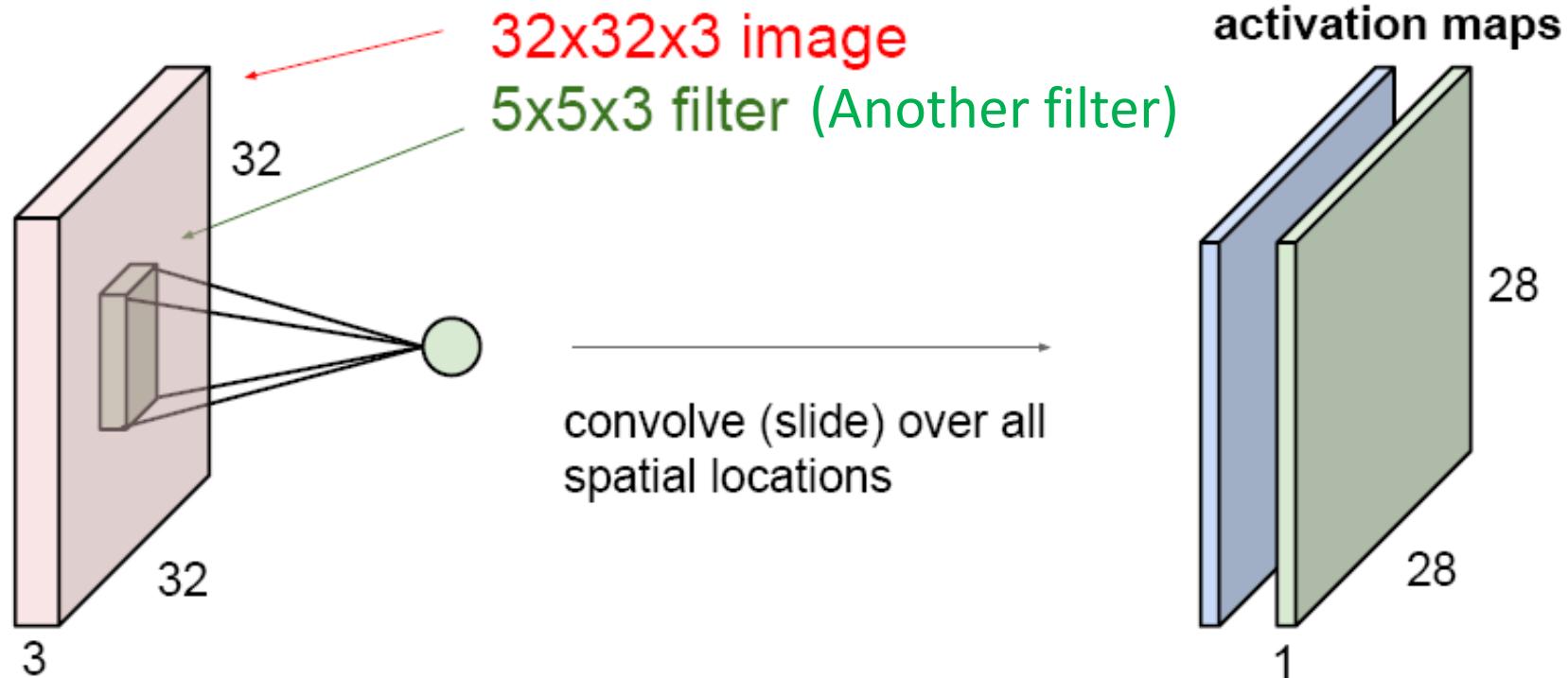
convolve (slide) over all
spatial locations

activation map
(feature map)



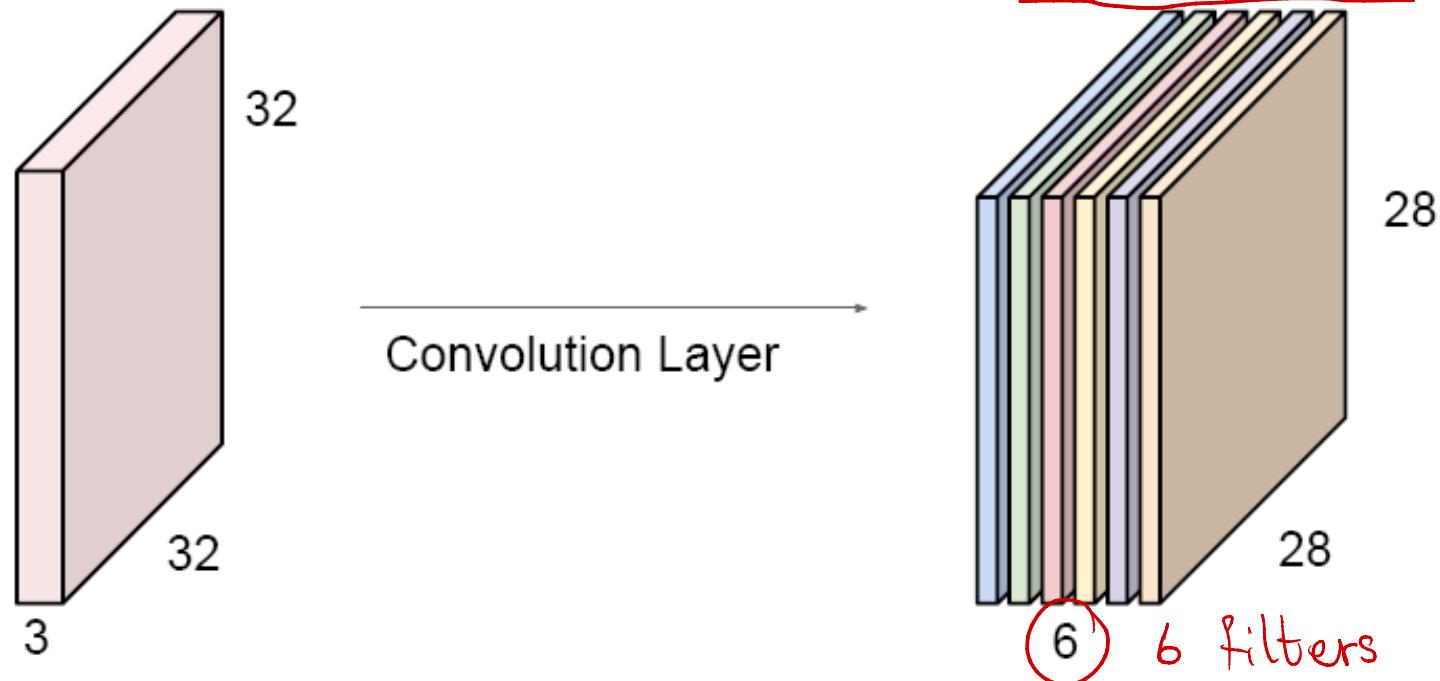
Convolution Layer

consider a second, green filter



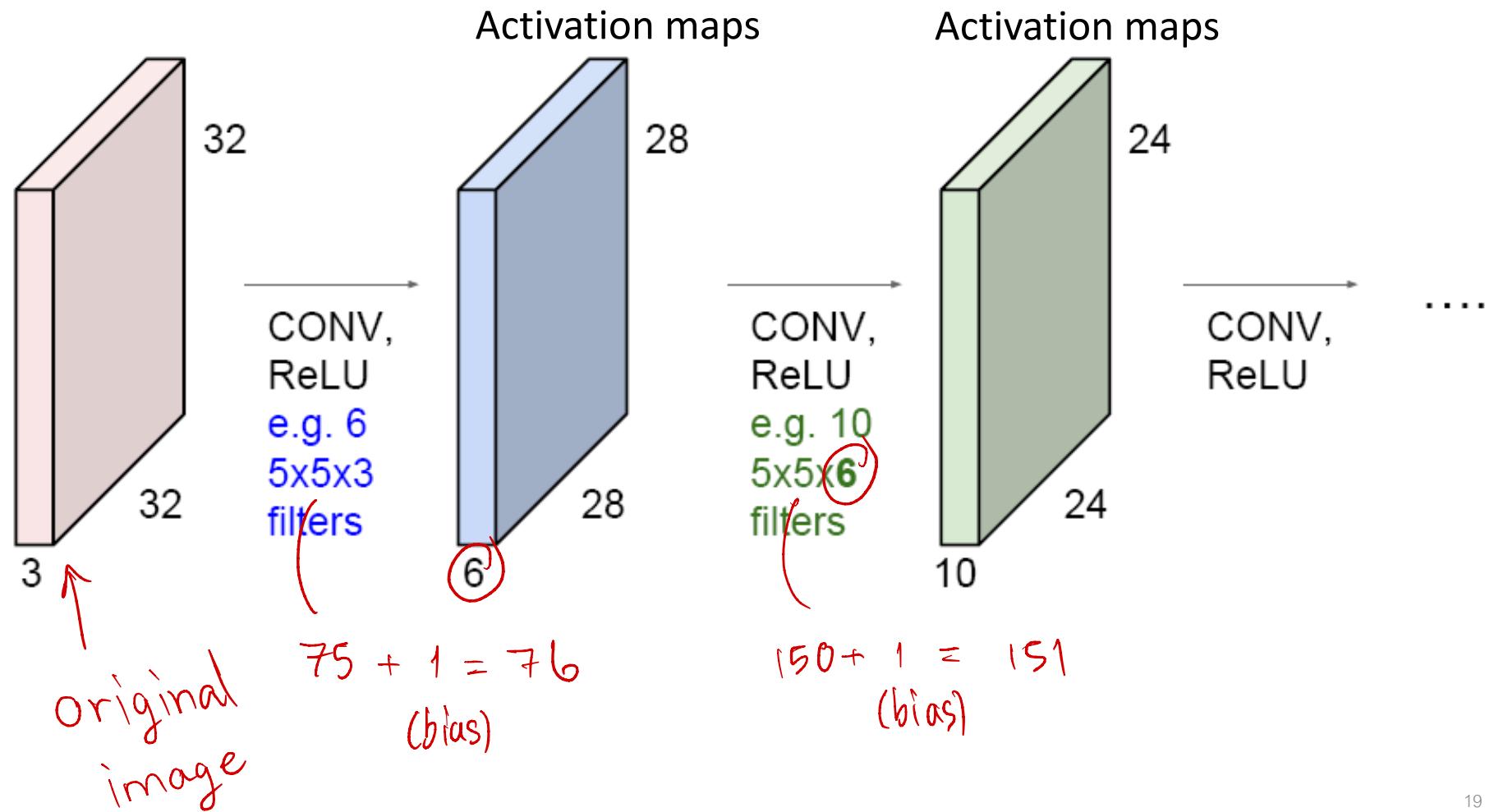
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

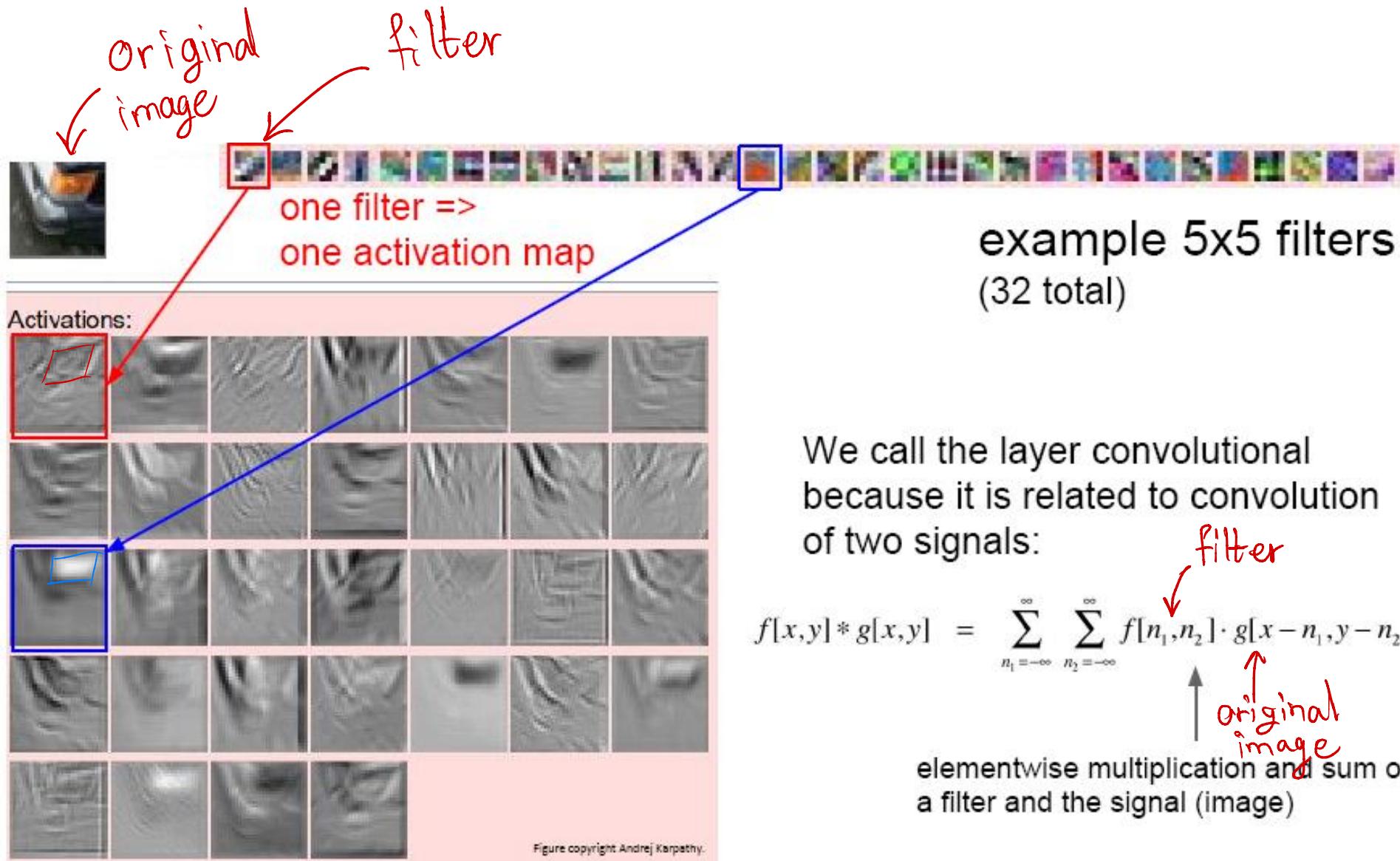
(Feature maps)
activation maps



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

Preview: ConvNet is a sequence of Convolution Layers



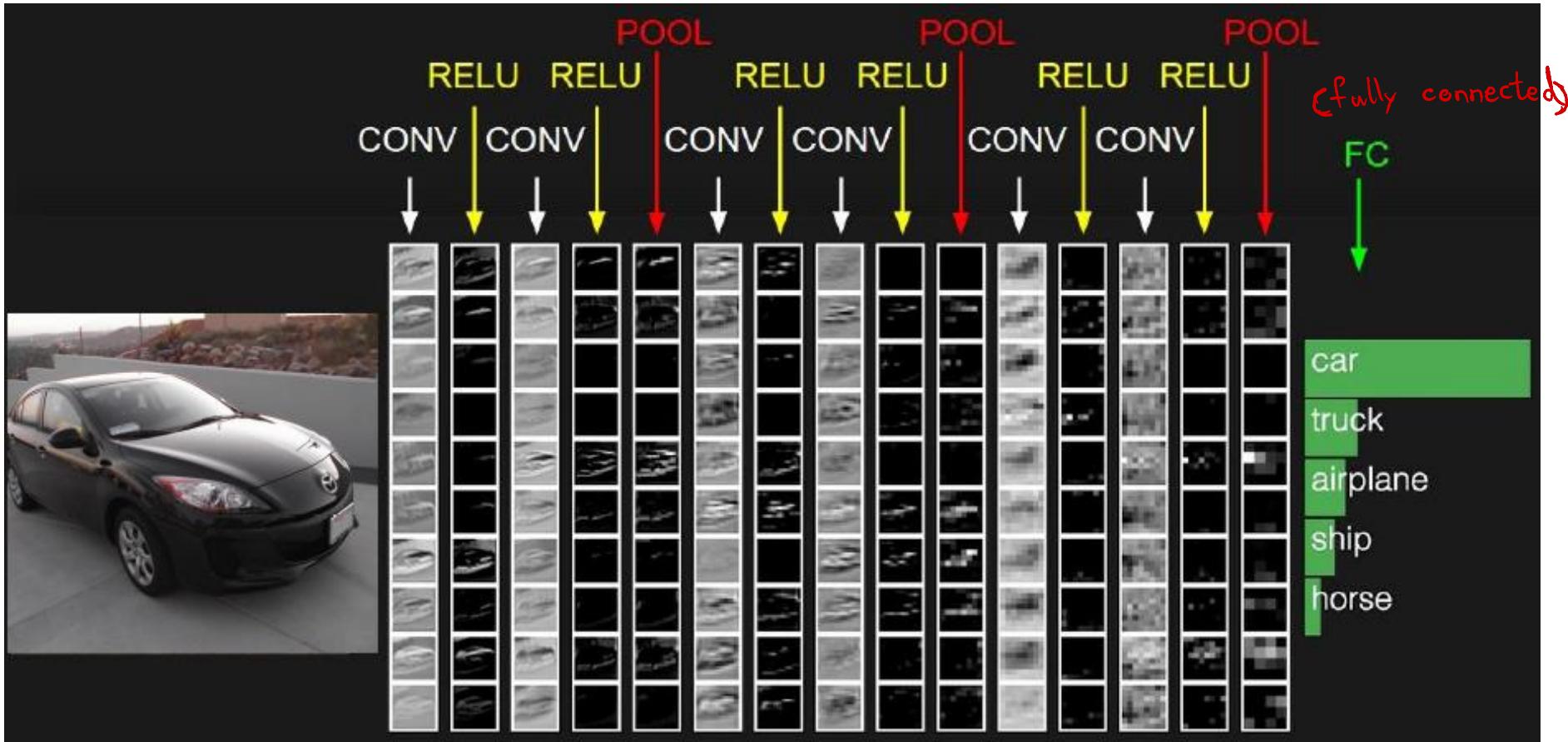


We call the layer convolutional because it is related to convolution of two signals:

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

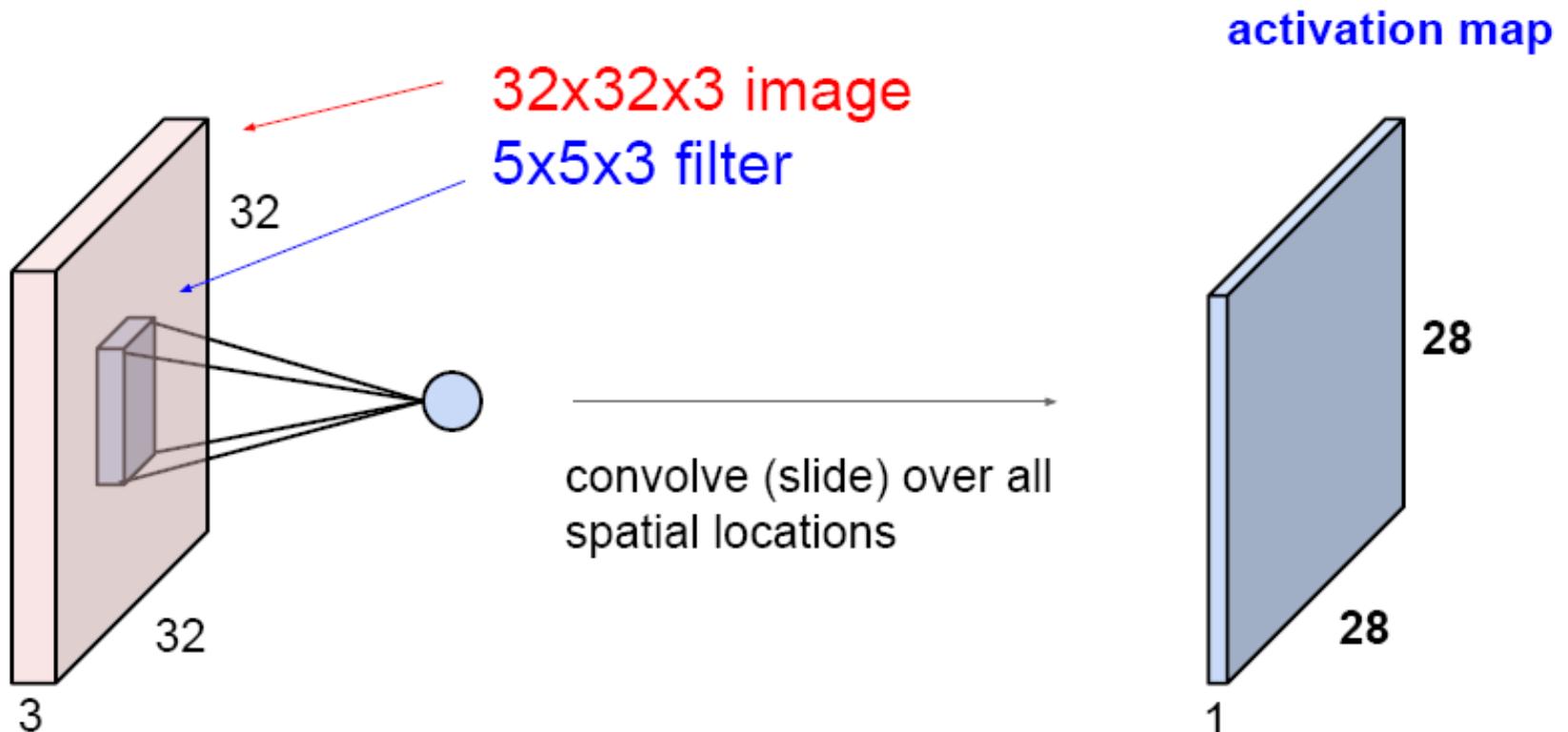
↑
original image
filter

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



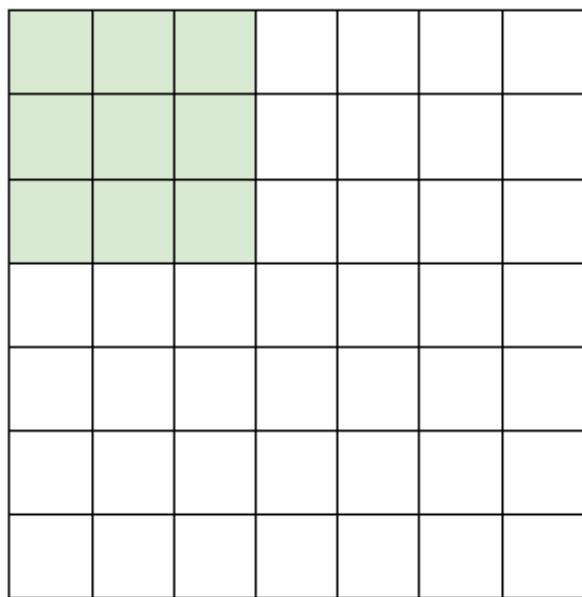
How to calculate the size of activation map ?

A closer look at spatial dimensions:



A closer look at spatial dimensions:

7



7x7 input (spatially)
assume 3x3 filter

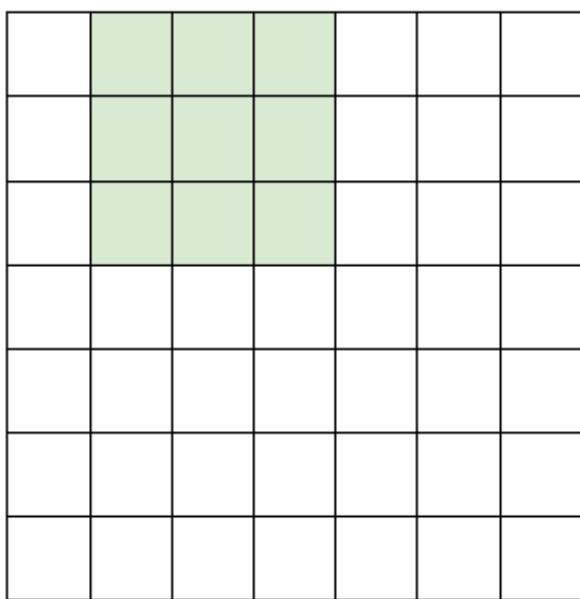


7

A closer look at spatial dimensions:

stride = 1

7

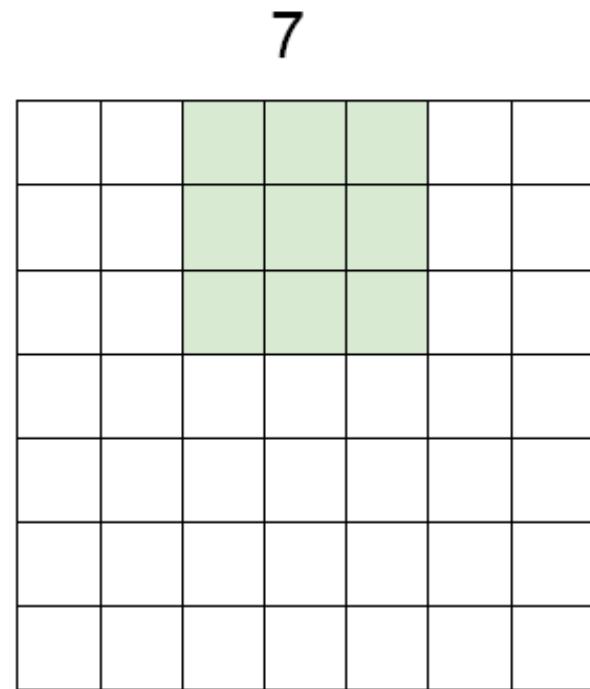


7x7 input (spatially)
assume 3x3 filter

00

7

A closer look at spatial dimensions:

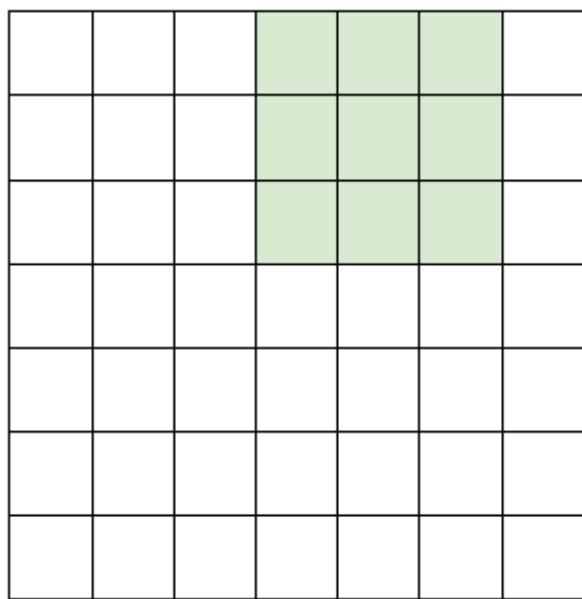


7x7 input (spatially)
assume 3x3 filter

000

A closer look at spatial dimensions:

7

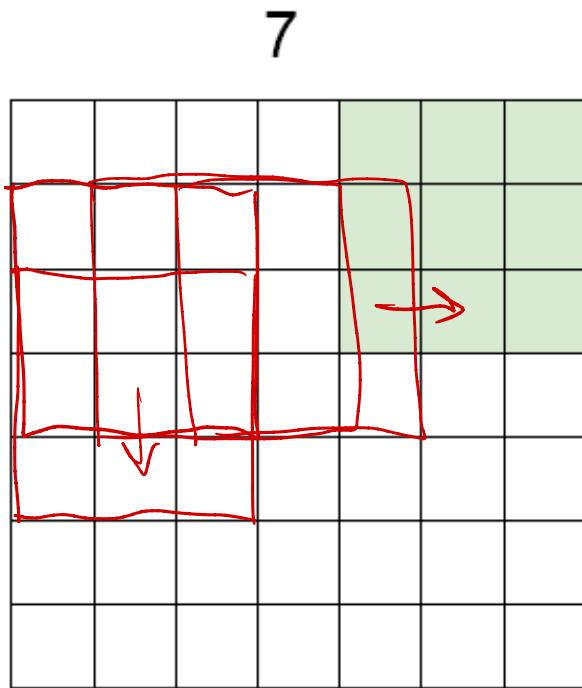


7x7 input (spatially)
assume 3x3 filter

0000

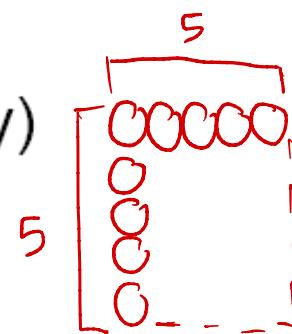
7

A closer look at spatial dimensions:

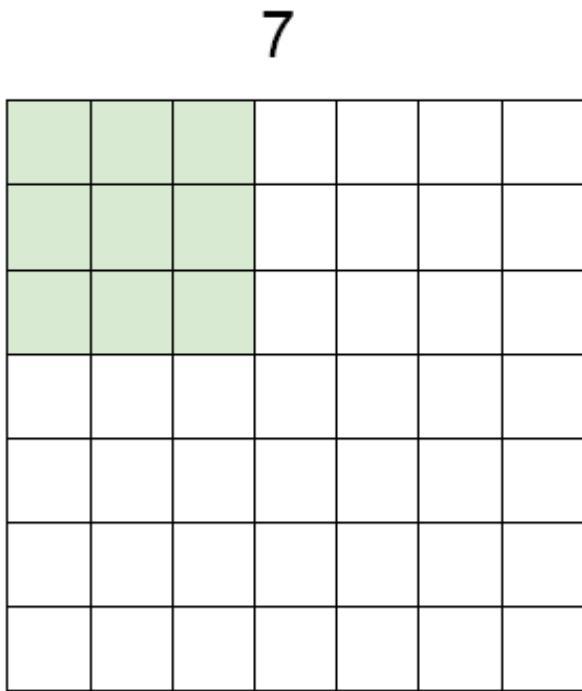


7x7 input (spatially)
assume 3x3 filter

=> 5x5 output

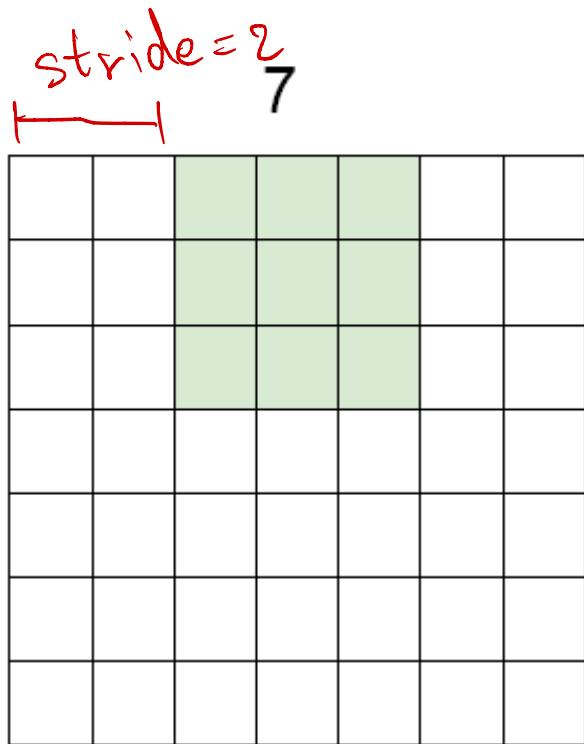


A closer look at spatial dimensions:



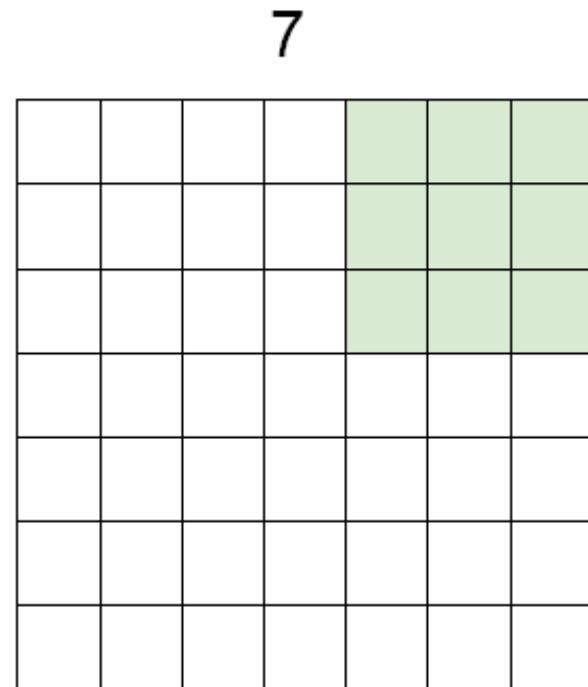
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



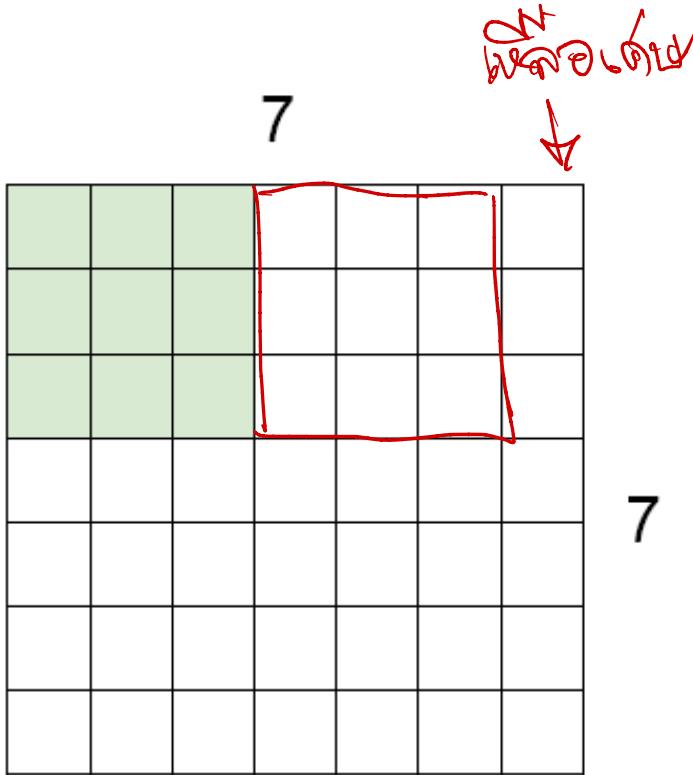
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



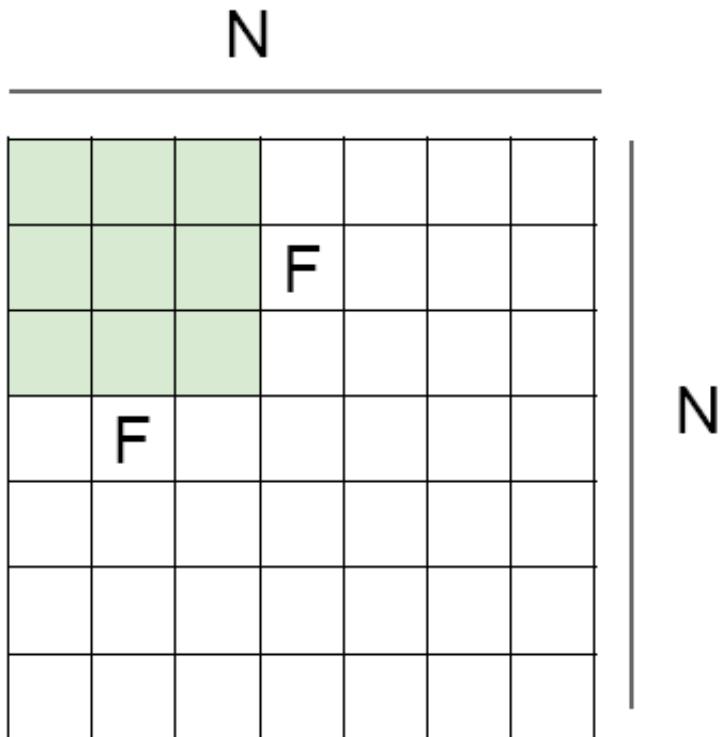
7x7 input (spatially)
assume 3x3 filter
applied with **stride 2**
=> 3x3 output!

A closer look at spatial dimensions:



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

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



Activation Map size
)

Output size:

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

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

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

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

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

“ក្នុងនេះគឺជាបរាយលស និងទទួលបានដំឡើ”

In practice: Common to zero pad the border

0	0	0	0	0	0	0			
0									
0									
0									
0									

$P=1$

e.g. input 7x7

3x3 filter, applied with **stride 1**

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

$$\text{Activation Map size} = (N-F+2P)/S + 1$$

7x7 output!

$$(7-3+2(1))/1 + 1$$

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

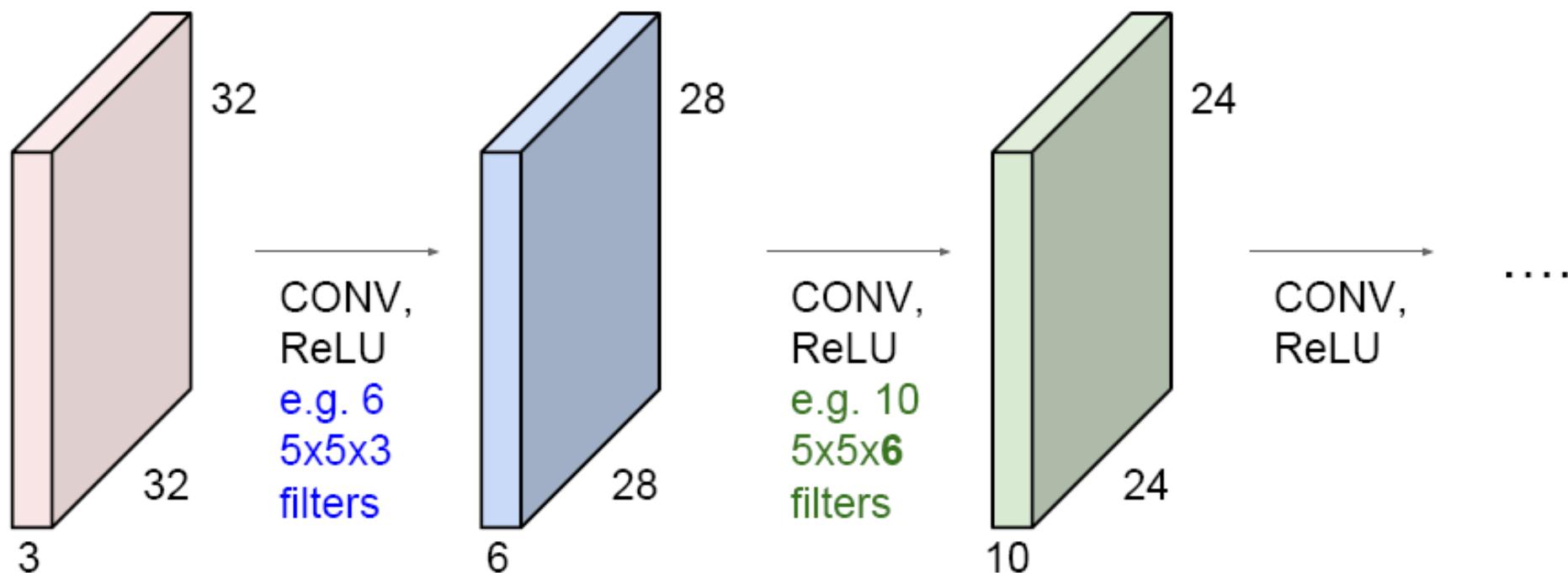
$F = 5 \Rightarrow$ zero pad with 2

$F = 7 \Rightarrow$ zero pad with 3

All of these produce 7x7 output.

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:

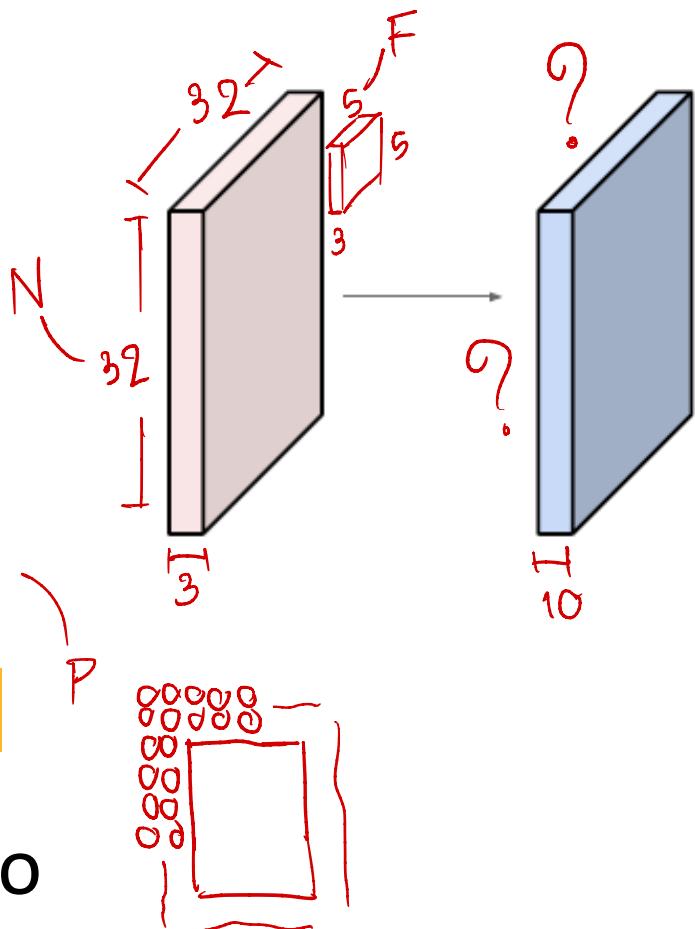
Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Output volume size: ? $(N-F+2P)/S + 1$

$$(32-5+2*2)/1+1 = 32 \text{ spatially, so}$$

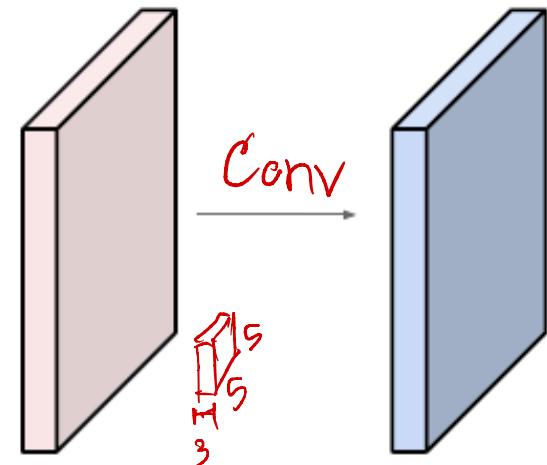
32x32x**10**



Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

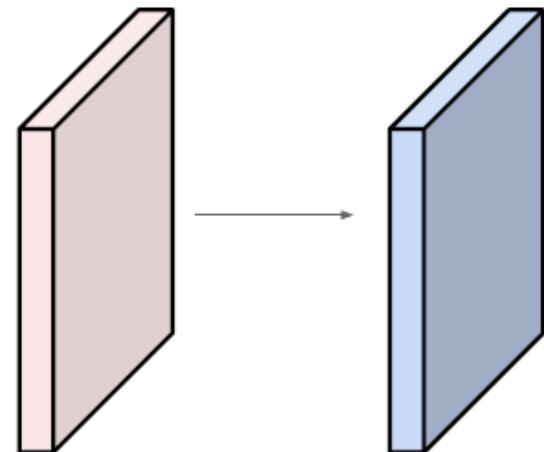


Number of parameters in this layer?

Examples time:

Input volume: **32x32x3**

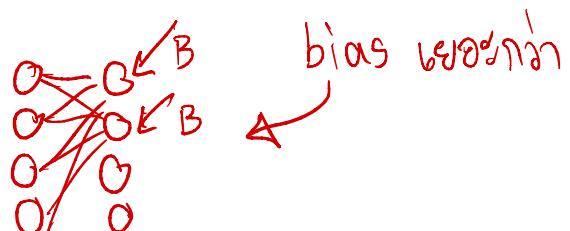
10 **5x5** filters with stride 1, pad 2



Number of parameters in this layer?

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

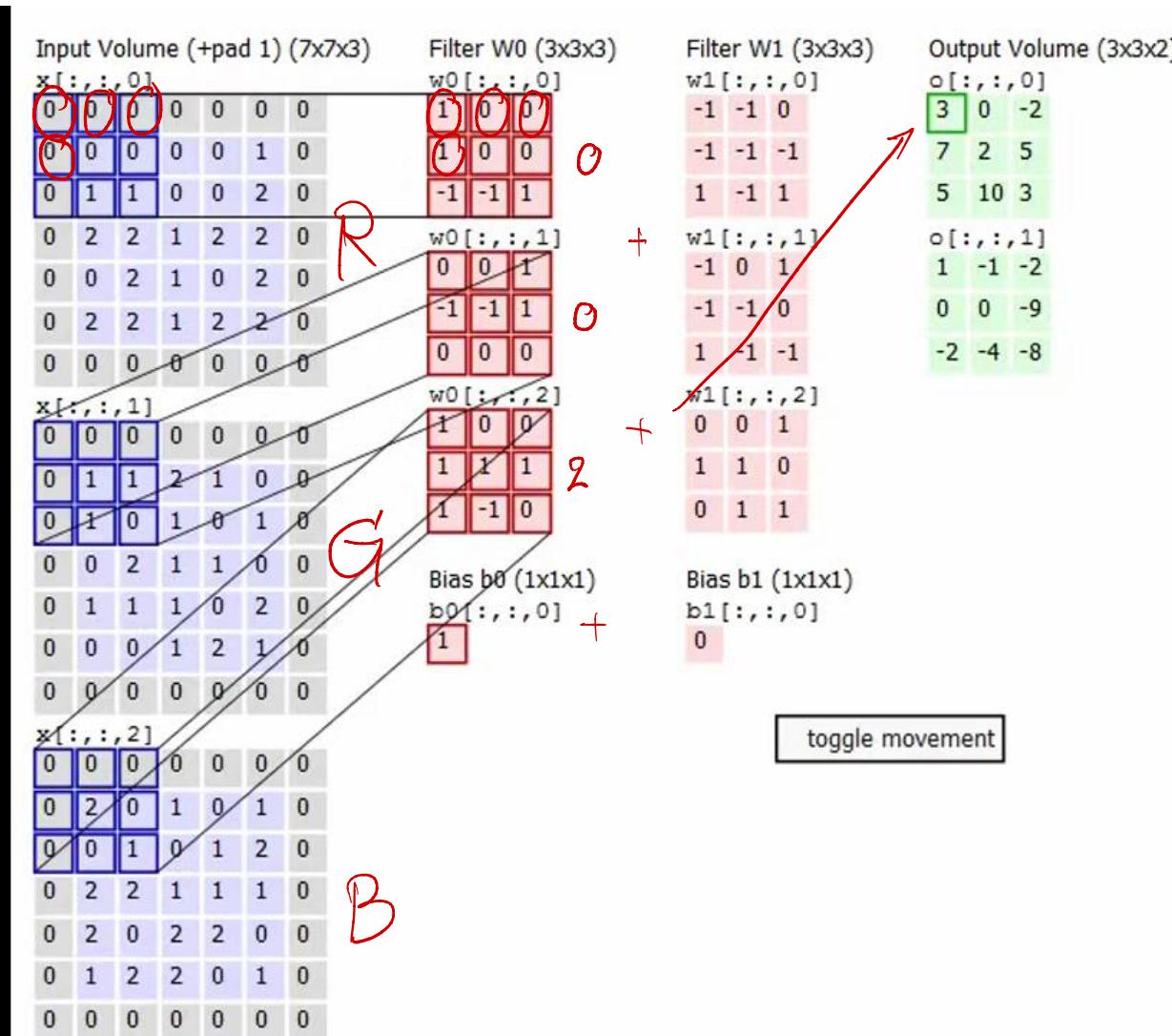
(+1 for bias)
 \equiv



Summary. To summarize, the Conv Layer:

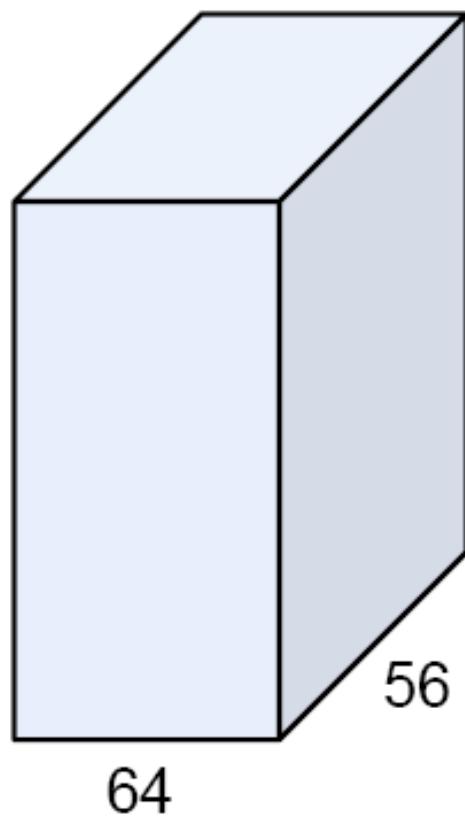
- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (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.

Watch a simulation at <http://cs231n.github.io/convolutional-networks/>



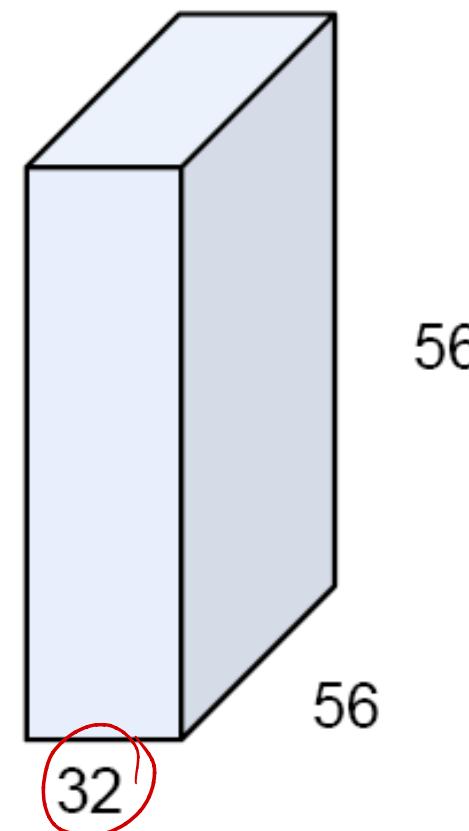
ເຕັມໄວ້ ລາຍ dimension ອີ 3 ຂົວການ

(btw, 1x1 convolution layers make perfect sense)



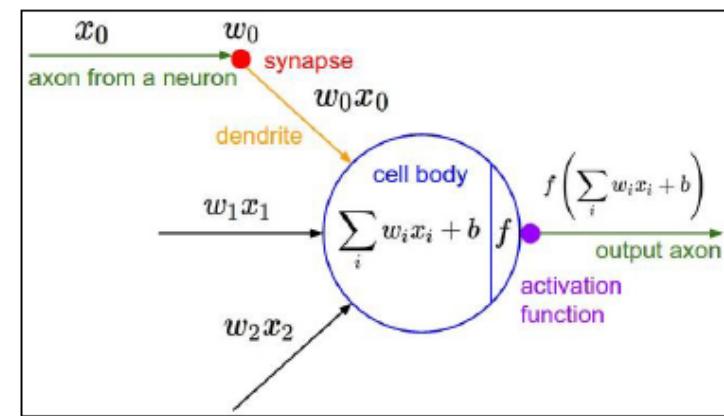
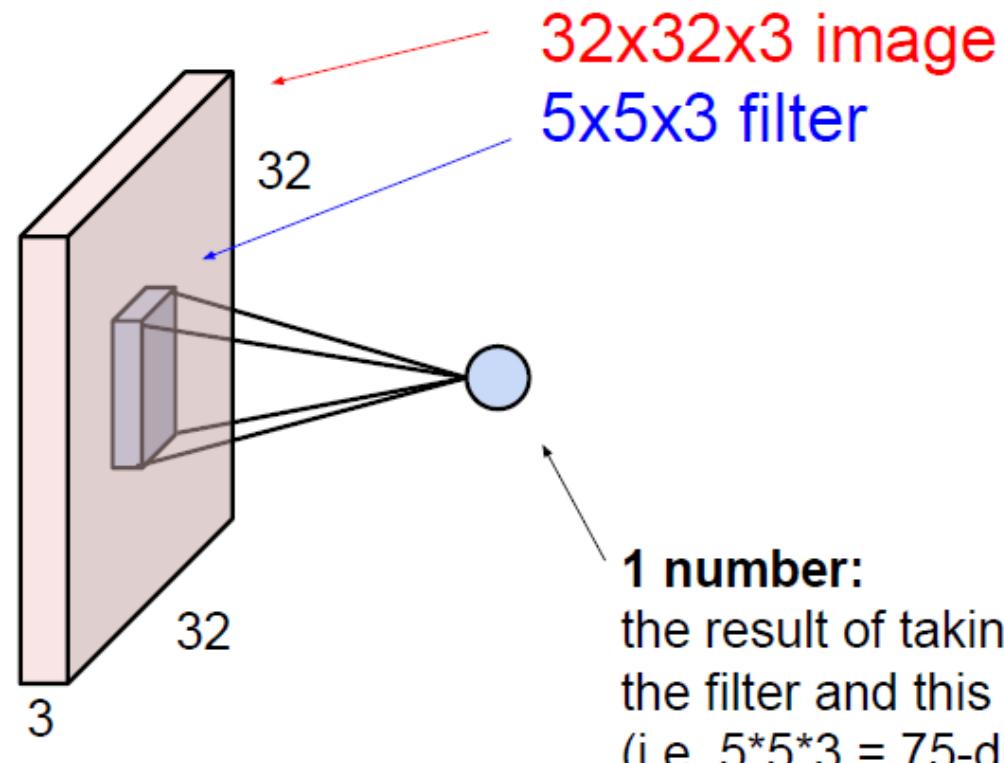
1x1 CONV
with 32 filters

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



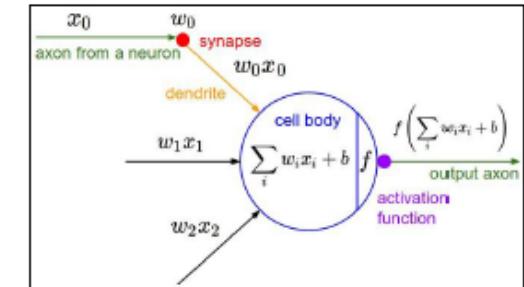
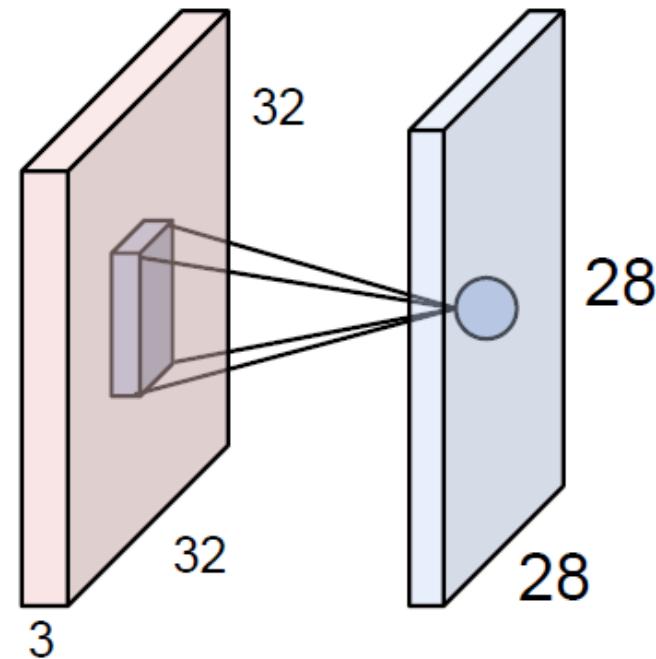
ເນັດວ່າອຸປະກອນໄດ້ຫຼືຍຸ?

The brain/neuron view of CONV Layer



It's just a neuron with local connectivity...

The brain/neuron view of CONV Layer

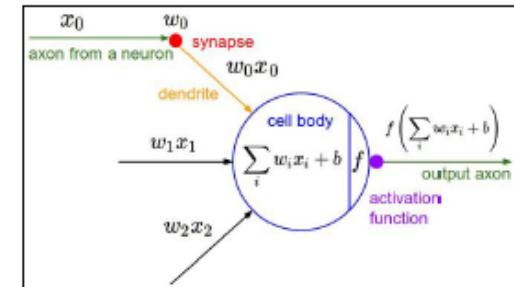
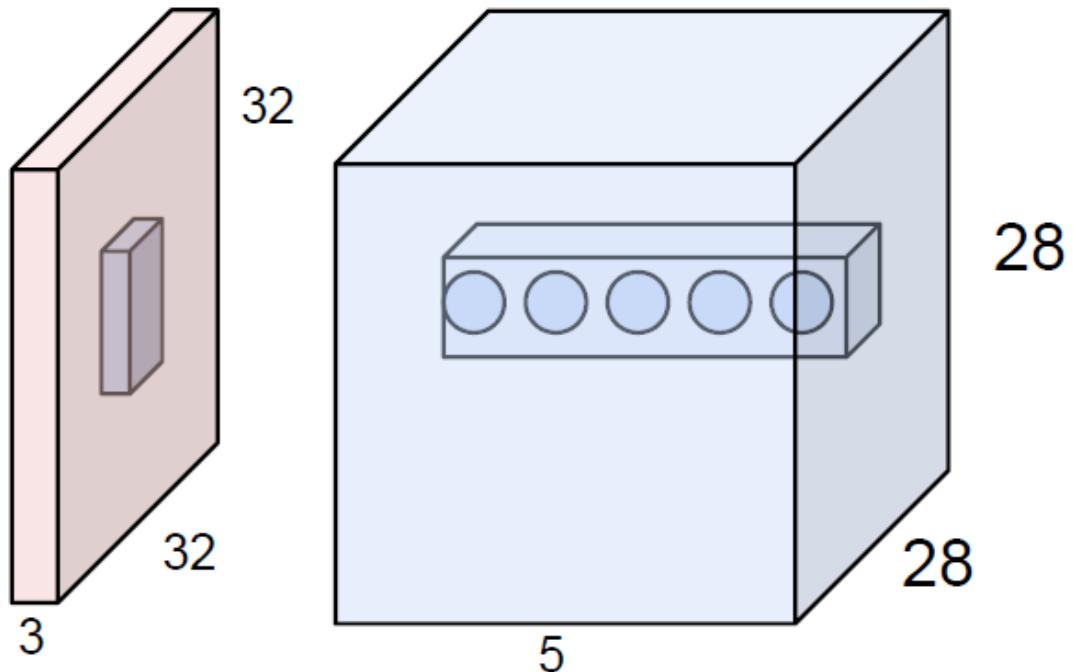


An activation map is a 28x28 sheet of neuron outputs:

1. Each is connected to a small region in the input
2. All of them share parameters

“5x5 filter” -> “5x5 receptive field for each neuron”

The brain/neuron view of CONV Layer



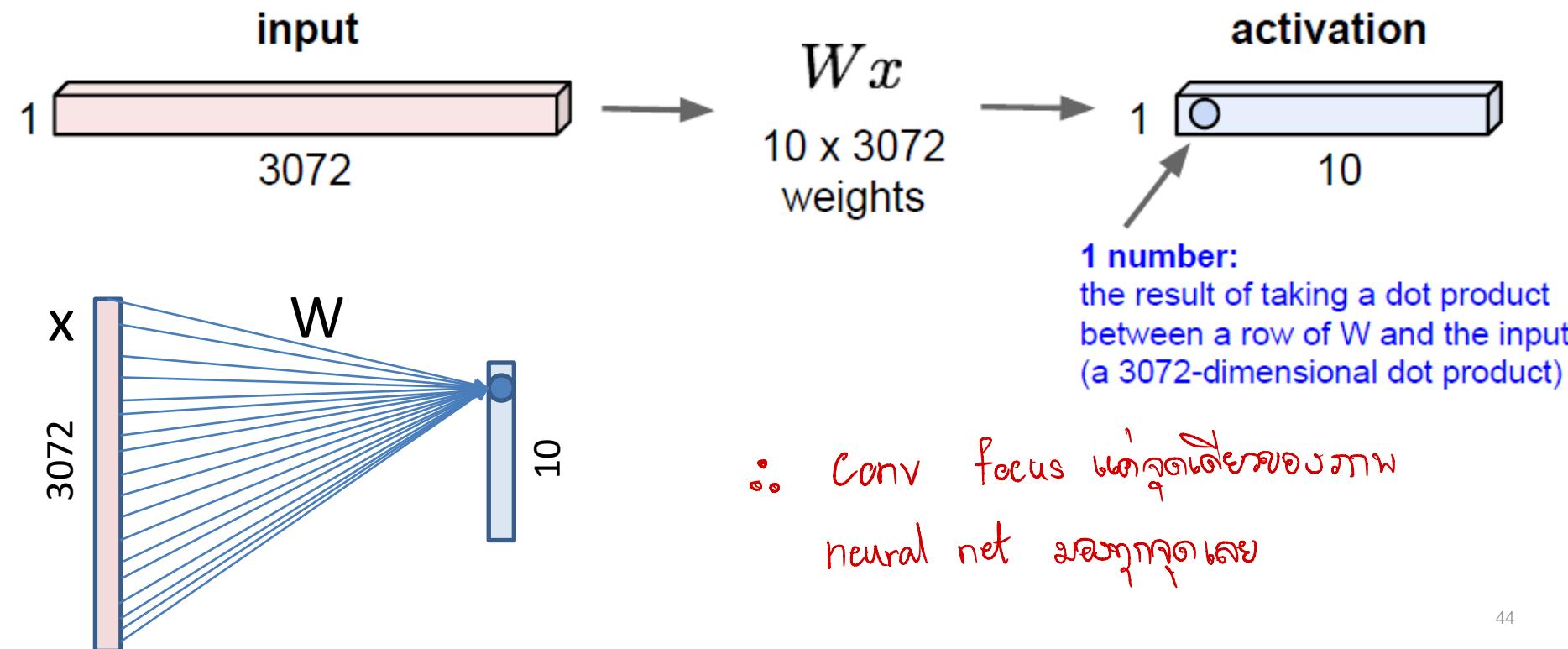
E.g. with 5 filters,
CONV layer consists of
neurons arranged in a 3D grid
($28 \times 28 \times 5$)

There will be 5 different
neurons all looking at the same
region in the input volume

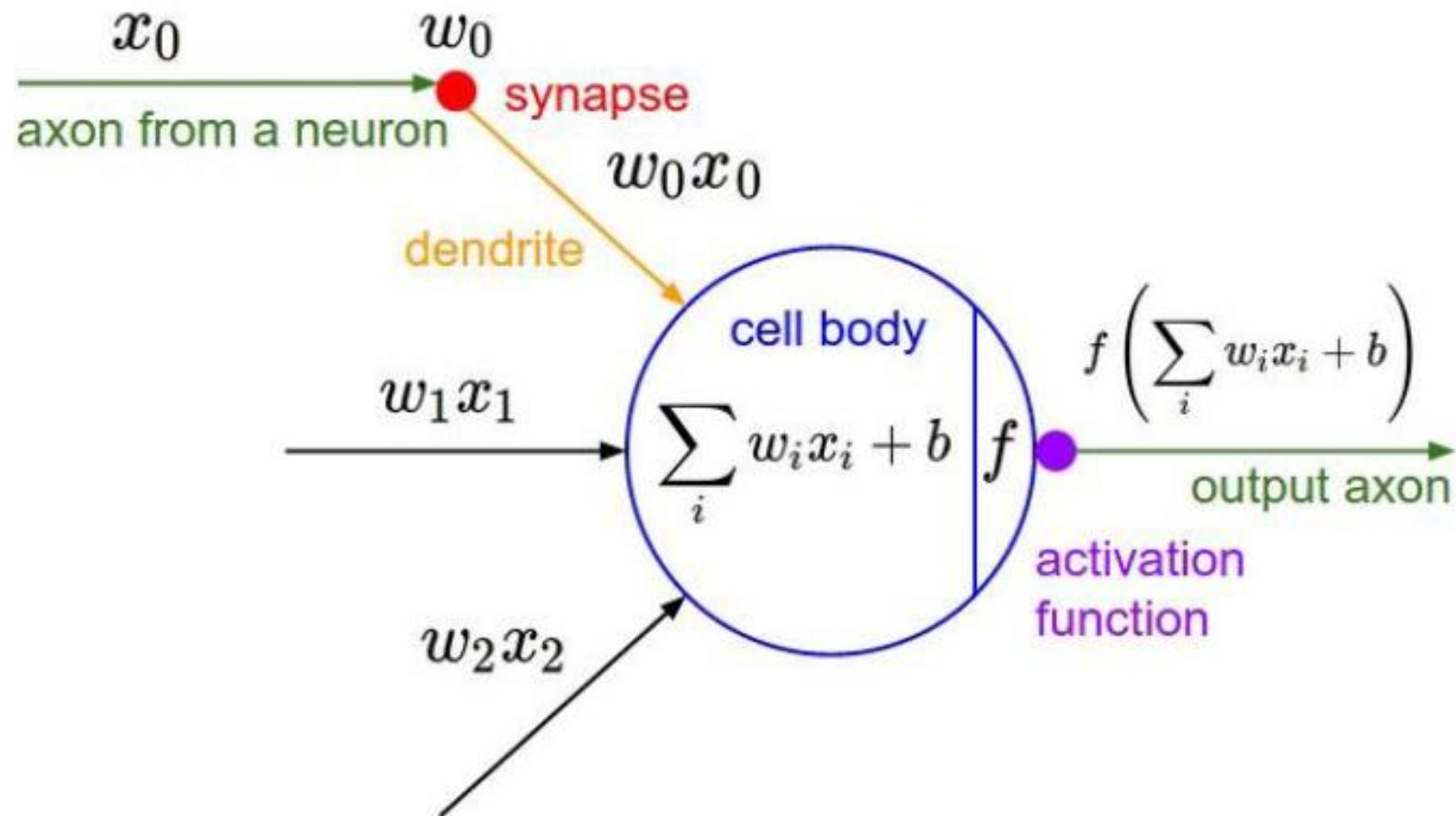
Reminder: Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

Each neuron looks at the full input volume

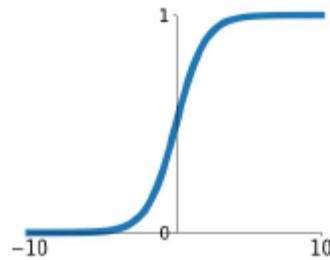


Activation of Convolution Layer



Sigmoid

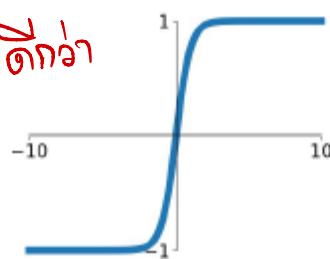
$$\sigma(x) = \frac{1}{1+e^{-x}}$$



tanh

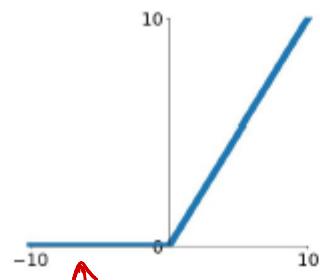
$$\tanh(x)$$

ReLU តាក់
 $\times 6$



ReLU

$$\max(0, x)$$

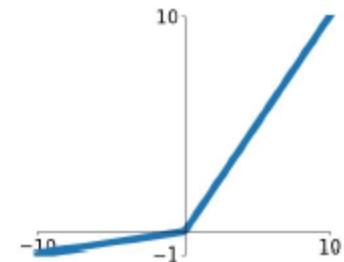


$$\frac{\partial \text{ReLU}}{\partial w} = 0$$

កំណើនត្រូវ
w មិនអែ

Leaky ReLU

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

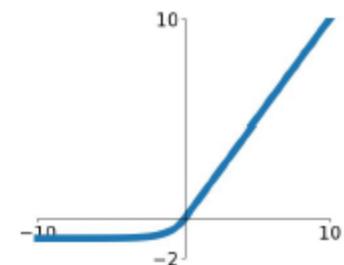


Maxout

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

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

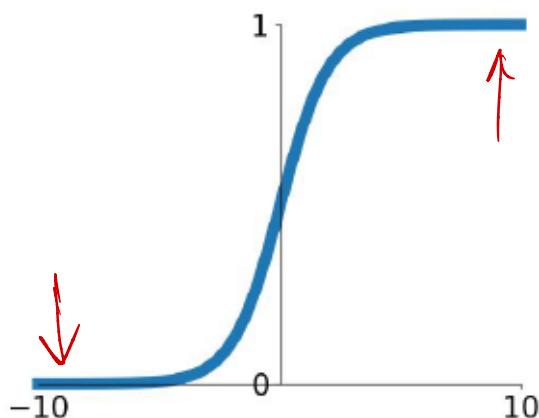


Why not
Sigmoid?

Activation Functions

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

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

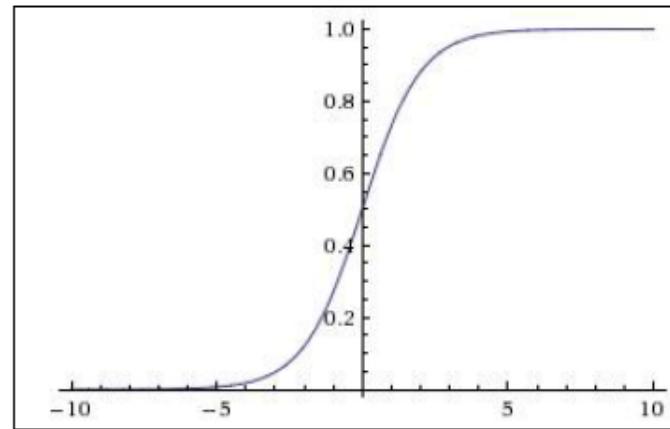
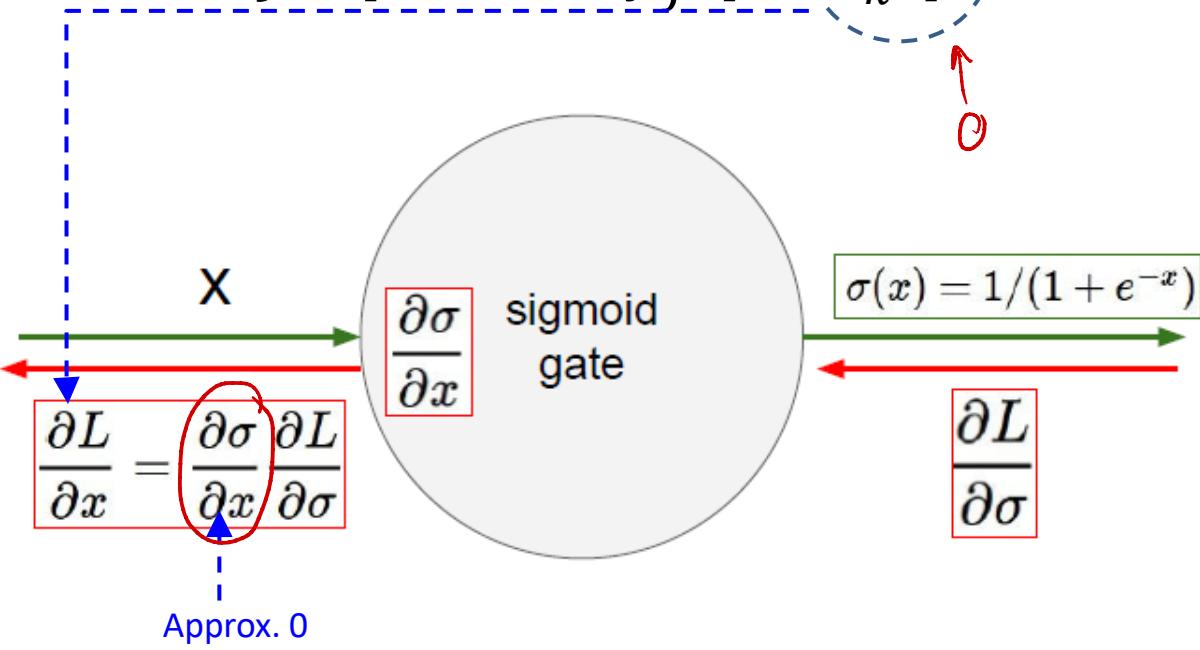


Sigmoid

1. Saturated neurons “kill” the gradients ≈ 0

ກຳອົນຫຼວງ ວ ຖະເຈົ້າ

$$\Delta w_{jk}(p) = \alpha \cdot y_j(p) \cdot \delta_k(p) ; \text{ from back-prop. algorithm}$$

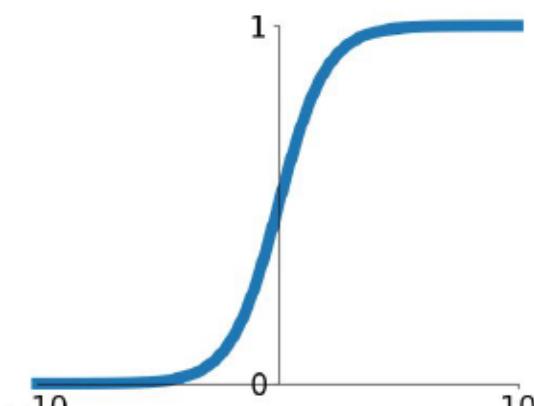


What happens when $x = -10$?

What happens when $x = 0$?

What happens when $x = 10$?

Activation Functions



Sigmoid

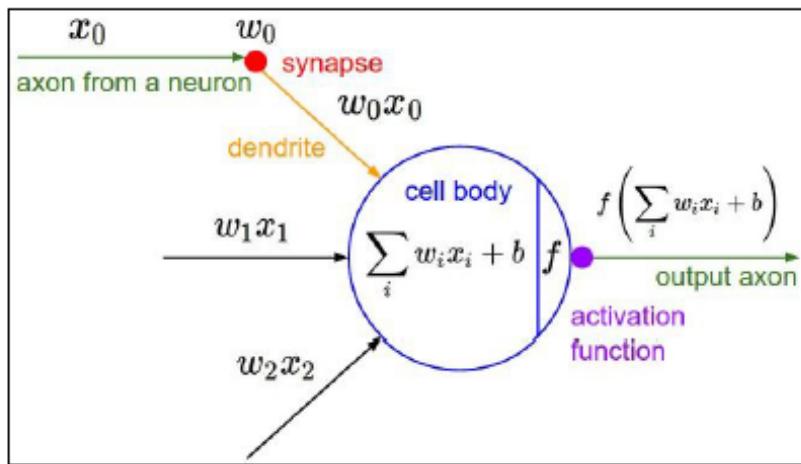
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered

Consider what happens when the input to a neuron (x) is always positive:

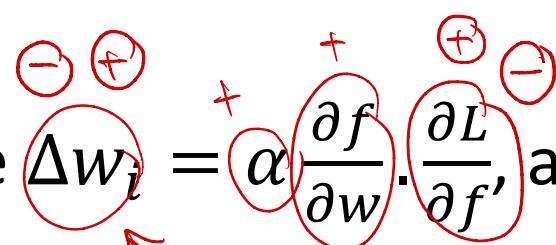


$$f \left(\sum_i w_i x_i + b \right)$$

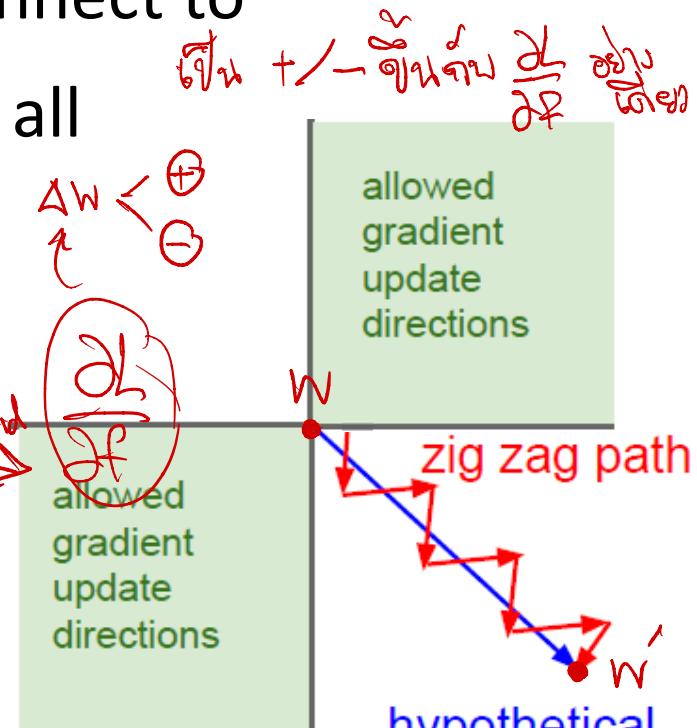
What can we say about the gradients on w ?

$$\frac{\partial f(\sum_i w_i x_i + b)}{\partial w} = \underbrace{\frac{1}{1 + e^{-(\sum_i w_i x_i + b)}}}_{\text{positive}} \times \underbrace{(1 - \frac{1}{1 + e^{-(\sum_i w_i x_i + b)}})}_{\text{positive}}$$

50


 Since $\Delta w_i = \alpha \frac{\partial f}{\partial w} \cdot \frac{\partial L}{\partial f}$, all Δw_i which connect to the same neuron, will be all positive or all negative (depending on the sign of $\frac{\partial L}{\partial f}$.)

$$f \left(\sum_i w_i x_i + b \right)$$



What can we say about the gradients on w ?
 Always all positive or all negative :(
 (this is also why you want zero-mean data!)

For example, suppose we have:

$$w = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \end{pmatrix} = \begin{pmatrix} 1 & -1 & 1 \end{pmatrix}$$

and we need to get it to:

$$w = \begin{pmatrix} -1 & 1 & -1 \end{pmatrix}$$

if we can have both positive and negative gradients we could do this in one step:

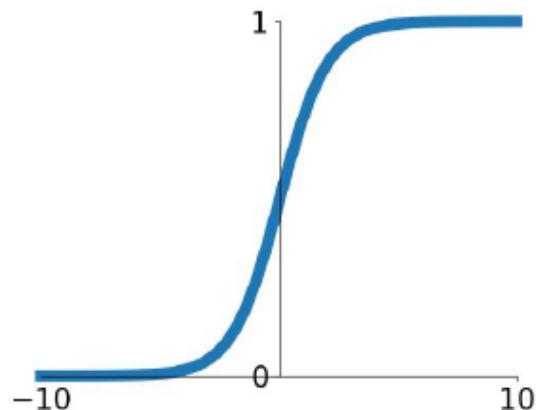
$$w = \begin{pmatrix} 1 & -1 & 1 \end{pmatrix} + \underbrace{\begin{pmatrix} -2 & 2 & -2 \end{pmatrix}}_{\Delta w} = \begin{pmatrix} -1 & 1 & -1 \end{pmatrix} \implies \text{Good. One-step update.}$$

However, the best we can do with only, positive or negative is two steps:

$$w = \begin{pmatrix} 1 & -1 & 1 \end{pmatrix} + \underbrace{\begin{pmatrix} -3 & -3 & -3 \end{pmatrix}}_{\Delta w_1} + \underbrace{\begin{pmatrix} 1 & 5 & 1 \end{pmatrix}}_{\Delta w_2} = \begin{pmatrix} -1 & 1 & -1 \end{pmatrix} \implies \text{Bad. ZIG-ZAG Weight Update}$$

“អ្នកលើនូវវិធានបង្ហាញ និង ការបង្ហាញ នៃសម្រាប់របៀប”

Activation Functions



Sigmoid

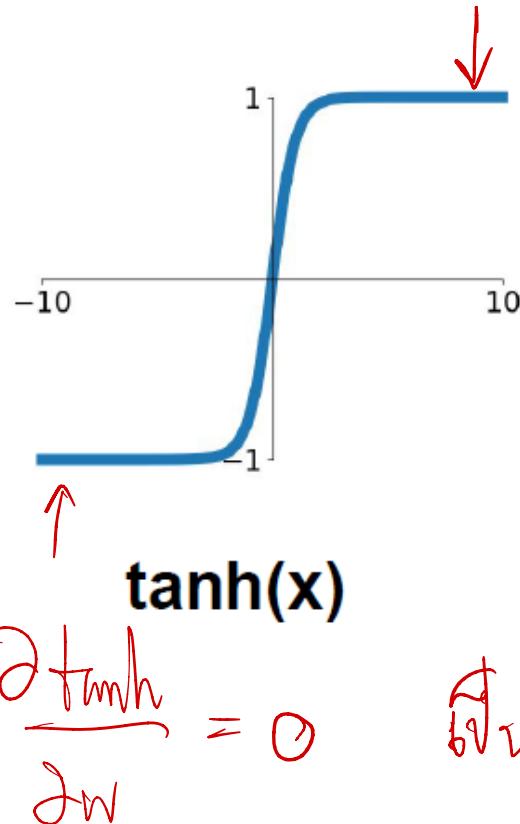
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

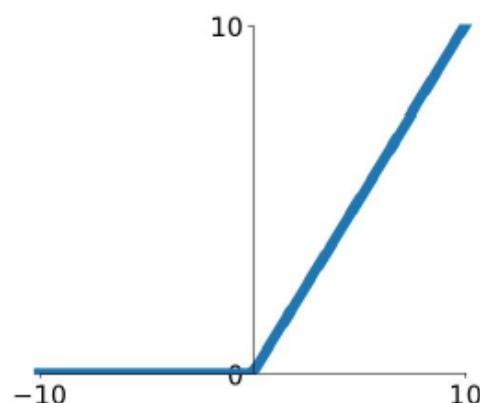
1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered
3. $\exp()$ is a bit compute expensive

Activation Functions



- Squashes numbers to range [-1, 1]
- zero centered (nice) ✓
- still kills gradients when saturated :(

Activation Functions



- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible than sigmoid

ReLU
(Rectified Linear Unit)

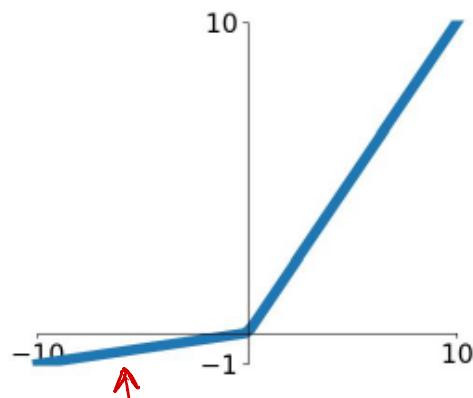
[Krizhevsky et al., 2012]

- Not zero-centered output
- An annoyance:

hint: what is the gradient when $x < 0$?

Activation Functions

[Mass et al., 2013]
[He et al., 2015]



Leaky ReLU

$$f(x) = \max(0.01x, x)$$

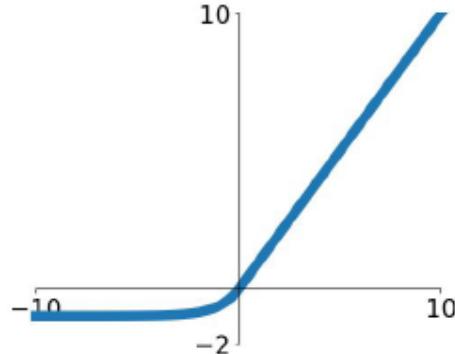
- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

ແຕ່ ກ.ລສັງ ສິພອນ ສະ ReLU

Activation Functions

[Clevert et al., 2015]

Exponential Linear Units (ELU)



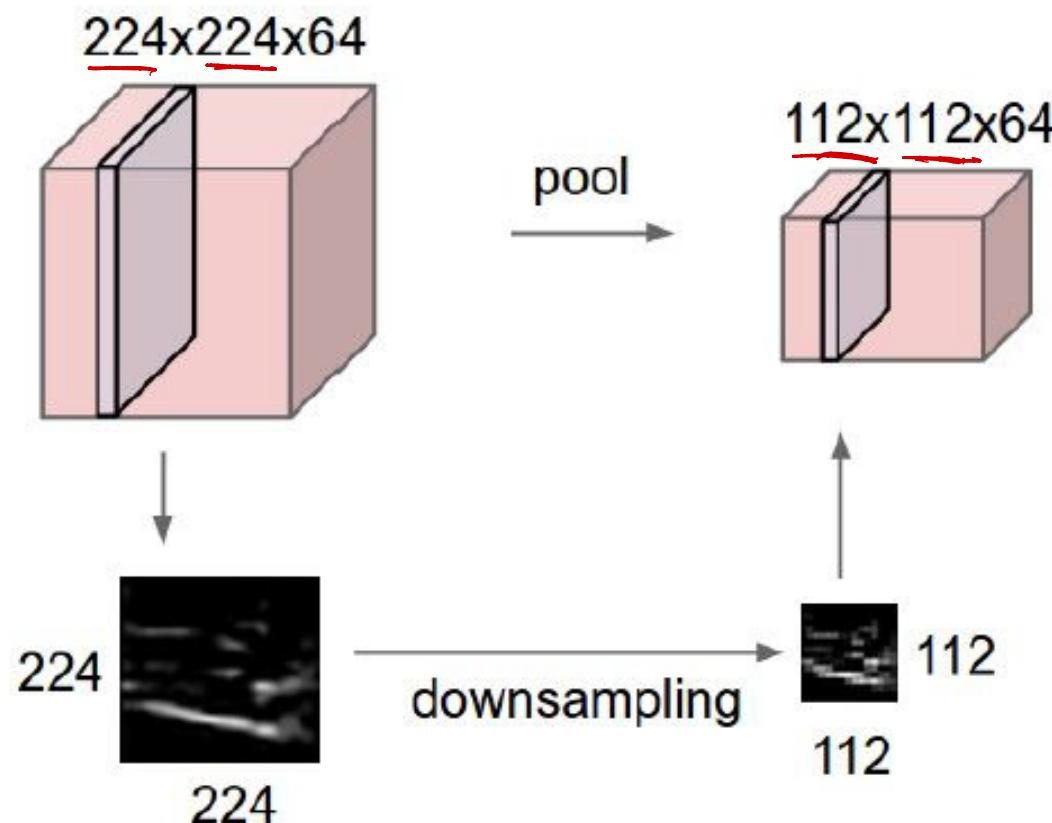
$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

- All benefits of ReLU
- Closer to zero mean outputs
- Negative saturation regime compared with Leaky ReLU adds some robustness to noise

- Computation requires `exp()`

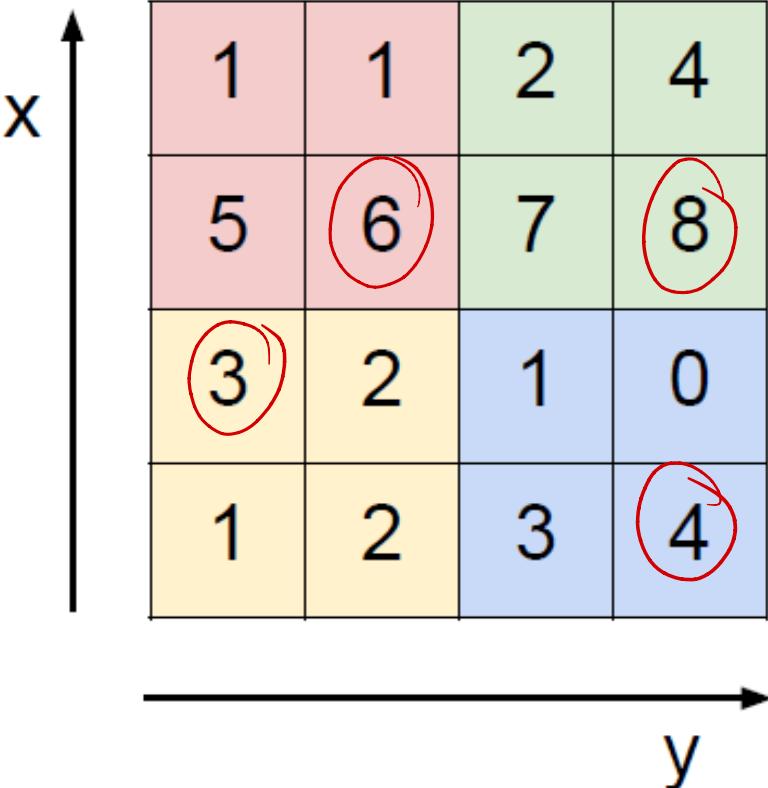
Pooling layer

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

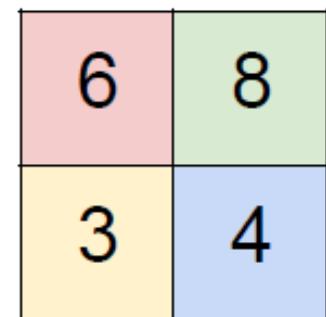


MAX POOLING

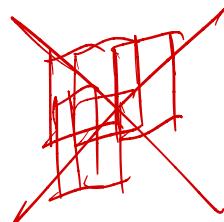
Single depth slice



max pool with 2x2 filters
and stride 2



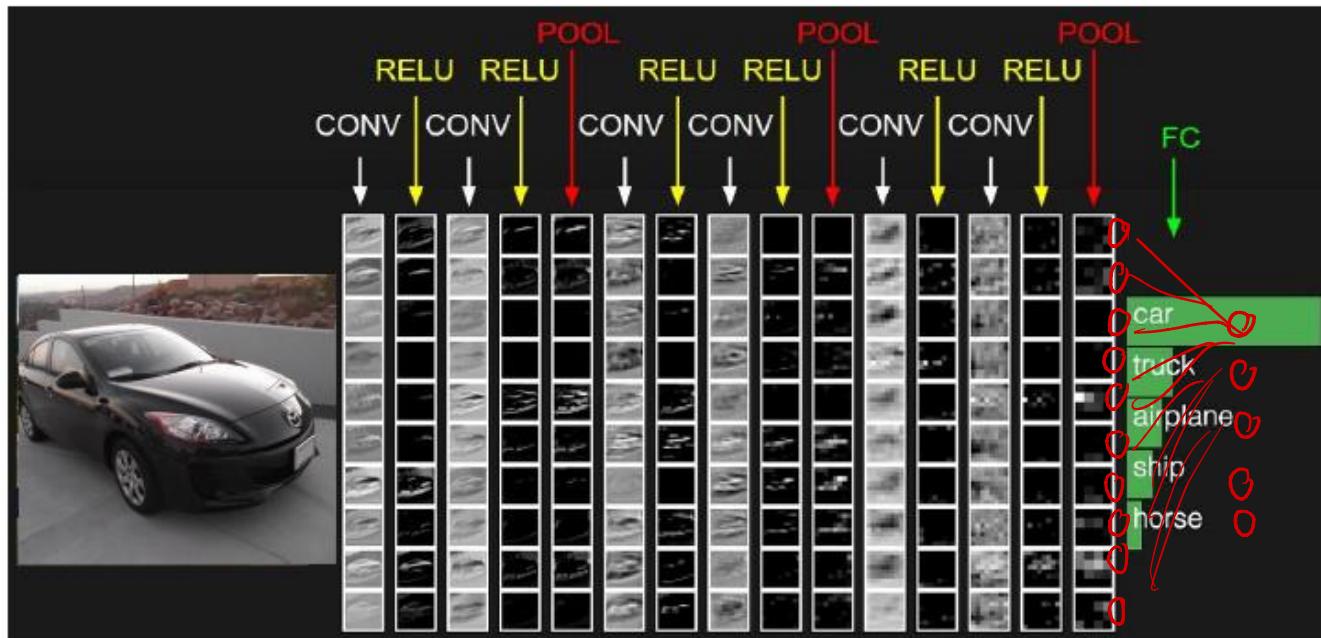
Two ~~ans~~ overlap ~~in~~



- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires two 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

Fully Connected Layer (FC layer)

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



Pre-trained ម៉ាស៊ីន Train បើមែន ចាត់ទូទៅ នៅវគ្គកំណែ
ជានេះ class ខោល