

# *Deep Learning Crash Course*



[www.deeplearningcrashcourse.org](http://www.deeplearningcrashcourse.org)

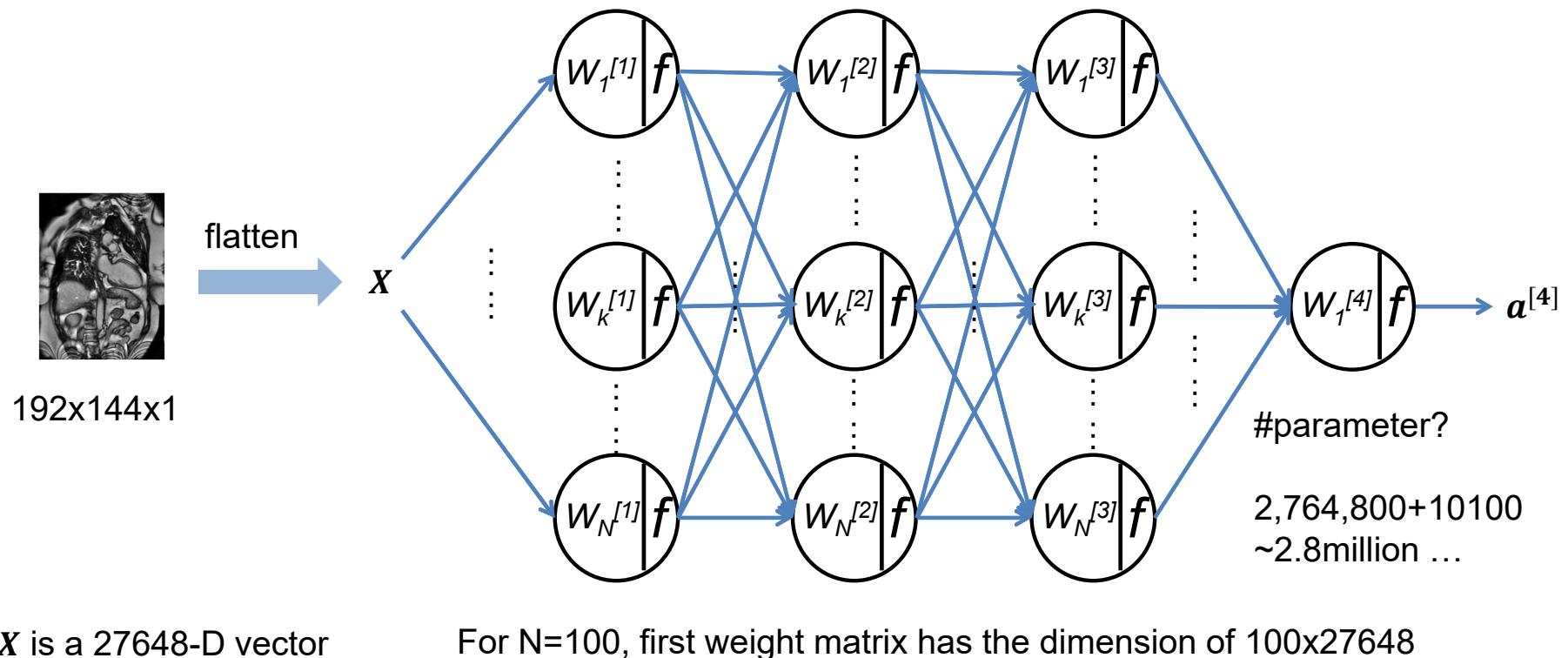
Hui Xue

Fall 2021

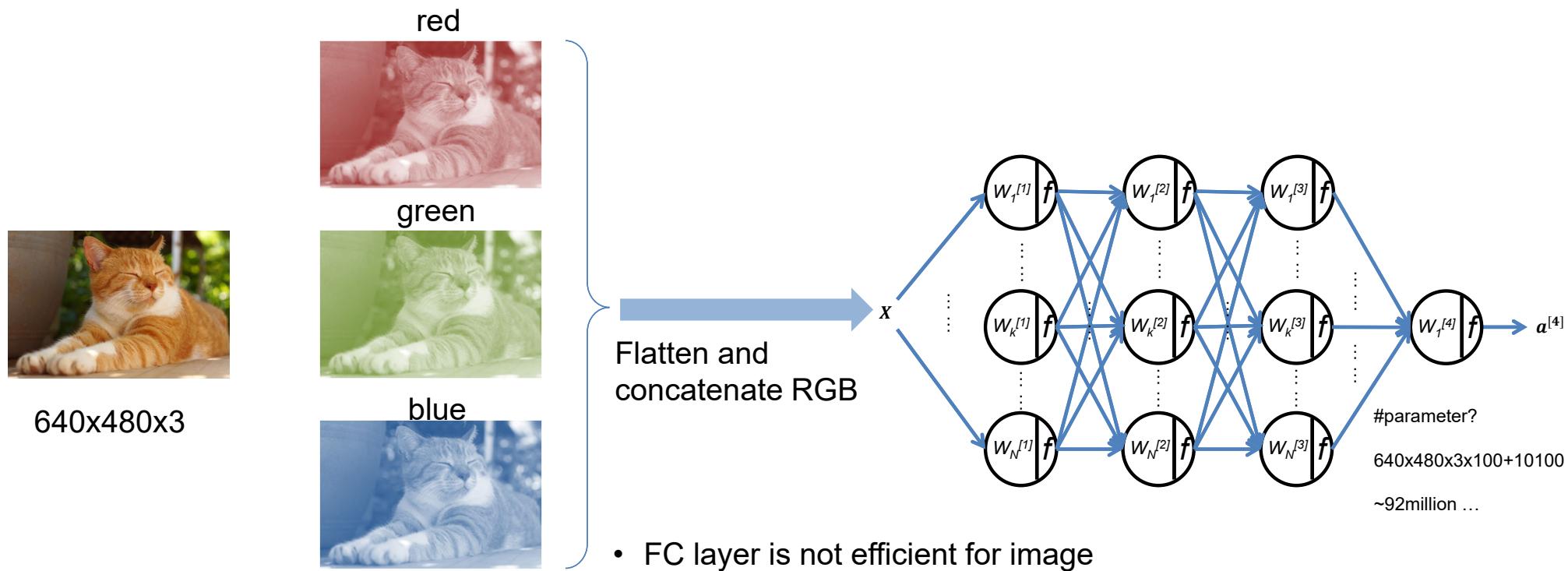
# Outline

- CNN basics
- Variation of convolution
- Batch Normalization
- CNN Examples: LeNet and Alex-Net

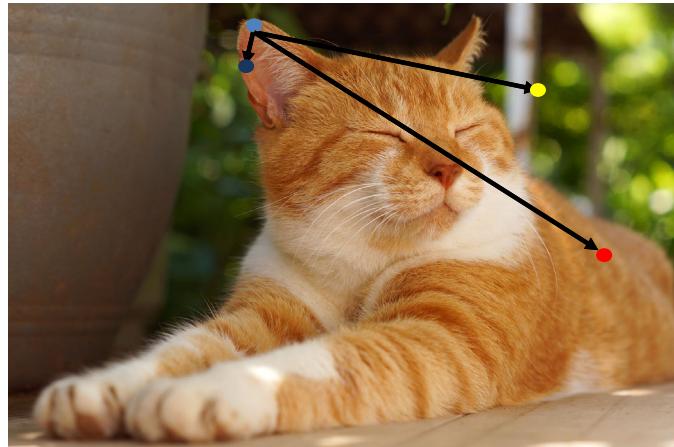
# Revisit image classification



# Revisit image classification



# FC layer is not efficient for image



- FC layer is not efficient for image
- Too many parameters -> hard to train, easy to get overfitting
- Treat all pixels equally; one pixel is connected to all other pixels in the image
- Does not utilize the spatial information among neighboring pixels

# Spatially invariant feature detection



Edge  
detection



1	0	-1
2	0	-2
1	0	-1

Gradient  
along row

1	2	1
0	0	0
-1	-2	-1

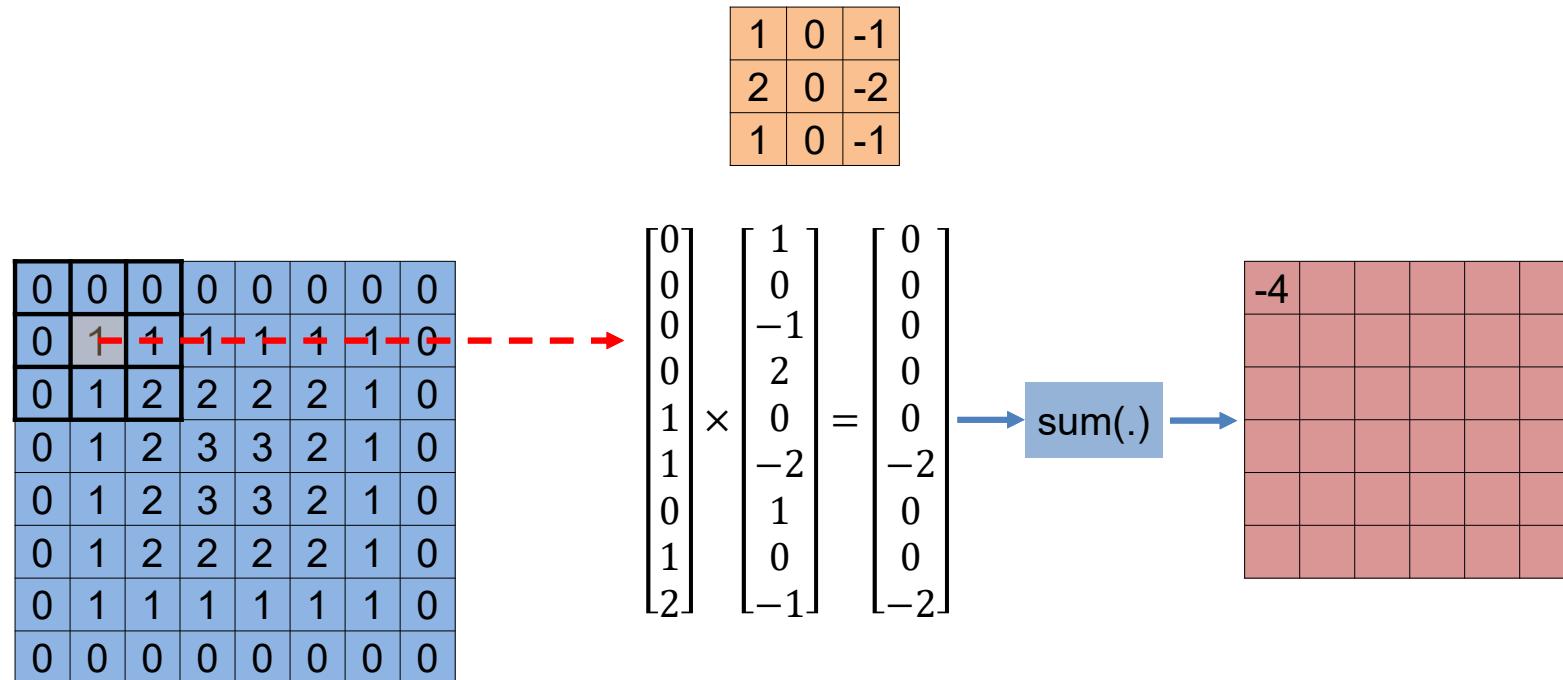
Gradient  
along column

Sobel  
filter

- Use the same filter “sliding over” the image
- Convolution
- No need to have different detectors for the same image features → FC method is not optimal

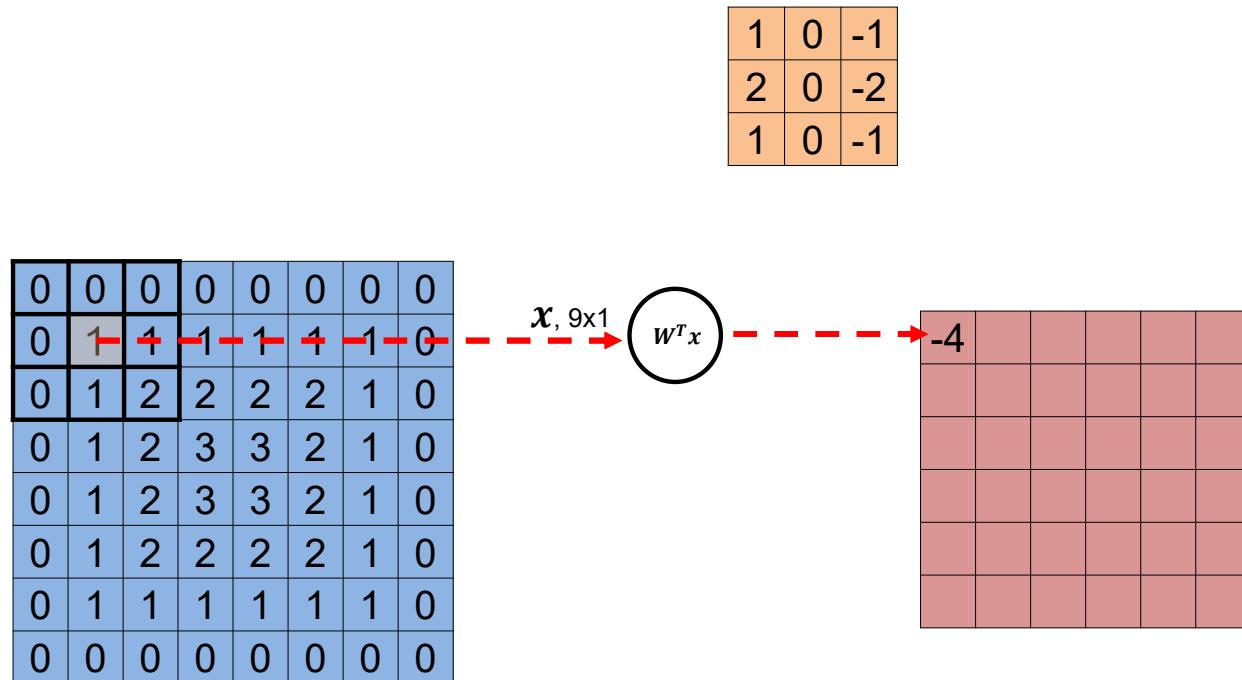
# CNN is a spatial-invariant, parameter-sharing FC

- Use the convolution operation for spatial invariant → *spatial invariant*
- Use the same filter for the whole image → *parameter sharing*



# CNN is a spatial-invariant, parameter-sharing FC

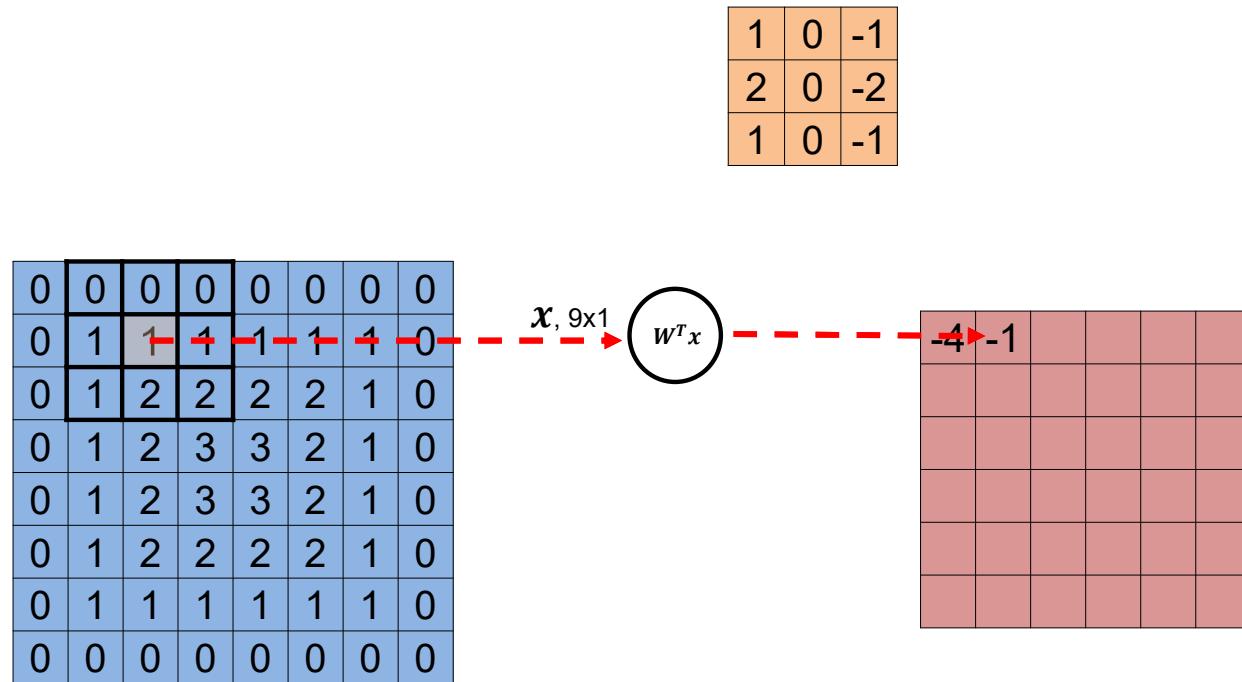
- Use the convolution operation for spatial invariant → spatial invariant
- Use the same filter for the whole image → parameter sharing



- Take a neighborhood and flatten
- Multiple with kernel and sum it up
- Can be seen as a local FC

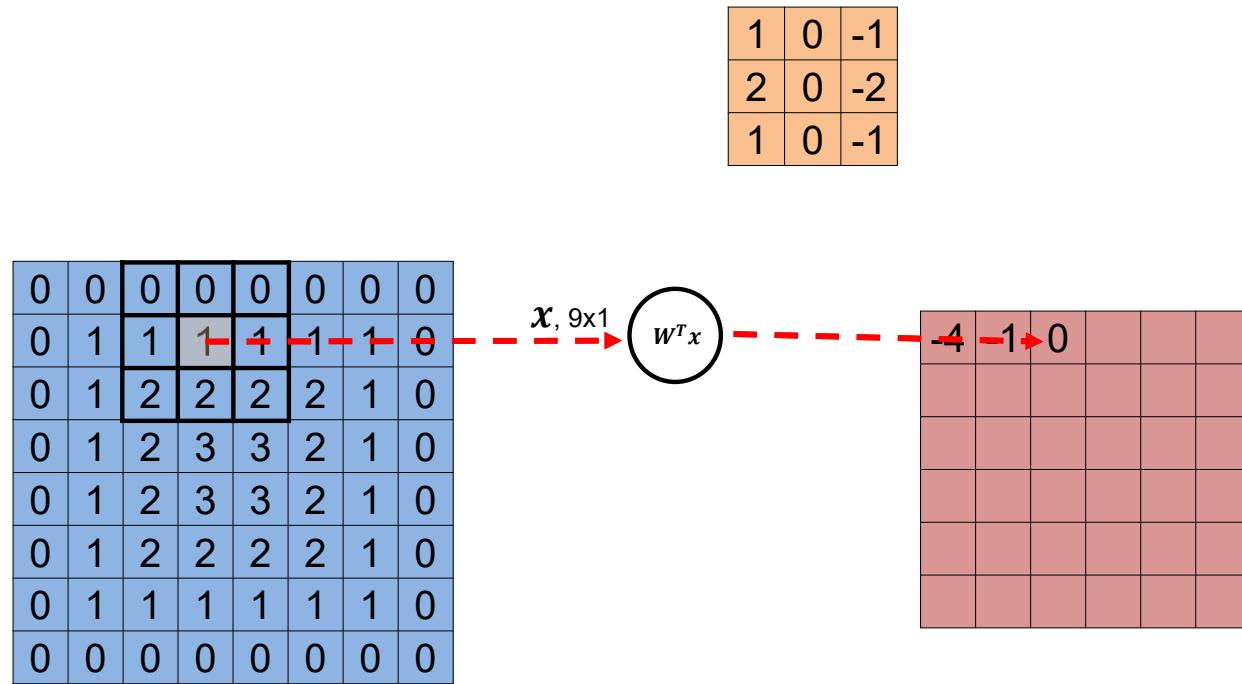
# CNN is a spatial-invariant, parameter-sharing FC

- Use the convolution operation for spatial invariant → spatial invariant
- Use the same filter for the whole image → parameter sharing



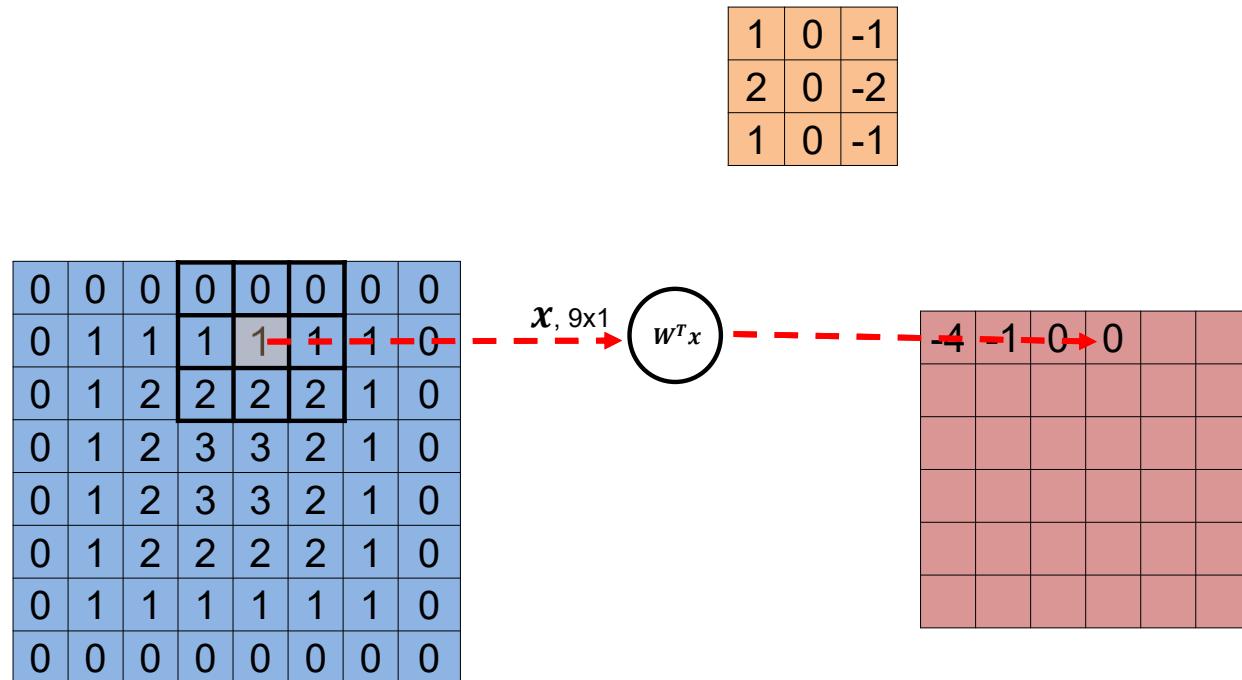
# CNN is a spatial-invariant, parameter-sharing FC

- Use the convolution operation for spatial invariant → spatial invariant
- Use the same filter for the whole image → parameter sharing



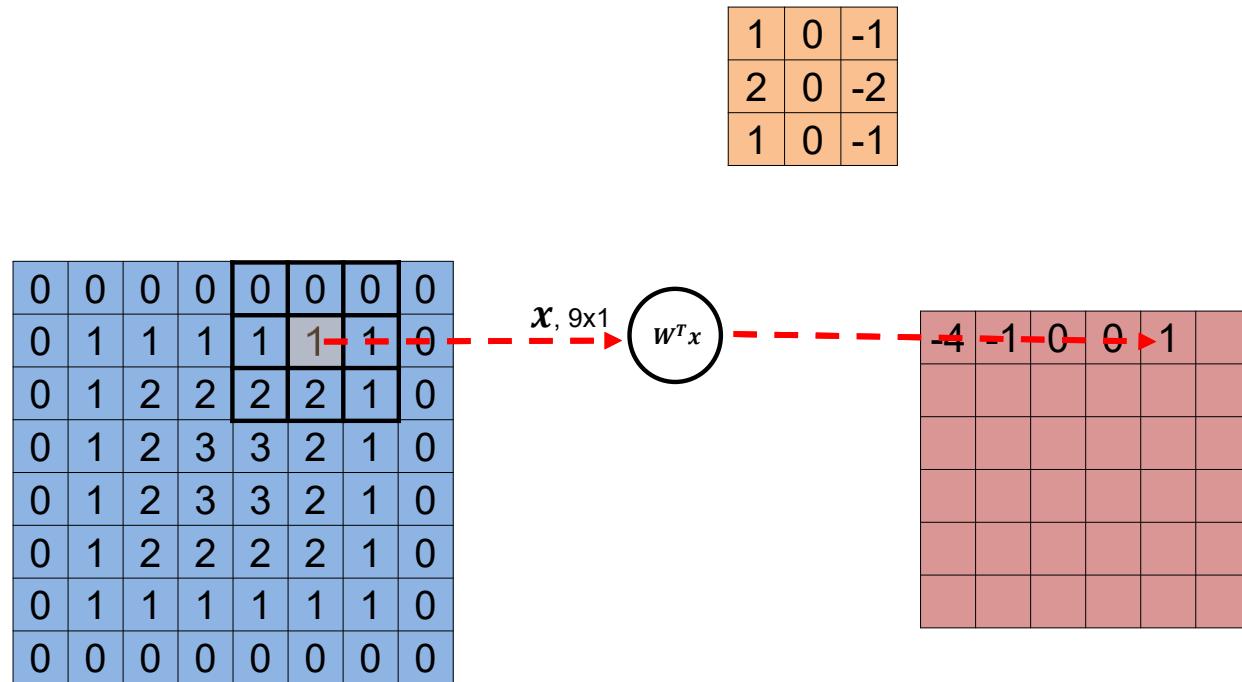
# CNN is a spatial-invariant, parameter-sharing FC

- Use the convolution operation for spatial invariant → spatial invariant
- Use the same filter for the whole image → parameter sharing



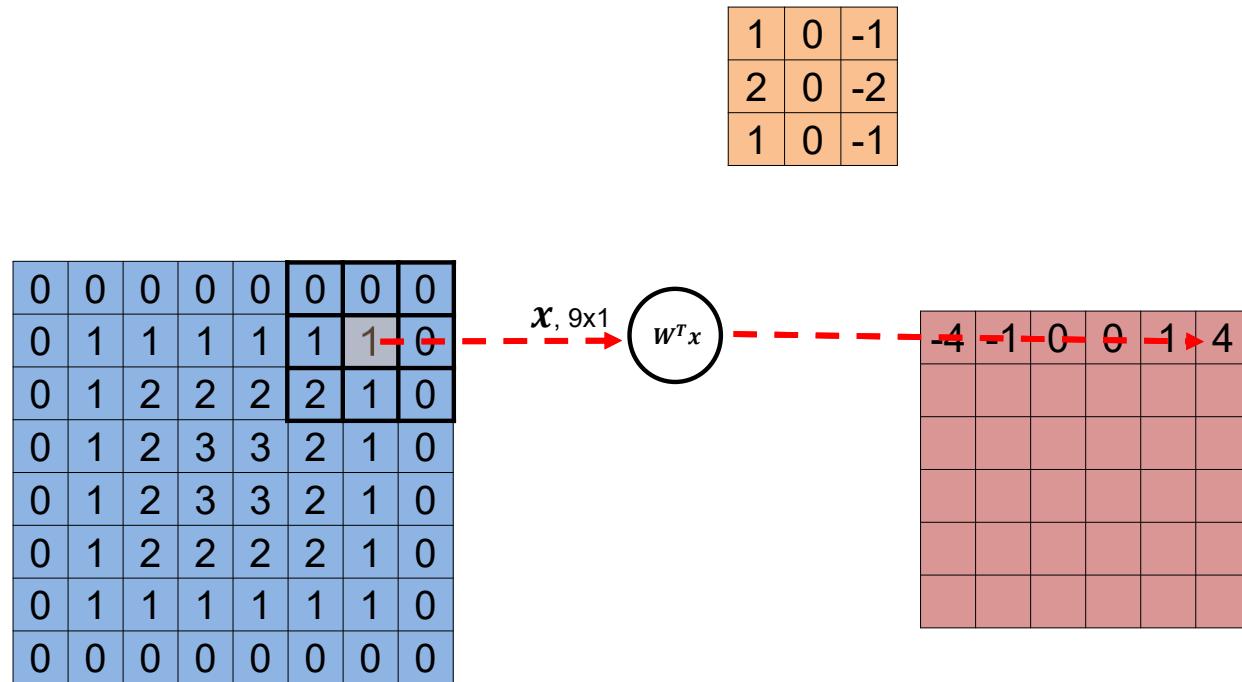
# CNN is a spatial-invariant, parameter-sharing FC

- Use the convolution operation for spatial invariant → spatial invariant
- Use the same filter for the whole image → parameter sharing



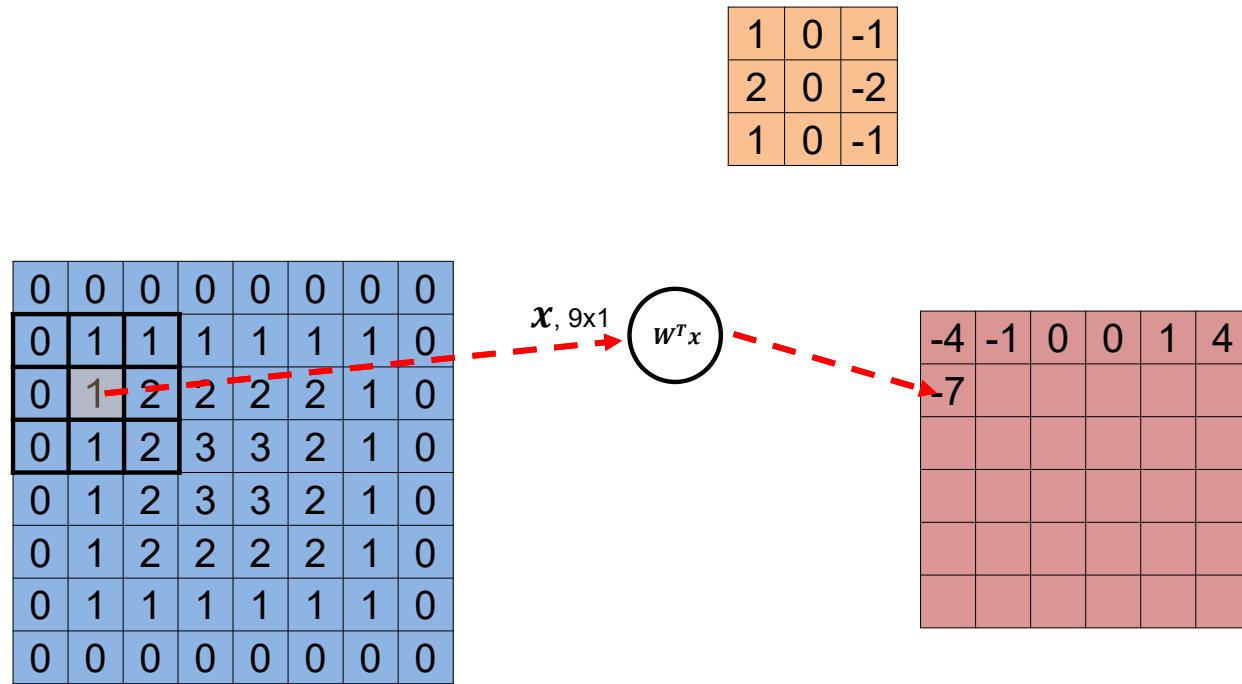
# CNN is a spatial-invariant, parameter-sharing FC

- Use the convolution operation for spatial invariant → spatial invariant
- Use the same filter for the whole image → parameter sharing



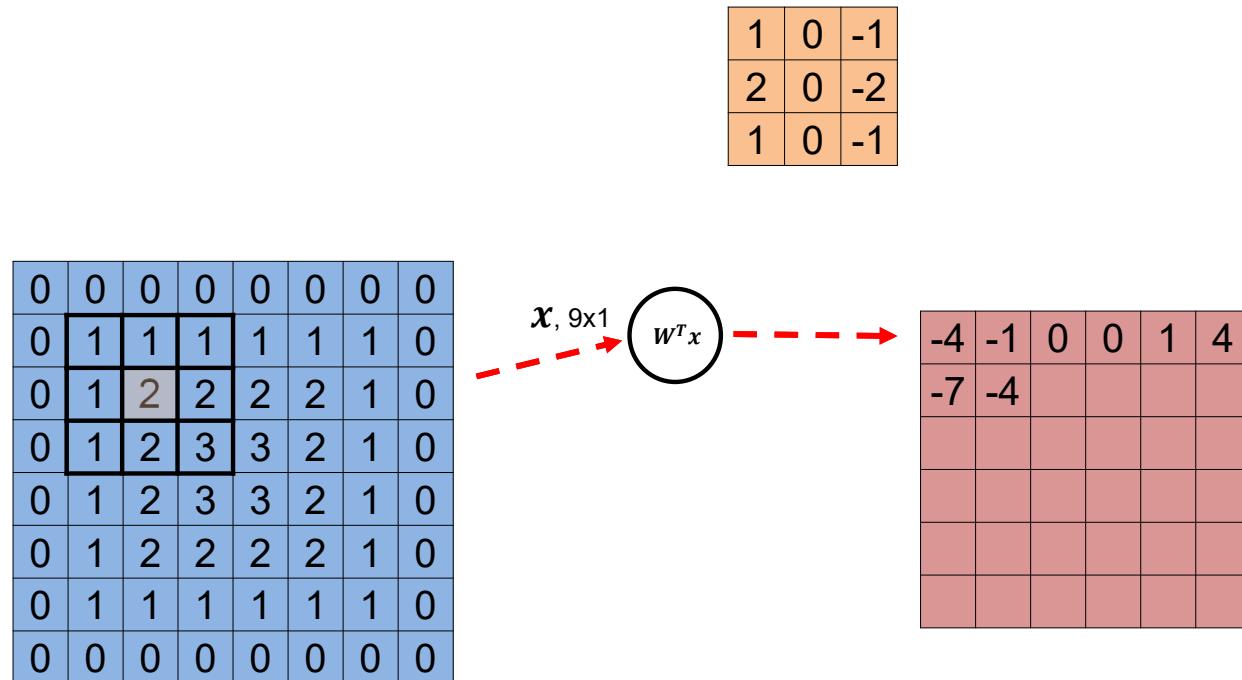
# CNN is a spatial-invariant, parameter-sharing FC

- Use the convolution operation for spatial invariant → spatial invariant
- Use the same filter for the whole image → parameter sharing



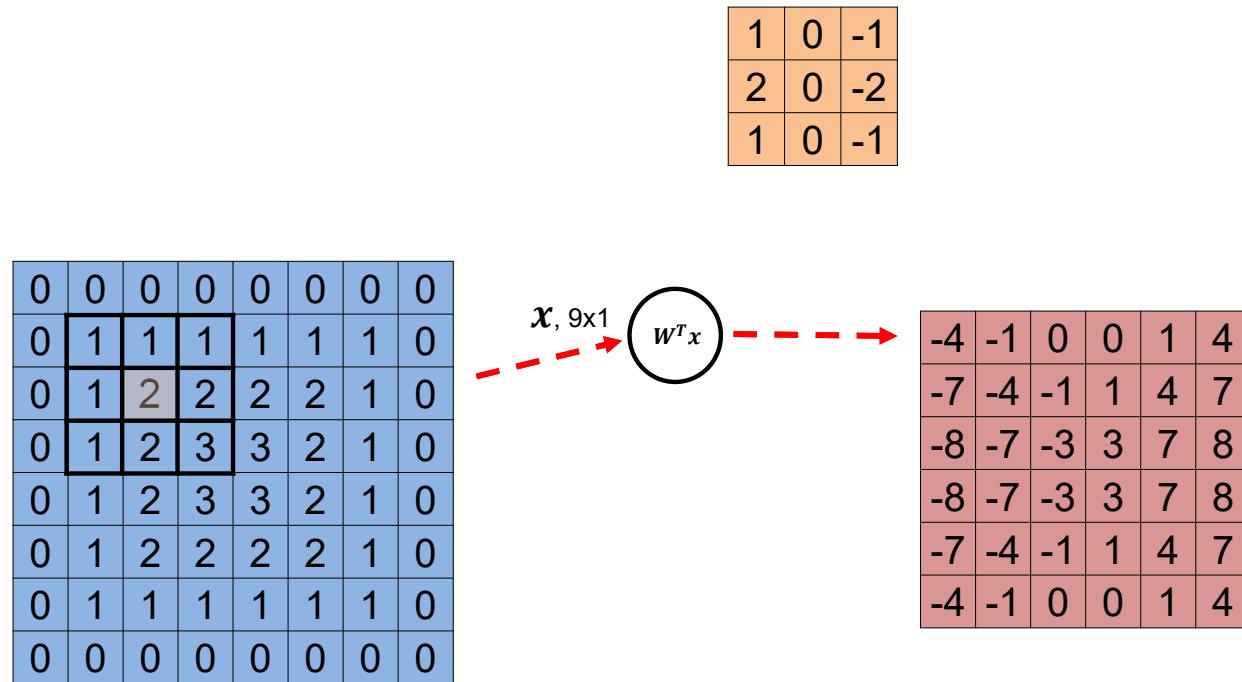
# CNN is a spatial-invariant, parameter-sharing FC

- Use the convolution operation for spatial invariant → spatial invariant
- Use the same filter for the whole image → parameter sharing



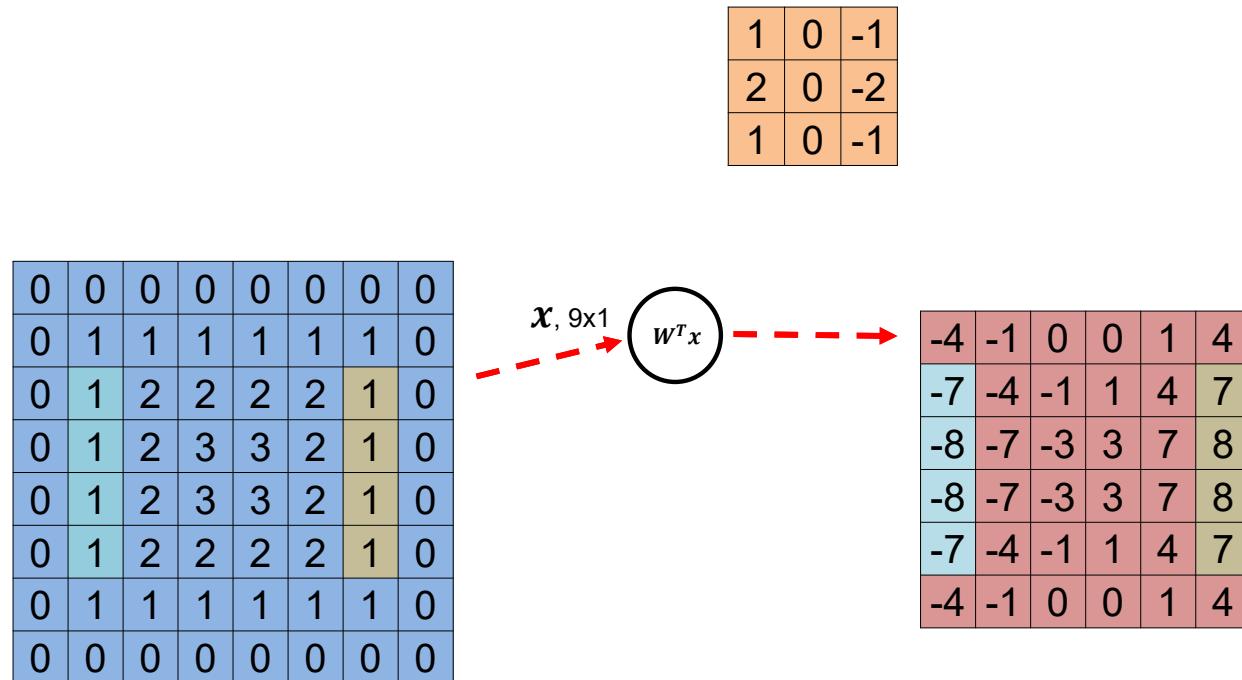
# CNN is a spatial-invariant, parameter-sharing FC

- Use the convolution operation for spatial invariant → spatial invariant
- Use the same filter for the whole image → parameter sharing



# CNN is a spatial-invariant, parameter-sharing FC

- Use the convolution operation for spatial invariant → spatial invariant
- Use the same filter for the whole image → parameter sharing



- Higher value in outputs indicate “vertical edges”
- With direction information (gradient direction)

# CNN layer with multiple channels

- Kernel has multiple channels too
- Sum over all channels

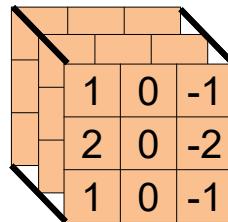
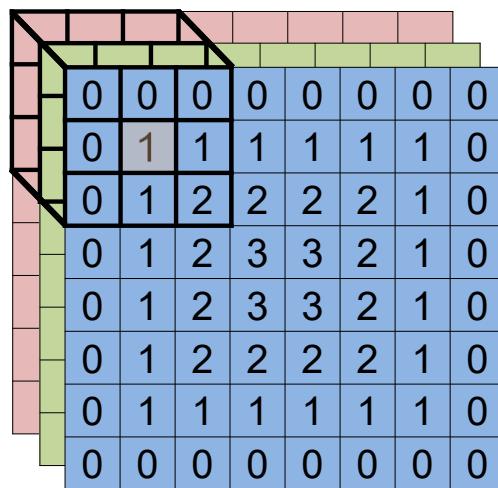
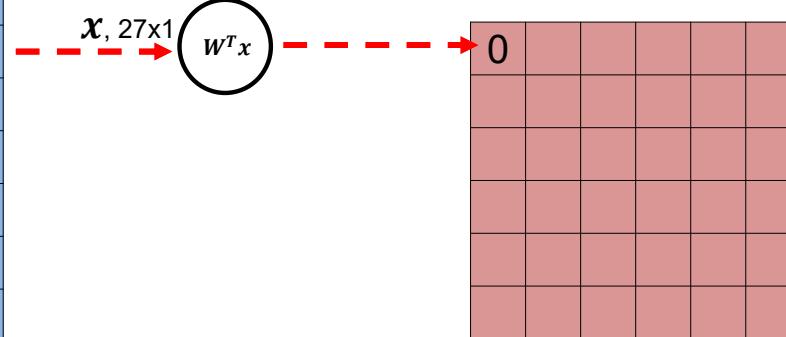


Image:  $H \times W \times C$   
Kernel:  $k_x \times k_y \times C$



# Padding

- Pixels on the boundary was not explored
- Reduced output size

Padding = 1

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	1	1	1	1	1	1	0	0
0	0	1	2	2	2	2	1	0	0
0	0	1	2	3	3	2	1	0	0
0	0	1	2	3	3	2	1	0	0
0	0	1	2	2	2	2	1	0	0
0	0	1	1	1	1	1	1	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

1	0	-1
2	0	-2
1	0	-1



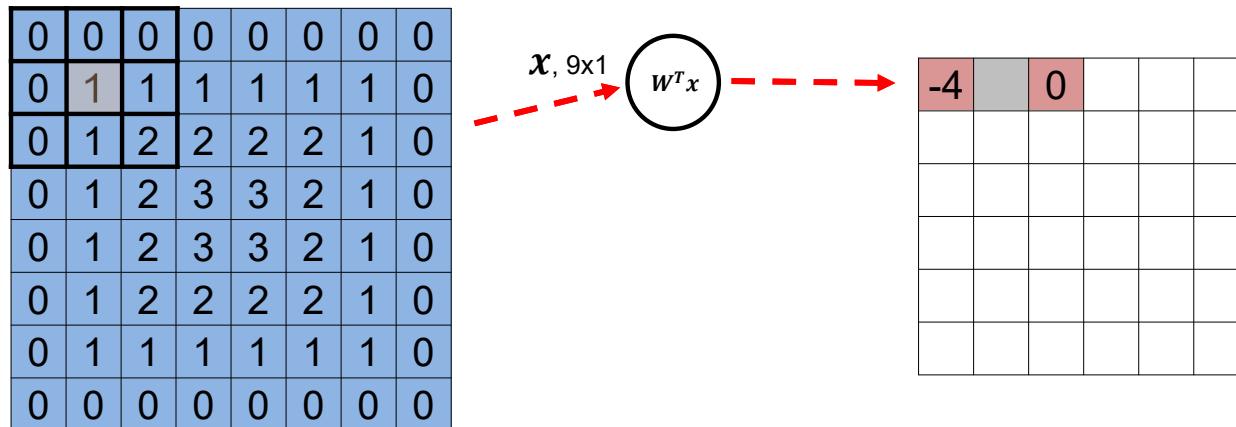
-1	-1	0	0	0	0	0	1	1
-3	-4	-1	0	0	1	4	-3	
-4	-7	-4	-1	1	4	7	-4	
-4	-8	-7	-3	3	7	8	-4	
-4	-8	-7	-3	3	7	8	-4	
-4	-7	-4	-1	1	4	7	-4	
-3	-4	-1	0	0	1	4	-3	
-1	-1	0	0	0	0	1	1	

- Often pad with 0
- Other choices : replica, mirror, periodic

# Stride

- When sliding the kernel over the image, the moving step-size is the stride
- Default is 1, but stride can be larger than 1
- Often used to reduce output size

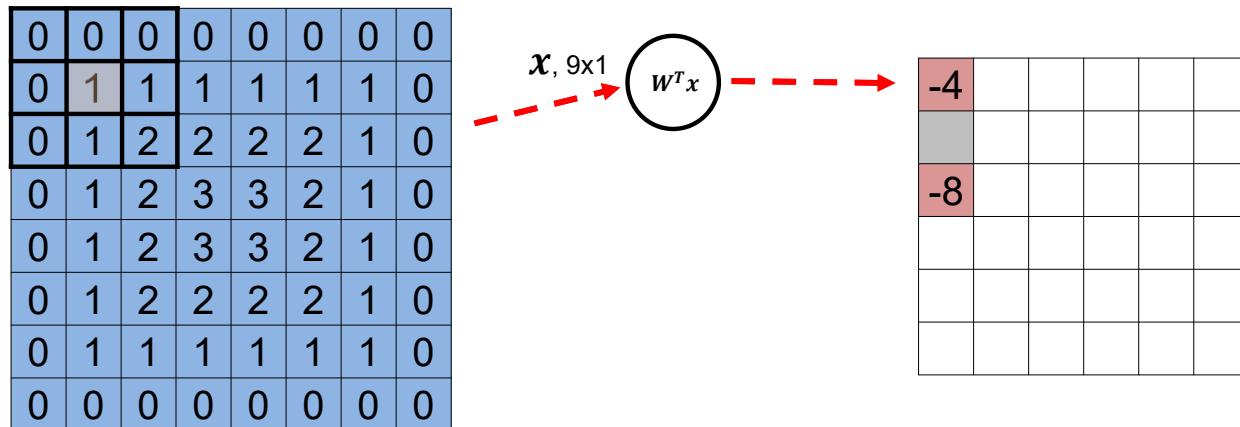
Stride=2 along row



# Stride

- When sliding the kernel over the image, the moving stepsize is the stride
- Default is 1, but stride can be larger than 1
- Often used to reduce output size

Stride=2 along column



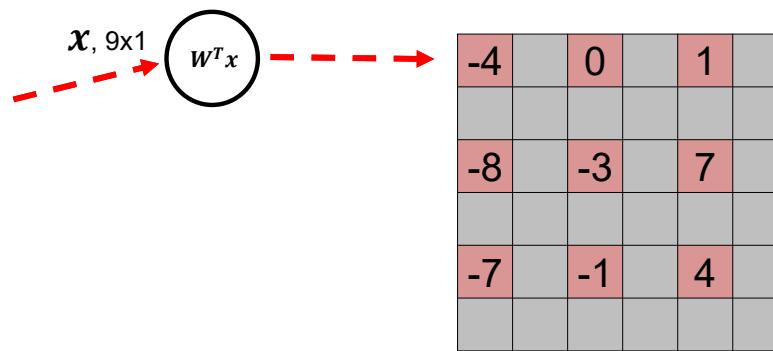
# Stride

- When sliding the kernel over the image, the moving stepsize is the stride
- Default is 1, but stride can be larger than 1
- Often used to reduce output size

Input 8x8  
Kernel 3x3  
Output 3x3

With padding=0, stride=2 for H and W

0	0	0	0	0	0	0	0	0
0	1	1	1	1	1	1	1	0
0	1	2	2	2	2	2	1	0
0	1	2	3	3	2	1	0	0
0	1	2	3	3	2	1	0	0
0	1	2	2	2	2	1	0	0
0	1	1	1	1	1	1	0	0
0	0	0	0	0	0	0	0	0



# Determine the output size

Image:  $H \times W \times C$

Kernel:  $k_x \times k_y \times C$

Padding size:  $p$

Stride:  $s$

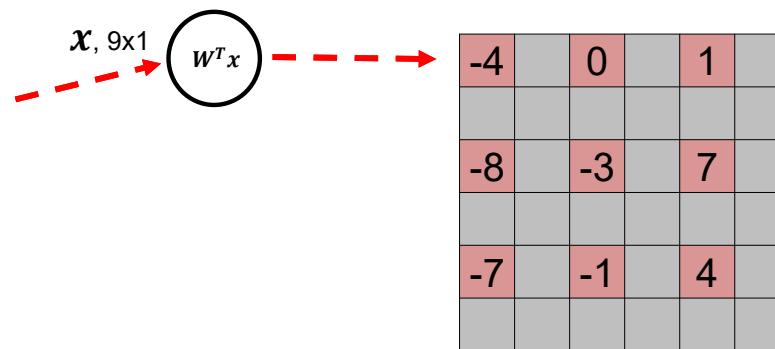
Output size:

$$H' = \left\lceil \frac{H + 2 \times p - k_x}{s} \right\rceil + 1$$

$$W' = \left\lceil \frac{W + 2 \times p - k_y}{s} \right\rceil + 1$$

$$\begin{aligned} & \left\lceil \frac{8 + 2 \times 0 - 3}{2} \right\rceil + 1 \\ &= \lceil 2.5 \rceil + 1 = 3 \end{aligned}$$

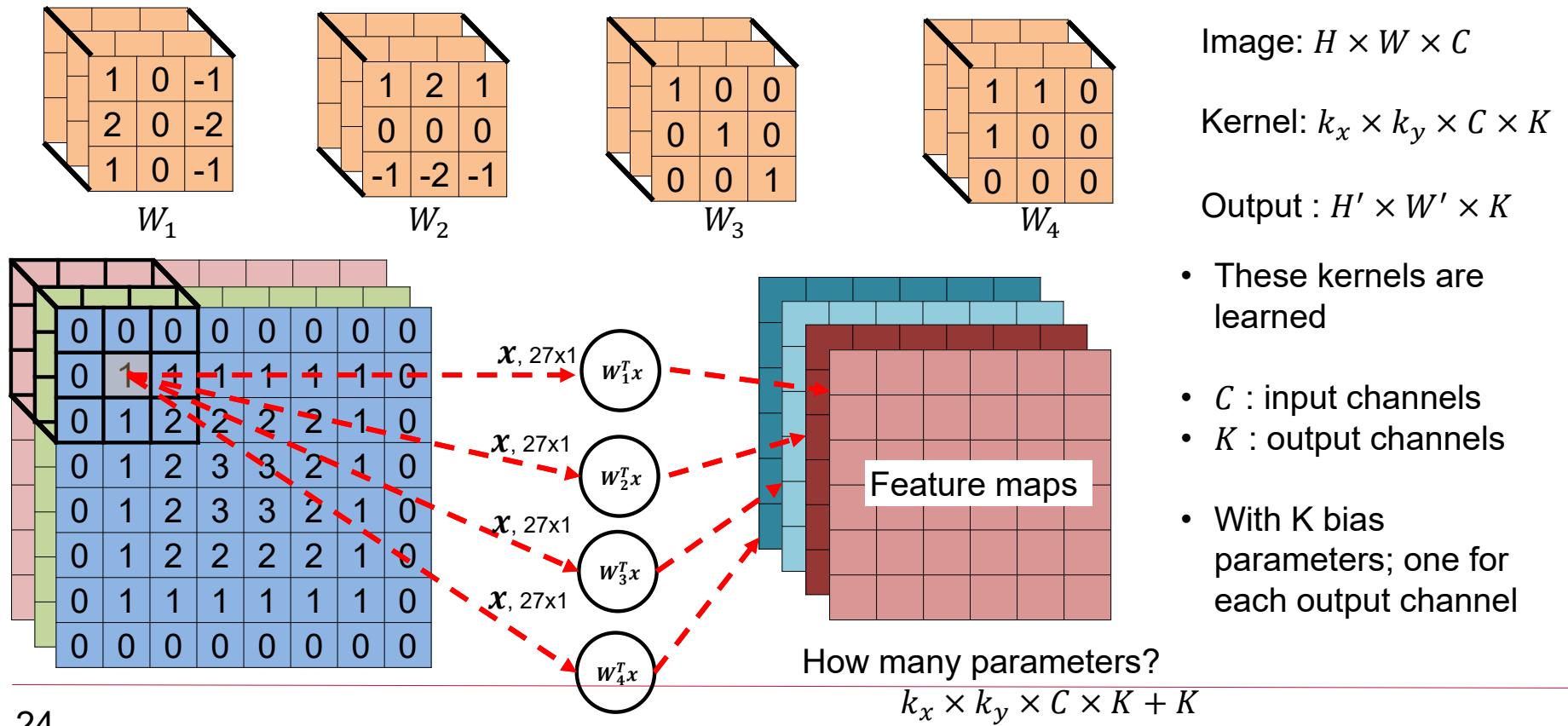
0	0	0	0	0	0	0	0	0
0	1	1	1	1	1	1	1	0
0	1	2	2	2	2	1	0	
0	1	2	3	3	2	1	0	
0	1	2	3	3	2	1	0	
0	1	2	2	2	2	1	0	
0	1	1	1	1	1	1	0	
0	0	0	0	0	0	0	0	0



- Stride can be used to reduce image size

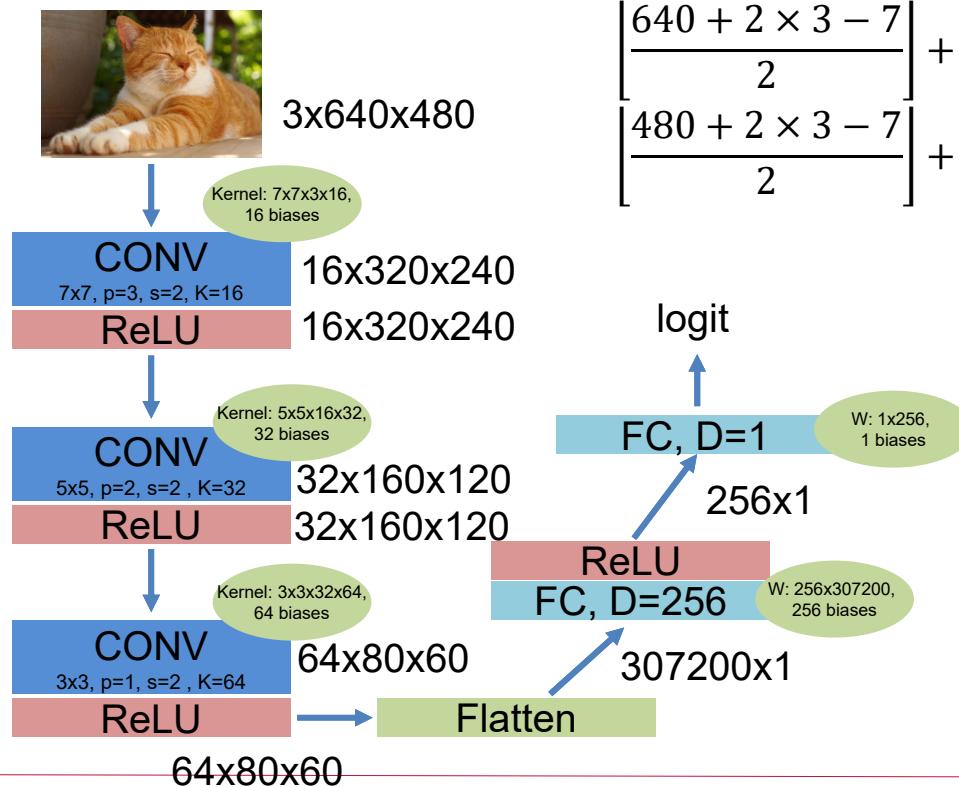
# Multiple kernels

- For one location in the image, multiple kernels can be used to capture different features



# Multiple CONV layers can be stacked

- Nonlinearity (e.g. ReLU) is applied to every element on feature maps



$$\begin{aligned} \frac{640 + 2 \times 3 - 7}{2} + 1 &= 320 \\ \frac{480 + 2 \times 3 - 7}{2} + 1 &= 240 \end{aligned}$$

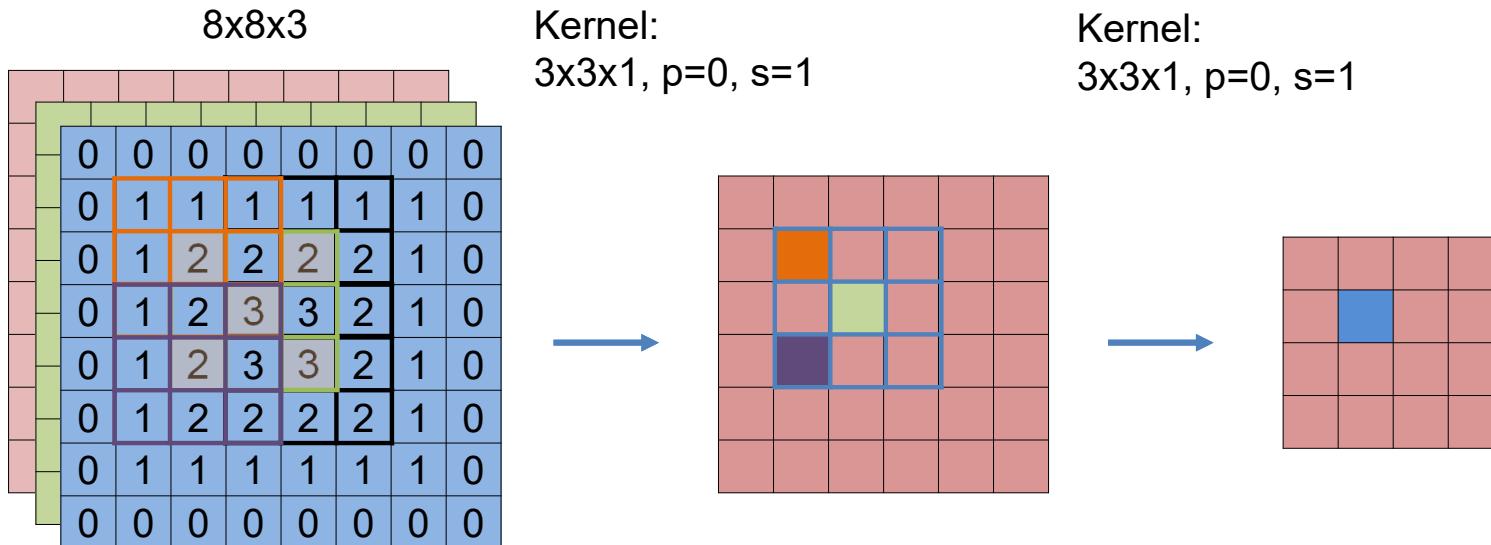
Total number of parameters:

$$\begin{aligned} &7 \times 7 \times 3 \times 16 + 16 + 5 \times 5 \times 16 \times 32 + 32 + 3 \\ &x 3 \times 32 \times 64 + 64 + 256 \times 307200 + 256 \\ &+ 1 \times 256 + 1 = 78,677,409, \\ &\sim 78 \text{million} \end{aligned}$$

One FC layer claims 99.96%

# Receptive field

- Every CONV filter looks at a small neighborhood  
May need to look at more global structure in the image?
- Receptive field : the region in the image to impact a feature after some CONV layers



- With deeper network, the receptive field increases

# 1D and 3D convolution

- Computation is the same as the 2D CONV

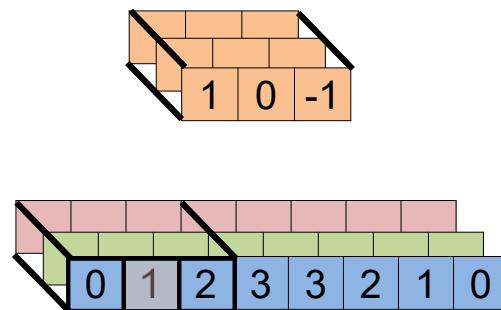


Image:  $L \times C$

Kernel:  $k \times C \times K$

Output :  $L' \times K$

Padding along two ends

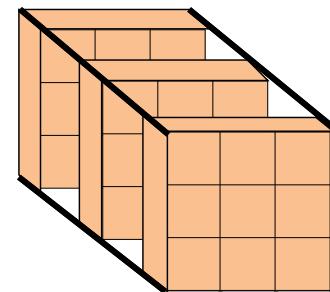
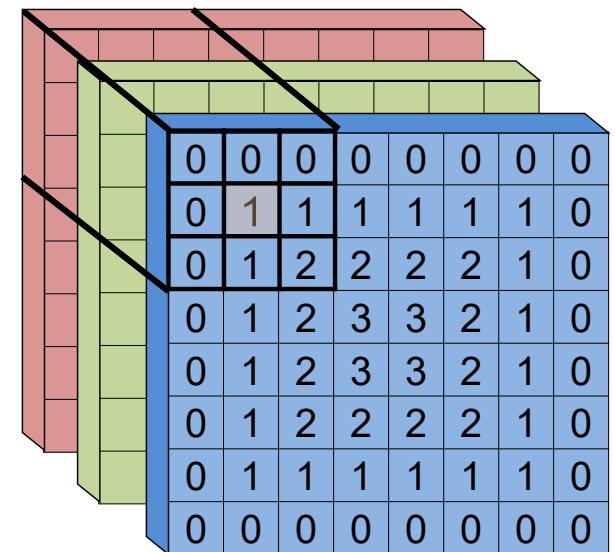


Image:  $H \times W \times D \times C$

Kernel:  $k_x \times k_y \times k_z \times C \times K$

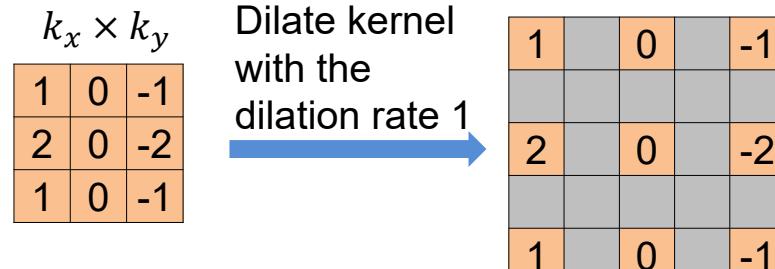
Output :  $H' \times W' \times D' \times K$

Padding along three dimensions



# Dilated convolution

- Increase receptive field by “skipping pixels” in input images



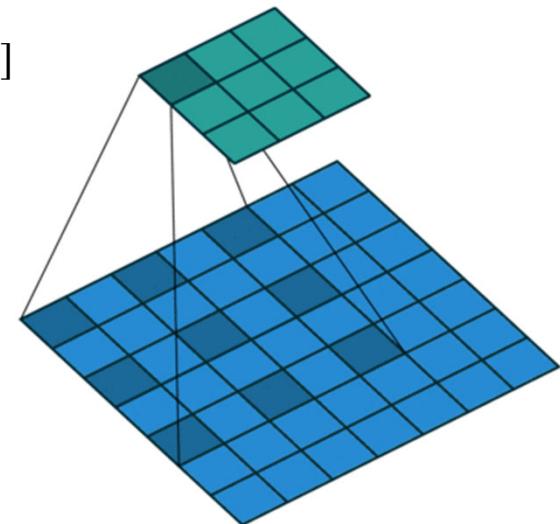
$$[k_x + d(k_x - 1)] \times [k_y + d(k_y - 1)]$$

0	0	0	0	0	0	0	0	0
0	1	1	1	1	1	1	0	0
0	1	2	2	2	2	1	0	0
0	1	2	3	3	2	1	0	0
0	1	2	3	3	2	1	0	0
0	1	2	2	2	2	1	0	0
0	1	1	1	1	1	1	0	0
0	0	0	0	0	0	0	0	0

$x, 9 \times 1$   $w^T x$

Then convolve with dilated kernel, but only count the  $k_x \times k_y$  locations

$$\begin{bmatrix} 1 & 1 & 1 & 1 & 3 & 2 & 1 & 2 & 2 \end{bmatrix} \times \begin{bmatrix} 1 & 0 & -1 & 2 & 0 & -2 & 1 & 0 & -1 \end{bmatrix}^T = -3$$



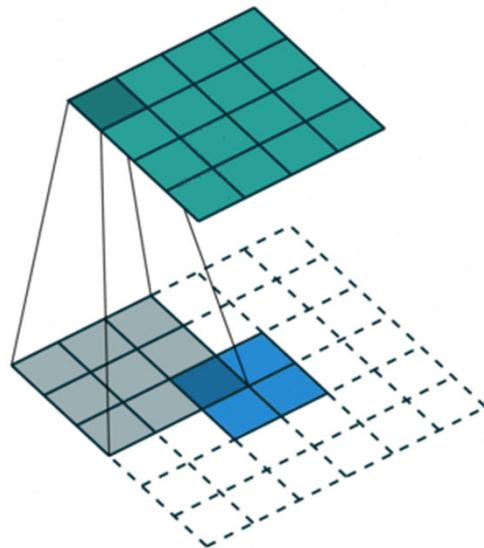
Input: 7x7, padding 0, stride 1  
dilation rate 1, kernel 3x3

$$\text{Output: } 3 \times 3, \left\lfloor \frac{7+2 \times 0 - [3+1 \times (3-1)]}{1} \right\rfloor + 1 = 3$$

Animation is from [https://github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)

# Transpose convolution

- Increase image size after convolution – often used to upsample feature maps



Input: 2x2, padding 0, stride 1  
Kernel: 3x3

Required Output: 4x4

What is the padding size p?

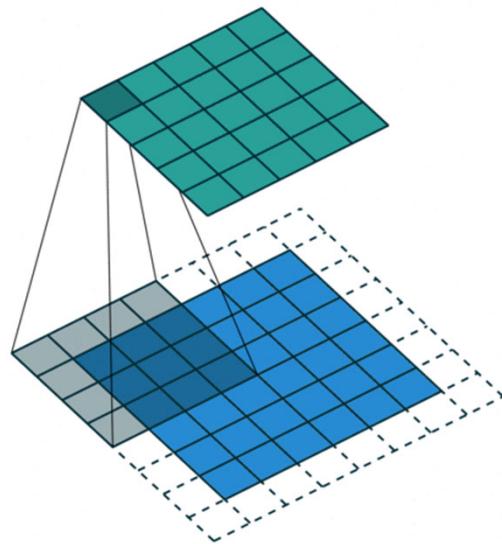
$$\left\lceil \frac{2 + 2 \times p - 3}{1} \right\rceil + 1 = 4$$

$$p = 2$$

- Determine padding size
- Sliding over the **output**
- For every location in output, take the corresponding pixels in input maps and apply the kernel

# Transpose convolution

- Increase image size after convolution – often used to upsample feature maps



Input: 6x6, padding 0, stride 1  
Kernel: 4x4

Required Output: 5x5

What is the padding size p?

$$\left\lceil \frac{6 + 2 \times p - 4}{1} \right\rceil + 1 = 5$$

$$p = 1$$

- Determine padding size
- Sliding over the **output**
- For every location in output, take the corresponding pixels in input maps and apply the kernel

# Convolution is shift-invariant, not for rotation and scale

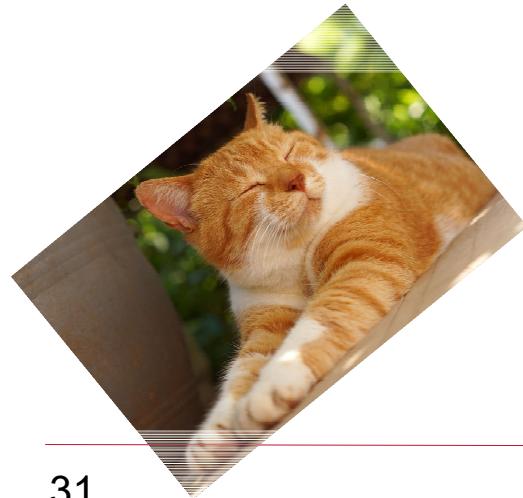
- CONV cannot deal with arbitrary spatial transformation



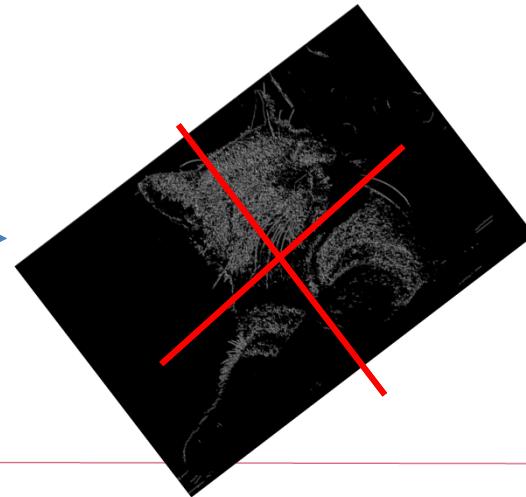
conv  
→



More general, the training and test sets can have different distribution



conv  
→



Training set, cats are roughly the same size

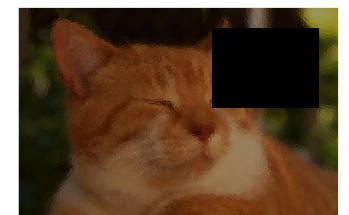
A test sample, small cat

# Data augmentation

- Introduce variation in the training data, so model can learn invariance



Randomly apply  
different  
transformations



And many more ...

Check <http://ai.stanford.edu/blog/data-augmentation/>

# Data augmentation

- **Do not** generate these random samples and save them
- Instead, implement them as a part of data loading



More examples can be found at [https://pytorch.org/vision/stable/auto\\_examples/plot\\_transforms.html#sphx-glr-auto-examples-plot-transforms-py](https://pytorch.org/vision/stable/auto_examples/plot_transforms.html#sphx-glr-auto-examples-plot-transforms-py)

# Pooling : another way to reduce computation

- Often used to reduce the image size
- Parameter free

0	0	0	0	0	0	0	0	0
0	1	1	1	1	1	1	1	0
0	1	2	2	2	2	1	0	
0	1	2	3	3	2	1	0	
0	1	2	3	3	2	1	0	
0	1	2	2	2	2	1	0	
0	1	1	1	1	1	1	1	0
0	0	0	0	0	0	0	0	0

Max pool  
→

1			

Input: 8x8  
Pool size: 2x2

Take every 2x2 neighborhood, compute its max (max pool) or average (averaging pool)

Typically, no padding, stride equals to the pool\_size

Output size

$$\left\lceil \frac{H - \text{pool\_size}}{\text{pool\_size}} \right\rceil + 1$$

# Pooling : another way to reduce computation

- Often used to reduce the image size
- Parameter free

0	0	0	0	0	0	0	0	0
0	1	1	1	1	1	1	1	0
0	1	2	2	2	2	2	1	0
0	1	2	3	3	2	1	0	
0	1	2	3	3	2	1	0	
0	1	2	2	2	2	2	1	0
0	1	1	1	1	1	1	1	0
0	0	0	0	0	0	0	0	0

Max pool  
→

1	1	1	1
1	3		

Input: 8x8  
Pool size: 2x2

Take every 2x2 neighborhood, compute its max (max pool) or average (averaging pool)

Typically, no padding, stride equals to the pool\_size

Output size

$$\left\lceil \frac{H - \text{pool\_size}}{\text{pool\_size}} \right\rceil + 1$$

# Pooling : another way to reduce computation

- Often used to reduce the image size
- Parameter free

0	0	0	0	0	0	0	0	0
0	1	1	1	1	1	1	1	0
0	1	2	2	2	2	1	0	
0	1	2	3	3	2	1	0	
0	1	2	3	3	2	1	0	
0	1	2	2	2	2	1	0	
0	1	1	1	1	1	1	0	
0	0	0	0	0	0	0	0	0

Max pool  
→

1	1	1	1
1	3	3	1
1	3	3	1
1	1	1	1

2x2 is the most typical choice

Input: 8x8  
Pool size: 2x2

Take every 2x2 neighborhood, compute its max (max pool) or average (averaging pool)

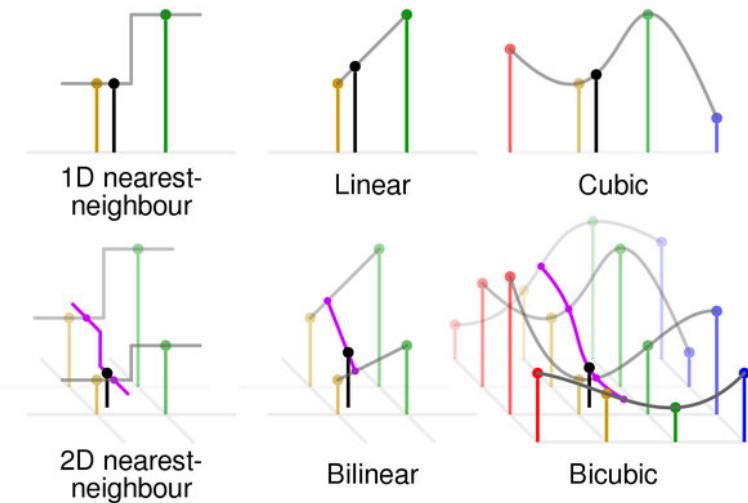
Typically, no padding, stride equals to the pool\_size

Output size

$$\left\lceil \frac{H - \text{pool\_size}}{\text{pool\_size}} \right\rceil + 1$$

# Use interpolation to resize images

- Another “parameter free” way to resize images



- Interpolation: use neighboring pixels to determine value of unknown location (weighted sum)
- Weights depend on pixel location; no learnable parameters

## TORCH.NN.FUNCTIONAL.INTERPOLATE

```
torch.nn.functional.interpolate(input, size=None, scale_factor=None, mode='nearest',  
align_corners=None, recompute_scale_factor=None)
```

[SOURCE]

Down/up samples the input to either the given `size` or the given `scale_factor`.

The algorithm used for interpolation is determined by `mode`.

Currently temporal, spatial and volumetric sampling are supported, i.e. expected inputs are 3-D, 4-D or 5-D in shape.

The input dimensions are interpreted in the form: `mini-batch x channels x [optional depth] x [optional height] x width`.

The modes available for resizing are: `nearest`, `linear` (3D-only), `bilinear`, `bicubic` (4D-only), `trilinear` (5D-only), `area`

- It is a layer in NN
- Differentiable by remembering the interpolation weights

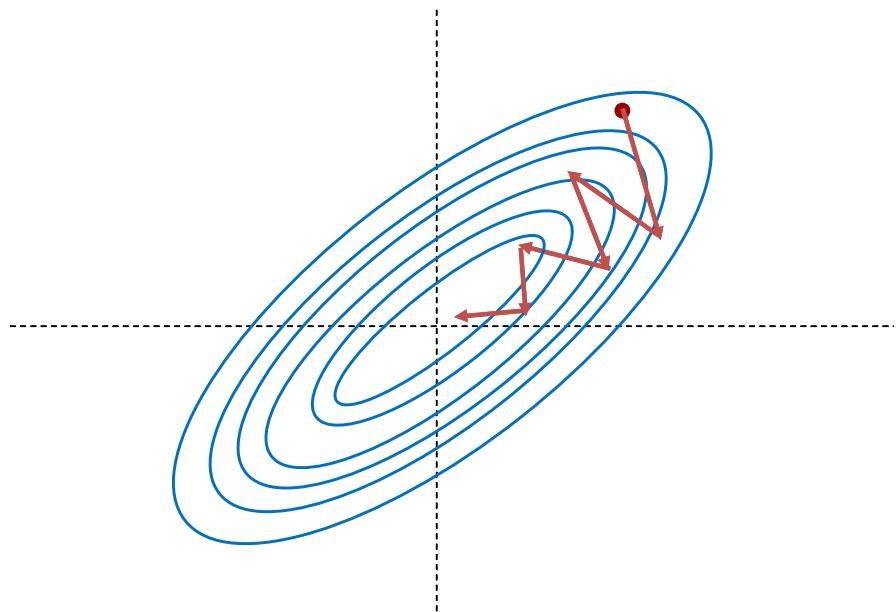
<https://www.quora.com/What-is-bicubic-interpolation-in-image-processing>

# Layers to adjust image size

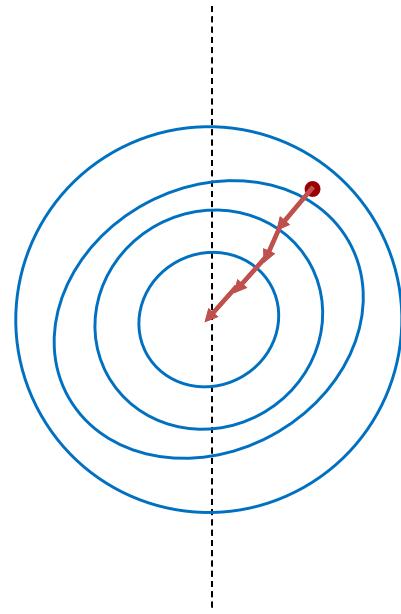
Layers	Downsample	Upsample	Pro/Cons
Stride CONV	Yes	No	Same as normal conv, fast to compute
Pooling	Yes	No	Fast to compute, reduce the model power, less popular nowadays
Interpolation	Yes	Yes	Parameter free, fast to compute
Transpose CONV	Yes, but often not used for this purpose	Yes	Has learnable parameters for upsampling

- Pooling still is used in classification tasks
- Full CONV architecture are more popular for other tasks, e.g. segmentation, super-resolution, denoising ...

# Batch Normalization



Features with varying magnitude



Features with normalized magnitude

$$y = w_1x_1 + w_2x_2 + w_3x_3$$

$$\frac{d\ell}{dw_1} = \frac{d\ell}{dy}x_1$$

$$\frac{d\ell}{dw_2} = \frac{d\ell}{dy}x_2$$

$$\frac{d\ell}{dw_3} = \frac{d\ell}{dy}x_3$$

- Features with unnormalized magnitude leads to harder optimization

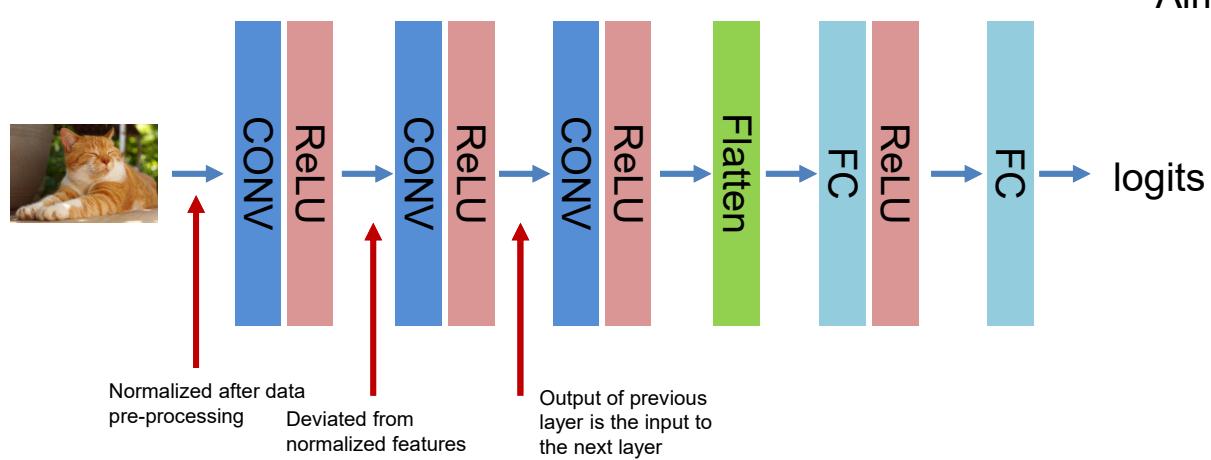
# Batch Normalization

Techniques to battle unnormalized features

- Data pre-processing, whitening ...
- Model parameter initialization (e.g. He and Xavier)

Model will be well-behaved at the beginning of training

- No guarantee for later epochs
- No guarantee for deeper layers
- Almost guaranteed to deviate



# Batch Normalization: if we need normalized features, add a layer to enforce it

Add a specific layer to normalize the output of previous layer

Given a mini-batch,  $X = \{x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad \dots \quad \dots \quad x^{(B)}\}$

Every sample  $x^{(k)}$  is a Nx1 vector or CxHxW image

- $\hat{x}^{(k)}$  is the Nx1 vector, same as  $x^{(k)}$
- $\mu$  and  $\delta$  are per-dimension mean and variance

- Compute mean of this batch
- Compute variance of this batch

$$\mu = \frac{1}{B} \sum_{k=1}^B x^{(k)}$$

$$\delta^2 = \frac{1}{B} \sum_{k=1}^B [x^{(k)} - \mu]^2$$

- Normalize the batch

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mu}{\sqrt{\delta^2 + \epsilon}}$$

← Element-wise operation

# Batch Normalization: if we need normalized features, add a layer to enforce it

Give a mini-batch,  $X = \{x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad \dots \quad \dots \quad x^{(B)}\}$

$x^{(k)}$  is a Nx1 vector or CxHxW image

$$\mu = \frac{1}{B} \sum_{k=1}^B x^{(k)} \quad \delta^2 = \frac{1}{B} \sum_{k=1}^B [x^{(k)} - \mu]^2$$

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mu}{\sqrt{\delta^2 + \epsilon}}$$

$\gamma$  and  $\beta$  are the same dimension as  $x^{(k)}$

$$y^{(k)} = \gamma \circ \hat{x}^{(k)} + \beta \quad \hat{x}^{(k)} = x^{(k)} \text{ if } \gamma = \delta \text{ and } \beta = \mu$$

- Normalization reduces the representation power of model
- To restore the representation power, introduce two learnable parameters

At training time,

Given a mini-batch, compute mean and variance

Compute  $\hat{x}$  and  $y$ , using  $\gamma$  and  $\beta$  and  $\mu$  and  $\delta$

Therefore, if batch size is too small (e.g. <16), BN may not work well

# Batch Normalization: if we need normalized features, add a layer to enforce it

At training time,

Given a mini-batch, compute  $\mu$  and  $\delta$  from this batch

Compute  $\hat{x}$  and  $y$ , using  $\gamma$  and  $\beta$  and  $\mu$  and  $\delta$

$$\mu = \frac{1}{B} \sum_{k=1}^B x^{(k)} \quad \delta^2 = \frac{1}{B} \sum_{k=1}^B [x^{(k)} - \mu]^2$$

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mu}{\sqrt{\delta^2 + \epsilon}}$$

$$y^{(k)} = \gamma \circ \hat{x}^{(k)} + \beta$$

At test time,

We don't have a mini-batch, often  $B=1$

- The assumption behind BN is the statistics from a batch can approximate the entire dataset
- Test and training datasets have the same distribution

$$\bar{\mu} = (1 - \alpha)\bar{\mu} + \alpha\mu$$

$$\bar{\delta} = (1 - \alpha)\bar{\delta} + \alpha\delta$$

$$\alpha = 0.1$$

or,  $\bar{\mu}$  and  $\bar{\delta}$  can be computed as simple average for all trained batches

# Spatial Batch Normalization for CONV

Batch Normalization

$B \times N$

↓  
Normalize along B

$B \times N$

$\gamma, \beta: [N \times 1]$  learnable parameters

Spatial Batch Normalization

$B \times C \times H \times W$

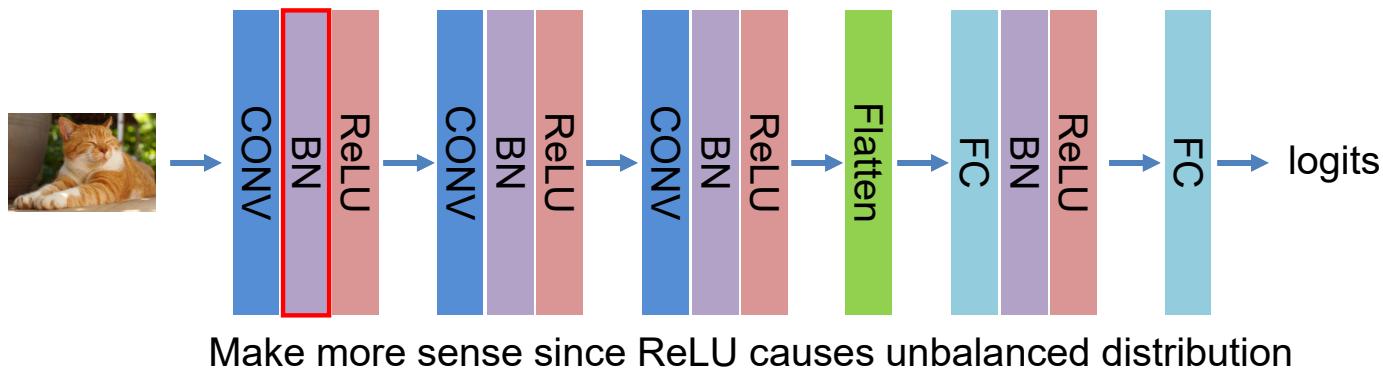
↓      ↓      ↓  
Normalize along  
 $B \times H \times W$

$B \times C \times H \times W$

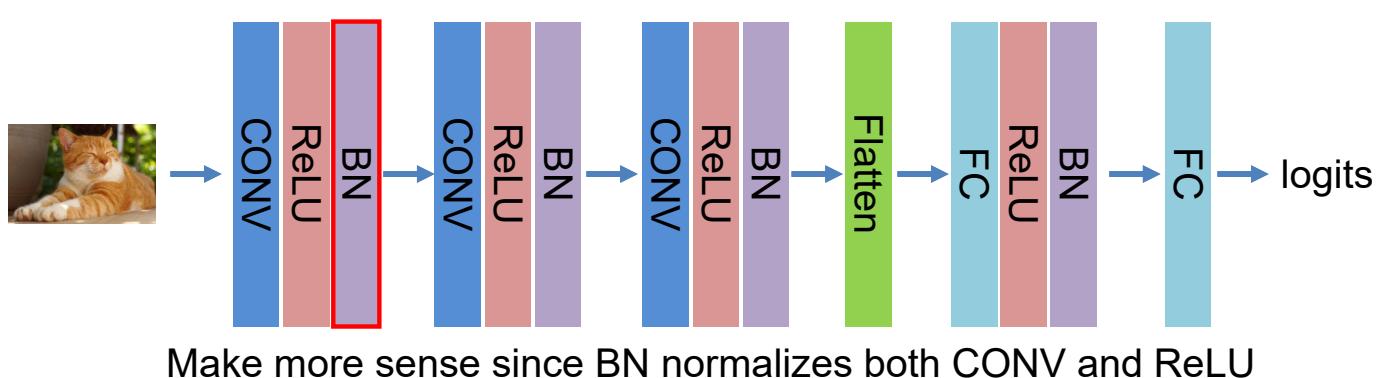
$\gamma, \beta: [C \times 1]$  learnable parameters

e.g. `torch.nn.BatchNorm2d`

# Model with BN layer



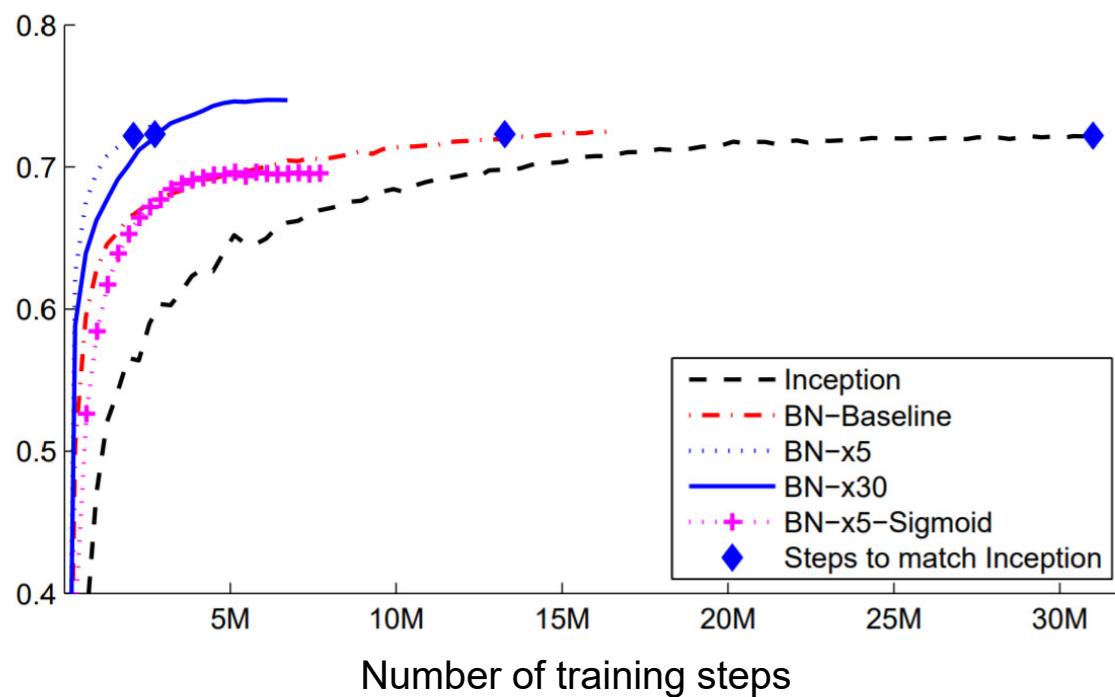
- Often lead to faster and more stable convergence
- Less sensitive to parameter initialization
- Less sensitive to hyper-parameter setting
- Allow flexible selection of nonlinearity functions



- Has regularization effect → introduce randomness in training
- Do not use with drop-out
- Can hide bugs in the code – debug by turning BN on/off

# From BN paper

Accuracy



- Test on ImageNet 2012 dataset with the inception network
- BN-x5, x30 is to increase learning rate by 5 times or 30 times → BN allows to use larger learning rate → faster convergence
- BN-x5-Sigmoid shows BN allows to use sigmoid activation function. Without BN, sigmoid activation does not lead to convergence
- BN-x30 got higher accuracy than BN-x5 → exploration

# Model evaluation and training

```
# evaluate model:  
model.eval()  
  
with torch.no_grad():  
    ...  
    out_data = model(data)  
    ...
```

```
# training step  
...  
model.train()  
...
```

- Do not forget to set model be in evaluation mode, so the moving averaged mean and standard deviation will not be updated
- When doing training, set model to be in training mode, so the moving averaged mean and standard deviation will be updated
- A common source of error

# More feature normalization

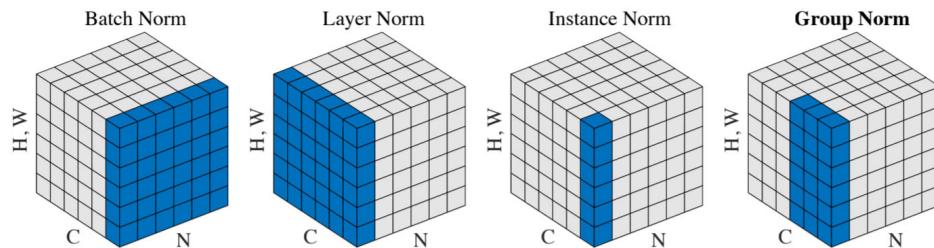


Figure 2. **Normalization methods.** Each subplot shows a feature map tensor, with  $N$  as the batch axis,  $C$  as the channel axis, and  $(H, W)$  as the spatial axes. The pixels in blue are normalized by the same mean and variance, computed by aggregating the values of these pixels.

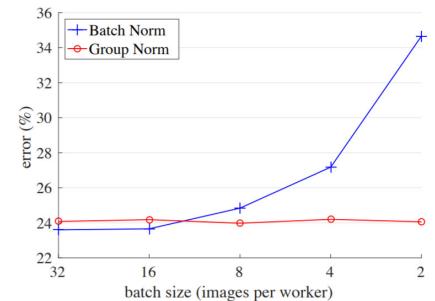
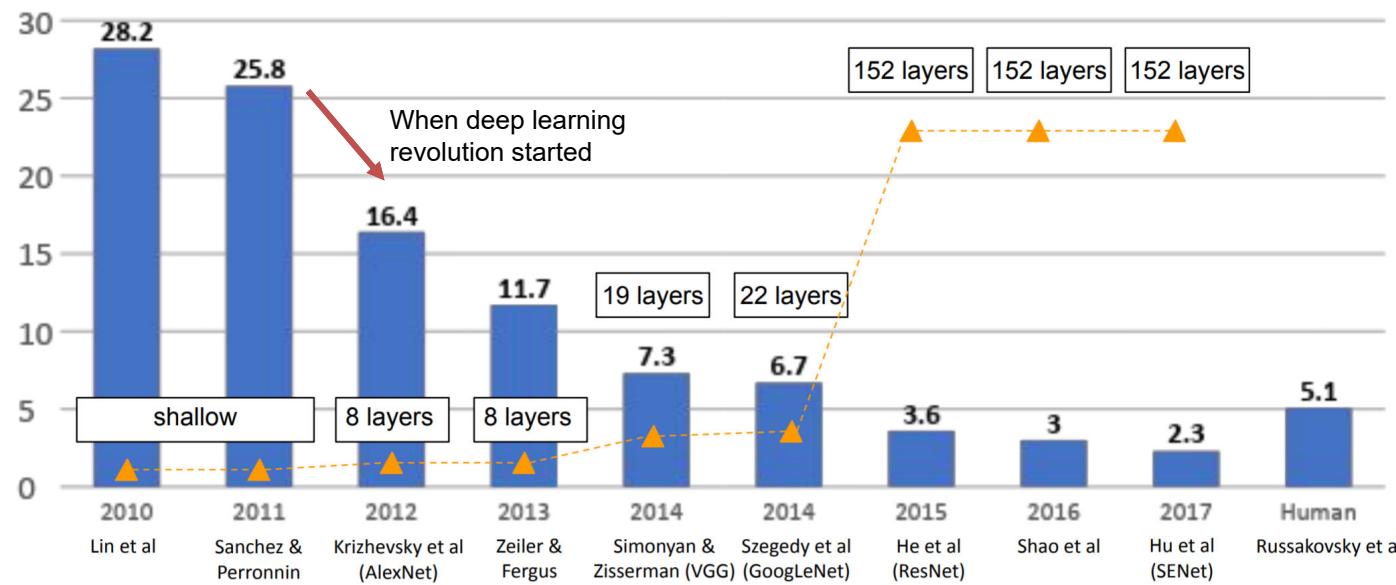


Figure 1. **ImageNet classification error vs. batch sizes.** This is a ResNet-50 model trained in the ImageNet training set using 8 workers (GPUs), evaluated in the validation set.

# CNN leads to the revolution of deep learning



- Deep learning revolution started in 2012 with AlexNet
  - Its origin were long back in time
  - LeNet 1989 or earlier
- CNN architecture development was well tracked by ImageNet winners
- What was learned was adapted to other applications

Figure is from [http://cs231n.stanford.edu/slides/2021/lecture\\_9.pdf](http://cs231n.stanford.edu/slides/2021/lecture_9.pdf)

# LeNet and MNIST datasets



```
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2  
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3  
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4  
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5  
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6  
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7  
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8  
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9
```

Modified National Institute of Standards and Technology database

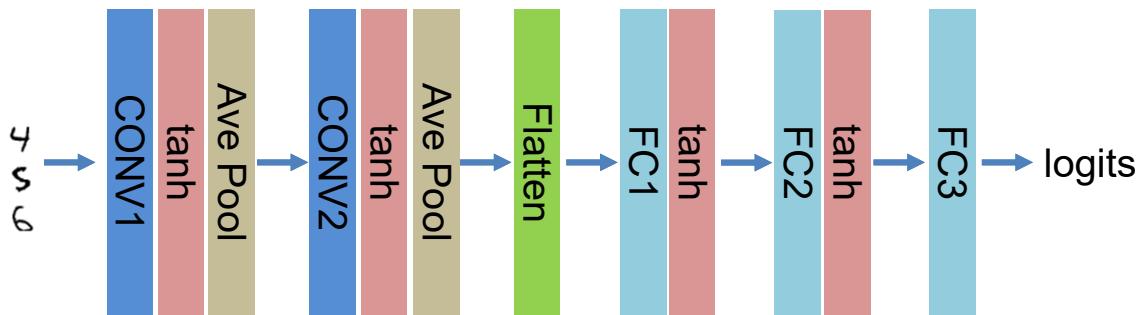
The MNIST database of handwritten digits, available from this page, has a training set of 60,000 examples, and a test set of 10,000 examples. It is a subset of a larger set available from NIST. The digits have been size-normalized (28x28) and centered in a fixed-size image.

- <http://yann.lecun.com/exdb/mnist/>

EMNIST has been published in 2017, which contains 240,000 training images, and 40,000 testing images of handwritten digits and characters.

Y. LeCun, L. Bottou, Y. Bengio and P. Haffner: Gradient-Based Learning Applied to Document Recognition, *Proceedings of the IEEE*, 86(11):2278-2324, November 1998

# LeNet-5



While the architecture of the best performing neural networks today are not the same as that of LeNet, **the network was the starting point for a large number of neural network architectures, and also brought inspiration to the field.**

Input : Bx1x28x28	Output
CONV1: 5x5, padding 2, stride 1, K = 6	Bx6x28x28
Ave pool: 2x2	Bx6x14x14
CONV2: 5x5, padding 0, stride 1, K = 16	Bx16x10x10
Ave pool: 2x2	Bx16x5x5
FC1: D=120	Bx120
FC2: D=84	Bx84
FC3: D=10	Bx10

~60K parameters

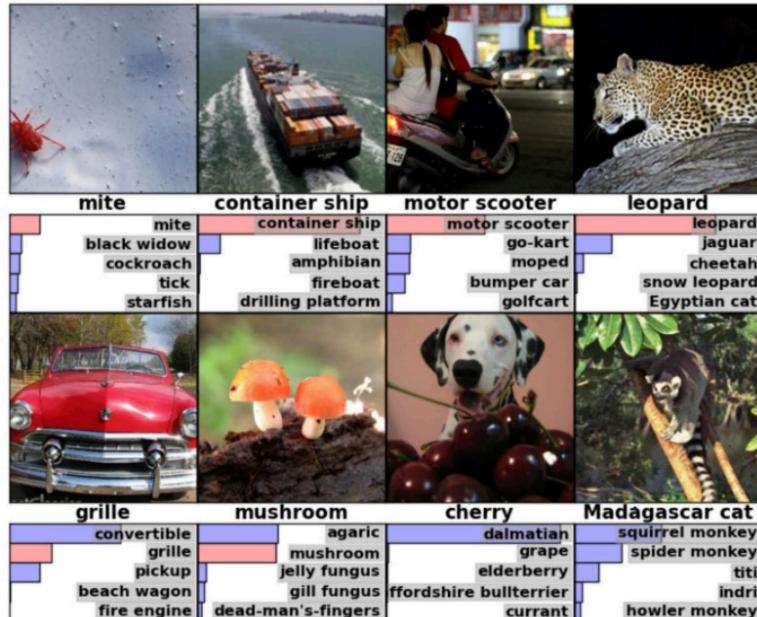
- <https://en.wikipedia.org/wiki/LeNet>

Used by the USPS, maybe the first deployed cnn?

- Backprop to train CONV
- Use ave pooling
- Use tanh
- Use flatten and FC
- Use softmax and CE loss
- Increase number of filters with reduced spatial resolution

LeCun, Y.; Boser, B.; Denker, J. S.; Henderson, D.; Howard, R. E.; Hubbard, W. & Jackel, L. D. (1989). Backpropagation applied to handwritten zip code recognition. Neural Computation, 1(4):541-551.

# ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



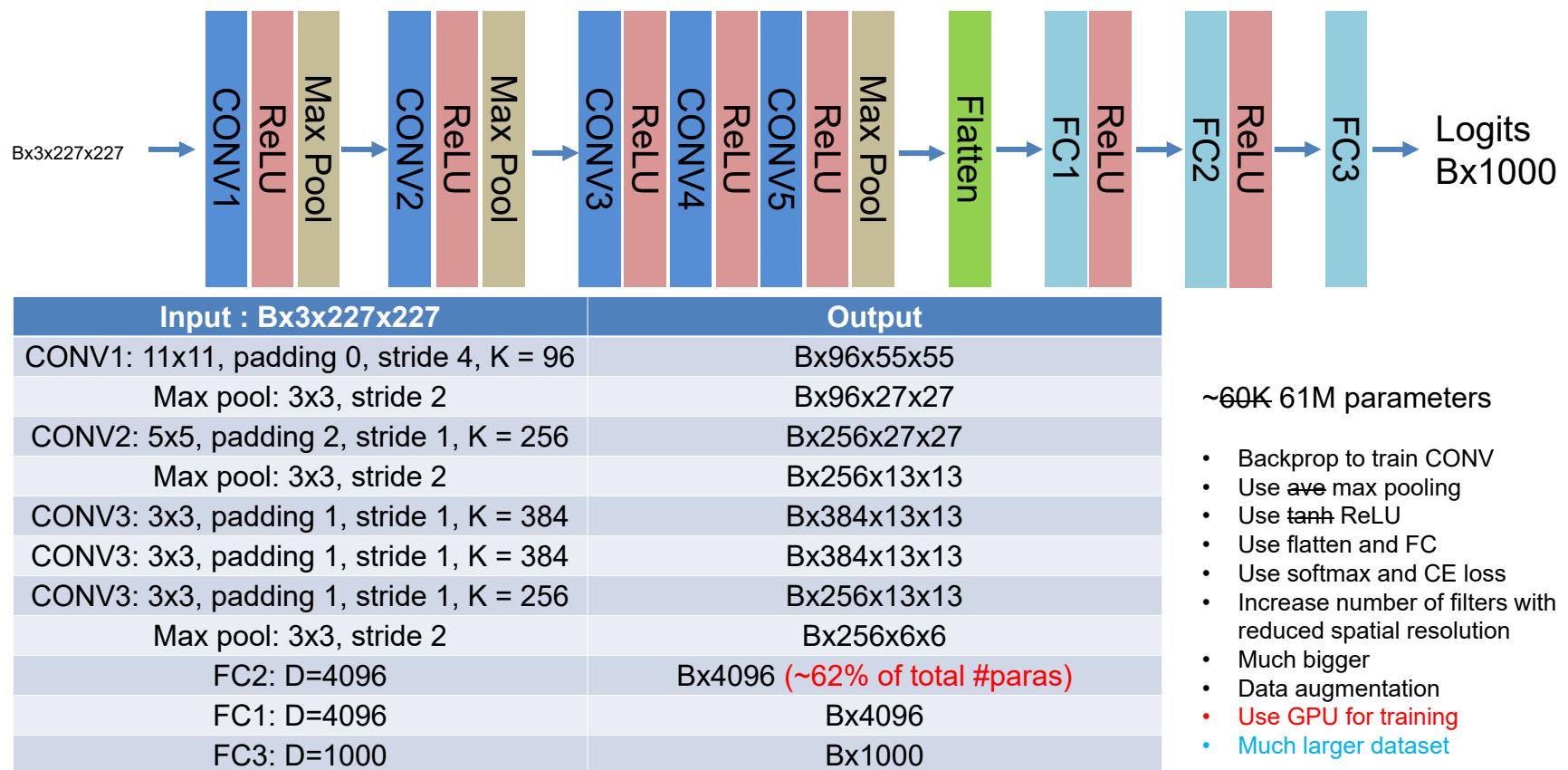
I. [Object localization](#) for 1000 categories.

II. [Object detection](#) for 200 fully labeled categories.

III. [Object detection from video](#) for 30 fully labeled categories.

- Running since 2010, currently still run in Kaggle  
<https://image-net.org/challenges/LSVRC/index.php>
- <https://www.kaggle.com/c/imagenet-object-localization-challenge/overview/description>
- <https://www.paperswithcode.com/sota/image-classification-on-imagenet?metric=Top%205%20Accuracy>
- ILSVRC used 1.2 million training images for 1000 classes, 50,000 validation images, and 150,000 testing images
- Include object classification, detection and video detection tasks
- ImageNet is a dataset of over 15 million labeled high-resolution images belonging to roughly 22,000 categories

# AlexNet



# To be continued ...

