

# Machine Learning (CE 40717)

Fall 2024

Ali Sharifi-Zarchi

CE Department  
Sharif University of Technology

October 25, 2024



## 1 Motivation

## 2 CNN vs. FCN

## 3 Convolution

## 4 Pooling

## 5 Receptive field

## 6 Inductive bias

## 7 Backpropagation

## 8 References

## 1 Motivation

## 2 CNN vs. FCN

## 3 Convolution

## 4 Pooling

## 5 Receptive field

## 6 Inductive bias

## 7 Backpropagation

## 8 References

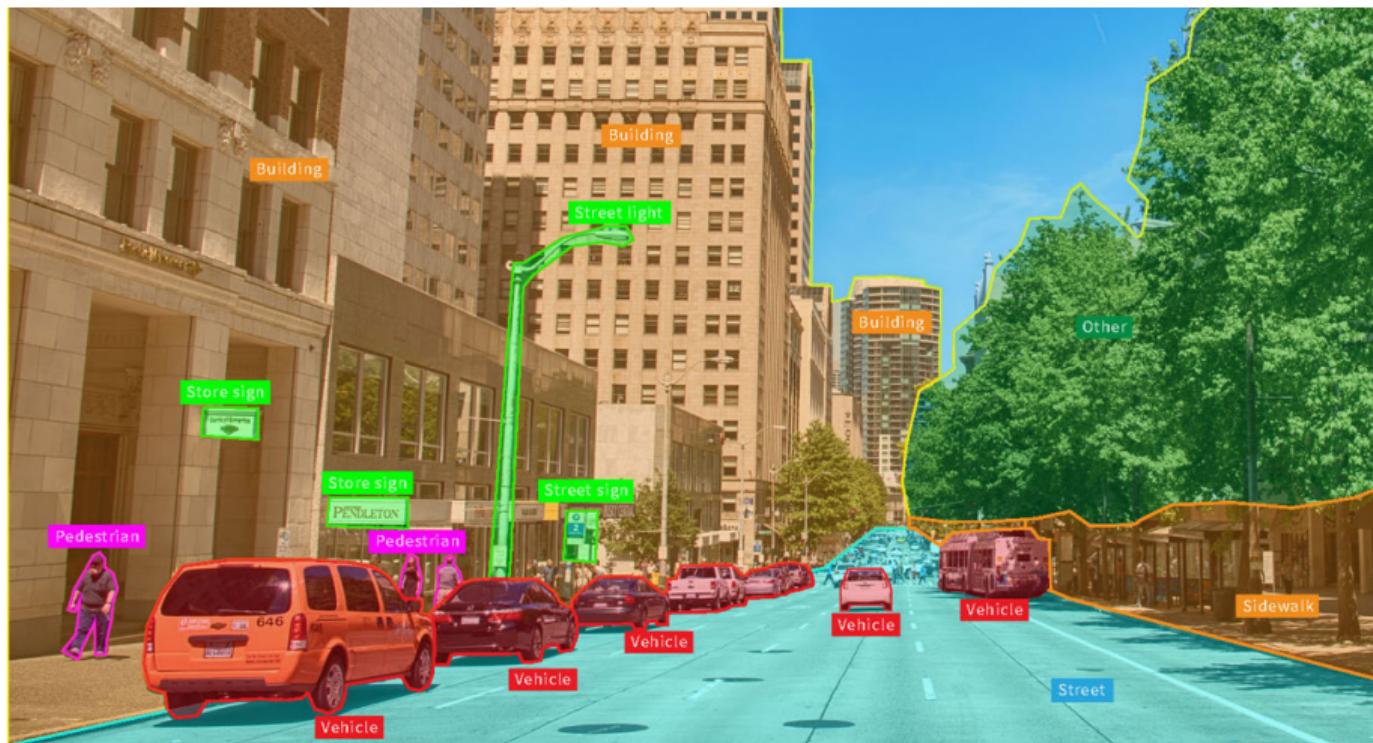
# How do you see?



## How can we help computers see?



# What is Computer Vision? (1)



## What is Computer Vision? (2)



Noisy Image

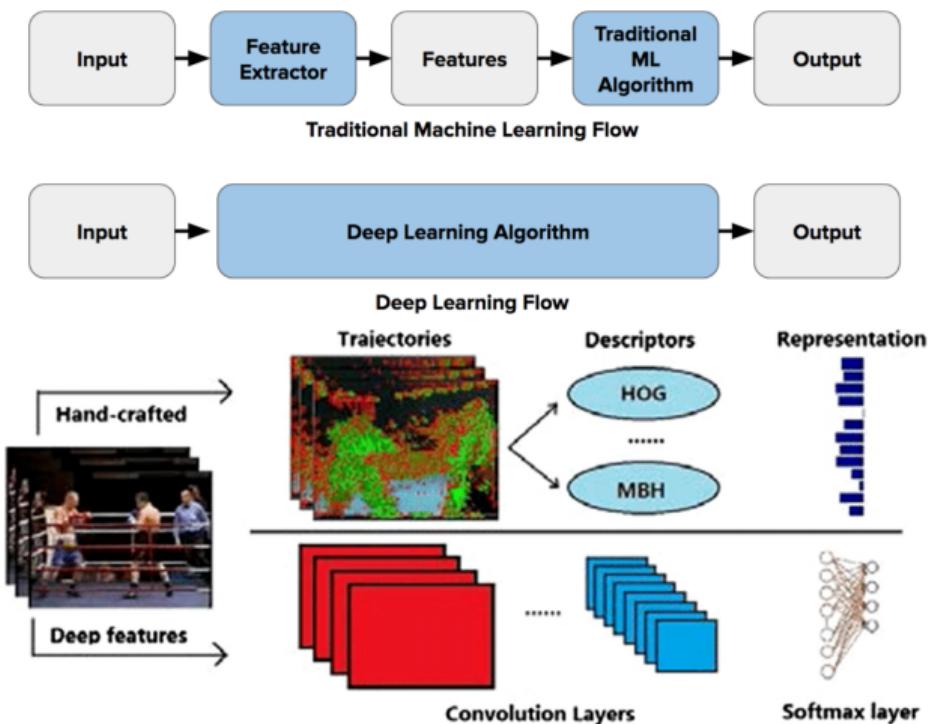


Denoised Image

# What is Computer Vision? (3)

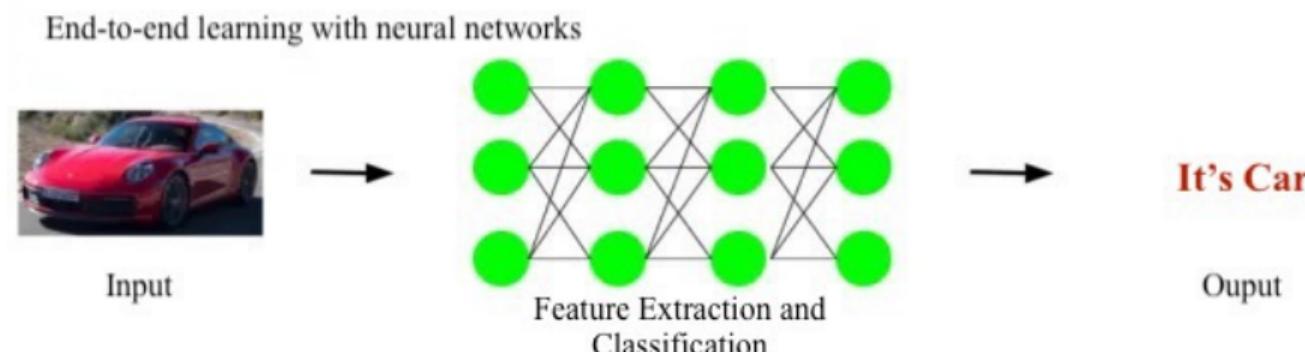
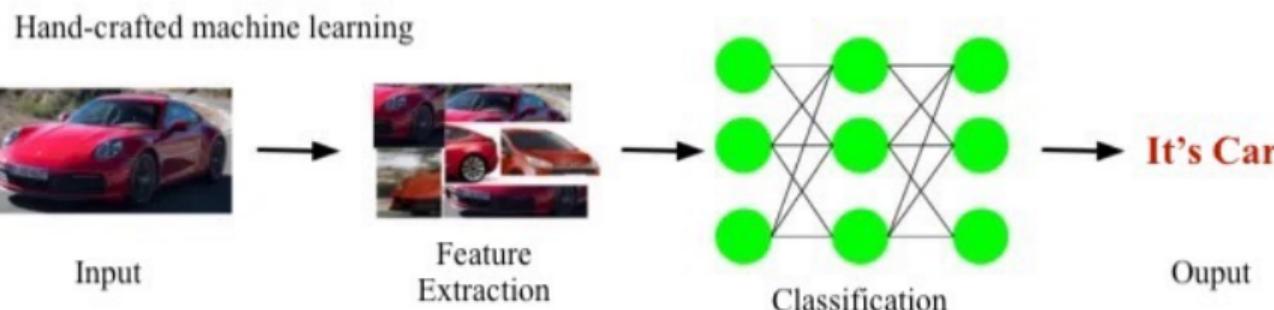


# Hand-Crafted Features: before deep learning features (1)

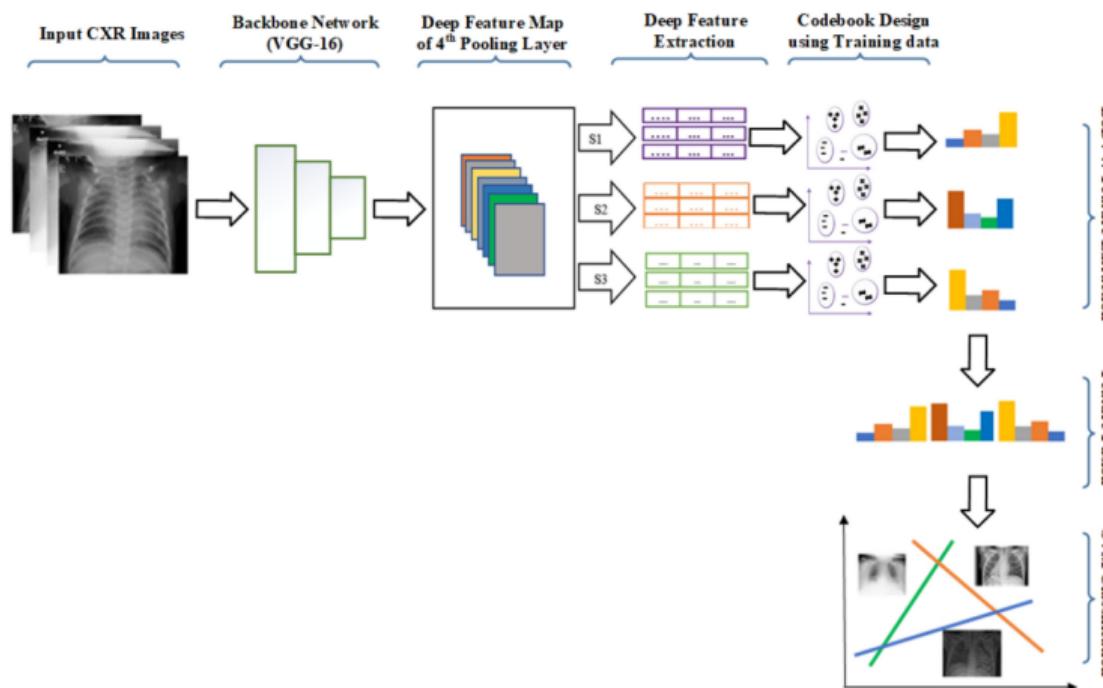


Two types of features used for action recognition, i.e., hand-crafted features and deep learning features

## Hand-Crafted Features: before deep learning features (2)

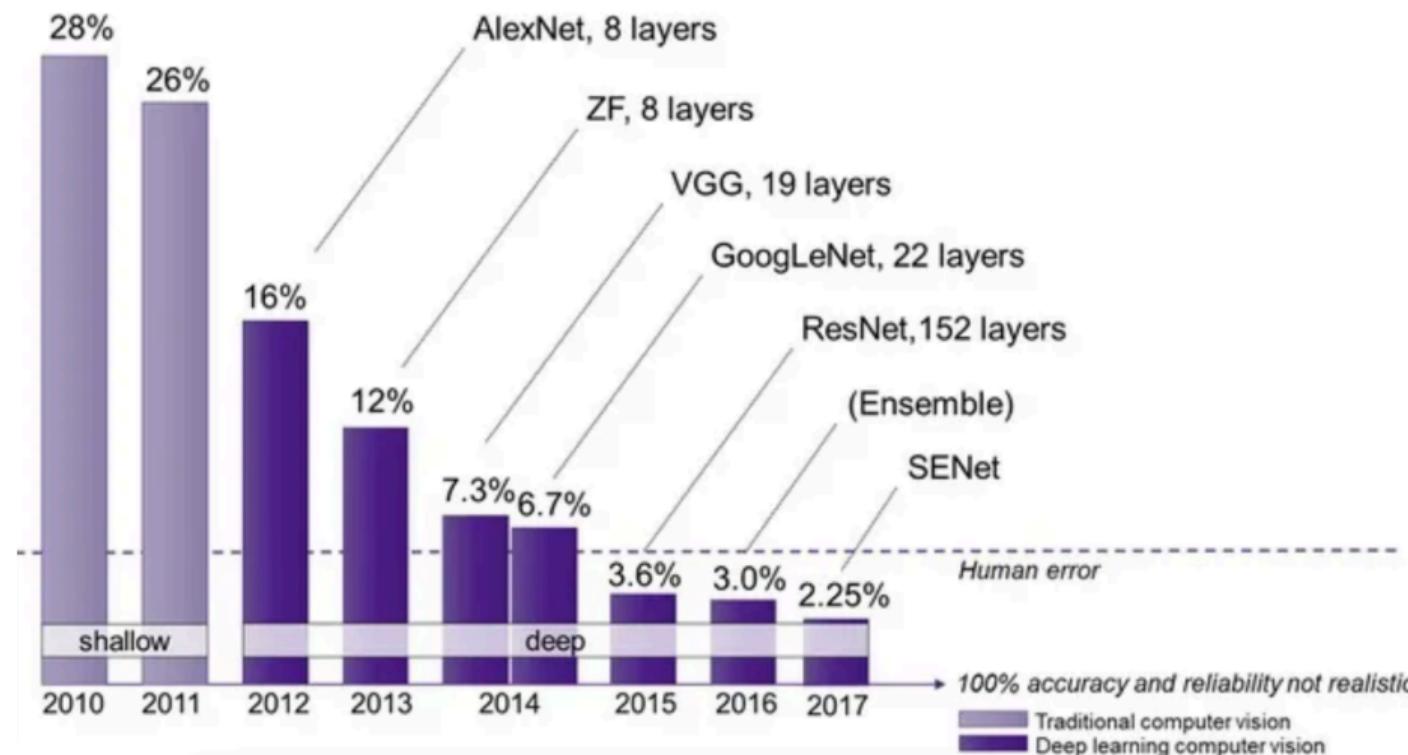


# Image Features



Fusion of multi-scale bag of deep visual words features of chest X-ray images to detect COVID-19 infection

# Effect of DeepLearning



Comparison between Deep Learning and Traditional Models

## 1 Motivation

## 2 CNN vs. FCN

## 3 Convolution

## 4 Pooling

## 5 Receptive field

## 6 Inductive bias

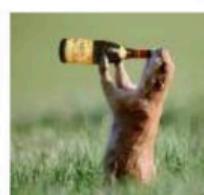
## 7 Backpropagation

## 8 References

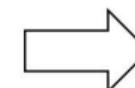
# FCLs on images

## Dense Vector Multiplication

32×32×3 image



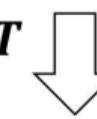
Reshape it into  
a vector

 $x$ 

100×3072

 $W$ 

$$Wx^T$$



100

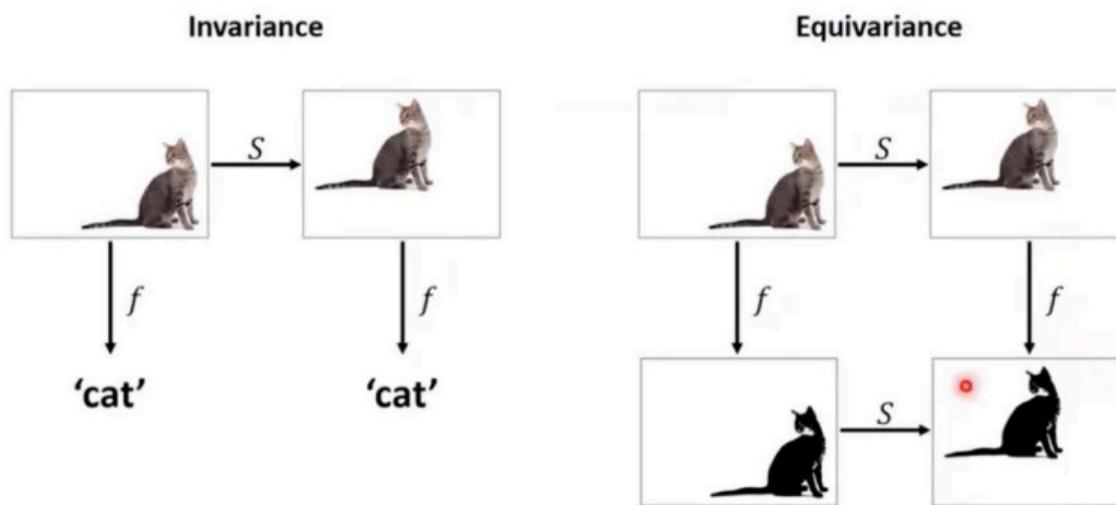
Each element contains the  
activation of 1 neuron



# Fully Connected Layers

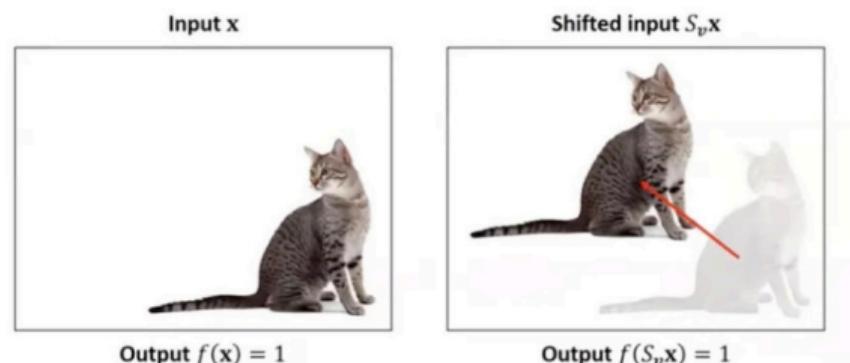
- Neurons in a single layer function completely independently and do not share any connections even for inputs.
- To be shift-invariant, many samples (in various time or locations) must be shown to them.
- **Regular Neural Nets dont scale well to full images**
  - Parameters would add up quickly!
  - Full connectivity is wasteful and the huge number of parameters would quickly lead to overfitting.

# Invariance vs Equivariance



- **Invariance:** The output remains the same when the input undergoes a transformation.
- **Equivariance:** The output changes in a predictable way when the input undergoes a transformation.

# Shift-Invariance



- *Cat detector*  $f: \mathbb{R}^d \rightarrow \mathbb{R}$
- *Shift operator*  $S_v: \mathbb{R}^d \rightarrow \mathbb{R}^d$  shifting the image by vector  $v$

**Shift invariance:**  $f(\mathbf{x}) = f(S_v\mathbf{x})$

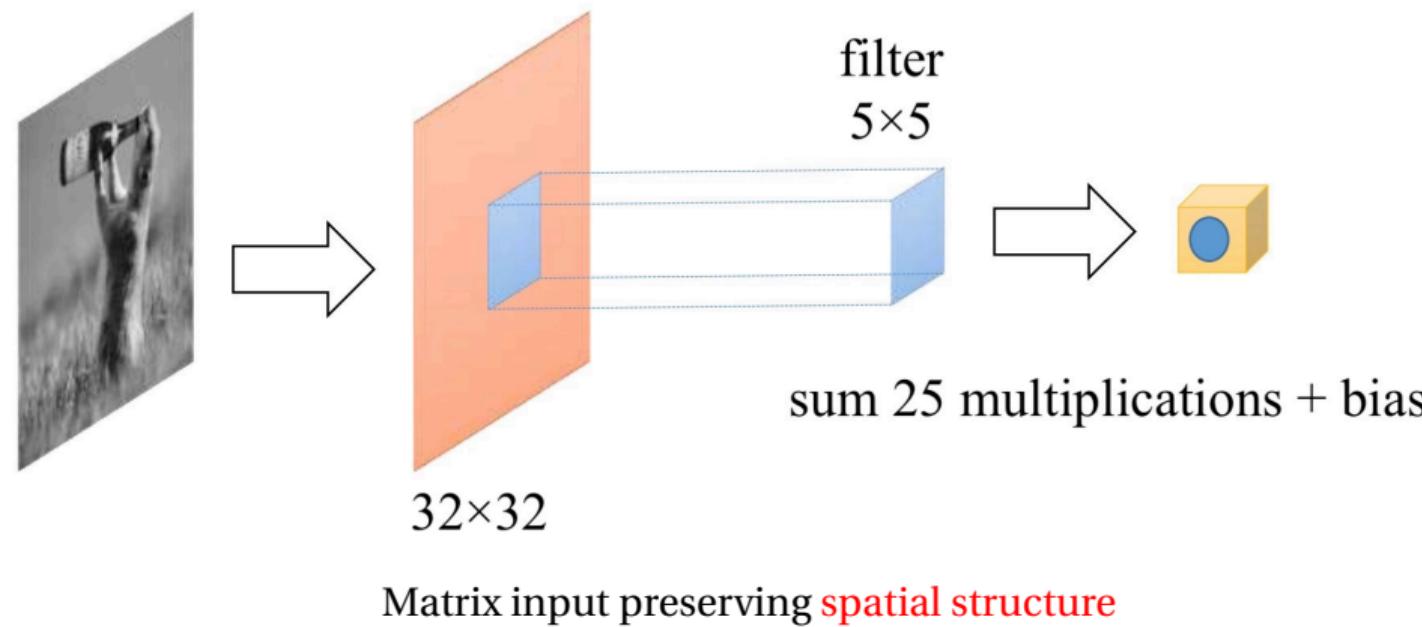
Will an MLP that recognizes the left image as a cat also recognize the shifted image on the right as a cat?

# A Problem

- In many problems the *location* of a pattern is not important
  - Only the presence of the pattern
- Conventional MLPs are sensitive to the location of the pattern
  - Moving it by one component results in an entirely different input that the MLP wont recognize
- **Requirement:** Network must be *shift invariant*

# Convolution on images

## Convolution (Refresher)



# Fully Connected or Convolution

**Question:** Can I just do classification on an flatten image (e.g., a 1080x1080x3 image) with fully connected layers?

**Answer:** No, using fully connected layers on such a large image is inefficient due to:

- Loss of spatial structure
- High parameter count

## Why Convolution:

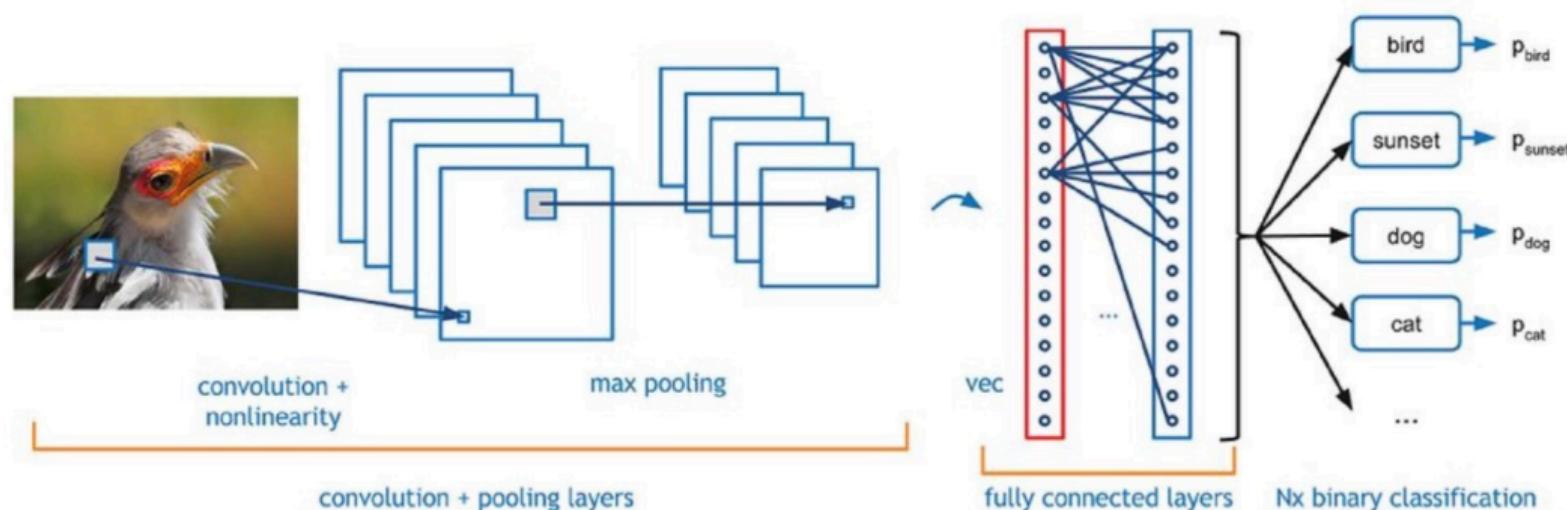
- **Exploit spatial structure:** Convolution sees local patterns by using filters over small patches of the image.
- **Reduce parameters:** By sharing weights across the image, convolution reduce the number of parameters.
- **Build features hierarchically:** Convolution starts with small patterns, then combines them to recognize more complex patters.

# What is CNN?

- It is a class of deep learning.
- Convolutional neural network (ConvNets or CNNs) is one of the main categories to do image recognition, image classifications, object detection, recognition of faces, etc.
- It is similar to the basic neural network. CNN also has learnable parameters like a neural network, i.e., weights, biases, etc.
- CNN is heavily used in computer vision.
- There are 3 basic components to define CNN:
  - The Convolution Layer
  - The Pooling Layer
  - The Output Layer (or) Fully Connected Layer

# Architecture of CNN

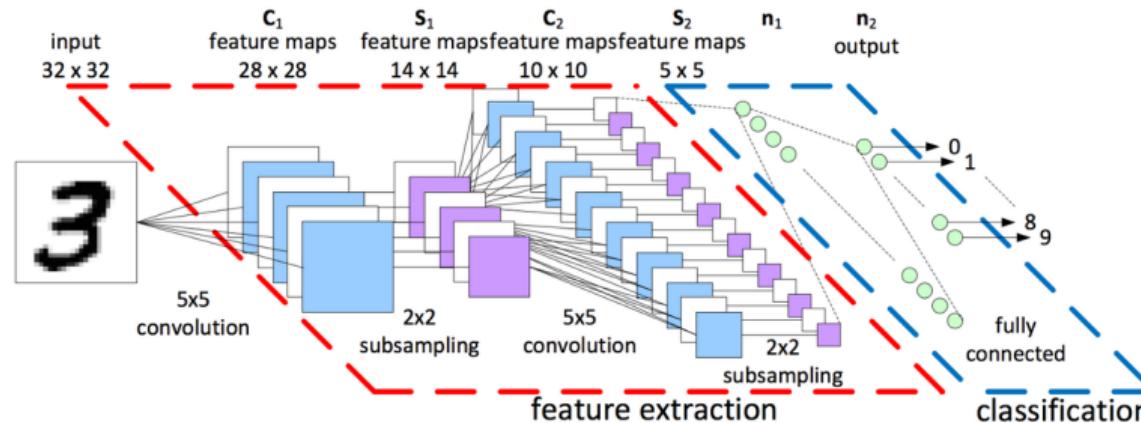
The basic idea of Convolution Neural Networks (CNN) Same idea as Back-propagation-neural networks (BPNN) but different implementation.



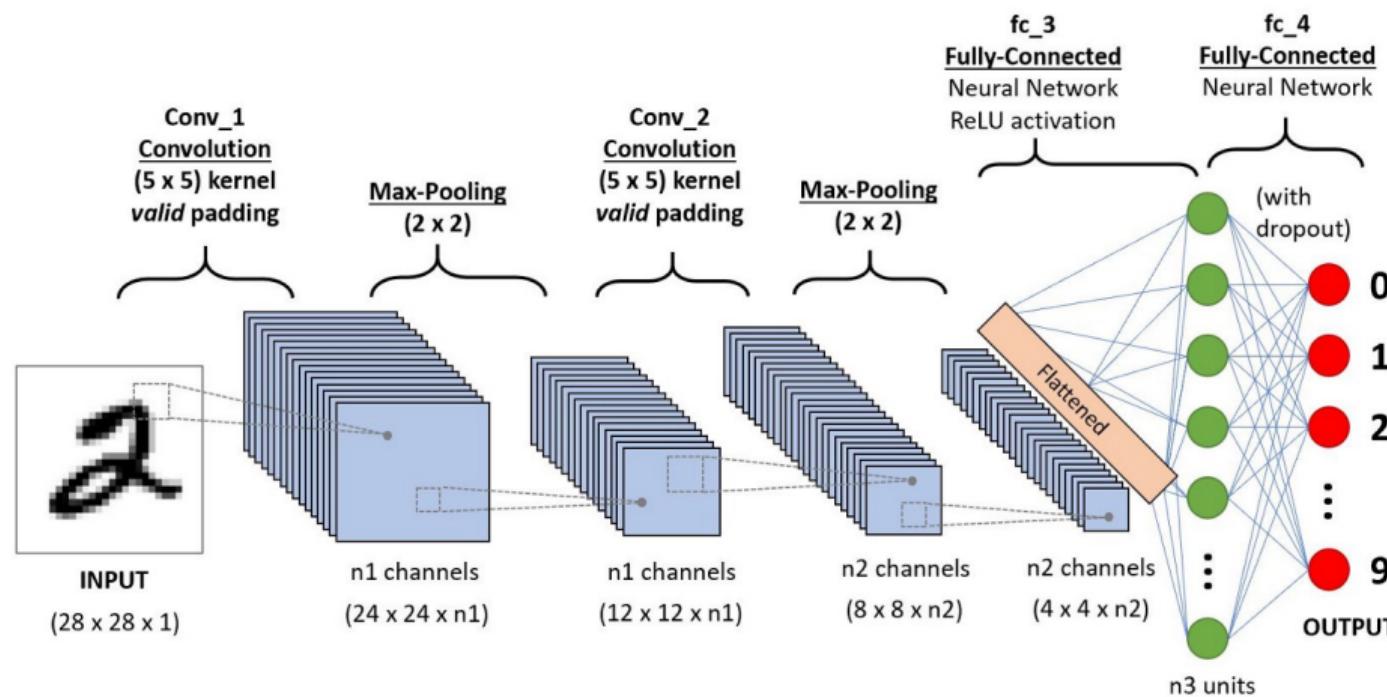
After vectorized (vec), the 2D arranged inputs become ID vectors. Then the network is just like a BPNN (Back propagation neural networks)

# The basic structure (1)

- Alternating Convolution (C) and subsampling layer (S)
- Subsampling allows the features to be flexibly positioned



## The basic structure (2)



# The basic structure (3)

## Three main types of layers

- **Convolutional Layer**

- Output of neurons are connected to local regions in the input.
- Applying the same filter on the whole image.
- CONV layers parameters consist of a set of learnable filters.

- **Pooling Layer**

- Performs a downsampling operation along the spatial dimensions.

- **Fully-Connected Layer**

- Typically used in the final stages of the network to combine high-level features and make predictions.

# CNN vs. FCN

- **Fully Connected Networks (FCN):**

- High number of parameters, leading to overfitting on large inputs (e.g., images).
- Lack of spatial awareness, making it sensitive to the exact positioning of patterns.
- Suitable for structured data but inefficient for image processing.

- **Convolutional Neural Networks (CNN):**

- Uses filters to process only small parts of the input at a time (locality).
- Weight sharing and local connectivity reduce the number of parameters.
- Can recognize patterns regardless of their location in the image (shift invariance).
- Efficient for image, video, and speech data.

## 1 Motivation

## 2 CNN vs. FCN

## 3 Convolution

## 4 Pooling

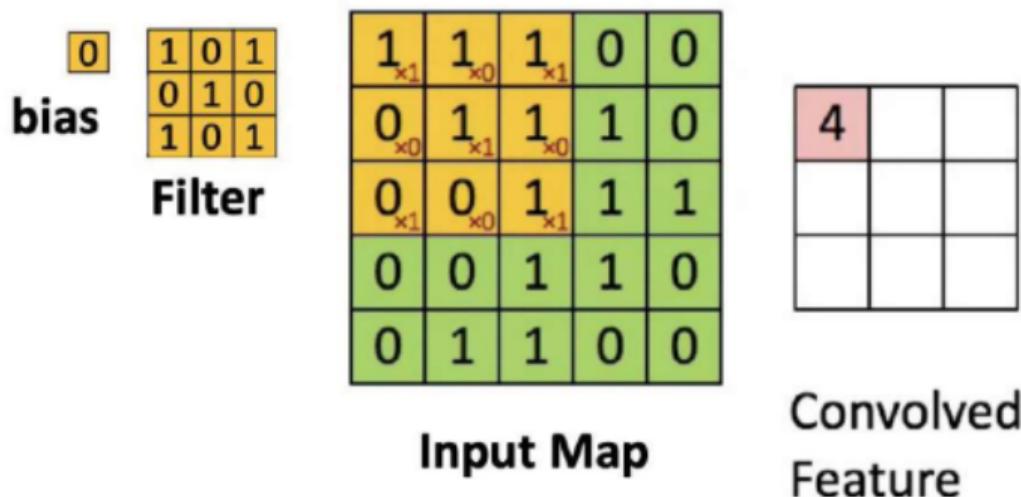
## 5 Receptive field

## 6 Inductive bias

## 7 Backpropagation

## 8 References

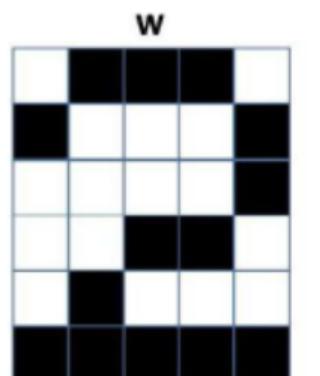
# What is a Convolve



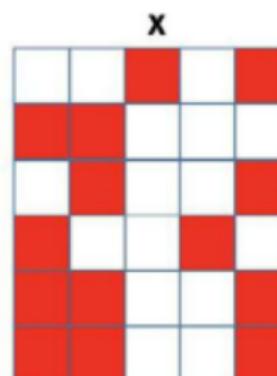
## Scanning an image with a "filter"

- A filter is really just a perceptron, with weights and a bias.
- At each location, the filter and the underlying map values are multiplied component-wise, and the products are added along with the bias.

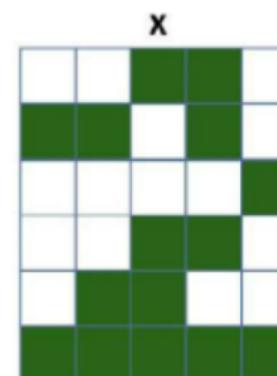
# Weights showing correlation



$$o = \begin{cases} 1 & \text{if } \sum_i w_i x_i \geq T \\ 0 & \text{else} \end{cases}$$



Correlation = 0.57

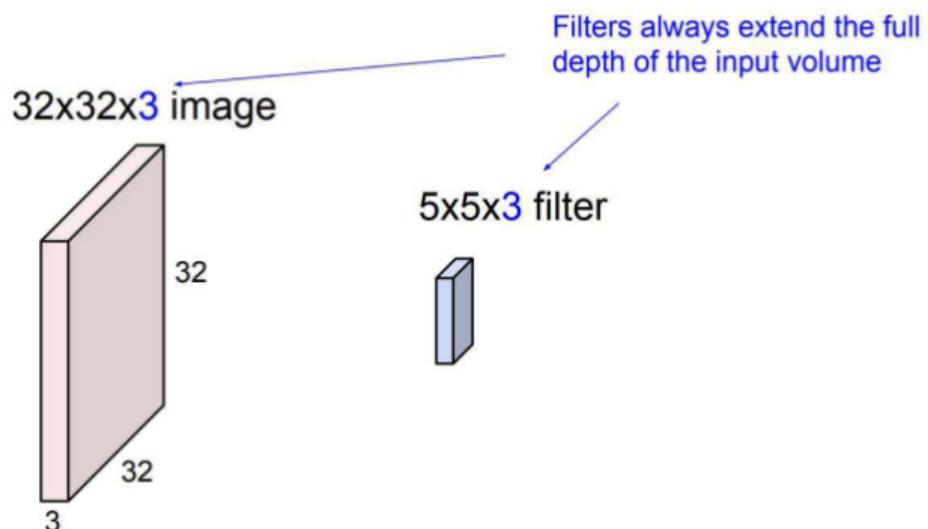


Correlation = 0.82



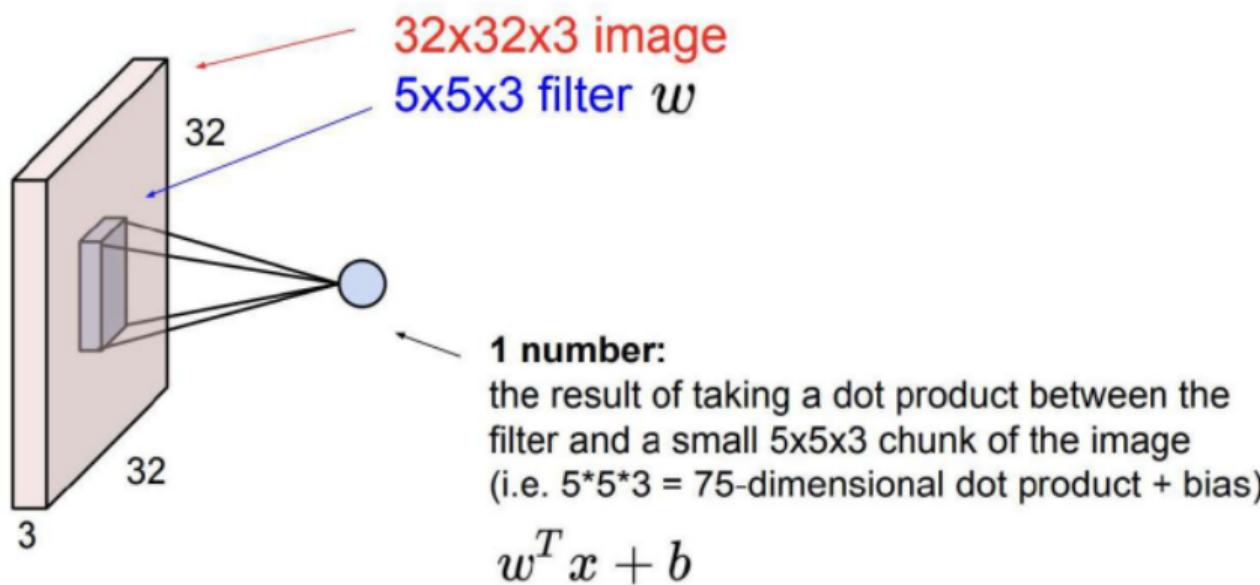
- The weights of the filter represent the appearance of the number "2".
- The **green** has a higher correlation with the filter compared to the **red**.
- The **green pattern** is more likely to represent the number "2".

# What is convolution



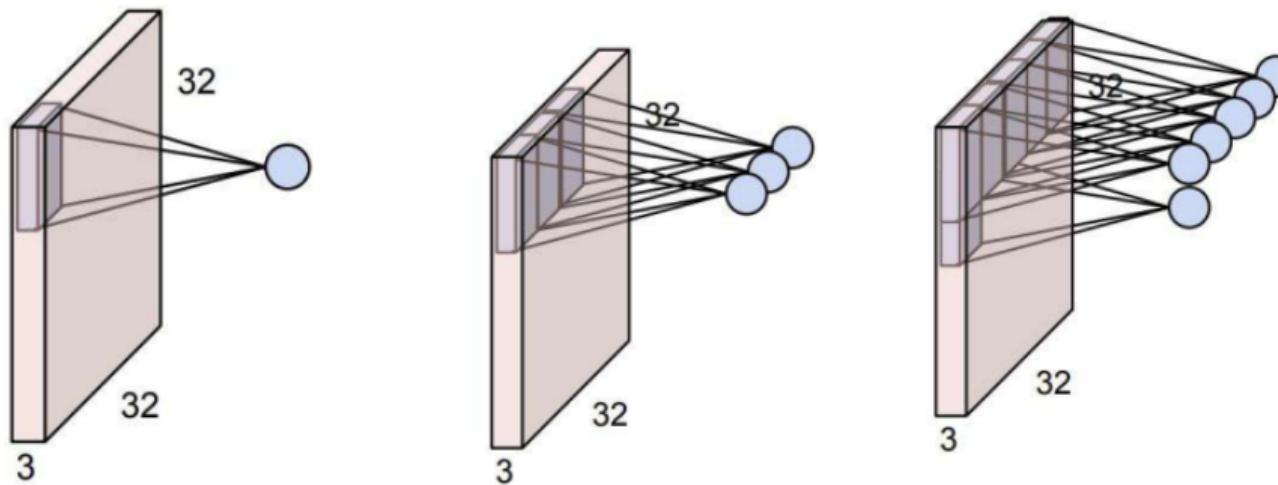
- **Convolve:** Slide over the image spatially, computing dot products.
- This allows us to preserve the spatial structure of the input.

# What is the output

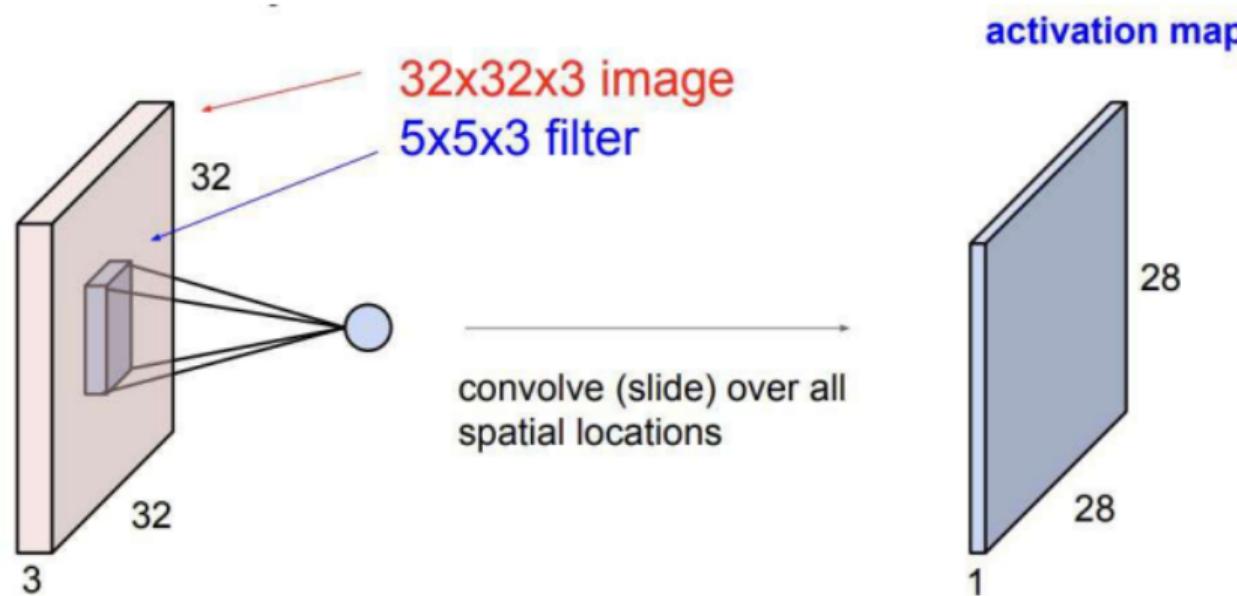


It's simply a neuron with local connectivity!

# Convolution process (1)

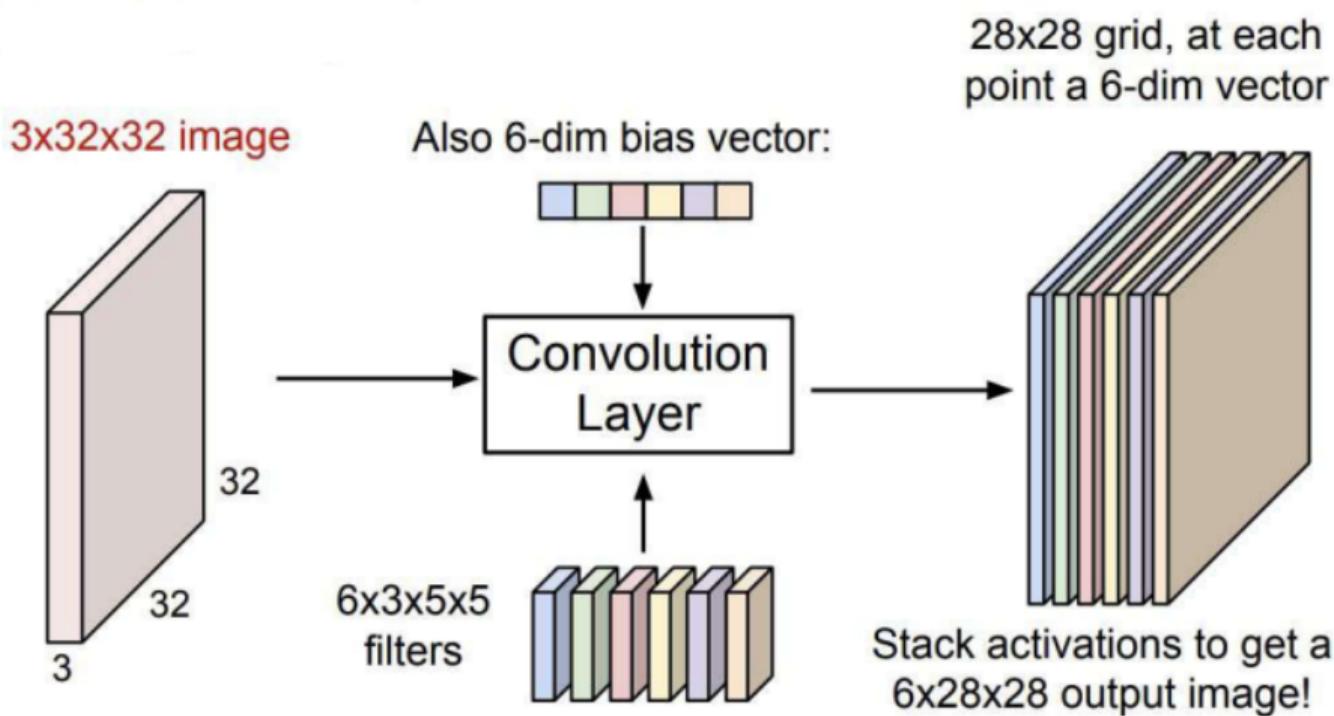


## Convolution process (2)

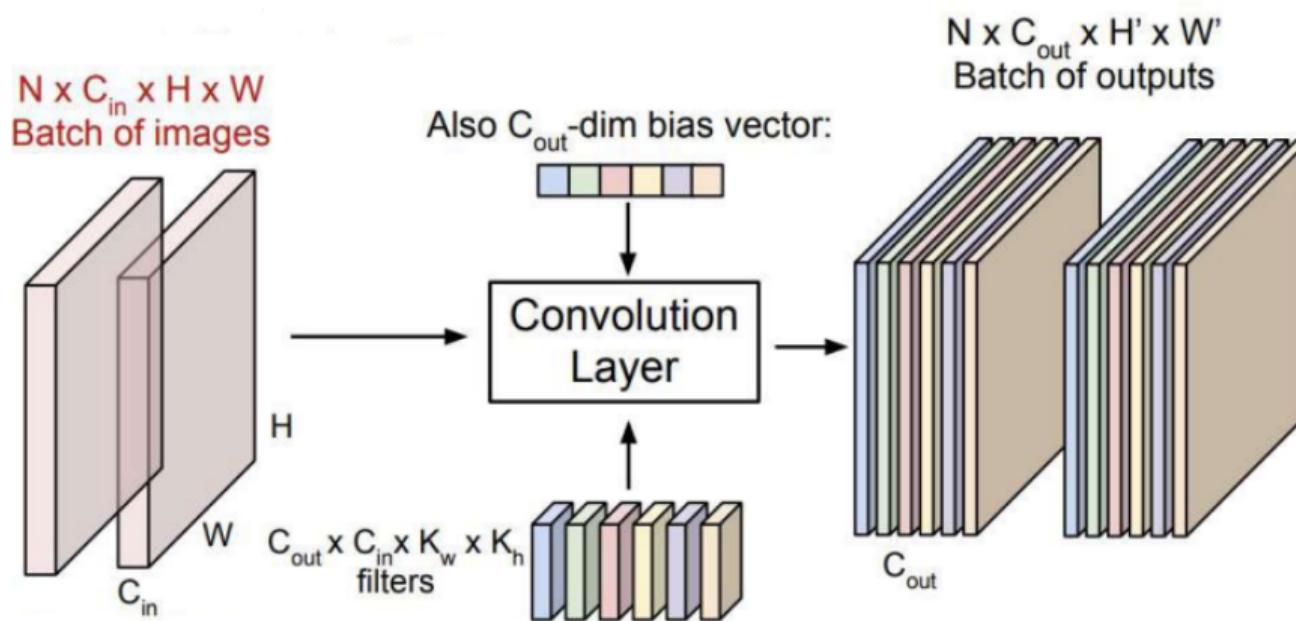


If we consider the **image** to be of size  $n \times m \times b$  and the **filter** to be of size  $n' \times m' \times b$ , we will have an **activation map** of size  $(n - n' + 1) \times (m - m' + 1) \times 1$  for each filter.

## Convolution process (3)

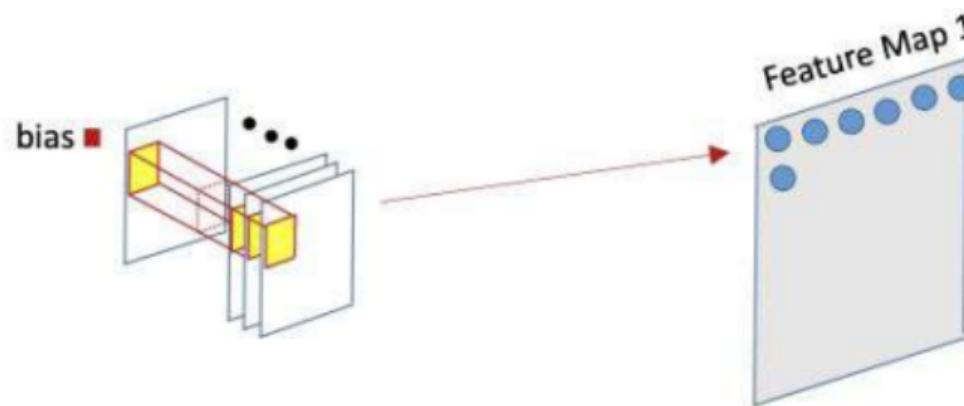


## Convolution process (4)



By batching the input images, we can expect a more generalized view as described above.

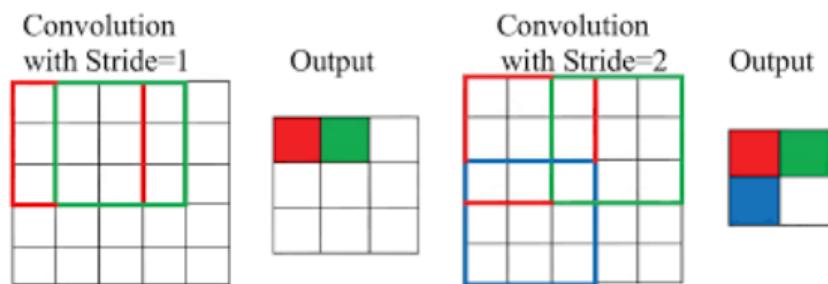
# A different view



$$z(i,j,s) = \underbrace{\sum_p \sum_{k=1}^L \sum_{l=1}^L w(k, l, p, s) a(i + l - 1, j + k - 1, p)}_{\text{One map}} + b(s)$$

$\underbrace{\qquad\qquad\qquad}_{\text{All maps}}$

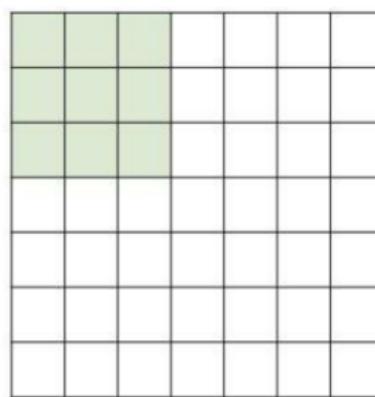
# Stride (1)



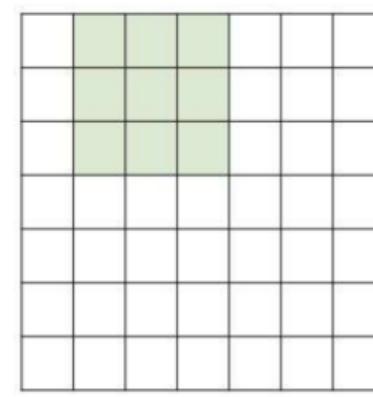
- The scans of the individual “filters” may advance by more than one pixel at a time
  - The “stride” may be greater than 1
  - Effectively increasing the granularity of the scan
    - Saves computation, sometimes at the risk of losing information
- This will result in a reduction of the size of the resulting maps
  - They will shrink by a factor equal to the stride
- This can happen at any layer

## Stride (2)

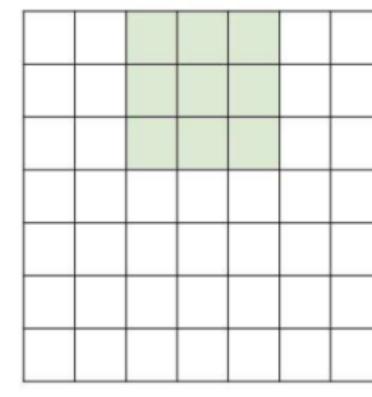
7



7



7

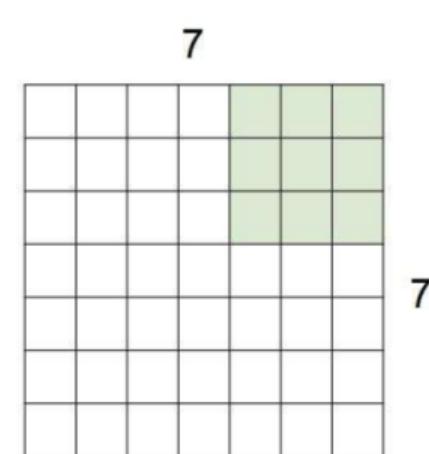
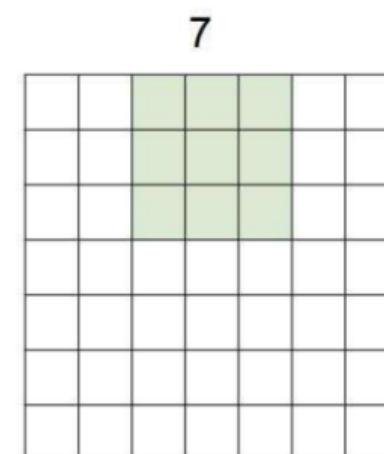
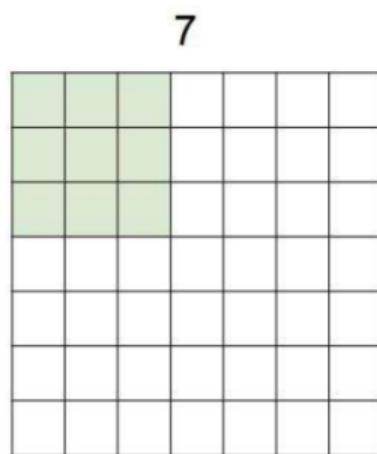


7

stride : 1

output dimension :  $5 \times 5$

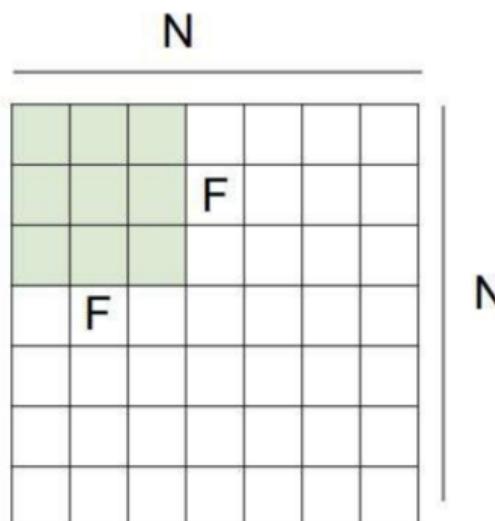
## Stride (3)



stride : 2

output dimension :  $3 \times 3$

## Stride (4)

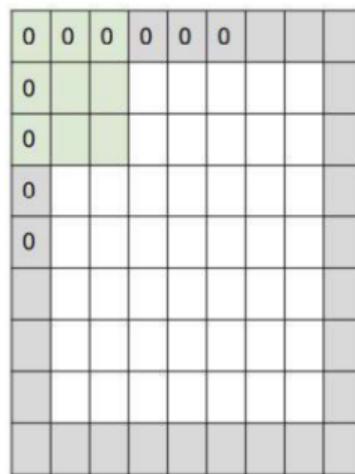


Output size:  
 $(N - F) / \text{stride} + 1$

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

- We can't choose any stride blindly!
- The output dimension should be an integer.
- This is why padding comes into the picture.

# Padding (1)



In practice, it is common to zero-pad the border. **Recall:**  
The original formula for convolution without padding:

$$\frac{N - F}{\text{stride}} + 1$$

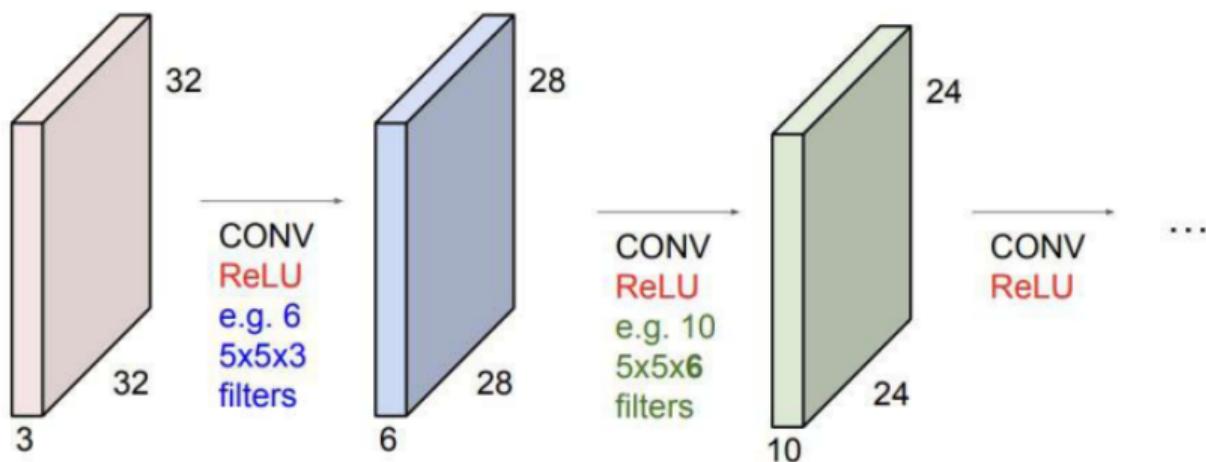
Now, we should adjust the formula considering padding  $P$ :

$$\frac{N + 2P - F}{\text{stride}} + 1$$

## Padding (2)

- Zero-padding is used not only for stride  $> 1$ , but also to prevent a reduction in output size even when  $S = 1$ .
- For stride  $> 1$ , zero padding is adjusted to ensure that the size of the convolved output is  $\lceil \frac{N}{S} \rceil$   
This is achieved by zero padding the image with  $P = S \lceil \frac{N}{S} \rceil - N$
- For an  $F$  width filter:
  - **Odd  $F$ :** Pad on both left and right with  $\frac{F-1}{2}$  columns of zeros
  - **Even  $F$ :** Pad one side with  $\frac{F}{2}$  columns of zeros, and the other with  $\frac{F}{2} - 1$  columns of zeros
  - The resulting image is width  $N + F - 1$
- The top/bottom zero padding follows the same rules to maintain map height after convolution

## ConvNet (1)



- **ConvNet:** sequence of convolution layers, interspersed with activation functions.

## ConvNet (2)

**Question:** What do convolutional filters at different levels of a ConvNet learn?

**Answer:**

- The filters in the early layers typically learn to detect simple features, such as edges, textures, and basic shapes.
- As we move deeper into the network, the filters learn more complex and abstract features, such as specific parts of objects (e.g., a nose, eyes, or other high-level patterns).

## Example-1

**Input volume:** 32x32x3

10 5x5 filters with stride 1, pad 2

**Output volume size:**

$$\left( \frac{32 + 2 \times 2 - 5}{1} \right) + 1 = 32 \text{ spatially, so } 32x32x10$$

## Example-2

**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)

$$\Rightarrow 76*10 = 760$$

# Parameter Setting (1)

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires four hyper-parameters:
  - 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:
  - $W_2 = \left(\frac{W_1 - F + 2P}{S}\right) + 1$
  - $H_2 = \left(\frac{H_1 - F + 2P}{S}\right) + 1$
  - $D_2 = K$

## Common settings:

- $K = \text{powers of 2}$  (e.g., 32, 64, ...)
- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = \text{whatever fits!}$
- $F = 1, S = 1, P = 0$

## Parameter Setting (2)

- 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 the  $d$ -th bias.

## 1 Motivation

## 2 CNN vs. FCN

## 3 Convolution

## 4 Pooling

## 5 Receptive field

## 6 Inductive bias

## 7 Backpropagation

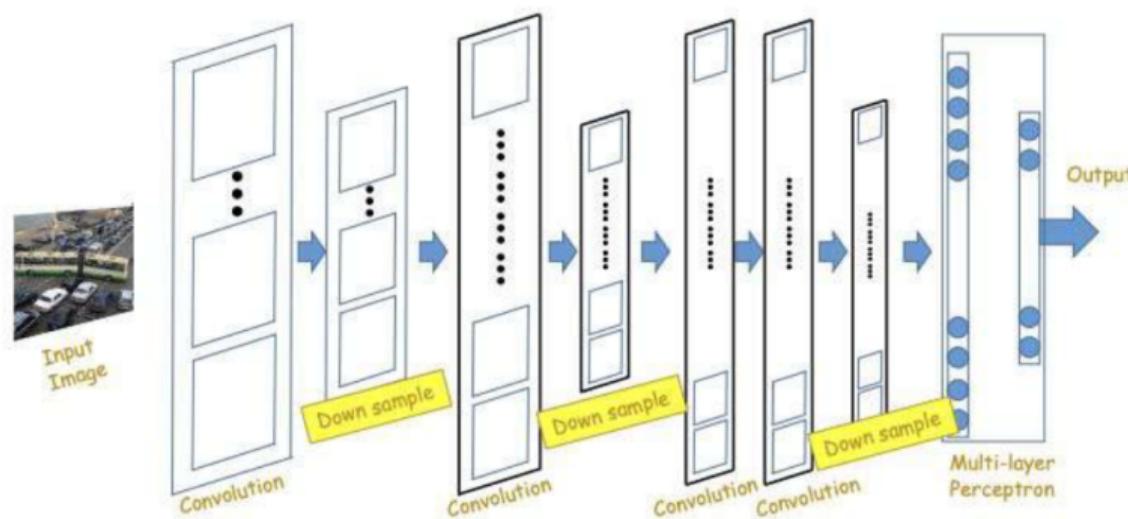
## 8 References

# Review

## Three main types of layers

- **Convolutional Layer**
  - Output of neurons are connected to local regions in the input.
  - Applying the same filter on the whole image.
  - CONV layers parameters consist of a set of learnable filters.
- **Pooling Layer**
  - Performs a downsampling operation along the spatial dimensions.
- **Fully-Connected Layer**
  - Typically used in the final stages of the network to combine high-level features and make predictions.

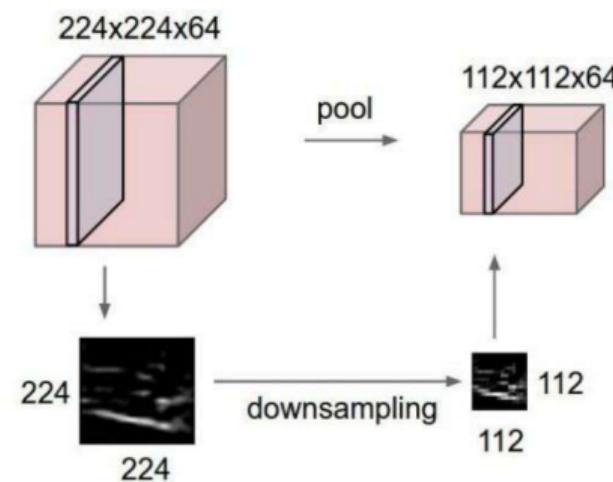
# Pooling (1)



Convolution activation layers are followed intermittently by pooling layers.

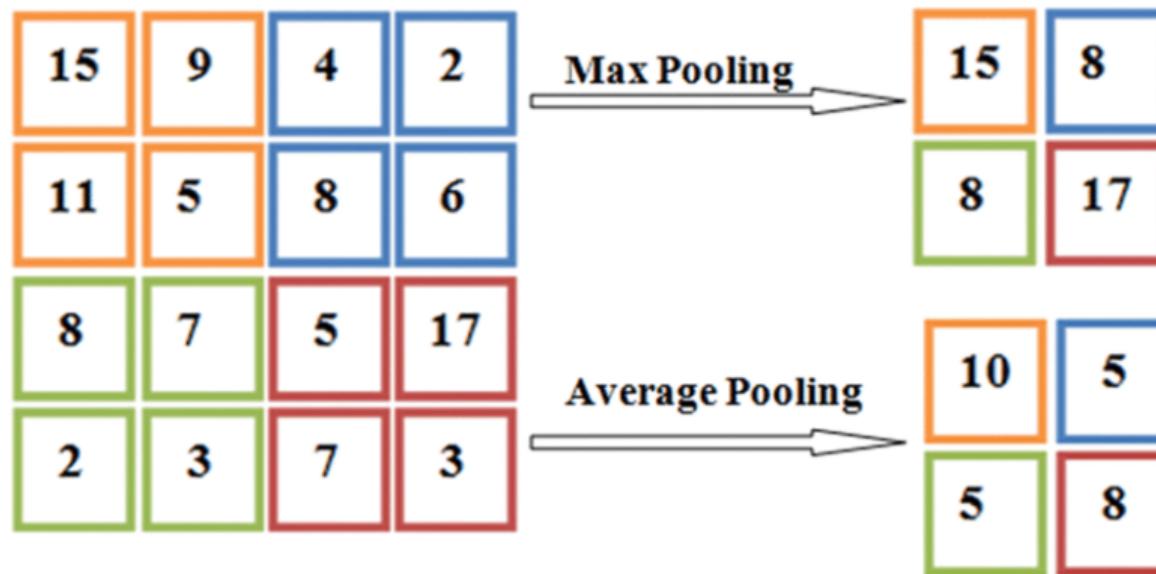
- Often, they alternate with convolution, though this is not necessary.

## Pooling (2)



- Reduce the spatial size of the representation.
  - To reduce the amount of parameters and computation in the network.
  - To control overfitting.
- Help the network to become invariant to small translations or distortions.
- Operates over each activation map independently.

## Pooling (3)



**Two main types of pooling:**

- **Max Pooling:** Selects the maximum value from each patch of the feature map.
- **Average Pooling:** Computes the average of the values in each patch of the feature map.

# Parameter Setting (1)

## Question:

How does a pooling layer affect the dimensions of the input image?

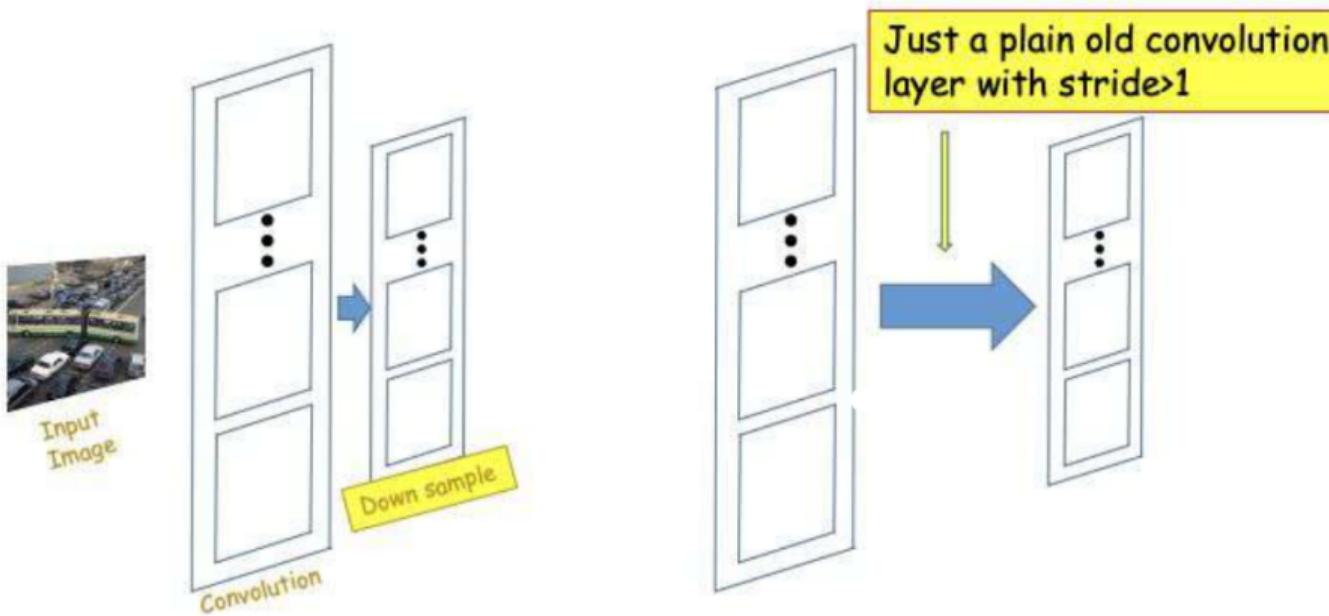
## Answer:

- An  $N \times N$  picture compressed by a  $P \times P$  pooling filter with stride  $S$  results in an output map of side  $\frac{(N-P)}{S} \lceil +1$ .
- Typically, pooling layers do not use zero-padding.

## Parameter Setting (2)

- Pooling 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:
  - $W_2 = \frac{(W_1 - F)}{S} + 1$
  - $H_2 = \frac{(H_1 - F)}{S} + 1$
  - $D_2 = D_1$
- 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.

## Alternative to Pooling



Downsampling be done by a simple convolution layer with stride larger than 1,  
Replacing the maxpooling layer with a conv layer.

## 1 Motivation

## 2 CNN vs. FCN

## 3 Convolution

## 4 Pooling

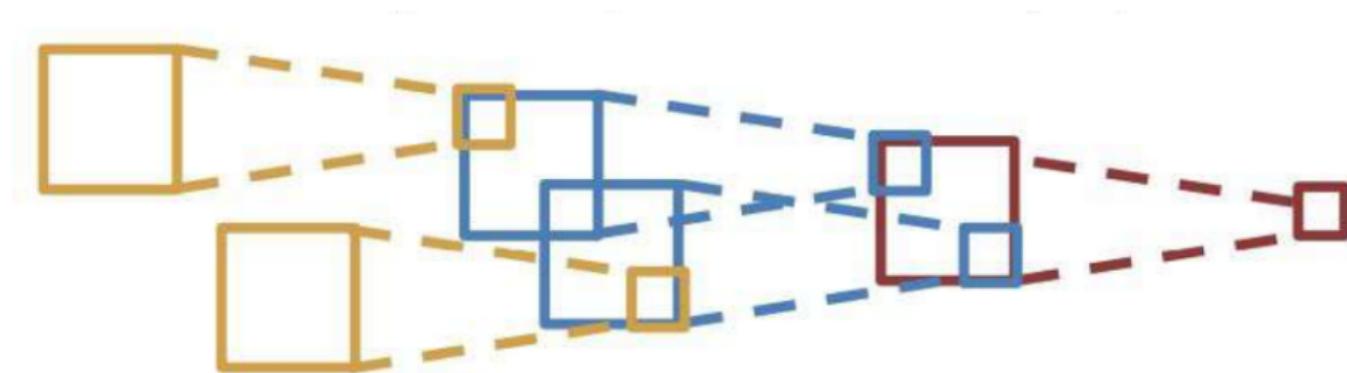
## 5 Receptive field

## 6 Inductive bias

## 7 Backpropagation

## 8 References

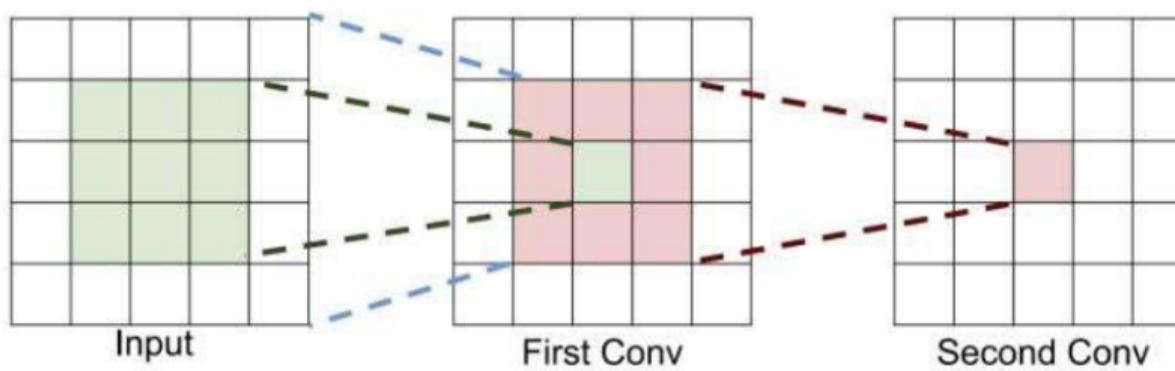
# Receptive Field (1)



**Receptive Field :** How big of a region in the **input or previous** layer does a neuron on the n-th conv-layer see?

For convolution with kernel size  $K$ , each element in the next layer depends on a  $K \times K$  **receptive field** in the previous layer.

## Receptive Field (2)



- Units in the deeper layers can be **indirectly** connected to most of the input image.
- Each successive convolution adds  $K - 1$  to the receptive field size. With  $L$  layers, the receptive field size is  $1 + L \cdot (K - 1)$ .
- Problem:** For large images, we need many layers for each output to "see" the whole image.
  - Solution:** Downsample inside the network using strides and pooling.

# Power of Small Filters

Suppose the input is  $H \times W \times C$  and we use convolutions with  $C$  filters to preserve depth (stride 1, padding to preserve  $H, W$ ).

## one CONV with 7 x 7 filters

Number of weights

$$= C \times (7 \times 7 \times C) = \mathbf{49}C^2$$

## three CONV with 3 x 3 filters

Number of weights

$$= 3 \times C \times (3 \times 3 \times C) = \mathbf{27}C^2$$

Both options result in a receptive field of 7, but using multiple smaller filters reduces the number of parameters, introduces more nonlinearity, and generally leads to a more efficient and expressive model.

## 1 Motivation

## 2 CNN vs. FCN

## 3 Convolution

## 4 Pooling

## 5 Receptive field

## 6 Inductive bias

## 7 Backpropagation

## 8 References

# Inductive Bias in CNNs

## Inductive Bias:

Refers to the assumptions a model incorporates to generalize from training data to unseen data.

## Key Features of Inductive Bias in CNNs:

- **Weight Sharing:**

A single filter is applied across different regions of the input, significantly reducing the number of parameters.

- **Locality:**

CNNs use small filters (e.g.,  $3 \times 3$ ) that focus on local regions, aligning well with image data where local structures are important.

- CNNs are more sample-efficient than FCNs due to their inductive biases.

- CNNs reduce sample complexity from  $O(d^2)$  in FCNs to  $O(\log^2 d)$ .

## 1 Motivation

## 2 CNN vs. FCN

## 3 Convolution

## 4 Pooling

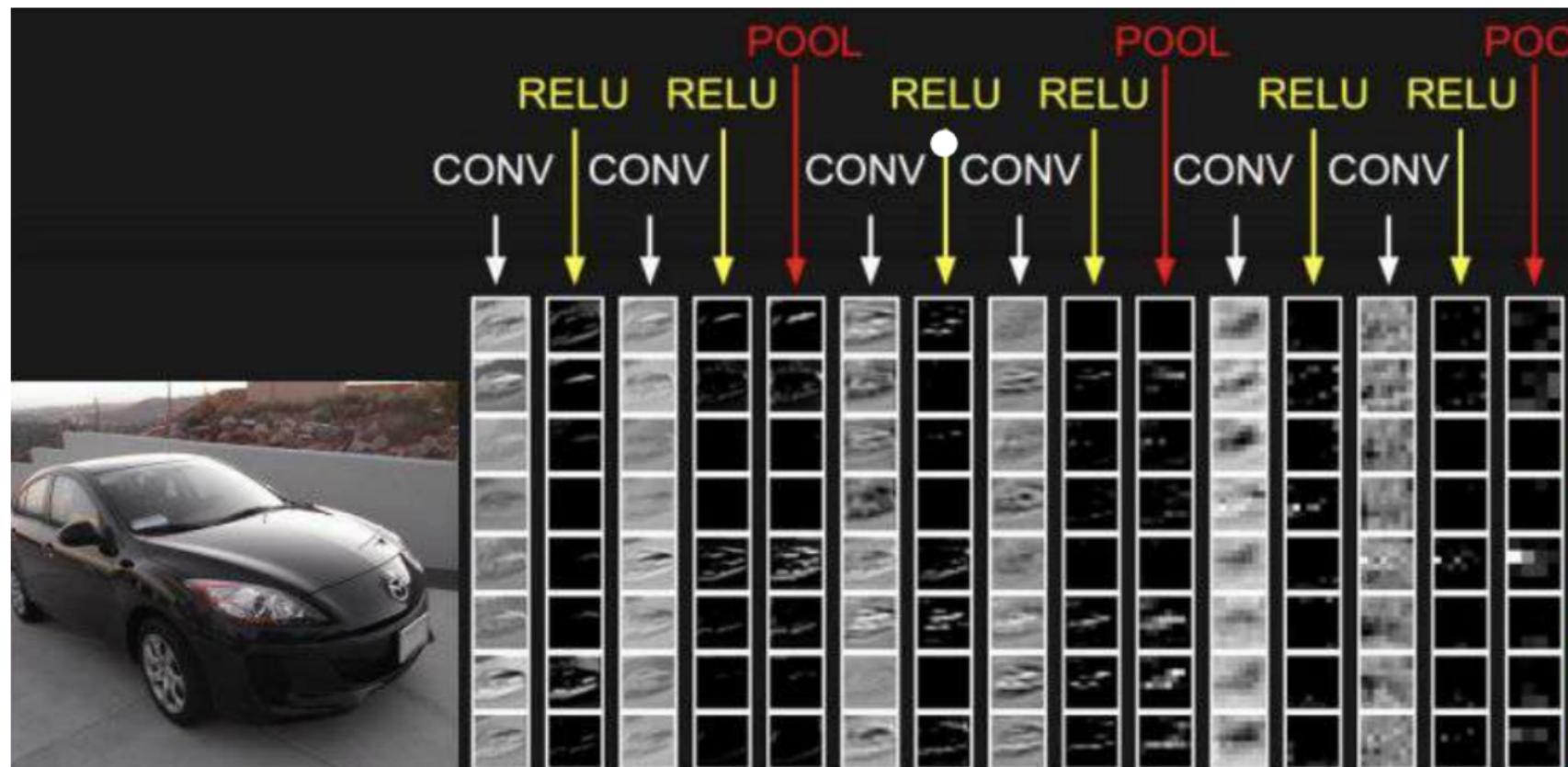
## 5 Receptive field

## 6 Inductive bias

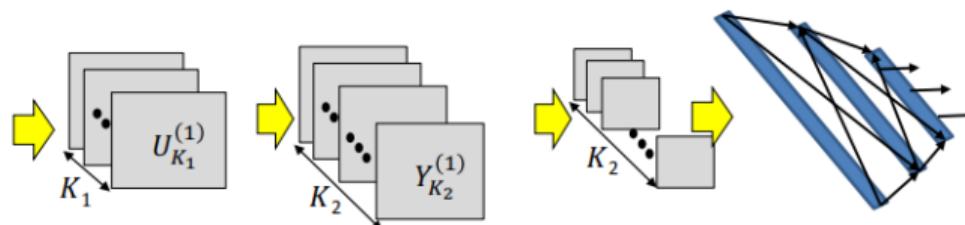
## 7 Backpropagation

## 8 References

# What we have

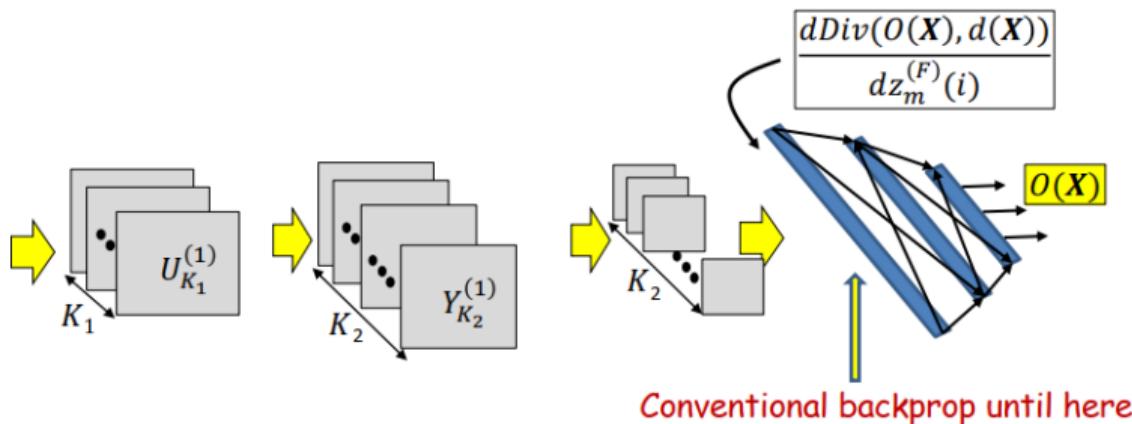


# Training (1)



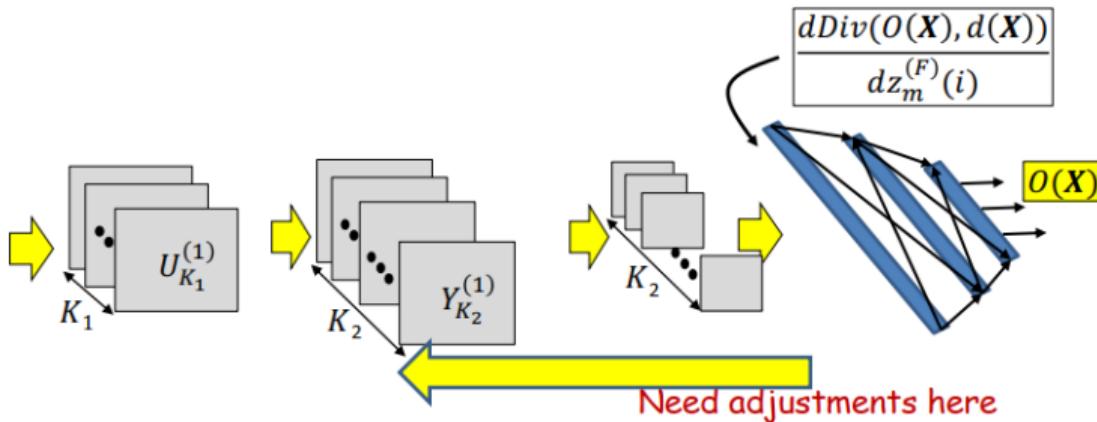
- Training is as in the case of the regular MLP
- Training examples of (Image, class) are provided
- A divergence between the desired output and the true output of the network in response to any input
- **Network parameters are trained through variants of gradient descent**
- **Gradients are computed through backpropagation**

## Training (2)



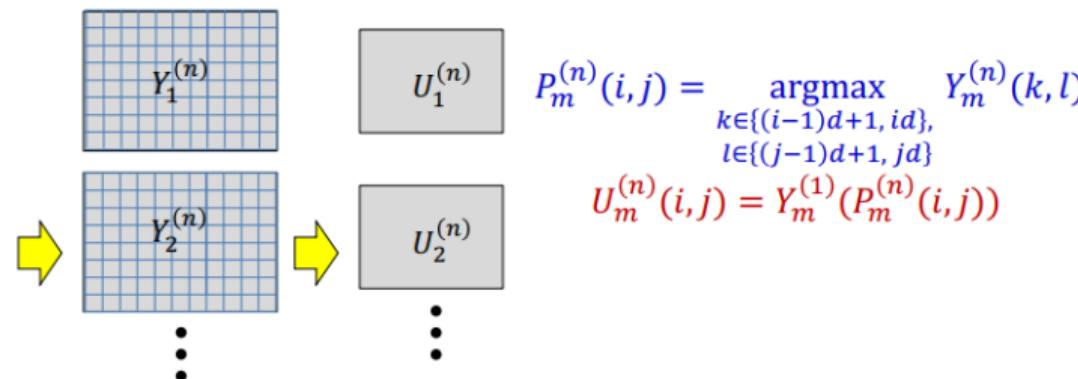
- Backpropagation continues in the usual manner until the computation of the derivative of the divergence w.r.t the inputs to the first “flat” layer

## Training (3)



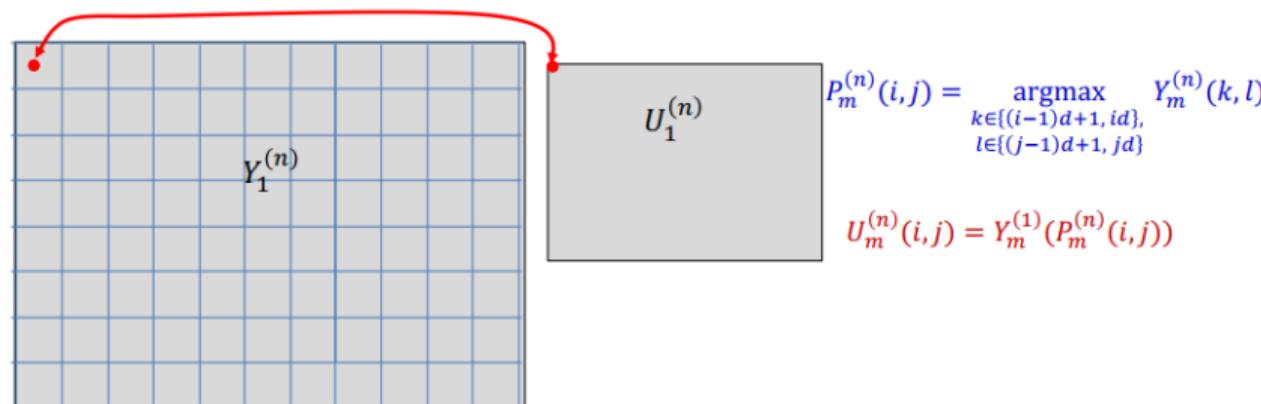
- Backpropagation from the flat MLP requires special consideration of:
  - The pooling layers (particularly Maxout)
  - The shared computation in the convolution layers

## Backpropagation maxout (1)



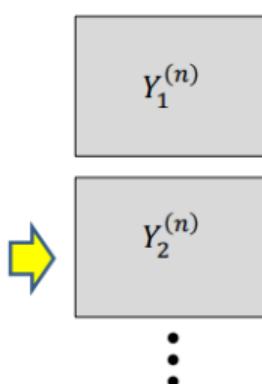
- The derivative w.r.t  $U_m^{(n)}(i,j)$  can be computed via backprop.
- But this cannot be propagated backwards to compute the derivative w.r.t.  $Y_m^{(n)}(k,l)$ .
- **Max and argmax are not differentiable.**

## Backpropagation maxout (2)



- $$\frac{dDiv(O(X), d(X))}{dY_m^{(n)}(k, l)} = \begin{cases} \frac{dDiv(O(X), d(X))}{dU_m^{(n)}(i, j)} & \text{if } (k, l) = P_m^{(n)}(i, j) \\ 0 & \text{otherwise} \end{cases}$$
- Approximation:** Derivative w.r.t the  $Y$  terms that did not contribute to the maxout map is 0.

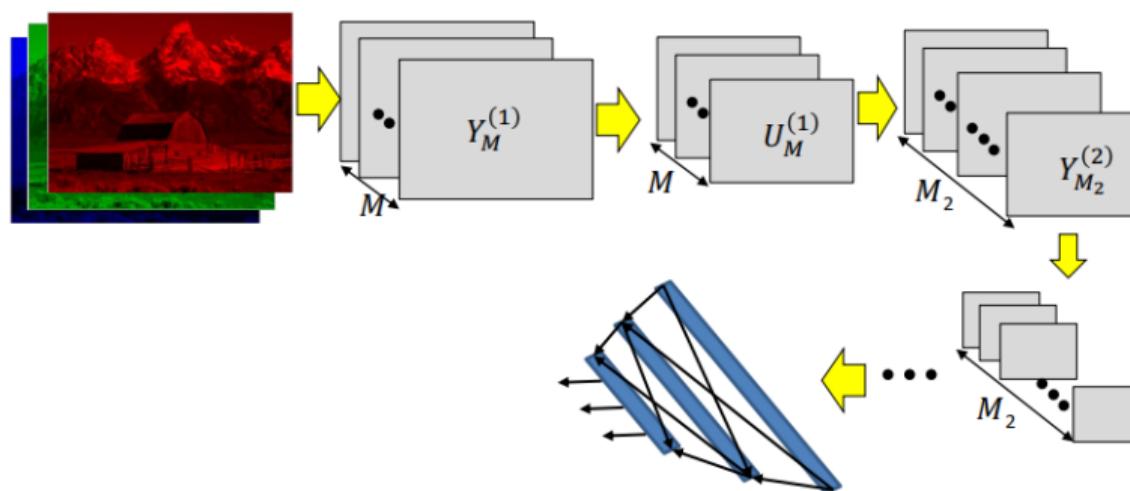
# Backpropagation weights



$$z_m^{(n)}(i, j) = \sum_{r=1}^{M_{n-1}} \sum_{k=1}^{L_n} \sum_{l=1}^{L_n} w_m^{(n)}(r, k, l) U_r^{(n-1)}(i + k, j + l)$$
$$Y_m^{(n)}(i, j) = f(z_m^{(n)}(i, j))$$

- **Note:** each weight contributes to *every* position in the map at the output of the convolutional layer.
- **Every position will contribute to the derivative of the weight**
  - Shared parameter updates

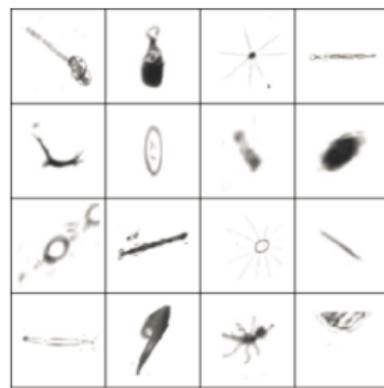
# Learning the network



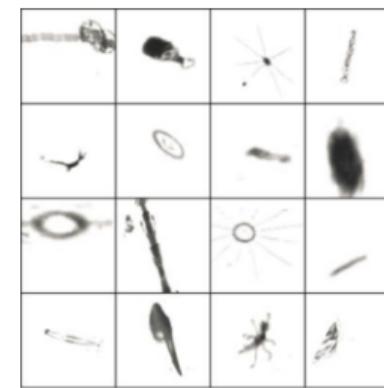
- Have shown the derivative of divergence w.r.t every intermediate output, and every free parameter (filter weights).
- Can now be embedded in the gradient descent framework to learn the network.

# A Problem

Original data



Augmented data



**Question:** How can we improve generalization? **Answer:** Rotation, Translation, Rescaling, Flipping, etc.

# Design Choices

- **Number of convolutional and downsampling layers**
  - And arrangement (order in which they follow one another)
- **For each convolution layer:**
  - Number of filters  $d^l$
  - Spatial extent of filter  $F^l \times F^l$ 
    - The "depth" of the filter is fixed by the number of filters in the previous layer  $d^{l-1}$
  - The stride  $S^l$
- **For each downsampling/pooling layer:**
  - Spatial extent of filter  $P^l \times P^l$
  - The stride  $S^l$
- **For the final MLP:**
  - Number of layers, and number of neurons in each layer

# Typical Architecture for CNNs

- $[(\text{CONV-RELU})^N \text{ POOL?}]^M (\text{FC-RELU})^K \text{ SOFTMAX}$ 
  - Where N is usually up to  $\sim 5$
  - M is large
  - $0 \leq K \leq 2$

**But recent advances such as ResNet and GoogleNet challenge this paradigm!**

# Conclusion

- Highly structured nature of image data
- Local correlations are important
- CNNs have local and shared connections
- CNNs make strong inductive biases
- To exploit the two-dimensional structure of image data to create inductive biases, we can use four interrelated concepts:
  - Hierarchy
  - Locality
  - Invariance

## 1 Motivation

## 2 CNN vs. FCN

## 3 Convolution

## 4 Pooling

## 5 Receptive field

## 6 Inductive bias

## 7 Backpropagation

## 8 References

content...