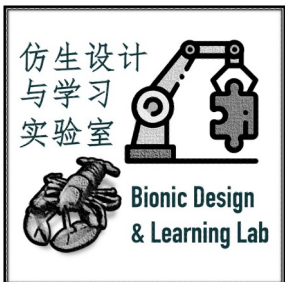


# Lecture 07

## Deep Networks II

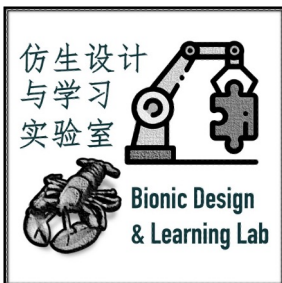


AncoraSIR.com

[Please refer to the course website for copyright credits]



# Convolutional Networks

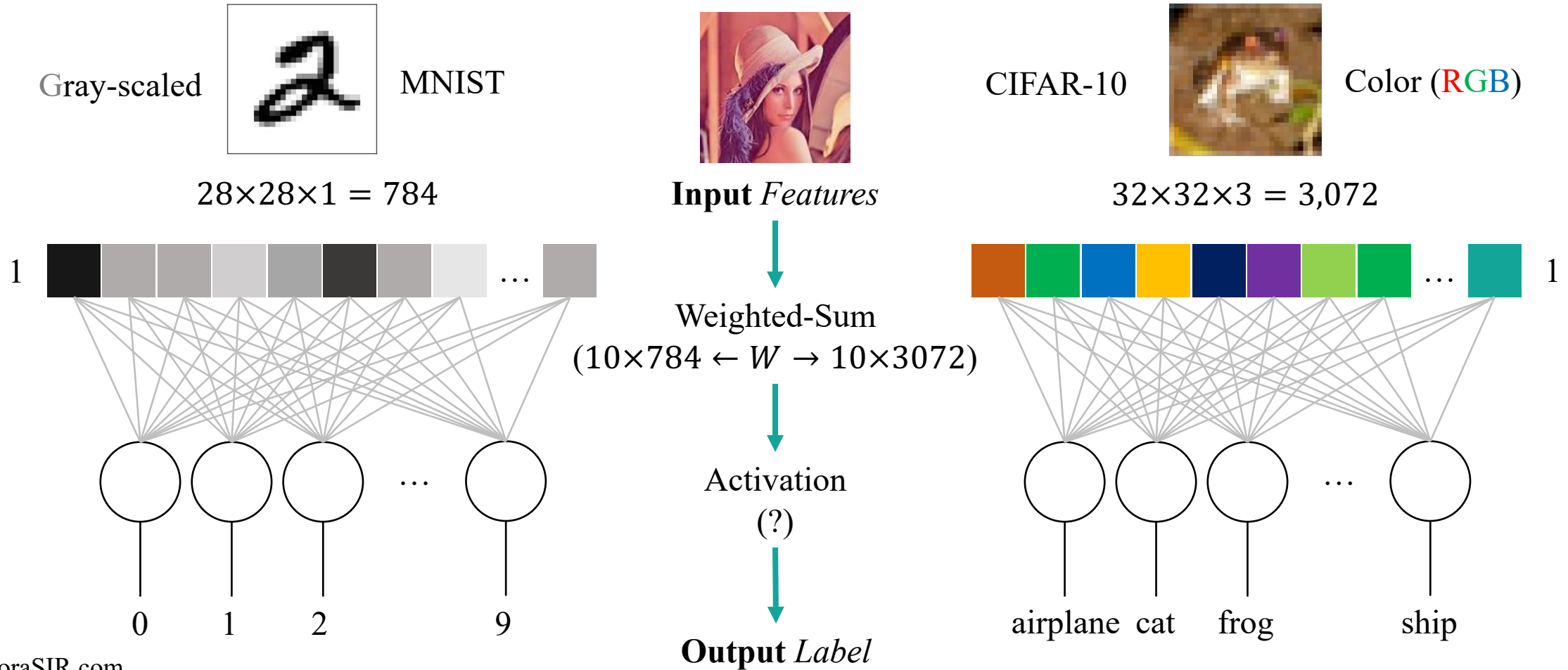


[AncoraSIR.com](http://AncoraSIR.com)

# A Design Challenge with Increasing Dimensions

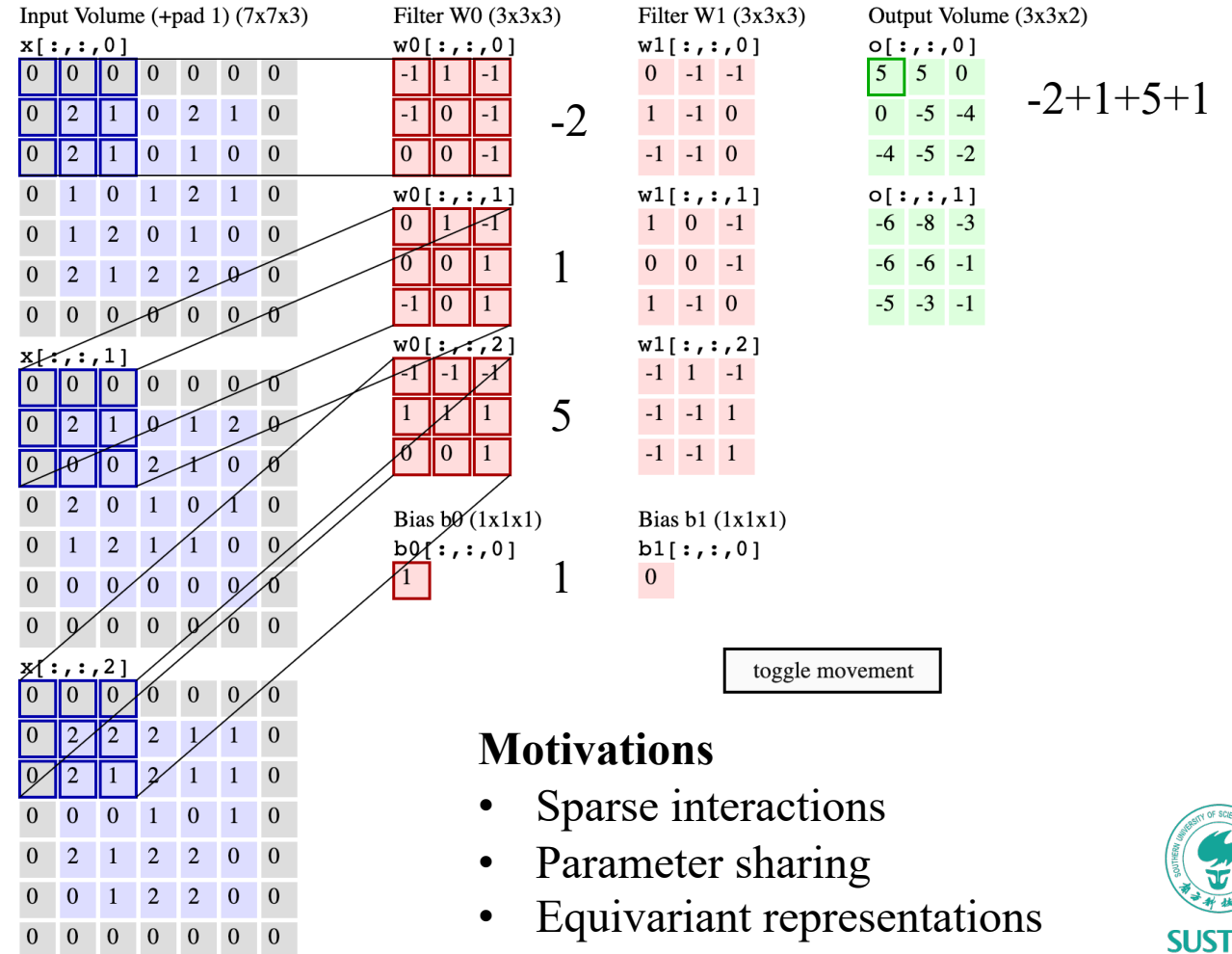
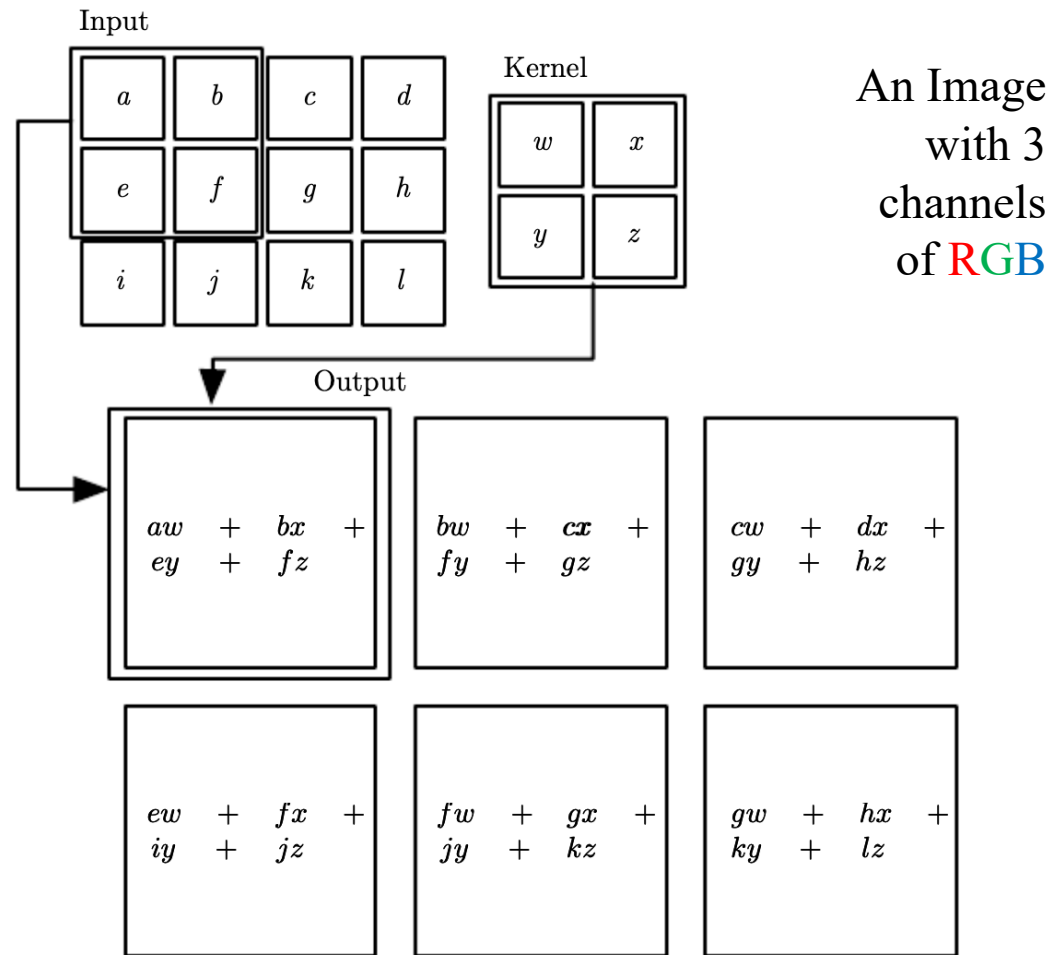
*Regular Neural Nets don't scale well to full images*

$$512 \times 512 \times 3 = 765,432$$



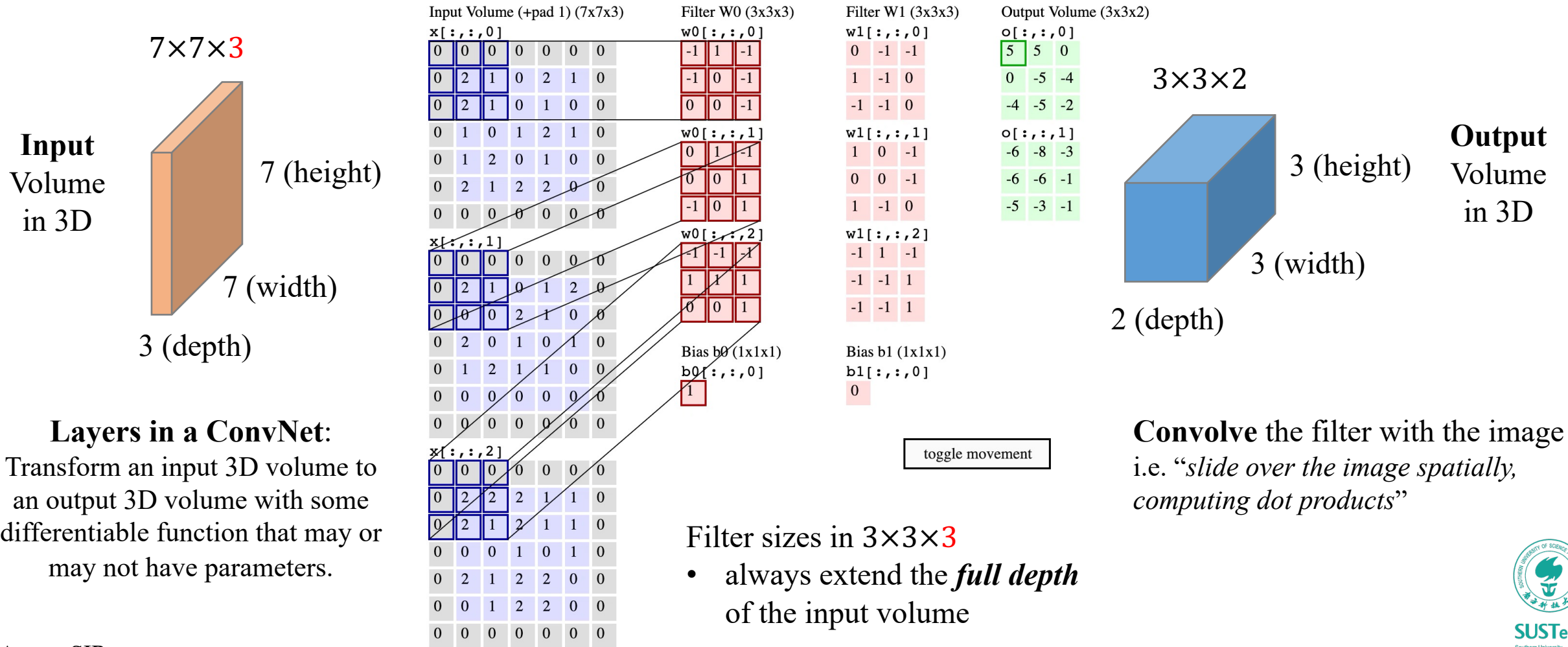
# Convolutional Operation

$$s(t) = \int x(a)w(t-a)da = (x * w)(t)$$



# Convolution in 3D Volumes

*Preserved spatial structure between the input and output volumes in width, height, number of channels*



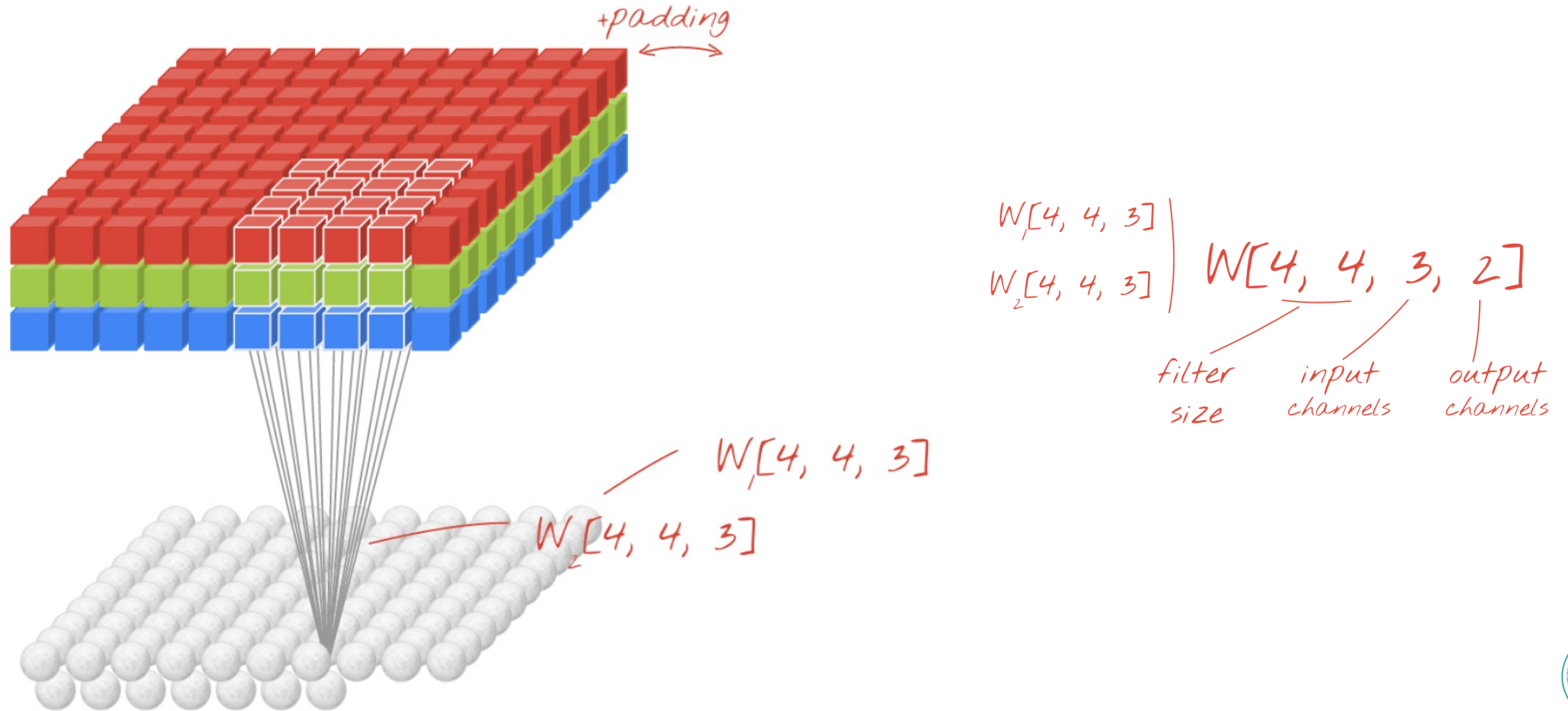
**Layers in a ConvNet:**  
 Transform an input 3D volume to an output 3D volume with some differentiable function that may or may not have parameters.

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

- Filter sizes in  $3 \times 3 \times 3$
- always extend the *full depth* of the input volume

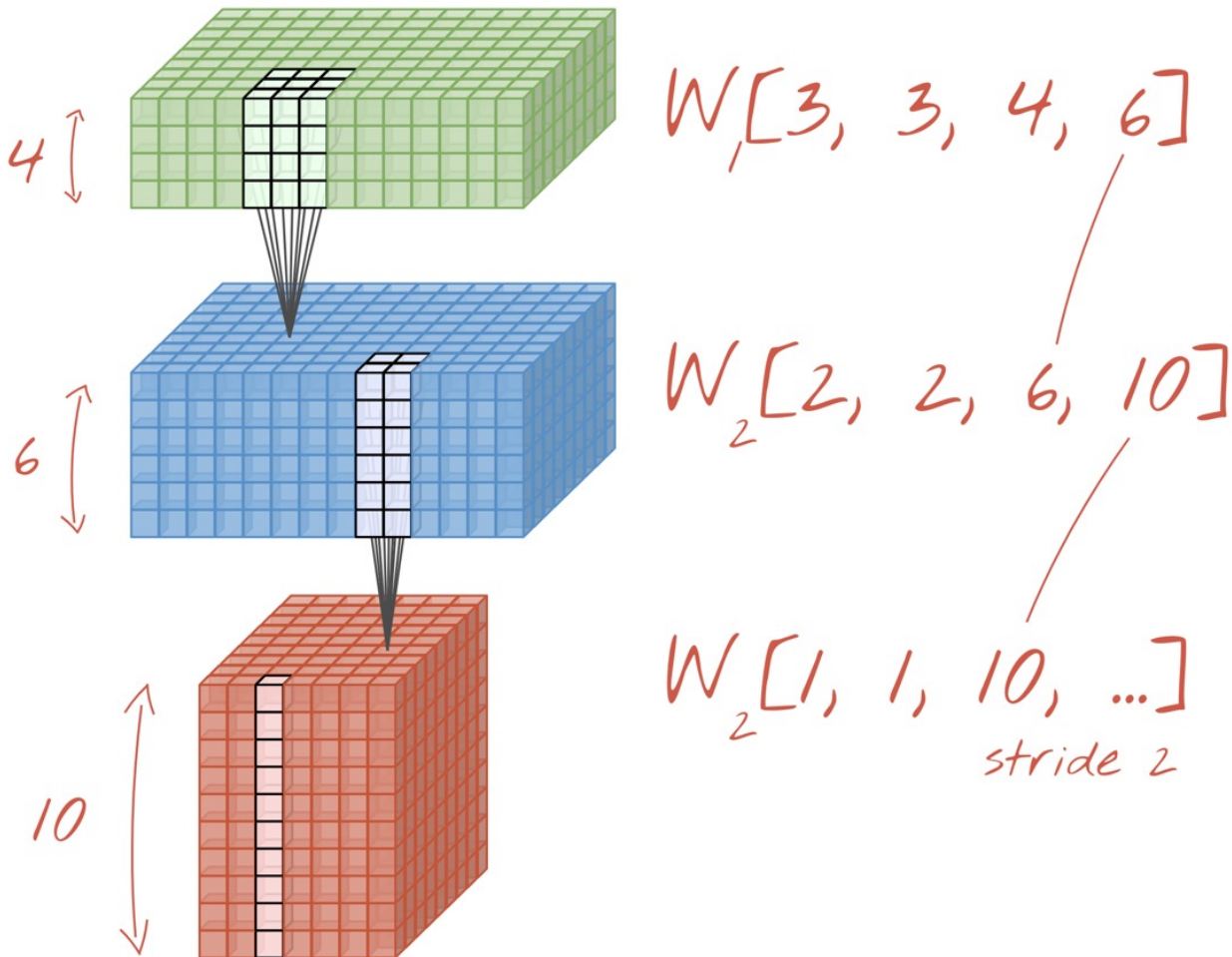
# The Design of a Convolutional Layer

*Defined by the filter (or kernel) size, the number of filters applied and the stride*

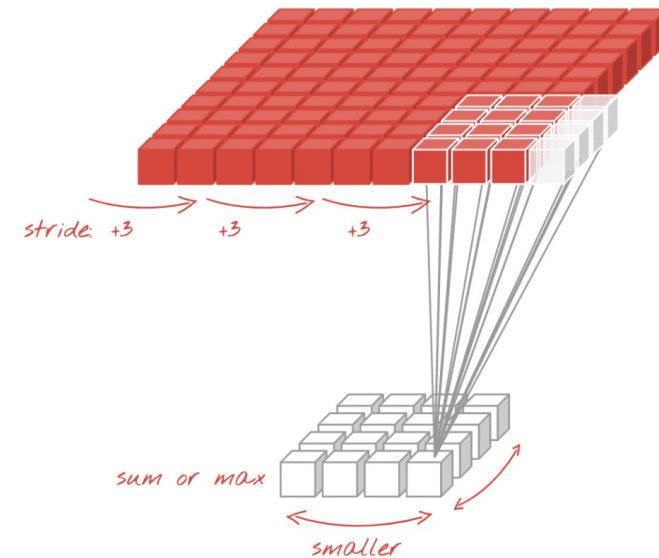


# Output Volume Size

Defined by the filter (or kernel) size, the number of filters applied and the stride



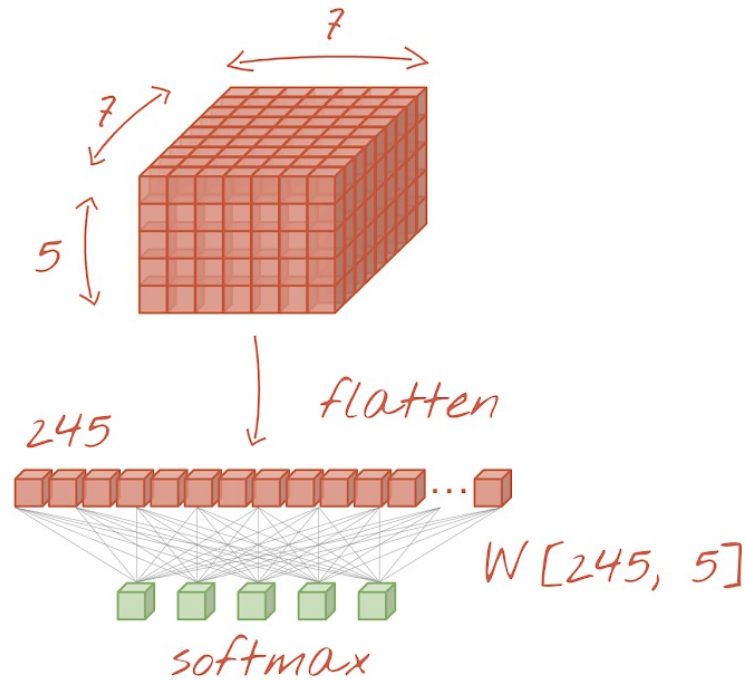
- Depth (number of channels):
  - adjusted by using more or fewer filters
- Width & Height:
  - adjusted by using a stride  $> 1$
  - (or with a max-pooling operation)



# The Last Layer

*From a Cubic Volume in 3D to predicted labels*

Fully connected layer

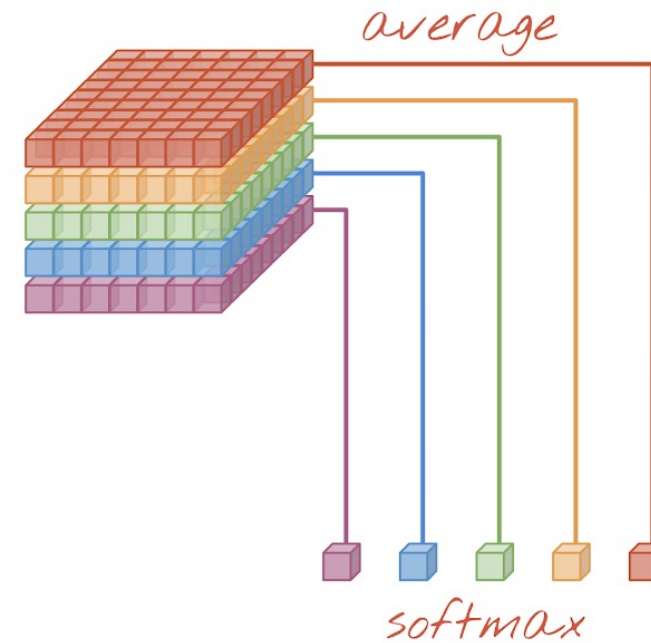


1225 weights

*cheaper* →

0 weights

Global average pooling



Much lighter in calculation

The average pooling explicitly discards all location data

Similar like a normal neural network

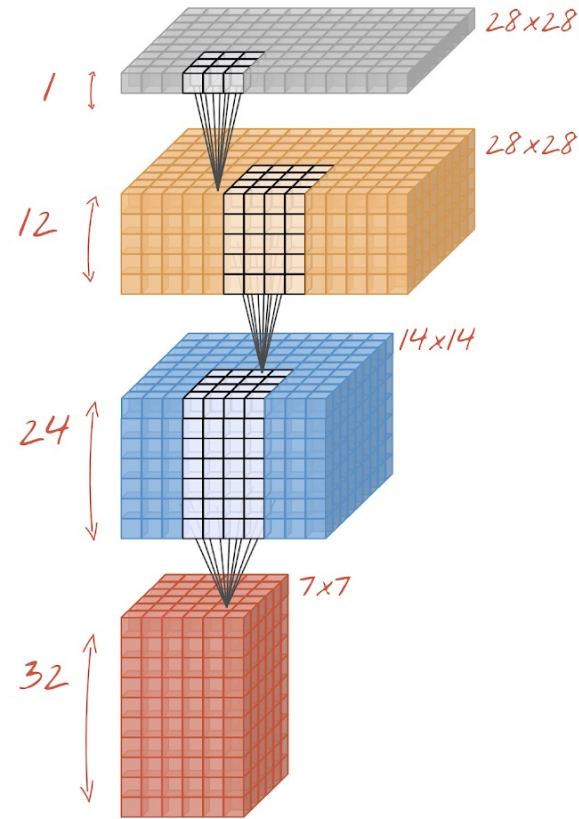
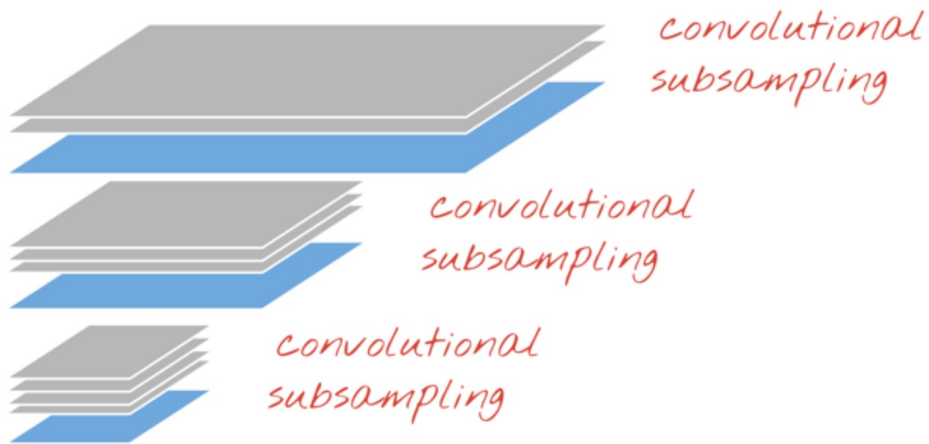
Expensive in #weights

But preserves the location data (x, y)



# Stacking Up a ConvNet

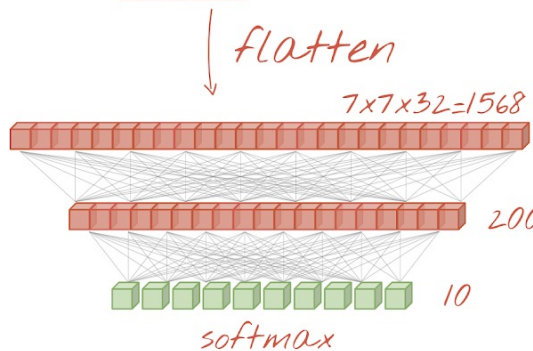
## Layer-by-layer



Convolutional 3x3 filters=12  
 $W_1[3, 3, 1, 12]$

Convolutional 6x6 filters=24  
 $W_2[6, 6, 12, 24]$  stride 2

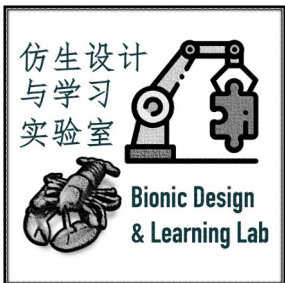
Convolutional 6x6 filters=32  
 $W_3[6, 6, 24, 32]$  stride 2



Dense layer  
 $W_4[1568, 200]$

Softmax dense layer  
 $W_5[200, 10]$

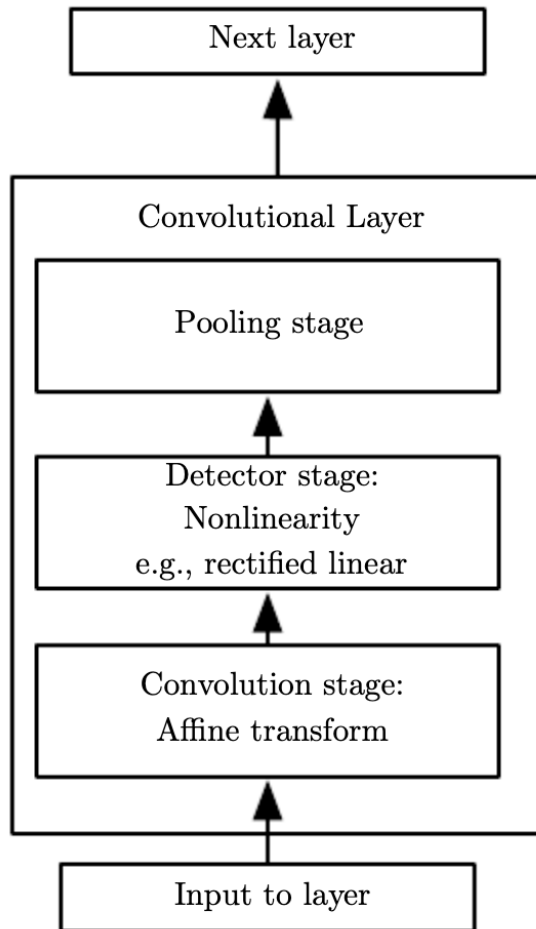
# Layers in ConvNets



[AncoraSIR.com](http://AncoraSIR.com)

# The Three Stages of a Typical ConvNet Layer

## *The Convolution, Detector and Pooling Stages*



- The maximum output within a rectangular neighborhood (max-pooling)
- The average of a rectangular neighborhood
- The L2 norm of a rectangular neighborhood
- A weighted average based on the distance from the central pixel

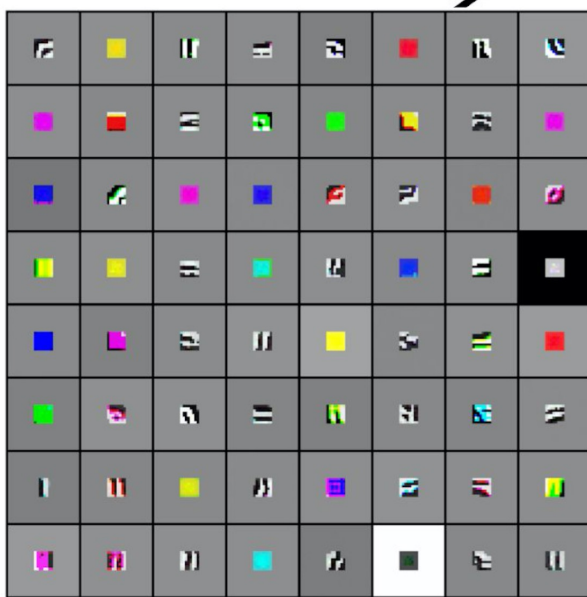
*Replace the output of the net at a certain location with a summary statistic of the nearby outputs (can be viewed as a further abstraction of the learned features)*

*Each linear activation is run through a nonlinear activation function, such as ReLU (can be viewed as activation function)*

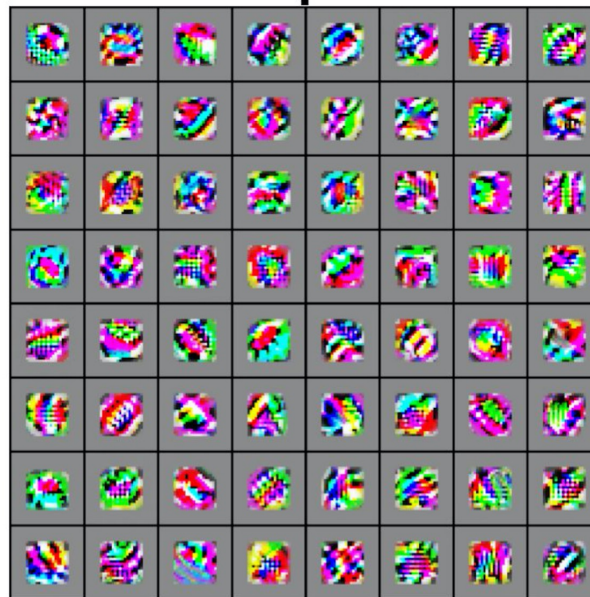
*Performs several convolutions in parallel to produce a set of linear activations (can be viewed as weighted-sum)*

# A Visualized Understanding of ConvNet

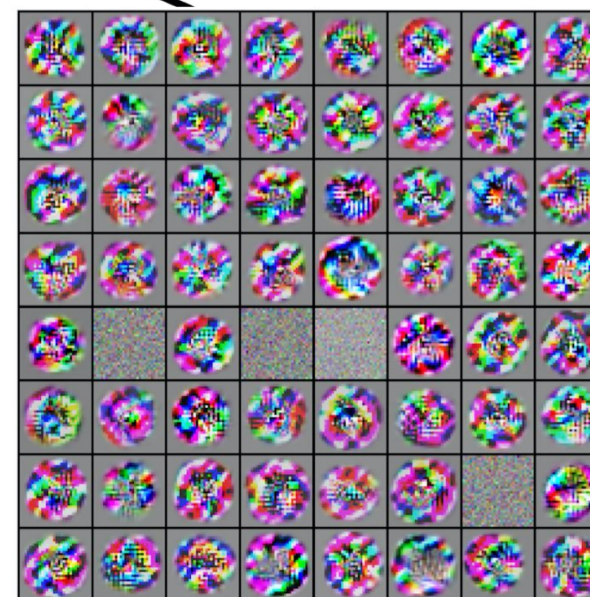
*Multi-layered abstraction of 3D features towards a linearly separable classification*



VGG-16 Conv1\_1



VGG-16 Conv3\_2



VGG-16 Conv5\_3

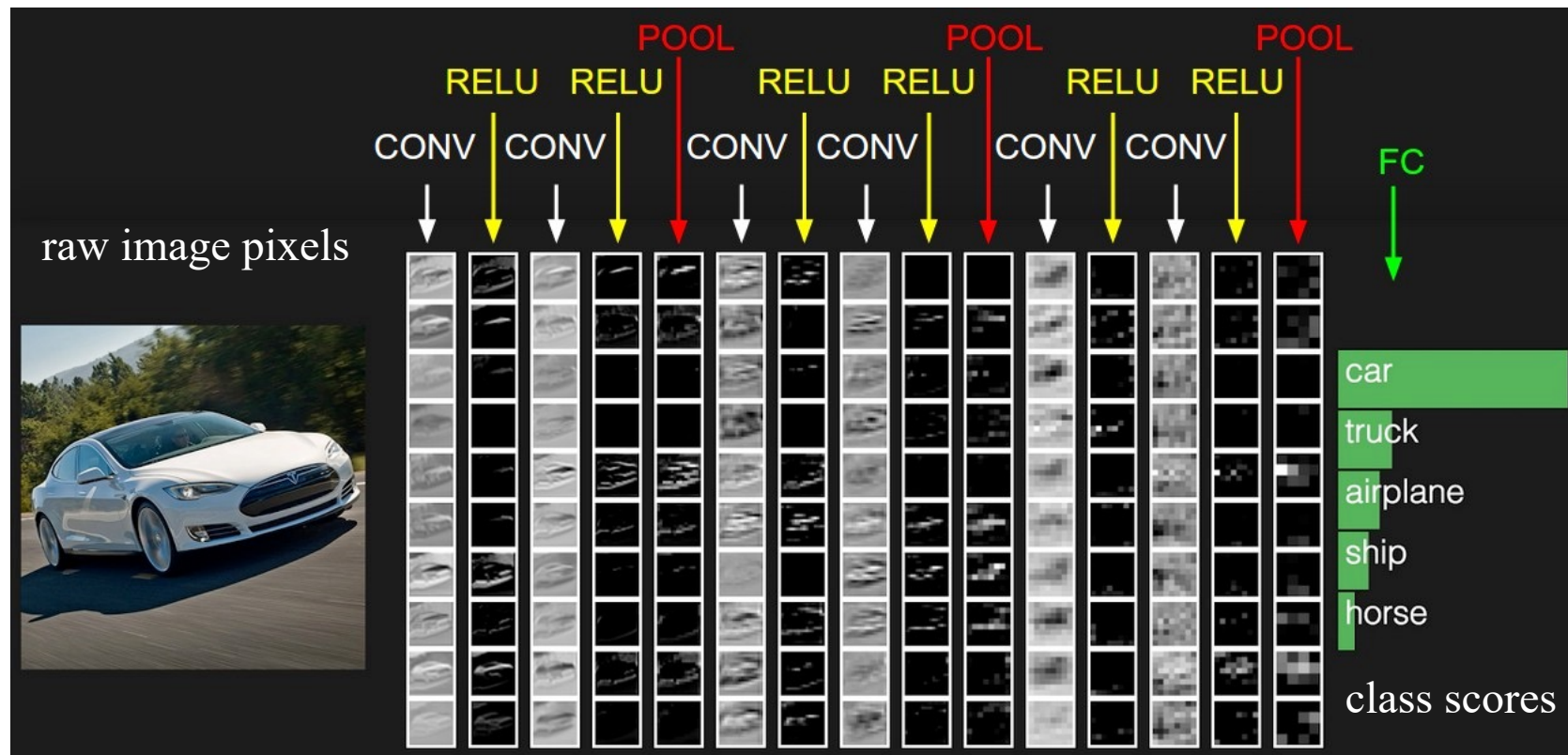
# A Simple ConvNet for CIFAR-10 Classification

*[INPUT - CONV - RELU - POOL - FC]*

**CONV** layer compute the output of neurons that are connected to local regions in the input, i.e.  $[32 \times 32 \times 12]$  with 12 filters.

**RELU** layer will apply an elementwise activation function, such as the  $\max(0, x)$  thresholding at zero. This leaves the size of the volume unchanged ( $[32 \times 32 \times 12]$ ).

**POOL** layer will perform a downsampling operation along the spatial dimensions (width, height), resulting in volume such as  $[16 \times 16 \times 12]$ .



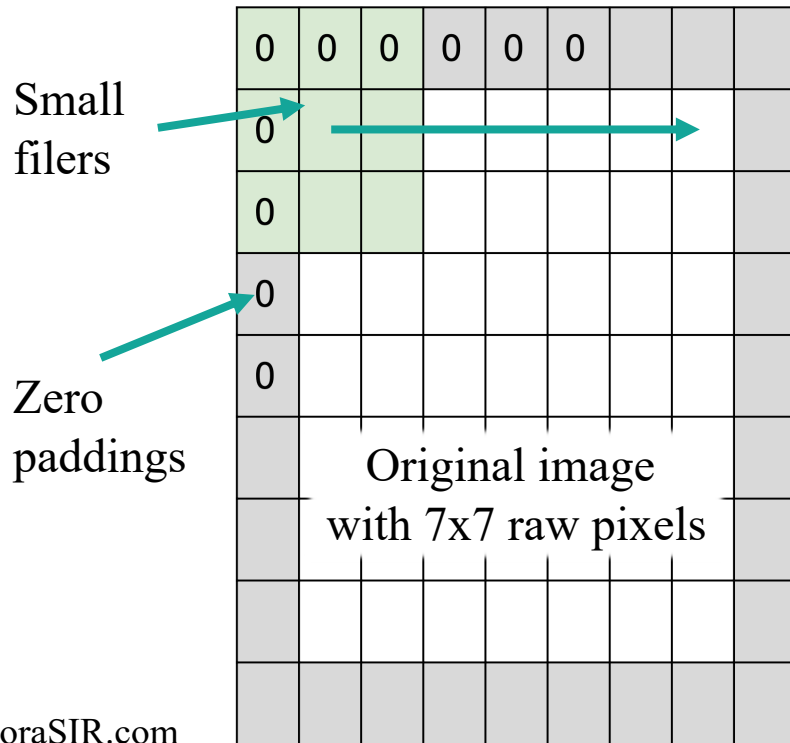
**INPUT** layer  $[32 \times 32 \times 3]$  will hold the raw pixel values of the image

**FC** (i.e. fully-connected) layer will compute the class scores, resulting in volume of size  $[1 \times 1 \times 10]$ , where each of the 10 numbers correspond to a class score

# Convolutional Layer

*Small filters that slide across the input volume*

- Small-size filers
  - e.g. 3x3 or at most 5x5, using a stride of  $S=1$ ,
  - Padding the input volume with zeros to avoid altering the spatial dimensions of the input.



INPUT features: 7x7

Filer size: 3x3

Stride: 1 (move step-by-step)

Padding: 1 pixel of 0 on all borders

OUTPUT features: 7x7

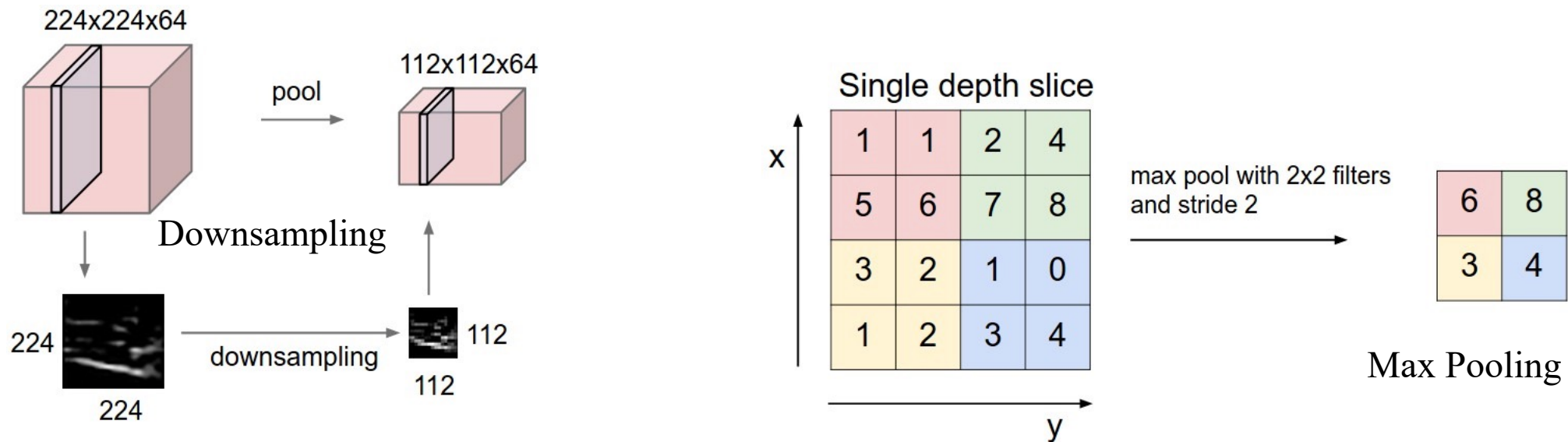
***What if without paddings on the border?***

- *The spatial dimensions of the input will be changed, causing information loss on the border*

# Pooling Layer

*Downsampling the spatial dimensions of the input volume*

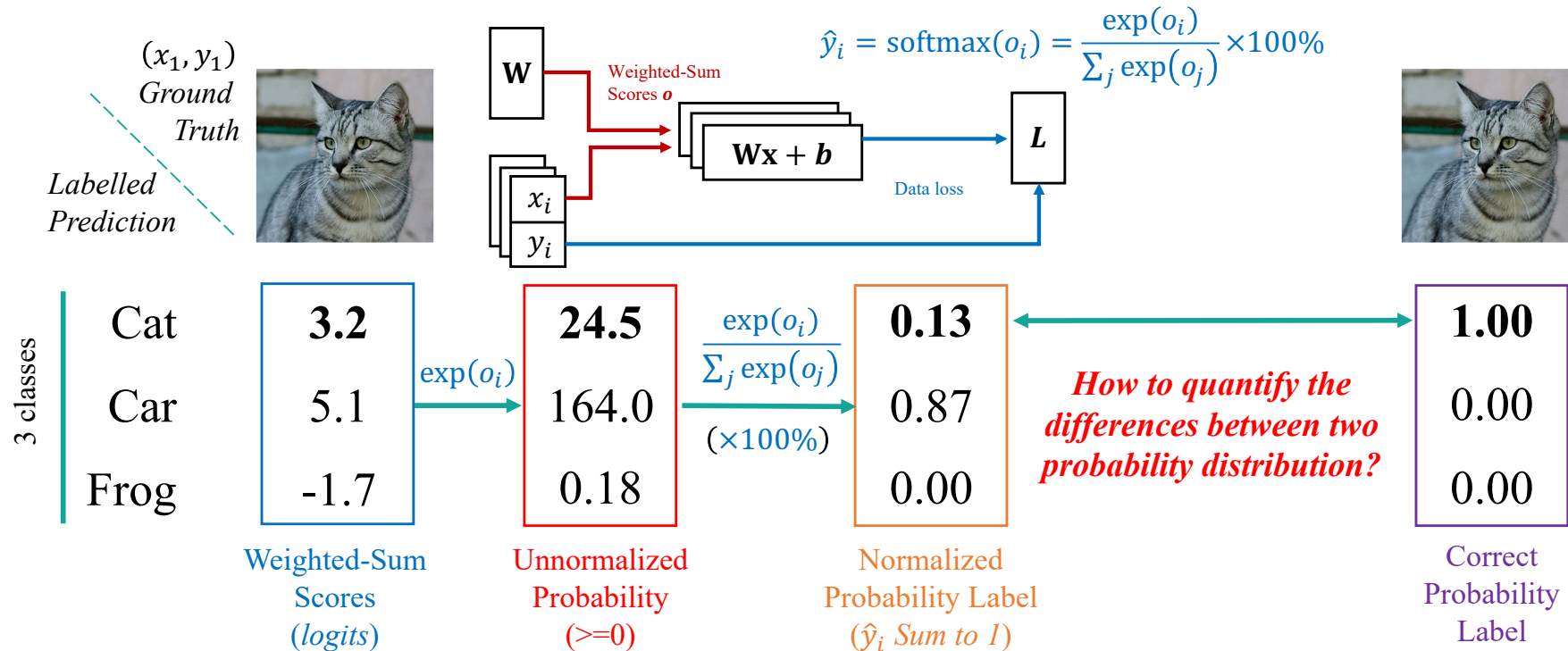
- A network-wise regularization
  - Progressively reduce the spatial size of the representation to reduce the amount of parameters and computation in the network, and hence to also control overfitting
  - Operates over each activation map independently
  - Usually, no need to zero padding (no convolutional operations)



# Fully-Connected Layer

*Full connections to all activations in the previous layer, as seen in regular Neural Networks*

- Contains neurons that connect to the entire input volume
- Softmax is a common choice





# ConvNet Architectures

*Common choice of hyperparameters of ConvNet designs*

- **INPUT**  $\rightarrow$   $[[\text{CONV} \rightarrow \text{RELU}] * N \rightarrow \text{POOL?}] * M \rightarrow [\text{FC} - \text{RELU}] * K \rightarrow \text{FC}$ 
  - the \* indicates repetition,
  - the POOL? indicates an optional pooling layer.
  - $N \geq 0$  (and usually  $N \leq 3$ ),  $M \geq 0$ ,  $K \geq 0$  (and usually  $K < 3$ )
- **INPUT** (that contains the image) should be divisible by 2 many times
  - 32 (e.g. CIFAR-10), 64, 96 (e.g. STL-10), or 224 (e.g. ImageNet), 384, and 512
- **CONV** should be using small filters using a stride of  $S=1$ 
  - 3x3 or at most 5x5 with zero padding of the input volume
- **POOL** downsamples the spatial dimensions of the input
  - Common setting is to use max-pooling with 2x2 receptive fields with a stride of 2

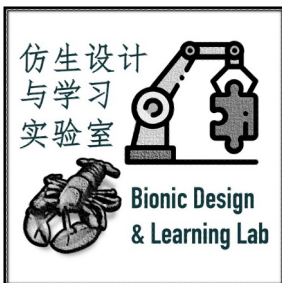
Bionic Design & Learning Lab  
@ SIR Group 仿生设计与学习实验室



Room 606  
7 Innovation Park  
南科创园7栋606室

Thank you~

[songcy@sustech.edu.cn](mailto:songcy@sustech.edu.cn)



AncoraSIR.com

