

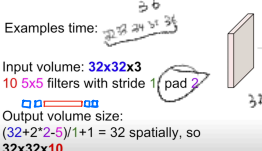
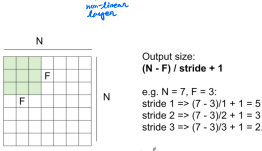
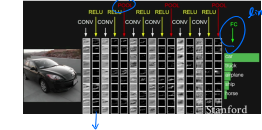
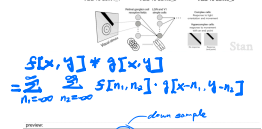
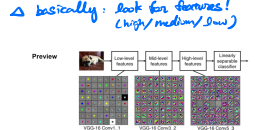
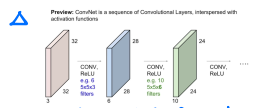
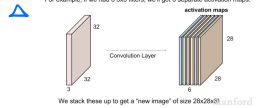
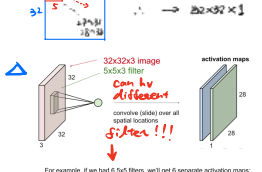
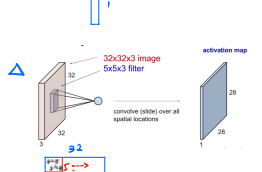
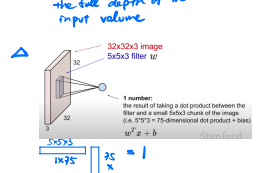
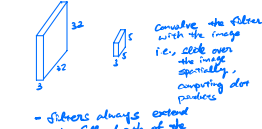
CNN (directional use)

recall fully-connected layer (linear)
 $32 \times 32 \times 3$ image

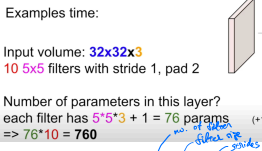
3072×1



convolutional layer
 $32 \times 32 \times 3$ image (preserve source)



Q: no. of parameters? (note weights)



you can do 1x1 conv

Example: CONV layer in TORCH

```

SpatialConvolution
Applies a 2D convolution over an input image composed of several input planes. The output tensor is returned in a 3D format (convolutional image + depth).

Parameters are the following:
- kernelSize: The number of spatial kernel planes in the image plane (height).
- inChannels: The number of input planes in the convolution layer (number).
- outChannels: The number of output planes in the convolution layer (number).
- stride: The kernel width of the convolution.
- padding: The kernel height of the convolution.
- dilation: The step of the convolution in the width dimension. Default is 1.
- groups: The additional zero added per width in the input planes. Default is 1, a group number in [1, outChannels].
- bias: The additional zero added per width in the input planes. Default is 1, a group number in [1, outChannels].

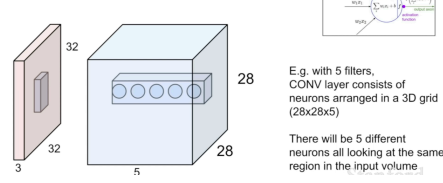
Note that depending on the size of your kernel, several of the last values or ones of the last image might be lost, it is up to the user to pad properly in images.

When two images in a 3D tensor (dimensions = height x width x channel) are passed as an input to this module, the output is a 3D tensor (dimensions = height x width x channel) whose size is:
height = floor((height - kernelSize) / stride) + 1
width = floor((width - kernelSize) / stride) + 1
channel = inChannels * outChannels / groups

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  $P$ 
  
```

5x5 filter \rightarrow 5x1 receptive field for each neuron

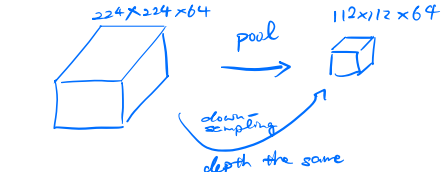
The brain/neuron view of CONV Layer



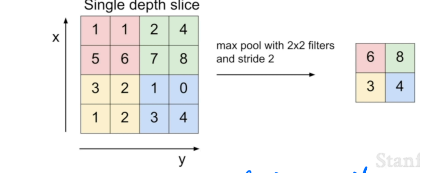
fully connected network (takes the entire input domain)

Pooling layer

makes the representations smaller & manageable
operates over each activation map independently



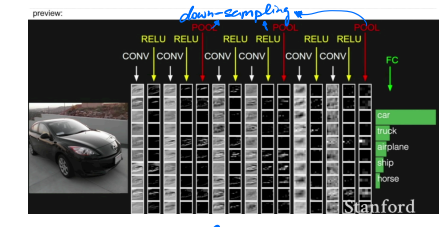
max pooling (instead of dot product, we do maximization)



(usually not over-lapping striding)

averaging pooling

Accepts a volume of size $W_1 \times H_1 \times D_1$
Requires three hyperparameters:
- their spatial extent F (filter size)
- 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
CNN output $N \times H \times D$ \rightarrow map \rightarrow result space

Input volume: $32 \times 32 \times 3$
10 5x5 filters with stride 1, pad 2

Output volume size:
 $(32 + 2 \cdot 2 - 5) / 1 + 1 = 32$ spatially, so
 $32 \times 32 \times 10$

Number of parameters in this layer?
each filter has $5^2 \cdot 3 + 1 = 76$ params (+1 for bias)
 $\Rightarrow 76 \cdot 10 = 760$

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:
 $W_2 = (W_1 - F) / S + 1$
 $H_2 = (H_1 - F) / S + 1$
 $D_2 = D_1$
- 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 $W_2 \times H_2 \times D_2$) is the result of performing a valid convolution of the d th slice over the input volume with stride S , since conv offset by d th bias.