

CS5242 : Neural Networks and Deep Learning

Lecture 8: Convolutional Neural Networks Implementation

Semester 1 2021/22

Xavier Bresson

<https://twitter.com/xbresson>

Department of Computer Science
National University of Singapore (NUS)



Outline

- Padding
- Stride
- Multi-dimensional input
- Convolution as linear operation
- Depth
- Convolutional layer with PyTorch
- Pooling layer with PyTorch
- LeNet5 architecture
- VGG architecture
- Case studies

Outline

- Padding
- Stride
- Multi-dimensional input
- Convolution as linear operation
- Depth
- Convolutional layer with PyTorch
- Pooling layer with PyTorch
- LeNet5 architecture
- VGG architecture
- Case studies

Padding

- Convolution operation **shrinks** the size of the image :

$$\begin{bmatrix} 3 & 2 & 3 & 3 & 3 & 1 \\ 0 & 1 & 3 & 3 & 1 & 2 \\ 2 & 2 & 0 & 0 & 2 & 0 \\ 3 & 0 & 1 & 3 & 3 & 0 \\ 2 & 3 & 2 & 0 & 0 & 1 \\ 2 & 2 & 2 & 3 & 1 & 0 \end{bmatrix}$$

Input image

6 x 6

*

$$\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 2 \end{bmatrix}$$

Convolutional filter

3 x 3

=

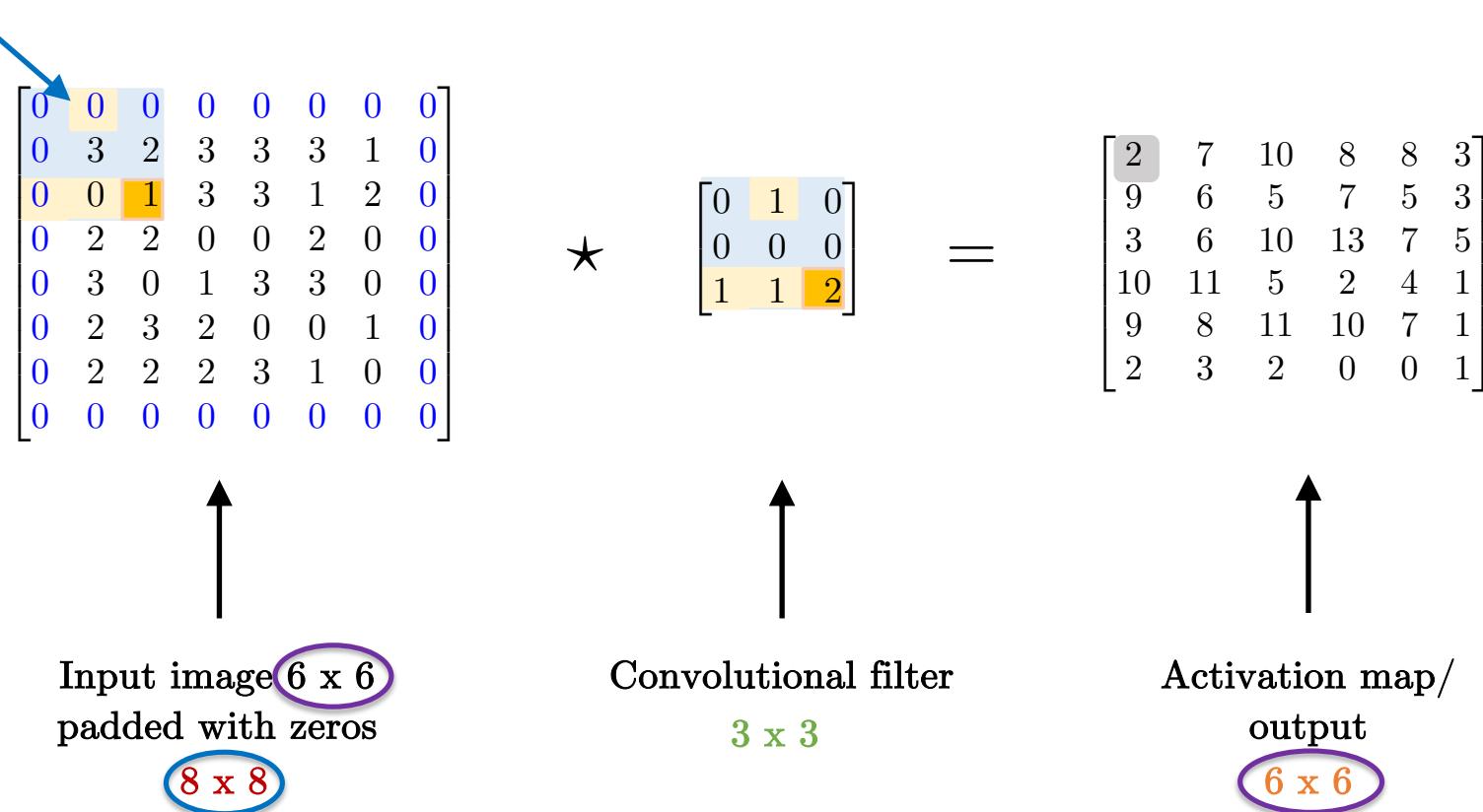
$$\begin{bmatrix} 6 & 5 & 7 & 5 \\ 6 & 10 & 13 & 7 \\ 11 & 5 & 2 & 4 \\ 8 & 11 & 10 & 7 \end{bmatrix}$$

Output image/
Activation map

4 x 4

Padding

- Shrinking the image size is an **issue** if we want to **stack many layers** (size will reduce to 1x1 quickly).
 - **Padding input image with zeros** allows to have the **same size** for the input and the output.



Padding

- Padding size depends on the filter size :
 - When using **5x5 filters** (radius=2), image must be padded with **2 layers** of zeros.

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \star \begin{bmatrix} 1 & 2 & 1 & 1 & 2 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 2 & 1 \\ 1 & 2 & 1 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 7 & 10 & 8 & 8 & 3 \\ 9 & 6 & 5 & 7 & 5 & 3 \\ 3 & 6 & 10 & 13 & 7 & 5 \\ 10 & 11 & 5 & 2 & 4 & 1 \\ 9 & 8 & 11 & 10 & 7 & 1 \\ 2 & 3 & 2 & 0 & 0 & 1 \end{bmatrix}$$

- When using **7x7 filters** (radius=3), image must be padded with **3 layers** of zeros.
- Etc.

Outline

- Padding
- **Stride**
- Multi-dimensional input
- Convolution as linear operation
- Depth
- Convolutional layer with PyTorch
- Pooling layer with PyTorch
- LeNet5 architecture
- VGG architecture
- Case studies

Strided convolution

- Strided convolution can be useful for large images (e.g. satellite images 10,000 x 10,000) :

Consider a stride = 2

$$\begin{bmatrix} 3 & 2 & 3 & 3 & 3 & 1 & 2 \\ 0 & 1 & 3 & 3 & 1 & 2 & 1 \\ 2 & 2 & 0 & 0 & 2 & 0 & 0 \\ 3 & 0 & 1 & 3 & 3 & 0 & 3 \\ 2 & 3 & 2 & 0 & 0 & 1 & 1 \\ 2 & 2 & 2 & 3 & 1 & 0 & 1 \\ 0 & 2 & 1 & 1 & 3 & 3 & 1 \end{bmatrix} \star \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 6 & 5 & 7 \\ 6 & 10 & 13 \\ 11 & 5 & 2 \end{bmatrix}$$

Strided convolution

Consider a stride = 2

$$\begin{bmatrix} 3 & 2 & 3 & 3 & 3 & 1 & 2 \\ 0 & 1 & 3 & 3 & 1 & 2 & 1 \\ 2 & 2 & 0 & 0 & 2 & 0 & 0 \\ 3 & 0 & 1 & 3 & 3 & 0 & 3 \\ 2 & 3 & 2 & 0 & 0 & 1 & 1 \\ 2 & 2 & 2 & 3 & 1 & 0 & 1 \\ 0 & 2 & 1 & 1 & 3 & 3 & 1 \end{bmatrix} \star \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 6 & 5 & 7 \\ 6 & 10 & 13 \\ 11 & 5 & 2 \end{bmatrix}$$

$$\begin{bmatrix} 3 & 2 & 3 & 3 & 3 & 1 & 2 \\ 0 & 1 & 3 & 3 & 1 & 2 & 1 \\ 2 & 2 & 0 & 0 & 2 & 0 & 0 \\ 3 & 0 & 1 & 3 & 3 & 0 & 3 \\ 2 & 3 & 2 & 0 & 0 & 1 & 1 \\ 2 & 2 & 2 & 3 & 1 & 0 & 1 \\ 0 & 2 & 1 & 1 & 3 & 3 & 1 \end{bmatrix} \star \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 6 & 5 & 7 \\ 6 & 10 & 13 \\ 11 & 5 & 2 \end{bmatrix}$$

Output size

- The **output size** of a convolution depends on :
 - Input size : n
 - Convolutional filter size : f
 - Padding : p
 - Stride : s

$$n \times n \quad \Rightarrow \quad \left\lfloor \frac{n + 2p - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n + 2p - f}{s} + 1 \right\rfloor$$

Input size Convolution operation
 Parameters: f, p, s Output size

Integer floor
value

Outline

- Padding
- Stride
- **Multi-dimensional input**
- Convolution as linear operation
- Depth
- Convolutional layer with PyTorch
- Pooling layer with PyTorch
- LeNet5 architecture
- VGG architecture
- Case studies

Multi-dimensional convolution

- N-D convolution using N-D filters with N-D inputs returns 1-D activation map.
- Example with 3-D input :

Patch of $3 \times 3 \times 3 =$
27 numbers

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

3-D input

$3 \times 5 \times 5$ tensor

Filter of $3 \times 3 \times 3 =$
27 numbers

0	1	0
0	3	0
1	1	2

*

3-D filter

$3 \times 3 \times 3$ tensor

$$(0)(0)+(1)(1)+(3)(0) + \dots + (1)(2)+(2)(-5)+(3)(-1)+\dots + (-2)(3)+(1)(2)+(2)(8)+\dots = 5$$

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

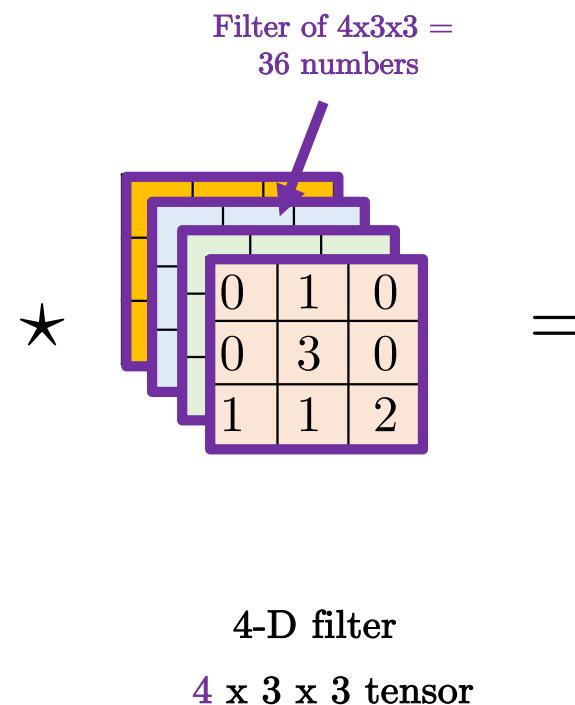
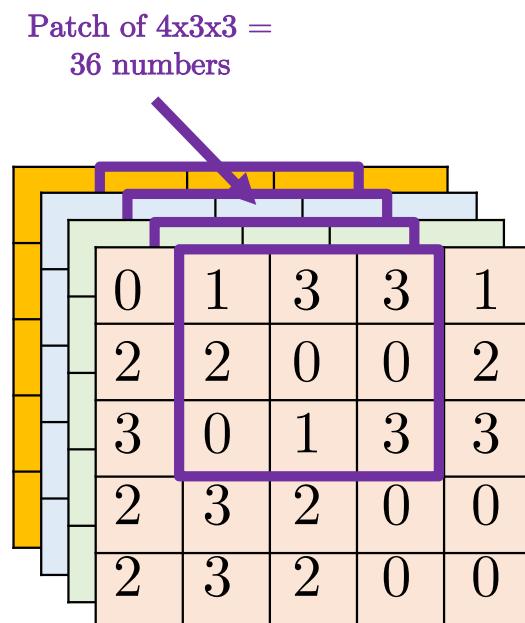
Padding=1

1-D activation map

5×5 tensor

Multi-dimensional convolution

- N-D convolution using N-D filters with N-D inputs returns 1-D activation map.
 - Example with 4-D input :



=

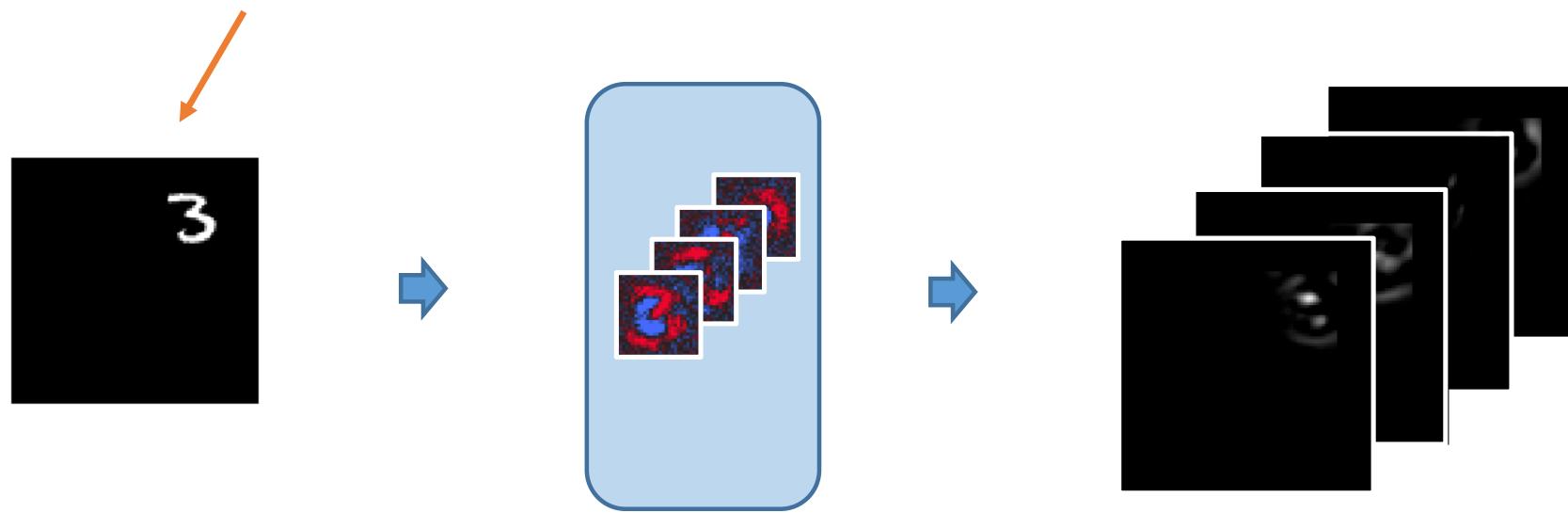
7	9	5	9	5
9	5	0	5	6
1	0	5	4	4
4	2	1	3	3
1	7	5	4	6

Padding=1

1-D activation map
 5×5 tensor

One-dimensional input

- 1-D input : Gray-scale images



The input has 1 channel → Each filter has 1 channel

128 x 128 tensor

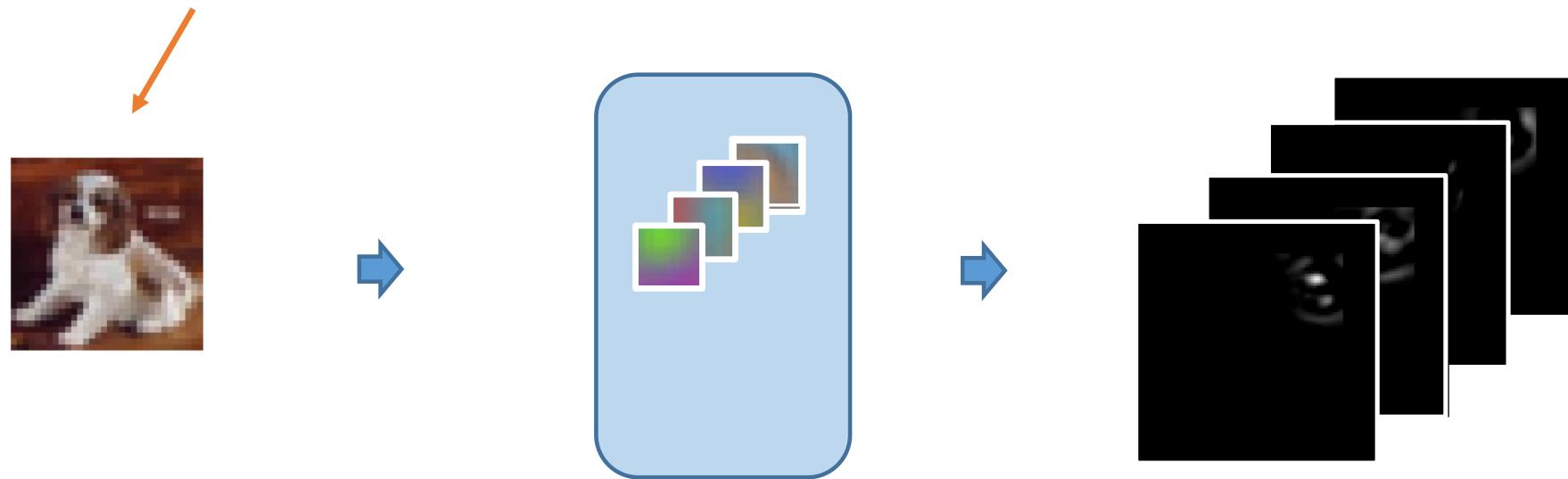
Four filters
4 x 28 x 28 tensor



Four activation maps
4 x 128 x 128 tensor

Three-dimensional input

- 3-D input : RGB images



The input has $N=3$ channels → Each filter has 3 channels

$3 \times 28 \times 28$ tensor

Four filters

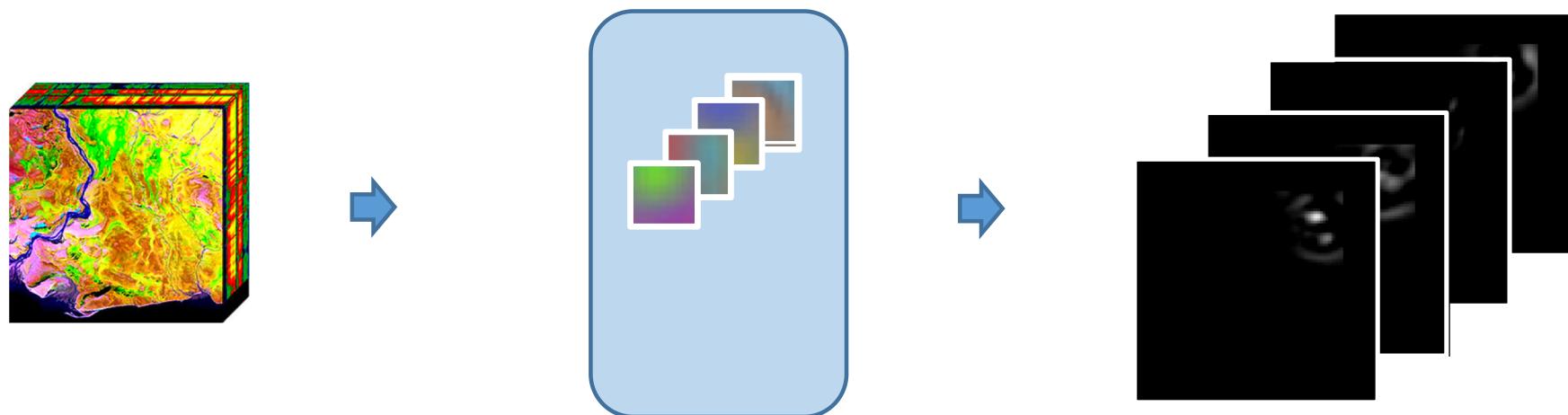
$4 \times 3 \times 3 \times 3$ tensor

Four activation maps

$4 \times 28 \times 28$ tensor

F-D activations map with N-D input

- Generate F activation maps with N-D inputs :
 - N channels, F filters/activation maps, $N \times f \times f$ filter size



The input has N channels



Each filter has N channels

$N \times n \times n$ tensor

F filters



F activation maps

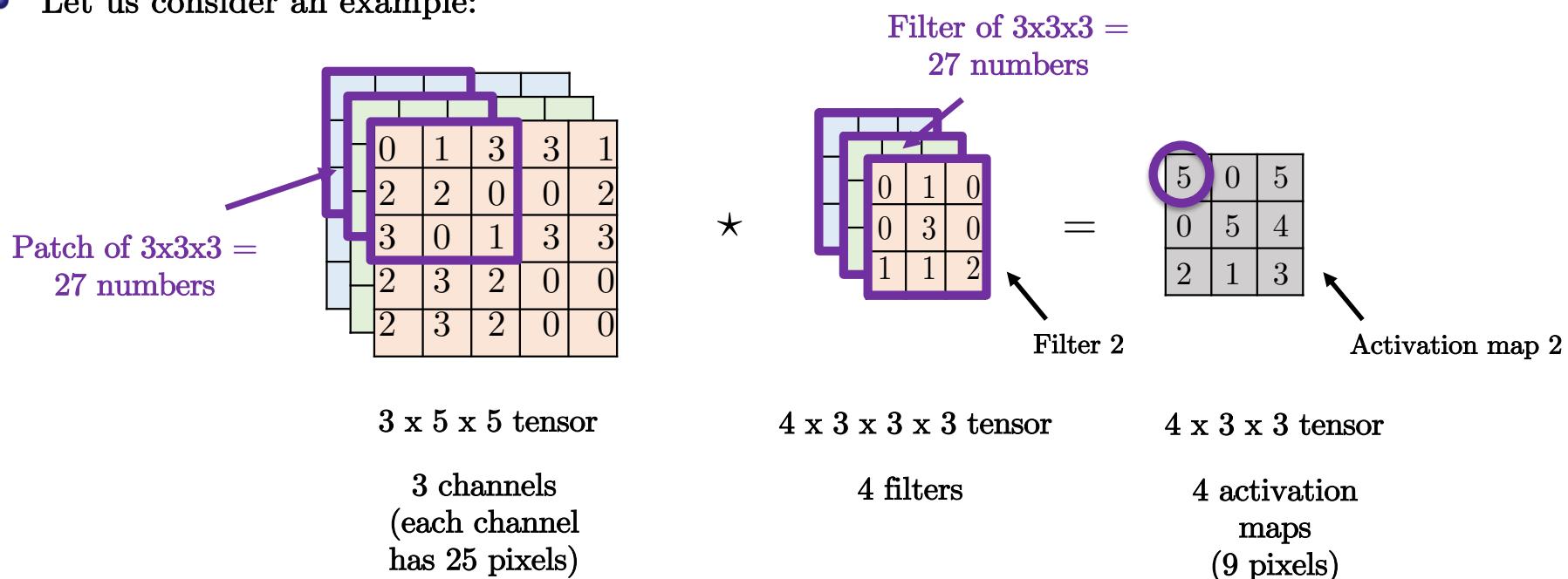
$F \times N \times f \times f$ tensor

Outline

- Padding
- Stride
- Multi-dimensional input
- **Convolution as linear operation**
- Depth
- Convolutional layer with PyTorch
- Pooling layer with PyTorch
- LeNet5 architecture
- VGG architecture
- Case studies

Convolution as linear operation

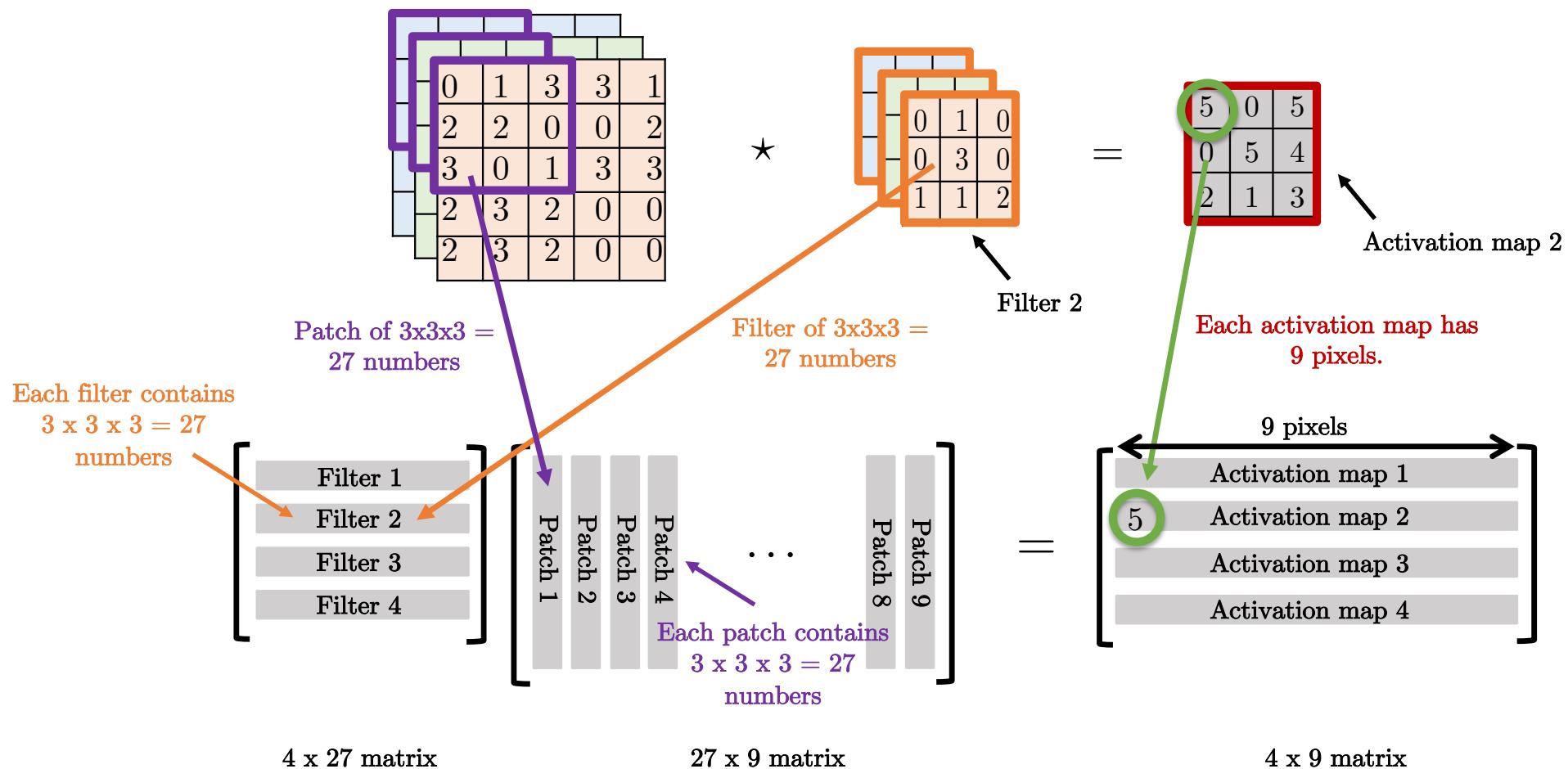
- Convolutional layers are not very different from regular linear layers.
 - They also define a linear operator with a special structure.
- Let us consider an example:



Key observation : Convolution operations are applied on patches !

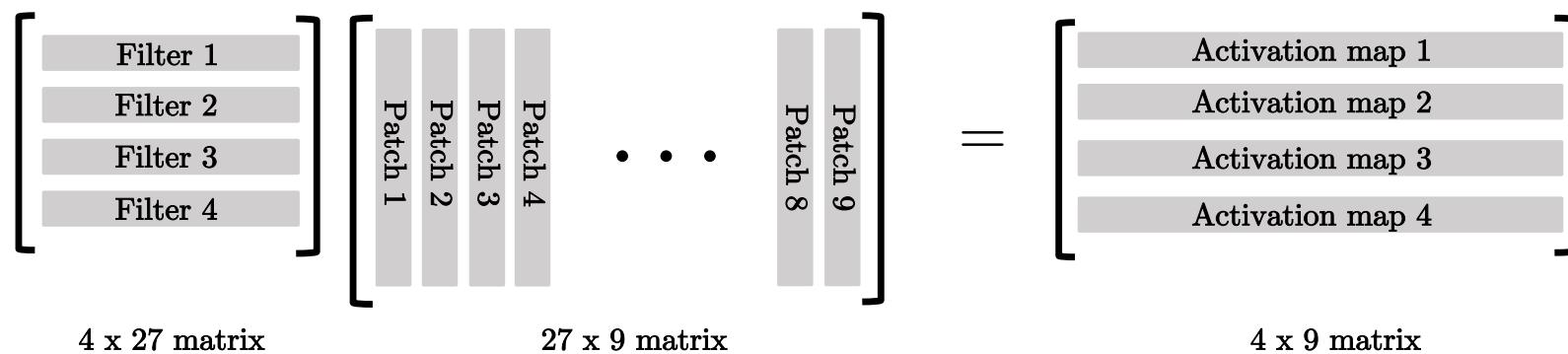
Convolution as linear operation

- Convolutions can be expressed as matrix-matrix multiplications :



Convolution as linear operation

- **Convolutions** can be expressed as **matrix-matrix multiplications** :



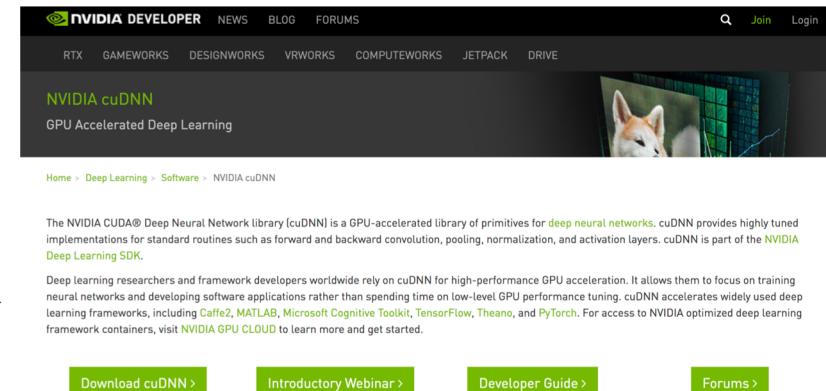
Convolution = matrix - matrix multiplication

Filters vectorized
(one filter per row)

Input pixels
rearranged as vectorized patches
(one patch per column)

Fast convolution

- Convolutions are simply matrix-matrix multiplications.
 - As matrix-matrix multiplications are standard linear algebra operations that have been highly parallelized.
 - CPU: MKL and BLAS libraries
 - GPU: CUDA library
- Fast convolutions on GPUs with cuDNN library.



<https://developer.nvidia.com/cudnn>

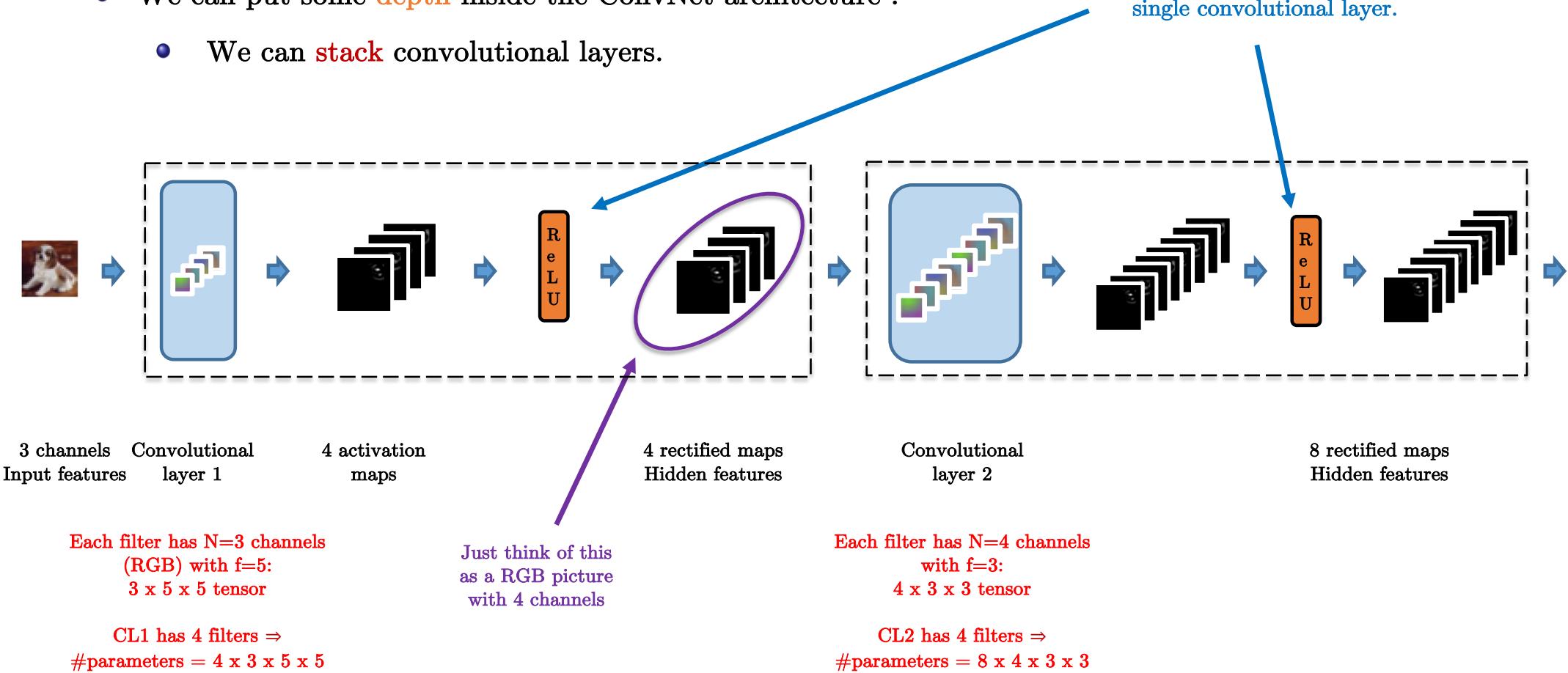
Outline

- Padding
- Stride
- Multi-dimensional input
- Convolution as linear operation
- **Depth**
- Convolutional layer with PyTorch
- Pooling layer with PyTorch
- LeNet5 architecture
- VGG architecture
- Case studies

Depth

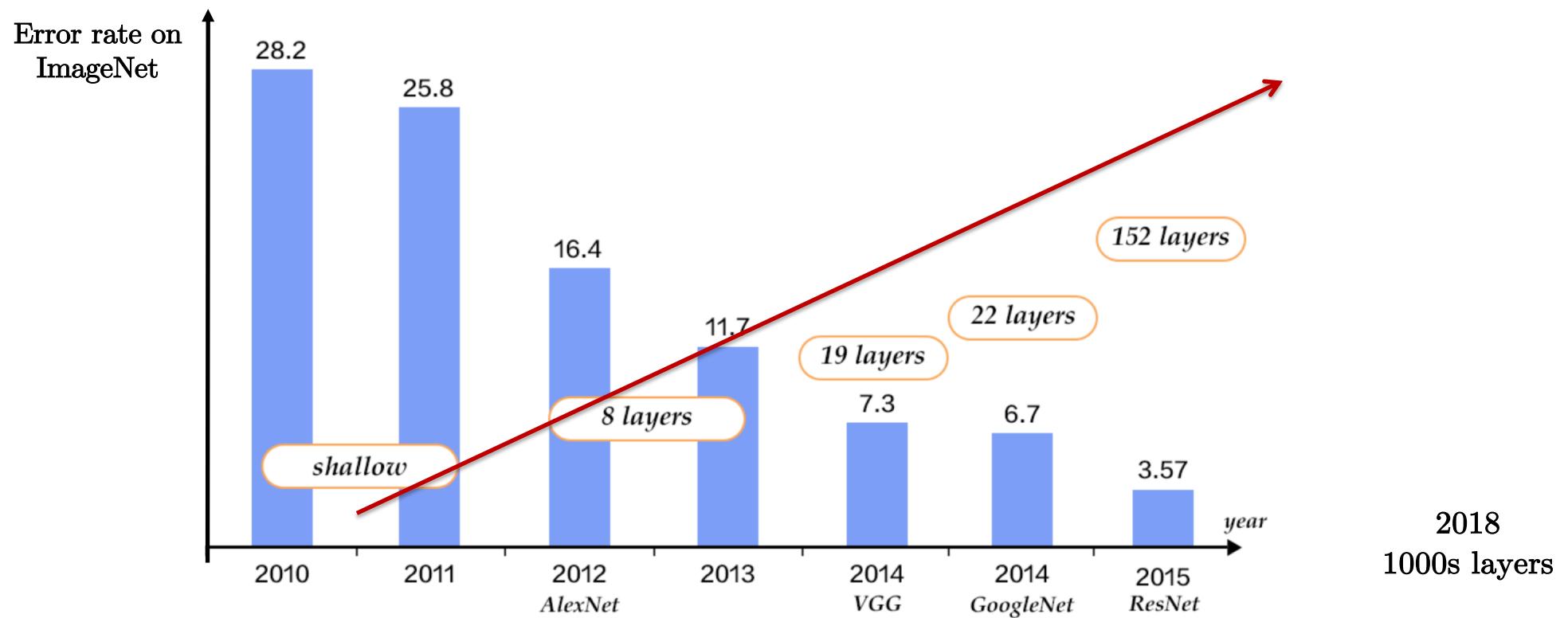
- We can put some **depth** inside the ConvNet architecture :
 - We can **stack** convolutional layers.

Non-linearity is needed between each layer to prevent the network to collapse into a single convolutional layer.



The deeper the better

- Revolution of depth :

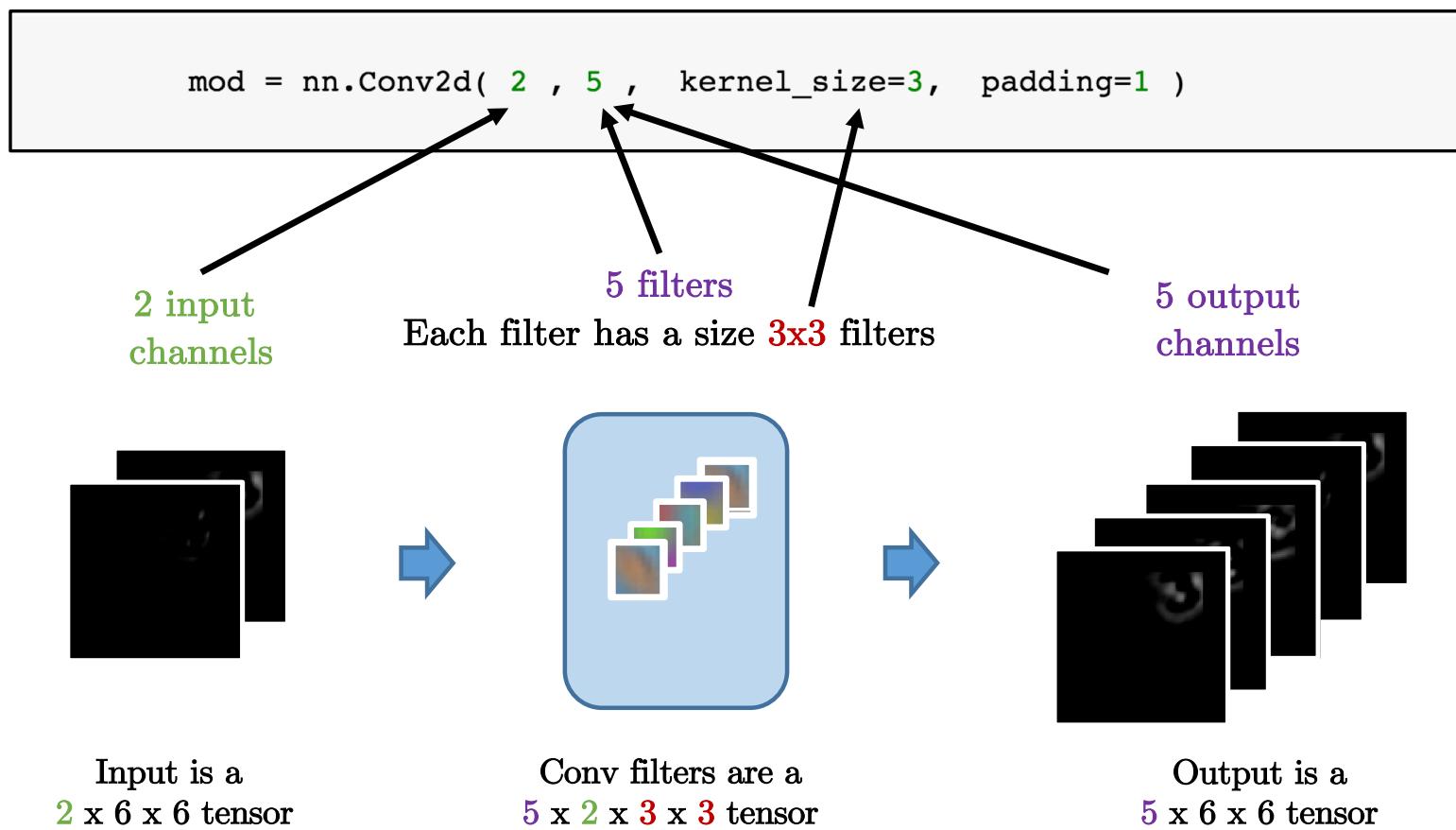


Outline

- Padding
- Stride
- Multi-dimensional input
- Convolution as linear operation
- Depth
- **Convolutional layer with PyTorch**
- Pooling layer with PyTorch
- LeNet5 architecture
- VGG architecture
- Case studies

Convolutional layer in PyTorch

- PyTorch function:



Lab 01

- Convolutional layer

The screenshot shows a Jupyter Notebook interface with the title "Lab 01 : Convolutional layer - demo". The notebook has three code cells:

- In [2]:**

```
import torch
import torch.nn as nn
```
- In [3]:**

```
mod = nn.Conv2d( 2 , 5 , kernel_size=3, padding=1 )
```
- In [4]:**

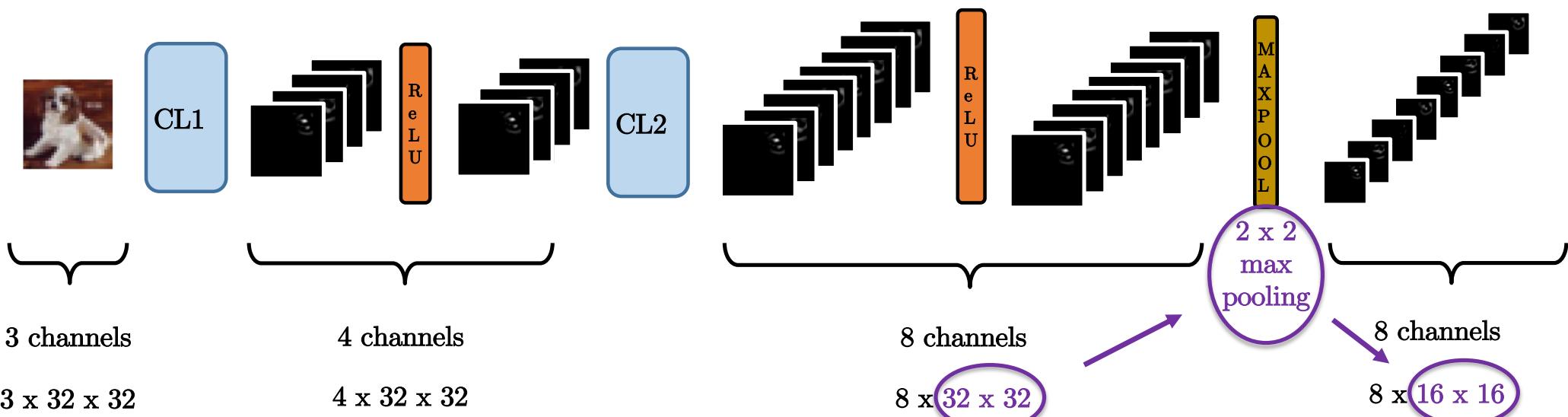
```
bs=1
x=torch.rand(bs,2,6,6)
print(x)
print(x.size())
tensor([[[[0.9647, 0.6641, 0.1846, 0.0741, 0.2426, 0.0286],
          [0.2926, 0.2690, 0.8945, 0.7917, 0.6449, 0.4339],
          [0.4199, 0.9387, 0.3030, 0.7605, 0.5784, 0.1366],
          [0.0920, 0.2539, 0.1089, 0.0882, 0.4449, 0.9903],
          [0.0181, 0.6897, 0.2516, 0.8399, 0.7515, 0.3681],
          [0.6694, 0.0813, 0.5470, 0.1661, 0.7311, 0.2856]]],
```

The notebook header shows "jupyter conv_layer_demo Last Checkpoint: a few seconds ago (autosaved)" and the Python 3 kernel.

Outline

- Padding
- Stride
- Multi-dimensional input
- Convolution as linear operation
- Depth
- Convolutional layer with PyTorch
- **Pooling layer with PyTorch**
- LeNet5 architecture
- VGG architecture
- Case studies

Max pooling



- MP (max pooling) layers are useful for :
 - Trading image resolution with features number (features capture data structures).
 - Controlling the computational complexity (image size /2 and #features x2).

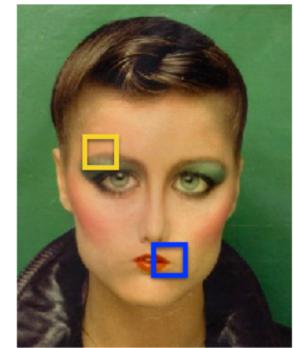
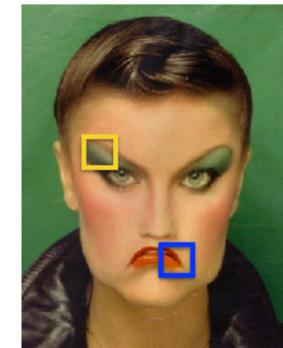
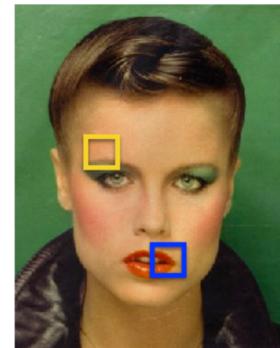
Max pooling

- MP (max pooling) layers are also useful for :

- Creating global invariance w.r.t. translation (same semantic object independent of location).



- Providing local invariance w.r.t. object deformations (robust to slight object changes).



Pooling layer in PyTorch

- PyTorch function:
 - Max Pooling with window size = 2 x 2

```
mod = nn.MaxPool2d(2,2)
```

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

6 x 6 activation map

Down-sampling
+ max pooling



7	1	5
9	1	6
9	0	1

3 x 3 activation map

Lab 02

- Pooling layer :

jupyter pool_layer_demo Last Checkpoint: a few seconds ago (autosaved) Logout

In [1]:

```
import torch
import torch.nn as nn
```

Make a pooling module

- inputs: activation maps of size n x n
- output: activation maps of size n/p x n/p
- p: pooling size

In [2]:

```
mod = nn.MaxPool2d(2,2)
```

Make an input 2 x 6 x 6 (two channels, each one has 6 x 6 pixels)

In [8]:

```
bs=1
x=torch.rand(bs,2,6,6)

print(x)
print(x.size())

tensor([[[[0.6796, 0.0749, 0.5749, 0.5729, 0.3750, 0.9115],
          [0.0087, 0.2287, 0.1944, 0.5500, 0.9770, 0.0556],
          [0.2348, 0.4120, 0.8975, 0.8404, 0.1589, 0.7432],
          [0.7057, 0.1542, 0.8852, 0.1988, 0.2206, 0.2997],
          [0.6186, 0.6267, 0.1895, 0.5063, 0.0490, 0.9198],
          [0.1485, 0.1590, 0.0135, 0.9169, 0.1770, 0.9111]],
```

Outline

- Padding
- Stride
- Multi-dimensional input
- Convolution as linear operation
- Depth
- Convolutional layer with PyTorch
- Pooling layer with PyTorch
- **LeNet5 architecture**
- VGG architecture
- Case studies

LeNet5 architecture

- LeCun, Bottou, Bengio, Haffner, 1998:

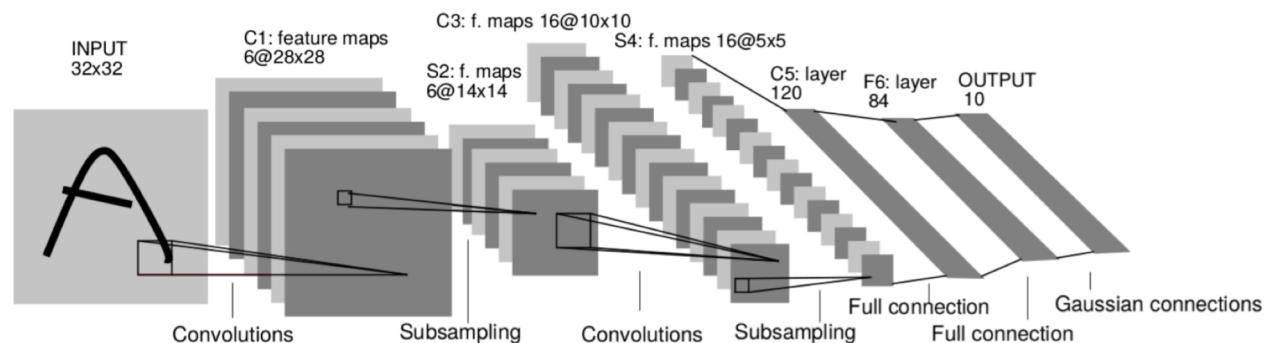
- 1st ConvNet architecture

- The dataset:

- MNIST

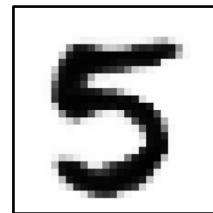
- The network:

- Convolutional/MP layers
 - Linear layers



MNIST dataset

- Number of classes: 10
- Training set: 50,000 labelled images
- Test set: 10,000 labelled images

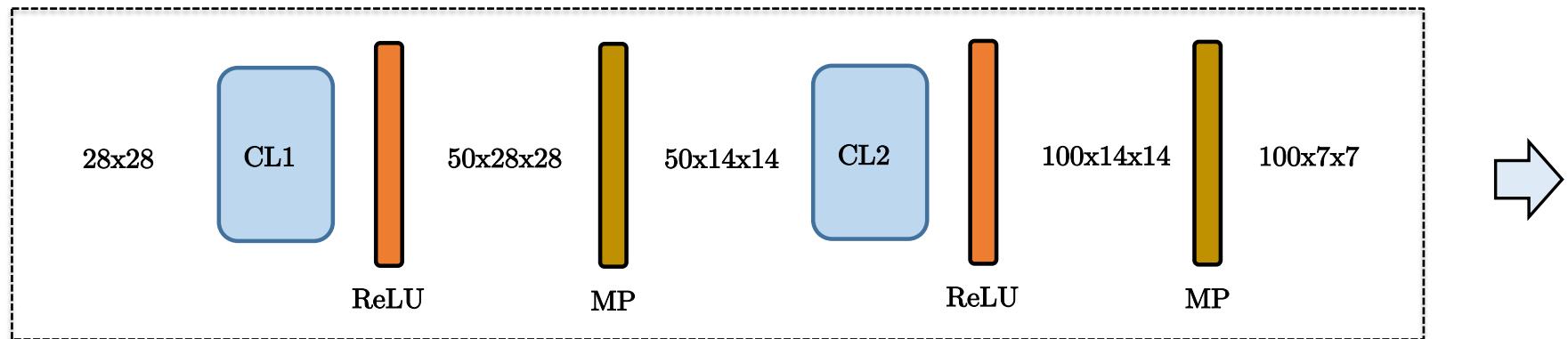


28 x 28 grayscale image

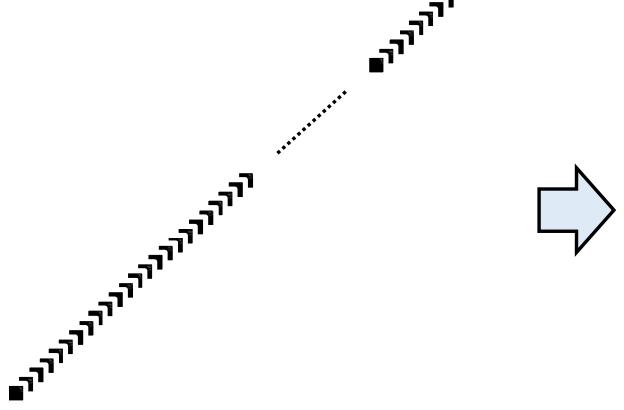
0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9

Two convolutional layers

- Convolutional/MP layers:

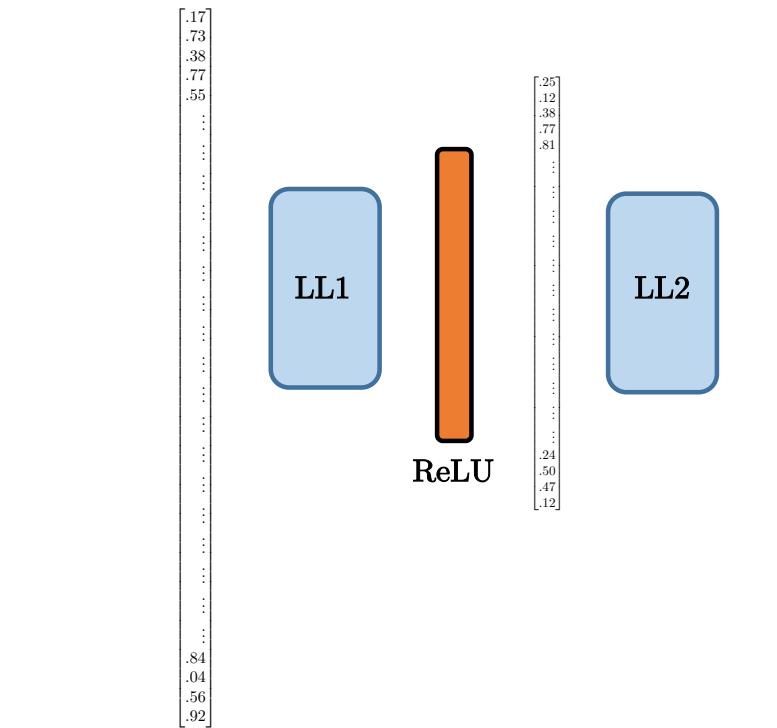


Linear layers



CL features

100 (small)
activation maps
(7×7 pixel each)



Vector with
4,900 entries

$$(100 \times 7 \times 7 = 4,900)$$

Vectors with
100 entries

10 Scores

Loss

Cross
Entropy
Criterion

0.03

Number of architecture parameters

- First CL:
 - $28 \times 28 \Rightarrow \text{CL1+ReLU} \Rightarrow 50 \times 28 \times 28$
 - $50 \times 28 \times 28 \Rightarrow \text{MP} \Rightarrow 50 \times 14 \times 14$
 - $\# \text{Parameters} = 50 \times 3 \times 3 + 50 = 500$
- Second CL:
 - $50 \times 14 \times 14 \Rightarrow \text{CL1+ReLU} \Rightarrow 100 \times 14 \times 14$
 - $100 \times 14 \times 14 \Rightarrow \text{MP} \Rightarrow 100 \times 7 \times 7$
 - $\# \text{Parameters} = 100 \times 50 \times 3 \times 3 + 100 = 45,100$
- First LL:
 - $100 \times 7 \times 7 \Rightarrow \text{Reshape} \Rightarrow 4,900$
 - $4,900 \Rightarrow \text{LL1} \Rightarrow 100$
 - $\# \text{Parameters} = 4900 \times 100 + 100 = 490,100$
- Second LL:
 - $100 \Rightarrow \text{LL2} \Rightarrow 10$
 - $\# \text{Parameters} = 100 \times 10 + 10 = 1,010$

Total for 4 Layers $\Rightarrow 536,710 / 0.54$ Millions of parameters

Lab 03

- LeNet5 on MNIST :

- Learning rate schedule:

- Start with lr=0.25
 - Divide by 1.5 every 5 epochs

The screenshot shows a Jupyter Notebook interface with the title "Lab 03 : LeNet5 architecture - exercise". The notebook contains several code cells:

- Cell 1:** Imports torch, nn, functional, optim, random, utils, and time.
- Cell 2:** A section titled "With or without GPU?" with a note: "It is recommended to run this code on GPU:
 - Time for 1 epoch on CPU : 96 sec (1.62 min)
 - Time for 1 epoch on GPU : 2 sec w/ GeForce GTX 1080 Ti"
- Cell 3:** Sets device to torch.device("cuda") or torch.device("cpu") and prints the device.
- Cell 4:** Downloads the MNIST dataset using check_mnist_dataset_exists and loads train and test data.

Results

- Results with a 4-layer ConvNet:
 - Number of parameters: 0.54 million
 - Error on training set: 0%
 - Error on test set: 0.78%
- Timing for one epoch:
 - CPU : 1.6 min
 - Tesla K80 GPU (Google Cloud) : 6 sec
 - GTX 1080 GPU (\$600) : 2 sec

Outline

- Padding
- Stride
- Multi-dimensional input
- Convolution as linear operation
- Depth
- Convolutional layer with PyTorch
- Pooling layer with PyTorch
- LeNet5 architecture
- **VGG architecture**
- Case studies

VGG architecture

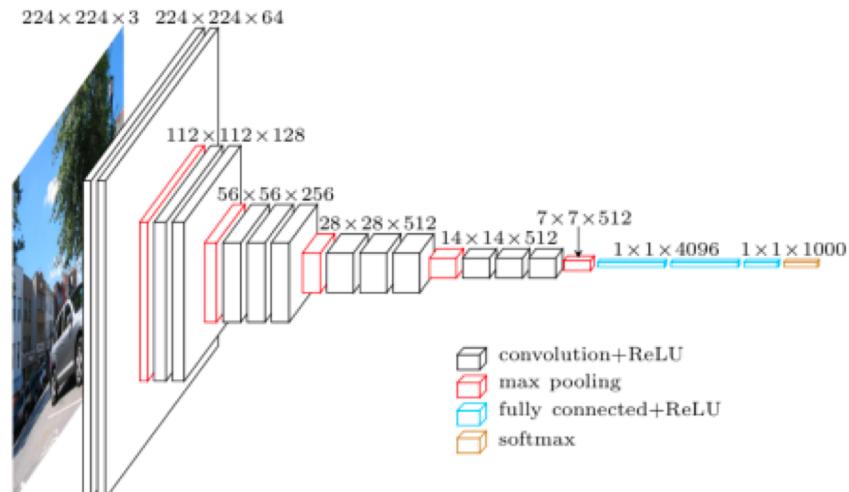
- Simonyan and Zisserman, 2014 :
 - 2nd place at the ImageNet challenge in 2014
 - All convolutional features are 3x3

- The dataset:

- CIFAR

- The network:

- Convolutional/MP layers
 - Linear layers

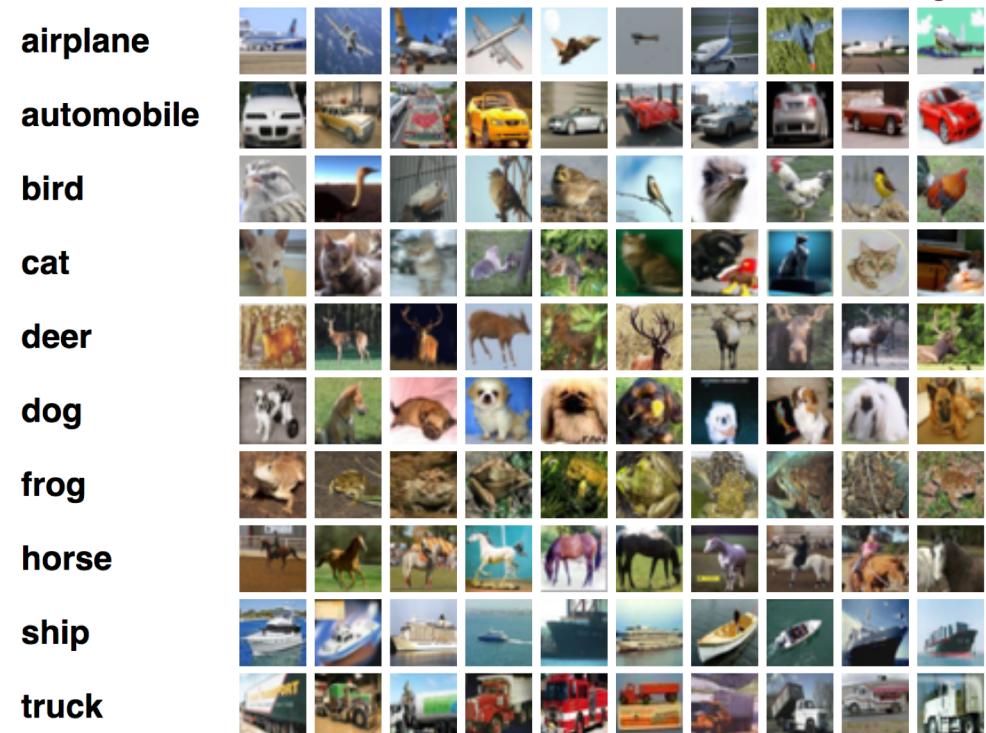


CIFAR dataset

- Number of classes: 10
- Training set: 50,000 labelled images
- Test set: 10,000 labelled images

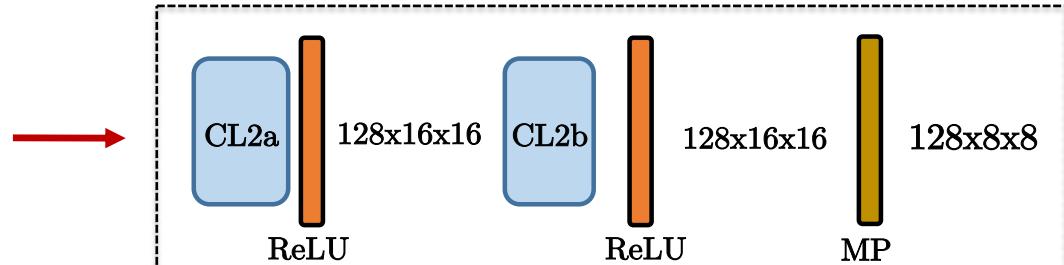
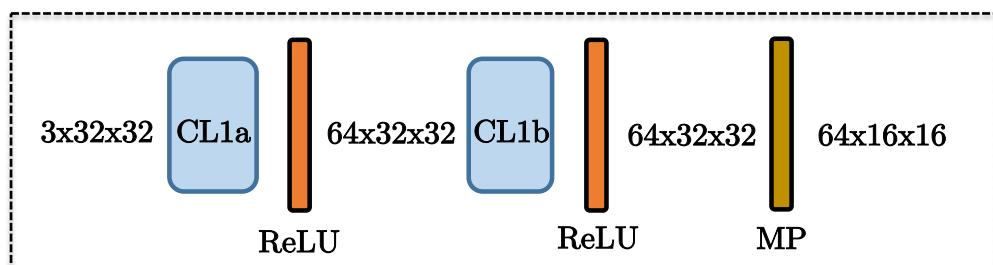


32 x 32 RGB image



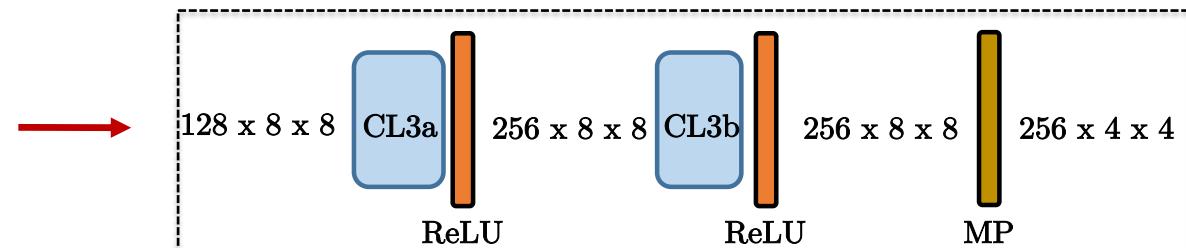
Four convolutional layers

- Convolutional/MP layers:
 - The number of features/channels increases by a factor 2.
 - The resolution of the activation maps decreases by a factor 2.

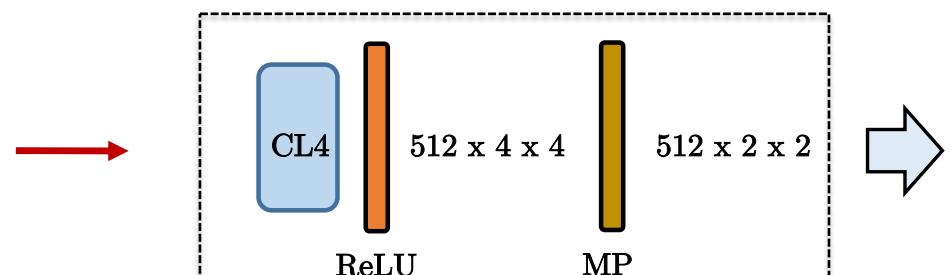


First block: from $3 \times 32 \times 32$ to $64 \times 16 \times 16$

Second block: from $64 \times 16 \times 16$ to $128 \times 8 \times 8$

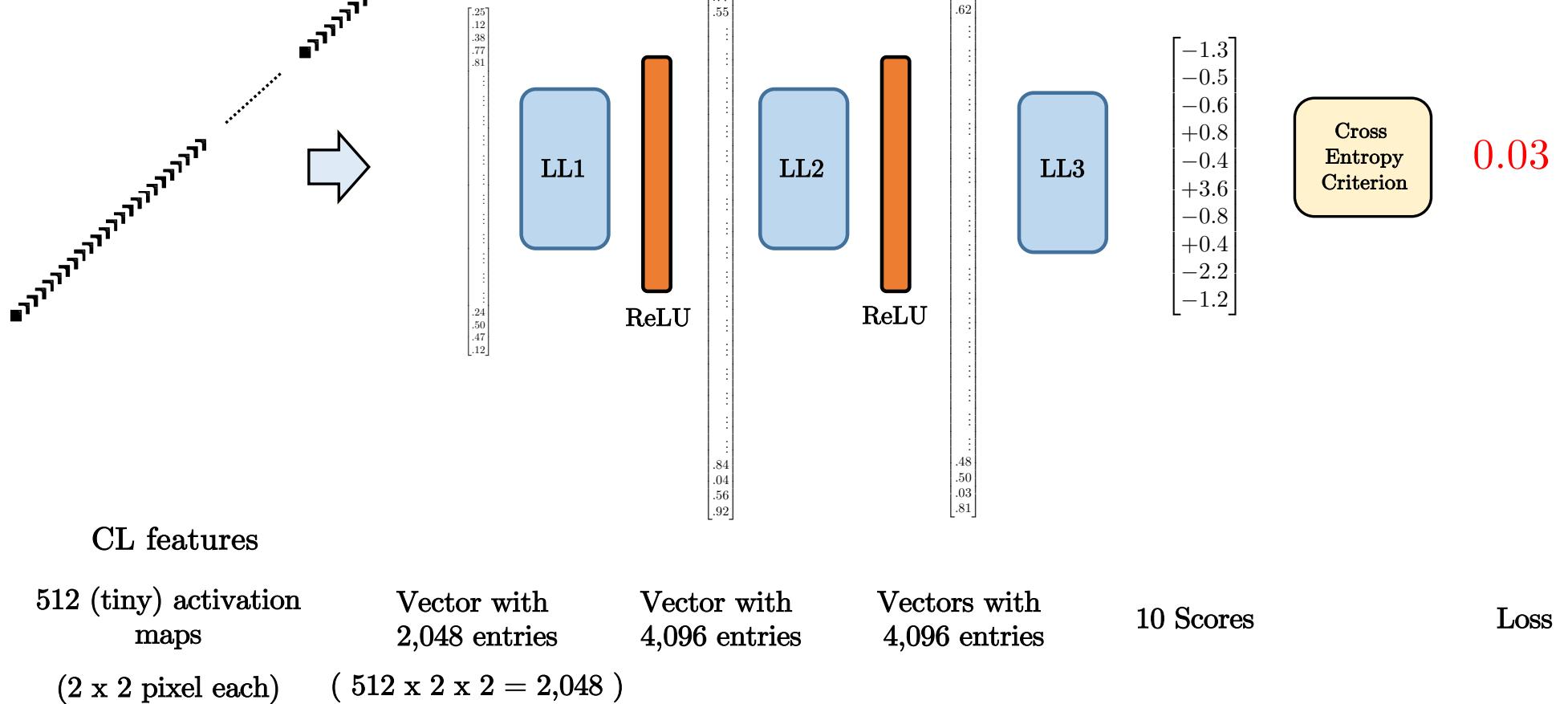


Third block: from $128 \times 8 \times 8$ to $256 \times 4 \times 4$



Fourth block: from $256 \times 4 \times 4$ to $512 \times 2 \times 2$

Linear layers



Number of architecture parameters

- First Block:
 - $3 \times 32 \times 32 \Rightarrow CL1a + ReLU \Rightarrow 64 \times 32 \times 32$
 - $64 \times 32 \times 32 \Rightarrow CL1b + ReLU \Rightarrow 64 \times 32 \times 32$
 - $64 \times 32 \times 32 \Rightarrow MP \Rightarrow 64 \times 16 \times 16$
 - $\#Parameters = 64 \times 3 \times 3 \times 3 + 64 + 64 \times 64 \times 3 \times 3 + 64 = 38K$
- Second Block:
 - $64 \times 16 \times 16 \Rightarrow CL2a + ReLU \Rightarrow 128 \times 16 \times 16$
 - $128 \times 16 \times 16 \Rightarrow CL2b + ReLU \Rightarrow 128 \times 16 \times 16$
 - $128 \times 16 \times 16 \Rightarrow MP \Rightarrow 128 \times 8 \times 8$
 - $\#Parameters = 128 \times 64 \times 3 \times 3 + 128 + 128 \times 128 \times 3 \times 3 + 128 = 221K$
- Third Block:
 - $128 \times 8 \times 8 \Rightarrow CL3a + ReLU \Rightarrow 256 \times 8 \times 8$
 - $256 \times 8 \times 8 \Rightarrow CL3b + ReLU \Rightarrow 256 \times 8 \times 8$
 - $256 \times 8 \times 8 \Rightarrow MP \Rightarrow 256 \times 4 \times 4$
- #Parameters = $256 \times 128 \times 3 \times 3 + 256 + 256 \times 256 \times 3 \times 3 + 256 = 885K$
- Fourth Block:
 - $256 \times 4 \times 4 \Rightarrow CL4 + ReLU \Rightarrow 512 \times 4 \times 4$
 - $512 \times 4 \times 4 \Rightarrow MP \Rightarrow 512 \times 2 \times 2$
 - $\#Parameters = 512 \times 216 \times 3 \times 3 + 512 + 512 \times 512 \times 3 \times 3 + 512 = 3.3M$
- LL:
 - $512 \times 2 \times 2 \Rightarrow \text{Reshape} \Rightarrow 2,048$
 - $2,048 \Rightarrow LL1 + ReLU \Rightarrow 4,096$
 - $4,096 \Rightarrow LL2 + ReLU \Rightarrow 4,096$
 - $4,096 \Rightarrow LL3 \Rightarrow 10$
 - $\#Parameters = 2048 \times 4096 + 4096 + 4096 \times 4096 + 4096 + 4096 \times 10 + 10 = 25.2M$

Total for 10 Layers \Rightarrow 30 Millions of parameters

Lab 04

- Learning rate schedule:

- Start with lr=0.25
- Divide by 2 at epoch 10
- Divide by 2 again at epoch 14
- Divide by 2 again at epoch 18

The screenshot shows a Jupyter Notebook interface with the title "Lab 04 : VGG architecture - exercise".

Code Snippet 1:

```
In [ ]: import torch
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim
from random import randint
import utils
import time
```

Text:

With or without GPU?

It is recommended to run this code on GPU:

- Time for 1 epoch on CPU : 841 sec (14.02 min)
- Time for 1 epoch on GPU : 9 sec w/ GeForce GTX 1080 Ti

Code Snippet 2:

```
In [ ]: device=torch.device("cuda")
device=torch.device("cpu")
print(device)
```

Text:

Download the CIFAR dataset

Code Snippet 3:

```
In [ ]: from utils import check_cifar_dataset_exists
data_path=check_cifar_dataset_exists()

train_data=torch.load(data_path+'cifar/train_data.pt')
train_label=torch.load(data_path+'cifar/train_label.pt')
test_data=torch.load(data_path+'cifar/test_data.pt')
test_label=torch.load(data_path+'cifar/test_label.pt')
```

Results

- Results with a **1-layer MLP** :

- Number of parameters: 0.03 million
- Error on training set: 59%
- Error on test set: 65%

- Results with a **3-layer MLP** :

- Number of parameters: 2 millions
- Error on training set: 0%
- Error on test set: 46%

- Results with a **10-layer ConvNet** :

- Number of parameters: **30 millions**
- Error on training set: 0%
- Error on test set: **20%**

Timing for one epoch:

- CPU : 14 min
- Tesla K80 GPU (Google Cloud): 1.2 min
- GTX 1080 GPU (\$600) : 9 sec

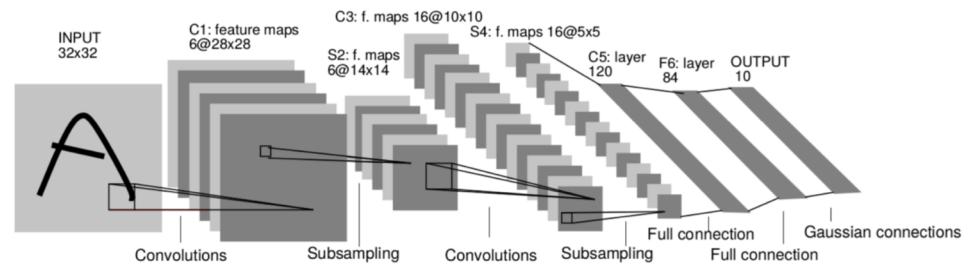
Outline

- Padding
- Stride
- Multi-dimensional input
- Convolution as linear operation
- Depth
- Convolutional layer with PyTorch
- Pooling layer with PyTorch
- LeNet5 architecture
- VGG architecture
- **Case studies**

ConvNet architectures

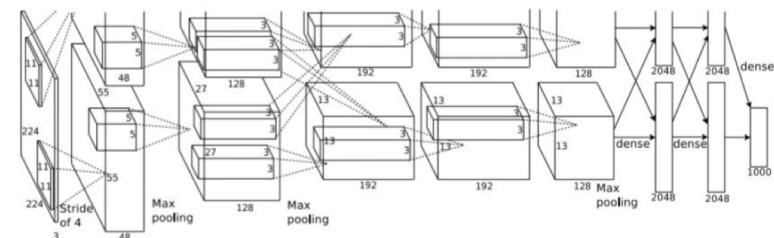
- **LeNet5** [LeCun, Bottou, Bengio, Haffner, 1998]:

- Input is 28x28
- Architecture is CL-MP-CL-MP-LL-LL
- Accuracy on MNIST is 99.6%
- #Parameters=0.6M



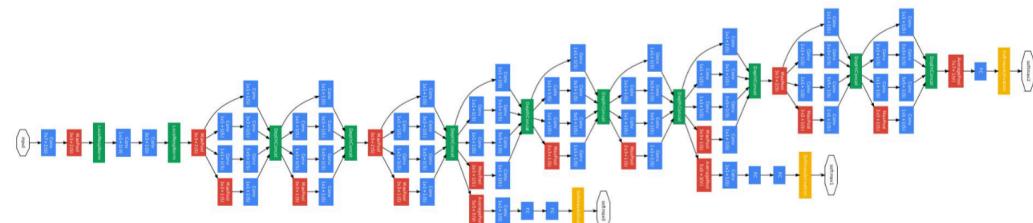
- **AlexNet** [Krizhevsky, Hinton, Sutskever, 2012]:

- Input is 227x227x3
- Architecture is 7CL-3MP-2FC
- Prediction error on ImageNet is 15.4%
- 1st rank
- #Parameters=62M



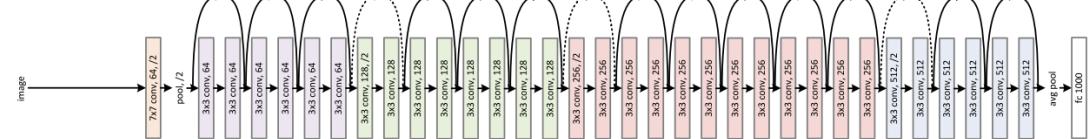
ConvNet architectures

- **VGGNet** [Simonyan, Zisserman, 2014]:
 - Input is 227x227x3
 - Architecture has 16 layers
 - Prediction error on ImageNet is 7.3%
 - 2nd rank
 - #Parameters=138M
- **GoogleNet/Inception** [Szegedy-et.al, 2014]:
 - Input is 227x227x3
 - Architecture has 22 layers
 - Prediction error on ImageNet is 6.7%
 - 1st rank
 - #Parameters=5M

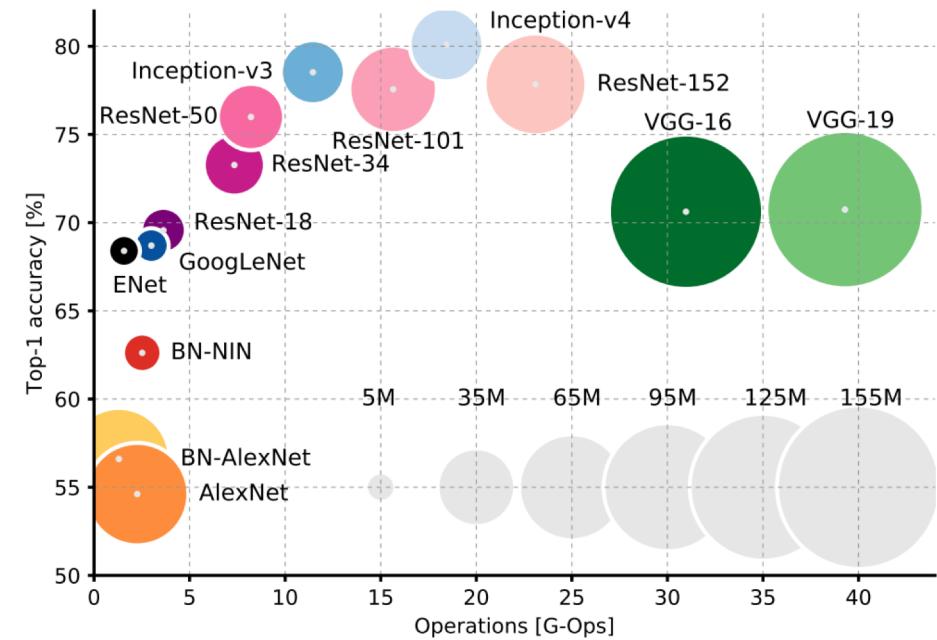
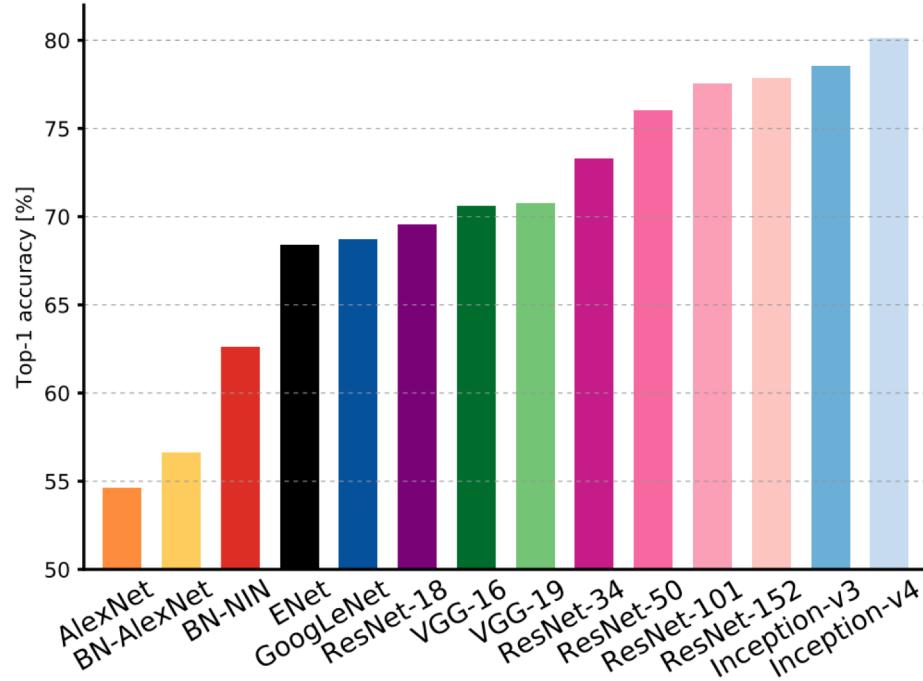


ConvNet architectures

- Microsoft/ResNet [He-et.al'15]:
 - Input is 227x227x3
 - Architecture has 152 layers
 - Prediction error on ImageNet is 3.6%
 - 1st rank
 - #Parameters=60M
- Etc.



ConvNet architectures





Questions?