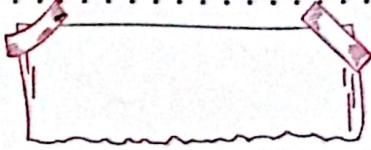


AS

\*PT - record  
layout



## Convolutional Neural Networks

Until now, we simply discarded each image's spatial structure by flattening them into one dimensional vectors, feeding the MLP parameters. Preferably, we would leverage our prior knowledge that nearby pixels are related to each other.

### Invariance

Not be overly concerned with the precise location of image.

CNN suitable for computer vision

1. Translation invariance

2. Localization principle

CNNs are special family of NN that contain convolutional layers.

Convolutional kernel or filter or simply the layer's weights that are often a learnable parameter

### Convolution

$$(f * g)(x) = \int f(z)g(x-z)$$

padding (outer), stride, pooling

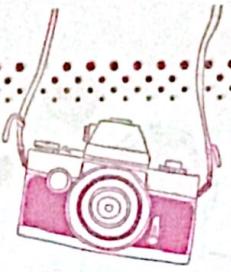
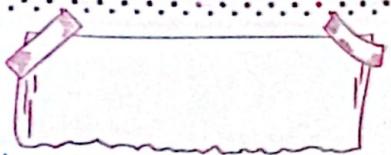
input                          Kernel                          output

[ 0 1 2 ]	*	[ 0 1 ]	=	[ 19 25 ]
[ 3 4 5 ]		[ 2 3 ]		[ 37 43 ]
[ 6 7 8 ]				

Meow  
Meow



My playlist today #



## Lenet

Consists of two parts:

- (i) Convolutional encoder with two convolutional layers
- (ii) Dense block with three fully connected layers

Basic units in each convolutional block:

- Convolutional layer
- Sigmoid activation function
- Average pooling operation

Note that while ReLU and max-pooling work better, these discoveries had not yet been made in the 1990s.

Each convolutional layer:

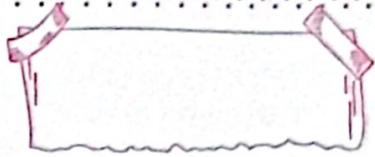
- $5 \times 5$  kernel
- sigmoid activation function

These layers map spatially arranged inputs to a number of two dimensional feature maps, typically increasing the number of channels.

In order to pass output from the convolutional block to dense block, we must flatten each example in a minibatch.



IDEA



## Deep Convolutional Neural Networks

Rather than training end-to-end (1 pixel to classification) systems, classical pipelines looked more like this:

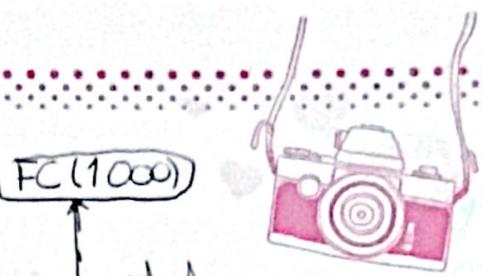
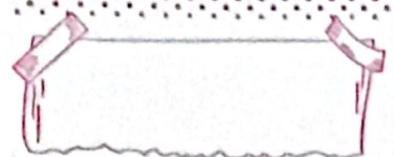
1. Obtain an interesting dataset. In the early days, these datasets required expensive sensors (at the time, 1 megapixel images were state-of-art)
2. Preprocess the dataset with hand-crafted features based on some knowledge of optics, geometry, other analytic tools, and occasionally on the serendipitous discoveries of likely graduate students
3. Feed the data through a standard set of feature extractors such as SIFT (Scale Invariant Feature Transform), the SURF (Speed up Robust Features) or any number of other hand-tuned pipeline
4. Dump the resulting representations into your favorite classifier, likely a linear model or a Kernel method, to train a class

Features ought to be learned - Alexnet

Higher layers of the representation to represent layer structures: eyes, noses, Even larger, final hidden state learns a compact representation of the image that summarizes its contents such that data belonging

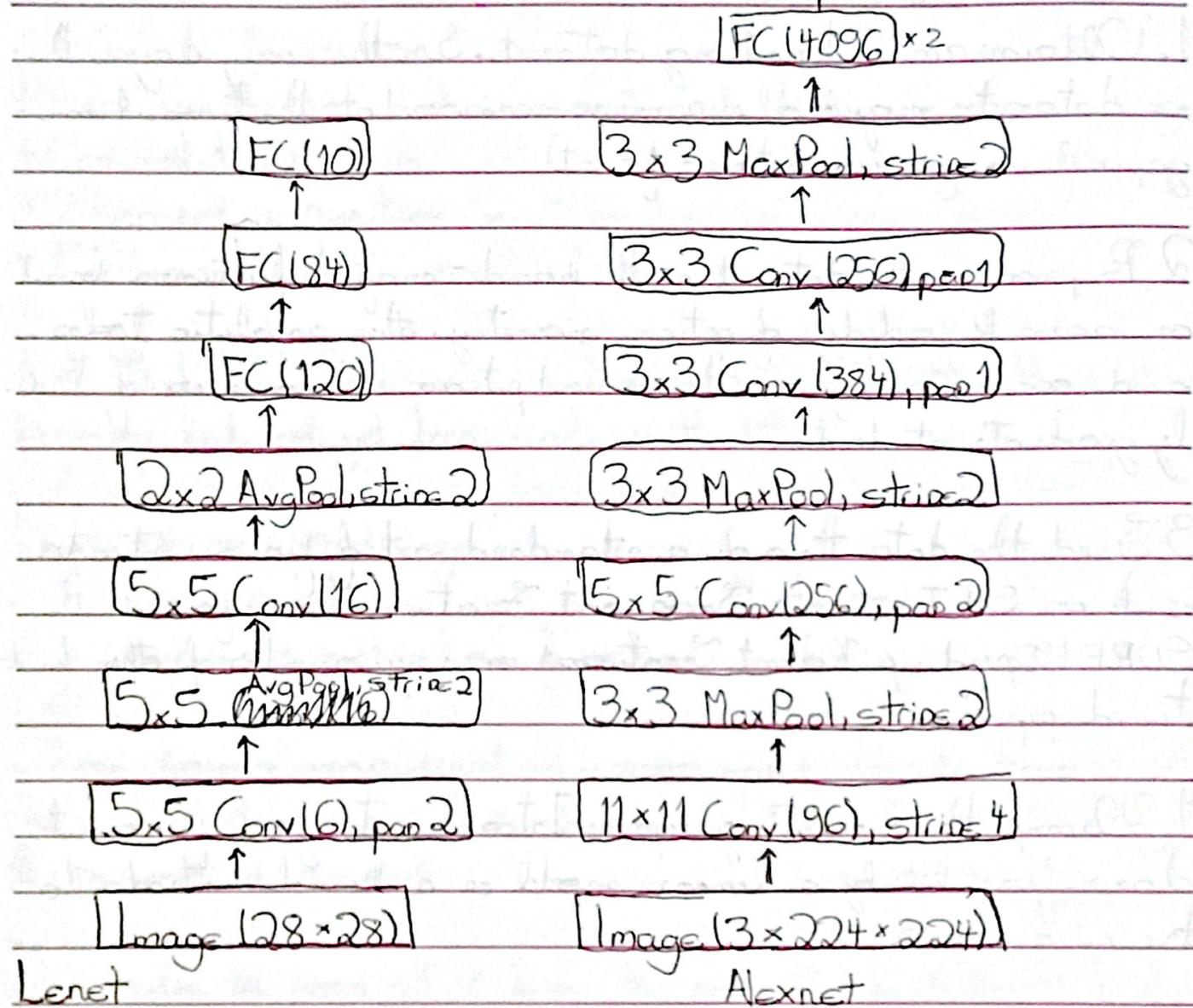


My playlist today: #



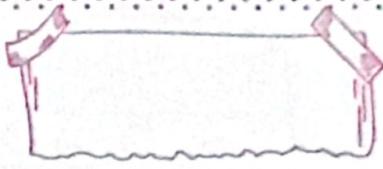
to different categories that can be easily separated

Mining ingredients: data, hardware  $\rightarrow$  GPU, cuda



Alexnet is much deeper than the comparatively small LeNet.  
Second, Alexnet used ReLu instead of sigmoid as its activation function.

AB



Since most images in ImageNet are ten times higher and wider than MNIST images, objects in ImageNet data tend to occupy more pixels. Consequently, a larger convolutional window is needed to capture the object. Moreover, Alexnet has ten times more convolutional channels than Lenet.

Alexnet controls the model complexity of the full connected layer by dropout while Lenet only uses weight decay.

To augment the data even further, the training loop of Alexnet added a great deal of image augmentation, such as flipping, clipping, and color changes. This makes the model more robust and larger sample size effectively reduces overfitting.

### Networks using blocks (VGG)

The idea of using blocks first emerged from the Visual Geometry Group at Oxford University, in their VGG network.

#### VGG blocks

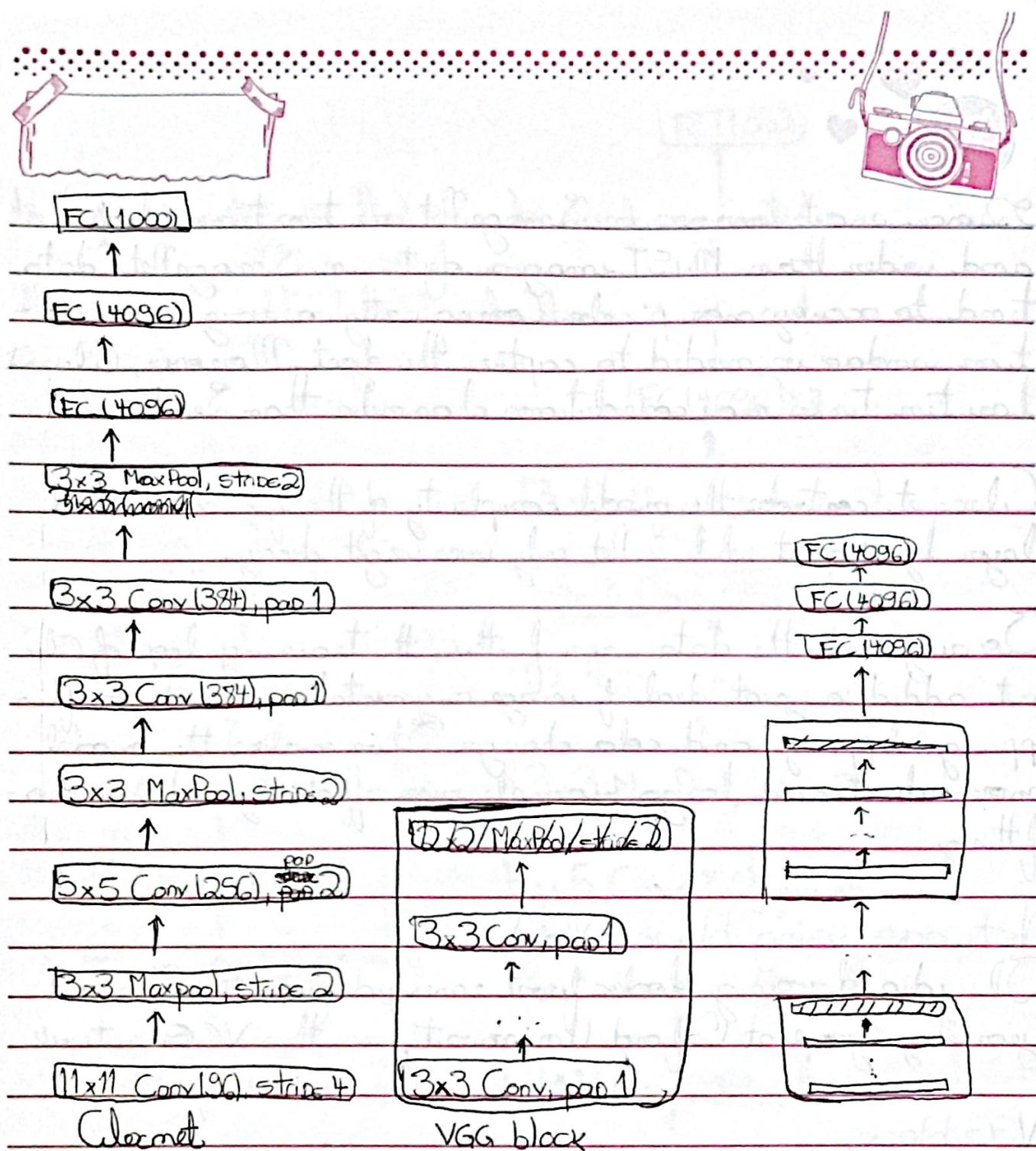
The basic idea of a classic CNN is a sequence of the following:

- i) A convolutional layer with padding to maintain resolution
- ii) A nonlinearity such as a ReLU
- iii) A pooling layer such as max

One VGG block consists a sequence of convolutional layers, followed by a maximum layer for spatial down-sampling.

My playlist today #





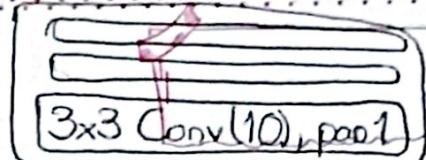
## Network in Network (NIN)

Use an MLP on the channels for each pixel separately

### NIN block

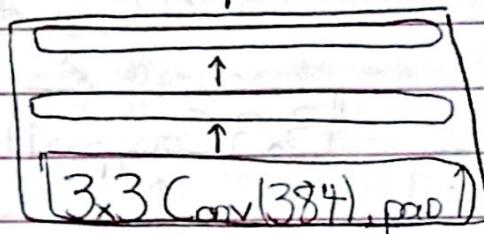
If we tie the weight across each spatial location, we could think of this as a  $1 \times 1$  convolutional layer or as a fully-connected layer.

Global Avg Pool

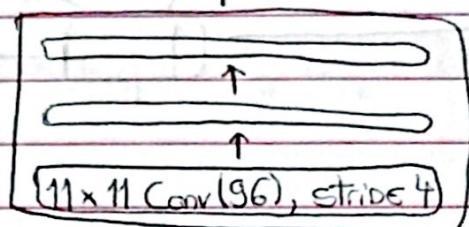


acting independently on each pixel location ↑

3x3 Maxpool, stride 2



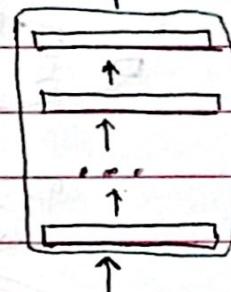
3x3 MaxPool, stride 2



FC (1000)

FC (4096)

FC (4096)

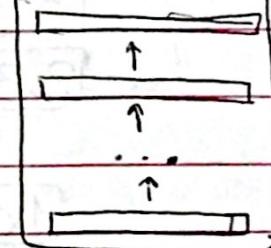


VGG block

3x3 MaxPool, stride 2

3x3 Conv, pool 1

3x3 Conv, pool 1



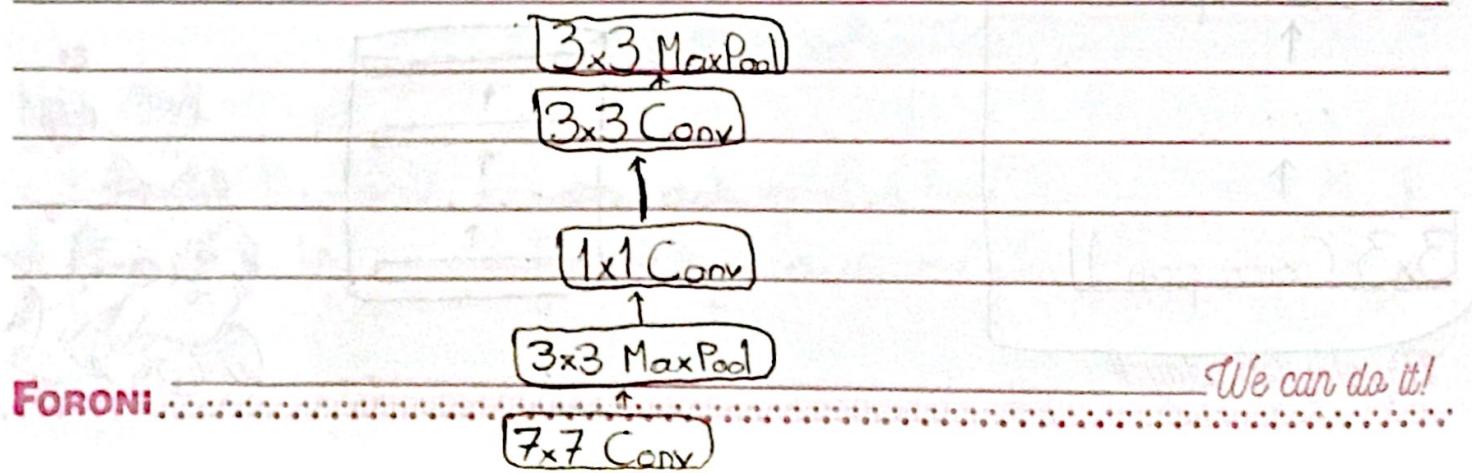
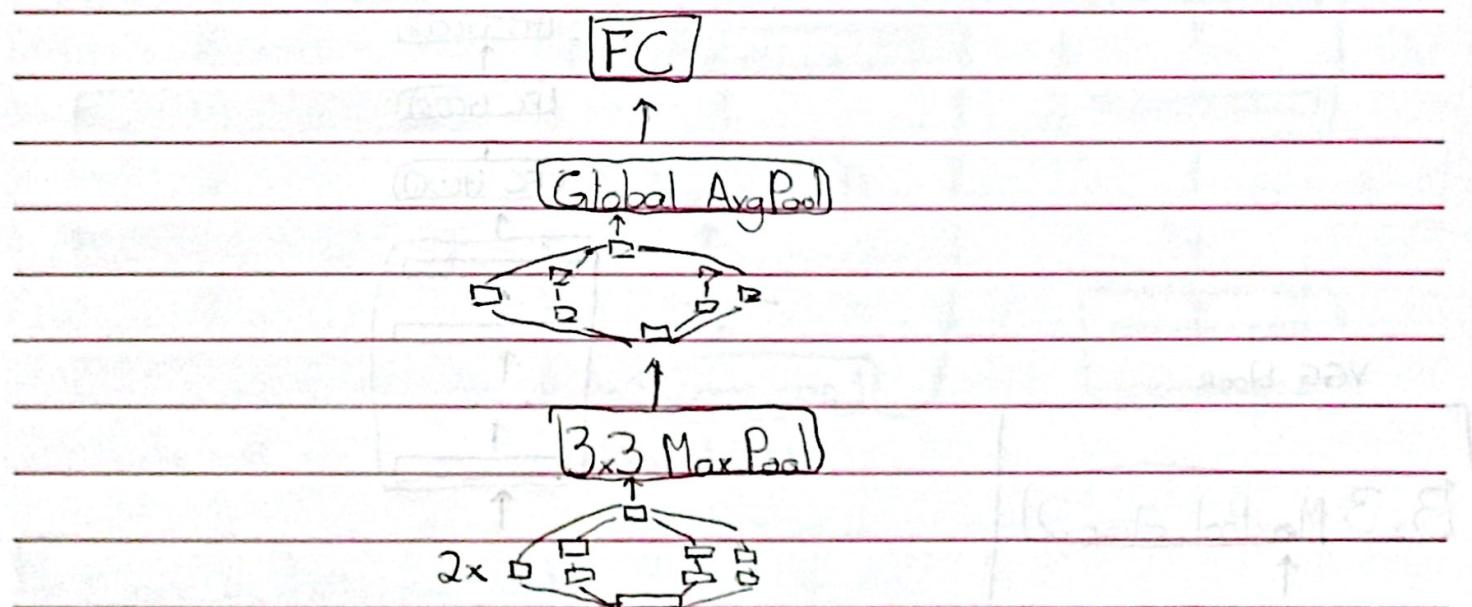
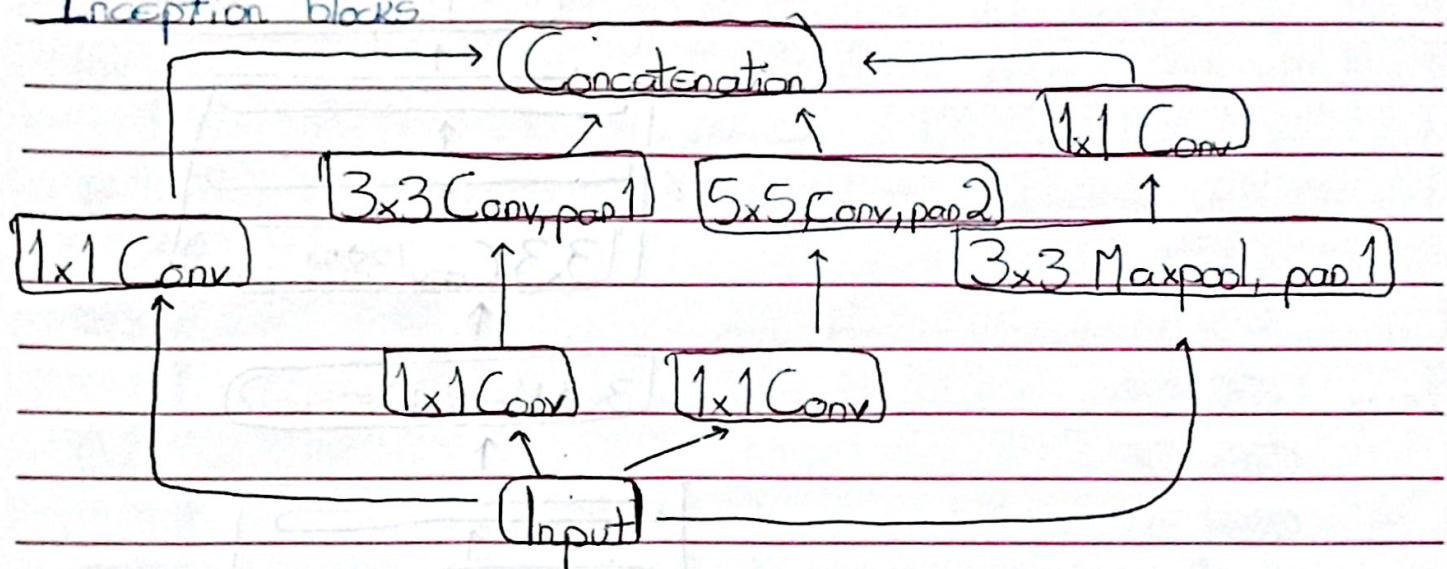
My playlist today. #

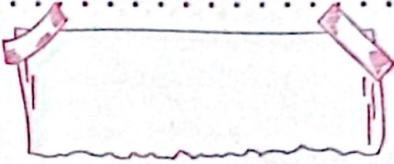




## Networks with parallel concatenation (GoogleLeNet)

### Inception blocks





## Batch Normalization

### ML Training challenges

1. Choices regarding data processing often make an enormous difference in results.
2. Typically as we train, the variables in intermediate layers may take values with varying magnitudes.
3. Deeper networks are complex and easily capable of overfitting.

Batch normalization is applied to individual layers and works as follows:

In each iteration, we first normalize the inputs by subtracting their mean and dividing by standard deviation, where both are estimated based on the statistics of the current minibatch.

Next, we apply a scale coefficient and a scale offset.

It is precisely due to this normalization based on batch statistics that batch normalization derives its name.

The choice of batch size affects batch normalization.

My playlist today #





Formally, denoting by  $x \in \mathcal{B}$  an input to batch normalization (BN) that is from a minibatch  $\mathcal{B}$

$$BN(x) = \gamma \circ x - \hat{\mu}_{\mathcal{B}} + \beta \quad \hat{\mu}_{\mathcal{B}} = \text{sample mean} \\ \hat{\sigma}_{\mathcal{B}} = \text{sample std}$$

$$\hat{\sigma}_{\mathcal{B}}^2 = \frac{1}{|\mathcal{B}|} \sum_{x \in \mathcal{B}} (x - \hat{\mu}_{\mathcal{B}})^2 + \epsilon \quad \hat{\mu}_{\mathcal{B}} = \frac{\sum_{x \in \mathcal{B}} x}{|\mathcal{B}|}$$

Note that we add a small constant  $\epsilon > 0$  to the variance estimate to ensure that we never attempt to divide by 0. You

Batch normalization layers function differently in training mode (normalizing by minibatch statistics) and in prediction mode (normalizing by dataset statistics)

### Residual Networks (ResNet)

Every additional layer should more easily contain the identity function as one of its elements. These considerations are

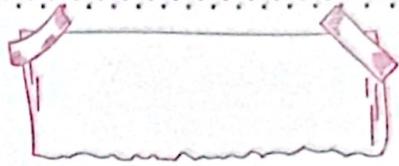
### Densely Connected Networks (DenseNet)

Taylor expansion for functions

Given a point  $x = 0$

$$f(0) + f'(0)x + f''(0)\frac{x^2}{2!} +$$

$3!$



The key point is that decomposes functions into

$$f(x) = x + g(x)$$

My playlist today #

