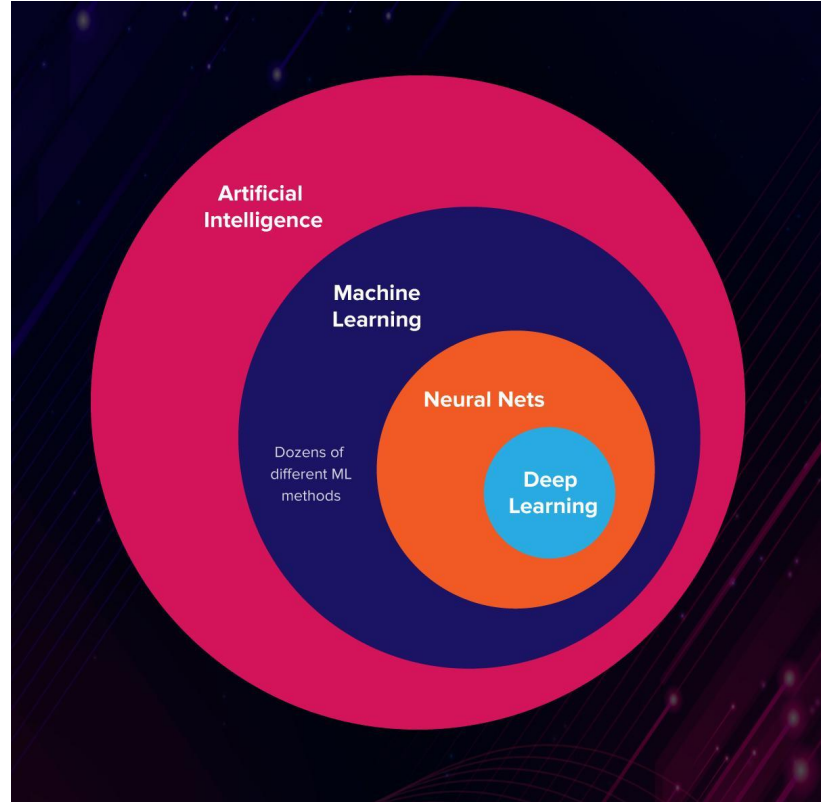# Intro to DL

Serafina Di Gioia
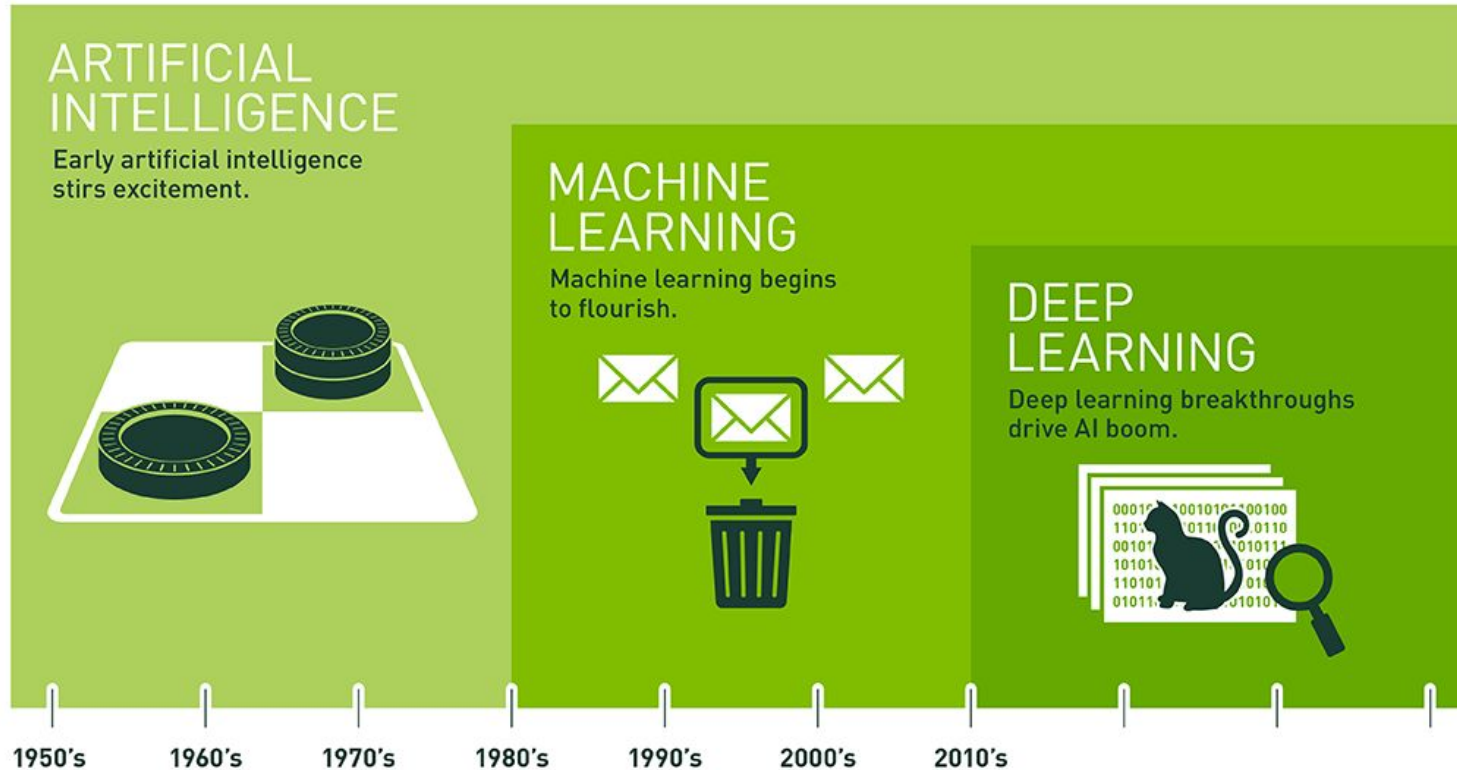PostDoctoral Researcher @ ICTP

# Where is DL in the picture?

**Deep learning:**

a type of machine learning based on artificial neural networks in which multiple layers of processing are used to extract progressively higher level features from data.

# From LISP to the DL revolution…



ARTIFICIAL INTELLIGENCE
Early artificial intelligence stirs excitement.

MACHINE LEARNING
Machine learning begins to flourish.

DEEP LEARNING
Deep learning breakthroughs drive AI boom.

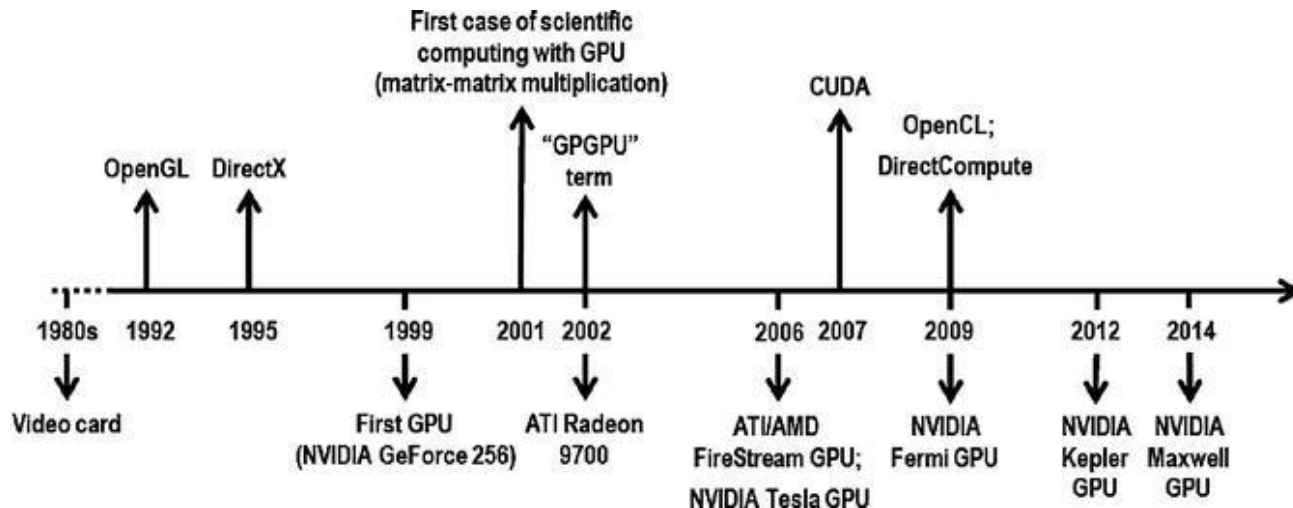1950's    1960's    1970's    1980's    1990's    2000's    2010's

Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.
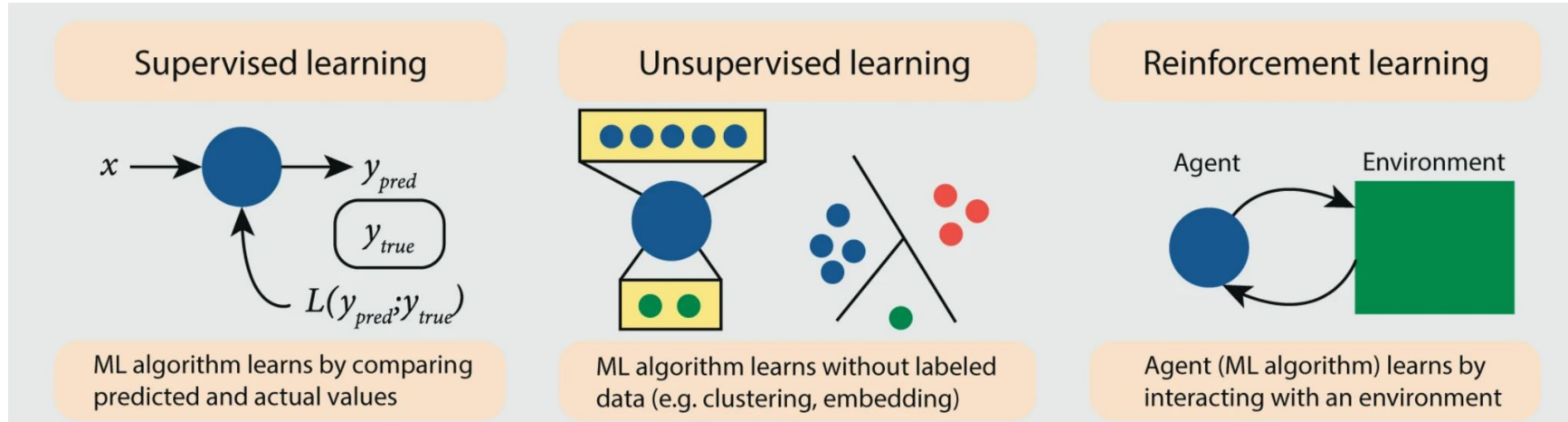
# Main ingredients for DL breakthrough

- large datasets available (e.g IMAGENET)
- GPUs development (in particular, CUDA introduction)
- increased involvement of developers from CV and scientific communities

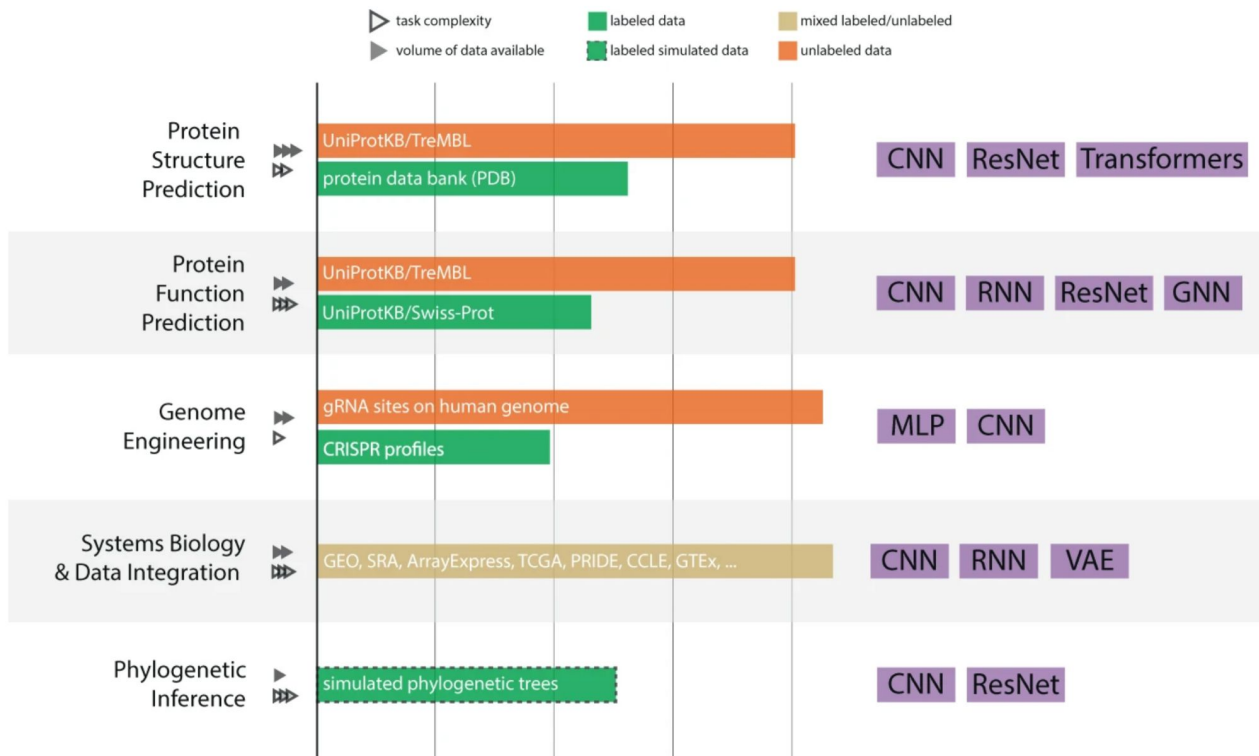The DL era starts few years after that CUDA came to light

# ML scenarios



**Supervised learning**

$x \rightarrow$ (node) $\rightarrow y_{pred}$

$y_{true}$

$L(y_{pred}; y_{true})$

ML algorithm learns by comparing predicted and actual values

**Unsupervised learning**

ML algorithm learns without labeled data (e.g. clustering, embedding)

**Reinforcement learning**

Agent     Environment

Agent (ML algorithm) learns by interacting with an environment

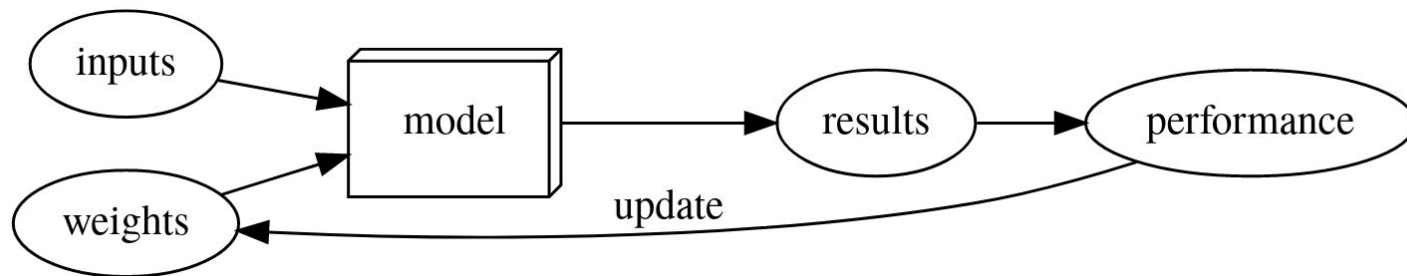| Learning type | Model building | Examples |
|---|---|---|
| Supervised | Algorithms or models learn from labeled data (task-driven approach) | Classification, regression |
| Unsupervised | Algorithms or models learn from unlabeled data (Data-Driven Approach) | Clustering, associations, dimensionality reduction |
| Semi-supervised | Models are built using combined data (labeled + unlabeled) | Classification, clustering |
| Reinforcement | Models are based on reward or penalty (environment-driven approach) | Classification, control |

# Typical sizes of DL data sets



The increasing complexity of the new datasets, typical of big data epoch motivates the need for GPU-based libraries
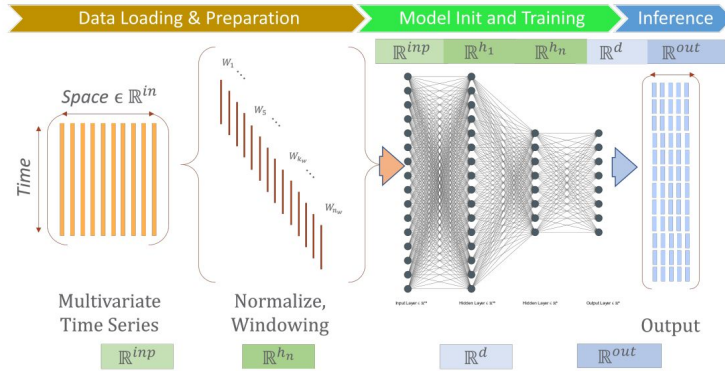
# Building blocks of ML algorithms

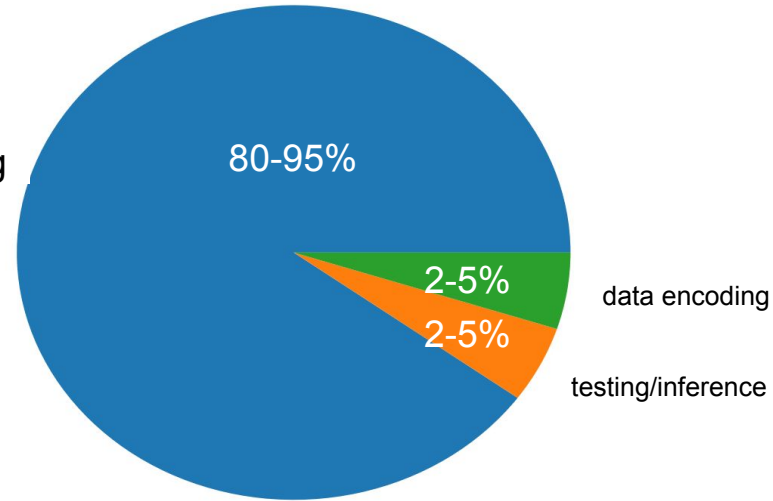ML algorithms have three main components

1. **decision process**: based on some input data, which can be labeled or unlabeled, your algorithm will produce an estimate about a pattern in the data. This estimate can be used to solve a prediction or classification task

2. **error function**: it evaluates the prediction of the model. If there are known examples, an error function can make a comparison to assess the accuracy of the model.

3. **Model Optimization Process**: If the model can fit better to the data points in the training set, then weights are adjusted to reduce the discrepancy between the known example and the model estimate. The algorithm will repeat this "evaluate and optimize" process, updating weights autonomously until a threshold of accuracy has been met.
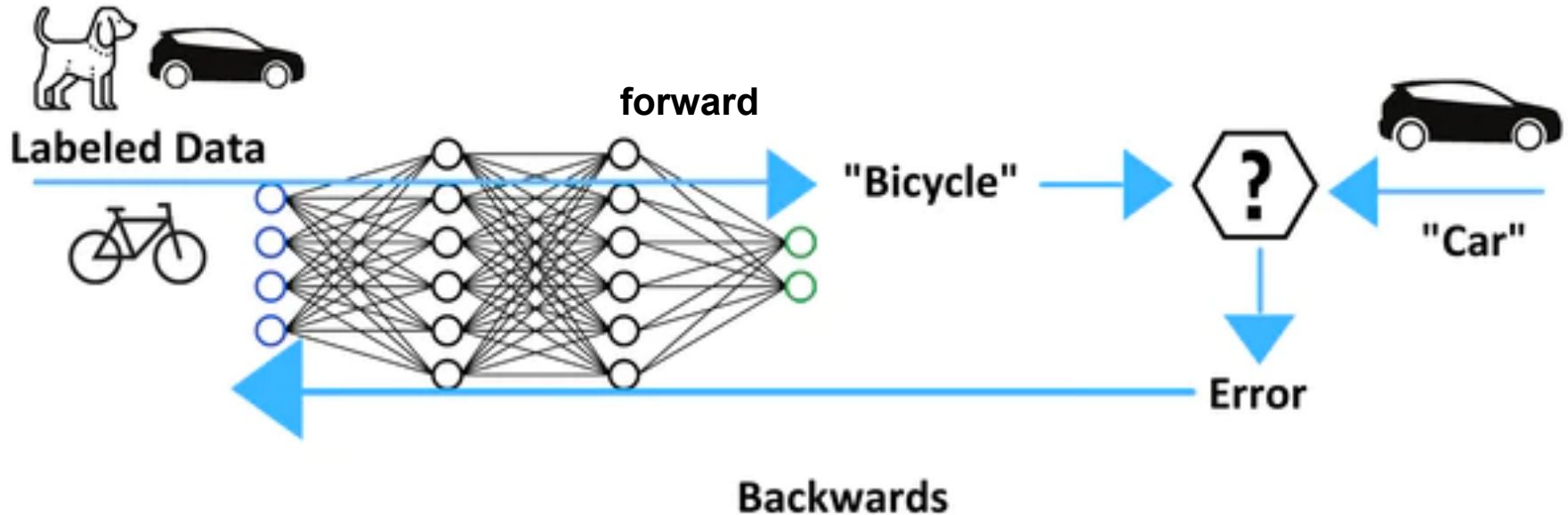
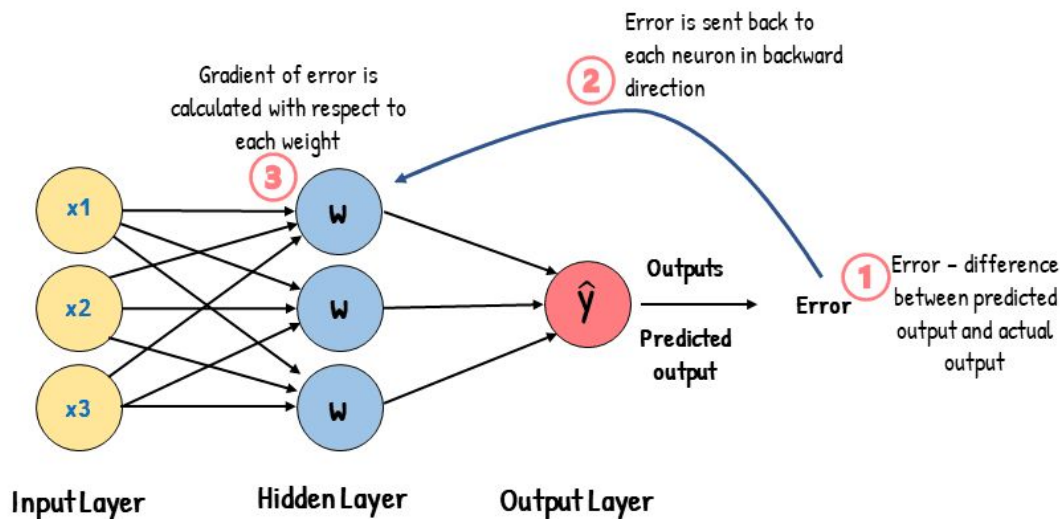# Analyzing the computational workload of DL models



training 80-95%

data encoding 2-5%

testing/inference 2-5%

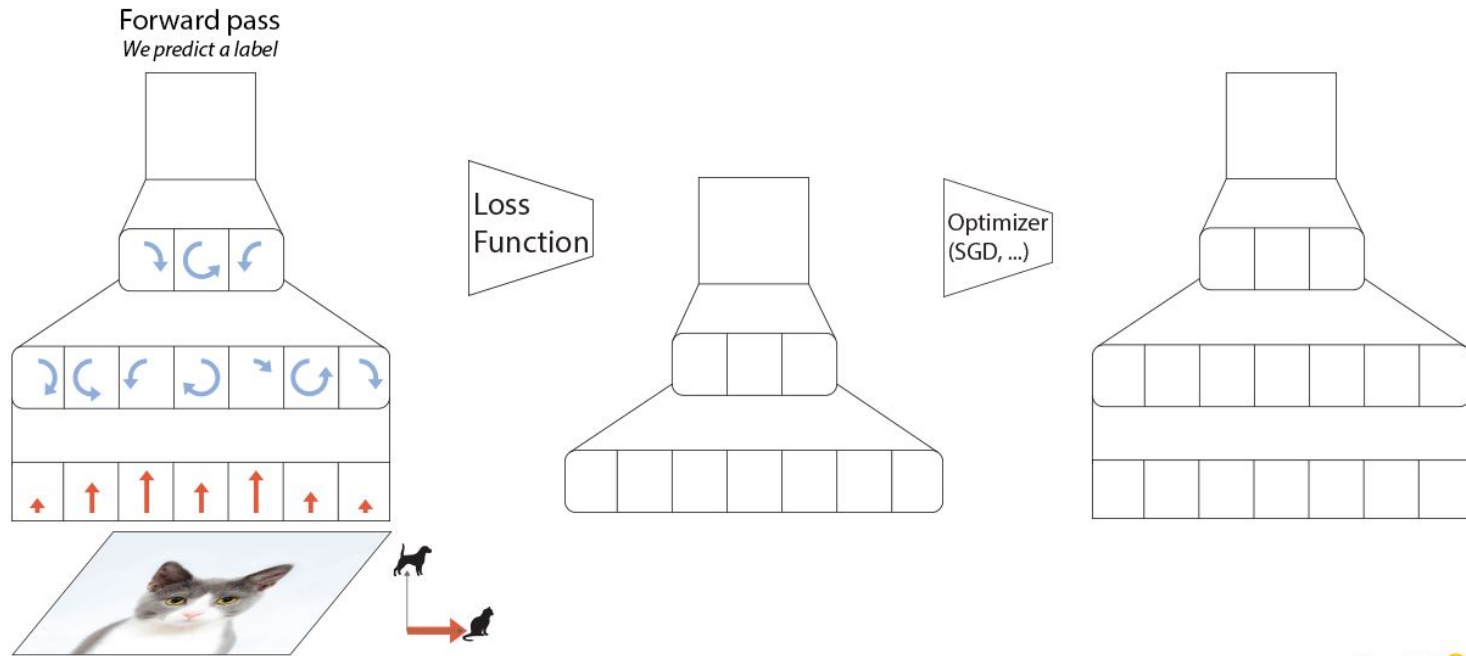# Training a DL model (in a supervised setting)

# Back-propagation

*The workhorse of DL is the Backpropagation algorithm (Rumelhart et al., 1986).*

It allows for efficient gradient computation by recursively applying the chain rule of calculus. It owes his name to the presence of a 'backward pass' of an error signal through the neural network.
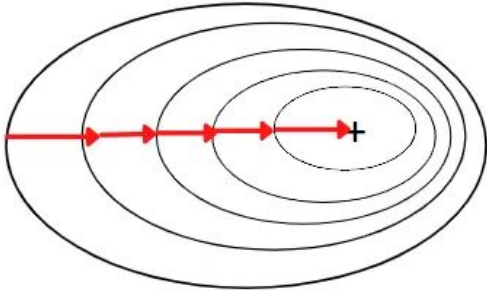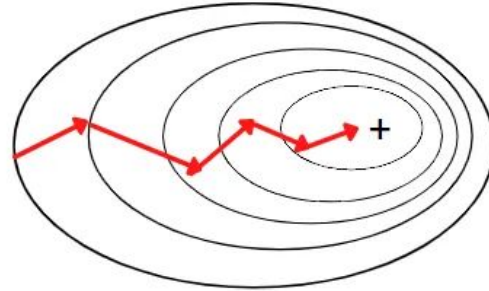
Forward pass
*We predict a label*

Loss Function

Optimizer (SGD, ...)

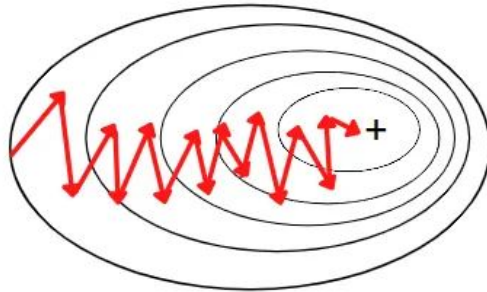@Thom_Wolf

# Gradient Descent variants

**Batch Gradient Descent**

**Mini-Batch Gradient Descent**

**Stochastic Gradient Descent** (SGD)

# Alternative Optimizers

Other optimizers have been proposed to enhance the speed and convergence of the training process:

• **SGD with Momentum** (Polyak, 1964): Speeds up gradient descent by adding a fraction of the previous update to the current one.

• **RMSprop** (Hinton, 2012): Adapts the learning rate for each parameter based on recent gradient magnitudes.

• **Adam** (Kingma and Ba, 2015): Combines momentum and adaptive learning rates for more efficient optimization.

# Hyperparameters importance

Hyperparameters are an overlooked but crucial factor in DL practise. They include:

- learning rate
- training steps/epochs
- batch size
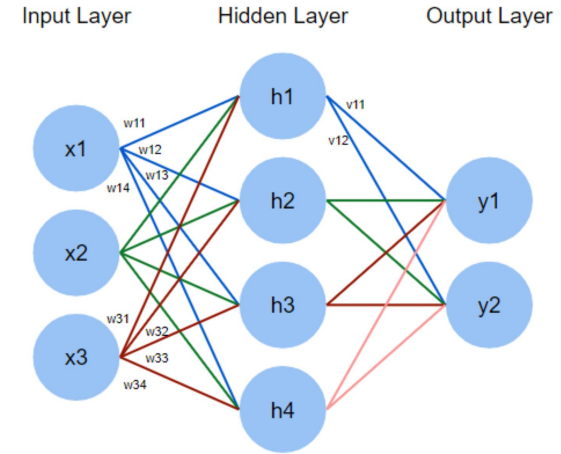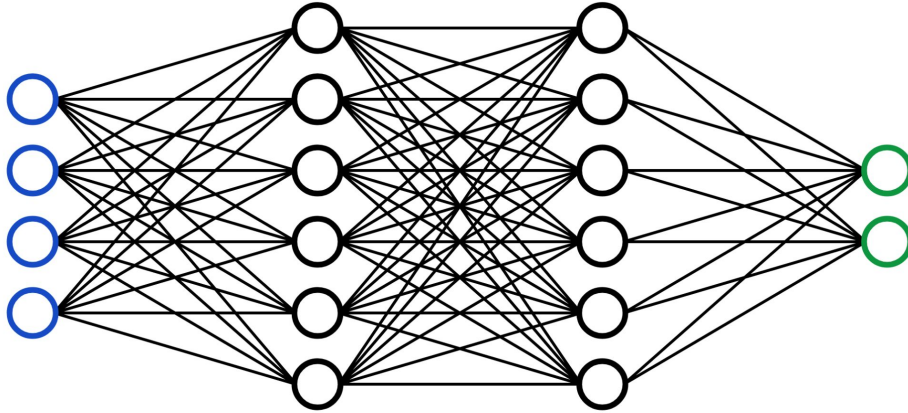- Optimizer choice-setting

**How to choose them?**

• Experience

• Trial and error

# Universal approximation theorem

*Universal approximation theorem (Hornik, 1991):*

*A feedforward neural network with at least one hidden layer can approximate any continuous function to any desired accuracy, given enough neurons and the right activation function.*
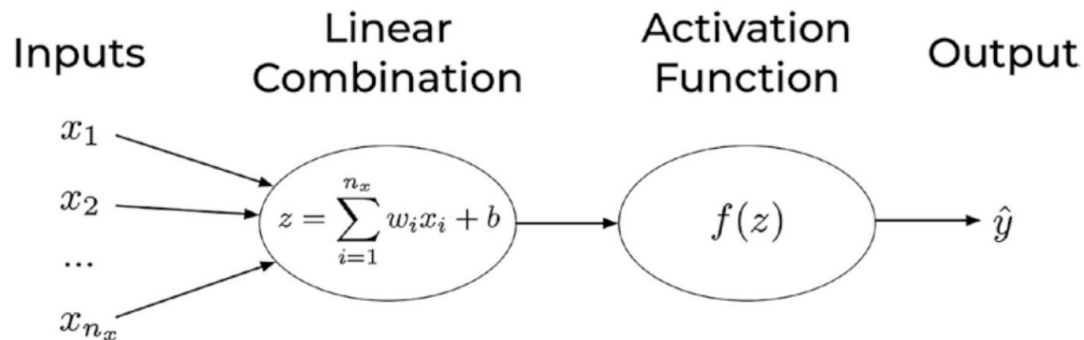
# Linear Neural Network (lNN)



$$
\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} * \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & w_{1,4} \\ w_{2,1} & w_{2,2} & w_{2,3} & w_{2,4} \\ w_{3,1} & w_{3,2} & w_{3,3} & w_{3,4} \end{bmatrix} = \begin{bmatrix} h'_1 & h'_2 & h'_3 & h'_4 \end{bmatrix}
$$

Alias of nn.linear() in Pytorch:   torch.mm(inputs, linear.weight.T).add(linear.bias)

# Artificial neuron

A neuron has two main components:

- The weights ($w_i$) (the bias $b$ is sometime included in the weights)
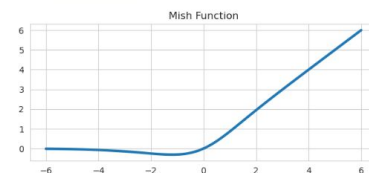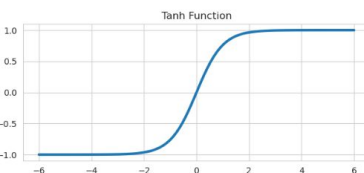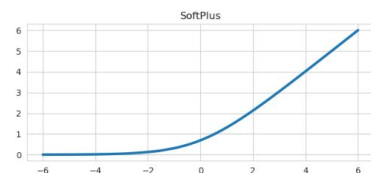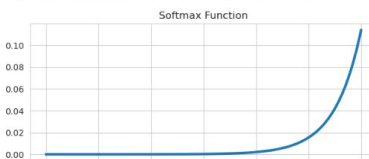- The *activation function* (the $f$



| Inputs | Linear Combination | Activation Function | Output |
|---|---|---|---|

$x_1$

$x_2$

...

$x_{n_x}$

$$z = \sum_{i=1}^{n_x} w_i x_i + b$$

$f(z)$

$\hat{y}$

**Which portion of a neural network model is responsible for the type of problem that can be solved?**

Two main components are responsible for the type of problem that can be solved:

- The output activation function

- The loss function

The optimiser is not related in any way to the type of problem solved (it does not depend on the type of the response variable).

# Types of activation functions

# Perceptron

Rosenblatt said:

"A perceptron is first and foremost a **brain model, not an invention for pattern recognition**. As a brain model, its utility is in enabling us to determine the physical conditions for the emergence of variou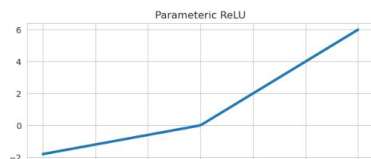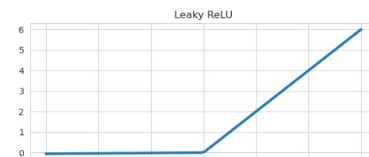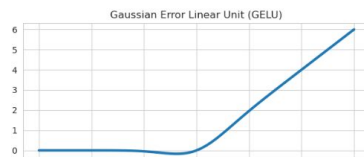s psychological properties. It is by no means a "complete" model, and we are fully aware of the simplifications which have been made from biological systems; but it is, at least, an analyzable model."



the perceptron can be described by:

- a linear function that aggregates the input signals
- a threshold-activation function that determines if the response neuron fires or not
- a learning procedure to adjust connection weights

# Learning procedure for the Perceptron

# Adaline vs Perceptron



**Perceptron training loop**

Aggregation function

Threshold function

$$z = b + \sum_{i=1}^{n} w_i x_i \quad f(z)$$

Output

$$\hat{y} = \in -1, 1$$

$$error = y - \hat{y}$$

$$\Delta w = \eta(error)x_k$$

$$w_{k+1} = w_k + \Delta w_k$$

Input units   Weight update   Error computation

**ADALINE training loop**

Aggregation function

Threshold function

$$\hat{y} = b + \sum_{i=1}^{n} w_i x_i \quad f(\hat{y})$$

Output

$$\hat{y}' = \in -1, 1$$

$$error = (\hat{y} - y)^2$$

$$w_{k+1} = w_k - \eta(2 * error_k)x_k$$

Input units   Weight update   Error computation

it introduces SGD in the training

# Multi-layer Perceptron



linear function

$$z_m = b + \sum_m x_n w_{mn}$$

matrix multiplication
input to hidden layer

sigmoid activation

matrix multiplication
hidden to output layer

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$a_m = \sigma(z_m) = \frac{1}{1 + e^{-z_m}}$$

$$\begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

$z_1 \mid a_1 = \sigma(z_1)$

$W_1$

$X_1$

$W_{11}$
$W_{12}$
$W_{13}$

$z_2 \mid a_2 = \sigma(z_2)$

$\mathbf{W_2}$

sigmoid activation

decision threshold

$W_{21}$
$W_{22}$

$z_1 \mid a_1 = \sigma(z_1)$

$\hat{y} = f(a_m) \begin{cases} +1, & \text{if } a > 0.5 \\ -1, & \text{otherwise} \end{cases}$

$X_2$

$W_{23}$

$z_3 \mid a_3 = \sigma(z_3)$

$\mathbf{W_{21}}$

input
layer

hidden
layer

output
layer

classification

Information flow forward pass neural network with one hidden layer

SMR 4067- AI and Climate Modeling
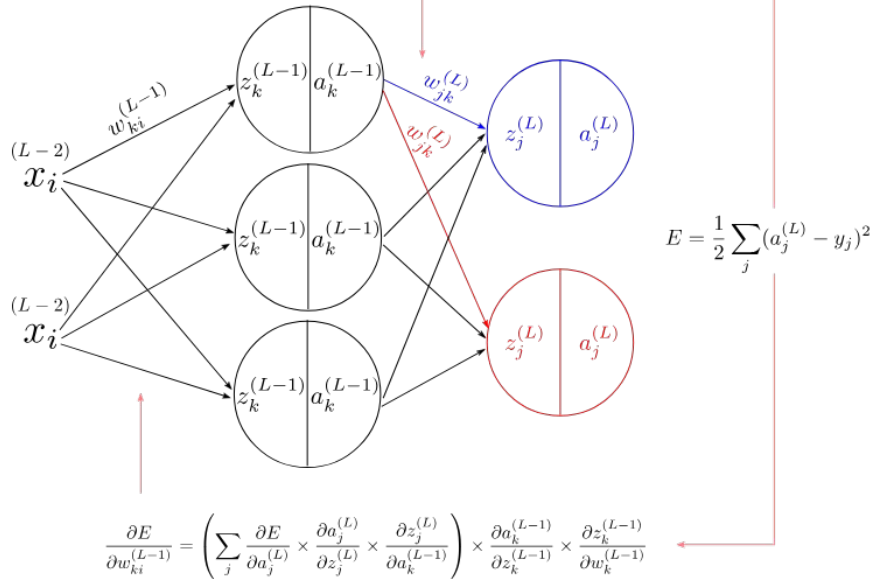
# Back-propagation in multi-layer perceptron



derivative of the error w.r.t. weights in (L)

$$\frac{\partial E_i}{\partial w_{jk}^{(L)}} = \frac{\partial E_i}{\partial a_j^{(L)}} \times \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \times \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}}$$

$$E = \frac{1}{2}\sum_j (a_j^{(L)} - y_j)^2$$

$$\frac{\partial E}{\partial w_{ki}^{(L-1)}} = \left( \sum_j \frac{\partial E}{\partial a_j^{(L)}} \times \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \times \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \right) \times \frac{\partial a_k^{(L-1)}}{\partial z_k^{(L-1)}} \times \frac{\partial z_k^{(L-1)}}{\partial w_k^{(L-1)}}$$

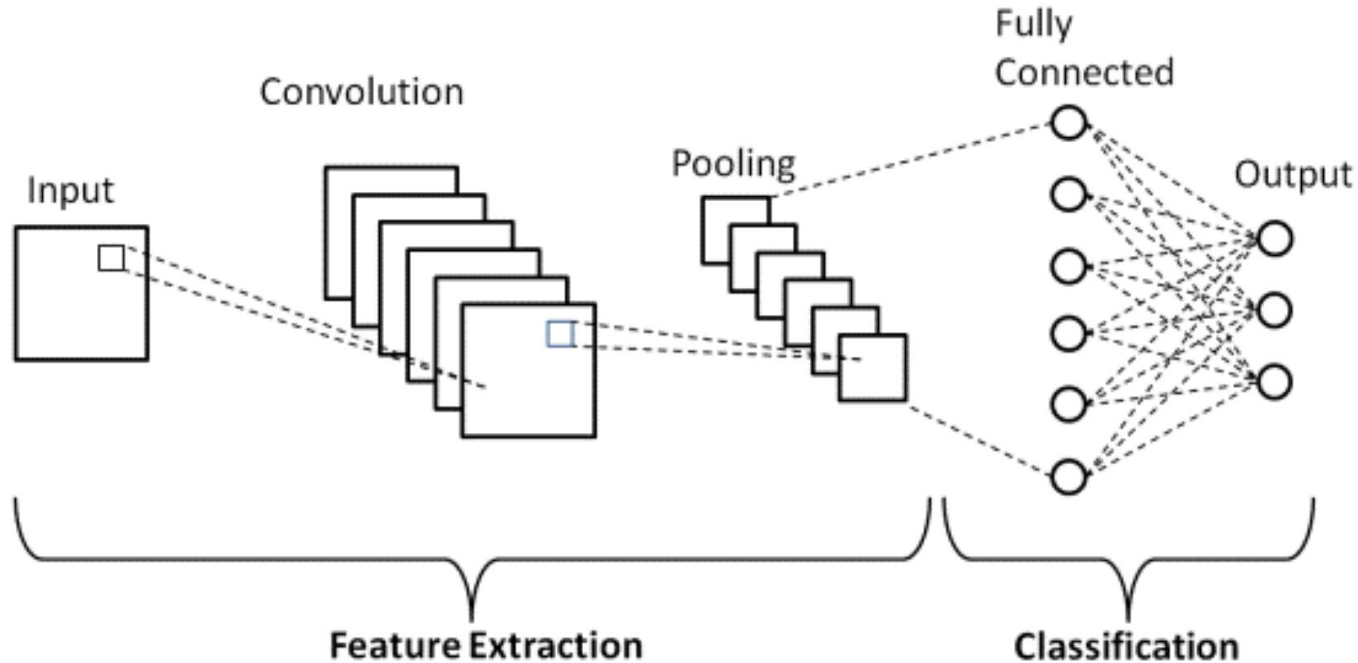derivative of the error w.r.t. weights in (L-1)

Formula for the weights update in the L-layer:

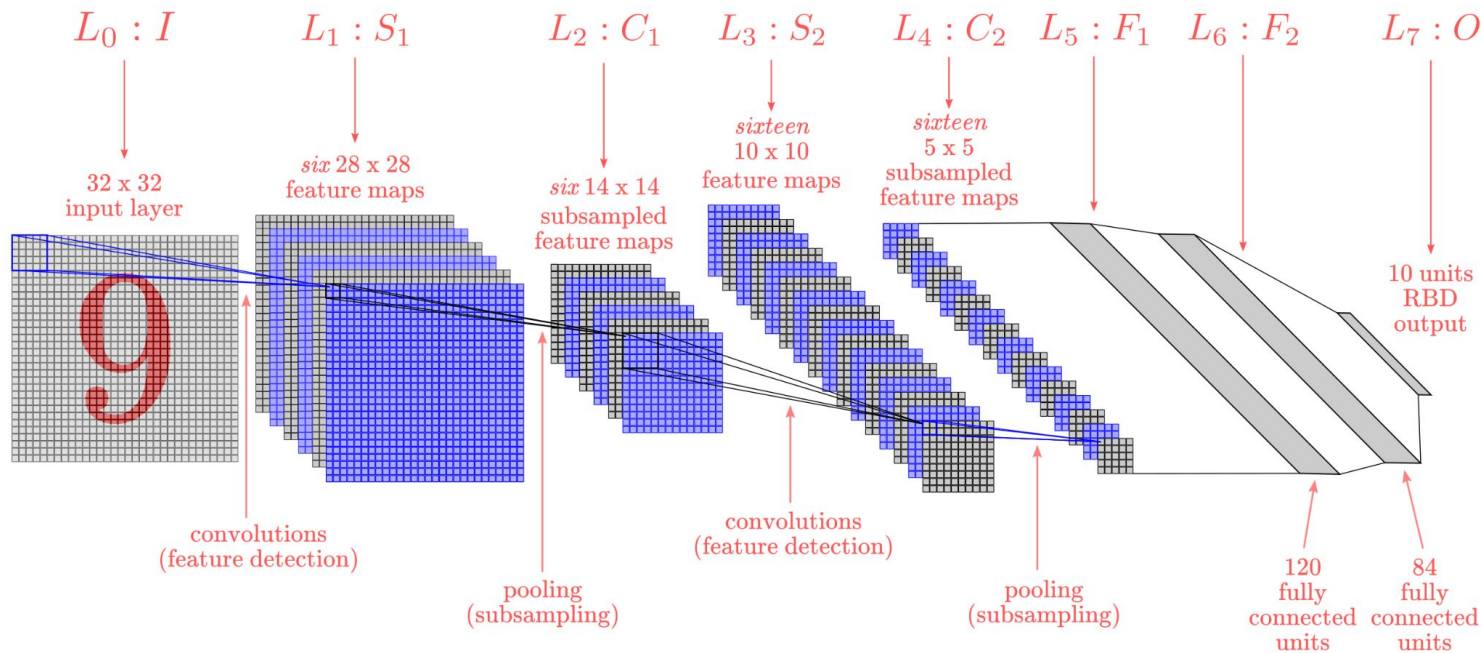$$w_{jk}^L = w_{jk}^L - \eta \times \frac{\partial E}{\partial w_{jk}^L}$$

Formula for the bias update in the L-layer:

$$b^{(L)} = b^{(L)} - \eta \times \frac{\partial E}{\partial b^{(L)}}$$

# Convolutional Neural Networks (CNN)

# LeNet-5



$$E(W) = \frac{1}{P} \sum_{p=1}^{X} (\hat{y}_{D^p}(X^p, W) + \log(e^{-j} + \sum_i e^{-\hat{y}(X^p, W)}))$$

# The Convolution step



$$F_{mn} = S(i,j) = (P * K)_{ij} = \sum_{m} \sum_{n} P_{i-m,j-n} * K_{m,n}$$

labels: convolution output (feature map at m,n position), convolution function, input function (matrix of pixel values), convolution operator, kernel function (matrix of weights), nested summation, sum over each row, sum over each col, col index for input, col index for input, row index for kernel, col index for kernel

Actually, several deep learning libraries like MXNet and Pytorch **DO NOT implement convolutions** but a closely related operation called **cross-correlation**

$$F_{mn} = S(i,j) = (P \star K)_{ij} = \sum_{m} \sum_{n} P_{i+m,j+n} \star K_{m,n}$$

# The convolution operation is simply a matrix multiplication

Let's take a look at basic element of CNN: convolution layer

Consider the case where we are applying (2,2) kernel



to a (3,3) matrix:

The convolution can be rewritten as

$$
\begin{bmatrix}
\alpha & \beta & 0 & \gamma & \delta & 0 & 0 & 0 & 0 \\
0 & \alpha & \beta & 0 & \gamma & \delta & 0 & 0 & 0 \\
0 & 0 & 0 & \alpha & \beta & 0 & \gamma & \delta & 0 \\
0 & 0 & 0 & 0 & \alpha & \beta & 0 & \gamma & \delta
\end{bmatrix}
*
\begin{bmatrix} A \\ B \\ C \\ D \\ E \\ F \\ G \\ H \\ J \end{bmatrix}
+
\begin{bmatrix} b \\ b \\ b \\ b \end{bmatrix}
=
\begin{bmatrix}
\alpha A + \beta B + 0C + \gamma D + \delta E + 0F + 0G + 0H + 0J + b \\
0A + \alpha B + \beta C + 0D + \gamma E + \delta F + 0G + 0H + 0J + b \\
0A + 0B + 0C + \alpha D + \beta E + 0F + \gamma G + \delta H + 0J + b \\
0A + 0B + 0C + 0D + \alpha E + \beta F + 0G + \gamma H + \delta J + b
\end{bmatrix}
=
\begin{bmatrix}
\alpha A + \beta B + \gamma D + \delta E + b \\
\alpha B + \beta C + \gamma E + \delta F + b \\
\alpha D + \beta E + \gamma G + \delta H + b \\
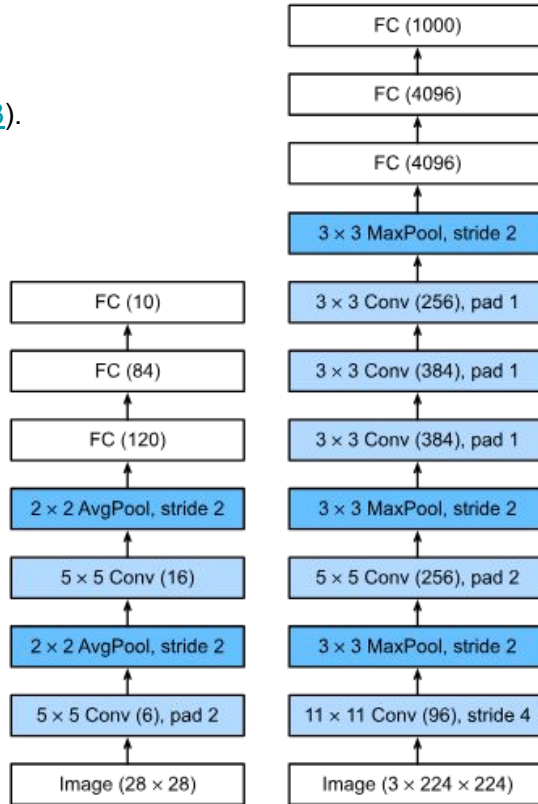\alpha E + \beta F + \gamma H + \delta J + b
\end{bmatrix}
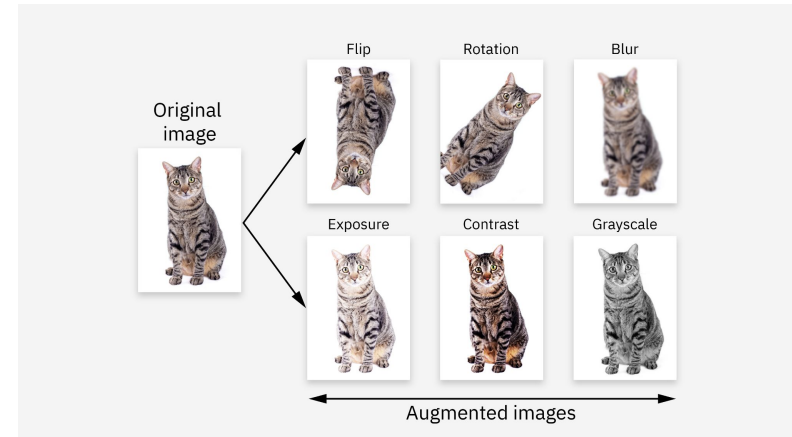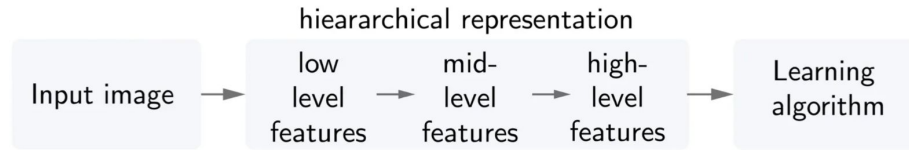=
\begin{bmatrix} P \\ Q \\ R \\ S \end{bmatrix}
$$

A  B  C  D  E  F  G  H  J

# AlexNet

(Russakovsky *et al.*, 2013).

FC (1000)

FC (4096)

FC (4096)

3 × 3 MaxPool, stride 2

| FC (10) | 3 × 3 Conv (256), pad 1 |

| FC (84) | 3 × 3 Conv (384), pad 1 |

| FC (120) | 3 × 3 Conv (384), pad 1 |

| 2 × 2 AvgPool, stride 2 | 3 × 3 MaxPool, stride 2 |

| 5 × 5 Conv (16) | 5 × 5 Conv (256), pad 2 |

| 2 × 2 AvgPool, stride 2 | 3 × 3 MaxPool, stride 2 |

| 5 × 5 Conv (6), pad 2 | 11 × 11 Conv (96), stride 4 |

| Image (28 × 28) | Image (3 × 224 × 224) |

The training procedure of AlexNet used for the first time data augmentation:



Original image

Flip     Rotation     Blur

Exposure     Contrast     Grayscale

Augmented images

# Hierarchical representation learning



Traditional pattern recognition
e.g., SIFT, HOG

hieararchical representation

Deep learning

# Limitations of CNNs

- they are sensible to ADVERSARIAL ATTACKS:



$$x \qquad + .007 \times \qquad \text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y)) \qquad = \qquad \begin{array}{c} \boldsymbol{x} + \\ \epsilon \text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y)) \end{array}$$

"panda"          "nematode"          "gibbon"
57.7% confidence     8.2% confidence     99.3 % confidence

- contain unrealistic features
- require heavy computation for the training

# GoogleNet (and the rise of inception blocks)

In 2014, *GoogLeNet* won the ImageNet
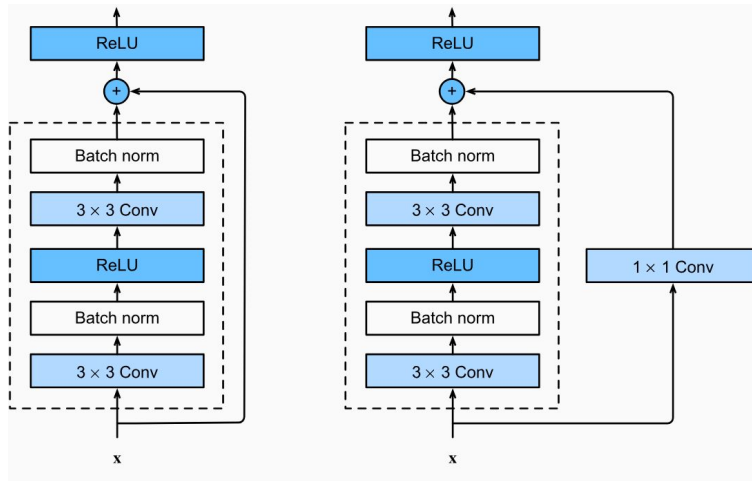Challenge (Szegedy *et al.*, 2015)
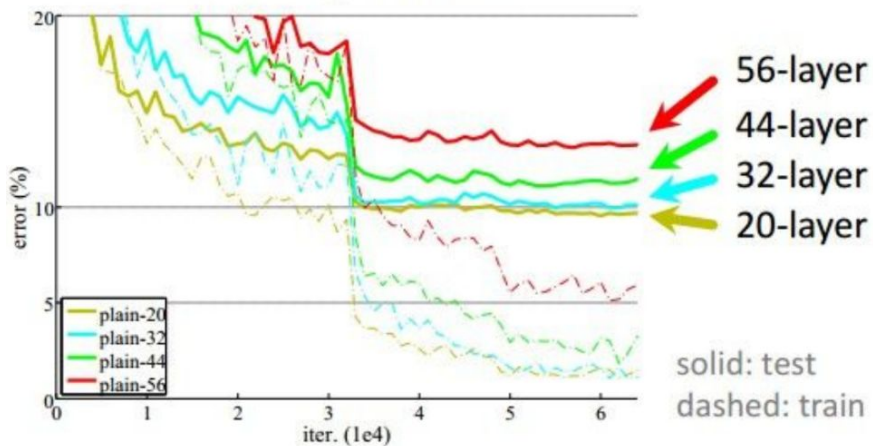


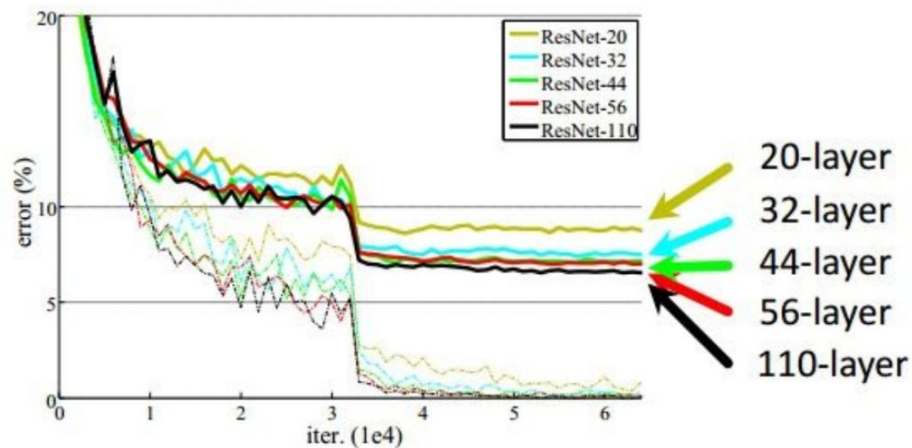inception
block

# ResNet
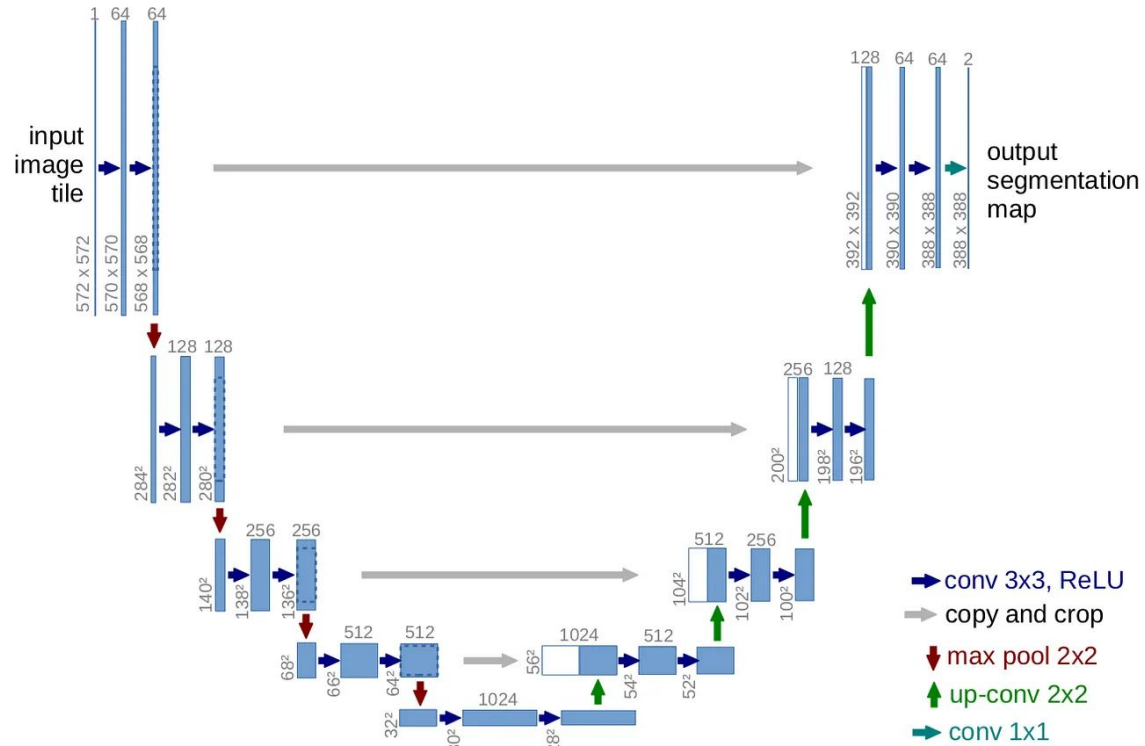
[He et al., 2015]


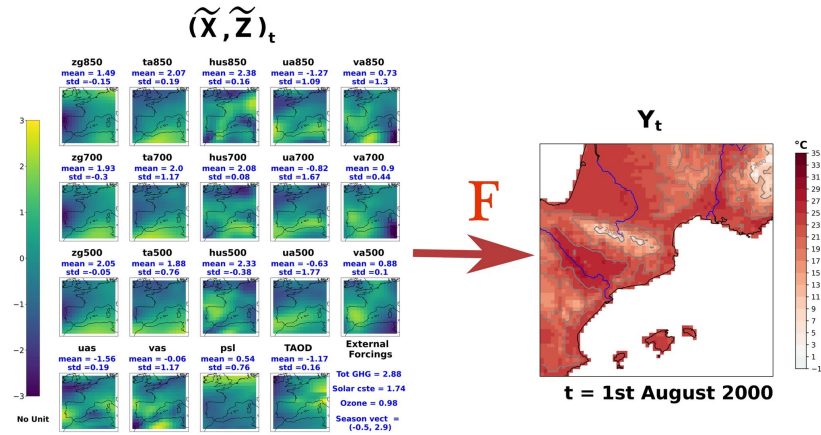
residual
block

# CIFAR experiments
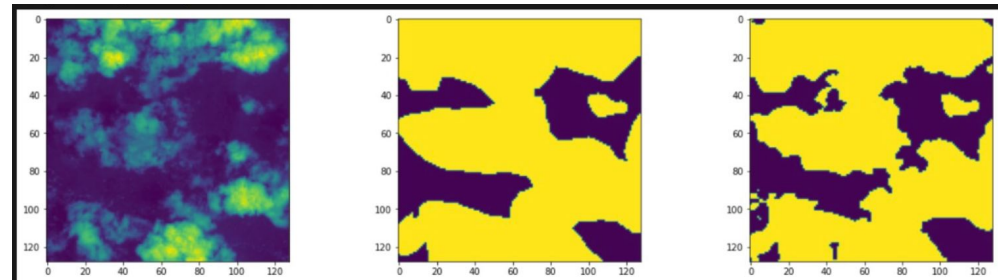
# UNet

# UNet applications in Climate studies

RCM emulator/downscaling



Doury et al. (2022-2024)

resol: 150 km -> 12.5 km

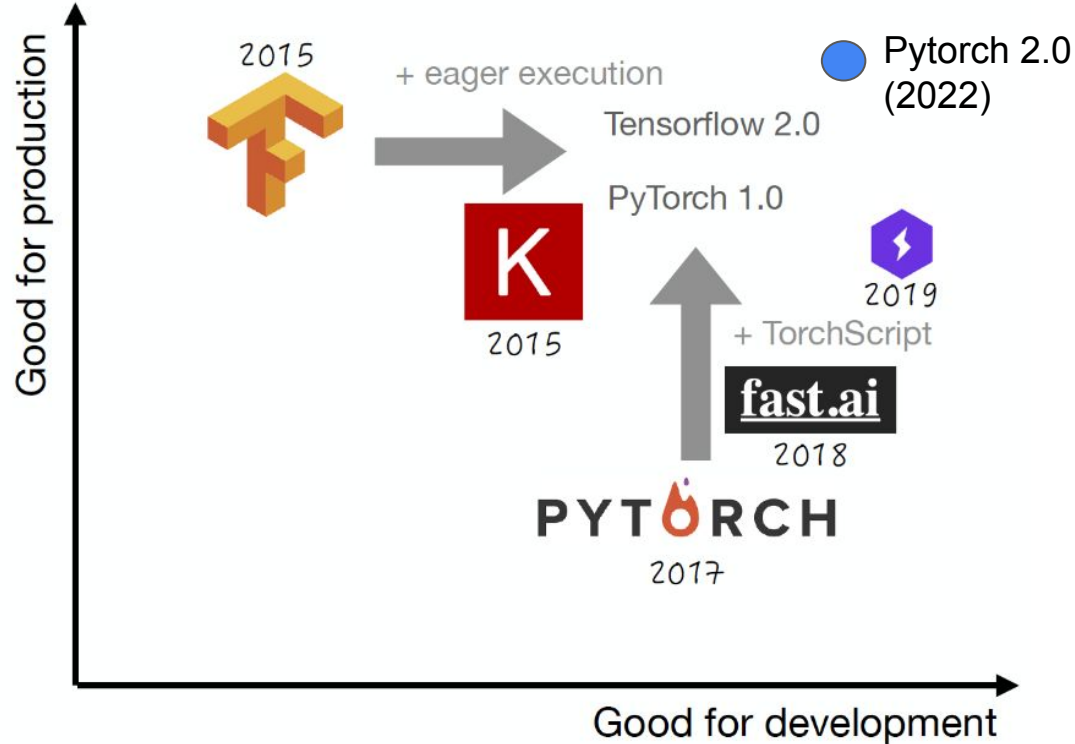semantic segmentation in satellite data (e.g. clouds)



De Souza 2023

# AI 4 climate



**Observation**
- Reconstruction
- Classification
- Super-resolution

**Key outcome:** optimally extended observational datasets (time, space and observable)

**Theory**
- Sub-grid parameterizations
- Reduced-order models

**Key outcome:** innovative approaches to capture sub-grid-scale phenomena

- Model validation
- Training datasets

- Hybrid models

- Initial value problems
- Data assimilation

- Boundary value problems
- Climate simulations

**Computation**
- Short-term weather forecasting
- Long-term climate prediction
- Model bias correction

**Key outcome:** breaking predictability barriers (such as weather and climate)
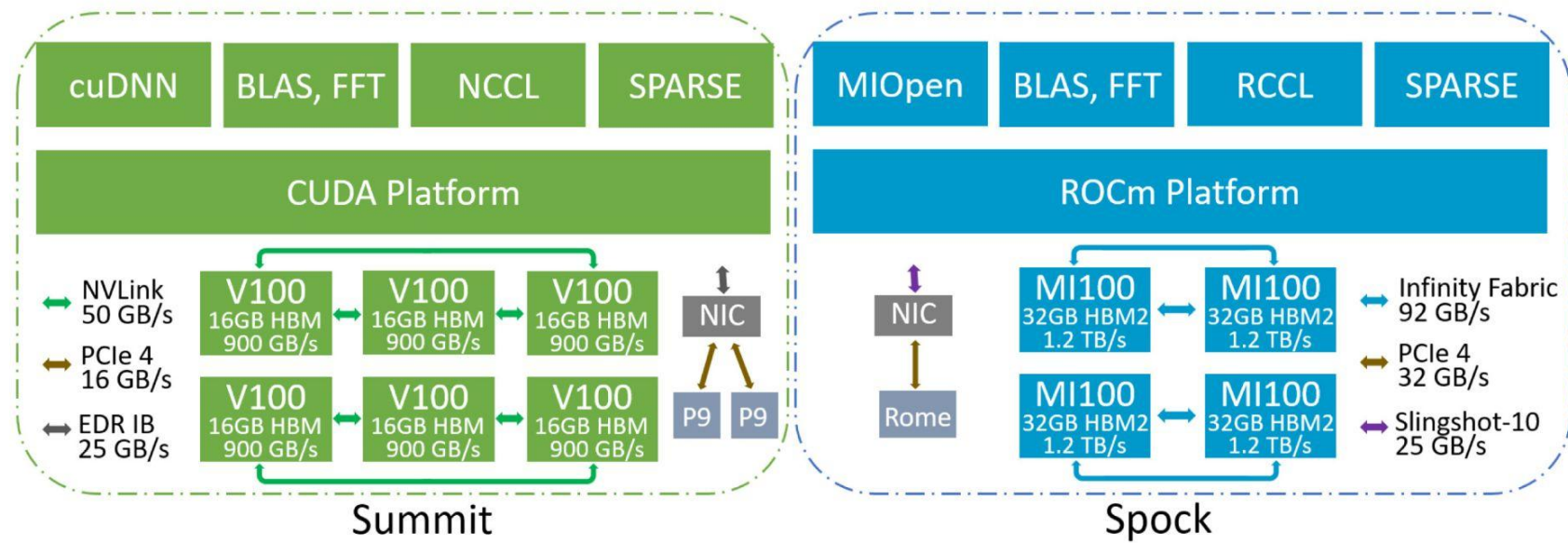
# Python most used libraries/frameworks for DL
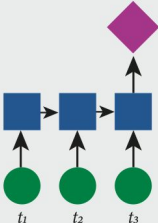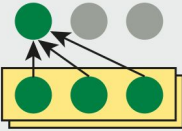
# What's behind Pytorch/Tensorflow?

# Diving into the world of DL models



| | Convolutional NN (CNN) | Recurrent NN (RNN) | Transformer | Autoencoder (AE) |
|---|---|---|---|---|
| **Goal** | Perform inference on data with local features | Perform inference on temporal data | Perform inference on sequential data | Embed high-dimensional data |
| **Key idea** | Learn shift-invariant filters | Learn temporal correlations via recurrent structure | Learn context based correlations via attention mechanism | Learn low-dimensional embedding of data |

| | Graph NN (GNN) | Generative Adversarial Network (GAN) |
|---|---|---|
| **Goal** | Capture graph based dependencies in the data | Generate samples from data distribution |
| **Key idea** | Perform message passing between nodes in a layer | Simultaneously train generator and discriminator |

**Denoising autoencoders (DAE)** are autoencoder models that learn low dimensional embeddings of noisy high dimensional data, i.e. inputs that differ by a small amount of noise give rise to a similar embedding vector.

**Attention mechanism** mimics cognitive attention by learning importance weights for the inputs based on the whole input context (e.g. in a task of translating codons to amino acids attention mechanism will learn to give higher weight to the first two nucleid acids). Attention is the key part of transformer models, but can also be applied in conjunction with other layer types.

**Convolutional layers** have dimension which indicates the dimension of learned filters. Thus, we can have a 1-dimensional convolutional layer for sequences, 2-dimensional layer for matrices, and so on.

**Graph convolutional network (GCN)** is a graph neural network with convolutional layers defined by the topology of the graph. Thus instead of passing neighboring sequence or matrix entries through a filter, graph defined neighborhoods are used.

# Why GPUs for DL?

1) **Neural networks are embarrassingly parallel algorithm**
2) **most of the operations performed in DL models can be rewritten as matrix multiplications**
3) **big datasets require to perform big matrix computation (**extremely slow on CPU with respect to GPU**)**
4) **well established libraries, with specific classes for ML objects (e.g. cuDNN, more recently tensorRT)**

**and we know that GPUs are very good in solving specific parallel tasks (e.g matrix multiplication) , thanks to**

- +1000 cores (>100K threads)
- **SIMD / SIMT**
- **high memory bandwidth**
- **newer GPUs have also tensor cores (particularly suited to tensor ops typical of NNs), and mixed precision**

**However, also GPUs have limitations:**

- ○ GPUs might not be as efficient for extreme sparse networks, due to the overhead of managing sparse data structures.
- ○ Some specialized sparse operations might not be as optimized as dense operations on GPUs.

# Different strategies for Multi-GPUs training

we can identify 5 different categories of parallelism

model parallelism
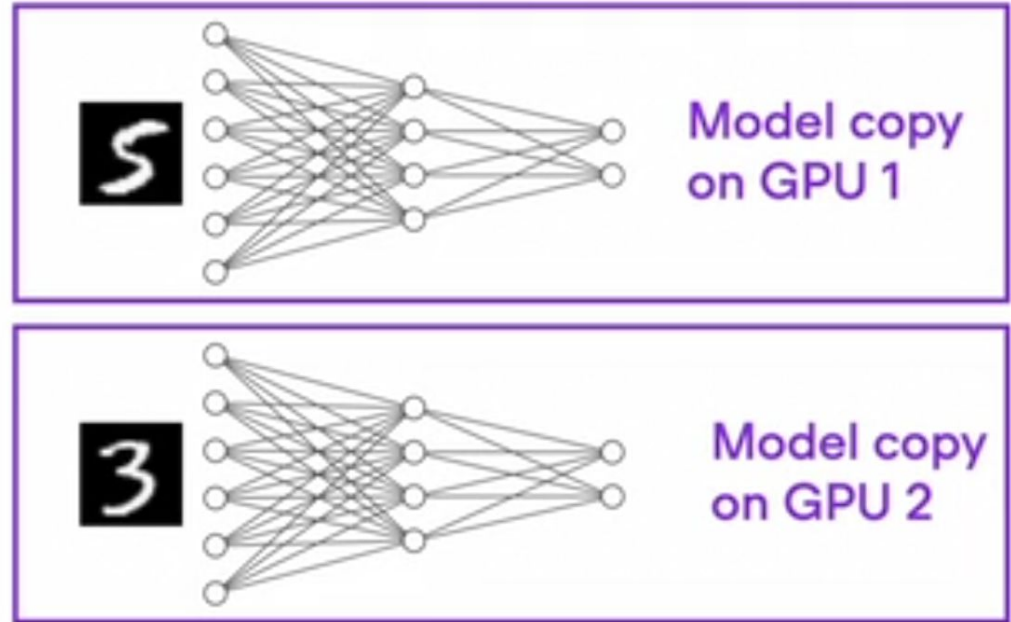
tensor parallelism

data parallelism

sequence parallelism
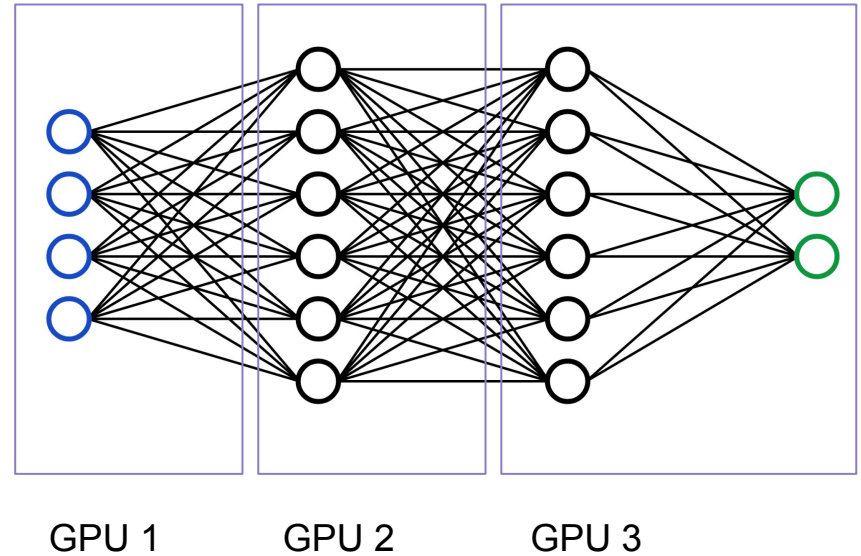
pipeline parallelism

# Data Parallelism

In this framework we split batches to train DL model into different GPUs



Model copy on GPU 1

Model copy on GPU 2

# Model parallelism

In this parallelism framework we choose to put different layers of the NN on different GPUs

to work around GPU memory limits



GPU 1          GPU 2          GPU 3
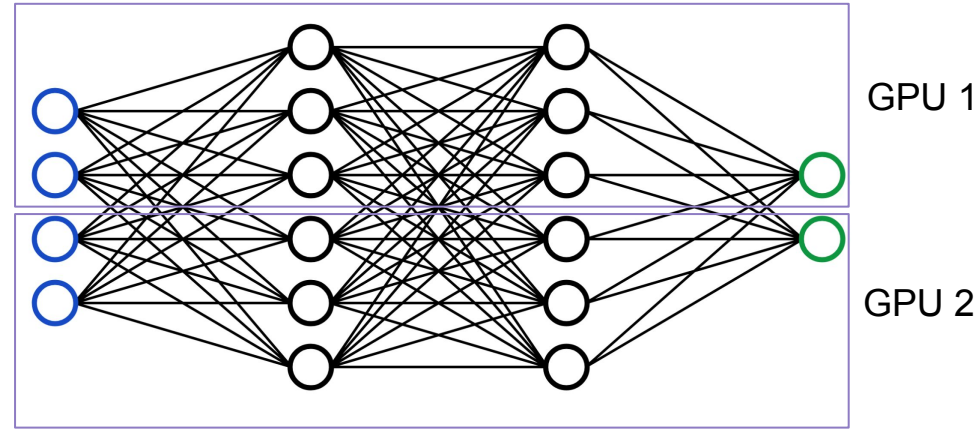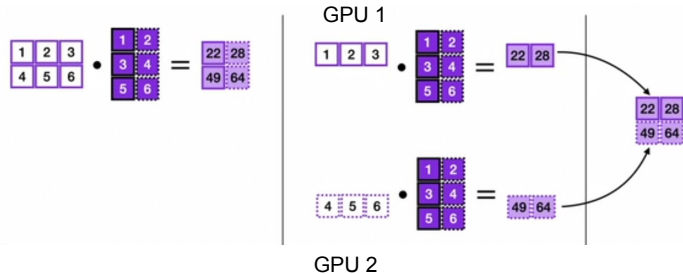
# Tensor parallelism

In this framework we split the tensor operation done at each layer among different GPUs
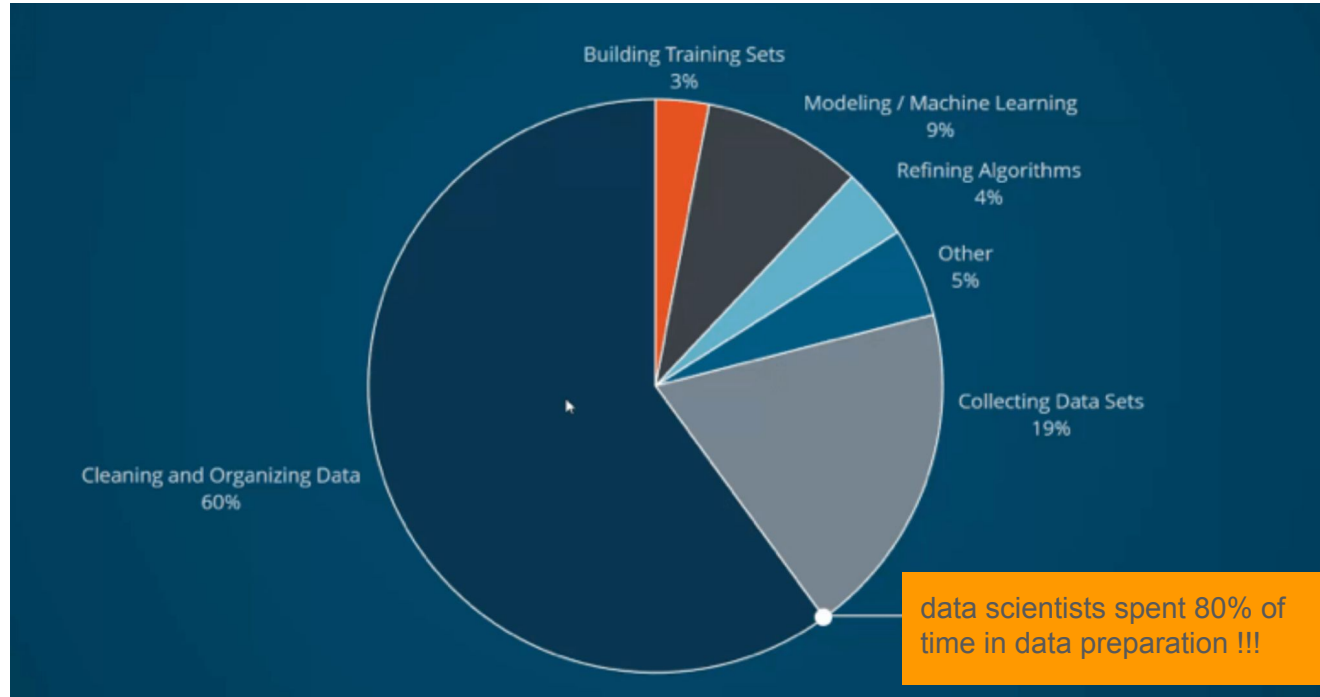
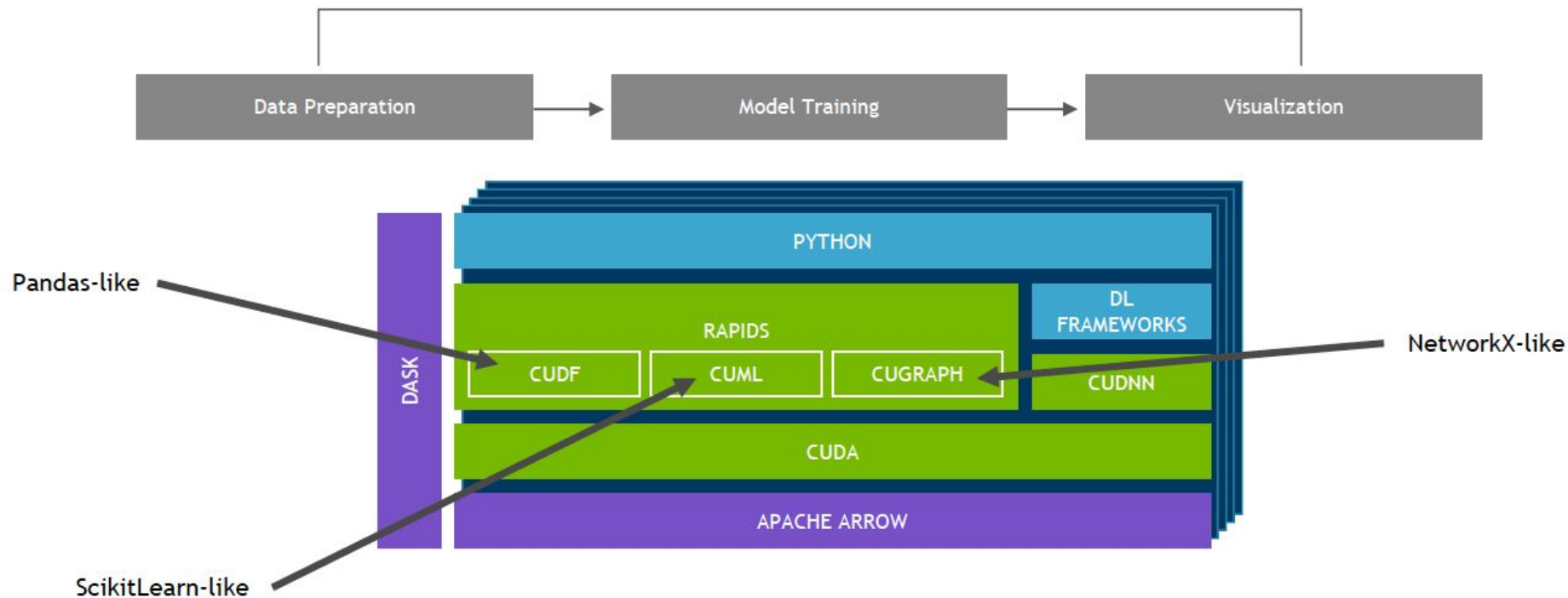similarly to what we would have done for matmul

# More complex strategies for DL training

Sequence parallelism and pipeline parallelism frameworks

are obtained combining the previous approaches,

and are typically applied to DL models dealing with
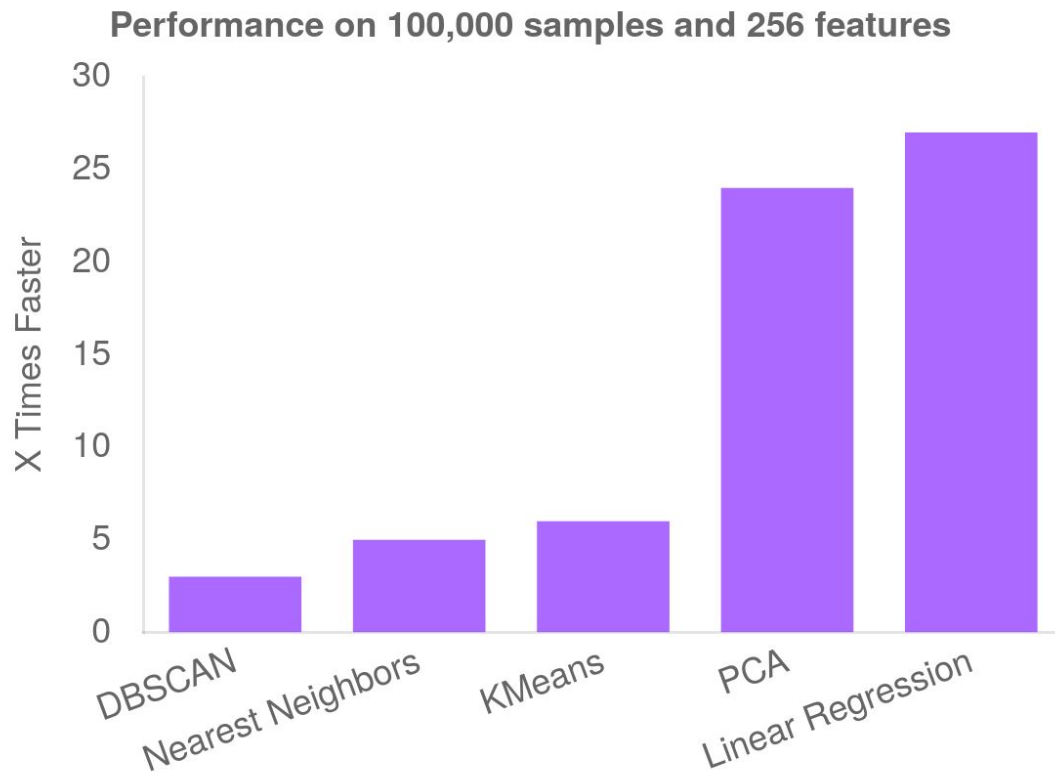
spatio-temporal data.

# Do we need GPUs also for other ML tasks?

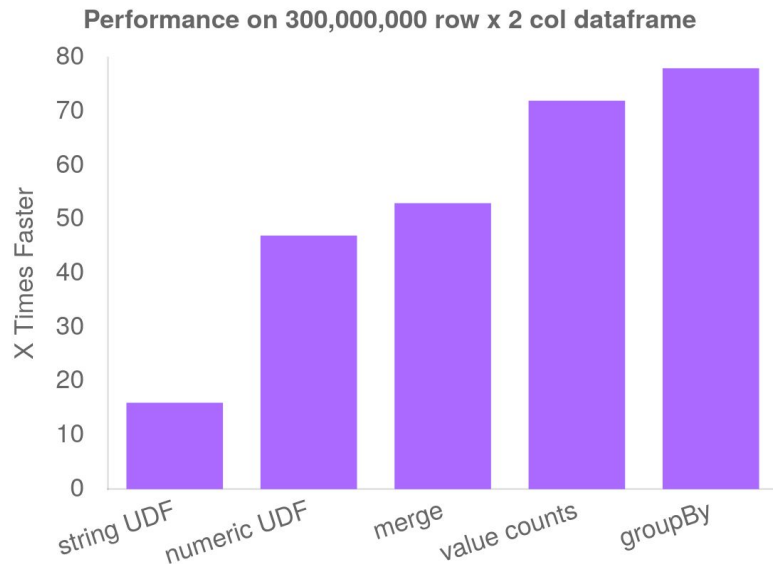# GPU-based libraries outside of Pytorch

# GPU for classical ML

**Performance on 100,000 samples and 256 features**



* Benchmark on AMD EPYC 7642 (using 1x 2.3GHz CPU core) w/ 512GB and NVIDIA A100 80GB (1x GPU) w/ scikit-learn v1.2 and cuML v23.02

# GPUs for data preprocessing



Performance on 300,000,000 row x 2 col dataframe

* Benchmark on AMD EPYC 7642 (using 1x 2.3GHz CPU core) w/
512GB and NVIDIA A100 80GB (1x GPU) w/ pandas v1.5 and cuDF
v23.02

# References

Bishop, C. M. (1995). Neural networks for pattern recognition. Oxford university press.

Rosenblatt, F. (1958). The perceptron: a probabilistic model for information storage and organization in the brain. Psychological review, 65(6), 386.

*Deep Learning: A Comprehensive Overview on Techniques, Taxonomy, Applications and Research Directions*

Junqi Yin et. al, 2021, *Comparative evaluation of deep learning workloads for leadership-class systems*

NVIDIA Booklet on GPU development, 2021