

# ML em HEP II

Escola INCT 2024

A. Sznajder

UERJ  
Instituto de Fisica

Novembro - 2024

## Outline

- 1 What is Machine Learning ?
- 2 Neural Networks
- 3 Deep Learning Revolution
- 4 Deep Architectures and Applications
  - Convolutional Networks (CNN)
  - Recurrent Networks (RNN)
  - Attention Mechanism and Transformers(TN)
  - Graph Neural Networks (GNN)
  - Unsupervised Learning and Autoencoders (AE,CAE,DAE)
  - Generative Networks And Density Estimation (VAE, GAN, NF, DM)

## Support Material

- **Textbook:**  
Aurélien Géron - Hands-On Machine Learning with Scikit-Learn, Keras, and Tensorflow
- **Jupyter Notebooks :**  
<https://github.com/INCT-CERN-Brasil/Escola2024>
- **Google Colab :**  
<https://colab.google>

# What is Machine Learning ?

## Machine Learning (ML)

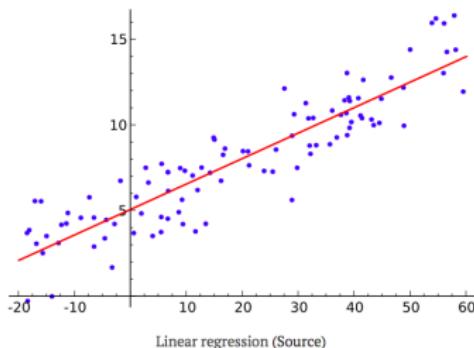
Machine learning (ML) is the study of computer algorithms capable of building a predictive model out of a data samples (**learn from examples**), without being explicitly programmed (**sequence of instructions**).



## Centuries Old Machine Learning <sup>1</sup>

Take some points on a 2D graph, and fit a function to them. What you have just done is generalized from a few  $(x, y)$  pairs (examples), to a general function that can map any input  $x$  to an output  $y$

### The Centuries Old Machine Learning Algorithm



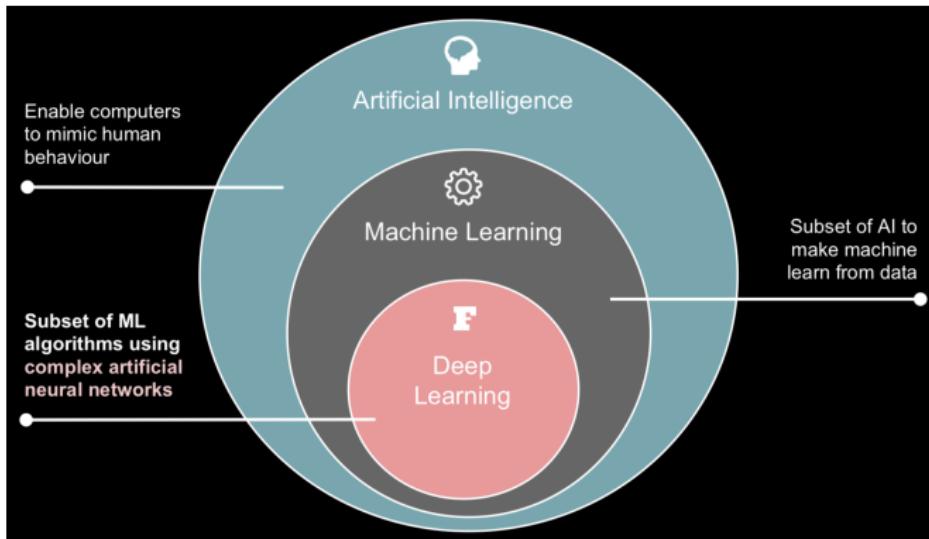
Linear regression (Source)

Linear regression is a bit too wimpy to solve the complex problems dealing with image, speech or text, but what it does is essentially what supervised ML is all about (**glorified regression**) !

<sup>1</sup><http://www.andreykurenkov.com/writing/ai/a-brief-history-of-neural-nets-and-deep-learning/>

# Artificial Intelligence

"Intelligence can be understood as the ability to process current information to inform future decisions"



None of current ML systems we have are real AI and the brain learns so efficiently that no ML method can match it<sup>2</sup>

<sup>2</sup>Yann LeCun , Epistemology of Deep Learning : <https://www.youtube.com/watch?v=gG5NCkMerHU&t=944s> , <https://www.youtube.com/watch?v=cWzi38-vDbE&t=768s>

# Artificial Intelligence

None of the systems we have nowadays are real AI. The brain learns so efficiently that no ML method can match it !

## AI versus Brain

- Brain has  $10^{14}$  parameters and we live only  $10^9$ s (lot more parameters than data)
- So, it must do lots of unsupervised and self-supervised learning and must predict what we observe !
- Supervised and reinforcement learning requires millions of examples(trials)
- Still missing a learning paradigm that builds predictive models through observation and action

Yann LeCun on the Epistemology of Deep Learning:

<https://www.youtube.com/watch?v=gG5NCkMerHU&t=944s>

<https://www.youtube.com/watch?v=cWzi38-vDbE&t=768s>

# Introduction to Machine Learning

Machine Learning(ML) can be approached from many different angles:

## 1) Tasks:

- Classification
- Regression
- Data Generation or Simulation

## 2) Data Structure

- Tabular
- Image
- Sequence
- Graph
- Sets

## 3) Learning Paradigms

- Supervised Learning (labeled)
- Un(Self)supervised Learning (unlabeled)
- Reinforcement Learning (reward)

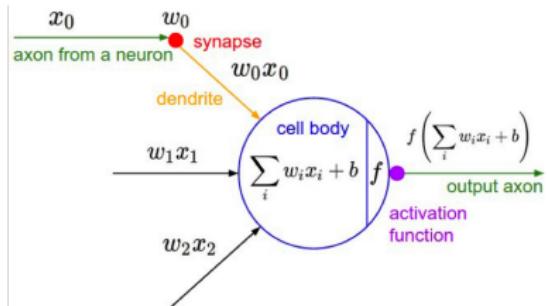
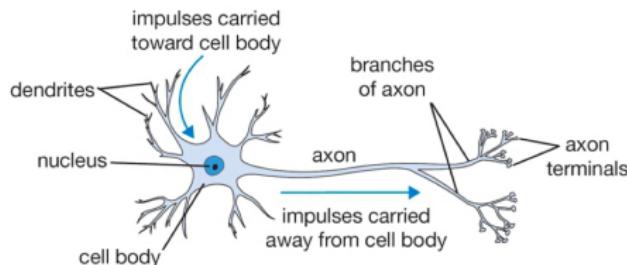
## 3) Neural Networks Architectures

- Multilayer Perceptron (MLP)
- Recurrent Networks (RNN,LSTM,GRU)
- Transformer Networks (TN)
- Convolutional Networks (CNN)
- Graph Networks (GCN,GAT,DGCN,IN)
- Point Nets and Deep Sets (PN,DS)
- Autoencoders (AE,DAE,VAE)
- Generative Adversarial Network (GAN)
- Normalizing Flows (NF)
- Diffusion Models (DM,SDM)

Obs: ML can be implemented by different algorithms (ex: SVM, BDT, PCA ...) but we will discuss only Neural Networks

# Neural Networks

Artificial Neural Networks (NN) are computational models vaguely inspired<sup>3</sup> by biological neural networks. A NN is formed by a network of basic elements called neurons, which receive an input, change their state according to the input and produce an output



Original goal of NN approach was to solve problems like a human brain. However, focus moved to performing specific tasks, deviating from biology. Nowadays NN are used on a variety of tasks: image and speech recognition, translation, filtering, playing games, medical diagnosis, autonomous vehicles, ...

<sup>3</sup> Design of airplanes was inspired by birds, but airplanes don't flap wings to fly !

# Artificial Neuron

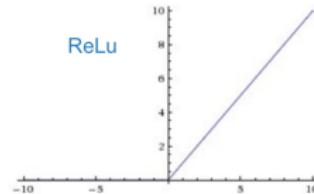
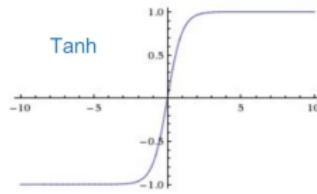
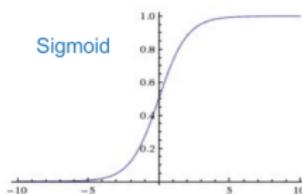
## Artificial Neuron Model

Each node of a NN receives inputs  $\vec{x} = \{x_1, \dots, x_n\}$  from other nodes or an external source and computes an output  $y$  according to the expression

$$y = F \left( \sum_{i=1}^n W_i x_i + B \right) = F(\vec{W} \cdot \vec{x} + B) \quad (1)$$

, where  $W_i$  are connection weights,  $B$  is the threshold and  $F$  the activation function <sup>4</sup>

There are a variety of possible activation function and the most common ones are



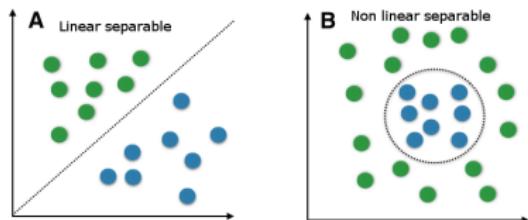
<sup>4</sup> Nonlinear activation is fundamental for nonlinear decision boundaries

# The Perceptron

The perceptron<sup>6</sup> algorithm is a binary linear classifier invented in 1957 by F.Rosenblatt. It's formed by a single neuron that takes input  $\vec{x} = (x_1, \dots, x_n)$  and outputs  $y = 0, 1$  according to

## Perceptron Model<sup>7</sup>

$$y = \begin{cases} 1, & \text{if } (\vec{W} \cdot \vec{x} + B) > 0 \\ 0, & \text{otherwise} \end{cases}$$



To simplify notation define  $W_0 = B$ ,  $\vec{x} = (1, x_1, \dots, x_n)$  and call  $\theta$  the Heaviside step function

## Perceptron Learning Algorithm

- ① Calculate the output error:  $y_j = \theta(\vec{W} \cdot \vec{x}_j)$  and  $Error = 1/m \sum_{j=1}^m |y_j - t_j|$
- ② Modify(update) the weights to minimize the error:  $\delta W_i = r \cdot (y_j - t_j) \cdot X_i$ , where  $r$  is the learning rate
- ③ Return to step 1 until output error is acceptable

<sup>6</sup><https://medium.com/towards-data-science/perceptrons-the-first-neural-network-model-8b3ee>

<sup>7</sup>Equation of a plane in  $\mathbb{R}^n$  is  $\vec{W} \cdot \vec{x} + B = 0$

# The Perceptron as a Universal Function Approximator

## Universal Approximation Theorem

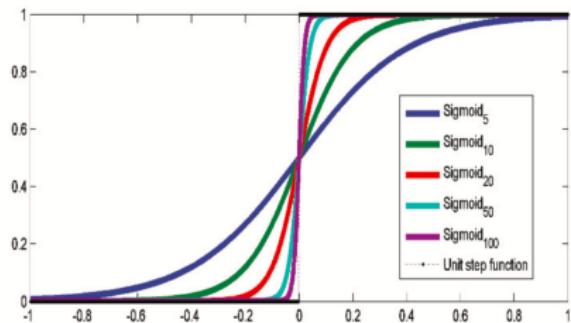
A single hidden layer feed forward neural network with a linear output unit can approximate any continuous function arbitrarily well, given enough hidden neurons<sup>9</sup>

The theorem doesn't tell us how many neurons or how much data is needed !

## Sigmoid → Step Function

For large weight  $W$  the sigmoid turns into a step function, while  $B$  gives its offset

$$y = \frac{1}{1 + e^{-(w \cdot x + b)}} \quad (2)$$



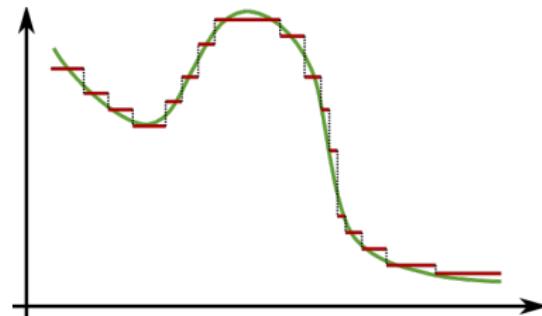
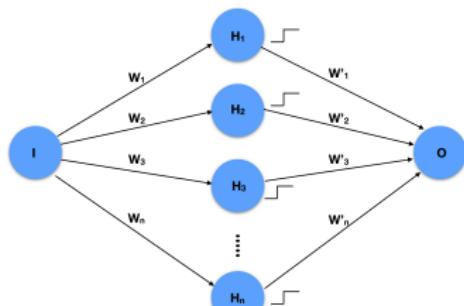
<sup>9</sup>Cybenko,G.(1989) Approximations by superpositions of sigmoidal functions, Math.ofCtrl.,Sig.,andSyst.,2(4),303  
Hornik,K.(1991) Approximation Capabilities of Multilayer Feedforward Networks, Neural Networks, 4(2), 251

# The Perceptron as a Universal Function Approximator

## Approximating $F(x)$ by Sum of Steps

A continuous function can be approximated by a finite sum of step functions. The larger the number of steps(nodes), the better the approximation.<sup>10</sup>

Consider a NN composed of a single input ,  $n$  hidden nodes and a single output. Tune the weights such that the activations approximate steps functions with appropriate threshold and add them together !



One can always tune weights such that any activation behaves approximately as a step function !

<sup>10</sup> M.Nielsen , <http://neuralnetworksanddeeplearning.com/chap4.html>  
 M.Cranmer, <https://www.youtube.com/watch?v=fk2r8y5TfNY>  
 S.Prince, Understanding Deep Learning 2023

# Multilayer Perceptron (MLP)

The Multilayer Perceptron(MLP) is a fully connected NN with at least 1 hidden layer and nonlinear activation function  $F$ <sup>11</sup>. It is the simplest deep NN with a compositional "inductive bias".<sup>12</sup>

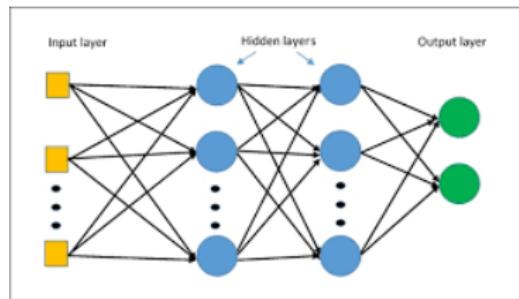
## Multilayer Perceptron Model

For a MLP with inputs nodes  $\vec{x}^{(0)}$ , one hidden layer of nodes  $\vec{x}^{(1)}$  and output layer of nodes  $\vec{x}^{(2)}$ , we have

$$\begin{cases} \vec{x}^{(1)} = \vec{F}^{(1)} (\vec{W}^{(1)}. \vec{x}^{(0)}) \\ \vec{x}^{(2)} = \vec{F}^{(2)} (\vec{W}^{(2)}. \vec{x}^{(1)}) \end{cases}$$

Eliminating the hidden layer variables  $\vec{H}$  we get

$$\Rightarrow \vec{x}^{(2)} = \vec{F}^{(2)} (\vec{W}^{(2)}. \vec{F}^{(1)} (\vec{W}^{(1)}. \vec{x}^{(0)})) \quad (3)$$



A MLP can be seen as a parametrized composite mapping  $F_{w,b} : \mathbb{R}^n \rightarrow \mathbb{R}^m$

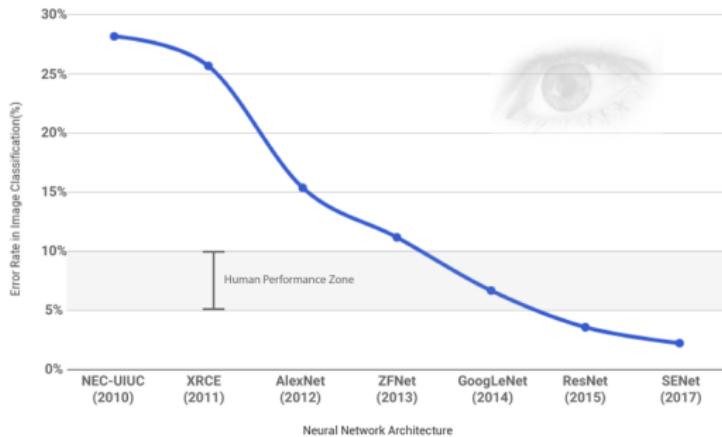
<sup>11</sup> A MLP with  $m$  layers using linear activation functions can be reduced to a single layer !

<sup>12</sup> The thresholds  $\vec{B}$  are represented as weights by redefining  $\vec{W} = (B, W_1, \dots, W_n)$  and  $\vec{x} = (1, x_1, \dots, x_n)$  (bias is equivalent to a weight on an extra input of activation=1 )

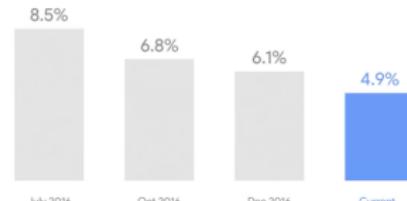
# Why Deep Learning ? and Why Now ?

## Why Deep Learning ?

### Image and Speech Recognition performance ( DNN versus Humans )



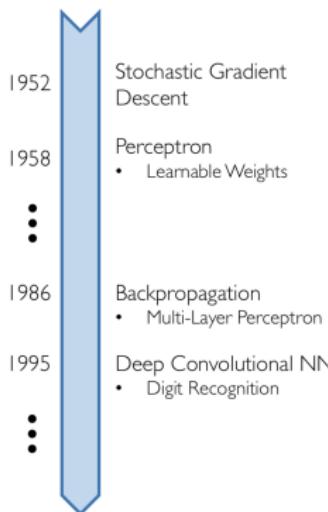
Speech Recognition  
Word Error Rate



<https://arxiv.org/pdf/1409.0575.pdf>

# Why Now ?

Neural networks date back decades , so why the current resurgence ?



The main catalysts for the current Deep Learning revolution have been:

- **Software:** TensorFlow, PyTorch, Keras and Scikit-Learn
- **Hardware:** GPU, TPU and FPGA
- **Large Datasets:** MNIST

# Training Datasets

Large and new open source datasets for machine learning research 26



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



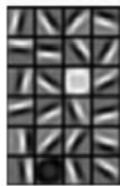
26 [https://en.wikipedia.org/wiki/List\\_of\\_datasets\\_for\\_machine\\_learning\\_research](https://en.wikipedia.org/wiki/List_of_datasets_for_machine_learning_research)

# Deep Learning - Need for Depth

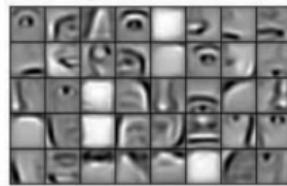
## Deep Neural Networks(DNN)

Depth allows the NN to factorize the data features, distributing its representation hierarchically across the layers

Layer 1



Layer 2



Layer 3



Edges

Object Parts

Object Models

⇒ DNN explores the compositional character of nature as an inductive bias !!!<sup>27</sup>

<sup>27</sup> Deep Learning , Y.LeCunn, J.Bengio, G.Hinton , Nature , vol. 521, pg. 436 , May 2015

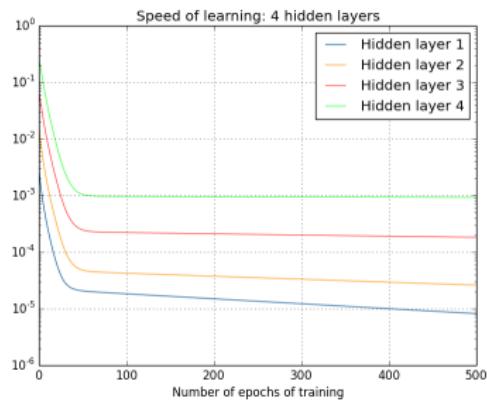
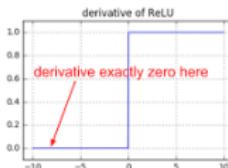
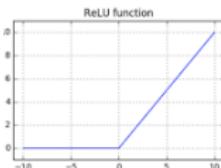
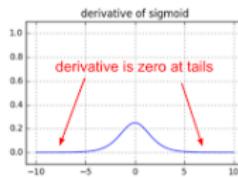
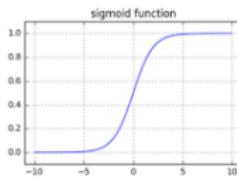
Deep Learning - Vanishing Gradient Problem

## Vanishing Gradient Problem<sup>28</sup>

Backpropagation computes gradients iteratively by multiplying the activation function derivate  $F'$  through  $n$  layers.

$$\delta_k^{(I)} = \left( \sum_m \delta_m^{(I+1)} W_{mk}^{(I+1)} \right) F'(z_k^{(I)})$$

For  $\sigma(x)$  and  $\text{Tanh}(x)$  the derivate  $F'$  is asymptotically zero, so weights updates gets vanishing small when backpropagated  $\Rightarrow$  Then, earlier layers learns much slower than later layers !!!

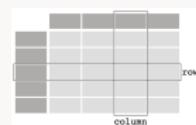


# Deep Architectures and Applications

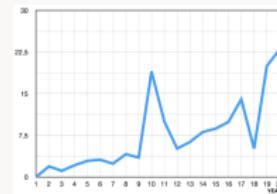
# Data Structures

Data comes structured in a variety of formats.

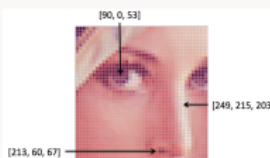
## Tabular Data



## Sequential Data



## Image Data



## Graph & Point Sets Data



# Image Structured Data

Image Representation

A computer sees an image as an array of numbers. The image is associated to a matrix containing numbers between 0 and 255, corresponding to each pixel brightness.



The image color **RGB** can be represented by expanding the depth of the representation, where each matrix represent a fundamental color intensity

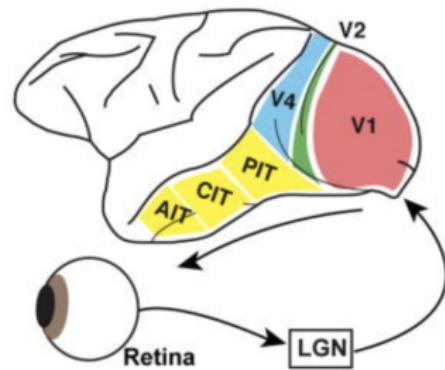
	165	167	265	96	
	74	52	211	207	165
50	144	170	217	41	140
57	100	32	241	23	106
206	118	124	27	55	20
210	206	105	169	65	200
25	170	190	197	4	11
115	104	74	111	85	124
32	69	728	203	74	

# Convolutional Neural Network

Full connectivity between neurons in a MLP makes it computationally too expensive to deal with high resolution images. A  $1000 \times 1000$  pixels image leads to  $O(10^6)$  weights per neuron !

## Convolutional Neural Network (CNN)

Inspired by the visual cortex<sup>30</sup>, where neurons respond to stimuli only in a restricted region of the visual field, a CNN mitigates the challenges of high dimensional inputs by restricting the connections between the input and hidden neurons. It connects only a small contiguous region, exploiting local image features.



A CNN uses inductive biases (local connectivity) built into the network layers to reduce de number of learnable parameters

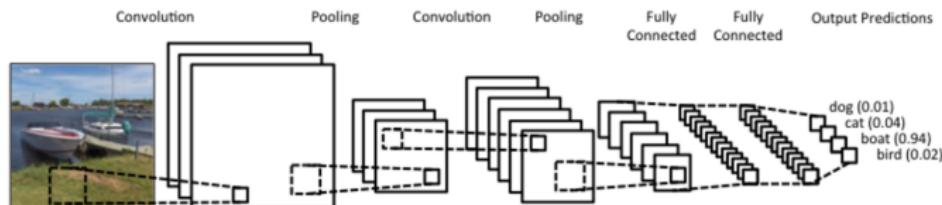
<sup>30</sup> How does the brain solve visual object recognition? , <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3306444>

# Convolutional Neural Network

## Convolutional Neural Network (CNN)

A typical CNN architecture is composed by a stack of distinct and specialized layers:

- ① Convolutional ( extract image feature maps )<sup>31</sup>
- ② Pooling (downsampling to reduce size)
- ③ Fully connected (MLP for image classification )



<sup>31</sup> <http://ufldl.stanford.edu/tutorial/supervised/FeatureExtractionUsingConvolution>

# CNN - Convolutional Layer

## Convolutional Layer

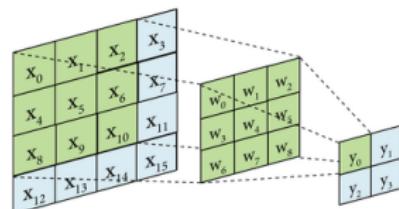
The convolutional layer<sup>32</sup> is the CNN core building block. A convolution can be seen as a sliding window transformation that applies a filter to extract local image features.

(Click on the figure)

## Discrete Convolution

The convolution has a set of learnable filters weights that are shared across the image

$$y_{ij}^{(l+1)} = \sum_{a=0}^{n-1} \sum_{b=0}^{m-1} w_{ab}^{(l)} x_{(i+a)(j+b)}^{(l)}$$



The same set of weights of a given filter are applied all across the image, reducing the number of learnable parameters (**shared weights**)

<sup>32</sup><http://ufldl.stanford.edu/tutorial/supervised/FeatureExtractionUsingConvolution>

<https://machinelearningmastery.com/padding-and-stride-for-convolutional-neural-networks>

## Convolutional Filters

A convolution filter can apply an effect ( sharpen, blurr ), as well as extract features ( edges, texture ) <sup>33</sup>

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & -1 & 5 & -1 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$



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



$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$



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



<sup>33</sup><https://docs.gimp.org/2.6/en/plug-in-convmatrix.html>

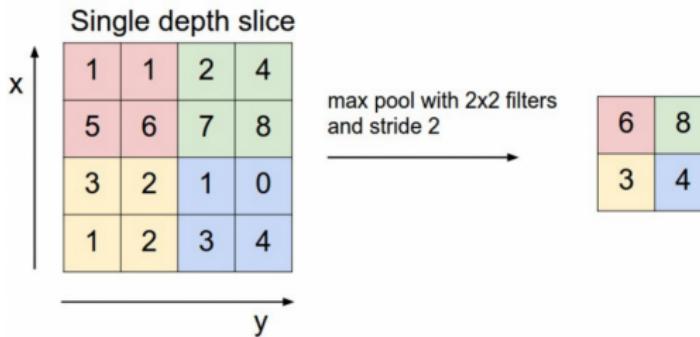
# CNN - Pooling Layer

## Pooling(Downsampling) Layer

The pooling (downsampling) layer<sup>34</sup> has no learning capabilities and serves a dual purpose:

- Decrease the representation size  $\Rightarrow$  reduce computation
- Make the representation approximately invariant to small input translations and rotations

Pooling layer partitions the input in non-overlapping regions and, for each sub-region, it outputs a single value (ex: max pooling, mean pooling)



<sup>34</sup><http://ufldl.stanford.edu/tutorial/supervised/Pooling>

# Typical CNN Architecture

## Typical Architecture

- CNN usually have pyramidal shape, decreasing in spatial extent and increasing in feature space
- Multiple layers of convolutional filters extract more-and-more abstract features

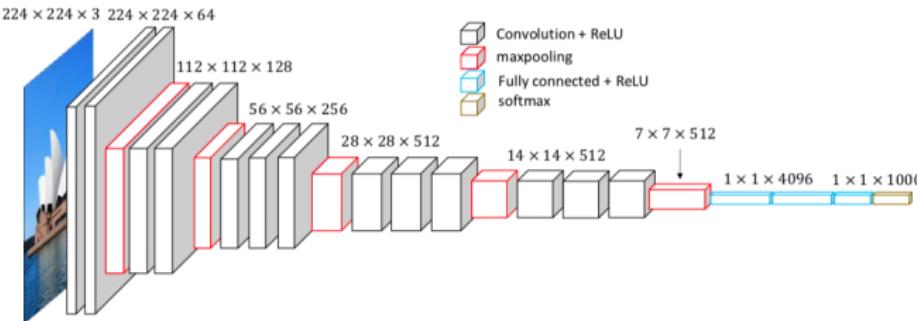
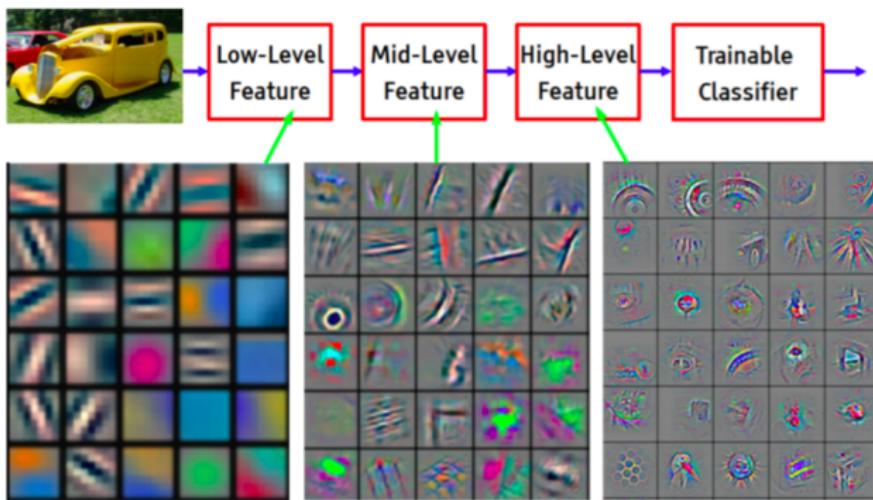


Figure 2: The architecture of VGG16 model .

## CNN Feature Visualization

Each CNN layer is responsible for capturing a different level of features as can be seen from ImagiNet<sup>35</sup>



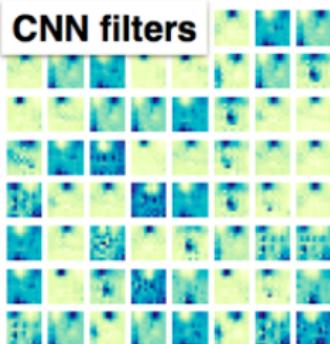
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

<sup>35</sup><https://arxiv.org/pdf/1311.2901.pdf>

CNN Application in HEP: Jet ID

# Jet images with convolutional nets

## CNN filters



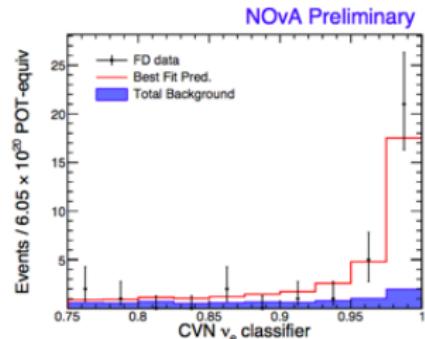
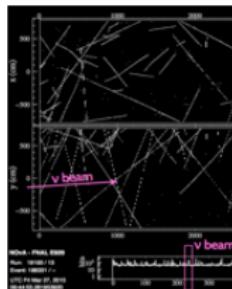
most activating images

L. de Oliveira et al., 2015



## CNN Application HEP: Neutrino ID

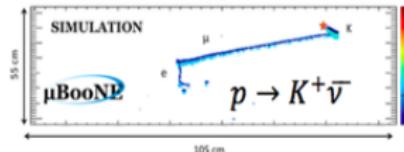
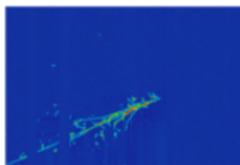
# Neutrinos with convolutional nets



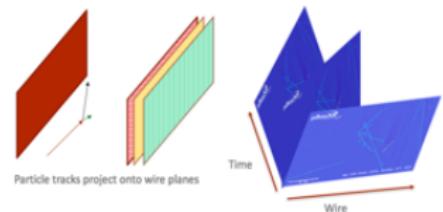
76% Purit  
73% Effici

An equivalent increased exposure of 30%

Aurisiano et al. 2016



**$\mu$ BooNE**



# Sequential Data

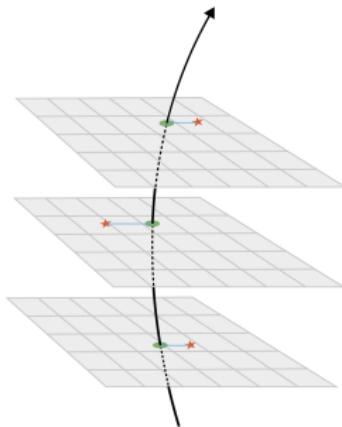
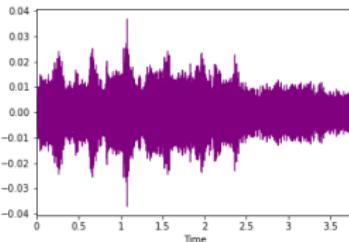
## Sequential Data

Sequential data is an interdependent data stream where data ordering contains relevant information

The food was good, not bad at all.

VS.

The food was bad, not good at all.



⇒ brain memorizes sequences for alphabet, words, phone numbers and not just symbols !

## Recurrent Neural Network(RNN)

Feed forward networks can't learn correlation between previous and current input !

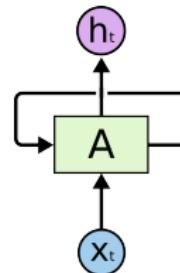
### Recurrent Neural Networks (RNN)

RNNs are networks that use feedback loops to process sequential data<sup>36</sup>. These feedbacks allows information to persist, which is an effect often described as **memory**.

#### RNN Cell ( Neuron <sup>37</sup> )

The hidden state depends not only on the current input, but also on the entire history of past inputs

- ① **Hidden State:**  $h^{[t]} = F(W_{xh}x^{[t]} + W_{hh}h^{[t-1]})$
- ② **Output:**  $y^{[t]} = W_{hy}h^{[t]}$



<sup>36</sup><https://eli.thegreenplace.net/2018/understanding-how-to-implement-a-character-based-rnn/>

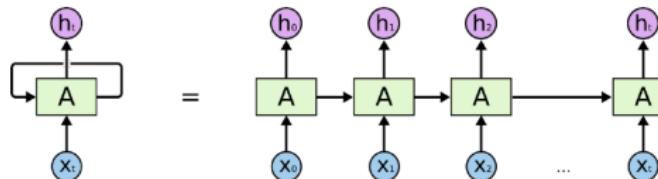
<sup>37</sup><https://www.analyticsvidhya.com/blog/2017/12/introduction-to-recurrent-neural-networks/>

## Recurrent Neural Network(RNN)

RNNs process each input sequence element at a time, maintaining in their hidden units a 'state vector' that contains information about all the past elements history.<sup>38</sup>

### RNN Unrolling

A RNN can be thought of as multiple copies of the same network, each passing a message to a successor. Unrolling is a visualization tool which views a RNN as a sequence of unit cells.



### Backpropagation Through Time (BPTT)

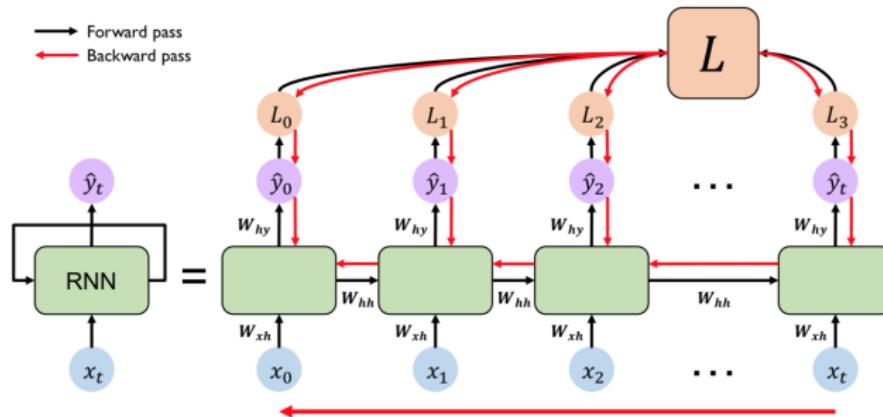
Backpropagation through time is just a fancy buzz word for backpropagation on an unrolled RNN

Unrolled RNN can lead to very deep networks  $\Rightarrow$  bias to capture only short term dependencies !

<sup>38</sup> <http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-hyperparameters-and-unrolled-rnn/>  
<https://towardsdatascience.com/rnn-recurrent-neural-networks-how-to-successfully-model-sequences-10f3a2a3a2>

## Recurrent Neural Network(RNN)

The RNN computational graph and information flow



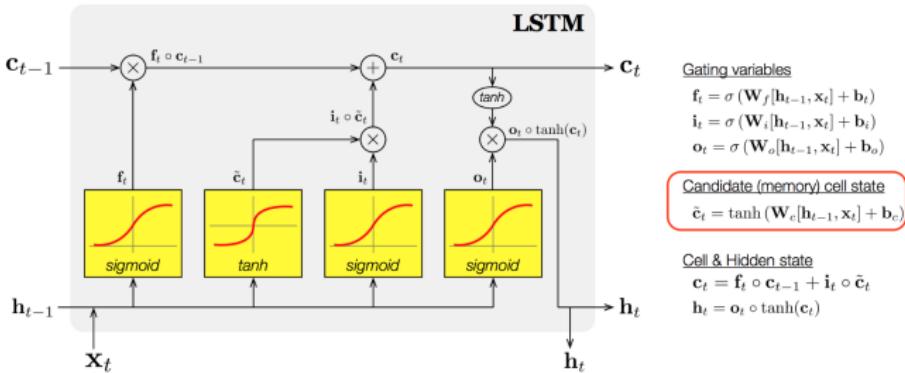
The same set of weights are used to compute the hidden state and output for all time steps (shared weights)

# Long Short Term Memory(LSTM)

## Long Short-Term Memory (LSTM) Network

LSTM<sup>40</sup> is a gated RNNs capable of learning long-term dependencies. It categorize data into **short** and **long** term, deciding its importance and what to remember or forget.

The LSTM unit cell has an **input** and a **forget gate**. The **input** defines how much of the newly computed state for the current input is accepted, while the **forget** defines how much of the previous state is accepted.



<sup>40</sup> <http://colah.github.io/posts/2015-08-Understanding-LSTMs>

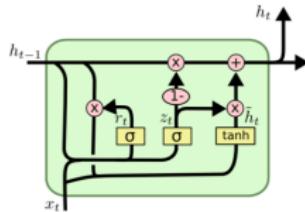
<https://towardsdatascience.com/lstm-recurrent-neural-networks-how-to-teach-a-network-to-remember-10133a2a2a1>

## Gated Recurrent Network (GRU)

### Gated Recurrent Network (GRU)

GRU<sup>41</sup> networks have been proposed as a simplified version of LSTM, which also avoids the vanishing gradient problem and is even easier to train.

In a GRU the **reset** gate determines how to combine the new input with the previous memory, and the **update** gate defines how much of the previous memory is kept



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

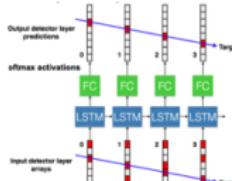
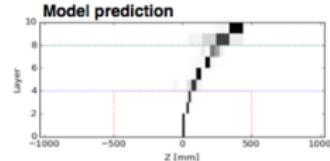
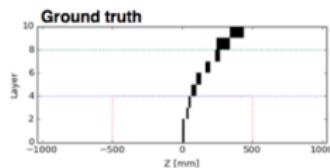
$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

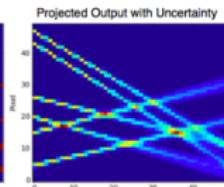
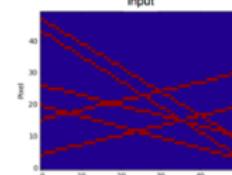
<sup>41</sup> <https://isaacchanghau.github.io/post/lstm-gru-formula>

# LSTM Application in HEP: Tracking

Tracking with recurrent neural networks ( LSTM ) 42



Time dimension  
(state memory)



## RNN Difficulties

### RNN Difficulties

- Large inputs leads to very deep RNN networks:
  - vanishing or exploding gradient problem
  - memory limitations
  - hard to train
- Process data in a sequential way  $\Rightarrow$  no parallel processing (GPU,TPU)

A solution to these problems is the attention mechanism , which allows to grab a whole data sequence in a parallel way and infer distant correlations. <sup>43</sup>

### Attention

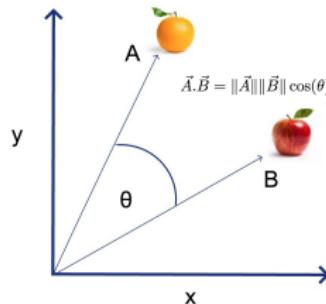
Attention is motivated by how we pay visual attention to different regions of an image, correlate words in sentences or focus on a voice ignoring background noise. It also allows us to pursue one thought at a time and remember an individual event rather than all events.

## Attention Mechanism

Geometric similarity between vectors (points) can be estimated by the "dot product"

### Attention

Attention mechanism generalizes geometric similarity by learning a metric to quantify similarity between arbitrary objects (words, images, ...)<sup>44</sup>



⇒ Attention can be interpreted as weights attributed to pixels in an image or words in a text, focusing on some part of the data while ignoring others.<sup>45</sup>

<sup>44</sup>A.Karpathy, Introduction to Transformers - <https://www.youtube.com/watch?v=XfpMkf4rD6E>

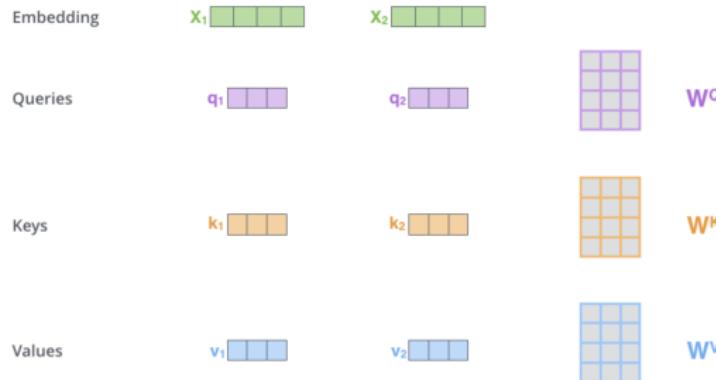
<sup>45</sup>C.Olah, Attention and Augmented RNN - <https://distill.pub/2016/augmented-rnns/> ◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏷ ⏸ ⏹ ⏺

## Transformers : Attention is All you Need

The Transformer architecture is a sequential data model that uses attention based MLPs outperforming RNNs.<sup>46</sup>

### Self-Attention

Self-attention is calculated from three vectors (Query, Key, Value) for each input vector (word embedding). They are defined by multiplying three weight matrices learned during training<sup>47</sup>



<sup>46</sup> A.Vaswani & all , Attention Is All you Need (2017) , <https://arxiv.org/abs/1706.03762>

<sup>47</sup> <https://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-m/>

## Transformer Networks

### Matrix Calculation of Self-Attention

Query, Key, and Value matrices are defined by packing the embeddings into a matrix  $X$  and multiplying by weight matrices  $W^Q$ ,  $W^K$  and  $W^V$ . The self-attention layer output is then softmaxed

$$\begin{array}{ccc}
 X & W^Q & Q \\
 \begin{matrix} \text{green} \\ \text{grid} \end{matrix} & \times & \begin{matrix} \text{purple} \\ \text{grid} \end{matrix} = \begin{matrix} \text{purple} \\ \text{grid} \end{matrix} \\
 \\ 
 X & W^K & K \\
 \begin{matrix} \text{green} \\ \text{grid} \end{matrix} & \times & \begin{matrix} \text{orange} \\ \text{grid} \end{matrix} = \begin{matrix} \text{orange} \\ \text{grid} \end{matrix} \\
 \\ 
 X & W^V & V \\
 \begin{matrix} \text{green} \\ \text{grid} \end{matrix} & \times & \begin{matrix} \text{blue} \\ \text{grid} \end{matrix} = \begin{matrix} \text{blue} \\ \text{grid} \end{matrix}
 \end{array}$$

$$\text{softmax} \left( \frac{Q \times K^T}{\sqrt{d_k}} \right) = Z$$

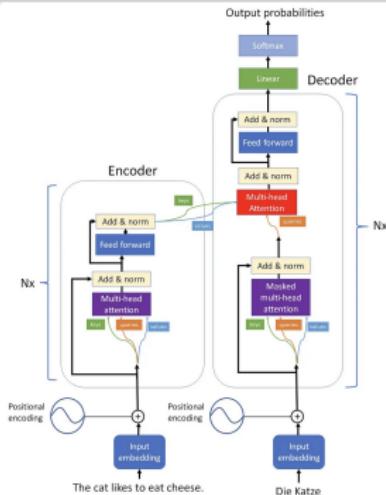
$$Z \times V$$

Transformers uses "multi-headed" attention where each "head" has its own set of learned weights ( like CNN channels ). Each head processes the inputs in parallel and independently, focusing on different aspects.

# Transformer Networks

## Transformer Architecture

The transformer architecture <sup>48</sup> has an encoder decoder structure. It uses attention and positional encoding to weight the importance of each part of the input data.



- The network has an encoder and decoder structure with independent weights.
- Self attention is used to weight the importance of each input token
- Cross attention combines keys and values from encoder (input sequence) with query from the decoder to predict the next output sequence token <sup>49</sup>
- Transformers can also deal with images by splitting it into patches ("image tokens")

<sup>48</sup> <https://towardsdatascience.com/transformers-explained-visually-part-1-overview-of-functionality-1f3a2a2a3a>  
<sup>49</sup> <https://vaclavkosar.com/ml/cross-attention-in-transformer-architecture>

## Transformer Networks

### ChatGPT

ChatGPT<sup>50</sup> is a generative pre-trained transformer (GPT) whose core function of a chatbot is to mimic a human conversationalist. ChatGPT is versatile and it can write and debug computer programs, compose music, fairy tales, poetry and student essays and it also answers test questions and play games.

<https://chat.openai.com/chat>

### Vision Transformers (ViT)

Computing relationships for every pixel pair in a typical image is computationally prohibitive. ViT splits an image into fixed size patches and flattens it to create lower-dimensional embeddings. After including positional embeddings the sequence is input to a Transformer. ViT attains state of the art performance in image classification, segmentation object detection and etc. ViT outperform the current state-of-the-art CNNs in terms of computational efficiency and accuracy.

<https://viso.ai/deep-learning/vision-transformer-vit>

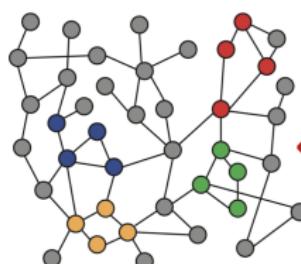
<sup>50</sup>A.Karpathy , Let's build GPT from scratch - <https://www.youtube.com/watch?v=kCc8FmEb1NY>



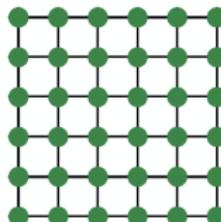
# Graph Structured Data

## Graphs

A graph is an abstract mathematical structure defined by its nodes(vertices) and edges(links):  $G(N, E)$ . The nodes list is unordered and can have variable length, while edges represent nodes relations. Each node and edge can have features ( ex: node position, edge length, ... ).



**Networks**



**Images**

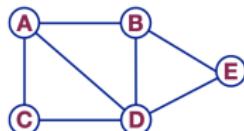


**Text**

Sequences ( linear graph ) and images ( regular grid graph ) can be seen as special types of graph structured data !

## Graph Representation

For a NN to operate on a given dataset it's necessary to encode it in a mathematical representation. A naive approach would be to represent graph structured data as an augmented adjacency matrix with node features information and use it as an MLP input.<sup>52</sup>



	A	B	C	D	E	Feat
A	0	1	1	1	0	1 0
B	1	0	0	1	1	0 0
C	1	0	0	1	0	0 1
D	1	1	1	0	1	1 1
E	0	1	0	1	0	1 0

### Problems:

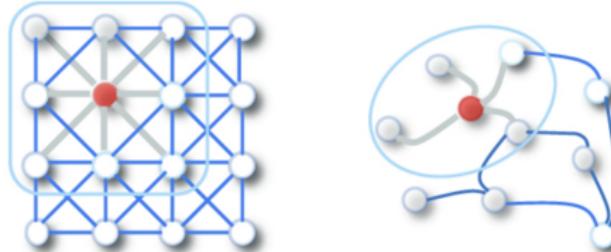
- Adjacency matrix depends on arbitrary node labeling (**need permutation invariance**)
- Adjacency matrix can become too large and sparse (**need compact representation**)
- Matrix size depend on number of nodes (**NN takes fixed input size**)

<sup>52</sup> <https://towardsdatascience.com/how-to-do-deep-learning-on-graphs-with-graph-convolutional-networks>  
<https://tkipf.github.io/graph-convolutional-networks>

## Graph Neural Network (GNN)

### From Convolution to Message Passing

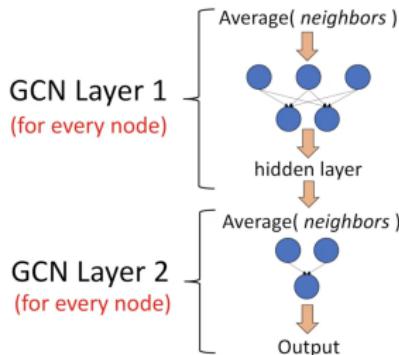
Inspired by convolution on images, which are grids of localized points, we can define a convolution on a graph where nodes have no spatial order. The node neighbourhood aggregation is often called *message passing*<sup>53</sup>.



<sup>53</sup> J.Gilmer & all. , Neural Message Passing for Quantum Chemistry , <https://arxiv.org/abs/1704.01212> ↗ ↘ ↙

## Graph Convolutional Network (GCN)

Taking the node features as inputs the network propagates <sup>55</sup> it through the hidden layers getting a latent representation ( embedding ) of the node features. This latent representation can have an arbitrary number of dimensions, like the number of filters of a CNN.



### Node Neural Network

$$h_v^0 = x_v \text{ (input features)}$$

$$h_v^k = \sigma \left( W_k' h_v^{k-1} + W_k \sum_{u \in N(v)} \frac{h_u^{k-1}}{|N(v)|} \right)$$

$$z_v = h_v^n \text{ (output)}$$

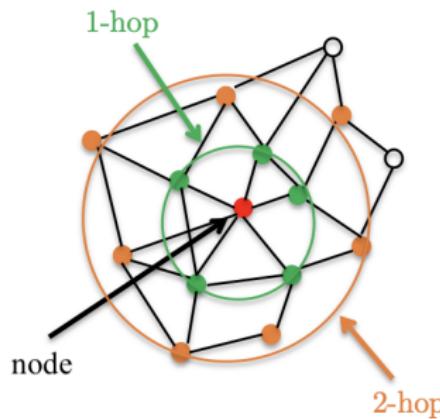
, where  $h_v^k$  is the node embedding in layer-k

<sup>55</sup> <https://towardsdatascience.com/hands-on-graph-neural-networks-with-pytorch-pytorch-geometry-113a2f3a2e0>  
<http://web.stanford.edu/class/cs224w/slides/08-GNN.pdf>

## Graph Neural Network (GNN)

### GNN Depth

GNNs with few hidden layers are enough to exhaust small graphs with a few hops .The larger the number of layers the farther away messages gets passed



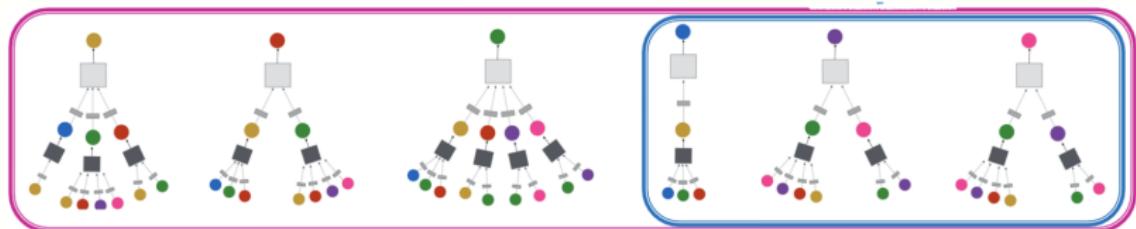
Obs: very deep GNNs tend to suffer from oversmoothing  $\Rightarrow$  nodes features becomes identical

## Graph Neural Network (GNN)

### Inductive Capability

GNN can be applied to graphs with arbitrary number of nodes, since weights  $W_k$  are shared across all nodes of a given layer-k.

For an example , a model trained on graphs with nodes A,B,C , can also evaluate graphs with nodes D,E,F.

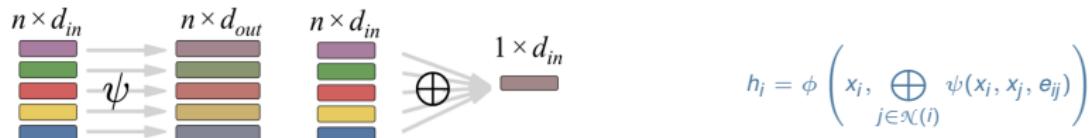


## Graph Network (GN)

A general graph can have edge features as well as node features.<sup>56</sup>

### Graph Network (GN)

The most general Graph Network(GN) can be modeled as a graph to graph mapping. An update (message passing) function  $\psi$  acts on nodes features  $x_i$  or edges features  $e_{ij}$ , followed by a permutation invariant aggregation  $\oplus$  and an activation function  $\phi$ .



The update function takes in a *fixed size*  $d_{in}$  input and returns a *fixed size*  $d_{out}$  output, while the aggregation function takes in *variable sized*  $n$  inputs and returns a *fixed size*  $d_{in}$  output.

Obs: GN formalism is not a specific model architecture and doesn't determine what the update and aggregate functions are.

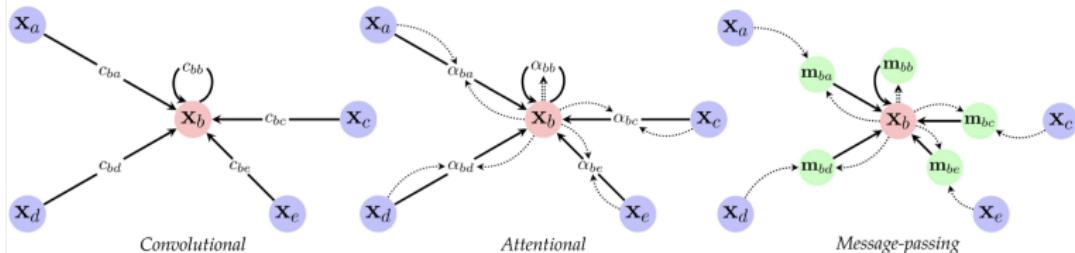
<sup>56</sup> [https://github.com/deepmind/graph\\_nets](https://github.com/deepmind/graph_nets)  
<https://arxiv.org/pdf/2007.13681>

# Graph Neural Network

## Three Particular Flavours of GNNs

GNNs can be further grouped<sup>57</sup> according to the following flavours:

- **Convolutional (GCN):** fixed weights coefficients  $c_{ij}$
- **Attention Based (GAT):** learnable attention weights calculated as  $\alpha_{ij} = a(x_i, x_j)$
- **Message Passing (GN):** general messages  $m_{ij} = \psi(x_i, x_j)$  sent across edges



$$\mathbf{h}_i = \phi \left( \mathbf{x}_i, \bigoplus_{j \in \mathcal{N}_i} c_{ij} \psi(\mathbf{x}_j) \right)$$

$$\mathbf{h}_i = \phi \left( \mathbf{x}_i, \bigoplus_{j \in \mathcal{N}_i} a(\mathbf{x}_i, \mathbf{x}_j) \psi(\mathbf{x}_j) \right)$$

$$\mathbf{h}_i = \phi \left( \mathbf{x}_i, \bigoplus_{j \in \mathcal{N}_i} \psi(\mathbf{x}_i, \mathbf{x}_j) \right)$$

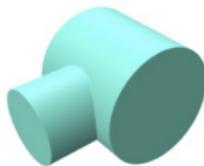
<sup>57</sup> <https://petar-v.com/talks/GNN-Wednesday.pdf>

## From Graphs to Sets

### Sets

Abstracting even more on graphs we encounter sets , which are collections without intrinsic order or relations<sup>58</sup>

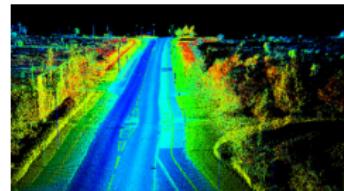
A MLP in principle can learn not to care about orderings, but training it for that would be a waste of computing resources. One can build a permutation invariant network by design, taking permutation symmetry as inductive bias in the architecture. DeepSets and PointNet are examples of networks capable of learning on point set structured data.



(a) Surface view.



(b) Point cloud.



Obs: data generated by 3D scanners and LIDAR sensors often come as point clouds with coordinates (x,y,z) as its features

<sup>58</sup> P.Velickovic, Graphs and Sets, [https://www.youtube.com/watch?v=E\\_Wweuk5iqA&list=PLn2-dEmQeTfQ8YVuHBOvAhUlnIPYxkeu3&index=5&t=13s](https://www.youtube.com/watch?v=E_Wweuk5iqA&list=PLn2-dEmQeTfQ8YVuHBOvAhUlnIPYxkeu3&index=5&t=13s)

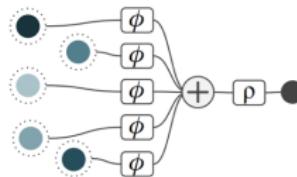
## DeepSets / PointNet Networks

Consider an input vector  $\vec{x}$  with  $N$  components. The set  $S_N$  of all permutations of its components. To construct a permutation invariant function we sum over all permutations  $\pi(\vec{x})$  of  $\vec{x}$ .<sup>59</sup>

### Janossy Pooling

$$f(\vec{x}) = \frac{1}{|S_n|} \sum_{\pi \in S_N} \phi(\pi(\vec{x})) \quad \Rightarrow \text{computationally expensive and scales as } N!$$

Complexity can be reduced using permutation invariants sum of  $k$ -tuples , where  $k < N$ . For  $k = 1$  we have the architectures known as DeepSets and PointNet

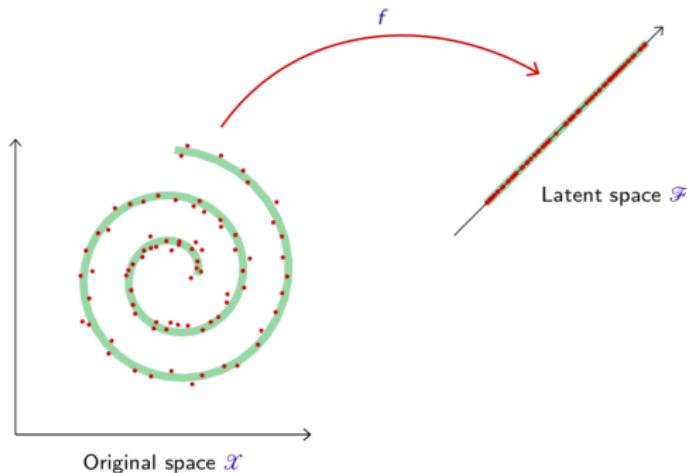


$\Rightarrow$  The sum aggregation of the  $\phi$  layers gives permutation invariance over inputs

<sup>59</sup> <https://fabianfuchsml.github.io/learningonsets>

## Autoencoder

Many applications such as data compression, denoising and data generation require to go beyond classification and regression problems. This modeling usually consists of finding “meaningful degrees of freedom”, that can describe high dimensional data in terms of a smaller dimensional representation



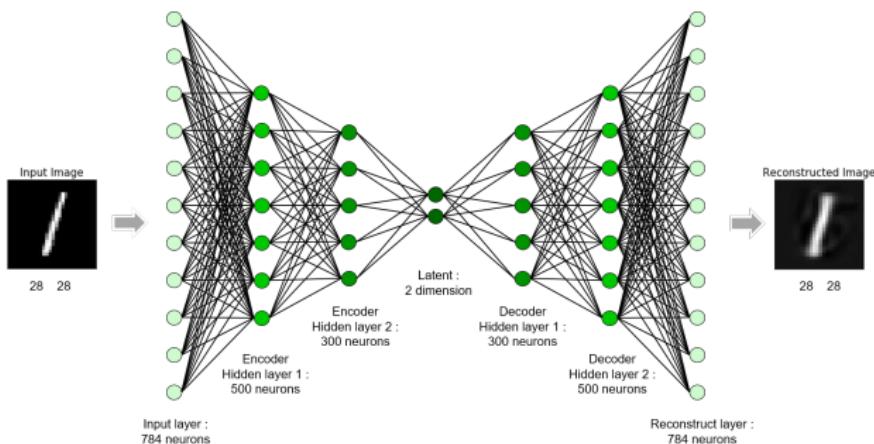
Traditionally, autoencoders were used for dimensionality reduction and denoising. Recently autoencoders are being used also in generative modeling

## Autoencoder(AE)

### Autoencoder(AE)

An AE is a neural network that is trained to attempt to copy its input to its output in an **self-supervised way**. In doing so, it learns a representation(encoding) of the data set features / in a low dimensional latent space.

It may be viewed as consisting of two parts: an encoder  $l = f(x)$  and a decoder  $y = g(l)$ .



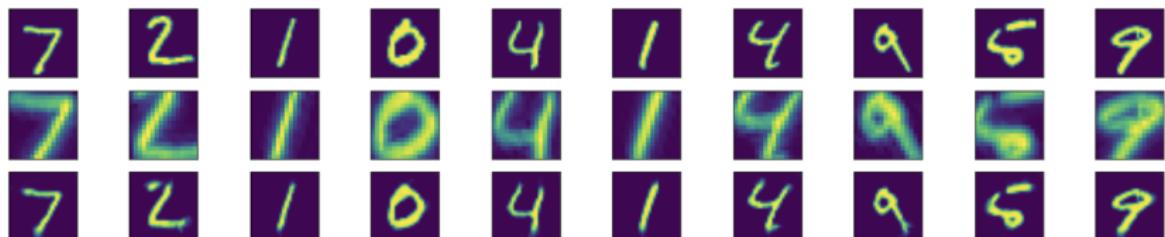
⇒ Understanding is data reduction

# Convolutional Autoencoder(CAE)

## Convolutional Autoencoder(CAE)

A CAE<sup>62</sup> is built out of convolutional layers assembled in a AE architecture and it's trained to attempt to reconstruct an output image from an input image after data compression. In doing so, it learns how to compress images.

Bellow we have MNIST digits ( 28x28 pixels ) used as input and the output images and the corresponding latent space compressed images ( 16x16 pixels )



⇒ CAE learns that image borders have no information and chops the central part.

<sup>62</sup> <https://towardsdatascience.com/introduction-to-autoencoders-b6fc3141f072>

## Denoising Autoencoder (DAE)

The DAE is an extension of a classical autoencoder where one corrupts the original image on purpose by adding random noise to its input. The autoencoder is trained to reconstruct the input from a corrupted version of it and then used as a tool for noise extraction

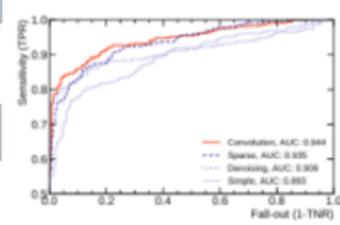
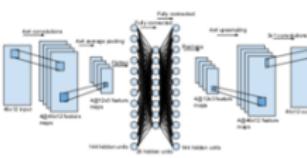
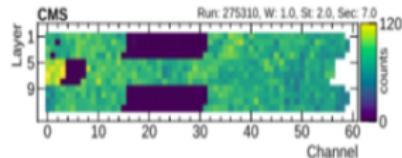
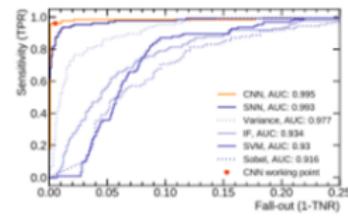
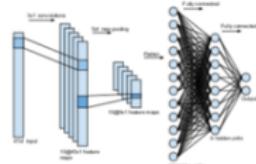
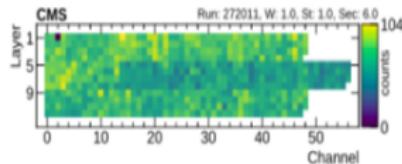


Autoencoder Application in HEP: DQM

AE can be used for anomaly detection by training on a single class , so that every anomaly gives a large reconstruction error

## Detector Quality Monitoring (DQM)

Monitoring the CMS data taking to spot failures ( anomalies ) in the detector systems

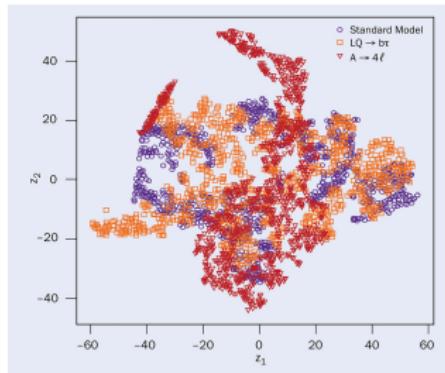


## Autoencoder Application in HEP: Anomaly Detection at LHC ( BSM )

AE can watch for signs of new physics at the LHC that have not yet been dreamt up by physicists <sup>63</sup>

### Anomaly Detection ( BSM Physics Search)

AE trained on a SM data sample , so that an anomalous event gives a large reconstruction error. The LHC collisions are compressed by an AE to a two-dimensional representation ( $z_1, z_2$ ). The most anomalous events populate the outlying regions.



Outliers could go into a data stream of interesting events to be further scrutinised

<sup>63</sup><https://cerncourier.com/a/hunting-anomalies-with-an-ai-trigger/>

## Variational Autoencoder(VAE)

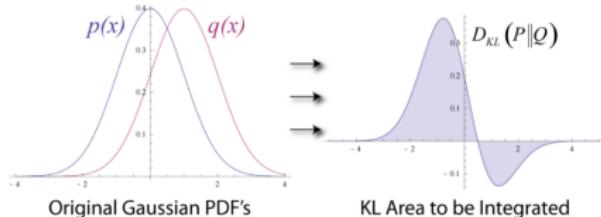
### Variational Autoencoder(VAE)

A VAE<sup>68</sup> is an autoencoder with a loss function penalty Kullback–Leibler (KL)divergence<sup>69 70</sup> that enforces a latent space representation that follows a normal distribution. To generate images with a VAE one just samples a latent vector from a unit gaussian and pass it through the decoder.

The KL divergence, or relative entropy, is a measures of how one probability distribution differs from a second, reference distribution (  $\Rightarrow$  information loss )

### Kullback–Leibler Divergence ( Relative Entropy )

$$D_{KL}(p||q) = \int_{-\infty}^{+\infty} dx p(x) \log \left( \frac{p(x)}{q(x)} \right)$$



<sup>68</sup> <http://kvfrans.com/variational-autoencoders-explained>

<sup>69</sup> [https://en.wikipedia.org/wiki/Evidence\\_lower\\_bound](https://en.wikipedia.org/wiki/Evidence_lower_bound)

<sup>70</sup> <https://www.youtube.com/watch?v=HxQ94L8n0vU>

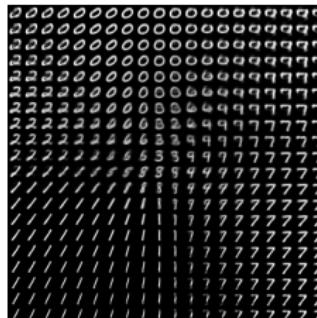
## Variational Autoencoder(VAE)

### VAE Loss

The VAE loss <sup>71</sup> is composed of a mean squared error term that measures the reconstruction accuracy and a KL divergence term that measures how close the latent variables match a gaussian.

$$L = \|x - f \circ g(x)\|^2 + D_{KL}(p(z|x) | q(z|x))$$

VAE generated numbers obtained by gaussian sampling a 2D latent space



⇒ KL loss is a regularization term that helps learning "well organized"latent space

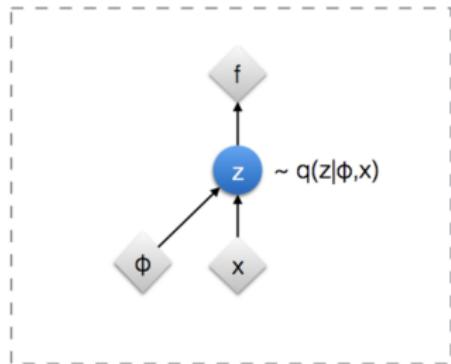
<sup>71</sup><https://tiao.io/post/tutorial-on-variational-autoencoders-with-a-concise-keras-implement/>

## Variational Autoencoder(VAE)

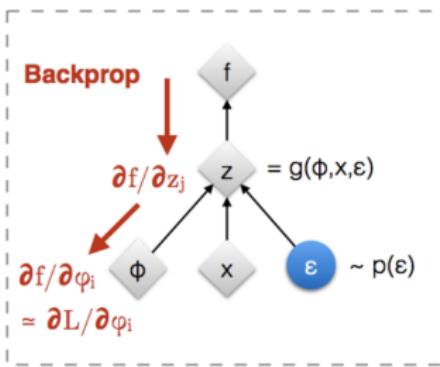
### VAE Reparametrization Trick

$$\begin{aligned} z &= \mu + \sigma * z_c = \mu + \sigma * (\mu_c + \sigma_c * \epsilon) \\ &= \underbrace{(\mu + \sigma * \mu_c)}_{\text{VAE mean}} + \underbrace{(\sigma * \sigma_c) * \epsilon}_{\text{VAE std}} \end{aligned}$$

Original form



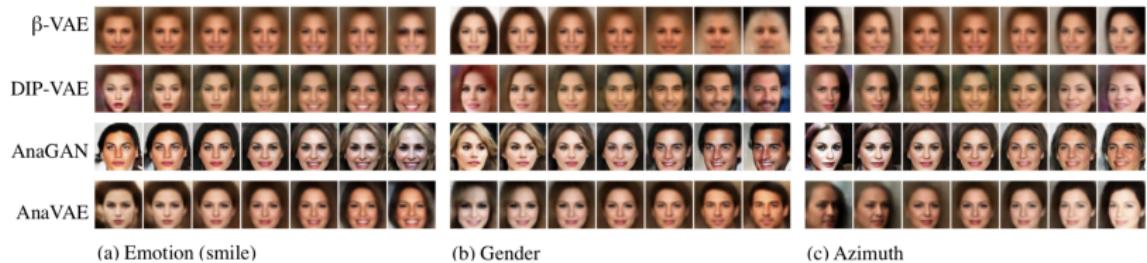
Reparameterised form



## Variational Autoencoder(VAE)

### Feature Disentanglement(Factorization) in VAE

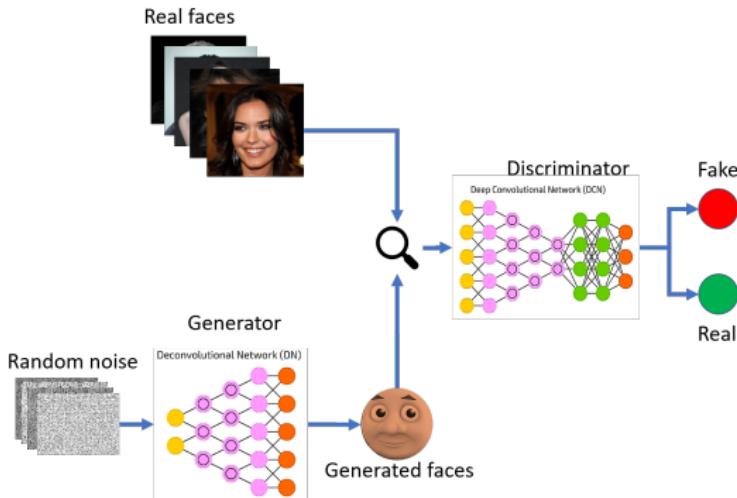
Feature disentanglement is isolating the source of variation in data <sup>72</sup>



# Generative Adversarial Network(GAN)

## Generative Adversarial Networks (GAN)

GANs are composed by two NN, where one generates candidates and the other classifies them. The generator learns a map from a latent space to data, while the classifier discriminates generated data from real data.

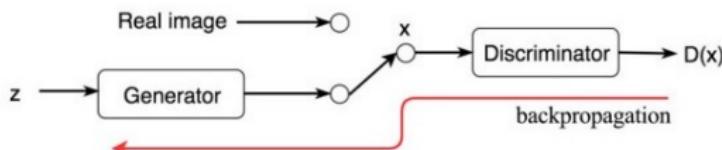


## Generative Adversarial Network(GAN)

GAN adversarial training <sup>73</sup> works by the two neural networks competing and training each other. The generator tries to "fool" the discriminator, while the discriminator tries to uncover the fake data.

### GAN Training

- ① The discriminator receives as input samples synthesized by the generator and real data . It is trained just like a classifier, so if the input is real, we want output=1 and if it's generated, output=0
- ② The generator is seeded with a randomized input that is sampled from a predefined latent space (ex: multivariate normal distribution)
- ③ We train the generator by backpropagating this target value all the way back to the generator
- ④ Both networks are trained in alternating steps and in competition



<sup>73</sup>[https://medium.com/@jonathan\\_hui/gan-whats-generative-adversarial-networks-and-its-applications-1f7d8c3e3d0](https://medium.com/@jonathan_hui/gan-whats-generative-adversarial-networks-and-its-applications-1f7d8c3e3d0)

## Generative Adversarial Network(GAN)

GANs are quite good on faking celebrities images<sup>75</sup> or Monet style paintings<sup>76</sup> !



Training Data



Sample Generator



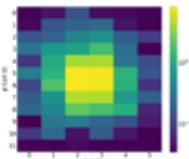
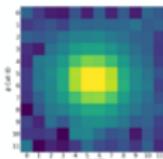
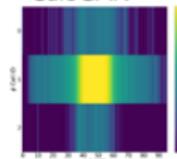
<sup>75</sup>[https://research.nvidia.com/publication/2017-10\\_Progressive-Growing-of](https://research.nvidia.com/publication/2017-10_Progressive-Growing-of)  
<sup>76</sup><https://github.com/junyanz/pytorch-CycleGAN-and-pix2pix>

# GAN Application in HEP: MC Simulation

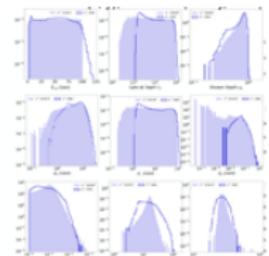
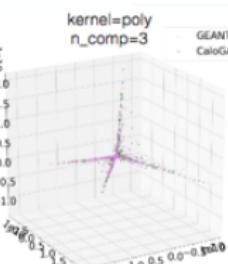
## CaloGAN

Simulating 3D high energy particle showers in multi-layer electromagnetic calorimeters with a GAN<sup>77</sup>

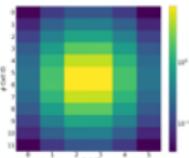
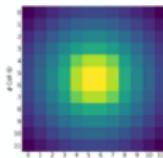
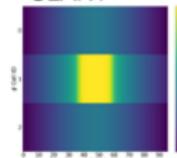
### CaloGAN



kernel=poly  
n\_comp=3



### GEANT



<sup>77</sup> <https://github.com/hep-lbdl/CaloGAN>

## Diffusion Model (DM)

A diffusion model works as a Denoising Autoencoder(DAE) trained to denoise images blurred with Gaussian noise.



Introducing Sora — OpenAI's text-to-video model

## Diffusion Model (DM)

### Diffusion Model

Diffusion models define a Markov chain of diffusion steps by slowly adding random noise to data and then learn to reverse the diffusion process in order to reconstruct desired data samples from the noise.<sup>82</sup>

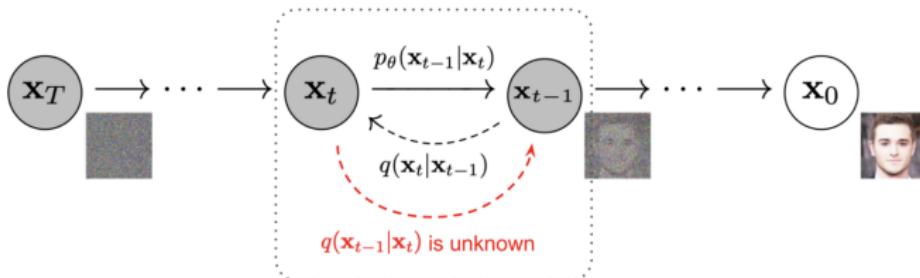


Fig. 2. The Markov chain of forward (reverse) diffusion process of generating a sample by slowly adding (removing) noise. (Image source: [Ho et al. 2020](#) with a few additional annotations)

82

<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

## Diffusion Model (DM)

### Forward Process

Given a data point  $x_0$ , the forward diffusion process adds small amount of Gaussian noise in steps (Markov Chain), producing a sequence of noisy samples. The step sizes are controlled by  $\beta_t$ . At the end we are left with a noisy image represented by an isotropic Gaussian distribution ( same variance along all dimensions).

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

⇒ The model latent space has the same dimensionality as the original data ( high dimensionality )

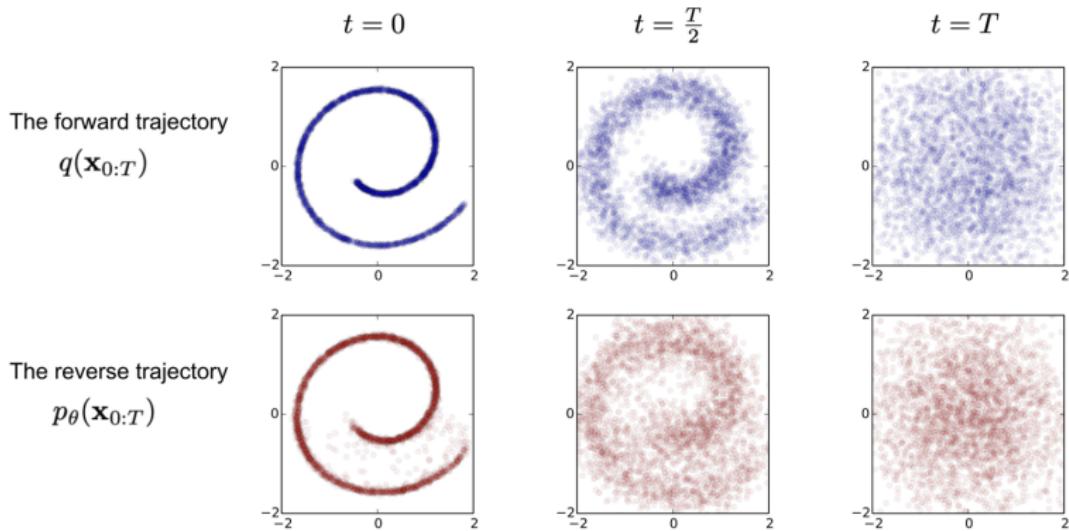
### Reverse/Backward Process

The task of removing the added noise in the forward process, again in an iterative fashion, is done using a neural network. If we can reverse the above process and sample from  $q(x_{t-1} | x_t)$ , we can recreate the original data sample from a Gaussian noise input  $X_T \sim \mathcal{N}(0, 1)$ .

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) \quad p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

## Diffusion Model (DM)

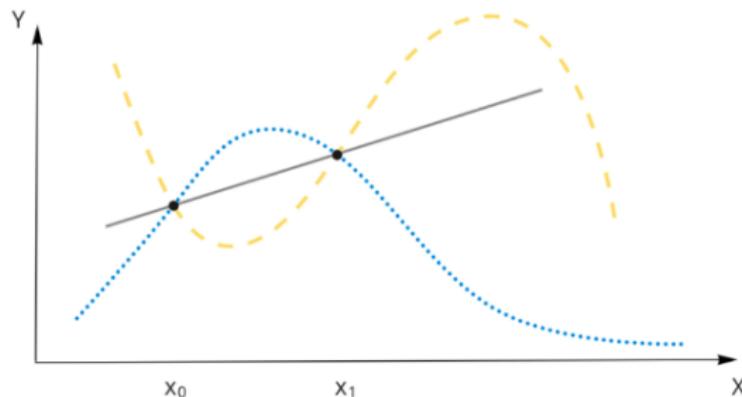
Example of training a diffusion model for modeling a 2D Swiss roll data <sup>83</sup>



<sup>83</sup><https://arxiv.org/abs/1503.03585>

## Domain Knowledge and Inductive Bias

As an example of how a set of observations can lead to different hypothesis depending on the inductive biases consider a fit given a set of data points



There's an infinite number of functions that could pass through the samples data points <sup>86</sup>

<sup>86</sup> <https://towardsdatascience.com/the-inductive-bias-of-ml-models-and-why-you-should-care-a>  
U. von Luxburg , Machine learning and inductive bias (Tubingen 2020) , <https://www.youtube.com/watch?v=WB8eYZSzyE&list=PL05umP7R6ij2XCvrRzLokX6EoHWaGA2cC&index=2&t=2372s>

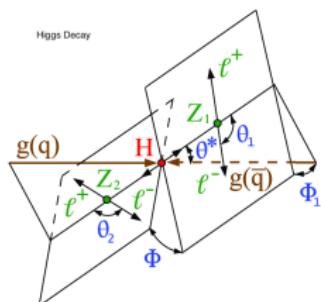
# Domain Knowledge and Inductive Bias

Domain knowledge can be infused into a model in different ways<sup>87</sup>

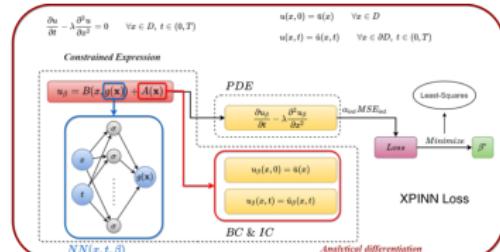
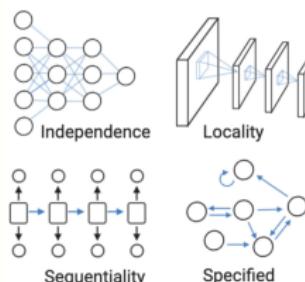
## Domain Knowledge

- Using engineered variables variables as input data
- Creating specialized architecture layers: CNN, RNN, GNN, DS
- Adding a penalty term to the loss function: PINN

⇒ It's also possible to use simultaneous combinations of the above methods



## Relational Inductive Biases



<sup>87</sup> <https://sgfin.github.io/2020/06/22/Induction-Intro/>  
<https://www.nature.com/articles/s41598-021-04590-0>

# The End !!!