# Deep Learning in Data Analysis: Introduction to Deep Learning in HEP
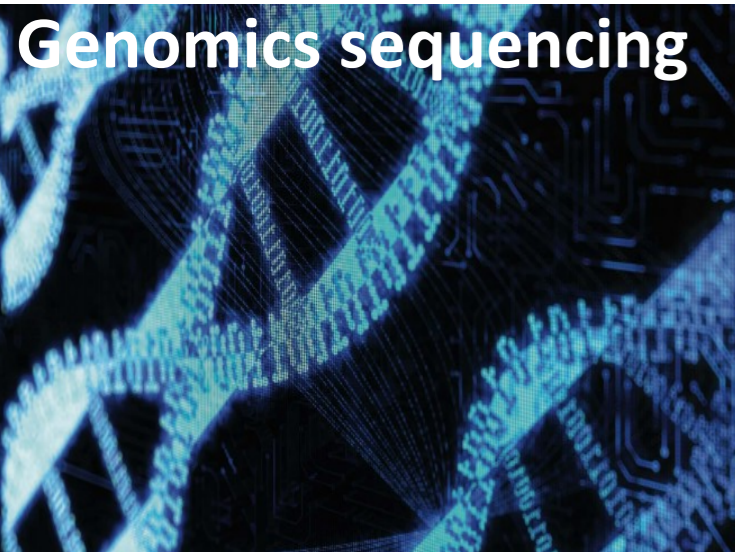
Lecture 1

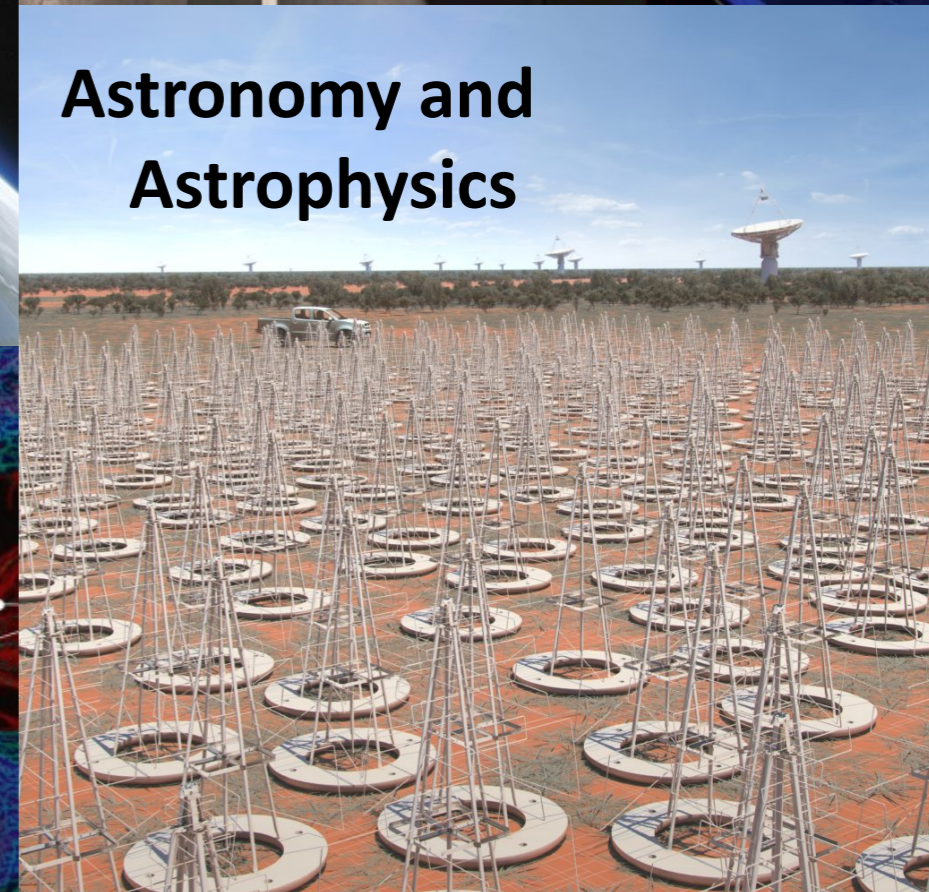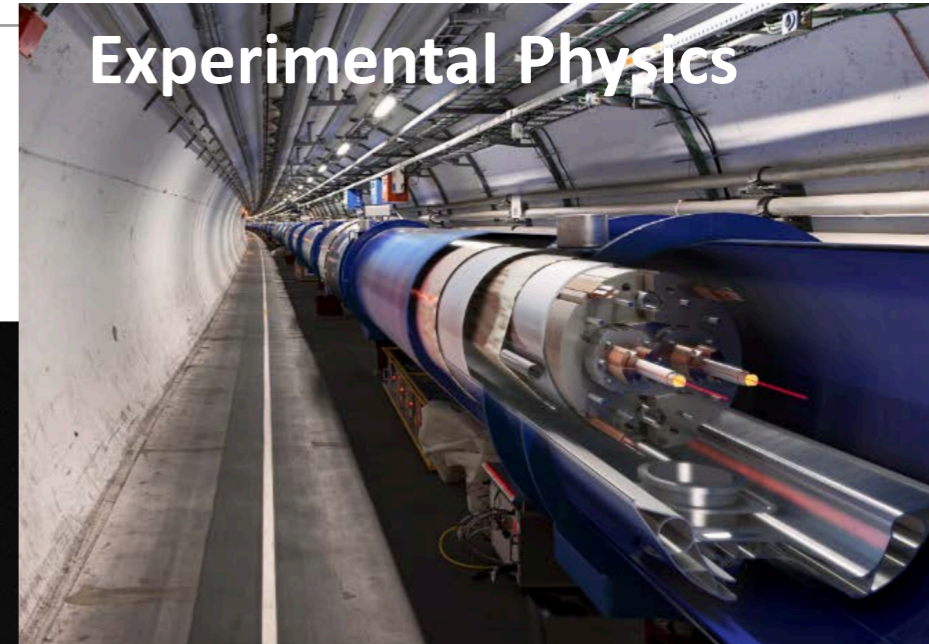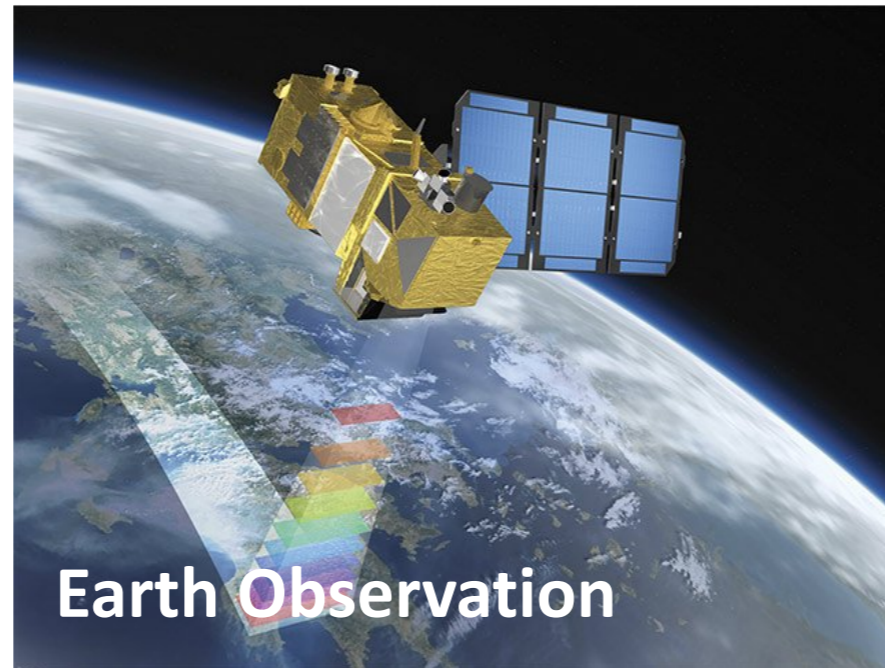Sofia Vallecorsa | Ilaria Luise

Thematic CERN School of Computing on Machine Learning
17th October 2024

# Outline

- Introduction

  - The need for depth  -  graphs complexity

- Computational challenges

- Generative Models

# Big Data in Science

Science produces more data than ever before and at an unmatched pace in history
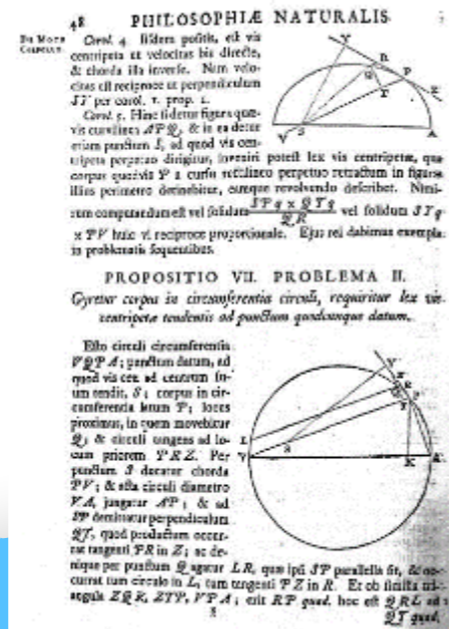
**Experimental Physics**

**Earth Observation**

**Astronomy and Astrophysics**

**Genomics sequencing**

**Biology and Microscopy imaging**

Sofia Vallecorsa, Ilaria Luise CERN - sofia.vallecorsa@cern.ch | ilaria.luise@cern.ch

# Four Paradigms of Scientific Research



**Today**

Data-driven science

**~50 years**

Simulations
Computational sciences

**500 years**

Generalization
Theoretical models

**4000 years**

Empirical observations

Sofia Vallecorsa, Ilaria Luise CERN - sofia.vallecorsa@cern.ch | ilaria.luise@cern.ch

# Data-driven science & AI

Is Artificial Intelligence just a **refined, faster** approach to computational science?

Machine Learning

Deep Learning

Artificial Intelligence

# Rediscovering physics

Udrescu, Silviu-Marian, and Max Tegmark. "**AI Feynman: A physics-inspired method for symbolic regression**." *Science Advances* 6.16 (2020): eaay2631.
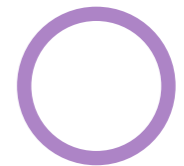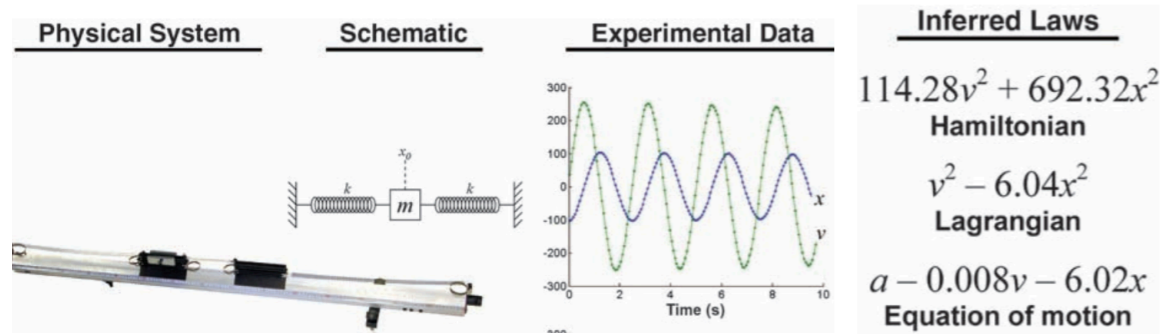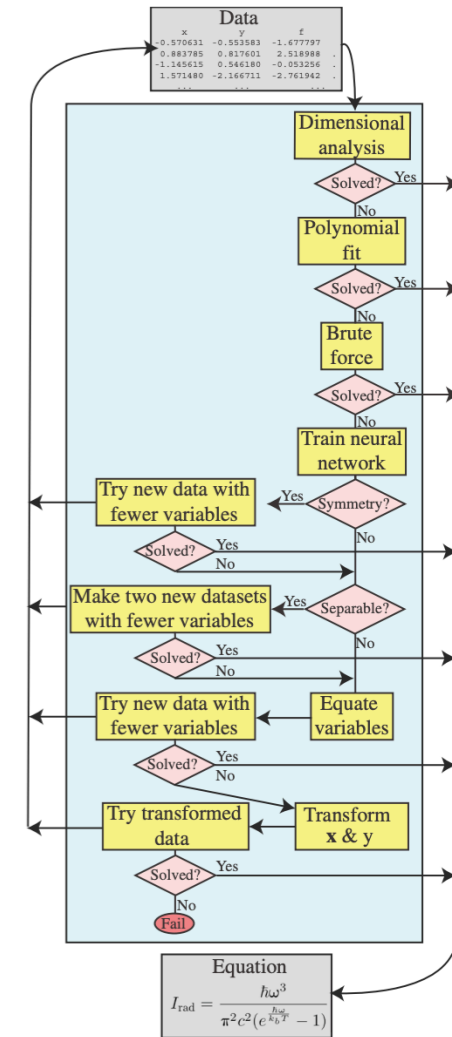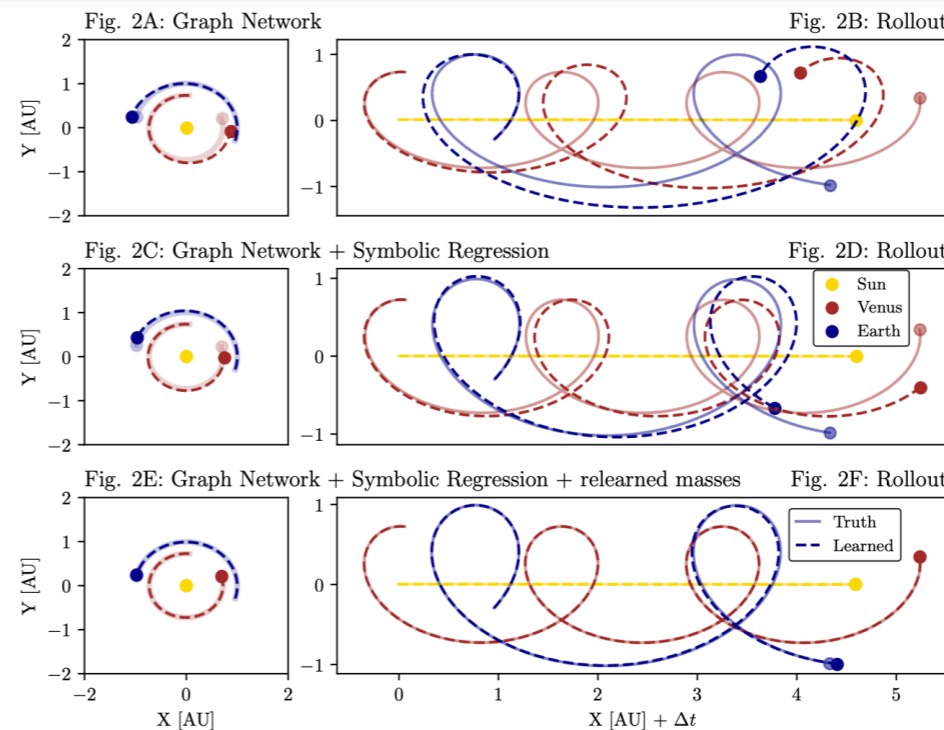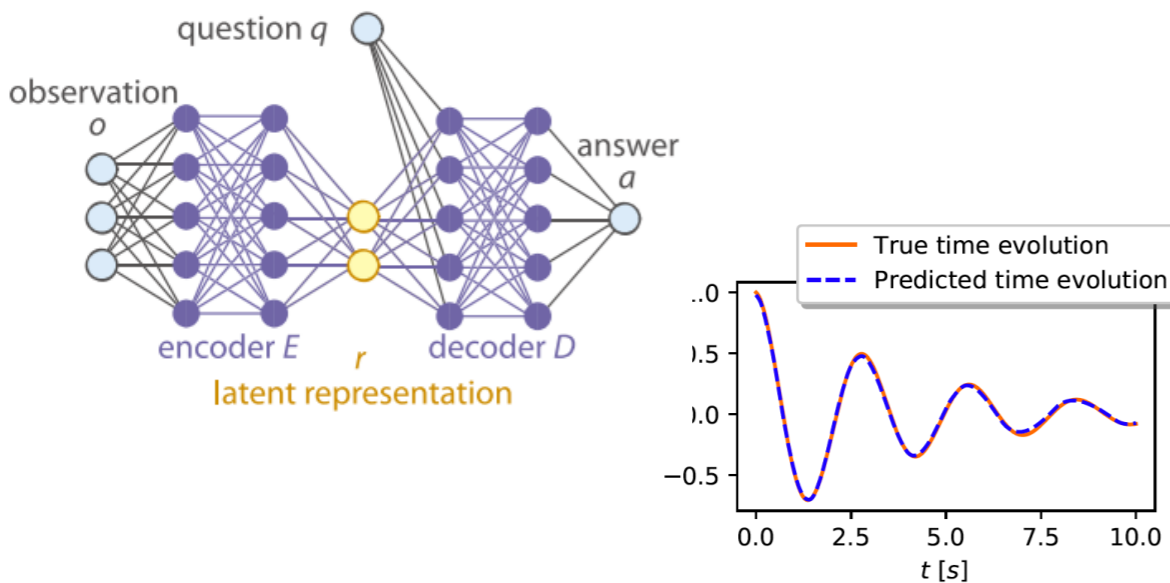
Schmidt, Michael, and Hod Lipson. "**Distilling free-form natural laws from experimental data**." *science* 324.5923 (2009): 81-85.



Physical System   Schematic   Experimental Data

**Inferred Laws**

$$114.28v^2 + 692.32x^2$$
**Hamiltonian**

$$v^2 - 6.04x^2$$
**Lagrangian**

$$a - 0.008v - 6.02x$$
**Equation of motion**

Lemos, Pablo, et al. "**Rediscovering orbital mechanics with machine learning**." *arXiv:2202.02306* (2022)



Fig. 2A: Graph Network    Fig. 2B: Rollout
Fig. 2C: Graph Network + Symbolic Regression    Fig. 2D: Rollout
Fig. 2E: Graph Network + Symbolic Regression + relearned masses    Fig. 2F: Rollout

Sun
Venus
Earth

Truth
Learned

Iten, Raban, et al. "**Discovering physical concepts with neural networks**." *Physical review letters* 124.1 (2020): 010508.



question q
observation o
answer a
encoder E       decoder D
r
latent representation



True time evolution
Predicted time evolution

Can we train AI to understand **physics itself** in order to achieve new discoveries ?

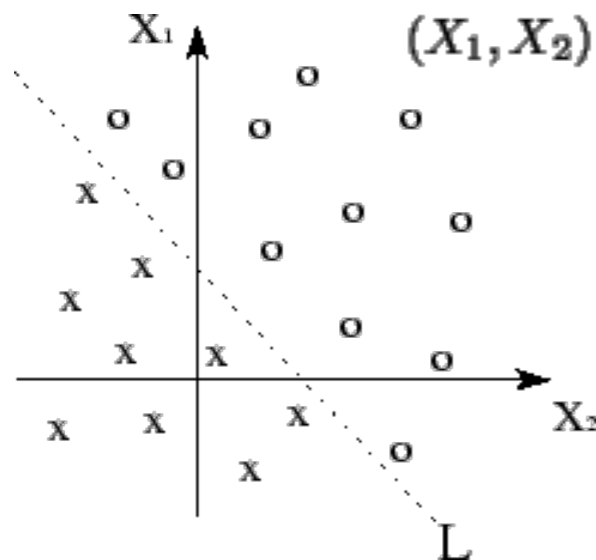Let's start at the beginning

# Universal approximator

NN with a single hidden layer containing a finite number of non-linear neurons approximate continuous functions to any desired degree of accuracy.

Hornik, Kurt; Tinchcombe, Maxwell; White, Halbert (1989). *Neural Networks*. **2**. Pergamon Press. pp. 359–366.

# The need for depth

A single layer perceptron can categorize "linearly separable" patterns (binary classification):

OR function is linearly separable

Exclusive OR is a non linearly separable pattern:



$$Y = X_1 \oplus X_2$$

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

(tutorial) http://www.theprojectspot.com/tutorial-post/introduction-to-artificial-neural-networks-part-1/7

# The need for depth (II)

Need a Multi-Layer architecture to solve the exclusive OR problem with a two-stages approach

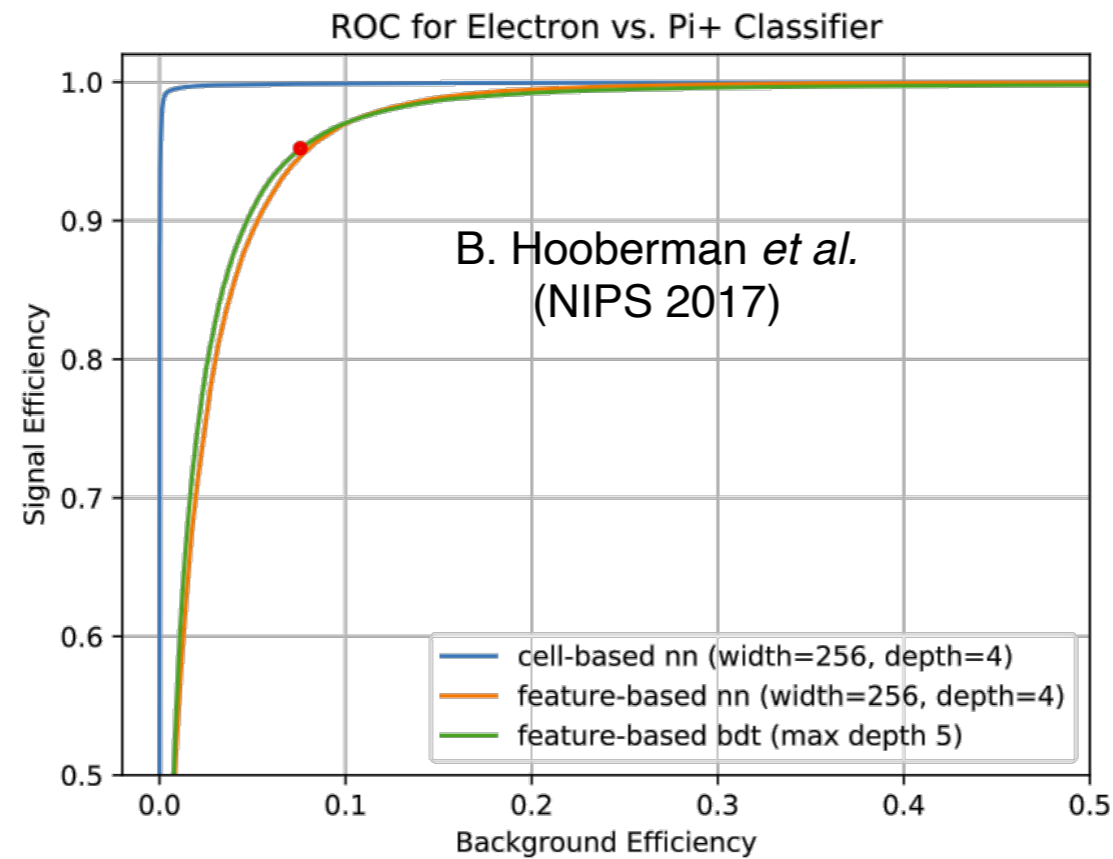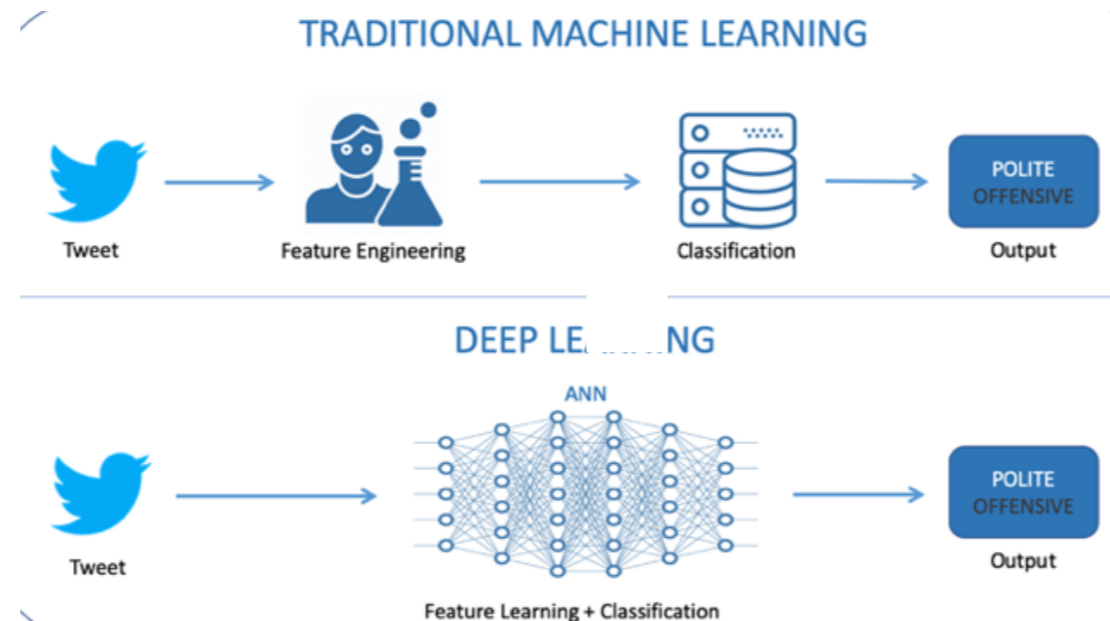# Deep Neural Networks

- "Deep learning allows computational models that are composed of multiple processing layers to learn **representations of data with multiple levels of abstraction**.

- …

- It discovers intricate structures in large data sets by using the back-propagation algorithm to indicate how a machine should change its internal parameters that are used to **compute the representation in each layer from the representation in the previous layer**…"

LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521,** 436–444 (2015).

# Machine Learning …Deep Learning



## TRADITIONAL MACHINE LEARNING

Tweet → Feature Engineering → Classification → Output (POLITE / OFFENSIVE)

## DEEP LEARNING

ANN

Tweet → Feature Learning + Classification → Output (POLITE / OFFENSIVE)

ROC for Electron vs. Pi+ Classifier

B. Hooberman *et al.*
(NIPS 2017)

Signal Efficiency

Background Efficiency

cell-based nn (width=256, depth=4)
feature-based nn (width=256, depth=4)
feature-based bdt (max depth 5)

Sofia Vallecorsa, Ilaria Luise CERN - sofia.vallecorsa@cern.ch | ilaria.luise@cern.ch

# A bit of history

# What you should already know…

- Stochastic Gradient Descent and Optimisers
- The problem of vanishing gradients
- Tips and tricks for training and the importance of data
- Regularisation strategies
- Example architectures : Convolutional Neural Networks
- 

- Small gradients slow down stochastic gradient descent
  - Limits ability to learn
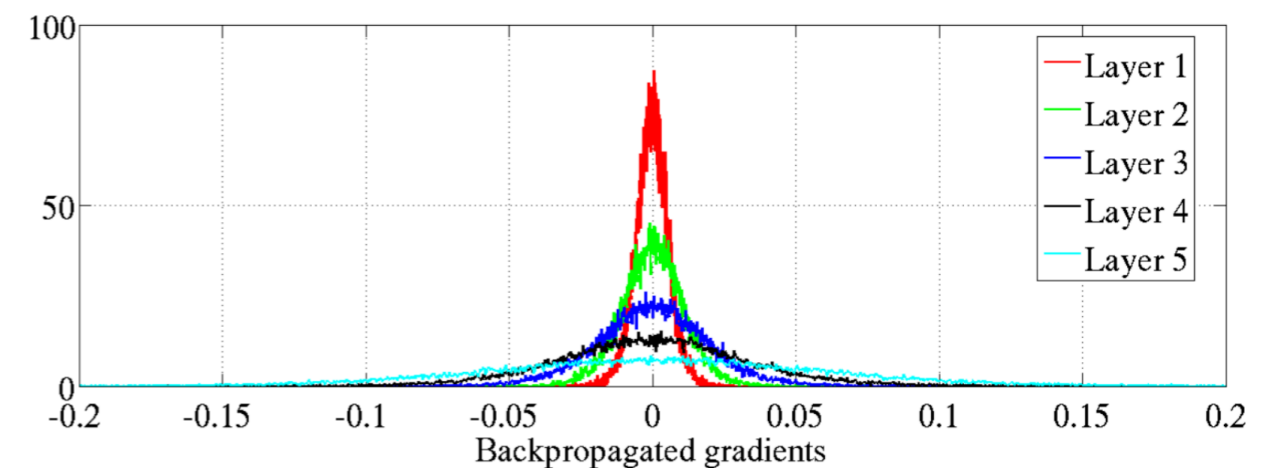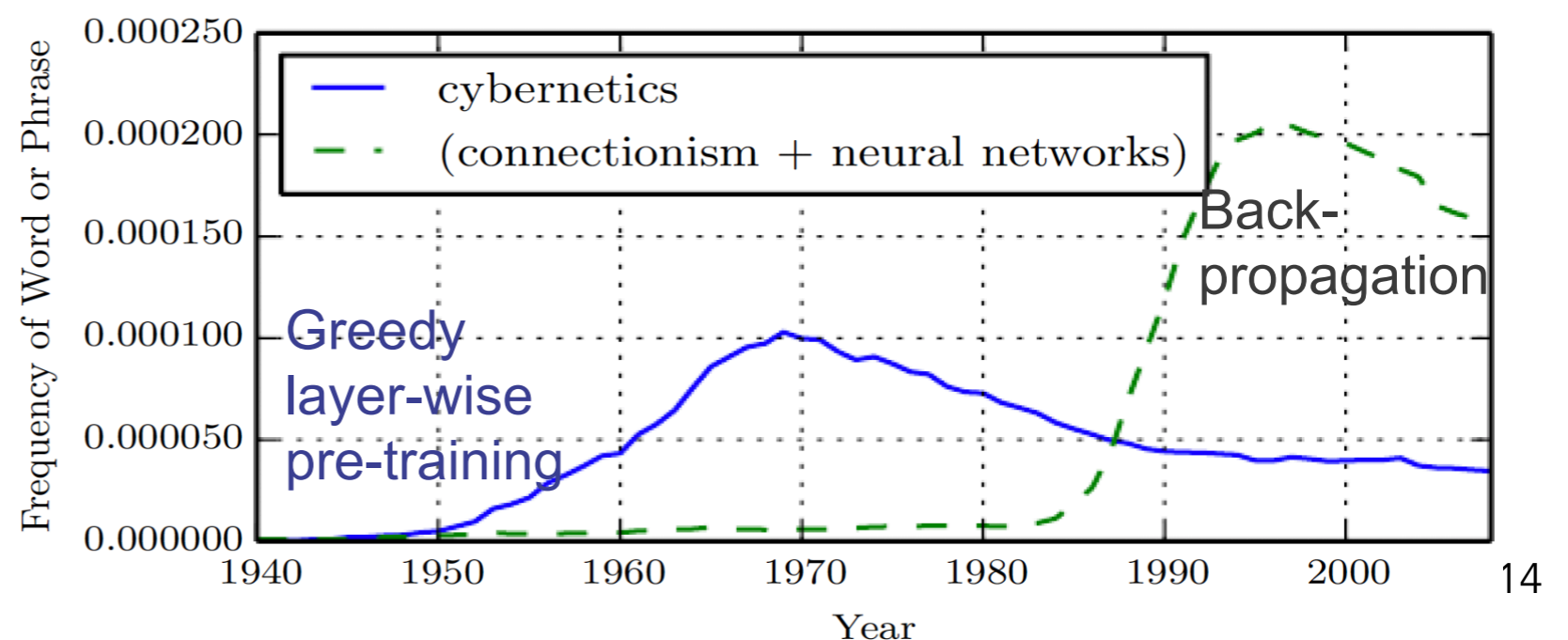- Gradients for layers far from the output vanish to zero.



Image from I. GoodFellow, Y. Bengio, A. Courville, "**Deep Learning**"

# AAAI 2020: Turing Award Keynote

*Deep Learning: more than just a deeper NN*



AAAI 20 keynotes Turing Award Winners (Geoff Hinton Yann Le Cunn, Yoshua Bengio): https://www.youtube.com/watch?v=UX8OubxsY8w

# Different primitives for different data representations

Perceptrons and MLP

Convolutions

Graphs

Recurrent Units (and LSTMs)

Point Cloud

...

# Then.. AI TakeOff….

Ilaria Luise, Sofia Vallecorsa CERN - ilaria.luise@cern.ch I sofia.vallecorsa@cern.ch

# Dataset sizes

| Domain | Data points |
|---|---|
| Vision | #Images (eg: a model trained on 3B images has a dataset size of 3B) |
| Language | #Words (eg: a model trained on 1T English tokens has a dataset size of ~750B words, the exact quantity depends on the tokenization) |

**Training datasets for language (left) and vision (right)**

18

Ilaria Luise,  Sofia Vallecorsa CERN - ilaria.luise@cern.ch I sofia.vallecorsa@cern.ch

# Machine learning at scale, for science

**Machine learning has been proven a very good tool to:**

- Extract information from (very large) datasets
- Efficiently analyse very large amounts of data
- Easily handle data from different sources
- Scalability to HPC environments

*Observation based datasets in physics are comparable or larger than these!*



*Can we use these tools for fully data-driven science?*

# Scientific opportunities

**Multi-scale dependencies:**

- **Model complex higher-order, statistical relationships between observations, fields, …**
- improve current simulations

**Compact representations:**

- **Condense dataset information in a compact representation**
- eg. condense the info in a few GB rather than TB

**Multi-source models:**

- **Enable multimodal and multi-source learning**
- eg. build models based on scientific data, GDP, birth rate etc..

**New discoveries:**

- **Explore the potential of unsupervised learning to extract new information directly from data**
- Learn unknown correlation patterns

Fundamental science

ML          HPC

20

# Transfer learning, pre-training, fine-tuning

*"**Transfer learning** and **domain adaptation** refer to the situation where what has been learned in one setting … is exploited to improve generalization in another setting"*

*Deep Learning, 2016.*



- Transferring knowledge to similar task

- Can be used to train large models

- Keep/modify pre-trained model?
  - CNN features are more generic in early layers and more dataset-specific in later layers

- Ex. **Flood detection in satellite images using U-Net**

- The **pre-training/fine-tuning** strategy has become key to the development of foundation models

(More on this later)



Nemni, Edoardo, et al., *Remote Sensing* 12.16 (2020): 2532.

# Accelerating training

- **Data parallelism**
  - Compute gradients on several batches independently
  - Update the model synchronously or asynchronously
- **Model Parallelism**, **Hybrid techniques**
- **Reduced precision** (INT8, BF16, …)
- Extreme parallelism using **combined strategies** and SGD algorithm optimization. Ex.
  - DeepSpeed and ZeRO-2 on Microsoft Azure



https://www.microsoft.com/en-us/research/blog/deepspeed-extreme-scale-model-training-for-everyone/

Sofia Vallecorsa, Ilaria Luise CERN - sofia.vallecorsa@cern.ch | ilaria.luise@cern.ch

# Sustainable AI

- AI inference more **energy efficient** than classical algorithms
- Energy cost of **AI training** can be high
- The community is defining **best practices**[1]
  - **Efficient AI architectures** can reduce computation by 3x–10x.
  - **AI-optimized processors vs general-purpose** can improve energy efficiency by 2x–5x[2].
  - **Cloud computing vs  on-prem** reduces energy usage by 1.4x–2x
- **Efficient training strategies**
  - Self-supervision, few-short learning, pre-training



Carbon Reduction 2017-2021 by Improving the 4Ms

Transformer model trained on P100 GPUs in an average data center in 2017

747 — Map⇒Oklahoma
83 — Mechanization⇒Google DC
57 — Machine⇒TPUv4
4 — Model⇒Primer



R. Cardoso, vCHEP2021 [2]

**2 Cardoso, Renato, et al. "Accelerating GAN training using highly parallel hardware on public cloud." EPJ Web of Conferences. Vol. 251. EDP Sciences, 2021.**

23

Sofia Vallecorsa, Ilaria Luise CERN - sofia.vallecorsa@cern.ch | ilaria.luise@cern.ch

# Generative Models



R. Feynman

# Generative models

The problem:

Assume data sample follows $p_{data}$ distribution

Can we draw samples from distribution $p_{model}$ such that $p_{model} \approx p_{data}$?

# Generative models

The problem:

Assume data sample follows $p_{data}$ distribution

Can we draw samples from distribution $p_{model}$ such that $p_{model} \approx p_{data}$?

**Maximum Likelihood Estimator:**
- Assume some form for $p_{model}$ (prior knowledge, parameterized by θ)
- draw samples from $p\theta_*$

$$\theta^* = \arg\max_{\theta} \sum_{\mathbf{x} \in \mathcal{D}} \log(p_{model}(\mathbf{x}; \theta))$$

Generative models don't look for mathematical expression of $p_{model}$

Train NN as a generator $\mathscr{G} : \mathbb{R}^m \rightarrow \mathbb{R}^n$

that maps samples from a tractable distribution supported in $\mathbb{R}^m$ to points in $\mathbb{R}^n$

31

26

# Latent Representation

Modeling Data and Meaningful Degrees of Freedom
53

Latent space $\mathscr{F}$

Original space $\mathscr{X}$

Common Approach to Deep Generative Modeling:
Choose known distribution for latent space and learn map to data space

Fleuret, Deep Learning Course

- Information content is preserved within a **hidden manifold with lower dimension**
- Can manipulate **latent space** (style specification, hypothesis testing directly in data, …)
- Can optimise latent representation according to a specific task (**guided compression**)
- Can help with **multi-modality**

**NB: Problems exhibiting complex symmetries may benefit from latent space representations connected to the specific underlying symmetry group!**

27

Ilaria Luise, Sofia Vallecorsa CERN - ilaria.luise@cern.ch I sofia.vallecorsa@cern.ch

# Deep Generative Models

Deep models allow **higher levels of abstractions** and **improve generalization** wrt to shallow models



(a) Autoregressive Models
(b) Variational Autoencoders (VAEs)
(c) Normalising Flows (NFs)
(d) Energy-based Models (EBMs)
(e) Generative Adversarial Networks (GANs)

Current Opinion in Structural Biology

See Danilo Rezende tutorial on Deep Generative Models

28

Ilaria Luise,  Sofia Vallecorsa CERN - ilaria.luise@cern.ch I sofia.vallecorsa@cern.ch

# Fully Observed Models

*Directly observe data without introducing new local (latent) variables*

$$p(x_{1,...,N}) = \prod_{i=1}^{N} p(x_i | x_{1,...,(i-1)})$$



Ex. Pixel Recurrent Neural Networks:



$$p(x_t | x_{1:t-1}) = p(x_t^{red} | x_{1:t-1}) p(x_t^{green} | x_{1:t-1}, x_t^{red}) p(x_t^{blue} | x_{1:t-1}, x_t^{red}, x_t^{green})$$

29

Sofia Vallecorsa, Ilaria Luise CERN - sofia.vallecorsa@cern.ch | ilaria.luise@cern.ch

# Auto-Encoders

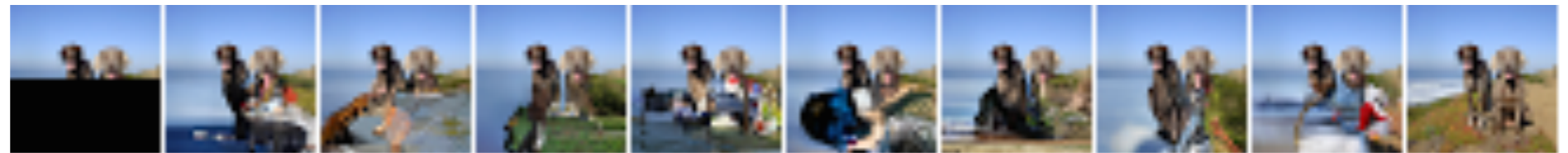Examples of latent variables models (and implicit..)

$$p(z)$$

$$z$$

## Ex. Auto-Encoder

$$x \in \mathbb{R}^{d_x} \quad z \in \mathbb{R}^{d_z} \quad \theta \in \mathbb{R}^{d_\theta}$$
$$\mathcal{D} = \{x_i\} \quad i \in \{1, ..., N\}$$



| | Encoder Network (conv) | Compress | Decoder Network (deconv) | |

$$x$$

$$p(x|z)$$

## Ex. Variational Auto-Encoder

**Explicit constraints** on encoded representations (learn the **latent variable distribution**)

Two components in the loss function (**reconstruction loss and KL divergence** to constrain latent to prior)

Ilaria Luise, Sofia Vallecorsa CERN - ilaria.luise@cern.ch I sofia.vallecorsa@cern.ch

# Likelihood-free learning

*Density estimation by comparison*

Sample-based comparison between **estimated** q(x) and **true distribution** p(x)

- Build **auxiliary model** to indicate how data simulated from the generative model differs from observed data.

- **Adjust model parameters** to better match the data distribution



D: Detective

R: Real Data

G: Generator (Forger)

I: Input for Generator

Ex. Generative Adversarial Networks
(I. Goodfellow 2014)

Sofia Vallecorsa, Ilaria Luise CERN - sofia.vallecorsa@cern.ch | ilaria.luise@cern.ch

# Diffusion models

- **Parametrized Markov Chains** trained using variational inference to produce samples matching the data after finite time.
  - Chain transitions are **reverse diffusions** (gradually adding noise to the data)
- Ex. DDPM (Diffusion Denoising Probabilistic Models) based on U-Net architecture, https://arxiv.org/pdf/2006.11239.pdf:
  - Iteratively add Gaussian noise to input image, eventually reaching pure noise
  - Generation process **inverts the diffusion:** start from pure noise sample, then iteratively de-noise it.

# Normalizing Flows

*Explicit density estimation*

**Bijective, differentiable** maps between two continuous variables

- Compositional

$$x = g(z) = g_n \bullet \cdots \bullet g_2 \bullet g_1(z)$$

- Simple prior density to complex target

$$\ln p(x) = \ln q(z) - \sum_i \ln \left| \det\left( \frac{\partial g_{i+1}}{\partial g_i} \right) \right|$$

**Generation**

**Inference**

$\phi(z)$

Invertible & Tractable Jacobian

$z$

$x$

$p(z)$

$g_1$ ... $g_i$

$g_{i+1}$ ... $g_n$

$p(x)$

Many simple layers composed to produce $\phi$

Easily sampled

Approximates desired dist.

Image credit: G. Kanwa

33

Sofia Vallecorsa, Ilaria Luise CERN - sofia.vallecorsa@cern.ch | ilaria.luise@cern.ch
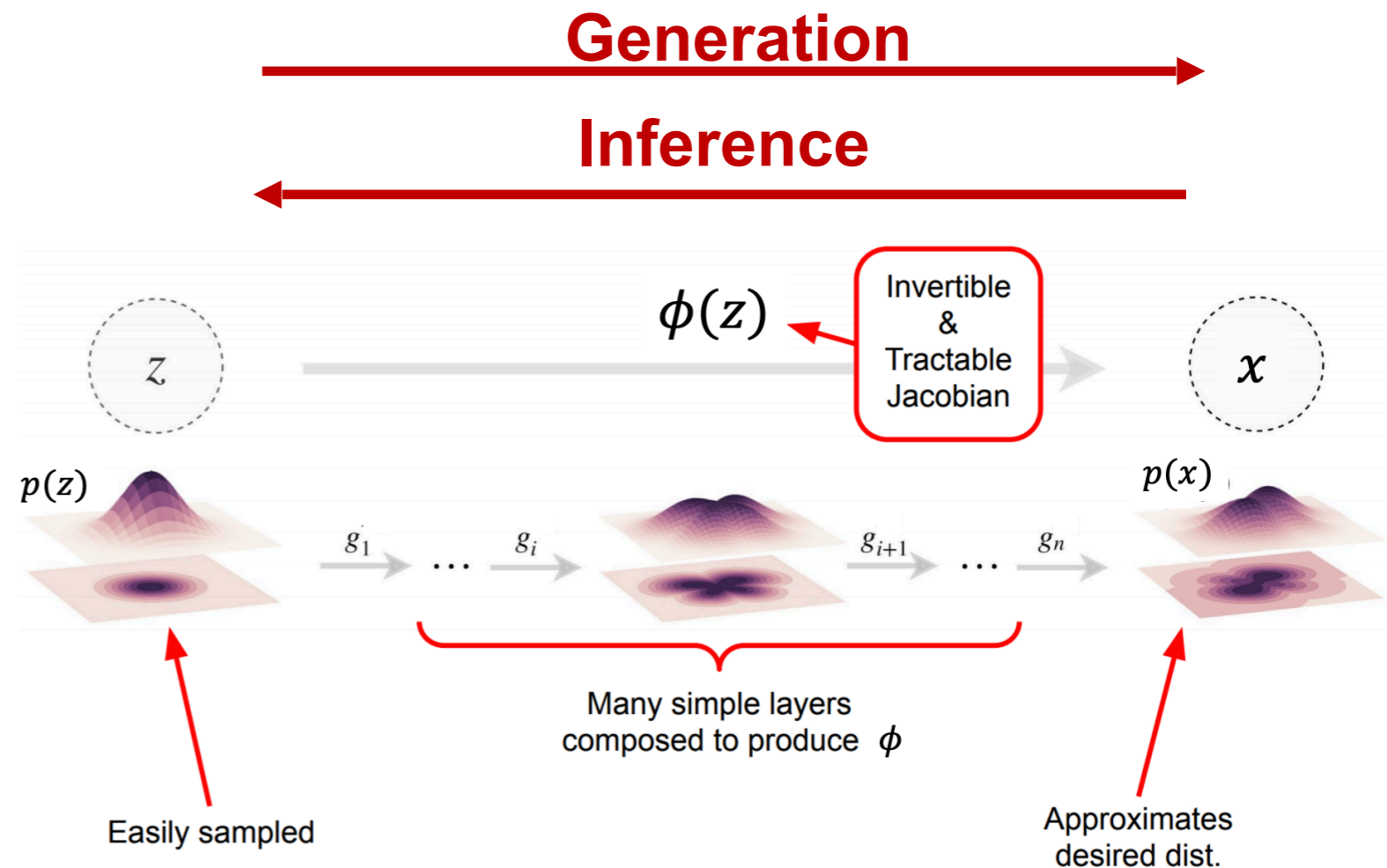
# Energy based models

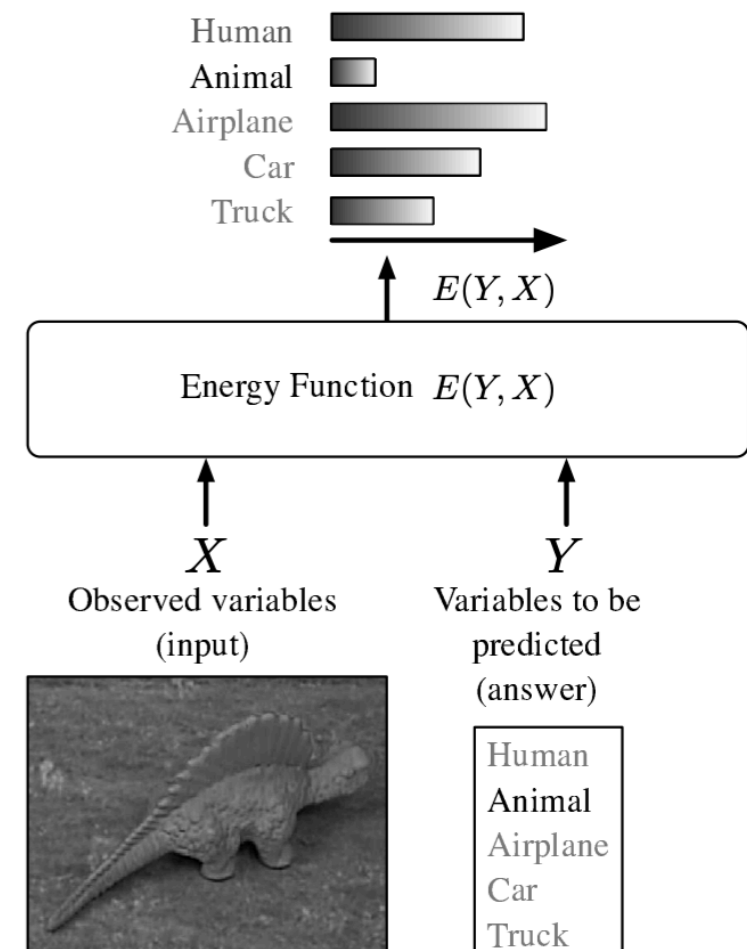*Model is an **energy function** measuring goodness of each (x,y) sample*

Inspired by **statistical mechanics**



- **Training**: Finding the best energy function E(W,X,Y).
- Minimise **loss functional** so that for any $X_i$, inference results in $Y_i$

$$\mathcal{L}(E, \mathcal{S}) = \frac{1}{P} \sum_{i=1}^{P} L(Y^i, E(W, \mathcal{Y}, X^i)) + R(W)$$

- **Inference:** strategy to find Y that minimizes E(X, Y) for classification, regression, generation
- Model combination can be tricky due to **E scale**
- Interpret E function as PDF (**Gibbs distribution**)

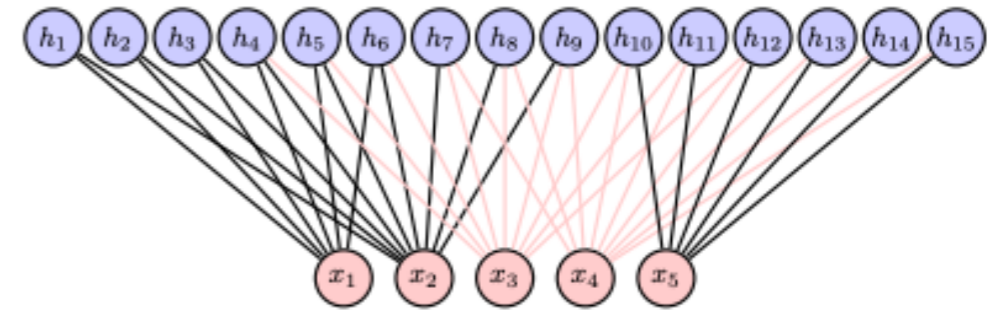$$P(Y|X) = \frac{e^{-\beta E(Y,X)}}{\int_{y \in \mathcal{Y}} e^{-\beta E(y,X)}}$$

http://yann.lecun.com/exdb/publis/pdf/lecun-06.pdf

34

# Variational Calculation and Boltzmann Machines

*Probability is a Boltzmann distribution*



Ex. Compute **expected value** of physical observable
- Statistical mechanics defines a probability function

- Minimize the **free energy** -ln $\mathcal{Z}$ (intractable in general) by defining its variational form for a normalized variational probability q(x)

$$\pi(x) = \frac{e^{-E(x)}}{\sum_x e^{-E(x)}} \qquad \mathcal{Z} = \sum_x e^{-E(x)}$$

- L is an upper bound for the physical free energy $- \ln \mathcal{Z}$

- Approximation is exact when variational distribution approaches the target

$$L = \sum_x q(x) \ln \frac{q(x)}{e^{-E(x)}} = \left\langle E(x) + \ln q(x) \right\rangle_{x \sim q(x)}$$

$$L + \ln Z = KL(q||\pi) \geq 0$$

Sofia Vallecorsa, Ilaria Luise CERN - sofia.vallecorsa@cern.ch | ilaria.luise@cern.ch

# Attention and Transformers

# A step back

## Recurrent States
29

- Input sequence $x \in S(\mathbb{R}^m)$ of *variable* length $T(x)$

- Recurrent model maintain a **recurrent state** $\boldsymbol{h}_t \in \mathbb{R}^q$ updated at each time step $t$. For $t = 1, \dots, T(x)$:

$$\boldsymbol{h}_{t+1} = \phi(\boldsymbol{x}_t, \boldsymbol{h}_t; \theta)$$

- Simplest model:

$$\phi(\boldsymbol{x}_t, \boldsymbol{h}_t; W, U) = \sigma(W\boldsymbol{x}_t + U\boldsymbol{h}_t)$$

- Predictions can be made at any time $t$ from the recurrent state

$$\boldsymbol{y}_t = \psi(\boldsymbol{h}_t; \theta)$$

Credit: F. Fleuret

## Recurrent Networks:



Credit: Fleuret

## LSTMs:



Credit: G. Louppe
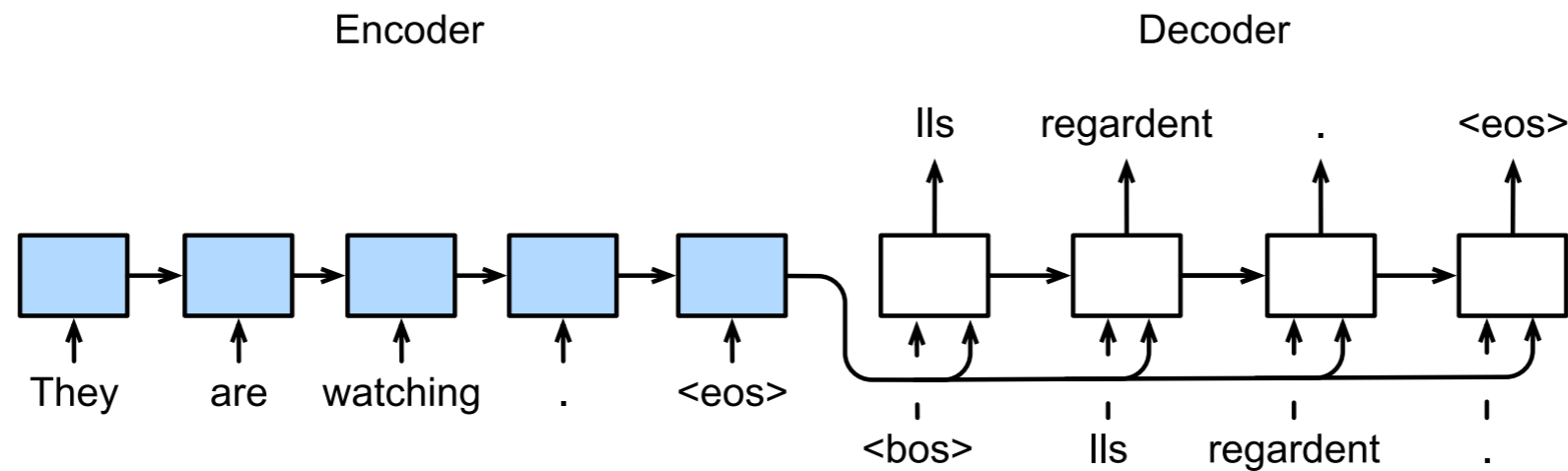
37

# Seq2seq models



Credit: d2l.ai

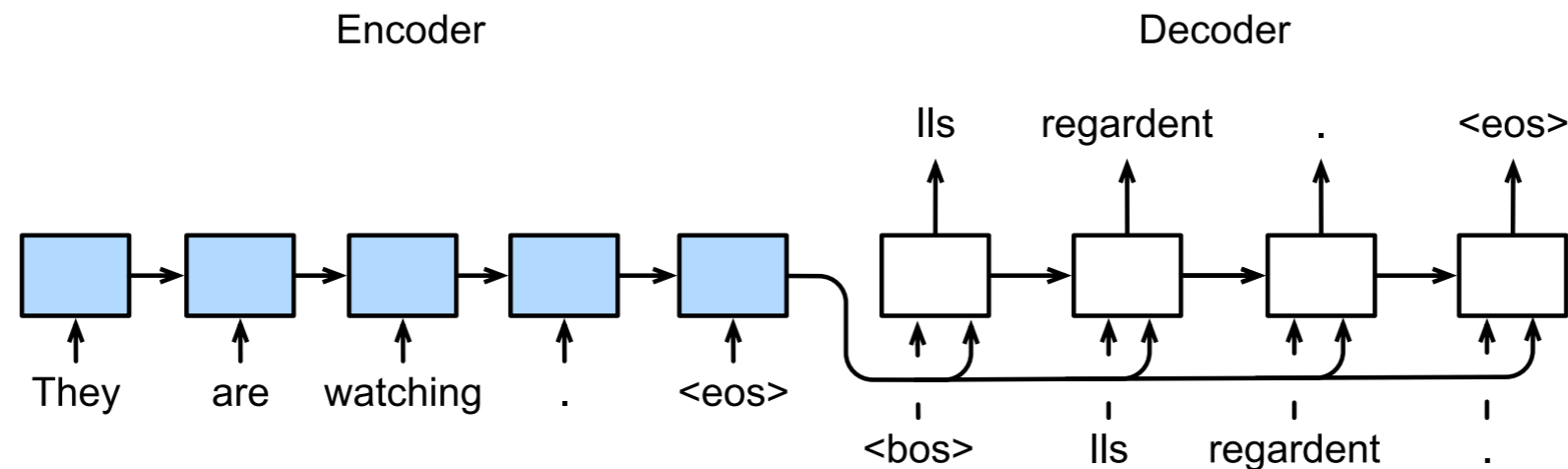**Seq2seq models analyse sequences**

Predict probability distributions of the next token given previous context

Encoder compresses the sequence in a fixed size vector

**Fixed size latent vector is a bottleneck**

Decoder **next-step generation is suboptimal** since latent vector contains the same information

Ilaria Luise,  Sofia Vallecorsa CERN - ilaria.luise@cern.ch I sofia.vallecorsa@cern.ch

# Information bottleneck requires attention

Encoder

Decoder

Ils    regardent    .    <eos>

They    are    watching    .    <eos>

<bos>    Ils    regardent    .

Credit: d2l.ai

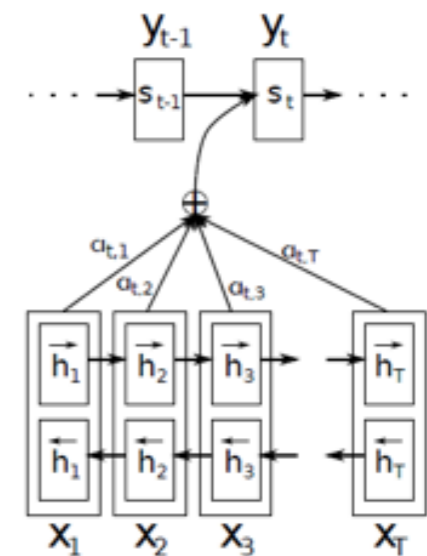Can we avoid compression and pass the decoder entire input?

Need a mechanism to **focus on most relevant** input tokens at each prediction step

Introduce **softmax to calculate probability** (maintain differentiable architecture)

Output is **independent of the order** of input examples (set instead of sequences)

Use **relationships between input elements** (as graph representation).

Attention mechanism as originally formulated in a bi-directional LSTM Auto-Encoder
https://arxiv.org/abs/1409.0473

$y_{t-1}$    $y_t$

$s_{t-1}$    $s_t$

$\alpha_{t,1}$    $\alpha_{t,T}$

$\alpha_{t,2}$    $\alpha_{t,3}$

$\overrightarrow{h_1}$    $\overrightarrow{h_2}$    $\overrightarrow{h_3}$    $\overrightarrow{h_T}$

$\overleftarrow{h_1}$    $\overleftarrow{h_2}$    $\overleftarrow{h_3}$    $\overleftarrow{h_T}$

$x_1$    $x_2$    $x_3$    $x_T$

Ilaria Luise,  Sofia Vallecorsa CERN - ilaria.luise@cern.ch I sofia.vallecorsa@cern.ch

# Attention mechanism

**A key-value database** (differentiable, entries are continues vectors):

$$Q = \{q_1, q_2, \ldots, q_m\} \quad \text{QUERIES}$$
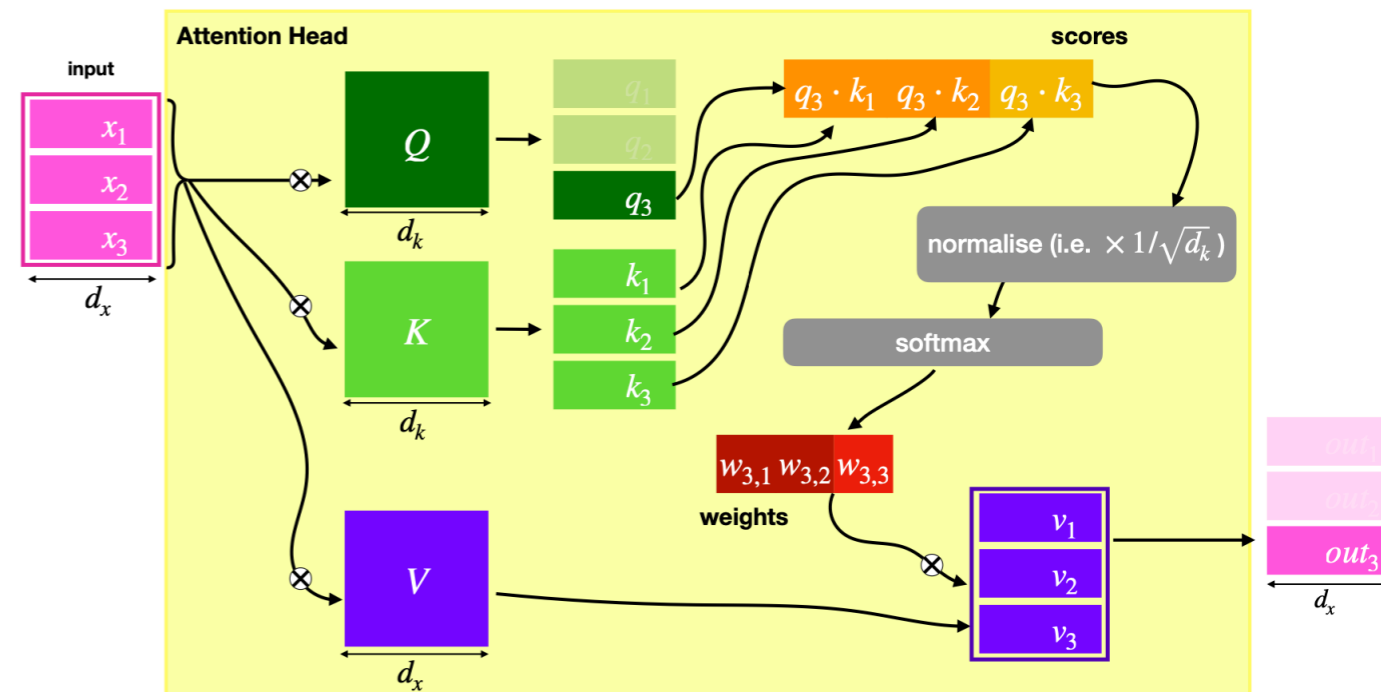$$K = \{k_1, k_2, \ldots, k_n\} \quad \text{KEYS}$$
$$V = \{v_1, v_2, \ldots, v_n\} \quad \text{VALUES}$$

A normalised **similarity** function between query-key pairs:

$$S_{ij} = \text{SIMILARITY}(q_i, k_j) \qquad A_{ij} = \text{NORMALIZE}(S_{ij}) = \frac{e^{S_{ij}}}{\sum_{l=1}^{n} e^{S_{il}}}$$

A **weighted average** over values {V}, based on similarity:

$$O_i = A_{ij} V^j$$



Credit: G. Weiss

**NB. Weights are probabilities (use softmax)**

$$\text{SIMILARITY}(q_i, k_j) = \frac{q_i \cdot k_j}{\sqrt{D}}$$

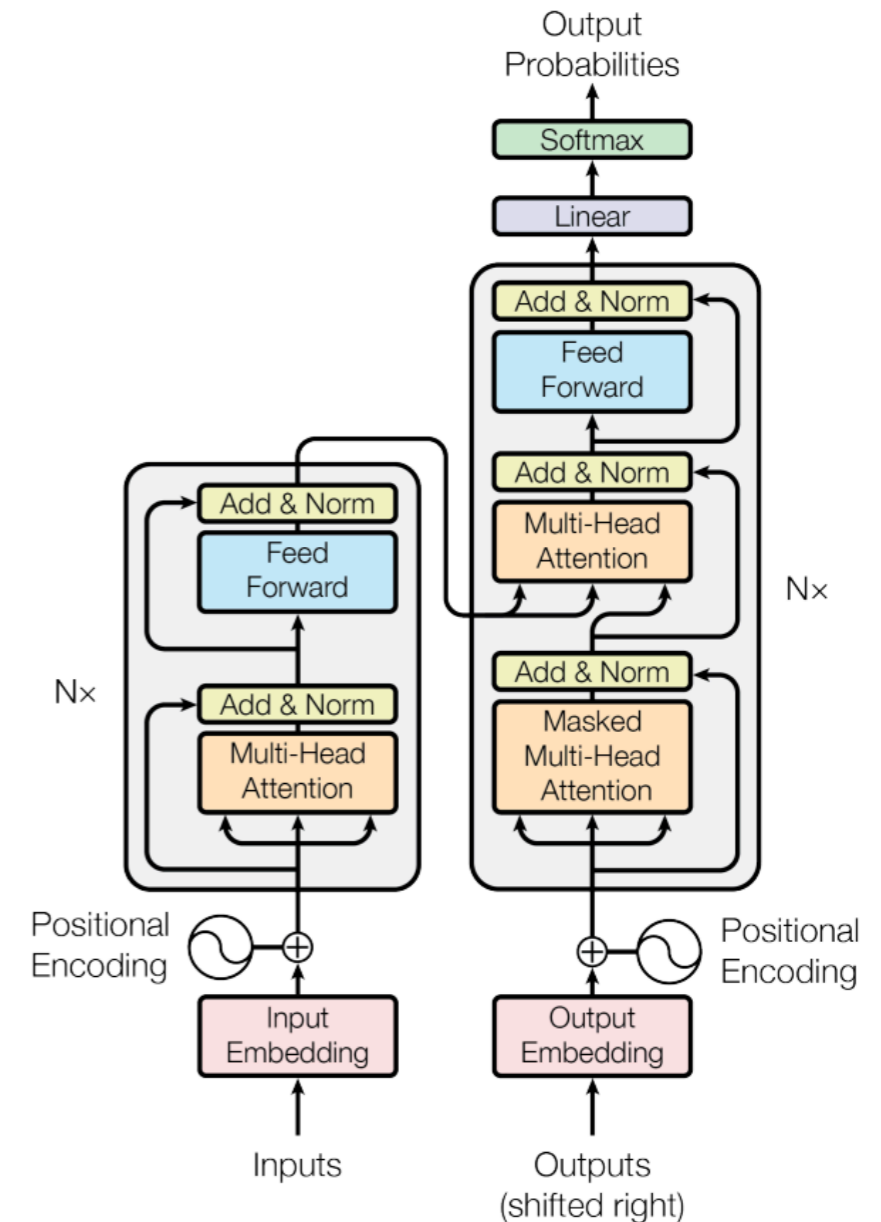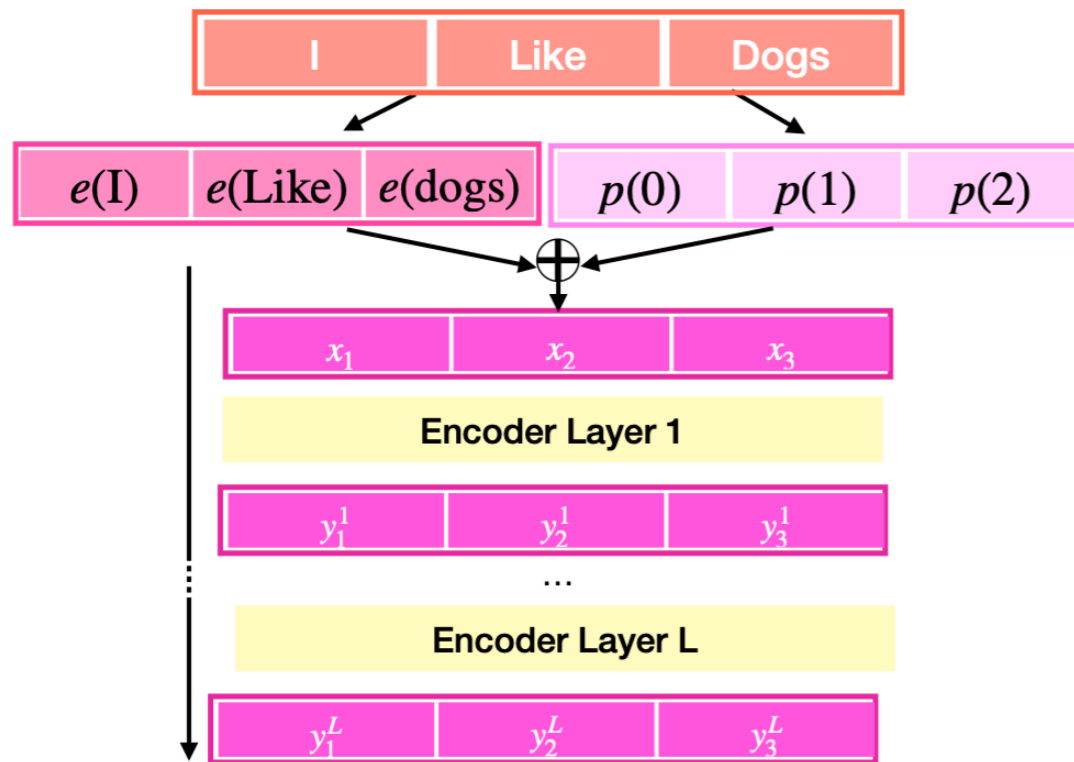**Self-attention** uses same input for values, keys and queries.
Focus on relationship between elements (adds context)

**Multi-head attention** splits input token in subgroups and processes them in parallel

**NB: Scaled dot-product is permutation equivariant**

40

Ilaria Luise, Sofia Vallecorsa CERN - ilaria.luise@cern.ch I sofia.vallecorsa@cern.ch

# Transformers

Transformer components include:

Multi Head **Attention**

**Normalisation** layers

Position Independent **Feed Forward Layers**

**Skip Connections**

**NB. All tokens are processed in parallel**

Vaswani et al., *Advances in Neural Information Processing Systems*, 2017, 5998–6008

41