



Machine Learning: A new toolbox for Theoretical Physics

Juan Rojo

VU Amsterdam & Theory group, Nikhef

D-ITP Advanced Topics in Theoretical Physics

11/12/2019

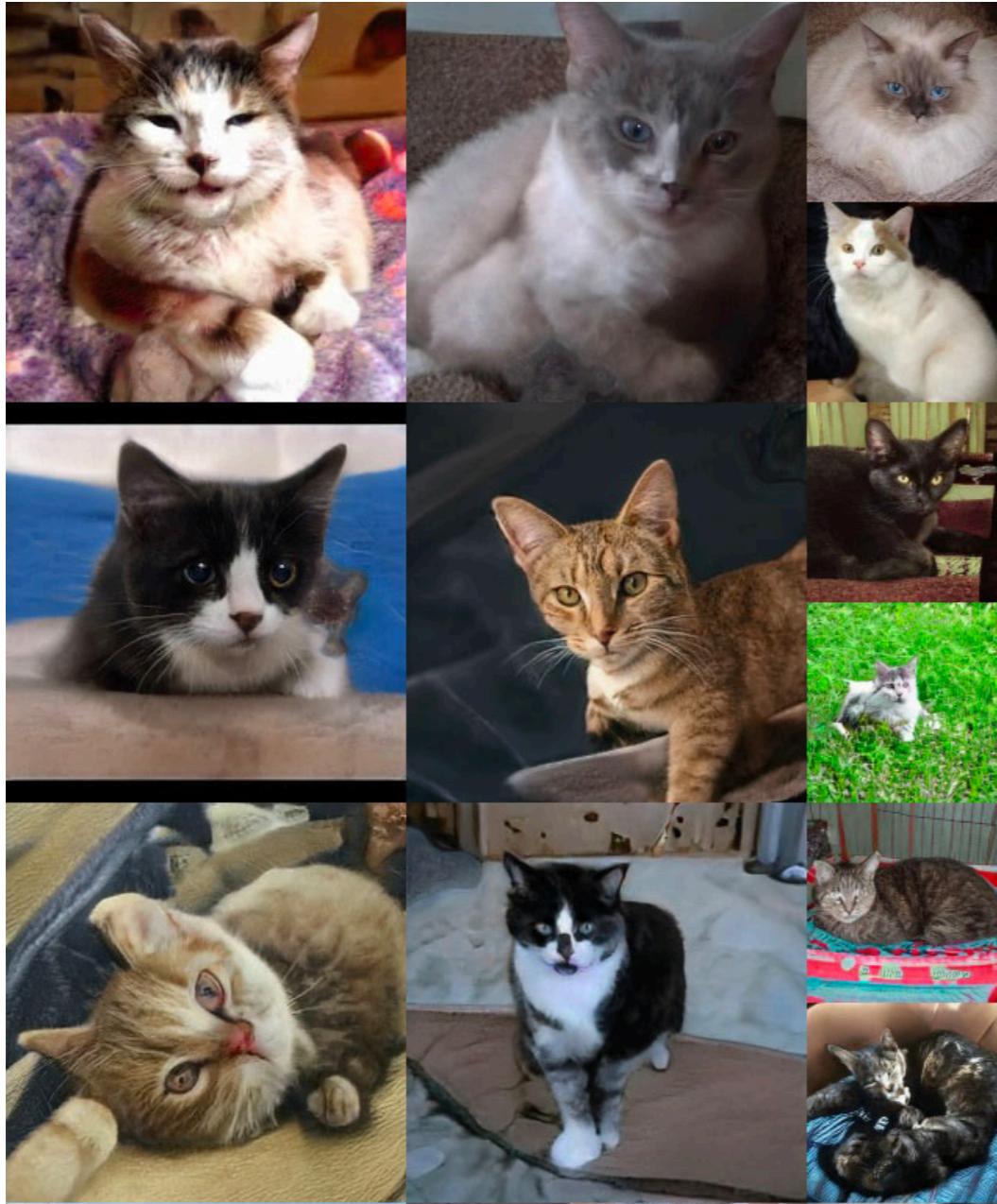
Today's lecture

- ✿ Convolutional Neural Networks
- ✿ Energy-based models and Boltzmann learning
- ✿ Generative Models and Adversarial Learning
- ✿ Applications of Machine Learning to Theoretical Physics

Convolutional Neural Networks

Convolutional Neural Networks

Like physical systems, many datasets and supervised learning tasks also possess additional **symmetries and structure** what can (and should) be exploited



e.g. we want to train a classifier to identify pictures of cats. What **high-level features** must one learn first?

Convolutional Neural Networks

Like physical systems, many datasets and supervised learning tasks also possess additional **symmetries and structure** what can (and should) be exploited



e.g. we want to train a classifier to identify pictures of cats. What **high-level features** must one learn first?

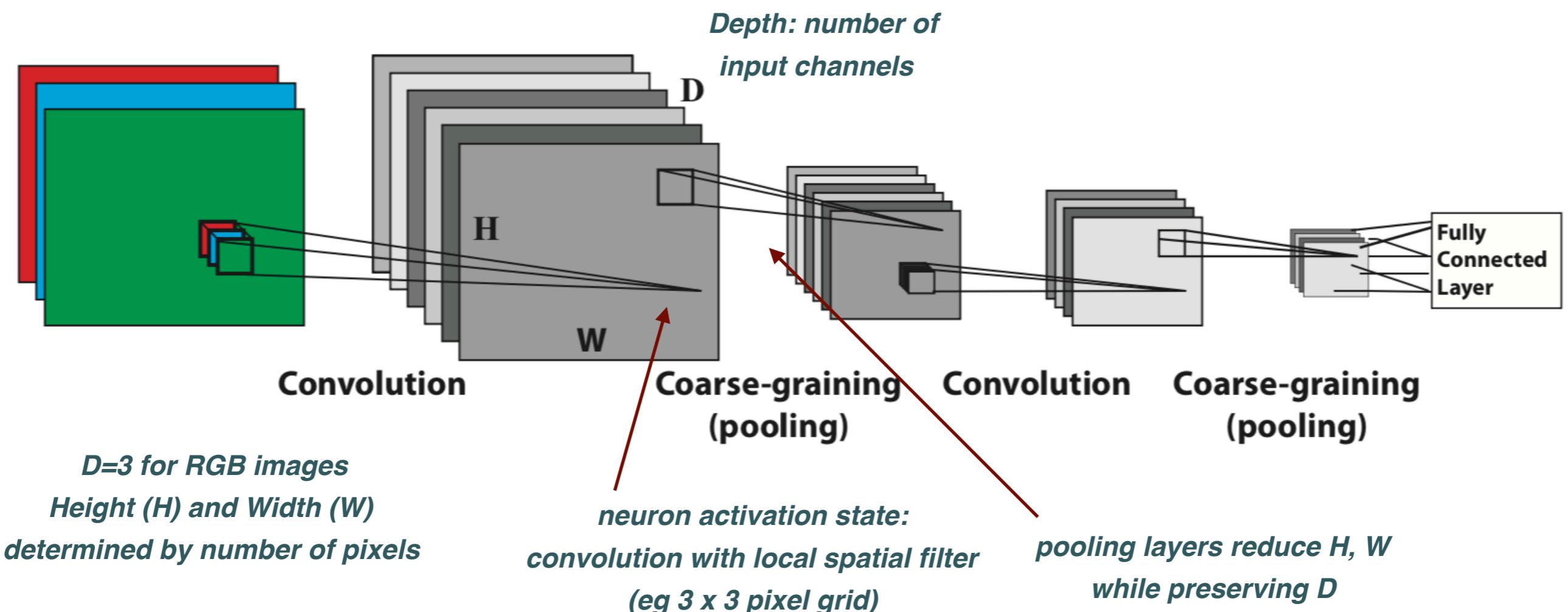
- ✿ *The features that define ``cat'' are local in the picture: whiskers, tail, paws ...: **locality***
- ✿ *Cats can be anywhere in the image: **translational invariance***
- ✿ *Relative position of features must be respected (eg whiskers and tail should appear in opposite sides of ``cat''): **rotational invariance***

Our classifier should exhibit all these high-level features

Convolutional Neural Networks

Convolutional Neural Networks (CNNs) are architectures that take **advantage of this additional high-level structures** that all-to-all coupled networks fail to exploit

A CNN is a translationally invariant neural network that respects locality of the input data



Convolutional Neural Networks

CNNs are composed by
two kinds of layers

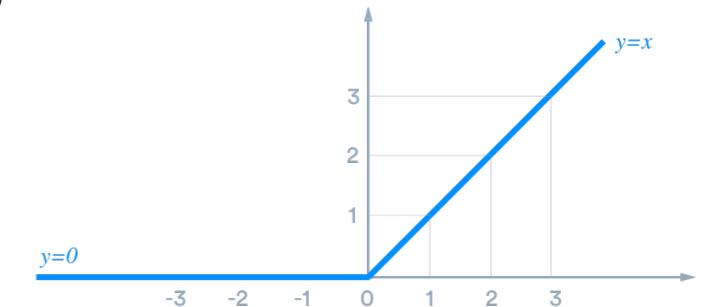
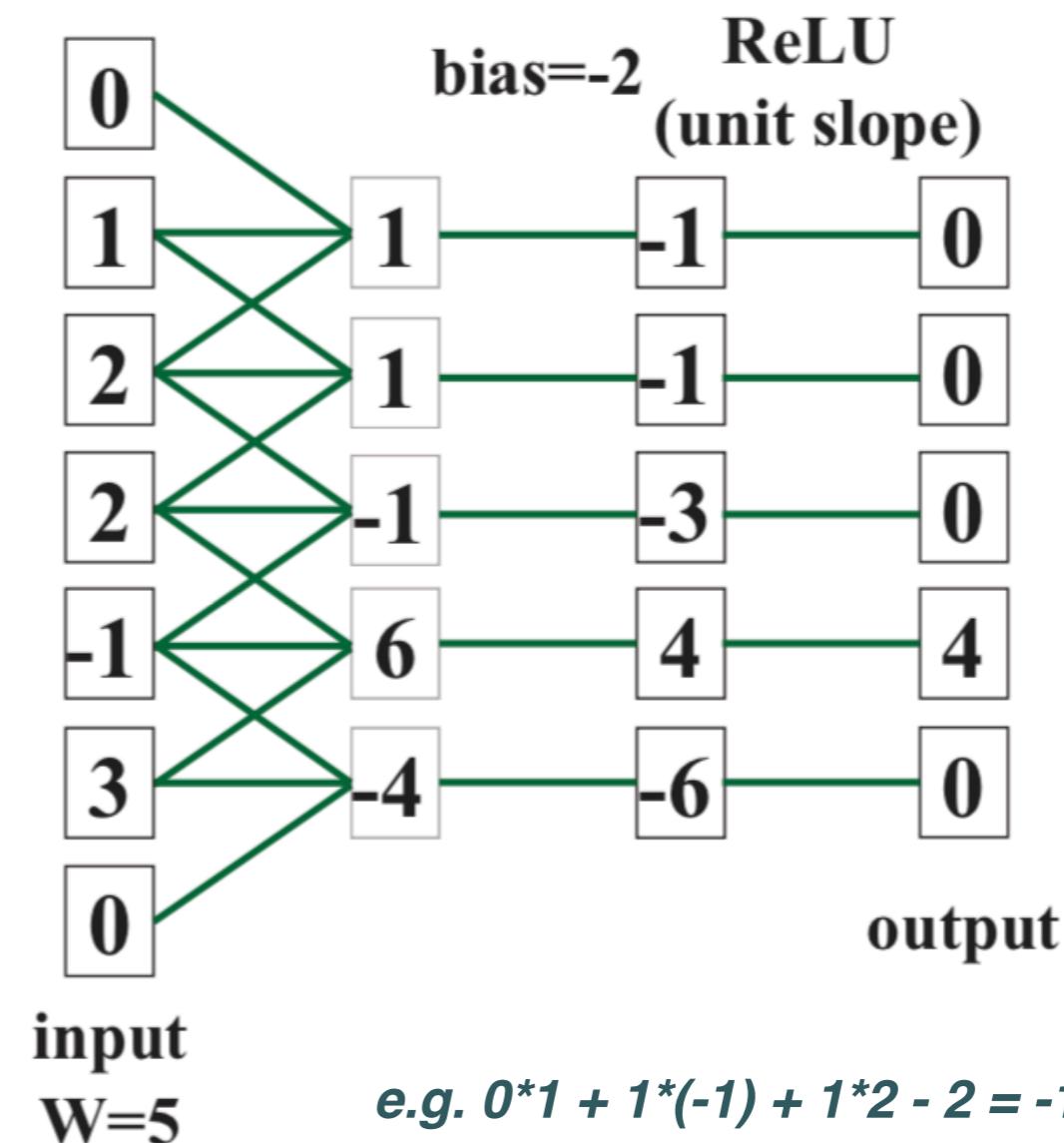
*example of
convolutional layer*

*note that convolution changes the depth,
but not the height and width of the network*

Convolution of input with **filters**

size-3 Filter

$F=3$
 $\text{weight}=[1, -1, 1]$



Convolutional Neural Networks

CNNs are composed by
two kinds of layers

Convolution layer of input with **filters**

Pooling layer that coarse-grains the input while
maintaining locality and spatial structure

e.g. **MaxPool**, the spatial dimensions are coarse-grained by replacing a small region by single neuron whose output is maximum value of the output in the region

in average pooling, one averages over output in region

12	20	30	0
8	12	2	0
34	70	37	4
112	100	25	12

max pooling

20	30
112	37

average pooling

13	8
79	20

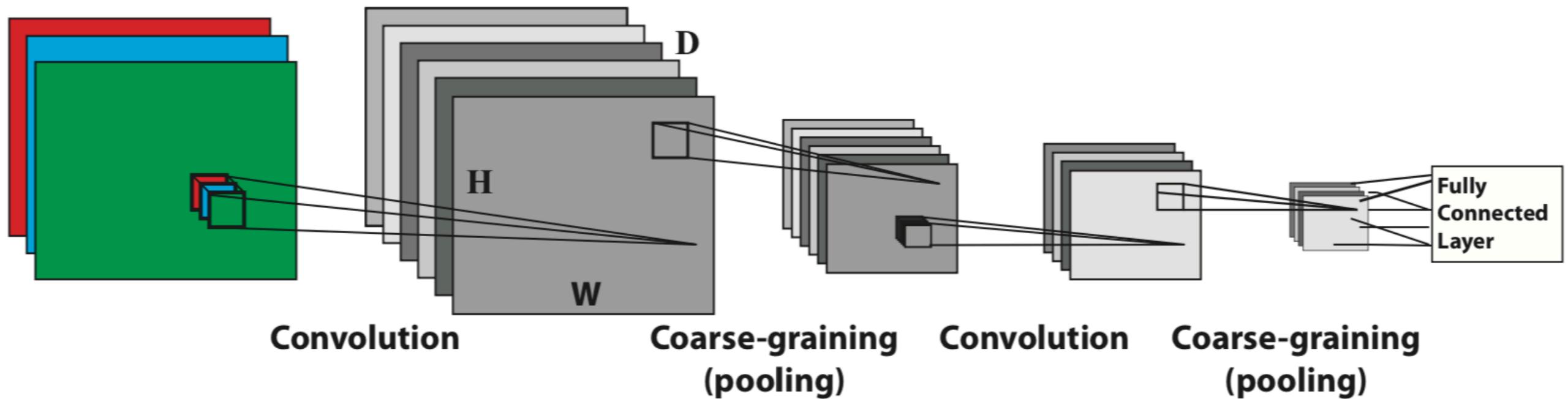
Convolutional Neural Networks

CNNs are composed by
two kinds of layers

Convolution layer of input with **filters**

Pooling layer that coarse-grains the input while
maintaining locality and spatial structure

the convolution and max-pool layers are followed by an **all-to-all connected layer and a high-level classifier**, so that one can train CNNs using the standard backpropagation algorithm



Convolutional Neural Networks

CNNs are composed by
two kinds of layers

Convolution layer of input with **filters**

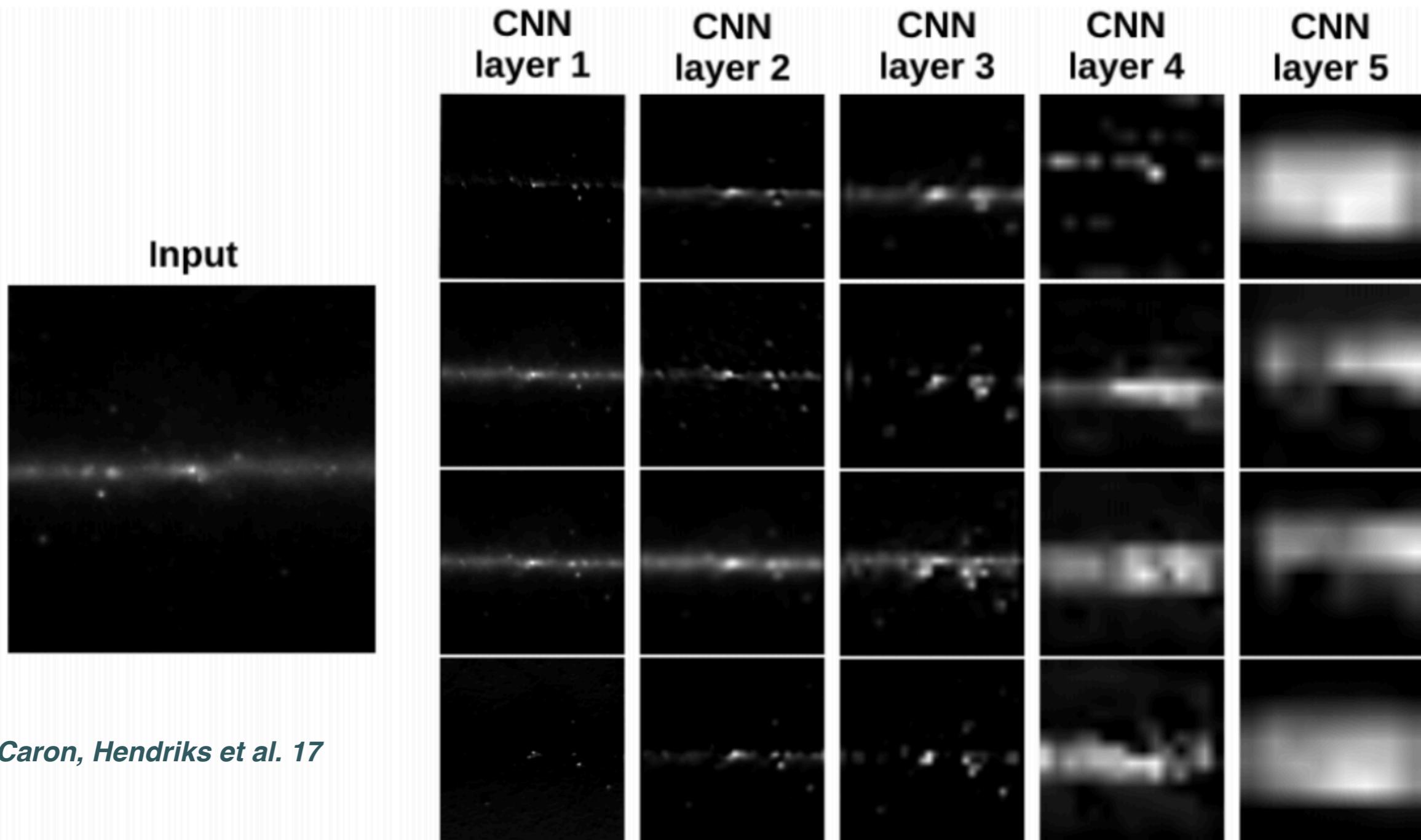
Pooling layer that coarse-grains the input while
maintaining locality and spatial structure

the convolution and max-pool layers are followed by an **all-to-all connected layer and a high-level classifier**, so that one can train CNNs using the standard backpropagation algorithm

note that only problems characterised by a **spatial locality** are amenable to CNNs:
for example the 2D Ising dataset can be studied with CNNs, but not the SUSY dataset

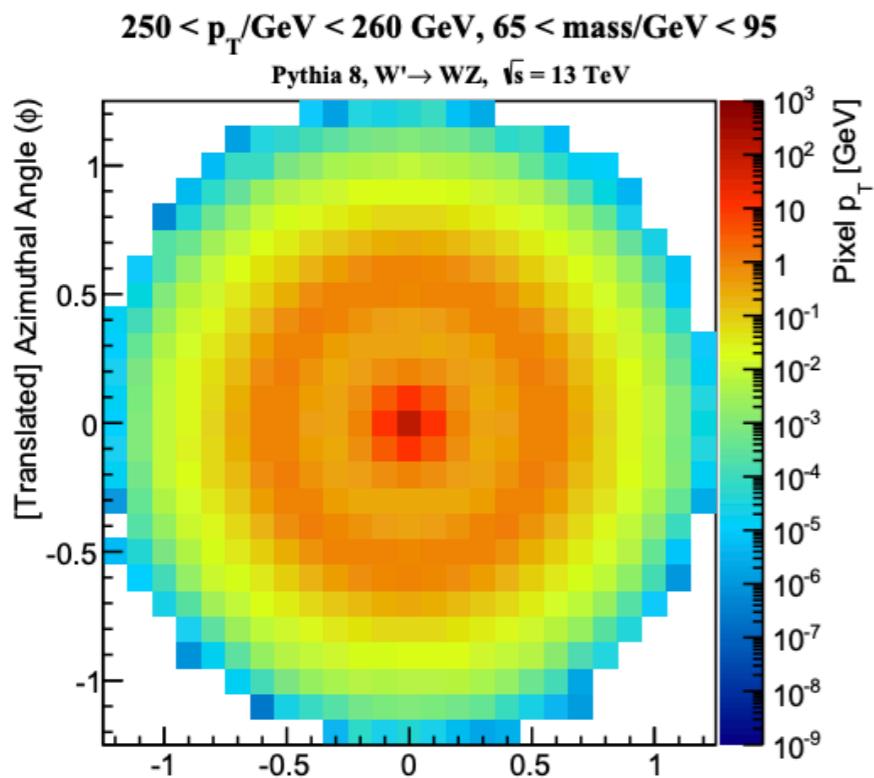
Application to Dark Matter searches

Use CNNs to discriminate between **point sources** (astrophysical origin) and **diffuse flux** (dark matter) in galactic centre images, to compare with predictions of DM models

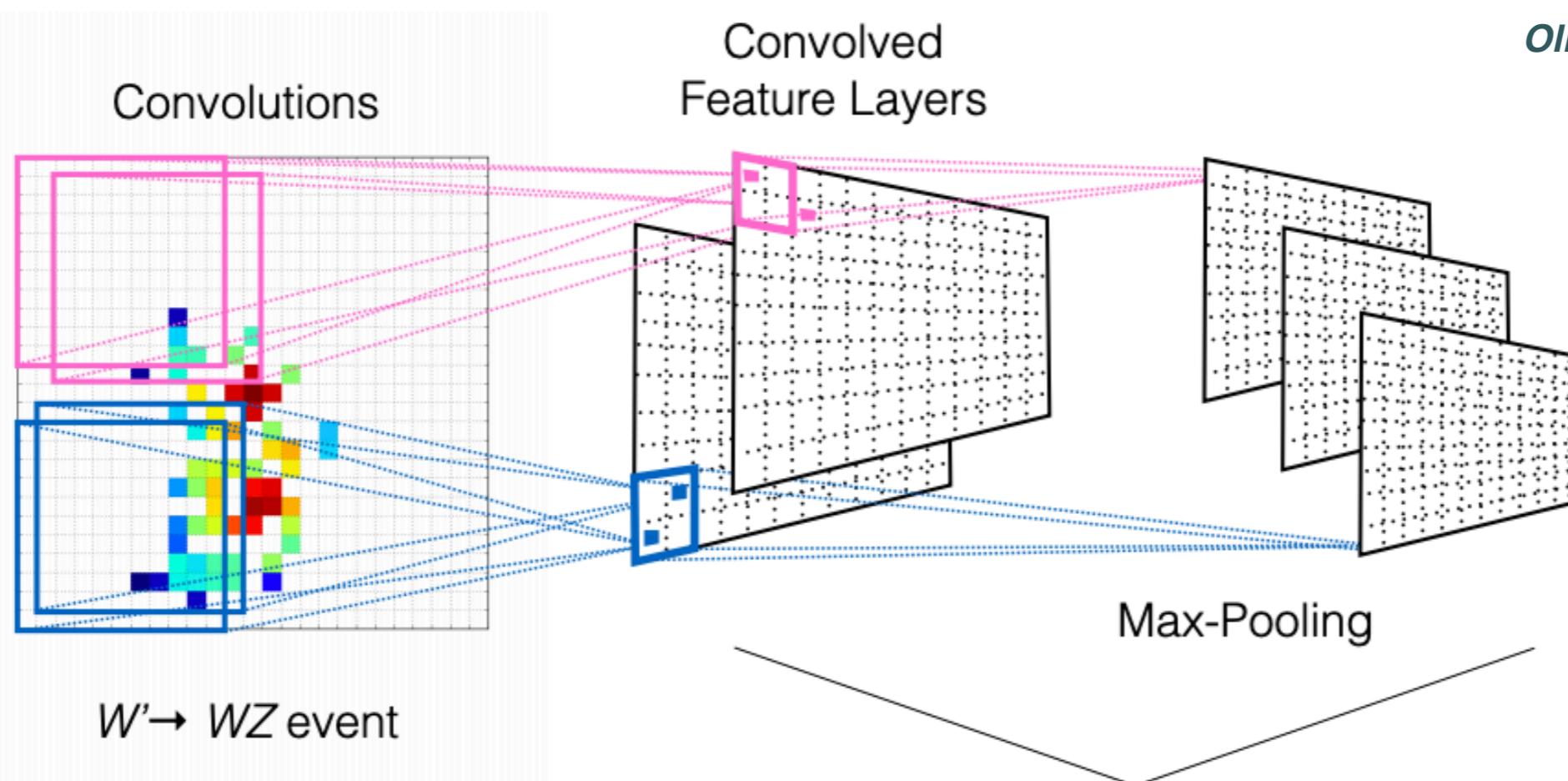
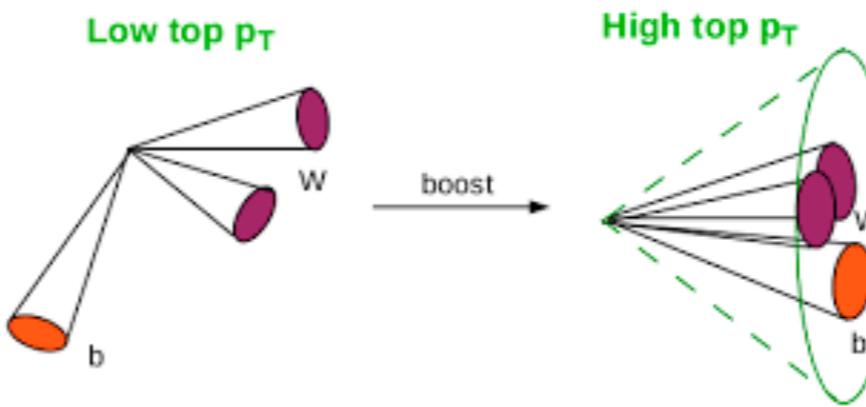


Caron, Hendriks et al. 17

Application to LHC physics



Train CNNs on jet **images** (energy deposits in detector) to discriminate between signal and background events



Energy-Based Models and Boltzmann Learning

Generative Models

Most ML models discussed here (Supervised NNs, logistic regression, ensemble models) are **discriminative**: designed to identify **differences between groups of data**

e.g. cats vs dogs discrimination

these models cannot carry some tasks such as **drawing new examples** from an unknown probability distribution: for this we need **generative models**

*e.g. learn how to draw new examples
of cat and dog images*

*e.g. generate new samples of a given phase
of the Ising model*

generative models are Machine Learning techniques that allows to learn **how to generate new examples** similar to those found in a training dataset

Here we consider **energy-based generative models**: close connection with statistical physics

hence the term Boltzmann Learning

Maximum Entropy Generative Models

basic concept is the **Shannon information-theoretic entropy**, which quantifies the statistical uncertainty one has about a random variable drawn from a probability distribution

$$S_p = \text{Tr}_x p(x) \log x$$

\uparrow
*sum/integral over all
possible values of variable*

assume we have a **set of models**, functions of x , whose average should coincide with some observed values. What should be their **underlying prob dist?**

$$\begin{array}{ll} \{f_i(x)\} & \langle f_i \rangle_{\text{obs}} \\ \textit{models} & \textit{observations} \end{array}$$

Principle of Maximum Entropy: choose the probability distribution with the largest uncertainty (Shannon entropy) subject to the observational constraints

$$\langle f_i \rangle_{\text{model}} = \int d\mathbf{x} f_i(\mathbf{x}) p(\mathbf{x}) = \langle f_i \rangle_{\text{obs}}$$

the selected distribution is the one that makes admits the most ignorance beyond the stated data

Maximum Entropy Generative Models

this condition can be expressed as a **Lagrange Multiplier problem** by minimising:

$$\mathcal{L}[p] = -S_p + \sum_i \lambda_i \left(\langle f_i \rangle_{\text{obs}} - \int d\mathbf{x} f_i(\mathbf{x}) p(\mathbf{x}) \right) + \gamma \left(1 - \int d\mathbf{x} p(\mathbf{x}) \right)$$

Shannon entropy *observational constraints* *normalisation*

whose solution gives us the **Maximum Entropy distribution**

$$p(\mathbf{x}) = \frac{\exp \left(\sum_i \lambda_i f_i(\mathbf{x}) \right)}{\int d\mathbf{x} \exp \left(\sum_i \lambda_i f_i(\mathbf{x}) \right)} = \frac{\exp \left(\sum_i \lambda_i f_i(\mathbf{x}) \right)}{Z}$$

partition function $Z = \sum_i \exp \left(\sum_i \lambda_i f_i(\mathbf{x}) \right)$

which is nothing but the **Boltzmann distribution in statistical mechanics**, and where the parameters of the distribution are fixed by the observations

$$\partial_{\lambda_i} \log Z = \langle f_i \rangle_{\text{data}}$$

Energy-based Generative Models

these MaxEnt models can be used to **infer the underlying probability distributions** from a finite set of observations, which subsequently can be used to generate new instances

training an energy-based generative model: using the data to infer the model parameters

$$E(\mathbf{x}; \boldsymbol{\theta}) = - \sum_i \theta_i f_i(\mathbf{x})$$

as in Supervised Learning, we need to specify a **cost function**, which however in the case of generative models is much subtler: what defines a good model?

the most useful method is to **maximise the log-likelihood** of the training set

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\{\boldsymbol{\theta}\})$$

$$\mathcal{L}(\{\boldsymbol{\theta}\}) = \langle \log p_{\boldsymbol{\theta}}(\mathbf{x}) \rangle_{\text{data}} = - \langle E(\mathbf{x}; \boldsymbol{\theta}) \rangle_{\text{data}} - \log Z(\{\boldsymbol{\theta}\})$$

where we have used that the **generative probability distribution** is of the Boltzmann form and that the partition function does not depend on the data

Boltzmann machines

The training of **energy-based generative models** proceeds usually via SGD

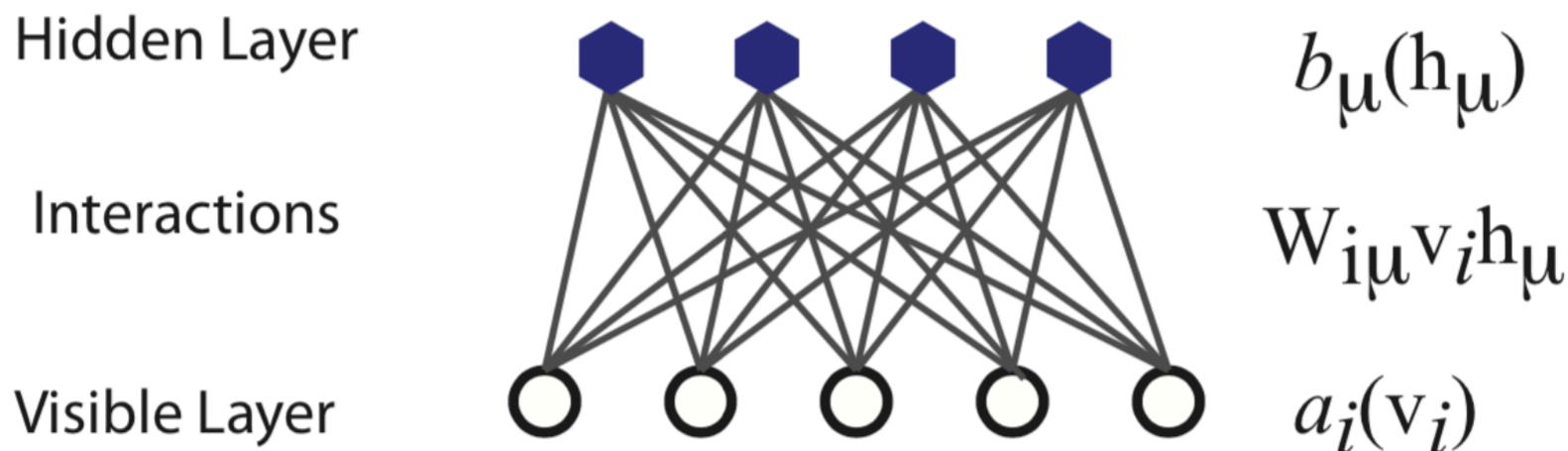
MaxEnt generative models are defined by the choice of the **energy**

$$p(\mathbf{x}) = \frac{\exp\left(\sum_i \lambda_i f_i(\mathbf{x})\right)}{\int d\mathbf{x} \exp\left(\sum_i \lambda_i f_i(\mathbf{x})\right)} = \frac{\exp\left(\sum_i \lambda_i f_i(\mathbf{x})\right)}{Z} = \frac{\exp(-E(\mathbf{x}; \lambda))}{Z}$$

one can construct various **other generative models** with different choices of the energy

e.g. *Restricted Boltzmann machines*

$$E(\mathbf{x}; \theta) = - \sum_i a_i(v_i) - \sum_{\mu} b_{\mu}(h_{\mu}) - \sum_{i\mu} W_{i\mu} v_i h_{\mu}$$



Generative Models & Adversarial Learning

The Kullback-Leibler divergence

The KL divergence, a measure of the similarity between two probability distributions $p(\mathbf{x})$ and $q(\mathbf{x})$ plays an important role in machine learning applications

$$D_{KL}(p \parallel q) = \int d\mathbf{x} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} \quad D_{KL}(q \parallel p) = \int d\mathbf{x} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

which can be symmetrised to construct a **squared metric** (distance)

$$D_{JS}(p \parallel q) = \frac{1}{2} \left(D_{KL}\left(p \middle\| \frac{p+q}{2}\right) + D_{KL}\left(q \middle\| \frac{p+q}{2}\right) \right)$$

Jensen-Shannon divergence

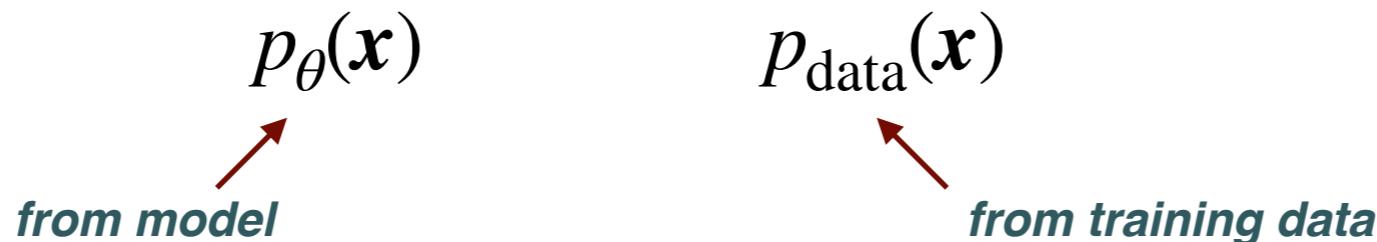
The KL-divergence is **positive-definite**, and only vanishes when $p(\mathbf{x})=q(\mathbf{x})$

$$D_{KL}(p \parallel q) \geq 0$$

in general the integral cannot be computed and one needs to sample the two probability distributions by means of a suitable binning

The Kullback-Leibler divergence

In **generative models** one deals with two probability distributions (data and model), which we would like to have as similar as possible



however subtleties about how we define **similarity** have large implications for the model training

maximising the **log-likelihood of the data under the model** is the same as **minimising the KL divergence** between the data distribution and the model distribution

$$\begin{aligned} D_{KL}(p_{\text{data}} || p_\theta) &= \int d\mathbf{x} p_{\text{data}}(\mathbf{x}) \log \frac{p_{\text{data}}(\mathbf{x})}{p_\theta(\mathbf{x})} \\ &= \int d\mathbf{x} p_{\text{data}}(\mathbf{x}) \log p_{\text{data}}(\mathbf{x}) - \int d\mathbf{x} p_{\text{data}}(\mathbf{x}) \log p_\theta(\mathbf{x}) \\ &= S[p_{\text{data}}] - \langle \log p_\theta \rangle_{\text{data}} \end{aligned}$$

The Kullback-Leibler divergence

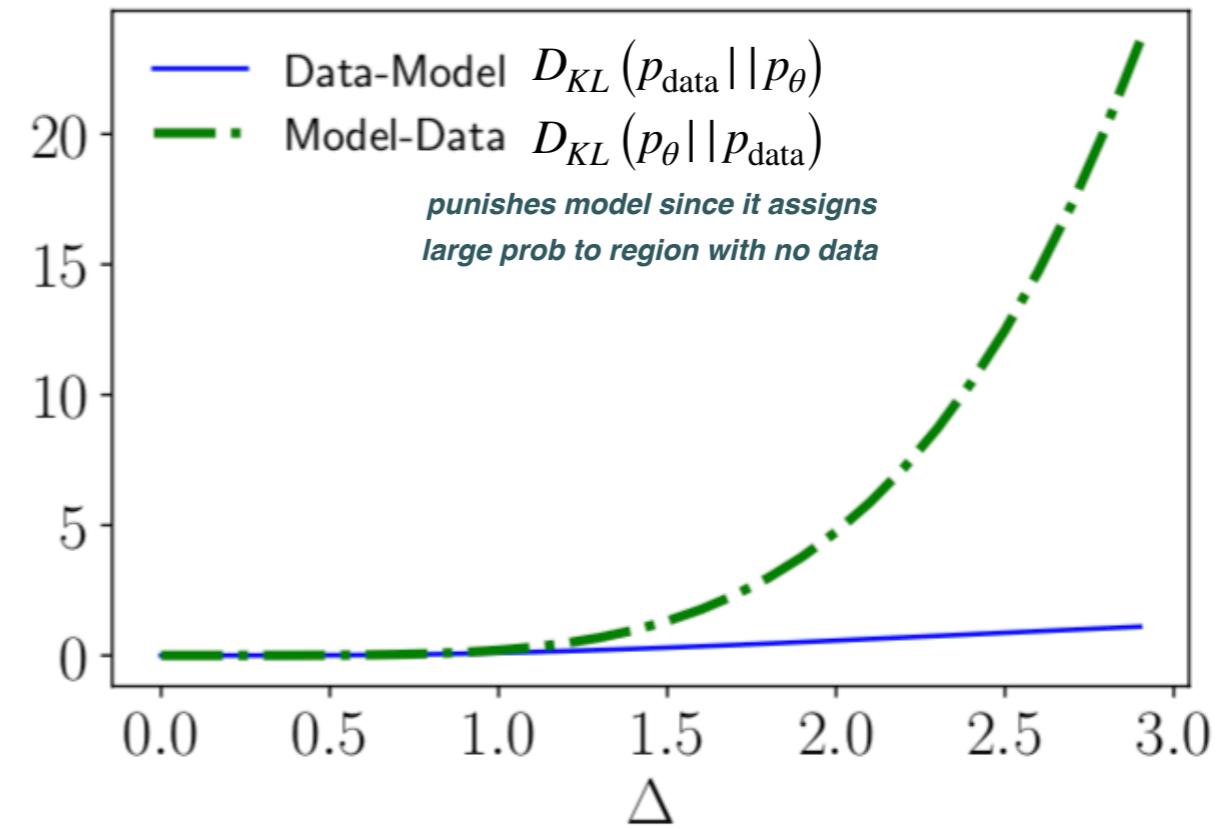
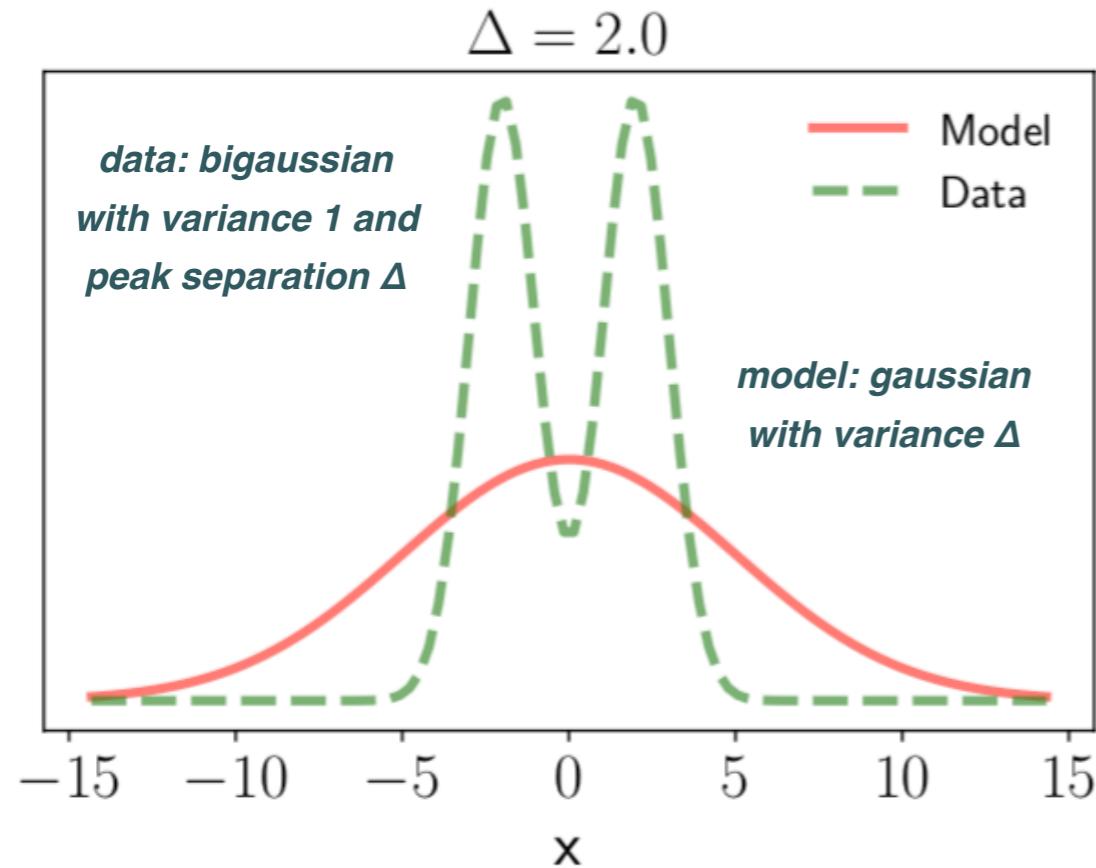
maximising the **log-likelihood** of the data under the model is the same as **minimising the KL divergence** between the data distribution and the model distribution

$$\langle \log p_\theta(x) \rangle_{\text{data}} = S[p_{\text{data}}] - D_{KL}(p_{\text{data}} \parallel p_\theta)$$

↑ *↑* *↑*
Log-likelihood of data under model *entropy of data:
independent of model parameters* *KL-divergence*

Similarity

Similarity between probability distributions is a subtle concept



$$D_{KL}(p_{\text{data}} \parallel p_{\theta}) \longrightarrow$$

misses important information when comparing the data and theory probability distributions

Adversarial Learning

(draw values of θ accordingly to p_θ and compare with data)

$$(1) \ D_{KL} (p_{\text{data}} || p_\theta) \longrightarrow \text{Calculable using sampling}$$

$$(2) \ D_{KL} (p_\theta || p_{\text{data}}) \longrightarrow \begin{aligned} &\text{Large when model over-weights low-} \\ &\text{density regions near real peaks} \\ &\text{but not calculable since } p_{\text{data}} \text{ unknown} \end{aligned}$$

In **Adversarial Learning** we achieve a similar goal as that of minimising **(2)** by training a **discriminator** to distinguish between real data points and samples from the model

By punishing the model for generating points that can be easily discriminated from the data, Adversarial Learning decreases the **weight of regions in the model space that are far away from data points**, regions that inevitably arise when maximising the likelihood

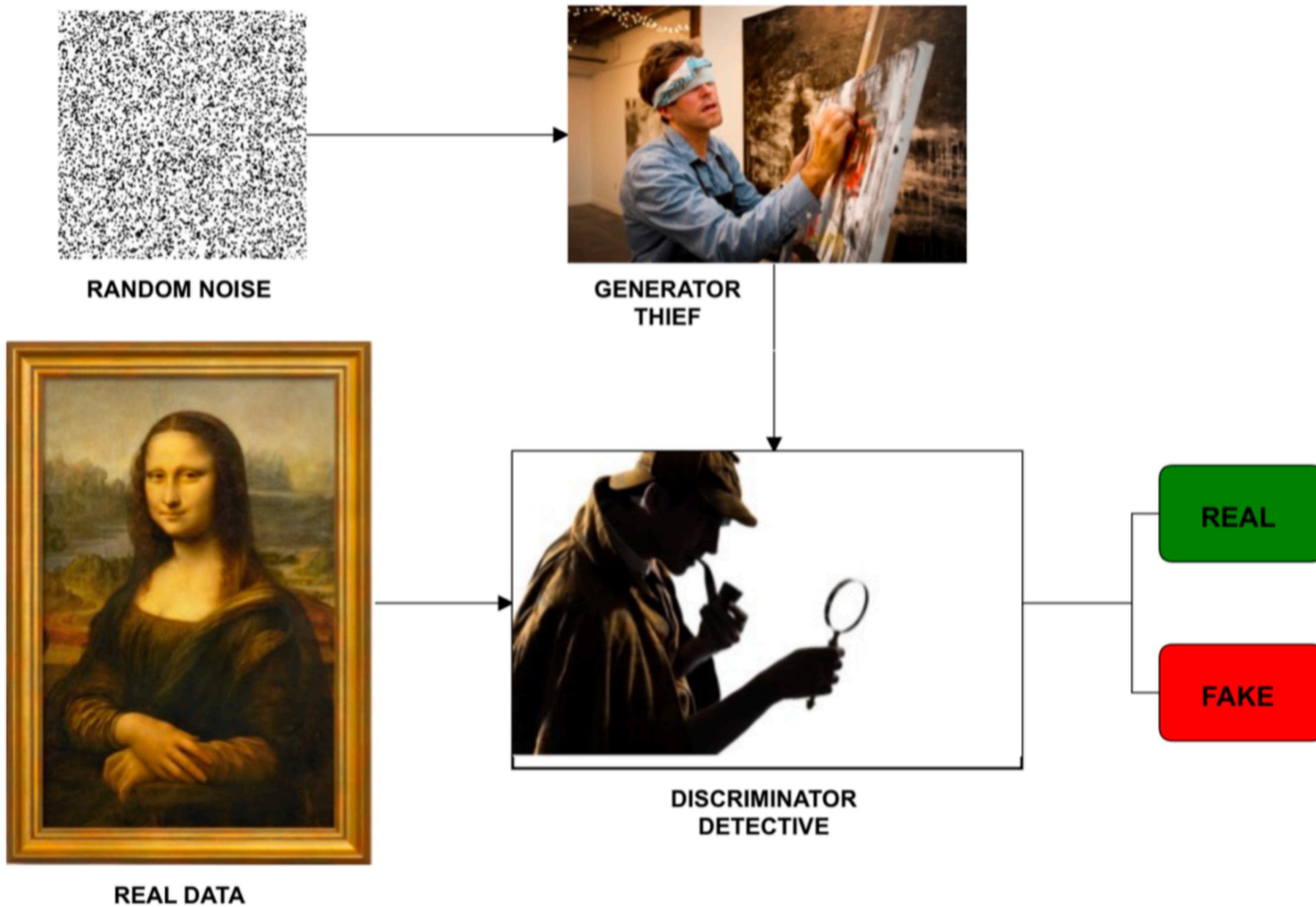
Generative adversarial networks

Generative Adversarial Networks (GANs) are deep neural network architectures, composed by two independent NNs which **compete against each other**

- (1) A **generator G** NN that creates (samples) pseudo-data by inferring the probability distribution associated to the training dataset
- (2) A **discriminator D** NN which determines the probability of a given sample arises from the actual training data rather than having been produced by **G**

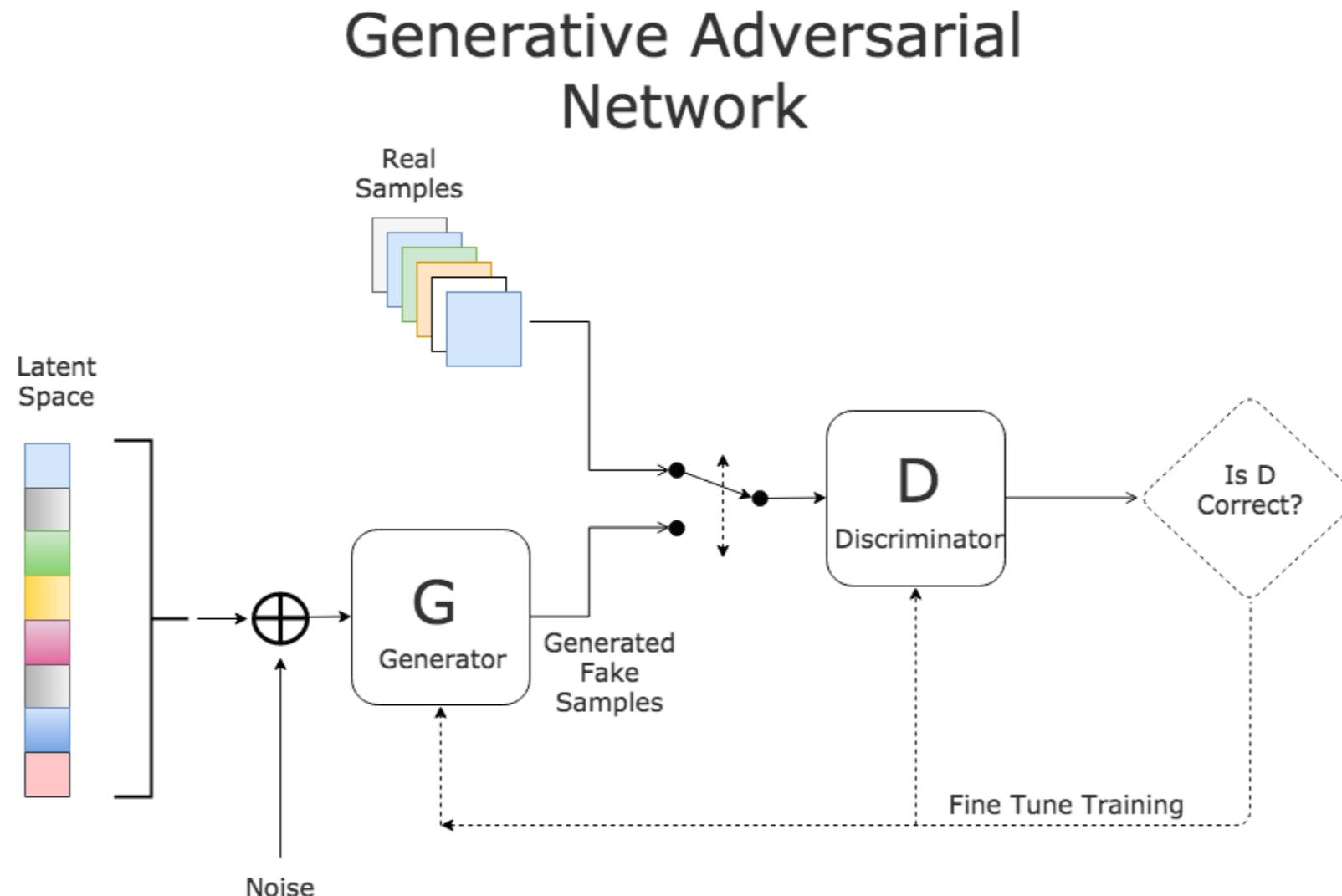
the generator network **G** should be trained to maximise the probability that the discriminator network **D** makes a mistake: that is, **G** should generate pseudo-data samples that are virtually **indistinguishable** from the actual data

Generative adversarial networks



Generative Adversarial Networks

- New architecture for an **unsupervised neural network training** (unlabelled samples)
- Based on two **independent nets** that work separately and act as **adversaries**:
 - the **Discriminator (D)** undergoes training and plays the role of classifier
 - the **Generator (G)** and is tasked to generate random samples that **resemble real samples** with a twist rendering them as fake samples.



GAN training

As with other NN architectures one uses GD to train GANs, but now one has to **update sequentially** the model parameters of both **G** and **D**

- Take a sample of N data points from the training set

$$\{\mathbf{x}_n\}_{n=1}^N \quad \mathbf{x}_n = (x_{n,1}, x_{n,2}, \dots, x_{n,p}) \quad p = \text{number of features per sample}$$

- Produce a sample of N pseudo-data points from generator **G** (at ite_0 this is random noise)

$$\{\mathbf{z}_n\}_{n=1}^N \quad \mathbf{z}_n = (z_{n,1}, z_{n,2}, \dots, z_{n,p})$$

- Evaluate the cost function: since we are dealing with binary classification (true/false) the appropriate cost function is the **cross-entropy**

$$C(\theta_D, \theta_G) = \frac{1}{N} \sum_{n=1}^N (\log D(\mathbf{x}_i) + \log(1 - D(G(\mathbf{z}_i))))$$

NN params of G output of D when input
NN params of D a real data sample a "fake" data sample produced by G

- Train **D** using GD to maximise its discrimination capability

GAN training

- Evaluate the cost function: since we are dealing with binary classification (true/false) the appropriate cost function is the **cross-entropy**

$$C(\theta_D, \theta_G) = \frac{1}{N} \sum_{n=1}^N (\log D(x_i) + \log(1 - D(G(z_i))))$$

- Train **D** using GD to maximise its discrimination capability

$$\mathbf{v}_t = \eta_t \nabla_{\theta_D} C(\theta_{D,t}, \theta_G), \quad \theta_{D,t+1} = \theta_D - \mathbf{v}_t$$

- At this point **D** can tell apart data from pseudo-data pretty well, so we need to train **G** to generate better (closer to the training set) pseudo-data samples

- Produce a sample of N pseudo-data points from the generator **G** $\{z_n\}_{n=1}^N$

$$C(\theta_D, \theta_G) = \frac{1}{N} \sum_{n=1}^N \log(1 - D(G(z_i)))$$

output of D (now with its parameters fixed)

$$\mathbf{v}_t = \eta_t \nabla_{\theta_G} C(\theta_{D,t}, \theta_G), \quad \theta_{G,t+1} = \theta_G - \mathbf{v}_t$$

GAN training

the generator and discriminator are sequentially trained and iterated until convergence is achieved, at this point **D** cannot tell apart the pseudo-data from **G** from the real data

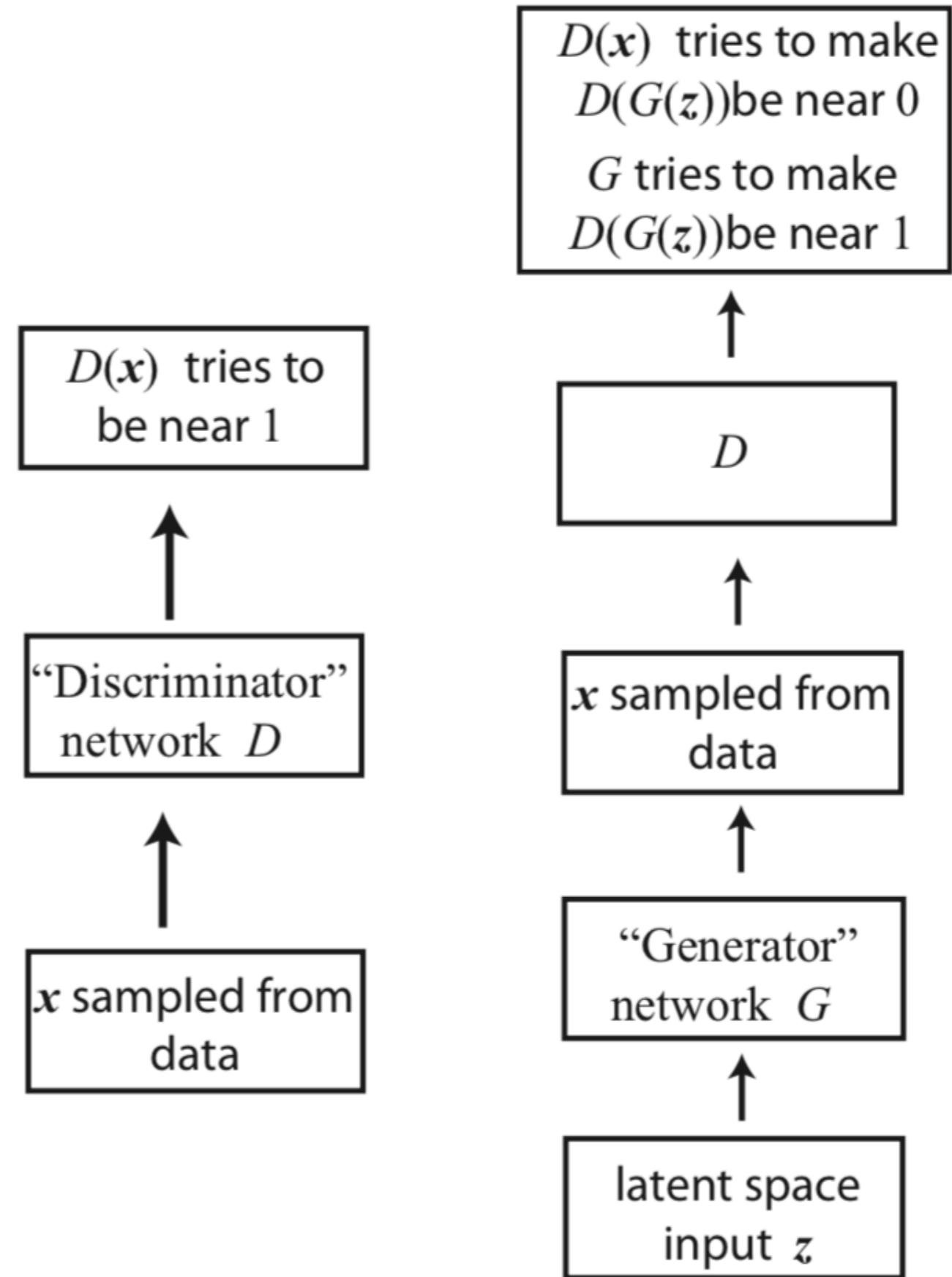


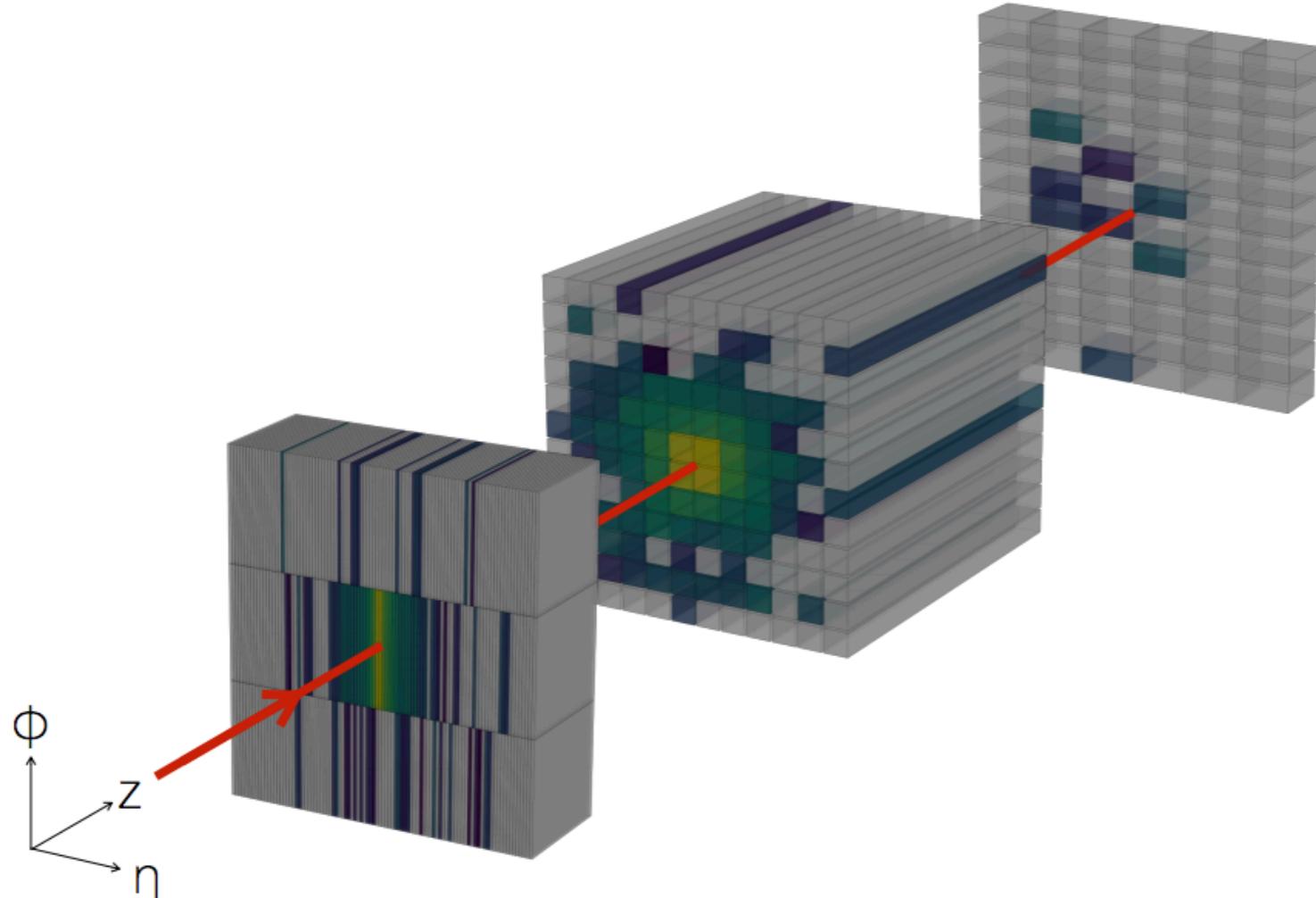
Image generation with GANs



<https://thispersondoesnotexist.com/>

GANs for detector simulation in HEP

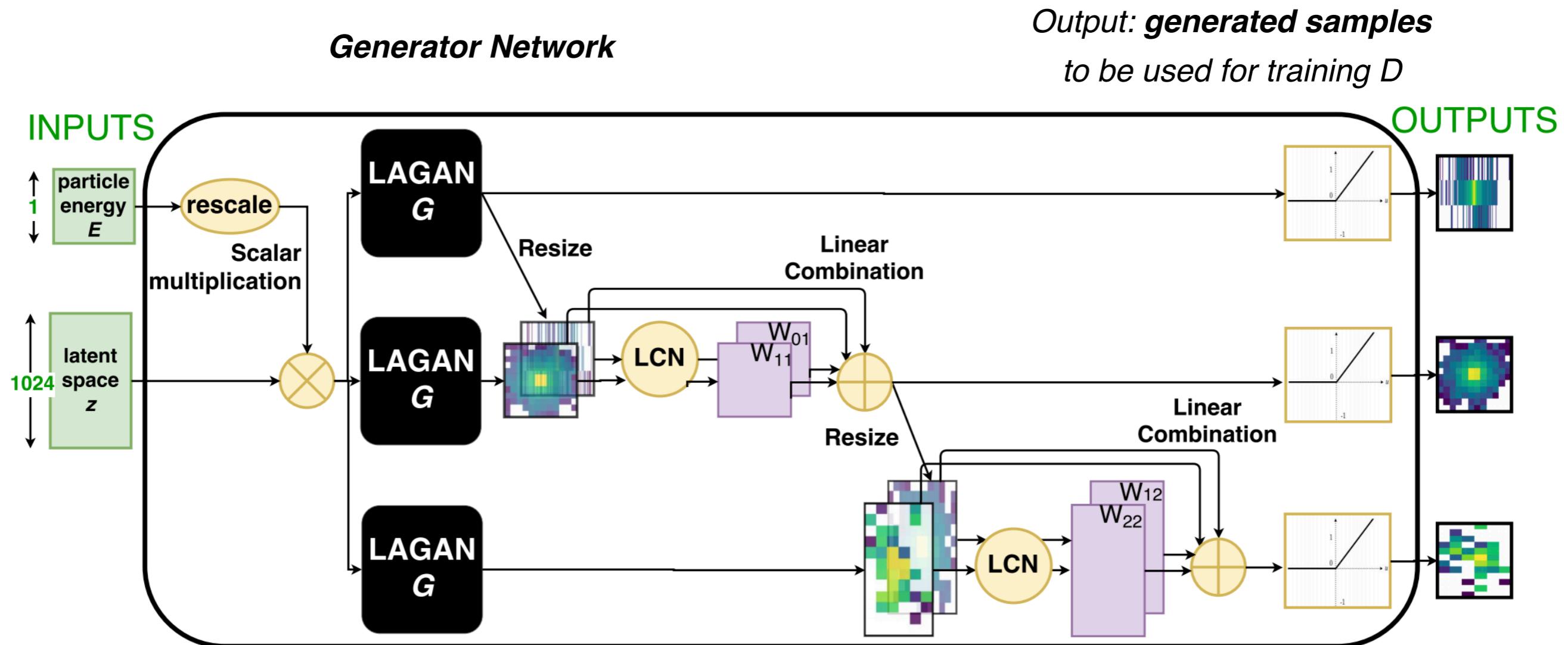
- Modelling accurately the response of detectors with the **propagation of high energy particles** is an essential task for present and future HEP experiment
- Detector simulation at the LHC** is a very CPU-intensive task, dominated by modelling of particle showers inside calorimeters
- Generative Adversarial Networks** can speed up detector simulation by orders of magnitude



*Task: to efficiently model the **propagation of high energy particles** (and their interaction) within the layers of electromagnetic and hadronic calorimeters*

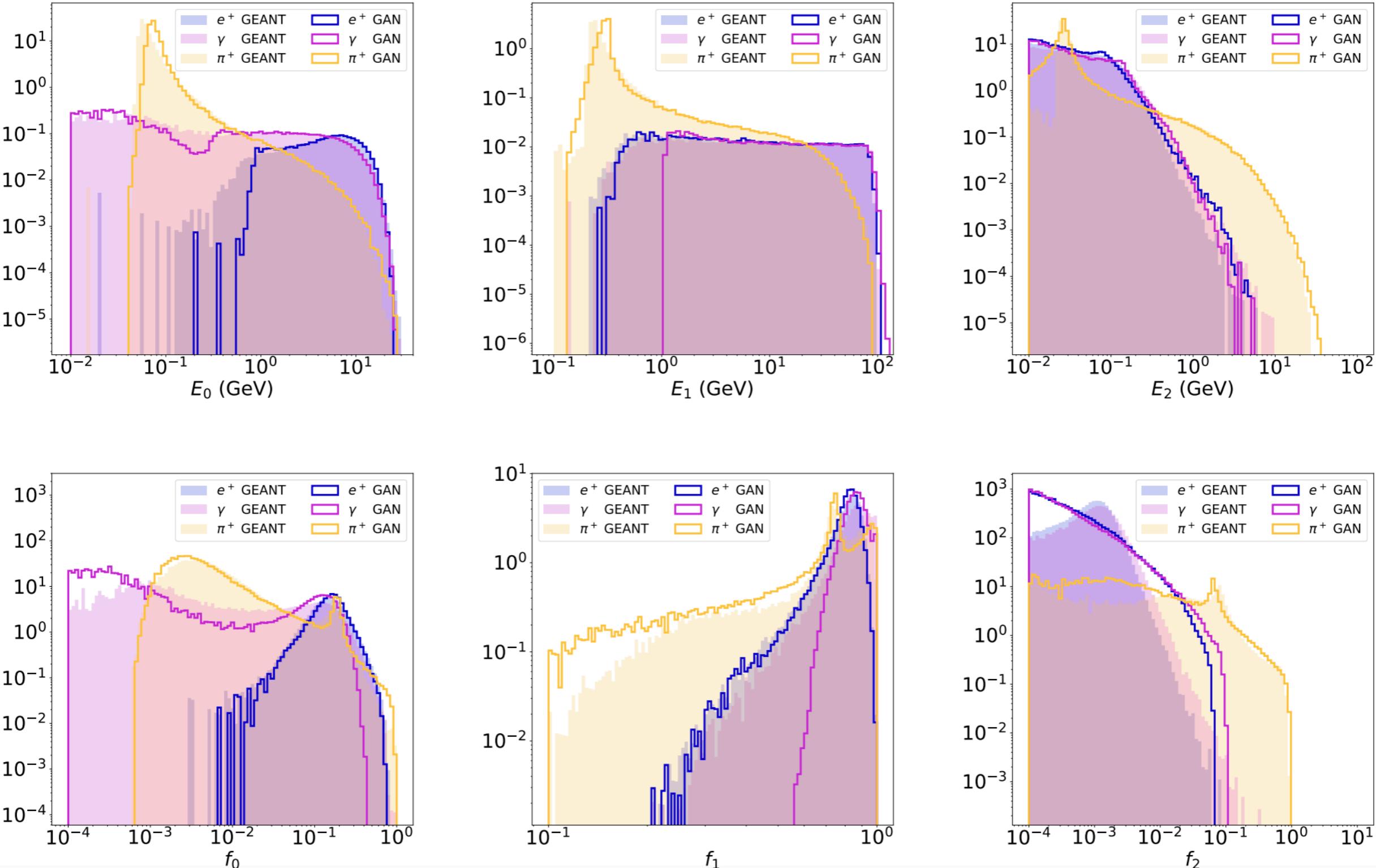
GANs for detector simulation in HEP

- Use GANs as a tool to **speed up full simulation of particle showers** in a HEP calorimeter
- The generator G learns a map from **a latent space** to space of **generated samples** for training
- Carefully understanding the **underlying physics of particle propagation** in a detector is crucial to optimise the training strategy, e.g. relationships between neighbouring detector layers



Paganini et al. 17

GANs for detector simulation in HEP



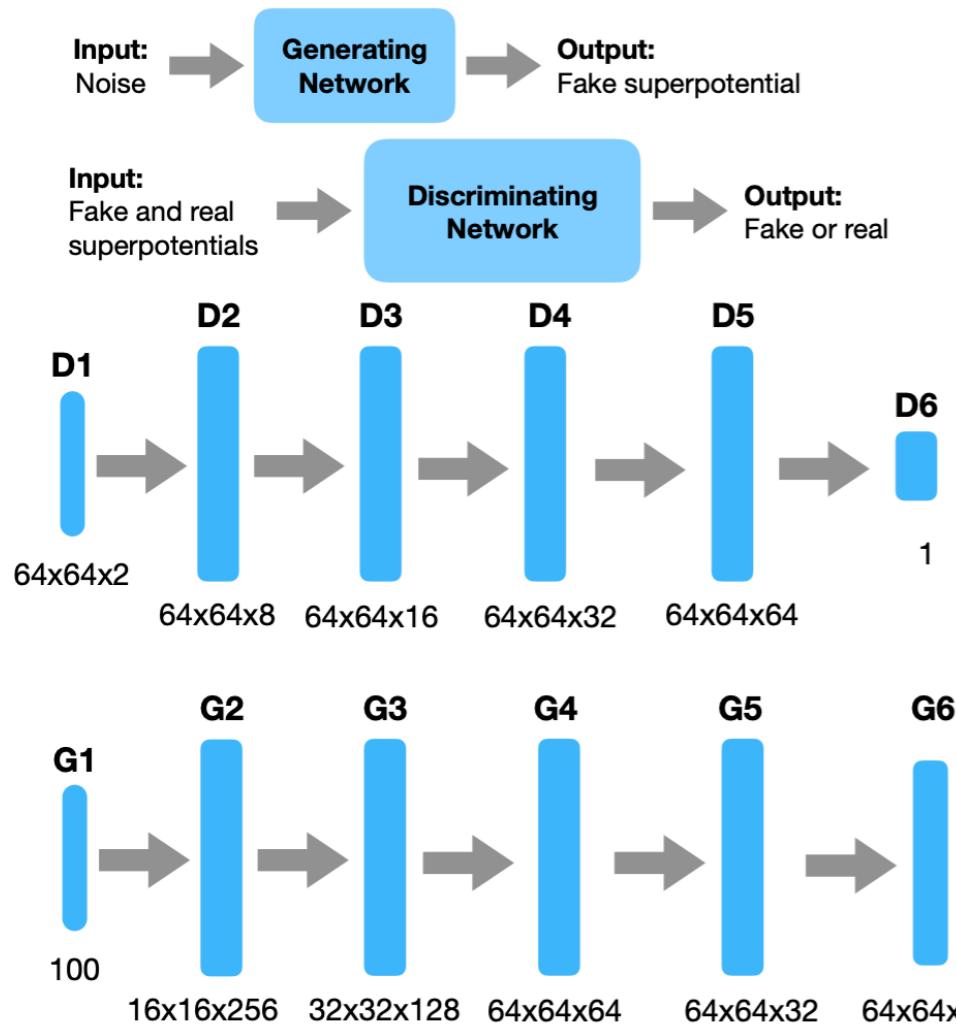
GANs for detector simulation in HEP

Simulator	Hardware	Batch Size	ms/shower
GEANT4	CPU	N/A	1772
CALOGAN	CPU	1	13.1
		10	5.11
		128	2.19
		1024	2.03
CALOGAN	GPU	1	14.5
		4	3.68
		128	0.021
		512	0.014
		1024	0.012

Speed-up by several orders of magnitude, specially when running in GPUs

(More) Applications of ML to Theoretical Physics

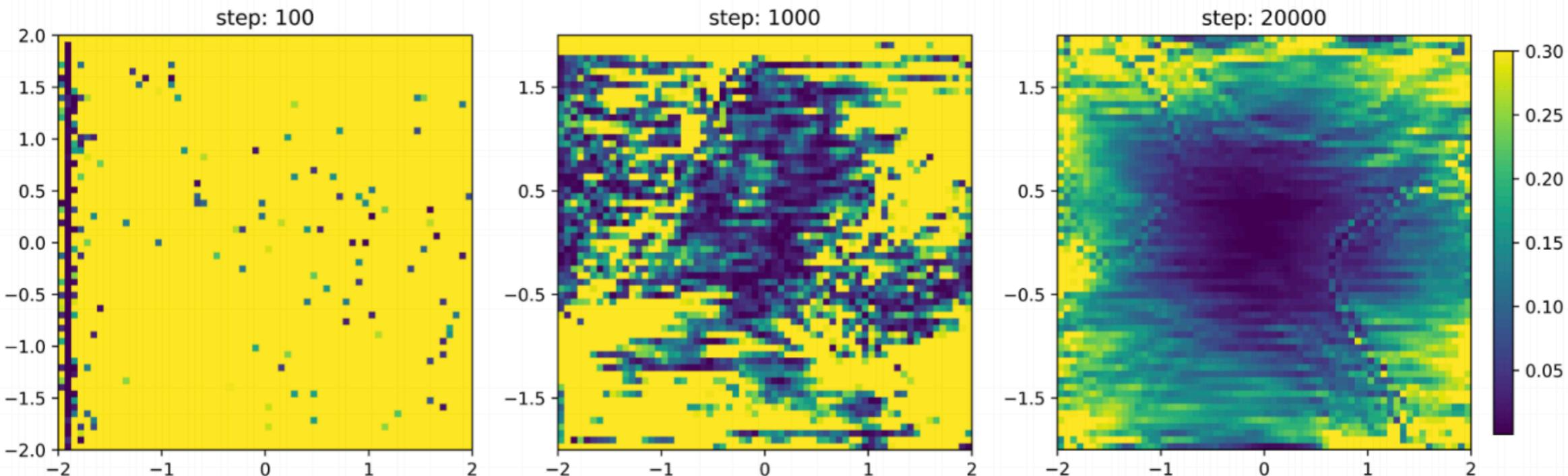
GANs for generating EFT models



- Starting from examples of **known SUSY QFTs** (with generating superpotentials) train GANs to generate other, equally consistent, theories
- After training the generator outputs new superpotentials that correspond to **new SUSY theories** not present in the training sample

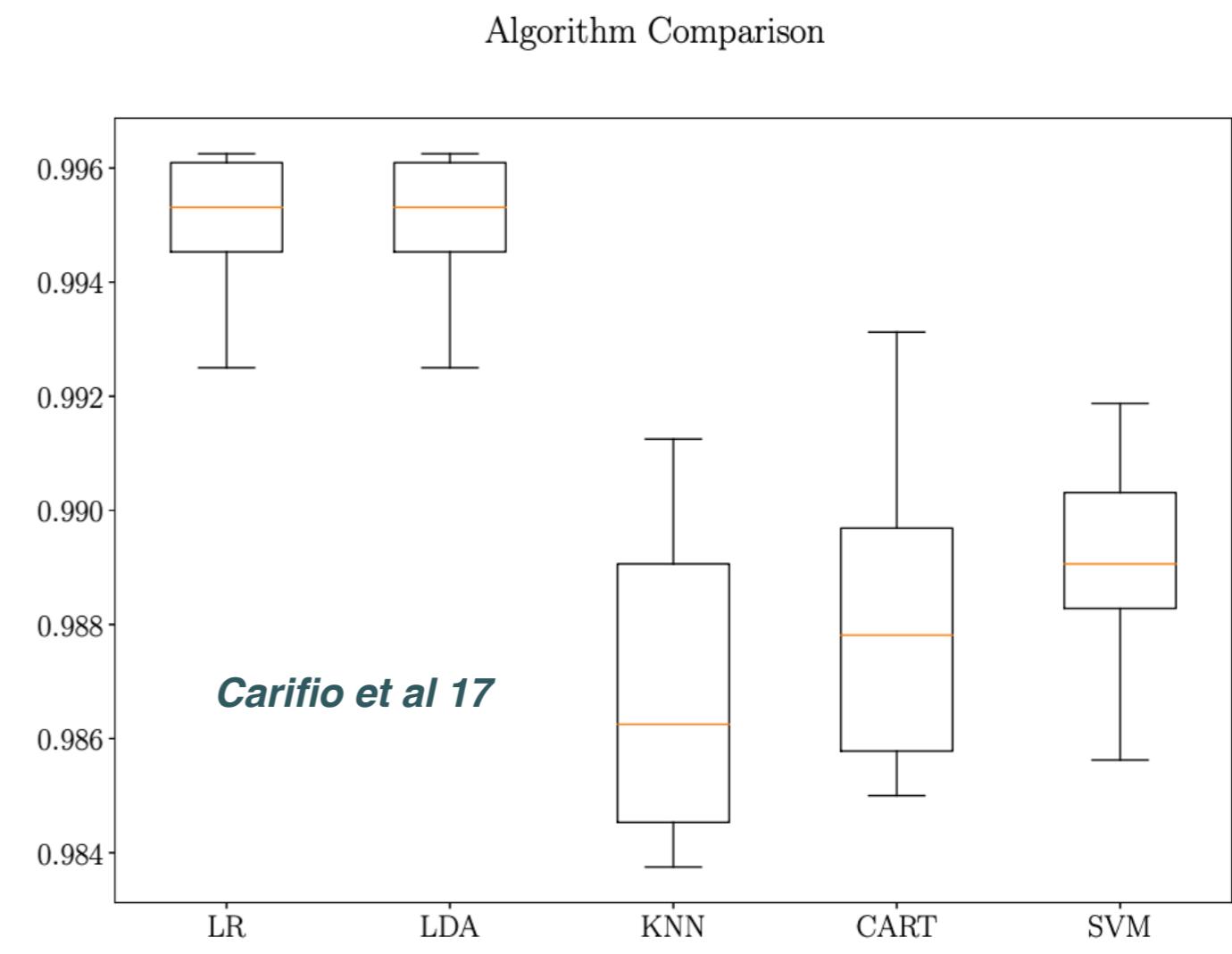
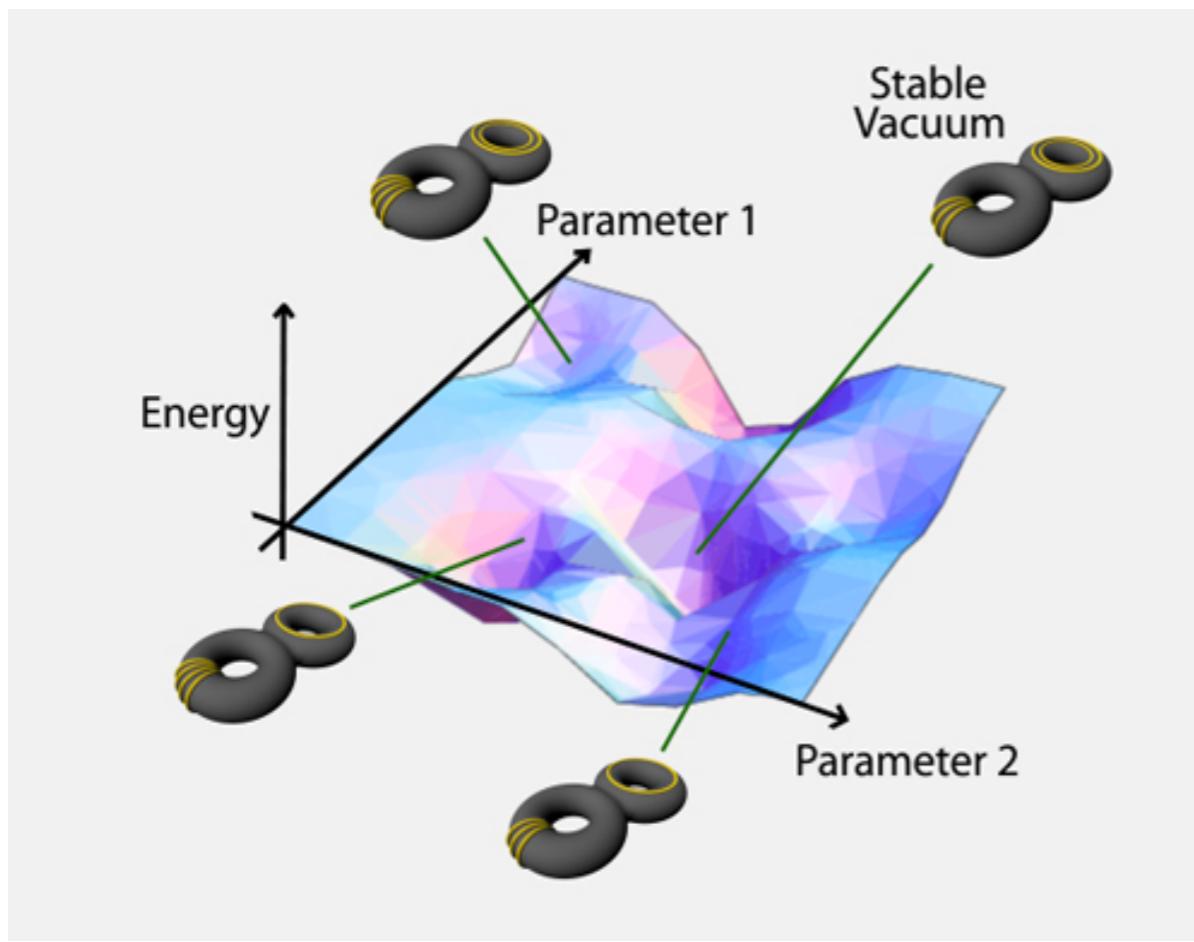
Erbin & Krippendorf 18

evolution of superpotential during training



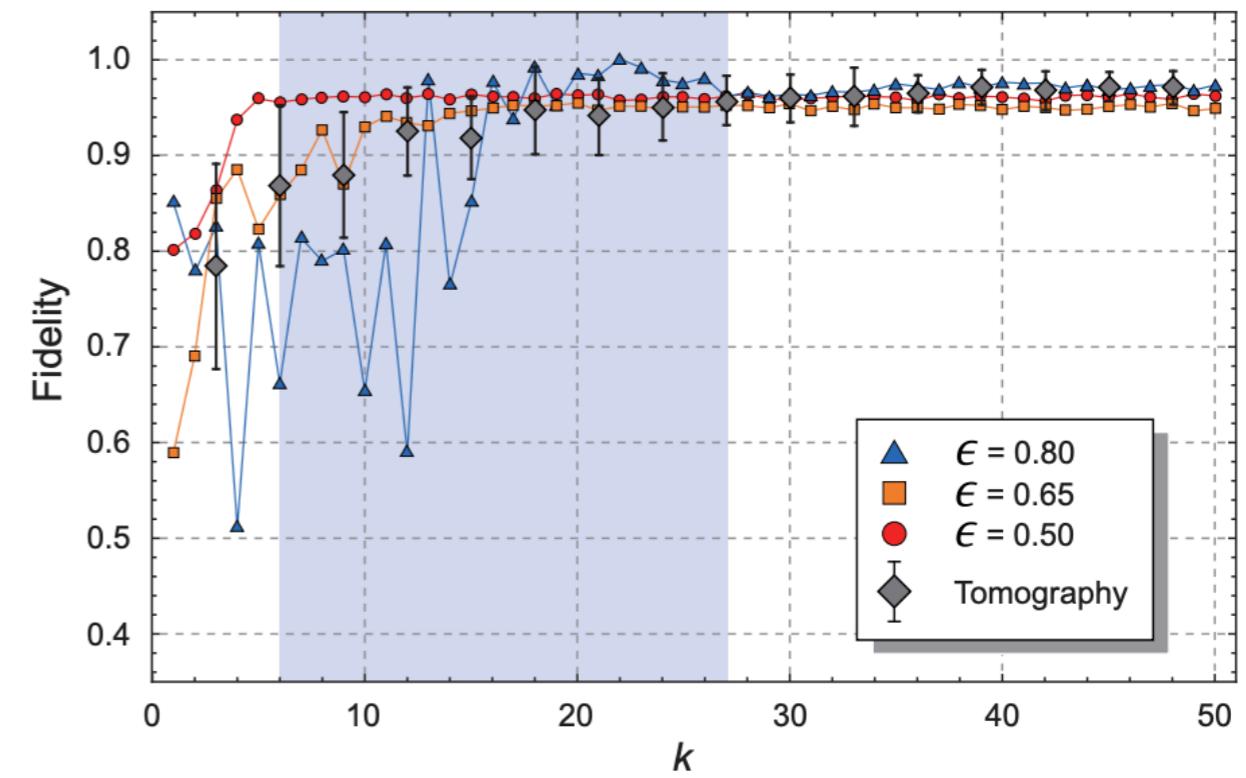
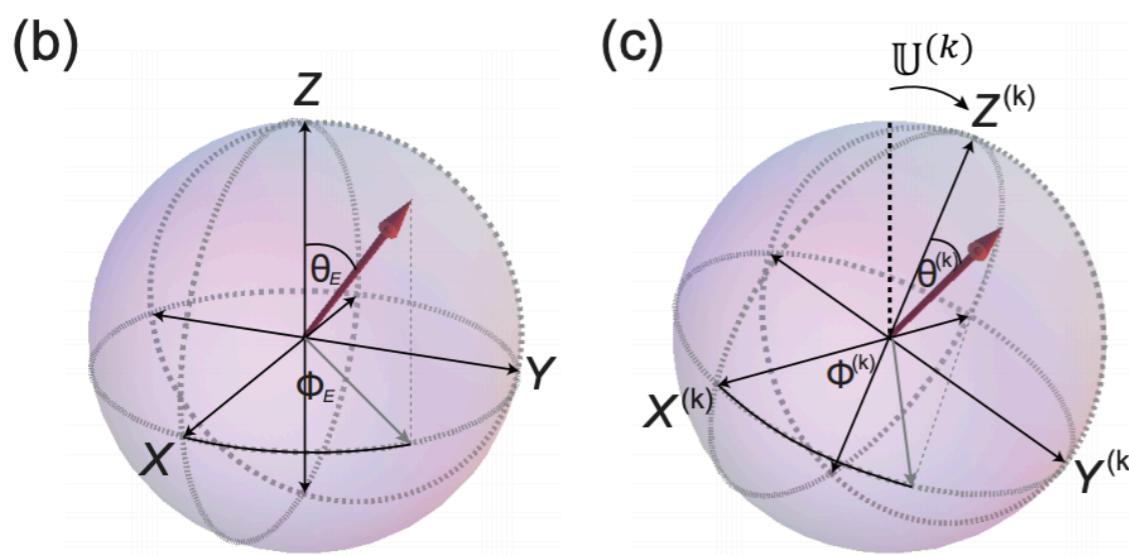
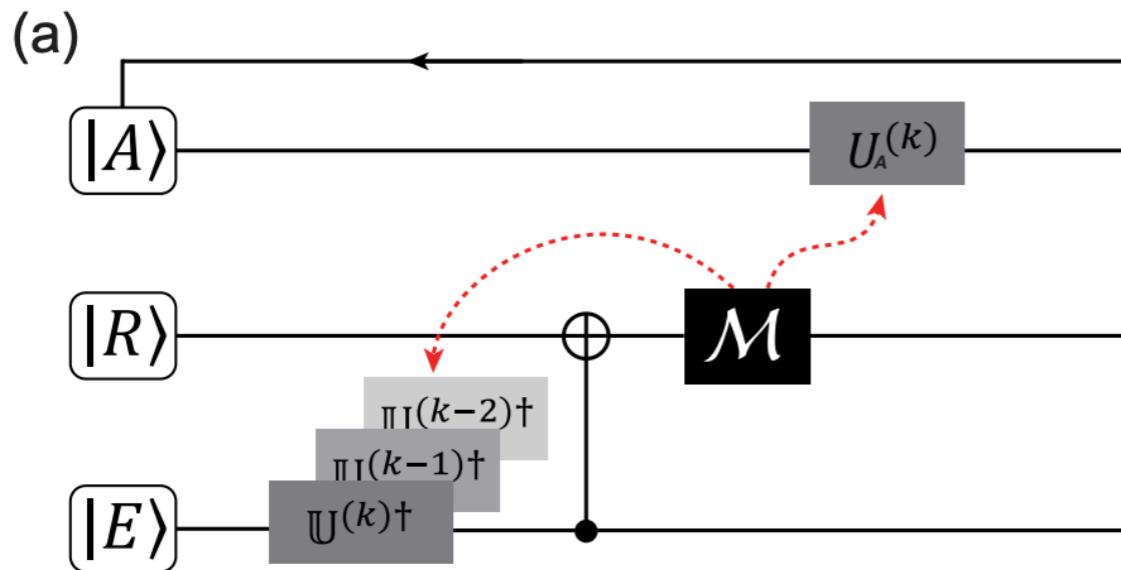
Machine Learning the String Landscape

- 💡 Numerical exploration of the string theory landscape
- 💡 **Formulation and verification of conjectures** about the properties of the landscape
- 💡 Using regression, K-means clustering, decision trees, Linear Discriminant Analysis



Qubit Reconstruction with Reinforcement Learning

to reconstruct an unknown photonic quantum state with a limited amount of copies, one employs a **semi-quantum reinforcement learning approach** to adapt one qubit state, an “agent”, to an unknown quantum state, an “environment”



Yu et al 19

More Examples

Machine learning and the physical sciences

Giuseppe Carleo

Center for Computational Quantum Physics, Flatiron Institute,
162 5th Avenue, New York, NY 10010, USA*

Ignacio Cirac

Max-Planck-Institut für Quantenoptik,
Hans-Kopfermann-Straße 1, D-85748 Garching, Germany

Kyle Cranmer

Center for Cosmology and Particle Physics, Center of Data Science,
New York University, 726 Broadway, New York, NY 10003, USA

Laurent Daudet

LightOn, 2 rue de la Bourse, F-75002 Paris, France

Maria Schuld

University of KwaZulu-Natal, Durban 4000, South Africa
National Institute for Theoretical Physics, KwaZulu-Natal, Durban 4000, South Africa,
and Xanadu Quantum Computing, 777 Bay Street, M5B 2H7 Toronto, Canada

Naftali Tishby

The Hebrew University of Jerusalem, Edmond Safra Campus, Jerusalem 91904, Israel

Leslie Vogt-Maranto

Department of Chemistry, New York University, New York, NY 10003, USA

Lenka Zdeborová

Institut de physique théorique, Université Paris Saclay, CNRS, CEA,
F-91191 Gif-sur-Yvette, France†

Contents

I. Introduction

- A. Concepts in machine learning
 - 1. Supervised learning and neural networks
 - 2. Unsupervised learning and generative modelling
 - 3. Reinforcement learning

II. Statistical Physics

- A. Historical note
- B. Theoretical puzzles in deep learning
- C. Statistical physics of unsupervised learning
 - 1. Contributions to understanding basic unsupervised methods
 - 2. Restricted Boltzmann machines
 - 3. Modern unsupervised and generative modelling
- D. Statistical physics of supervised learning
 - 1. Perceptron and GLMs
 - 2. Physics results on multi-layer neural networks

- 3. Information Bottleneck
- 4. Landscapes and glassiness of deep learning
- E. Applications of ML in Statistical Physics
- F. Outlook and Challenges

III. Particle Physics and Cosmology

- A. The role of the simulation
- B. Classification and regression in particle physics
 - 1. Jet Physics
 - 2. Neutrino physics
 - 3. Robustness to systematic uncertainties
 - 4. Triggering
 - 5. Theoretical particle physics
- C. Classification and regression in cosmology
 - 1. Photometric Redshift
 - 2. Gravitational lens finding and parameter estimation
 - 3. Other examples
- D. Inverse Problems and Likelihood-free inference
 - 1. Likelihood-free Inference
 - 2. Examples in particle physics
 - 3. Examples in Cosmology
- E. Generative Models
- F. Outlook and Challenges

IV. Many-Body Quantum Matter

- A. Neural-Network quantum states
 - 1. Representation theory
 - 2. Learning from data
 - 3. Variational Learning
- B. Speed up many-body simulations
- C. Classifying many-body quantum phases
 - 1. Synthetic data
 - 2. Experimental data
- D. Tensor networks for machine learning
- E. Outlook and Challenges

V. Quantum computing

- A. Quantum state tomography
- B. Controlling and preparing qubits
- C. Error correction

VI. Chemistry and Materials

- A. Energies and forces based on atomic environments
- B. Potential and free energy surfaces
- C. Materials properties
- D. Electron densities for density functional theory
- E. Data set generation
- F. Outlook and Challenges

VII. AI acceleration with classical and quantum hardware

- A. Beyond von Neumann architectures
- B. Neural networks running on light
- C. Revealing features in data
- D. Quantum-enhanced machine learning
- E. Outlook and Challenges