



# HISS Summer School

## Intro and Overview

Prof. Gregor Kasieczka  
[gregor.kasieczka@uni-hamburg.de](mailto:gregor.kasieczka@uni-hamburg.de)

# What's special about Physics?

Different fields, unified by **digitisation challenges**:

- **Immense data volumes** with high complexity
- Weak signals overshadowed by much larger background processes
- Large data-throughput requiring custom solutions

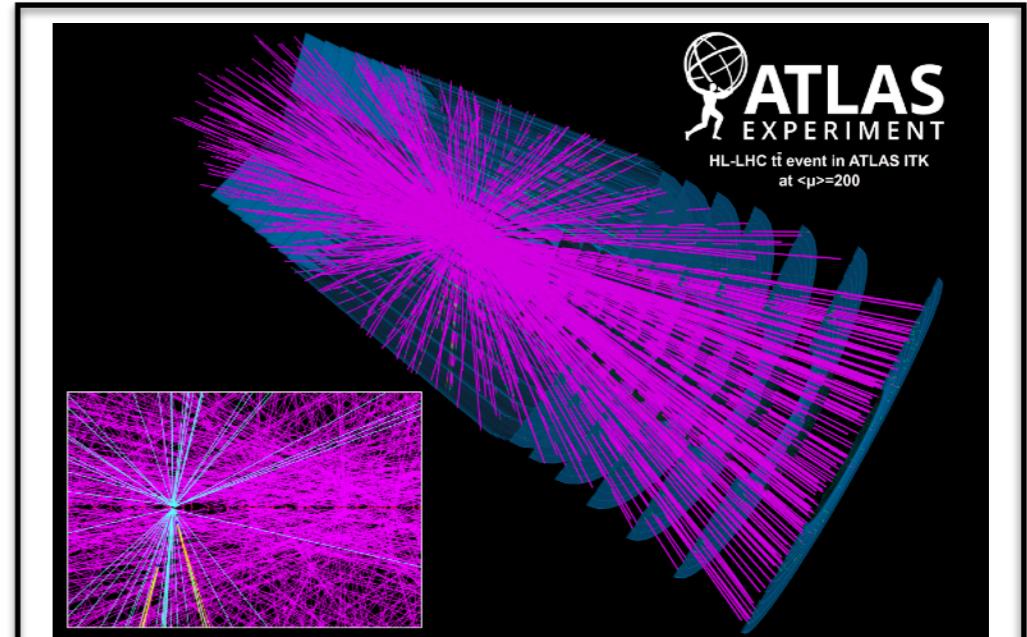


The planned Square Kilometre Array (SKA) telescope will produce one exabyte of raw data/day (~10 PB after compression)

# What's special about Physics?

Different fields, unified by **digitisation challenges**:

- Immense data volumes **with high complexity**
- Weak signals overshadowed by much larger background processes
- Large data-throughput requiring custom solutions

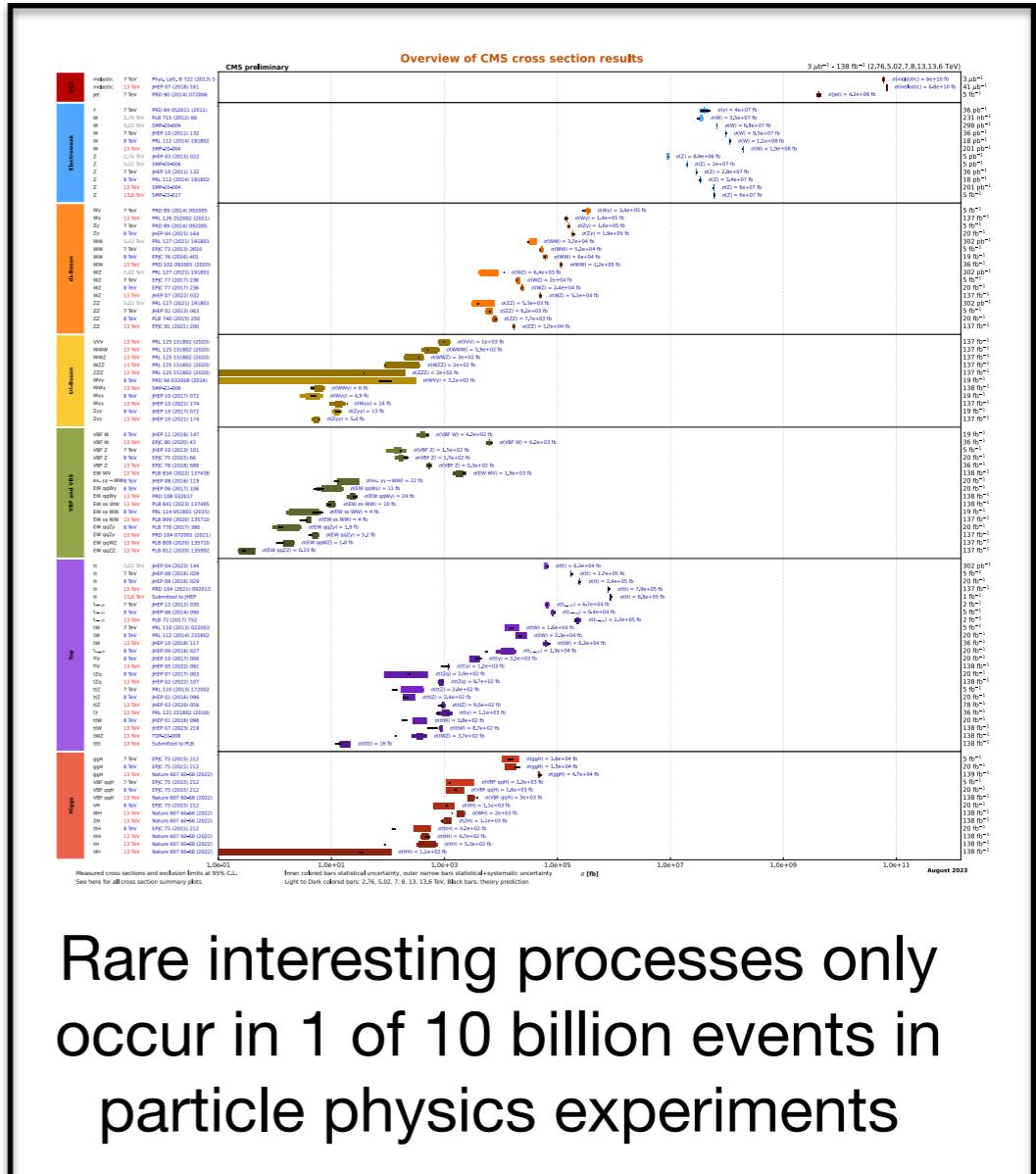


At the High Luminosity LHC, overlapping collisions mean that thousands of trajectories need to be reconstructed in parallel

# What's special about Physics?

Different fields, unified by digitisation challenges:

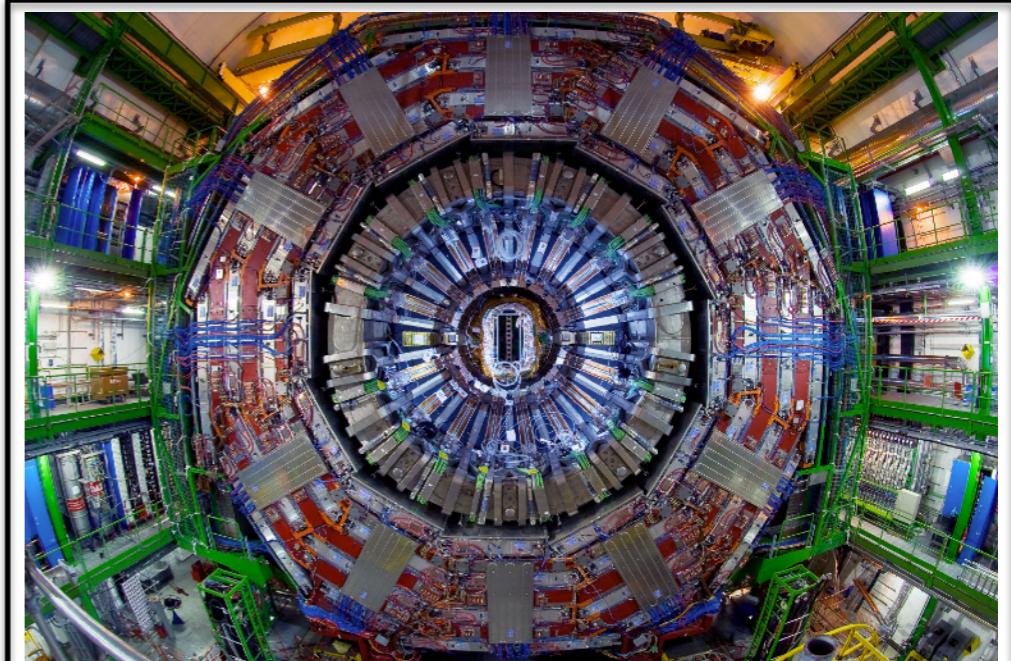
- Immense data volumes with high complexity
- Weak signals overshadowed by much larger background processes
- Large data-throughput requiring custom solutions



# What's special about Physics?

Different fields, unified by **digitisation challenges**:

- Immense data volumes with high complexity
- Weak signals overshadowed by much larger background processes
- **Large data-throughput requiring custom solutions**



Trigger systems filter 40 million particle collisions/second at experiments at the Large Hadron Collider (LHC) at CERN

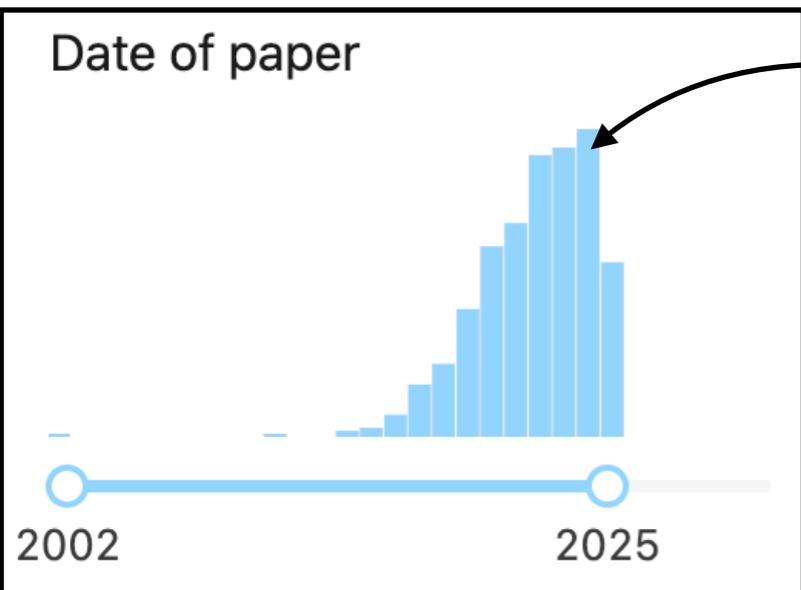
# State of AI in Physics

literature ▾

('machine learning' or 'deep learning') in hep-ex



Date of paper



200+ papers in 2024

Similar for nucl-ex,  
astro-ph, hep-ph, ...

Rapid **rise of AI** in fundamental physics

Transition from **concepts to applications**

Recognition of AI work as **Nobel-prize worthy**

Nobel Prize in Chemistry 2024



The Nobel Prize in Physics 2024



Ill. Niklas Elmehed © Nobel Prize Outreach  
John J. Hopfield  
Prize share: 1/2



Ill. Niklas Elmehed © Nobel Prize Outreach  
Geoffrey Hinton  
Prize share: 1/2

# Plan

- Lecture 1:
  - Basics of Neural Networks
  - Classification example: Top Tagging
- Lecture 2:
  - Generative Models
  - Anomaly Detection

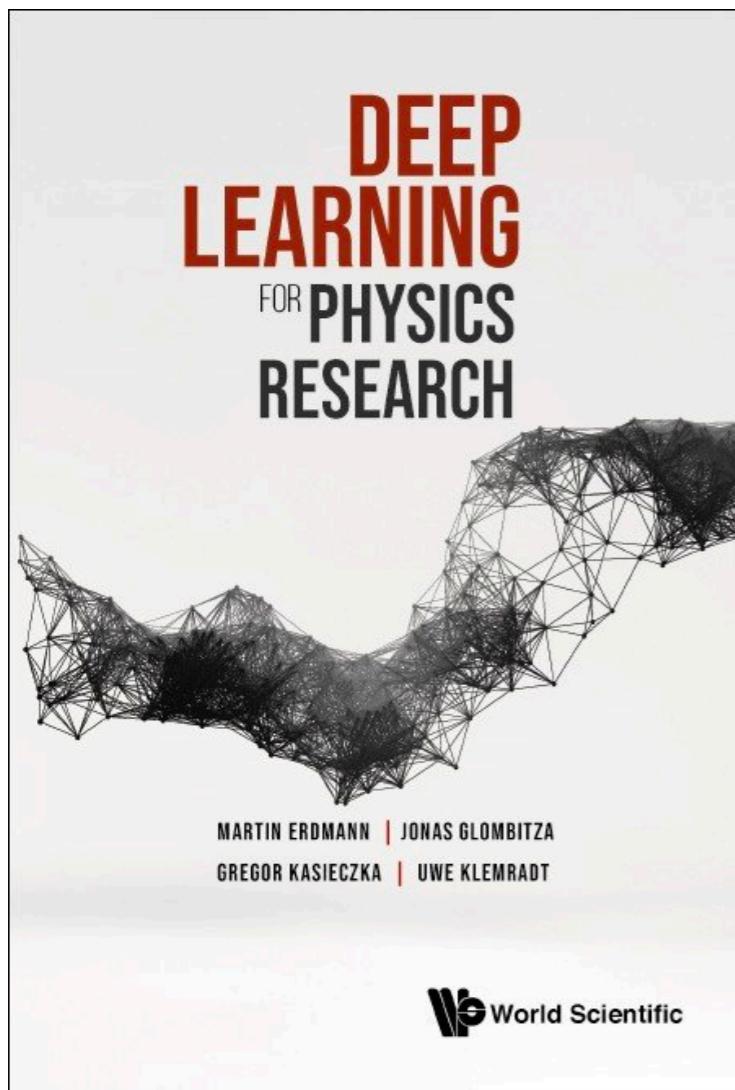
# Resources

## Deep Learning

An MIT Press book

Ian Goodfellow and Yoshua Bengio and Aaron Courville

Free online: <https://www.deeplearningbook.org/>



[https://www.worldscientific.com/  
worldscibooks/  
10.1142/12294#t=aboutBook](https://www.worldscientific.com/worldscibooks/10.1142/12294#t=aboutBook)  
(available in UHH library)

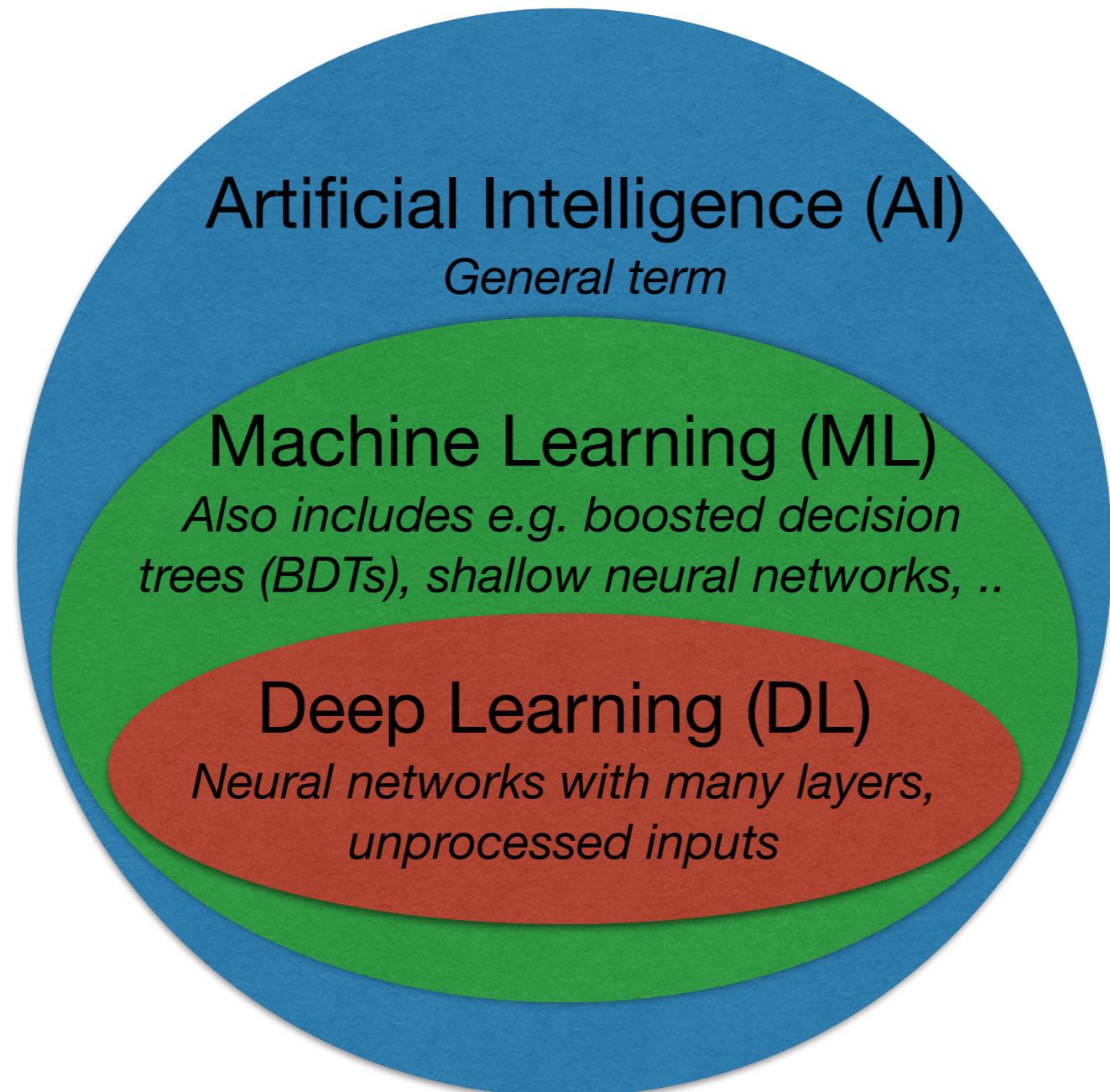
## Lecture Slides



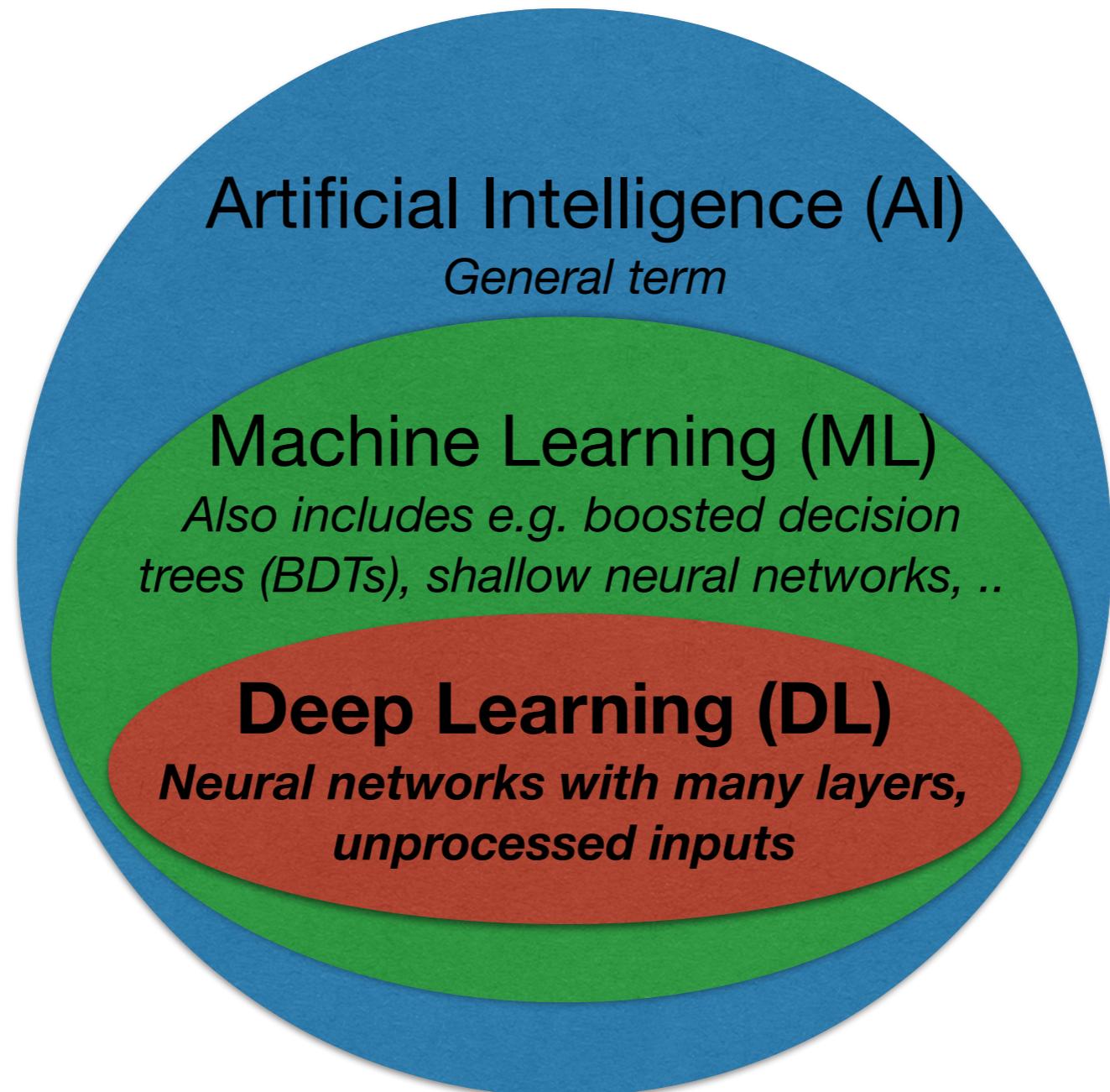
Prof. Gregor Kasieczka  
[gregor.kasieczka@uni-hamburg.de](mailto:gregor.kasieczka@uni-hamburg.de)

# **Learning like a machine**

# Terminology



# Terminology



# Basic idea

- Classical approach:
  - Write a sequence of instructions to solve a specific task
  - e.g.:
    - Tracking algorithm
    - Jet clustering
    - Calculation of physical observables

```
21     if (trackerTopology.pxfSide(detId) == 1) {
22         trackingLayer._side = TrackingLayer::Side::NEG_ENDCAP;
23     } else if (trackerTopology.pxfSide(detId) == 2) {
24         trackingLayer._side = TrackingLayer::Side::POS_ENDCAP;
25     } else {
26         throw cms::Exception("FastSimulation/Tracking")
27             << "Tracker hit for seeding in FPix seems neither on positive nor on negative disk side: "
28             << trackerTopology.print(detId).c_str();
29     }
30     trackingLayer._layerNumber = trackerTopology.pxfDisk(detId);
31 }
32 //TIB
33 else if (subdet == StripSubdetector::TIB) {
34     trackingLayer._subDet = TrackingLayer::Det::TIB;
35     trackingLayer._side = TrackingLayer::Side::BARREL;
36     trackingLayer._layerNumber = trackerTopology.tibLayer(detId);
37 }
38 //TID
39 else if (subdet == StripSubdetector::TID) {
40     trackingLayer._subDet = TrackingLayer::Det::TID;
41     if (trackerTopology.tidSide(detId) == 1) {
42         trackingLayer._side = TrackingLayer::Side::NEG_ENDCAP;
43     } else if (trackerTopology.tidSide(detId) == 2) {
44         trackingLayer._side = TrackingLayer::Side::POS_ENDCAP;
45     } else {
46         throw cms::Exception("FastSimulation/Tracking")
47             << "Tracker hit for seeding in TID seems neither on positive nor on negative disk side: "
48             << trackerTopology.print(detId).c_str();
49     }
50     trackingLayer._layerNumber = trackerTopology.tidWheel(detId);
51 }
52 //TOB
53 else if (subdet == StripSubdetector::TOB) {
```

# Basic idea

- Machine learning approach:
  - Rephrase task as a minimisation problem..
  - ..and “simply” solve:

$$\theta^* = \operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\mathcal{L}(f_{\theta}(\mathbf{x}), \mathbf{x})]$$

- Will now go step-by-step to understand the underlying ideas, focusing on neural networks.

# General Strategy

Loss function  $\mathcal{L}$

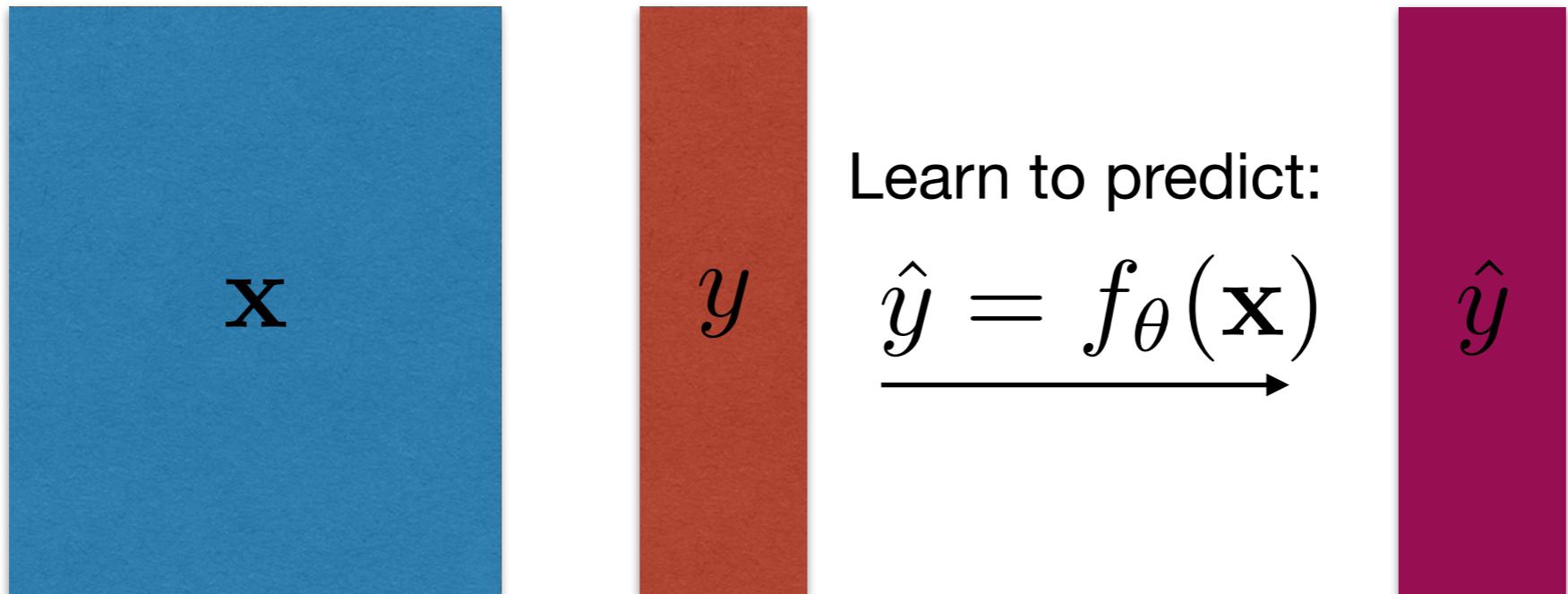
- Define an optimisation target (loss function)

# Loss function: Supervised

## Supervised Learning:

Attempt to infer some target (*truth label*):  
classification, regression (often also clustering/inference)

Use training data with known labels  
(often from Monte Carlo simulation)



observable features  
such as kinematics,  
tracks,...

truth label  
(e.g. true energy)

predicted energy

**Regression: Minimize mean squared error:**

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

# Classification Tasks

*Distinguish a pair of classes (binary) or several (multi-class).*

*Images*



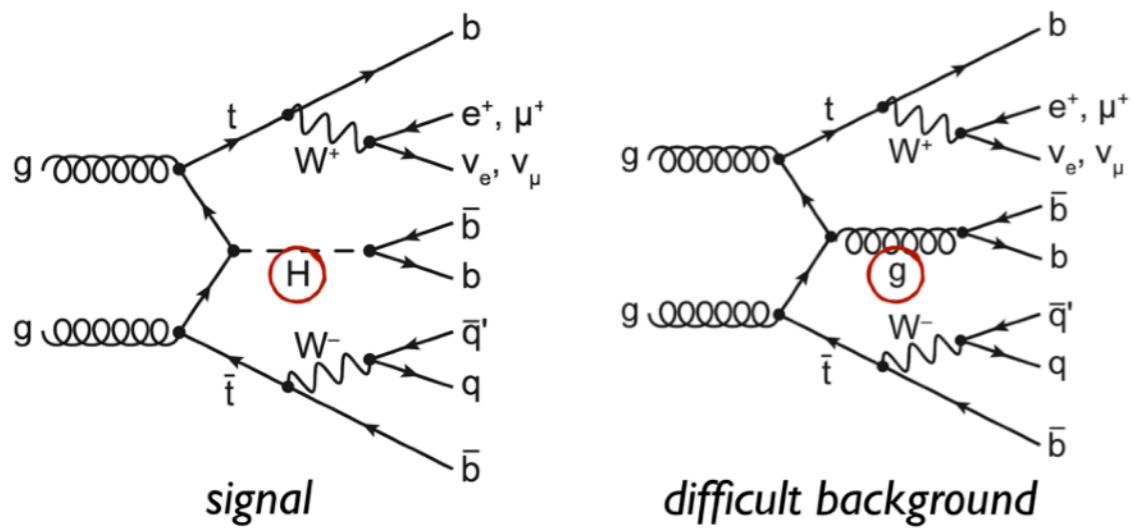
**Loss function: Cross entropy**

$$\mathcal{L} = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

# Classification Tasks

*Distinguish a pair of classes (binary) or several (multi-class).*

## Event Topologies



## Particles

Standard-Modell der Elementarteilchen									
Drei Generationen der Materie (Fermionen)			Wechselwirkungen (Bosonen)						
Masse	Ladung	Spin							
$\approx 2.2 \text{ GeV}/c^2$	$\frac{2}{3}$	$\frac{1}{2}$	u	c	t	g	H		
$\approx 4.7 \text{ MeV}/c^2$	$-\frac{1}{3}$	$\frac{1}{2}$	d	s	b	$\gamma$			
$\approx 0.511 \text{ MeV}/c^2$	-1	$\frac{1}{2}$	e	$\mu$	$\tau$	Z			
$<1.0 \text{ eV}/c^2$	0	$\frac{1}{2}$	$\nu_e$	$\nu_\mu$	$\nu_\tau$	W			
QUARKS									
LEPTONEN									
SKALARBOSONEN									
EICHBOSONEN									
VEKTORBOSONEN									

**Loss function: Cross entropy**

$$\mathcal{L} = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

# Cross Entropy

- Example: Binary classification
  - 2 Classes with labels:  $y_k = 0, 1$
  - Build the NN such that we the output is a probability:  $\hat{y}(x_k) \in [0,1]$
- -> Bernoulli Trial with probabilities p and q = 1 – p
- Likelihood:



Jacob Bernoulli 1654-1705

$$L = \prod_k^N [y_k \cdot \hat{y}(x_k) + (1 - y_k) \cdot (1 - \hat{y}(x_k))]$$

- We want to minimize the average negative log-likelihood:

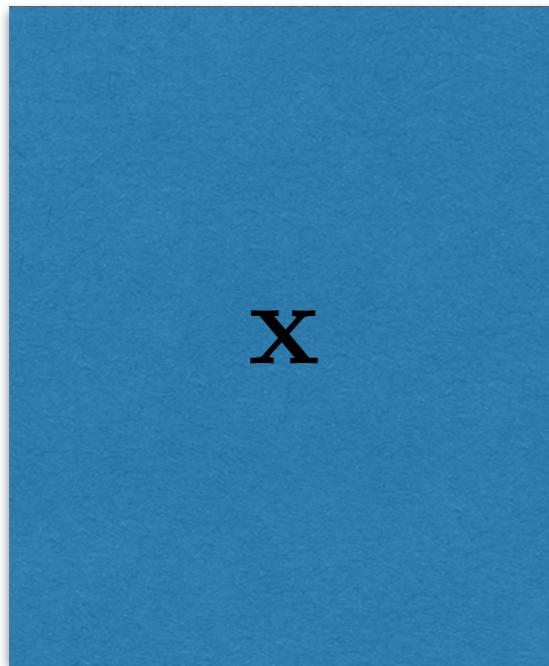
$$-\frac{1}{N} \log(L) = -\frac{1}{N} \sum_k^N [y_k \cdot \log(\hat{y}(x_k)) + (1 - y_k) \cdot \log(1 - \hat{y}(x_k))]$$

# Loss function: Unsupervised

## Unsupervised Learning:

No target, learn the probability distribution (directly from data)

Can use for sampling, anomaly detection, unfolding, ...



Learn to predict:

$$\hat{p}(\mathbf{x}) = f_{\theta}(\mathbf{x}) \longrightarrow p(\mathbf{x})$$

True probability density

**Distribution learning: Maximise likelihood (minimize log-likelihood):**  
(either directly or with approximations)

$$\mathcal{L} = -\log (\hat{p}(\mathbf{x}))$$

**\*There also exists a number of other less-than-supervised approaches (weakly supervised learning, semi-supervised learning, ...) Not so important for now.**

# General Strategy

Loss function  $\mathcal{L}$  ✓

Neural network  $f_\theta$

Parameters  $\theta$

Opt. Parameters  $\theta^*$

Data  $\mathbf{x}$

Data distribution  $p(\mathbf{x})$

$$\theta^* = \operatorname{argmin}_\theta \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\mathcal{L}(f_\theta(\mathbf{x}), \mathbf{x})]$$

# General Strategy

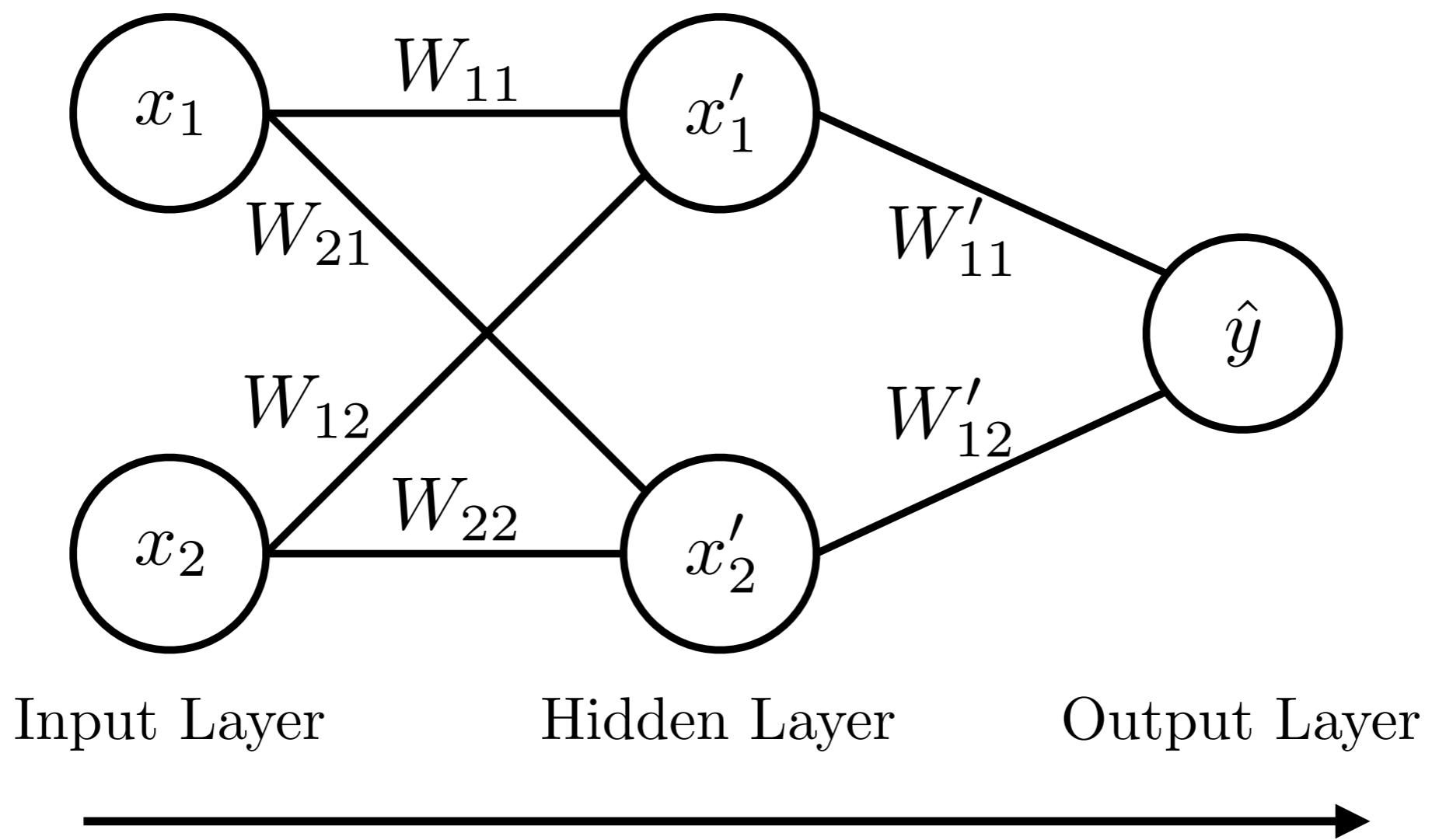
Loss function  $\mathcal{L}$   
Neural network  $f_\theta$   
Parameters  $\theta$

- Define an optimisation target (loss function)
- Choose a non-linear, expressive, parametrised function (e.g. a neural network)

# Neural networks

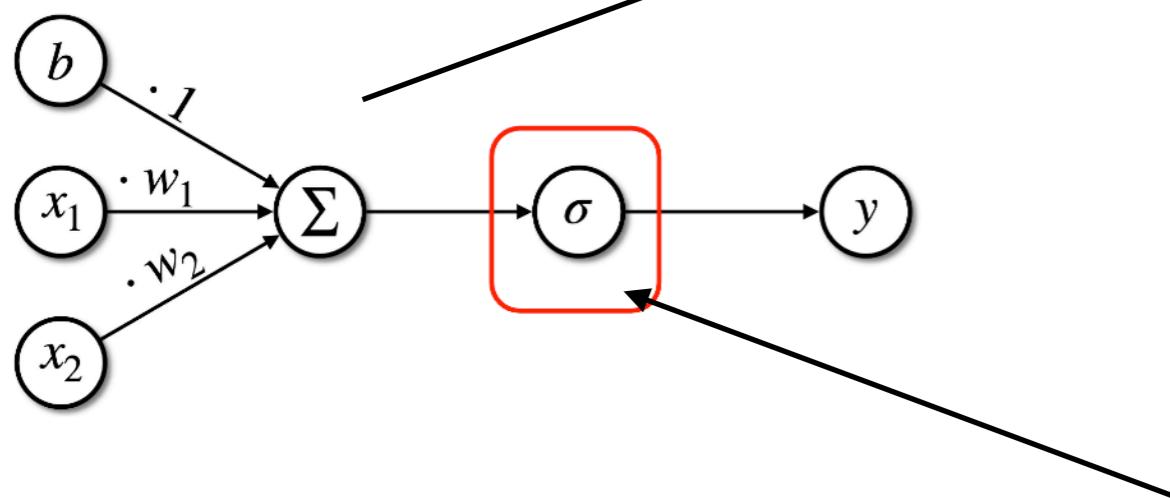
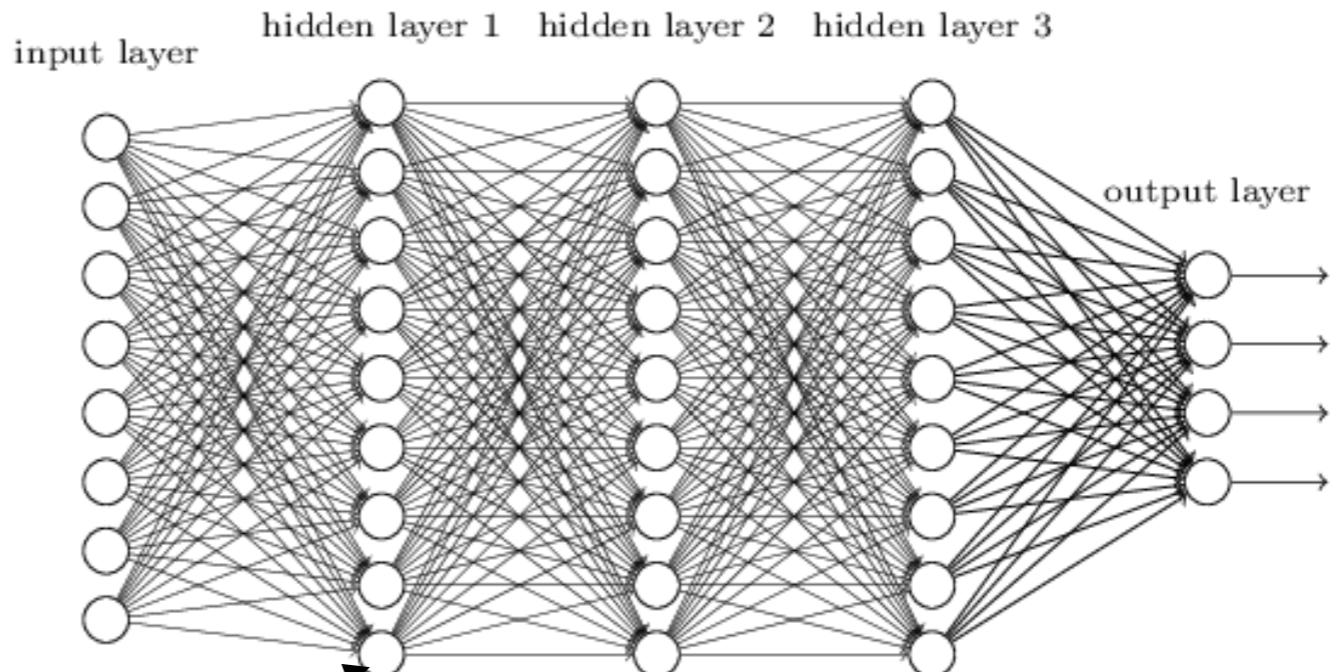
We need an expressive function  
(universal approximator)  
with tuneable parameters and useful  
implicit biases

$$x'_1 = f(W_{11} \cdot x_1 + W_{12} \cdot x_2)$$



# Neural Networks

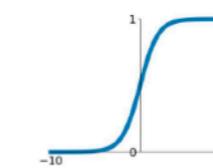
We need an expressive function  
 (universal approximator)  
 with tuneable parameters:  
**Neural networks**



## Activation Functions

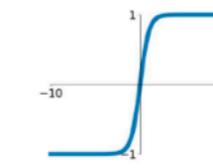
### Sigmoid

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



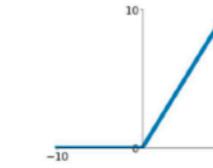
### tanh

$$\tanh(x)$$



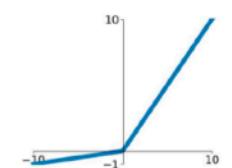
### ReLU

$$\max(0, x)$$



### Leaky ReLU

$$\max(0.1x, x)$$

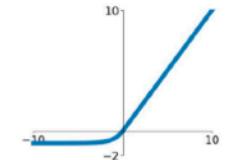


### Maxout

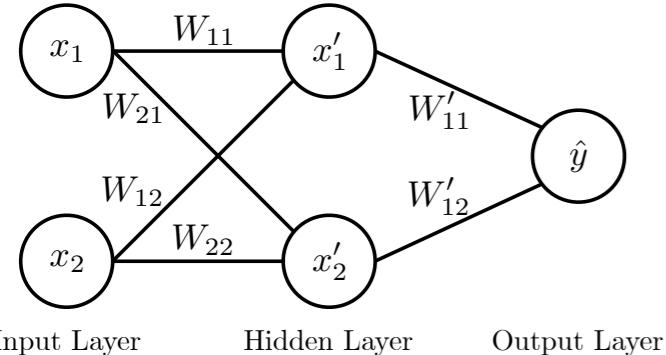
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

### ELU

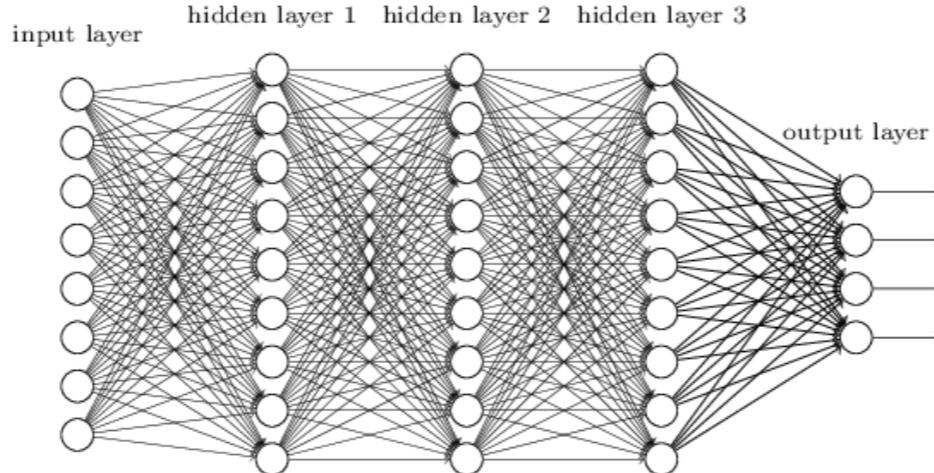
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Complexity



**6 weights**



**300 weights**

**Deep Learning:**  
**Complex network + low level inputs**

**25 million weights:**  
*2016 state of the art for  
image classification*

stage	output	ResNet-50	ResNeXt-50 (32×4d)
conv1	112×112	7×7, 64, stride 2	7×7, 64, stride 2
conv2	56×56	3×3 max pool, stride 2 [ 1×1, 64 3×3, 64 1×1, 256 ] ×3	3×3 max pool, stride 2 [ 1×1, 128 3×3, 128, C=32 1×1, 256 ] ×3
conv3	28×28	[ 1×1, 128 3×3, 128 1×1, 512 ] ×4	[ 1×1, 256 3×3, 256, C=32 1×1, 512 ] ×4
conv4	14×14	[ 1×1, 256 3×3, 256 1×1, 1024 ] ×6	[ 1×1, 512 3×3, 512, C=32 1×1, 1024 ] ×6
conv5	7×7	[ 1×1, 512 3×3, 512 1×1, 2048 ] ×3	[ 1×1, 1024 3×3, 1024, C=32 1×1, 2048 ] ×3
	1×1	global average pool 1000-d fc, softmax	global average pool 1000-d fc, softmax
# params.		<b>25.5×10<sup>6</sup></b>	<b>25.0×10<sup>6</sup></b>
FLOPs		<b>4.1×10<sup>9</sup></b>	<b>4.2×10<sup>9</sup></b>

**175 billion weights: 2020**  
**GPT-3 text model**  
**(GPT-4 ~1.8 trillion  
weights)**

Model Name	n <sub>params</sub>	n <sub>layers</sub>	d <sub>model</sub>	n <sub>heads</sub>	d <sub>head</sub>	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	$6.0 \times 10^{-4}$
GPT-3 Medium	350M	24	1024	16	64	0.5M	$3.0 \times 10^{-4}$
GPT-3 Large	760M	24	1536	16	96	0.5M	$2.5 \times 10^{-4}$
GPT-3 XL	1.3B	24	2048	24	128	1M	$2.0 \times 10^{-4}$
GPT-3 2.7B	2.7B	32	2560	32	80	1M	$1.6 \times 10^{-4}$
GPT-3 6.7B	6.7B	32	4096	32	128	2M	$1.2 \times 10^{-4}$
GPT-3 13B	13.0B	40	5140	40	128	2M	$1.0 \times 10^{-4}$
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128	3.2M	$0.6 \times 10^{-4}$

# General Strategy

Loss function  $\mathcal{L}$  ✓

Neural network  $f_\theta$  ✓

Parameters  $\theta$  ✓

Opt. Parameters  $\theta^*$

Data  $\mathbf{x}$

Data distribution  $p(\mathbf{x})$

- Define an optimisation target (loss function)
- Choose a non-linear, expressive, parametrised function (e.g. a neural network)
- Use training data to optimise parameters, then apply to new examples

$$\theta^* = \operatorname{argmin}_\theta \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\mathcal{L}(f_\theta(\mathbf{x}), \mathbf{x})]$$

# How do networks learn?

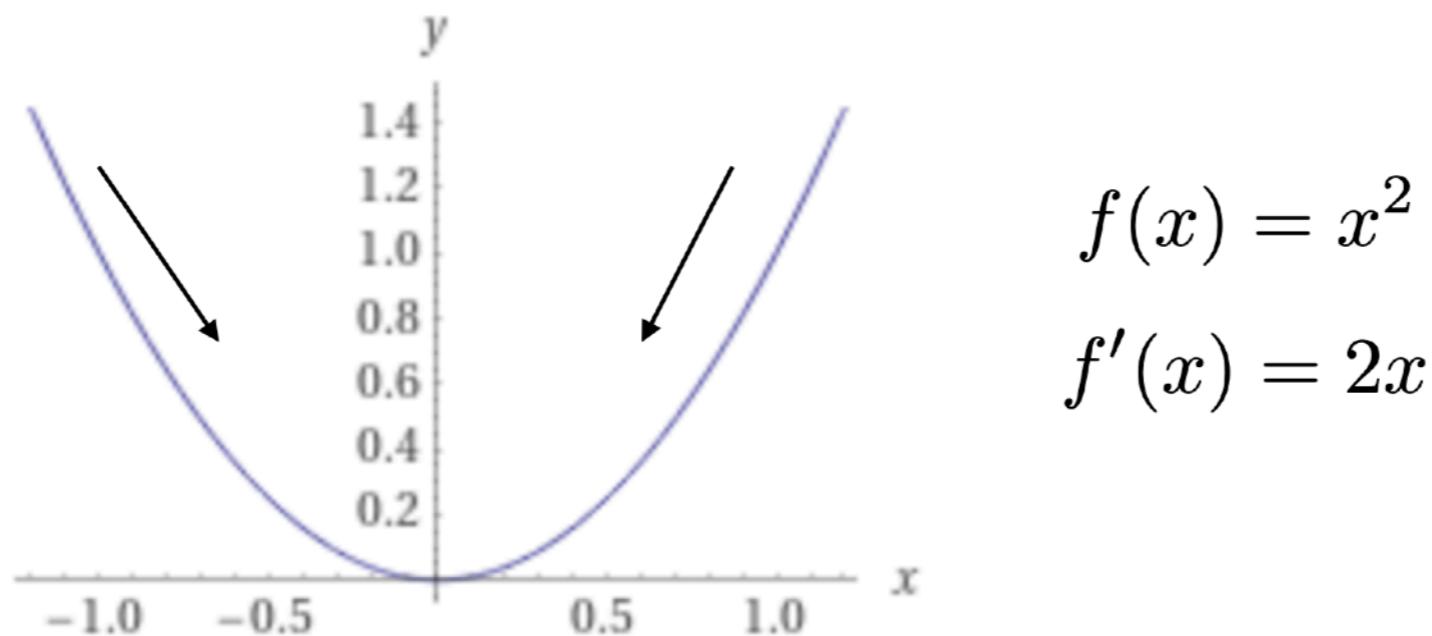
- *Backpropagation + Gradient descent*
- Important: Loss function needs to be differentiable
  - (Or find a differentiable approximation)
- Pass input ( $x_1, x_2, \dots$ ) to networks
- From output calculate loss function  
Find gradient of loss function with respect to weights
- Use gradient to find new weights

$$\theta_{t+1} = \theta_t - \eta \frac{\partial \mathcal{L}}{\partial \theta_t} = \theta_t - \eta \nabla \mathcal{L}$$

*Learning rate*

- In practice, handled by optimise algorithm (e.g. Adam)

# How do networks learn?



$$x < 0 \implies -f'(x) > 0$$

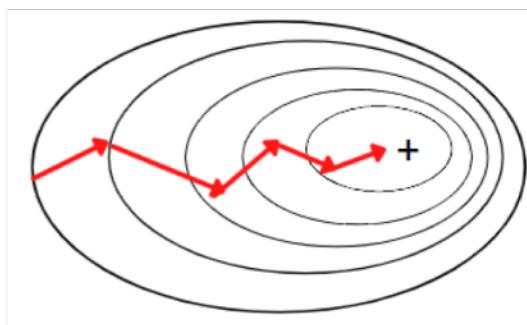
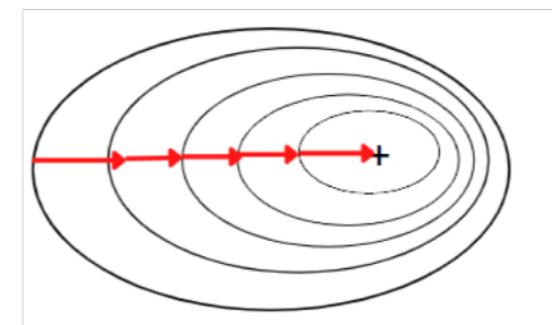
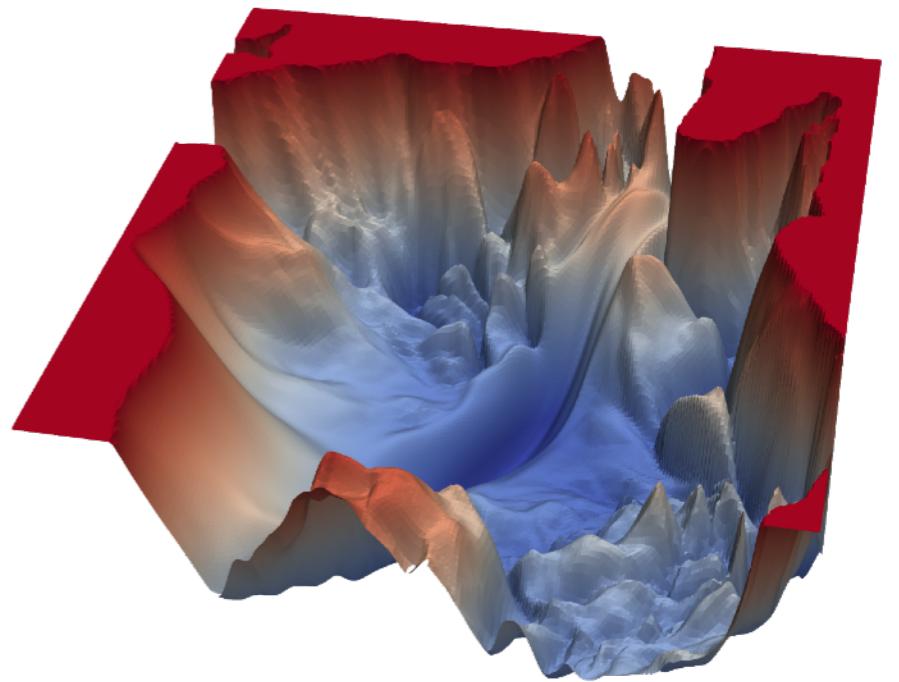
“Go right”

$$x > 0 \implies -f'(x) < 0$$

“Go left”

# Batches and Epochs

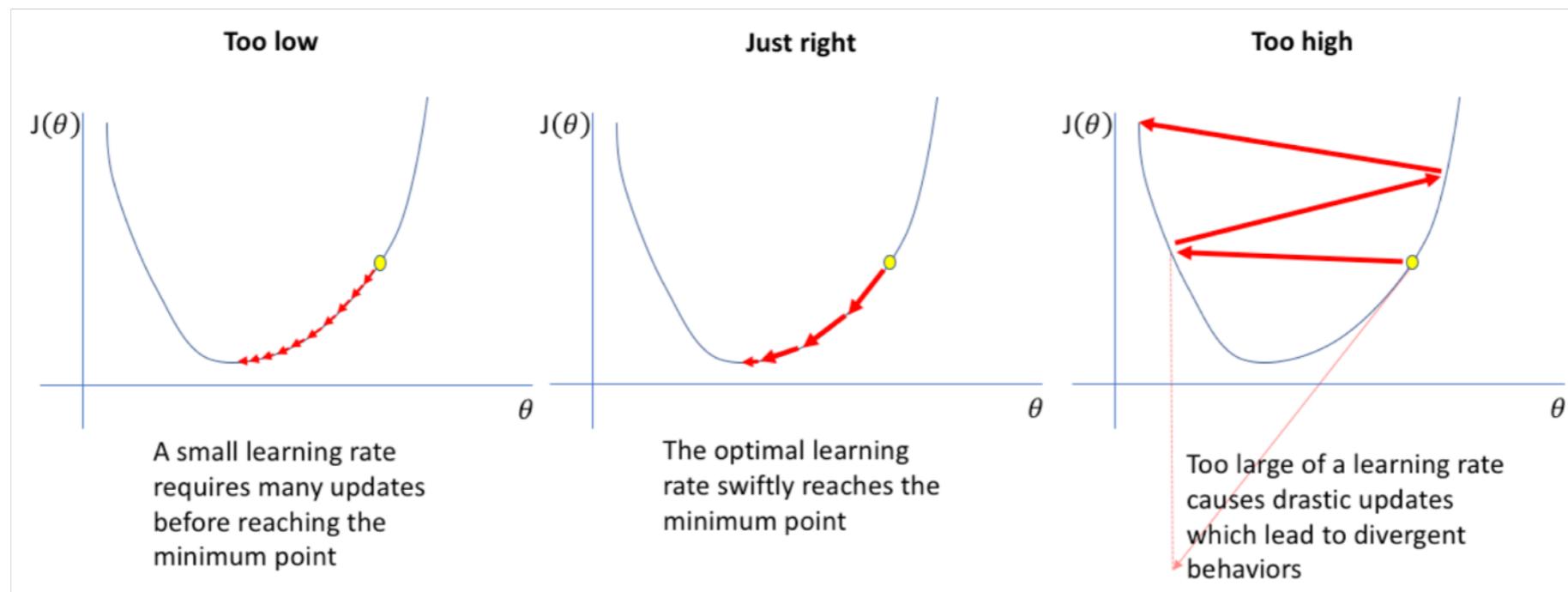
- The loss depends on  $O(1k-1b)$  parameters
- To minimize it these parameters need to be adjusted
- For large networks the loss landscape can get very complicated with many local minima
- Better convergence
  - Multiple passes (epochs) over training data
  - Split training data in batches with a minimization call and model update after each batch



# Learning Rate

- Many popular minimization algorithms have a tunable parameter called Learning rate  $\alpha$
- Should be chosen such that the training smoothly converges

$$\theta \rightarrow \theta - \alpha \frac{\partial L(\theta)}{\partial \theta}$$



<https://www.jeremyjordan.me/nn-learning-rate/>

# Adam

$$m_w^{(t+1)} \leftarrow \beta_1 m_w^{(t)} + (1 - \beta_1) \nabla_w L^{(t)}$$

$$v_w^{(t+1)} \leftarrow \beta_2 v_w^{(t)} + (1 - \beta_2) (\nabla_w L^{(t)})^2$$

$$\hat{m}_w = \frac{m_w^{(t+1)}}{1 - \beta_1^t}$$

$$\hat{v}_w = \frac{v_w^{(t+1)}}{1 - \beta_2^t}$$

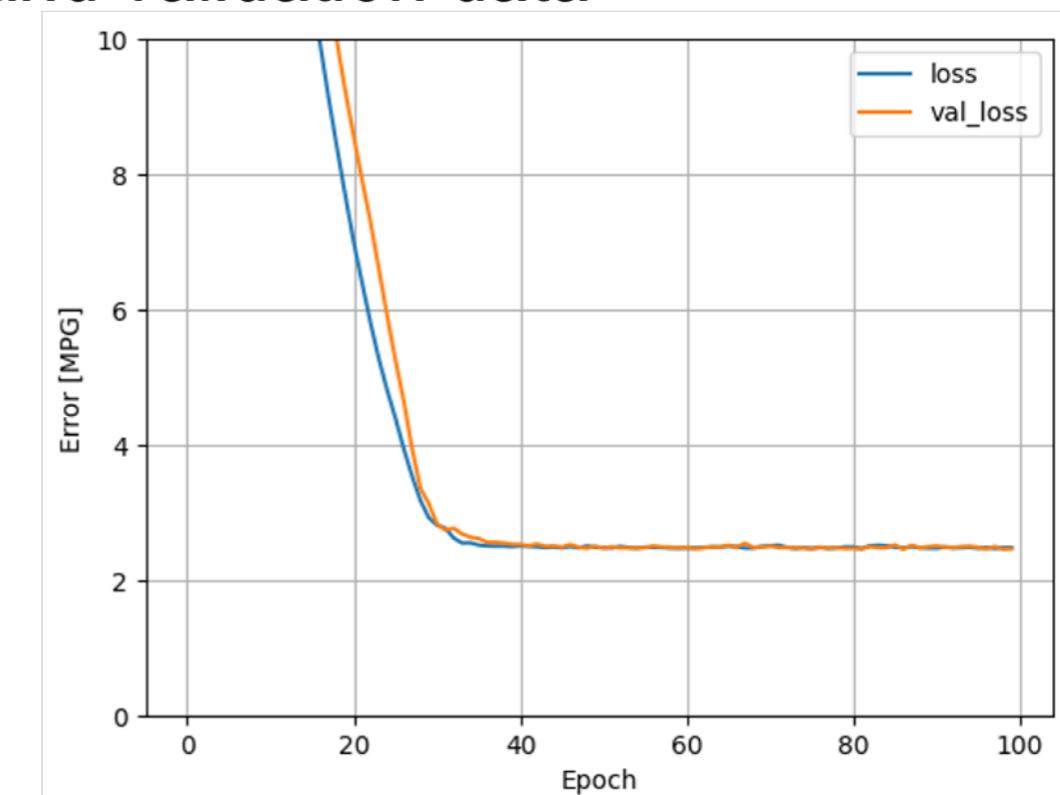
$$w^{(t+1)} \leftarrow w^{(t)} - \eta \frac{\hat{m}_w}{\sqrt{\hat{v}_w} + \epsilon}$$

Maintain history of previous changes:  
momentum

Exponential moving average of  
gradient and gradient squares

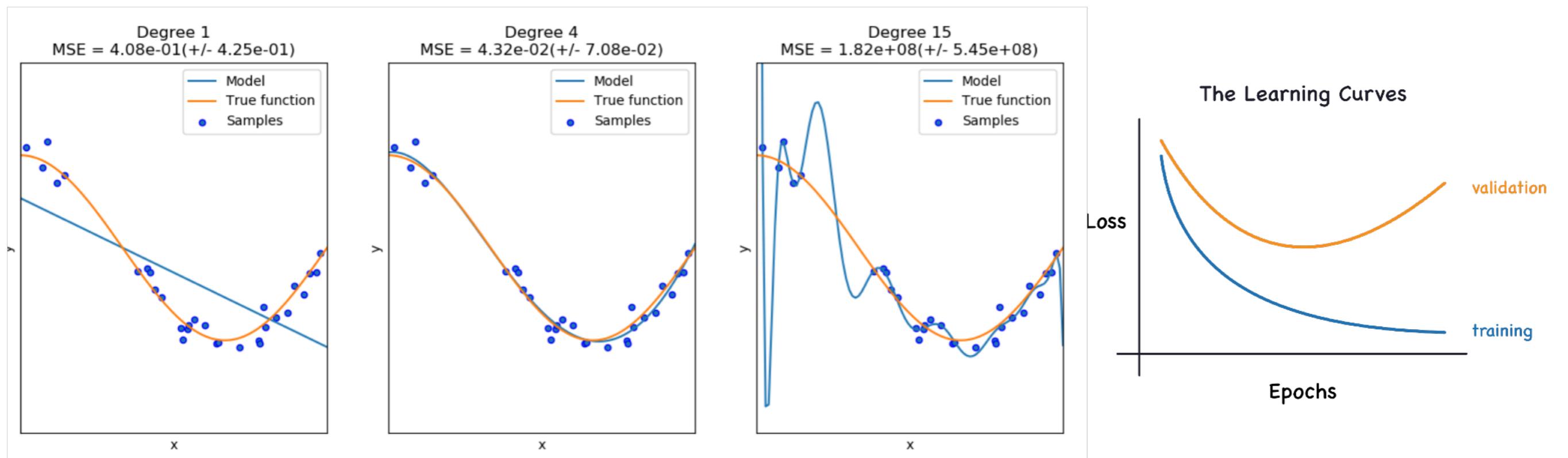
# The loss curve

- Once the training is done it is time to analyze the training process and the performance of the model
- For the training process it is a good idea to look at the loss as a function of the number of epochs – for the training data and validation data
- Ideally the training converges with similar loss values for the validation and training data



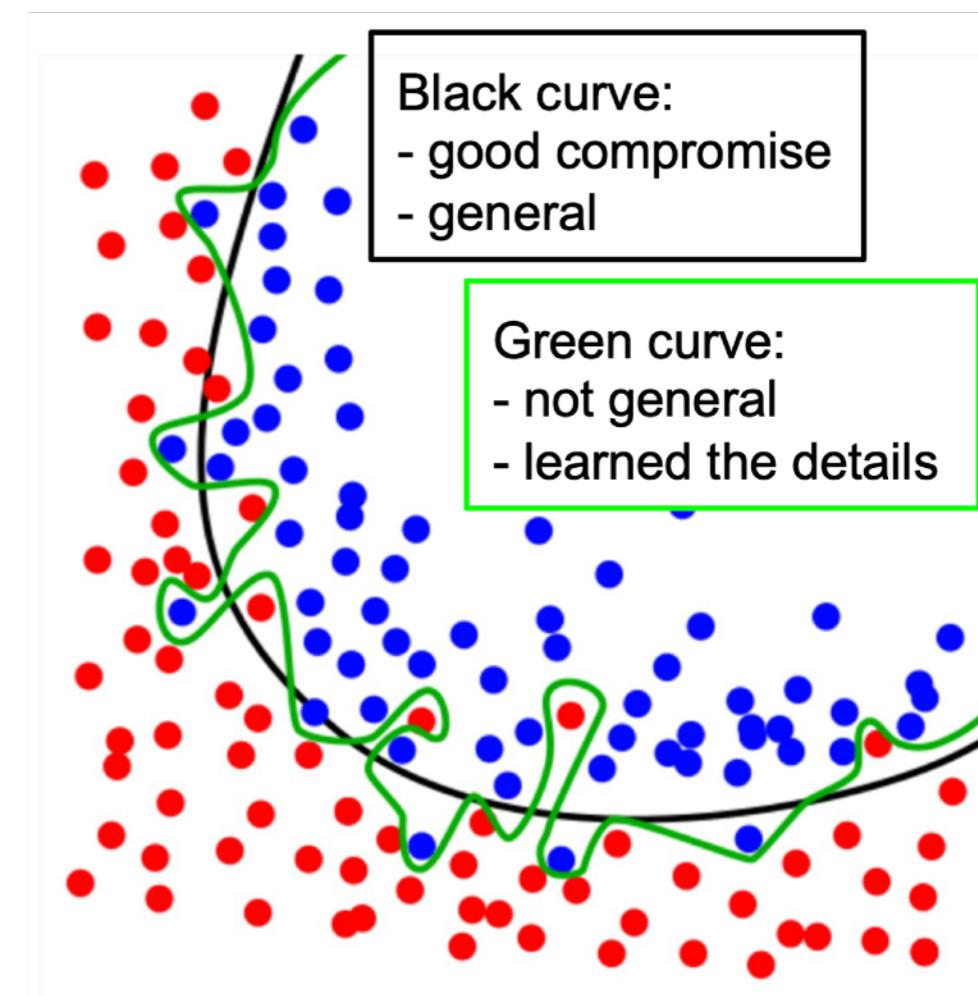
# Overfitting

- If the model has too many parameters it can overfit the data
- This leads to a worse performance on data with different statistical fluctuations than the training data
- -> the training and validation losses diverge



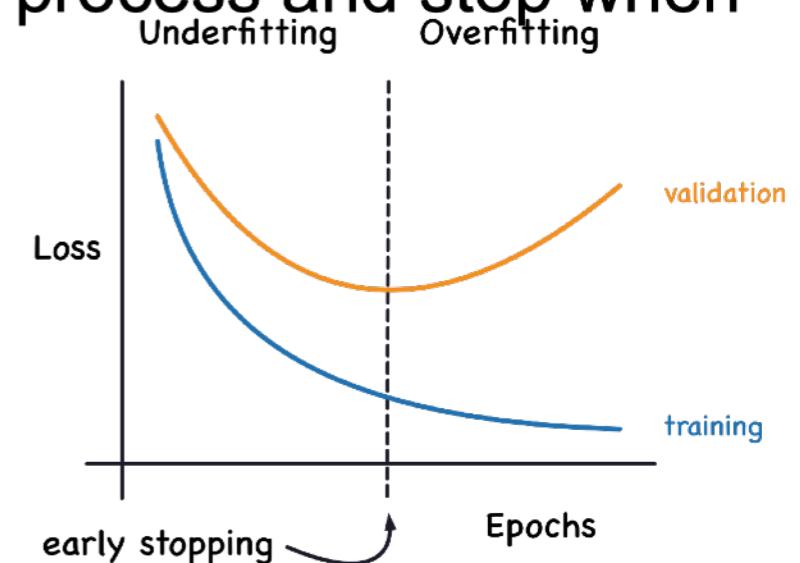
# Overfitting

- Also an issue in classification tasks (which are after all a form of regression of discrete values)
- Overfitting can be interpreted as remembering the training samples
- We can see that the performance would change if the dots move around even a bit

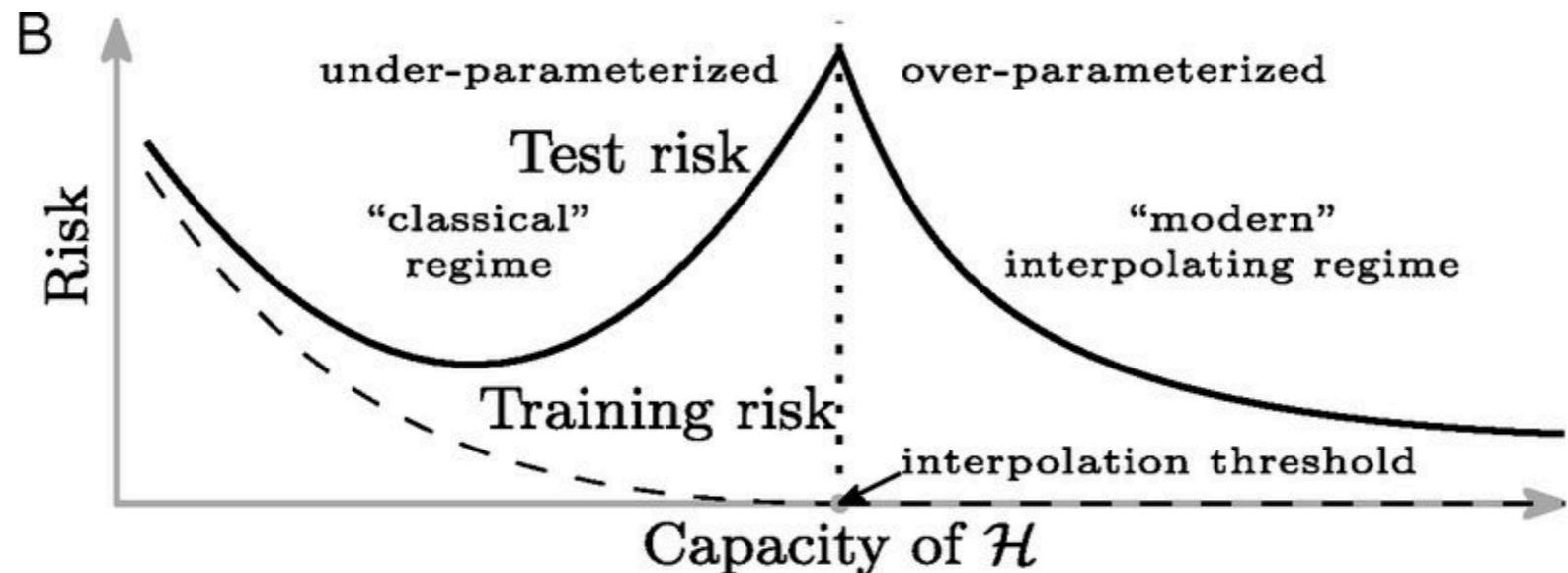


# Overfitting - Mitigation

- There are multiple ways to mitigate overfitting such as dropout, regularization or early stopping
- Dropout:
  - Add layers that randomly discard a fraction of inputs in training phase
- Early stopping:
  - Monitor training and validation loss during the training process and stop when they diverge
- Regularisation:
  - Add penalty for large weight values in loss function



# Side note: Overfitting



- Recent observation in computer science
  - Extremely high capacity does not guarantee overfitting
- Even more recent:
  - Grokking: Sudden generalisation of a trained model
- Active research into foundations, does not stop us from using it

# General Strategy

Loss function  $\mathcal{L}$  ✓

Neural network  $f_\theta$  ✓

Parameters  $\theta$  ✓

Opt. Parameters  $\theta^*$  ✓

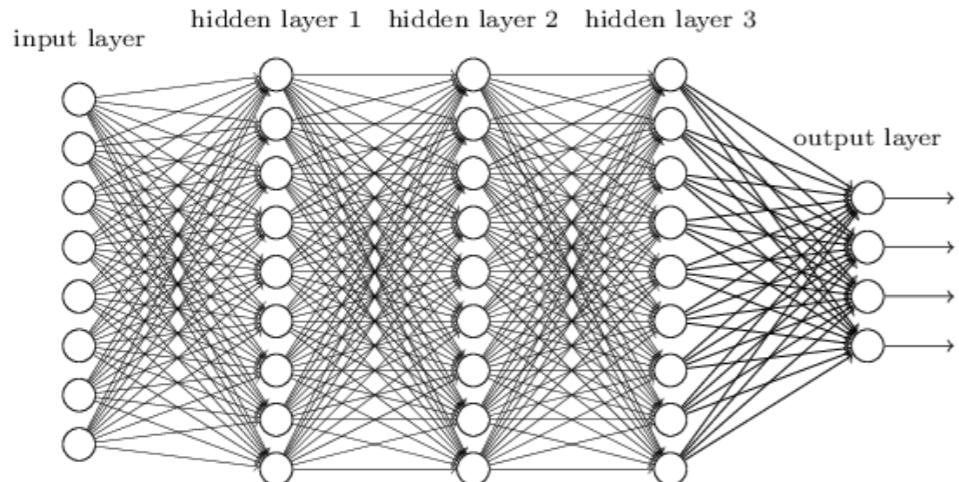
Data  $\mathbf{x}$

Data distribution  $p(\mathbf{x})$

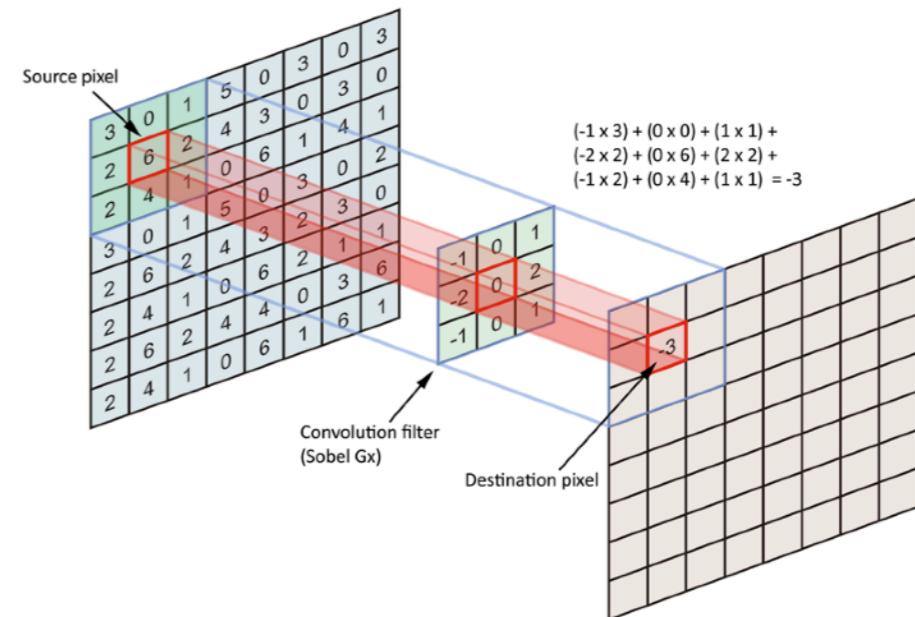
$$\theta^* = \operatorname{argmin}_\theta \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\mathcal{L}(f_\theta(\mathbf{x}), \mathbf{x})]$$

# Data Representation

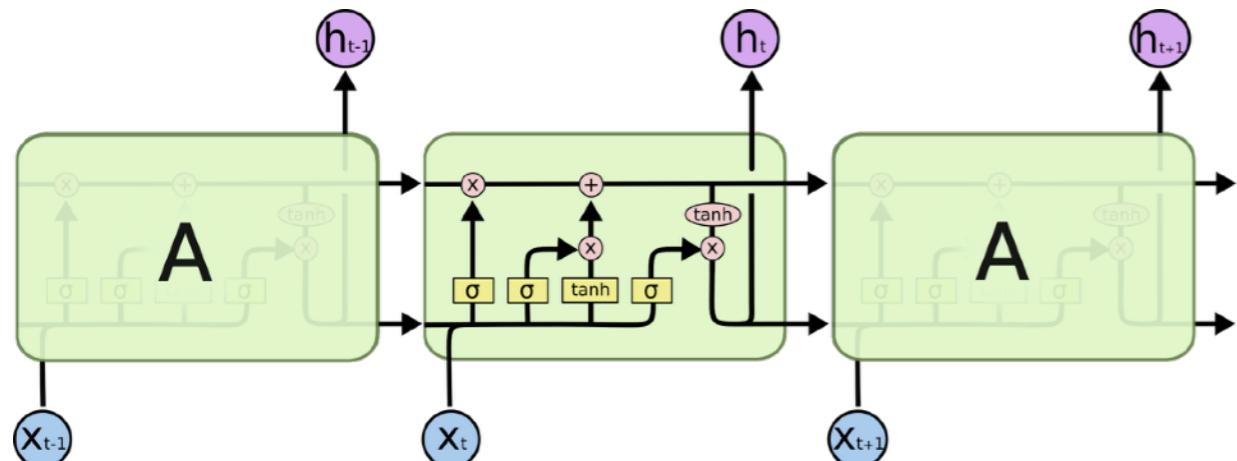
High-level/no structure:  
Fully Connected



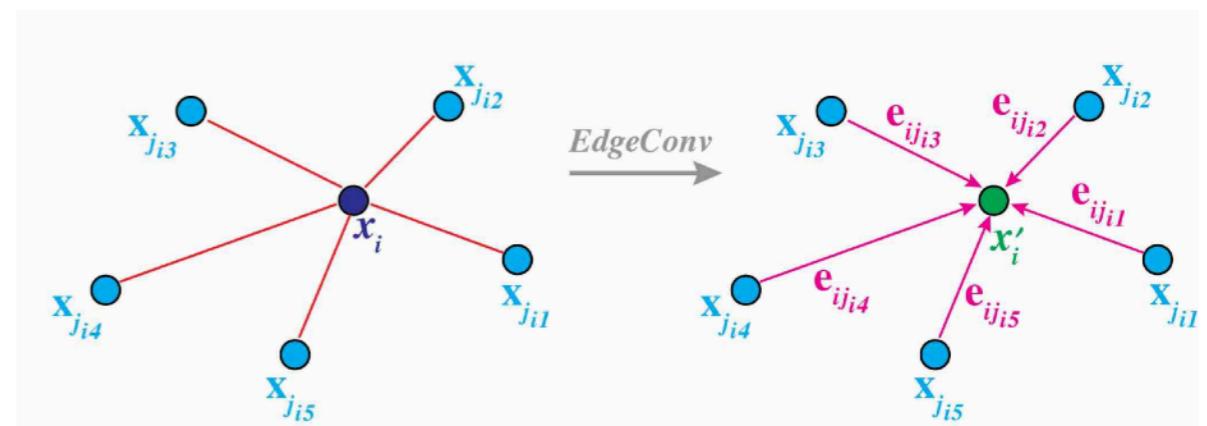
Regular grid: Convolution



Sequence/Time Series: Recurrent



Point cloud: Sets & Graphs



# General Strategy

Loss function  $\mathcal{L}$  ✓

Neural network  $f_\theta$  ✓

Parameters  $\theta$  ✓

Opt. Parameters  $\theta^*$  ✓

Data  $\mathbf{x}$  ✓

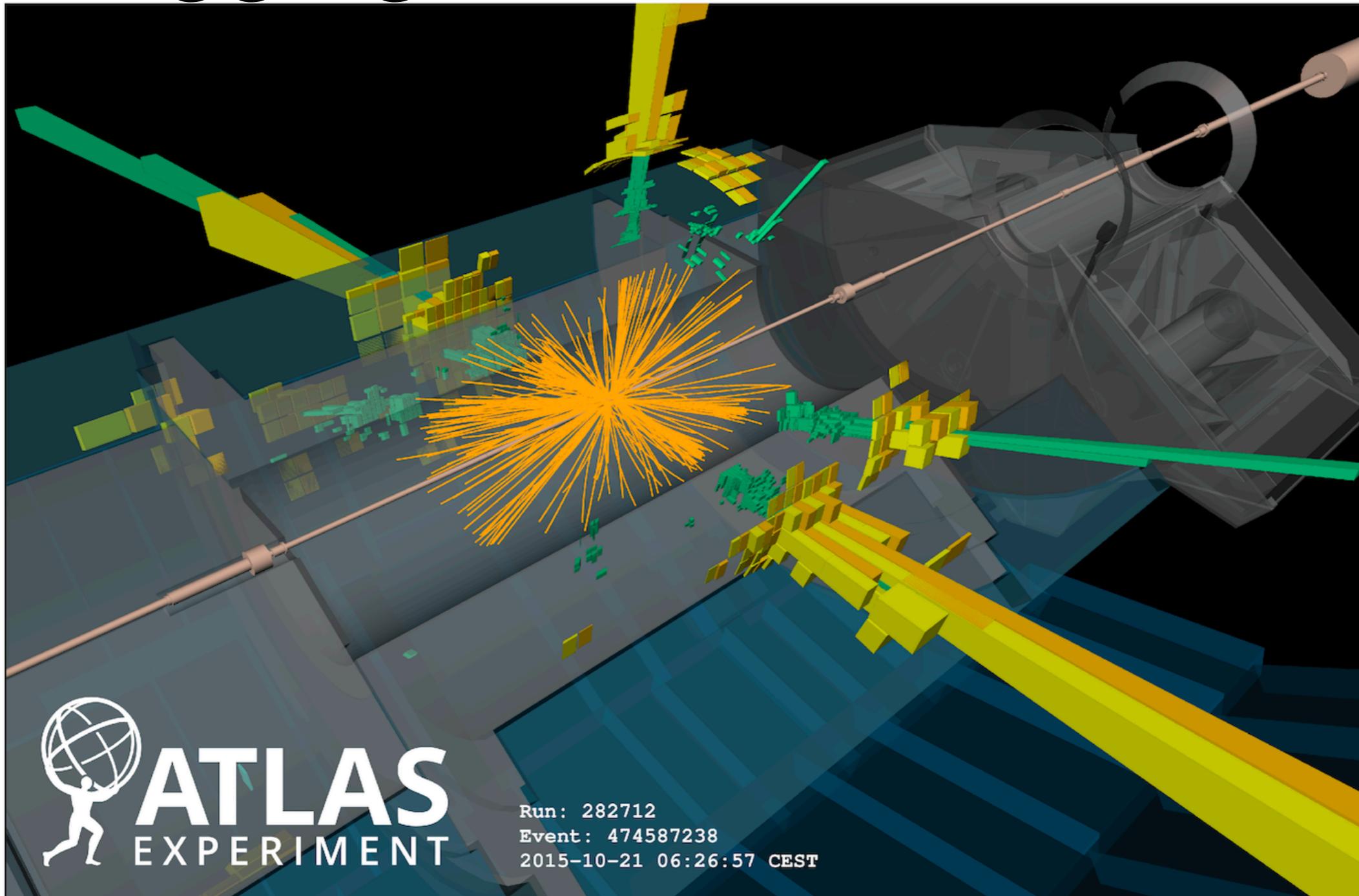
Data distribution  $p(\mathbf{x})$  ✓

$$\theta^* = \operatorname{argmin}_\theta \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\mathcal{L}(f_\theta(\mathbf{x}), \mathbf{x})]$$

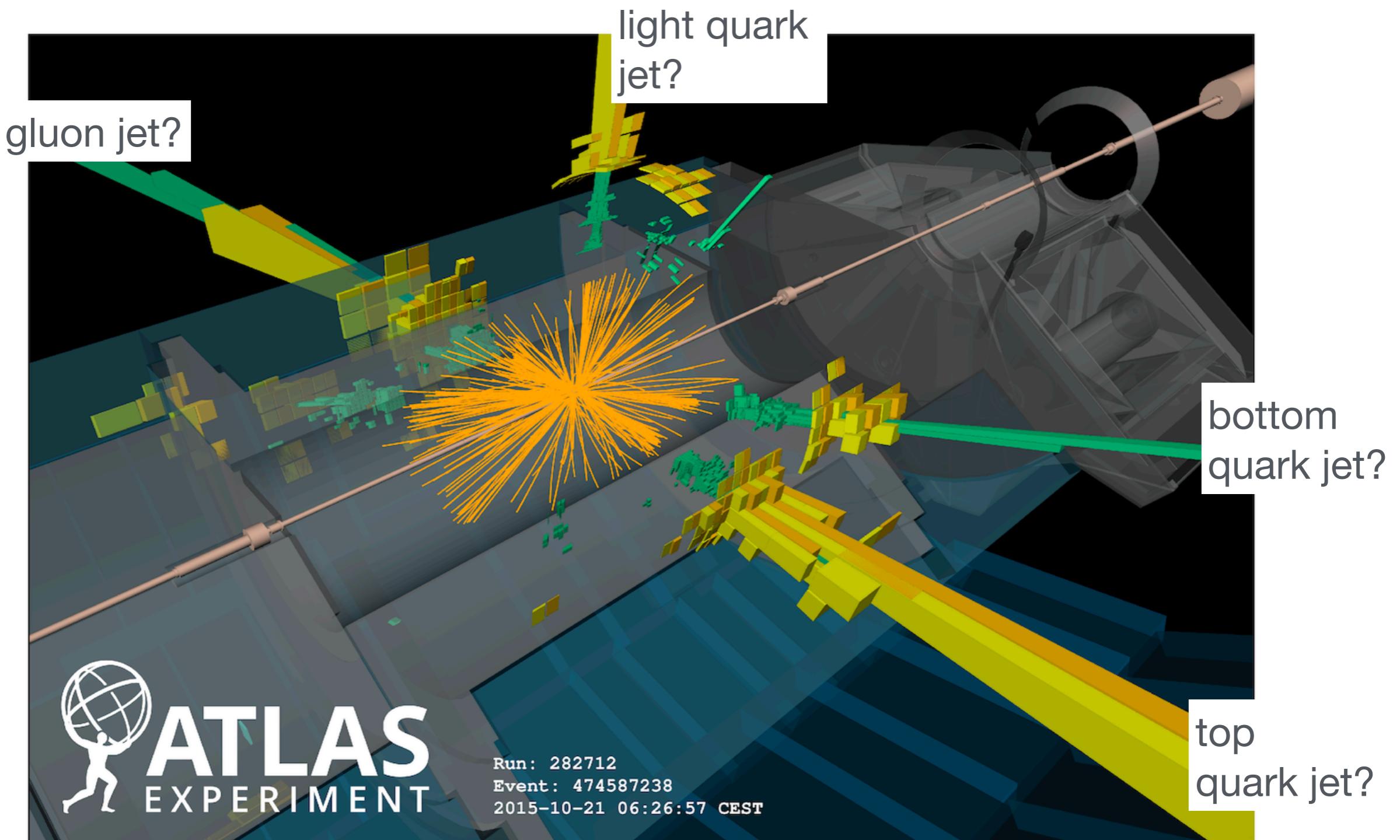
# Summary so far

- Key idea of deep learning:
  - Define a optimisation target and use an expressive function (neural network) to minimise it using gradient descent
- Let's look at a typical problem in particle physics...

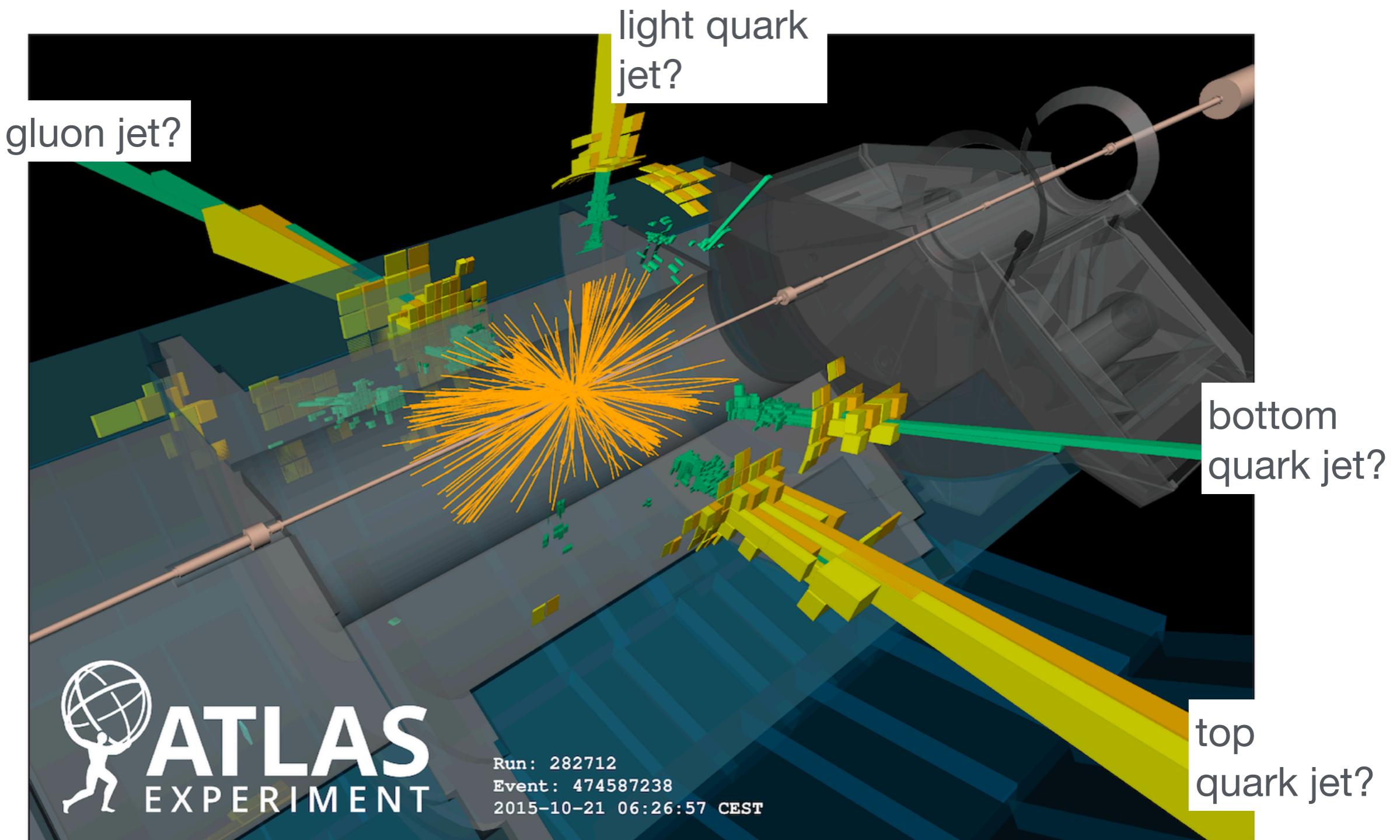
# Jet Tagging



A jet is a  
collimated shower of particles in the detector



We want to know  
which particle produced a jet



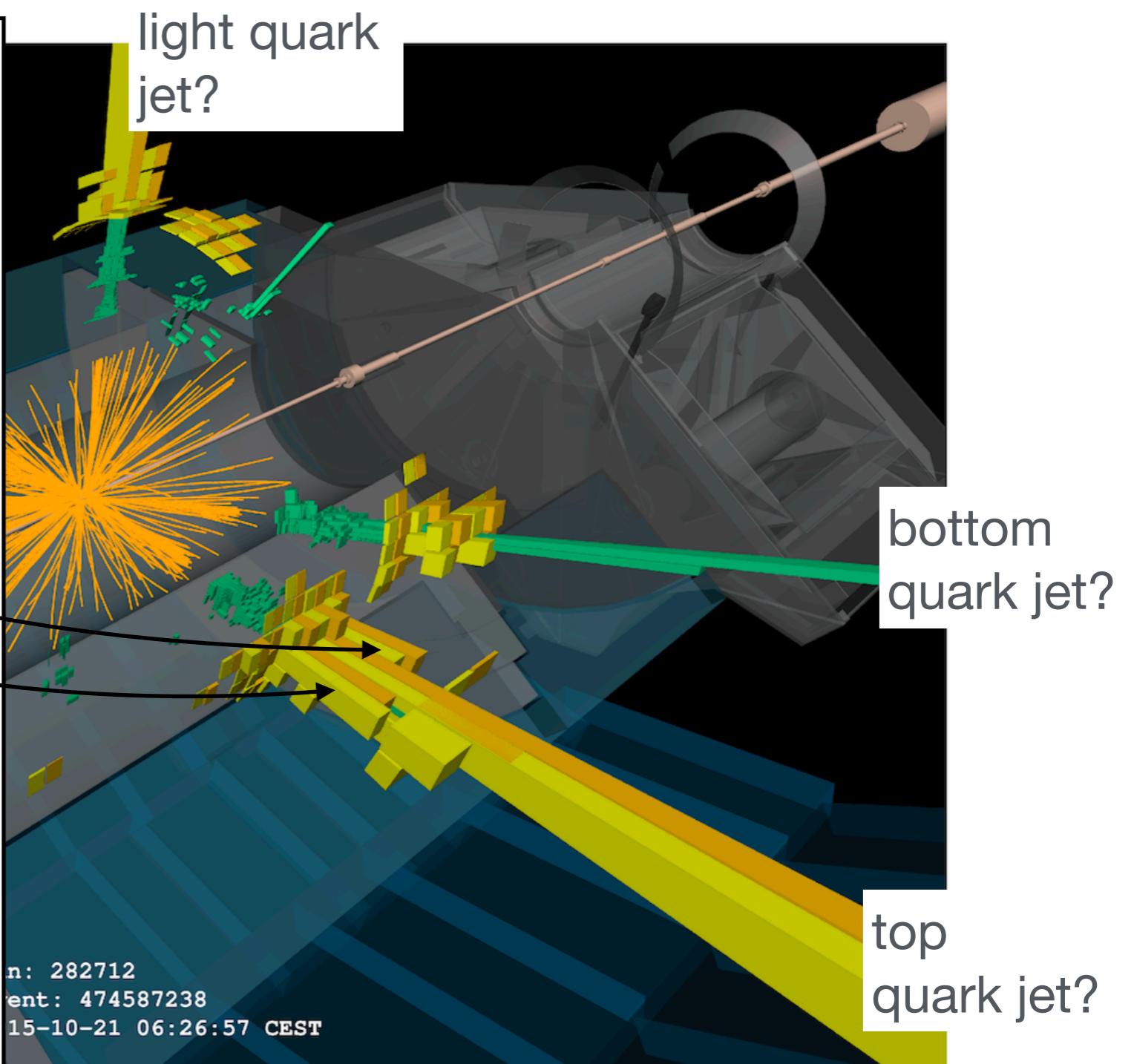
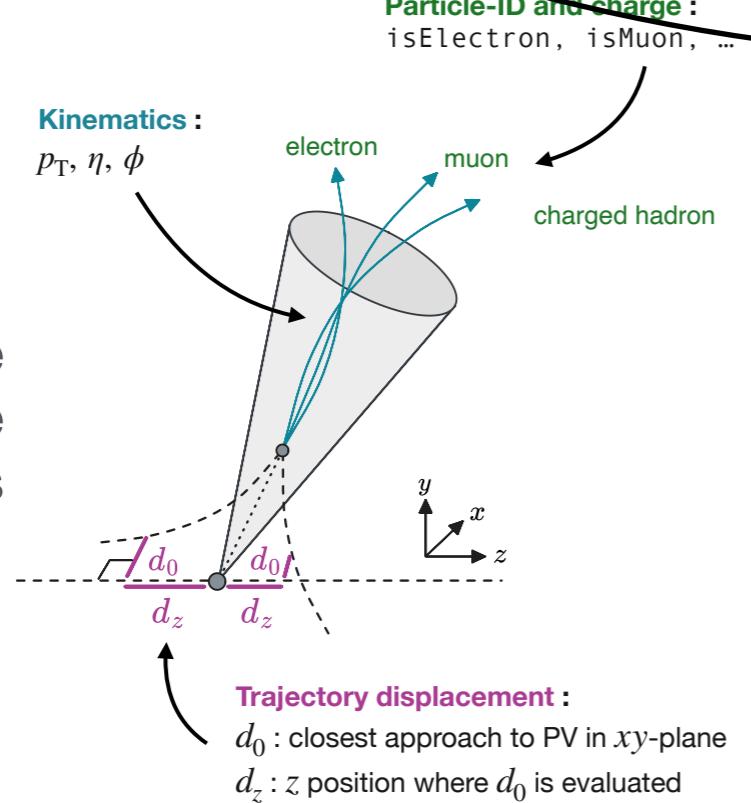
How to build ML algorithms for **complex, heterogenous** data?

Data most naturally viewed as **point cloud**:

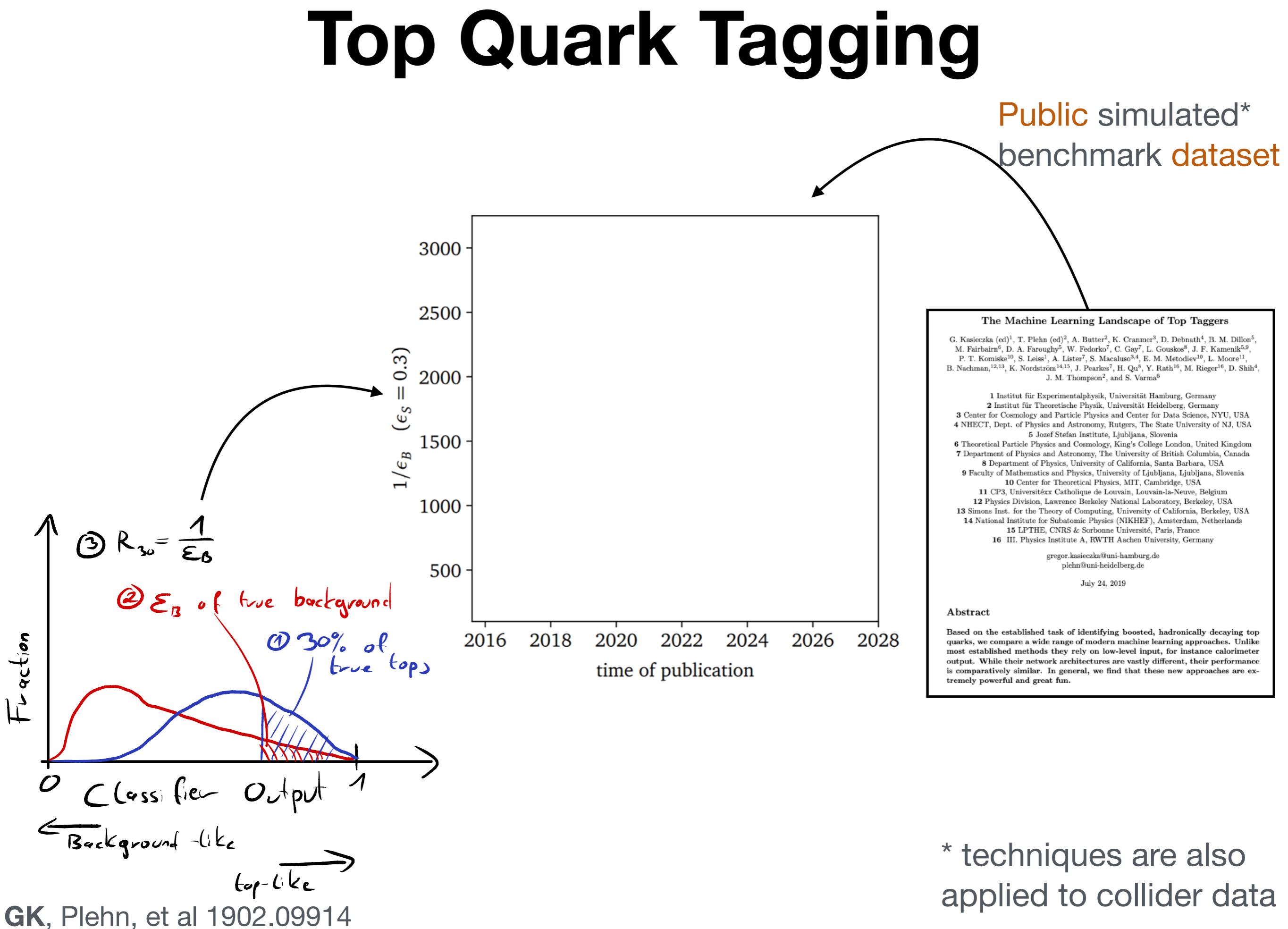
Each **input** (e.g. jet, event, ...) is a **set of k-dimensional vectors** (individual particles, hits, ...)

$$J_i = \{\vec{p_1}, \dots, \vec{p_n}\}$$

Example per-particle features



# Top Quark Tagging



# Landscape Dataset

- Open dataset for the development of better tagging algorithms for particle physics
- 2 million simulated examples
- Input: momentum sorted list of 200 particles/jet with 3 features/particle ( $p_x$ ,  $p_Y$ ,  $p_z$ )
- Perfect class labels:  
top jet or light quark/gluon jet
- Supervised learning problem

arXiv:1902.09914v3 [hep-ph] 23 Jul 2019

**SciPost Physics** **Submission**

**The Machine Learning Landscape of Top Taggers**

G. Kasieczka (ed)<sup>1</sup>, T. Plehn (ed)<sup>2</sup>, A. Butter<sup>2</sup>, K. Cranmer<sup>3</sup>, D. Debnath<sup>4</sup>, B. M. Dillon<sup>5</sup>, M. Fairbairn<sup>6</sup>, D. A. Faroughy<sup>5</sup>, W. Fedorko<sup>7</sup>, C. Gay<sup>7</sup>, L. Gouskos<sup>8</sup>, J. F. Kamenik<sup>5,9</sup>, P. T. Komiske<sup>10</sup>, S. Leiss<sup>1</sup>, A. Lister<sup>7</sup>, S. Macaluso<sup>3,4</sup>, E. M. Metodiev<sup>10</sup>, L. Moore<sup>11</sup>, B. Nachman,<sup>12,13</sup>, K. Nordström<sup>14,15</sup>, J. Pearkes<sup>7</sup>, H. Qu<sup>8</sup>, Y. Rath<sup>16</sup>, M. Rieger<sup>16</sup>, D. Shih<sup>4</sup>, J. M. Thompson<sup>2</sup>, and S. Varma<sup>6</sup>

**1** Institut für Experimentalphysik, Universität Hamburg, Germany  
**2** Institut für Theoretische Physik, Universität Heidelberg, Germany  
**3** Center for Cosmology and Particle Physics and Center for Data Science, NYU, USA  
**4** NHECT, Dept. of Physics and Astronomy, Rutgers, The State University of NJ, USA  
**5** Jozef Stefan Institute, Ljubljana, Slovenia  
**6** Theoretical Particle Physics and Cosmology, King's College London, United Kingdom  
**7** Department of Physics and Astronomy, The University of British Columbia, Canada  
**8** Department of Physics, University of California, Santa Barbara, USA  
**9** Faculty of Mathematics and Physics, University of Ljubljana, Ljubljana, Slovenia  
**10** Center for Theoretical Physics, MIT, Cambridge, USA  
**11** CP3, Université Catholique de Louvain, Louvain-la-Neuve, Belgium  
**12** Physics Division, Lawrence Berkeley National Laboratory, Berkeley, USA  
**13** Simons Inst. for the Theory of Computing, University of California, Berkeley, USA  
**14** National Institute for Subatomic Physics (NIKHEF), Amsterdam, Netherlands  
**15** LPTHE, CNRS & Sorbonne Université, Paris, France  
**16** III. Physics Institute A, RWTH Aachen University, Germany

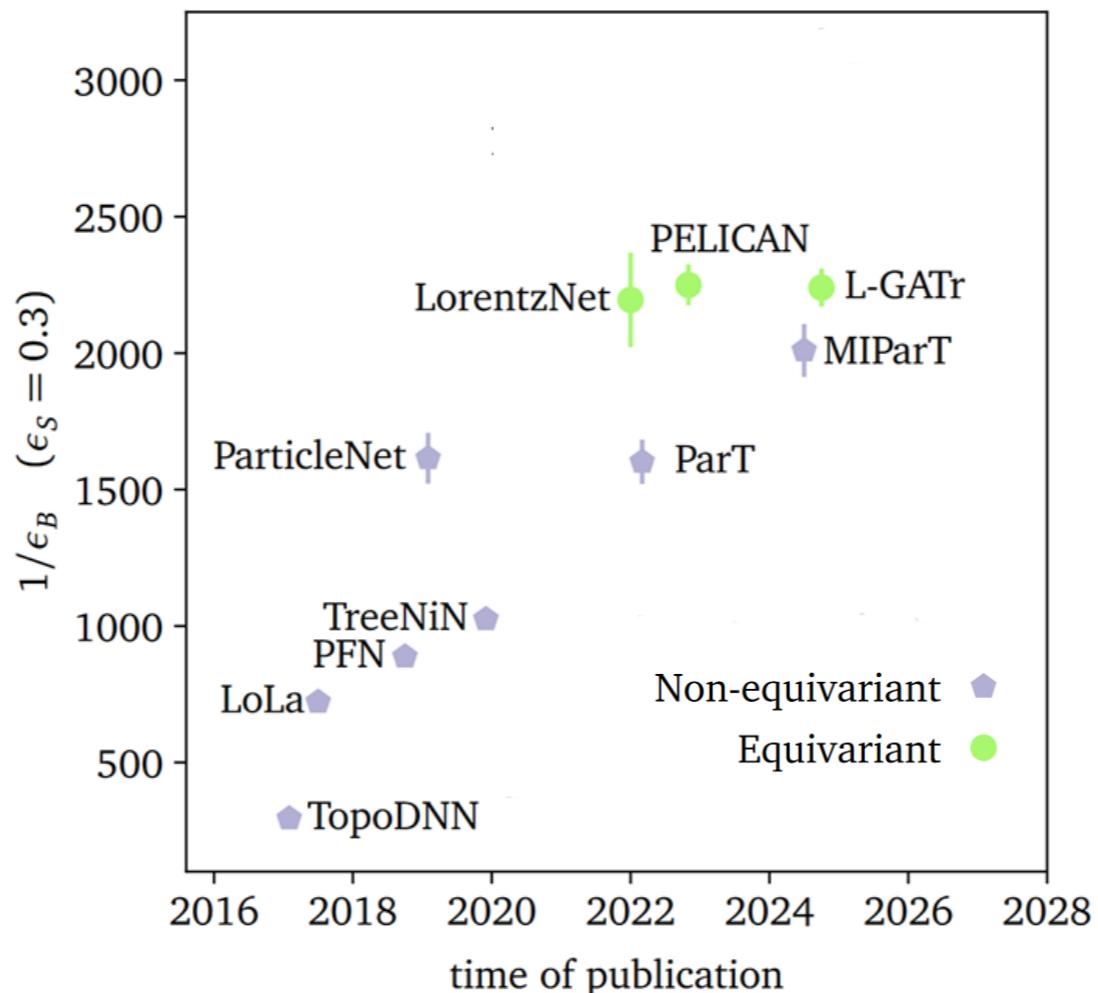
gregor.kasieczka@uni-hamburg.de  
plehn@uni-heidelberg.de

July 24, 2019

**Abstract**

Based on the established task of identifying boosted, hadronically decaying top quarks, we compare a wide range of modern machine learning approaches. Unlike most established methods they rely on low-level input, for instance calorimeter output. While their network architectures are vastly different, their performance is comparatively similar. In general, we find that these new approaches are extremely powerful and great fun.

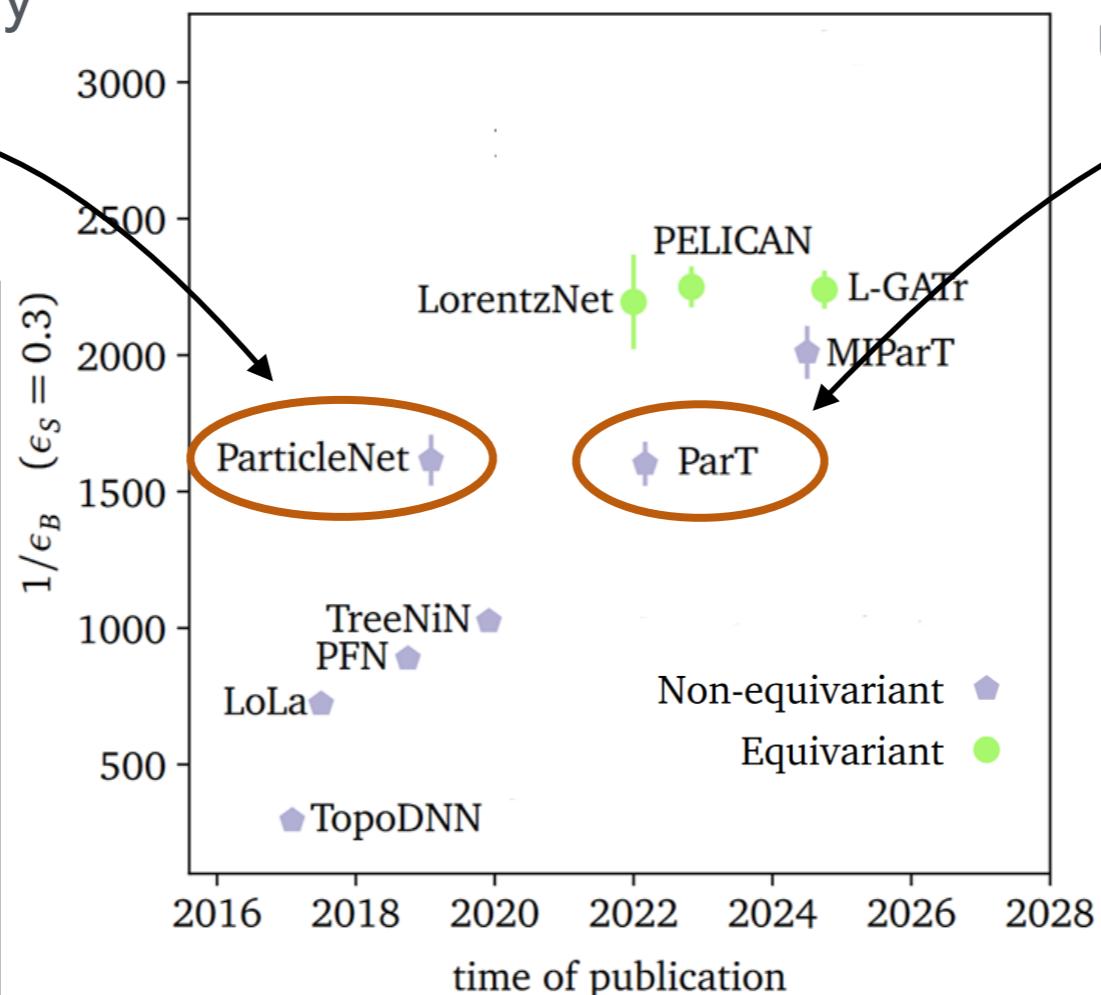
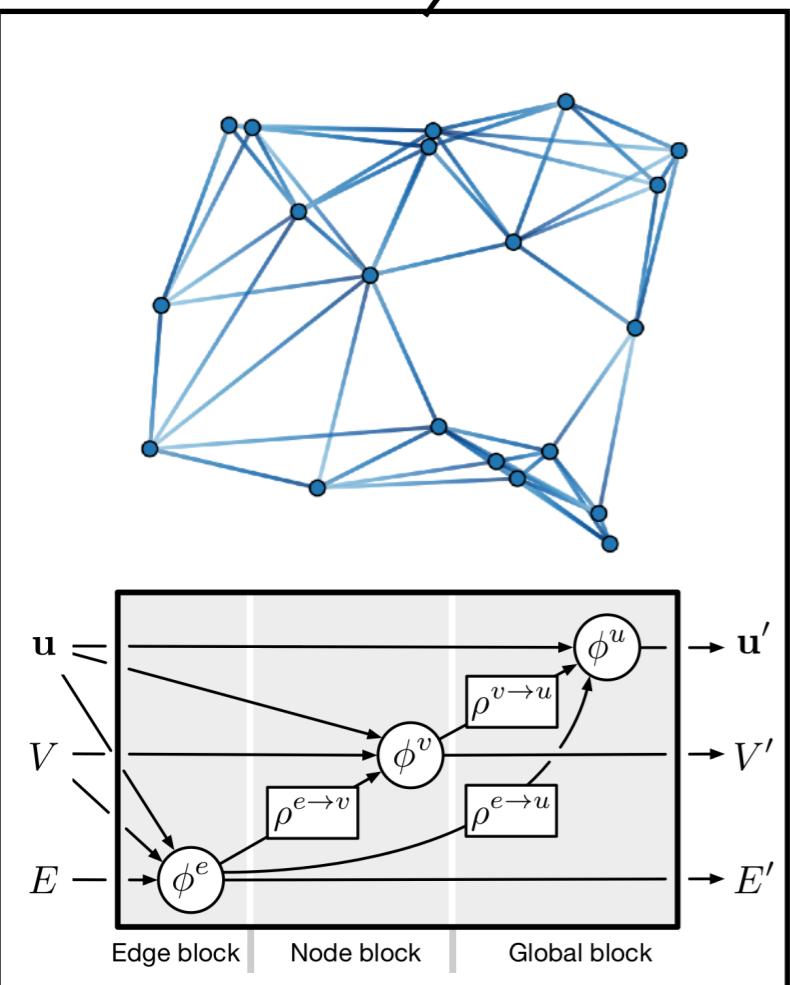
# Top Quark Tagging



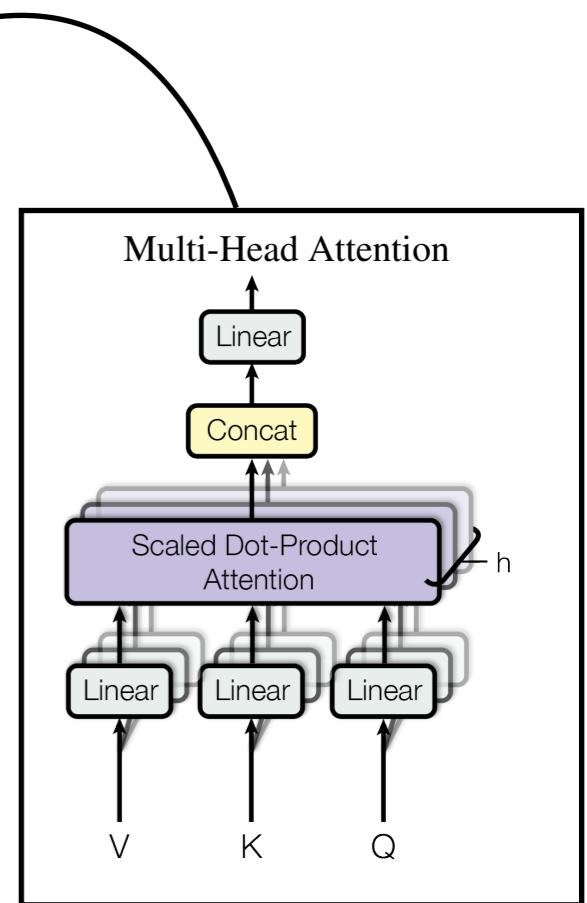
(\*Subset of methods, graph from  
V Breso at ML4Jets)

# Top Quark Tagging

View data as graph based on geometry & use message passing

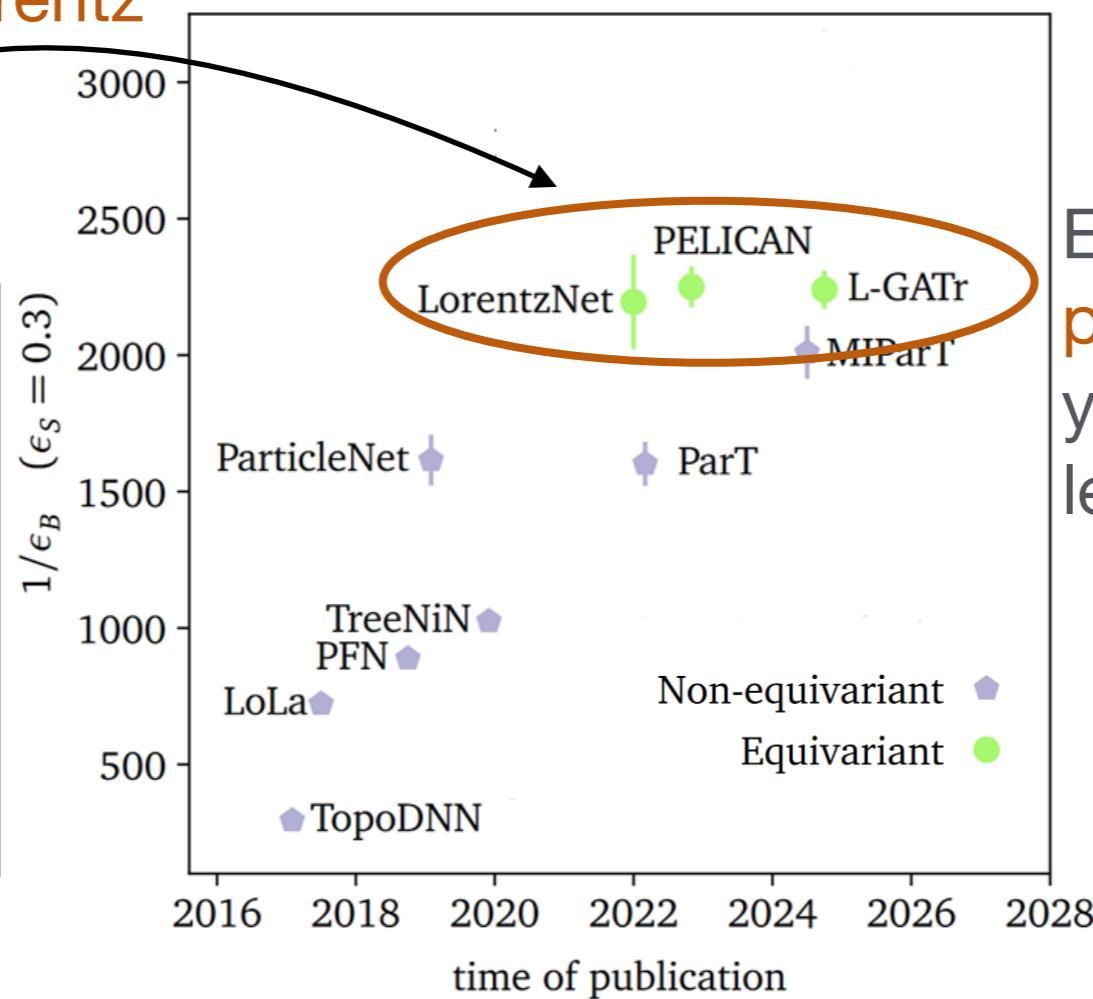
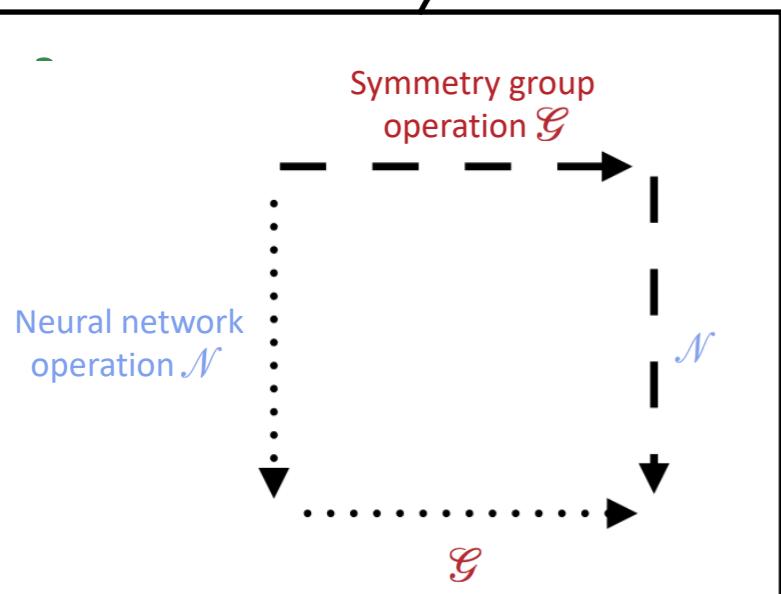


Learn relations from data using attention



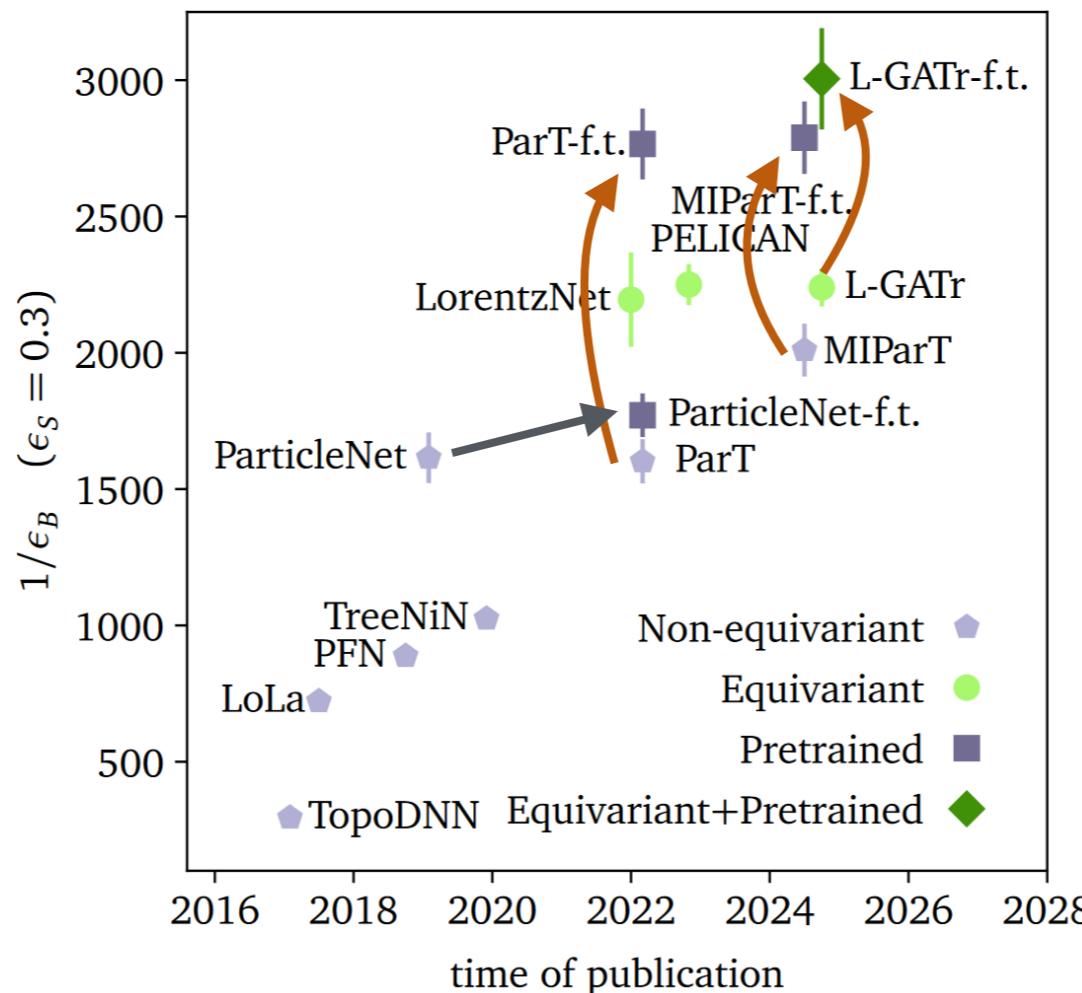
# Top Quark Tagging

Constrain functions to be equivariant under the Lorentz group



Explicitly injecting physics knowledge yields more efficient learning

# Top Quark Tagging



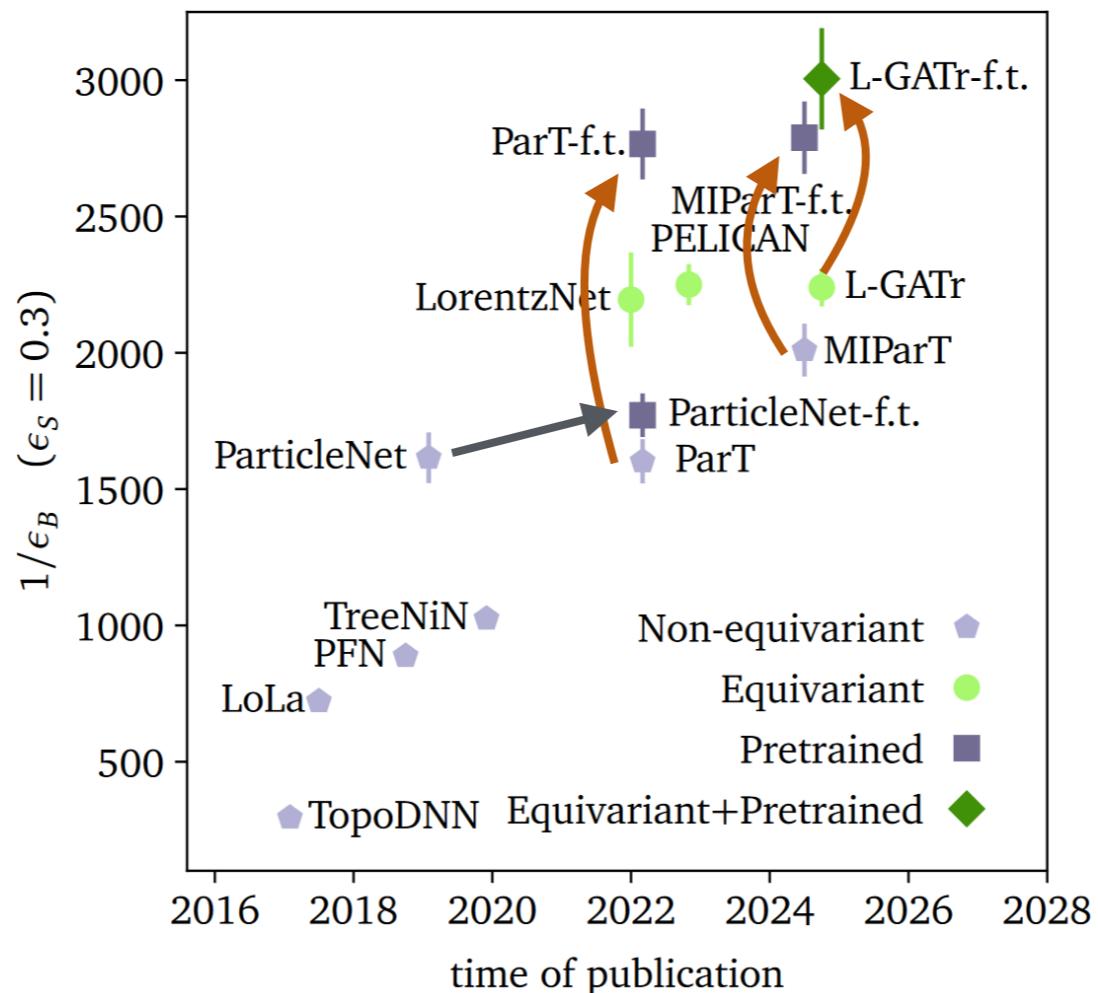
**Fine tuning:**  
Pre-train on one dataset  
and recycle weights  
for training on new data

**Yields substantial  
boost in performance  
(for **transformer**-based  
models)**

Can combine with  
equivariance

# Top Quark Tagging

- Usefulness of open benchmark data
- Good representation of data pays off
- Transformers + physics rules
- Boost by fine-tuning across datasets



Important issues on application side

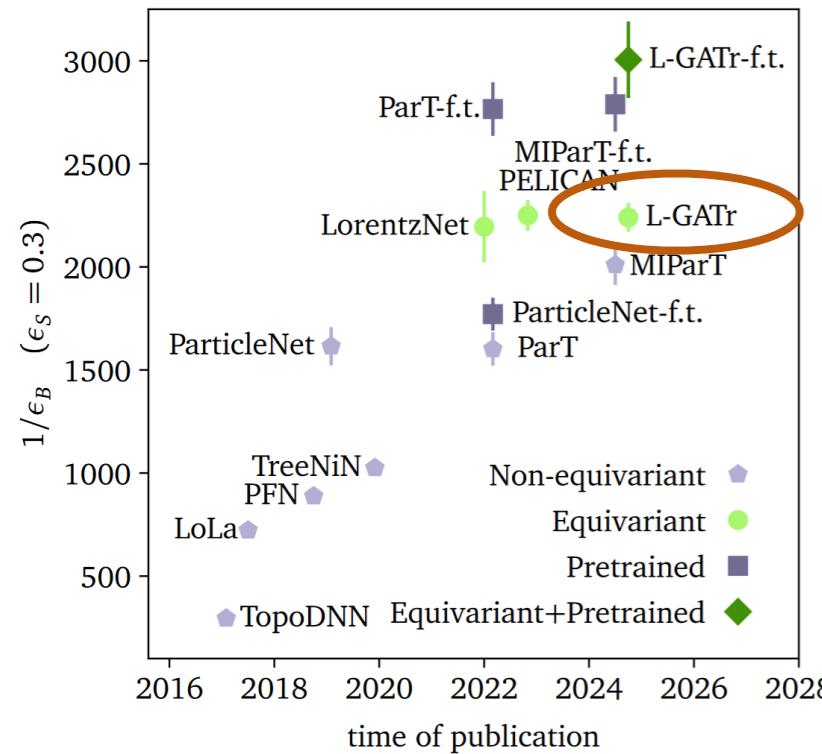
- Domain shift
- Calibrateable
- Compute cost

Modern tagging algorithms are widely used in searches for new particles

**End of Part I.**

# **Bonus Material**

# Symmetries



Reminder: Space-time described by Lorentz group (rotations, boosts)

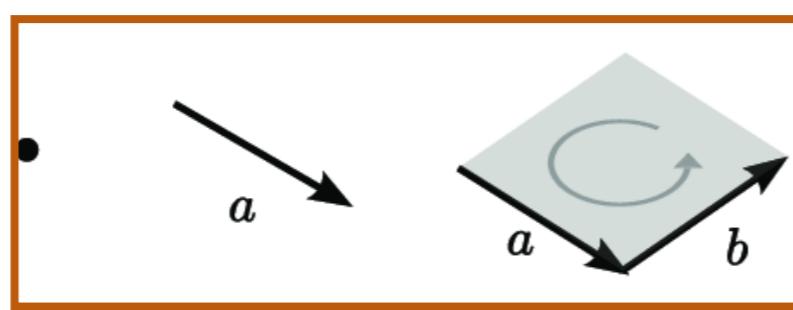
$$(t, x, y, z) \mapsto t^2 - x^2 - y^2 - z^2$$

Equivariance:  $\text{L-GATr}(\Lambda(x)) = \Lambda(\text{L-GATr}(x))$

Multivector:

$$x = x^S 1 + x_\mu^V \gamma^\mu + x_{\mu\nu}^B \sigma^{\mu\nu} + x_\mu^A \gamma^\mu \gamma^5 + x^P \gamma^5$$

“Grades” with well defined behaviour under Lorentz transformations



Scalar

Vector

Bi-Vector

Axial-Vector

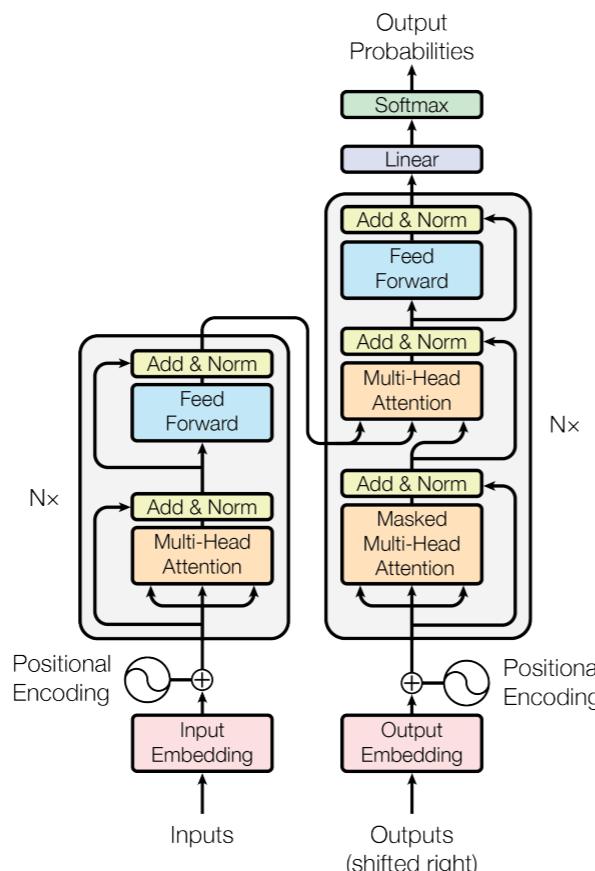
Pseudo-Scalar

$$\gamma^0 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \end{pmatrix}, \quad \gamma^1 = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & -1 & 0 & 0 \\ -1 & 0 & 0 & 0 \end{pmatrix},$$

$$\gamma^2 = \begin{pmatrix} 0 & 0 & 0 & -i \\ 0 & 0 & i & 0 \\ 0 & i & 0 & 0 \\ -i & 0 & 0 & 0 \end{pmatrix}, \quad \gamma^3 = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1 \\ -1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}.$$

# Symmetries

Layer type	Transformer	L-GATr
Linear( $x$ )	$vx + w$	$\sum_{k=0}^4 v_k \langle x \rangle_k$
Attention( $q, k, v$ ) <sub>ic</sub>	$\sum_{j=1}^{n_t} \text{Softmax}_j \left( \sum_{c'=1}^{n_c} \frac{q_{ic'} k_{jc'}}{\sqrt{n_c}} \right) v_{jc}$	$\sum_{j=1}^{n_t} \text{Softmax}_j \left( \sum_{c'=1}^{n_c} \frac{\langle q_{ic'}, k_{jc'} \rangle}{\sqrt{16n_c}} \right) v_{jc}$
LayerNorm( $x$ )	$x \left[ \frac{1}{n_c} \sum_{c=1}^{n_c} x_c^2 + \epsilon \right]^{-1/2}$	$x \left[ \frac{1}{n_c} \sum_{c=1}^{n_c} \sum_{k=0}^4 \left  \langle \langle x_c \rangle_k, \langle x_c \rangle_k \rangle \right  + \epsilon \right]^{-1/2}$
Activation( $x$ )	GELU( $x$ )	GELU( $\langle x \rangle_0$ ) $x$
GP( $x, y$ )	—	$xy$



Norm of  
inner product  
by grade

Inner product  
(requires some work but  
can compute Algebra-  
rules in advances)

$$x_\mu x^\mu = \eta_{\mu\nu} x^\mu x^\nu = (ct)^2 - \mathbf{x} \cdot \mathbf{x} \stackrel{\text{def}}{=} (c\tau)^2$$

$$\eta^{\mu\nu} = \eta_{\mu\nu} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \end{pmatrix}.$$

# **Classification**

# Cross Entropy?

$$-\frac{1}{N} \log(L) = -\frac{1}{N} \sum_k^N [y_k \cdot \log(\hat{y}(x_k)) + (1 - y_k) \cdot \log(1 - \hat{y}(x_k))]$$

- Note:

$$y_k \cdot \log(\hat{y}(x_k)) + (1 - y_k) \cdot \log(1 - \hat{y}(x_k)) = \sum_i q_i \log(p_i)$$

- Now remember entropy

$$\text{entropy} = \int p(x) \log(p(x)) dx$$

$$\text{cross-entropy} = \int q(x) \log(p(x)) dx$$

Gibbs inequality:

Cross entropy is minimal if  
 $p=q$

For multiple classes: Categorical cross entropy

# Multiple classes

- For more than 2 classes it is usually a good idea to use **One-Hot-Encoding**  
-> For  $m$  classes the output is a  $m$ -dimensional vector with components  

$$\hat{y}_c(x_k) \in [0,1]$$
- One neuron per class in output layer
- All outputs should sum to one
- Loss function: Extension of cross entropy to multiple classes: **Categorical cross entropy**

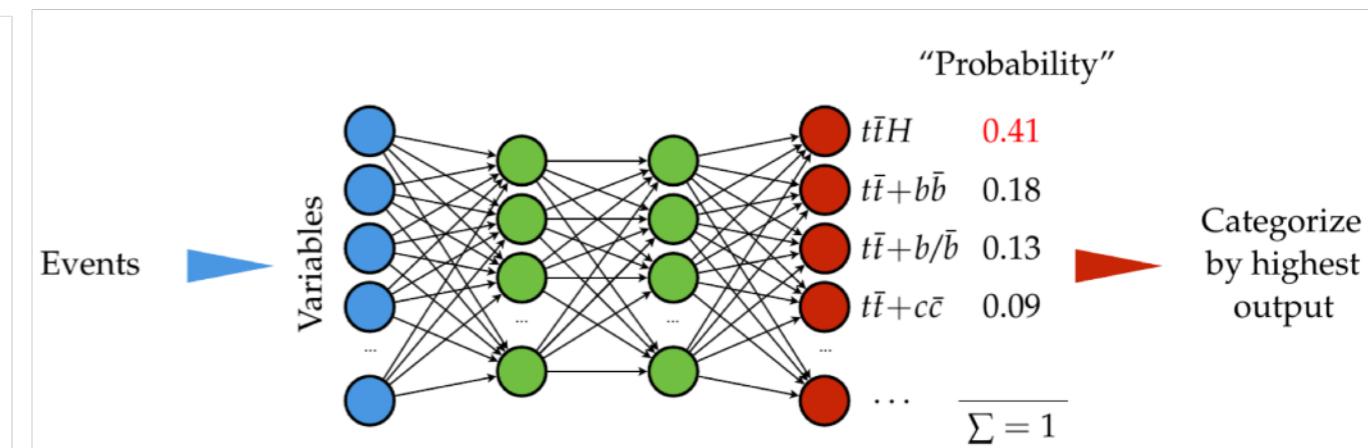
Label Encoding

Food Name	Categorical #	Calories
Apple	1	95
Chicken	2	231
Broccoli	3	50

→

One Hot Encoding

Apple	Chicken	Broccoli	Calories
1	0	0	95
0	1	0	231
0	0	1	50



<https://medium.com/@michaeldelsole/>

# Activation functions for Classification

- Need activations in the output layer that allow probability interpretation:

- $\hat{y}(x_k) \in [0,1]$
- $\sum_{\text{classes}} \hat{y}(x_k) = 1$

- Binary classification:

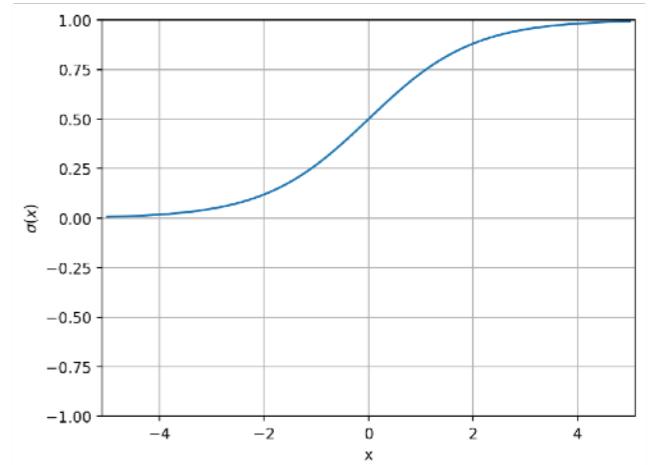
- Only 1 output node needed
- Activation: Sigmoid

- Multiclassification:

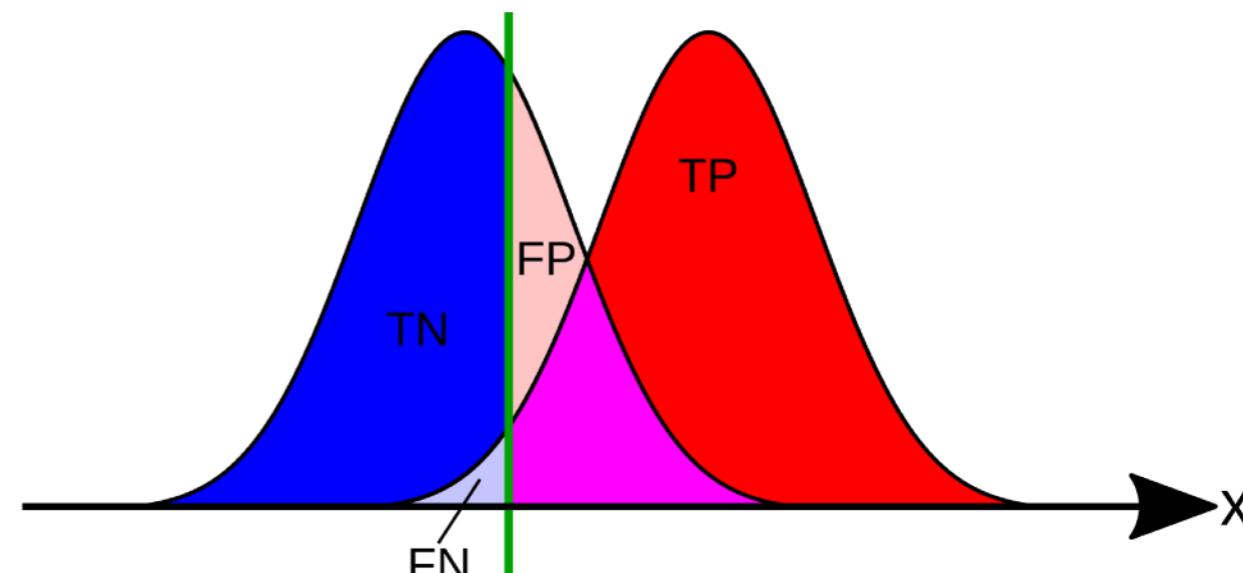
- One output node per class i
- Softmax activation:

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

$$\sigma(x) = \frac{e^{x_i}}{\sum_j^{\text{classes}} e^{x_j}}$$

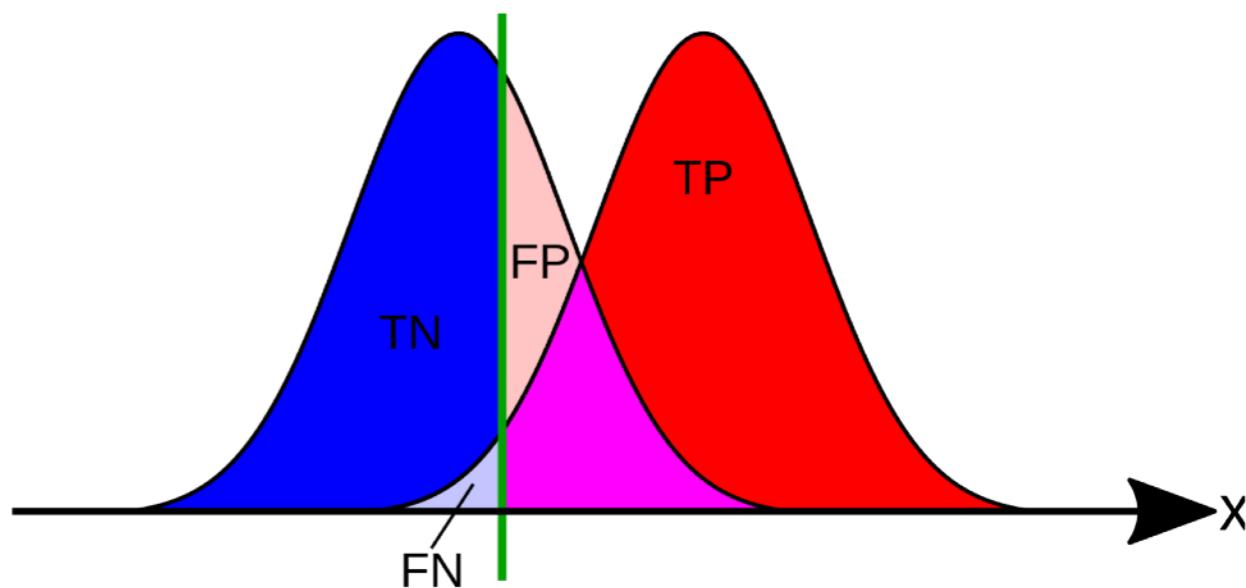


# Binary Classification Outputs

- The NN outputs a score for belonging to class „signal“ (high value) or class „background“ (low value)
  - Often you need to make a decision based on this score -> Choose a threshold
  - Key quantities:
    - True Positives (**TP**)
    - False Positives (**FP**)
    - True Negatives (**TN**)
    - False Negatives (**FN**)
- 
- Purity/Precision =  $\text{TP}/(\text{TP}+\text{FP})$ : How pure is the sample predicted as signal?
  - True Positive Rate/Sensitivity/Efficiency/Recall =  $\text{TP}/(\text{TP}+\text{FN})$ : What fraction of the signal is classified correctly?
  - False Positive Rate =  $\text{FP}/(\text{FP}+\text{TN})$
  - Accuracy =  $(\text{TP}+\text{TN})/(\text{TP}+\text{FP}+\text{TN}+\text{FN})$ : What fraction of all data is classified correctly?

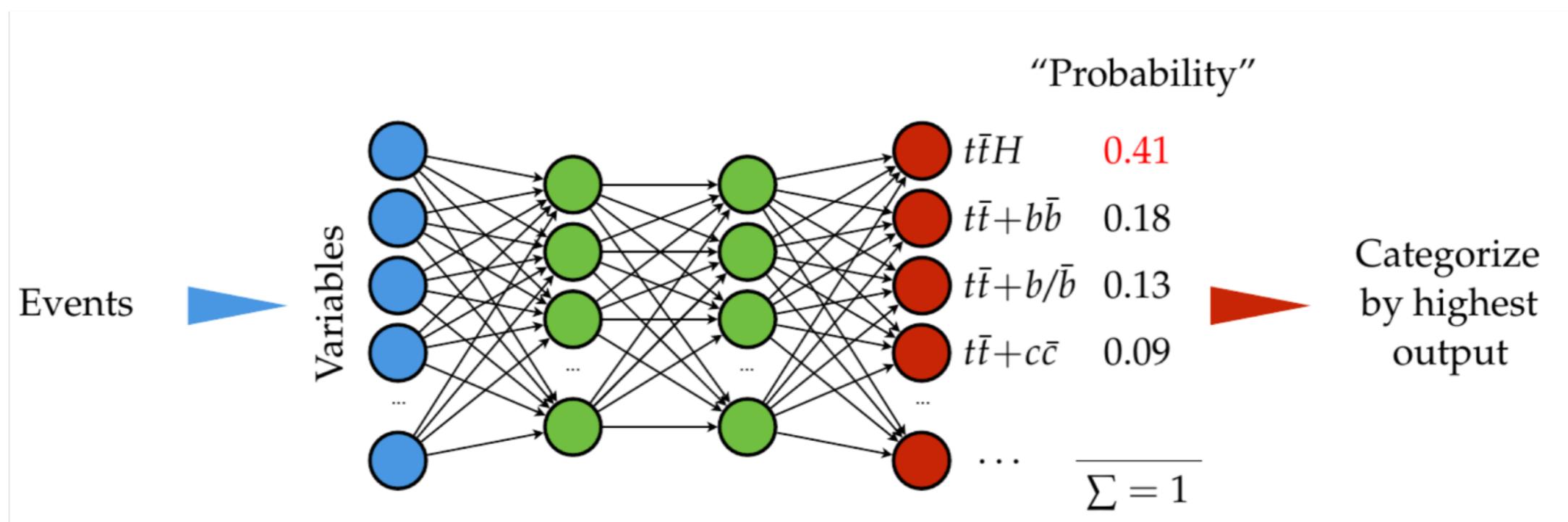
# Choosing a threshold

- Which threshold to choose?
- Depends on your problem:
  - Think about what questions you want to answer
  - The cost of being wrong
- Do you need a high purity? (“We need to be sure that all animals we tag as harmless are actually harmless!”)
- Do you need a high Efficiency? („We need to detect as many positive Covid cases as possible“)



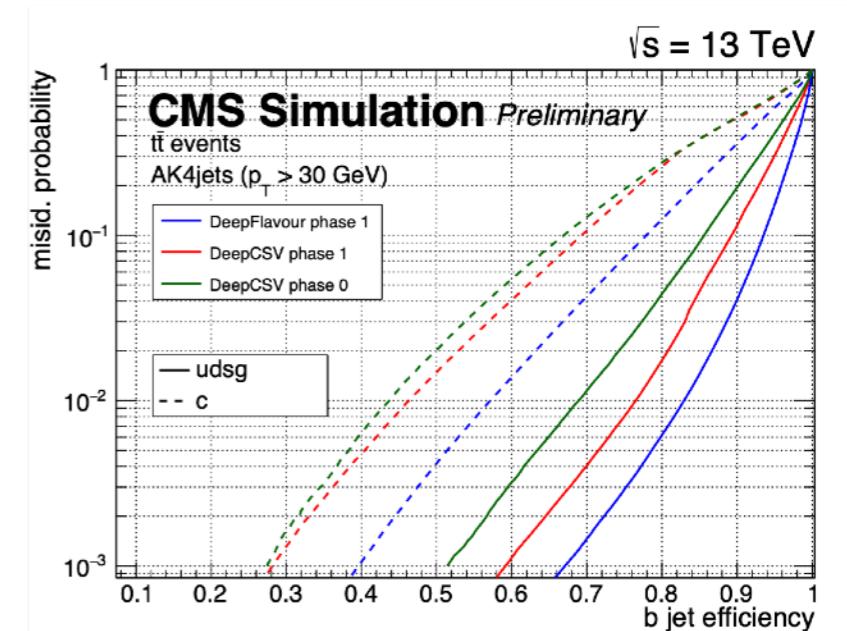
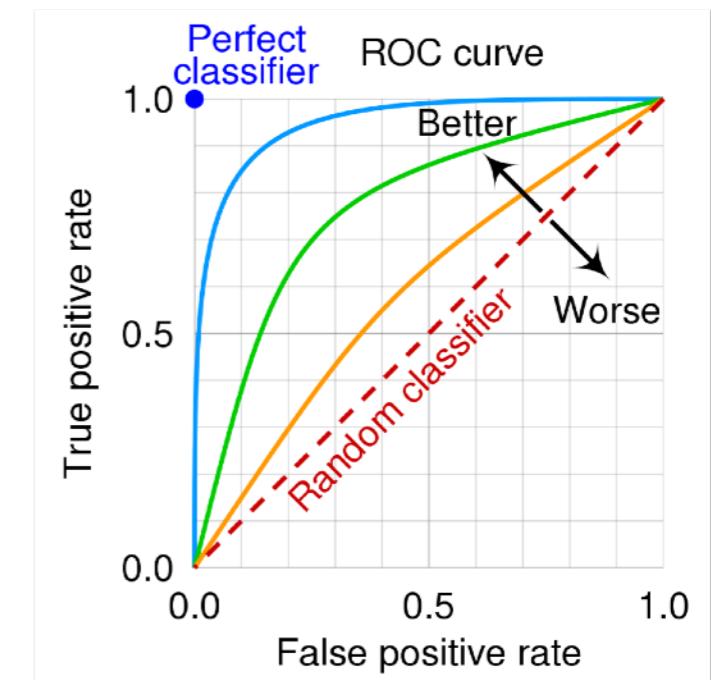
# Multiclassification

- In multiclassification we have one output neuron for each class
- If we use softmax as final activation function, the values of each output neuron sum to 1 and can be interpreted as probability
- Common approach: Classify a sample according to the neuron with the highest output value



# Receiver Operator Characteristic (ROC)

- Overall performance of a classifier is usually visualized using the ROC curve:
  - Scan all possible thresholds and evaluate True Positive Rate and False Positive Rate
  - Plot one as a function of the other
- Area under the curve (AUC) often quoted as performance metric
- ROCs can aide in choosing a classifier and threshold
- Note: Sometimes ROCs of different classifiers cross each other -> they are better at different TPR regions



# Confusion matrices

- With a chosen threshold we can plot confusion (migration) matrices
- Shows migration from one class to another
- Very usefull for multiclassification:
  - Which background class looks most signal like?
  - Which classes are difficult to separate?
- Can be normalized to:
  - True classes (efficiency matrix)
  - Predicted class (purity matrix)

		Actual Values	
		Positive	Negative
Predicted Values	Positive	True Positive	False Positive
	Negative	False Negative	True Negative

		CMS Simulation (13 TeV)				
		Reconstructed decay mode				
		$h^{\pm}h^{\pm}\pi^0$	$h^{\pm}h^{\mp}$	$h^{\pm}h^{\pm}(\pi^0s)$	$h^{\pm}\pi^0s$	$h^{\pm}$
None		0.11	0.25	0.10	0.17	0.38
$h^{\pm}h^{\pm}\pi^0$		0.00	0.01	0.05	0.36	0.11
$h^{\pm}h^{\mp}$		0.00	0.01	0.61	0.27	0.07
$h^{\pm}h^{\pm}(\pi^0s)$		0.00	0.02	0.19	0.13	0.03
$h^{\pm}\pi^0s$		0.09	0.57	0.02	0.06	0.36
$h^{\pm}$		0.80	0.14	0.03	0.01	0.04
		$h^{\pm}$	$h^{\pm}\pi^0s$	$h^{\pm}h^{\pm}\pi^0$	$h^{\pm}h^{\pm}h^{\mp}\pi^0s$	$Other$
		Generated decay mode				