

Machine Learning

Lesson 2

School on Data Science in Fundamental Physics, IGFAE/USC, Spain

Dr. Pietro Vischia

pietro.vischia@cern.ch

[@pietrovischia](https://twitter.com/pietrovischia)

<https://vischia.github.io/>



**Supported by project
RYC2021- 033305-I
funded by**



MINISTERIO
DE CIENCIA
E INNOVACIÓN

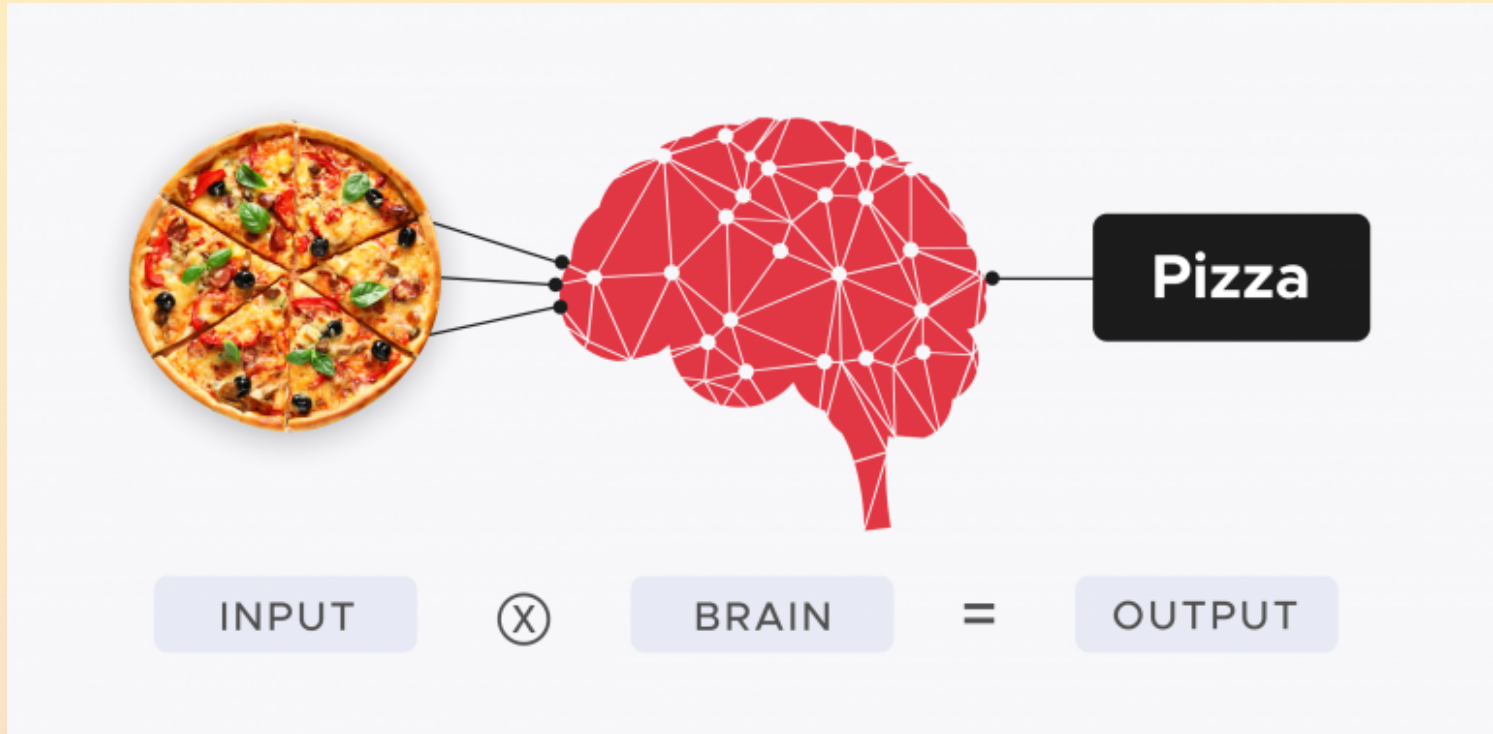


If you are reading this as a web page: have fun! If you are reading this as a PDF:
please visit

https://www.hep.uniovi.es/vischia/persistent/2024-06-03to07_MachineLearningAtDataScienceSchoolIGFAE_vischia_2.html

to get the version with working animations

Brain activity...



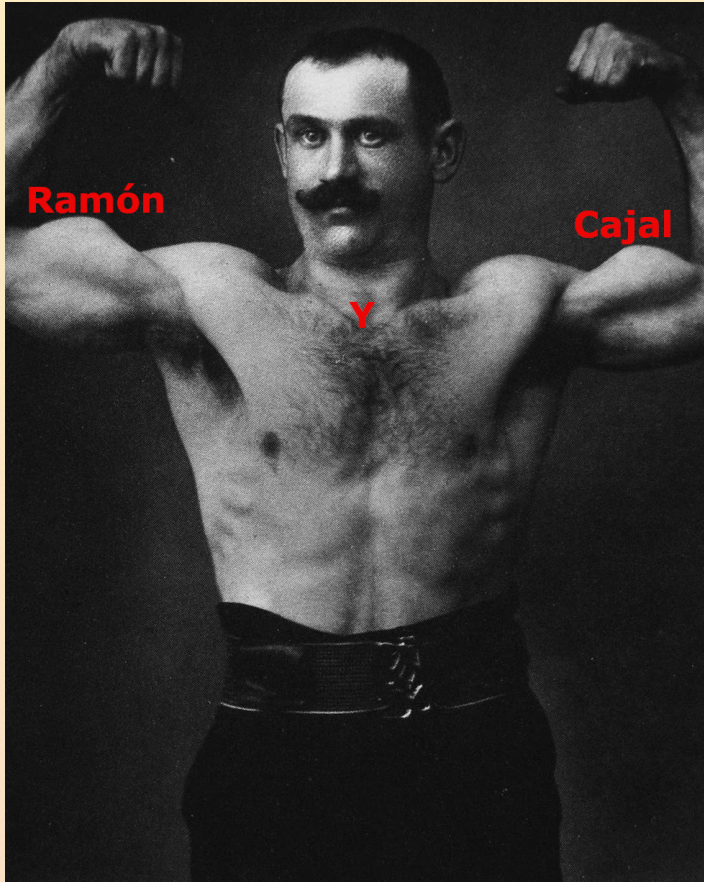
...approximated...



...using computers

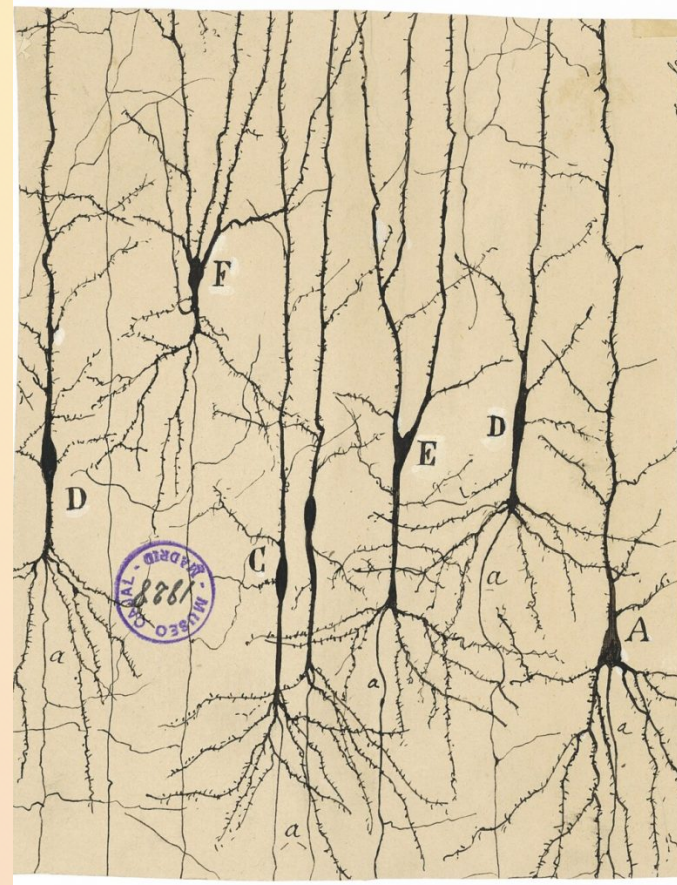
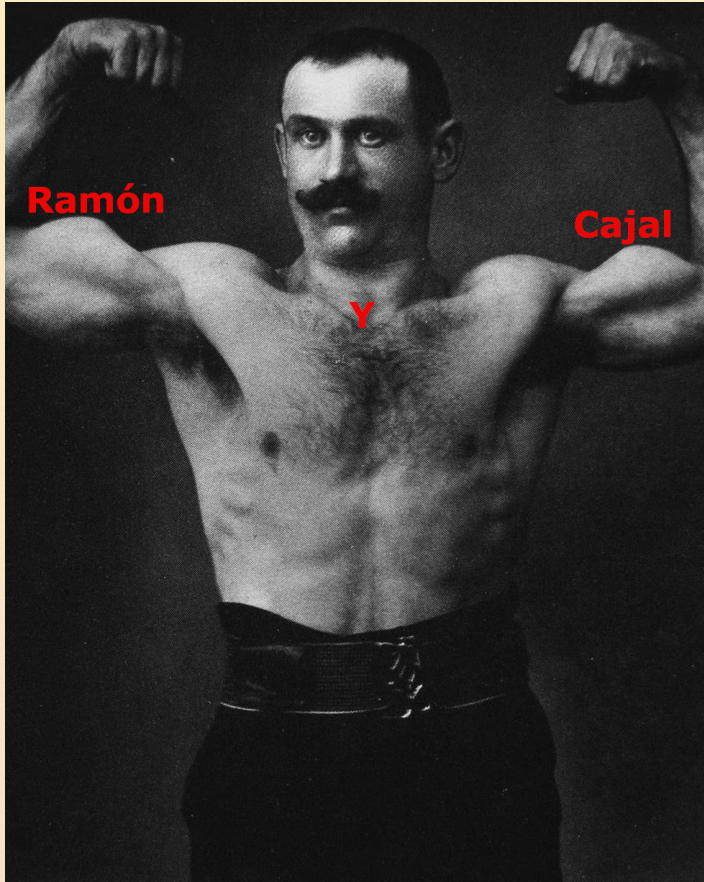


Santiago Ramón y Cajal



- "The Spanish father of culturism"

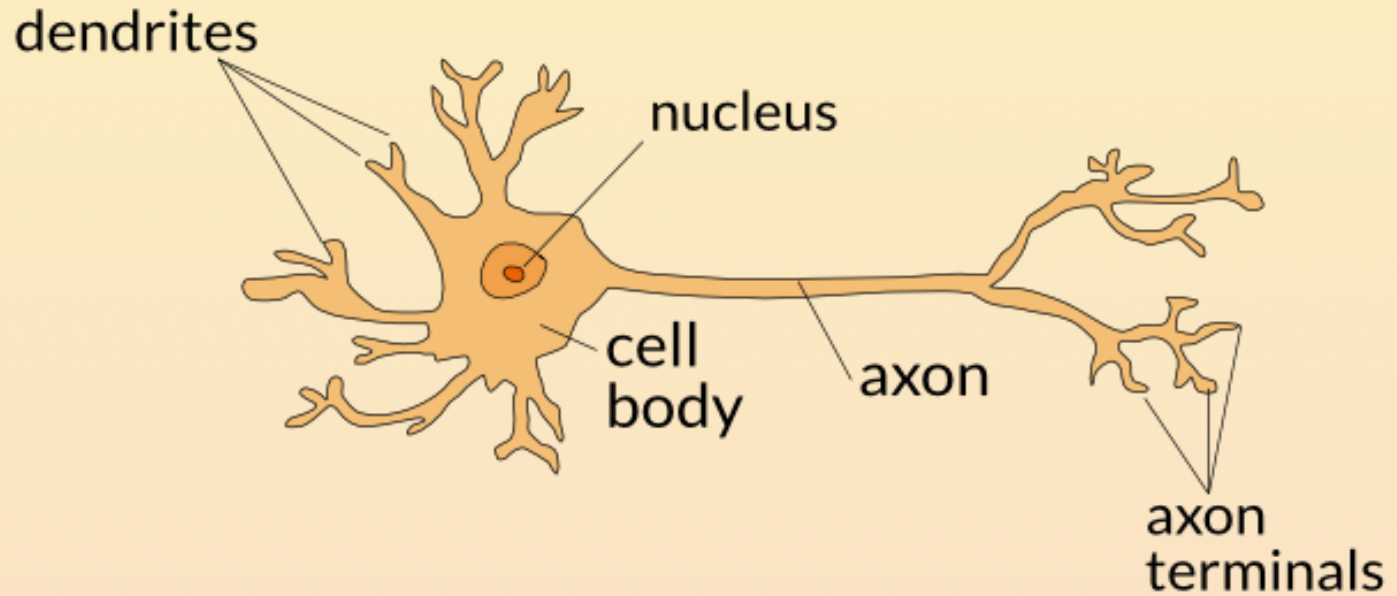
Santiago Ramón y Cajal



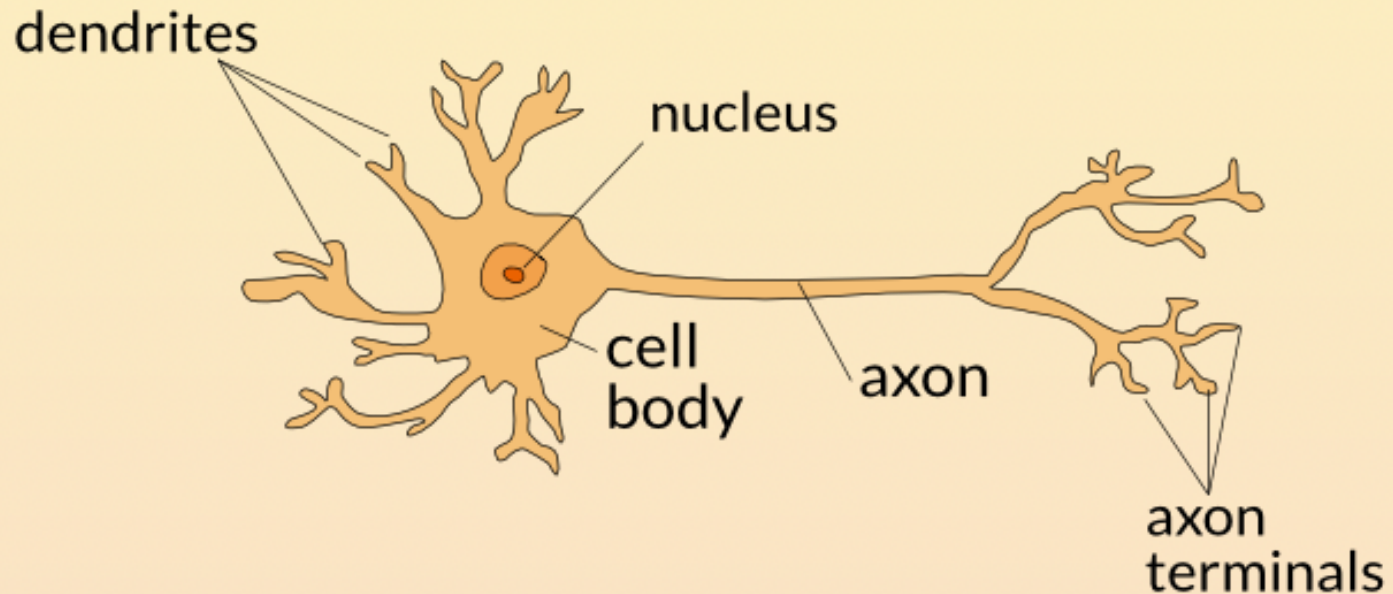
- "The Spanish father of culturism"

- The 1906 Nobel Prize in Medicine

Real neurons



Real neurons



$$I = C \frac{dV}{dt} + G_{Na} m^3 h (V - V_{Na}) + G_K n^4 (V - V_K) + G_L (V - V_L)$$

Computationally heavy



Simplified Neurons

Bulletin of Mathematical Biology Vol. 52, No. 1/2, pp. 99–115, 1990.
Printed in Great Britain.

0092-8240/90\$3.00 + 0.00
Pergamon Press plc
Society for Mathematical Biology

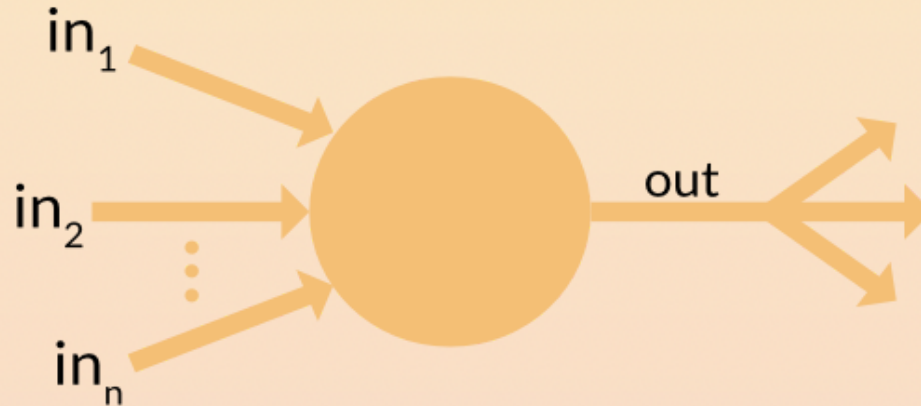
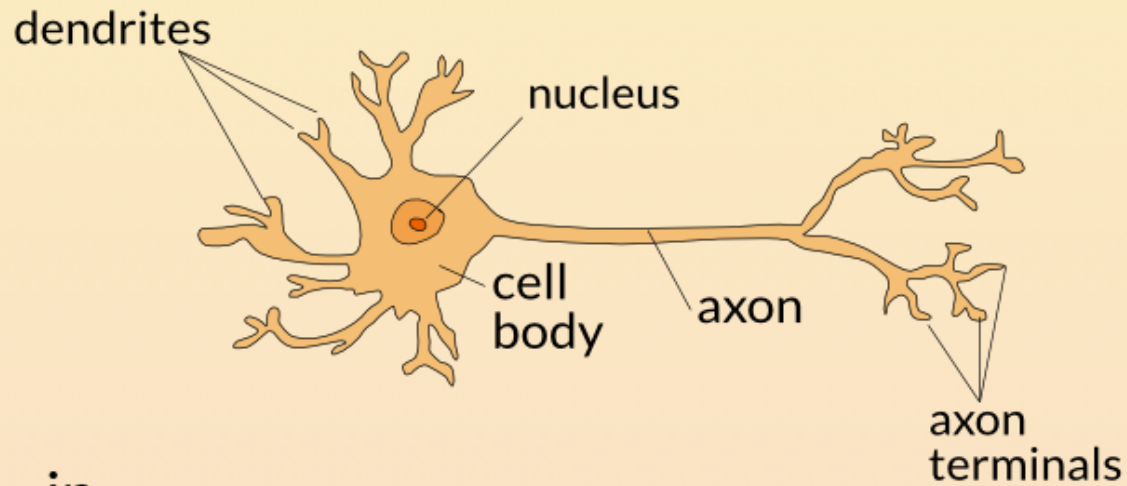
A LOGICAL CALCULUS OF THE IDEAS IMMANENT IN NERVOUS ACTIVITY*

- WARREN S. MCCULLOCH AND WALTER PITTS
University of Illinois, College of Medicine,
Department of Psychiatry at the Illinois Neuropsychiatric Institute,
University of Chicago, Chicago, U.S.A.

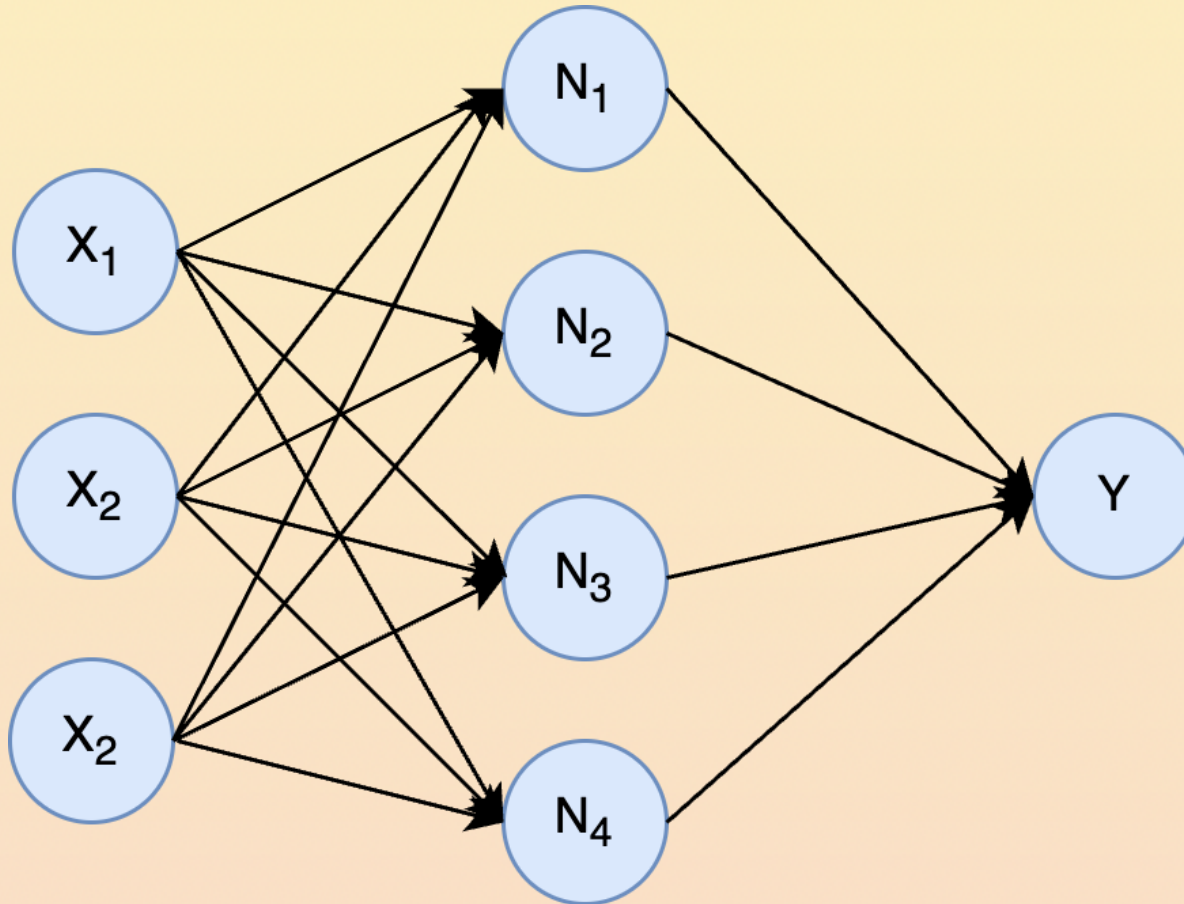
Because of the “all-or-none” character of nervous activity, neural events and the relations among them can be treated by means of propositional logic. It is found that the behavior of every net can be described in these terms, with the addition of more complicated logical means for nets containing circles; and that for any logical expression satisfying certain conditions, one can find a net behaving in the fashion it describes. It is shown that many particular choices among possible neurophysiological assumptions are equivalent, in the sense that for every net behaving under one assumption, there exists another net which behaves under the other and gives the same results, although perhaps not in the same time. Various applications of the calculus are discussed.

Perceptrons

$$y = f\left(b_i + \sum w_i x_i\right)$$



Artificial Neural Networks



Perceptron step by step

- Linear combination of the inputs

$$\sum_j w_j x_j$$

- Activation function (imitates the activation of real neurons, where a voltage peak, a **spike**, is generated when the injected potential passes a threshold)

$$g\left(w_0 + \sum_j w_j x_j\right)$$

- Bias term, an order-zero term in the inputs (translation)

$$\hat{y} = g\left(w_0 + \sum_j w_j x_j\right)$$

- In matrix form, a row-column product

$$\hat{y} = g\left(w_0 + \mathbf{X}^T \mathbf{W}\right)$$

Activation Function

- Originally, the activation function was linear

$$g(z) = 1 \text{ if } z = w_0 + \sum w_i x_i \geq 0$$

$$g(z) = -1 \text{ if } z = w_0 + \sum w_i x_i < 0$$

- Activation function is the only chance of estimating nonlinear functions
- Otherwise, a neural network would be just a fancier linear regression model

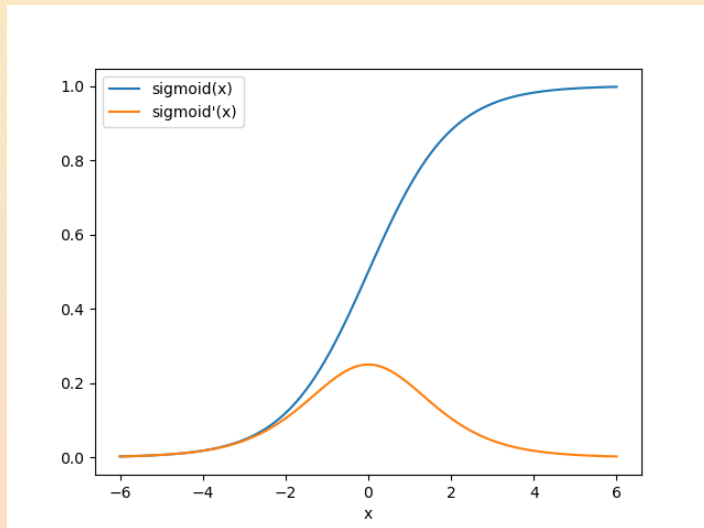
$$\hat{y} = g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ 2 \end{bmatrix}\right) = g(1 + 3x_1 + 2x_2)$$

Popular activation functions

Sigmoid

$$g(z) = \frac{1}{1 + e^{-z}}$$

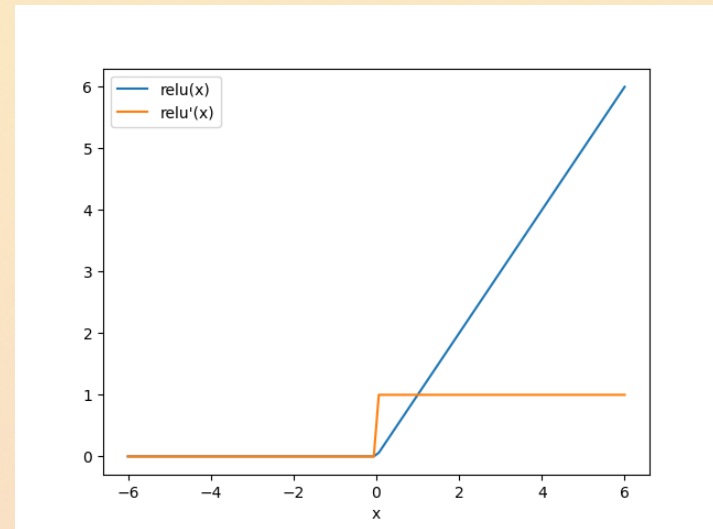
$$g'(z) = g(z)(1 - g(z))$$



Rectified Linear Unit (ReLU)

$$g(z) = \max(0, z)$$

$$g'(z) = 1 \text{ if } z > 0, 0 \text{ otherwise}$$



Multidimensional outputs

- \hat{y}_1, \hat{y}_2 , each one with the same formula as a single neuron \rightarrow just with an additional index

$$\hat{y}_i = g(w_{0,i} + \sum_j w_{j,i} x_j)$$

Neural network with one internal layer

- Between input and hidden layer: $\mathbf{W}^{(1)}$
- Between hidden layer and output layer: $\mathbf{W}^{(2)}$
- Output of the hidden layer neurons:

$$z_i = g(w_{0,i}^{(1)} + \sum_j w_{j,i}^{(1)} x_j)$$

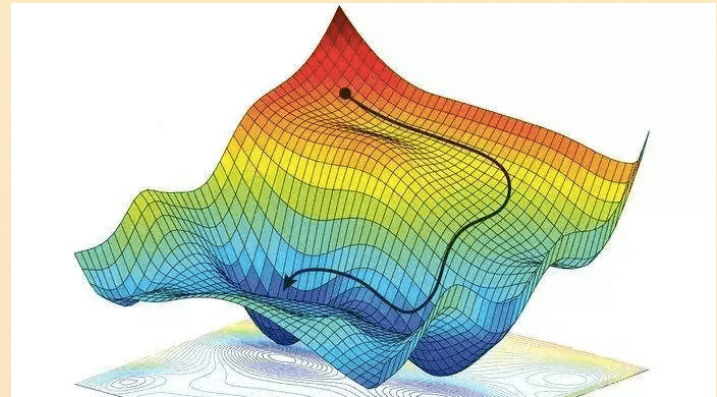
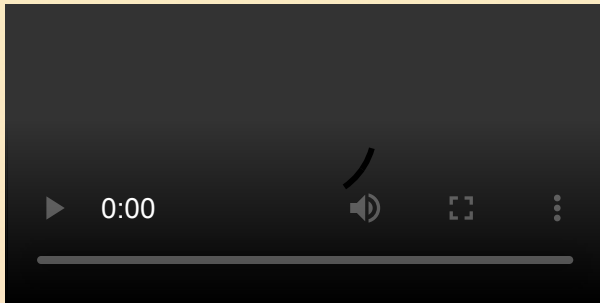
- Output of the network

$$\hat{y}_i = g(w_{0,i}^{(2)} + \sum_j w_{j,i}^{(2)} z_j)$$

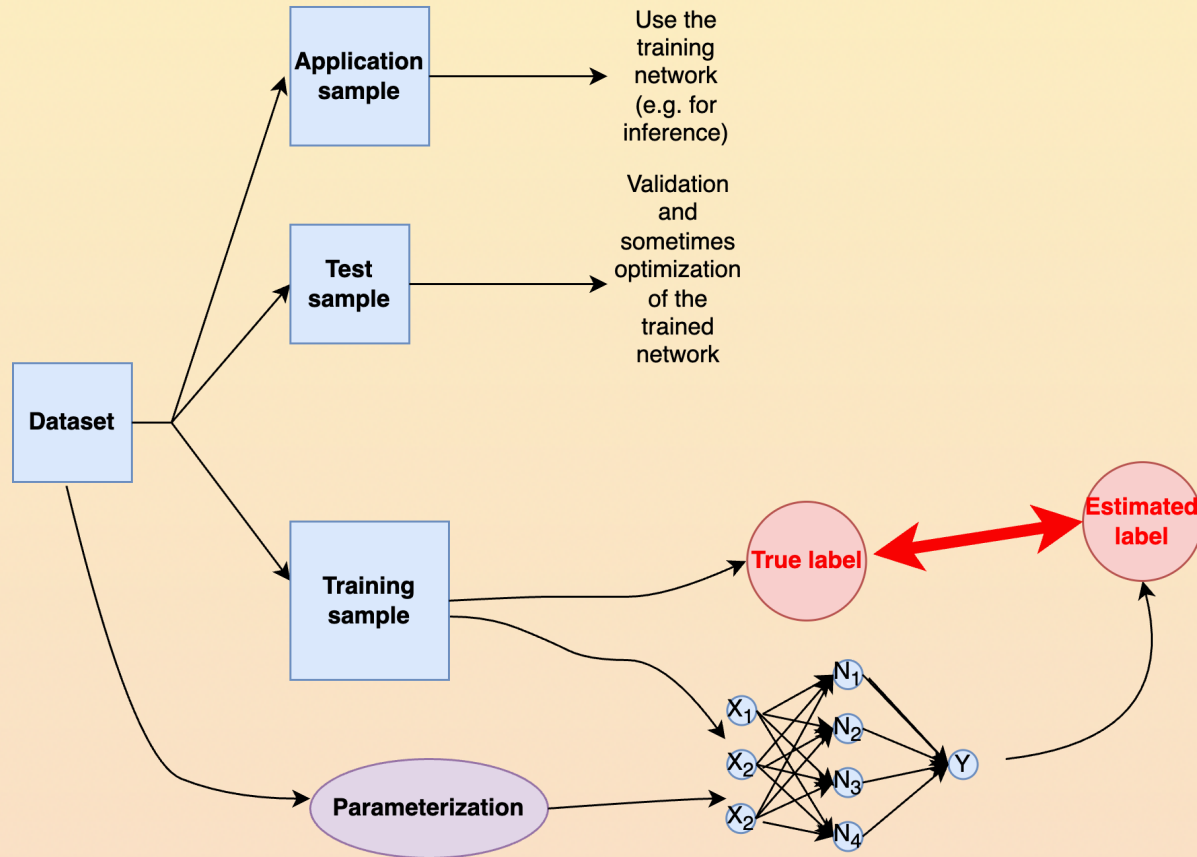
- The generalization to multiple outputs \hat{y}_i is also trivial

Gradient Descent

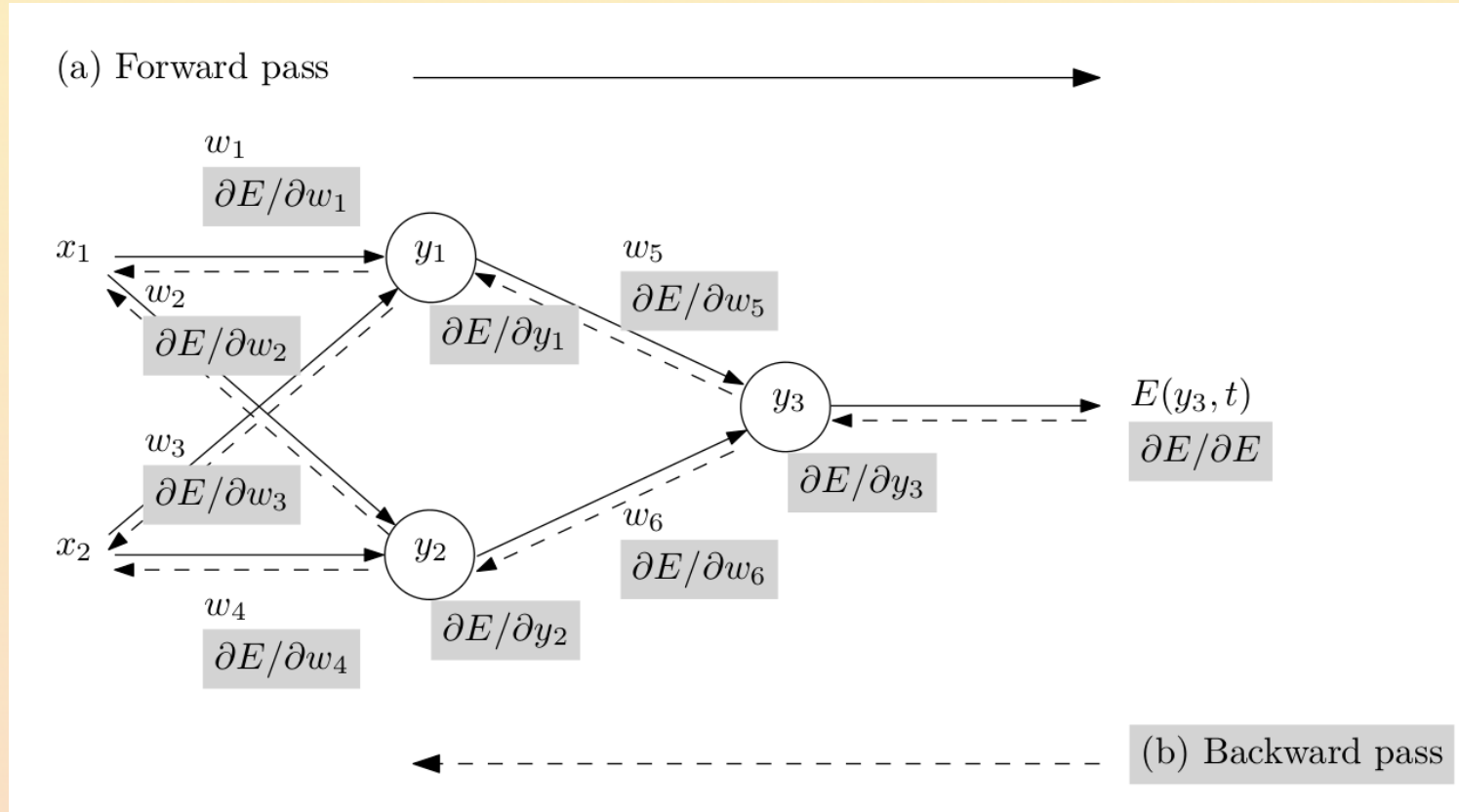
- Search for the "minimum error" as a function of the values of the training paramters (weights)



Learning



Backpropagation



Backpropagation

- Empirical Loss Function
- Cost function
- Empirical risk

$$\mathbf{J}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{*(i)})$$

- Minimized by:

$$\mathbf{W}^0 = \operatorname{argmin}_{\mathcal{W}} \mathbf{J}(\mathbf{W}) = \operatorname{argmin}_{\mathcal{W}} \mathbf{J}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{*(i)})$$

Loss function comes from inference

- Decision-theoretic approach (C.P. Robert, "The Bayesian Choice")
 - \mathcal{X} : observation space
 - Θ : parameter space
 - \mathcal{D} : decision (action) space
- **Statistical inference** take a decision $d \in \mathcal{D}$ related to parameter $\theta \in \Theta$ based on observation $x \in \mathcal{X}$, under $f(x|\theta)$
 - Typically, d consists in estimating $h(\theta)$ accurately

$$U(\theta, d) = \mathbb{E}_{\theta, d} \left[U(r) \right]$$

Loss function comes from inference

- Loss function: $L(\theta, d) = -U(\theta, d)$
 - Represents intuitively the loss or error in which you incur when you make a bad decision (a bad estimation of the target function)
 - Lower bound at 0: avoids "infinite utility" paradoxes (St. Petersburg paradox, martingale-based strategies)
- Generally impossible to uniformly minimize in d the loss for θ unknown
 - Need for a practical prescription to use the loss function as a comparison criterion in practice

Frequentist loss, Bayesian loss

- Frequentist loss (risk) is integrated (averaged) on \mathcal{X} : $R(\theta, \delta) = \mathbb{E}_{\theta} \left[L(\theta, \delta(x)) \right]$
 - $\delta(\cdot)$ is an `\textbf{estimator}` of θ (e.g. MLE)
 - Compare estimators, find the best estimator based on long-run performance for all values of unknown θ
 - Issues: based on long run performance (not optimal for x_{obs}); repeatability of the experiment; no total ordering on the set of estimators
- Bayesian loss: is integrated on Θ : $\rho(\pi, d|x) = \mathbb{E}^{\pi} \left[L(\theta, d)|x \right]$
 - π is the prior distribution
 - Posterior expected loss averages the error over the posterior distribution of θ conditional on x_{obs}
 - Can use the conditionality because x_{obs} is known!
 - Can also integrate the frequentist risk; integrated risk $r(\pi, \delta) = \mathbb{E}^{\pi} \left[R(\theta, \delta) \right]$ averaged over θ according to π (total ordering)

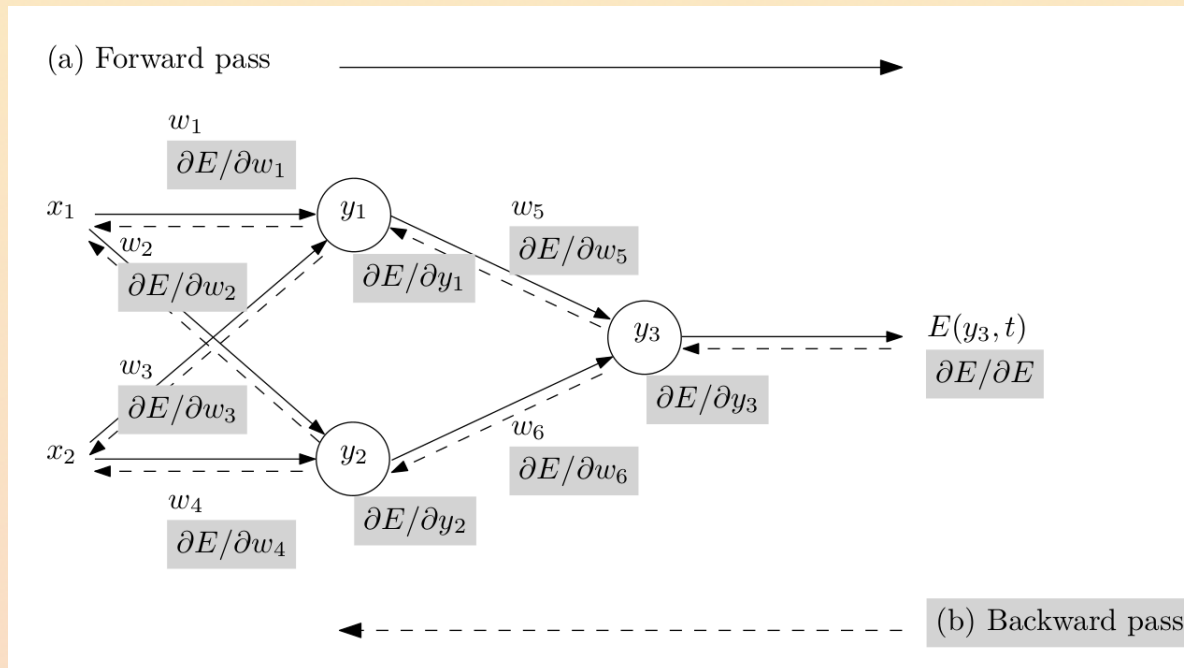
Training loop

- Initialize the weights (for instance $w \sim \text{Gaus}(0, \sigma^2)$)
- Loop until convergence
 - Compute network output (**forward step**)
 - Compute gradients $\frac{\partial \mathbf{J}(\mathbf{W})}{\partial \mathbf{W}}$
 - Update the weights according to the learning rate $\mathbf{W} \leftarrow \mathbf{W} + \eta \frac{\partial \mathbf{J}(\mathbf{W})}{\partial \mathbf{W}}$
 - Return weights

Backpropagation

$$\mathbf{J}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{*(i)}), \quad \mathbf{W}^0 = \operatorname{argmin}_{\mathbf{W}} \mathbf{J}(\mathbf{W})$$

$$\mathbf{W} \leftarrow \mathbf{W} + \eta \frac{\partial \mathbf{J}(\mathbf{W})}{\partial \mathbf{W}}$$



Jacobian and Hessian

$$\mathbf{J}_{ij} = \frac{\partial f_i}{\partial x_j}$$

$$\mathbf{J}(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}) & \dots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}) \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}) & \dots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}) \end{bmatrix}$$

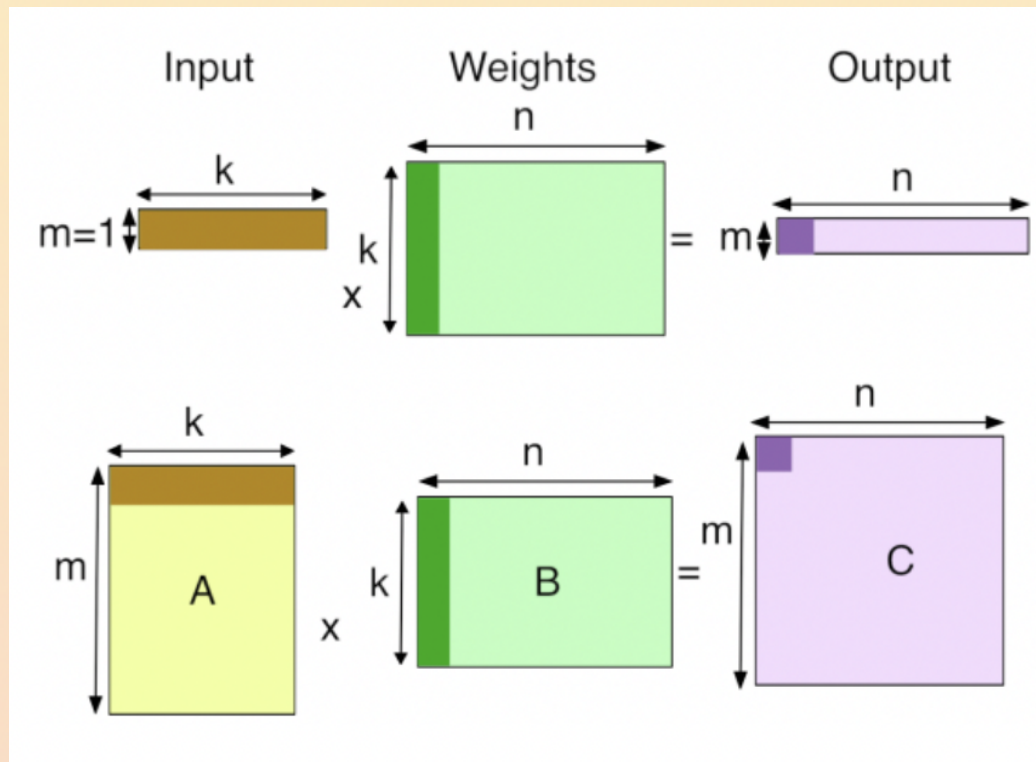
$$\mathbf{H}_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

$$\mathbf{H}(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{x}) & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{x}) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_2^2}(\mathbf{x}) & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(\mathbf{x}) \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_n \partial x_2}(\mathbf{x}) & \dots & \frac{\partial^2 f}{\partial x_n^2}(\mathbf{x}) \end{bmatrix}$$

- $\mathbf{H}(f(\mathbf{x})) = \mathbf{J}(\nabla f(\mathbf{x}))$
(describes local curvature)

Matrix multiplication

- Neural network weights expressible as matrices
- Generalize matrix calculus to tensors ([tensorflow](#))
- Optimize for efficient tensor calculus (e.g. GPU→TPU, computational tricks)



Example: Google's TPUs

- Systolic flow
 - Hide four-stage process within the matrix multiplication operation
 - E.g. decoupled access/execution when reading weights
 - Trick flow control into thinking inputs are read and update results at once

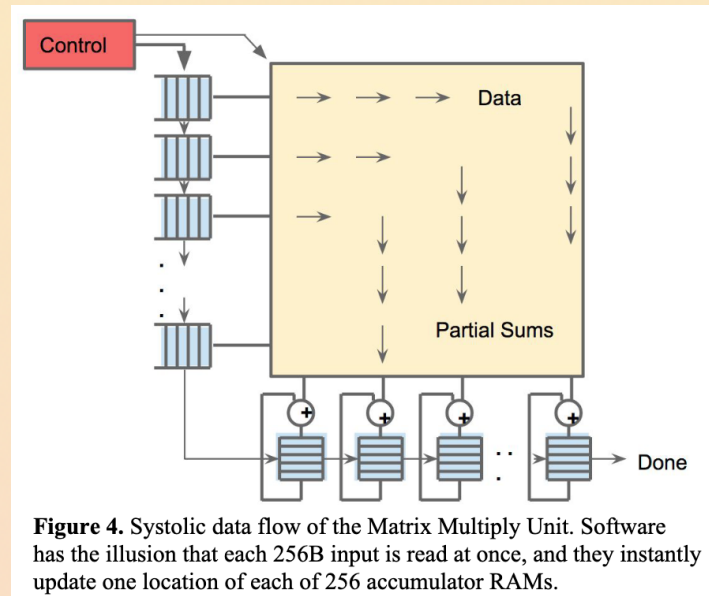
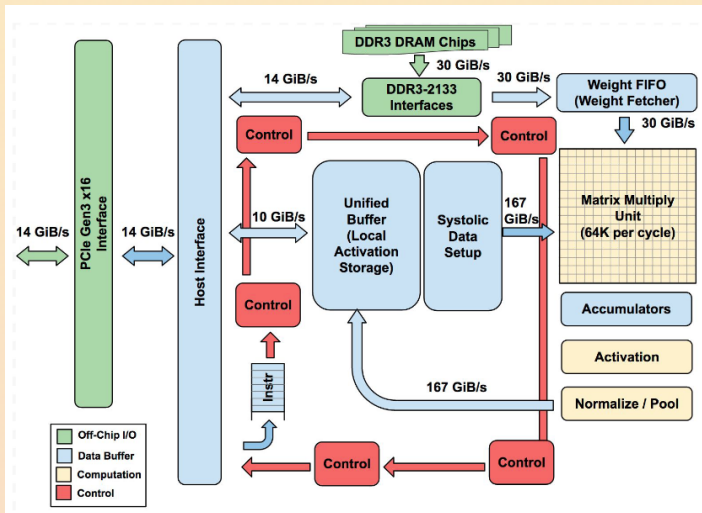
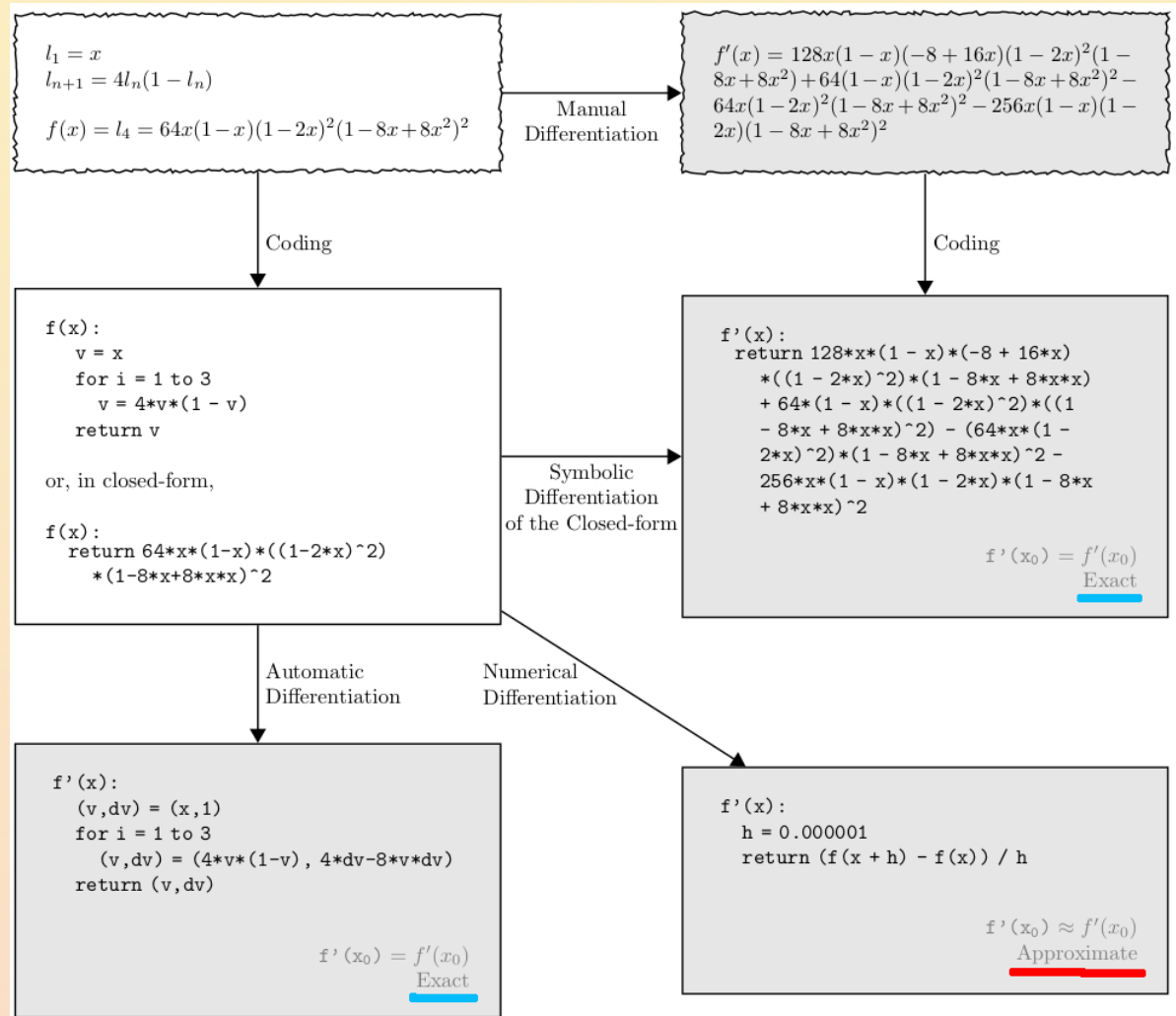


Figure 4. Systolic data flow of the Matrix Multiply Unit. Software has the illusion that each 256B input is read at once, and they instantly update one location of each of 256 accumulator RAMs.

Derive



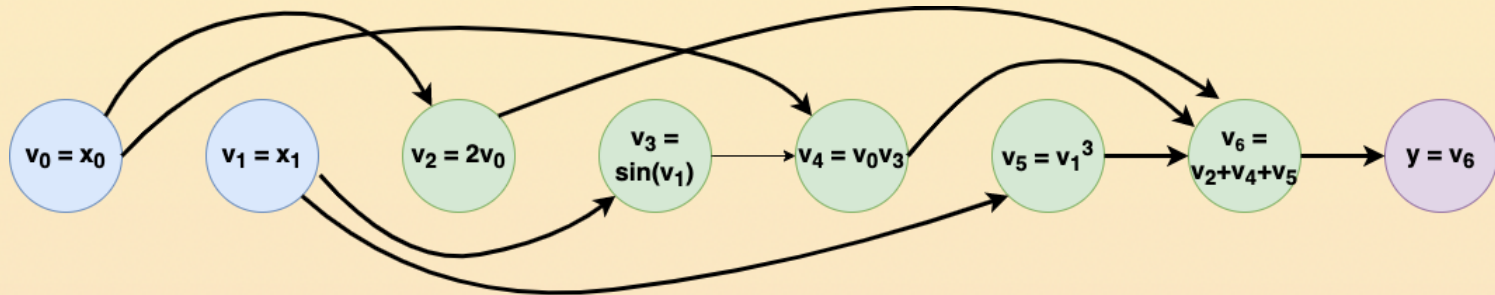
Automatic differentiation

has many names

- Automatic differentiation
- Algorithmic differentiation
- AD
- Autodiff
- Algodiff
- Autograd

Automatic differentiation

$$z(x, y) = 2x + x \sin(y) + y^3$$



Forward mode

- To the extreme, $f : \mathbb{R} \rightarrow \mathbb{R}^m$
- Evaluates $(\frac{\partial f_1}{\partial x}, \dots, \frac{\partial f_m}{\partial x})$
- Computational cost of calculating $\mathbf{J}_f(\mathbf{x})$ for $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ in $\mathbb{R}^n \times \mathbb{R}^m$

$$\mathcal{O}(n \text{ time}(f))$$

Reverse mode

- To the extreme, $f : \mathbb{R}^n \rightarrow \mathbb{R}$
- Evaluate $\nabla f(\mathbf{x})(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n})$

$$\mathcal{O}(m \text{ time}(f))$$

Forward and reverse (==backprop) modes

Primal: independent to dependent

Adjoint (derivatives): dependent to independent

$$y(\mathbf{x}) = 2x_0 + x_0 \sin(x_1) + x_1^3$$

Fwd Primal Trace Atomic operation	Value in (1, 2)	Fwd Tangent Trace (set $\dot{x}_0 = 1$ to compute $\frac{\partial y}{\partial x_0}$) Atomic operation	Value in (1, 2)
$v_0 = x_0$ $v_1 = x_1$	1 2	$\dot{v}_0 = \dot{x}_0$ $\dot{v}_1 = \dot{x}_1$	1 0
$v_2 = 2v_0$ $v_3 = \sin(v_1)$ $v_4 = v_0 v_3$ $v_5 = v_1^3$ $v_6 = v_2 + v_4 + v_5$	2 0.9093 0.9093 8 10.9093	$\dot{v}_2 = 2\dot{v}_0$ $\dot{v}_3 = \dot{v}_1 \cos(v_1)$ $\dot{v}_4 = \dot{v}_0 v_3 + v_0 \dot{v}_3$ $\dot{v}_5 = 3\dot{v}_1 v_1^2$ $\dot{v}_6 = \dot{v}_2 + \dot{v}_4 + \dot{v}_5$	2×1 $0 \times$ -0.41 $1 \times$ $0.9093 +$ 1×0 $3 \times 0 \times 4$ $2 +$ $0.9093 +$ 0
$y = v_6$	10.9093	$\dot{y} = \dot{v}_6$	2.9093

Fwd Primal Trace Atomic operation	Value in (1, 2)	Rev Adjoint Trace (set $\bar{y} = 1$ to compute $\frac{\partial v}{\partial y}$) Atomic operation	Value in (1, 2)
$v_0 = x_0$ $v_1 = x_1$	1 2	$\bar{x}_0 = \bar{v}_0$ $\bar{x}_1 = \bar{v}_1$	2.9093 11.5839
$v_2 = 2v_0$ $v_3 = \sin(v_1)$ $v_4 = v_0 v_3$ $v_5 = v_1^3$ $v_6 = v_2 + v_4 + v_5$	2 0.9093 0.9093 8 10.9093	$\bar{v}_0 = \bar{v}_0 + \bar{v}_2 \partial v_2 / \partial v_0$ $\bar{v}_0 = \bar{v}_0 + \bar{v}_4 \partial v_4 / \partial v_0$ $\bar{v}_1 = \bar{v}_1 + \bar{v}_3 \partial v_3 / \partial v_1$ $\bar{v}_1 = \bar{v}_1 + \bar{v}_5 \partial v_5 / \partial v_1$ $\bar{v}_2 = \bar{v}_2 + \bar{v}_6 \partial v_6 / \partial v_2$ $\bar{v}_3 = \bar{v}_3 + \bar{v}_4 \partial v_4 / \partial v_3$ $\bar{v}_4 = \bar{v}_4 + \bar{v}_6 \partial v_6 / \partial v_4$ $\bar{v}_5 = \bar{v}_5 + \bar{v}_6 \partial v_6 / \partial v_5$	$\bar{v}_0 + \bar{v}_2 \times 2 = 2.9093$ $\bar{v}_4 \times v_3 = 0.9093$ $\bar{v}_1 + \bar{v}_3 \times \cos(v_1) = 11.5839$ $\bar{v}_5 \times 3v_1^2 = 12$ $\bar{v}_6 \times 1 = 1$ $\bar{v}_4 \times v_0 = 1$ $\bar{v}_6 \times 1 = 1$ $\bar{v}_6 \times 1 = 1$
$y = v_6$	10.9093	$\bar{y} = \bar{y}$	1

Designed to be simple in software

```
import torch, math
x0 = torch.tensor(1., requires_grad=True)
x1 = torch.tensor(2., requires_grad=True)
p = 2*x0 + x0*torch.sin(x1) + x1**3
print(p)
p.backward()
print(x0.grad, x1.grad)
```

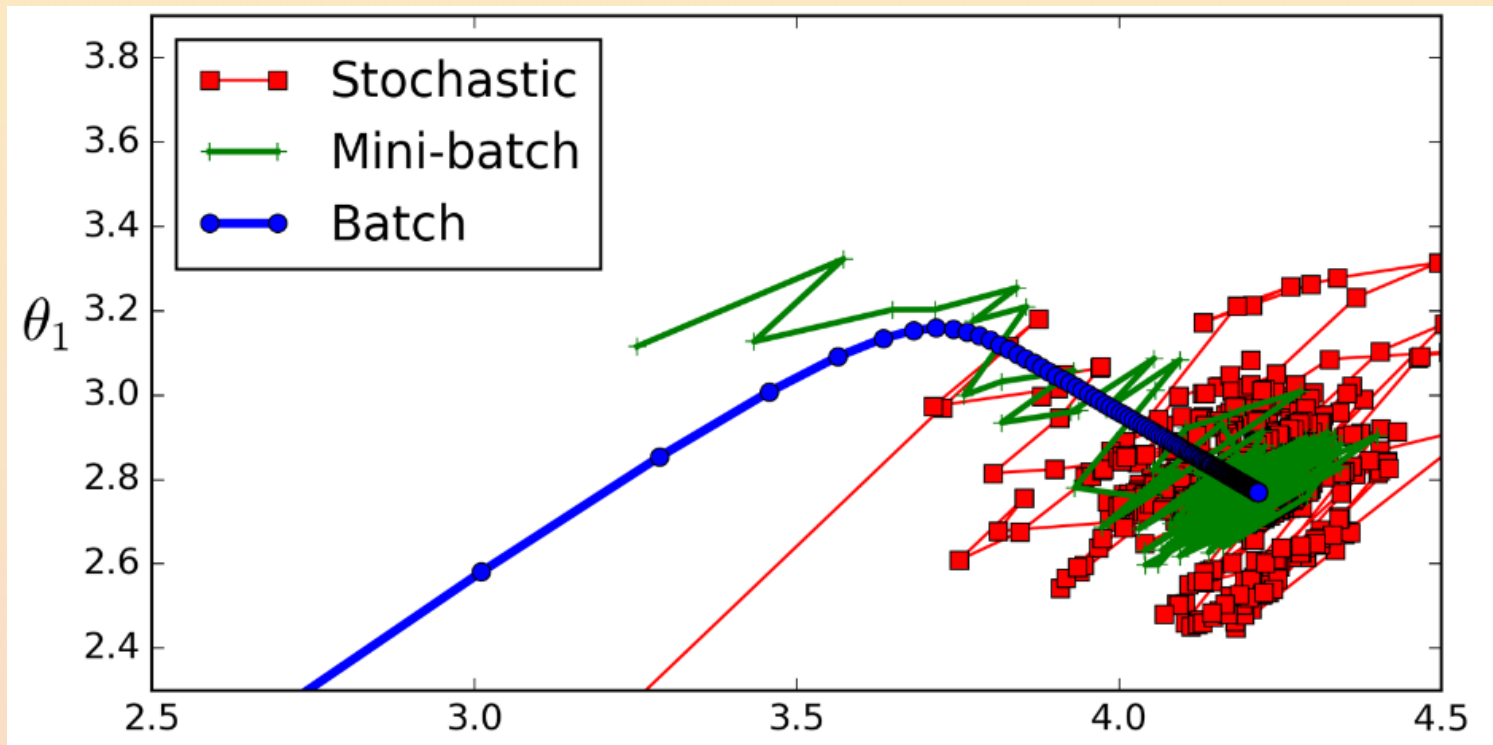
yielding

```
Primal: tensor(10.9093, grad_fn=<AddBackward0>)
Adjoint: tensor(2.9093) tensor(11.5839)
```

- Torch (and similar software) will correctly differentiate **only when the atomic operations are supported within it**
 - Common operations are overloaded (`__mul__` rewritten by `torch._mult_`)
 - Operations from libraries (`math.sin()`) must be replaced by their differentiation-aware equivalents (`torch.sin()`)

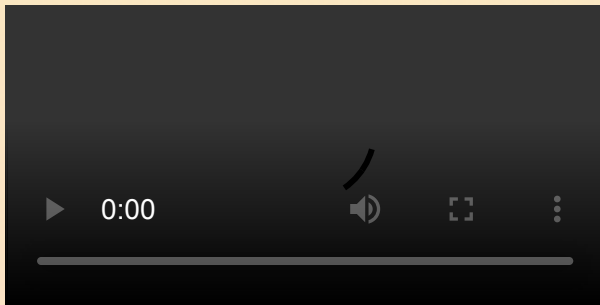
Sampling scheme

- Batch: compute on the whole training set (for large sets becomes too costly)
- Stochastic: compute on one sample (large noise, difficult to converge)
- Mini-batch: use a relatively small sample of data (tradeoff)



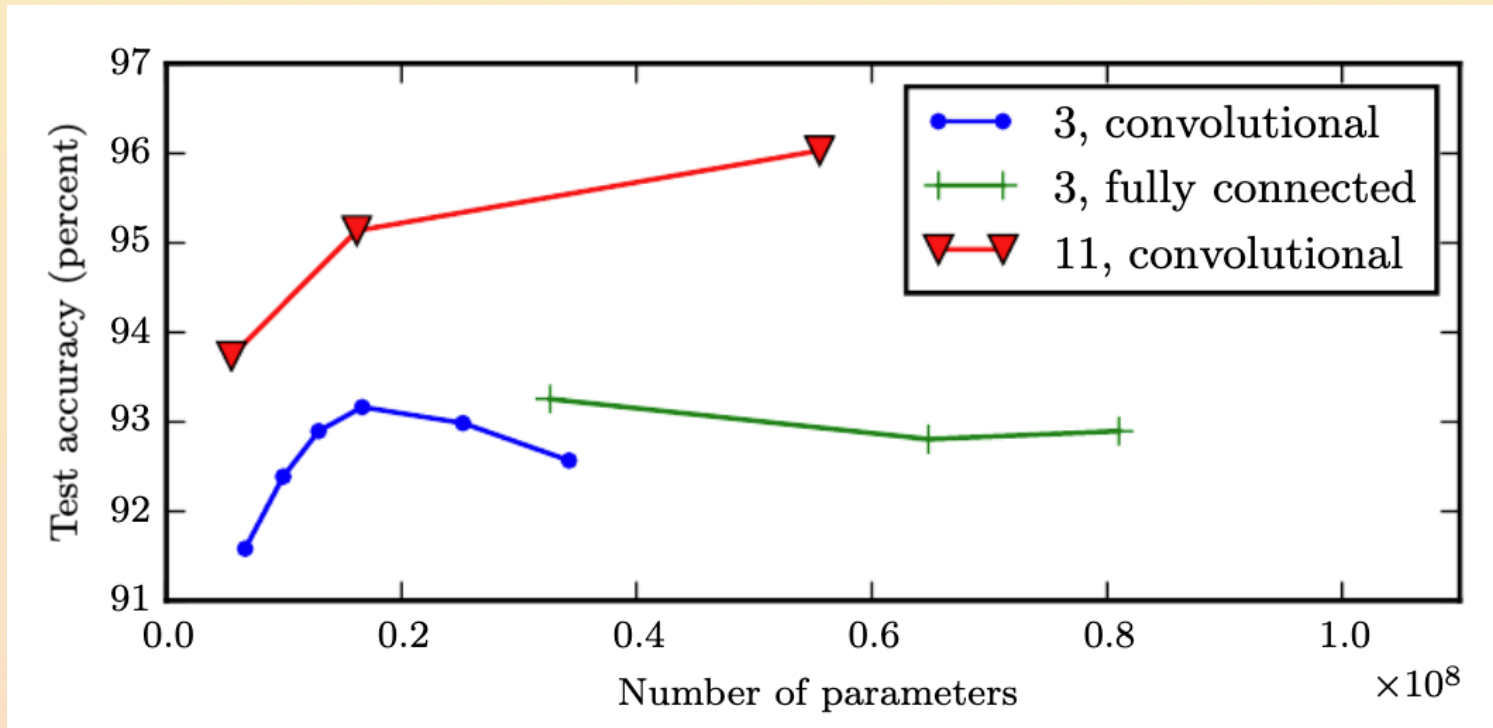
Descent strategies

- Mostly nonconvex optimization: very complicated problem, convergence in general not guaranteed
- Nesterov momentum: big jumps followed by correction seem to help!
- Adaptive moments: gradient steps decrease when getting closer to the minimum (avoids overshooting)



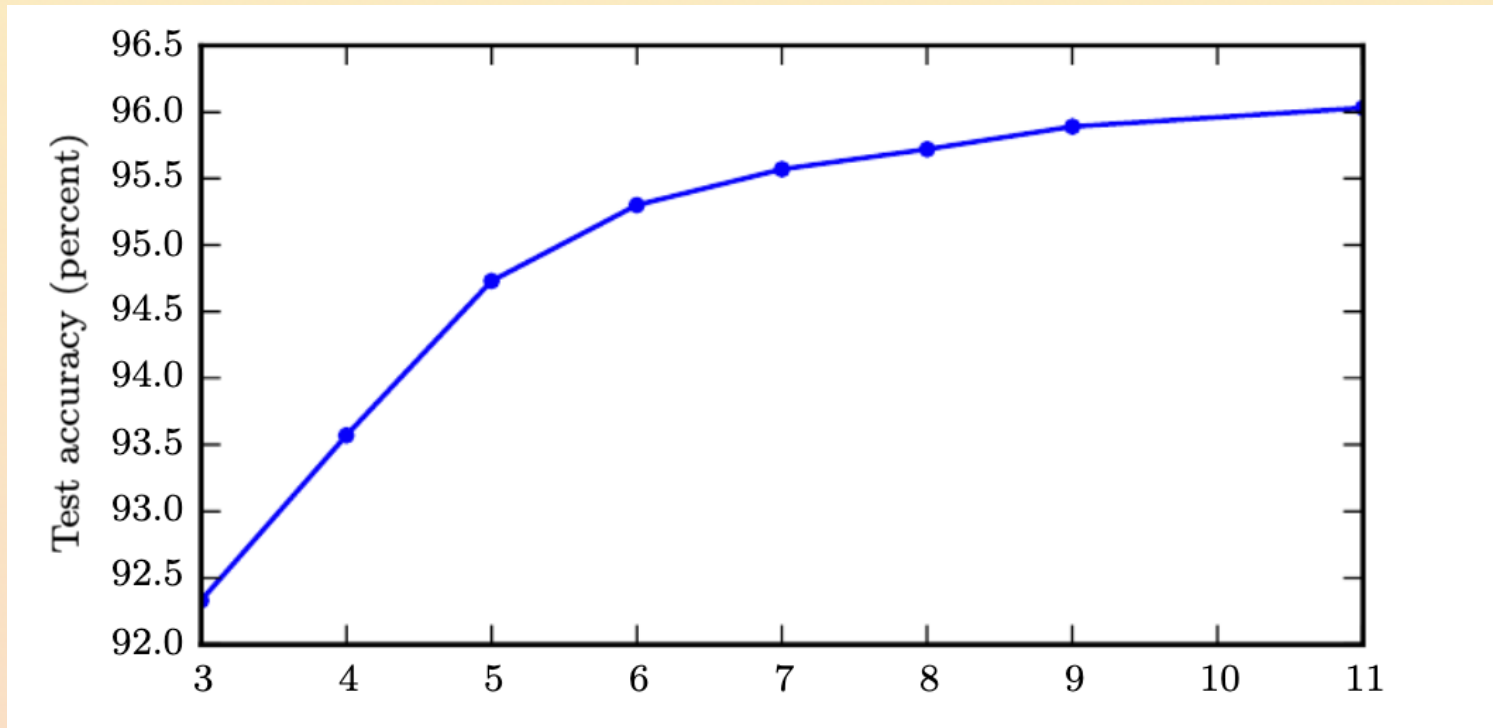
Number of parameters

- Empirical studies: increasing number of parameters doesn't help beyond a certain point



Depth

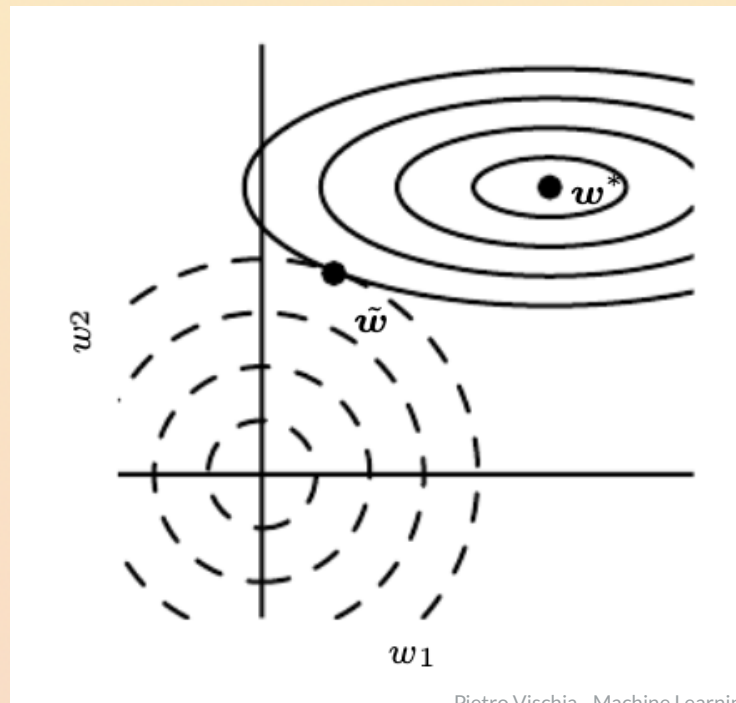
- Empirical studies: increasing depth tends to always result in some improvement



Regularization: weight decay

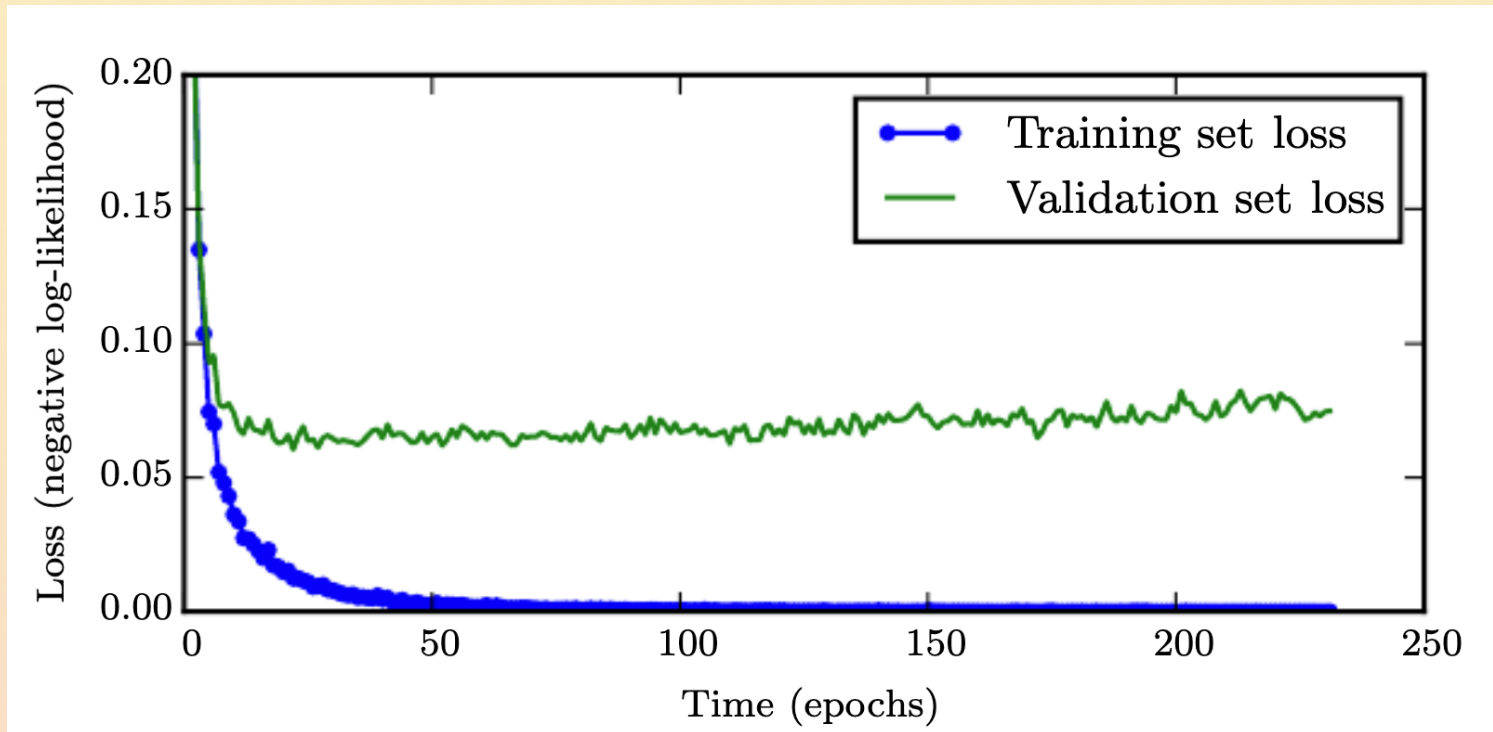
- Another way of regularizing is via **weight decay** (L^2 regularization)
 - Tradeoff between good fitting (small MSE) and small norm (smaller slope, or fewer features with large weights)
 - In statistics, "ridge regression", "Tikhonov regularization"

$$J(\mathbf{w}) = MSE_{train} + \lambda \mathbf{w}^T \mathbf{w}$$



Early stopping...

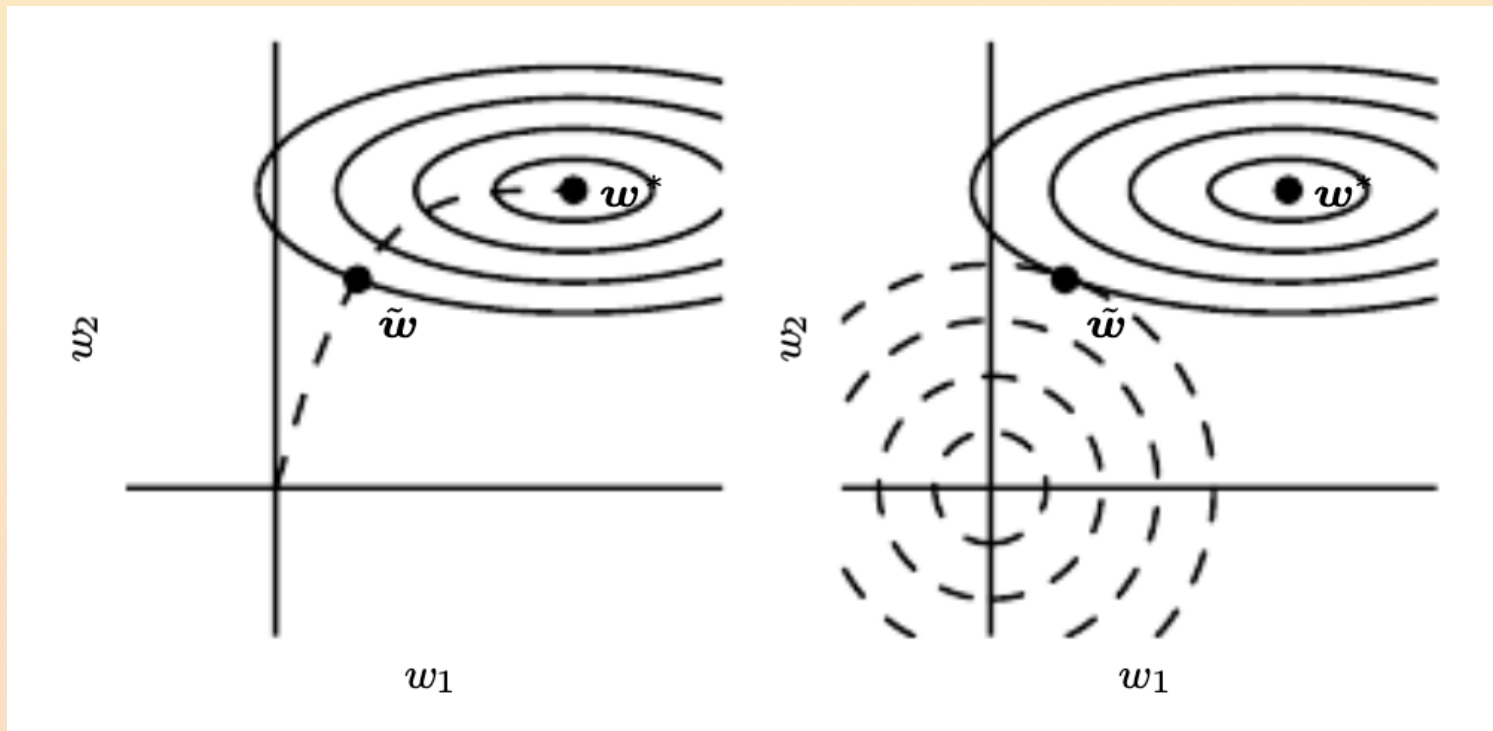
- Train until the validation set loss starts increasing, and pick the model corresponding to the minimum validation loss



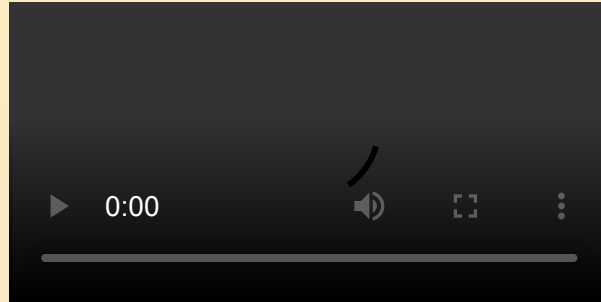
...is a form of regularisation

- Early stopping limits the reachable phase space, and is therefore analogous to L2 regularization (weight decay)

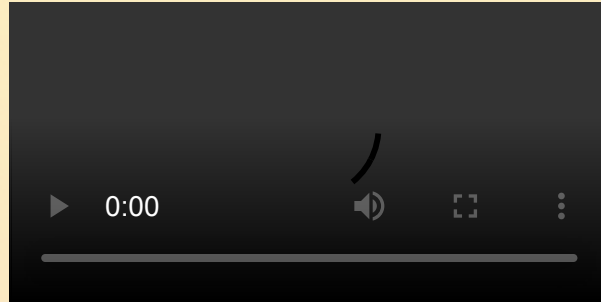
$$J(\mathbf{w}) = MSE_{train} + \lambda \mathbf{w}^T \mathbf{w}$$



Impressive results



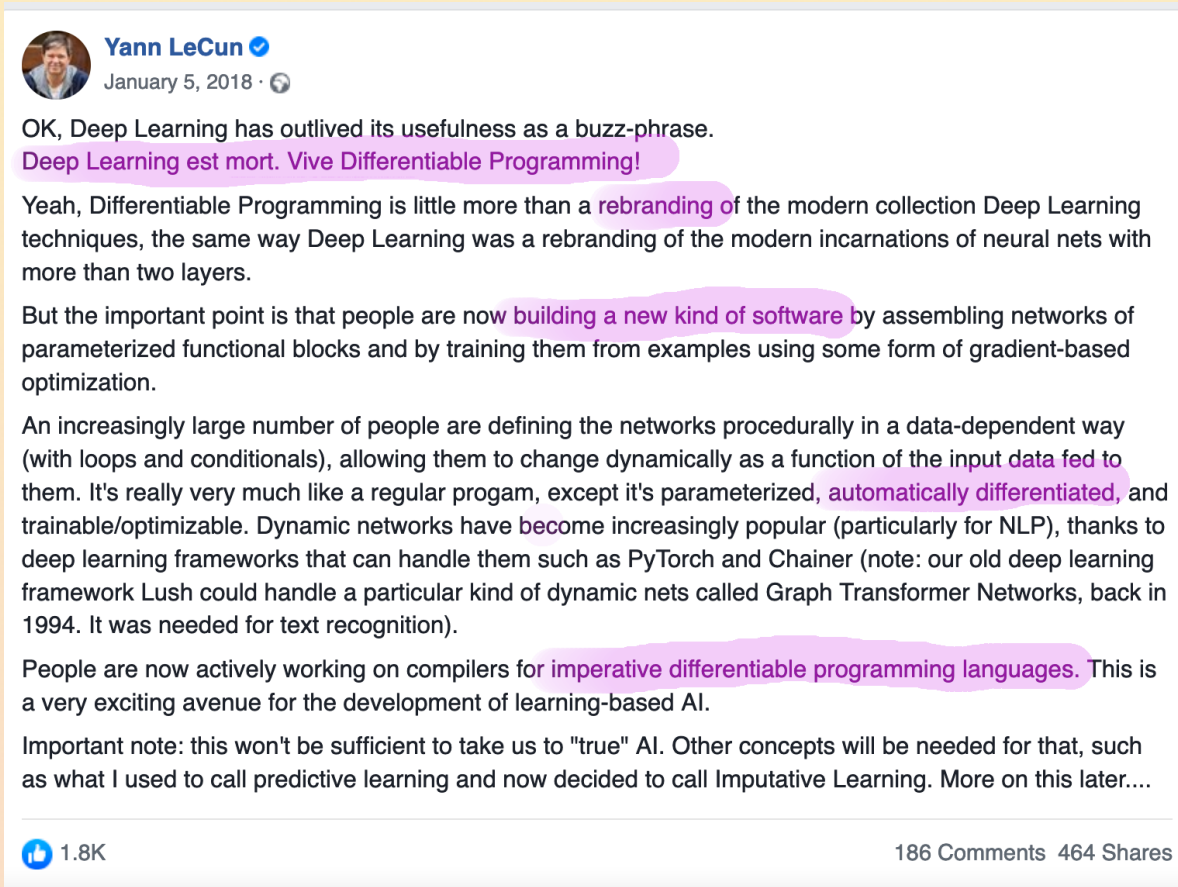
Impressive results



- Busco colaboraciones para aplicaciones médicas de inteligencia artificial

Differentiable Programming

Execute **differentiable functions (programs)** via **automatic differentiation**



Yann LeCun ✓
January 5, 2018 · 🌐

OK, Deep Learning has outlived its usefulness as a buzz-phrase.
Deep Learning est mort. Vive Differentiable Programming!

Yeah, Differentiable Programming is little more than a **rebranding** of the modern collection Deep Learning techniques, the same way Deep Learning was a rebranding of the modern incarnations of neural nets with more than two layers.

But the important point is that people are now **building a new kind of software** by assembling networks of parameterized functional blocks and by training them from examples using some form of gradient-based optimization.

An increasingly large number of people are defining the networks procedurally in a data-dependent way (with loops and conditionals), allowing them to change dynamically as a function of the input **data fed to** them. It's really very much like a regular program, except it's parameterized, **automatically differentiated**, and trainable/optimizable. Dynamic networks have **become** increasingly popular (particularly for NLP), thanks to deep learning frameworks that can handle them such as PyTorch and Chainer (note: our old deep learning framework Lush could handle a particular kind of dynamic nets called Graph Transformer Networks, back in 1994. It was needed for text recognition).

People are now actively working on compilers for **imperative differentiable programming languages**. This is a very exciting avenue for the development of learning-based AI.

Important note: this won't be sufficient to take us to "true" AI. Other concepts will be needed for that, such as what I used to call predictive learning and now decided to call Imputative Learning. More on this later...

👍 1.8K 186 Comments 464 Shares

Many ways of inserting our biases

- Regularization corresponds to inserting our bias into the algorithm
 - "I know that the solution should not wiggle", "I know that the curvature must not be too large"
- Two models A and B performing the same classification task
 - $\hat{y}^{(A)} = f(\mathbf{w}^{(A)}, \mathbf{x})$
 - $\hat{y}^{(B)} = f(\mathbf{w}^{(B)}, \mathbf{x})$
- If the inputs distributions are somehow different, but we know (or want that) the output are related, we can assume that the weights should be similar

$$J(\mathbf{w}) = MSE_{train} + \lambda \|\mathbf{w}^{(A)} - \mathbf{w}^{(B)}\|_2^2$$

Parameter sharing

- Simply require parameters are **equal**
 - If they are equal, you can store only one number in memory (sometimes dramatic memory footprint reduction)



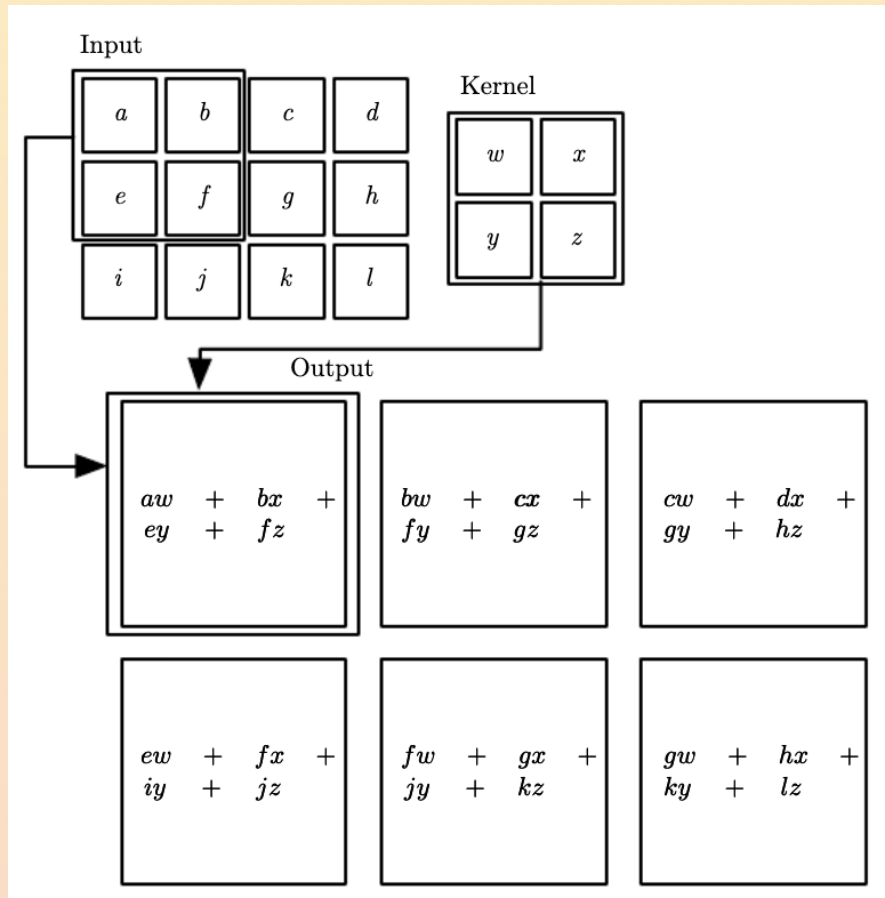
Convolution: a form of averaging

$$s(t) = \int x(a)w(t - a)da$$

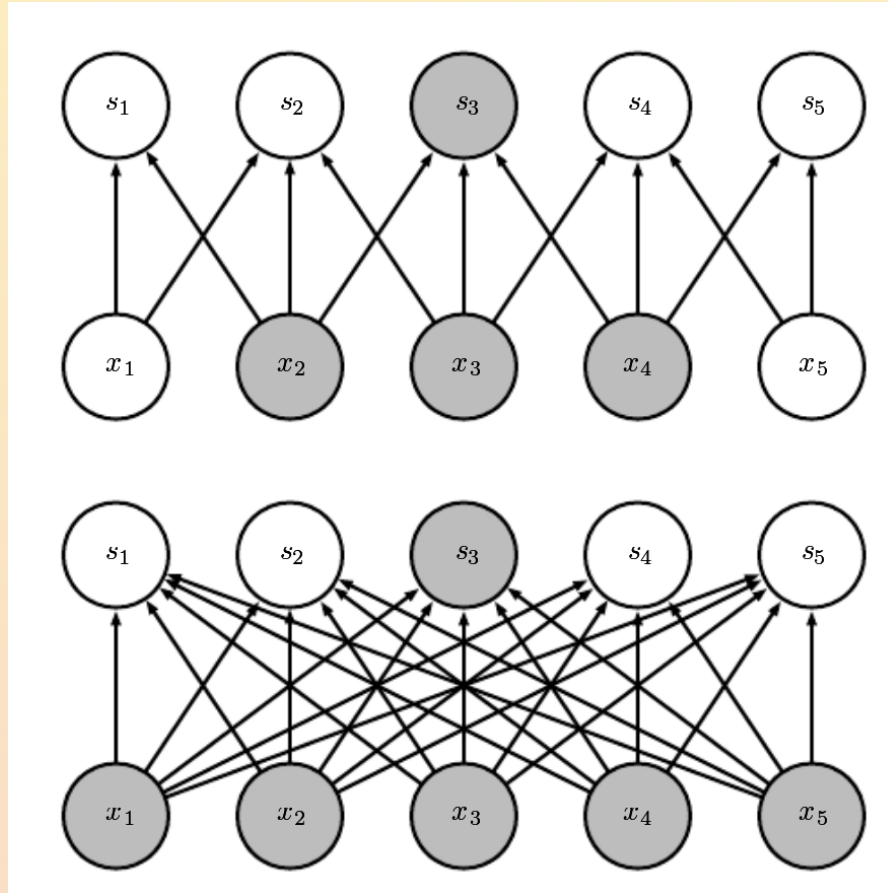
- When discretized, integral becomes a sum
 - x input
 - w kernel: specifies how far does the averaging goes
 - s feature map

Convolution: a form of averaging

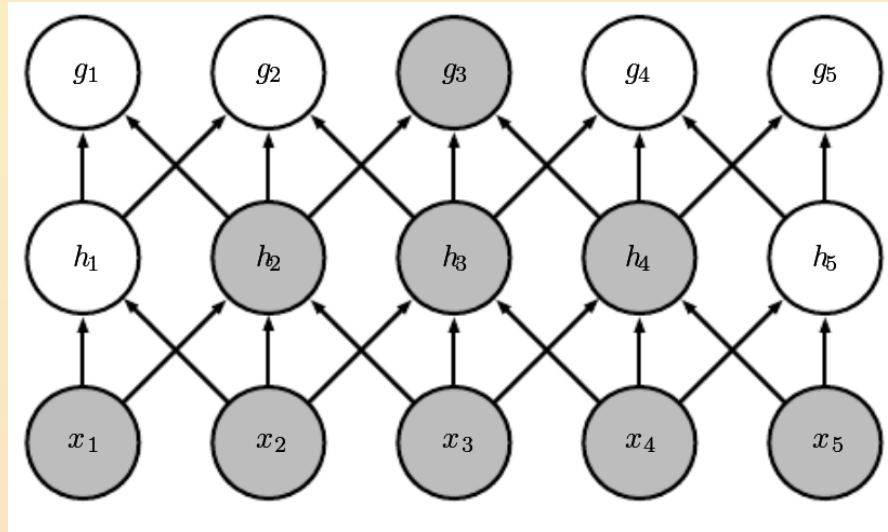
$$S(i, j) = (K * I)(i, j) = \sum_m \sum_n I(i - m, j - n)K(m, n)$$



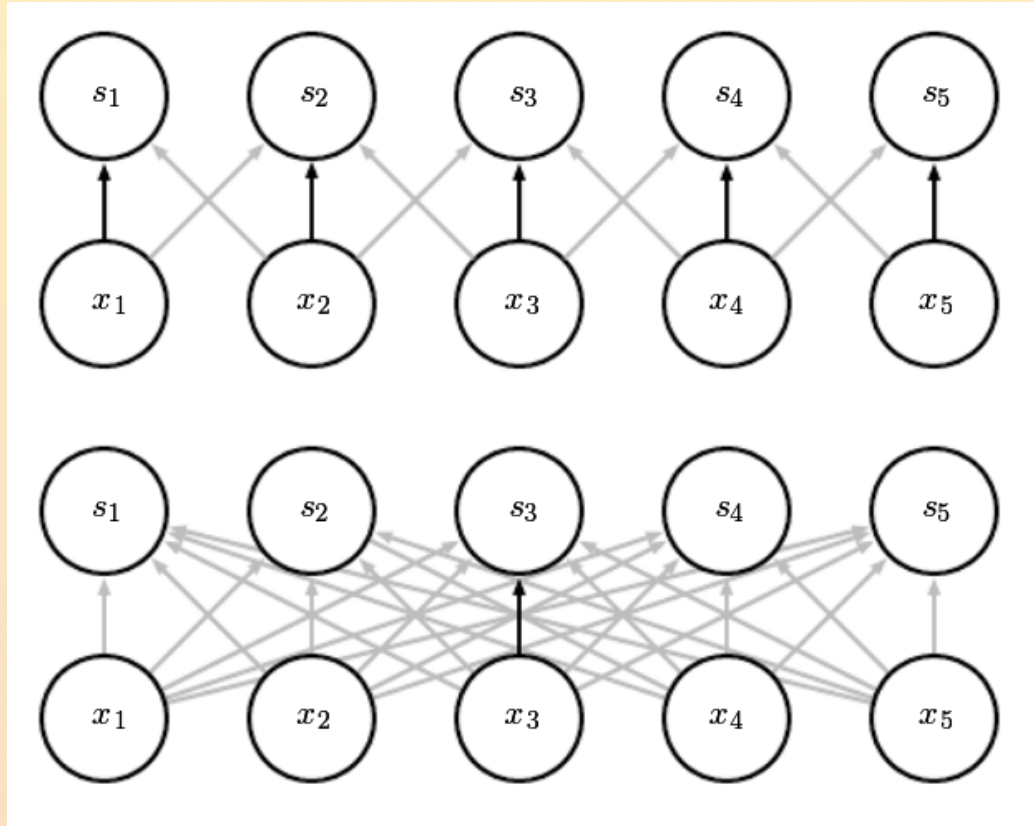
Receptive field



Receptive field: deeper = larger

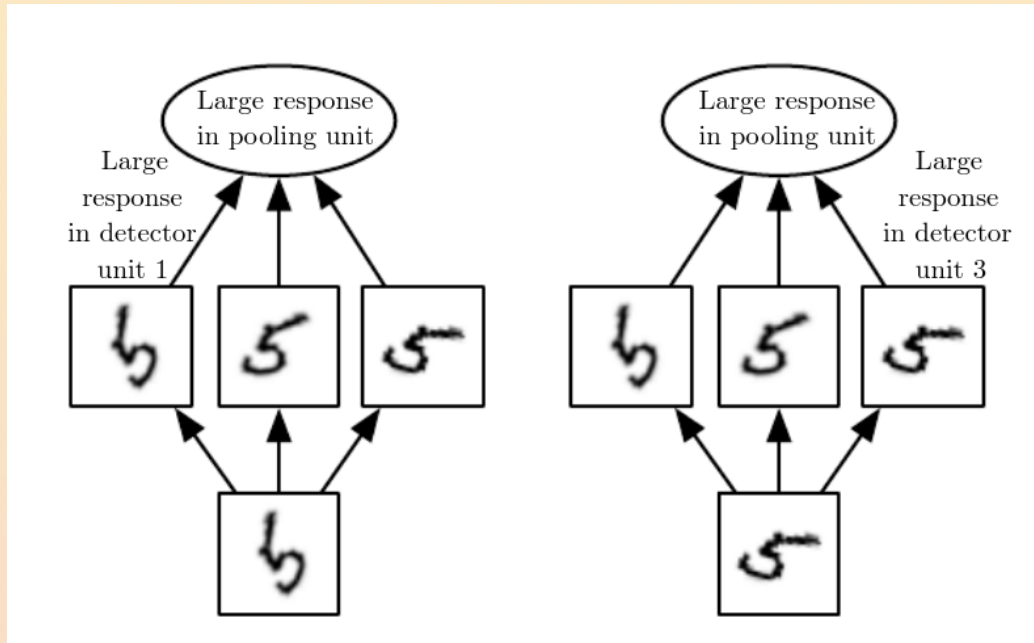


Parameter sharing

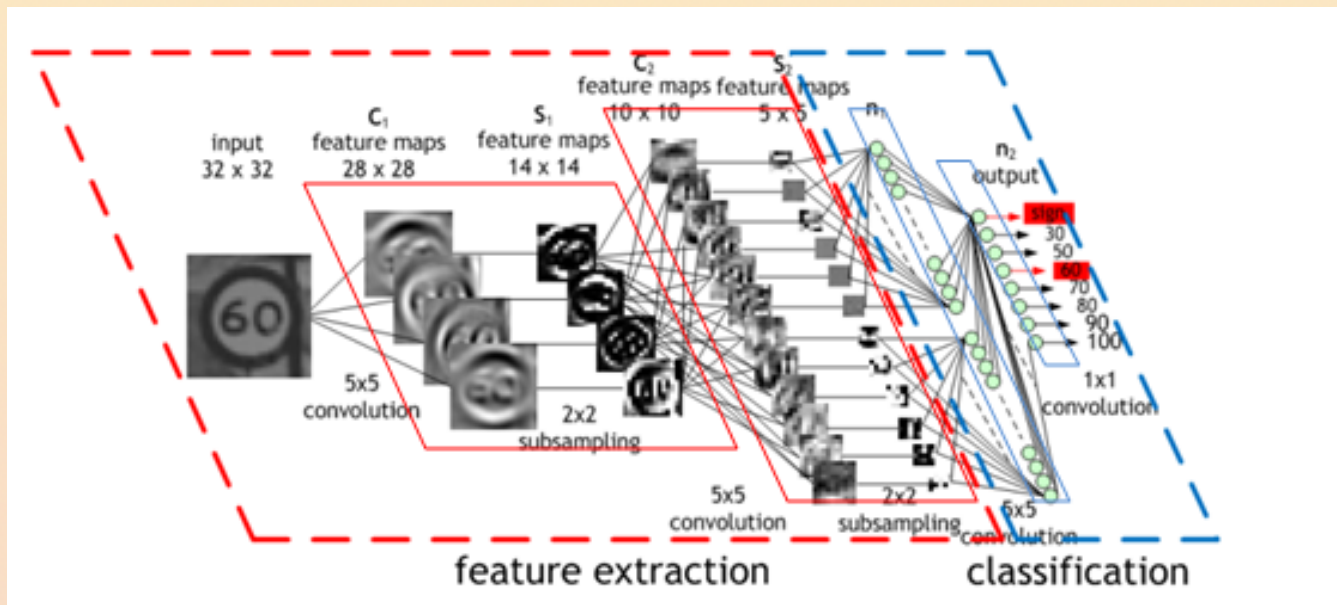


Convolutional network

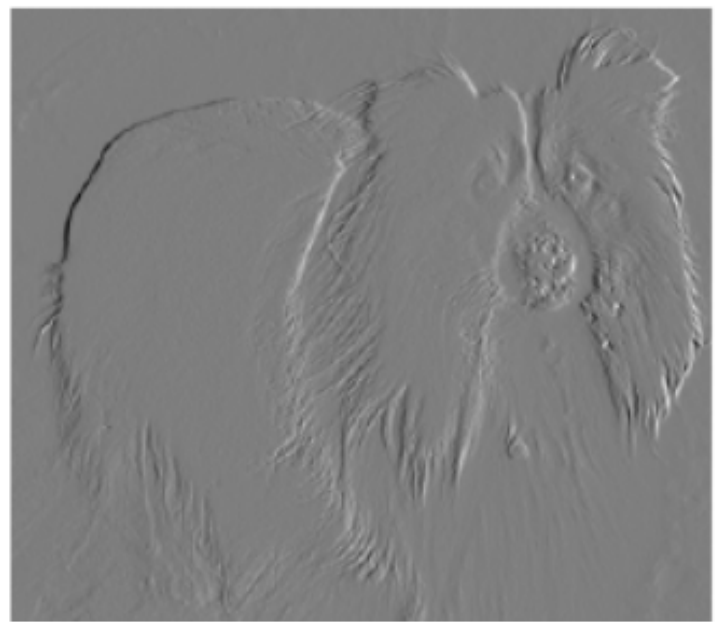
- Convolution → nonlinear activation → pooling
- Pooling: replace output at a location with a summary statistic
 - e.g., max pooling = report the maximum output in a neighbourhood
 - Helps with invariance for translations



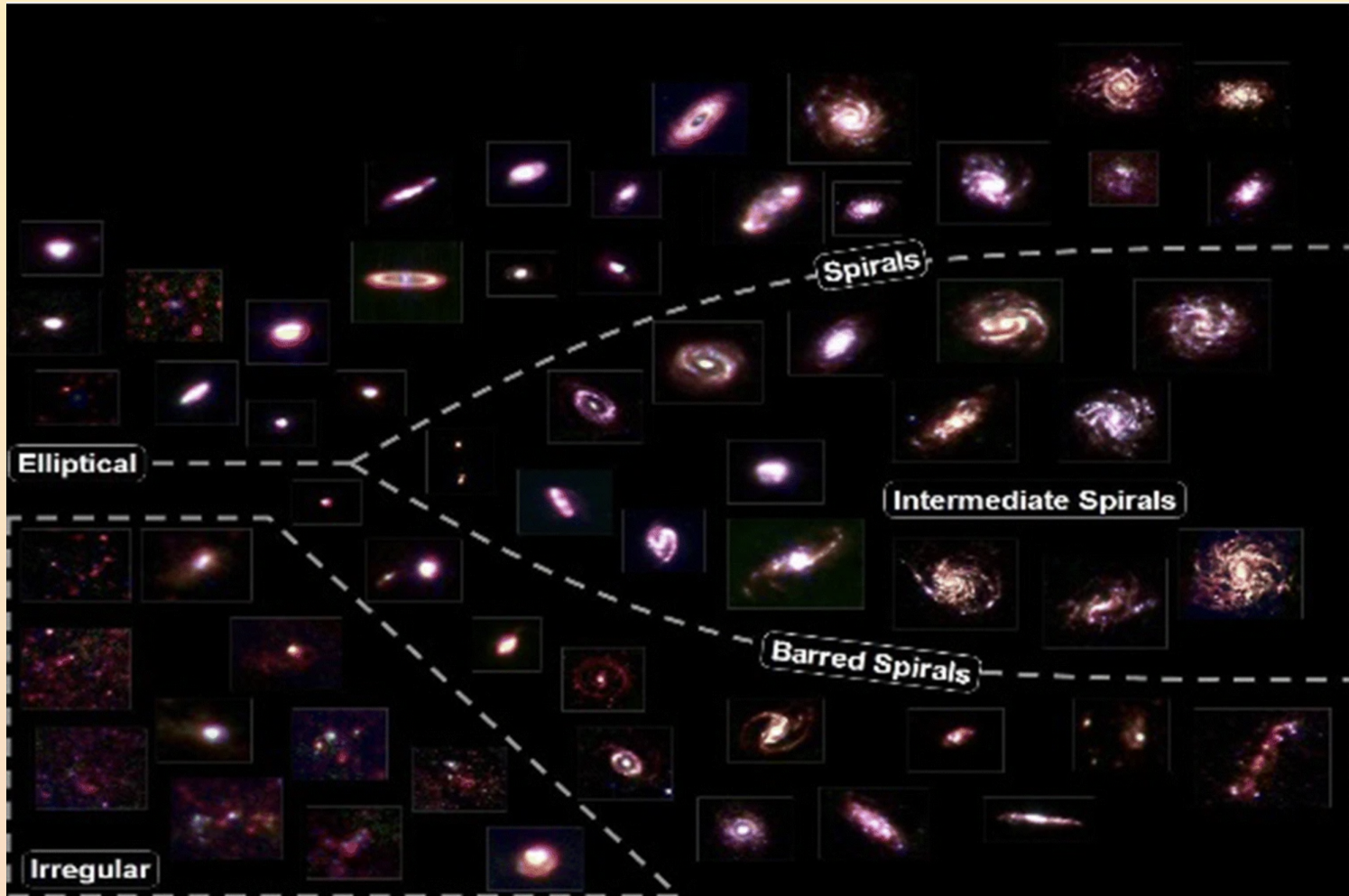
Convolutional networks



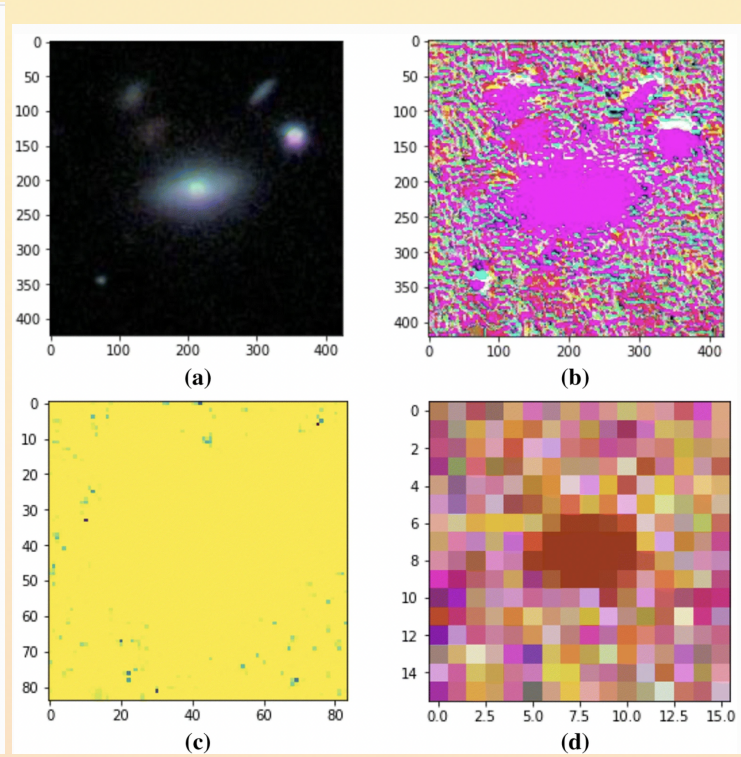
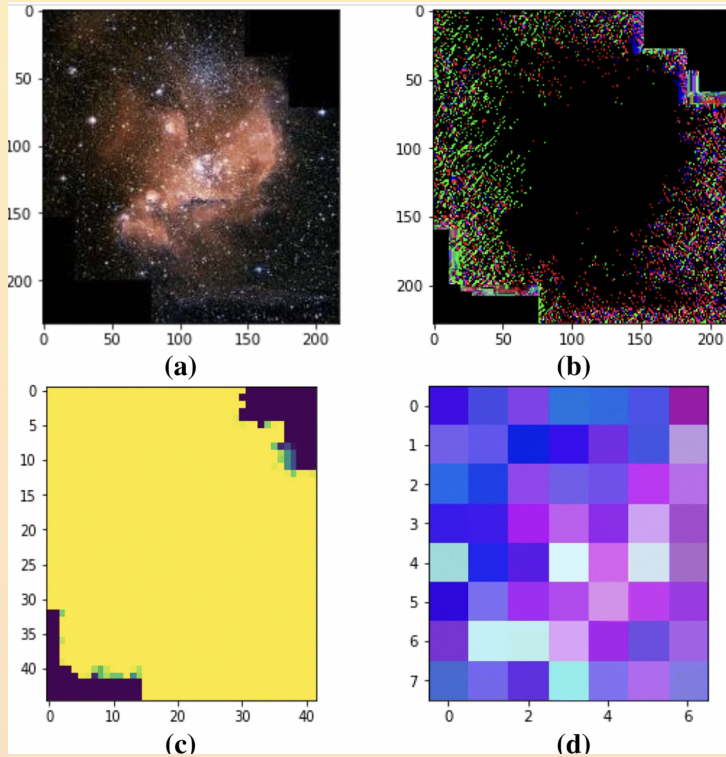
Intermediate representations



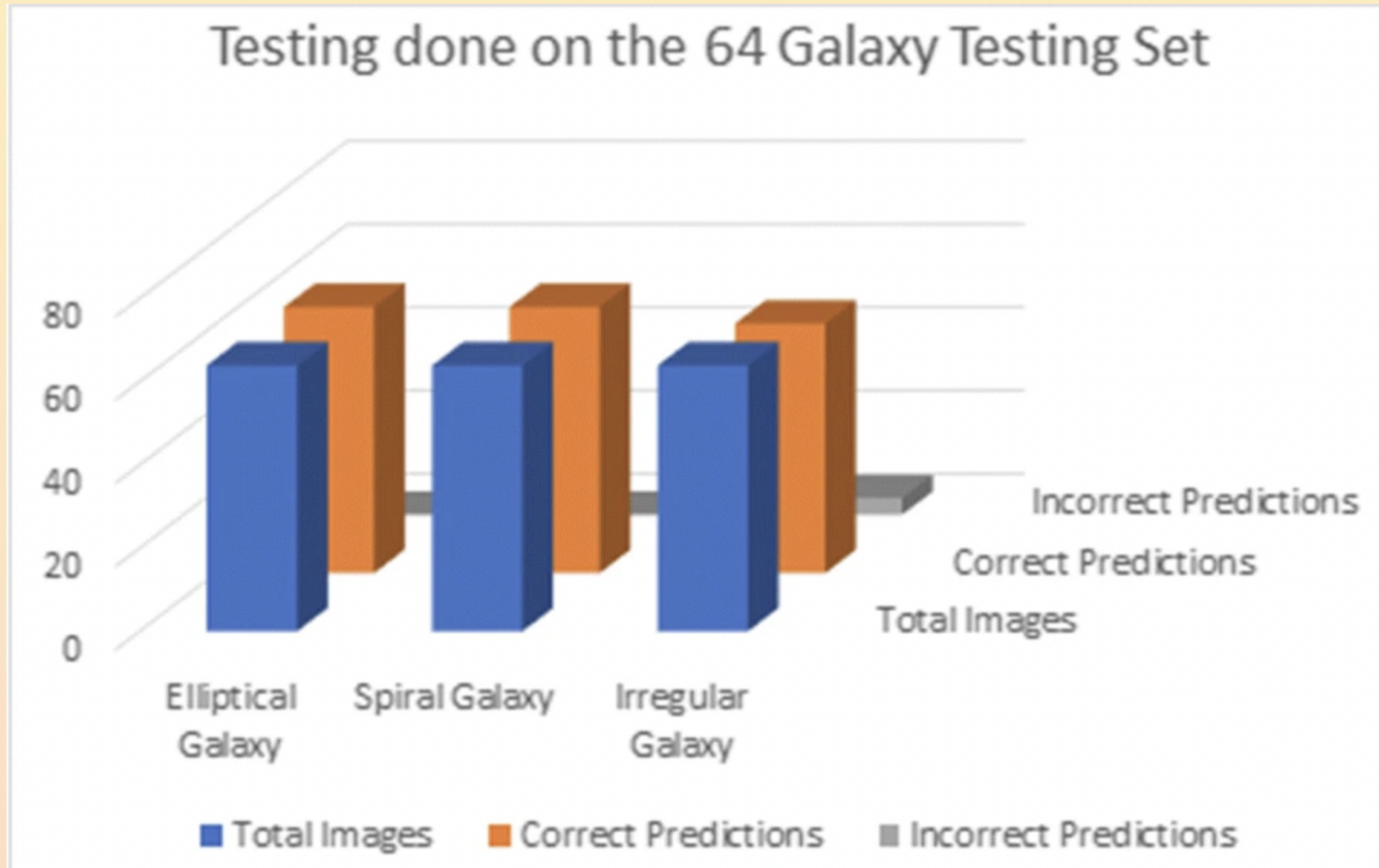
Morphology of galaxies



Representations of galaxies...



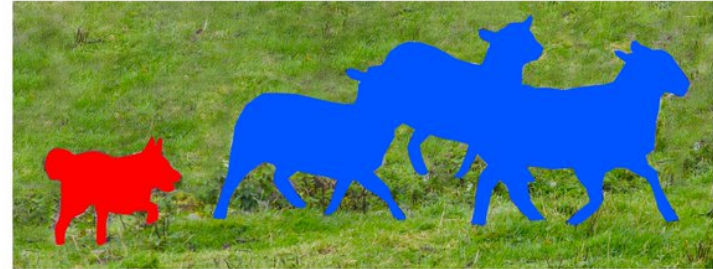
...work pretty well



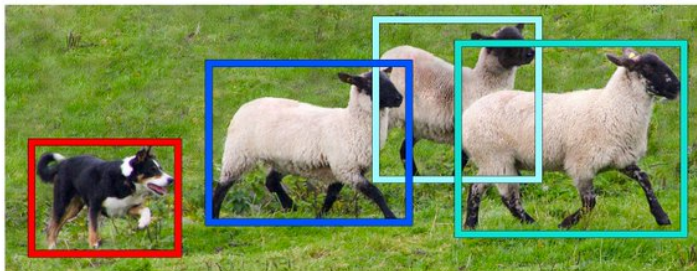
Semantic representations



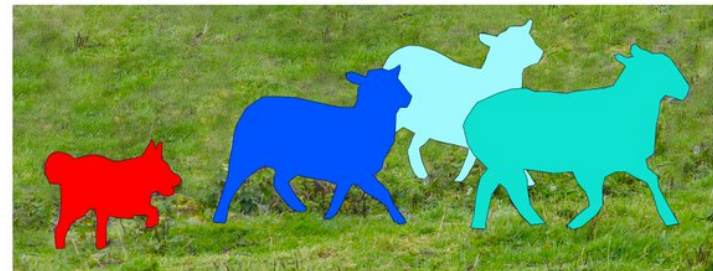
Image Recognition



Semantic Segmentation



Object Detection



Instance Segmentation

What about time?

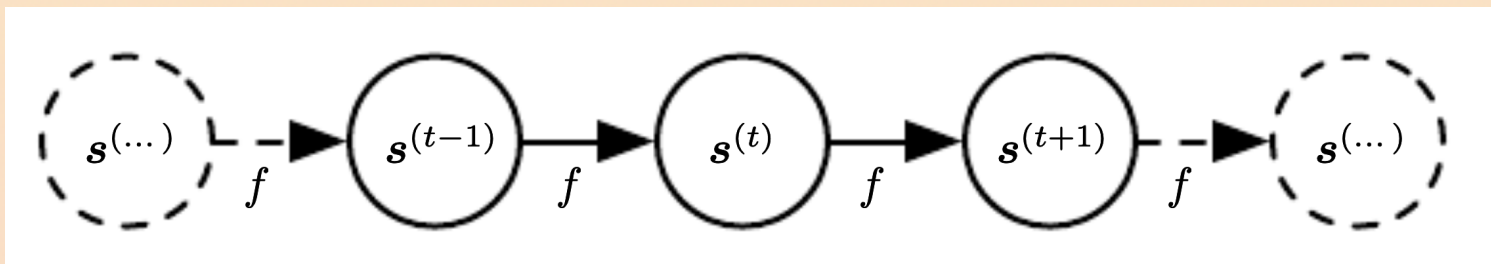
- Convolutional network: process grid of values (e.g. images)
- Recurrent networks: process a sequence of values indicised by a "time" component
 - Language is a sequence
- Parameter sharing crucial to generalize:
 - lengths unseen in training
 - different positions in the sentences
- Without parameter sharing, a network would have to learn all the language rules at each step of the sequence
 - Very impractical
- Both scale very well (thanks to parameter sharing)

Convolutional networks for sequences?

- Could "link" the steps of the sequence via the convolution
- Use the same kernel at each time step
- Shallow: it links only neighbouring time steps

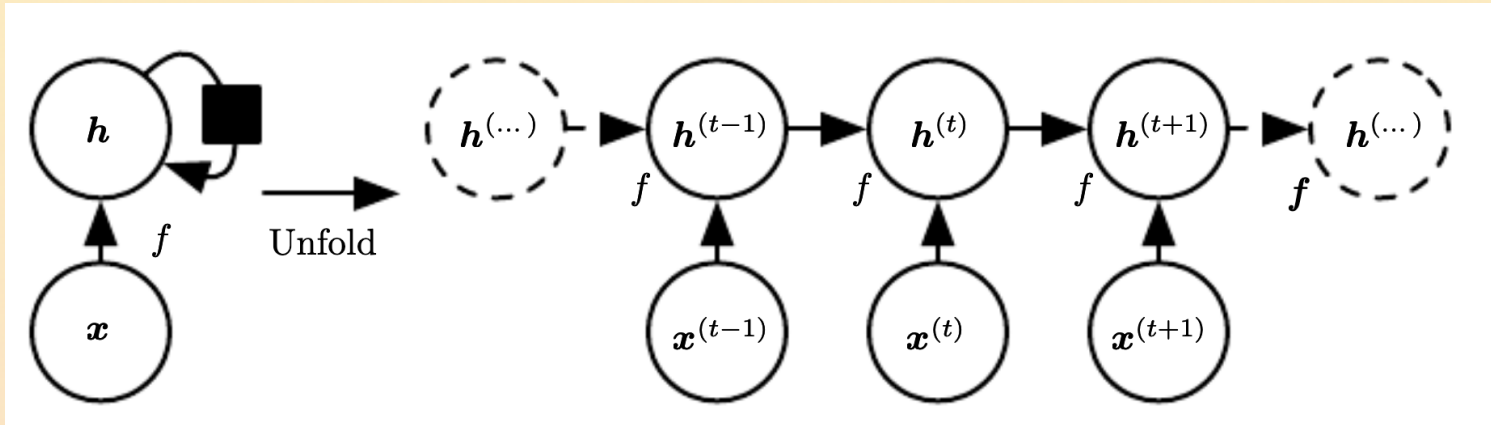
Recurrent network

- Use the same parameter at the same step, $s^{(t)} = f(s^{(t-1)}, \theta)$
 - Very deep structure



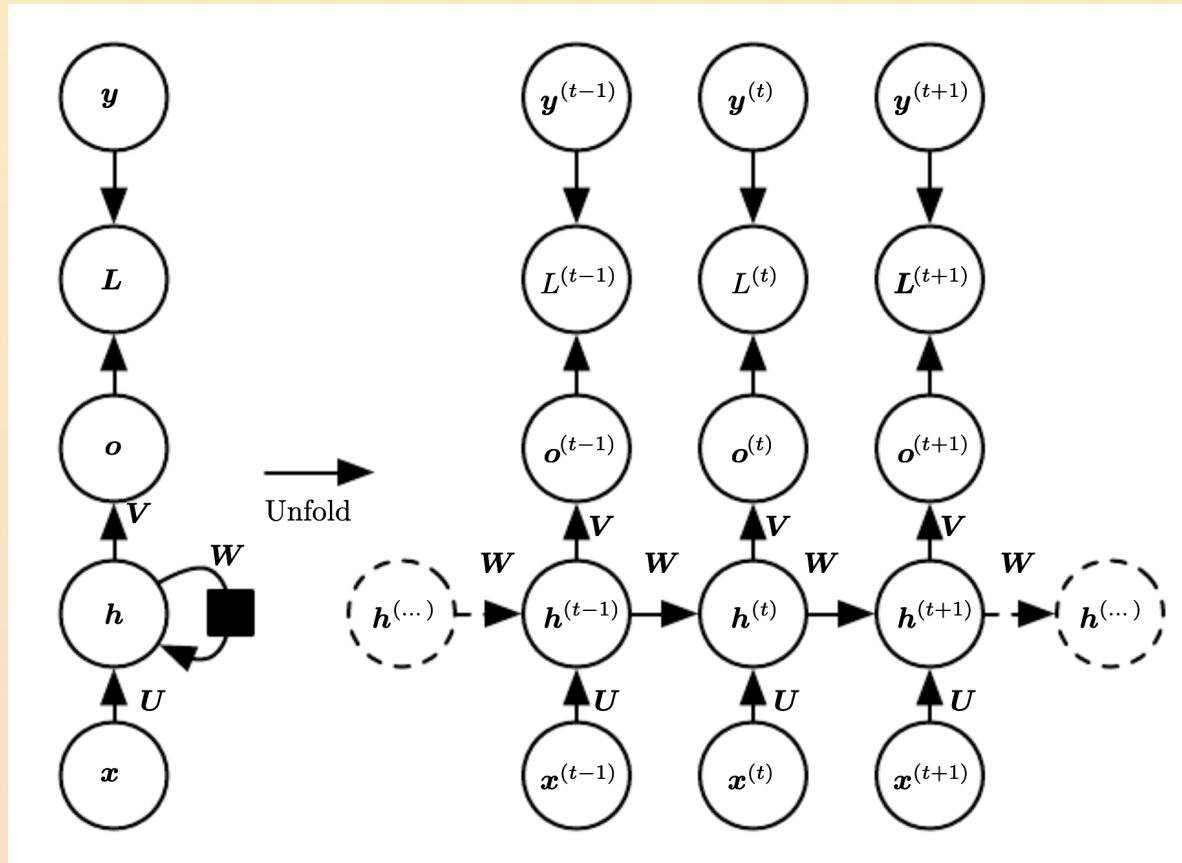
Unfold the graph

$$s^{(3)} = f(s^{(2)}, \theta) = f(f(s^{(1)}, \theta), \theta)$$



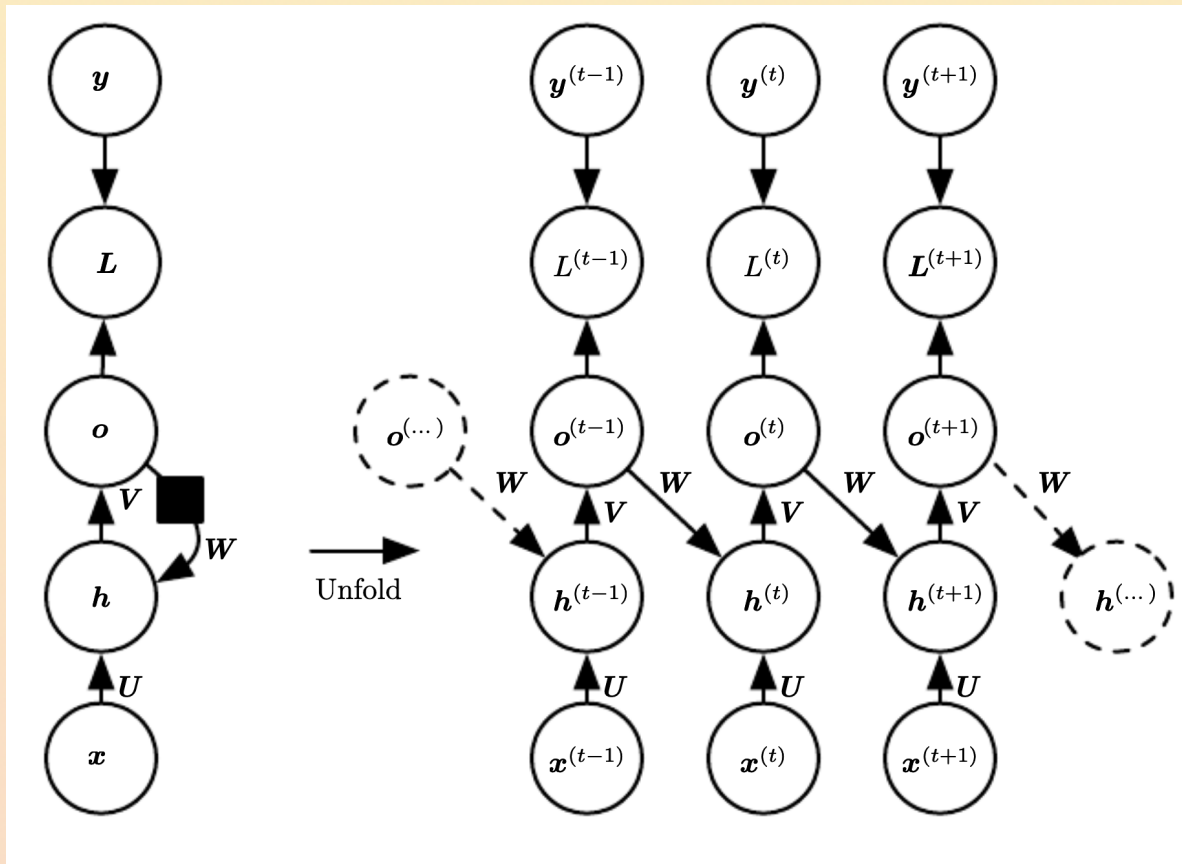
Vast zoology

= An output at each time step, recurrent connections between hidden units



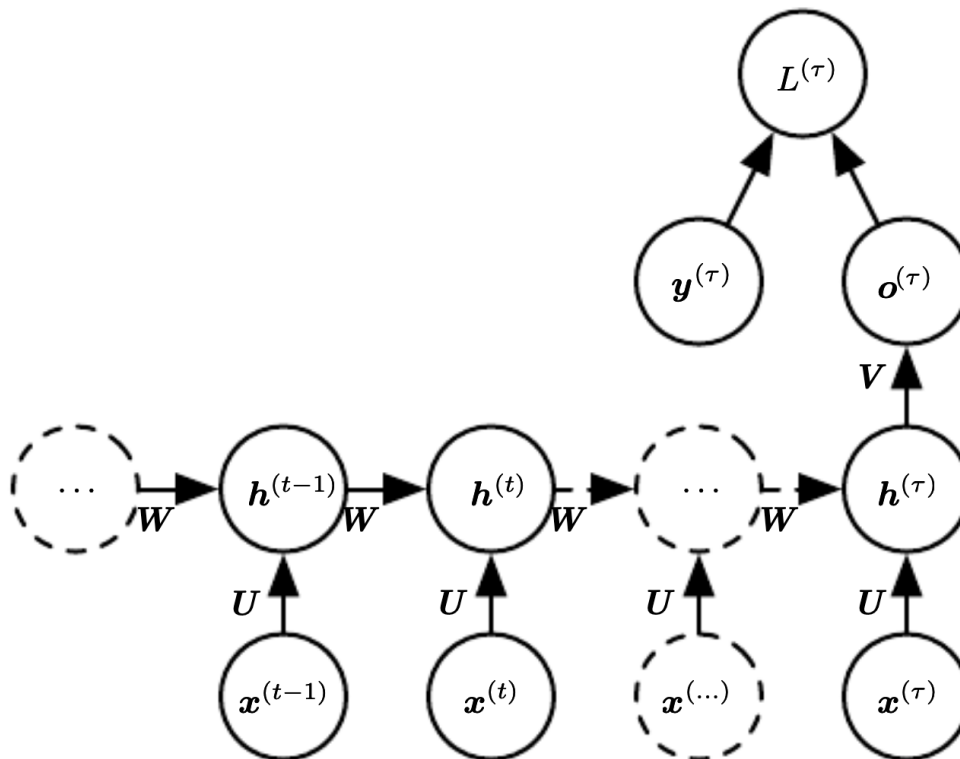
Vast zoology

- An output at each time step, recurrent connections only from the output at one time step to the hidden units at the next time step

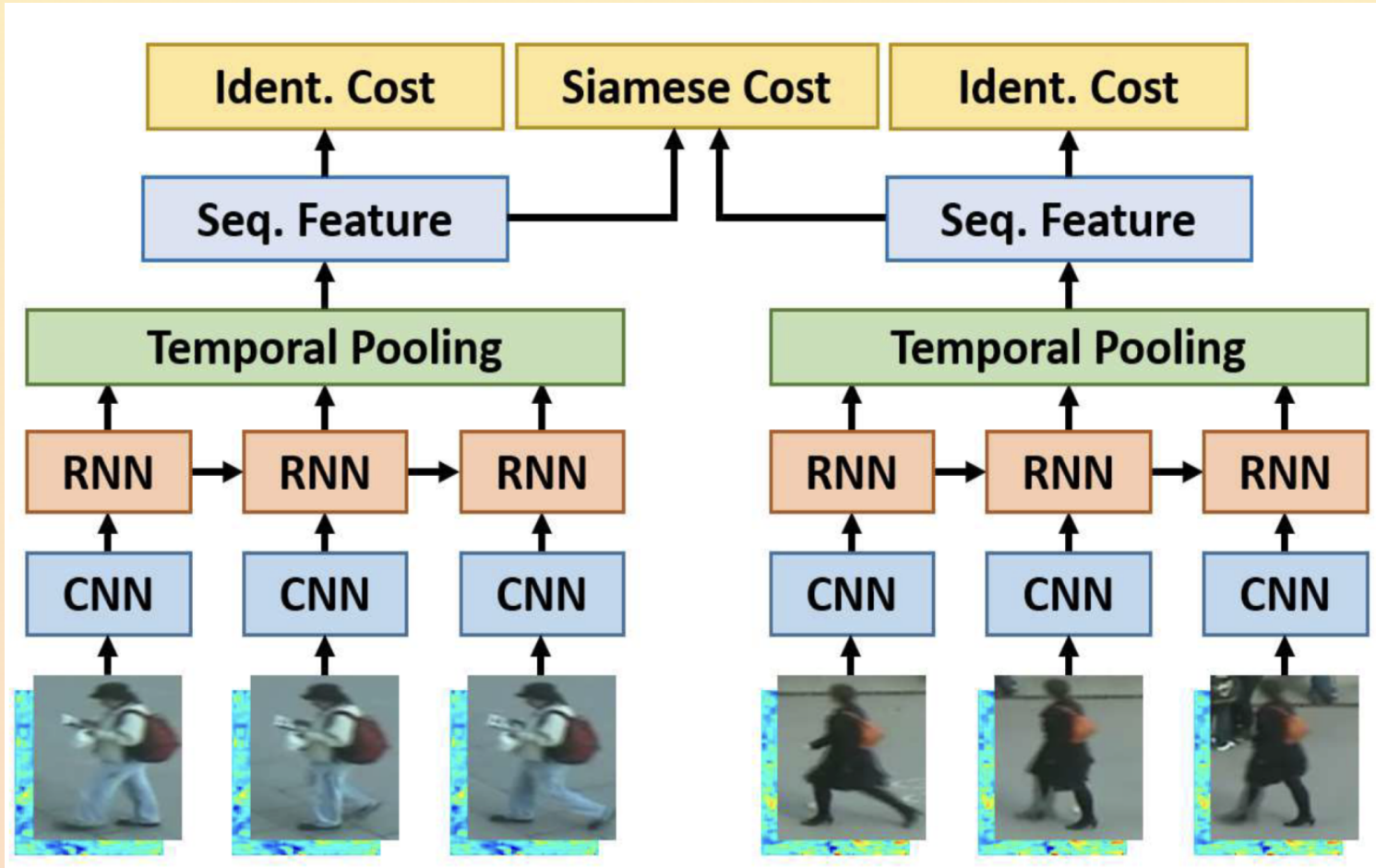


Vast zoology

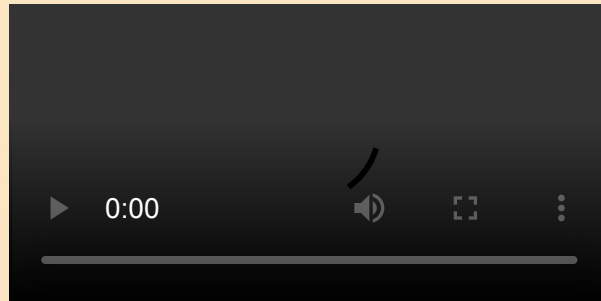
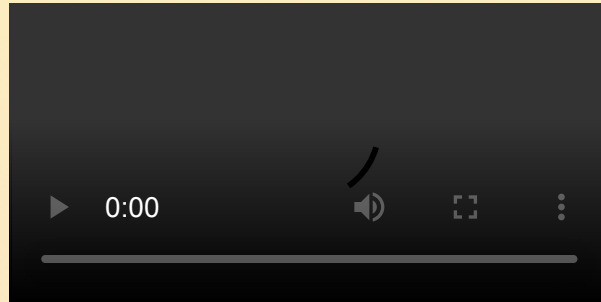
- Recurrent connections between hidden units, that read an entire sequence and then produce a single output



Sequences of images

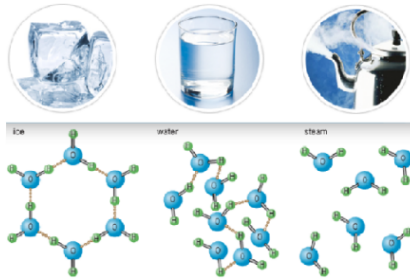


Real-time segmentation

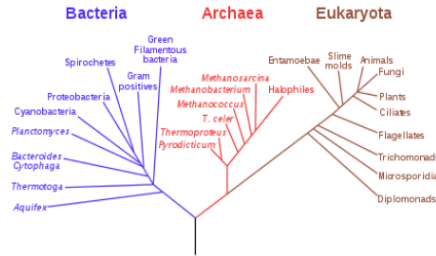


Graphs Represent Structure

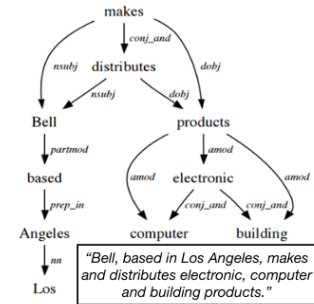
Molecules



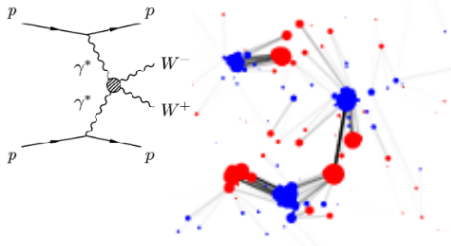
Biological species



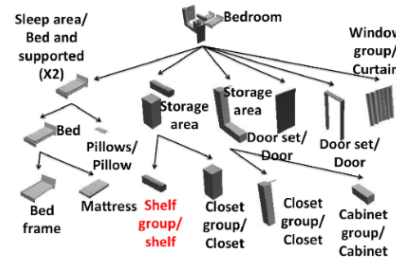
Natural language



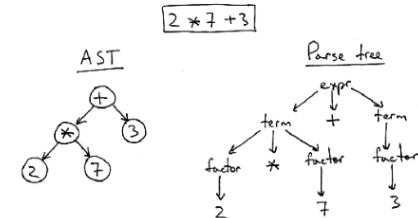
Sub-atomic particles



Everyday scenes

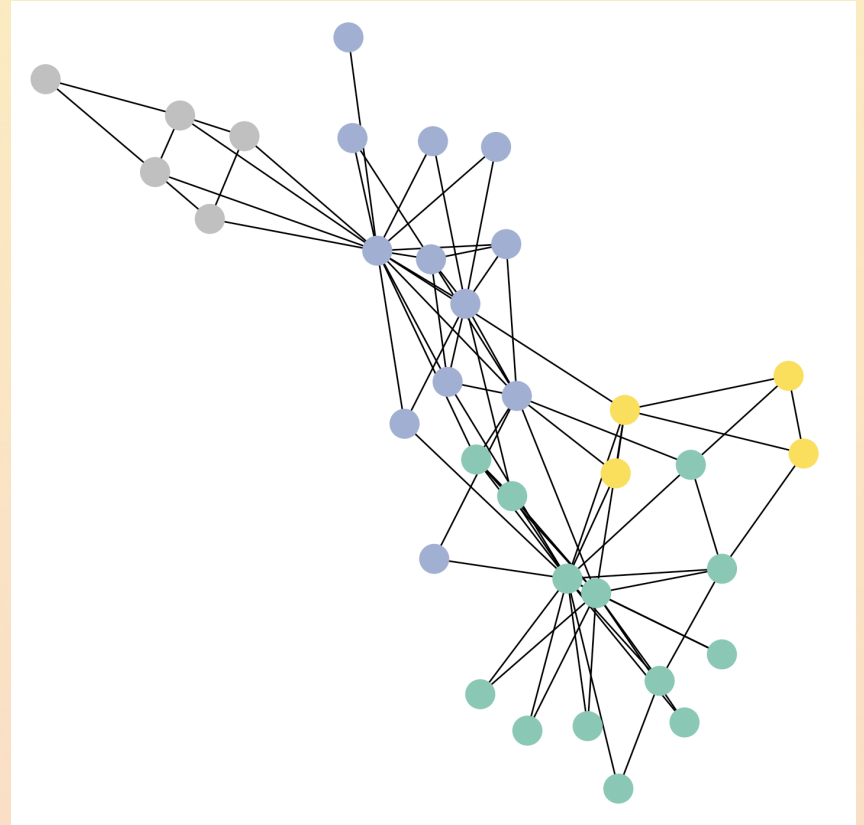


Code



Graph networks

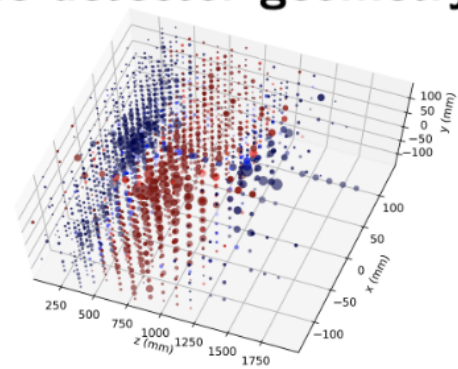
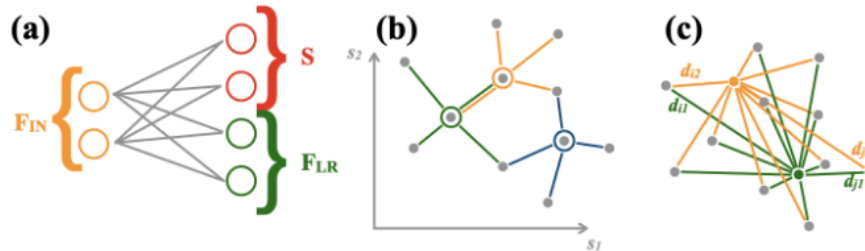
- Represent data as **point clouds**
- Connect data points with weight-dependent connections
- Train the network to find which weights are strongest
 - Learning the connectivity structure of the data



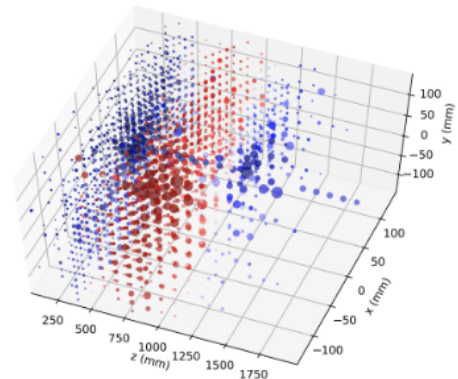
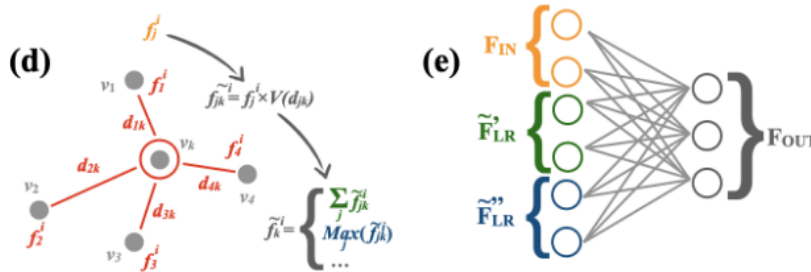
CMS High-granularity calorimeter

- 6 million cells with $\sim 3mm$ spatial resolution, over $600m^2$ of sensors
- Non-projective geometry

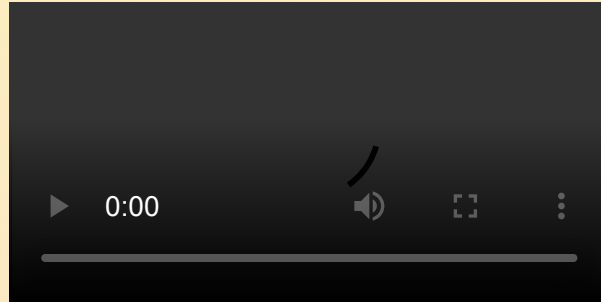
Learning representations of irregular particle-detector geometry with distance-weighted graph networks



(a) Truth



Graphs for water simulation



Plug the Physics into the AI

How are symmetries implemented?

- Data augmentation
Li, Dobriban '20

- Loss function penalties
- conserved quantities

- Architectural design

- Approximate symmetries (CNN)

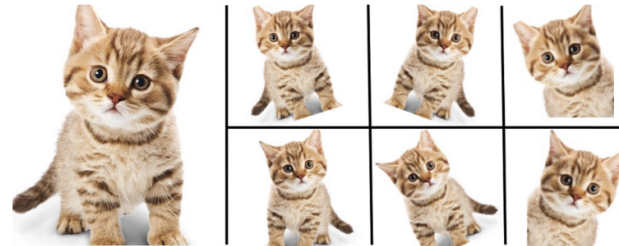
- Exact symmetries

- Weight sharing (message passing)
- Parameterization of symmetry preserving functions (group convolutions)

- Symmetries as constraints Finzi et al '21

- Irreducible representations Kondor, Thomas '18, Fuchs '20
Smidt

- Steerable CNNs Cohen '17, Welling..



Enlarge your Dataset

Credit: Bharath Raj

Cohen, Welling '19 Ravanbakhsh
Rose YU '21, '20. Weiler '21

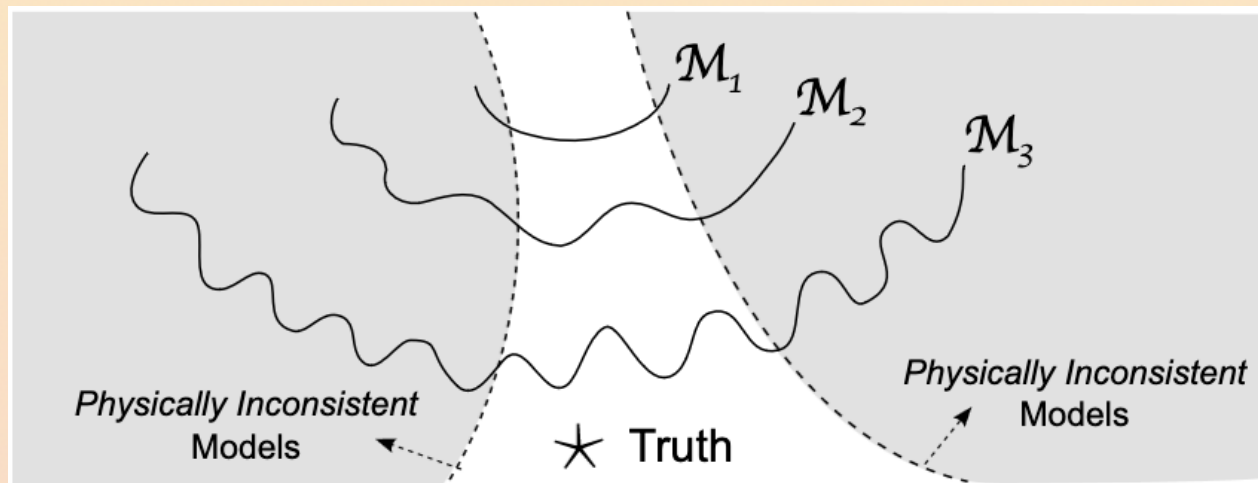
Kondor '18
Maron '18
Cohen '18

Plug the physics into the AI: constraints

$$\hat{y} = f(\mathbf{x}, \theta)$$

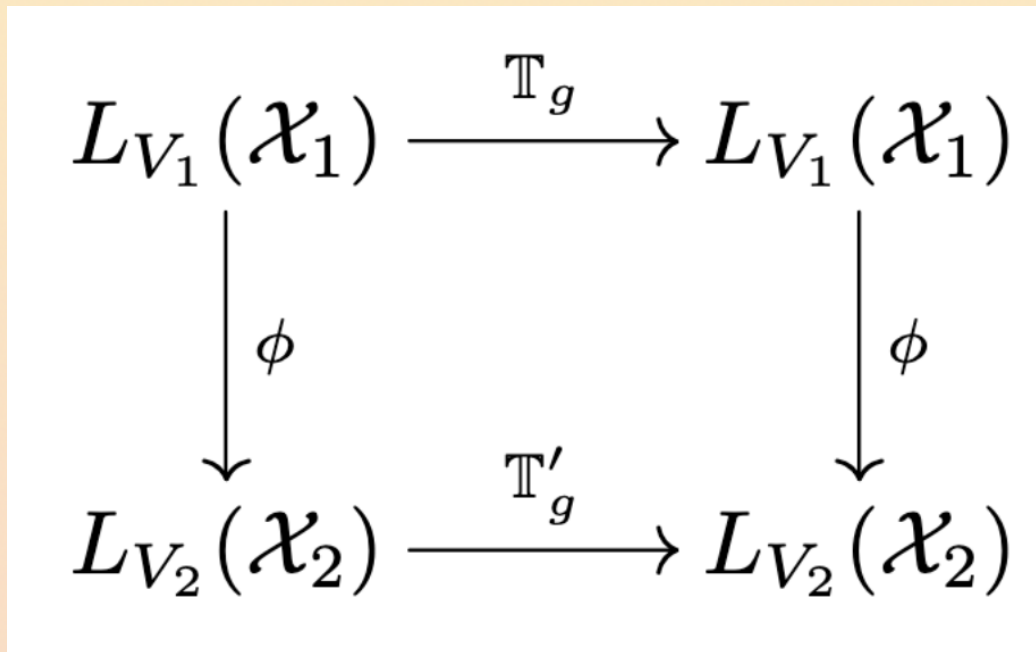
- Encode physics knowledge (e.g. inconsistency of models) inside the loss function as a penalty term

$$\mathbf{J}(\mathbf{w}) = Loss(y, \hat{y}) + \lambda \|\mathbf{w}\|_2^2 + \gamma \Omega(\hat{y}, \Phi)$$



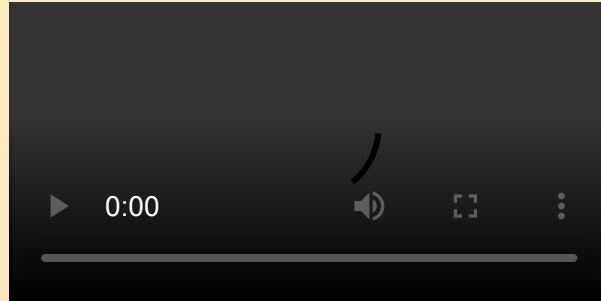
Plug the physics into the AI: network structure

- Equivariance under group transformation can e.g. enforced by convolutional layers
- Some implementations [available in pytorch](#)



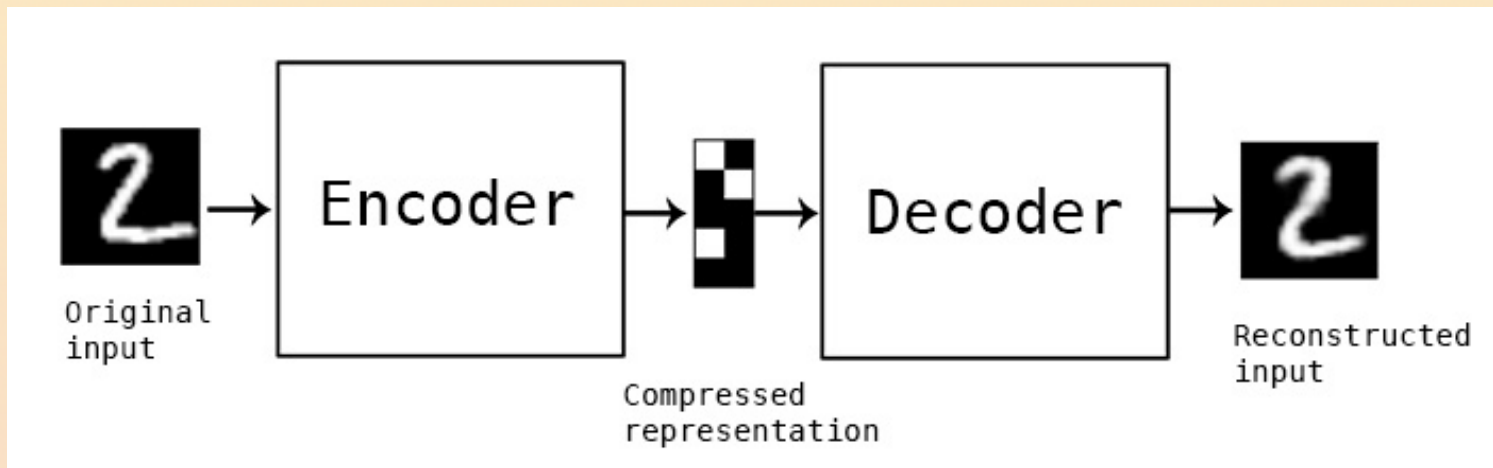
Plug the Physics into the AI

- Physics-aware differential equations solving

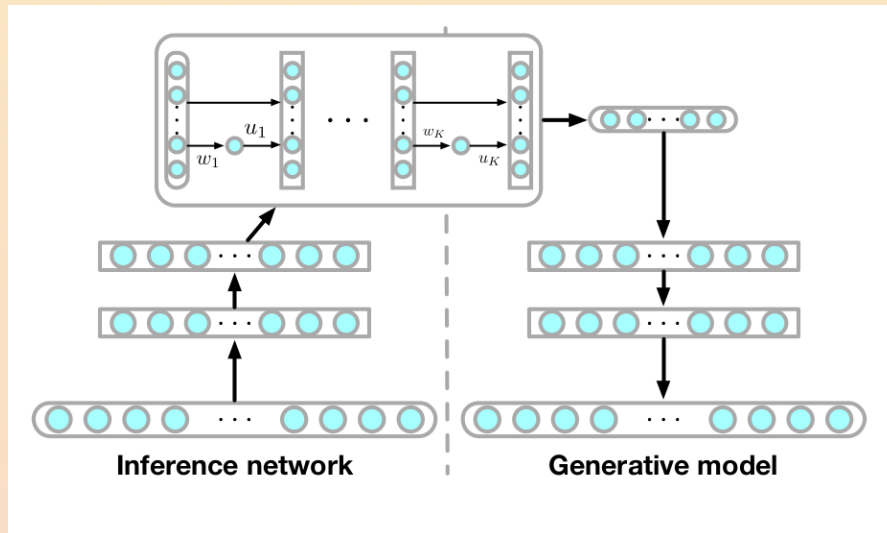
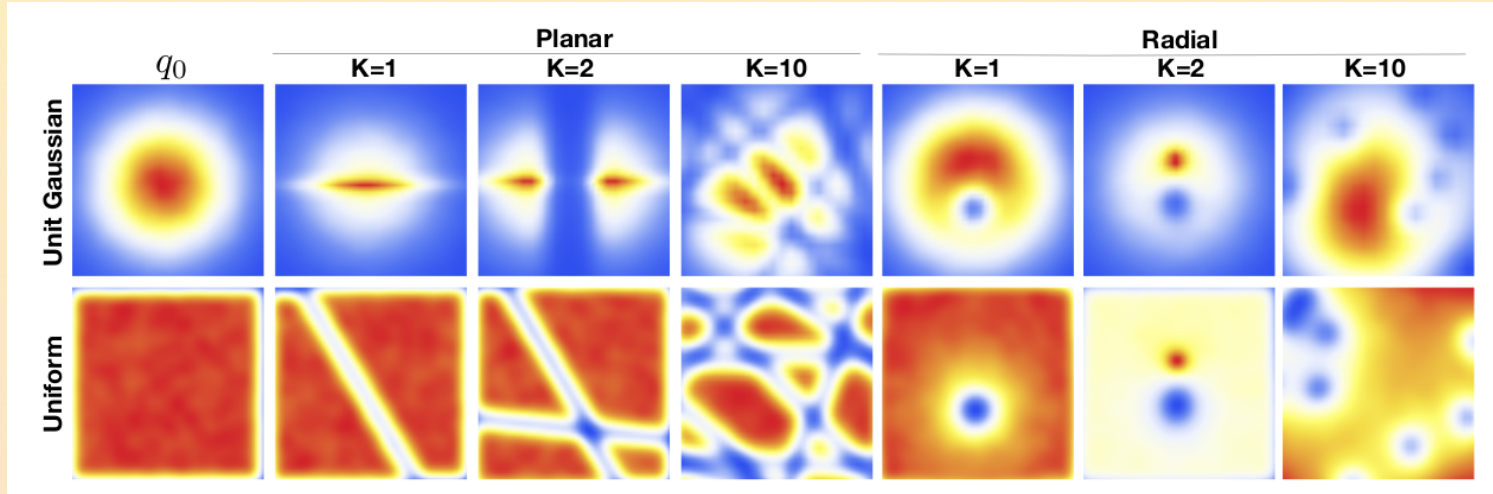


Autoencoders

- Learn the data itself passing by a lower-dimensional intermediate representations
 - Capture data generation features into a lower-dimensional space
- Can use for anomaly detection
- Can sample from the latent space to obtain random samples (generative AI)



Invertible networks



Solve inverse problems ("unfolding")

- Correct detector observation noise to recover source distribution

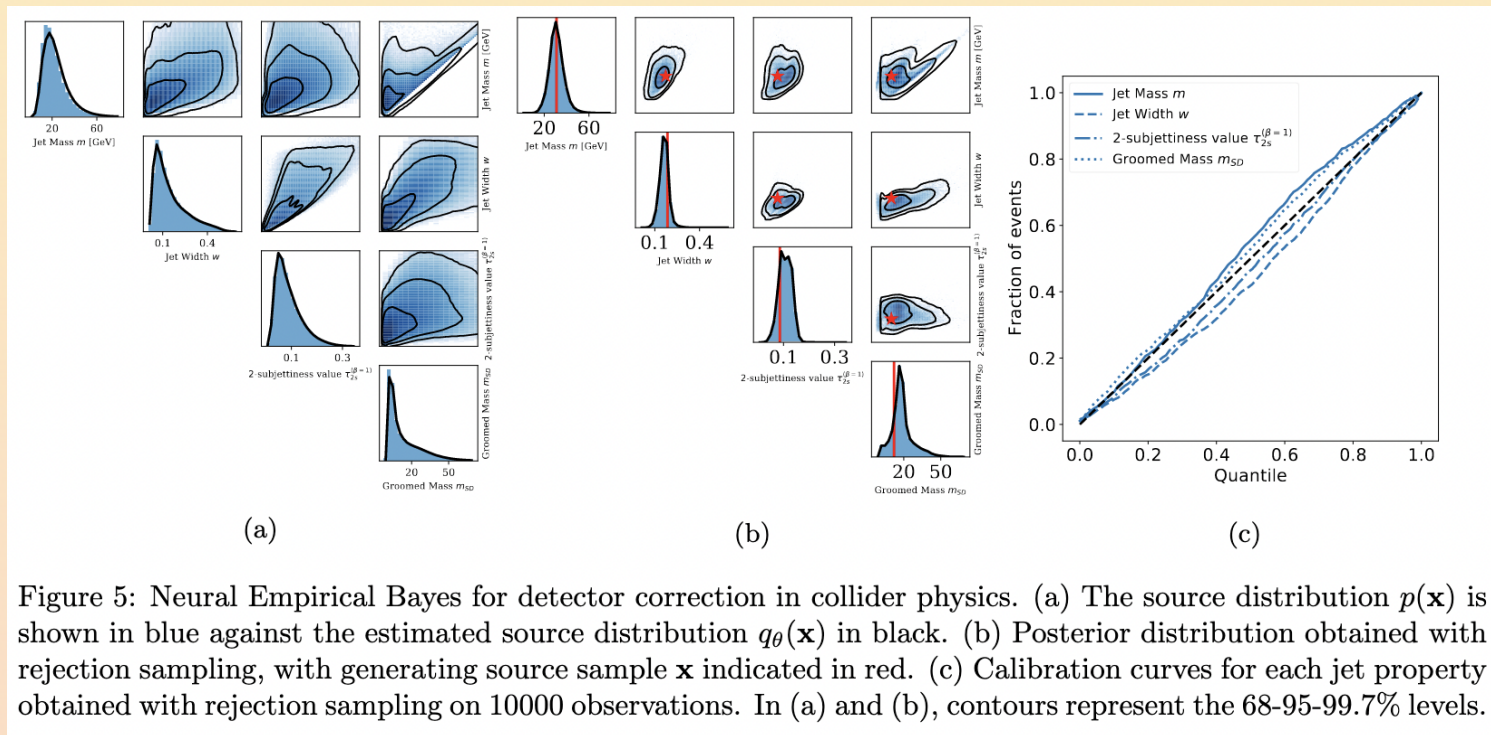
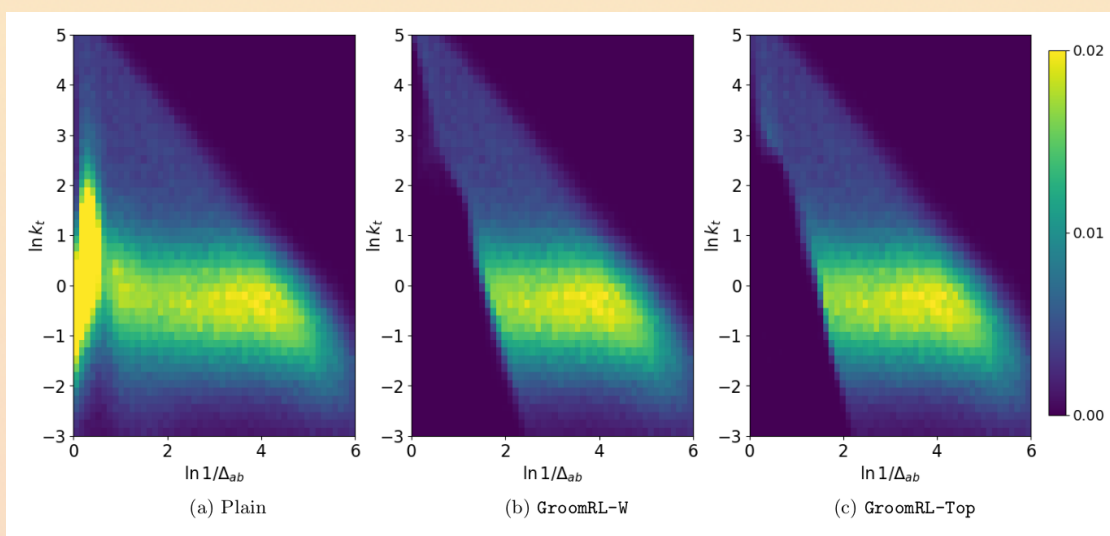
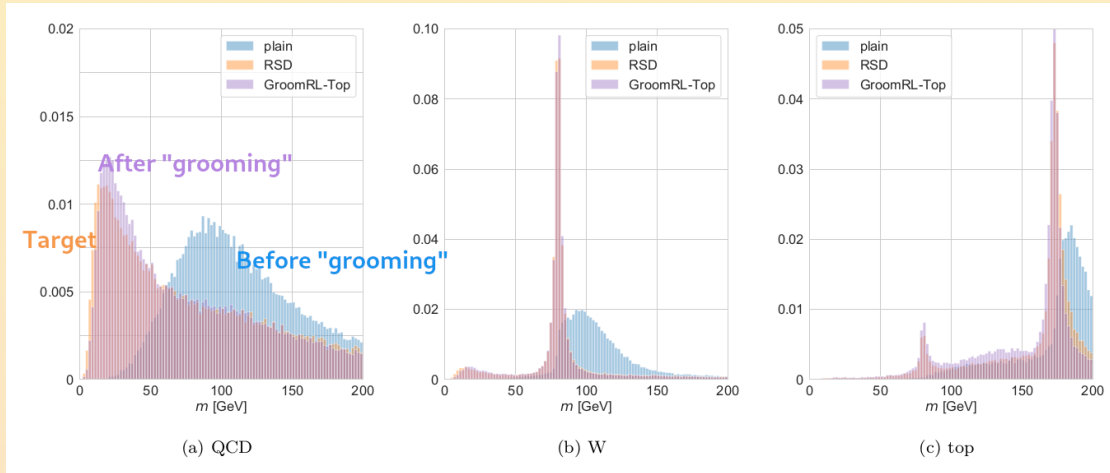


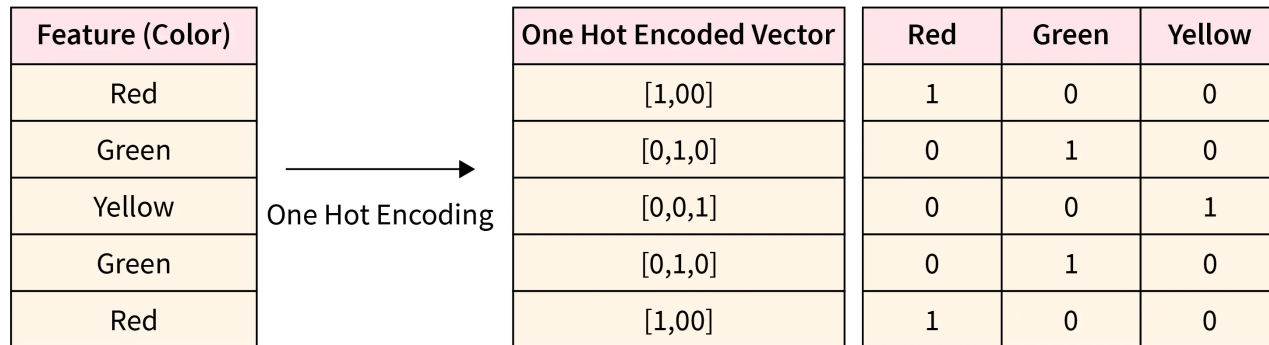
Figure 5: Neural Empirical Bayes for detector correction in collider physics. (a) The source distribution $p(\mathbf{x})$ is shown in blue against the estimated source distribution $q_{\theta}(\mathbf{x})$ in black. (b) Posterior distribution obtained with rejection sampling, with generating source sample \mathbf{x} indicated in red. (c) Calibration curves for each jet property obtained with rejection sampling on 10000 observations. In (a) and (b), contours represent the 68-95-99.7% levels.

Interpretability



Encode sequences

- **One-hot encoding** for unordered sequences
 - Works e.g. for text



SCALER
Topics

Encode sequences

- "Yellow" $[0, 0, 1]$ can be predicted as "0 for red and 0 for green"
 - One-hot-encoded features highly correlated ("multicollinearity")
- Dummy variable trap
 - Drop one of the "dimensions"

Feature (Color)		Red	Green
Red		1	0
Green		0	1
Yellow		0	0
Green		0	1
Red		1	0

One Hot Encoding

Red	Green
1	0
0	1
0	0
0	1
1	0

Yellow Column dropped to avoid the Dummy Variable Trap

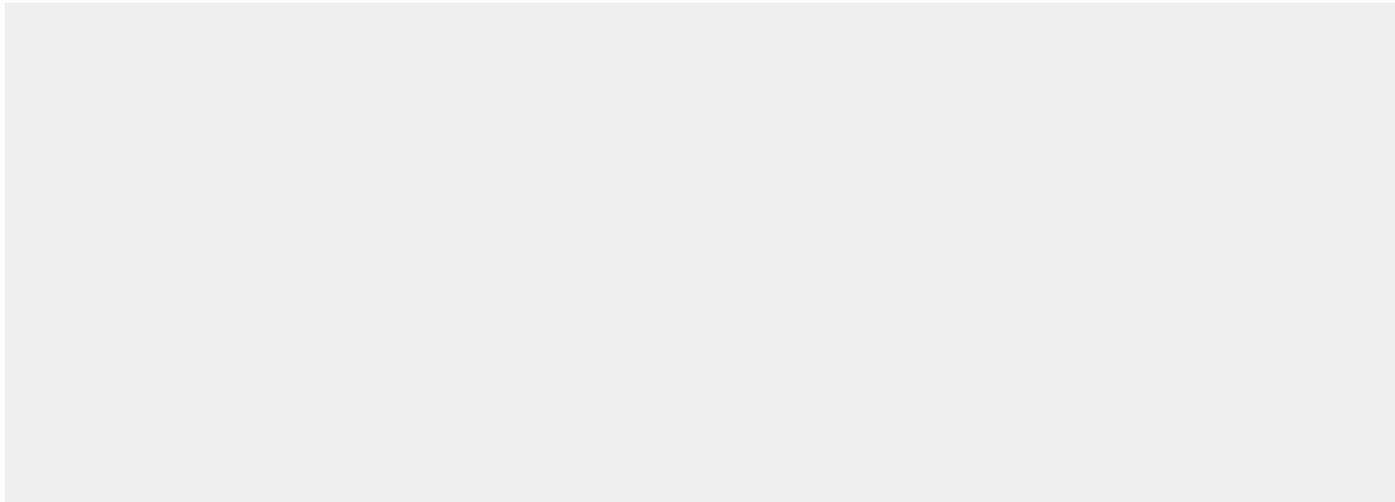
Self-attention

- Capture dependencies and relationships within inputs
 - Mostly in natural language processing and computer vision
- N inputs, N outputs
 - Allow inputs to interact with each other and find out which ones to pay attention to
 - Output is an aggregate of interactions and attention scores
- Useful for:
 - Long-range dependencies: understand complex patterns and dependencies
 - Contextual understanding: assign appropriate weights to important elements in the sequence
 - Parallel computation: can be computed in parallel → efficient and scalable for large datasets.

Self-attention

- Inputs (green) must be represented as: key (orange), query (red), value (purple)
 - Initially, by random reweighting of inputs themselves

Self-attention



input #1

1	0	1	0
---	---	---	---

input #2

0	2	0	2
---	---	---	---

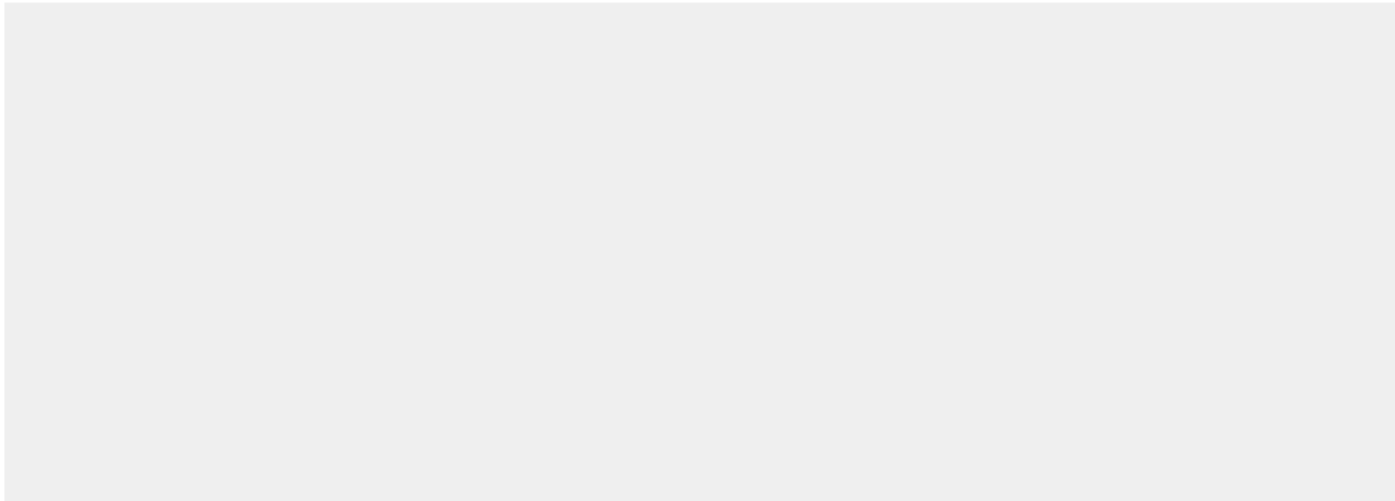
input #3

1	1	1	1
---	---	---	---

Self-attention

- Calculate attention score
 - Multiply (dot product) each query with all keys
 - For each query: N keys \rightarrow N attention scores

Self-attention



input #1
1 0 1 0

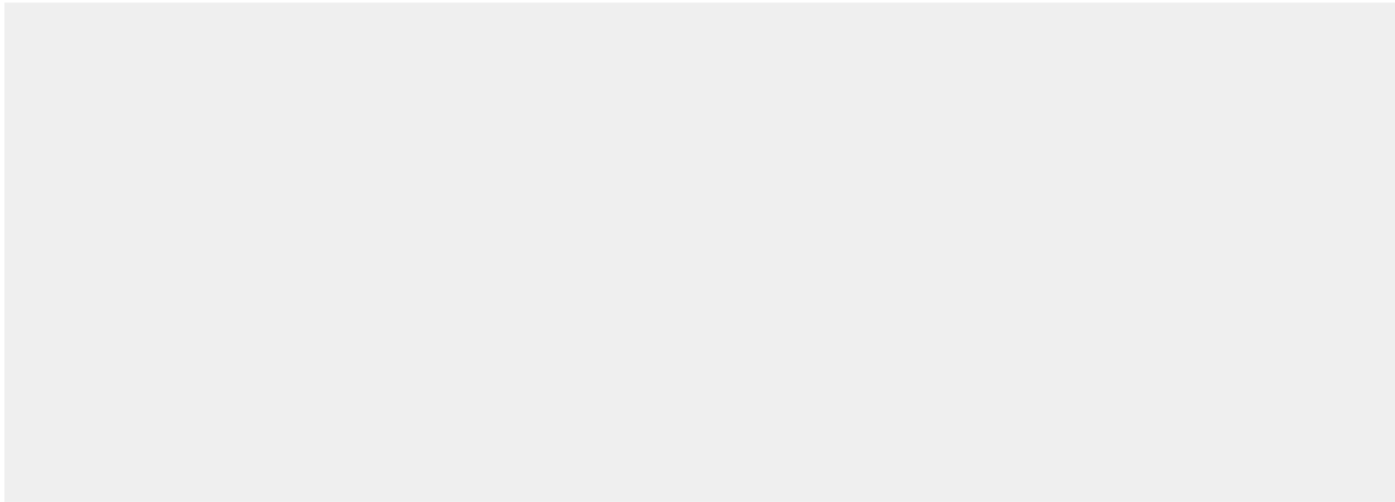
input #2
0 2 0 2

input #3
1 1 1 1

Self-attention

- Activation function (softmax) of attention scores

Self-attention



input #1

1	0	1	0
---	---	---	---

input #2

0	2	0	2
---	---	---	---

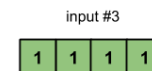
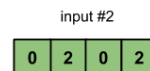
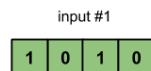
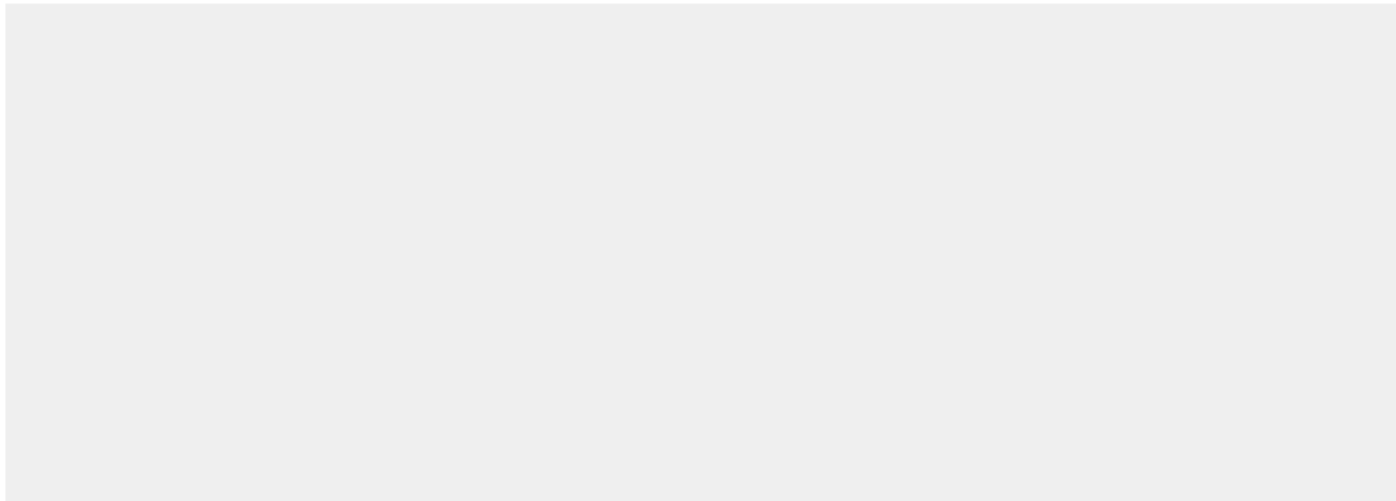
input #3

1	1	1	1
---	---	---	---

Self-attention

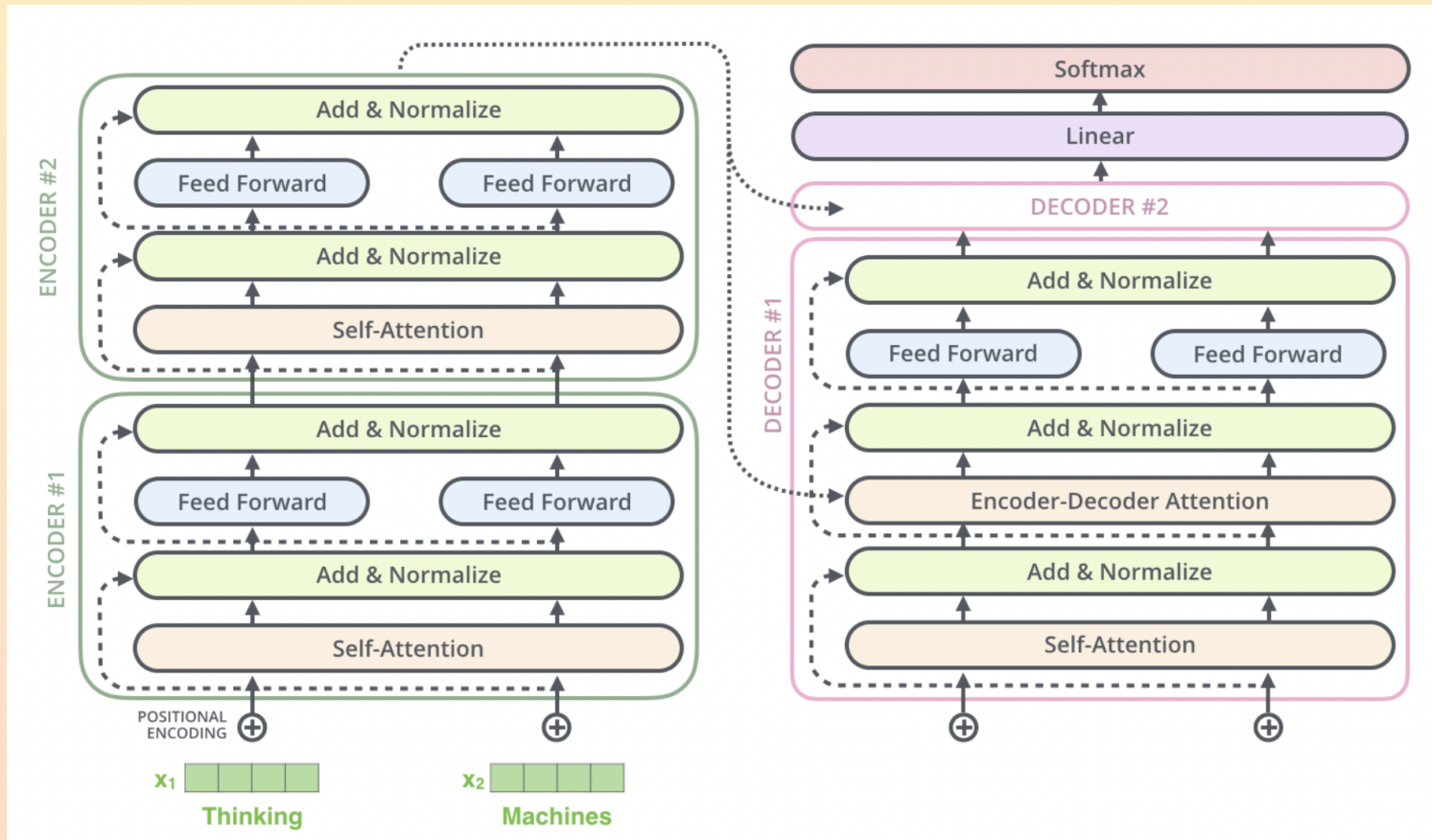
- Calculate alignment vectors (yellow), i.e. weighted values
 - Multiply each attention score (blue) by its value (purple)
- Sum alignment vectors to get input for output 1, repeat for 2 and 3

Self-attention

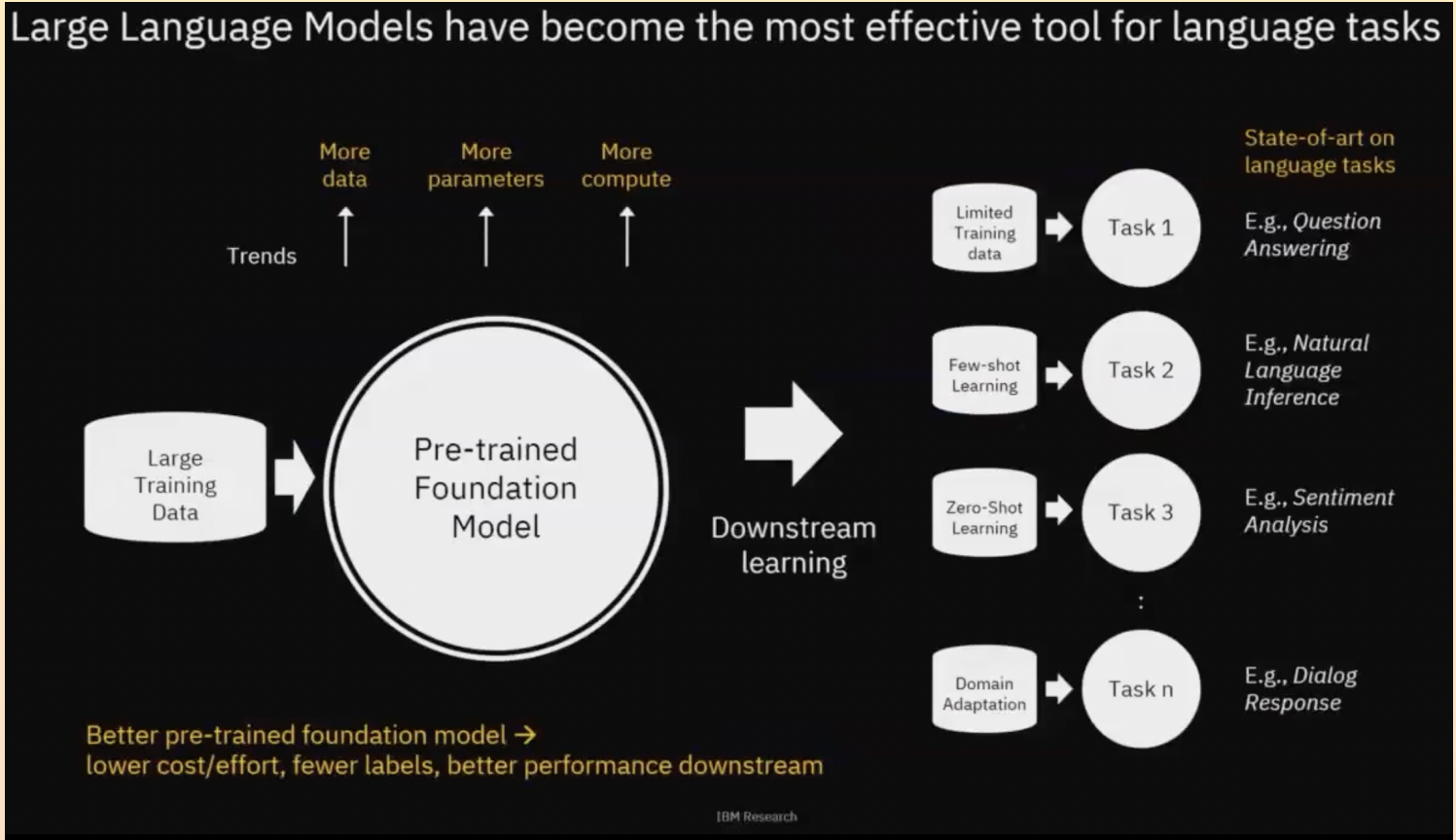


Transformers

- The engine behind GPT3

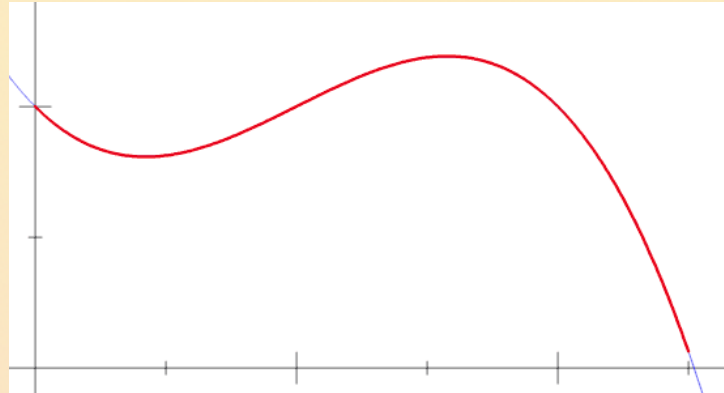


Foundation Models



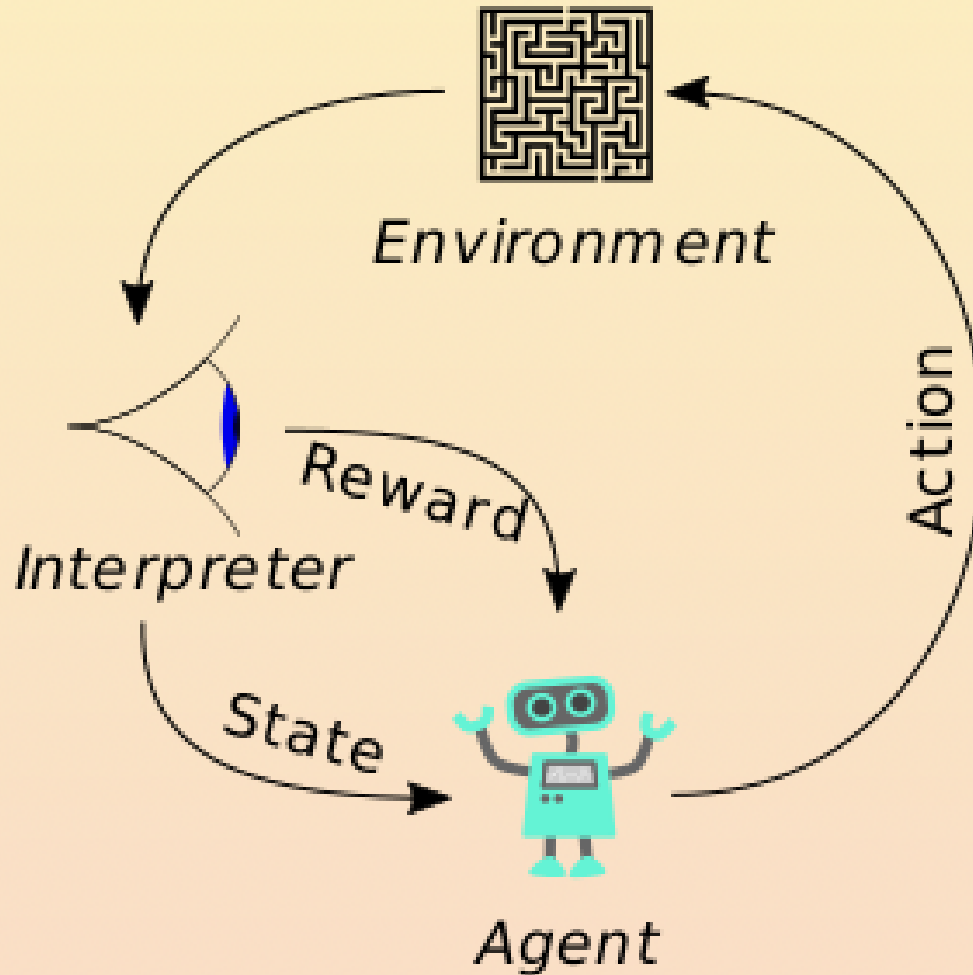
Translate Problems into Solutions

- Symbolic integration: find the analytic formula for the area of the curve

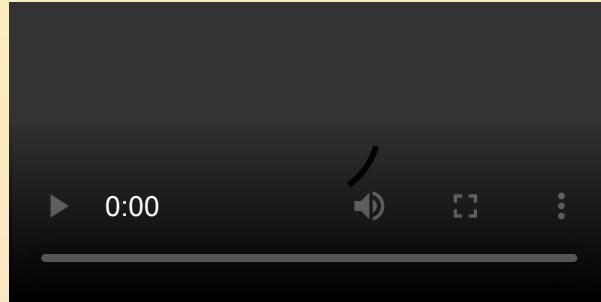


	Integration (BWD)	ODE (order 1)	ODE (order 2)
Mathematica (30s)	84.0	77.2	61.6
Matlab	65.2	-	-
Maple	67.4	-	-
Beam size 1	98.4	81.2	40.8
Beam size 10	99.6	94.0	73.2
Beam size 50	99.6	97.0	81.0

Reinforcement learning



From Videogames...

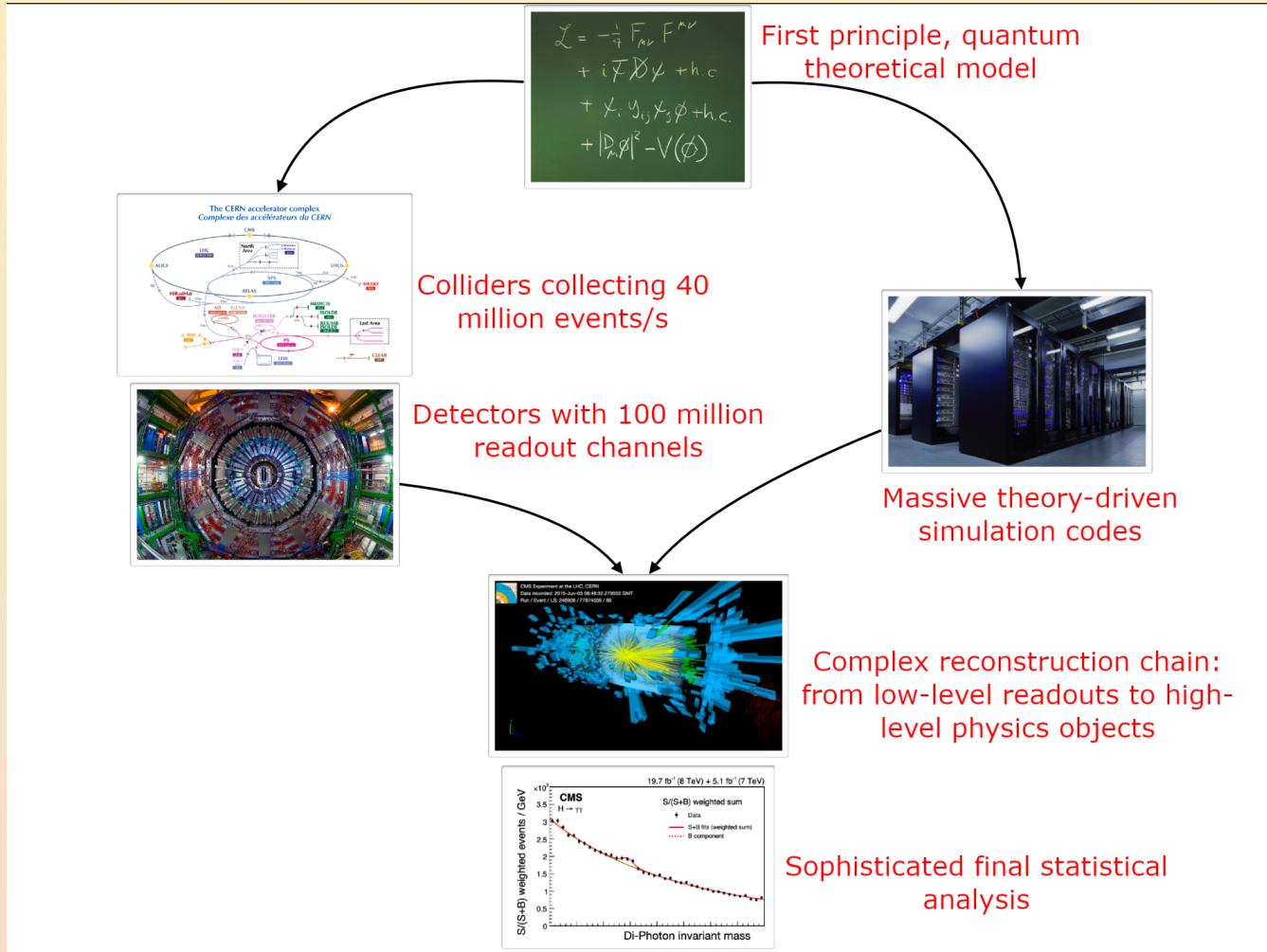


...to Physics

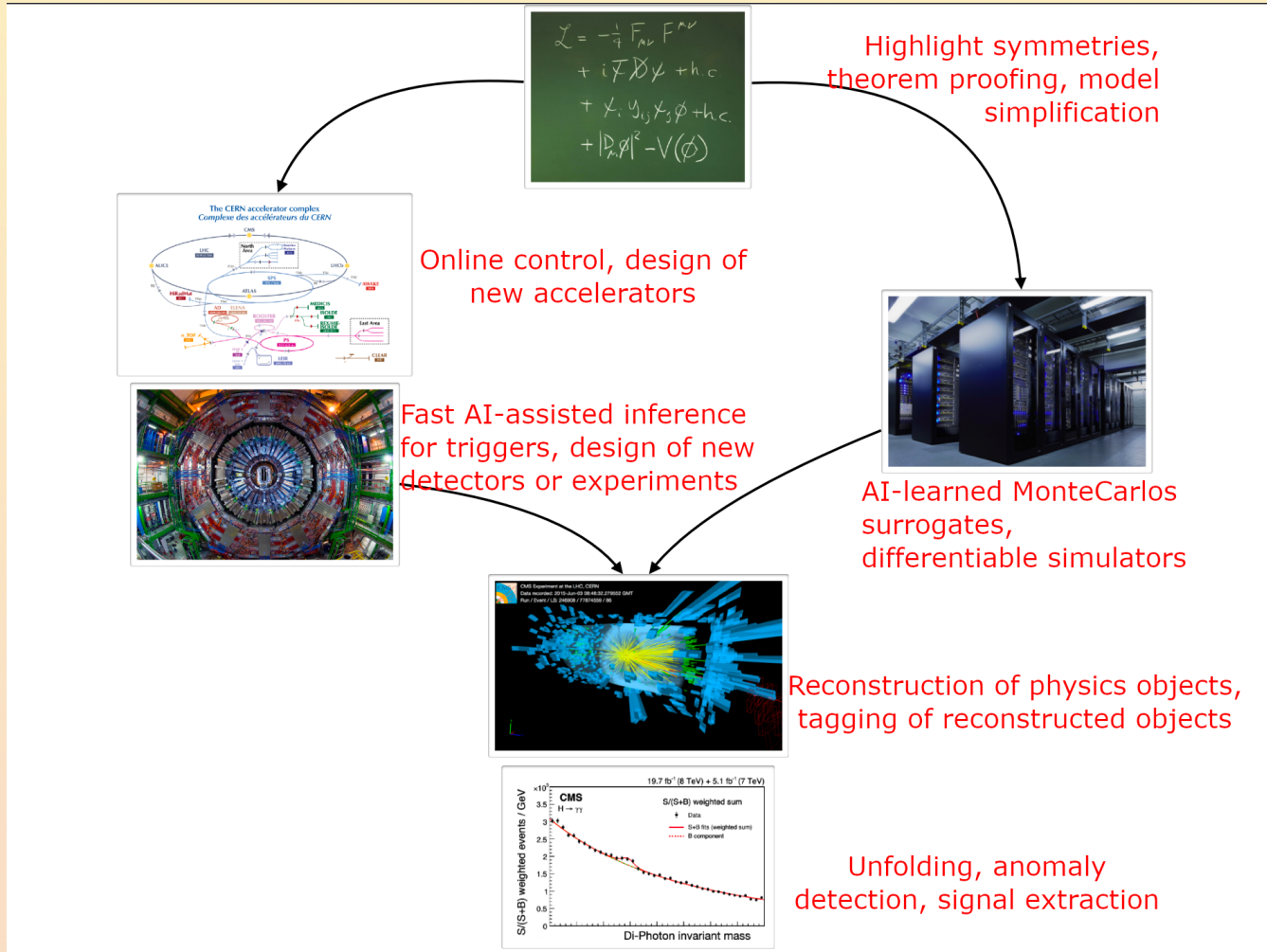
- Reward models consistent with the observed quark properties

charges	$\mathcal{Q} = \left(\begin{array}{c ccc ccc ccc c c} & Q_1 & Q_2 & Q_3 & u_1 & u_2 & u_3 & d_1 & d_2 & d_3 & H & \phi \\ \hline q & 6 & 4 & 3 & -2 & 2 & 4 & -3 & -1 & -1 & 1 & 1 \end{array} \right)$										
$\mathcal{O}(1)$ coeff.	$(a_{ij}) \simeq \begin{pmatrix} -1.975 & 1.284 & -1.219 \\ 1.875 & -1.802 & -0.639 \\ 0.592 & 1.772 & 0.982 \end{pmatrix} \quad (b_{ij}) \simeq \begin{pmatrix} -1.349 & 1.042 & 1.200 \\ 1.632 & 0.830 & -1.758 \\ -1.259 & -1.085 & 1.949 \end{pmatrix}$										
VEV, Value	$v_1 \simeq 0.224, \quad \mathcal{V}(\mathcal{Q}) \simeq -0.598$										
charges	$\mathcal{Q} = \left(\begin{array}{c ccc ccc ccc c c} & Q_1 & Q_2 & Q_3 & u_1 & u_2 & u_3 & d_1 & d_2 & d_3 & H & \phi \\ \hline 1 & 2 & 0 & -1 & -3 & 1 & -3 & -5 & -4 & 1 & 1 \end{array} \right)$										
$\mathcal{O}(1)$ coeff.	$(a_{ij}) \simeq \begin{pmatrix} -0.601 & 1.996 & 0.537 \\ -0.976 & -1.498 & -1.156 \\ 1.513 & 1.565 & 0.982 \end{pmatrix} \quad (b_{ij}) \simeq \begin{pmatrix} 0.740 & -1.581 & -1.664 \\ -1.199 & -1.383 & 0.542 \\ 0.968 & 0.679 & -1.153 \end{pmatrix}$										
VEV, value	$v_1 \simeq 0.158, \quad \mathcal{V}(\mathcal{Q}) \simeq -0.621$										

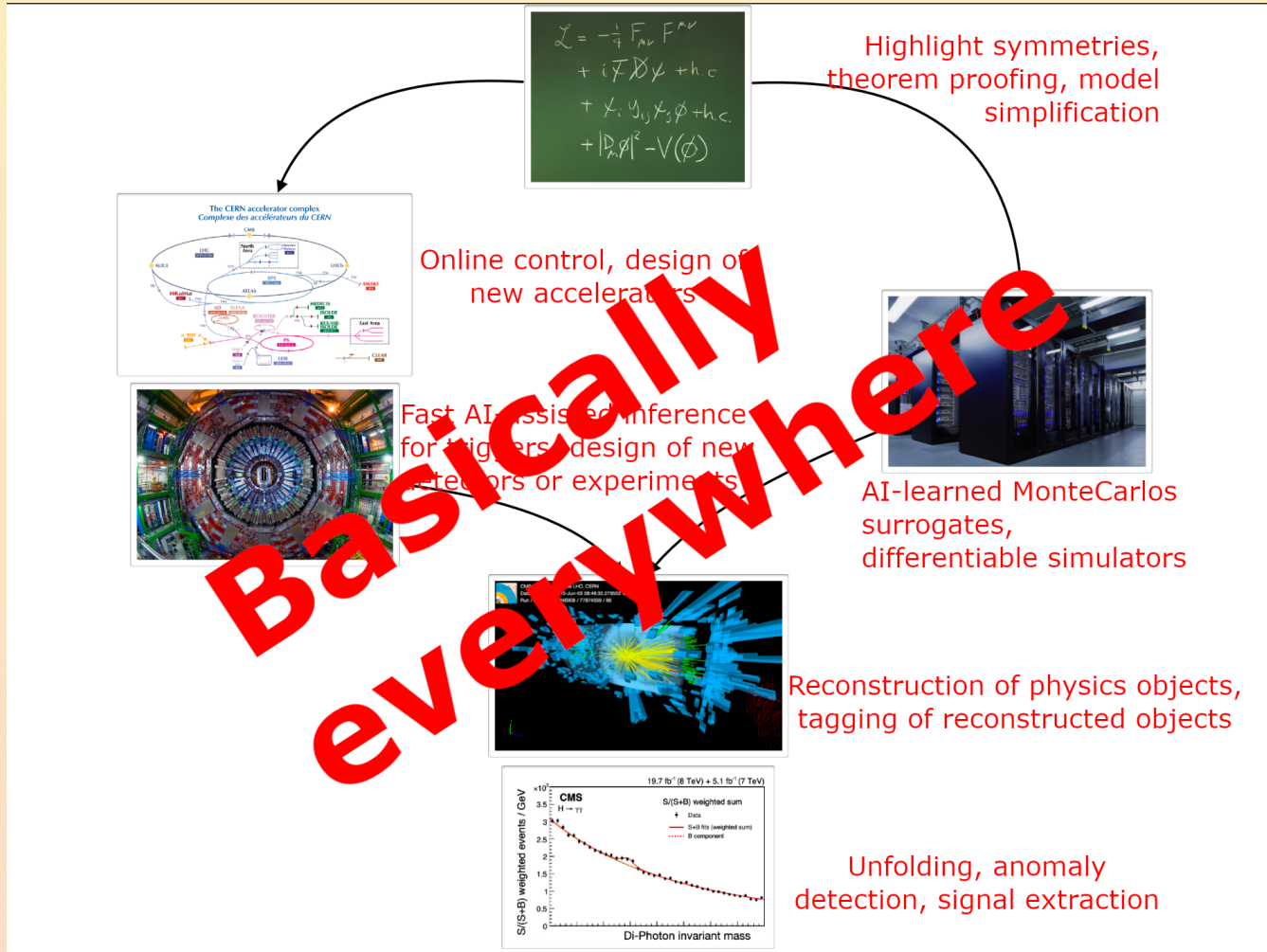
What do we do



Where we can plug AI




Where we can plug AI




Most theory papers are symbolic


- AI-assisted theorem proving



Lean
[de Moura et al., 2015]



Coq
[Barras et al., 1997]



Isabelle
[Nipkow et al., 2002]

<https://machine-learning-for-theorem-proving.github.io/> (NeurIPS 2023)

- LLMs to solve mathematical problems

Article | [Open access](#) | Published: 14 December 2023

Mathematical discoveries from program search with large language models

[Bernardino Romera-Paredes](#) , [Mohammadamin Barekatin](#), [Alexander Novikov](#), [Matej Balog](#), [M. Pawan](#)

- Simplify polylogarithms (no classical algorithm available, LLMs 91% success!)

Dutch:
naamsveranderingsdocumentenbriefgeheel

$$\begin{aligned}
 f(x) = & 9 \left(-\text{Li}_3(x) - \text{Li}_3\left(\frac{2ix}{-i + \sqrt{3}}\right) - \text{Li}_3\left(-\frac{2ix}{i + \sqrt{3}}\right) \right) \\
 & + 4 \left(-\text{Li}_3(x) + \text{Li}_3\left(\frac{x}{x+1}\right) + \text{Li}_3(x+1) - \text{Li}_2(-x) \ln(x+1) \right) \\
 & - 4 \left(\text{Li}_2(x+1) \ln(x+1) + \frac{1}{6} \ln^3(x+1) + \frac{1}{2} \ln(-x) \ln^2(x+1) \right)
 \end{aligned}$$

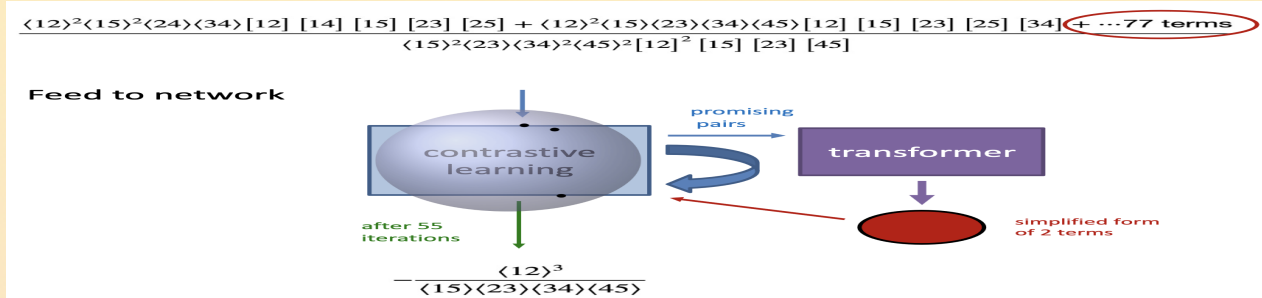
translate

English: dossier

$$f(x) = -\text{Li}_3(x^3) - \text{Li}_3(x^2) + 4\zeta_3$$

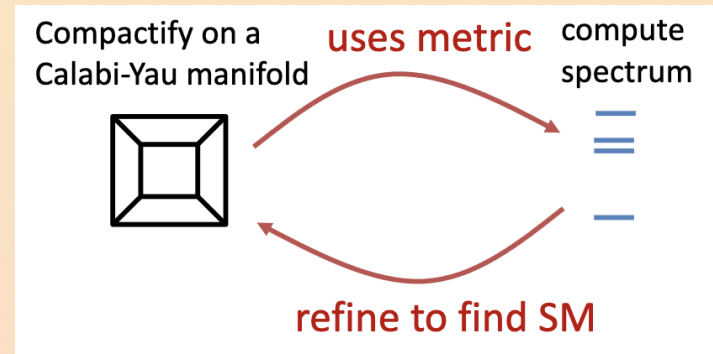
Most theory papers are symbolic

- 5-point MHV amplitude w/ Feynman diagrams: from 1990 tokens to 27 tokens



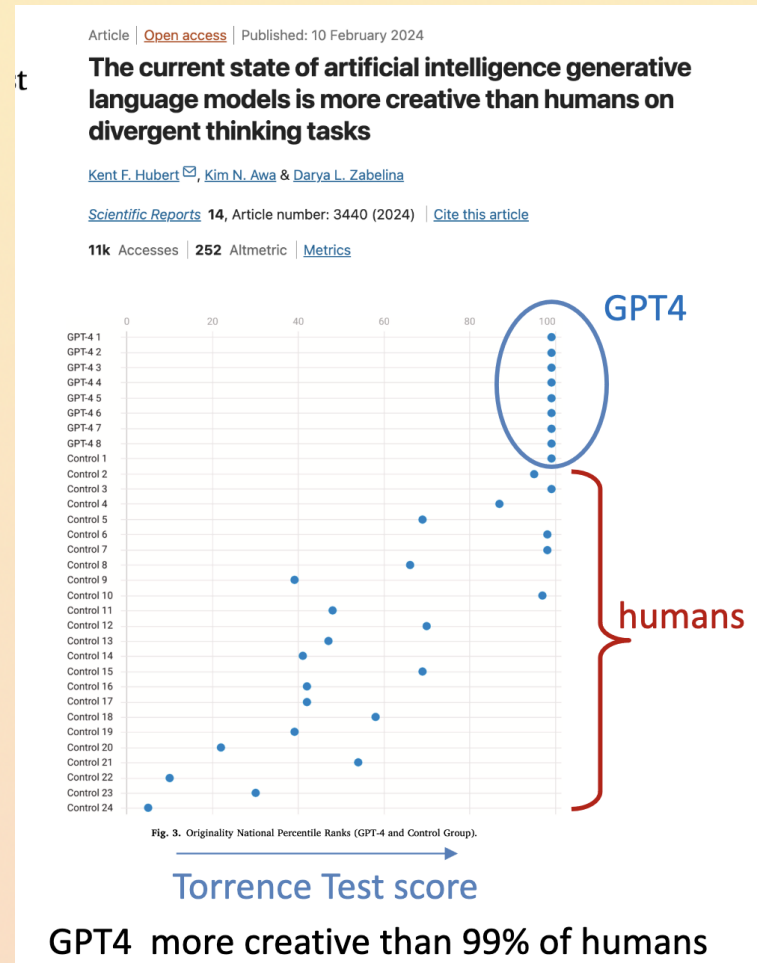
Solve string theory 🤪

- Find nontrivial Calabi-Yau metrics (1910.08605)
- Look for fixed points of metric flows (2310.19870)
- Predict rank of gauge group (1707.00655, prediction later proven)



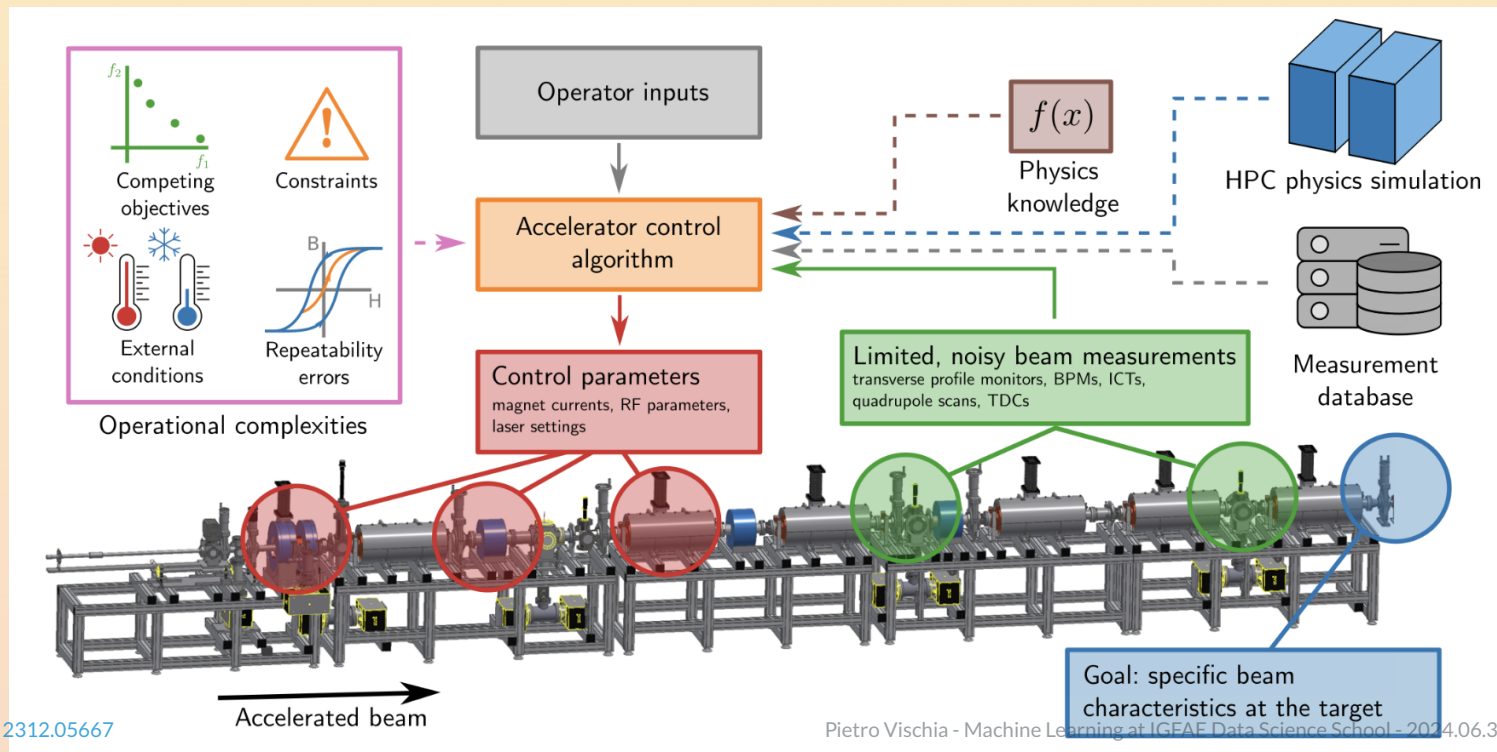
Beyond symbolic manipulation

- Can AI find interesting questions?
- Can AI models teach themselves to be good physicists using data?
- **If AI understands physics (can calculate everything) but we do not, do we consider it an acceptable "understanding"?**



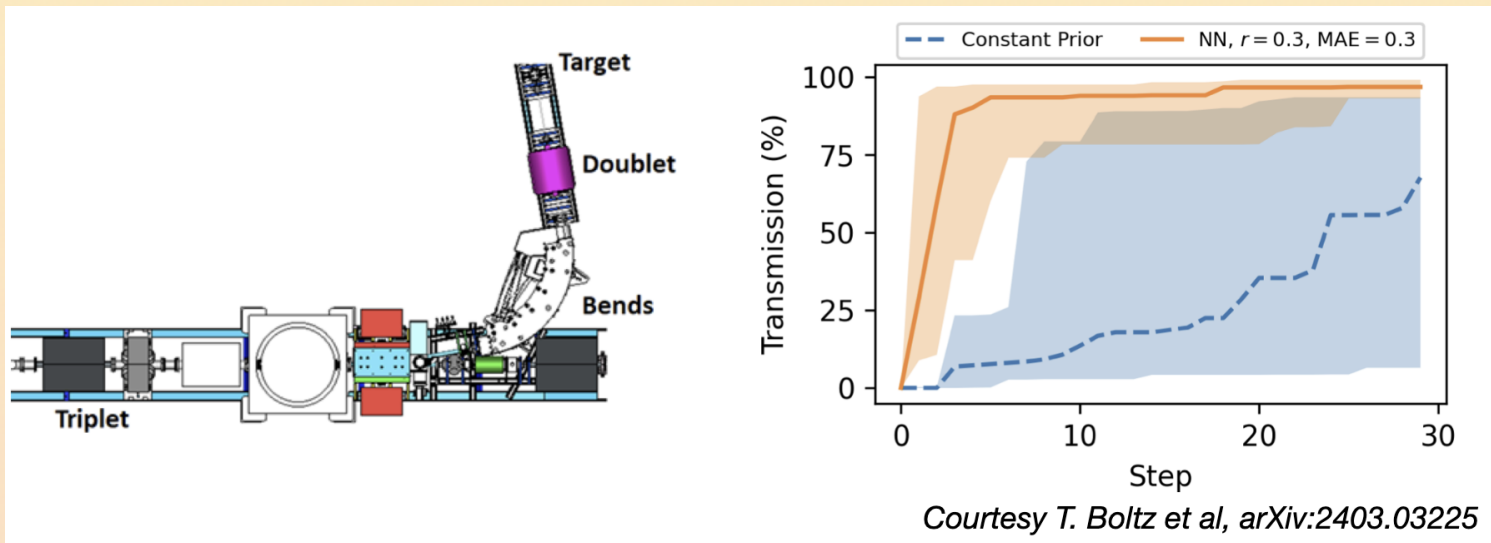
Accelerate accelerators

- Daily operation and control have huge impact on resources and efficiency
 - Beam scheduling: changing supercycle requires 20-100 clicks (2-25min) about 60 times/day
 - 15% of the yearly cost of SPS fixed target cycle employed for "waste" cycles to mitigate hysteresis problems
- What if we could make them fully automatic (like e.g. Space telescopes)?



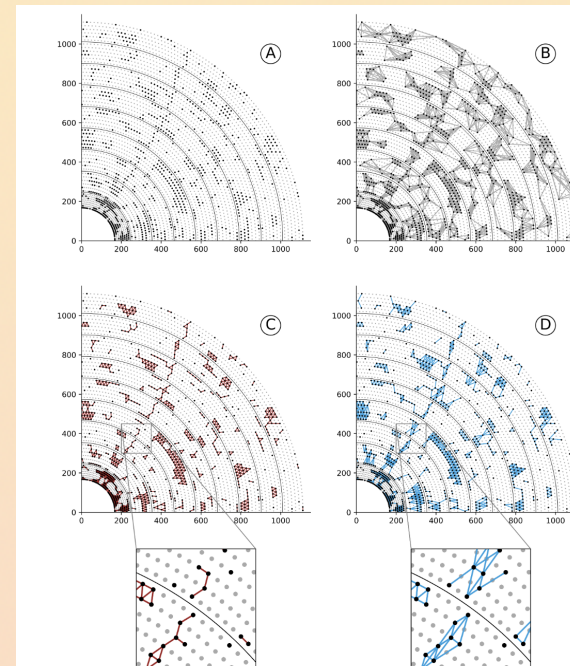
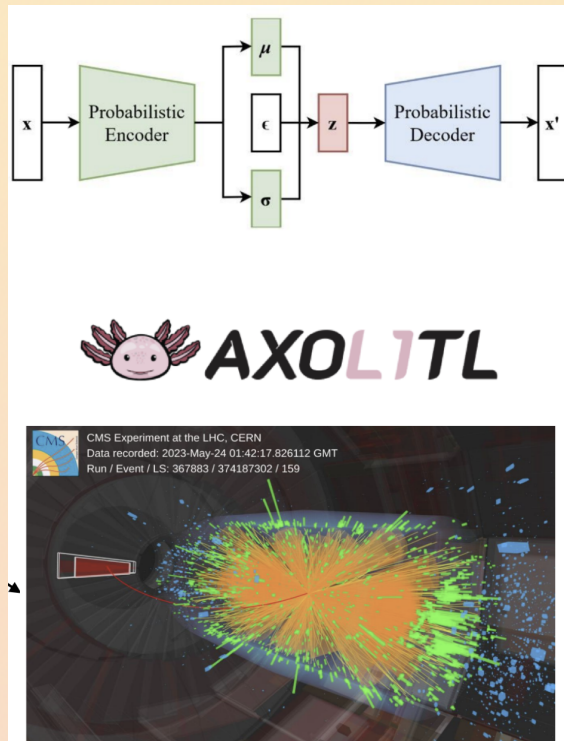
Accelerate accelerators

- Hierarchical, AI-controlled autonomous systems
- Optimize transmission to target in a system with 5 DoF, using Bayesian Optimization



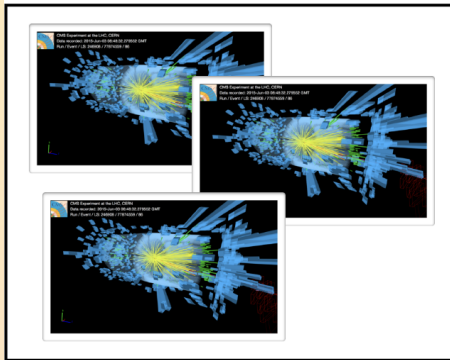
Trigger

- See talks by A. Zabi and S. Folgueras
- Pack AI models into the L1 trigger → improve selection criteria
 - At ICTEA!
- Can do e.g. anomaly detection, and online graph building

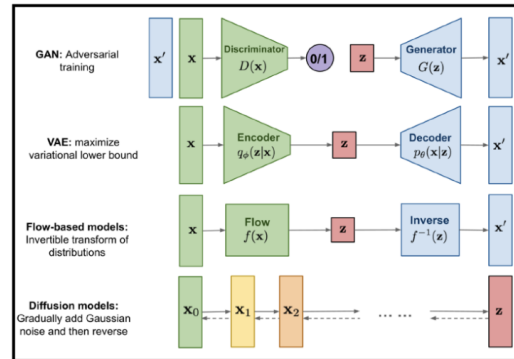


Simulation: two solutions

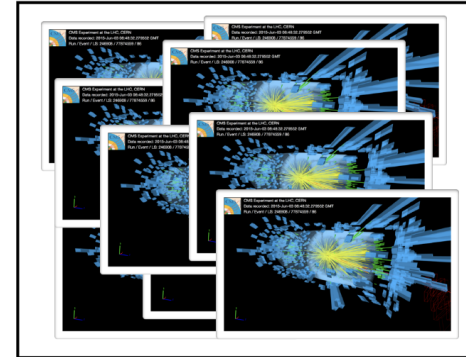
1. Use classical simulation or collider data as input



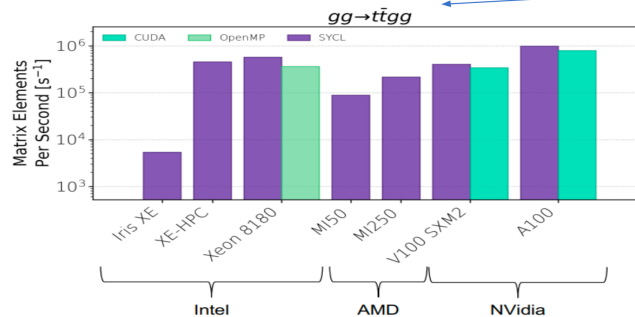
2. Train generative surrogate



3. Oversample



- Very recently, Madgraph5_aMC@NLO authors deployed a version of their code that can run on GPUs.
- This version significantly improves computation times (see [this talk](#)).



Talk by C. Vico Villalba

So our idea is: can we do this on hardware based accelerators?

- **FPGAs** are:
 - **Highly** parallelizable
 - In some cases not as fast as GPU.
 - But less power consuming.
 - Hardware based! really versatile.

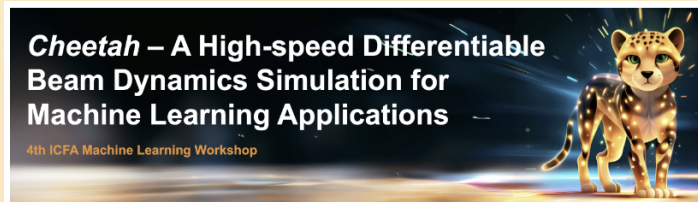


See talk by P. Leguina.

Simulation: long term solution

- Make everything differentiable, exploiting **differentiable programming**

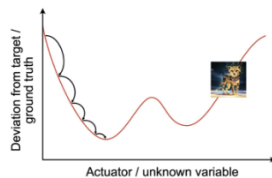
Cheetah



Gradient-based Tuning

Transverse beam tuning at ARES

- Tune magnet settings or lattice parameters using the **gradient of the beam dynamics model** computed through **automatic differentiation**.
- Seamless **integration with PyTorch** tools tuning neural networks.
- Becomes very useful for **high-dimensional tuning tasks** (see neural network training).



```
ares_ea.AREA0ZM1.k1 = nn.Parameter(0.0)
ares_ea.AREA0ZM2.k1 = nn.Parameter(0.0)
ares_ea.AREA0ZM1.angle = nn.Parameter(0.0)
ares_ea.AREA0ZM3.k1 = nn.Parameter(0.0)
ares_ea.AREA0M1.angle = nn.Parameter(0.0)

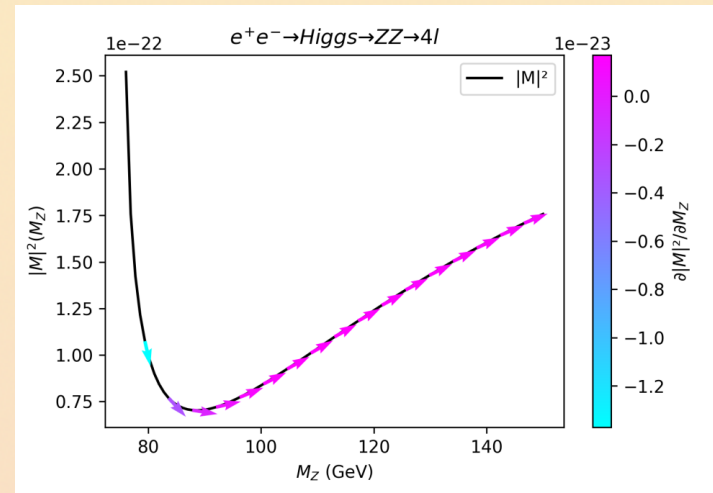
optimizer = Adam(ares_ea.parameters())

for _ in range(42):
    outgoing = ares_ea.track(incoming)
    loss = loss_fn(outgoing)

    loss.backward()
    optimizer.step()
    optimizer.zero_grad()
```

MadJax

(differentiable matrix element computation)



Towards full differentiability

- Matrix Element: differentiable ([MadJax](#))
- Integration in Madgraph: multi-channel integration speed up ([MadNIS](#))
- Parton shower: mostly differentiable ([2208.02274](#), and recent work by Kagan+Heinrich)
- Detector simulation: GATE/GEANT4 numerically differentiable (in small ranges) ([2202.05551](#))
- Operator overloading for GEANT4 ([2405.07944](#))

AD in GEANT4

- Probabilistic generator $f : \Theta \times \Omega \rightarrow Y$, $(\theta, \omega) \mapsto y$
- Find optimal inputs θ to maximize output y
 - Need $\frac{\partial y}{\partial \theta}$, pathwise derivative $\frac{\partial}{\partial \theta} f(\theta, \omega)$
 - **Derivgrind**: insert AD logic into the program (a sort of debugger): cannot support tricky cases
 - **CoDiPack**: operator overloading (e.g. replace `double` type): can run out of memory when storing the real-arithmetic evaluation graph (tape)
 - **Clad**: compiler-based source transformation tools: could use smaller tapes, more advanced optimization

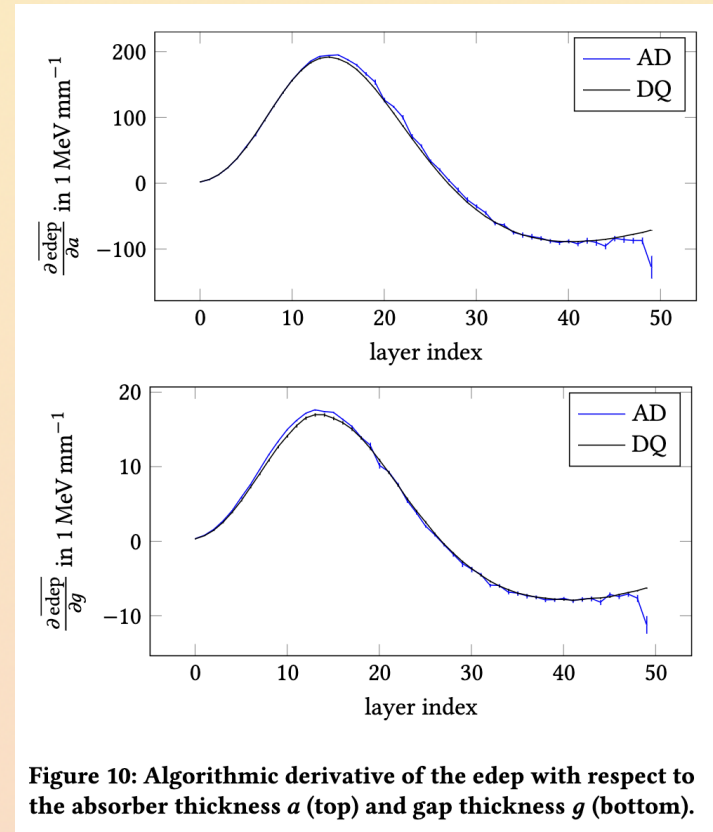
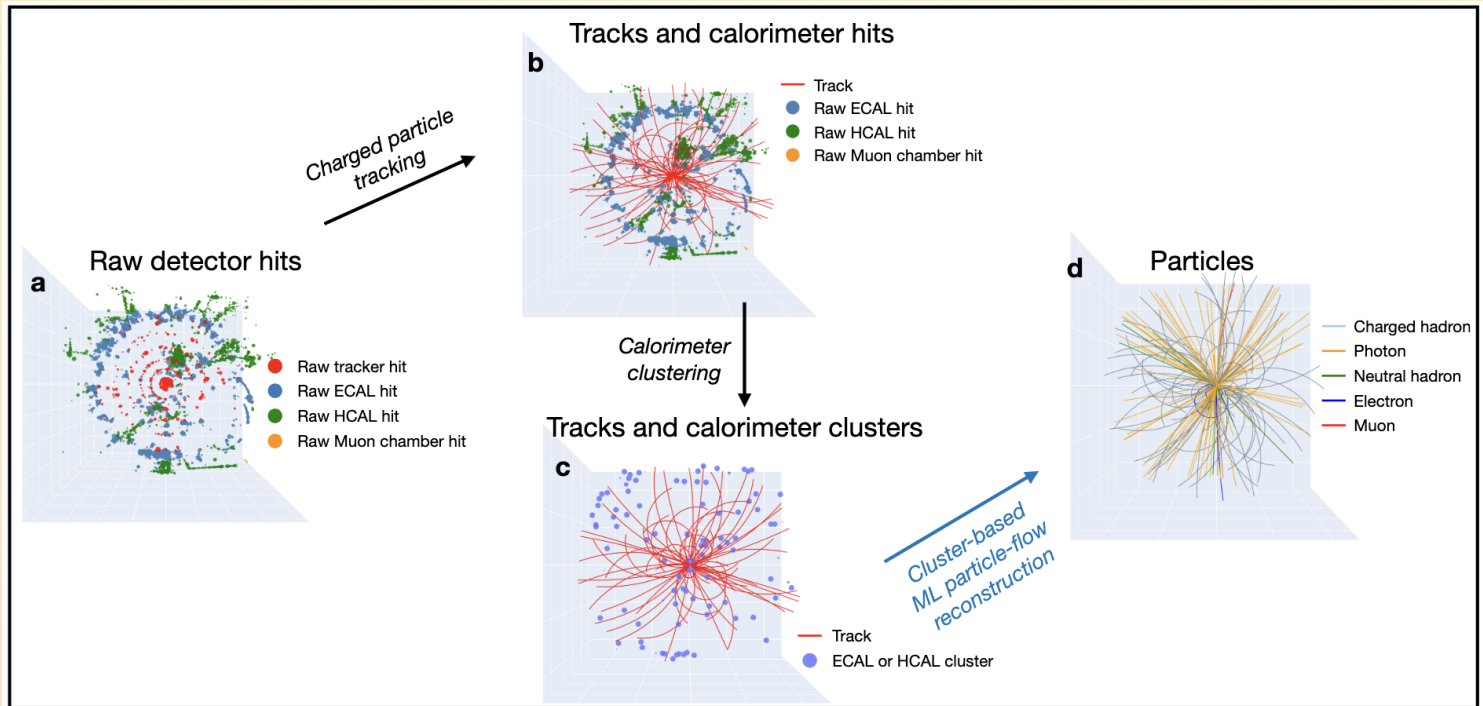


Figure 10: Algorithmic derivative of the edep with respect to the absorber thickness a (top) and gap thickness g (bottom).

Reconstruction...

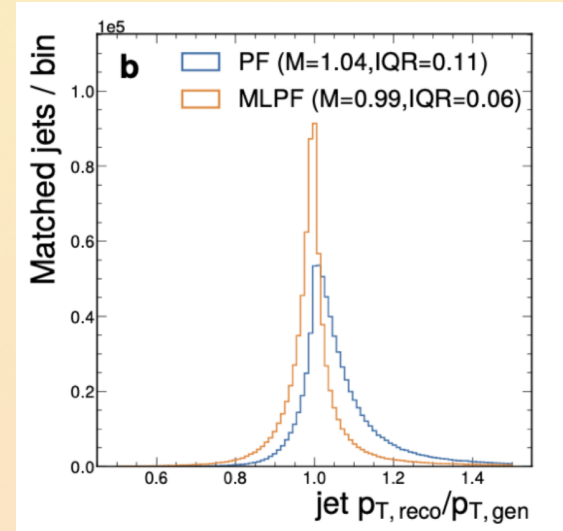
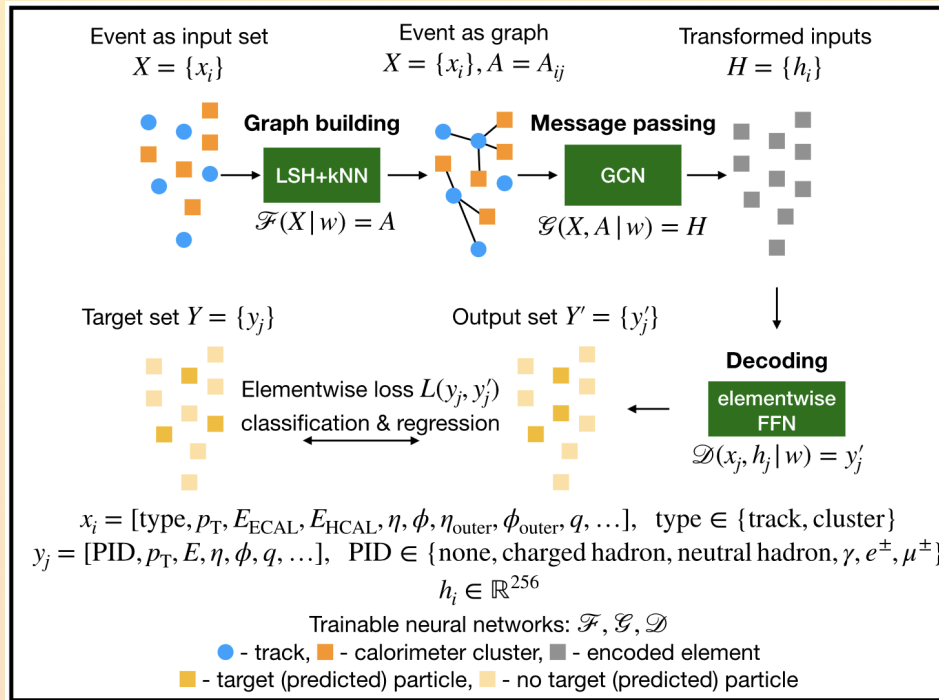


Reconstruction maps **low-level** detector read-outs

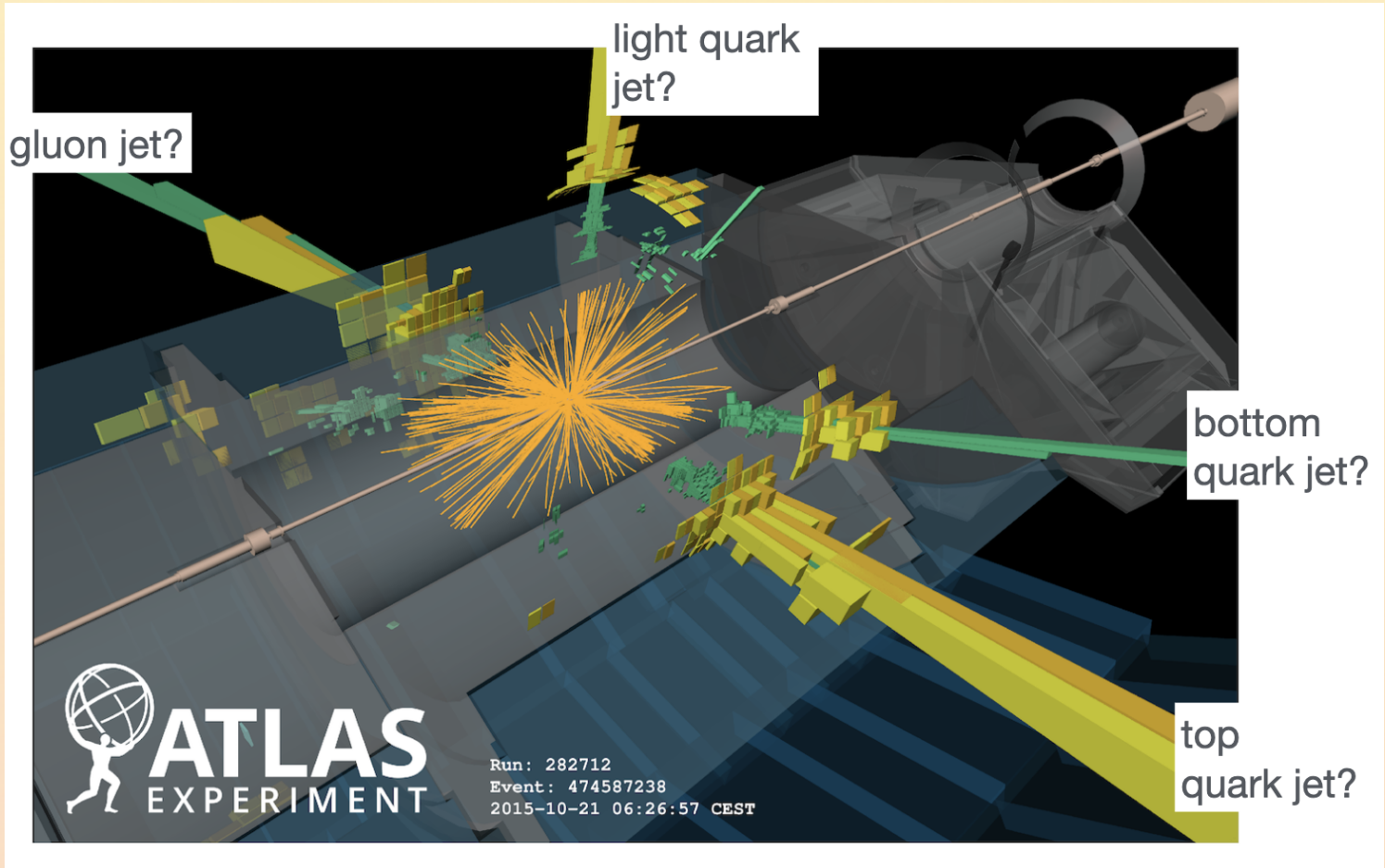
to **physical particles**

(Which in-turn are the basis of higher-level interpretation)

...with AI

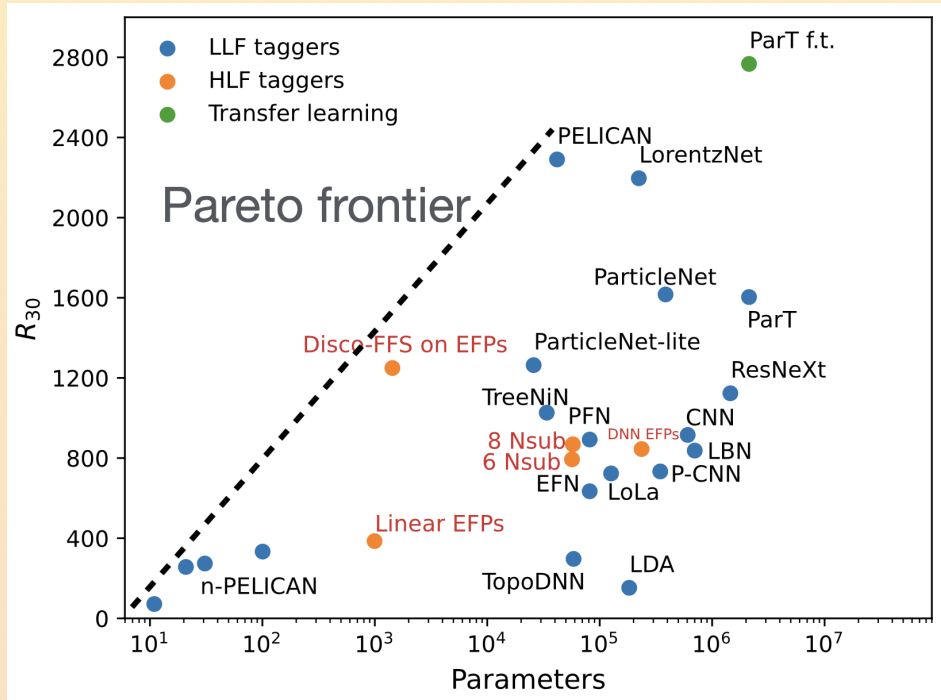


Identification...

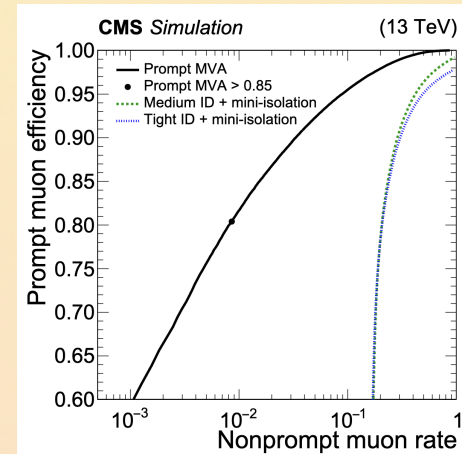


...with AI

Vast landscape of taggers

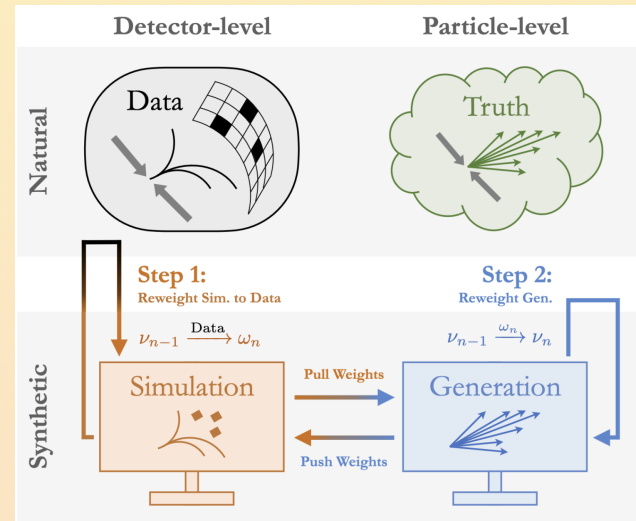


CMS Muon ID: made in ICTEA!

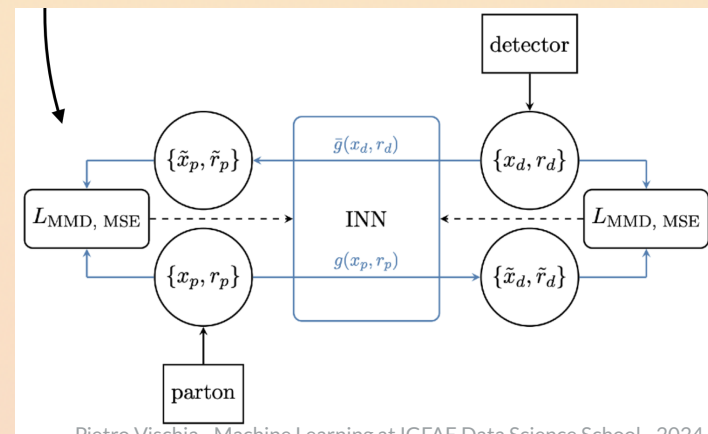


Inference: unfolding

- Use classifiers to learn appropriate weights



- Morph distributions one into the other using diffusion models



Inference: anomaly detection

Gaussian processes

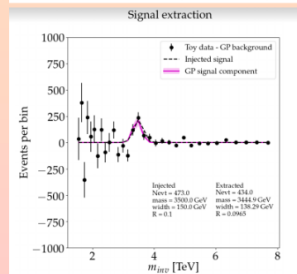
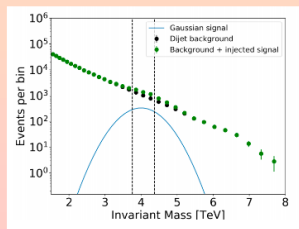
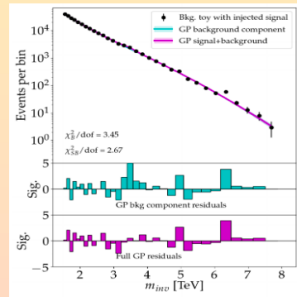
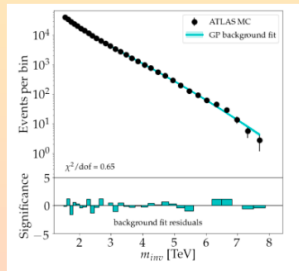
- Multivariate gaussian associated to a set of random variables ($N_{dim} = N_{random\ variables}$)
 - Kernel as a similarity measure between bin centers (counts) and a averaging function

$$\mu(x) = 0, \quad (9)$$

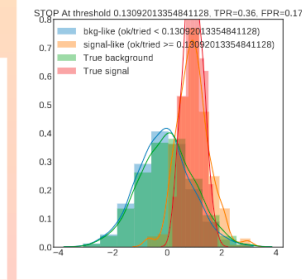
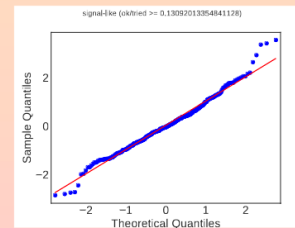
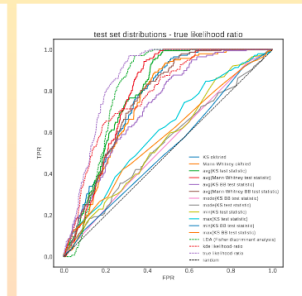
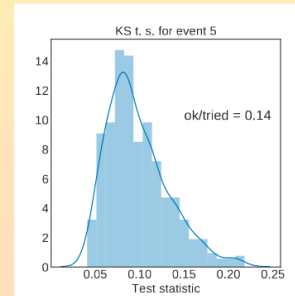
$$\Sigma_B(x, x') = A \exp\left(\frac{d - (x + x')}{2a}\right) \sqrt{\frac{2l(x)l(x')}{l(x)^2 + l(x')^2}} \exp\left(\frac{-(x - x')^2}{l(x)^2 + l(x')^2}\right), \quad (10)$$

$$\Sigma_S(x, x') = C \exp\left(-\frac{1}{2}(x - x')^2/k^2\right) \exp\left(-\frac{1}{2}((x - m)^2 + (x' - m)^2)/t^2\right), \quad (11)$$

- Signal is not parameterized
- Hyperparameters fixed by the B-only fit
- S: residual of B-subtraction

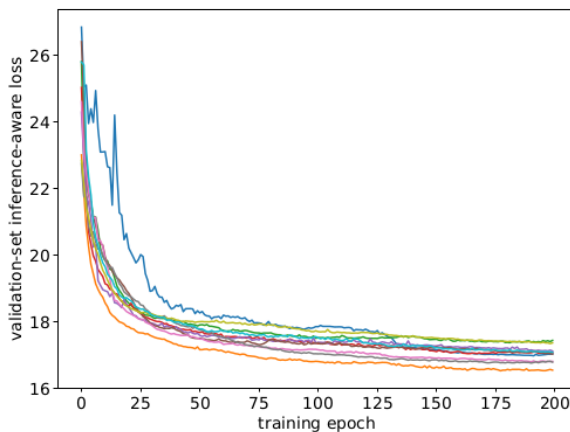
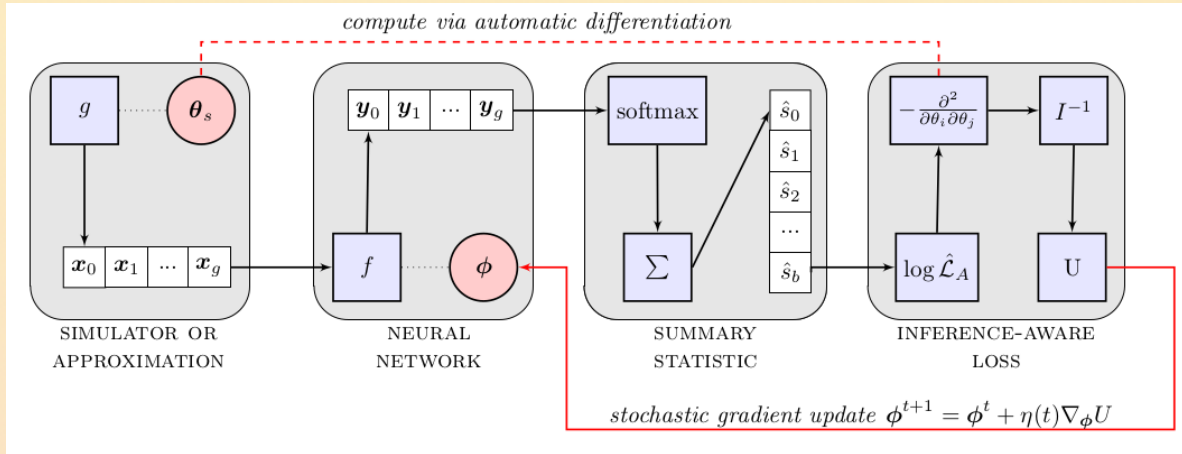


- Data: mixture model with small S
- Classification based on sample properties
 - Compare bootstrapped samples with reference (pure B)
 - Use Metodiev theorem to translate inference into signal fraction
- Validate with LR y LDA **At ICTEA!**
 - Promising results

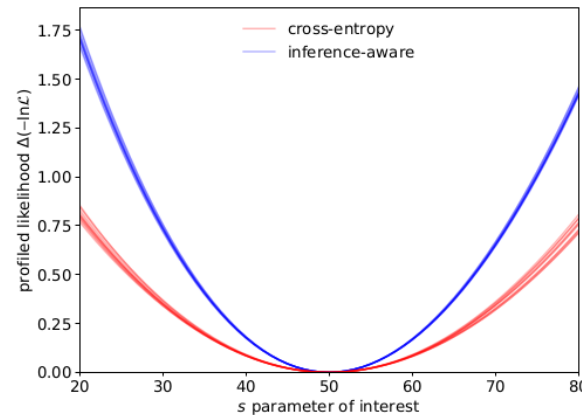


Vischia-Dorigo arXiv:1611.08256, doi:10.1051/epjconf/201713711009, and P.

Go to INFERNO: syst-aware inference opt.



(a) inference-aware training loss



(b) profile-likelihood comparison

The likelihood ratio trick

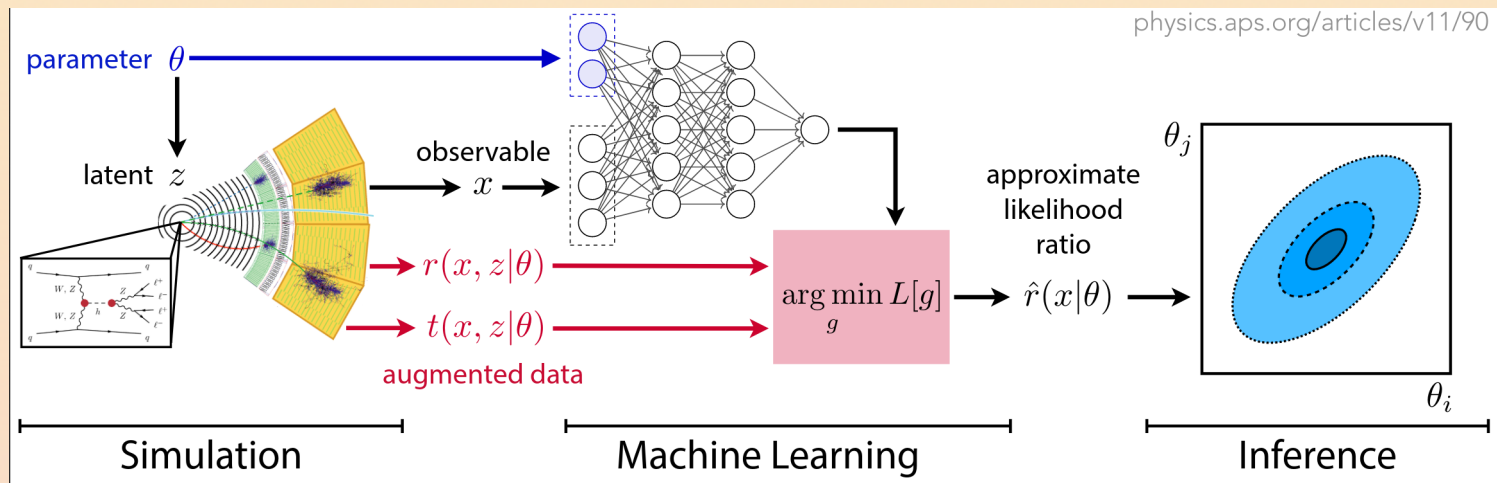
- Train surrogate to discriminate $x_i \sim p(x|\theta_0)$ from $x_i \sim p(x|\theta_1)$
- Cross-entropy

$$L_{\text{XE}} = -\mathbb{E}[1(\theta = \theta_1) \log \hat{s}(x|\theta_0, \theta_1) + 1(\theta = \theta_0) \log(1 - \hat{s}(x|\theta_0, \theta_1))]$$

- Minimized, $s(x|\theta_0, \theta_1) = p(x|\theta_1)/(p(x|\theta_0) + p(x|\theta_1))$

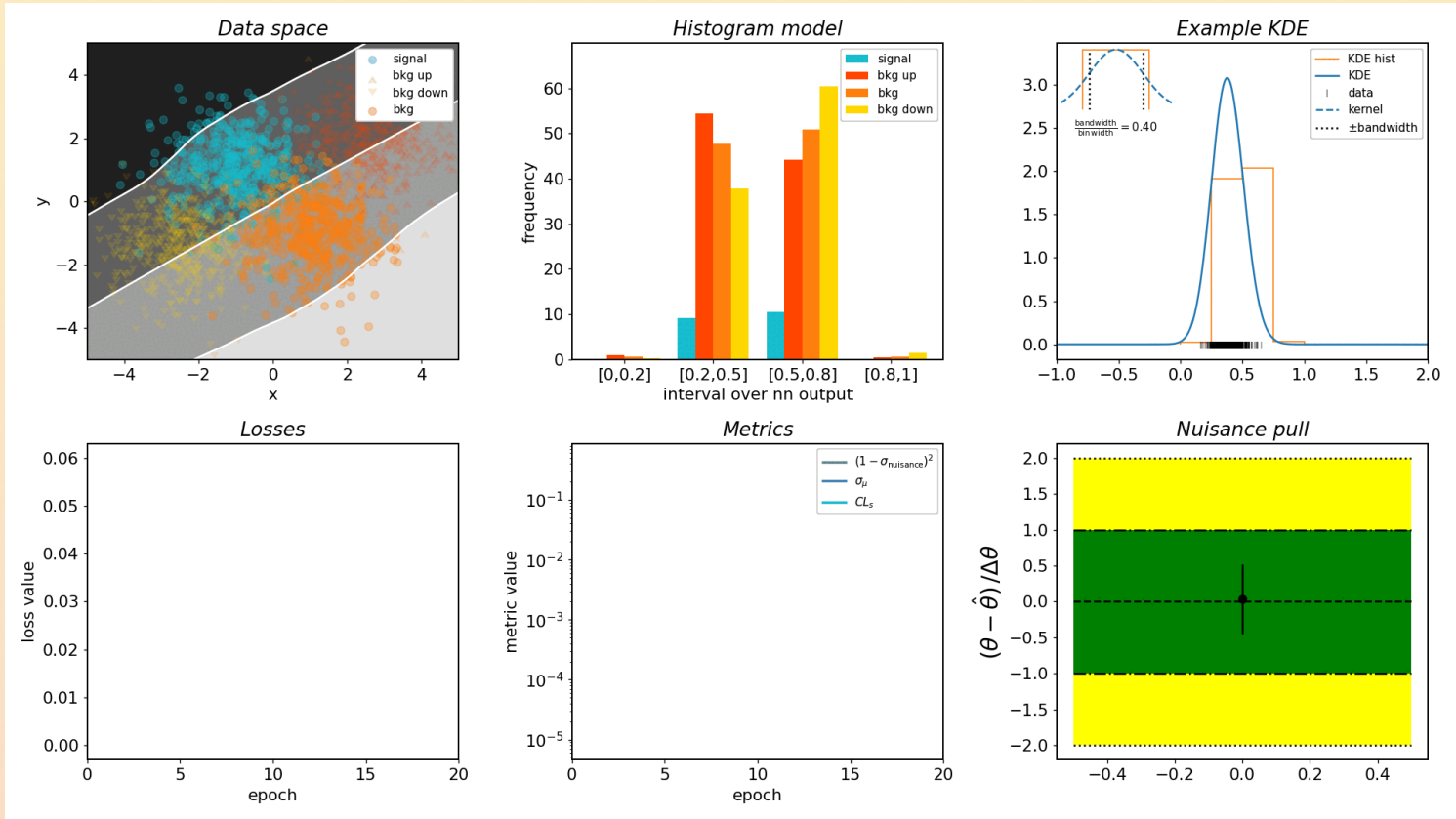
- Invert, to estimate the likelihood ratio:

$$\hat{r}(x|\theta_0, \theta_1) = (1 - \hat{s}(x|\theta_0, \theta_1))/\hat{s}(x|\theta_0, \theta_1)$$



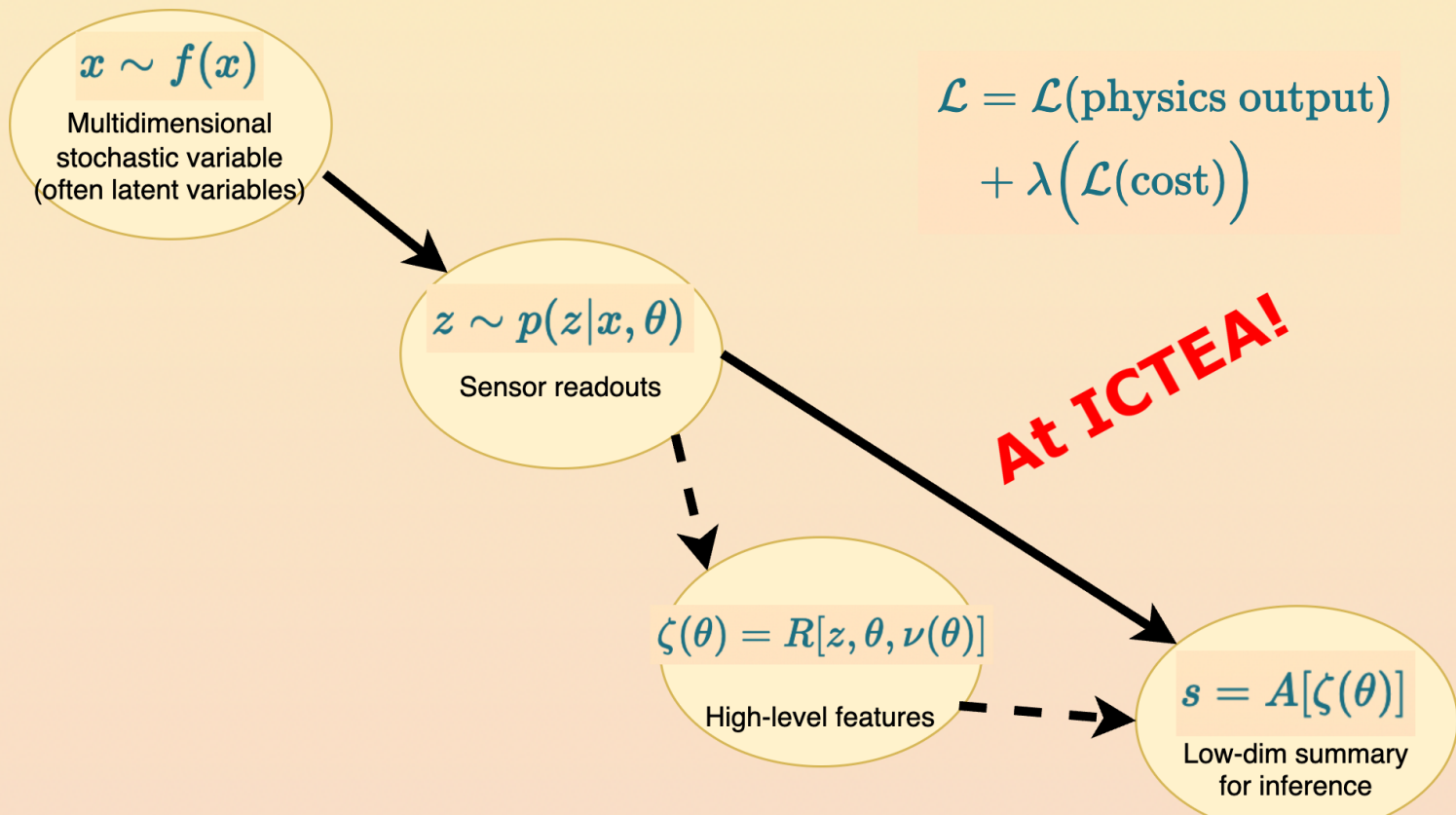
Measurement-aware analysis opt.

neos



Measurement-aware detector opt.!

- Joint optimization of design parameters w.r.t. inference made with data
- MODE White Paper, [10.1016/j.revip.2023.100085 \(2203.13818\)](https://arxiv.org/abs/2203.13818), 117-pages document, physicists + computer scientists

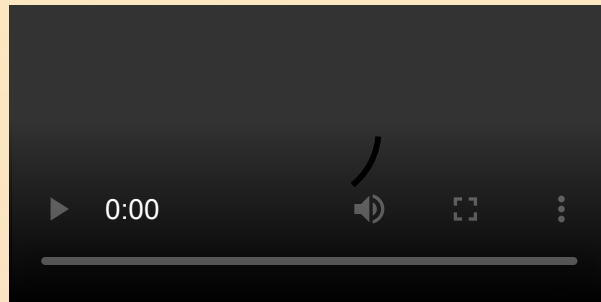


Prototype for muon tomography

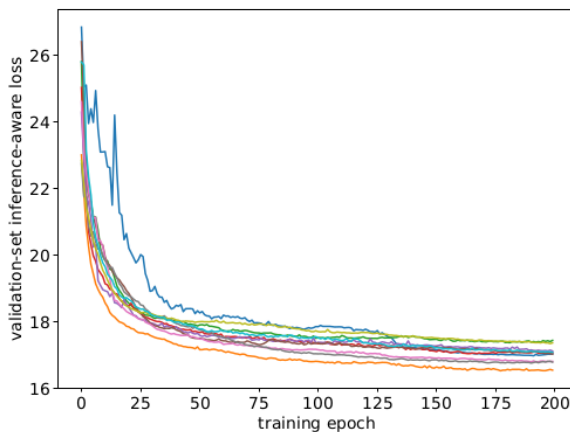
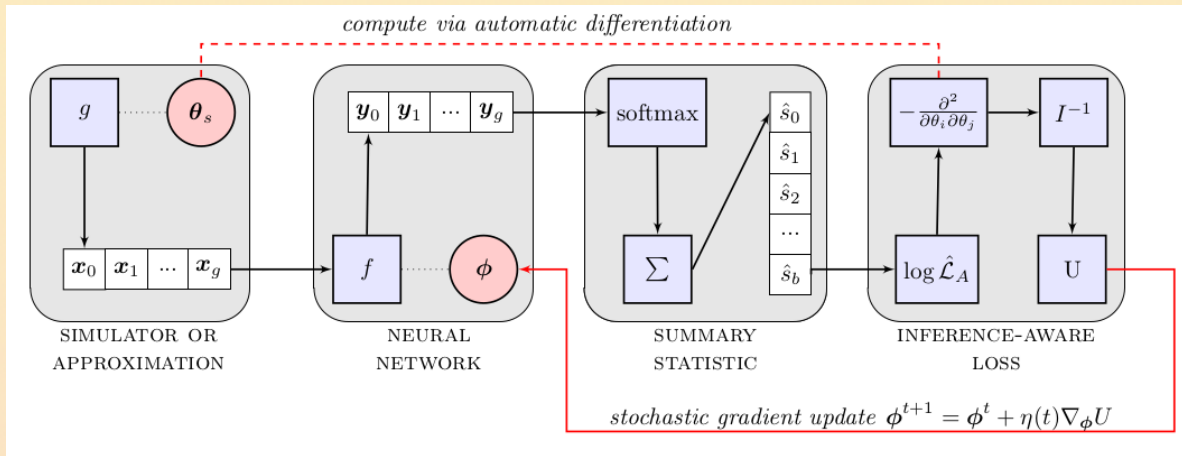
TomOpt: Differential optimisation for task- and constraint-aware design of particle detectors in the context of muon tomography

Giles C. Strong, Maxime Lagrange, Aitor Orio, Anna Bordignon, Florian Bury, Tommaso Dorigo, Andrea Giammanco, Mariam Heikal, Jan Kieseler, Max Lamparth, Pablo Martínez Ruíz del Árbol, Federico Nardi, Pietro Vischia, Haitham Zaraket

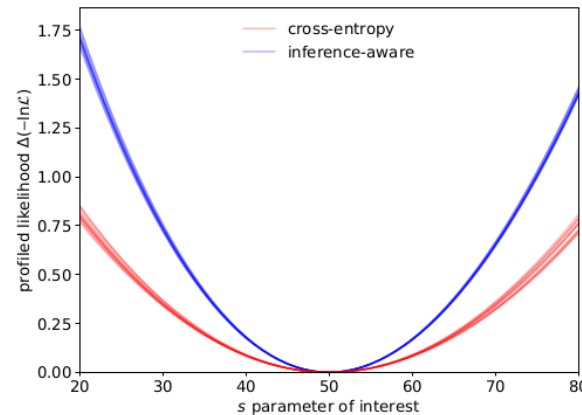
We describe a software package, TomOpt, developed to optimise the geometrical layout and specifications of detectors designed for tomography by scattering of cosmic-ray muons. The software exploits differentiable programming for the modelling of muon interactions with detectors and scanned volumes, the inference of volume properties, and the optimisation cycle performing the loss minimisation. In doing so, we provide the first demonstration of end-to-end-differentiable and inference-aware optimisation of particle physics instruments. We study the performance of the software on a relevant benchmark scenarios and discuss its potential applications.



Go to INFERNO: syst-aware inference opt.



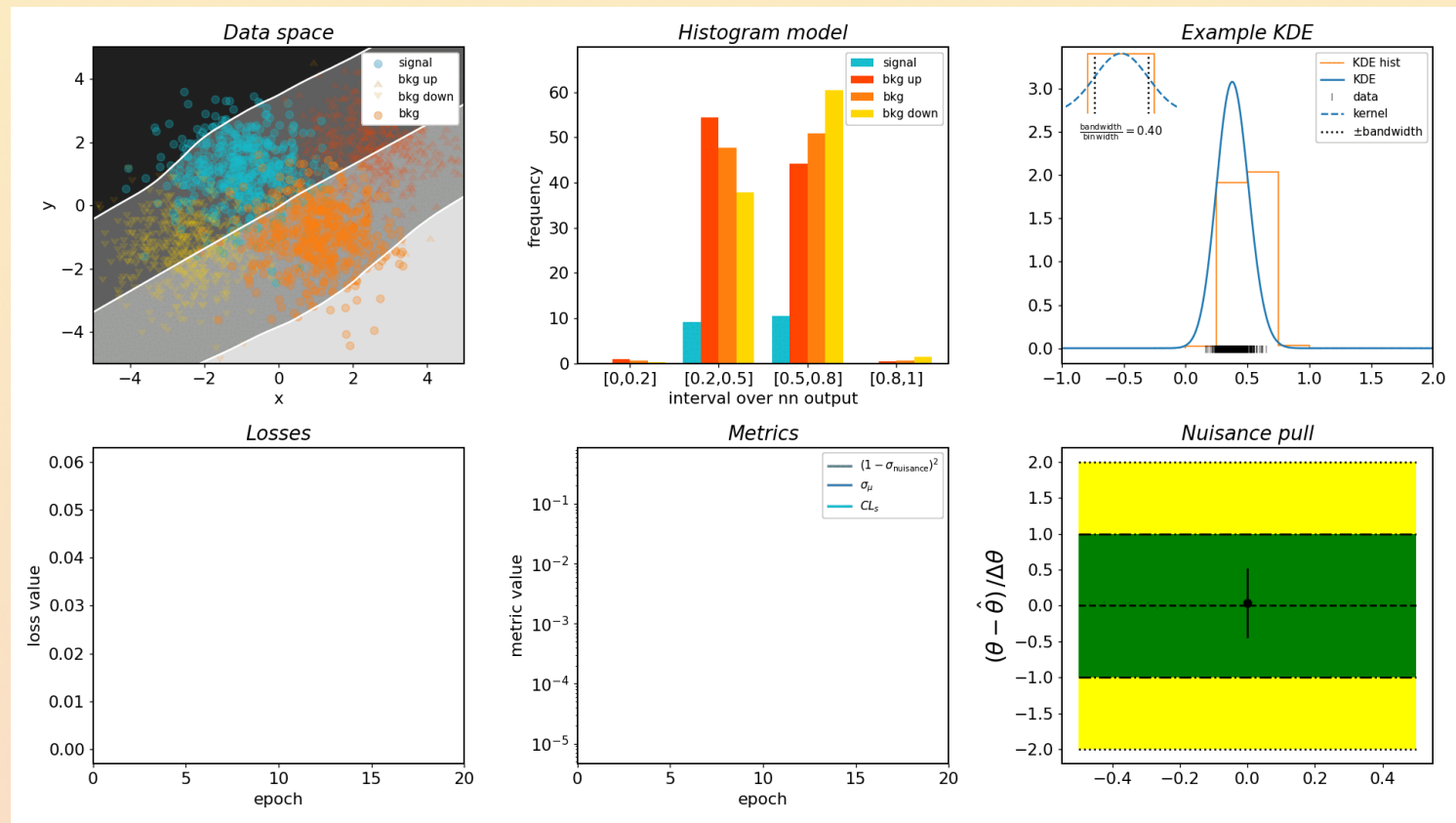
(a) inference-aware training loss



(b) profile-likelihood comparison

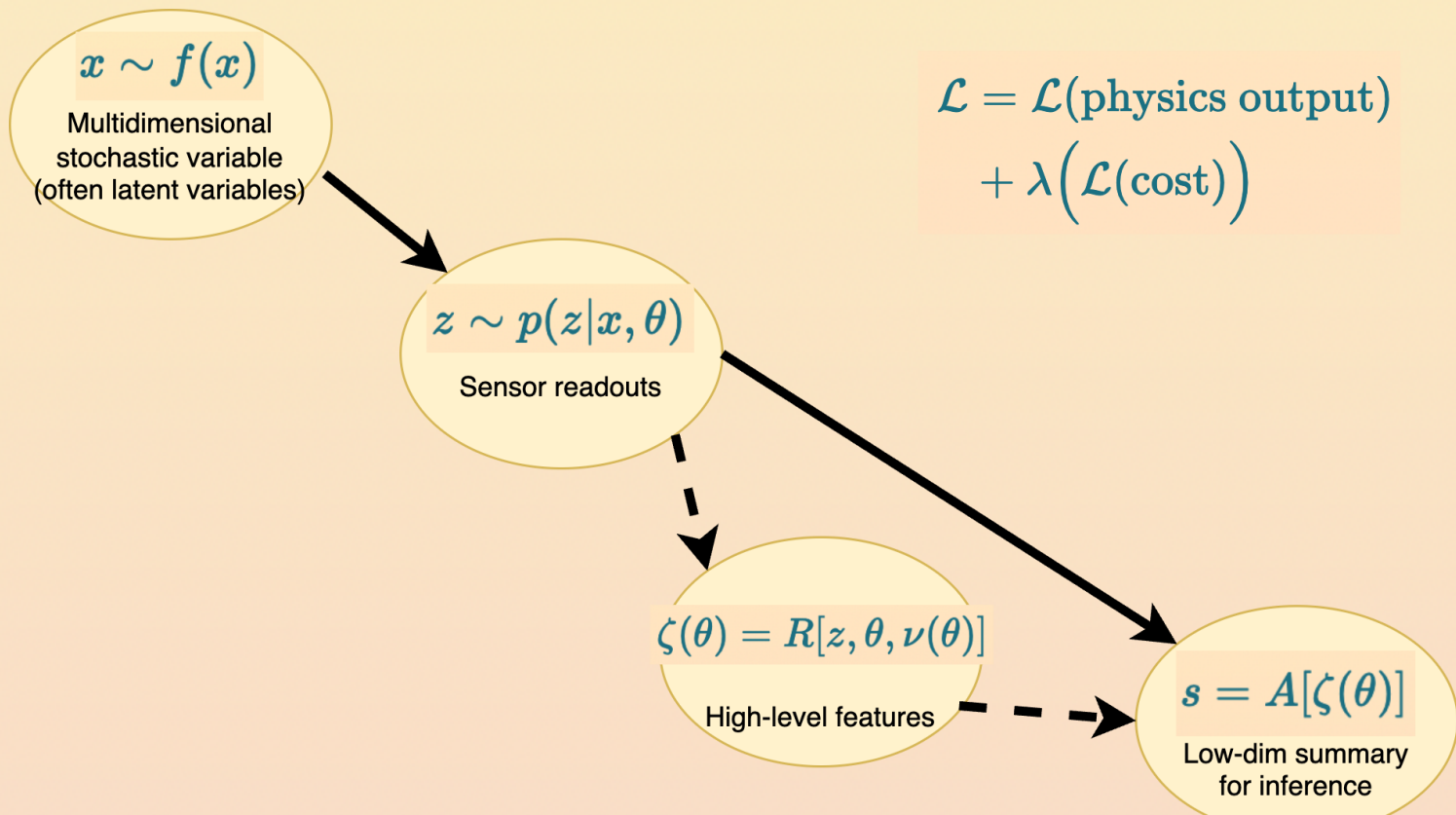
Measurement-aware analysis opt.

neos



Measurement-aware detector opt.!

- Joint optimization of design parameters w.r.t. inference made with data
- MODE White Paper, [10.1016/j.revip.2023.100085 \(2203.13818\)](https://arxiv.org/abs/2203.13818), 117-page document, physicists + computer scientists



Guarantee feasibility within constraints

- Monetary cost
- Case-specific technical constraints

$$\mathcal{L}_{\text{cost}} = c(\theta, \phi)$$

- θ : local, specific to the technology used (e.g. active components material)
- ϕ : global, describing overall detector conception (e.g. number, size, position of detector modules)
- Fixed costs can be added separately to the loss function

In general

Depends on z and nuisances

Cost of the layout with parameters θ

Closed form

$$\hat{\theta} = \arg \min_{\theta} \int L[A(\zeta), c(\theta)] p(z|x, \theta) f(x) dx dz ,$$

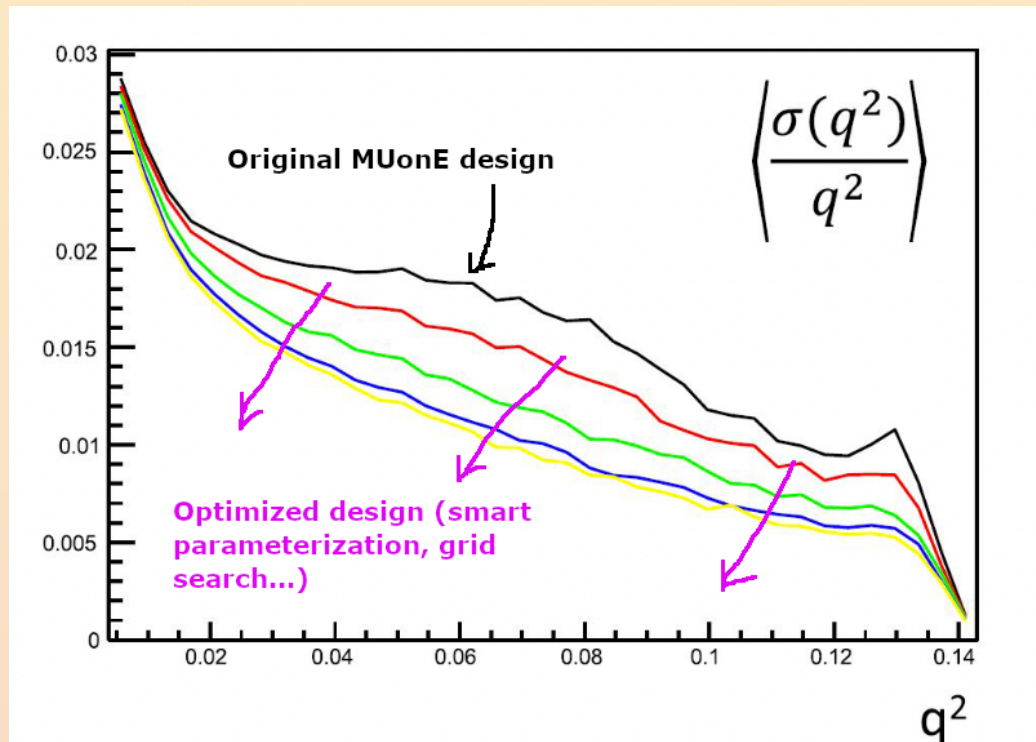
Weight desirable goals while obeying cost constraints

Thrive in asymmetries



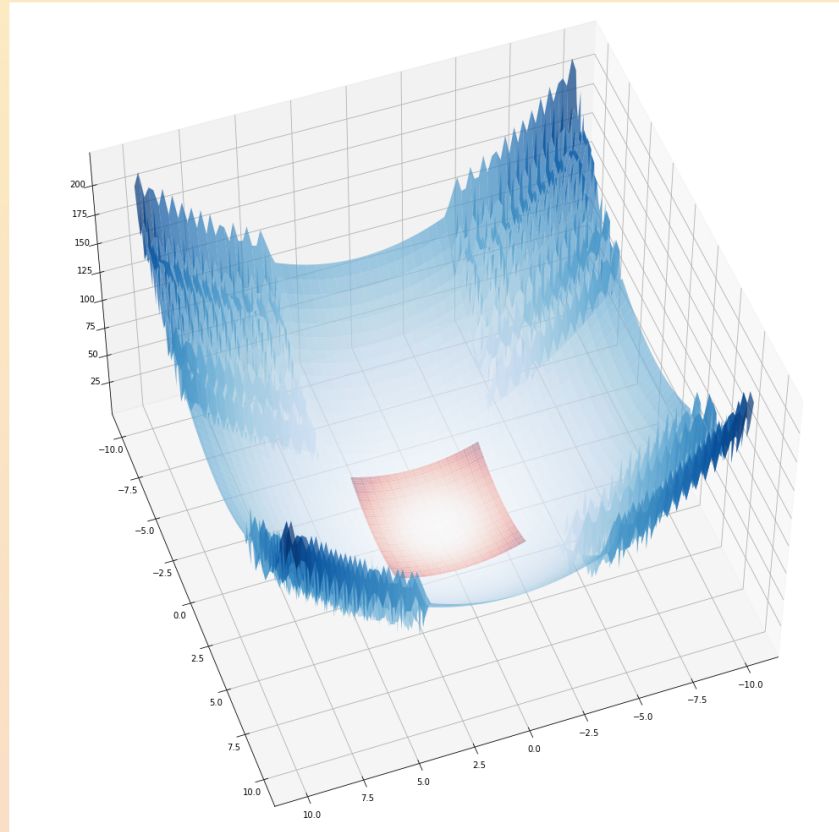
Large gains to be had

- MUonE: proposed 150 GeV muon beam experiment to be built at CERN
 - Measure precisely the q^2 differential cross section in electron-muon scattering
 - 40 tracking stations and a calorimeter
- **Dramatic improvement** in the resolution on q^2 even from a simple grid search



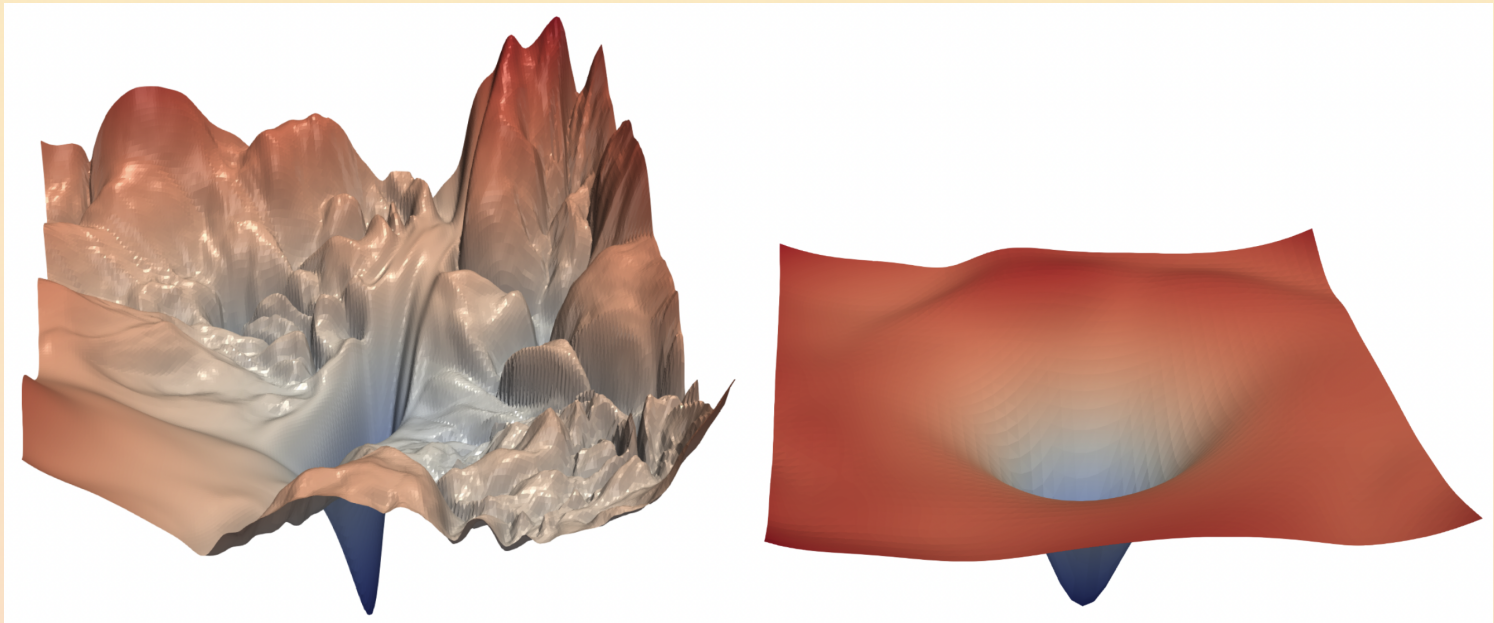
Assist the physicist with a landscape of solutions

- Cannot parameterize everything
- **The optimal solution:** unrealistic
- Provide feasible solutions near optimality
- The physicist will fine tune

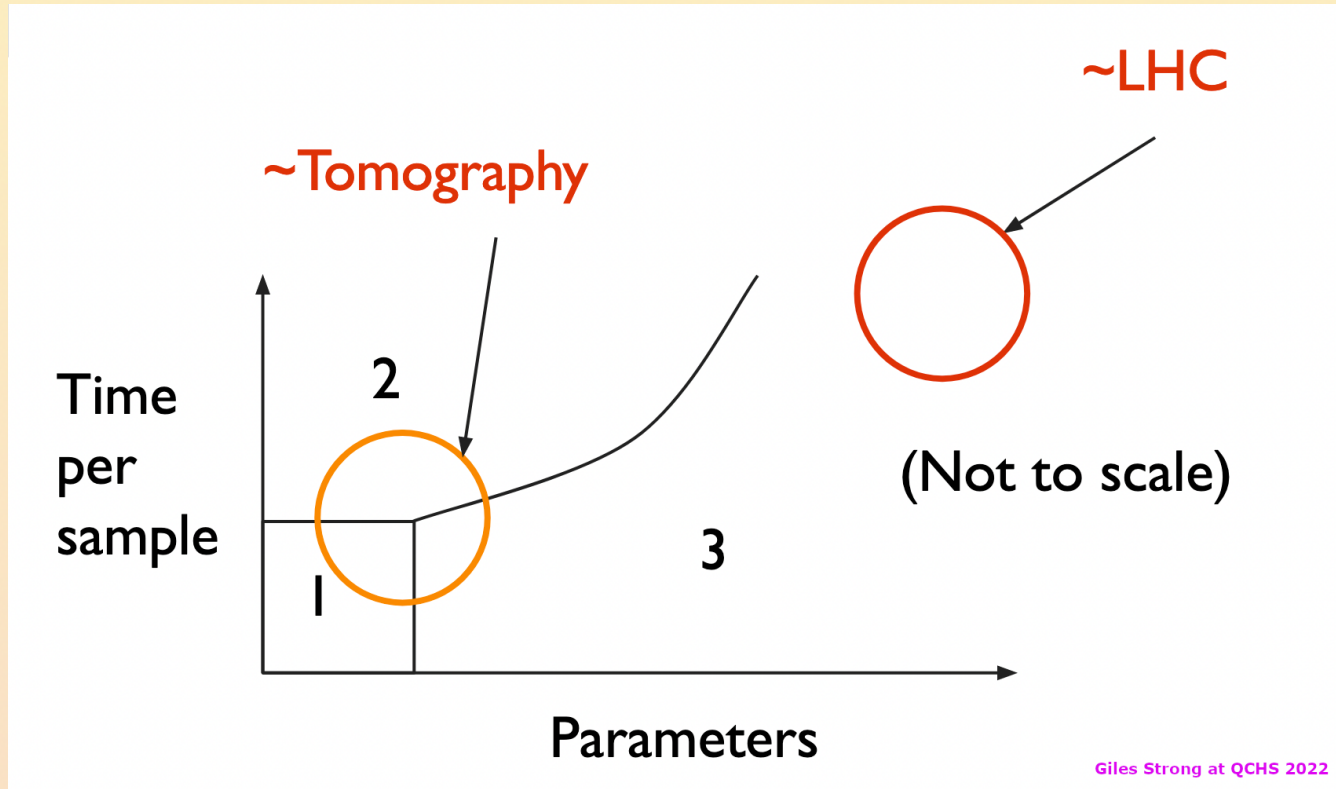


How far from optimality?

- Can we define in a general way an acceptable **increase in loss**?
 - Tradeoff performance/cost
- For sure we can regularize the loss landscape to select our scale of interest



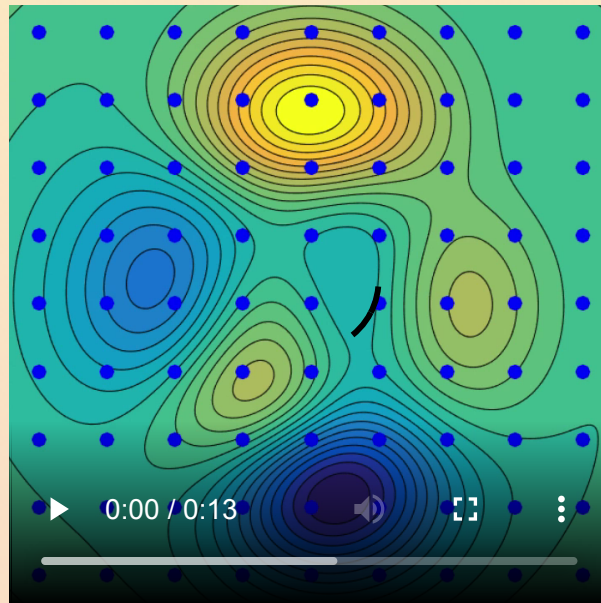
Method of choice depends on scale



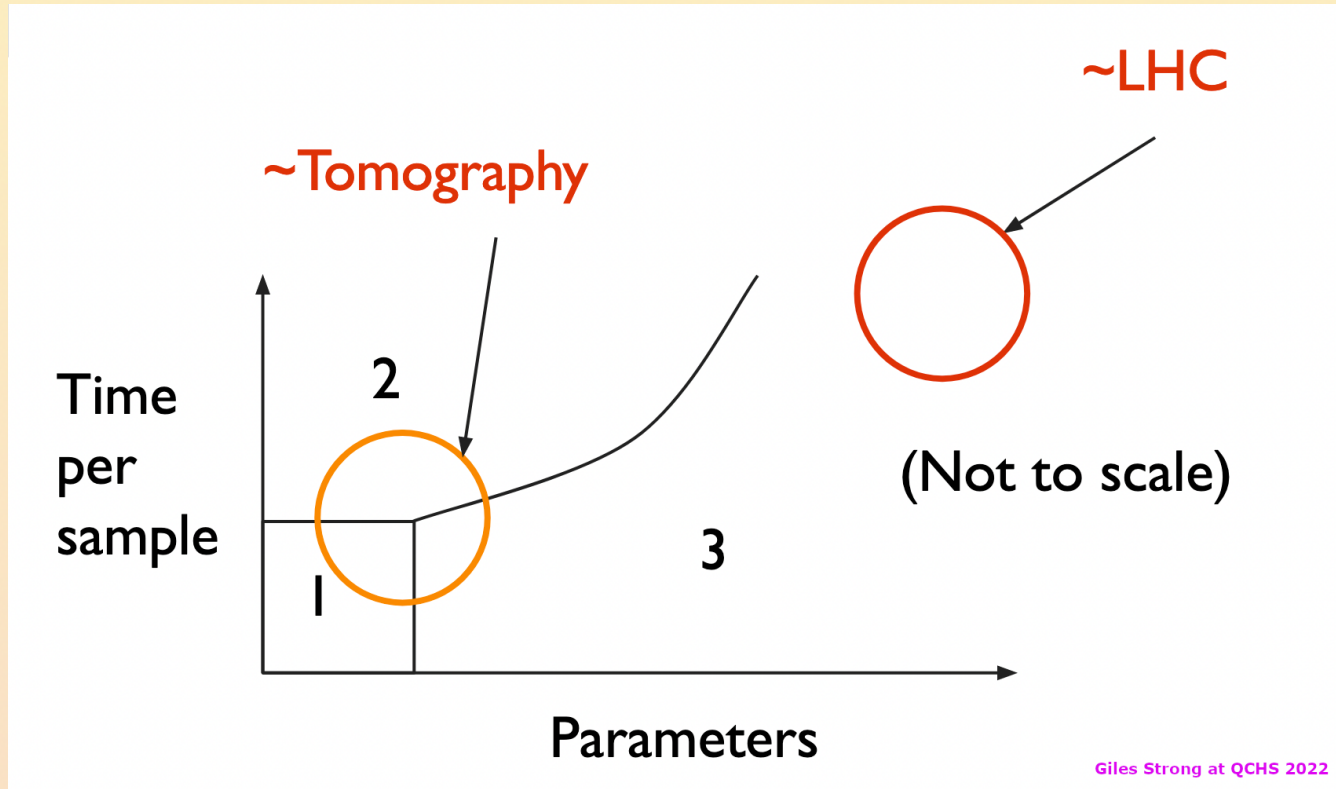
1. Grid/random search
2. Bayesian opt, simulated annealing, genetic algos, ...
3. Gradient-based optimization (Newton, BFGS, gradient descent, ...)

Experimental design: present and future

- Gradient descent applied to experiment design works!!!
 - Discreteness and stochasticity mostly solvable or avoidable
- What now?



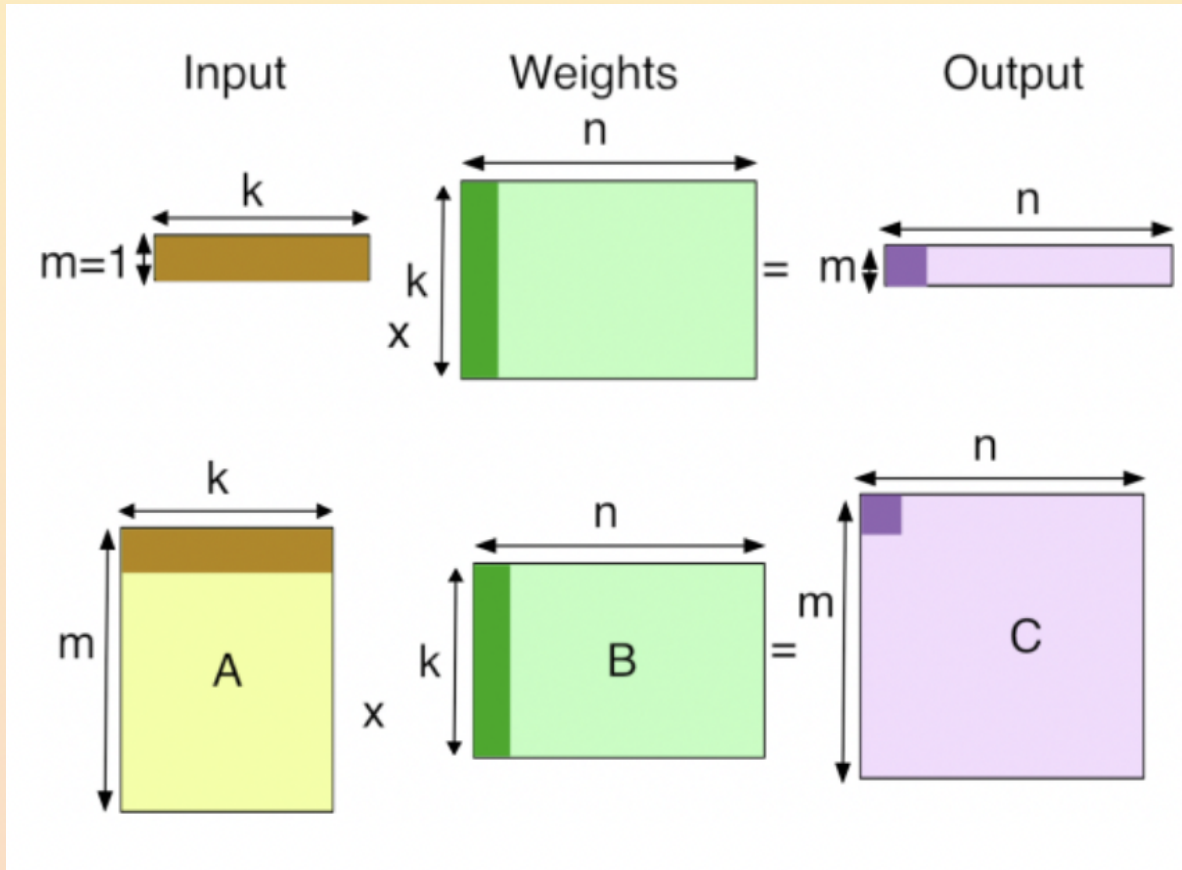
Method of choice depends on scale



1. Grid/random search
2. Bayesian opt, simulated annealing, genetic algos, ...
3. Gradient-based optimization (Newton, BFGS, gradient descent, ...)

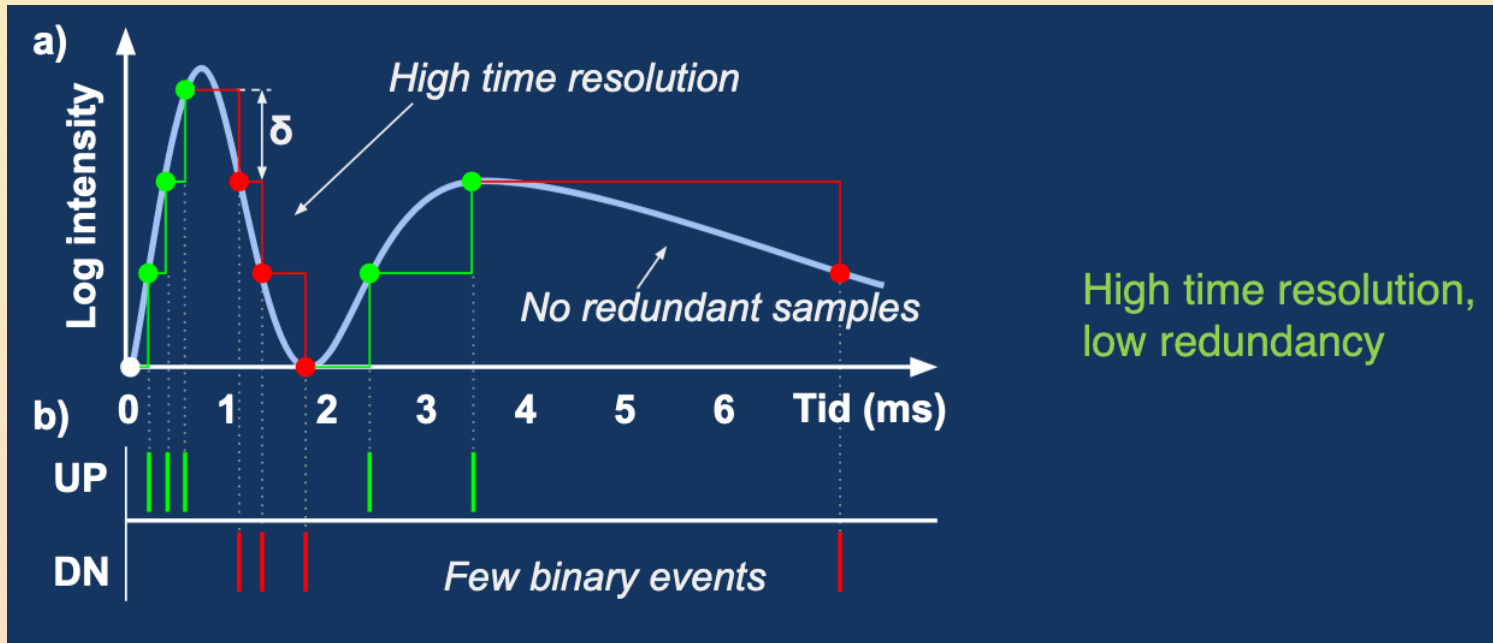
From perceptron-based networks...

- Matrix multiplication



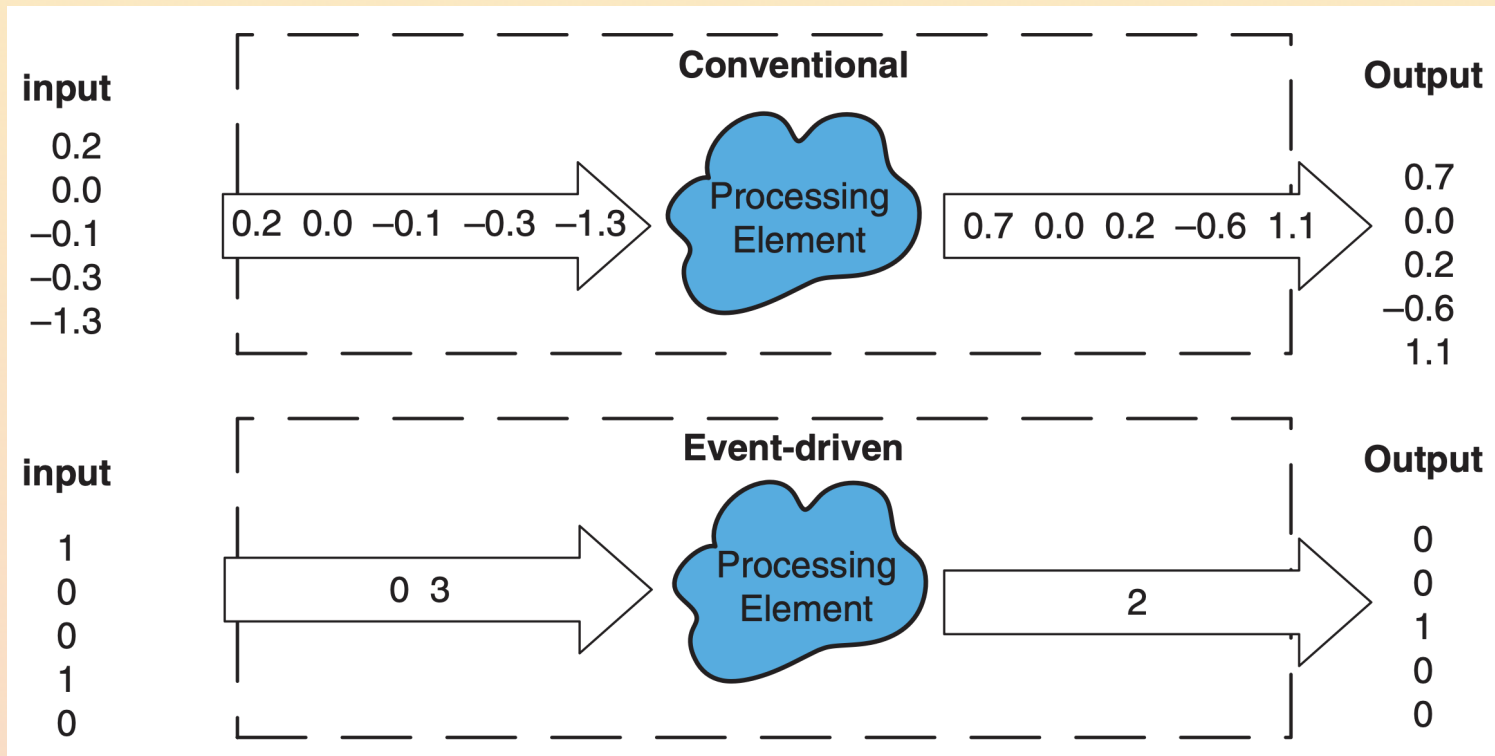
...to spiking neural networks

- Event-driven computations
 - "when a spike occurs, compute something"



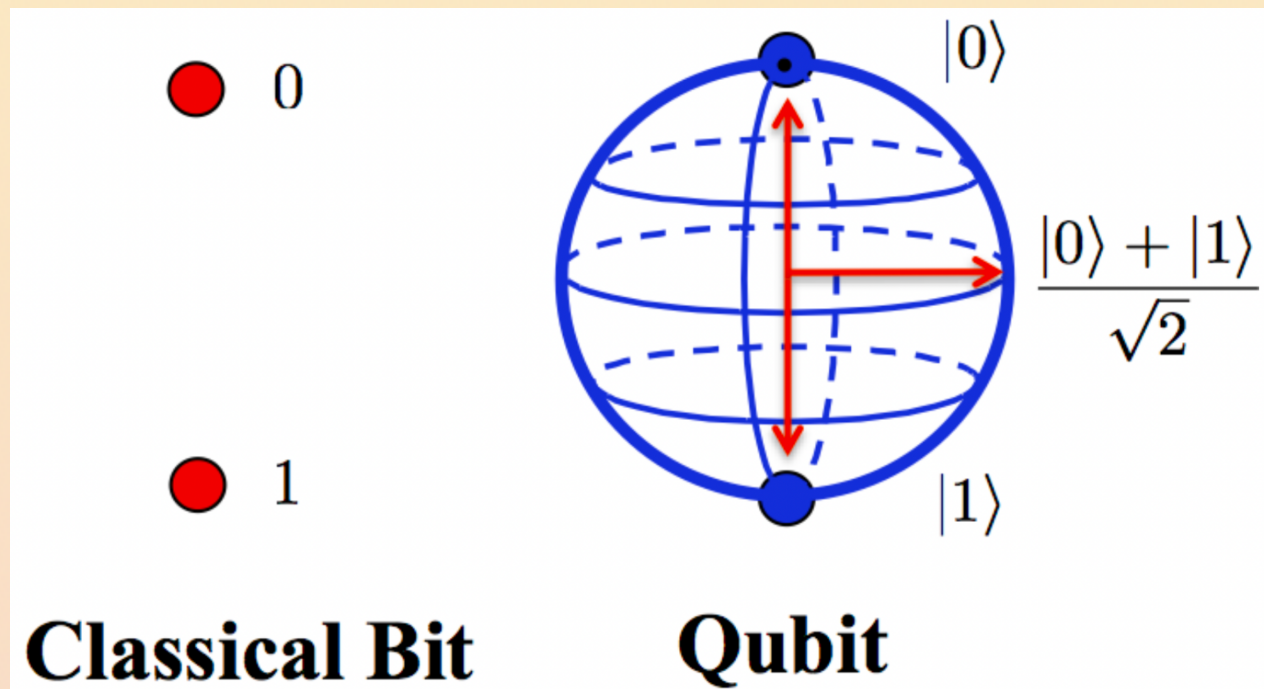
The energy advantage

- Perceptron-based networks: matrix multiplication
 - Sparsity doesn't affect much the throughput and energy consumption
- Spiking neural networks: event-driven computations
 - Sparser inputs require less computations, therefore less time and energy



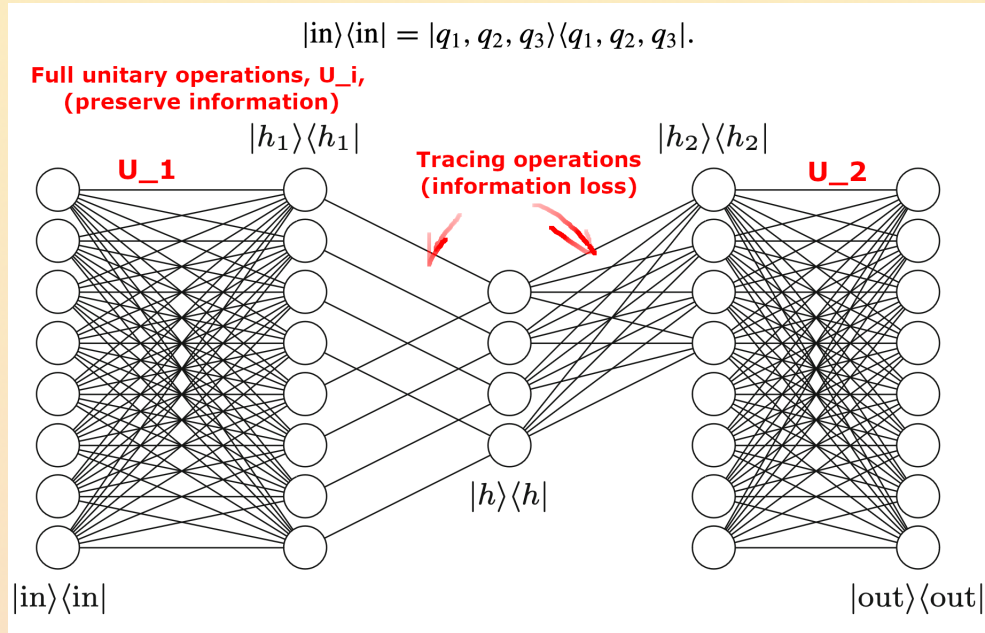
Encode information with Qubits

- Random bit (Bernoulli random variable) whose description is not governed by classical probability theory but by quantum mechanics
- Not only "because it can take real values in $[0, 1]$ ": complex numbers as coefficients α and β create **interference**
 - Interference is not reproducible with classical bits



Represent neural networks

- Qubit operations can represent rather naturally neural networks



- Gradient descent exploits intrinsic **analytic differentiability** of quantum circuits

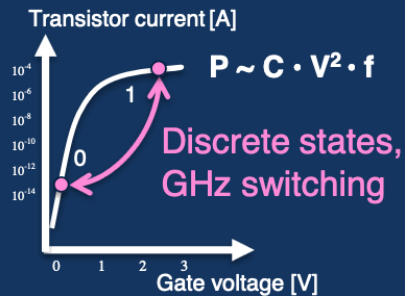
$$\begin{aligned}
 \partial_\mu \langle \psi(x, \theta) | \sigma_z | \psi(x, \theta) \rangle &= \langle 0 | \dots \partial_\mu e^{-i\mu\sigma} \dots \sigma_z \dots e^{i\mu\sigma} \dots | 0 \rangle \\
 &\quad + \langle 0 | \dots e^{-i\mu\sigma} \dots \sigma_z \dots \partial_\mu e^{i\mu\sigma} \dots | 0 \rangle \\
 &= \langle 0 | \dots (-i\sigma) e^{-i\mu\sigma} \dots \sigma_z \dots e^{i\mu\sigma} \dots | 0 \rangle \\
 &\quad + \langle 0 | \dots e^{-i\mu\sigma} \dots \sigma_z \dots (i\sigma) e^{i\mu\sigma} \dots | 0 \rangle \\
 &= \langle 0 | \dots (1 - i\sigma) e^{-i\mu\sigma} \dots \sigma_z \dots (1 + i\sigma) e^{i\mu\sigma} \dots | 0 \rangle \\
 &\quad + \langle 0 | \dots (1 + i\sigma) e^{-i\mu\sigma} \dots \sigma_z \dots (1 - i\sigma) e^{i\mu\sigma} \dots | 0 \rangle
 \end{aligned}$$

Need for new paradigm

- If you are interested in Neuromorphic computing or Quantum computing, drop me a line!

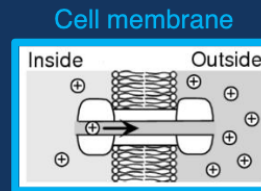
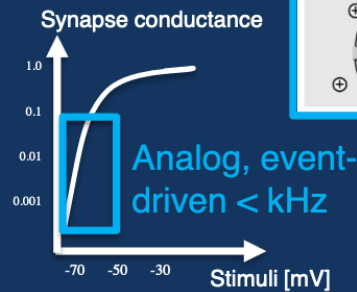
Conventional computers

mimic logical and analytical thinking

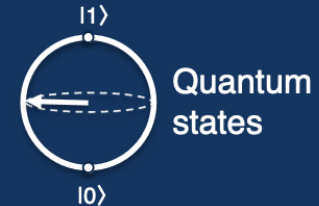


Neuromorphic processors

mimic the senses, learning and perception



Quantum processors
use quantum superpositions for probabilistic inference



Technology readiness?

Thank you!

