

Modern Machine Learning in HEP

8th BCD ISHEP Cargèse School

27-31 March 2023

Julien Donini – Université Clermont Auvergne / LPC



Machine learning is statistics + algorithms + computing power

HEP: used since the 80's : old style **NN**, then **Boosted Decision Trees** “era”

Modern Machine learning since about 10 years

→ **Breakthrough** ideas supported by statisticians, computer scientists, etc

→ Increasing computing **power** to run efficiently **complex** algorithms

Powerful tool which enables us to do better but also **new kinds of physics**

2014 Kaggle Higgs Challenge (<https://www.kaggle.com/c/higgs-boson>)

- Improve measurements of **Higgs** decaying to pair of tau leptons
- 1800 participating teams (physicists, statisticians, computer scientists).

Solution	Score
Gabor Melis (DNN pooling)	3.806
MultiBoost	3.405
TMVA boosted trees	3.200
Naive Bayesian classifier	2.060
1D cut-based selection	1.535

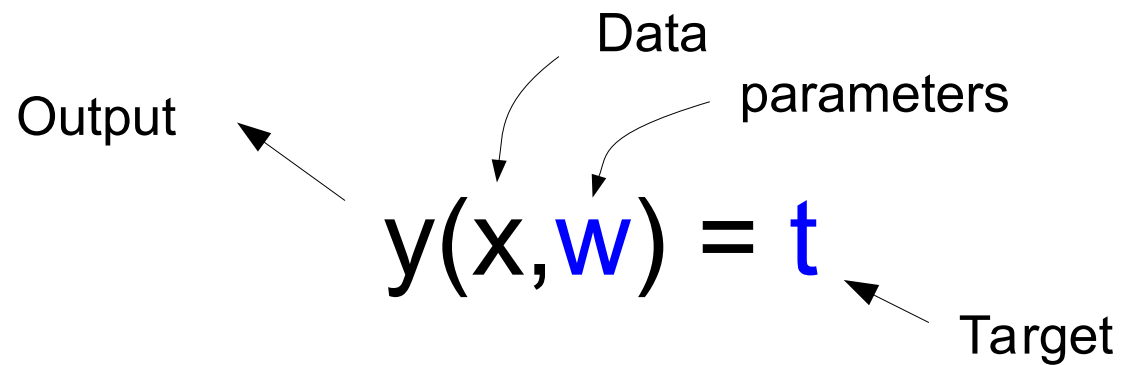


Winning solution performance equivalent to having **6 times more data** !

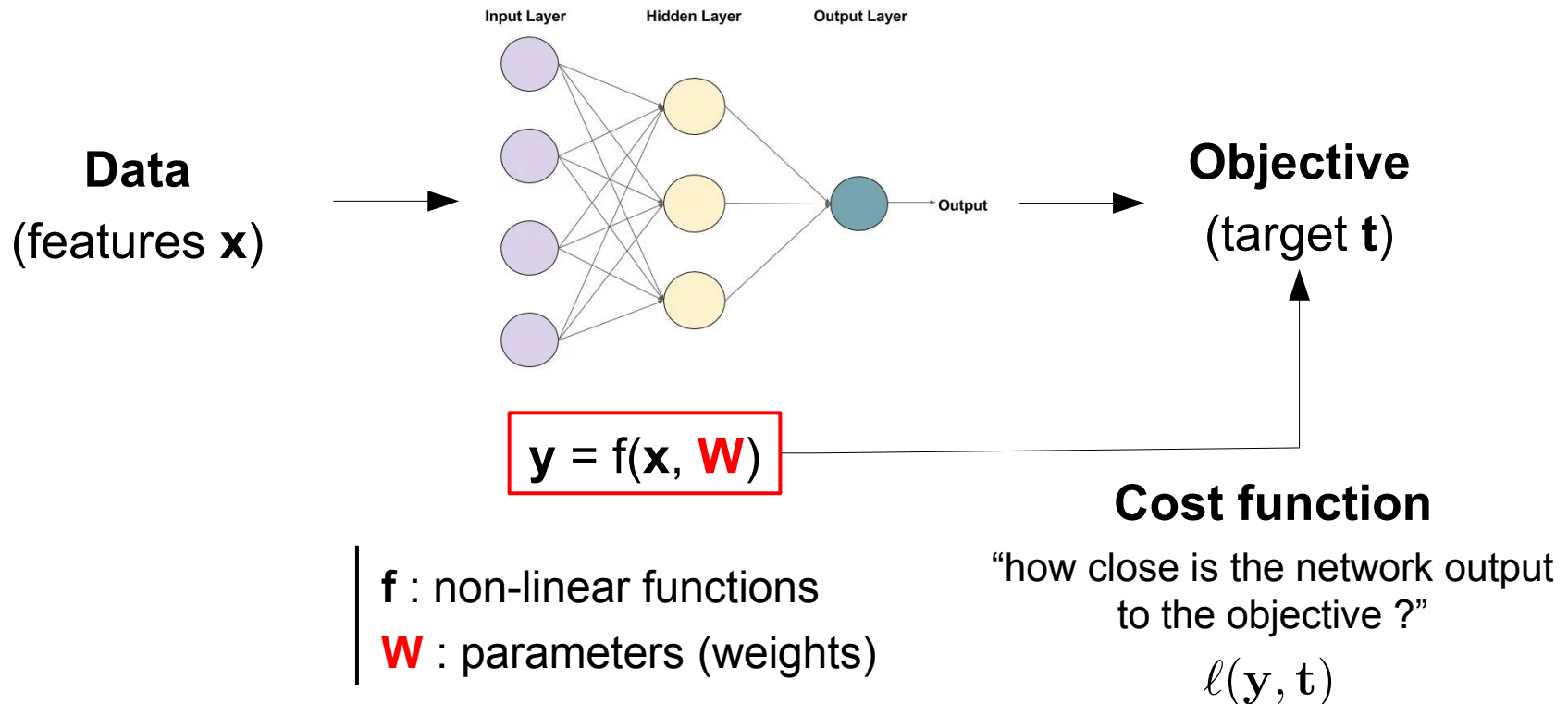
Part1: Modern ML for HEP

- Introduction to ML
- Example of advanced methods
- Differentiable programming

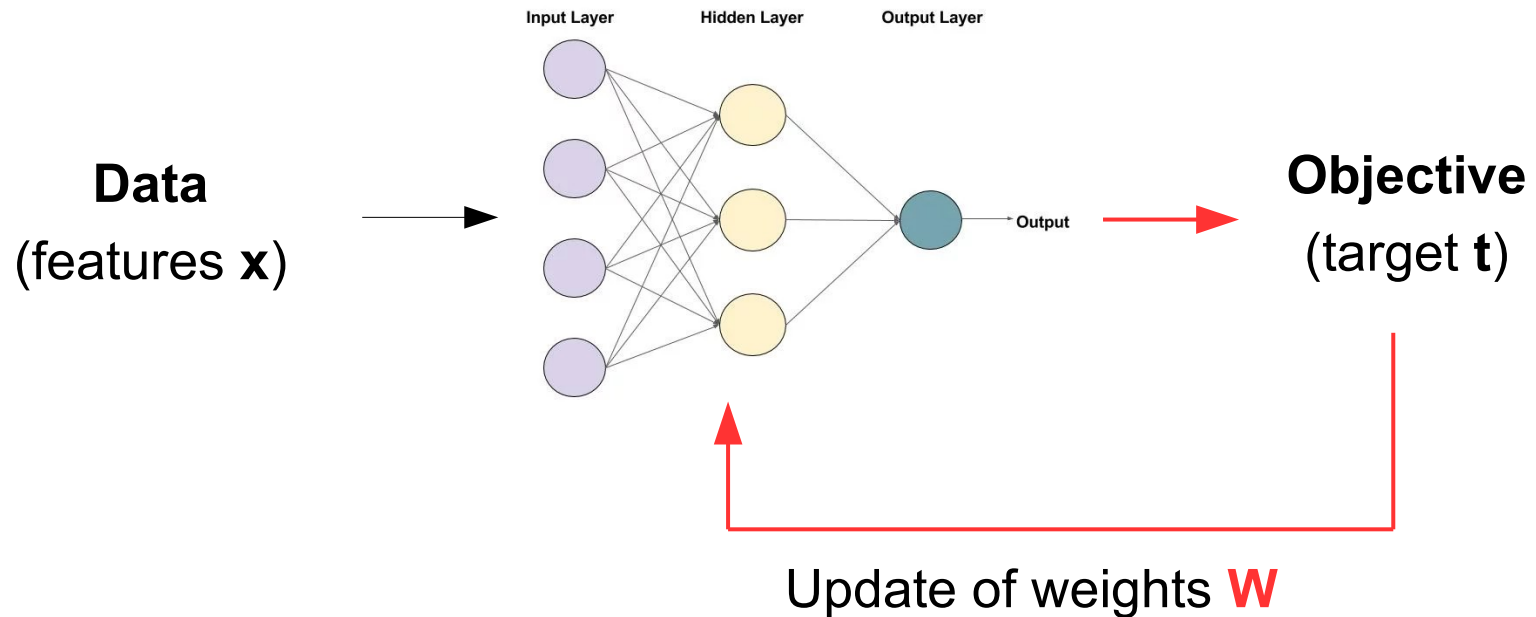
Part 2: Neural Networks and variational inference



Training Neural Networks



Training Neural Networks



$$\mathbf{W} \rightarrow \mathbf{W} - \eta \sum_N \frac{\partial \ell(\mathbf{y}, \mathbf{t})}{\partial \mathbf{W}}$$

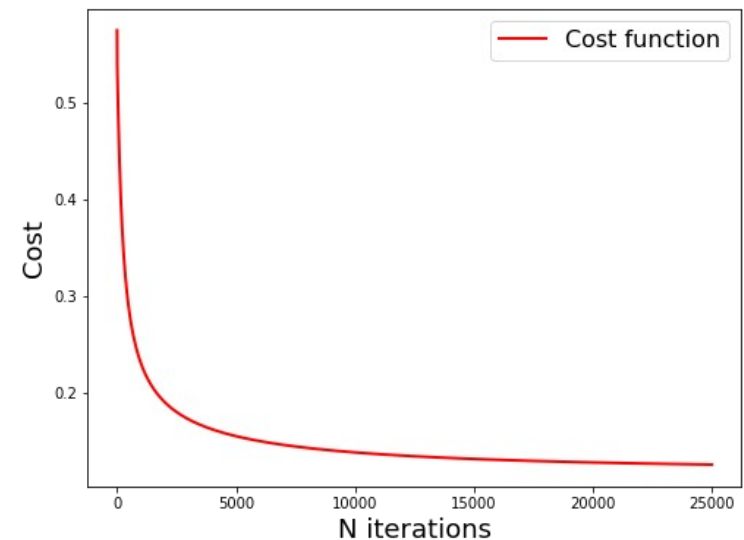
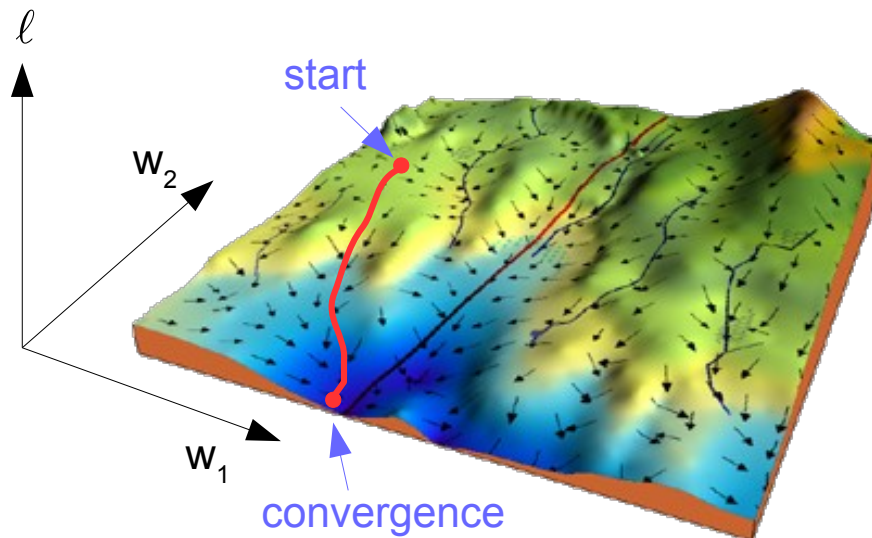
Gradient descent

Start from initial set of weights \mathbf{w} and subtract gradient of ℓ iteratively:

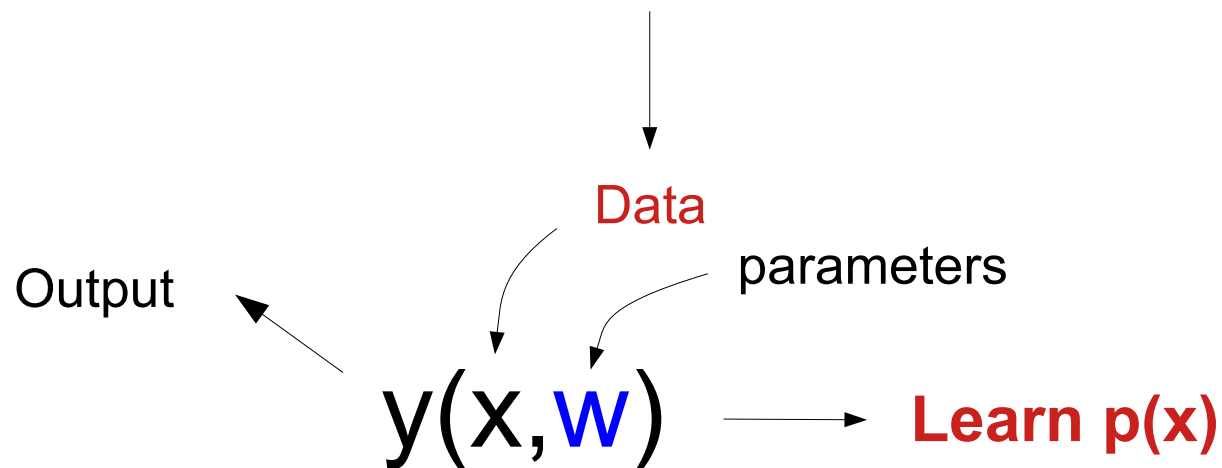
$$\mathbf{W}^k \rightarrow \mathbf{W}^{k+1} = \mathbf{W}^k - \eta \sum_N \frac{\partial \ell(\mathbf{W}^k)}{\partial \mathbf{W}}$$

k : iteration, η : learning speed

Repeat until convergence.



Data \mathbf{X} drawn from some unknown **distribution** $p(\mathbf{x})$



Learn $p(\mathbf{x})$ itself → density estimation, eg Normalizing Flows

Conditional densities $p(\mathbf{x} | \mathbf{y})$ → conditional density estimation

Sampling from $p(\mathbf{x})$ → generative modeling, eg GANs, VAEs, ...

Ratios of densities $p_1(\mathbf{x})/p_2(\mathbf{x})$ → classification, eg CNNs, RNNs, GNNs, ...

A Living Review of ML for HEP

Huge collection of references : <https://iml-wg.github.io/HEPML-LivingReview/>

Covered topics

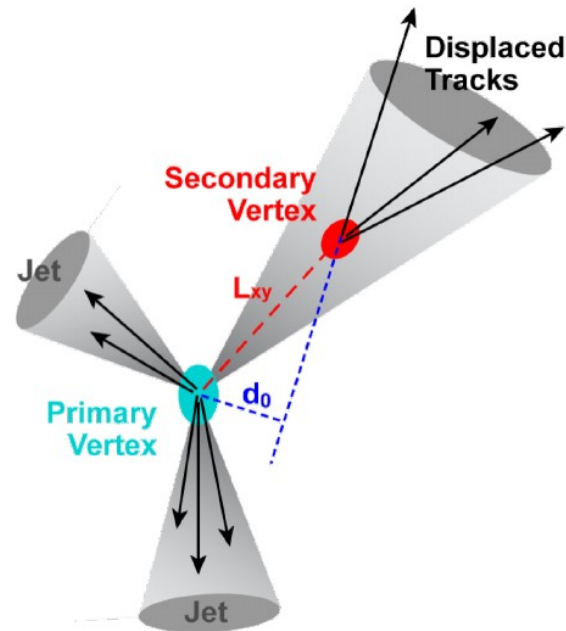
- **ML reviews**: modern, historical, ...
- **Classification**: jet images, graphs, flavor tagging, ...
- **Regression**: calibration, parameter estimation, matrix element, ...
- **Generative** models/density estimation: GAN, normalizing flows, ...
- **Anomaly** detection: BSM searches, hardware faults, real time detection...
- Simulation-based **Inference**: parameter estimation, unfolding, ...
- **Uncertainty** Quantification: interpretation, mitigation, estimation, ...
- ...

Many more “Proof-of-concept” than applications “in production”



Flavor tagging

BSM searches



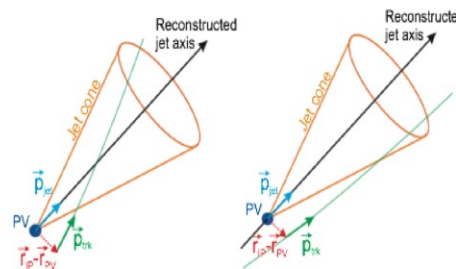
Goal: Discriminate b-jets from c-jets, light-jets, tau

B-tagging algorithms utilize the long lifetime and displaced decays of b-hadrons to look for secondary vertices and displaced tracks

Flavor-tagging: Baseline Algorithms

- **IP3D / IP2D:**

- Impact parameter algorithm
- Exploit (in)compatibility of track with PV



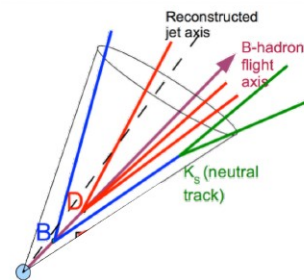
- **SV:**

- **Inclusive Secondary vertexing**
- Determination of single inclusive weak b-hadron decay vertex



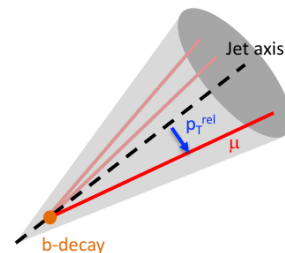
- **JetFitter:**

- **PV \rightarrow B \rightarrow D decay chain finding**
- More detailed determination of decay vertex topology



- **SMT:** Soft-muon tagger BDT utilizes:

- Muon kinematics and impact parameter
- Track quality, to reject fakes and decays in flight

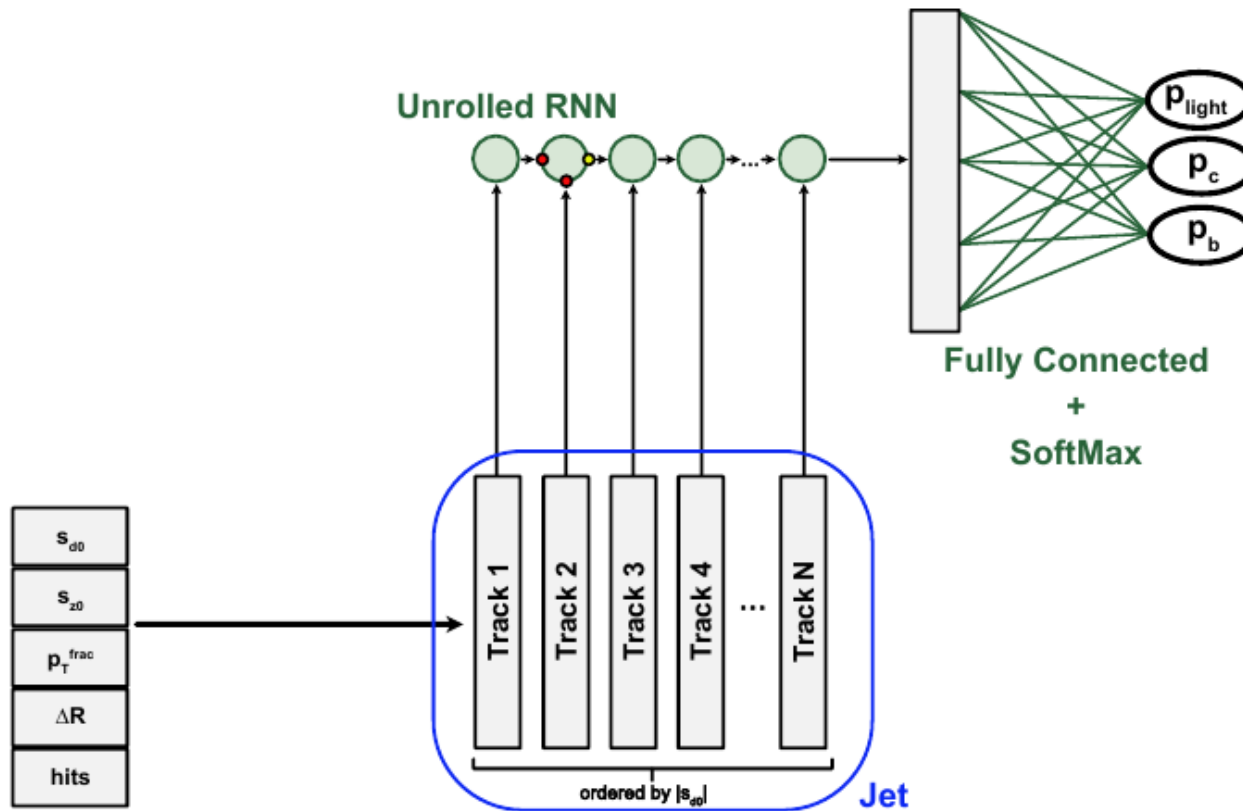


[slide M. Kagan]

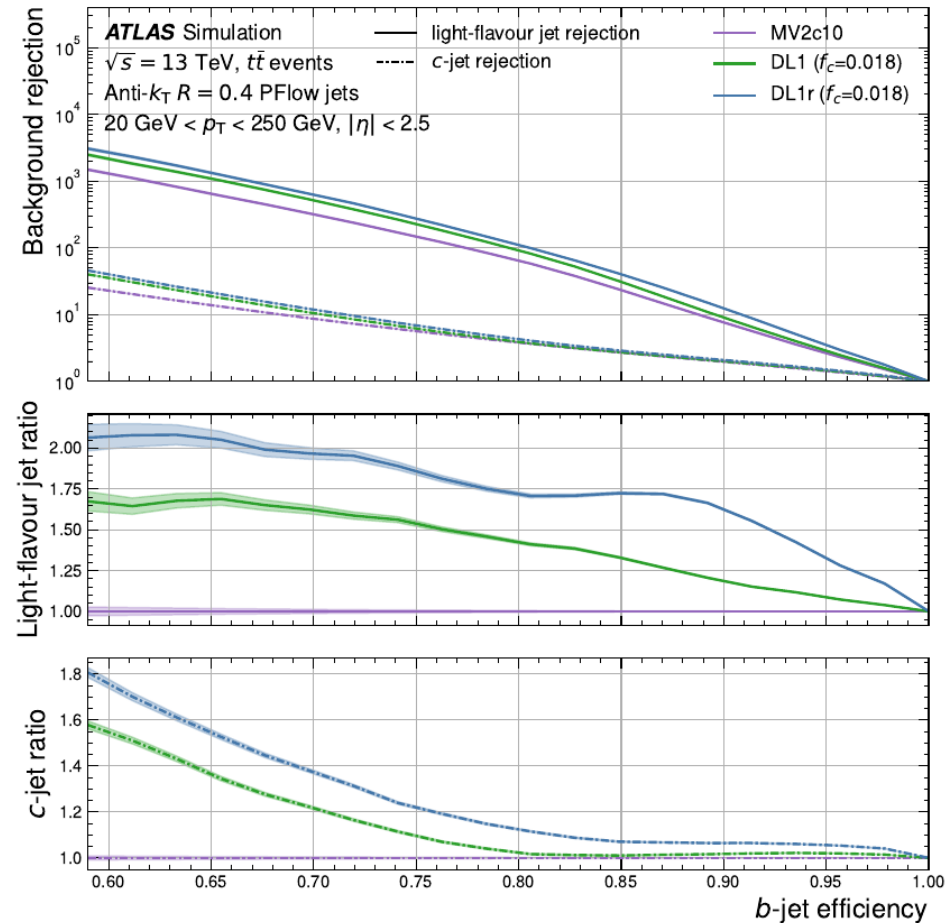
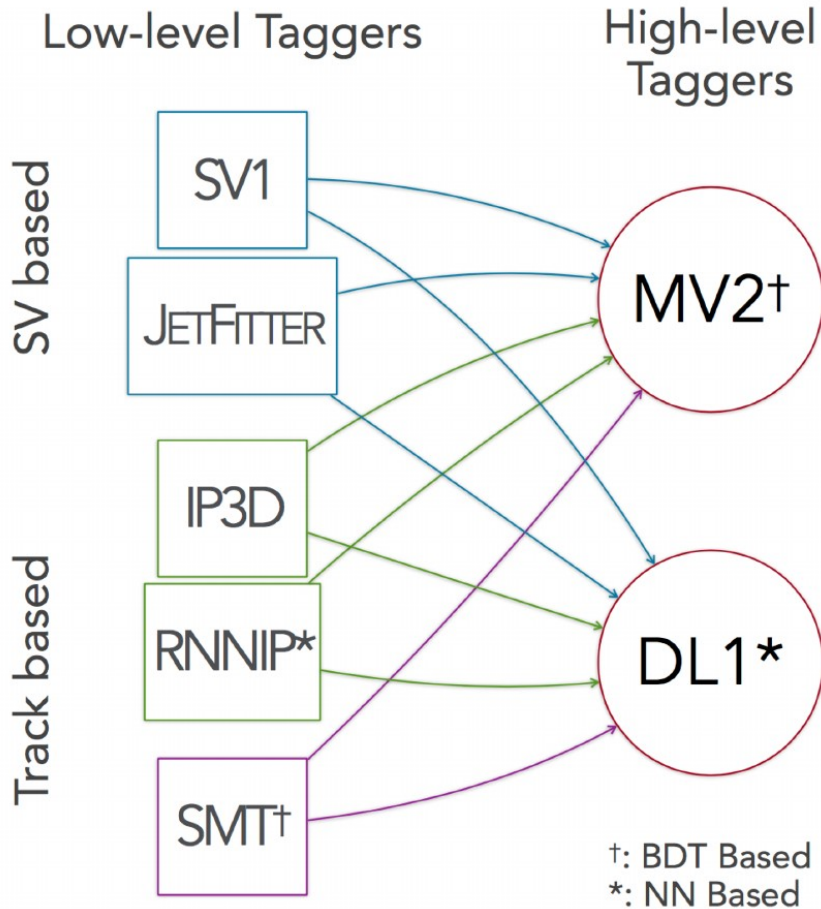
Flavor-tagging: Recurrent NN

Recurrent NN (sentence classification, NLP, ...)

- Treat tracks as a sequence, ordered by impact parameter significance



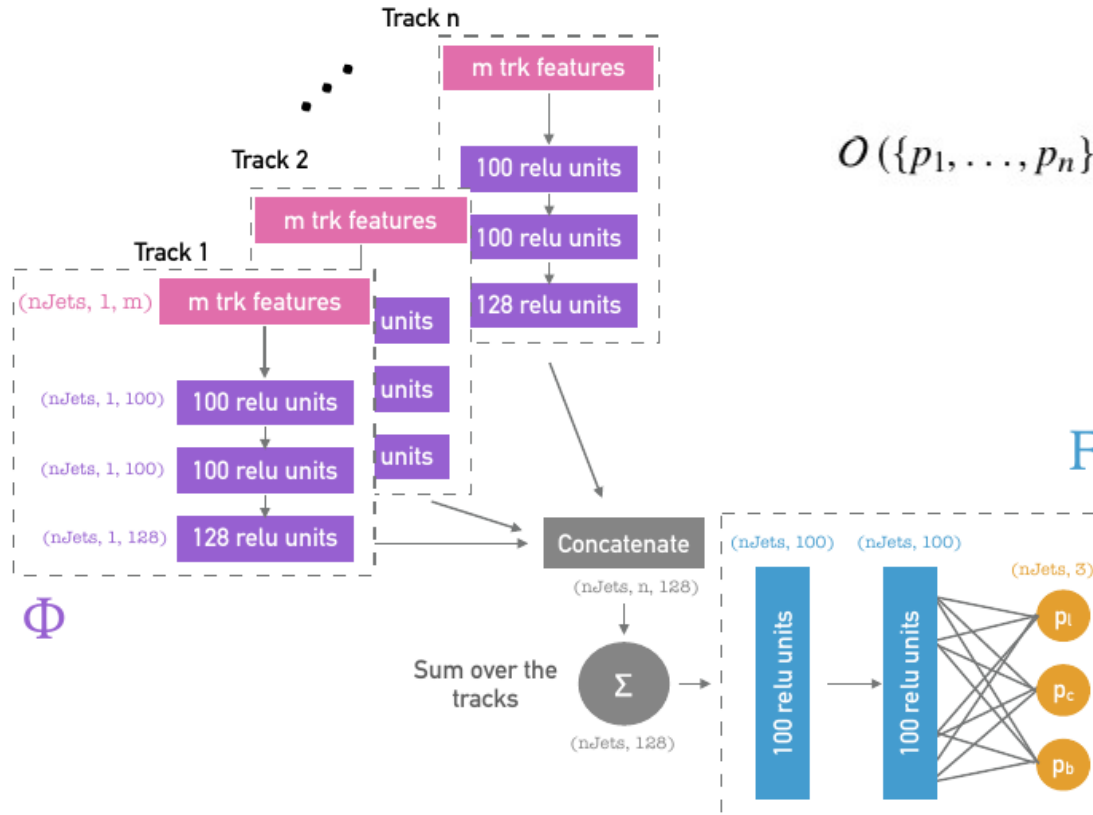
High-Level Taggers



<https://arxiv.org/abs/2211.16345>

Improved algorithms: Deep Sets

The **Deep Sets** architecture treats each tracks as a set without any specific order = maintain benefit of RNN without requiring ordering



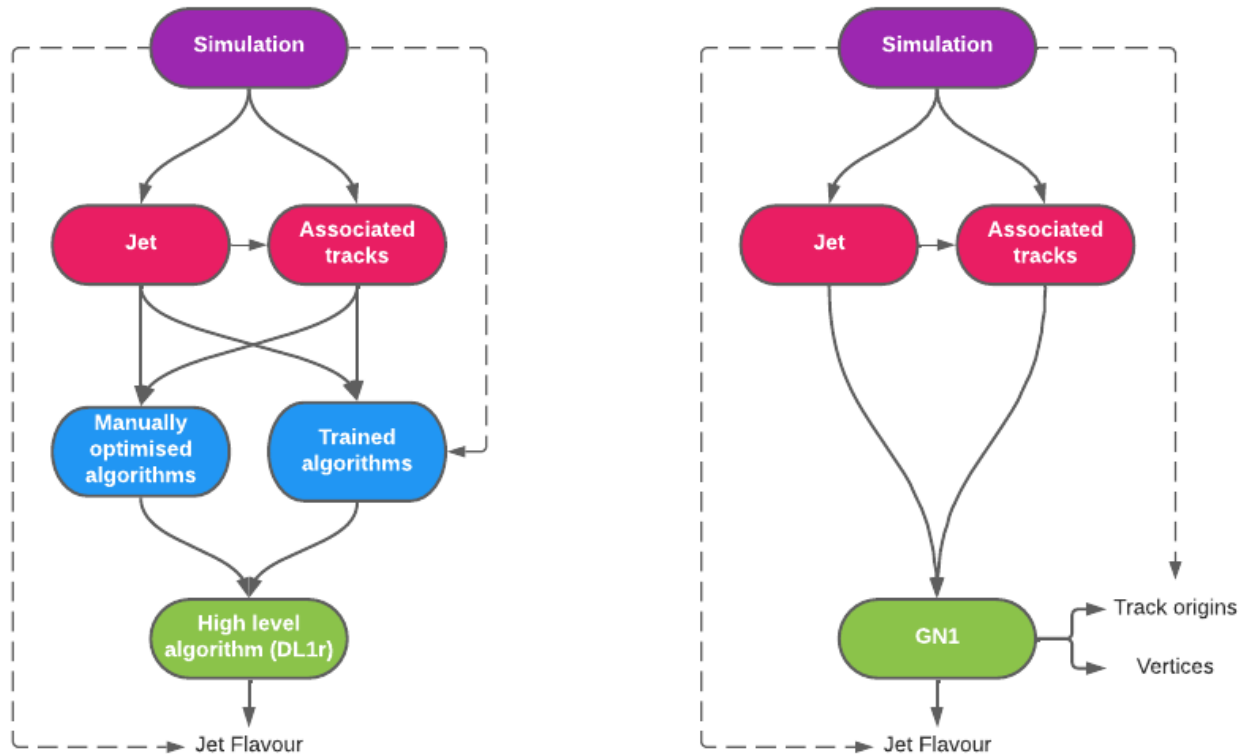
$$O(\{p_1, \dots, p_n\}) = F\left(\sum_{i=1}^n \Phi(p_i)\right)$$

Faster training and improved performances (~ factor 2 bkgd rejection).

<https://cds.cern.ch/record/2718948>

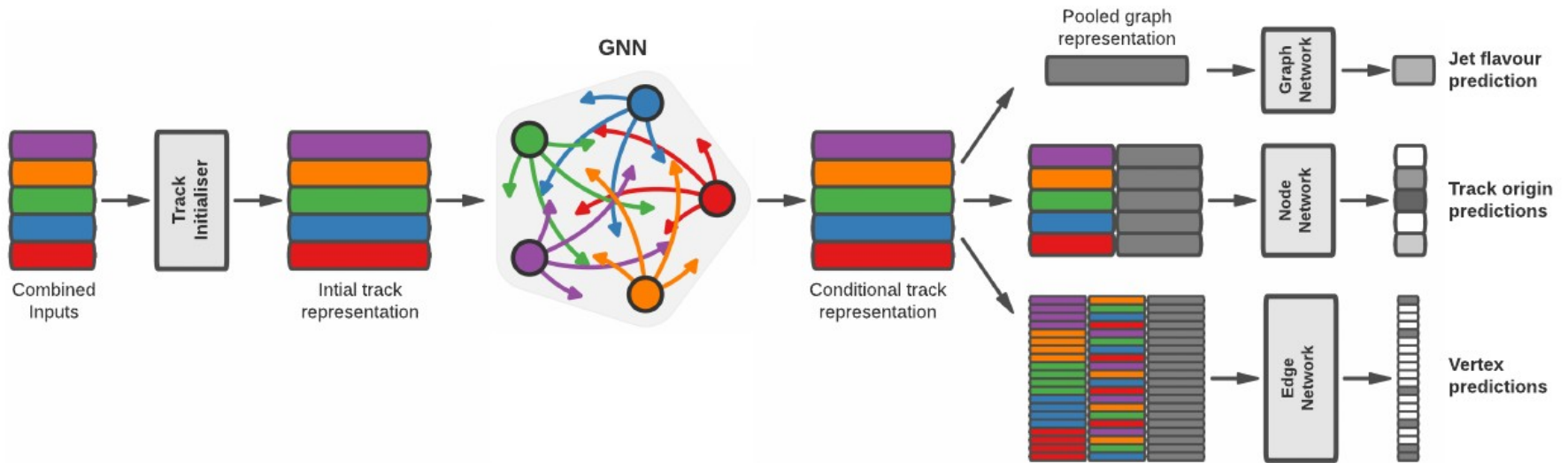
Improved algorithms: GNN

GNN* directly operates on tracks to perform b-tagging. It also performs vertexing and track classification, removing the need for low-level algorithms



* GNN ?? Read e.g. here <https://arxiv.org/abs/2007.13681>

Improved algorithms: GNN

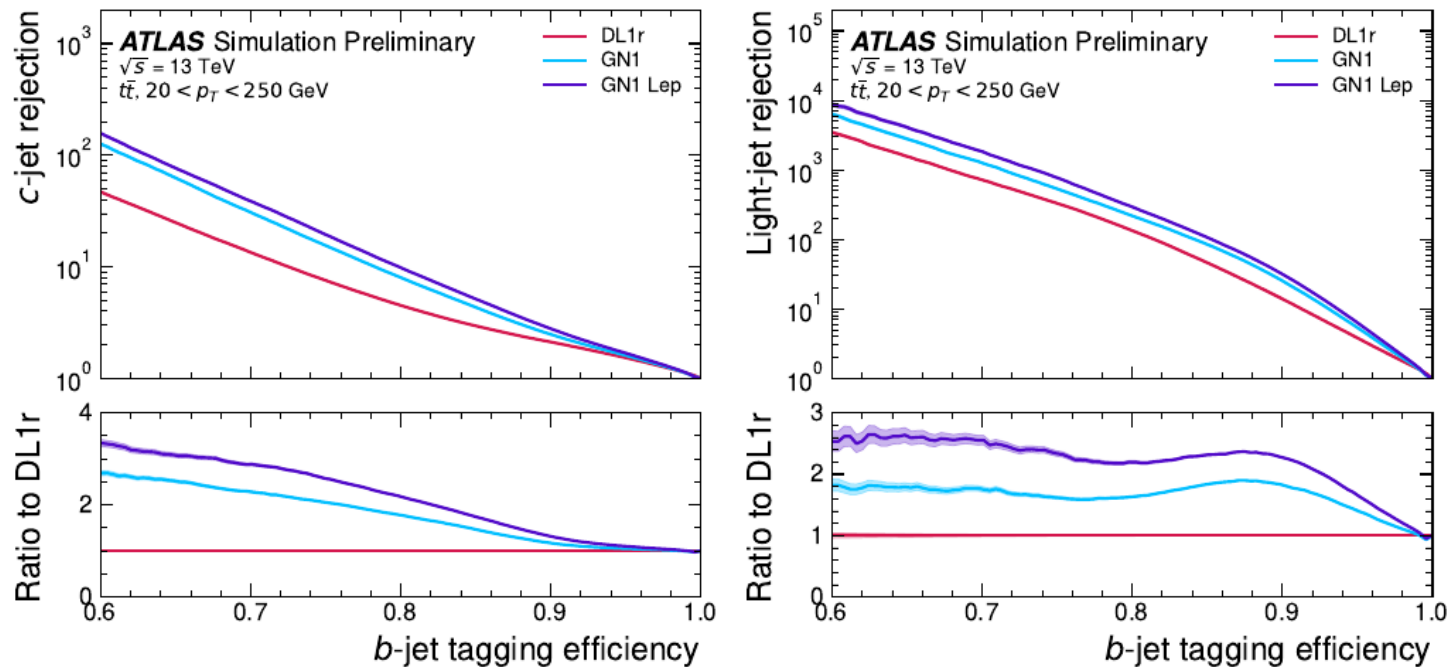


The model combines jet- and track-level information into a combined input, It is then fed into a per-track initialisation network, which outputs a latent representation of each track. These representations are used to populate the node features of a fully connected graph network.

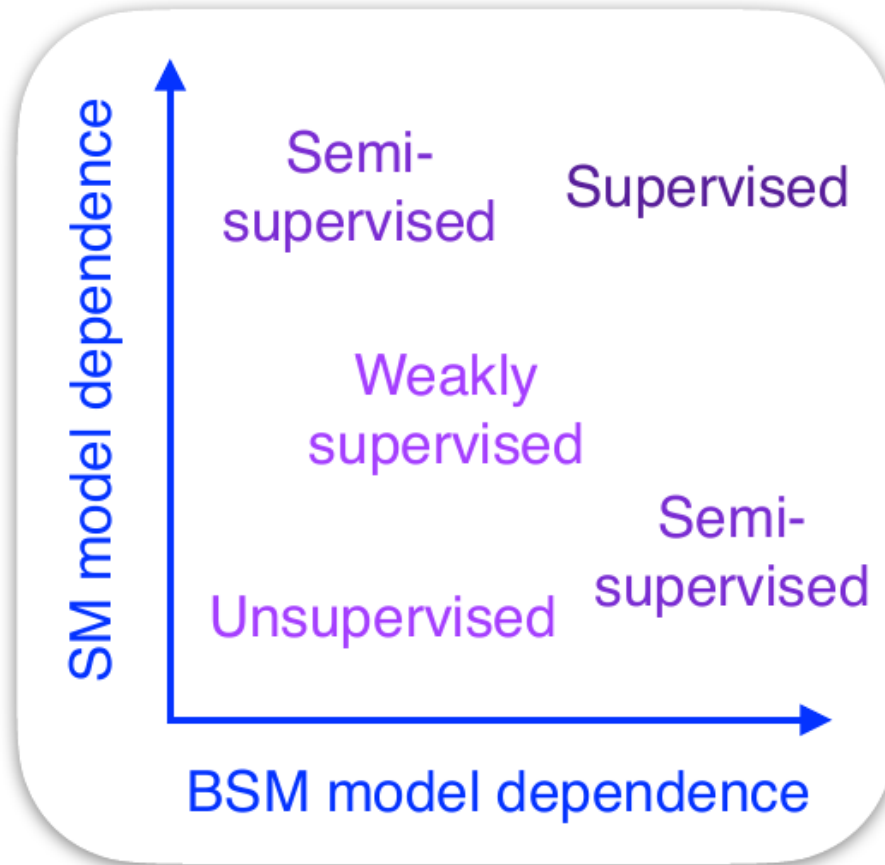
After the graph network, the resulting node representations are used to predict the jet flavour, the track origins, and the track-pair vertex compatibility.

Improved algorithms: GNN

Performance improvement: factor 2-3 with respect to DL1+RNN



Story not over: other algorithmic improvements planned:
GNN with transformers/attention, etc.



<https://arxiv.org/abs/2112.03769>

Anomaly detection as BSM Searches

Search for **unknown signal** (=anomaly) in data

Derive **background** model directly from **data**

Select region of phase-space potentially enriched in signal

Scan **multiple signatures** and variables

	e	μ	τ	q/g	b	t	γ	Z/W	H	BSM \rightarrow SM ₁ \times SM ₁				BSM \rightarrow SM ₁ \times SM ₂			BSM \rightarrow complex			
										q/g	$\gamma/\pi^{0,s}$	b	...	tZ/H	bH	...	$\tau qq'$	eqq'	$\mu qq'$...
e	[37, 38]	[39, 40]	[39]	\emptyset	\emptyset	\emptyset	[41]	[42]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	[43, 44]	\emptyset	
μ		[37, 38]	[39]	\emptyset	\emptyset	\emptyset	[41]	[42]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	[43, 44]	
τ			[45, 46]	\emptyset	[47]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	[48, 49]	\emptyset	
q/g				[29, 30, 50, 51]	[52]	\emptyset	[53, 54]	[55]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	
b					[29, 52, 56]	[57]	[54]	[58]	[59]	\emptyset	\emptyset	\emptyset	\emptyset	[60]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	
t						[61]	\emptyset	[62]	[63]	\emptyset	\emptyset	\emptyset	\emptyset	[64]	[60]	\emptyset	\emptyset	\emptyset	\emptyset	
γ							[65, 66]	[67-69]	[68, 70]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	
Z/W								[71]	[71]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	
H									[72, 73]	[74]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	
BSM \rightarrow SM ₁ \times SM ₁	q/g									\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	
	$\gamma/\pi^{0,s}$										[75]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	
	b											[76, 77]	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	
	\vdots													\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	

Searches for two-body resonances (arXiv:1907.06659)

Two open-dataset challenges fostered many novel ideas for anomaly detection

The LHC Olympics 2020

A Community Challenge for Anomaly
Detection in High Energy Physics



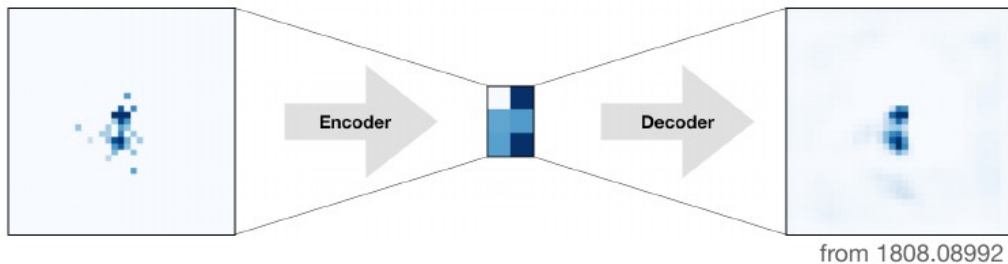
Gregor Kasieczka (ed),¹ Benjamin Nachman (ed),^{2,3} David Shih (ed),⁴ Oz Amram,⁵ Anders Andreassen,⁶ Kees Benkendorfer,^{2,7} Blaz Bortolato,⁸ Gustaaf Brooijmans,⁹ Florencia Canelli,¹⁰ Jack H. Collins,¹¹ Biwei Dai,¹² Felipe F. De Freitas,¹³ Barry M. Dillon,^{8,14} Ioan-Mihail Dinu,⁵ Zhongtian Dong,¹⁵ Julien Donini,¹⁶ Javier Duarte,¹⁷ D. A. Faroughy,¹⁰ Julia Gonski,⁹ Philip Harris,¹⁸ Alan Kahn,⁹ Jernej F. Kamenik,^{8,19} Charanjit K. Khosa,^{20,30} Patrick Komiske,²¹ Luc Le Pottier,^{2,22} Pablo Martín-Ramiro,^{2,23} Andrej Matevc,^{8,19} Eric Metodiev,²¹ Vinicius Mikuni,¹⁰ Inês Ochoa,²⁴ Sang Eon Park,¹⁸ Maurizio Pierini,²⁵ Dylan Rankin,¹⁸ Veronica Sanz,^{20,26} Nilai Sarda,²⁷ Uroš Seljak,^{2,3,12} Aleks Smolkovic,⁸ George Stein,^{2,12} Cristina Mantilla Suarez,⁵ Manuel Szwec,²⁸ Jesse Thaler,²¹ Steven Tsan,¹⁷ Silviu-Marian Udrescu,¹⁸ Louis Vasilin,¹⁶ Jean-Roch Vlimant,²⁹ Daniel Williams,⁹ Mikael Yunus¹⁸

<https://arxiv.org/abs/2101.08320>

The Dark Machines Anomaly Score Challenge:
Benchmark Data and Model Independent Event
Classification for the Large Hadron Collider

T. Aarrestad^a M. van Beekveld^b M. Bona^c A. Boveia^e S. Caron^d J. Davies^c
A. De Simone^{f,g} C. Doglioni^h J. M. Duarteⁱ A. Farbin^j H. Gupta^k L. Hendriks^d
L. Heinrich^a J. Howarth^l P. Jawahar^{m,a} A. Jueidⁿ J. Lastow^h A. Leinweber^o
J. Mamuzic^p E. Merényi^q A. Morandini^r P. Moskvitina^d C. Nellist^d J. Ngadiuba^{s,t}
B. Ostdiek^{u,v} M. Pierini^a B. Ravina^l R. Ruiz de Austri^p S. Sekmen^w
M. Touranakou^{x,a} M. Vaskeviciute^l R. Vilalta^y J.-R. Vlimant^f R. Verheyen^z
M. White^o E. Wulff^h E. Wallin^h K.A. Wozniak^{a,a} Z. Zhang^d

<https://arxiv.org/abs/2105.14027>



Autoencoders

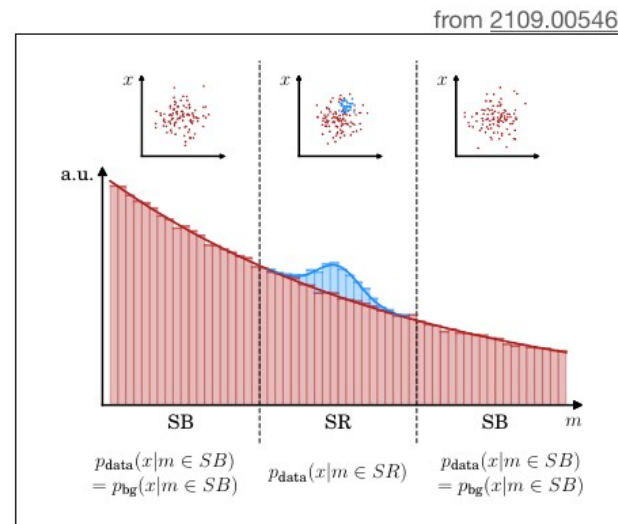
Fully unsupervised

Sensitive to outliers (low $p(x)$)

Farina, Nakai & **DS** [1808.08992](#)

Heimel et al [1808.08979](#)

and many more!!



Enhanced bump hunts

Weakly supervised

Sensitive to overdensities (high $p_{data}(x)/p_{bg}(x)$)

CWoLa Hunting [Collins, Howe & Nachman [1805.02664](#), [1902.02634](#)]

ANODE [Nachman & **DS** [2001.04990](#)]

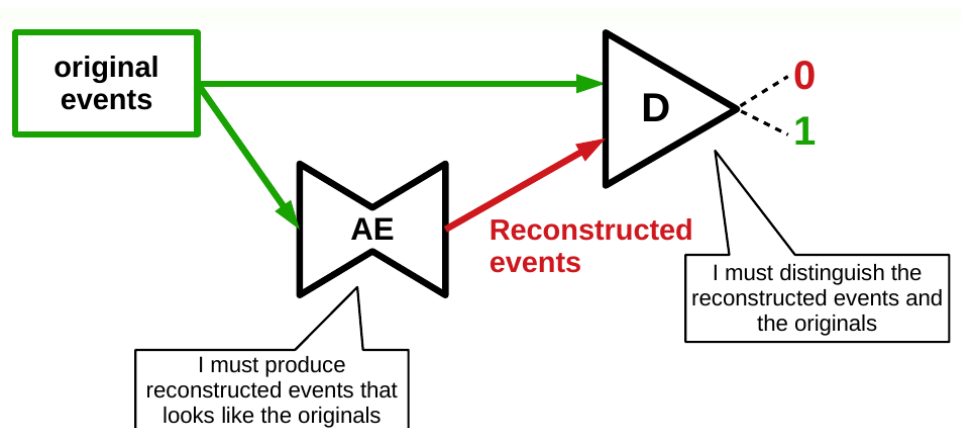
CATHODE [Hallin et al [2109.00546](#)]

CURTAINS [Raine et al [2203.09470](#)]

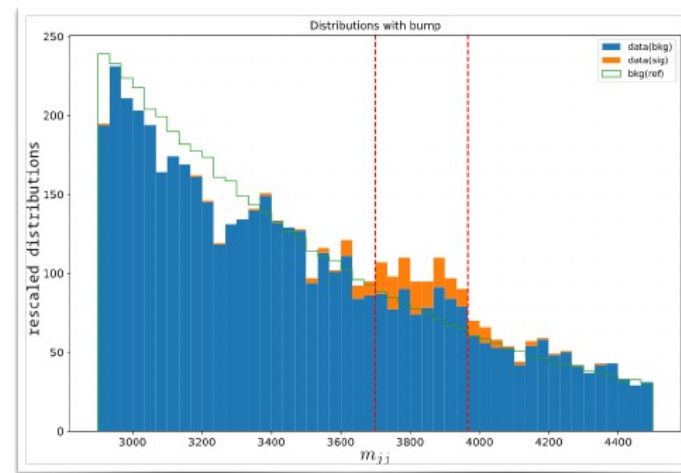
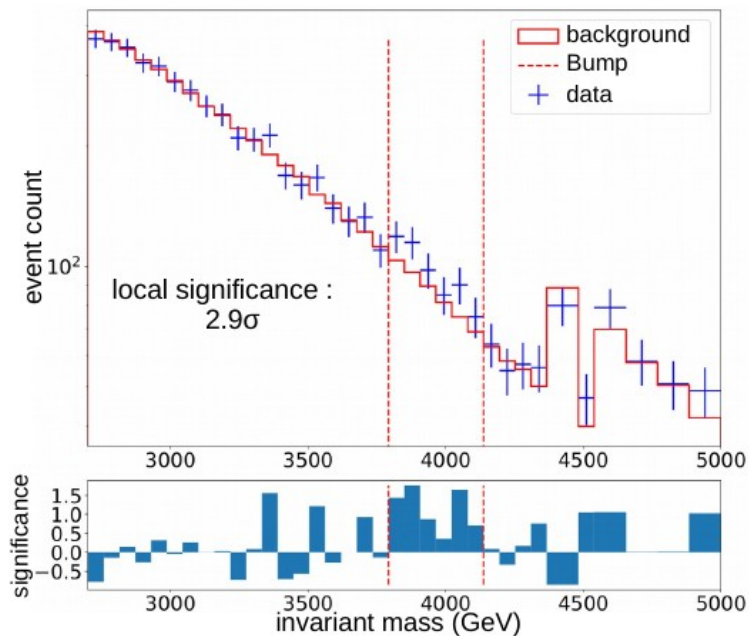
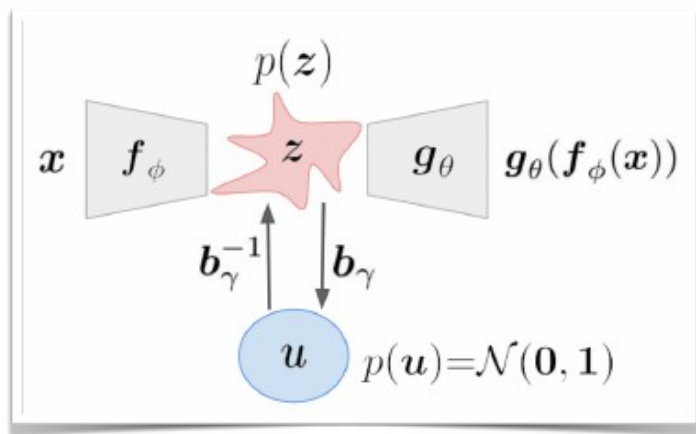
and more...

[slide D. Shih]

Adversarial auto-encoder (L. Valsin et al.)

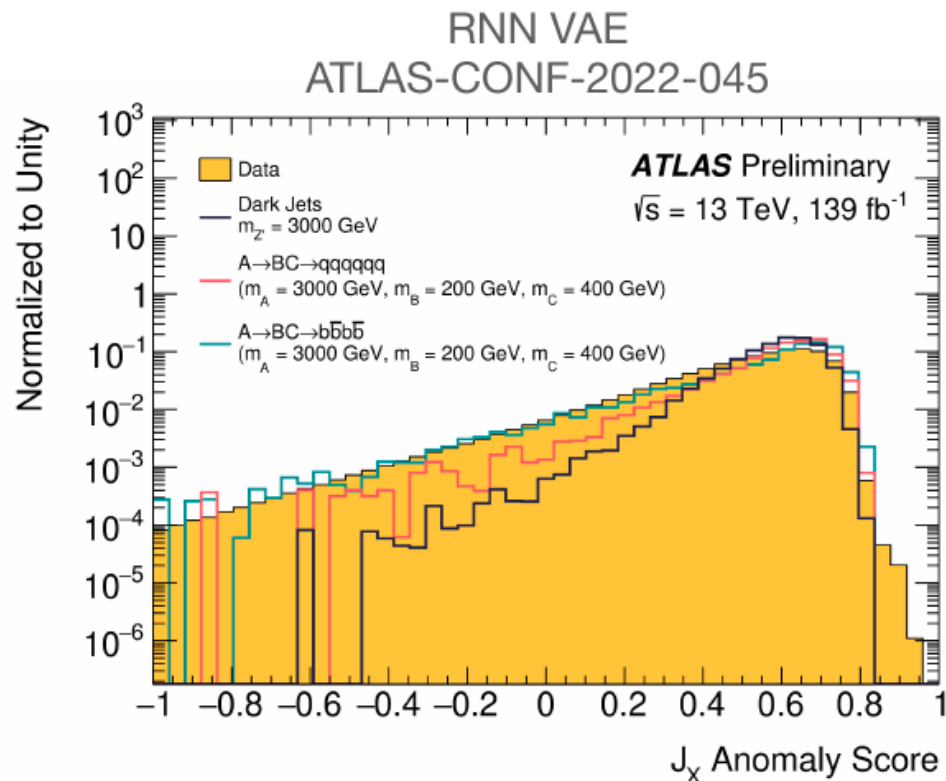
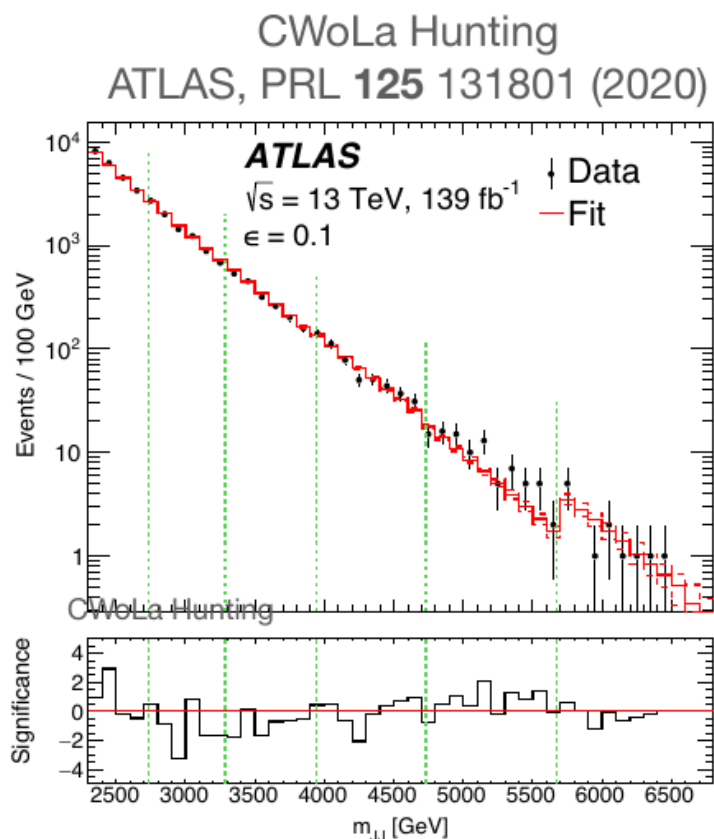


Probabilistic auto-encoder (I. Dinu)



Work documented in community paper [Rep. Prog. Phys. 84 124201](#)

Proof-of-concept are becoming actual LHC searches:

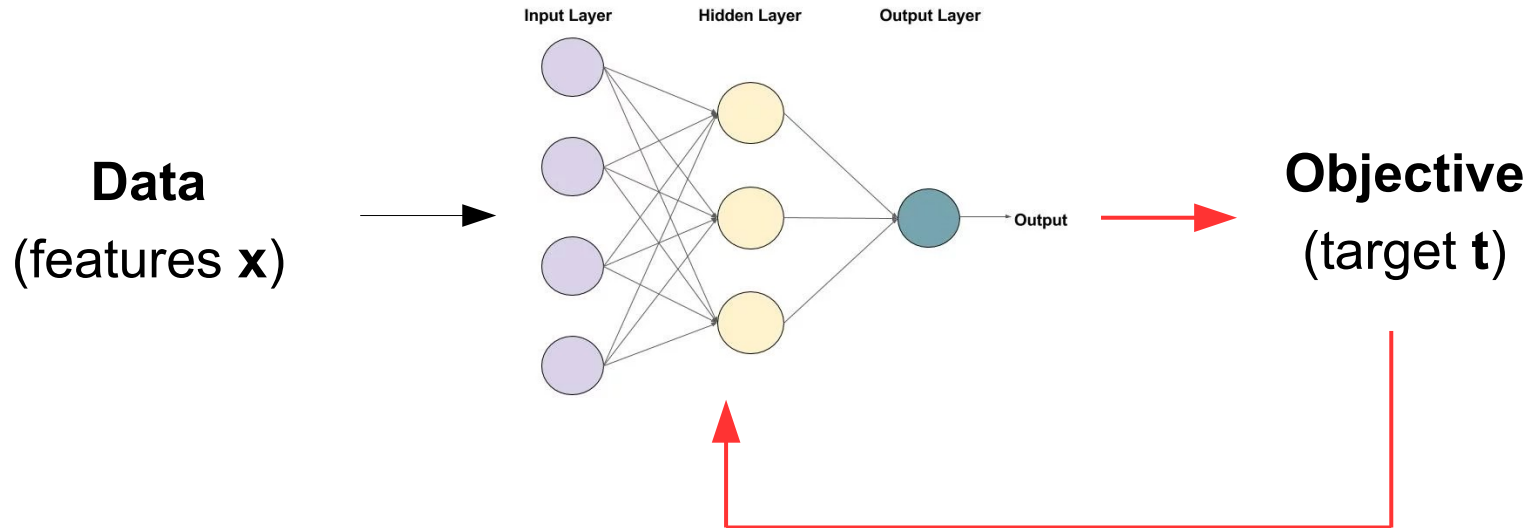


[slide D. Shih]

Automatic differentiation

$$\frac{\partial}{\partial a} \ln f_{a, \sigma^2}(\xi_1) = \frac{(\xi_1 - a)}{\sigma^2} f_{a, \sigma^2}(\xi_1) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(\xi_1 - a)^2}{2\sigma^2}\right)$$
$$\int_{\mathbb{R}_n} T(x) \cdot \frac{\partial}{\partial \theta} f(x, \theta) dx = M\left(T(\xi) \cdot \frac{\partial}{\partial \theta} \ln L(\xi, \theta)\right)$$
$$\int_{\mathbb{R}_n} T(x) \cdot \left(\frac{\partial}{\partial \theta} \ln L(x, \theta)\right) \cdot f(x, \theta) dx = \int_{\mathbb{R}_n} T(x) \cdot \left(\frac{\frac{\partial}{\partial \theta} f(x, \theta)}{f(x, \theta)}\right) \cdot f(x, \theta) dx$$
$$\frac{\partial}{\partial \theta} \int_{\mathbb{R}_n} T(x) f(x, \theta) dx = \int_{\mathbb{R}_n} \frac{\partial}{\partial \theta} T(x) f(x, \theta) dx = \int_{\mathbb{R}_n} \frac{\partial}{\partial \theta} T(x) f(x, \theta) dx$$

Remember this slide ?



Update of weights \mathbf{W}

$$\mathbf{W} \rightarrow \mathbf{W} - \eta \sum_N \frac{\partial \ell(\mathbf{y}, \mathbf{t})}{\partial \mathbf{W}}$$



The **loss** measures how **close** the network **output** is to the **objective**

$$\text{Error} = \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{y}(\mathbf{x}_i, \mathbf{w}), \mathbf{t}_i)$$

Diagram illustrating the loss function formula. The formula is $\text{Error} = \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{y}(\mathbf{x}_i, \mathbf{w}), \mathbf{t}_i)$. Arrows point from labels to parts of the formula: "Training events" points to the summation symbol, "NN output" points to $\mathbf{y}(\mathbf{x}_i, \mathbf{w})$, and "Target value" points to \mathbf{t}_i .

Examples:

$$\left\{ \begin{array}{l} \ell(y_i, t_i) = (\mathbf{y}(\mathbf{x}_i, \mathbf{w}) - \mathbf{t}_i)^2 \\ \ell(y_i, t_i) = t_i \ln(\mathbf{y}(\mathbf{x}_i, \mathbf{w})) + (1 - t_i) \ln(1 - \mathbf{y}(\mathbf{x}_i, \mathbf{w})) \end{array} \right.$$

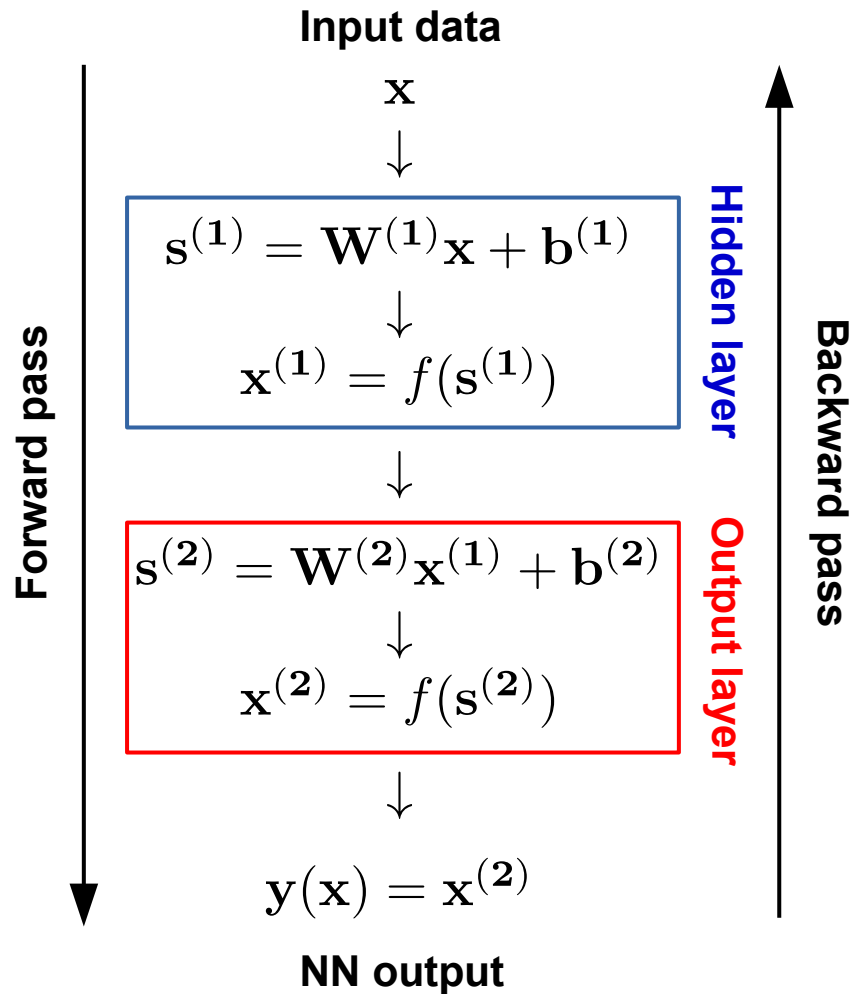
Mean square error

Cross entropy

The loss must be **minimized** → we need to compute its **gradient**

But **y** itself is a function of functions, with a lot of non-linearities → **complex !**

Example NN with 2 layers



Use **chain rule** to compute derivatives of the loss $\ell(\mathbf{y}, \mathbf{t})$

$$\begin{aligned}\frac{\partial \ell}{\partial \mathbf{W}^{(2)}} &= \frac{\partial \ell}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{s}^{(2)}} \frac{\partial \mathbf{s}^{(2)}}{\partial \mathbf{W}^{(2)}} \\ &= \frac{\partial \ell}{\partial \mathbf{y}} \frac{\partial f(\mathbf{s}^{(2)})}{\partial \mathbf{s}^{(2)}} \mathbf{x}^{(1)}\end{aligned}$$

$$\begin{aligned}\frac{\partial \ell}{\partial \mathbf{W}^{(1)}} &= \frac{\partial \ell}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{s}^{(2)}} \frac{\partial \mathbf{s}^{(2)}}{\partial \mathbf{x}^{(1)}} \frac{\partial \mathbf{x}^{(1)}}{\partial \mathbf{s}^{(1)}} \frac{\partial \mathbf{s}^{(1)}}{\partial \mathbf{W}^{(1)}} \\ &= \frac{\partial \ell}{\partial \mathbf{y}} \frac{\partial f(\mathbf{s}^{(2)})}{\partial \mathbf{s}^{(2)}} \frac{\partial \mathbf{s}^{(2)}}{\partial \mathbf{x}^{(1)}} \frac{\partial f(\mathbf{s}^{(1)})}{\partial \mathbf{s}^{(1)}} \mathbf{x}\end{aligned}$$

Automatic (algorithmic) differentiation (AD)

- **Numerical derivative evaluations** rather than derivative **expressions**
- Composition of **operations** for which **derivatives are known**
- **No need to rearrange** the code in a closed-form expression
- Accurate at **machine precision**

For each function a **computational graph** is constructed

→ **evaluation** of the function (forward pass)

→ calculation of **gradient** (backward pass)

Computational graph (example)

Forward pass

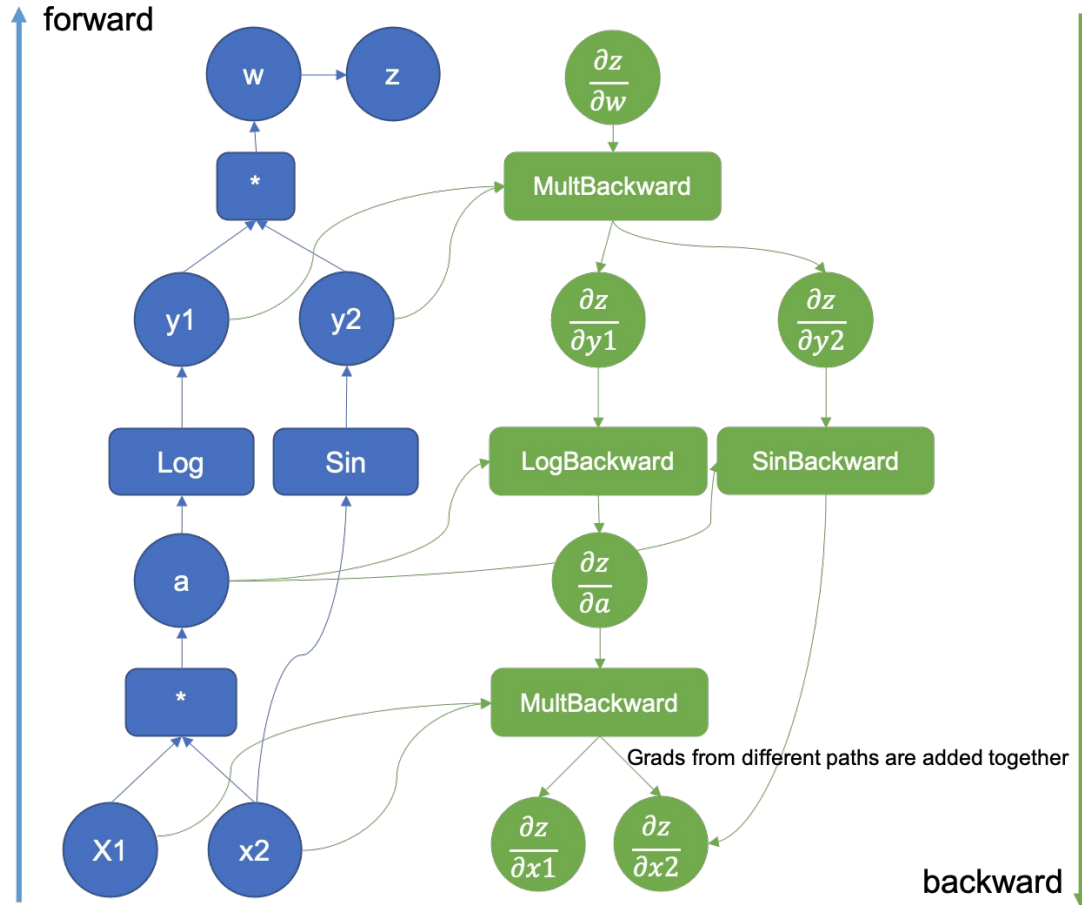
$$\vec{y} = \begin{pmatrix} \log(x_1 x_2) \\ \sin(x_2) \end{pmatrix}$$

$$z = y_1 y_2$$

Backward pass

$$\begin{pmatrix} \frac{\partial z}{\partial x_1} \\ \frac{\partial z}{\partial x_2} \end{pmatrix} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} \end{pmatrix}^T \begin{pmatrix} \frac{\partial z}{\partial y_1} \\ \frac{\partial z}{\partial y_2} \end{pmatrix}$$

$$= \begin{pmatrix} y_2 & x_1 \\ y_2 + y_1 \cos(x_1) & \end{pmatrix}$$



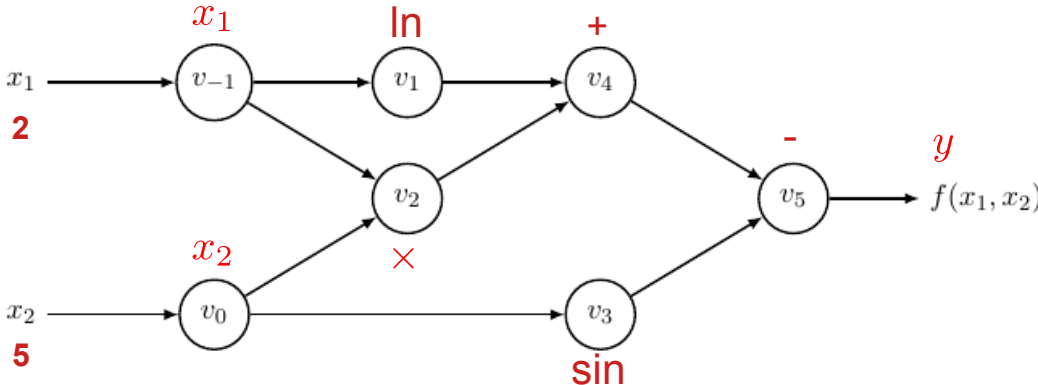
(see: <https://pytorch.org/blog/overview-of-pytorch-autograd-engine/>)

Backpropagation

Example

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$

Propagates derivatives backwards from output $\bar{v}_i = \frac{\partial y}{\partial v_i}$



Forward Primal Trace	Reverse Adjoint (Derivative) Trace
$v_{-1} = x_1 = 2$	$\bar{x}_1 = \bar{v}_{-1} = 5.5$
$v_0 = x_2 = 5$	$\bar{x}_2 = \bar{v}_0 = 1.716$
$v_1 = \ln v_{-1} = \ln 2$	$\bar{v}_{-1} = \bar{v}_{-1} + \bar{v}_1 \frac{\partial v_1}{\partial v_{-1}} = \bar{v}_{-1} + \bar{v}_1 / v_{-1} = 5.5$
$v_2 = v_{-1} \times v_0 = 2 \times 5$	$\bar{v}_0 = \bar{v}_0 + \bar{v}_2 \frac{\partial v_2}{\partial v_0} = \bar{v}_0 + \bar{v}_2 \times v_{-1} = 1.716$
$v_3 = \sin v_0 = \sin 5$	$\bar{v}_{-1} = \bar{v}_2 \frac{\partial v_2}{\partial v_{-1}} = \bar{v}_2 \times v_0 = 5$
$v_4 = v_1 + v_2 = 0.693 + 10$	$\bar{v}_0 = \bar{v}_3 \frac{\partial v_3}{\partial v_0} = \bar{v}_3 \times \cos v_0 = -0.284$
$v_5 = v_4 - v_3 = 10.693 + 0.959$	$\bar{v}_2 = \bar{v}_4 \frac{\partial v_4}{\partial v_2} = \bar{v}_4 \times 1 = 1$
$y = v_5 = 11.652$	$\bar{v}_1 = \bar{v}_4 \frac{\partial v_4}{\partial v_1} = \bar{v}_4 \times 1 = 1$
	$\bar{v}_3 = \bar{v}_5 \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \times (-1) = -1$
	$\bar{v}_4 = \bar{v}_5 \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \times 1 = 1$
	$\bar{v}_5 = \bar{y} = 1$

Reverse mode example, evaluated at $(x_1, x_2) = (2, 5)$. Both $\frac{\partial y}{\partial x_1}$ and $\frac{\partial y}{\partial x_2}$ are computed on the same reverse pass starting from the output

$$\bar{v}_5 = \bar{y} = \frac{\partial y}{\partial y} = 1$$

[1502.05767]

Can we push this idea further ?

Automatic differentiation is used to optimize complex networks

Could we use it to optimize **complex problems** ?

Yes: differentiable programming



Yann LeCun

January 5 · 🌐 (2018)

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

Gradient-based **optimization** methods

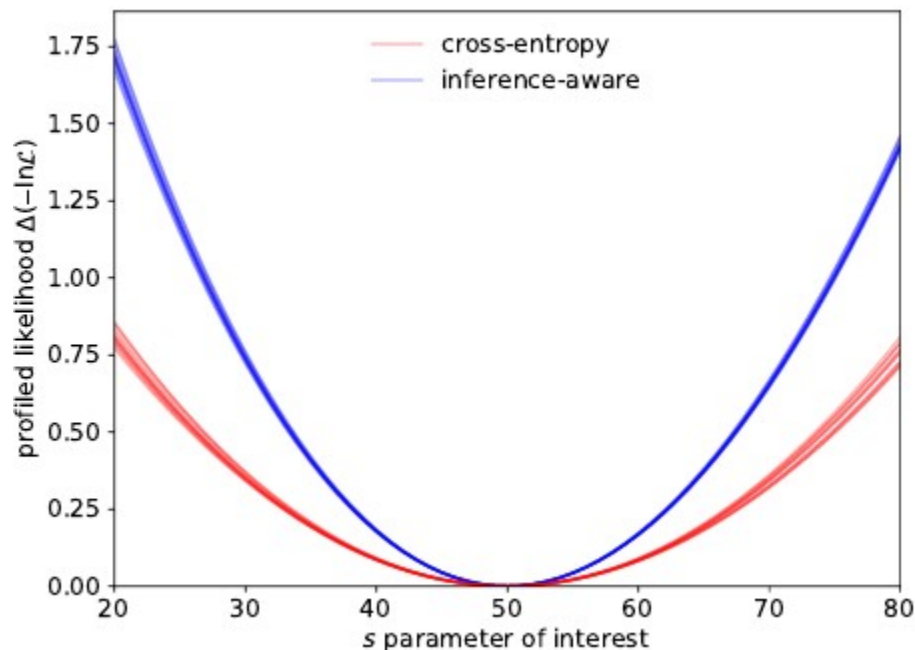
- Code composed of **differentiable and parameterized building blocks**
- **Software** optimized via **automatic differentiation**

ML for Precision Measurements

Inference Aware Neural Optimization

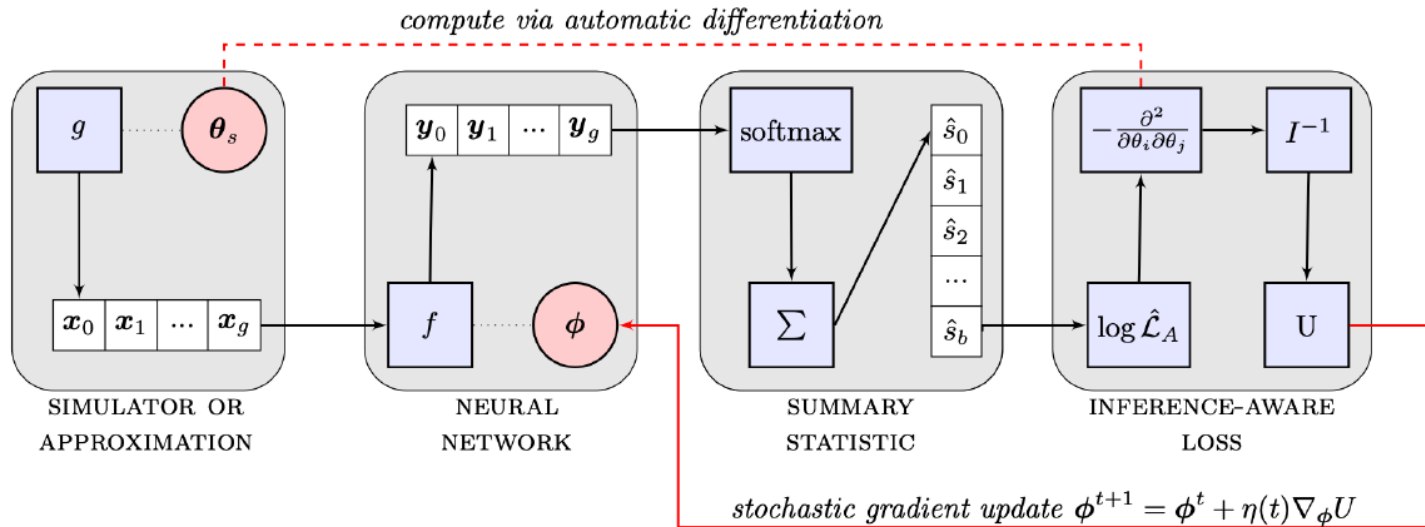
[1806.04743, de Castro, Dorigo]

- **Include nuisance parameters** in the loss function and directly minimize **precision of parameters of interest** (e.g. signal strength measurement)



Profiled likelihood around the expectation value for the parameter of interest for **inference-aware** models and **cross-entropy** loss based models.

INFERNO Algorithm

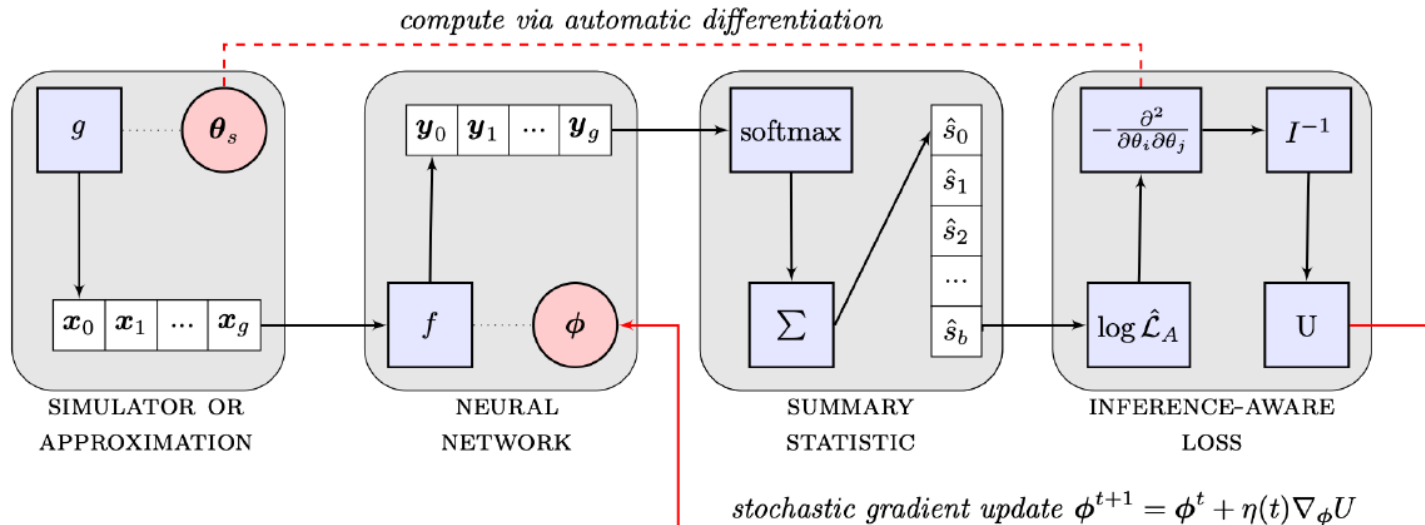


NN output summary statistics from input data

Loss function: **uncertainty** on parameter of interest (U)

Obtained by computing full **hessian** of the **Likelihood** with respect to all **nuisance** parameters

INFERNO Algorithm



Algorithm 1 Inference-Aware Neural Optimisation.

Input 1: differentiable simulator or variational approximation $g(\theta)$.

Input 2: initial parameter values θ_s .

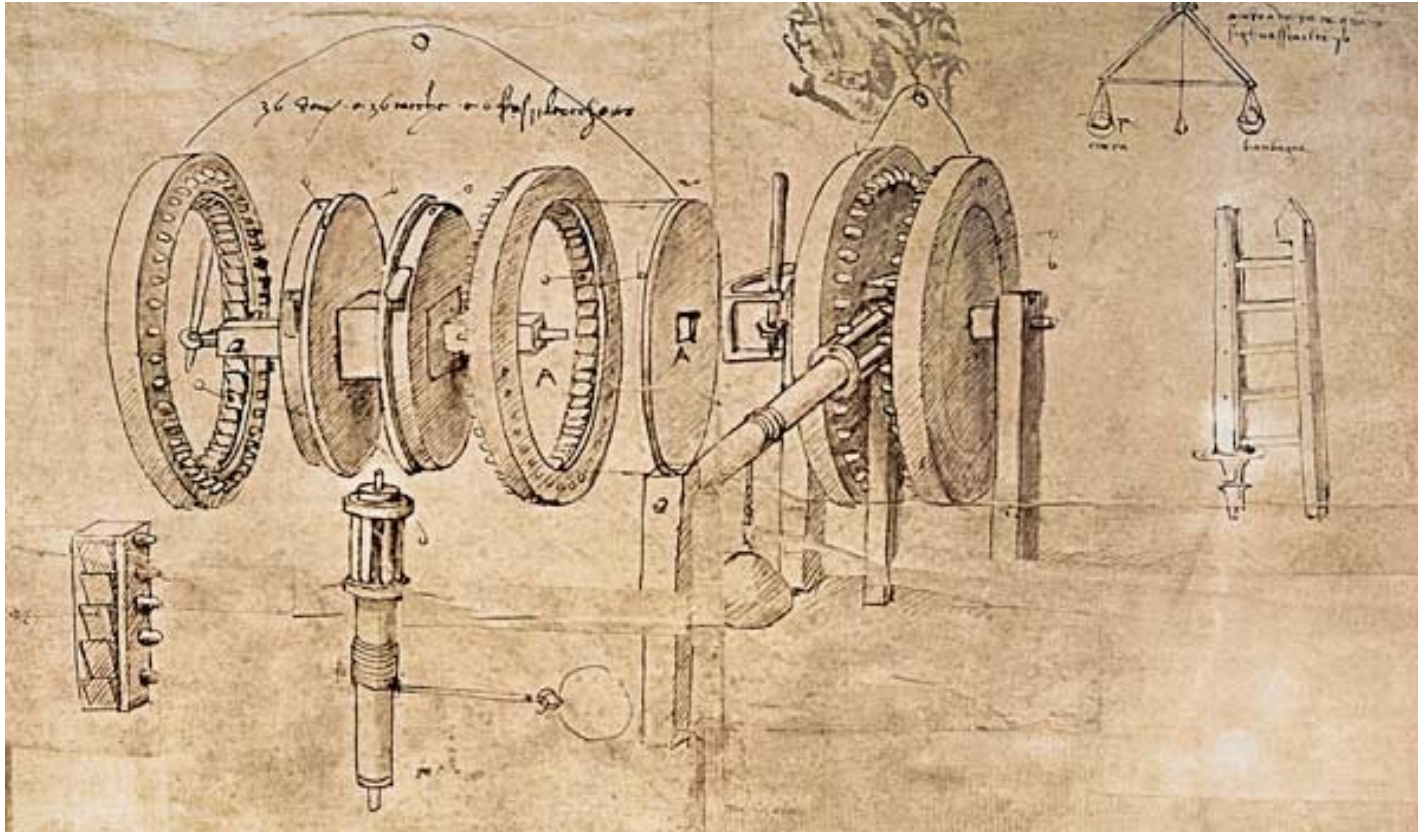
Input 3: parameter of interest $\omega_0 = \theta_k$.

Output: learned summary statistic $s(D; \phi)$.

- 1: **for** $i = 1$ to N **do**
 - 2: Sample a representative mini-batch G_s from $g(\theta_s)$.
 - 3: Compute differentiable summary statistic $\hat{s}(G_s; \phi)$.
 - 4: Construct Asimov likelihood $\mathcal{L}_A(\theta, \phi)$.
 - 5: Get information matrix inverse $I(\theta)^{-1} = \mathbf{H}_{\theta}^{-1}(\log \mathcal{L}_A(\theta, \phi))$.
 - 6: Obtain loss $U = I_{kk}^{-1}(\theta_s)$.
 - 7: Update network parameters $\phi \rightarrow \text{SGD}(\nabla_{\phi} U)$.
 - 8: **end for**
-

For a detailed explanation of the algorithm see [G. Strong blog post](#)

Can automatic differentiation be applied to detector optimization ?



Optimization of Detector Design

Design of detectors traditionally **relies** on individual **optimization of subdetector**

- **Track first, destroy later**

- First detect ionization tracks in tracker, then measure energy deposits from destructive interaction with thick calorimeters

- **Per-subdetector optimization**

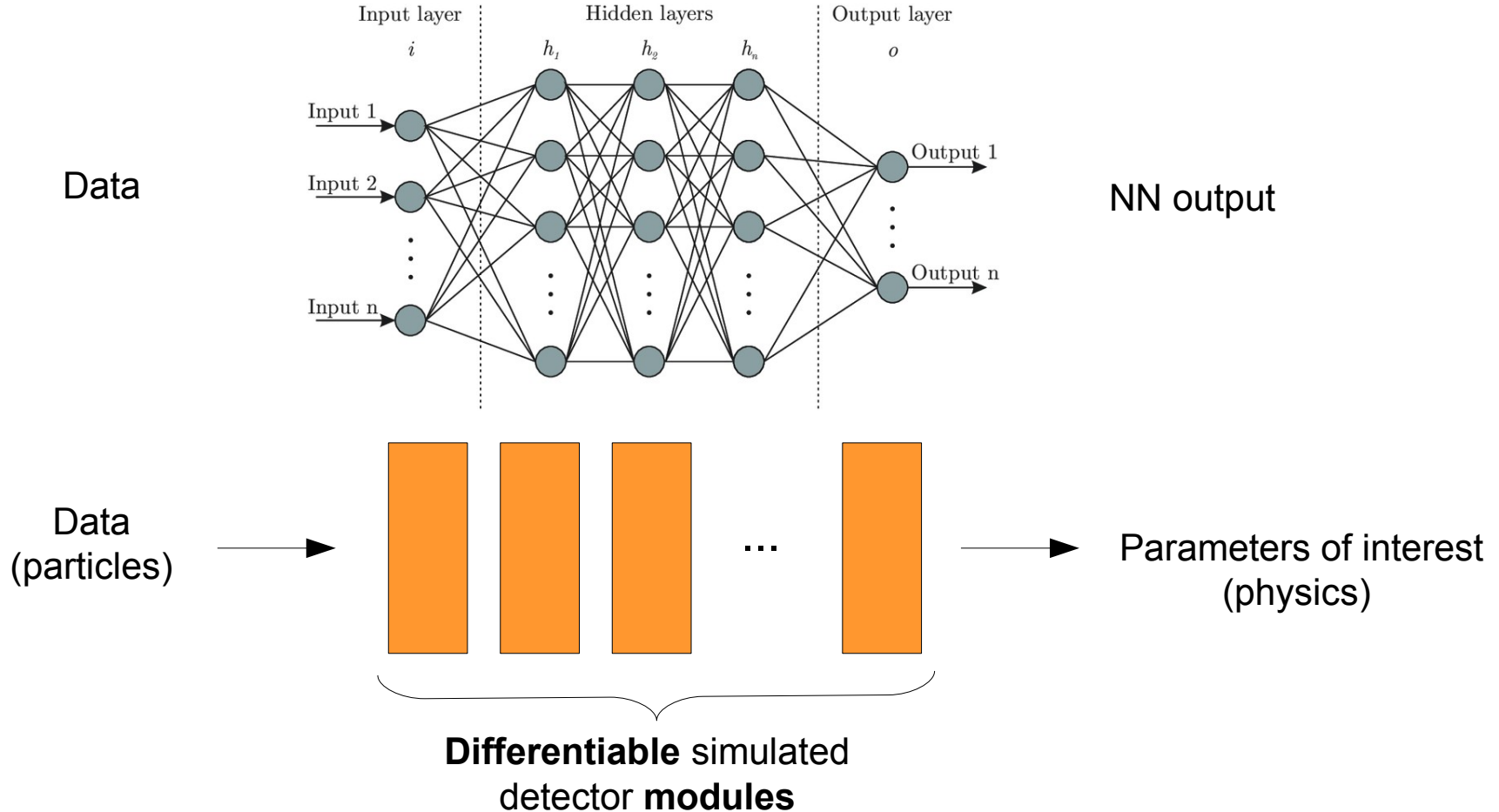
- subdetector-specific figures of merit (e.g. momentum resolution)

- **Impact on physics goals** typically considered in a second step

Optimization of a **joint problem** \neq different from **individual optimization**

$$\operatorname{argmax}_{x,y}(\mathcal{L}(x,y)) \neq \left[\operatorname{argmax}_x(\int \mathcal{L}(x,y)dy), \operatorname{argmax}_y(\int \mathcal{L}(x,y)dx) \right]$$

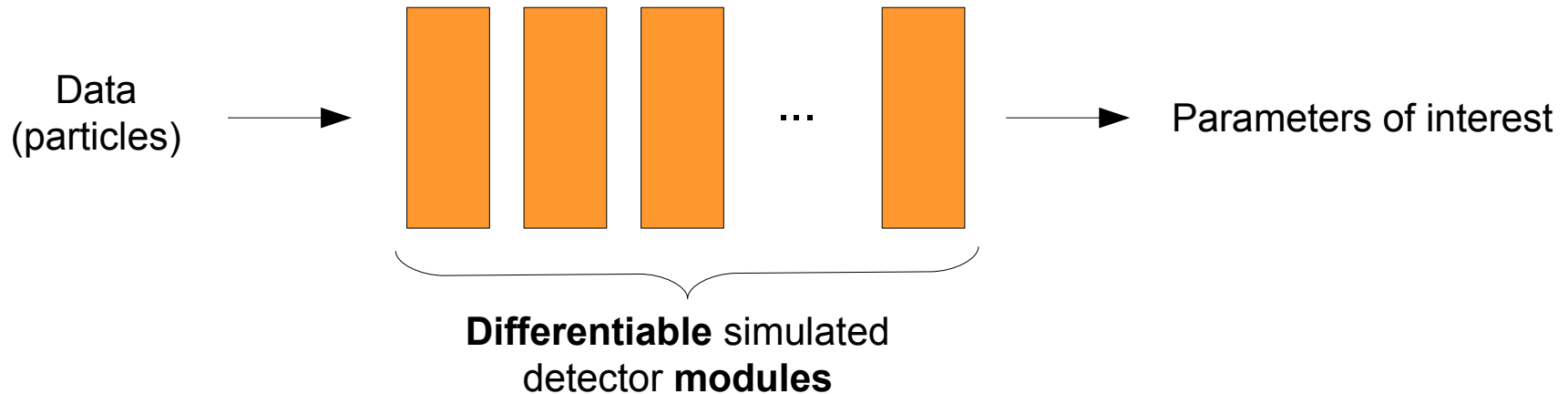
ML for Detector Optimization



Minimization of objective function through automatic differentiation

ML for Detector Optimization

What if simulator is **not differentiable** ? Try differentiable **surrogate models**



GEANT4 → 10^{10} events **SLOW but ACCURATE**

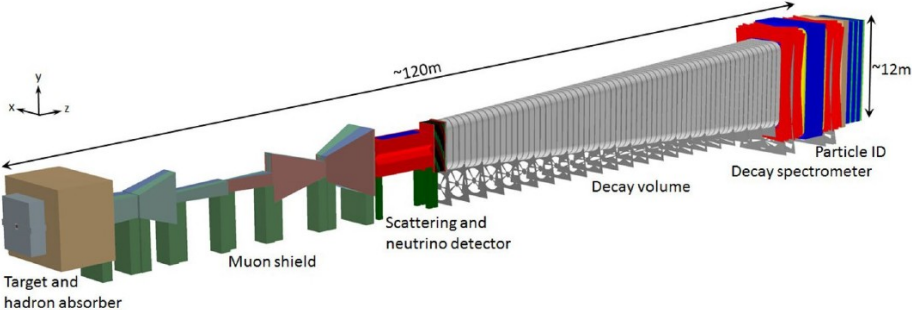
GEANT4 → 10^5 events → Surrogate model → 10^{10} events

(GAN, VAE, Normalizing Flow, ...)
Learn underlying distribution of GEANT4 events

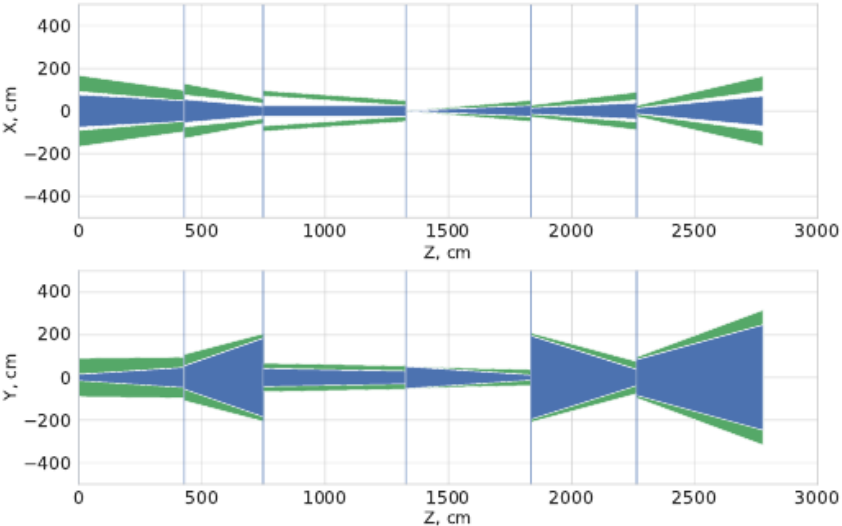
Differentiable

Muon shielding in SHIP

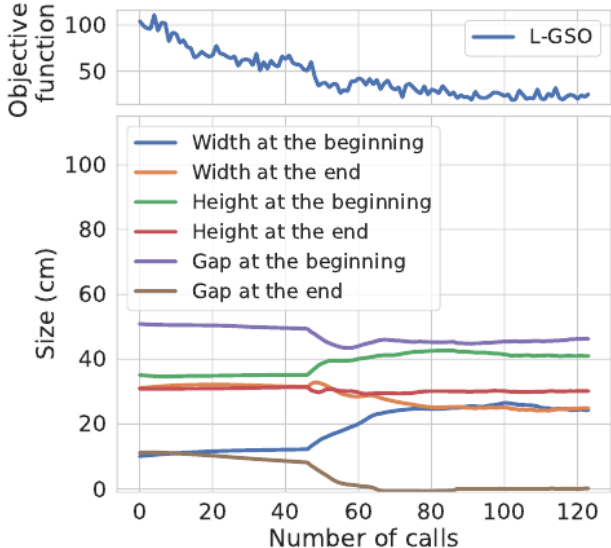
Minimize muon background fluxes in the SHIP steel magnet by varying its geometry



Local generative surrogate solution is shorter and has lower mass than other proposal, hence improving efficacy of the experiment and reducing its cost



Geometry of the magnet
42 parameters to optimize



Evolution of 6 parameters
during optimization

[2002.04632]

Machine-Learning Optimized Design of Experiments

MODE Collaboration

<https://mode-collaboration.github.io>



A. G. Baydin⁵, A. Belias¹⁰, A. Boldyrev⁴, K. Cranmer⁸, P. de Castro Manzano¹, T. Dorigo^{1,14}, C. Delaere², D. Derkach⁴, J. Donini³, F. Fanzago¹, A. Giammanco², C. Glaser¹¹, L. Heinrich¹², J. Kieseler⁷, C. Krause¹³, M. Lagrange², M. Lamparth¹², G. Louppe⁶, L. Layer¹, F. Nardi^{3,14}, P. Martinez Ruiz del Arbol⁹, F. Ratnikov⁴, P. Stowell¹⁵, G. Strong¹, M. Tosi^{1,14}, A. Ustyuzhanin⁴, P. Vischia², H. Yarar¹, H. Zaraket¹⁶

1 INFN, Italy

2 Université Catholique de Louvain, Belgium

3 Université Clermont Auvergne, France

4 Laboratory for big data analysis of the HSE, Russia

5 University of Oxford

6 Université de Liege

7 CERN

8 New York University

9 IFCA, Spain

10 GSI, Germany

11 Uppsala Universitet, Sweden

12 TU Munchen, Germany

13 Rutgers University, US

14 Università di Padova, Italy

15 Durham University, UK

16 Lebanese University, Lebanon

Sponsored by



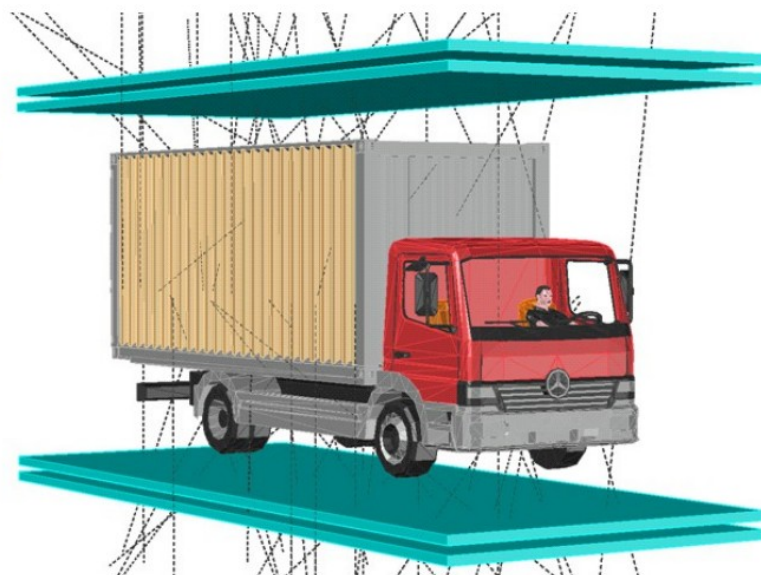
Differentiable programming for muography

Tomography: exploit **atmospheric muon** flux to **map** the interior of **objects**

Muon absorption

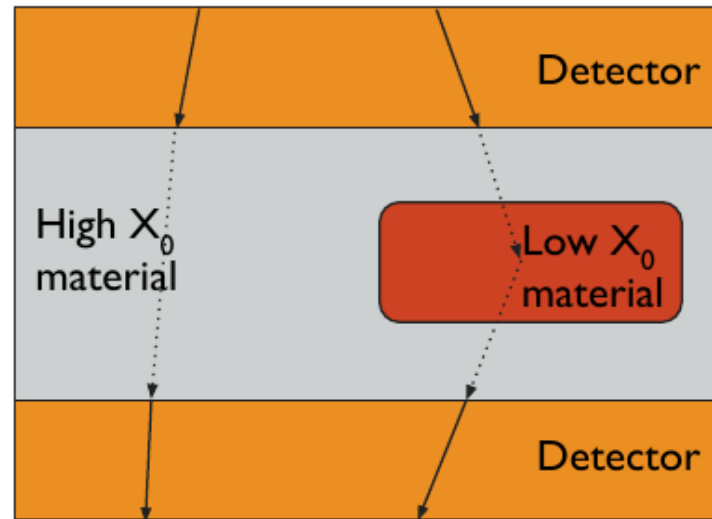


Muon scattering



[images : A. Giammanco]

Volume with unknown **composition** sandwiched between **detectors**



High X_0 = low
scattering

Low X_0 = high
scattering

Infer X_0 (radiation length) of volume by measuring **muon scattering**

How should detectors be positioned for best performances ?

- i.e Muon detection **accuracy**, **resolution** on X_0 , ...
- But also: **cost**, **size**, ...

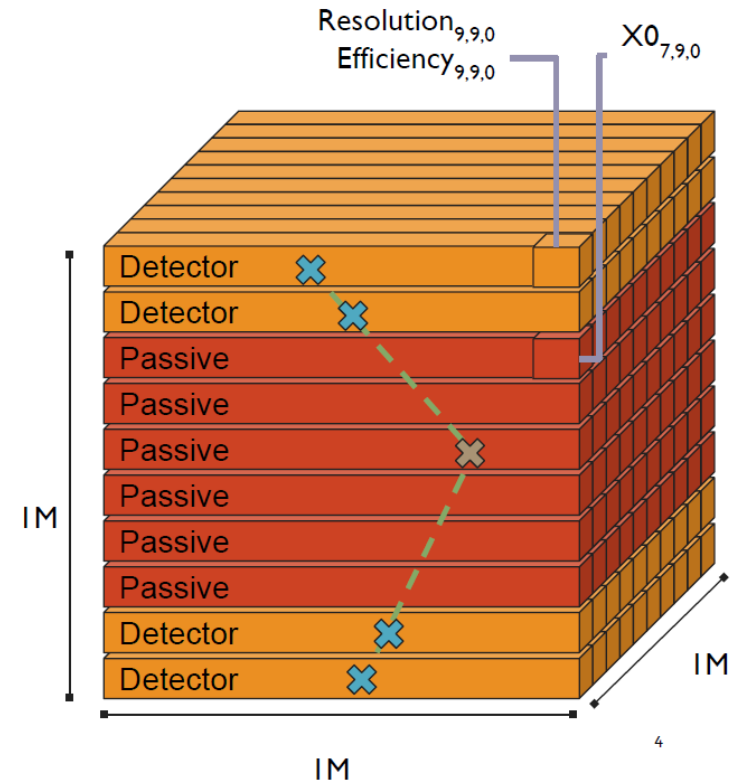
[see G. Strong [talk](#)]

TomOpt: Tomography Optimization

Muon scan of volume of unknown density

Volume of unknown density (e.g. one including a high-Z block of $0.5 \times 0.1 \times 0.1 \text{ m}^3$ somewhere inside a $0.6 \times 1 \times 1 \text{ m}^3$ of low-Z material)

The system «learns» how to compromise cost and precision to optimize the inference on the Z map, and where detector elements are less useful



Python package for differential optimisation of muon-tomography detectors

G.Strong, T.Dorigo, F.Fanzago, A.Giammanco, M.Lagrange, M.Lamparth, F.Nardi, and P.Vischia

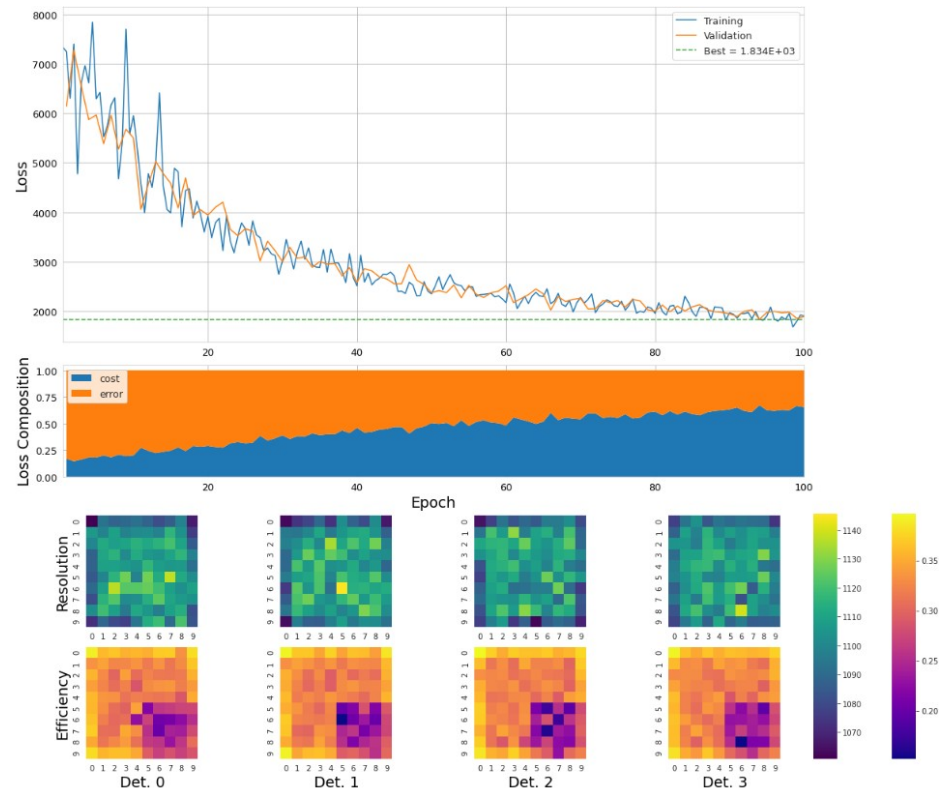
TomOpt: Tomography Optimization

Result of a run of 100 epochs training, followed by a prediction with 100k muons

Shows training of a differentiable model of a schematic muon tomography apparatus.

The **loss** is a combination of **detector cost** (itself a function of sensors efficiency and resolution) and **RMSE on rad length estimate**

Still a long way to go..., but an important milestone



Above, top to bottom: loss, loss composition, resolution map, and efficiency map of detection elements after minimization.

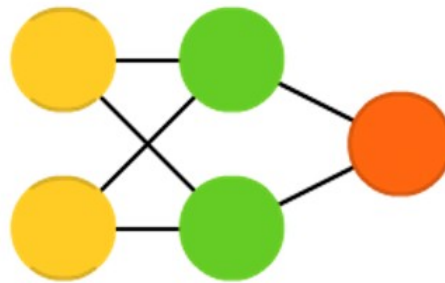
[code and results Giles Strong]

Modern ML is not doing what we did before, but more quickly. It is really doing physics that could not be **possible** otherwise.

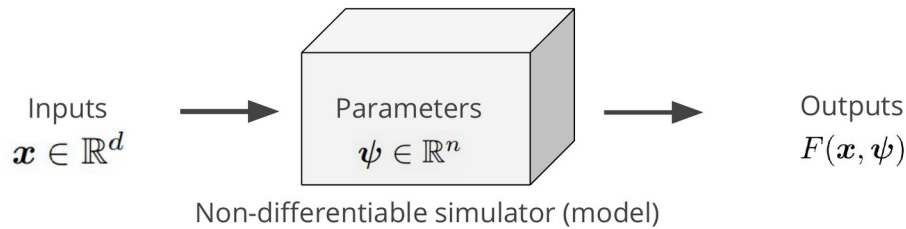
Huge **potential** for new physics searches, triggering, fast simulation, instrumentation, theory, etc.

Very **active** field in HEP in recent years with lots of ideas and developments.


Exciting opportunities for (young) physicists !



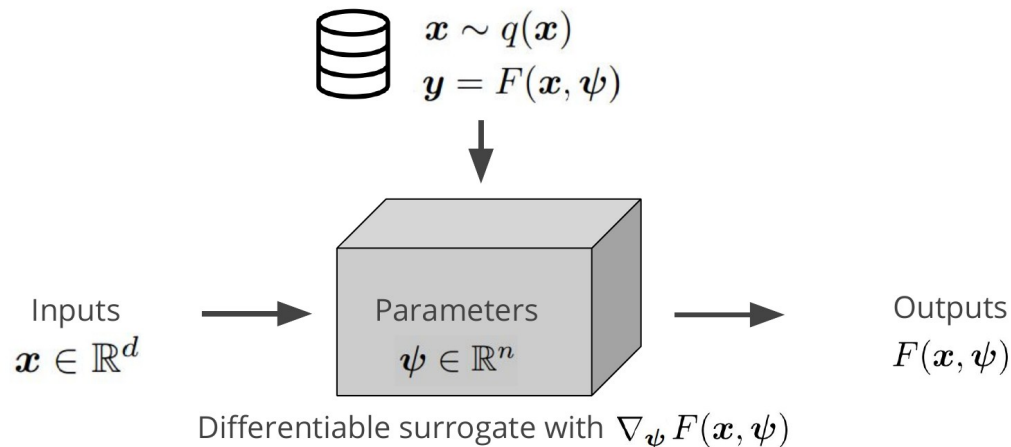
Surrogates for differentiability



- Run simulator many times
- Generate a (large) dataset of input - output pairs capturing simulator's behavior


$$\mathbf{x} \sim q(\mathbf{x})$$
$$\mathbf{y} = F(\mathbf{x}, \psi)$$

- Use the dataset to learn a differentiable approximation of the simulator (e.g., a deep generative model)



[slides G. Baydin]

Surrogates for differentiability

Algorithm 1 Local Generative Surrogate Optimization (L-GSO) procedure

Require: number N of ψ , number M of \mathbf{x} for surrogate training, number K of \mathbf{x} for ψ optimization step, trust region U_ϵ , size of the neighborhood ϵ , Euclidean distance d

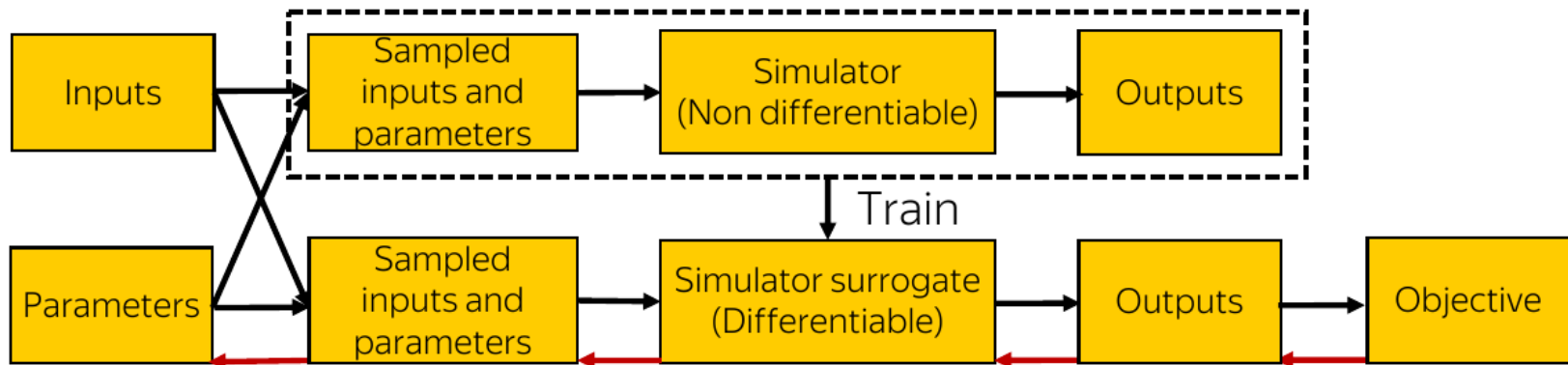
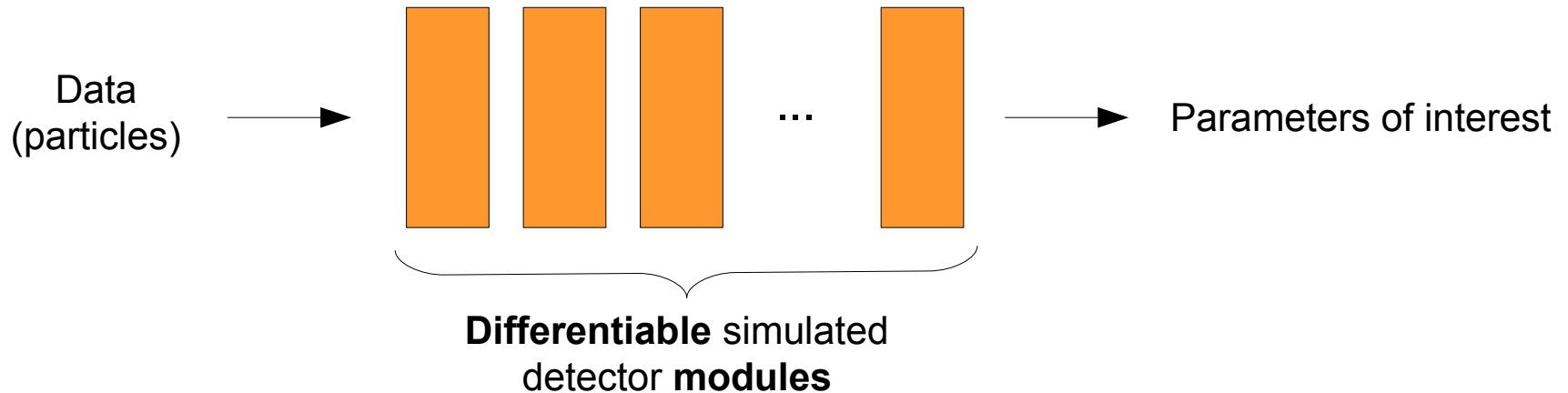
- 1: Choose initial parameter ψ
 - 2: **while** ψ has not converged **do**
 - 3: Sample ψ_i in the region U_ϵ^ψ , $i = 1, \dots, N$
 - 4: For each ψ_i , sample inputs $\{\mathbf{x}_j^i\}_{j=1}^M \sim q(\mathbf{x})$
 - 5: Sample $M \times N$ training examples from simulator $\mathbf{y}_{ij} = F(\mathbf{x}_j^i; \psi_i)$
 - 6: Store $\mathbf{y}_{ij}, \mathbf{x}_j^i, \psi_i$ in history H
 $i = 1, \dots, N; j = 1, \dots, M$
 - 7: Extract all $\mathbf{y}_l, \mathbf{x}_l, \psi_l$ from history H ,
iff $d(\psi, \psi_l) < \epsilon$
 - 8: Train generative surrogate model
 $S_\theta(\mathbf{z}_l, \mathbf{x}_l; \psi_l)$, where $\mathbf{z}_l \sim \mathcal{N}(0, 1)$
 - 9: Fix weights of the surrogate model θ
 - 10: Sample $\bar{\mathbf{y}}_k = S_\theta(\mathbf{z}_k, \mathbf{x}_k; \psi)$, $\mathbf{z}_k \sim \mathcal{N}(0, 1)$,
 $\mathbf{x}_k \sim q(\mathbf{x})$, $k = 1, \dots, K$
 - 11: $\nabla_\psi \mathbb{E}[\mathcal{R}(\bar{\mathbf{y}})] \leftarrow \frac{1}{K} \sum_{k=1}^K \frac{\partial \mathcal{R}}{\partial \bar{\mathbf{y}}_k} \frac{\partial S_\theta(\mathbf{z}_k, \mathbf{x}_k; \psi)}{\partial \psi}$
 - 12: $\psi \leftarrow \text{SGD}(\psi, \nabla_\psi \mathbb{E}[\mathcal{R}(\bar{\mathbf{y}})])$
 - 13: **end while**
-

$$\begin{aligned} \psi^* &= \arg \min_{\psi} \mathbb{E}[\mathcal{R}(\mathbf{y})] = \arg \min_{\psi} \int \mathcal{R}(\mathbf{y}) p(\mathbf{y}|\mathbf{x}; \psi) q(\mathbf{x}) d\mathbf{x} d\mathbf{y} \\ &\approx \arg \min_{\psi} \frac{1}{N} \sum_{i=1}^N \mathcal{R}(F(\mathbf{x}_i; \psi)) \end{aligned}$$

$$\nabla_{\psi} \mathbb{E}[\mathcal{R}(\mathbf{y})] \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\psi} \mathcal{R}(S_{\theta}(\mathbf{z}_i, \mathbf{x}_i; \psi))$$

ML for Detector Optimization

What if simulator is **not differentiable** ? Try differentiable **surrogate models**



Black-Box Optimization with Local Generative Surrogates, S. Shirobokov, V. Belavin, M. Kagan, A. Ustyuzhanin, A. G. Baydin, <https://arxiv.org/abs/2002.04632>

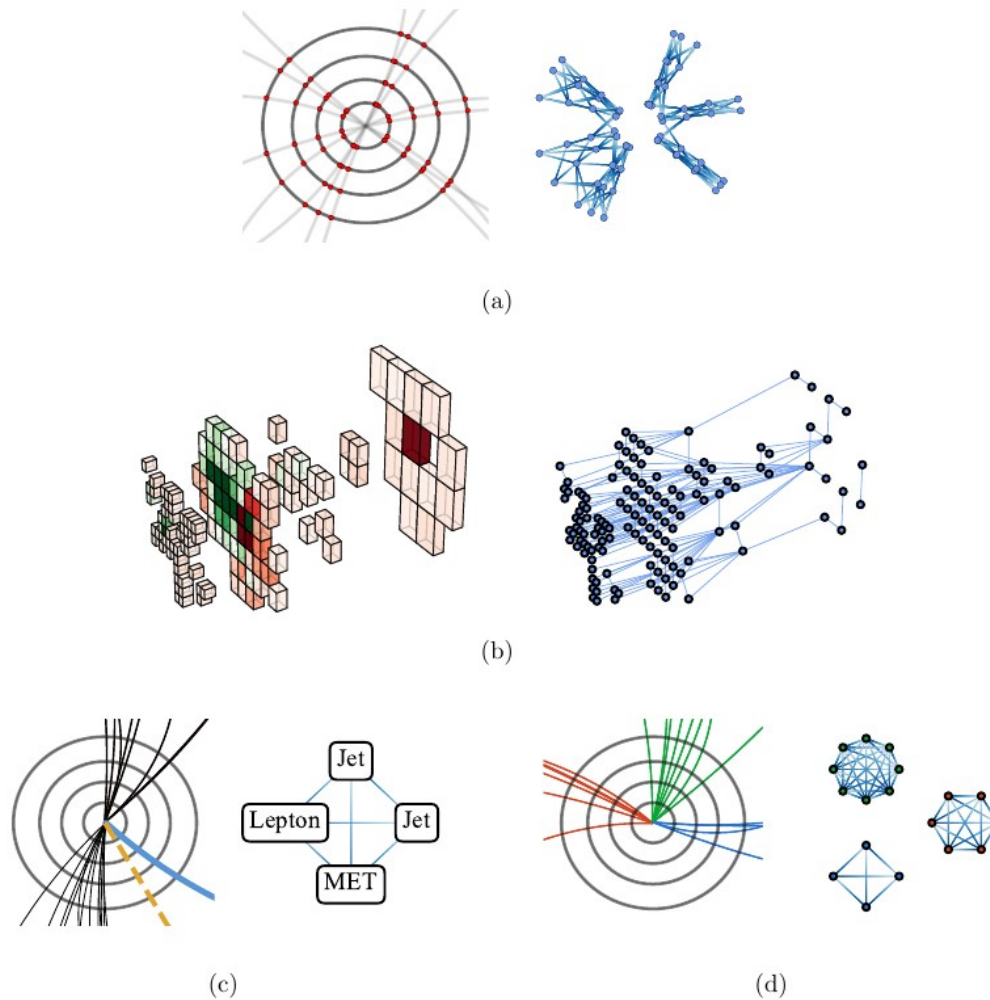


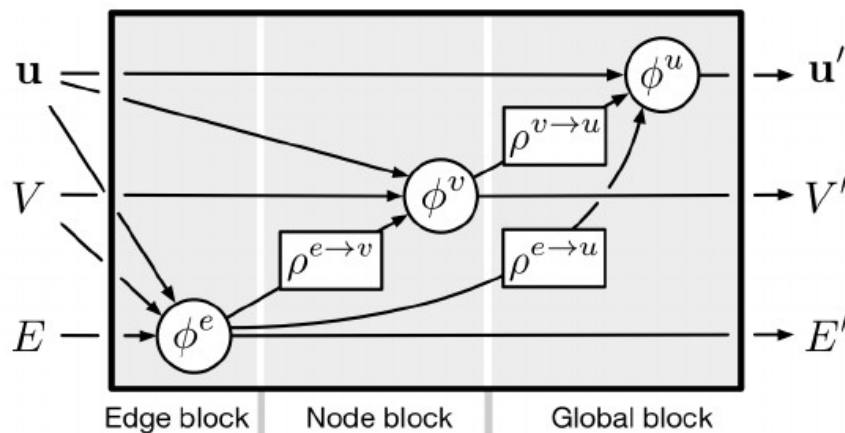
Figure 2. HEP data lend itself to being represented as a graph for many applications: (a) clustering tracking detector hits into tracks, (b) segmenting calorimeter cells, (c) classifying events with multiple types of physics objects, (d) jet classification based on the particles associated to the jet.

A graph can be represented by, $G = (u, V, E)$, with N_v vertices and N_e edges. The u represents graph-level attributes. The set of nodes is V and the set of edges is E .

A GN's stages of processing are as follows.

$\mathbf{e}'_k = \phi^e(\mathbf{e}_k, \mathbf{v}_{r_k}, \mathbf{v}_{s_k}, \mathbf{u})$	$\bar{\mathbf{e}}'_i = \rho^{e \rightarrow v}(E'_i)$	▷ Edge block
$\mathbf{v}'_i = \phi^v(\bar{\mathbf{e}}'_i, \mathbf{v}_i, \mathbf{u})$	$\bar{\mathbf{e}}' = \rho^{e \rightarrow u}(E')$	▷ Vertex block
$\mathbf{u}' = \phi^u(\bar{\mathbf{e}}', \bar{\mathbf{v}}', \mathbf{u})$	$\bar{\mathbf{v}}' = \rho^{v \rightarrow u}(V')$	▷ Global block

A GN block contains 6 internal functions: 3 *update functions* (ϕ^e , ϕ^v , and ϕ^u) and 3 *aggregation functions* ($\rho^{e \rightarrow v}$, $\rho^{e \rightarrow u}$, and $\rho^{v \rightarrow u}$).



[arXiv:2007.13681]