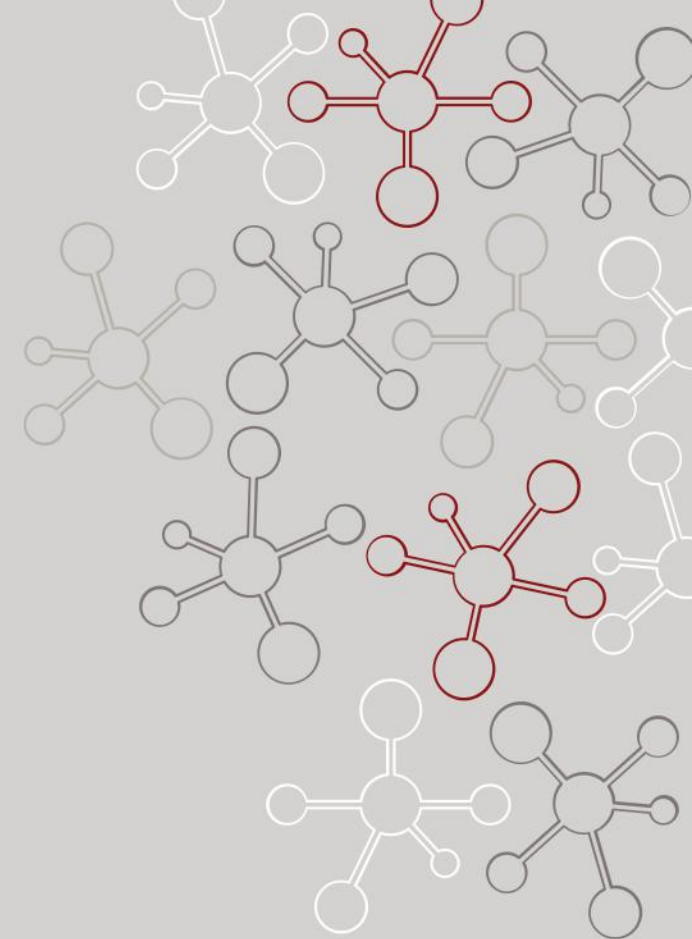# Using Gradients to Get More Out of High Energy Physics

Michael Kagan, SLAC
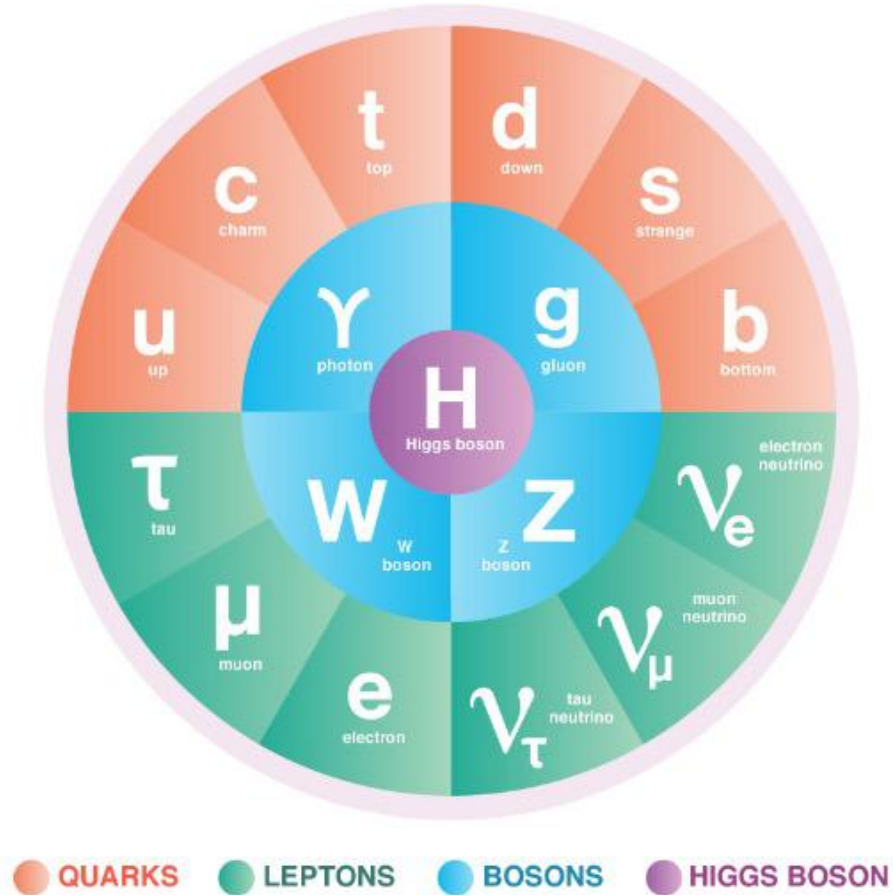
*CMU STAMPS Seminar*
November 10, 2023

Image source: Symmetry Magazine

# High Energy Physics – Big Questions

## Why is the Higgs so light?

Image source: Symmetry Magazine

## What is Dark Matter?
## What is Dark Energy?

Image source: NASA/CXC/CFA/ M.MARKEVITCH

## Why is there more matter than anti-matter in the universe?

Image source: Symmetry Magazine

# Studying Physics at the Smallest Scales

The Large Hadron Collider at CERN

LHCb

ATLAS

CMS

ALICE

$E = mc^2$

$pp(qq) \rightarrow X \rightarrow M\,M$

27 km ring, 100m underground



46m

25m

ATLAS

Run: 303079
Event: 197351611
2016-07-01 05:01:26 CEST

# Studying Collisions



muon

Bottom Quark

Bottom Quark

neutrino

electron

ATLAS EXPERIMENT

Run Number: 160958, Event Number: 9038972

Date: 2010-08-08 11:01:12 BST

O(20) Fundamental physics parameters $\theta$

O(10) particles

O(100) particles

O($10^8$) detector elements

Deep knowledge of data generation process

**Likelihood intractable,
but can simulate with high-fidelity simulators**

Data        Parameters of interest

$$p(x|\theta) = \int p(x|z)p(z|\theta)dz$$

Latent variables
O($10^8$) Dimensional

# Data Analysis "Inverts" this Process

O(20) Fundamental physics parameters $\theta$

O(10) particles

O(100) particles

O($10^8$) detector elements

1909.02845

Statistical Inference

Reconstruct & Select interesting events

Reconstruct particles

# Data Analysis Workflow

Experiment

$\theta$

Simulation

$\phi$

Data Analysis

$\phi$

Result $\hat{\theta}$

**Summarize**: Reduce 100M → 1 informative number

**Statistical Inference**: Compare simulation & data

EPJC 80 (2020) 942

ATLAS
H → ZZ* → 4l
√s = 13 TeV, 139 fb⁻¹

- Data
- Higgs (125 GeV)
- Z(Z*)
- tXX, VVV
- Z+jets, tt̄
- Uncertainty

Events/2.5 GeV

$m_{4l}$ [GeV]

Standard Model Production Cross Section Measurements — Status: February 2022

ATLAS Preliminary, $\sqrt{s} = 5,7,8,13$ TeV

# But we want to get the most out of our data!



Improve Confidence Intervals

$\theta_2$

$\theta_1$



Which collider? Which Detectors?

$$\min_{\phi} \mathbb{E}[f(x,\phi)] = \min_{\phi} \int f(x,\phi) p_{\phi}(x|\theta) dx$$

$$\min_{\phi} \mathbb{E}[f(x,\phi)] = \min_{\phi} \int f(x,\phi)p_{\phi}(x|\theta)dx$$

Statistical Analysis / Design Objective

Realizations of measurements:

E.g. Simulations

Analysis params / Design params

Probability to see a measurement, e.g.
- Scattering prob.
- Detection prob.

# What are we optimizing

$$\min_\phi \mathbb{E}[f(x, \phi)] = \min_\phi \int f(x, \phi) p_\phi(x|\theta) dx$$

$$\approx \min_\phi \frac{1}{N} \sum_{x_i \sim p_\phi(x|\theta)} f(x_i, \phi)$$

**Deep learning looks very similar**, optimizing an objective over parameters of a model using set of samples

**Stochastic gradient descent (SGD):** go-to optimization algorithm for training deep neural networks with even $O(10^{11})$ parameters

$$\theta \leftarrow \theta - \nabla_\theta L(x, \theta)$$



To deal with hyper-planes in a 14-dimensional space, visualize a 3D space and say 'fourteen' to yourself very loudly.
-G. Hinton

# Differentiable Programming

Derivatives for gradient-based optimization come from running **differentiable code** via **automatic differentiation (AD)**

Image credit: L. Heinrich



Image credit: Wikipedia

$f(x_1, x_2)$

$\bar{f} = \bar{w}_5 = 1$ (seed)

$w_5$

$\bar{w}_4 = \bar{w}_5 \frac{\partial w_5}{\partial w_4} = \bar{w}_5 \cdot 1$    $\bar{w}_3 = \bar{w}_5 \frac{\partial w_5}{\partial w_3} = \bar{w}_5 \cdot 1$

$w_4$    $w_3$

$\bar{w}_1^a = \bar{w}_4 \cos(w_1)$    $\bar{w}_2 = \bar{w}_3 \frac{\partial w_3}{\partial w_2} = \bar{w}_3 w_1$

$\bar{w}_1^b = \bar{w}_3 w_2$

$x_1$    $x_2$

$\bar{x}_1 = \bar{w}_1^a + \bar{w}_1^b = \cos(x_1) + x_2$    $\bar{x}_2 = \bar{w}_2 = x_1$

Backward propagation of derivative values

# Differentiable Programming in HEP

Mix **learnable ML modules** with domain-specific computations,
e.g. **physics code**, and use the full pipeline as a jointly optimizable entity



Image credit: L. Heinrich

# A Growing Body of Work in HEP

**Differentiable Matrix Elements**
[Heinrich, **MK**, 2203.00057]

**Differentiable Accelerator Simulation**
[Roussel, Edelen, 2211.09077]

**Differentiable LAr TPC simulation**
[Gasiorowski, et al., 2309.04639]

**Differentiable Parton Distribution Functions** [Ball, et al., 2109.02671]

**Differentiable Inference**
[Feickert, Heinrich, Stark, 2211.15838]

**Differentiable Analysis + Inference**
[Simpson, Heinrich, 2203.05570]

# Example: Jet Classification



**Jet** = Unordered set of particles

Each particles has a list of features:

**Particle** = {momentum, direction, position, … }



Image Credit: 1909.12285

# Differentiable Vertexing

Vertex ($\bigstar$)

$$v^* = \arg\min_v \chi^2(v, \alpha)$$

# Differentiable Vertexing

Vertex (★)

$$v^* = \arg\min_v \chi^2(v, \alpha)$$

Gradients from implicit differentiation

$$\mathcal{G} \equiv \left.\frac{\partial \chi^2(v, \alpha)}{\partial v}\right|_{v^*} = 0$$

$$\frac{\partial v^*}{\partial \alpha} = -\left(\frac{\partial \mathcal{G}}{\partial v}\right)^{-1}\left(\frac{\partial \mathcal{G}}{\partial \alpha}\right)$$

Smith, Ochoa, Inacio, Shoemaker, **MK**, 2310.12804

Smith, Ochoa, Inacio, Shoemaker, **MK**, 2310.12804

Smith, Ochoa, Inacio, Shoemaker, **MK**, 2310.12804

# What happens if there is stochasticity inside the model?

# Automatic Differentiation & ML

AD is great for differentiating deterministic functions like

$$y = f(x)$$

ML requires differentiating expectation values

$$\frac{\partial}{\partial \phi} \mathbb{E}_{p(x)}[f(x, \phi)] = \frac{\partial}{\partial \phi} \int f(x, \phi) p(x) dx$$

No param dependence in the distribution $p(\cdot)$

AD is great for differentiating deterministic functions like

$$y = f(x)$$

ML requires differentiating expectation values

$$\frac{\partial}{\partial \phi} \mathbb{E}_{p(x)}[f(x, \phi)] = \frac{\partial}{\partial \phi} \int f(x, \phi) p(x) dx$$

"Easy" Case

$$= \int \frac{\partial f(x, \phi)}{\partial \phi} p(x) dx = \mathbb{E}_{p(x)} \left[ \frac{\partial f(x, \phi)}{\partial \phi} \right]$$

$z \sim p_\phi(z|x)$

param dependence in $q(\cdot)$

$$L = \mathbb{E}_{q_\phi(z|x)}[\log p(x|z)] - \mathbb{E}_{q_\phi(z|x)}\left[\log \frac{q_\phi(z|x)}{p(z)}\right]$$

Kingma, Welling, 1312.6114
Rezende, Mohamed, Wierstra, 1401.4082

$$L = \mathbb{E}_{p(\epsilon)}[\log p(x|z(\epsilon,\phi))] - \mathbb{E}_{p(\epsilon)}\left[\log \frac{q_\phi(z(\epsilon,\phi)|x)}{p(z(\epsilon,\phi))}\right]$$

Kingma, Welling, 1312.6114
Rezende, Mohamed, Wierstra, 1401.4082

Separate parameters from stochasticity

$$x \sim p_\theta(x) \quad \rightarrow \quad \text{rewrite } x = g(\epsilon, \theta) \text{ with } \epsilon \sim p(\epsilon)$$

Example:

$$x \sim \mathcal{N}(\mu, \sigma) \quad \rightarrow \quad x = g(\epsilon, \mu, \sigma) = \epsilon * \sigma + \mu \quad \text{with } \epsilon \sim \mathcal{N}(0,1)$$

$$\frac{d}{d\theta} \mathbb{E}_{p_\theta(x)}[f(x)] = \frac{d}{d\theta} \mathbb{E}_{p(\epsilon)}\left[f\big(g(\epsilon, \theta)\big)\right] = \mathbb{E}_{p(\epsilon)}\left[\frac{df}{dg}\frac{dg}{d\theta}\right]$$

# Is that all we need?

# A Problem: Discrete Random Variables & Choices

Discrete random variables and discrete choices are all over HEP

Branching / Showering Processes

Clustering Algorithms

```
def f(x):
    theta = (sin(2.*x))**2
    b = bernoulli(theta)
    g = x*x
    return g+b
```

Bernoulli parameter $\theta$ depends on $x$

$$f(x) = x^2 + b \qquad b \sim Bern(\theta = \sin^2(2x))$$

# Programs with Discrete Randomness

```
def f(x):
    theta = (sin(2.*x))**2
    b = bernoulli(theta)
    g = x*x
    return g+b
```

$$f(x) = x^2 + b \qquad b \sim Bern(\theta = \sin^2(2x))$$

$$\mathbb{E}_b[f(x)] = x^2 + \sin^2 2x$$

$$\nabla_x \mathbb{E}_b[f(x)] = 2x + 4\sin(2x)\cos(2x)$$



Even if a program contains discrete randomness,
expected value can be smooth and have a well-defined derivative

# Programs with Discrete Randomness

```python
def f(x):
    theta = (sin(2.*x))**2
    b = bernoulli(theta)
    g = x*x
    return g+b
```

$$f(x) = x^2 + b \qquad b \sim Bern(\theta = \sin^2(2x))$$



$$\mathbb{E}_b[f(x)] = x^2 + \sin^2 2x$$

$$\nabla_x \mathbb{E}_b[f(x)] = 2x + 4\sin(2x)\cos(2x)$$

AD Gradient:  $grad(f_i(x)) \rightarrow 2x_i$

Standard AD tools don't know how to handle discrete randomness that depends on the parameter of differentiation → **We need another approach**

# Derivatives for Discrete Randomness

**Do Some Work,
Get Better Derivatives**

**Approximate
Derivatives**

**Numerical
Derivatives**

**Don't Use
Derivatives**

| Score Functions | Stochastic AD | Smoothing / Relaxations | Surrogates | Finite Differences | Gradient-Free Methods |
|---|---|---|---|---|---|

*Related work*:
Smooth perturbation analysis

*Example*:
Differentiable ranking & sorting

*HEP Example*:
Surrogates for SHiP magnet design

$$\frac{f(x + \epsilon) - f(x)}{\epsilon}$$

Bayesian Opt.,
Genetic Algs, …



2002.08871



**MK**, et al. 2002.04632



BO for
EIC Design

Figure credit: C. Fanelli

$$\nabla_\theta \mathbb{E}_{p_\theta(x)}[f(x)] = \int \nabla_\theta p_\theta(x) f(x) dx$$

$$\nabla_\theta \mathbb{E}_{p_\theta(x)}[f(x)] = \int p_\theta(x) f(x) \nabla_\theta \log p_\theta(x) \, dx$$

because:

$$\nabla_\theta \log p_\theta(x) = \frac{1}{p_\theta(x)} \nabla_\theta p_\theta(x)$$

$$\nabla_\theta \mathbb{E}_{p_\theta(x)}[f(x)] = \mathbb{E}_{p_\theta(x)}[f(x)\nabla_\theta \log p_\theta(x)]$$

Gradient estimator used in Reinforcement Learning

Works with discrete $x$ and even non-differentiable $f(\cdot)$

Requires tracking probabilities $\log p_\theta(x)$ throughout program



AlphaStar Vinyals et al. 2019

$$\frac{d}{d\theta} \mathbb{E}_{p_\theta}[f(x)] = \mathbb{E}_{p_\theta}[\delta + \beta(y - x)]$$

Standard AD          Weight          Alternative value of rv

Recently, *Arya et al.* extended fwd-mode AD to discrete-stochastic environments

Importantly, this includes a *composition rule* for how to combine weights $\beta$ step-by-step along the computation chain

## Automatic Differentiation of Programs with Discrete Randomness

**Gaurav Arya**
Massachusetts Institute of Technology, USA
aryag@mit.edu

**Moritz Schauer**
Chalmers University of Technology, Sweden
University of Gothenburg, Sweden
smoritz@chalmers.se

Can use the inversion method to reparameterize discrete random variables

$$x \sim Bernoulli(p) \qquad \rightarrow$$

$$\omega \sim Uniform[0,1]$$

$$x = \begin{cases} 1 & if\ \omega > 1 - p \\ 0 & Otherwise. \end{cases}$$

Can use the inversion method to reparameterize discrete random variables

$$x \sim Bernoulli(p) \qquad \rightarrow$$

$$\omega \sim Uniform[0,1]$$

$$x = \begin{cases} 1 & if \ \omega > 1 - p \\ 0 & Otherwise. \end{cases}$$

$$\mathbb{E}[b] = \int 1_{[\omega > 1-p]} p(\omega) d\omega = \int_{1-p}^{1} d\omega$$

# Intuition

Can use the inversion method to reparameterize discrete random variables

$$x \sim Bernoulli(p) \qquad \rightarrow$$

$$\omega \sim Uniform[0,1]$$

$$x = \begin{cases} 1 & if \ \omega > 1 - p \\ 0 & Otherwise. \end{cases}$$

$$\mathbb{E}[b] = \int 1_{[\omega > 1-p]} p(\omega) d\omega = \int_{1-p}^{1} d\omega$$

Standard AD on Monte Carlo expectation of this program would still be wrong

$$grad_p \left( \frac{1}{N} \sum_i [1 \ if \ (\omega_i > 1 - p) \ else \ 0] \right) = 0$$

Can use the inversion method to reparameterize discrete random variables

$$x \sim Bernoulli(p) \quad \rightarrow$$

$$\omega \sim Uniform[0,1]$$

$$x = \begin{cases} 1 & if \ \omega > 1 - p \\ 0 & Otherwise. \end{cases}$$

$$\mathbb{E}[b] = \int 1_{[\omega > 1-p]} p(\omega) d\omega = \int_{1-p}^{1} d\omega$$

Param. dependence in integration bounds

Correct derivative must account for boundary dependence → *Leibniz Rule*

# Intuition

Can use the inversion method to reparameterize discrete random variables

$$x \sim Bernoulli(p) \qquad \rightarrow$$

$$\omega \sim Uniform[0,1]$$

$$x = \begin{cases} 1 & if\ \omega > 1 - p \\ 0 & Otherwise. \end{cases}$$



Output

1

$X(p+\varepsilon)$    $X(p)$

$Uniform[0,1]$

$0 \quad 1-p-\varepsilon \quad 1-p \quad\quad 1$

Arya et al., [2210.08572](2210.08572)

Can use the inversion method to reparameterize discrete random variables

$$x \sim Bernoulli(p) \qquad \rightarrow$$

$$\omega \sim Uniform[0,1]$$

$$x = \begin{cases} 1 & if \ \omega > 1 - p \\ 0 & Otherwise. \end{cases}$$



Output

$X(p + \epsilon) - X(p)$

$1$

$Uniform[0,1]$

$0 \quad 1 - p - \varepsilon \quad 1 - p \quad 1$

$$\frac{d}{d\theta} \mathbb{E}_{p_\theta}[f(x)] = \mathbb{E}_{p_\theta}[\delta + \beta(y - x)]$$

The weight $\beta$ accounts for the derivative of the probability of a jump in program

Equivalently, the weight accounts for the boundary derivative

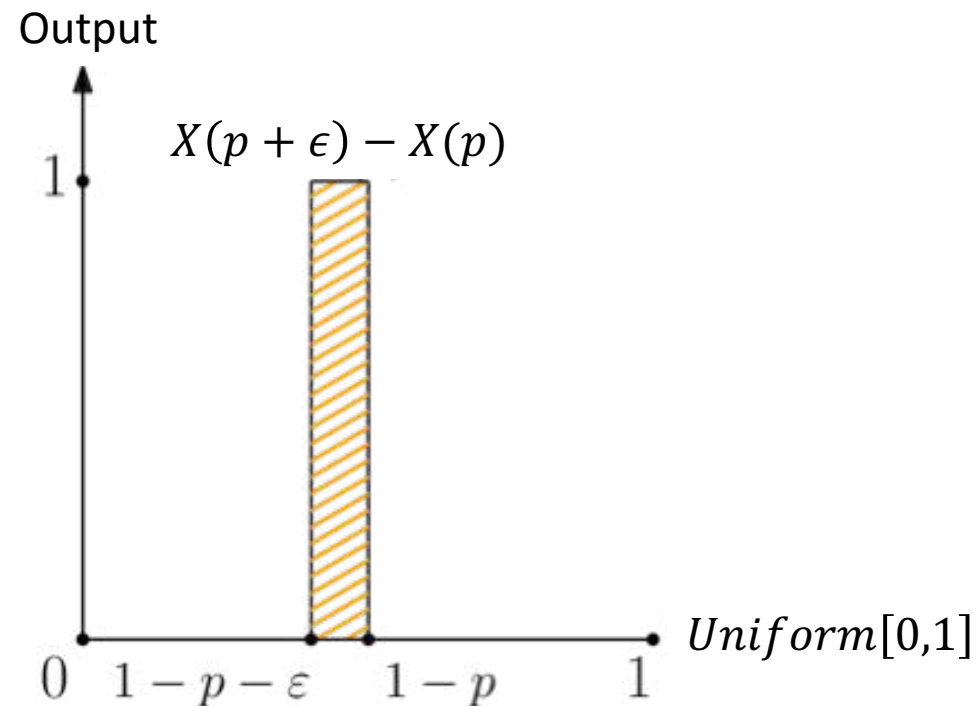In many cases: $\beta = \frac{\partial_\theta CDF_\theta(X(\theta))}{PDF_\theta(X(\theta))}$



Output

$X(p + \epsilon) - X(p)$

$Uniform[0,1]$

$0 \quad 1 - p - \varepsilon \quad 1 - p \quad 1$

Arya et al., 2210.08572

# Automatic Differentiation

Primal

Pruning follows one alternative path

$y \sim Y|X(p) = x$
$x = X(p)(\omega_1)$

Smoothing flattens onto primal

State space

Steps

Uncoupled path

$X(p+\varepsilon)(\omega_2)$

Correlated paths → low variance

$\mathcal{O}(1)$ unbiased forward mode AD

Arya et al., 2210.08572

# Are these methods useful?

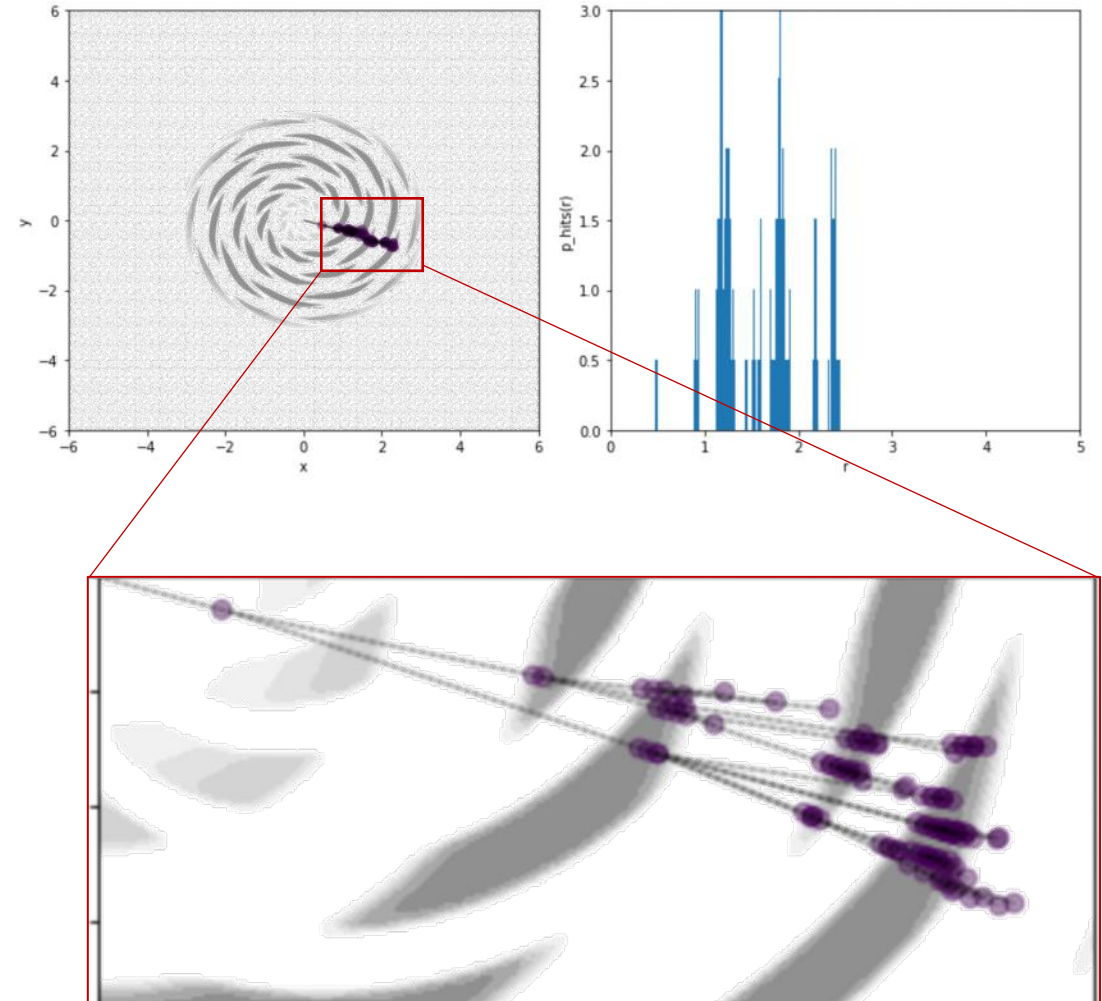# Toy Shower

*Simplified particle shower:*

Including Energy loss and splitting

*Design parameter:*

Radial distance of material
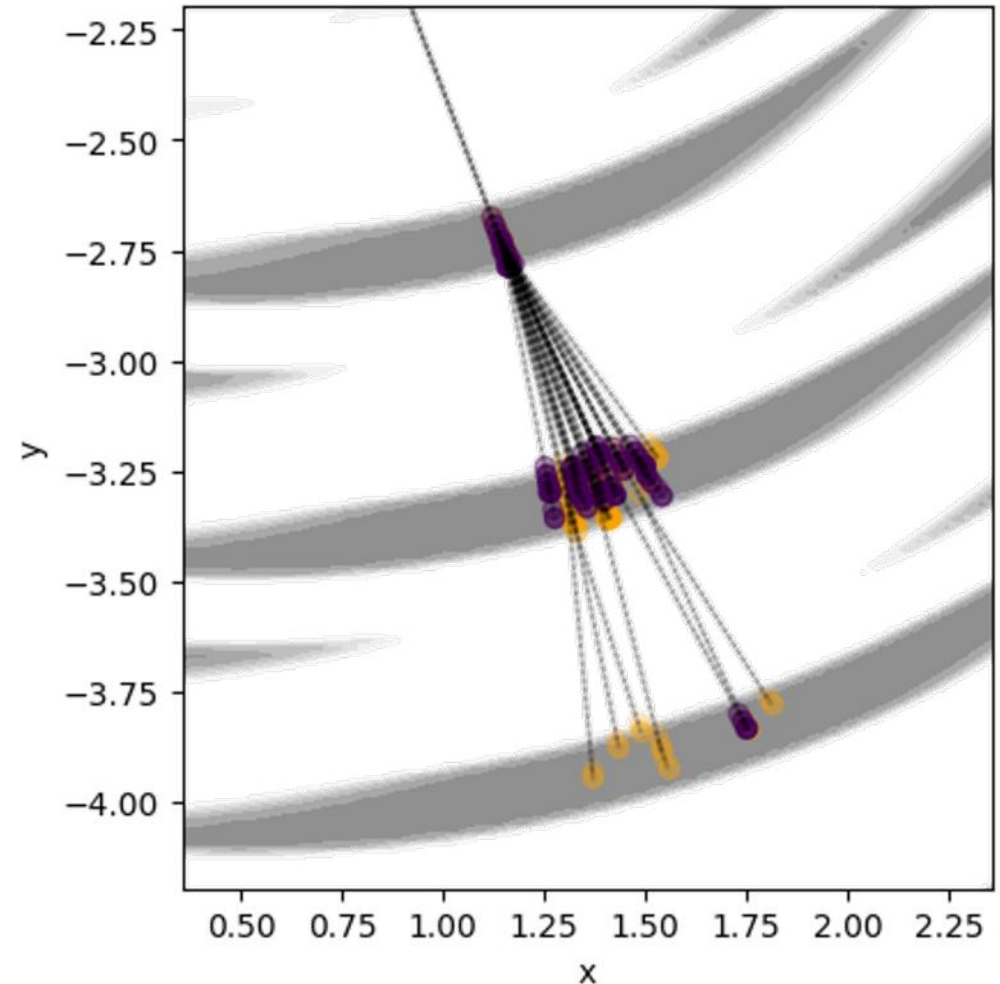
*Design goal:*

Specify average shower depth

# Toy Shower

*Simplified particle shower:*
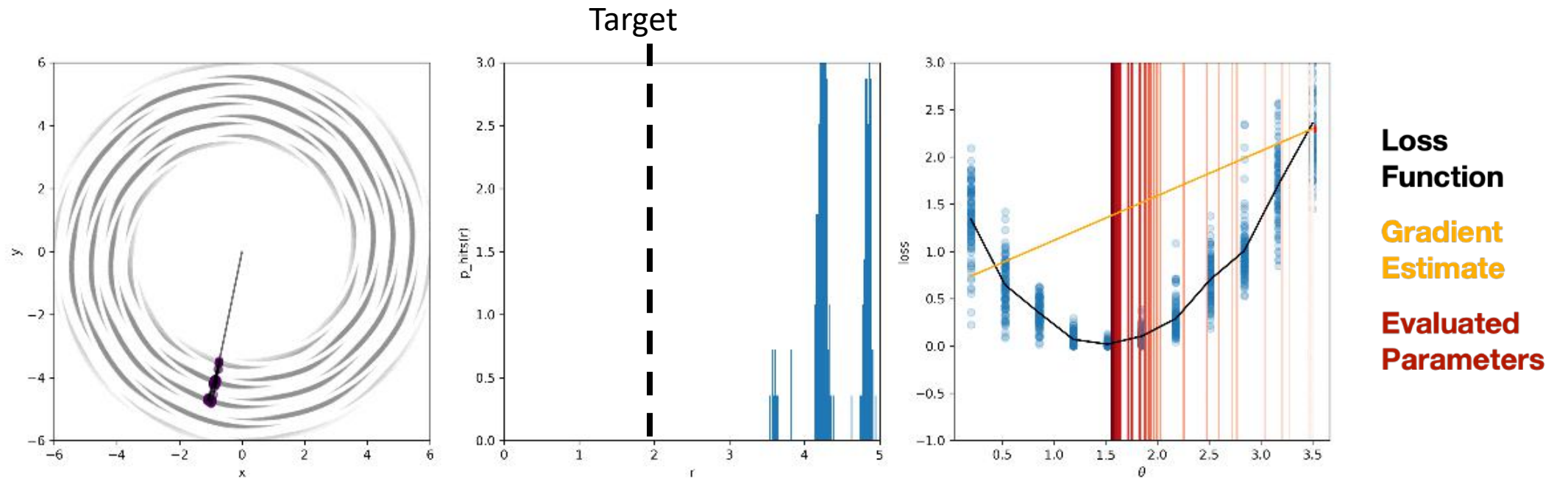
Including Energy loss and splitting

*Design parameter:*

Radial distance of material

*Design goal:*

Specify average shower depth

Dedicated implementation of
Stochastic AD
→ Can generate "alternative showers"

# Example Optimization



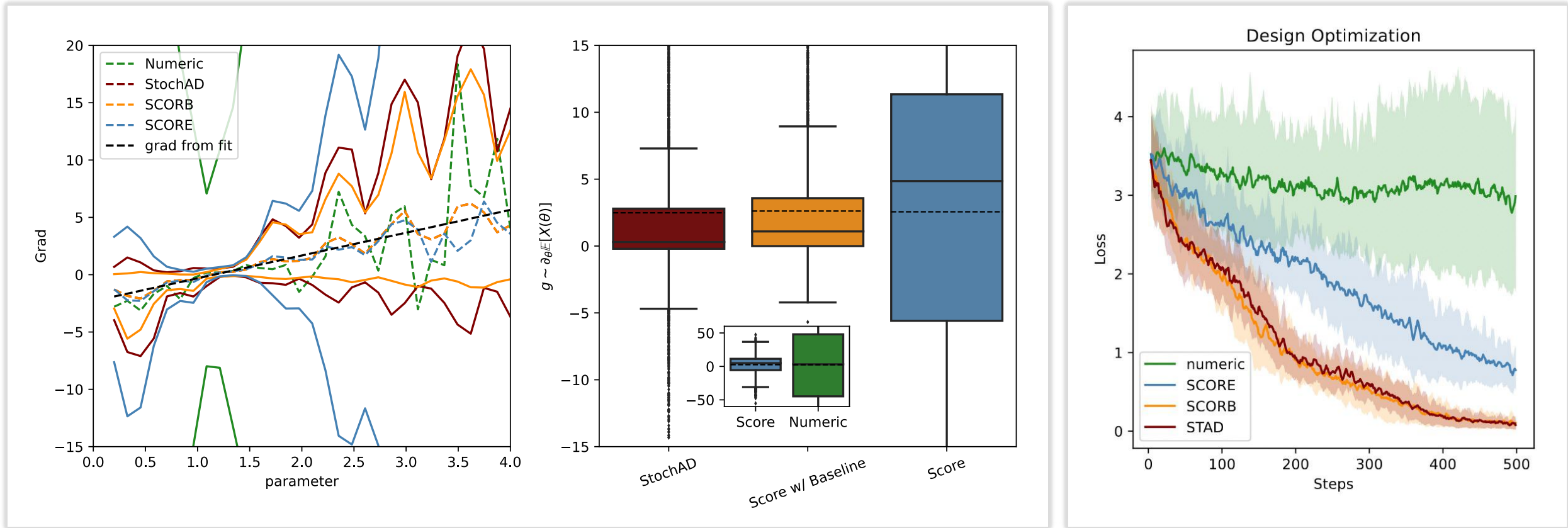Target

Loss Function

Gradient Estimate

Evaluated Parameters

Gradients are noisy but in right direction (on average) → optimization works!

Both score function and Stochastic AD have reasonable variance of gradients



**MK**, Heinrich, 2308.16680

# Summary

LHC and future colliders present unique opportunities
Need to make the most of it!



Differentiable programming provides powerful method to
Optimize our data analysis and simulation pipelines
Embed physics knowledge in our ML tools



Not everything is easily differentiable,
lots of challenges along the way...
especially for programs with discrete randomness