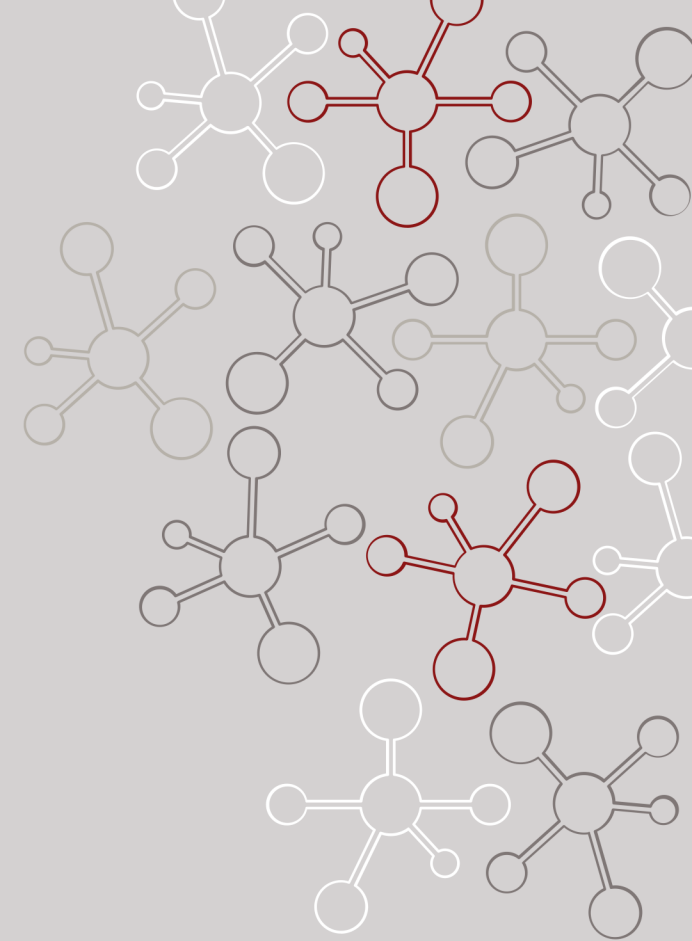


# Estimating Gradients of Programs with Discrete Randomness:

With Applications in Detector Design Optimization

Michael Kagan, SLAC  
Lukas Heinrich, TUM

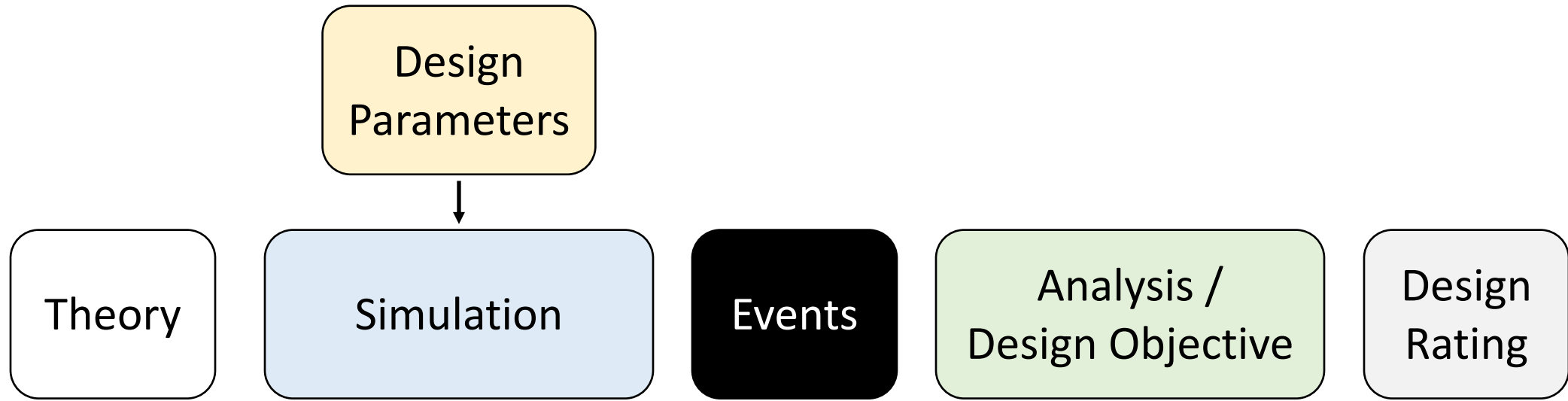
3<sup>rd</sup> MODE Workshop  
July 26, 2023



# Detector Design Optimization

2

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



# Detector Design Optimization

3

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


Design Objective

Realizations of  
measurements:  
E.g. Simulations

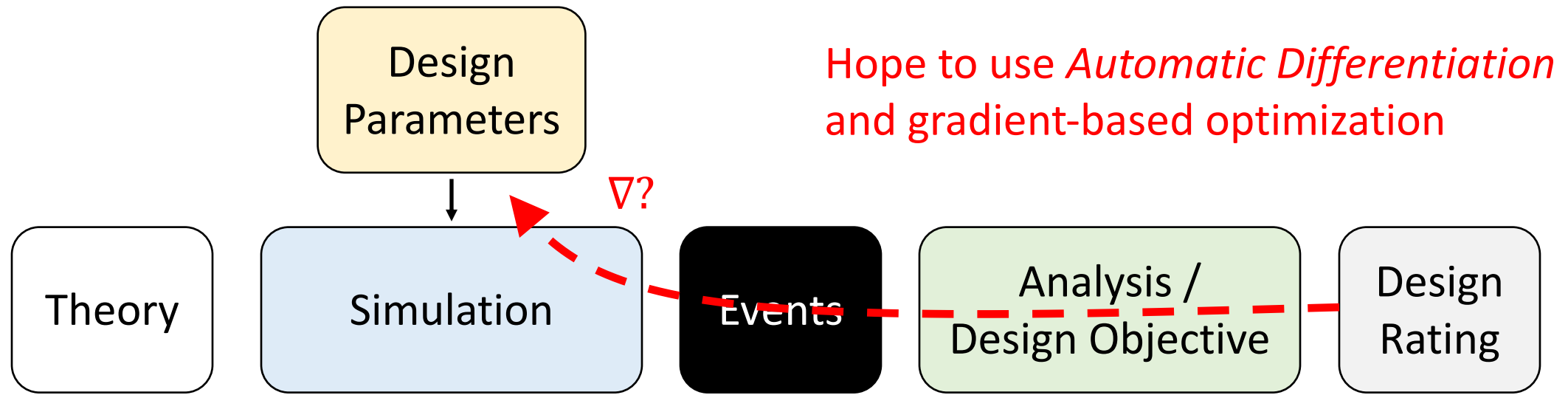
Design Parameters

Probability to see a  
measurement:  
E.g. Interaction and  
detection probability

# Detector Design Optimization

4

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

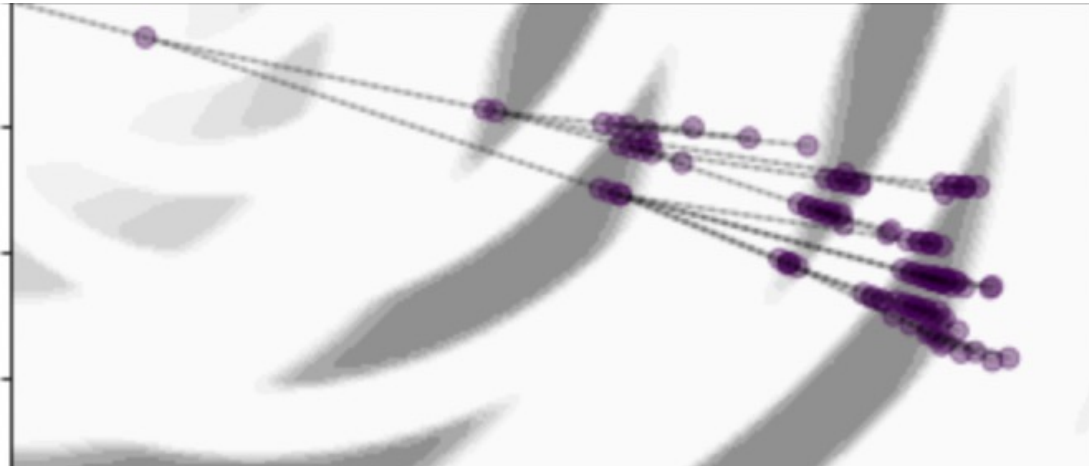


# A Problem: Discrete Random Variables & Choices

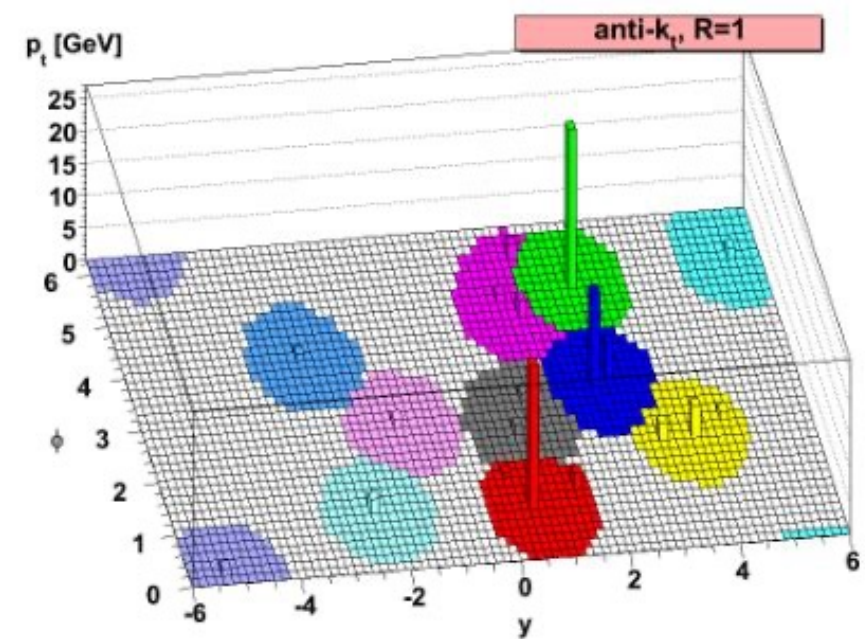
5

Discrete random variables and discrete choices are all over HEP

Branching / Showering Processes



Clustering Algorithms



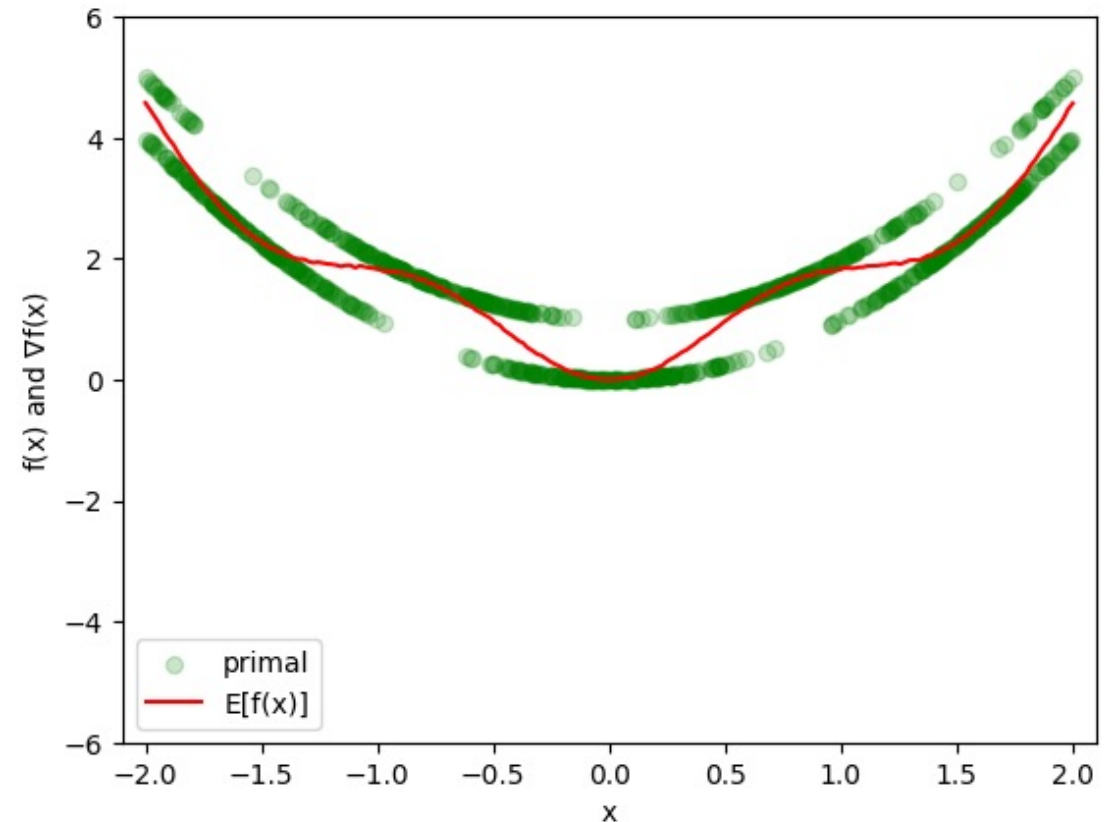
# Programs with Discrete Randomness

6

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

Bernoulli parameter  $\theta$   
depends on  $x$

$$f(x) = x^2 + b \quad b \sim \text{Bernoulli}(\theta = \sin^2(2x))$$



# Gradients of Expected Values

7

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

Simple example:

```
def f(theta):  
    b = bernoulli(theta)  
    return b
```

$$\mathbb{E}[f(\theta)] = \mathbb{E}_{b \sim \text{Bern}(\theta)}[b] = \theta$$

$$\nabla_{\theta} \mathbb{E}[f(\theta)] = 1$$

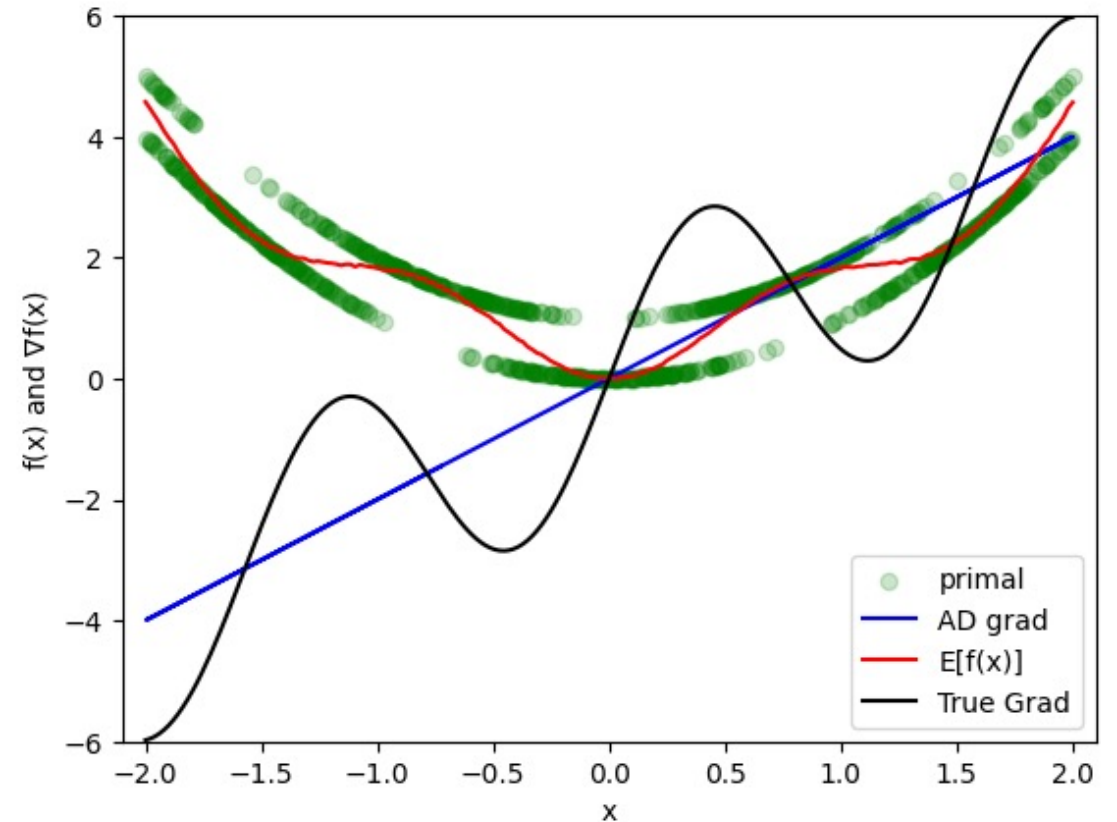
# AD in Programs with Discrete Randomness

8

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

AD Gradient:  $\text{grad}(f_i(x))$

True gradient:  $\nabla_x E[f(x)]$





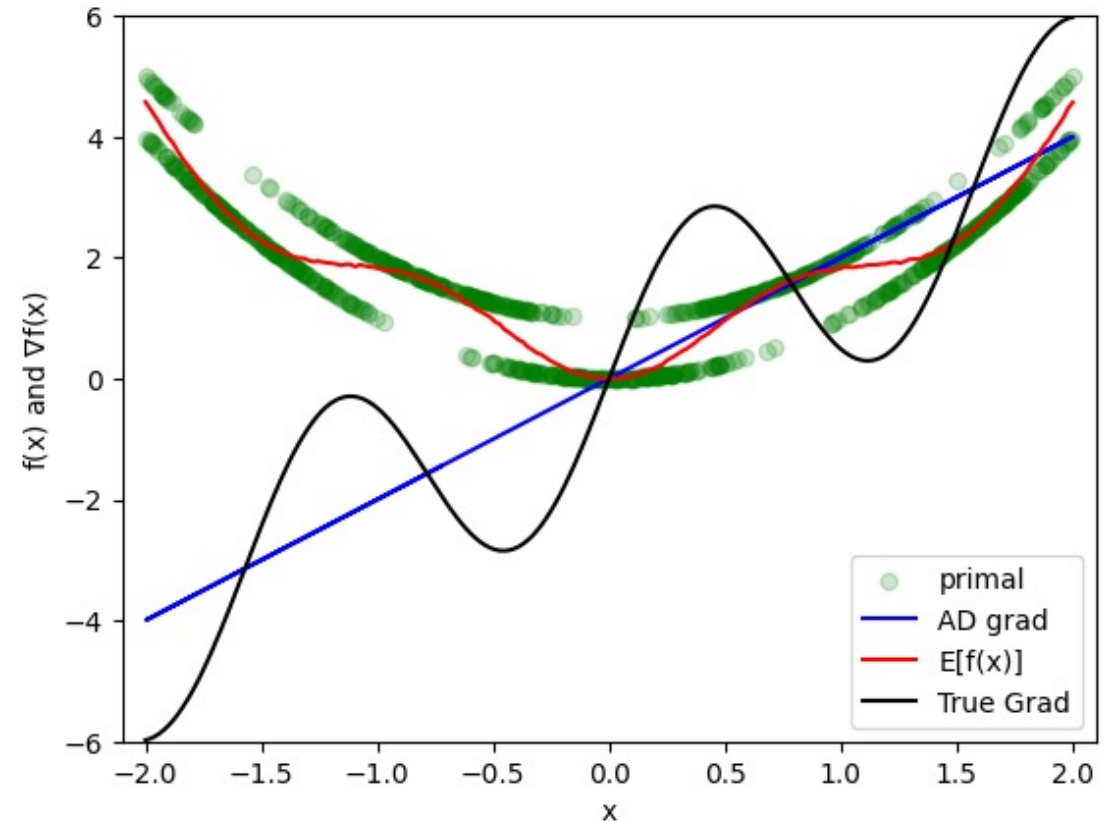
# AD in Programs with Discrete Randomness

9

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

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

10

Do Some Work,  
Get Better Derivatives

Score  
Functions

Stochastic AD

Approximate  
Derivatives

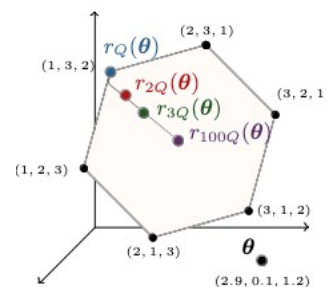
Smoothing /  
Relaxations

Surrogates

Don't Use  
Derivatives

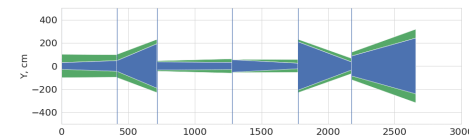
Gradient-Free  
Methods

*Example:*  
Differentiable  
ranking and sorting



[2002.08871](#)

*HEP Example:*  
Surrogates for  
SHiP magnet design



MK, et al. [2002.04632](#)

Bayesian Optimization,  
Genetic Algorithms, ...

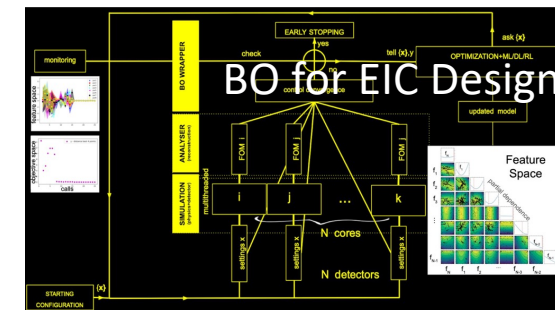


Figure credit: C. Fanelli

# Score Functions / REINFORCE

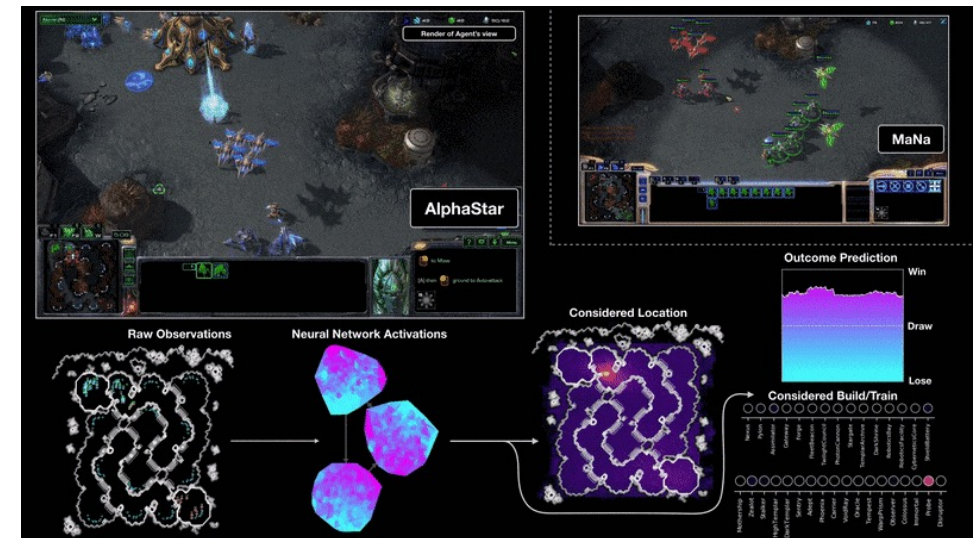
11

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

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

[2210.08572](#)

---

## 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

## Quick Aside: Reparameterization Trick

13

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

Common method for continuous rv's is the *reparameterization trick*

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

Separate parameters from stochasticity

Example:

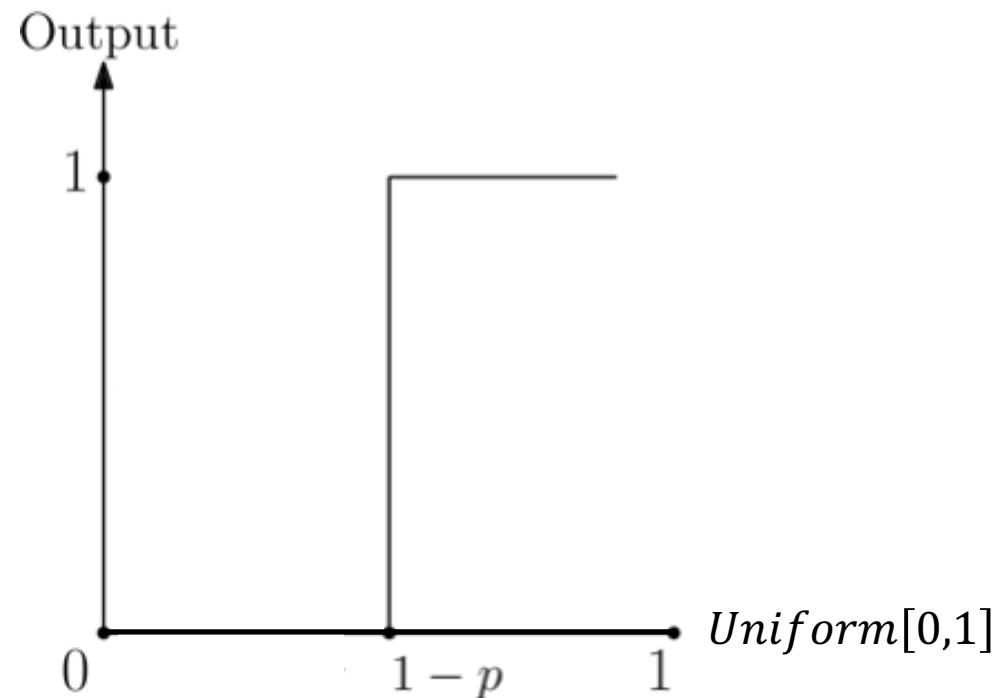
$$x \sim \mathcal{N}(\mu, \sigma) \rightarrow x = \epsilon * \sigma + \mu \text{ with } \epsilon \sim \mathcal{N}(0,1)$$

# Intuition

14

Can use the inversion method to reparameterize discrete random variables

$$x \sim \text{Bernoulli}(p) \quad \rightarrow \quad \begin{aligned} \omega &\sim \text{Uniform}[0,1] \\ x &= \begin{cases} 1 & \text{if } \omega > 1 - p \\ 0 & \text{Otherwise.} \end{cases} \end{aligned}$$



Can use the inversion method to reparameterize discrete random variables

$$x \sim \text{Bernoulli}(p) \quad \rightarrow \quad \begin{aligned} \omega &\sim \text{Uniform}[0,1] \\ x &= \begin{cases} 1 & \text{if } \omega > 1 - p \\ 0 & \text{Otherwise.} \end{cases} \end{aligned}$$

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

Can use the inversion method to reparameterize discrete random variables

$$x \sim \text{Bernoulli}(p) \quad \rightarrow \quad \begin{aligned} \omega &\sim \text{Uniform}[0,1] \\ x &= \begin{cases} 1 & \text{if } \omega > 1 - p \\ 0 & \text{Otherwise.} \end{cases} \end{aligned}$$

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

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

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



# Intuition

17

Can use the inversion method to reparameterize discrete random variables

$$x \sim \text{Bernoulli}(p) \quad \rightarrow \quad \begin{aligned} \omega &\sim \text{Uniform}[0,1] \\ x &= \begin{cases} 1 & \text{if } \omega > 1 - p \\ 0 & \text{Otherwise.} \end{cases} \end{aligned}$$

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

Param. dependence  
in integration bounds

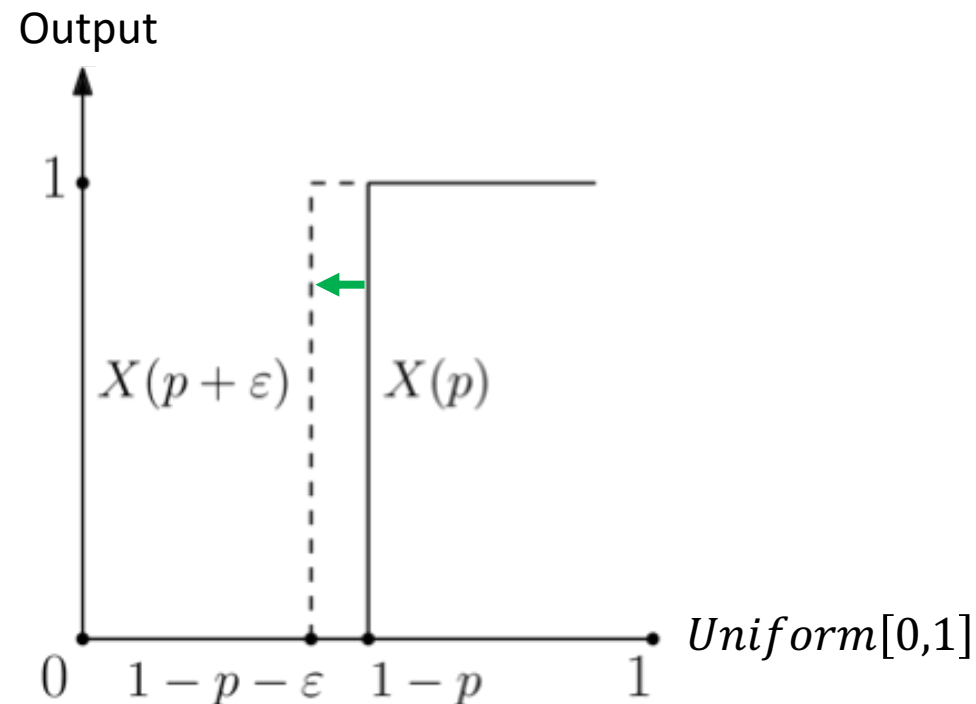
Correct derivative must account for boundary dependence  $\rightarrow$  *Leibniz Rule*

# Intuition

18

Can use the inversion method to reparameterize discrete random variables

$$x \sim \text{Bernoulli}(p) \quad \rightarrow \quad \begin{aligned} &\omega \sim \text{Uniform}[0,1] \\ &x = \begin{cases} 1 & \text{if } \omega > 1 - p \\ 0 & \text{Otherwise.} \end{cases} \end{aligned}$$

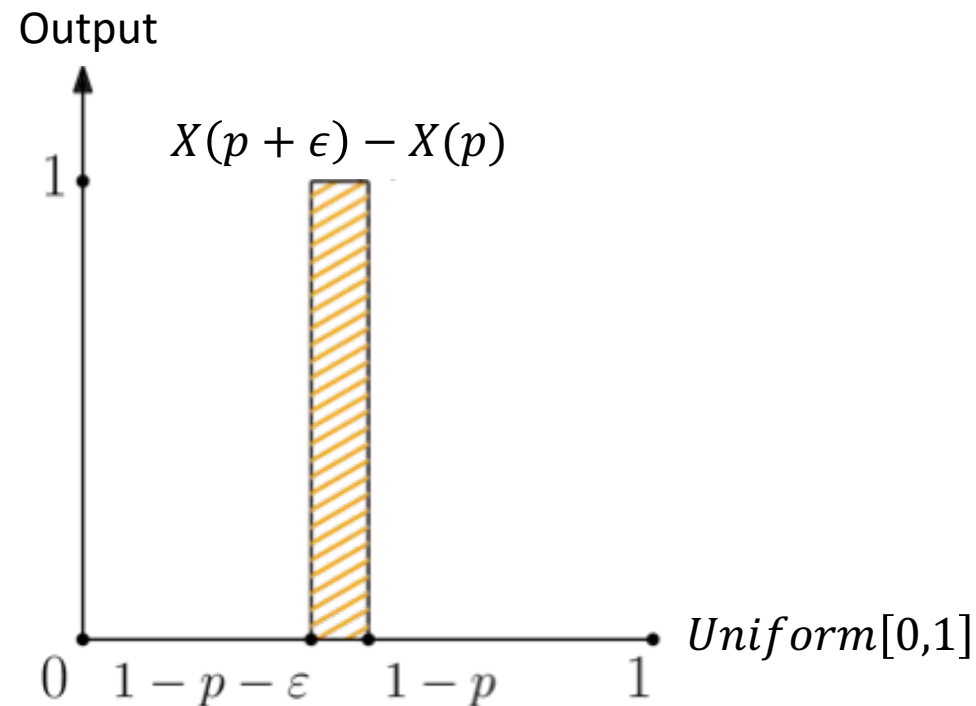


# Intuition

19

Can use the inversion method to reparameterize discrete random variables

$$x \sim \text{Bernoulli}(p) \quad \rightarrow \quad \begin{aligned} \omega &\sim \text{Uniform}[0,1] \\ x &= \begin{cases} 1 & \text{if } \omega > 1 - p \\ 0 & \text{Otherwise.} \end{cases} \end{aligned}$$

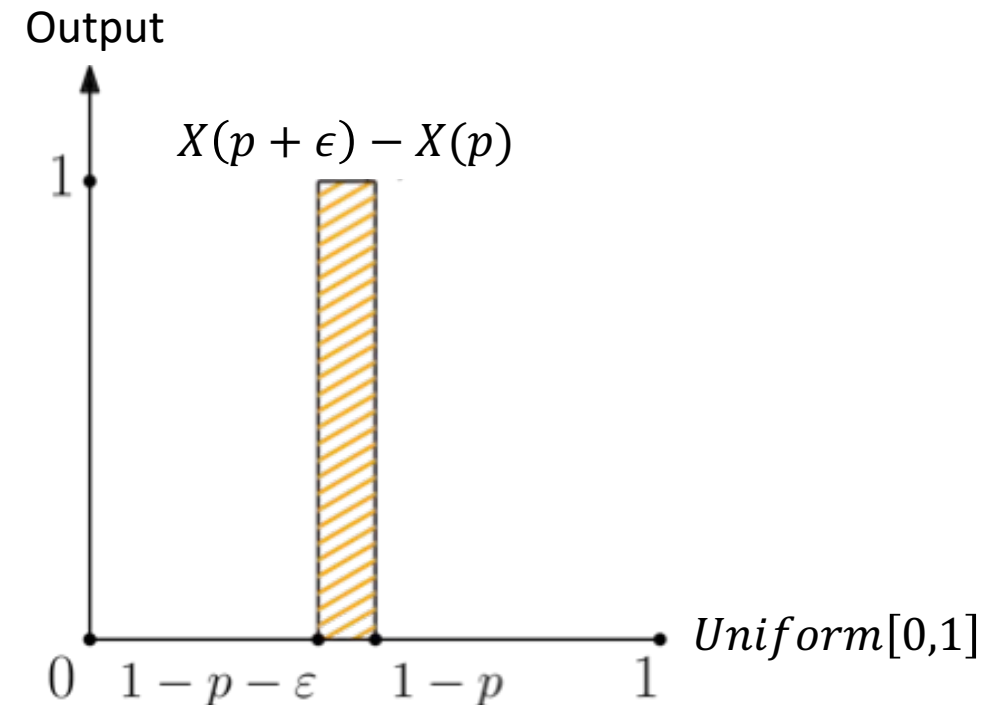


$$\frac{d}{d\theta} \mathbb{E}_{p_\theta}[f(x)] = \mathbb{E}_{p_\theta}[\delta + \beta(y - x)]$$

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

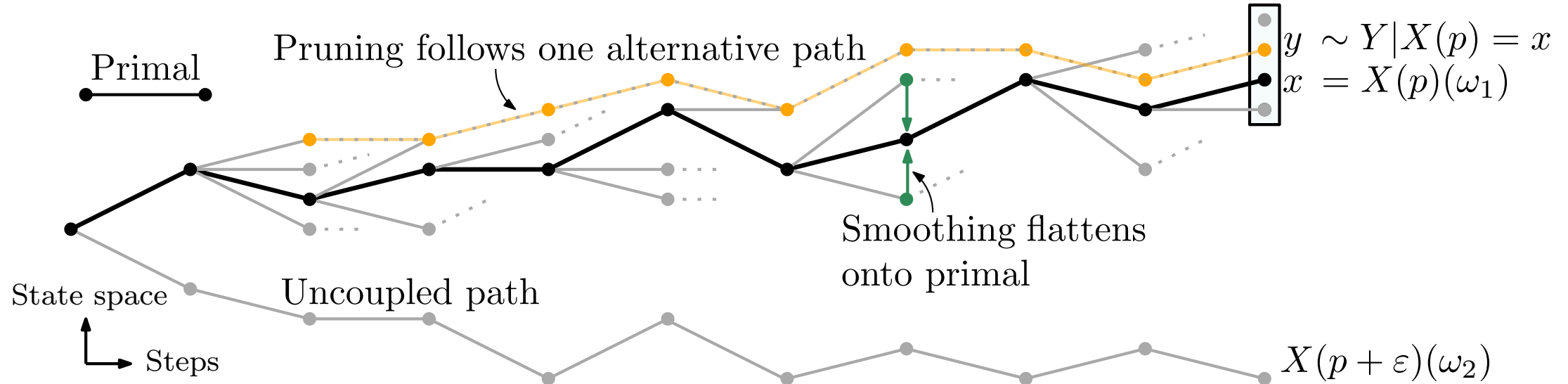
Equivalently, the weight accounts for the boundary derivative

In many cases:  $\beta = \frac{\partial_\theta \text{CDF}_\theta(X(\theta))}{\text{PDF}_\theta(X(\theta))}$



# Automatic Differentiation

21



Correlated paths  $\rightarrow$  low variance

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

Are these methods useful?

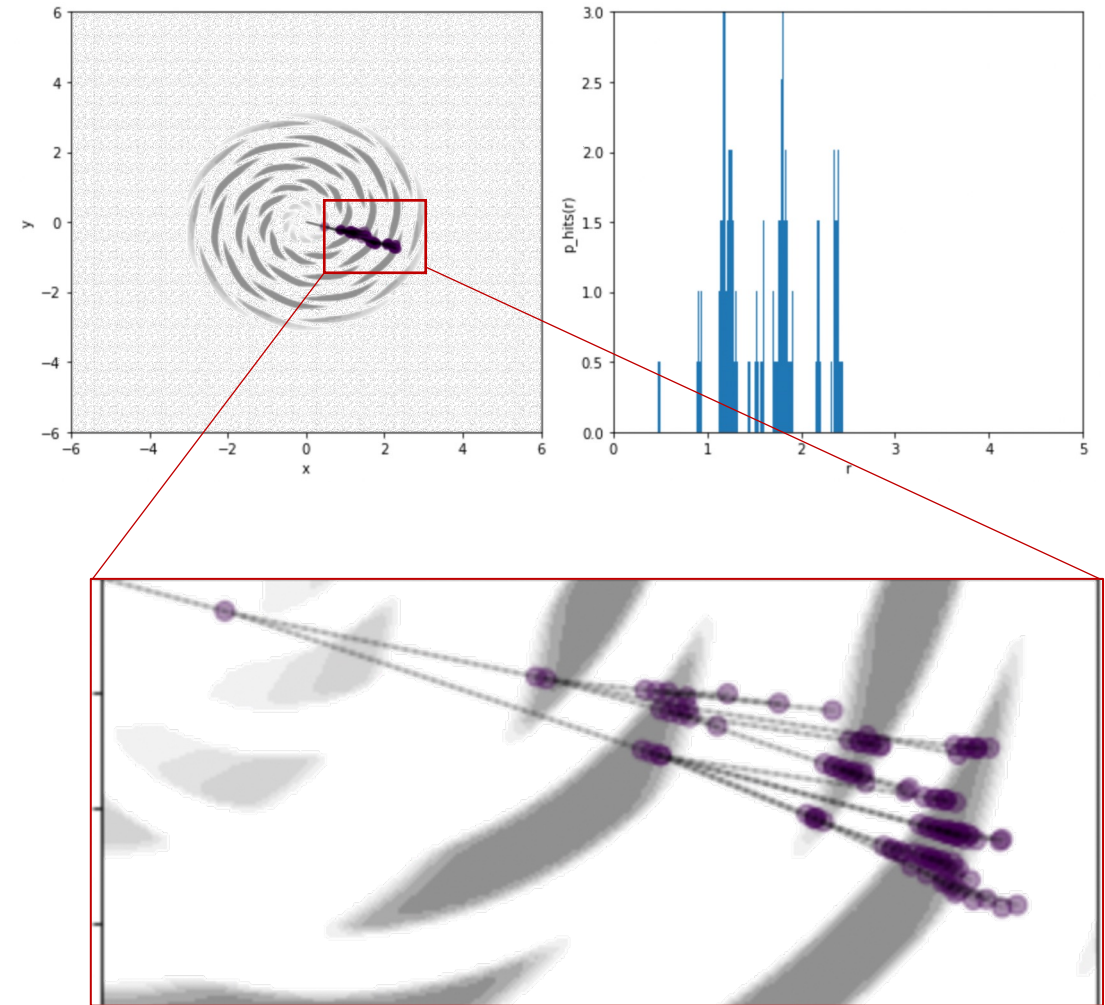
# Toy Shower

23

*Simplified particle shower:*  
Including Energy loss and splitting

*Design parameter:*  
Radial distance of material

*Design goal:*  
Specify average shower depth



# Toy Shower

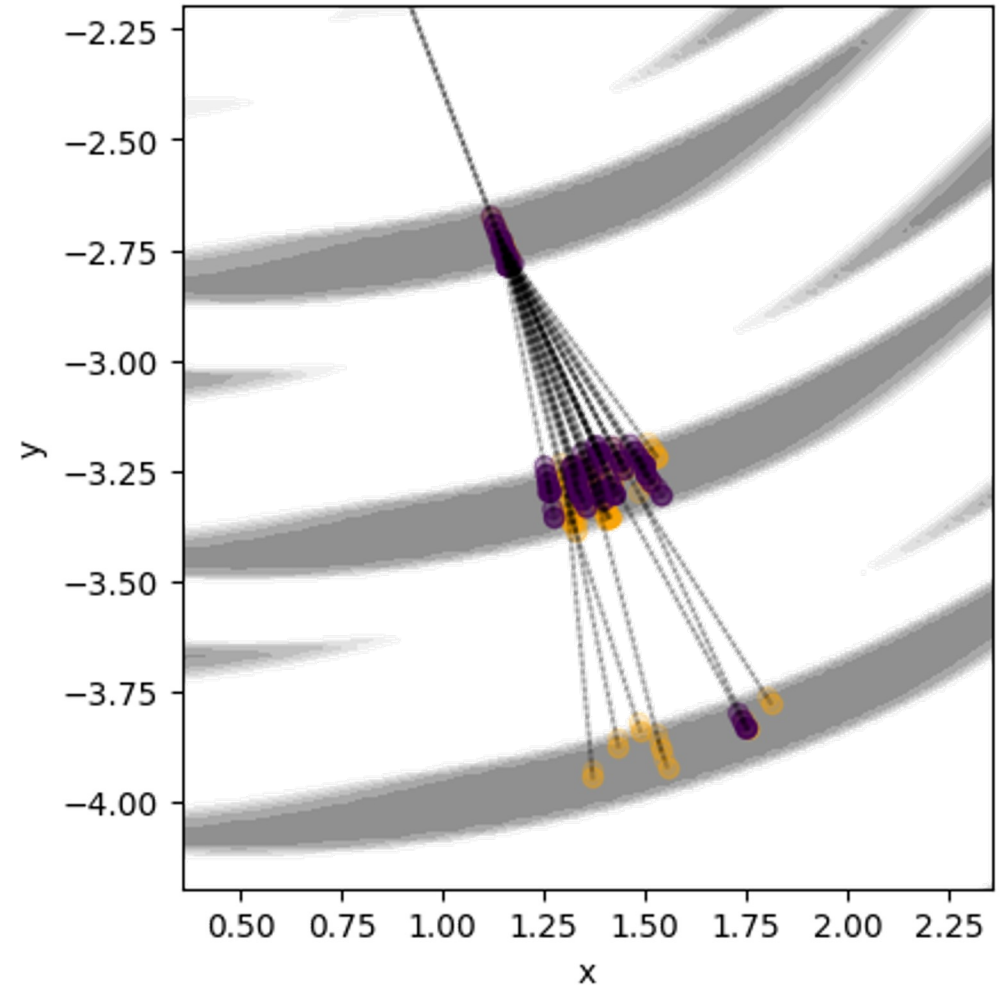
24

*Simplified particle shower:*  
Including Energy loss and splitting

*Design parameter:*  
Radial distance of material

*Design goal:*  
Specify average shower depth

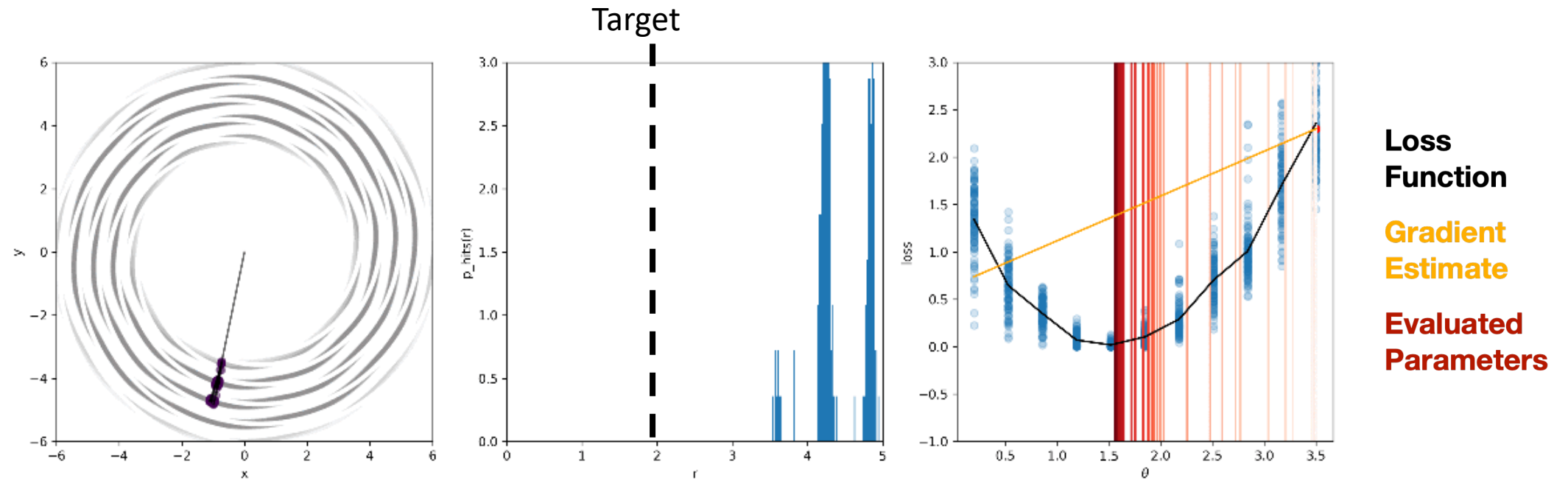
Dedicated implementation of  
Stochastic AD  
→ Can generate “alternative showers”





# Example Optimization

25



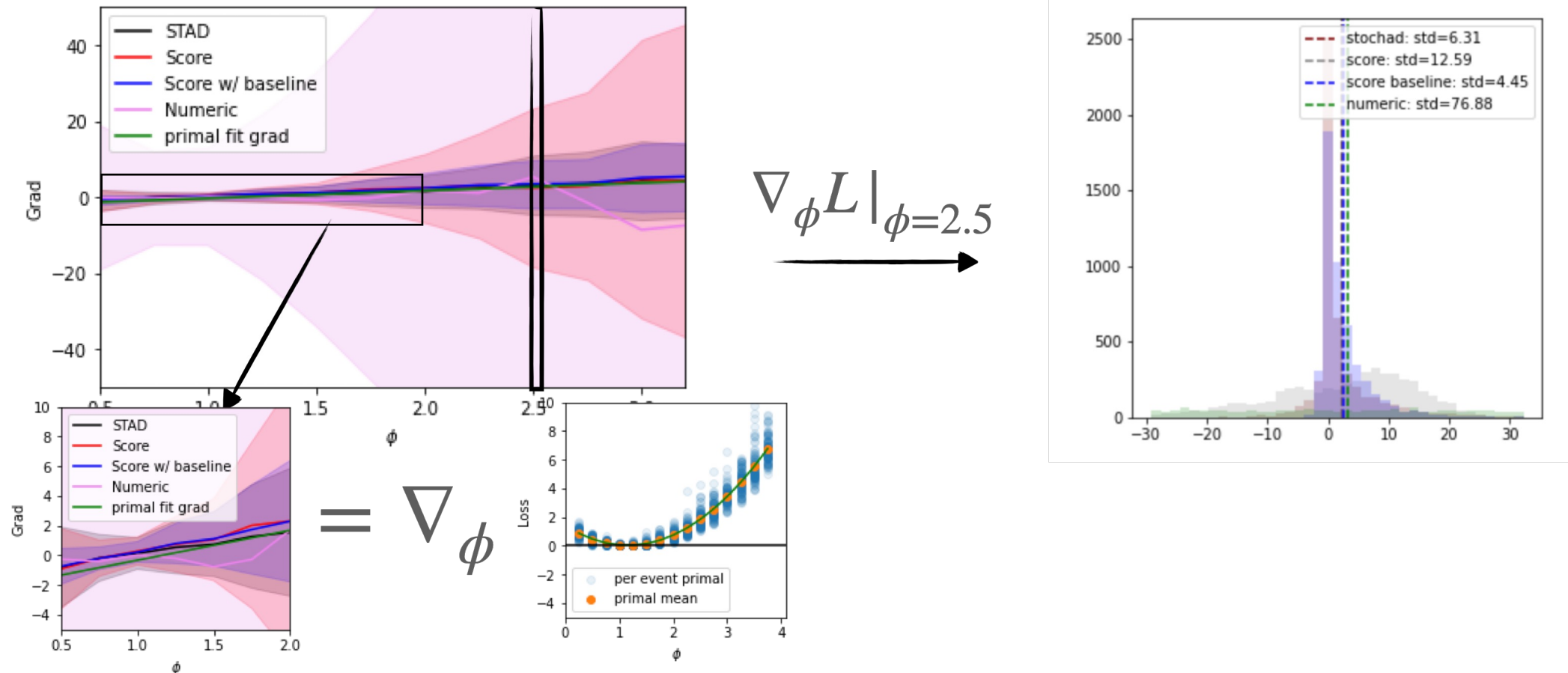
Gradients are noisy but in right direction (on average)  $\rightarrow$  optimization works!

# Comparisons

26

Both score function and Stochastic AD have reasonable variance gradients

Much to explore on how to *couple* primal & alternative programs to lower variance



# Summary

27

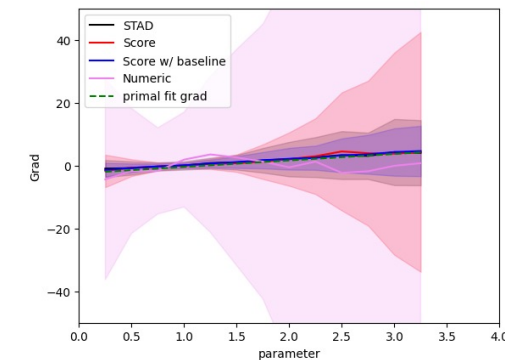
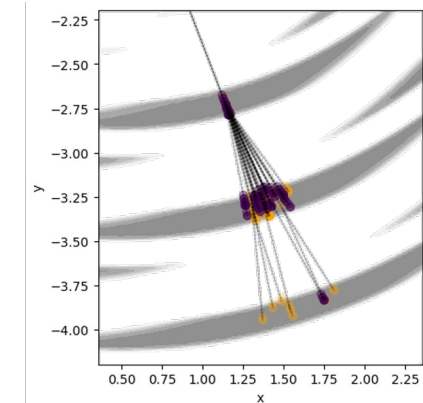
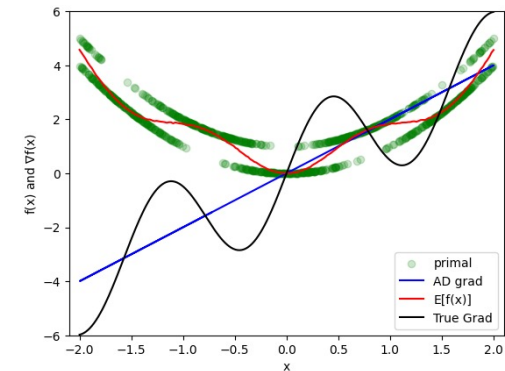
Programs with discrete stochasticity are all over HEP

Discrete stochasticity is a problem for differentiation, but expected value of these programs (which is what we want anyway) may be differentiable

Variety of tools can potentially handle this, like score function or recent exciting work on Stochastic AD.

In a toy shower, we could successfully differentiate the program and perform optimization...

- Proof of principle... we may be able to scale this up!
- Still much work on reducing gradient variance



# Backup

# Differentiable Programming

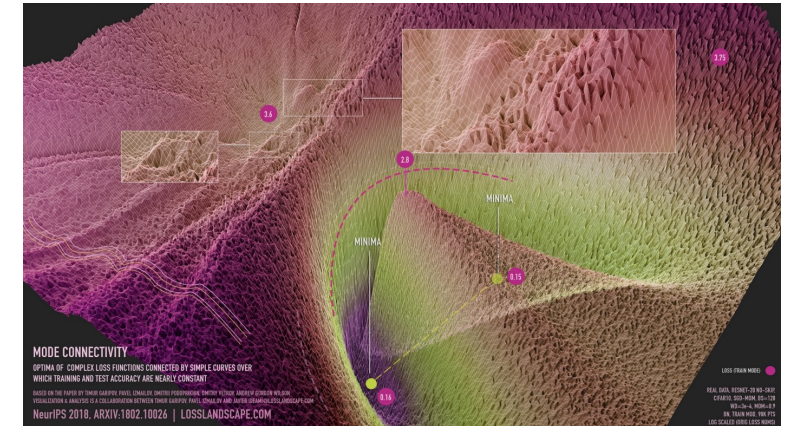
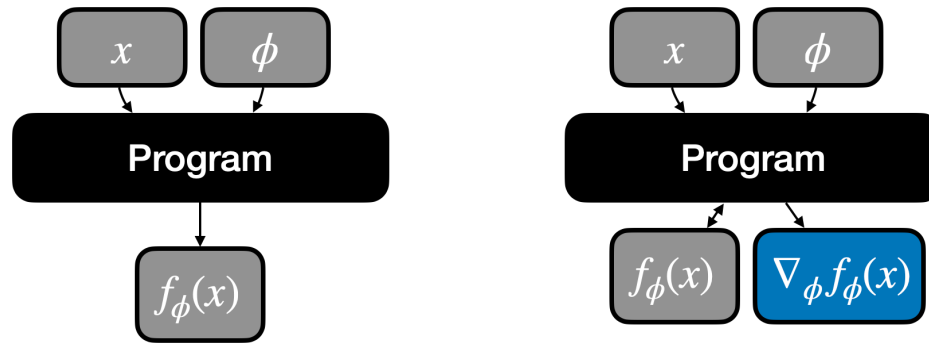
29

Automatic differentiation super-charges code:

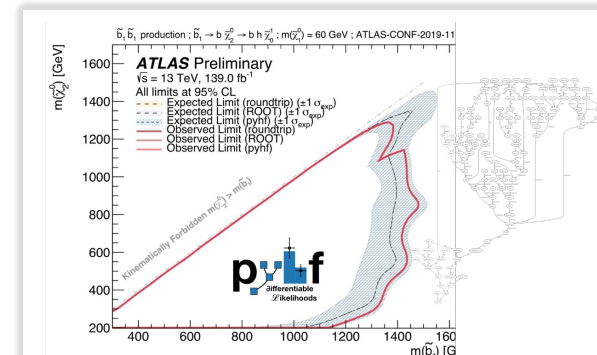
- Can compute gradients of numeric programs



Major enabler of ML

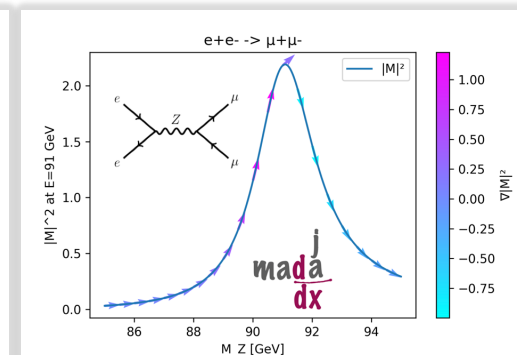


Differentiable programming applied to programs that are not (fully) NNs  
→ powerful way to combine physics & ML



Differentiable Inference

10.21105/joss.02823



Differentiable Matrix Elements

2203.00057

# AD in Programs with Discrete Randomness

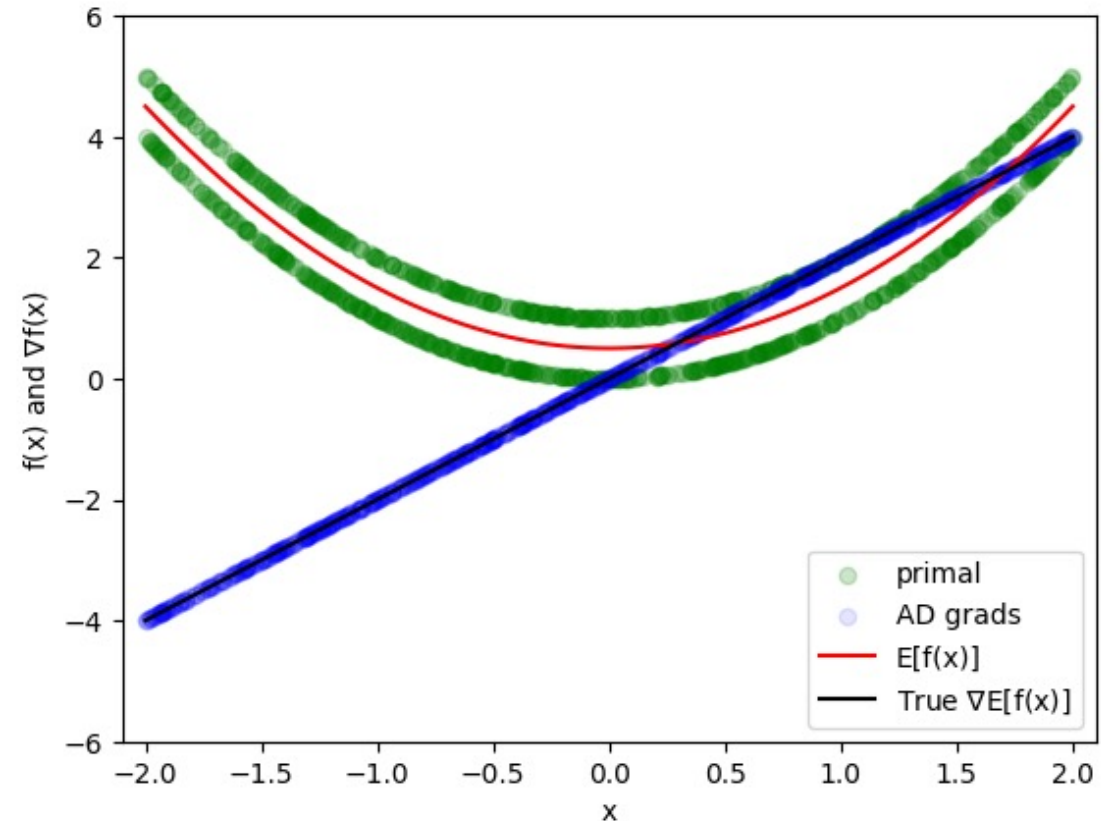
30

```
def f(x):  
    theta = 0.5  
    b = bernoulli(theta)  
    g = x*x  
    return g+b
```

AD Gradient:  $\text{grad}(f(x_i))$

Expected value:  $E[f(x)]$

True gradient:  $\nabla_x E[f(x)]$

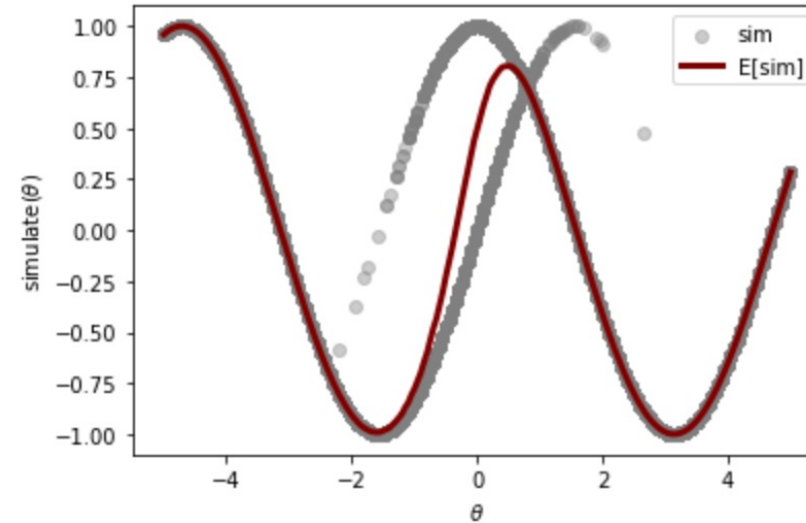


# Programs with Discrete Randomness

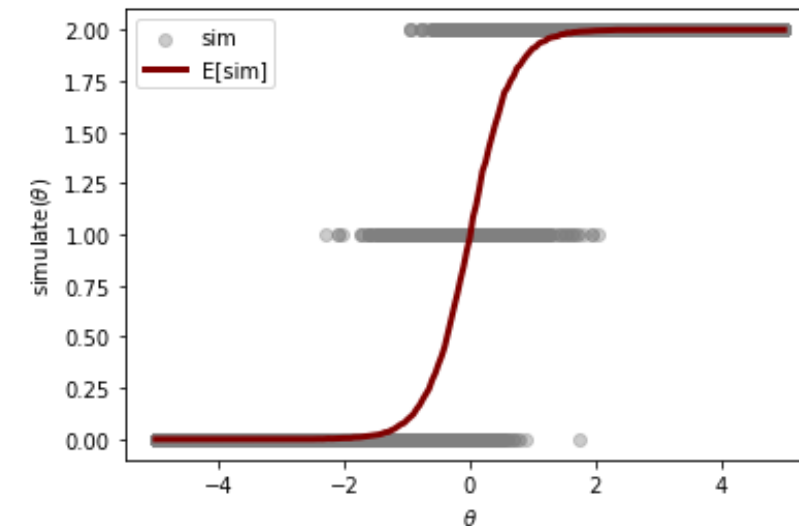
31



```
def simulate(theta):  
    p = sigmoid(theta)  
    x = bernoulli(p) #0 or 1  
    if x == 0:  
        eval = sin(theta)  
    else:  
        eval = cos(theta)  
    return eval
```



```
def program(theta):  
    p = sigmoid(theta)  
    x1 = bernoulli(p) #0/1  
    x2 = bernoulli(p) #0/1  
    eval = x1 + x2 # 0, 1 or 2  
    return eval
```



```
def simulate(t,  $\phi$ ){  
    t' = propagate(t)  
    dointeract ~ p(interact | material(x, y, z |  $\phi$ ))  
    if dointeract:  
        hits.append([x, y, z])  
        dosplit ~ p(split | t')  
        if dosplit:  
            t1, t2 = split(t')  
            simulate(t1), simulate(t2)  
        else:  
            E = (1- $\alpha$ ) E // energy loss  
            simulate(t')  
}
```