

Simulation-Based Inference

ÖAW AI Winter School

Lukas Heinrich

Introduction

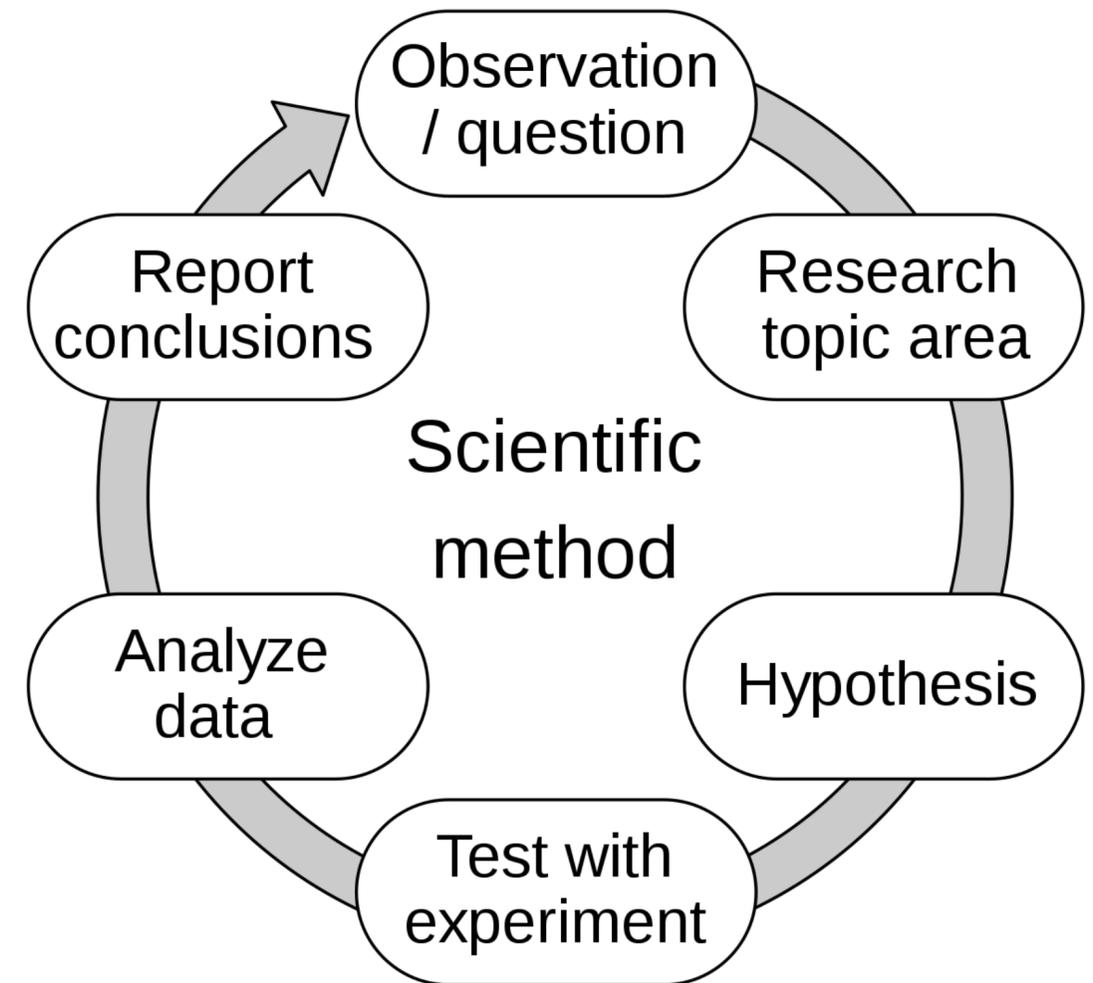
Let's go back to basics and ask ourselves, and look at the

Scientific Method

and how we analyze data.

Scientific Method

In many ways the scientific method relies on the **comparison of hypothesized outcomes with reality**

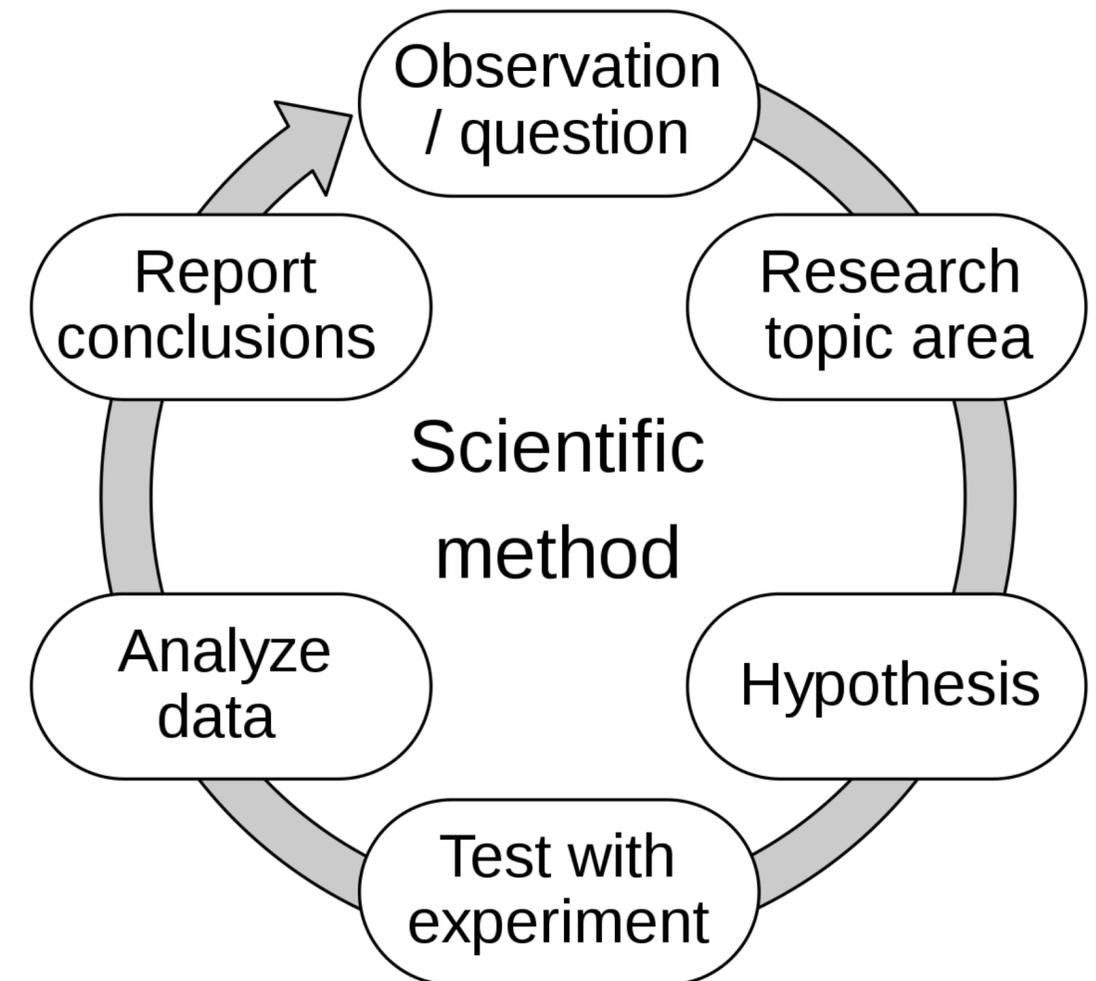


Scientific Method

In many ways the scientific method relies on the comparison of hypothesized outcomes with reality

it's often not the full story, but it's the idealized picture

- sometimes the data is taken without a hypothesis leading to “unexpected results”, etc...

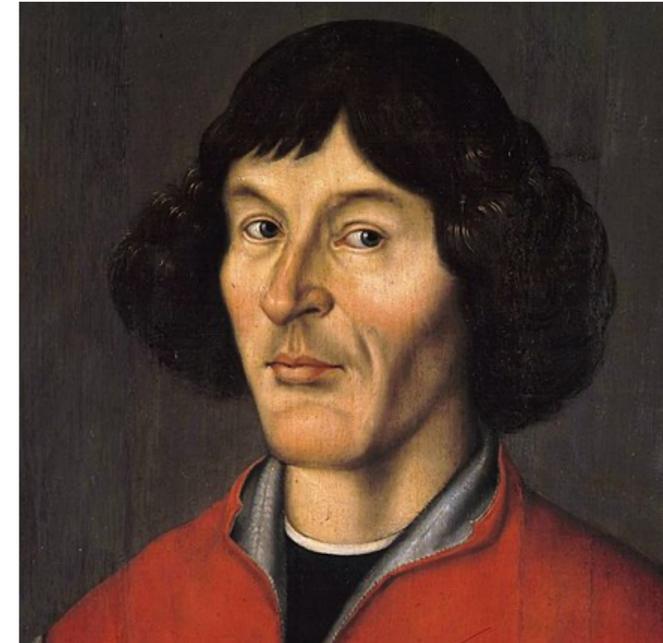
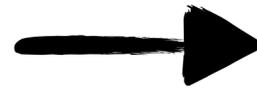


Successes



Claudius Ptolemy

2nd century



Nikolaus Copernicus

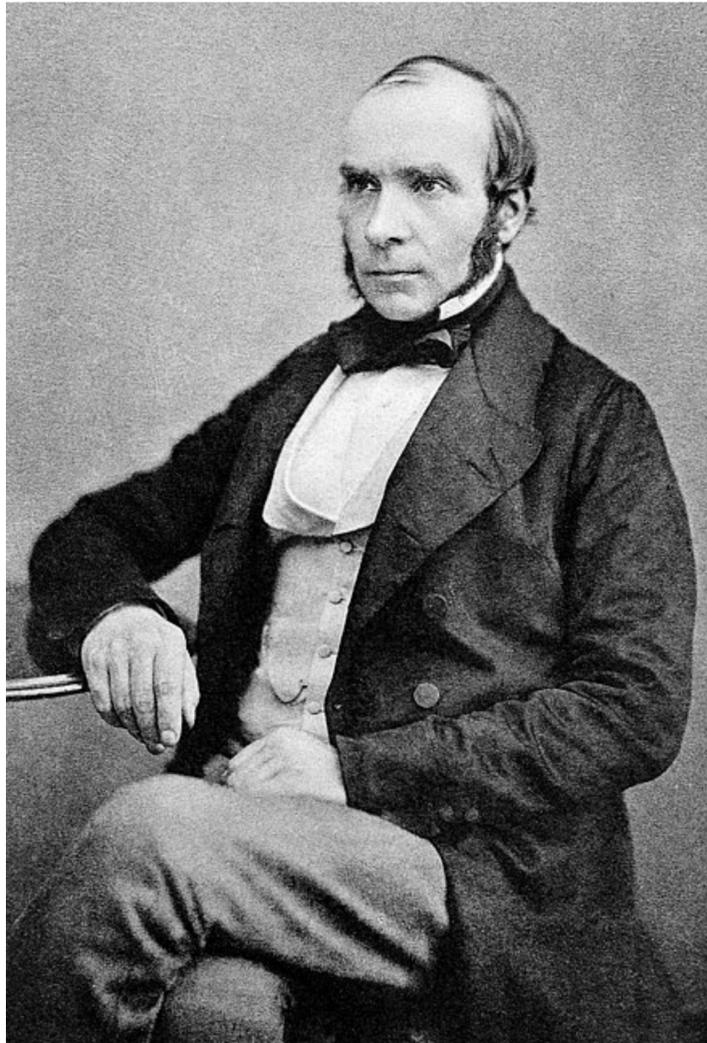
15th century



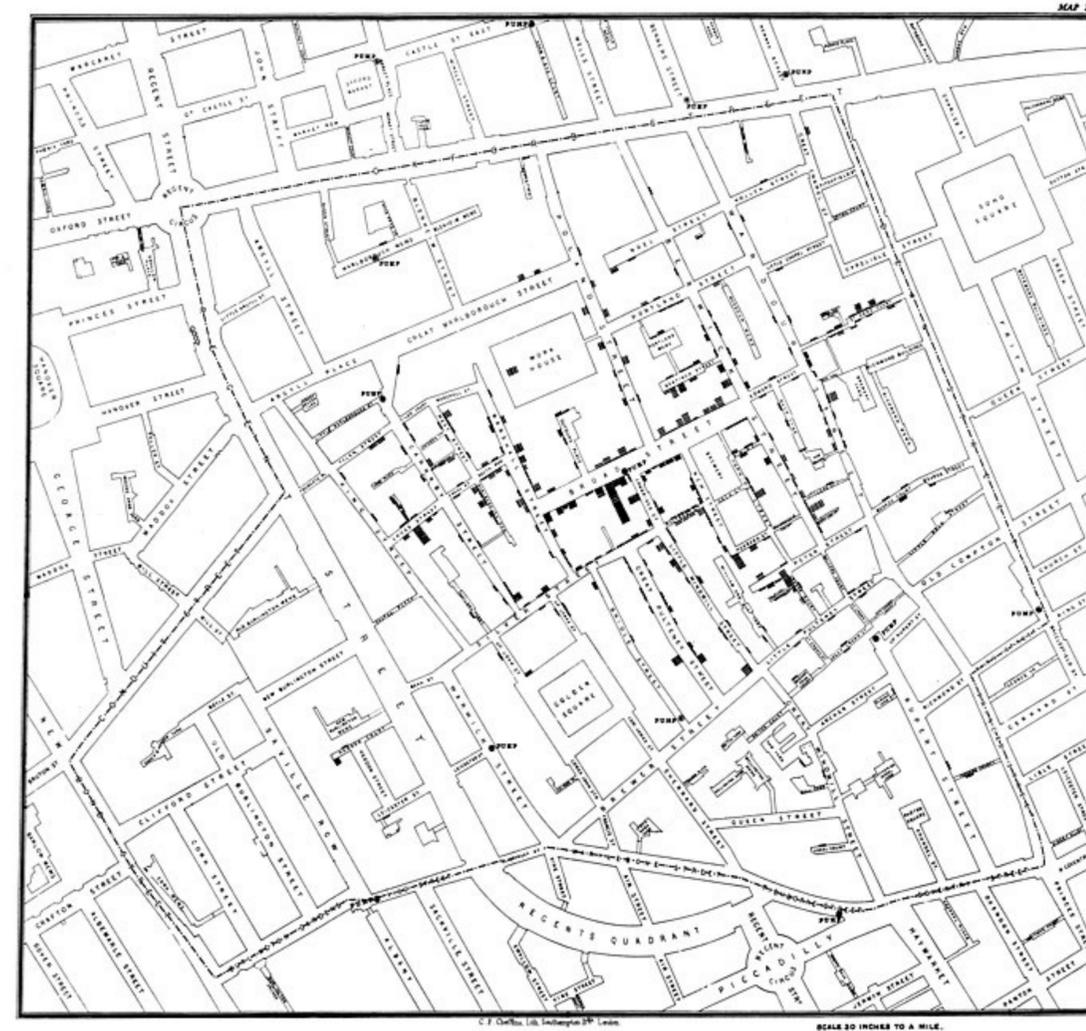
Johannes Kepler

17th century

Successes



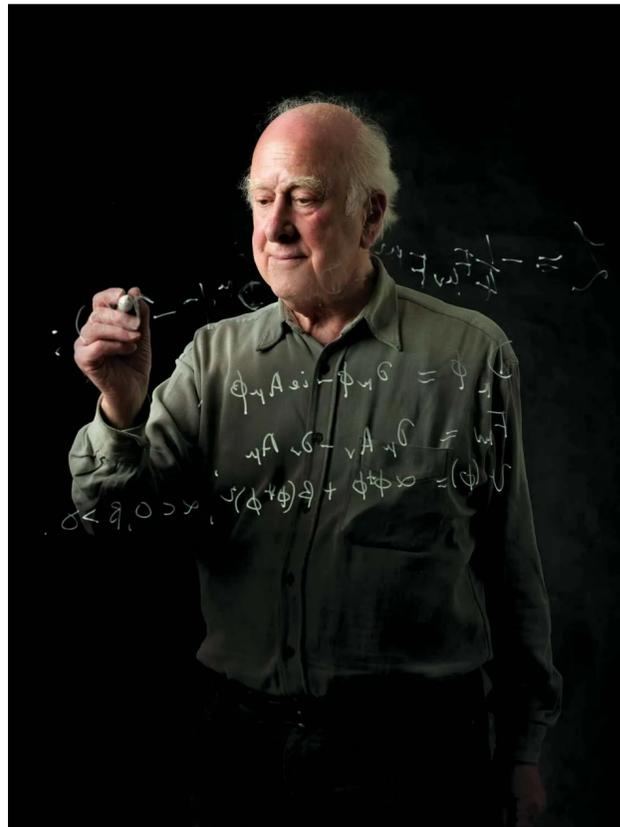
John Snow



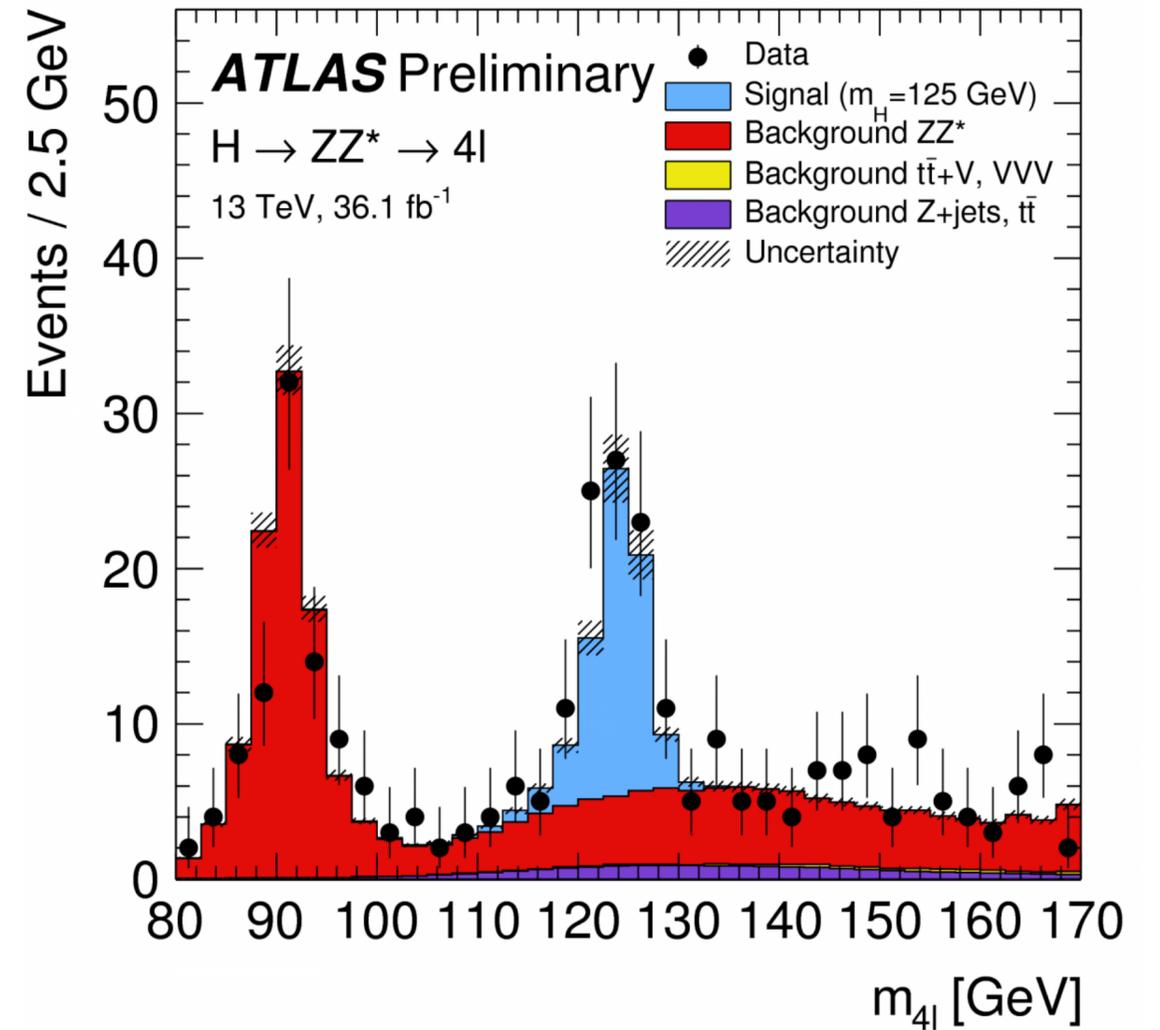
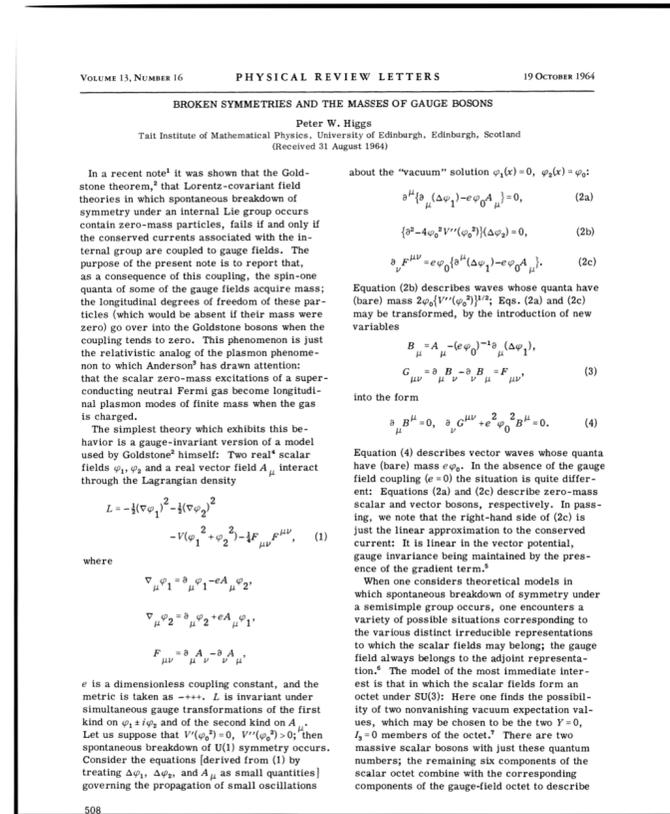
Cholera Outbreak in London

Germ Theory

Successes



Peter Higgs
1964



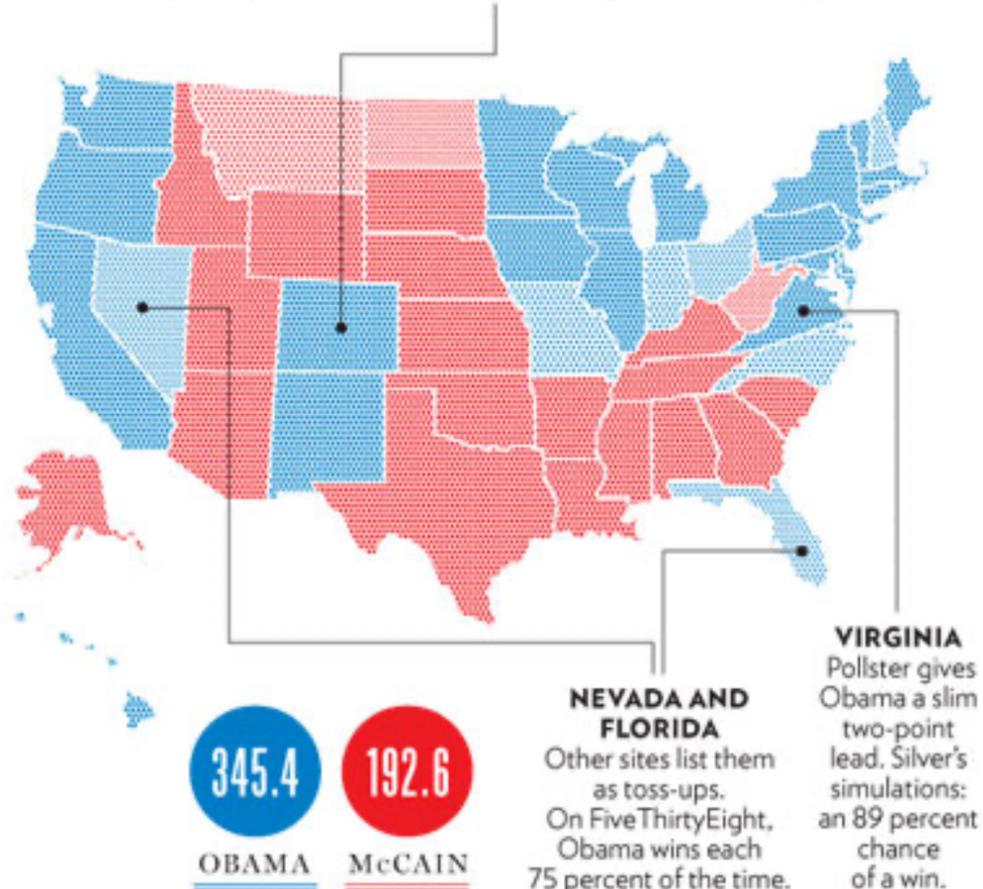
ATLAS and CMS Experiments at
the LHC 2012

Successes

THE ELECTORAL MAP ON NOVEMBER 4

Most electoral-vote maps give a snapshot of where we are now. *FiveThirtyEight's* map, which incorporates weighted averages and simulated scenarios, is a projection of where we'll be on Election Day. Here's how his map looked on October 8, the day after the town-hall debate.

COLORADO RealClearPolitics calls it a toss-up, with McCain trailing Obama by just four points. But Silver's projection gives Obama a seven-point victory—which, when run through a simulation that accounts for the greater certainty of polls as the election nears, means Obama has an 87 percent likelihood of taking the state.



The New York Times

Finding Fame With a Prescient Call for Obama

Give this article



By **Stephanie Clifford**

Nov. 9, 2008

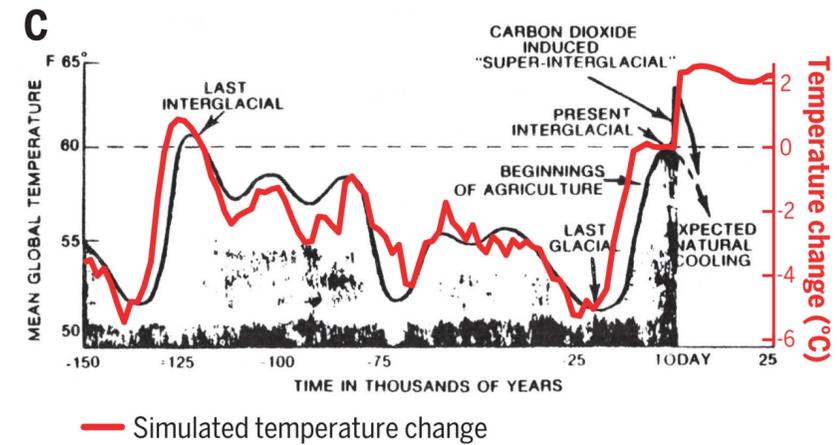
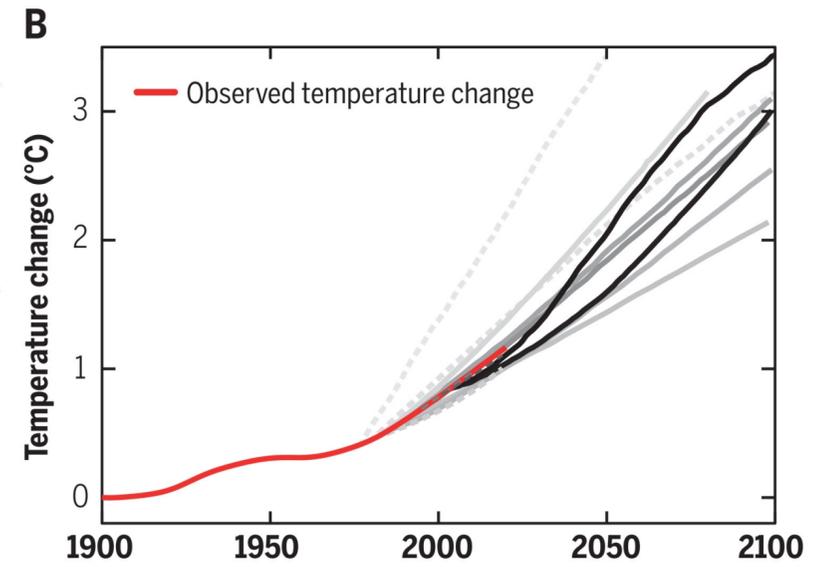
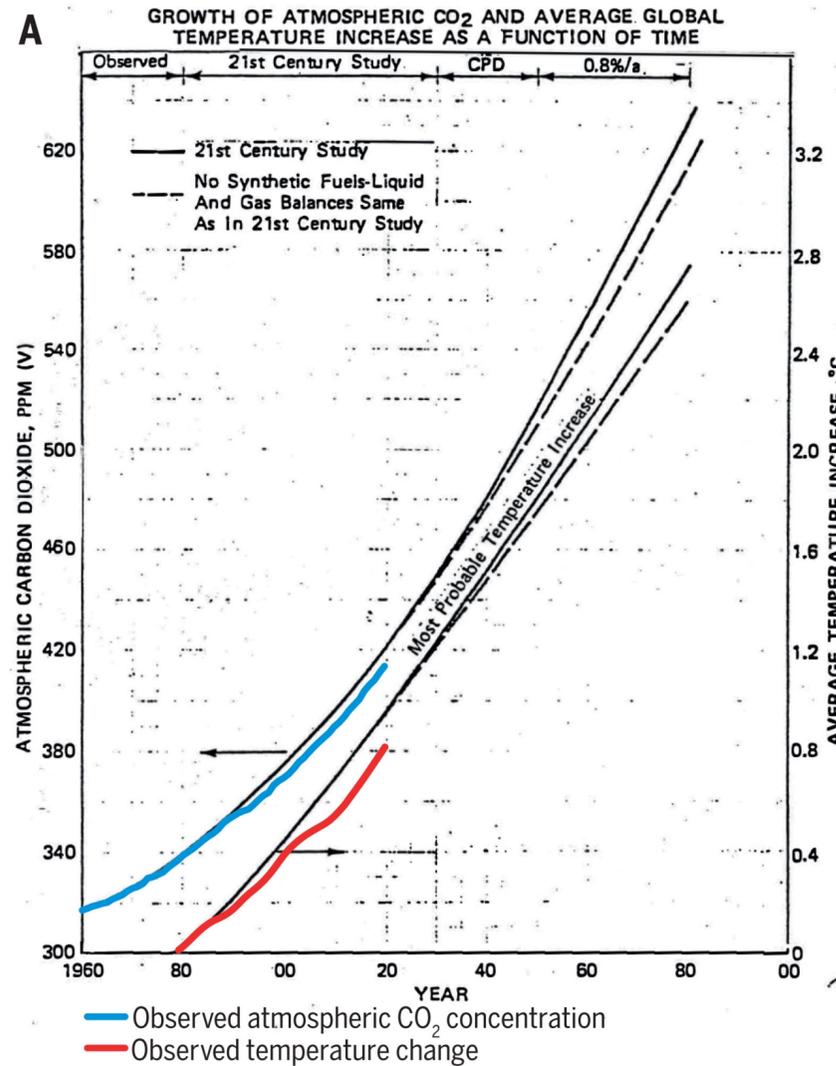
At 9:46 p.m., blogging on his site *FiveThirtyEight.com*, [Nate Silver](#) called the presidential election for Barack Obama. The television networks followed suit about an hour and 15 minutes later after most polls in Western states closed.

Of course, Mr. Silver had a head start: he had forecast that Senator Obama would beat Senator John McCain back in March.

In an election season of unlikely outcomes, Mr. Silver, 30, is perhaps the most unlikely media star to emerge. A baseball statistician who began analyzing political polls only last year, he introduced his site, [FiveThirtyEight.com](#), in March, where he used his own formula to predict federal and state results and run Election Day possibilities based on a host of factors.

Other sites combine polls, notably RealClearPolitics and Pollster, but *FiveThirtyEight*, which drew almost five million page views on Election Day, has become one of the breakout online stars of the

Successes

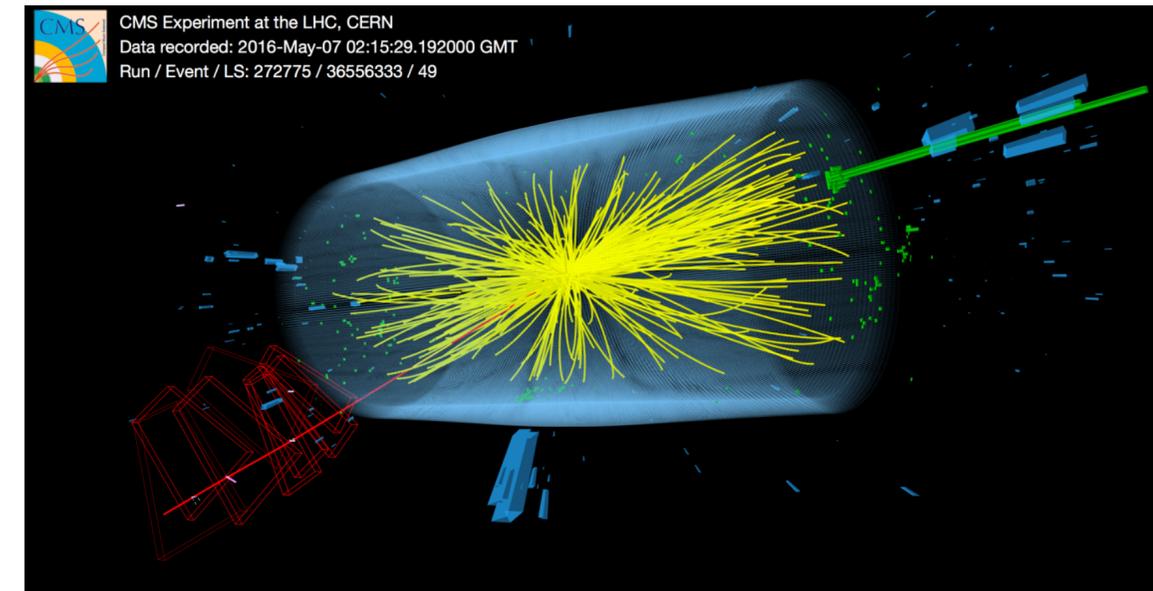
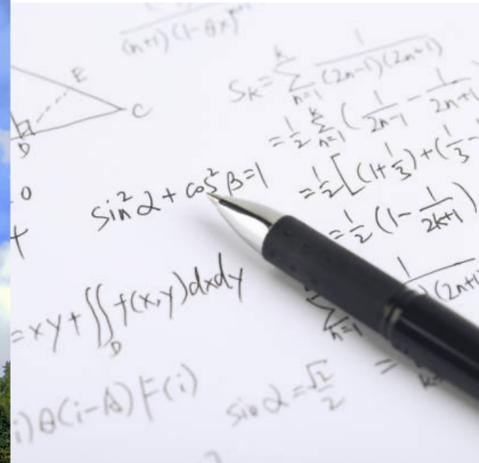
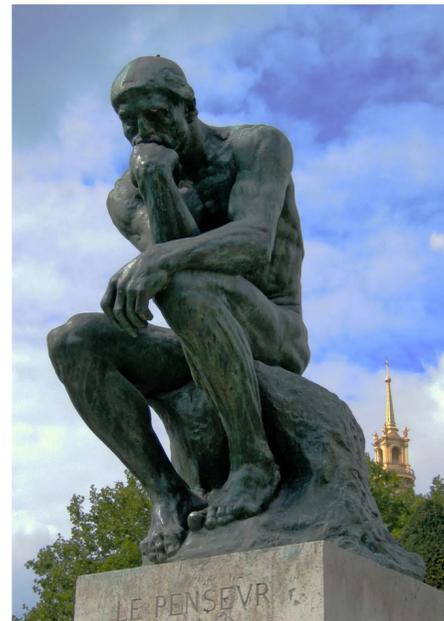


Basic Ingredient

Key aspect in science is not only to formulate the theory but also to **work out what it predicts in the real world**

$$\begin{aligned} \mathcal{L} = & -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} \\ & + i \bar{\psi} \not{D} \psi + h.c. \\ & + \bar{\psi} i \gamma_{ij} \psi \phi + h.c. \\ & + |\mathcal{D}_\mu \phi|^2 - V(\phi) \end{aligned}$$

Theory

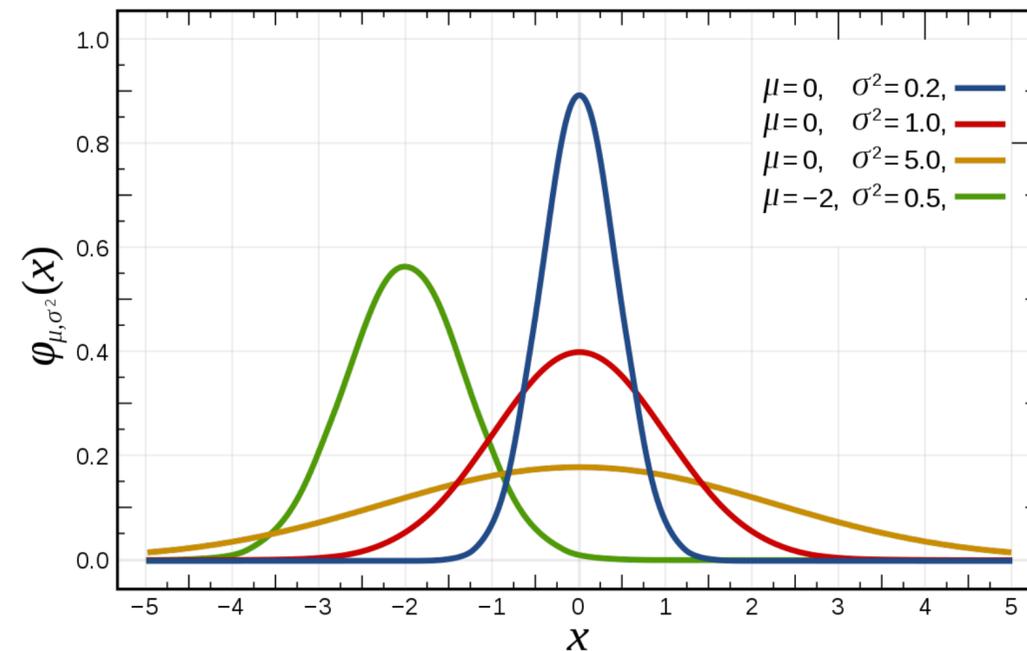


Testable Predictions

But things are never exactly known

We almost **never have exact predictions**, due to inherent or effective randomness of the measurement process

Our predictions are usually statistical in nature



$$p(\text{data} \mid \text{theory})$$

Bayes Theorem

concrete empirical
prediction lives here



$$p(\text{theory} \mid \text{data}) = \frac{p(\text{data} \mid \text{theory})}{p(\text{data})} p(\text{theory})$$



updated beliefs
about the state
of nature

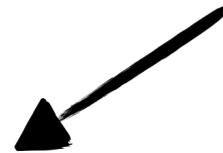


prior beliefs
about the state
of nature

Tractable Models

Standard statistics (i.e. inference) is concerned with the situation, in which you (think you) know the model

$$p(\text{theory} \mid \text{data})$$



Frequentist Statistics

Bayesian Statistics

e.g. Maximum Likelihood Estimates

e.g. Posterior Estimates

Closed-Form Solvable Models

E.g. for Bayesian Statistics, in certain cases we can solve Bayes' Theorem by hand

Advantage:

very fast

Disadvantage:

rarely reflects our true beliefs

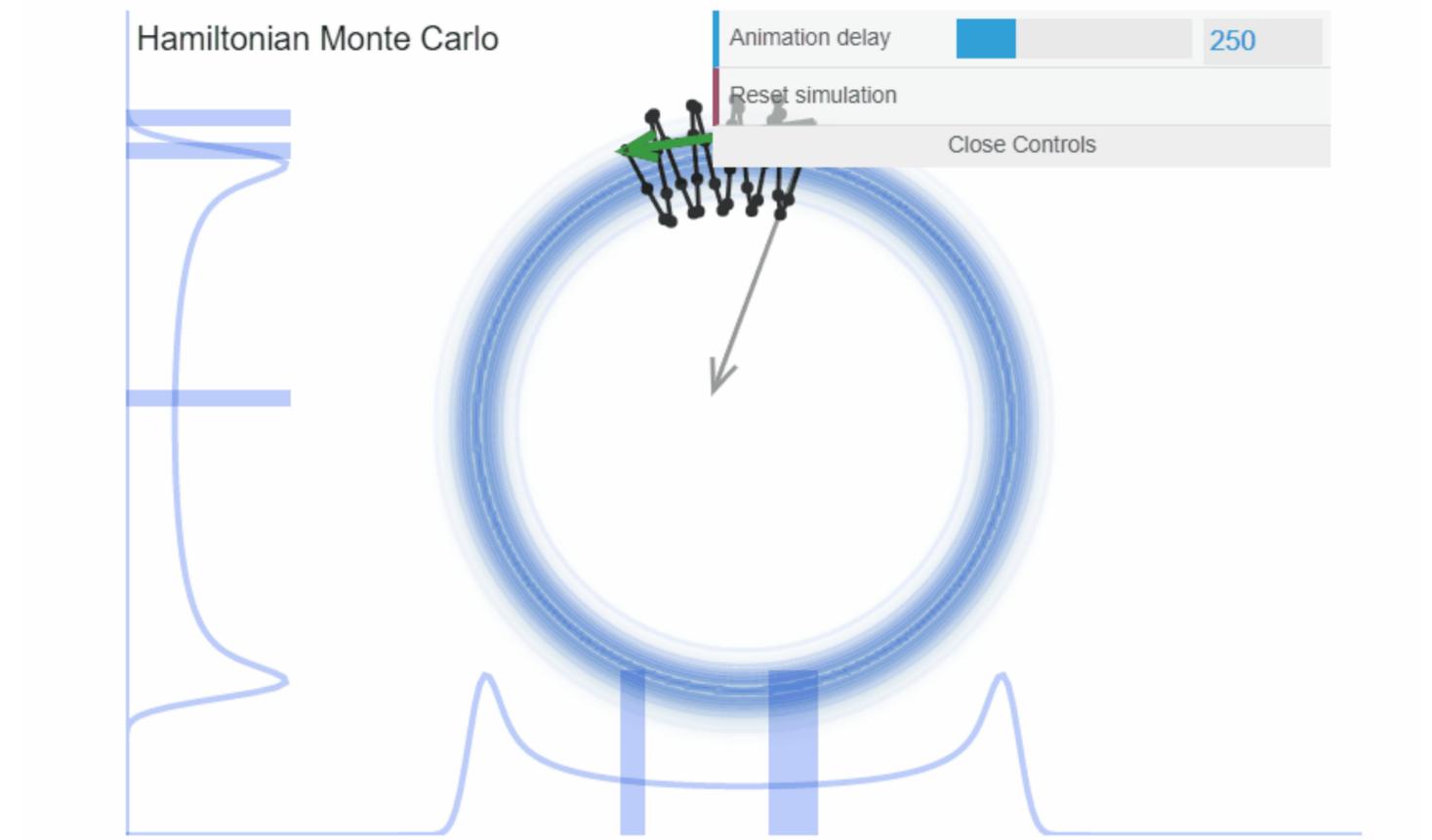
Likelihood	Model parameters	Conjugate prior distribution	Prior hyperparameters	Posterior hyperparameters ^[note 1]	Interpretation of hyperparameters	Posterior predictive ^[note 2]
Bernoulli	p (probability)	Beta	$\alpha, \beta \in \mathbb{R}$	$\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i$	α successes, β failures ^[note 3]	$p(\tilde{x} = 1) = \frac{\alpha'}{\alpha' + \beta'}$
Binomial with known number of trials, m	p (probability)	Beta	$\alpha, \beta \in \mathbb{R}$	$\alpha + \sum_{i=1}^n x_i, \beta + \sum_{i=1}^n N_i - \sum_{i=1}^n x_i$	α successes, β failures ^[note 3]	BetaBin($\tilde{x} \alpha', \beta'$) (beta-binomial)
Negative binomial with known failure number, r	p (probability)	Beta	$\alpha, \beta \in \mathbb{R}$	$\alpha + rn, \beta + \sum_{i=1}^n x_i$	α total successes, β failures ^[note 3] (i.e., $\frac{\beta}{r}$ experiments, assuming r stays fixed)	BetaNegBin($\tilde{x} \alpha', \beta'$) (beta-negative binomial)
Poisson	λ (rate)	Gamma	$k, \theta \in \mathbb{R}$	$k + \sum_{i=1}^n x_i, \frac{\theta}{n\theta + 1}$	k total occurrences in $\frac{1}{\theta}$ intervals	NB($\tilde{x} k', \frac{1}{\theta' + 1}$) (negative binomial)
			α, β ^[note 4]	$\alpha + \sum_{i=1}^n x_i, \beta + n$	α total occurrences in β intervals	NB($\tilde{x} \alpha', \frac{\beta'}{1 + \beta'}$) (negative binomial)
Categorical	\mathbf{p} (probability vector), k (number of categories; i.e., size of \mathbf{p})	Dirichlet	$\boldsymbol{\alpha} \in \mathbb{R}^k$	$\boldsymbol{\alpha} + (c_1, \dots, c_k)$, where c_i is the number of observations in category i	α_i occurrences of category i ^[note 3]	$p(\tilde{x} = i) = \frac{\alpha_i'}{\sum_i \alpha_i'} = \frac{\alpha_i + c_i}{\sum_i \alpha_i + n}$
	\mathbf{p}					

Numeric Solutions (e.g. MCMC)

With increasing computational power, more we started to learn how to solve more realistic problems

MCMC Methods: for known but complex $p(x | \theta)$ **we cannot compute $p(\theta | x)$ but we can sample from it!**

$$\theta \sim p(\theta | x)$$



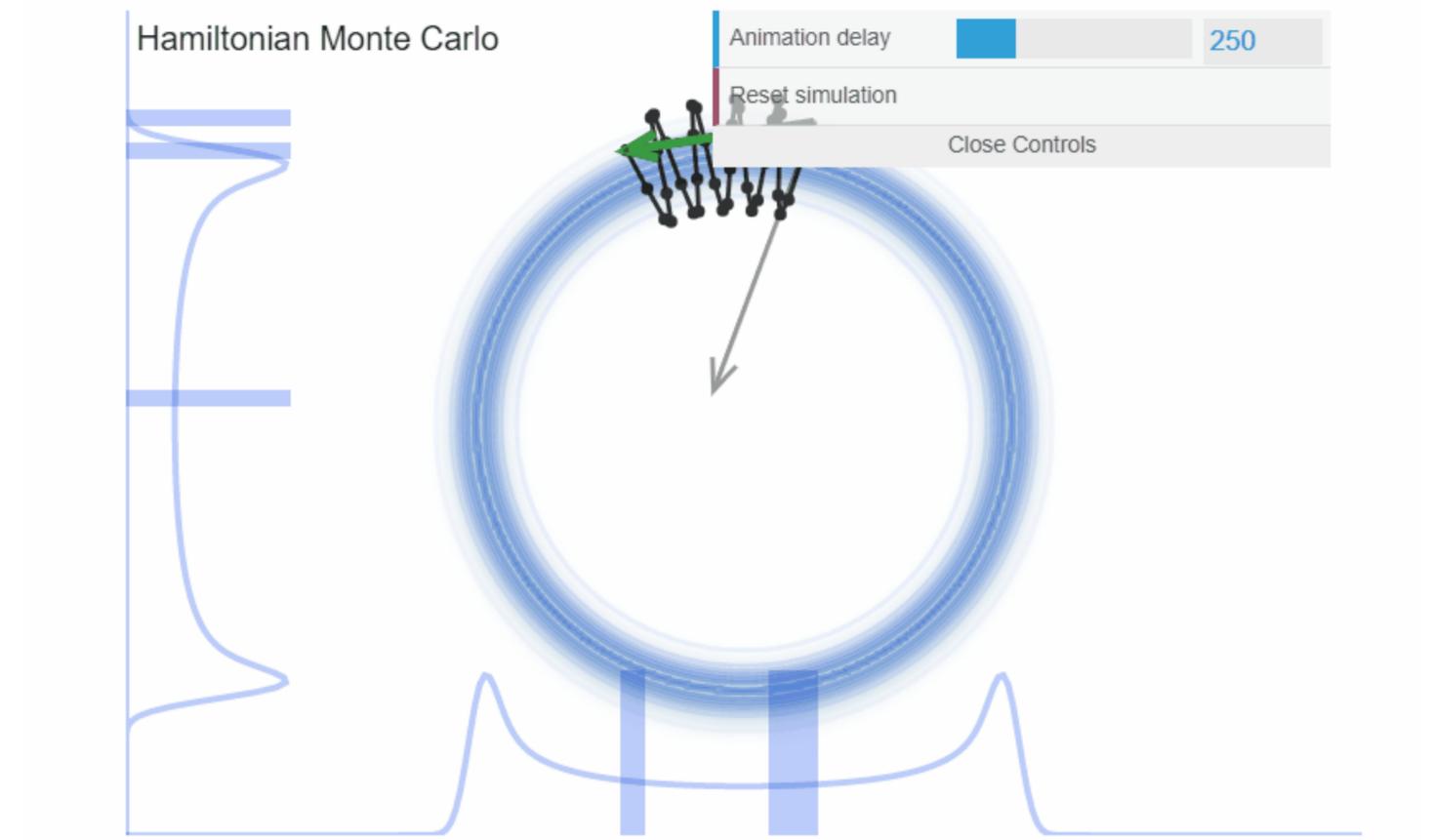
Gradients

We just discussed automatic differentiation. MCMC is even more powerful if you can not only compute $l(\theta) = \log p(x | \theta)$

but also its gradient...

$$\nabla_{\theta} \log p(x | \theta)$$

Hamiltonian Monte Carlo

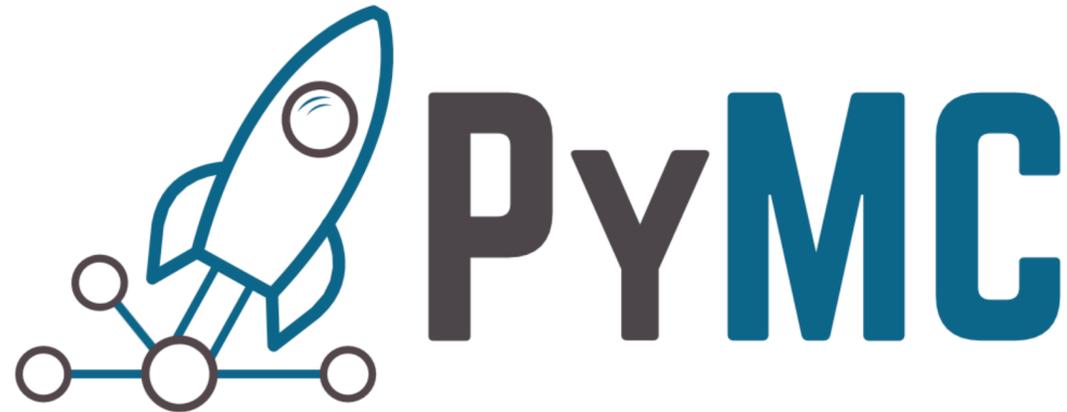


Gradients

We just discuss
more powerful
but also its gra

$$\nabla_{\theta} \log$$

Hamiltonian



PyMC is a probabilistic programming library for Python that allows users to build Bayesian models with a simple Python API and fit them using Markov chain Monte Carlo (MCMC) methods.

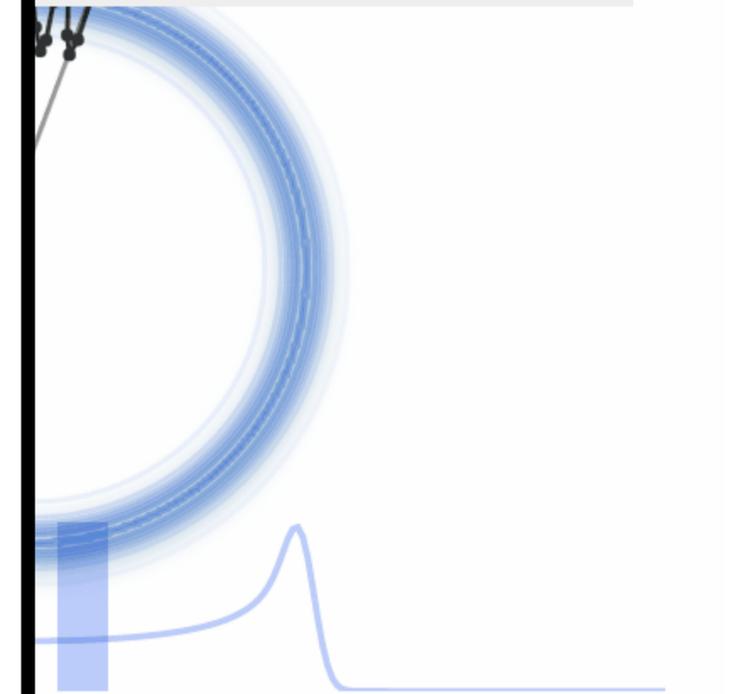
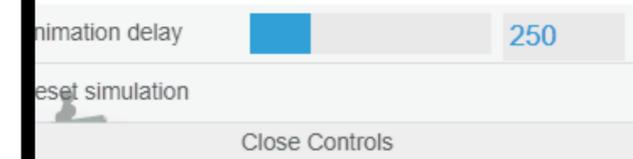
Features

PyMC strives to make Bayesian modeling as simple and painless as possible, allowing users to focus on their problem rather than the methods.

Here is what sets it apart:

- **Modern:** Includes state-of-the-art inference algorithms, including MCMC (NUTS) and variational inference (ADVI).
- **User friendly:** Write your models using friendly Python syntax. [Learn Bayesian modeling](#) from the many [example notebooks](#).
- **Fast:** Uses [Aesara](#) as its computational backend to compile to C and JAX, [run your models on the GPU](#), and benefit from complex graph-optimizations.

MCMC is even
 $\log p(x | \theta)$



Bad News

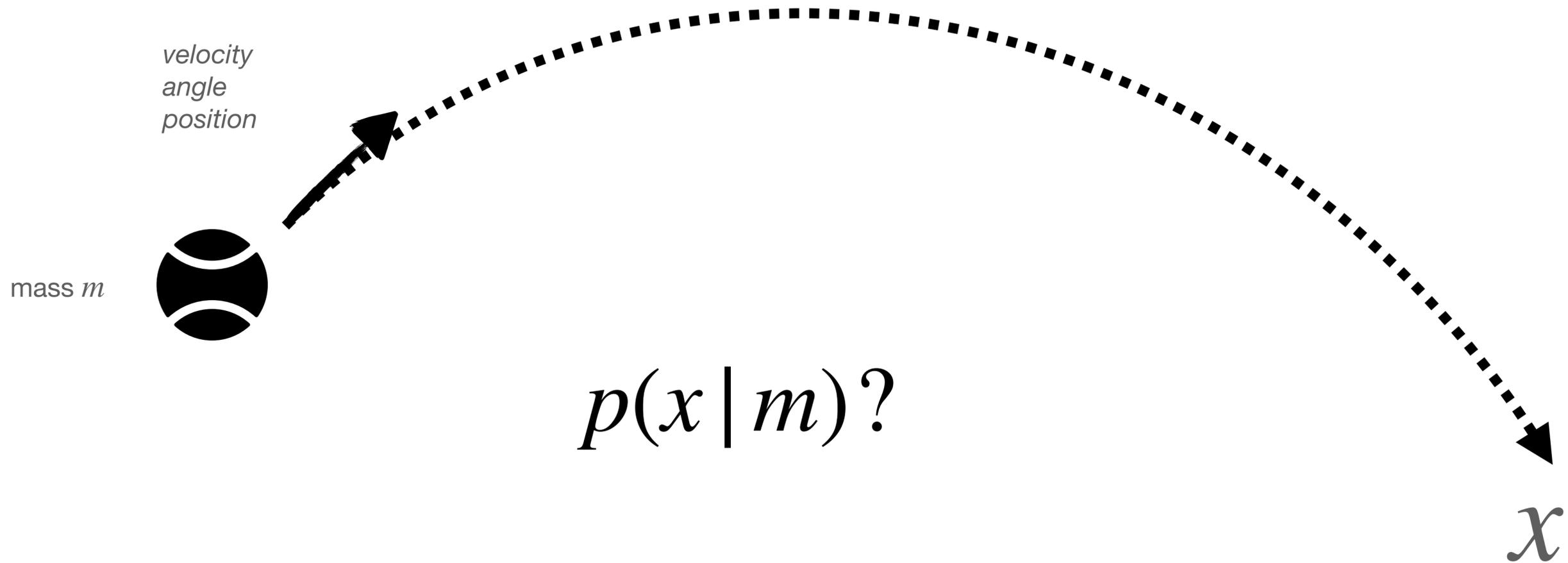
But what if ... we didn't know how to compute

$$p(\text{data} \mid \text{theory})$$

(this happens much more often than we'd like)

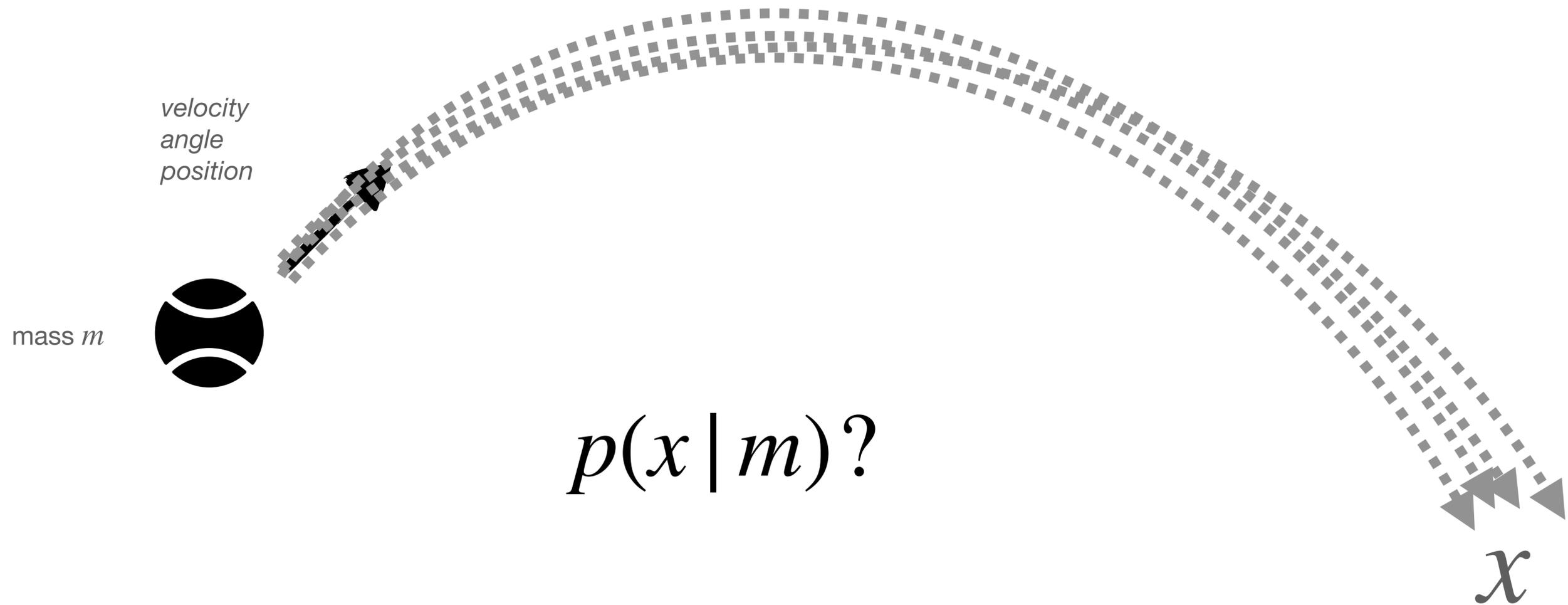
Example

What's the probability for a ball to measure a ballistic ball launched at an angle at a final position x ?



Example

What's the probability for a ball to land at position x ?



Example

Need to integrate over all possible unobserved paths for a given final outcome

$$p(x | m) = \int p(x | f_{\text{physics}}(v, \alpha, \dots)) p(v, \alpha, \dots | m) ?$$

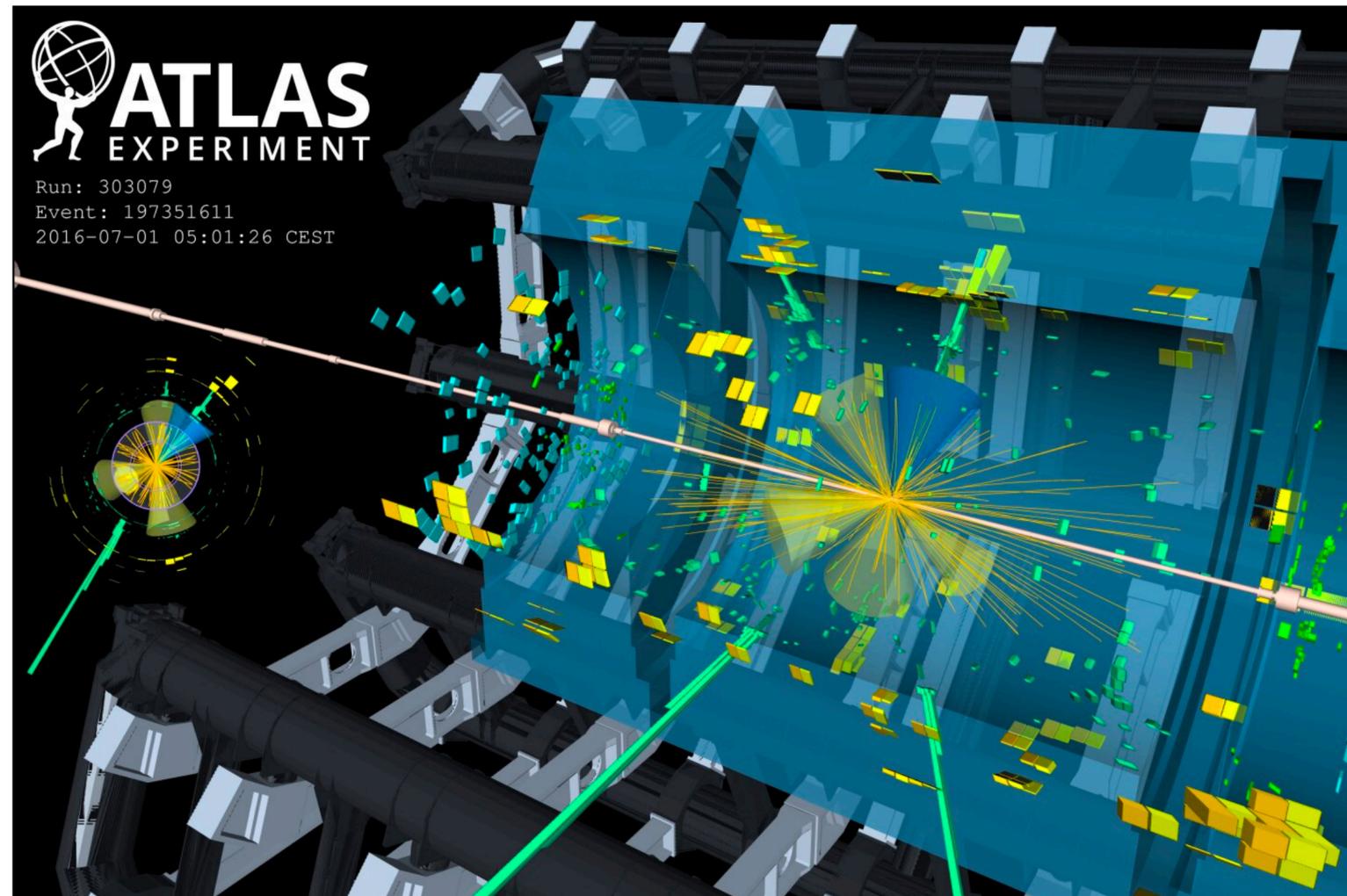
measured position
is likely distributed
over theoretically
calculated one

predicted trajectory
of ball based on physics

probabilities to launch
at given speed / velocity

A slightly more complicated example...

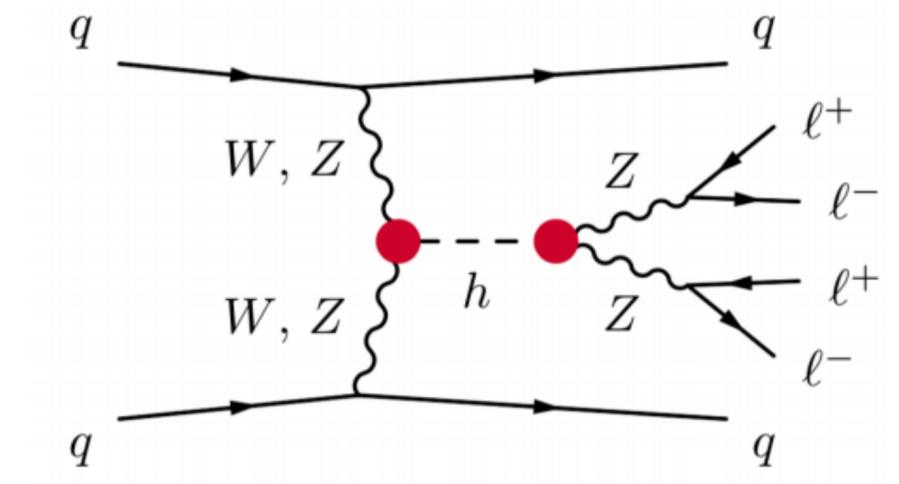
Given a Higgs mass, what do the 100M sensors in ATLAS read?



Why?

Because there is a huge gap between theory and data we need to bridge.

*Core Interaction
(Matrix Elements)*

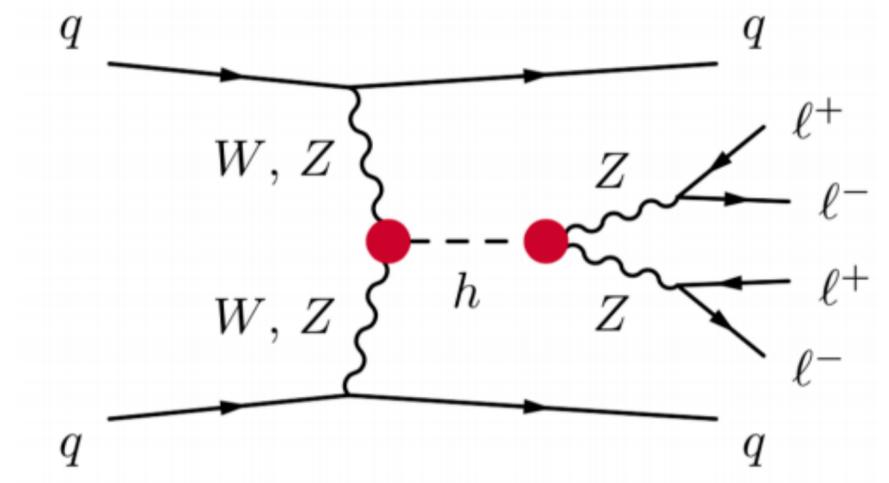
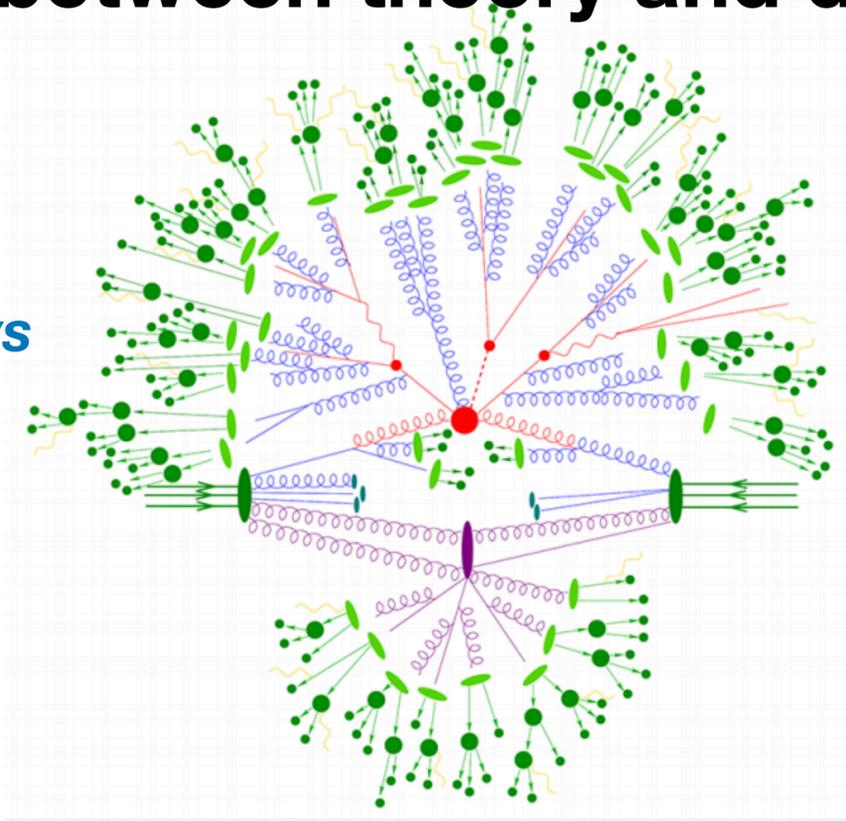


$$p(z_p | \theta)$$

Why?

Because there is a huge gap between theory and data we need to bridge.

*Secondary, Tertiary, n-ary decays
(hadronization, parton shower)*

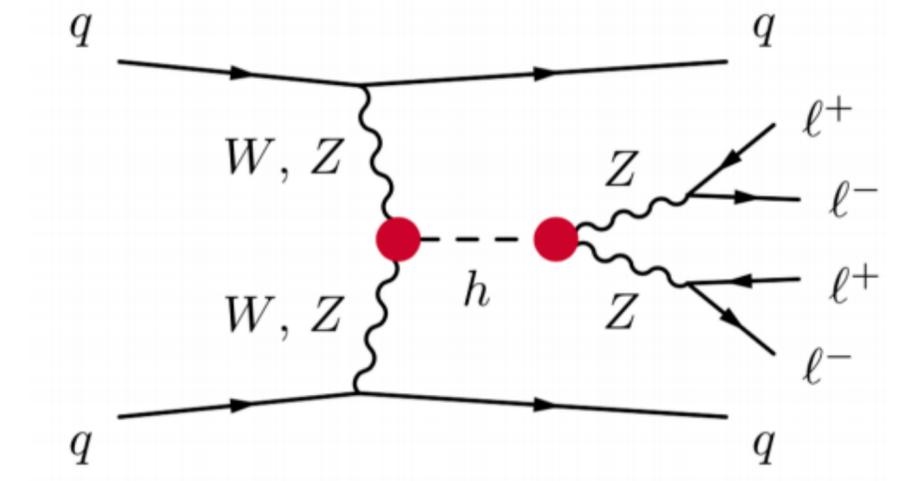
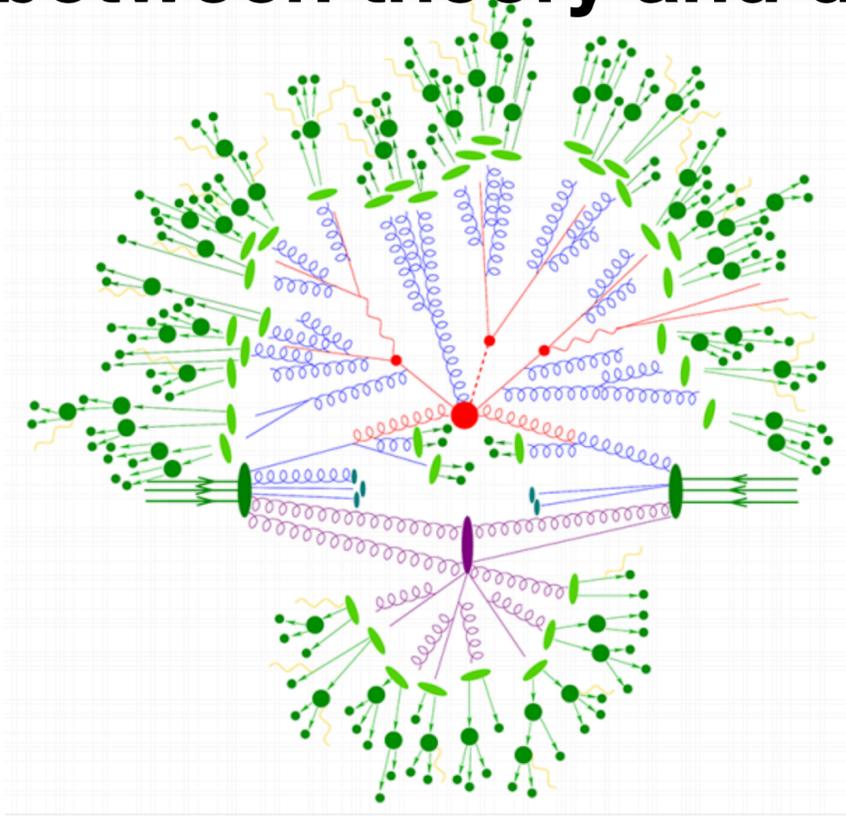
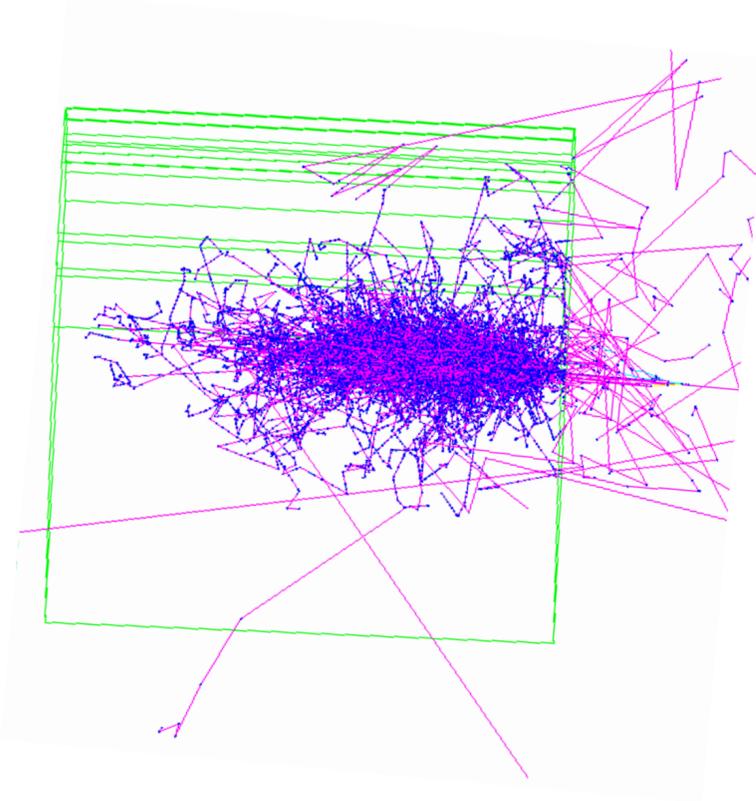


$$p(z_h | z_p) p(z_p | \theta)$$

Why?

Because there is a huge gap between theory and data we need to bridge.

*sensor interaction
(particle transport)*

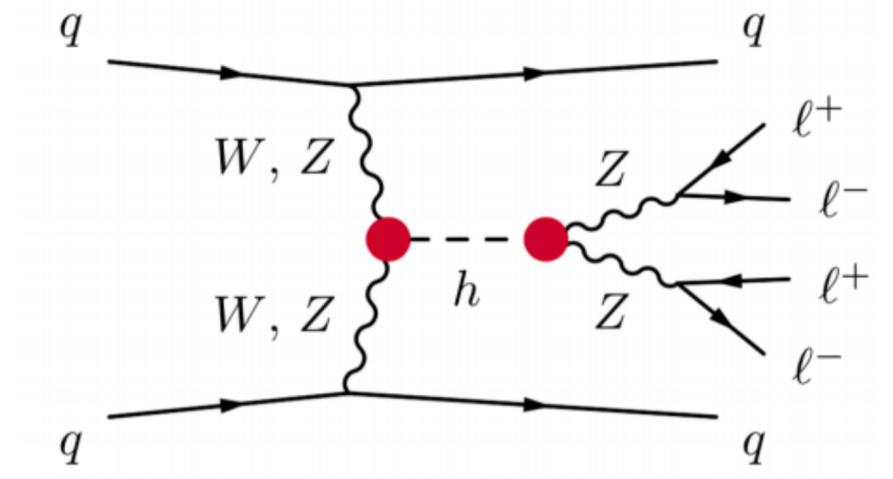
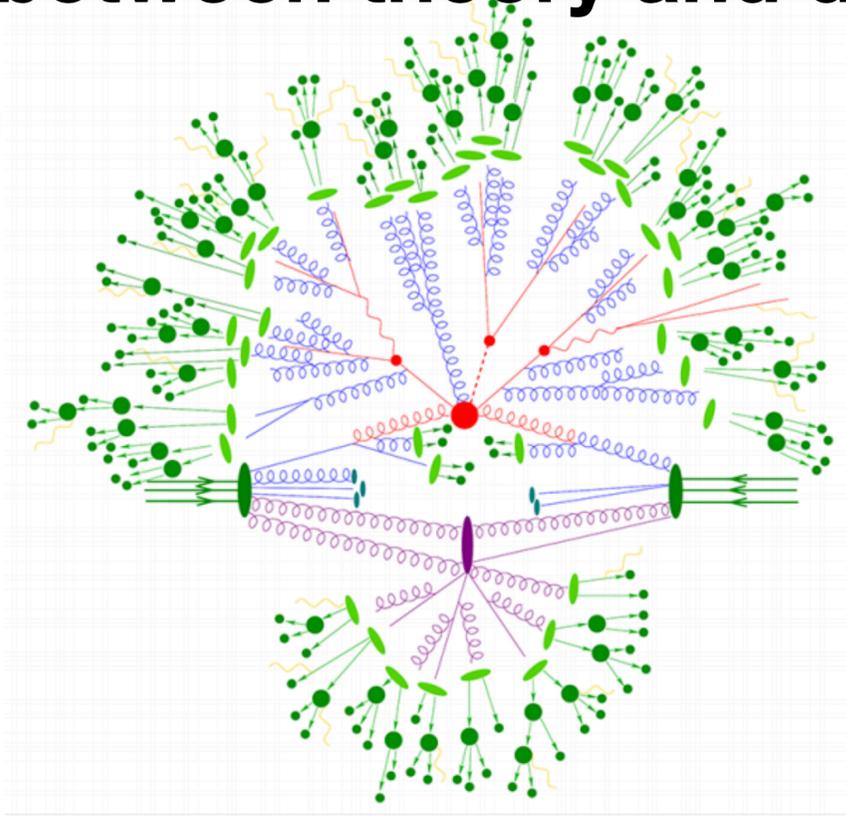
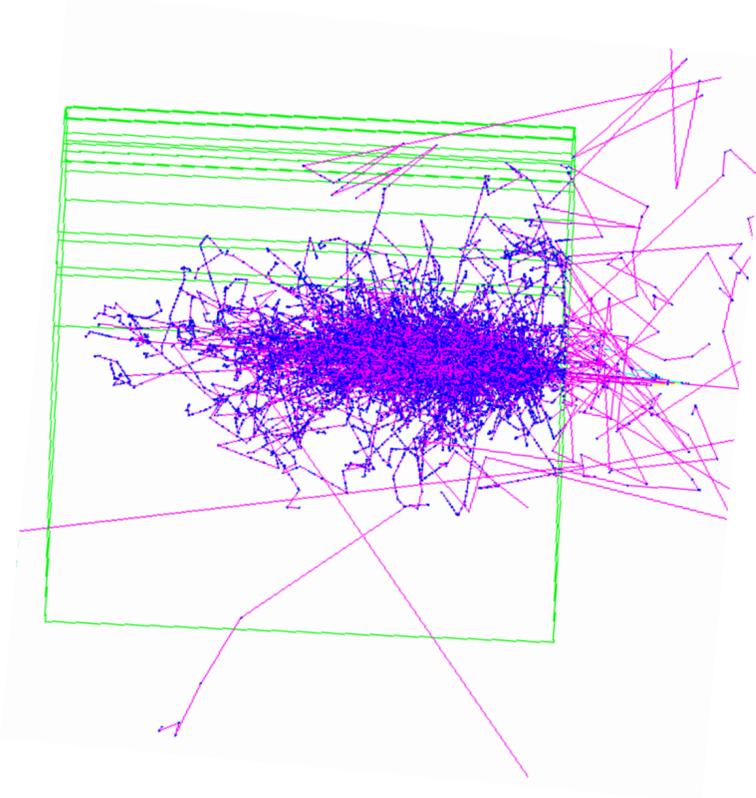


$$p(x | z_h) p(z_h | z_p) p(z_p | \theta)$$

Why?

Because there is a huge gap between theory and data we need to bridge.

*sensor interaction
(particle transport)*



*all of this is quantum-mechanical and
thus random, we only observe x at the end
→ integrate*

$$p(x | \theta) = \int dz \ p(x | z_h) p(z_h | z_p) p(z_p | \theta)$$

Example

More generally, often $p(x | \theta)$ is intractable because of hidden unobserved dynamics / stochastic processes

- **histories of how the data came to be**
- **intractable even if we know (in principle) all the details (!)**



Example

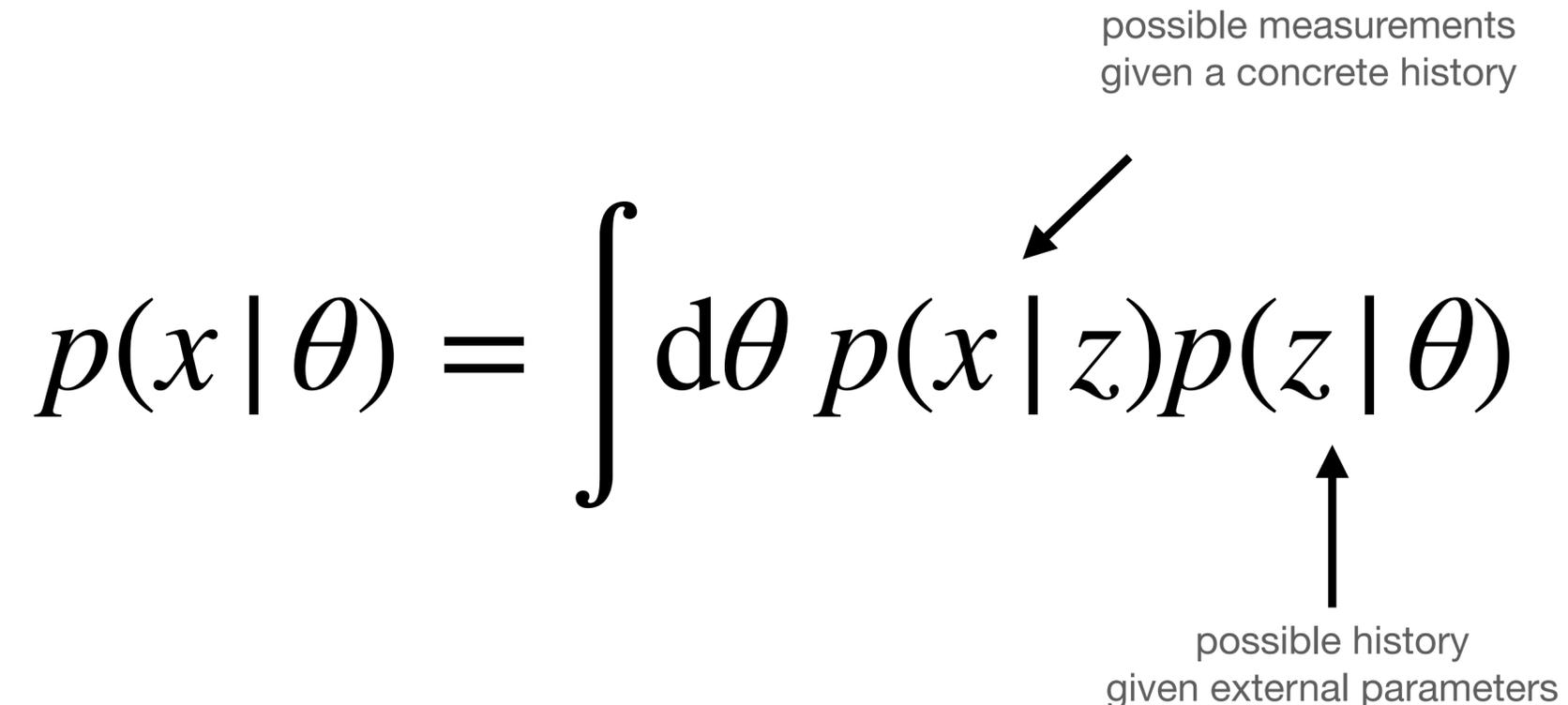
More generally, often $p(x | \theta)$ is intractable because of hidden unobserved dynamics / stochastic processes

Mathematically: intractable integral

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

possible measurements
given a concrete history

possible history
given external parameters

The diagram shows the equation $p(x | \theta) = \int d\theta p(x | z)p(z | \theta)$. An arrow points from the text "possible measurements given a concrete history" to the variable z in the term $p(x | z)$. Another arrow points from the text "possible history given external parameters" to the variable θ in the term $p(z | \theta)$.

Now what?

If we cannot calculate the values of $p(x | \theta)$ how could we possibly do statistics ? Are we doomed?

No! Let's take a closer look at $p(x | \theta)$

The Two Faces of $p(x)$

Info about $p(x)$? What is $p(x)$? It's two things at once:



The Two Faces of $p(x)$

Info about $p(x)$? What is $p(x)$? It's two things at once:

A formula



$$\mathbb{R}^2 \rightarrow \mathbb{R}$$

$$p_{\mu, \Sigma}(x) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma (x - \mu)\right)$$

*Evaluating the Probability
for a given sample*

The Two Faces of $p(x | \theta)$

Info about $p(x)$? What is $p(x)$? It's two things at once:

A process

$$\text{Die} \rightarrow \mathbb{R}^2$$

$$x \sim p(x | \theta)$$

*Generating new samples
from randomness*



A formula

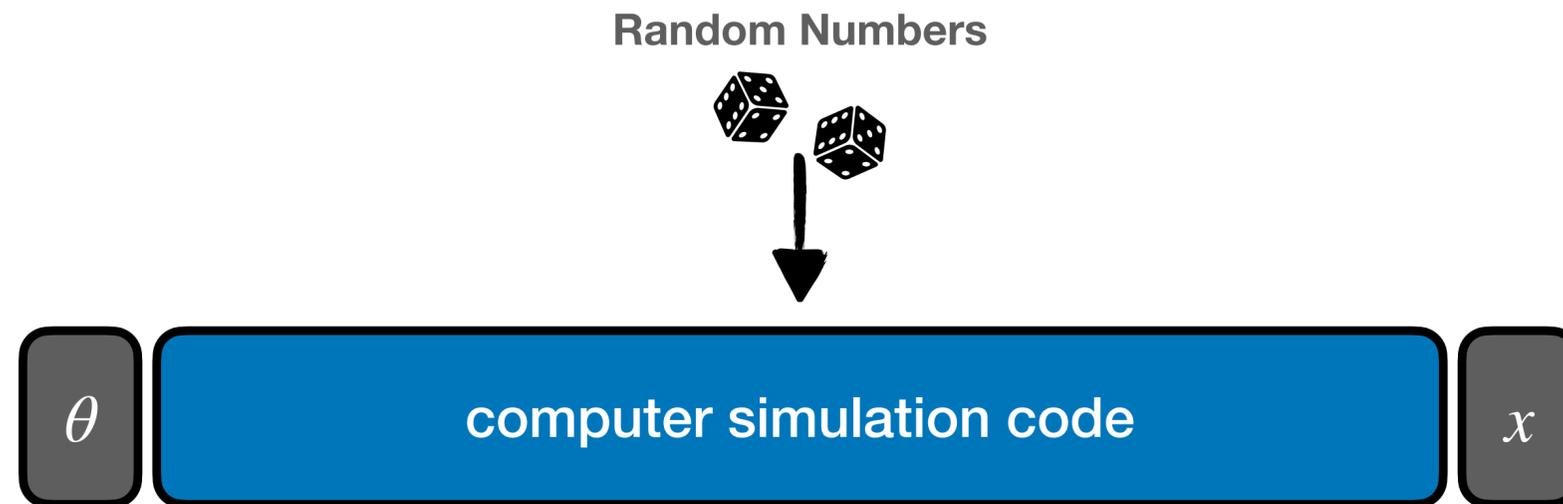
$$\mathbb{R}^2 \rightarrow \mathbb{R}$$

$$p_{\mu, \Sigma}(x) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma (x - \mu)\right)$$

*Evaluating the Probability
for a given sample*

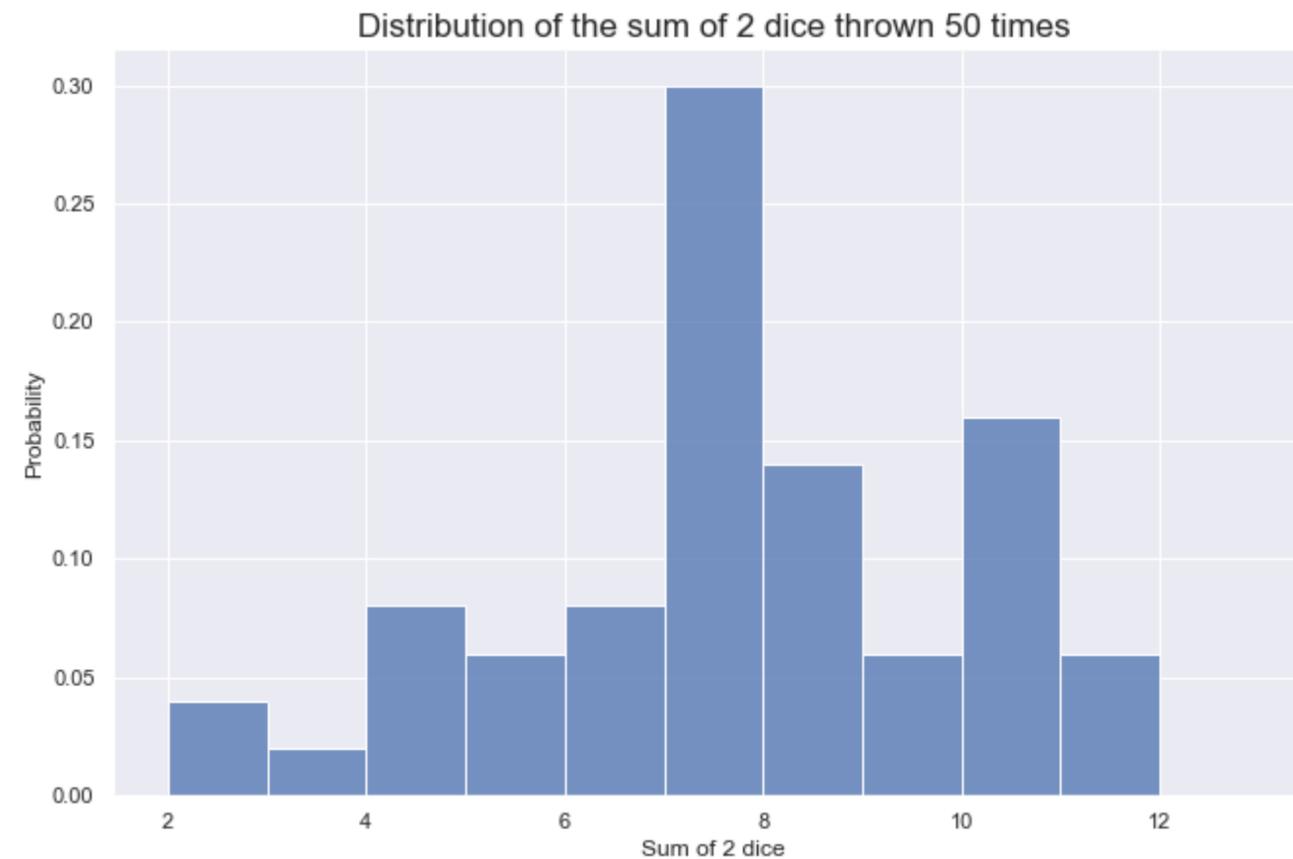
Writing Samplers of Physical Experiments

If we understand the mechanics of the data-generating process it's feasible to write “stochastic simulators” that implicitly implement $p(x | \theta)$



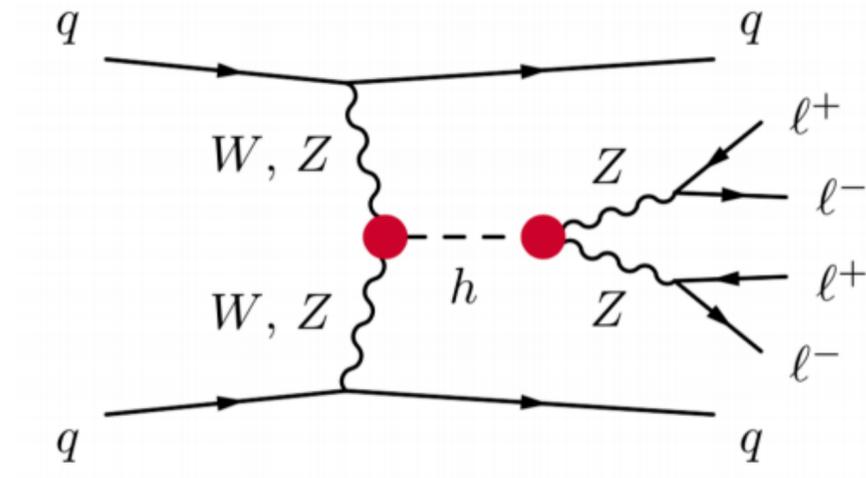
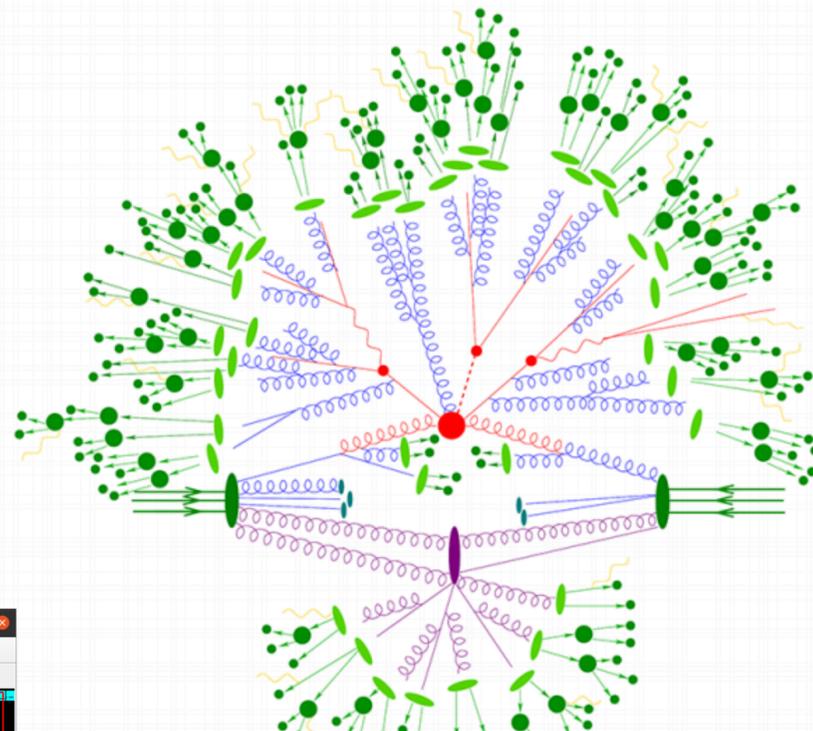
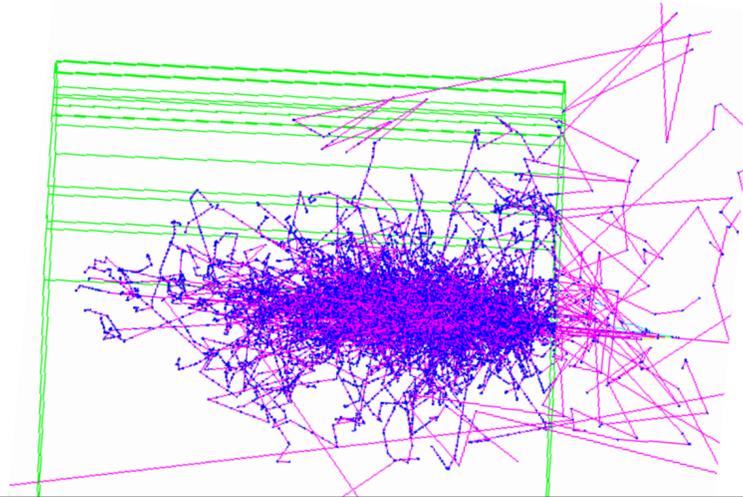
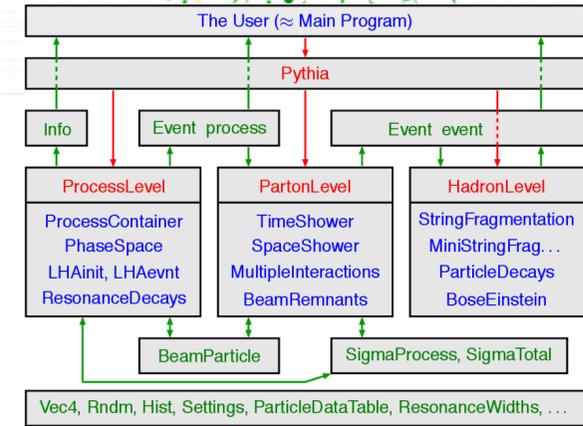
Example:

Without calculating anything, we can easily “approximately” determine $p(x | \theta)$ by analyzing **samples from** $x \sim p(x | \theta)$



Simulators are Everywhere

Stochastic Simulation Programs are how we can model complex physical systems. **The Science lives in Software**

Center for Particle Physics and Phenomenology - CP3

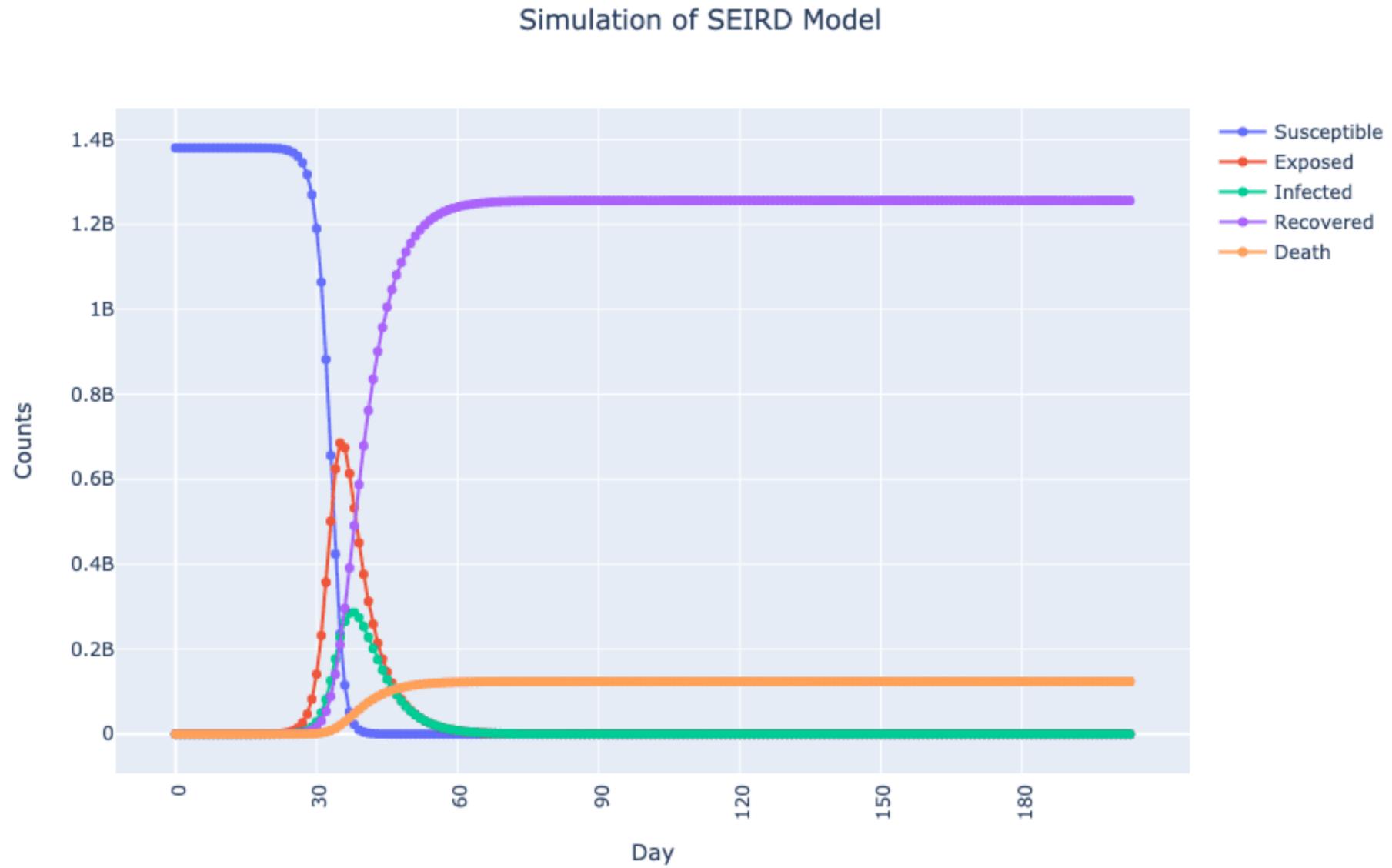
[The MadGraph5_aMC@NLO homepage](#)

UCL Milan Launchpad Github
by the MG5aMC@NLO Development team

Generate Process [aMC@NLO](#) [Tools](#) [Downloads](#) [Wiki](#) [Answers](#) [Bug reports](#) [Citation](#)

Server Maintenance in progress

Simulators are Everywhere



Epidemiology Simulation

Simulation-based Inference

Simulation-based Inference:

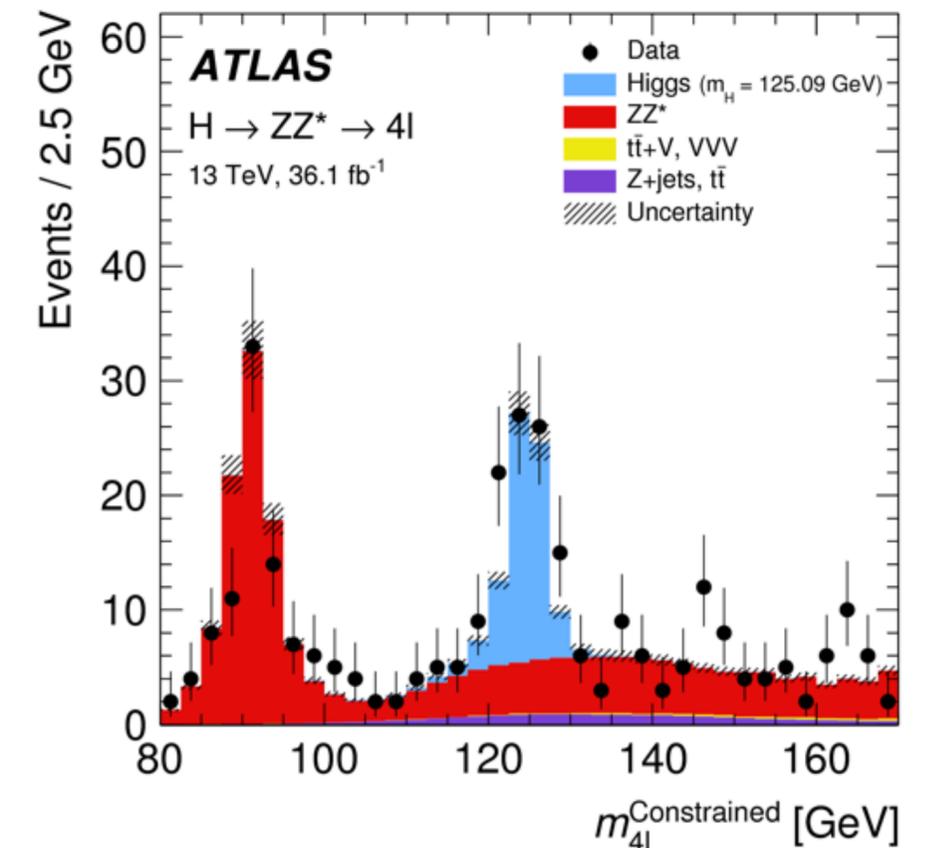
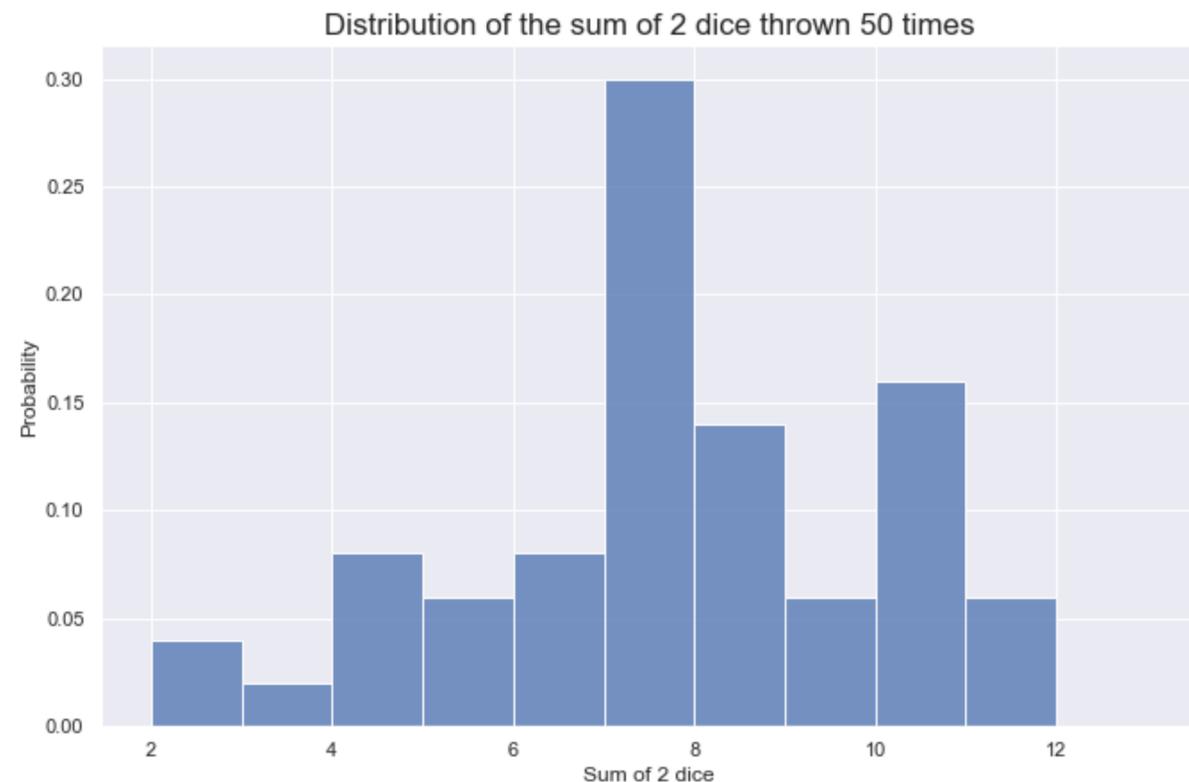
statistical techniques to apply when

- you don't have a tractable likelihood $L(\theta) = p(x | \theta)$
- but do have a simulator that can give you faithful samples $x \sim p(x | \theta)$

Simulation-based Inference

Simplest method: density estimation

If we don't have $p(x | \theta)$ then approximate $\hat{p}(x | \theta)$ it as best as possible - then do standard statistics



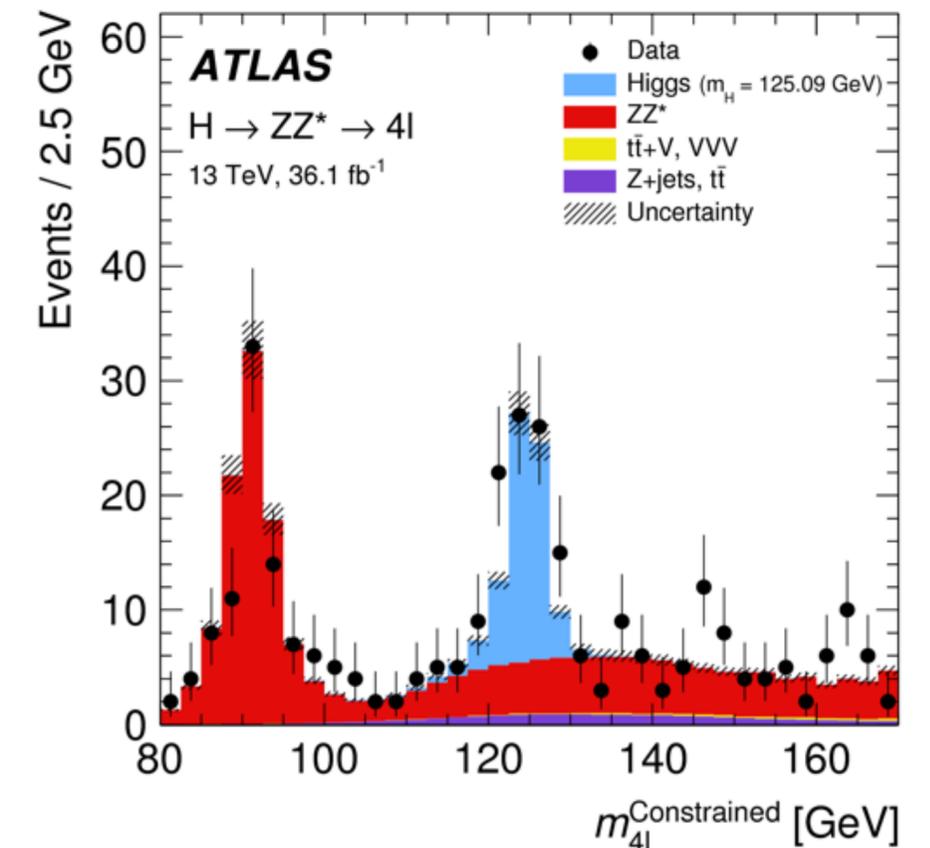
Simulation-based Inference

Simplest method: density estimation

If we don't have $p(x | \theta)$ then approximate $\hat{p}(x | \theta)$ it as best as possible - then do standard statistics

Advantage: no change in methodology

Disadvantage: density estimation in high dimensions is very difficult



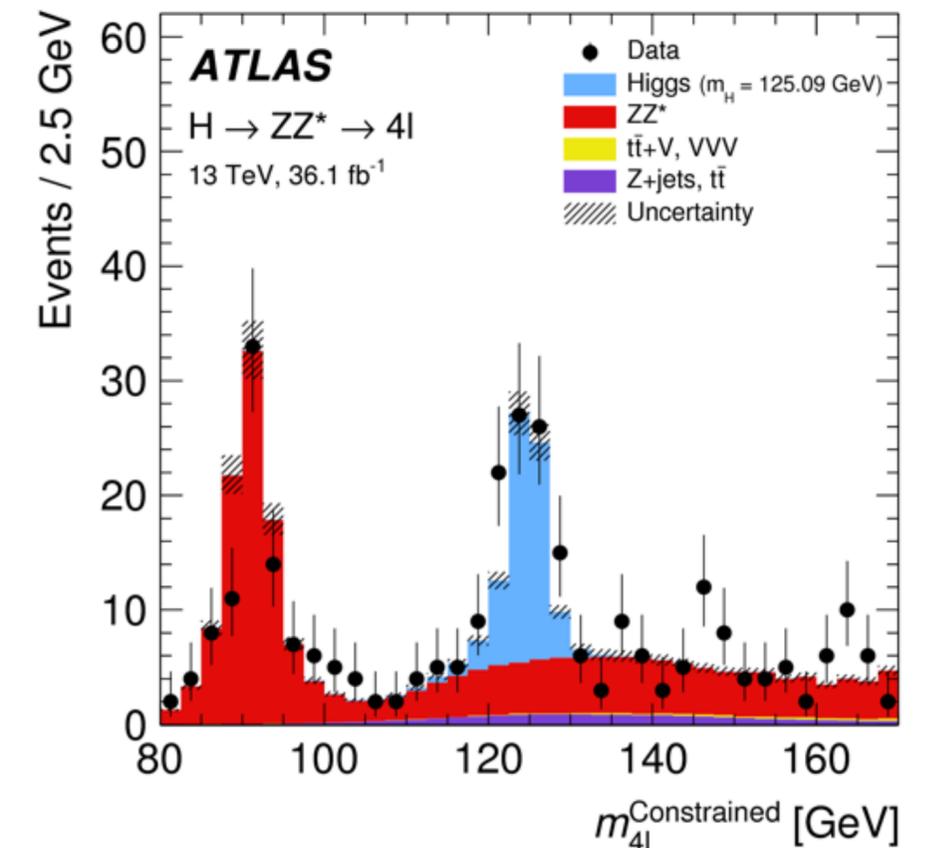
Simulation-based Inference

Simplest method: density estimation

Solution for high-dimensions: find good **summary statistics**

powerful functions $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m \quad m \ll n$

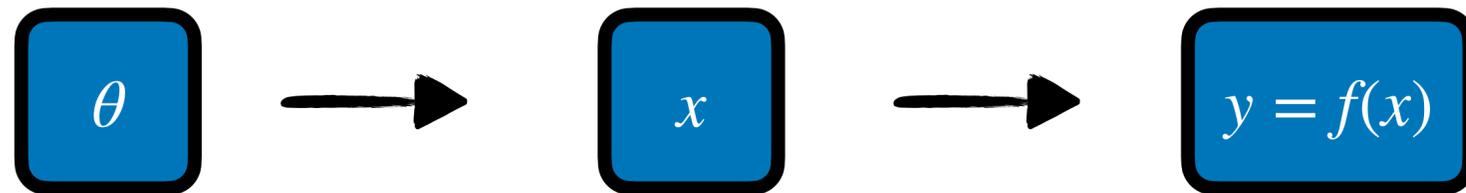
That retain the relevant information and “filter out” the noise as best as possible



Data Processing Inequality

This approach relies crucially on the quality of $f(x)$]

**At best it is as good as the real data “sufficient statistics”,
but usually a lossy process**

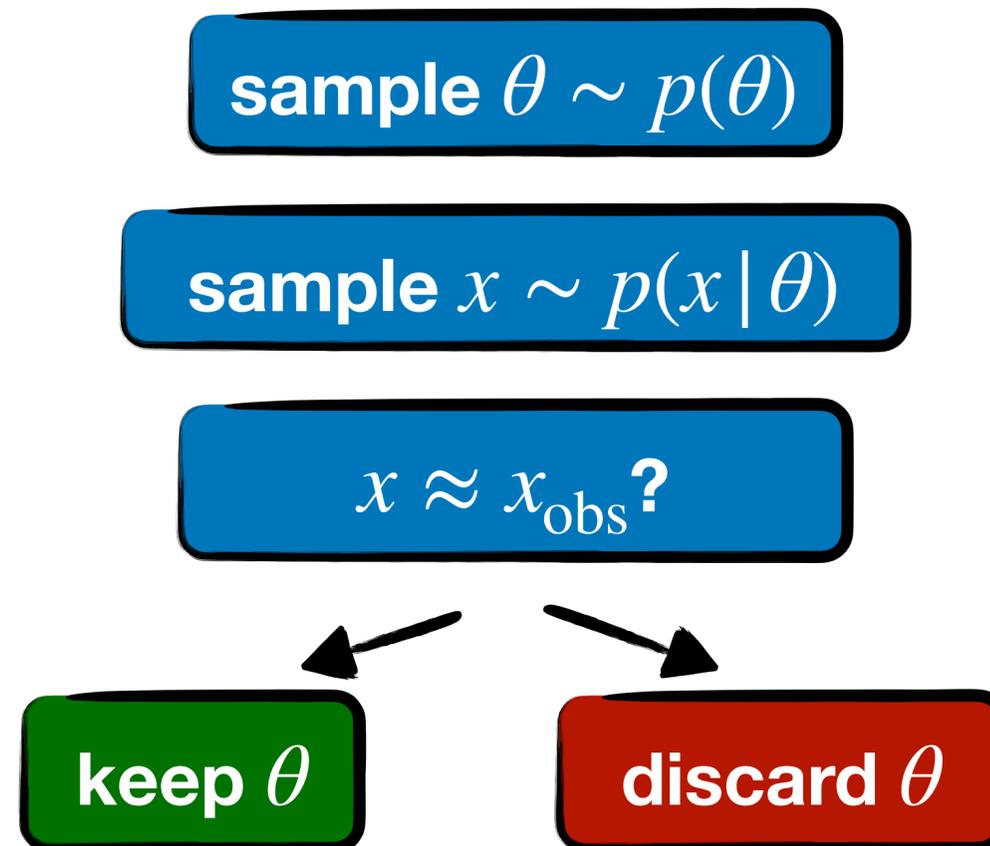


$$I(x; \theta) \geq I(y; \theta)$$

Brute-Force SBI

Alternatively, you can brute force e.g. posterior estimation with samples using **Approximate Bayesian Computation**

Simple Recipe to sample $\theta \sim p(\theta | x)$



Brute-Force SBI

Alternatively, you can brute force e.g. posterior estimation with samples using **Approximate Bayesian Computation**

Simple Recipe to sample $\theta \sim p(\theta | x)$

$$p_\epsilon(\theta | x_o) = \int \delta(|x - x_o| < \epsilon) p(x | \theta) p(\theta)$$

Advantage: no density estimation

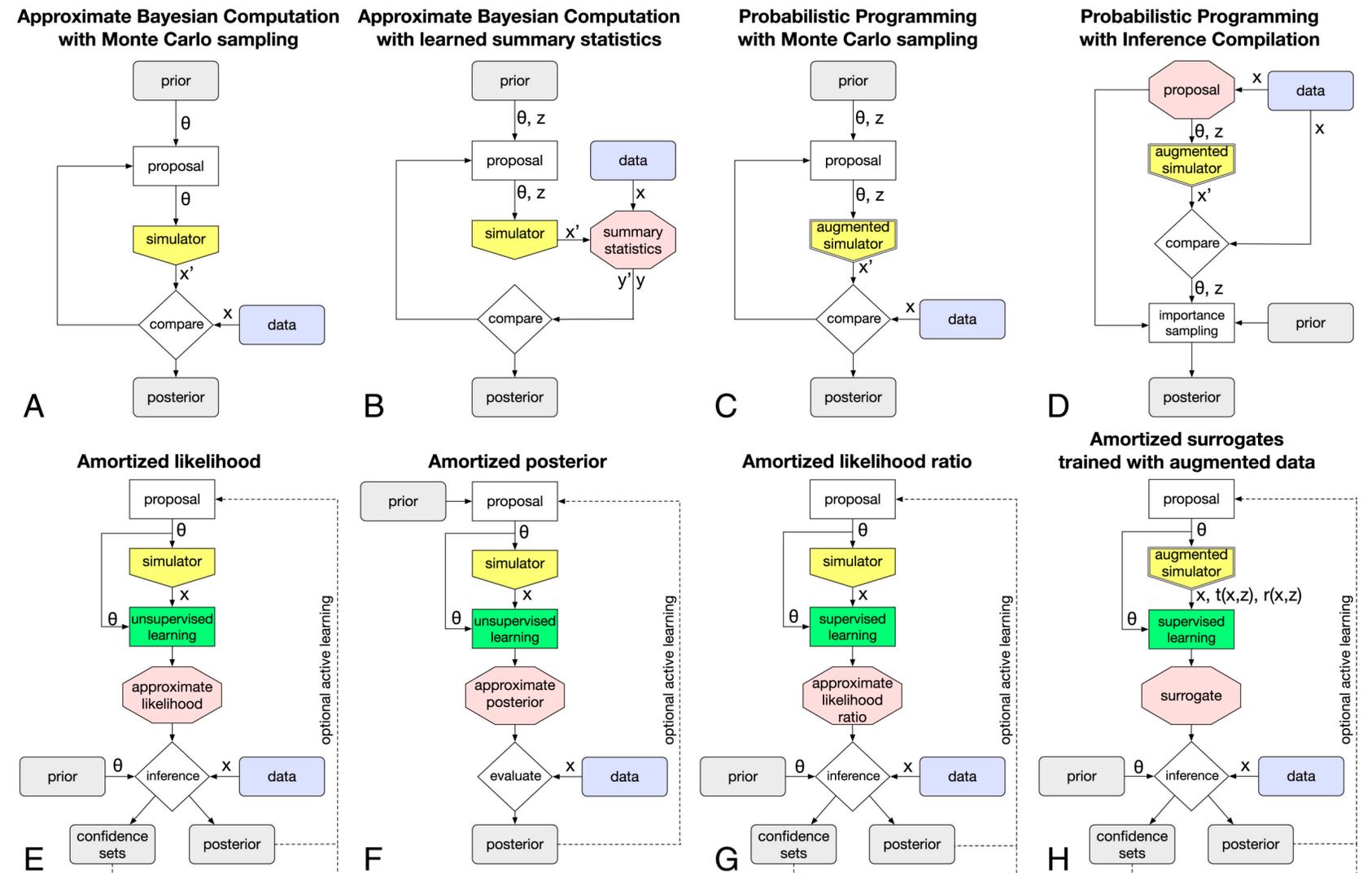
Disadvantage: very inefficient or very approximate

Simulation-based Inference with ML

More Recently there has been a big push to utilize Machine Learning for Simulation-based Inference

An entire family of methods forming, that combine **domain-simulators** with ML for higher fidelity, efficiency, etc..

Can only scratch the surface

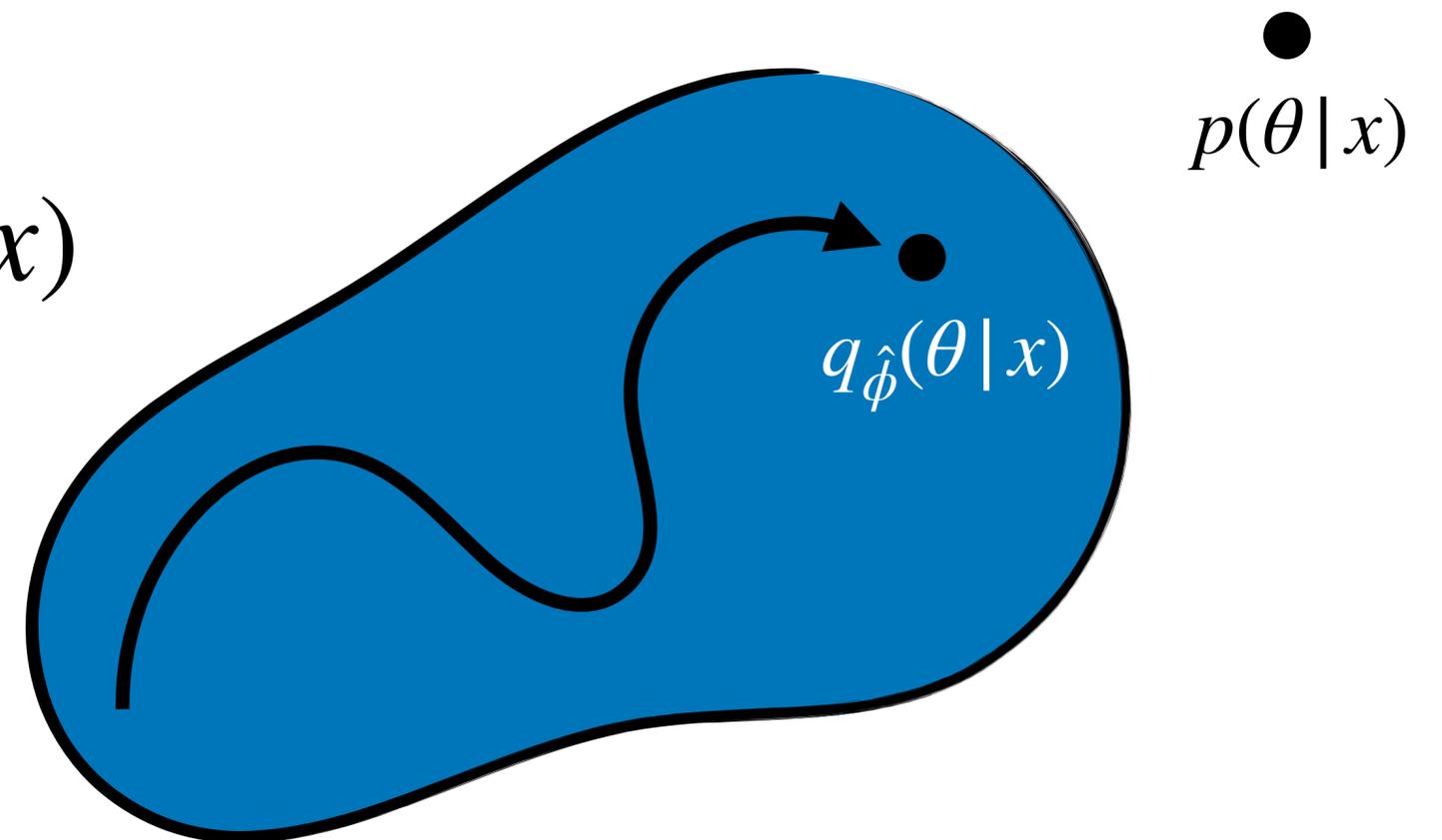


Simplest Approach: Variational Inference

Learning the posterior $p(\theta | x)$ is exactly what supervised machine learning does.

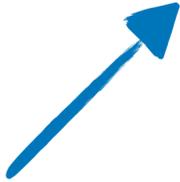
It does to by proposing a family of candidate posteriors, parametrized by ϕ : $q_{\phi}(\theta | x)$

Goal: Optimize ϕ such that $q_{\phi}(\theta | x)$ is a good approx. for all x



Simplest Approach: Variational Inference

Want to minimize

$$L(x) = D_{\text{KL}}(p(\theta | x) || q_{\phi}(\theta | x)) = \mathbb{E}_{\theta|x} \log \frac{p(\theta | x)}{q_{\phi}(\theta | x)}$$


Nice try: no way to sample from $p(\theta | x)$

Power of Amortization

Amazingly if we try to simultaneously solve for all x the problem becomes tractable

$$L' = \mathbb{E}_x L(x) = \mathbb{E}_x \mathbb{E}_{\theta|x} [-\log q_\phi(\theta | x)] = \mathbb{E}_{x,\theta} [-\log q(\theta | x)]$$

just pairs of samples of data and label (θ, x)



This is exactly the “Cross-Entropy” loss in ML!

i.e. ML will converge to the true $p(\theta | x)$ in the limit of infinite data and this purely based on samples!

Limitations

The main limitation of “standard ML” is that the families of candidate posteriors $q_\phi(\theta | x)$ is usually very limited

- Classification: $q_\phi(\theta | x) = \text{Bernoulli}(\theta | p = f_\phi(x))$
- Regression: $q_\phi(\theta | x) = \text{Normal}(\theta | \mu = f_\phi(x), \sigma = \sigma_0)$

True posteriors are likely much more complicated

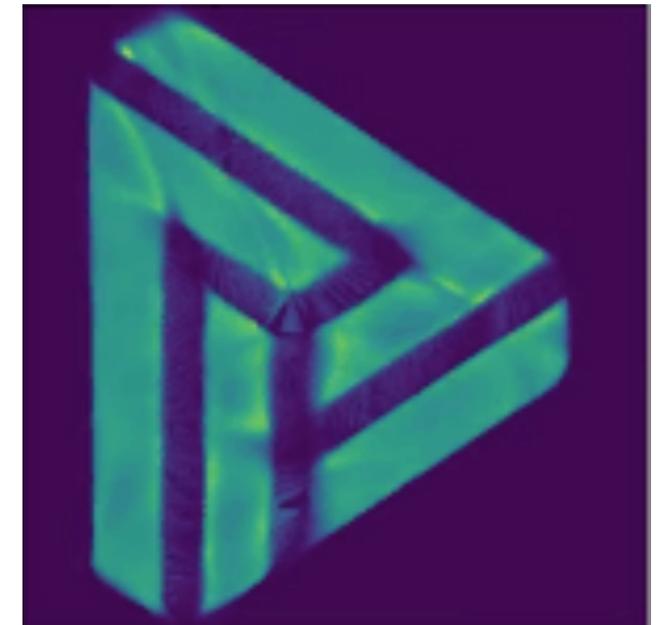
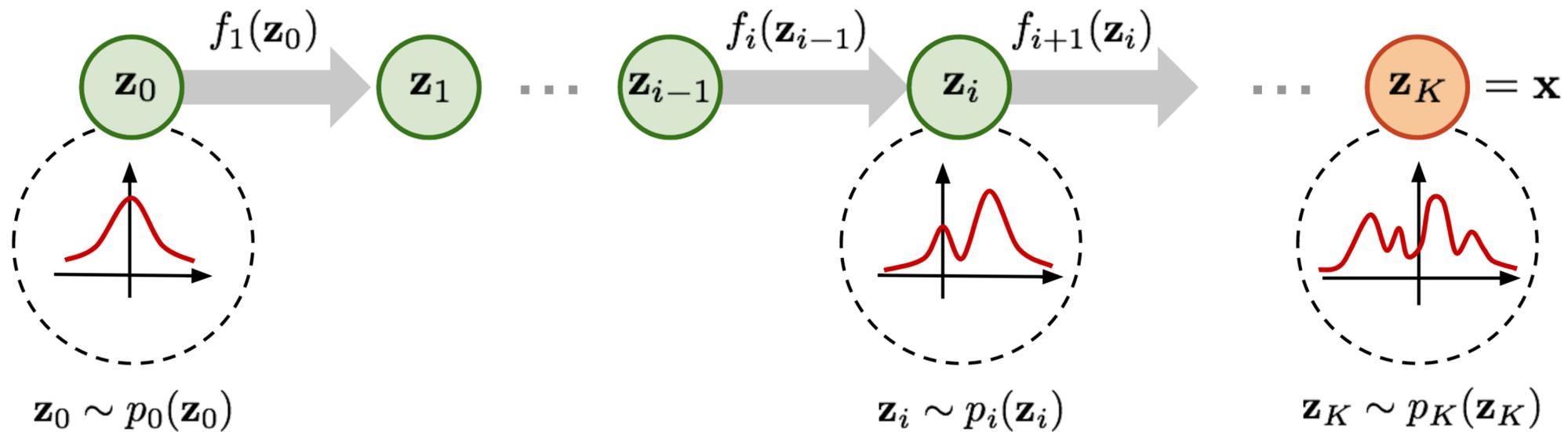
Complexifying $q_{\phi}(\theta | x)$

We can go beyond simple “scaffolding” densities (NNs used to just predict parameters of well known density families) using normalizing flows.

Neural Posterior Estimation

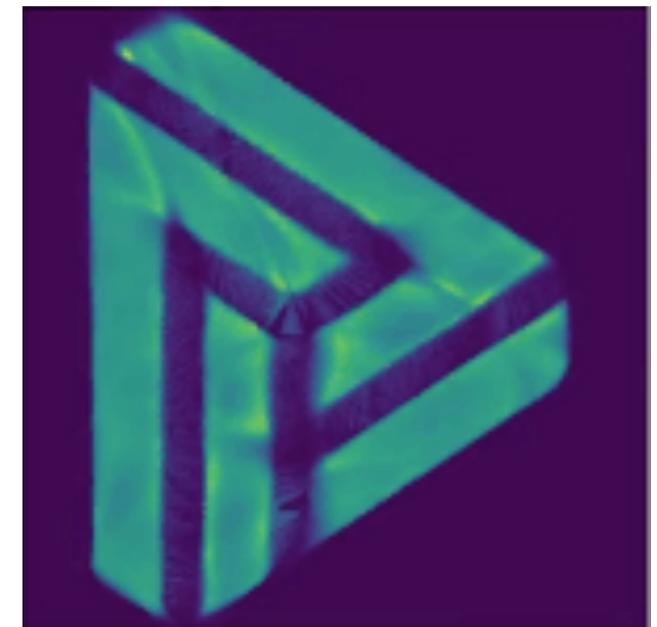
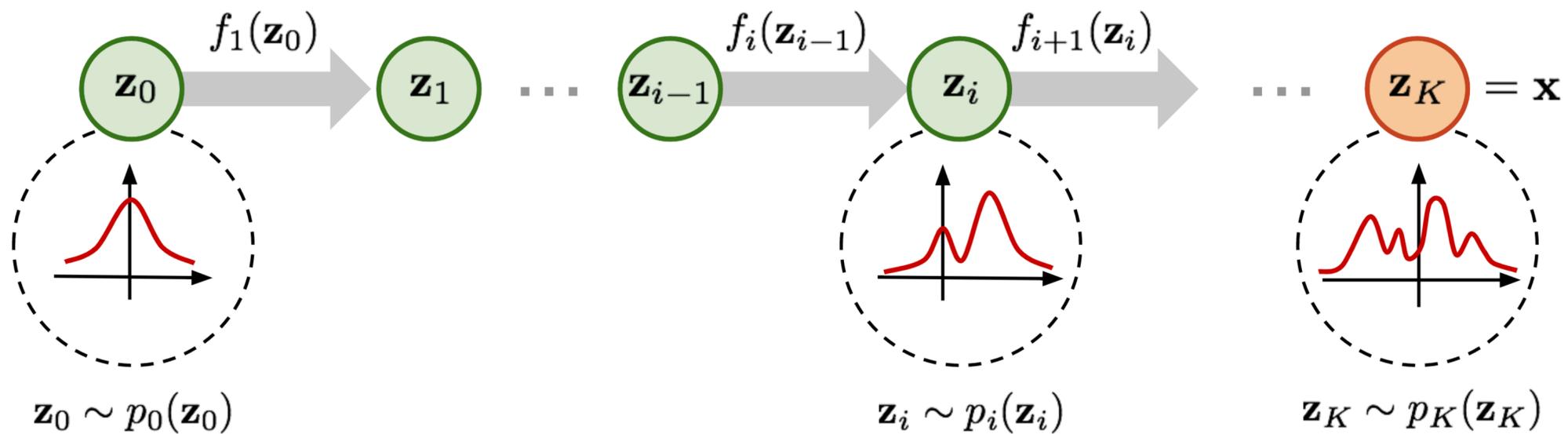
Normalizing Flows

Normalizing Flows are an expressive density modelling method based on a sequence of learnable bijections



Normalizing Flows

Training standard supervised-learning but with normalizing flows gives an immediate recipe for “simulation-based inference”



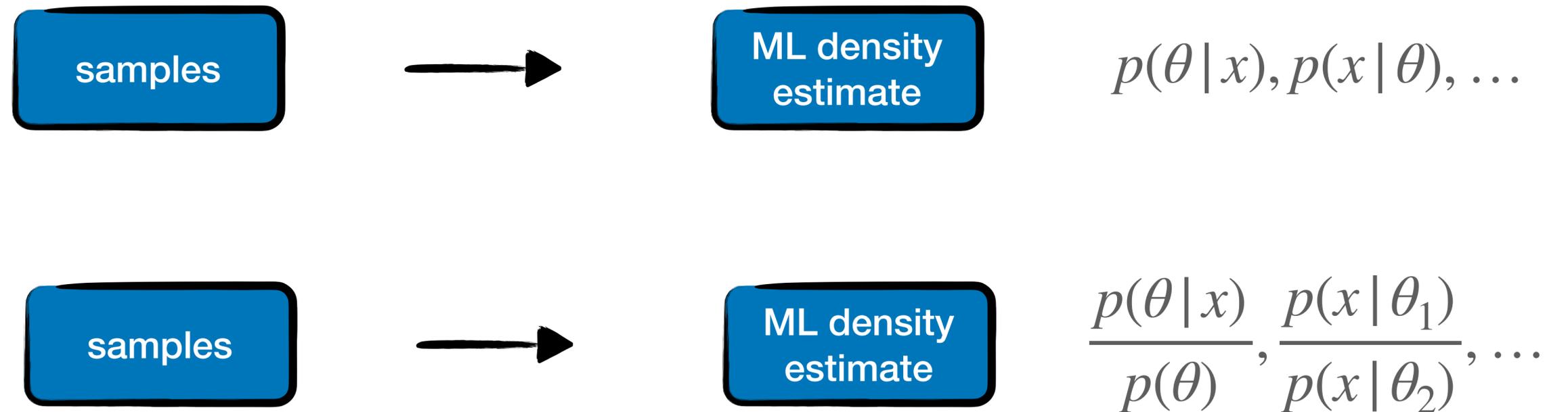
Disadvantage: density estimation is still a hard problem

Density Ratio Estimation

An alternative to density estimation is ...

Density Ratio Estimation

Insight: a lot of standard statistics doesn't really need a fully normalized density. Un- (or arbitrarily) normalized is fine.



Going back to $f(x)$

We actually know what the best possible summary statistic is for e.g. deciding between two hypotheses

**It's the likelihood ratio, due to the
Neyman-Pearson Lemma**

$$t(x) = \frac{p(x | \theta_0)}{p(x | \theta_1)}$$

Can we find it without knowing $p(x | \theta)$?

Yes we can!

Let's consider binary classification between θ_0 and θ_1

We know: BXE minimization leads to a the true posterior (in the asymptotic limit)

$$p(\theta_1 | x) = \frac{p(x | \theta_1)}{p(x)} p(\theta_1)$$

Yes we can!

Let's consider binary classification between θ_0 and θ_1

We know: BXE minimization leads to a the true posterior (in the asymptotic limit)

$$p(\theta_1 | x) = \frac{p(x | \theta_1)}{p(x | \theta_1)p(\theta_1) + p(x | \theta_0)p(\theta_0)} p(\theta_1)$$

Yes we can!

Let's consider binary classification between θ_0 and θ_1

We know: BXE minimization leads to a the true posterior (in the asymptotic limit)

$$p(\theta_1 | x) = \frac{1}{1 + \frac{p(x | \theta_0)}{p(x | \theta_1)}}$$

Yes we can!

Let's consider binary classification between θ_0 and θ_1

We know: BXE minimization leads to a the true posterior (in the asymptotic limit)

$$p(\theta_1 | x) = \frac{1}{1 + \frac{p(x | \theta_0)}{p(x | \theta_1)}}$$

“The Likelihood Ratio Trick”

binary classification \leftrightarrow learn l'hood ratio

 1/Likelihood Ratio

Yes we can!

Let's consider binary classification between θ_0 and θ_1

We know: BXE minimization leads to a the true posterior (in the asymptotic limit)

$$p(\theta_1 | x) = \frac{1}{1 + \exp(-\log(LR))}$$

Yes we can!

Let's consider binary classification between θ_0 and θ_1

We know: BXE minimization leads to a the true posterior (in the asymptotic limit)

$$p(\theta_1 | x) = \text{sigmoid}(\log \text{LR})$$

The pre-activation in the last NN layer is actually an estimate of the log-likelihood ratio!

Not very surprising

We actually know what the best possible summary statistic is for e.g. deciding between two hypotheses

**It's the likelihood ratio, due to the
Neyman-Pearson Lemma**

$$t(x) = \frac{p(x | \theta_0)}{p(x | \theta_1)}$$

**So training a ML algorithm to decide,
will learn the same thing**

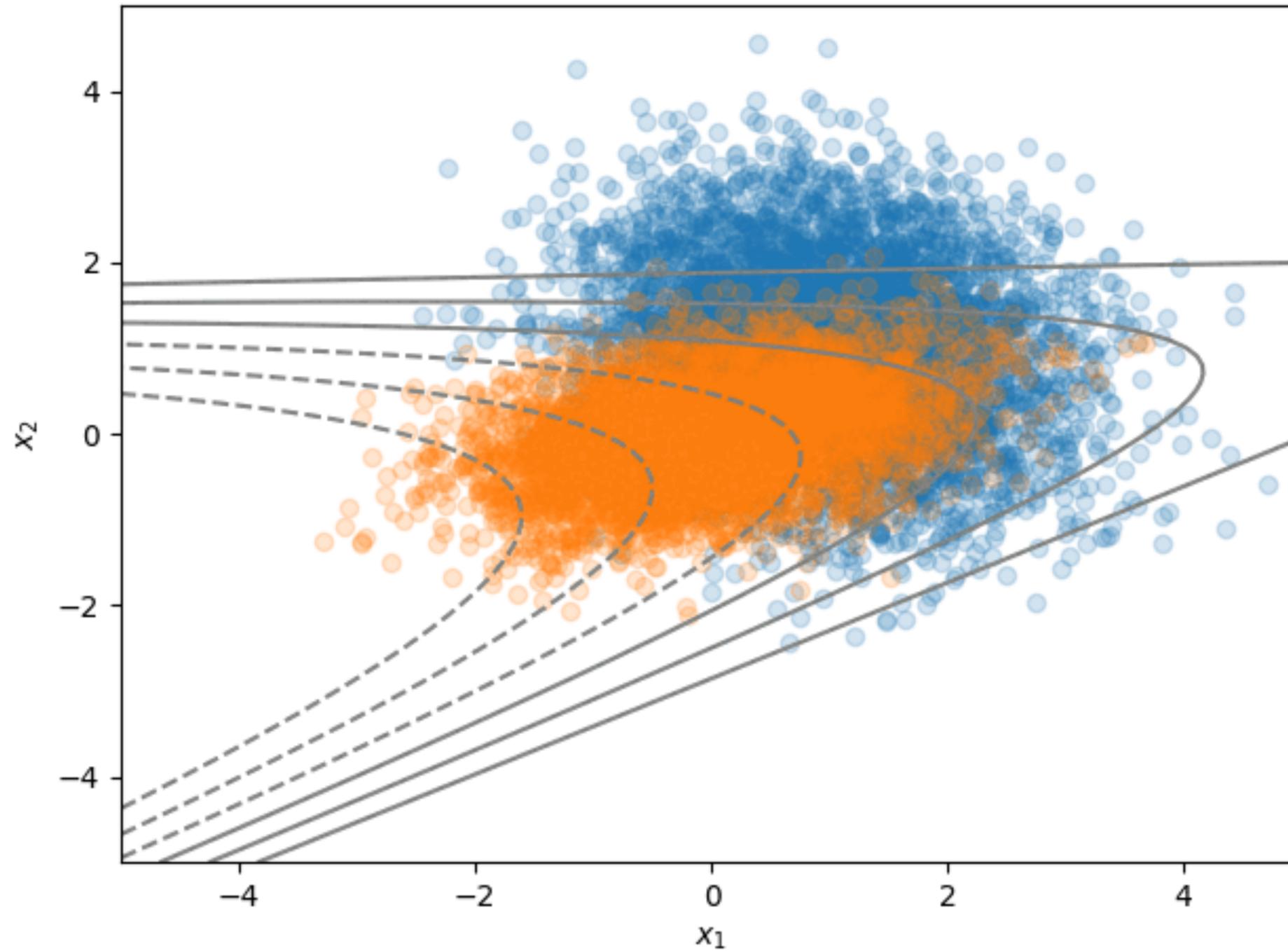
Recovering LR from a Classifier

We can recover the Likelihood Ratio from the ML model by just rearranging the formula

$$y = q_{\phi}(\theta_1 | x) = f_{\phi}(x) \approx \frac{1}{1 + \frac{1}{LR}}$$

$$LR \approx \frac{y}{1 - y}$$

Likelihood Ratio Trick - Visualized



The Value of Density Ratios

With the Likelihood-Ratio trick in our bag, we have a nice way to learn the likelihood ratio for a high-dimensional parameter space

We now want $A(x, \theta) = \frac{p(x | \theta)}{p(x | \theta_0)}$ instead of $A(x) = \frac{p(x | \theta_1)}{p(x | \theta_0)}$

For most purposes $A(x, \theta)$ is just as useful as $p(x | \theta)$

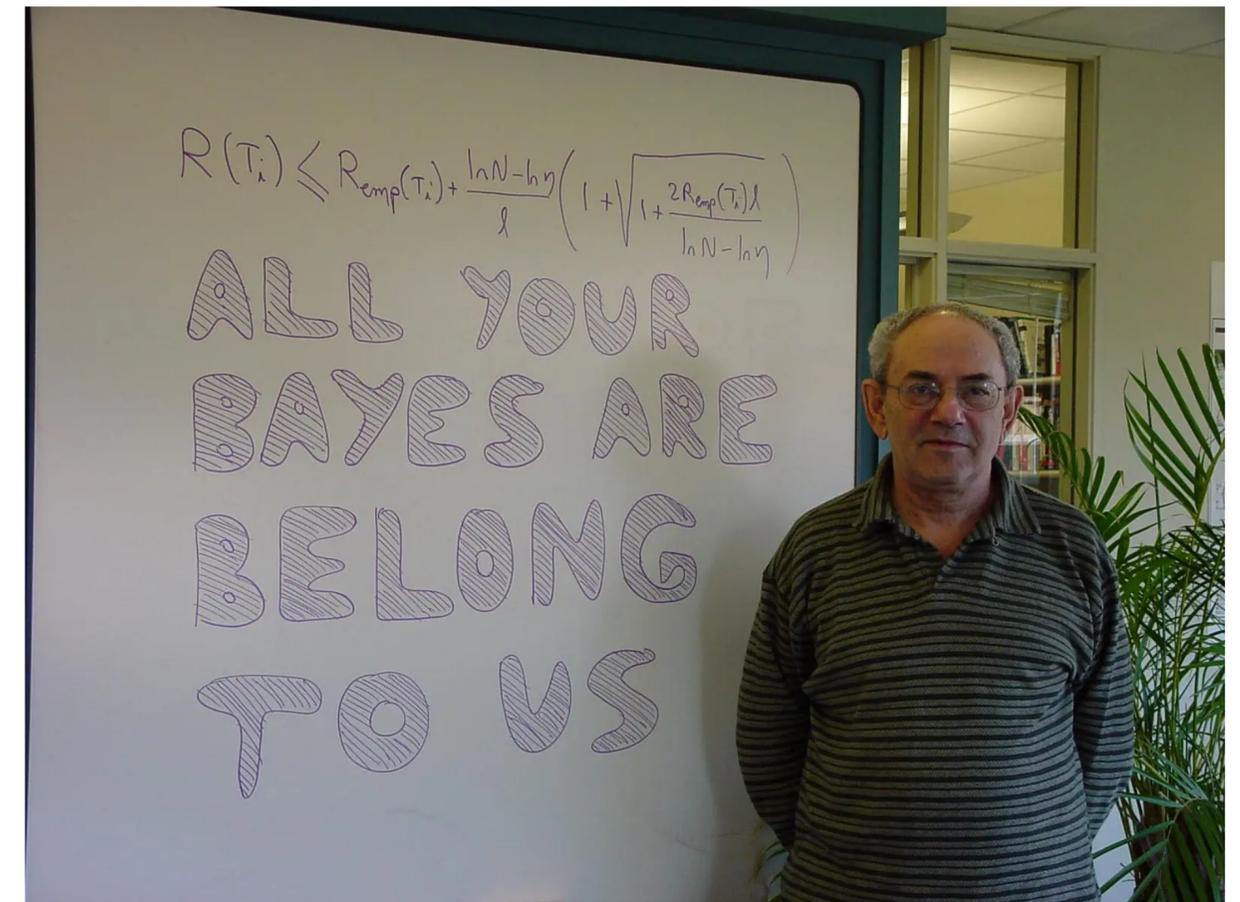
- **MLE estimation, MCMC, HMC, you name it...**

The Value of Density Ratios

Insight:

"When solving a problem of interest, do not solve a more general problem as an intermediate step" -Vladimir Vapnik

(with density estimation we did try to solve a more general problem)



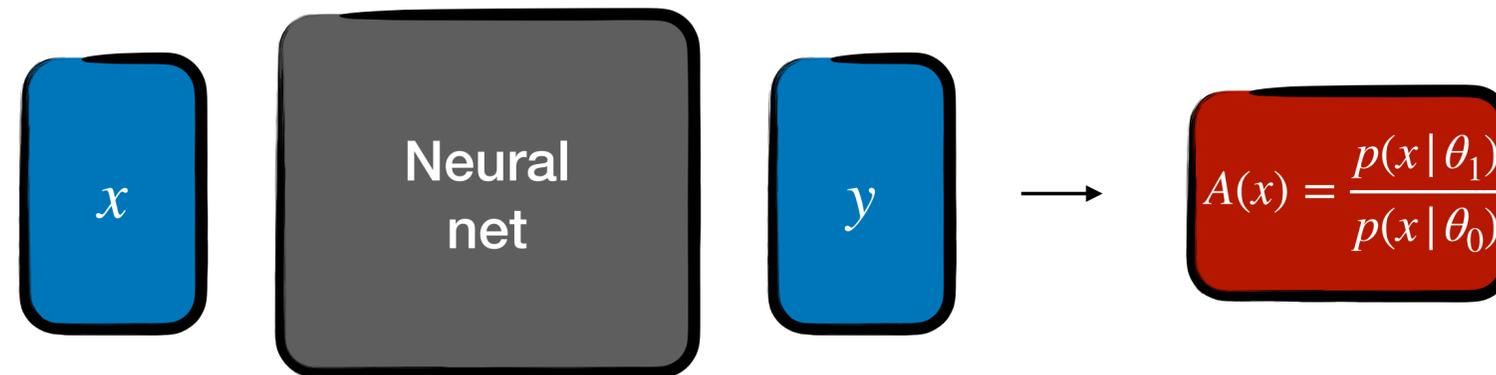
Parametrized Neural Network

We can use our classification-as-LR estimation trick to learn our target function $A(x, \theta)$ by

- **conditioning (or parametrizing) out neural network**
- **expanding our classification training regime**

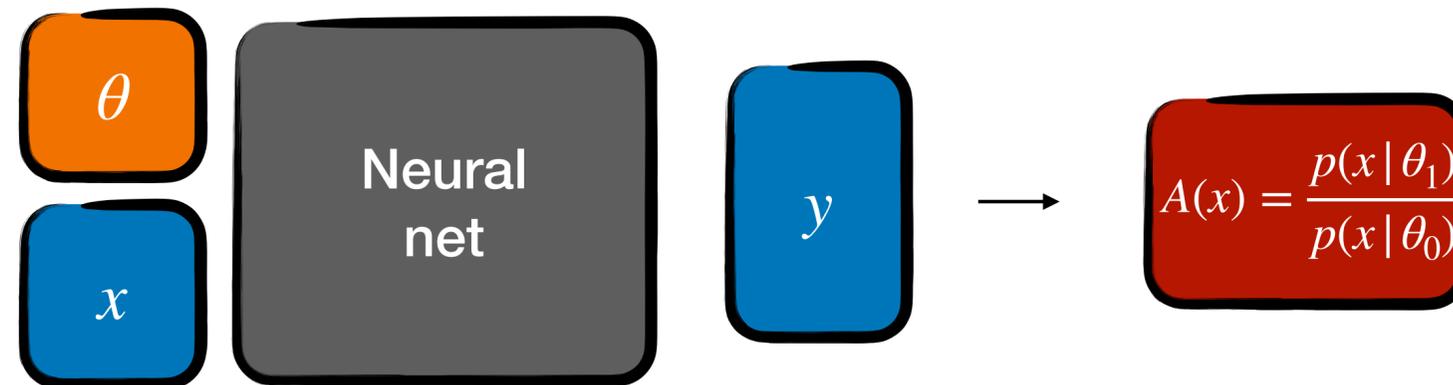
Conditioning the Neural Net

Normally the neural network just gets the data it's supposed to classify



Conditioning the Neural Net

But now we extend it by also giving it an additional input θ to turn $A(x) \rightarrow A(x, \theta)$



The input θ is signaling to the NN, which classification problem we are solving at any given time

Parametrized Training

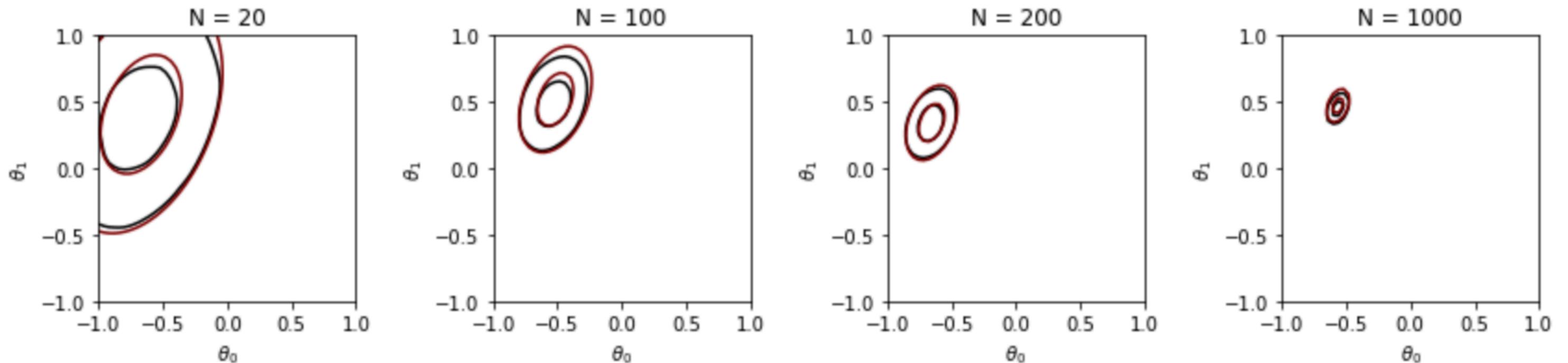
We can now have the following training loop

- sample a $\theta \sim p(\theta)$
- set up a classification problem for θ
 - sample “class 0” from simulator at $\theta_1: x \sim p(x | \theta_1)$
 - sample “class 1” from simulator at $\theta: x \sim p(x | \theta)$
 - do one train step in binary classification by running the network at value θ (will converge to likelihood ratio)
- rinse & repeat

We will do this as an exercise!

Example Results

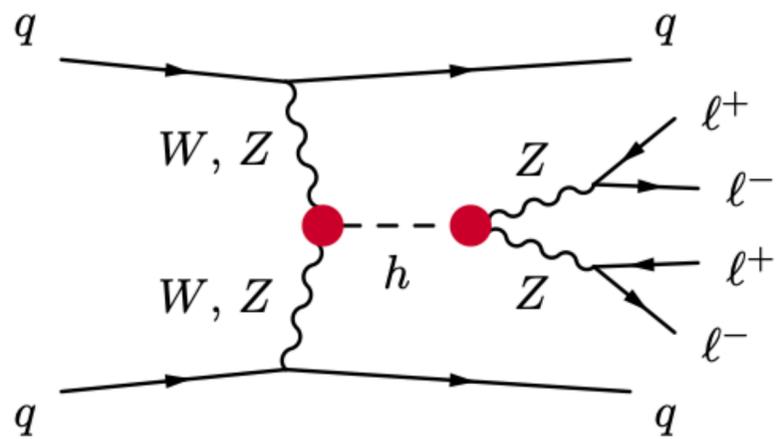
With sufficient data, we can learn $A(x, \theta) = p(x | \theta) / C$ to an impressive precision.



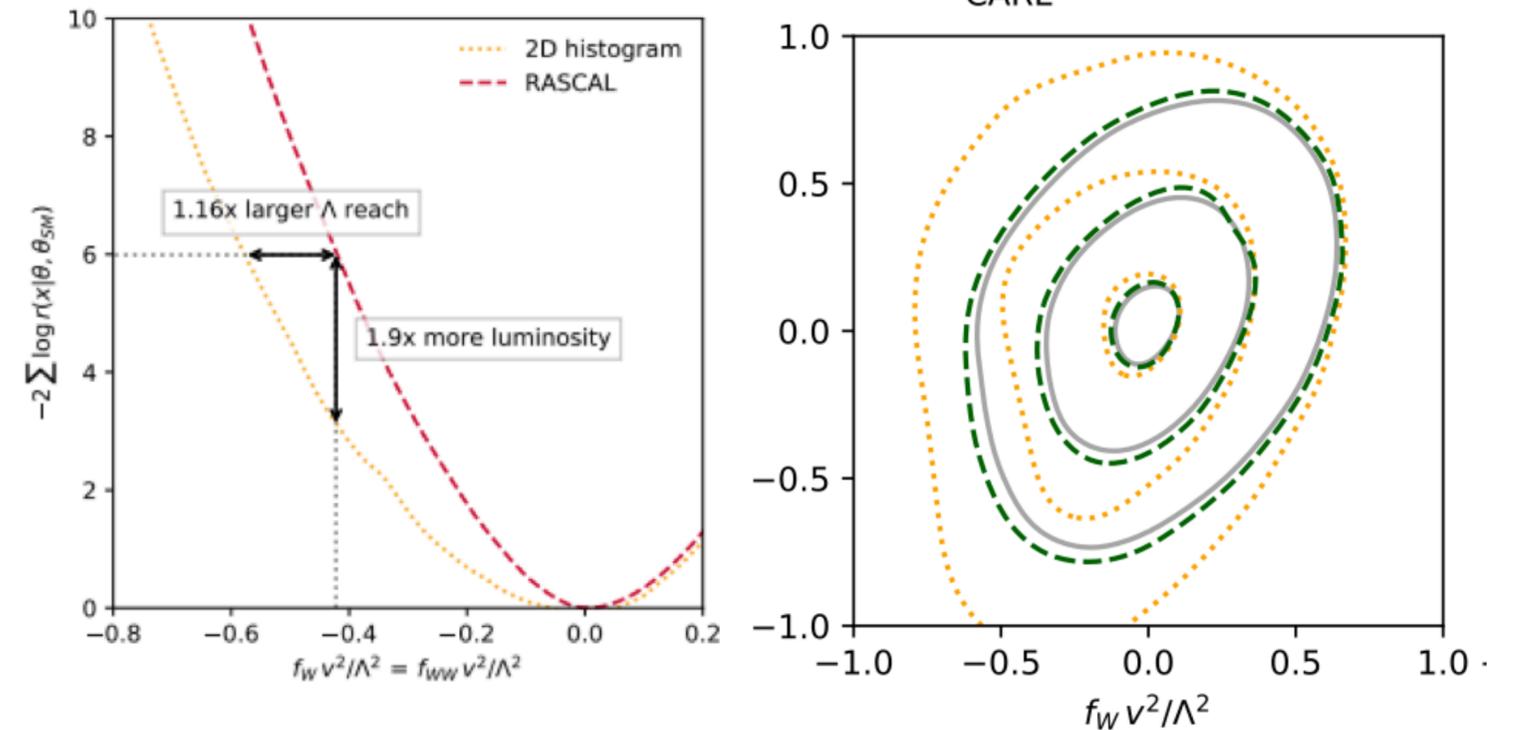
Note: no(!) data processing inequality, we use the full data!

Impact, e.g. at the LHC

The difference between summary statistics (hand-chosen $f(x)$) and “optimal observables” (i.e. LR estimates) can be equivalent to many years of LHC running



Studying the Higgs Boson in detail
at the LHC



Extensions

The procedure to learn density ratios through classification is very general. Can be applied to many different scenarios

E.g.

$$p(\theta | x) = \frac{p(x | \theta)}{p(x)} p(\theta)$$

Extensions

The procedure to learn density ratios through classification is very general. Can be applied to many different scenarios

E.g.

$$p(\theta | x) = \frac{p(x | \theta) p(\theta)}{p(x) p(\theta)} p(\theta)$$

Extensions

The procedure to learn density ratios through classification is very general. Can be applied to many different scenarios

E.g.

Look Ma! A Density Ratio


$$p(\theta | x) = \frac{p(x, \theta)}{p(x)p(\theta)} p(\theta)$$

The key difference between the posterior and the prior is captured by the density ratio of the **joint distribution** vs the **product of the marginals**

Extensions

We can sample from both!

- **joint density** $(x, \theta): \theta \sim p(\theta) \rightarrow x \sim p(x | \theta)$
- **marginal** $(x, \theta): \theta \sim p(\theta) \quad \theta' \sim p(\theta) \rightarrow x \sim p(x | \theta')$

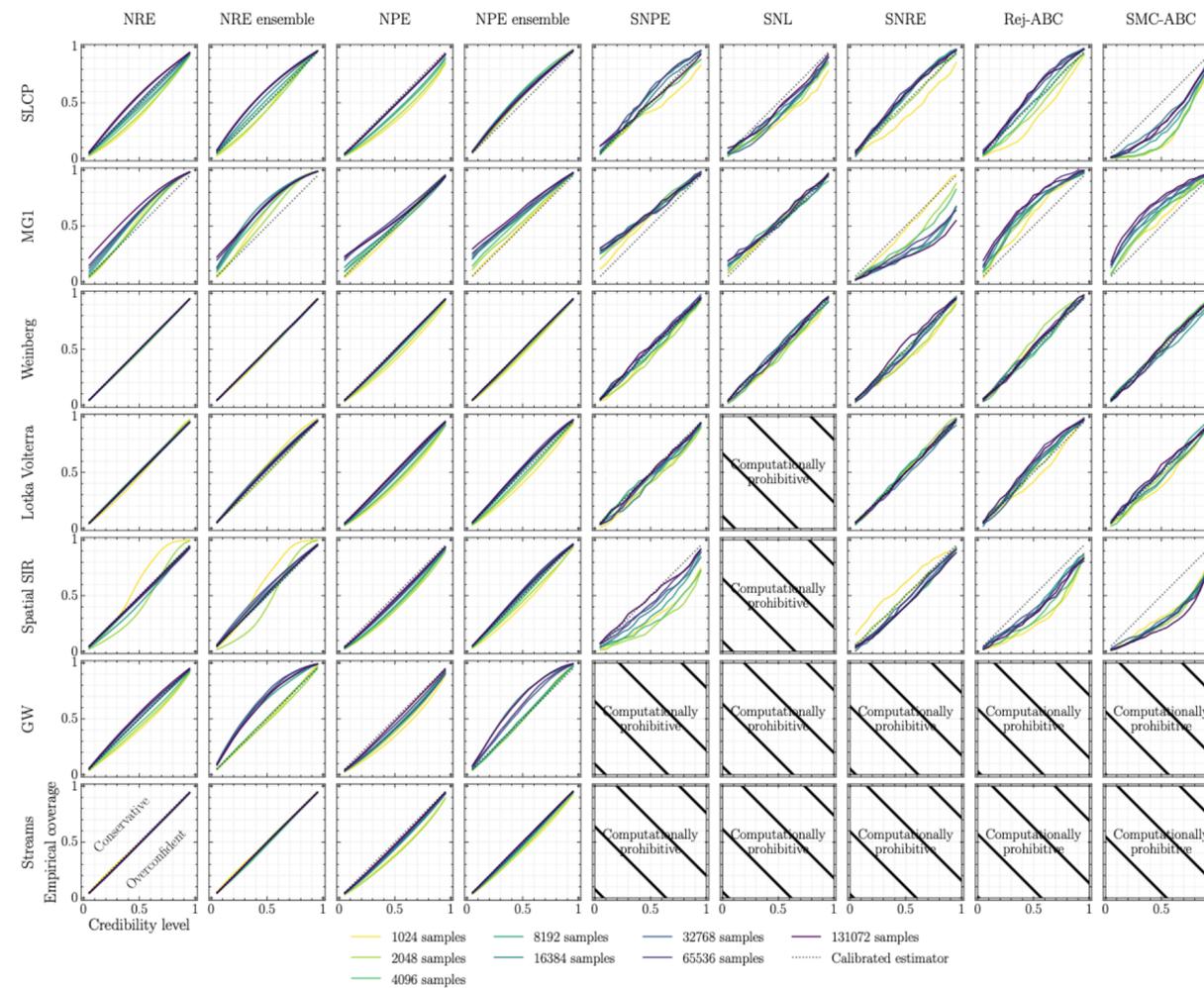
And with this we can setup the same classification regime to learn

$$\frac{p(x, \theta)}{p(x)p(\theta)}$$

Neural Ratio Estimation (NRE) (for Bayes)

Benchmarks

Methods for Simulation-based inference behave well in different regimes. Extensive Benchmarking to guide you.



Software Tools for SBI

sbil : simulation-based inference

sbil : A Python toolbox for simulation-based inference.

```
prior = BoxUniform(low=zeros(2), high=2*ones(2)) # Box prior [0,2]x[0,2]
def simulator(theta): return theta + 0.1*randn_like(theta) # Gaussian in 2D
posterior = infer(simulator, prior, method='SNPE', num_simulations=500)
```

```
Running 500 simulations.: 100% |██████████| 500/500 [00:00<00:00, 57141.55it/s]
Training neural network. Epochs trained: 90
```

```
samples = posterior.sample((1000,), x=observed)
```

Table of contents

Motivation and approach

Publications

Posterior estimation (SNPE)

Likelihood-estimation (SNLE)

Likelihood-ratio-estimation (SNRE)

Utilities

Inference can be run in a single line of code:

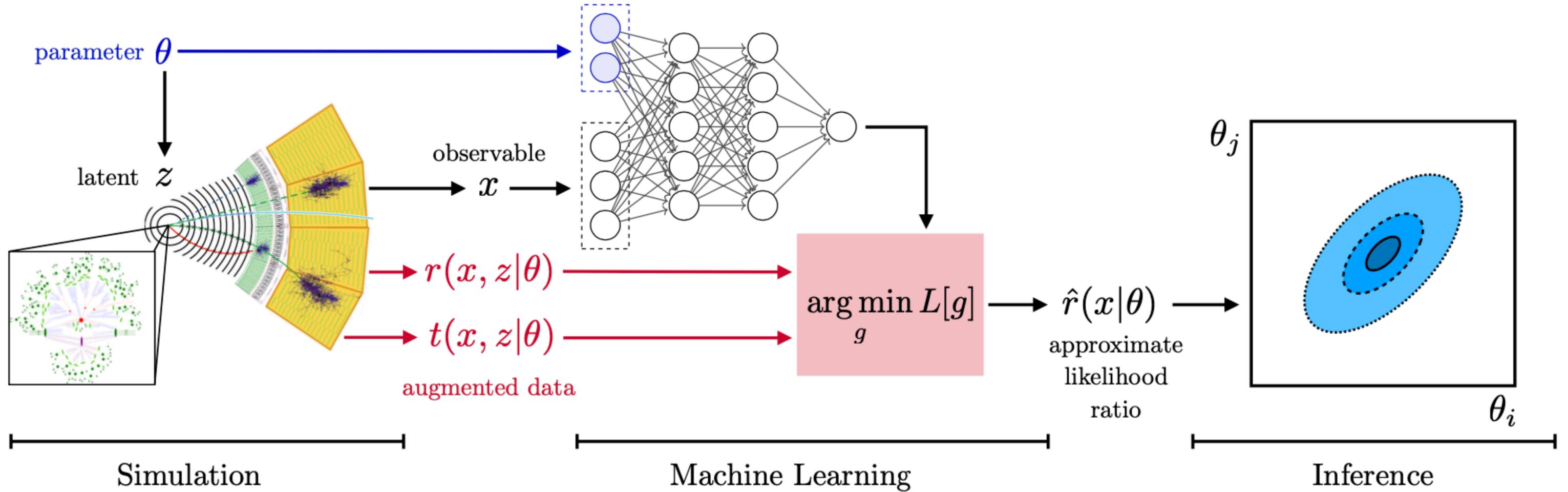
```
posterior = infer(simulator, prior, method='SNPE', num_simulations=1000)
```

- To learn about the general motivation behind simulation-based inference, and the inference methods included in `sbil`, read on below.
- For example applications to canonical problems in neuroscience, browse the recent research article [Training](#)

Display a menu [density estimators to identify mechanistic models of neural dynamics.](#)

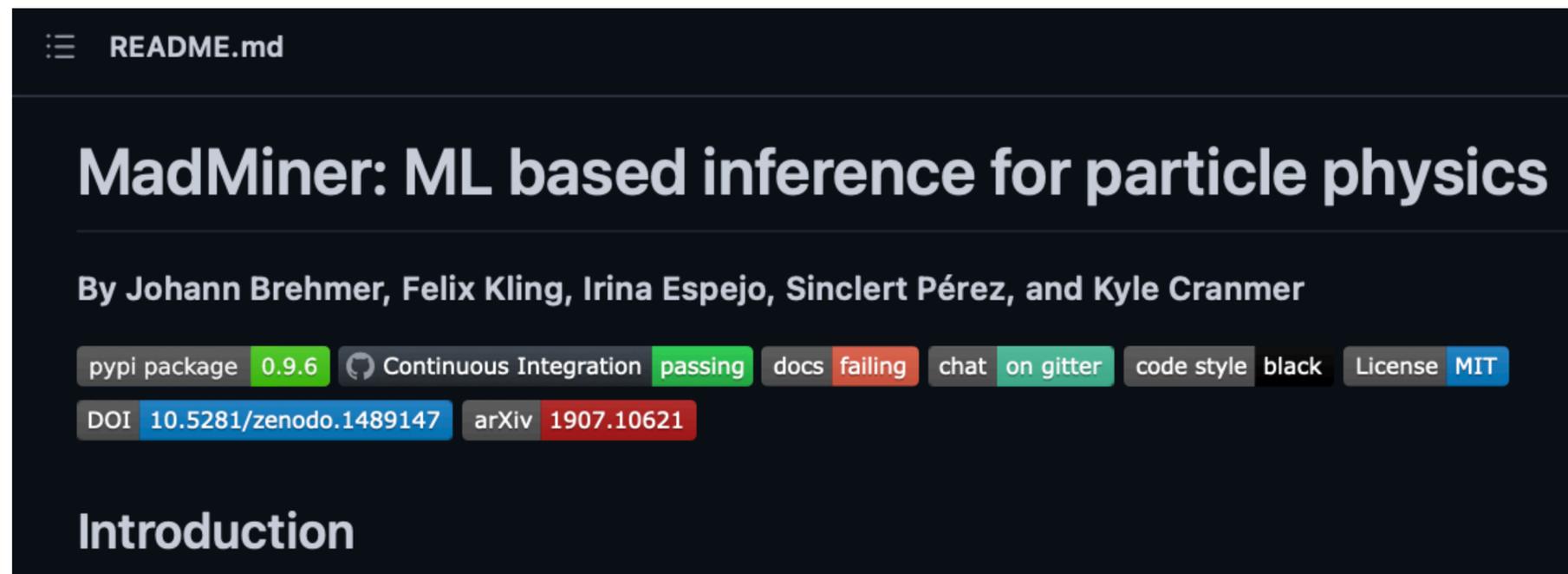
```
prior = BoxUniform(low=zeros(2), high=2*ones(2)) # Box prior [0,2]x[0,2]
def simulator(theta): return theta + 0.1*randn_like(theta) # Gaussian in 2D
```

In Particle Physics



In Particle Physics

MadMiner: a package for end-to-end SBI in particle physics



☰ README.md

MadMiner: ML based inference for particle physics

By Johann Brehmer, Felix Kling, Irina Espejo, Sinclert Pérez, and Kyle Cranmer

pypi package 0.9.6 Continuous Integration passing docs failing chat on gitter code style black License MIT

DOI 10.5281/zenodo.1489147 arXiv 1907.10621

Introduction

<https://madminer.readthedocs.io/en/latest/?badge=latest>

Bonus: Data Efficiency & Mining Gold

Sounds too good to be true?

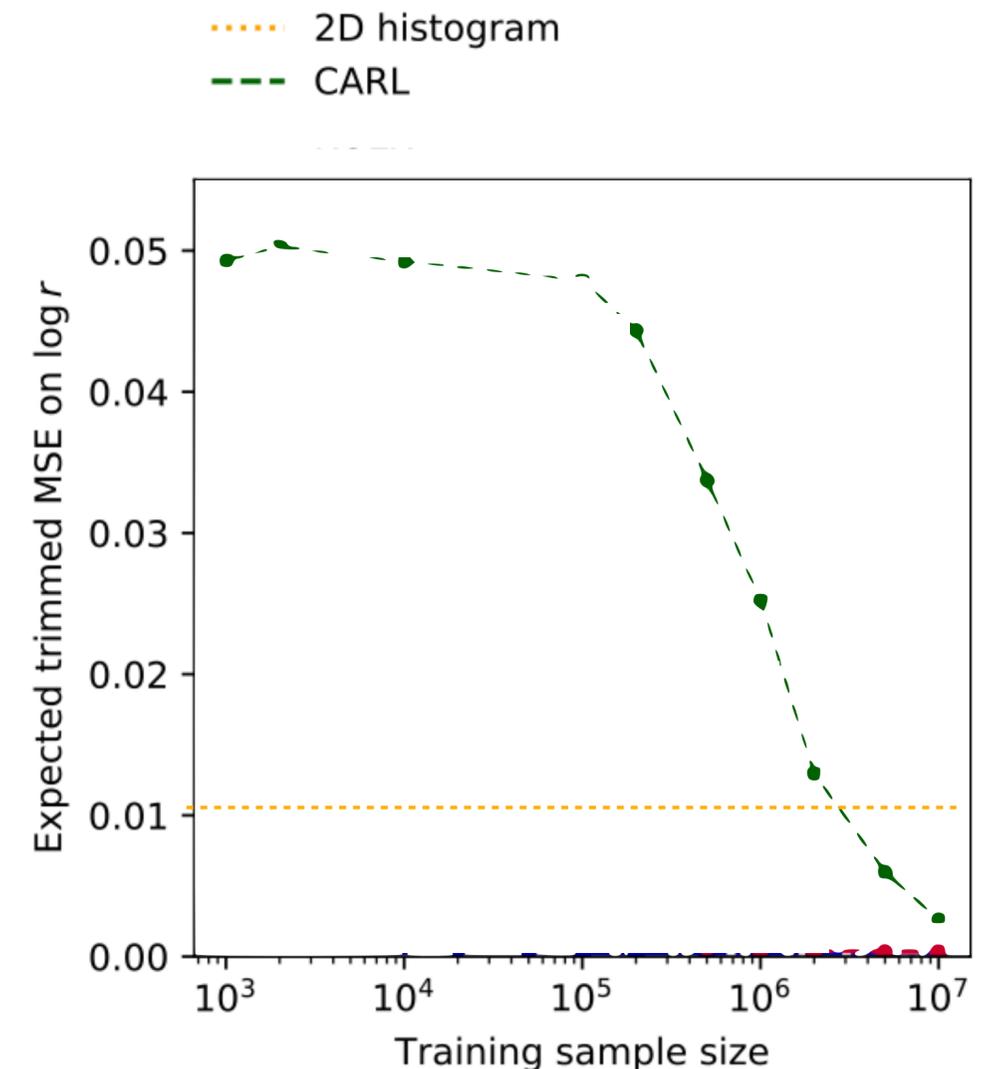
This is nice in theory. But of course there is a catch.

This requires an very large amount of MC to become better than the traditional approach.

Higgs Example:

- 10 Million samples per theory point
- $O(10)$ - $O(100)$ possible theories

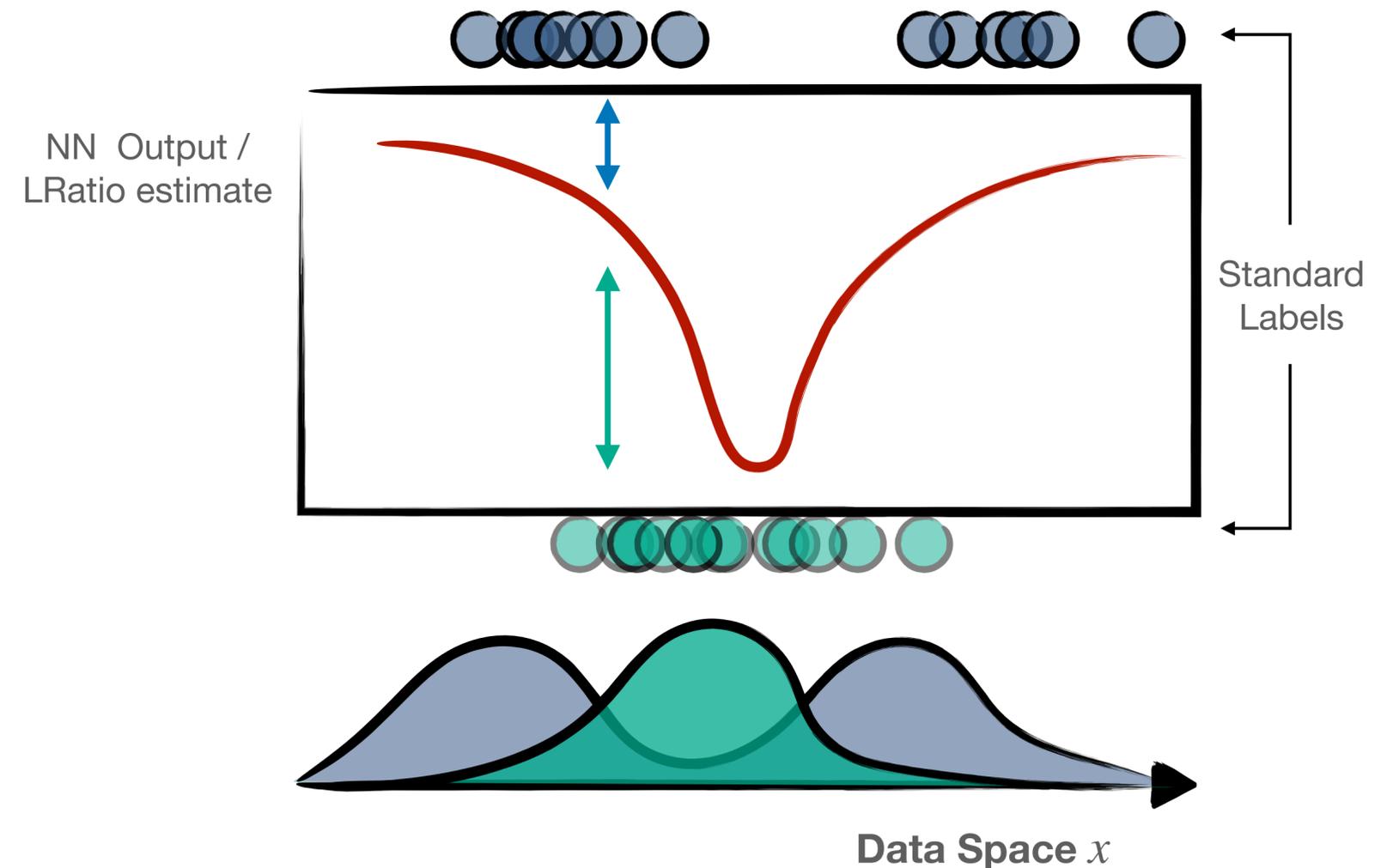
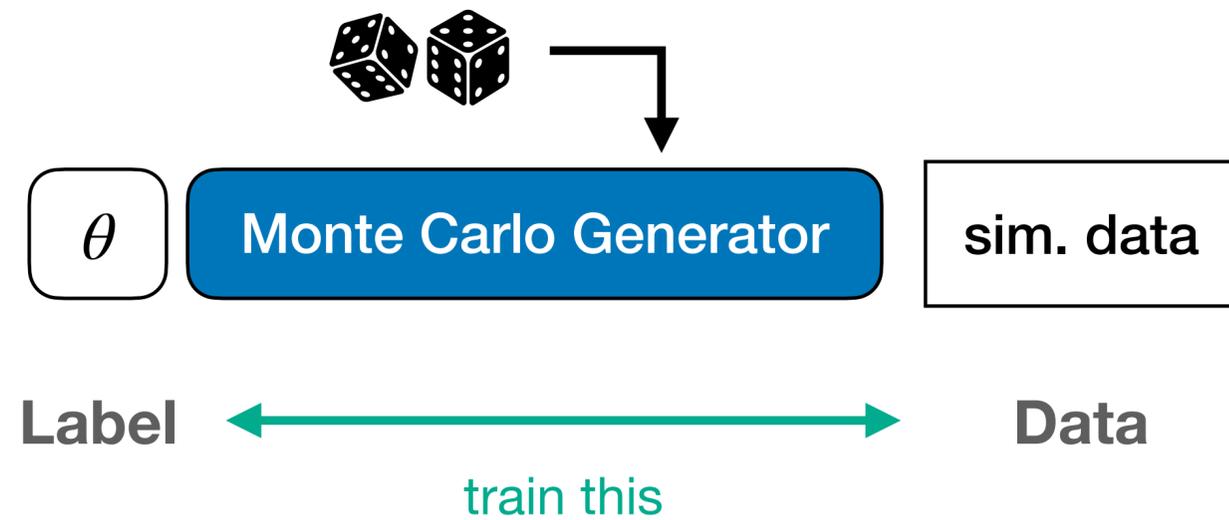
Can be computationally prohibitive



Increasing Data Efficiency

The reason training converges slowly is again distance between theory and data.

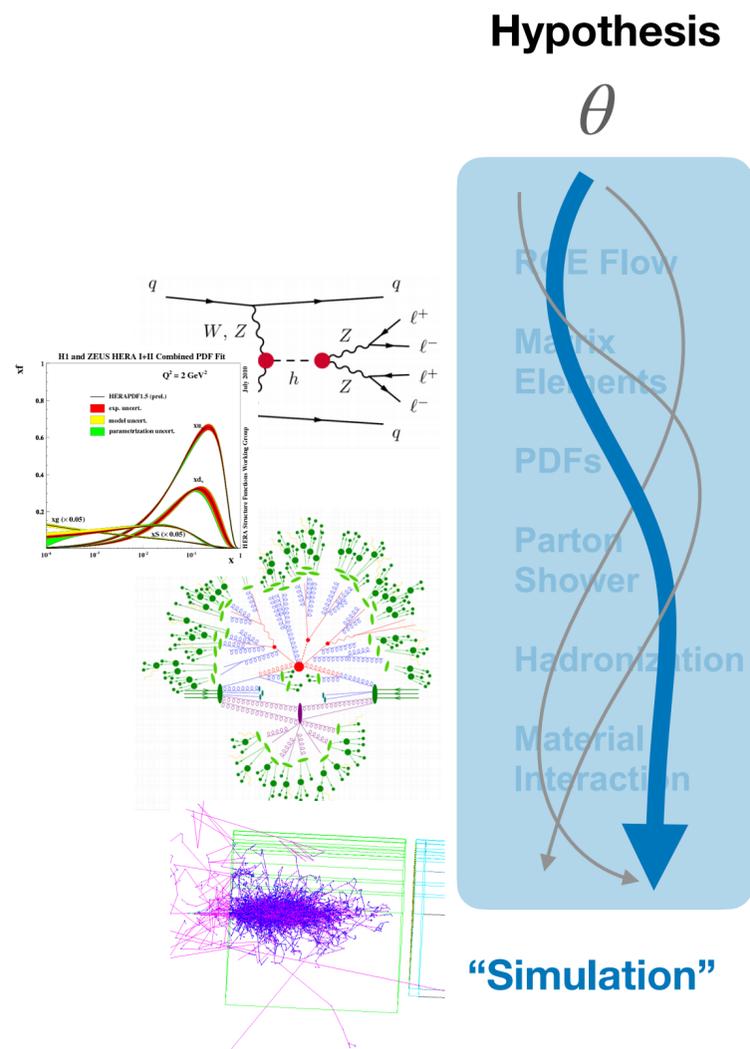
Training = “tug of war” between samples



More tricks

To increase data efficiency we can use more tricks. Remember:

- We can calculate $p(x, z)$ for any given path through our MC.
- it's just the integral that kills us



$$p(x | \theta) = \int dz \frac{p(x | z_h) p(z_h | z_p) p(z_p | \theta)}{p(x, z | \theta)}$$

NOT OK

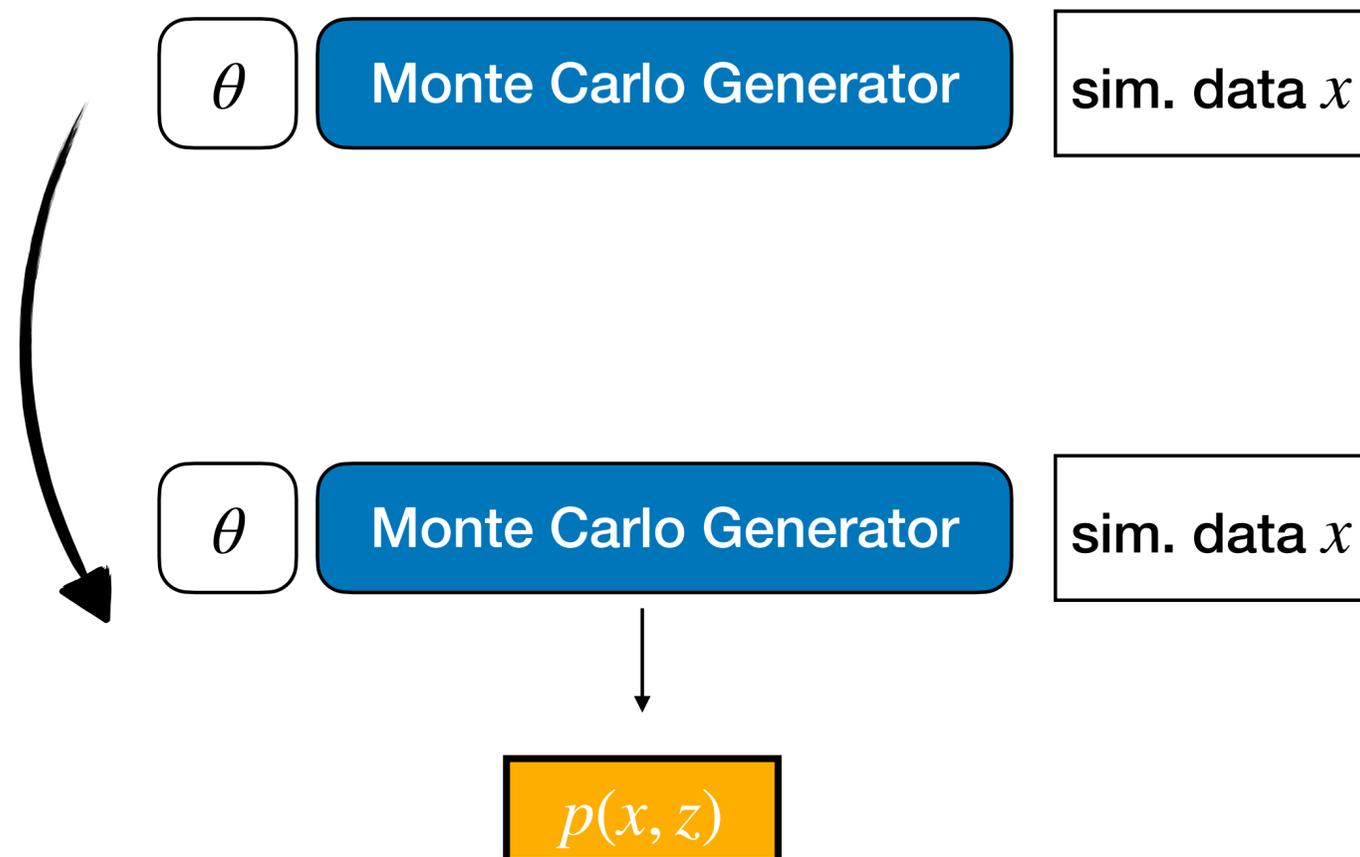
OK

Mining Gold

By adapting our MC to not only spit out the simulated event **but also the “path probability”** $p(x, z | \theta)$ we are extracting valuable information for ML training

→ exploit our physics knowledge in ML training

“Mining Gold”



Interesting Property

Turns out: the true likelihood ratio, and the path-dependent ratio are closely related

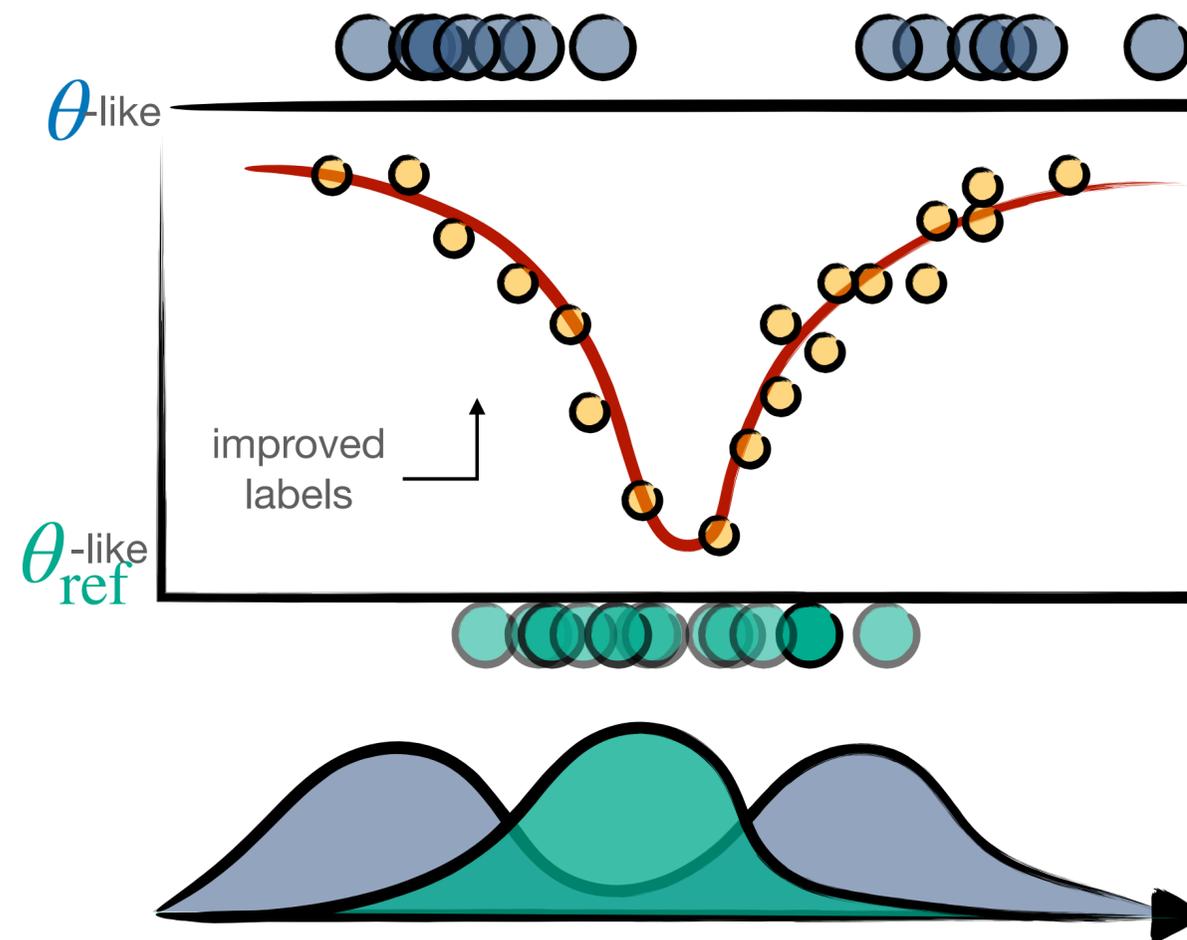
$$\text{NOT OK } \frac{p(x | \theta)}{p(x | \theta_0)} \quad \text{OK } \frac{p(x, z | \theta)}{p(x, z | \theta_0)}$$

But...

$$\frac{p(x | \theta)}{p(x | \theta_0)} = \mathbb{E}_{p(z|x)} \frac{p(x, z | \theta_1)}{p(x, z | \theta_0)}$$

Interesting Property

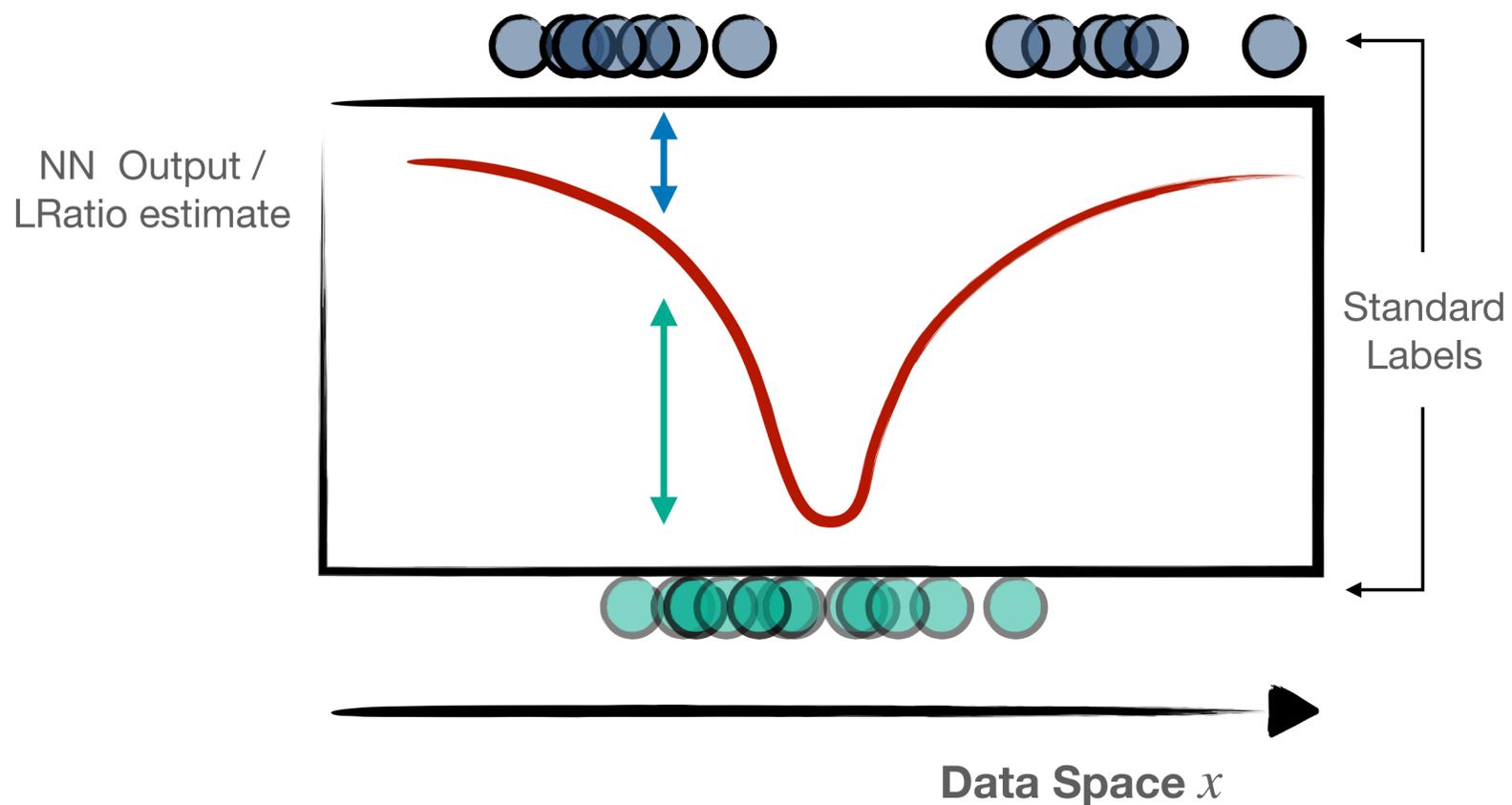
The paths through the simulation chain add “theory noise” but in expectation average out to be the true likelihood ratio!



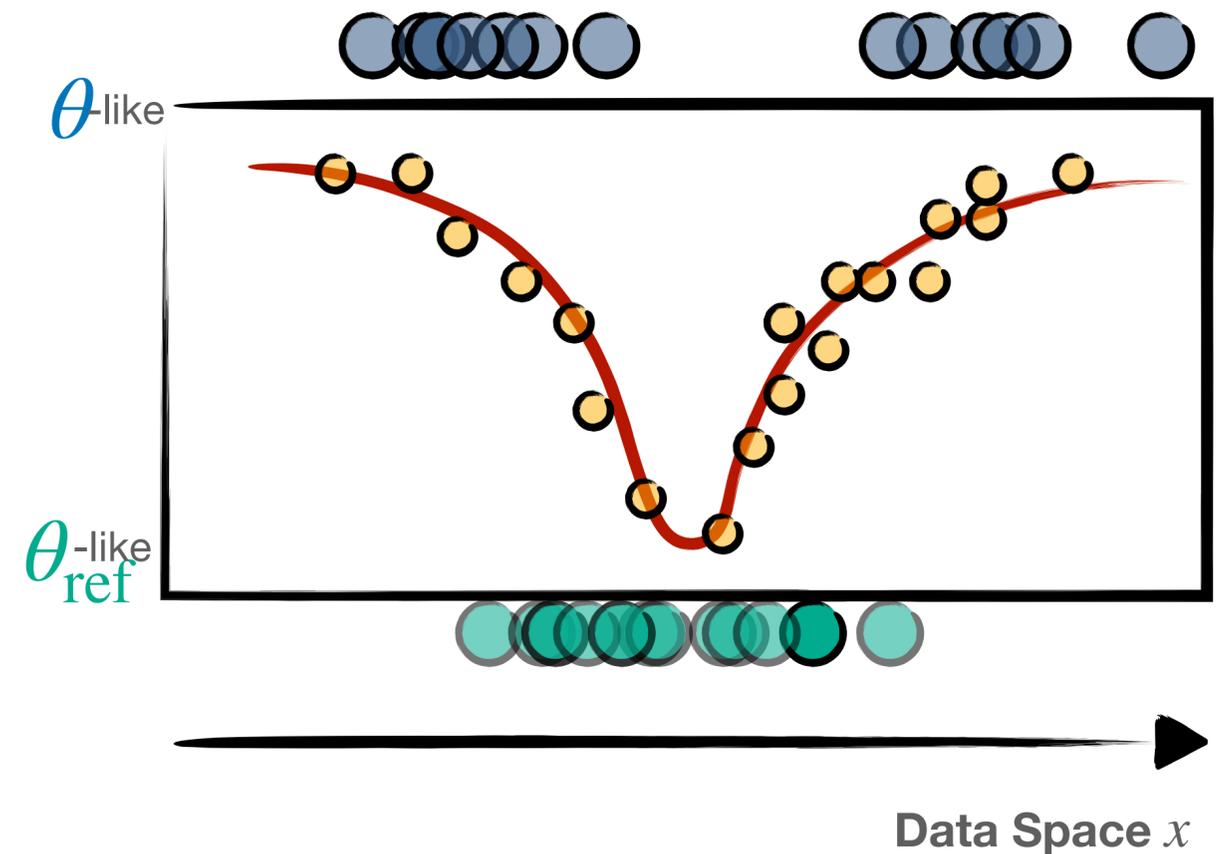
Better Training

That means: by supercharging our MC generators w/ “Mining Gold” we can extract much much better training data.

Instead of “Tug of War” ...



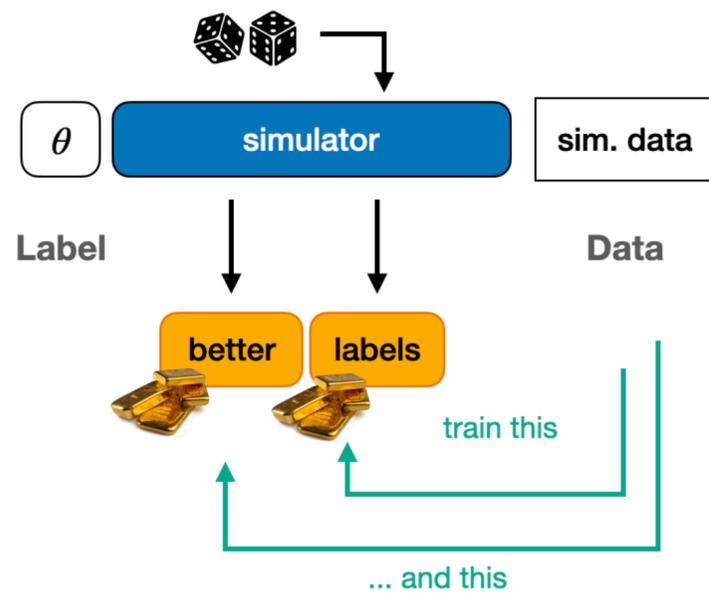
... it's just simple regression!



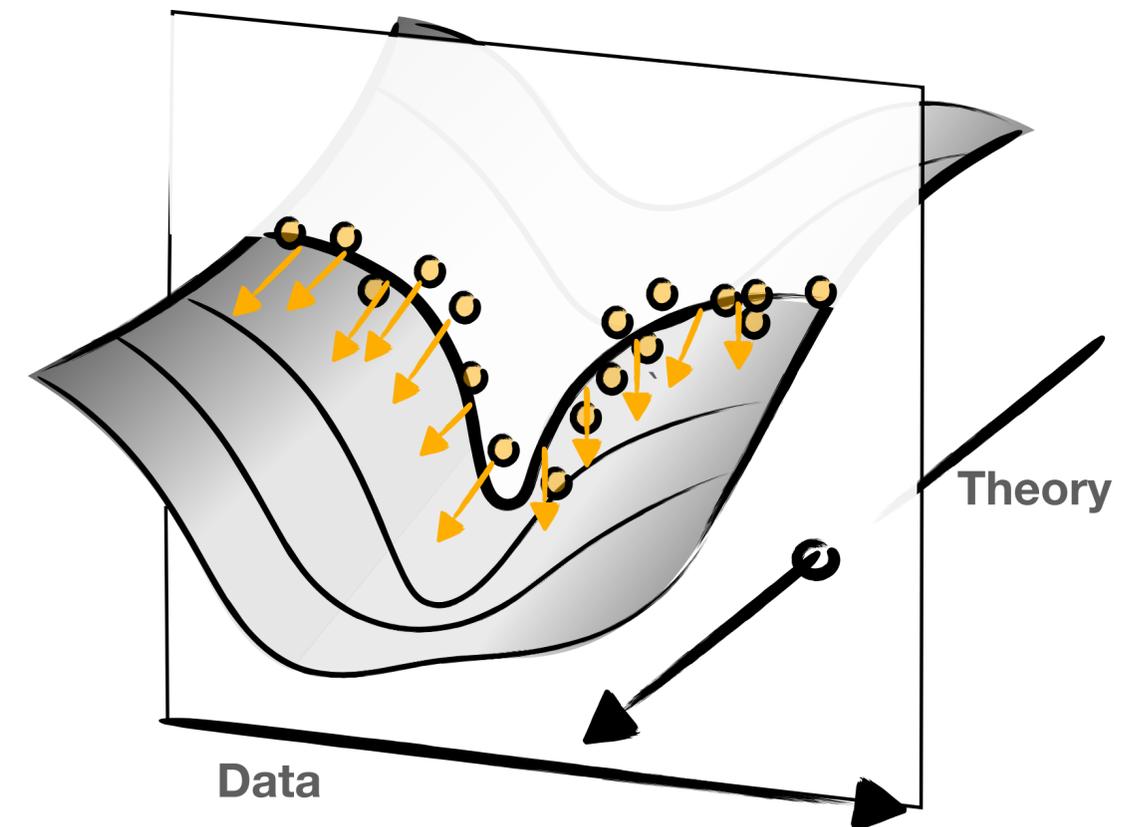
Gradients for the win

Idea: **derivatives of our noisy labels** w.r.t theory parameters are **additional powerful labels** we can extract from the simulator...

$$\nabla \log p(x, z | \theta)$$

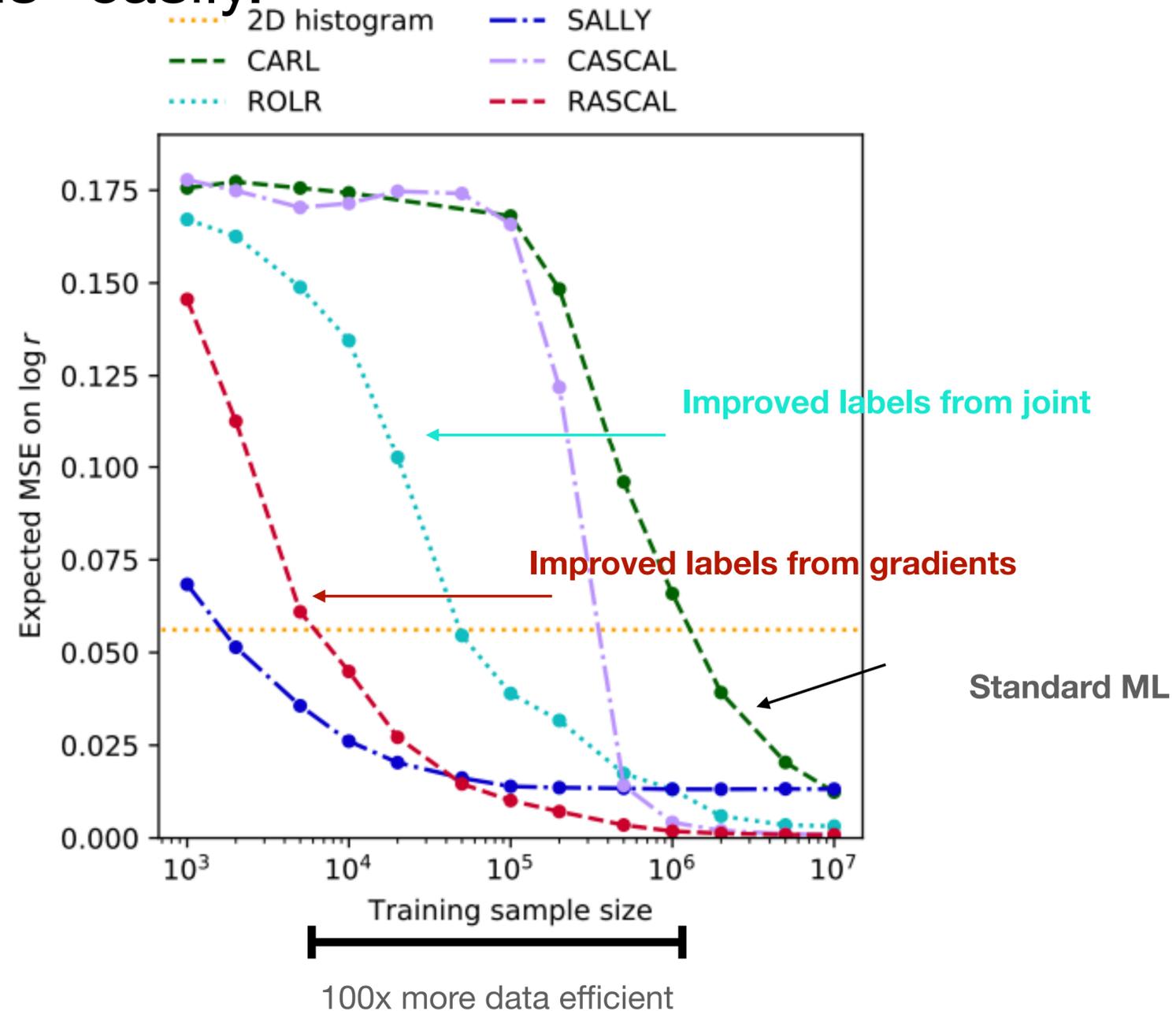


Needs a Differentiable Simulator!



Impressive Gains in Data Efficiency

This gives us the ability to train with 100x-1000x less training data. Beat “traditional analysis” easily.

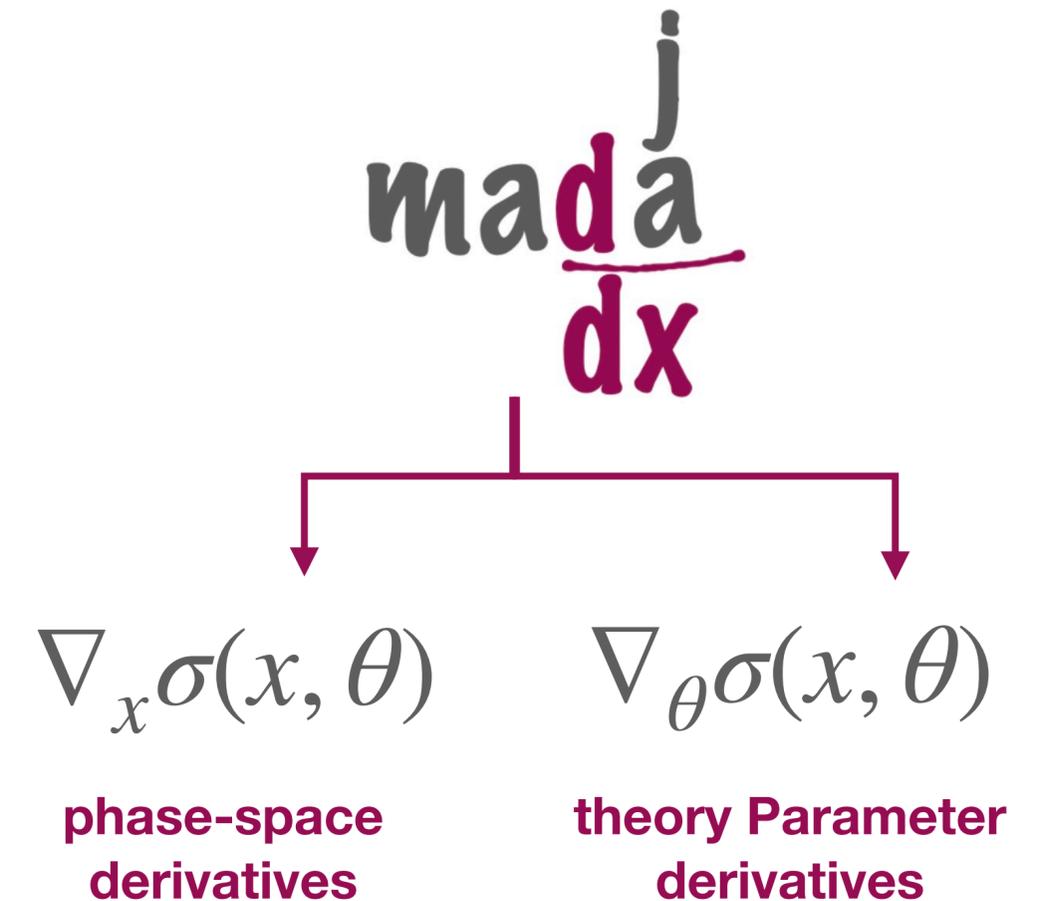
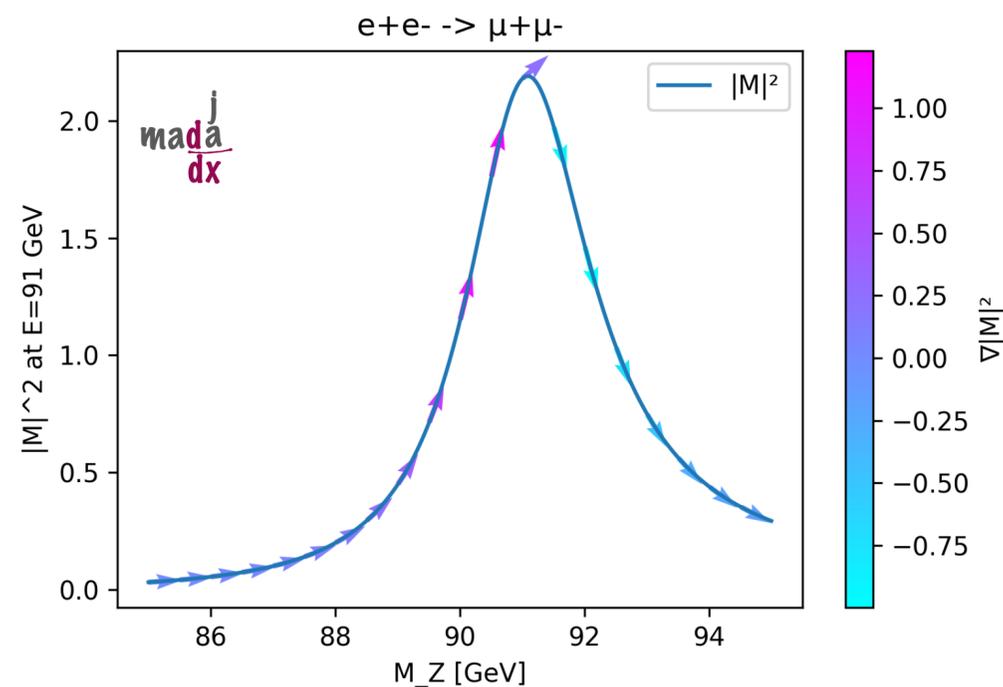


Differentiable Simulators

For these more advanced methods, it's key to have differentiable simulators in order to extract the max. amount of information from the generating process.

As we've seen: a wide & active field.

In particle physics pursued e.g. through **MadJax**



Summary

Advanced scientific predictions are rarely computable using pen & paper.

The rise of numerical and stochastic simulation

For these situations a new breed of statistical methods is emerging that fuses ML with domain simulators to extract reliable statistical inferences despite lack of likelihood

End.