

Constraining Effective Field Theories with Machine Learning

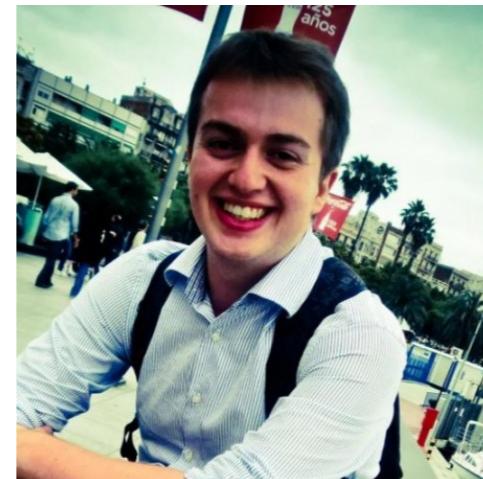
Johann Brehmer

NYU

HET seminar, BNL
June 27, 2018



Kyle Cranmer



Gilles Louppe



Juan Pavez

Constraining Effective Field Theories with Machine Learning

arXiv:1805.00013

A Guide to Constraining Effective Field Theories with Machine Learning

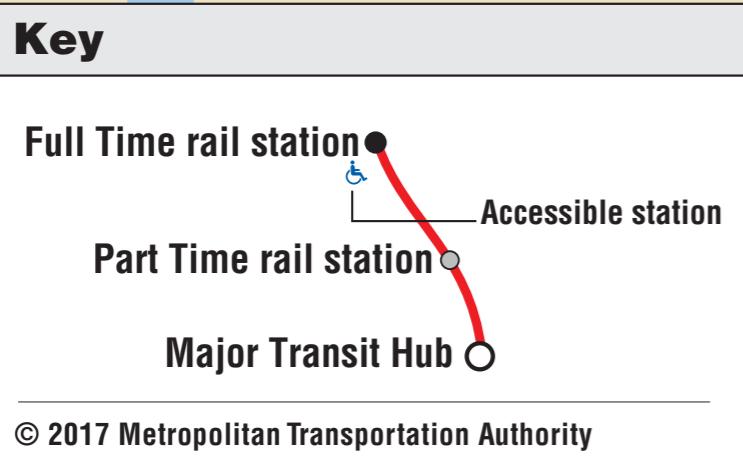
arXiv:1805.00020

Mining gold from implicit models to improve likelihood-free inference

arXiv:1805.12244

Thanks to Kyle and Gilles for ~~letting me steal~~ inspiring many slides!

Long Island Rail Road



Likelihood-
free inference

Established
methods

Constraining
EFTs with ML

Mining gold
from the
simulator



Likelihood-free inference

The Galton board



FIG. 7.

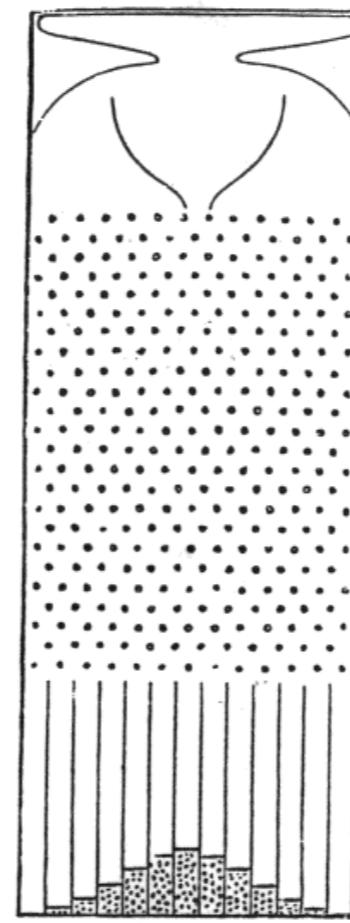


FIG. 8.

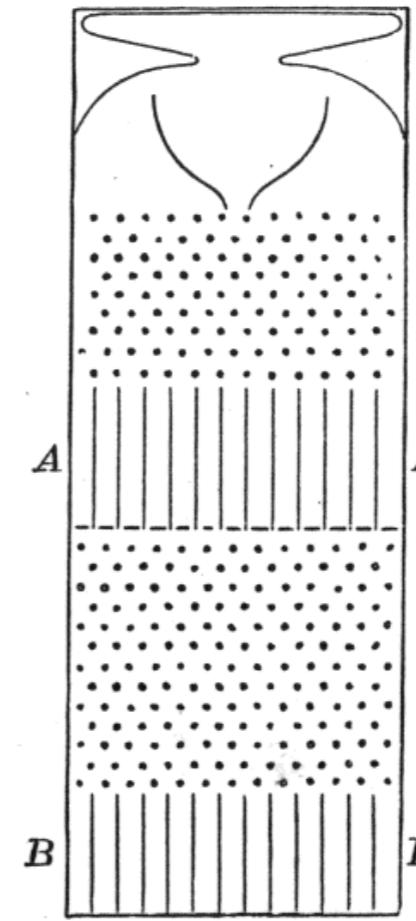
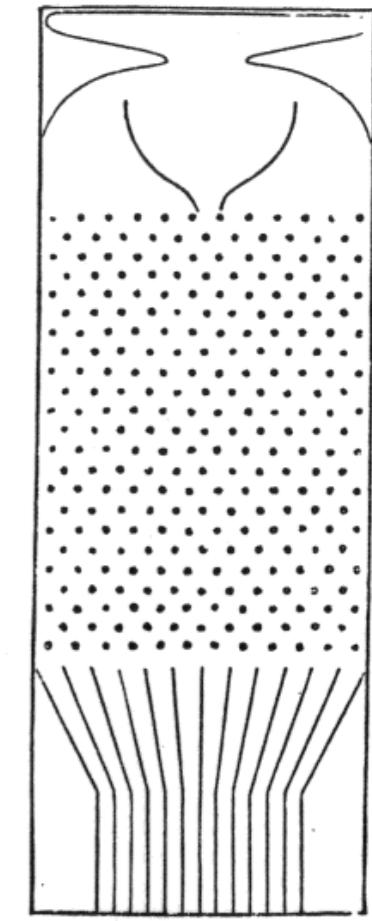
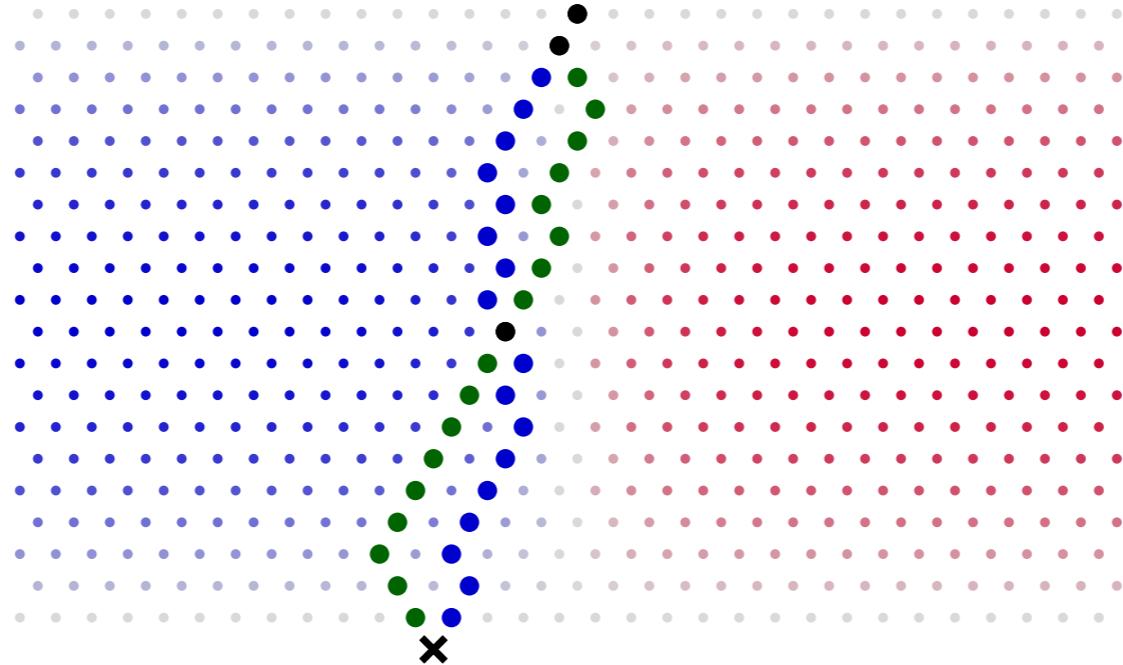


FIG. 9.



[Source: F. Galton 1889]

Probabilities



Probability of ending in bin x :

$$p(x) = \int dz p(x, z)$$

Sum over
all trajectories
("latent variables")

Probability of
each path z
from start to x

The generalized Galton board

- What if probability to go left at a nail is not always 0.5, but some (known) function of some parameters θ ?
- **Prediction**: given θ , generate observations $\{x_i\}$... just drop balls!
- **Inference**: given observations $\{x_i\}$, what are the most likely values for θ ?

The generalized Galton board

- What if probability to go left at a nail is not always 0.5, but some (known) function of some parameters θ ?
- **Prediction**: given θ , generate observations $\{x_i\}$... just drop balls!
- **Inference**: given observations $\{x_i\}$, what are the most likely values for θ ?

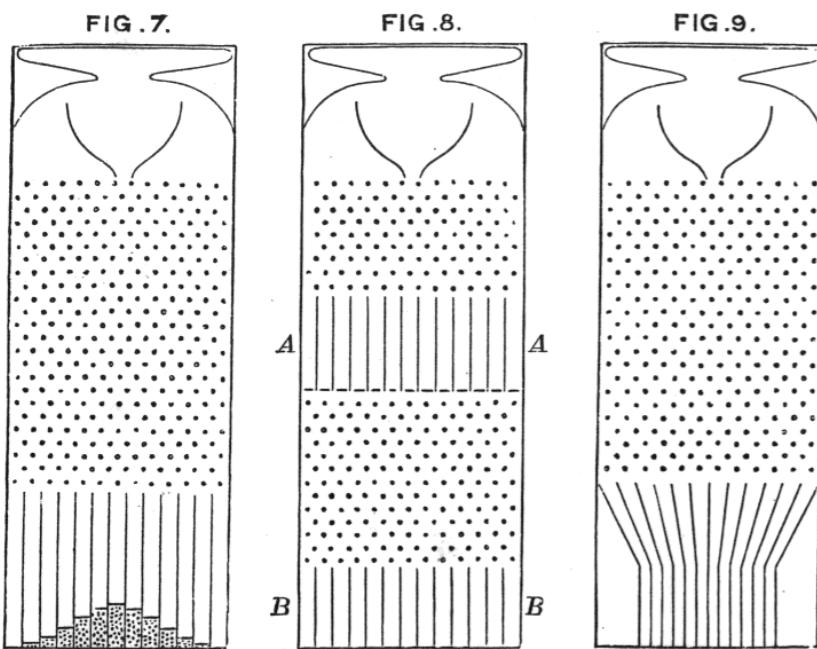
“Easy” problem if we can evaluate likelihood $p(x_i|\theta)$. But

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

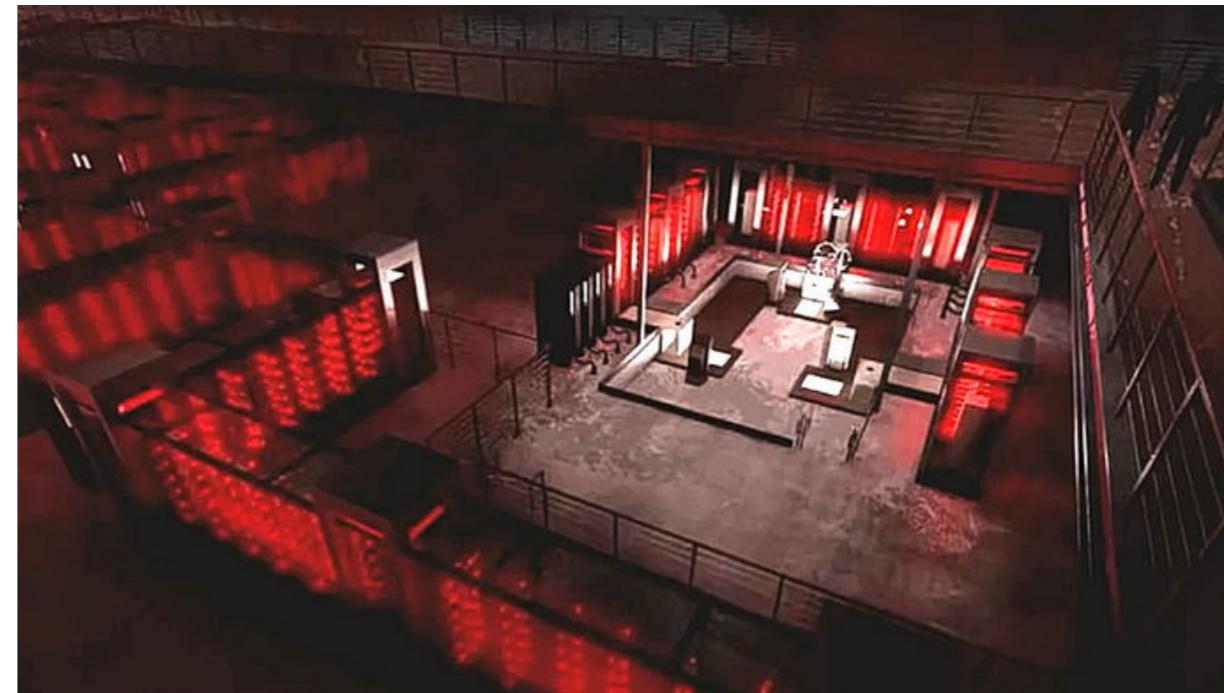
The number of possible **paths** z can be huge, and it becomes impossible to calculate the integral

⇒ Likelihood not tractable, only implicitly defined through “simulator”

Galton board: metaphor for simulator-based science



[Source: F. Galton 1889]



[Source: HBO 2018]

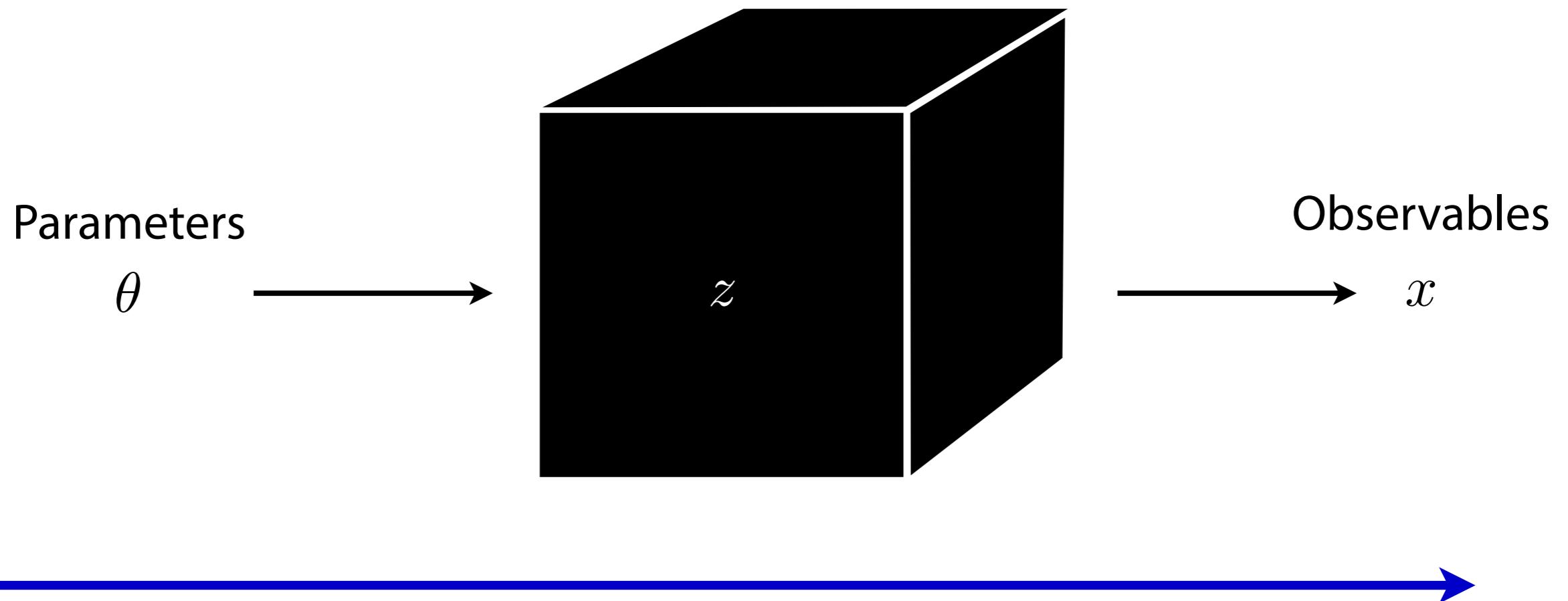
Galton board device → Computer simulation

Parameters θ → Parameters of interest θ

Bins x → Observables x

Path z → Latent variables z
(stochastic execution trace through simulator)

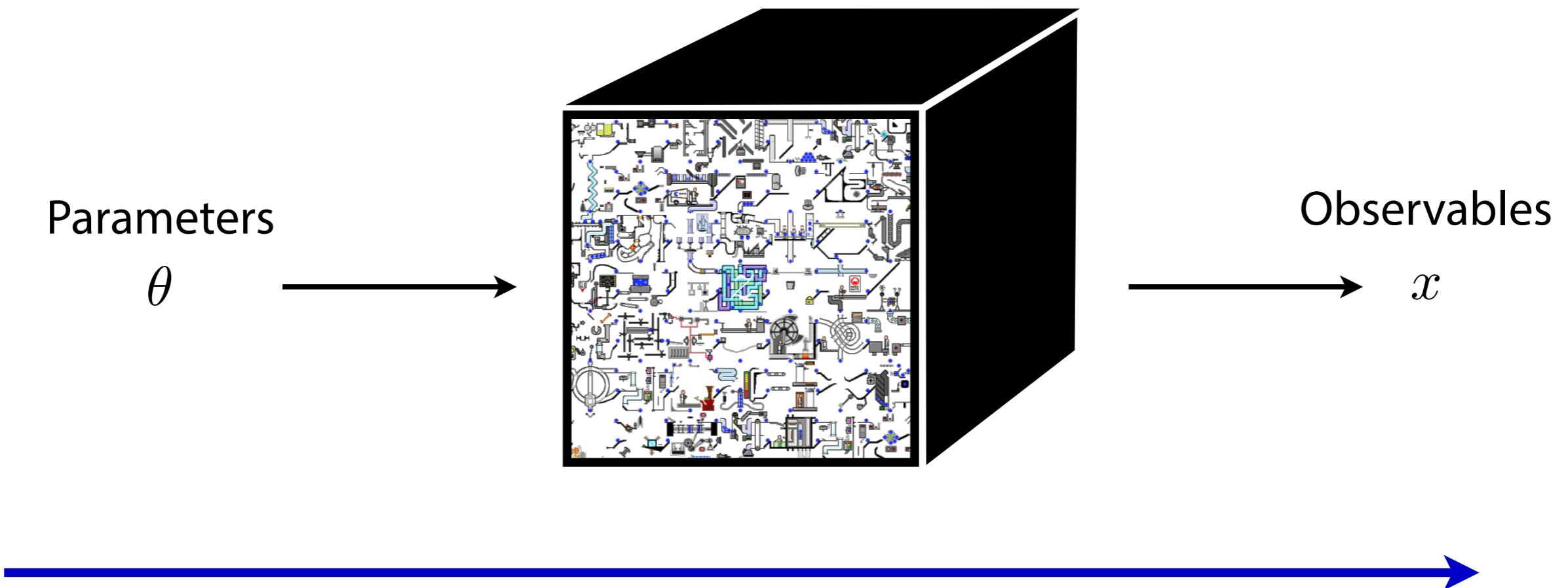
Likelihood-free inference



Prediction (simulation):

- Well-understood mechanistic model
- Simulator can generate samples

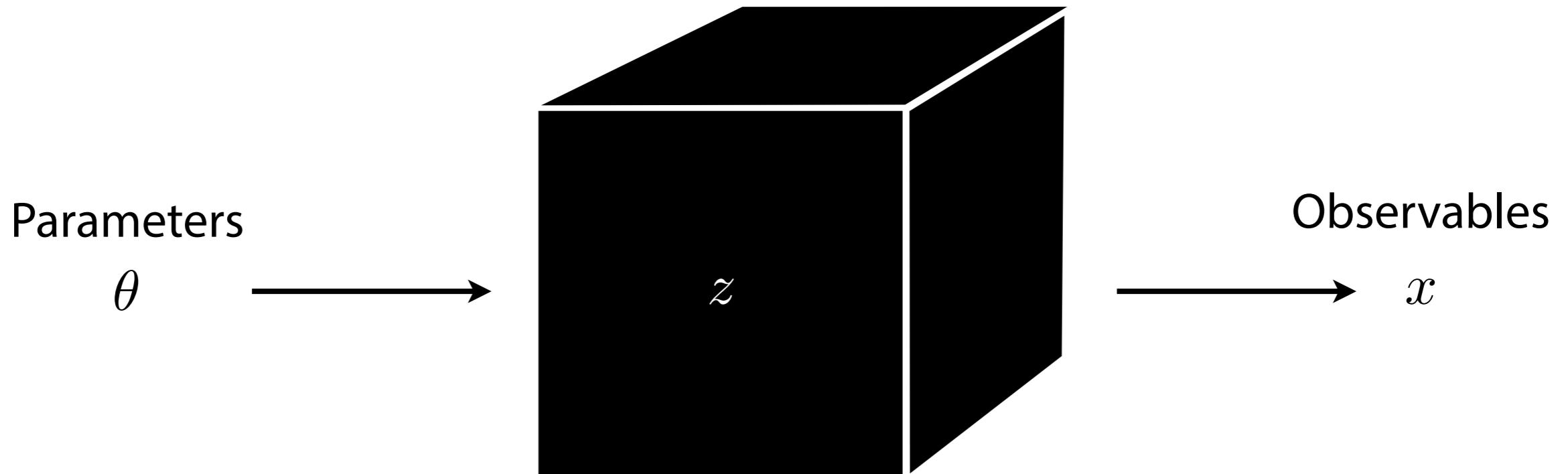
Likelihood-free inference



Prediction (simulation):

- Well-understood mechanistic model
- Simulator can generate samples

Likelihood-free inference



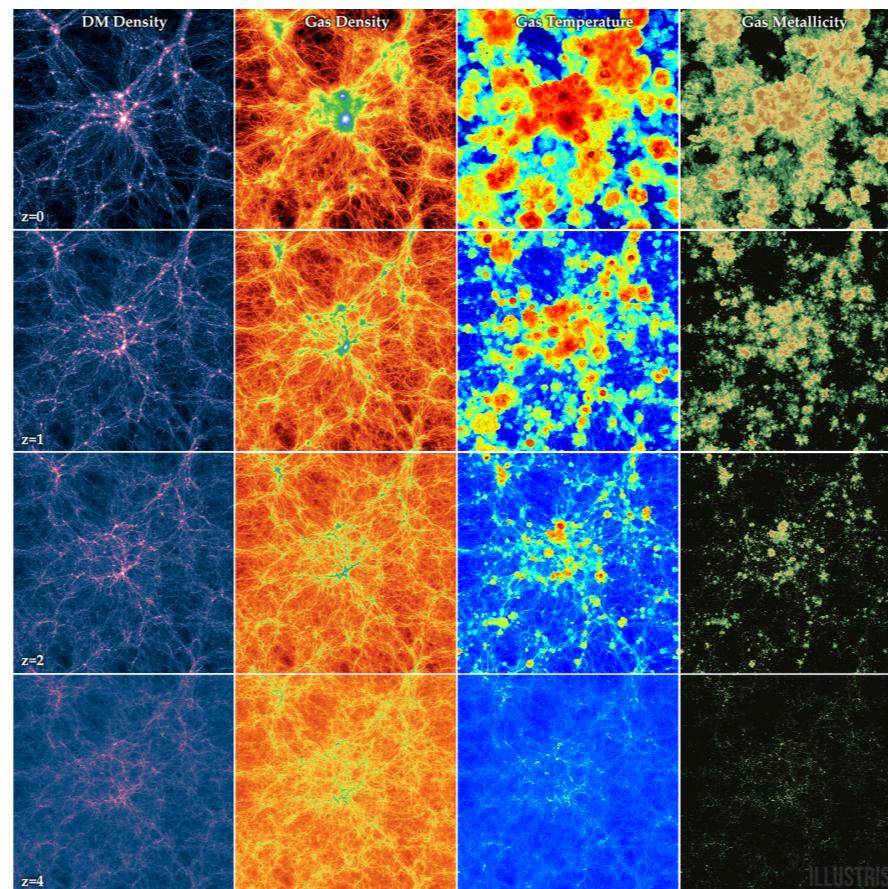
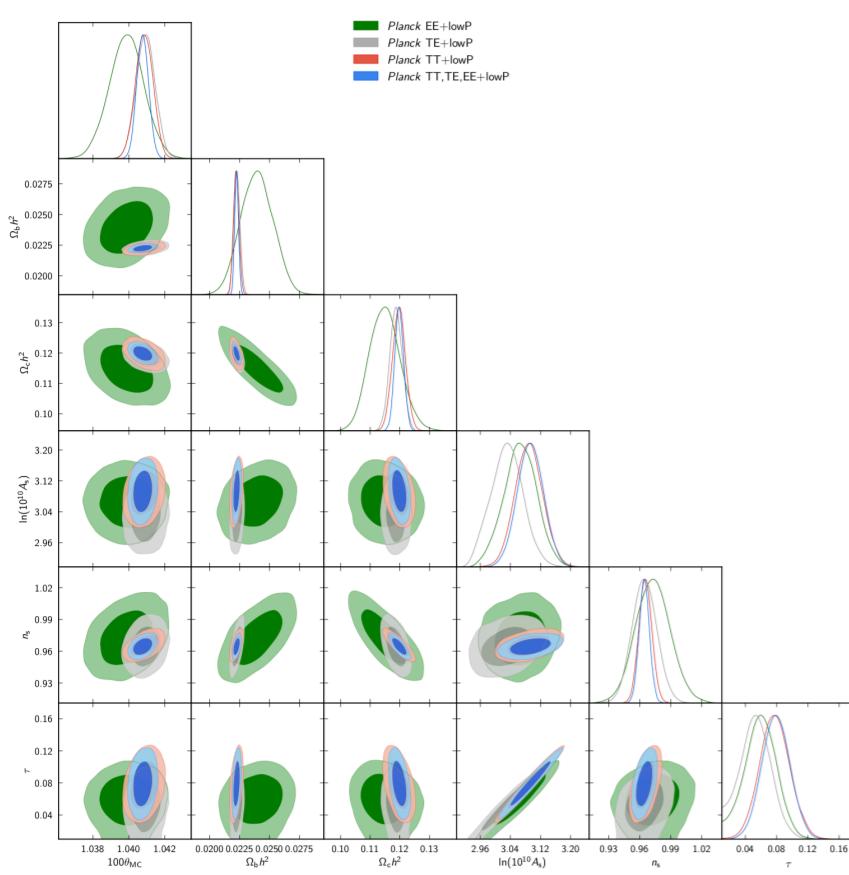
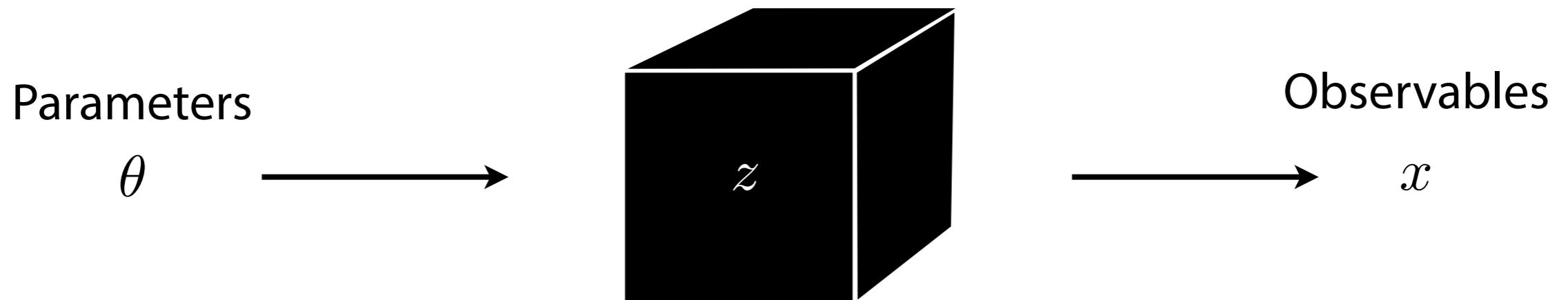
Prediction (simulation):

- Well-understood mechanistic model
- Simulator can generate samples

Inference:

- Likelihood function $p(x|\theta)$ is intractable
- Goal: estimator $\hat{p}(x|\theta)$

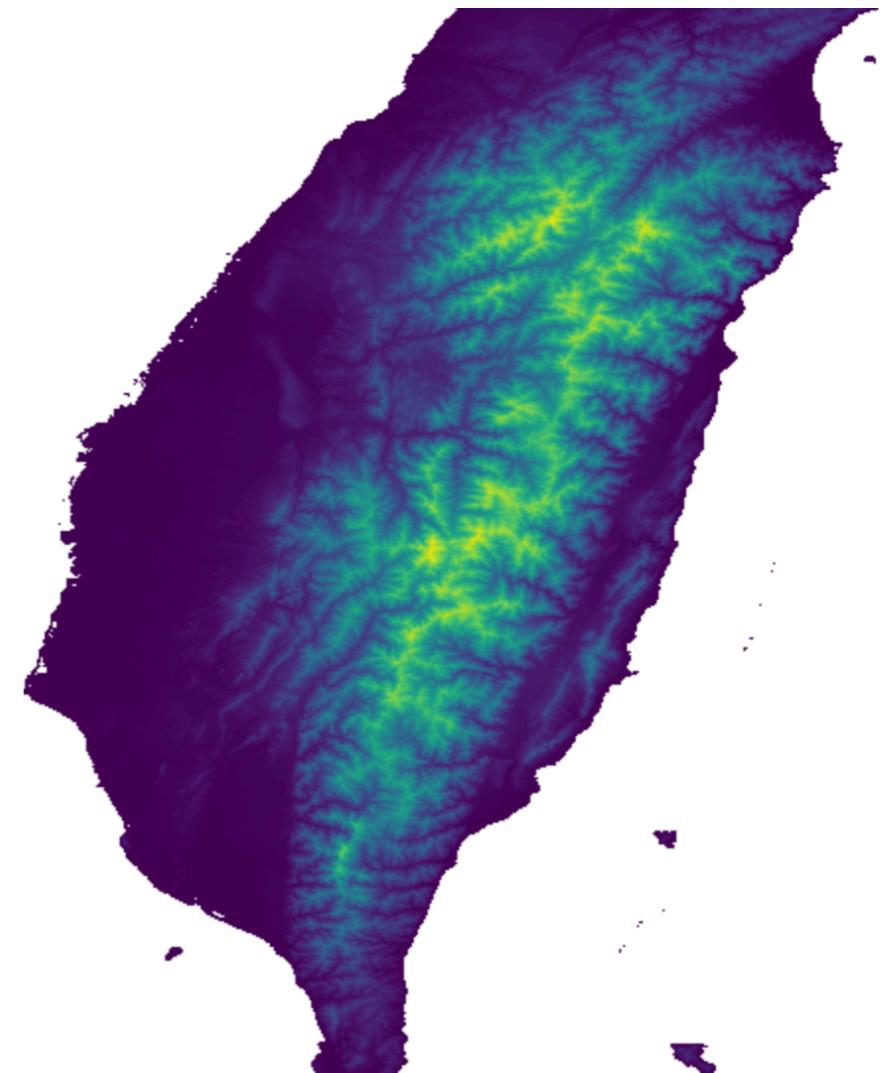
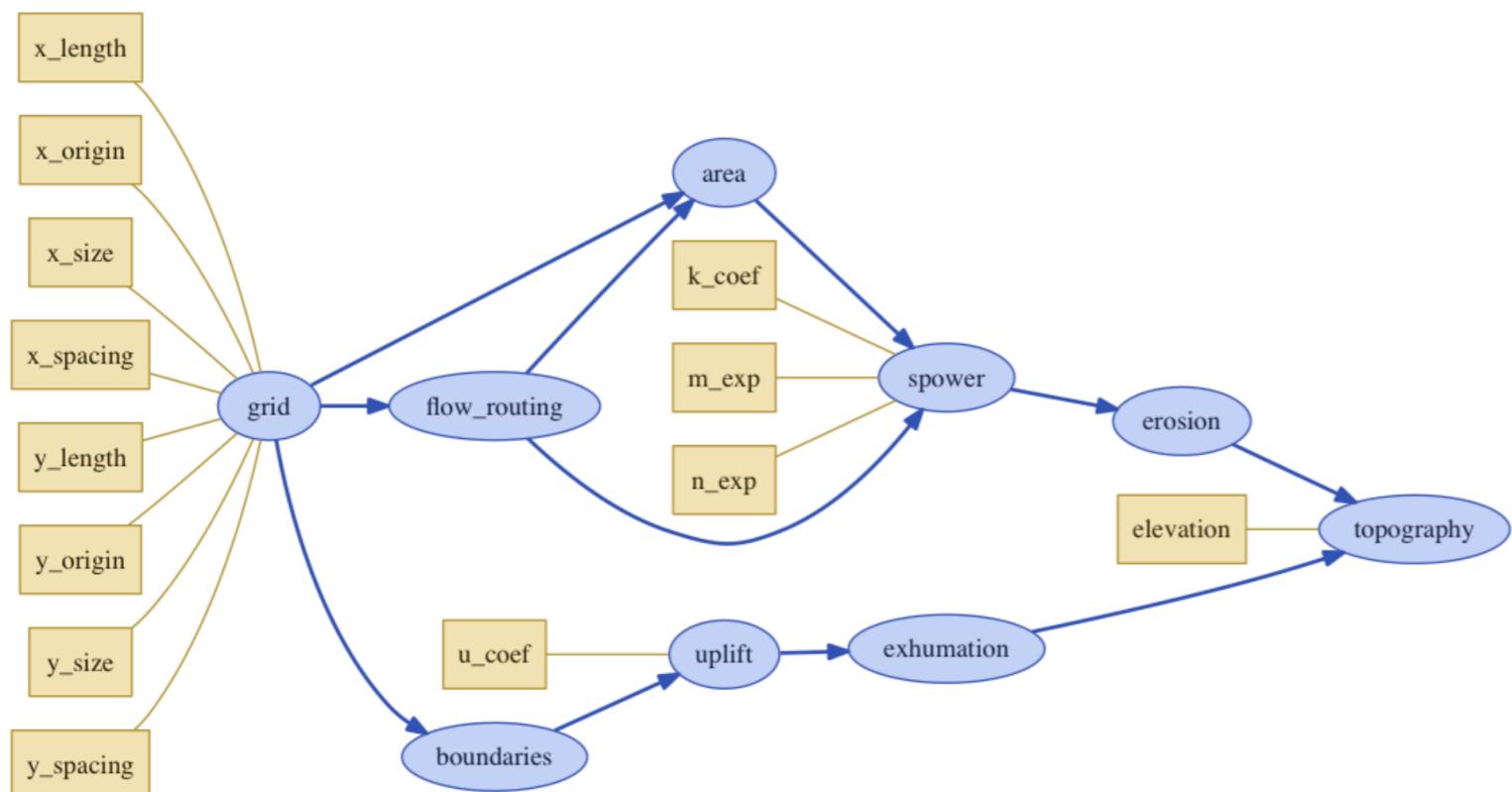
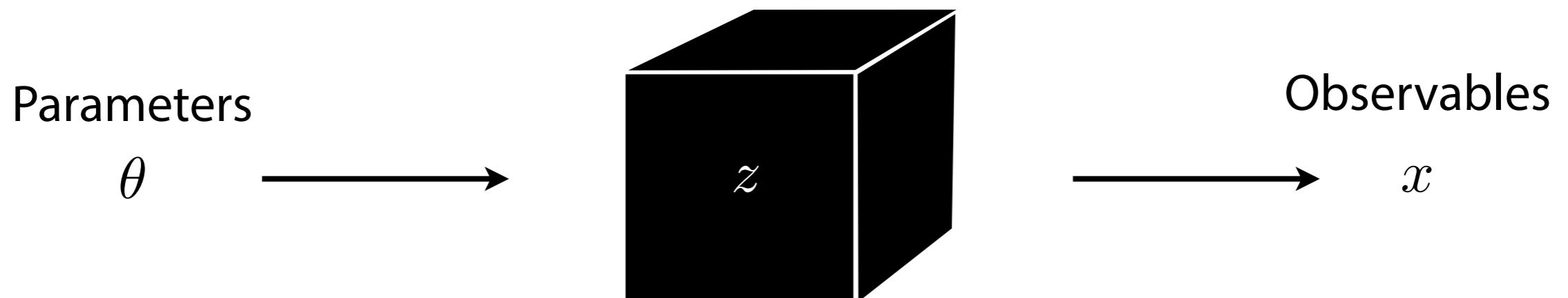
Cosmological N-body simulations



[Source: Planck 1502.01589]

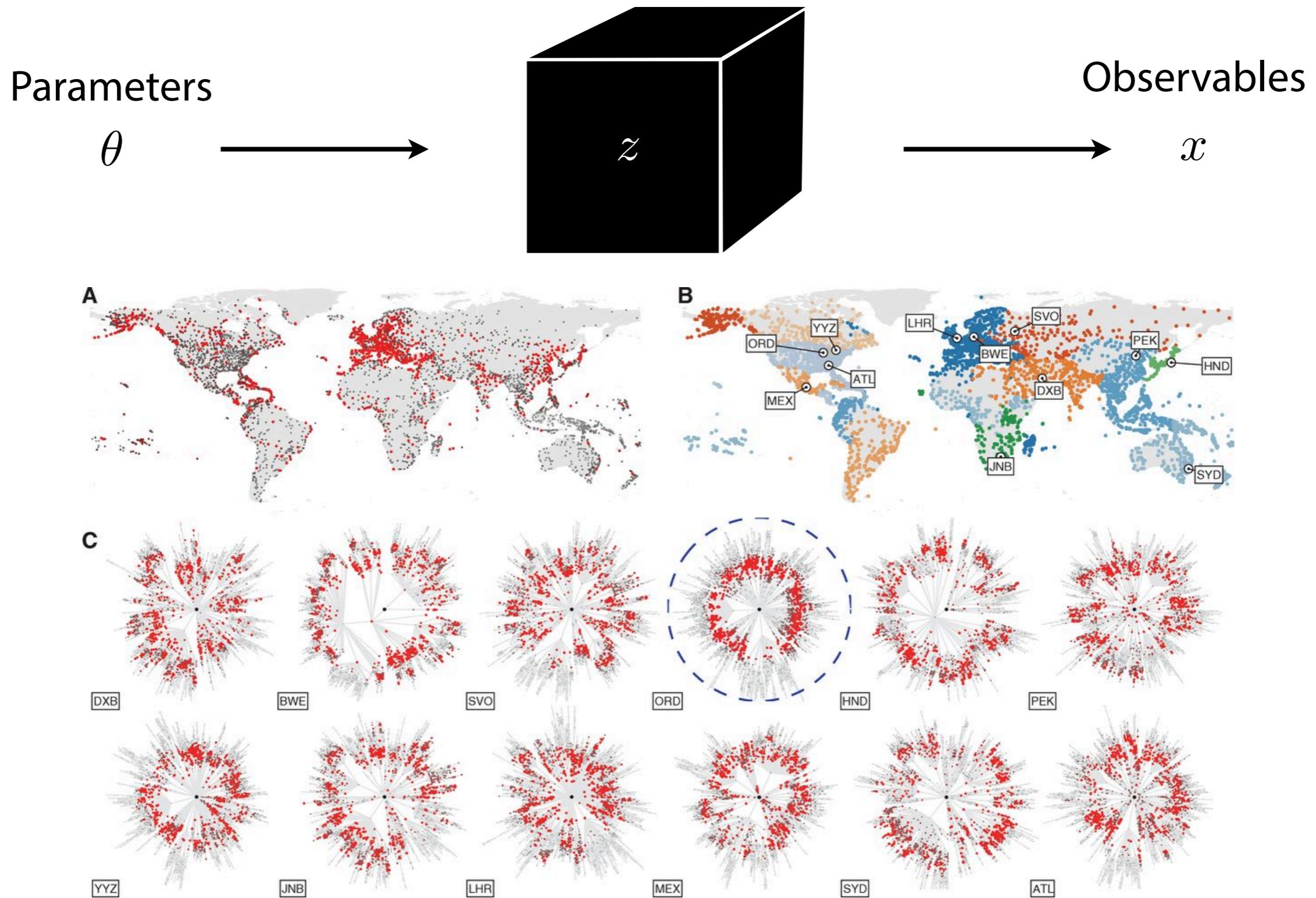
[Source: Illustris 1405.2921]

Computational topography



[Source: Benoit Bovy]

Epidemiology



[Source: D. Brockmann, D. Helbing 2013]

Particle physics

Parameters
of interest

Theory
parameters

θ



Evolution

Particle physics

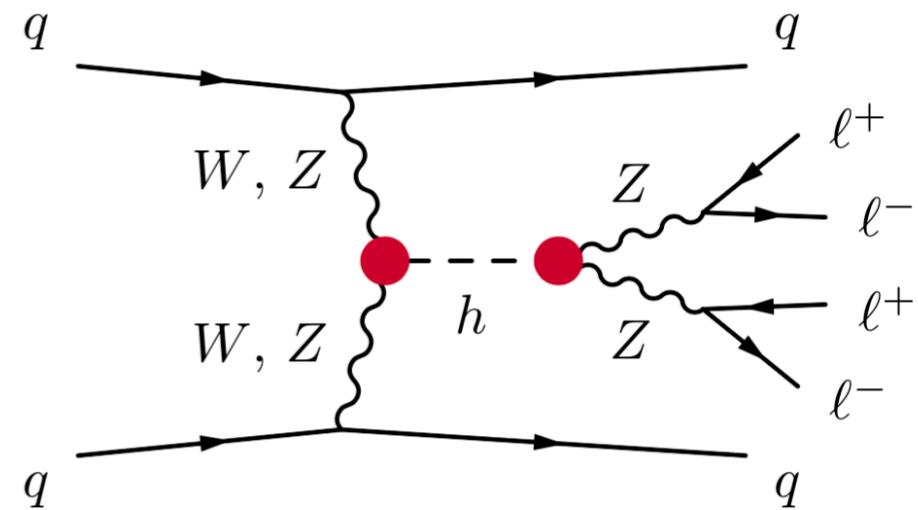
Latent variables

Parameters
of interest

Parton-level
momenta

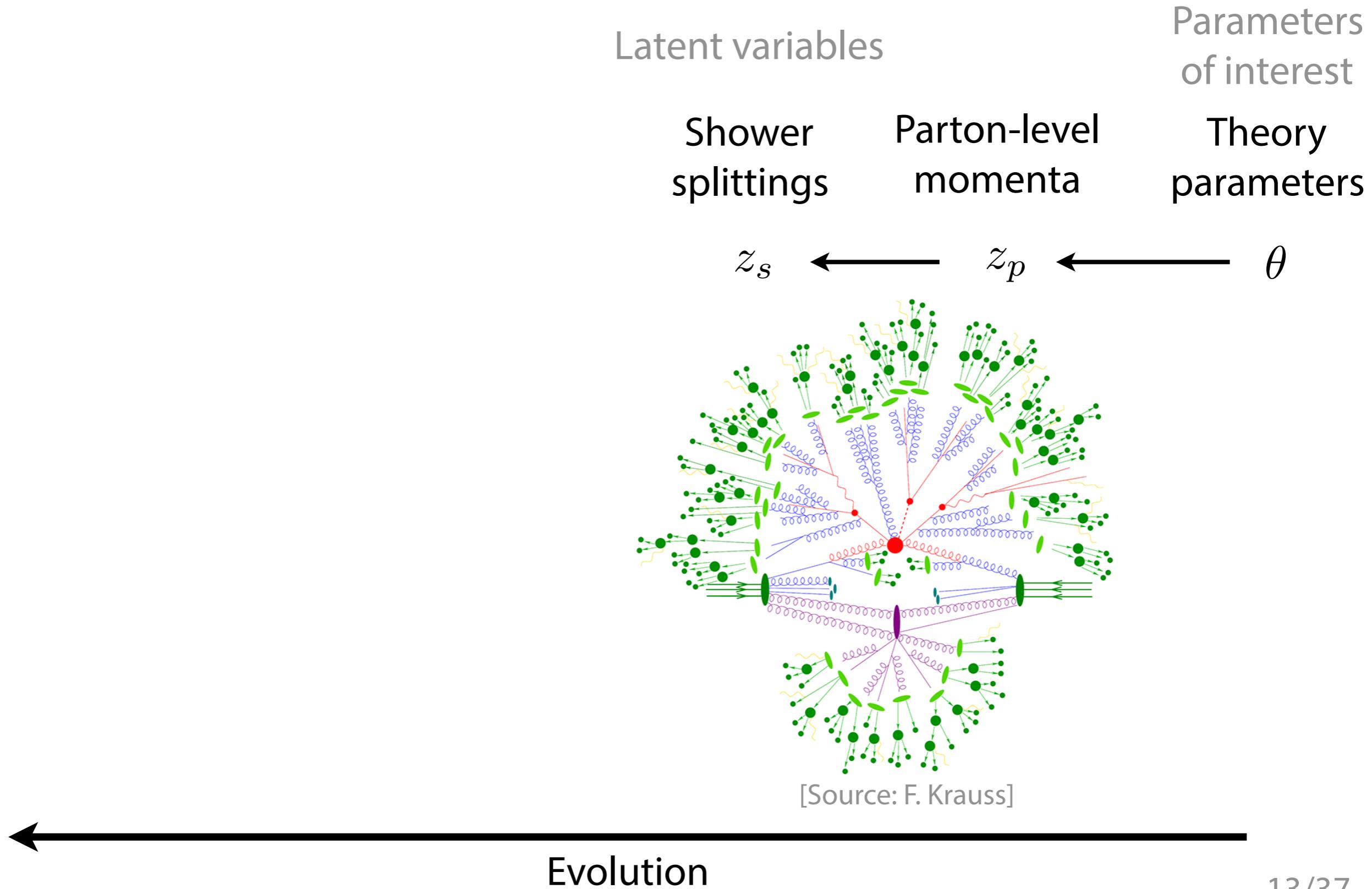
Theory
parameters

$$z_p \leftarrow \theta$$

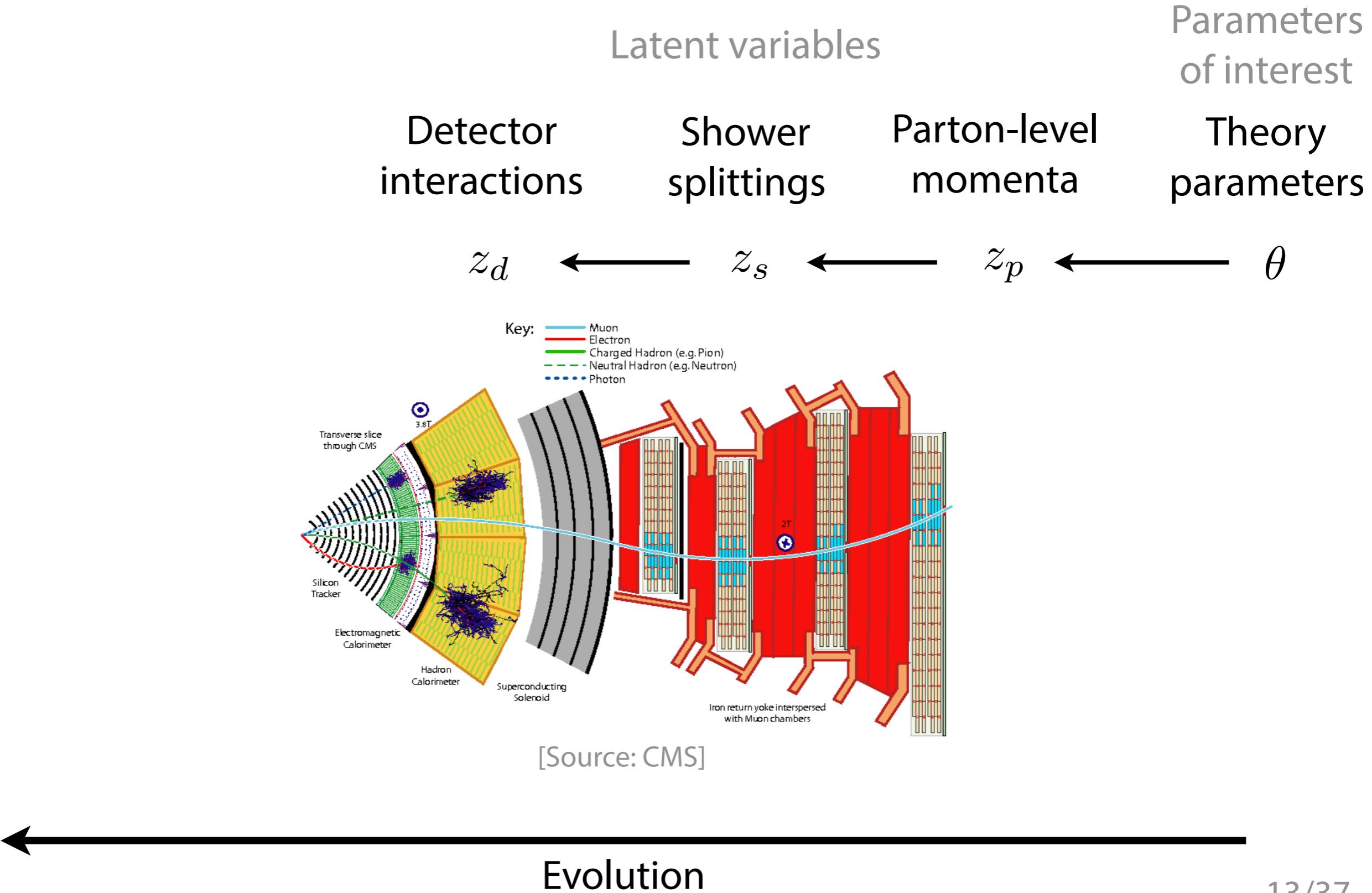


Evolution

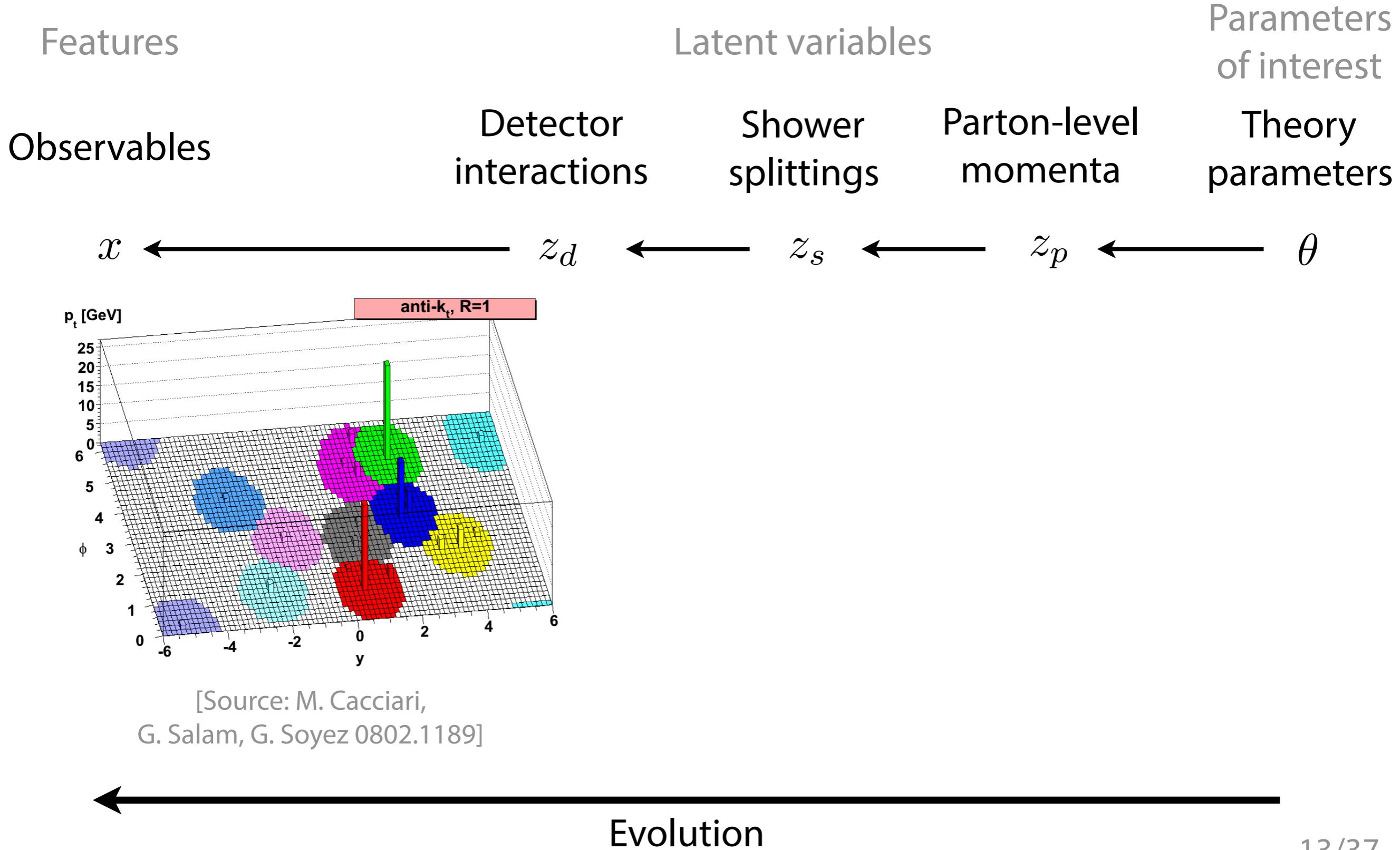
Particle physics



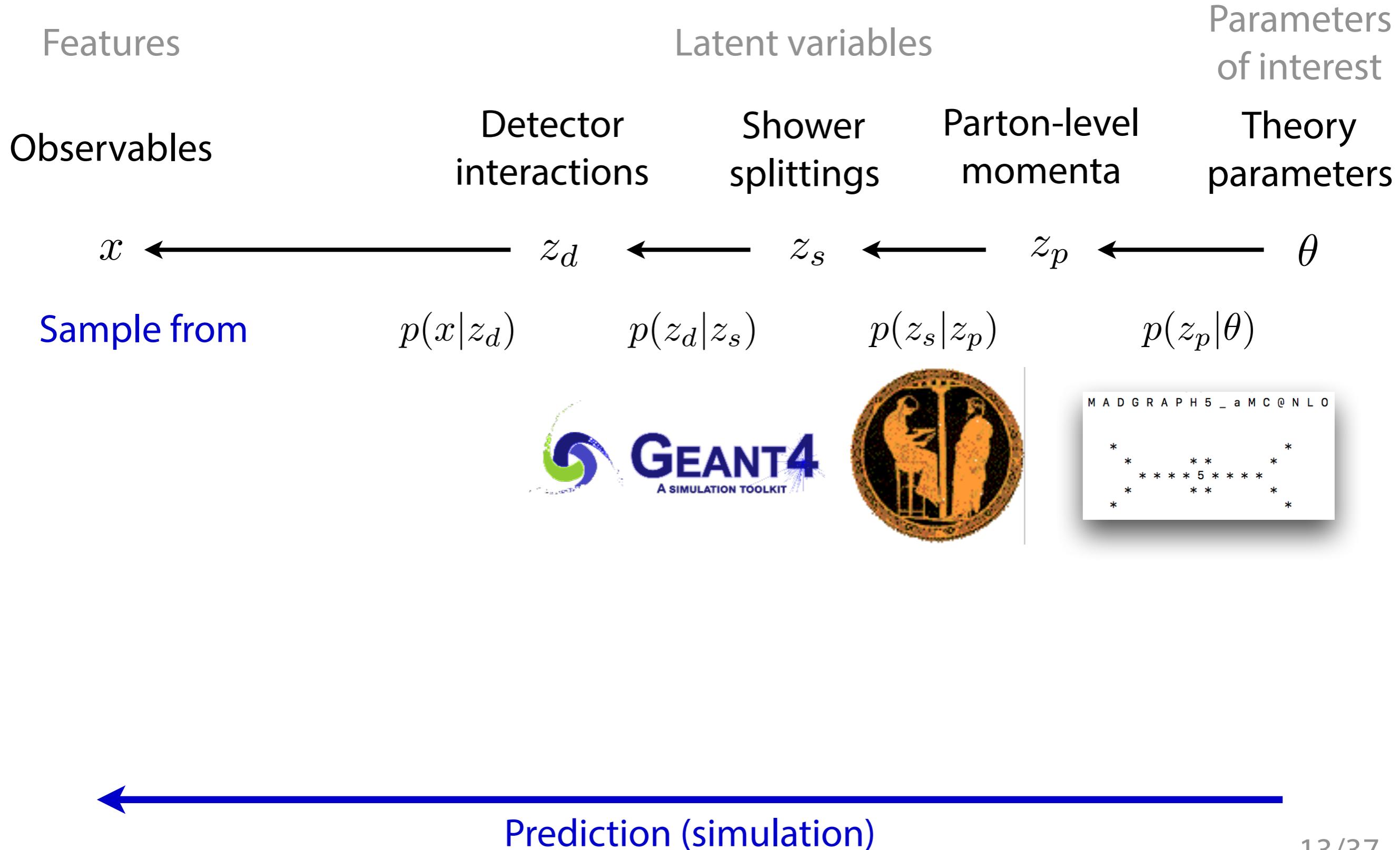
Particle physics



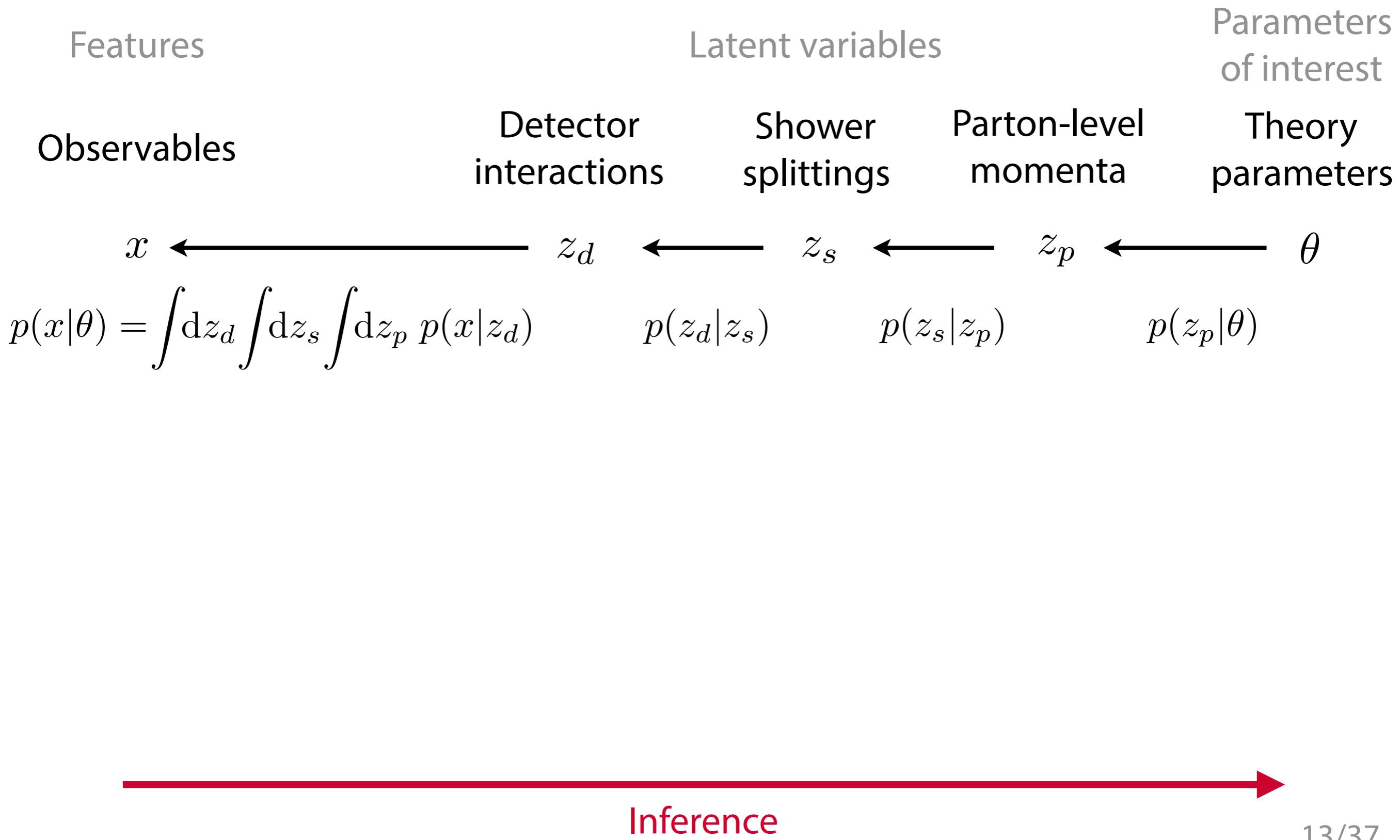
Particle physics



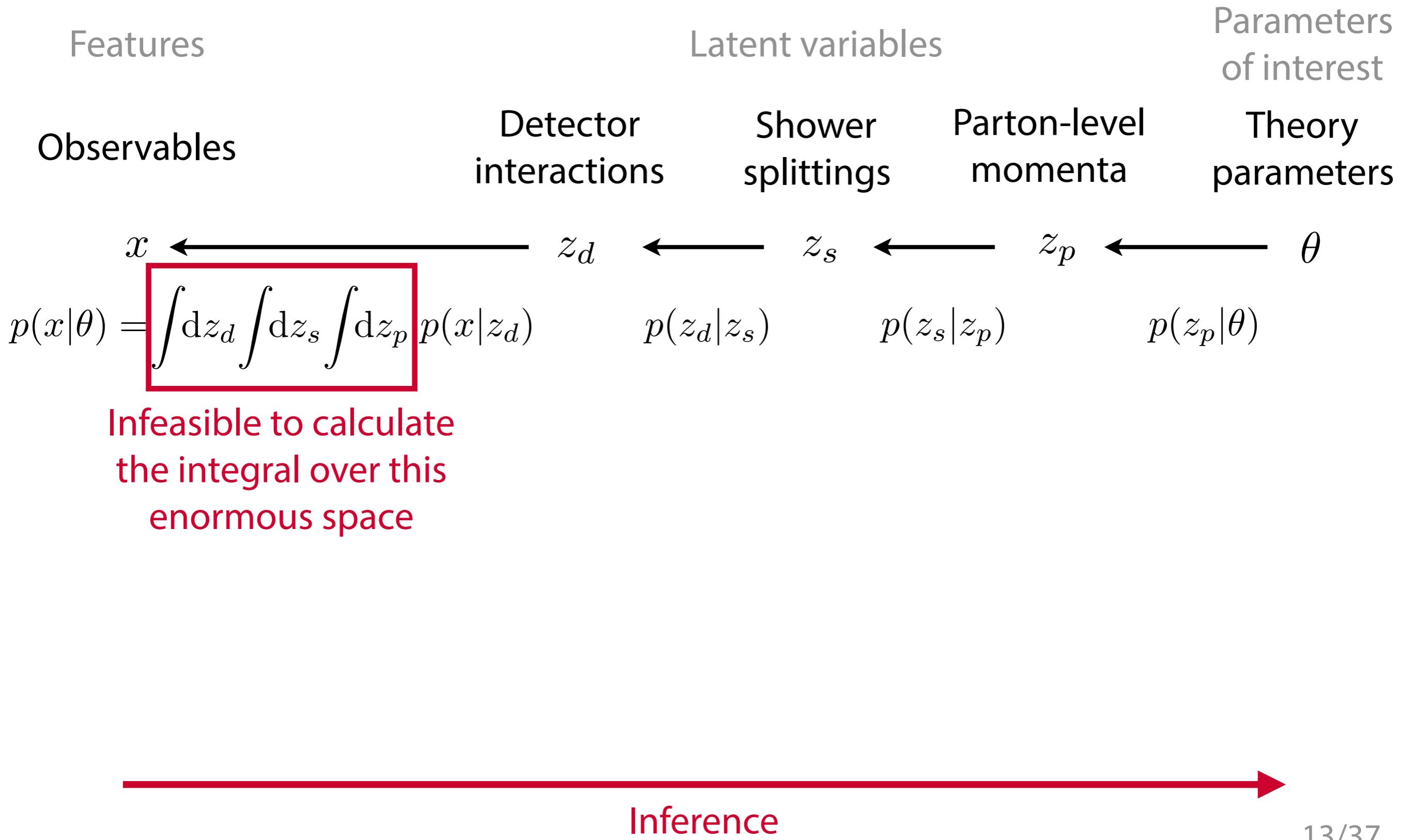
Particle physics



Particle physics

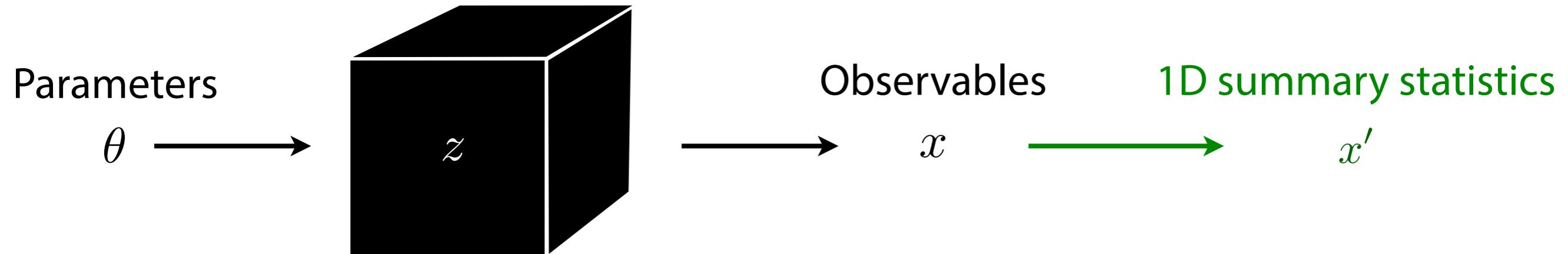


Particle physics

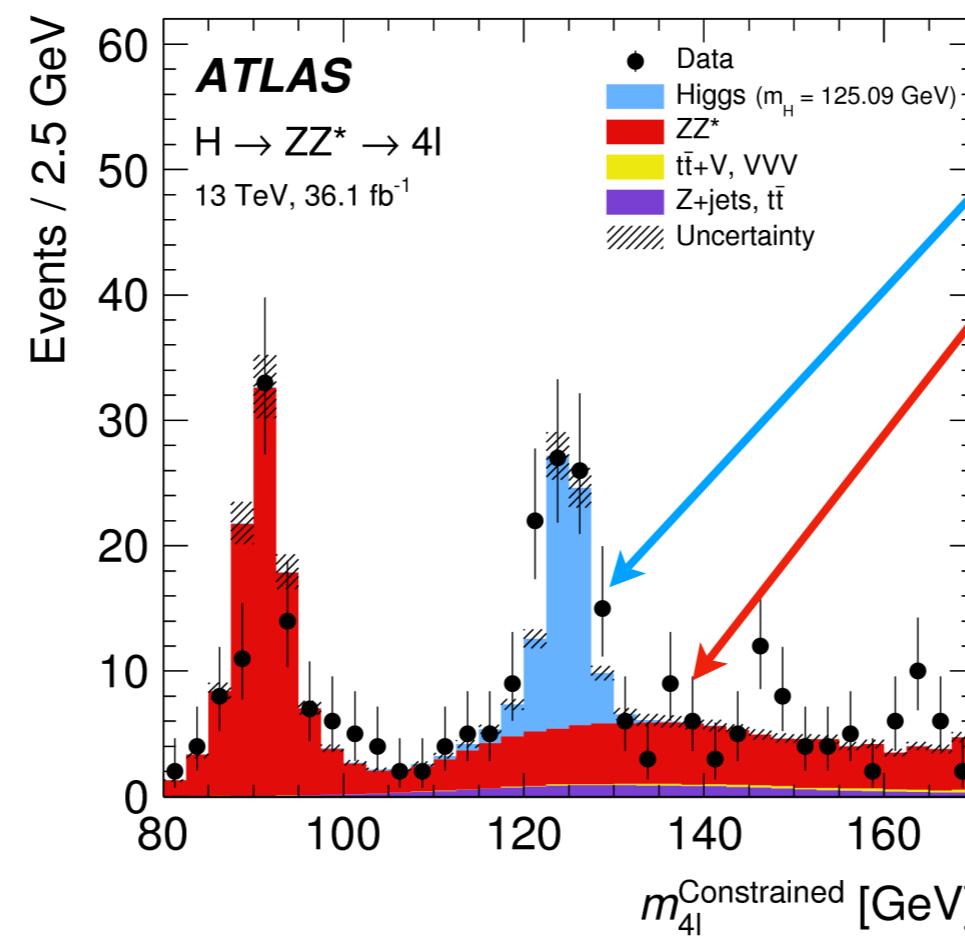
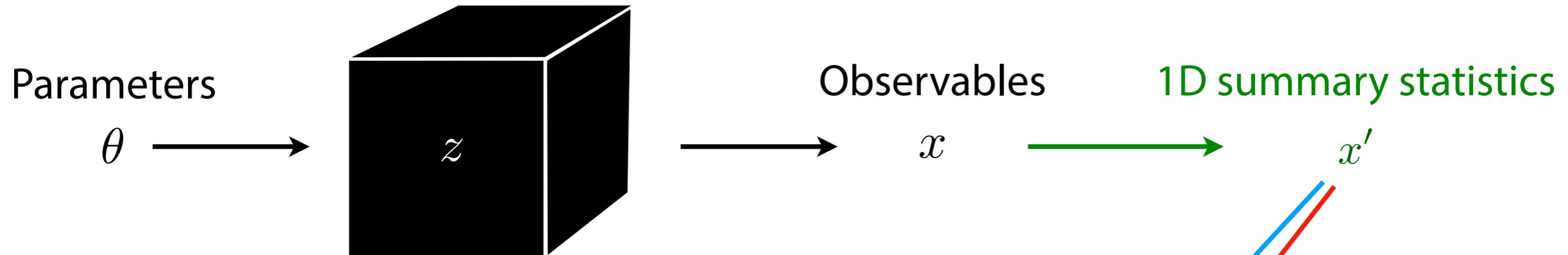


**Why has that not stopped
us so far?**

The physicist's way



The physicist's way



[Source: ATLAS 1712.02304]

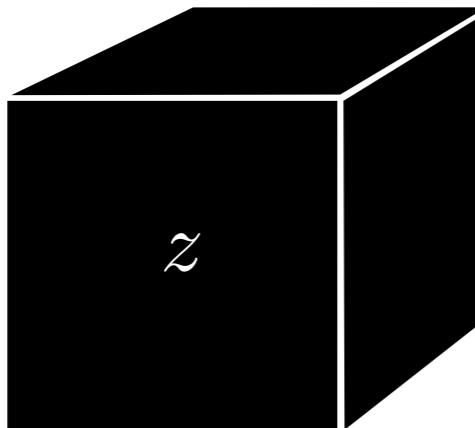
Run simulator for
different θ ,
fill histograms

x'

The physicist's way

Parameters

$$\theta$$



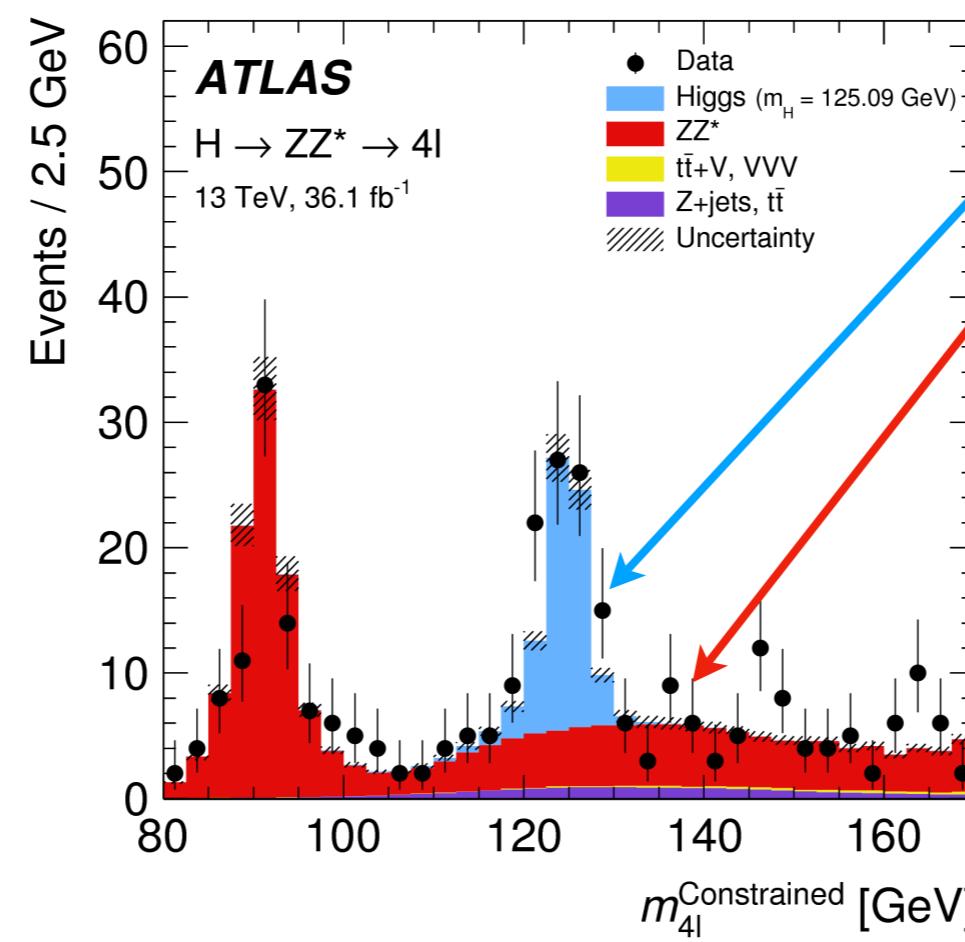
Observables

$$x$$

1D summary statistics

$$x'$$

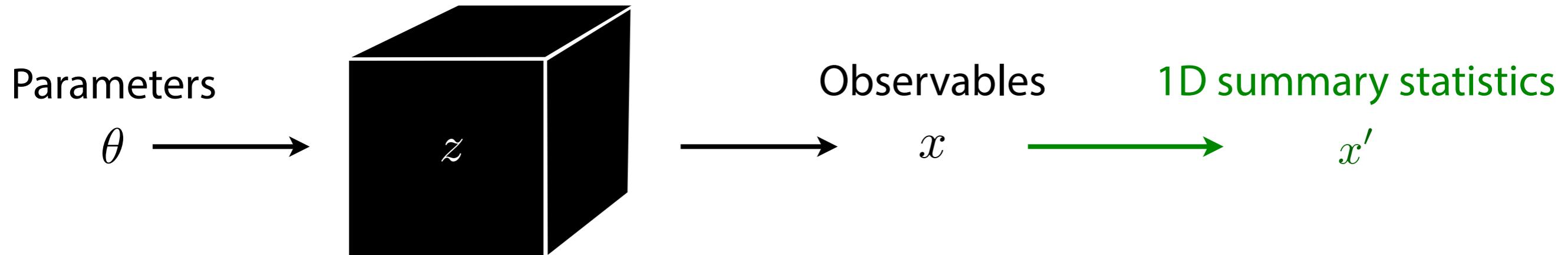
$$\hat{p}(x|\theta) = p(x'|\theta) =$$



Run simulator for
different θ ,
fill histograms

$$x'$$

Does the histogram method scale?

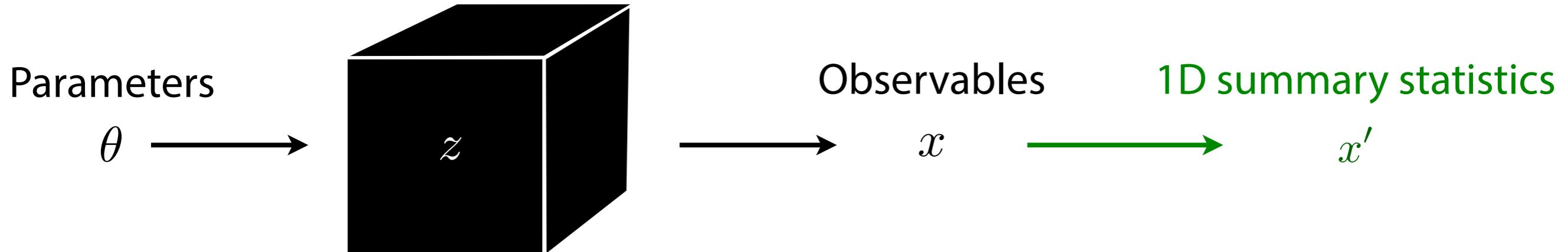


- Choosing x' is difficult and problem-dependent
- Often there is no single good variable — compressing to any x' loses information!

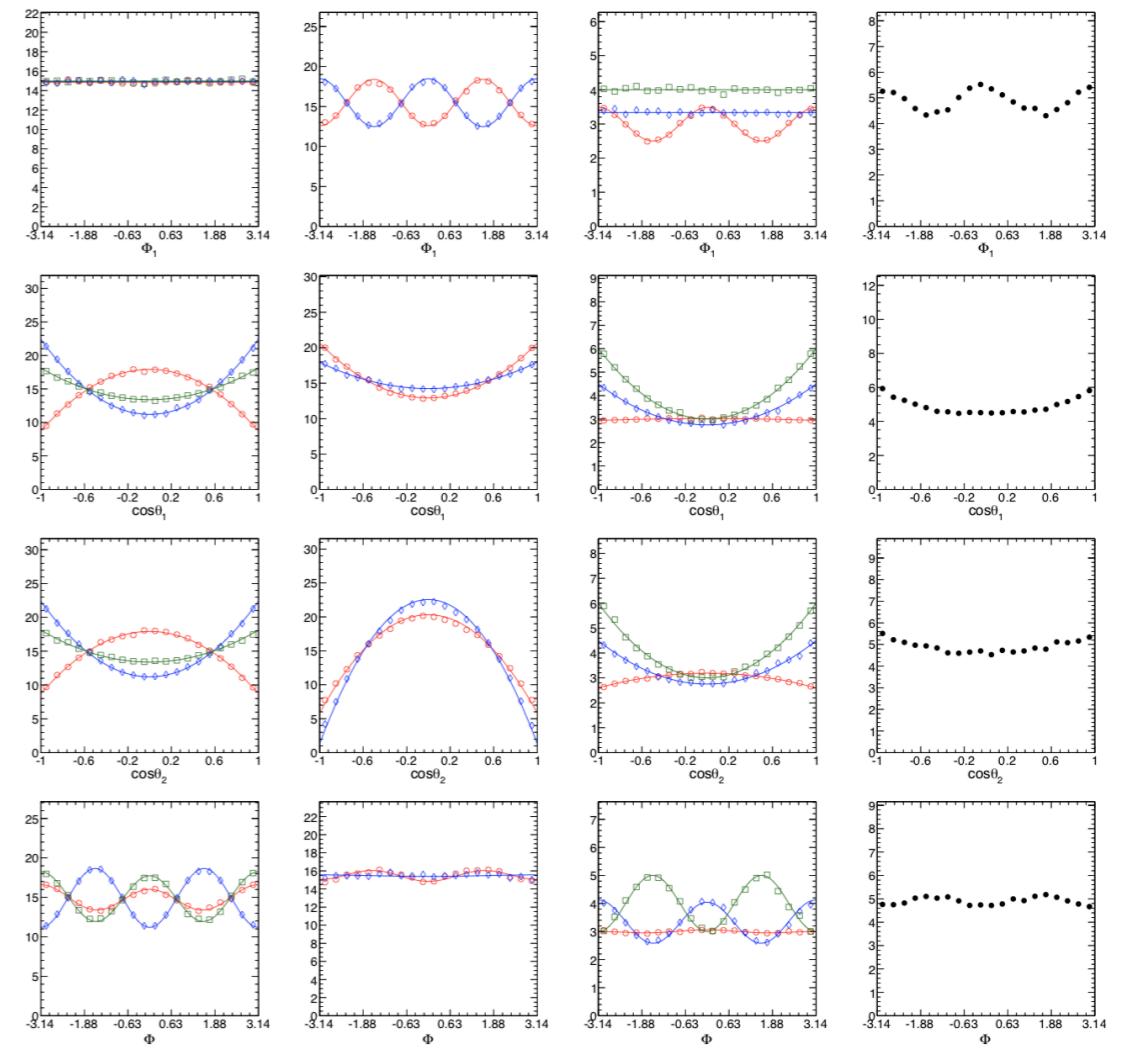
[JB, K. Cranmer, F. Kling, T. Plehn 1612.05261;
JB, F. Kling, T. Plehn , T. Tait 1712.02350]

- Ideally: analyze high-dimensional x including all correlations

Does the histogram method scale?



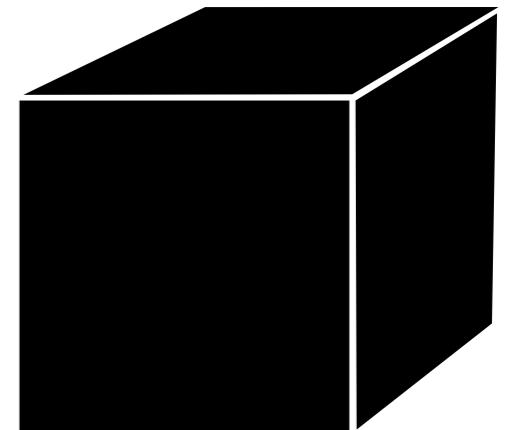
- Choosing x' is difficult and problem-dependent
 - Often there is no single good variable — compressing to any x' loses information!
- [JB, K. Cranmer, F. Kling, T. Plehn 1612.05261;
JB, F. Kling, T. Plehn , T. Tait 1712.02350]
- Ideally: analyze high-dimensional x including all correlations



[Source: Bolognesi et al. 1208.4018

An incomplete list of established methods

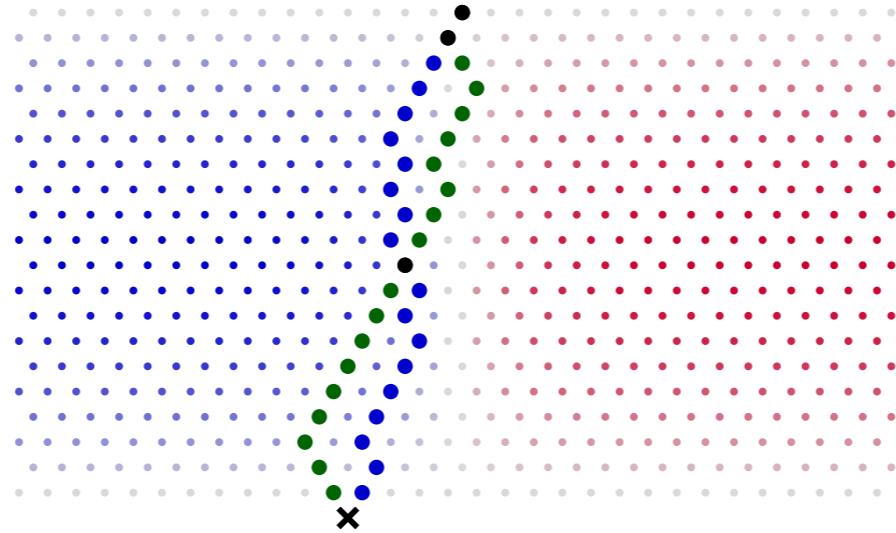
- Histograms of observables
Summary statistics
- Approximate Bayesian Computation
Summary statistics
- Machine Learning techniques
Density networks, CARL, autoregressive models,
normalizing flows, ...
- Matrix Element Method / Optimal Observables
Neglect or approximate shower + detector,
explicitly calculate integral



$$\hat{p}(x|\theta) = \int dz_p \tilde{p}(x|z_p) p(z_p|\theta)$$

Mining gold from the simulator

Back to the Galton board



- Likelihood

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

is usually intractable because of the integral over **all possible paths z**

- But: we can calculate the **probability of each individual path**

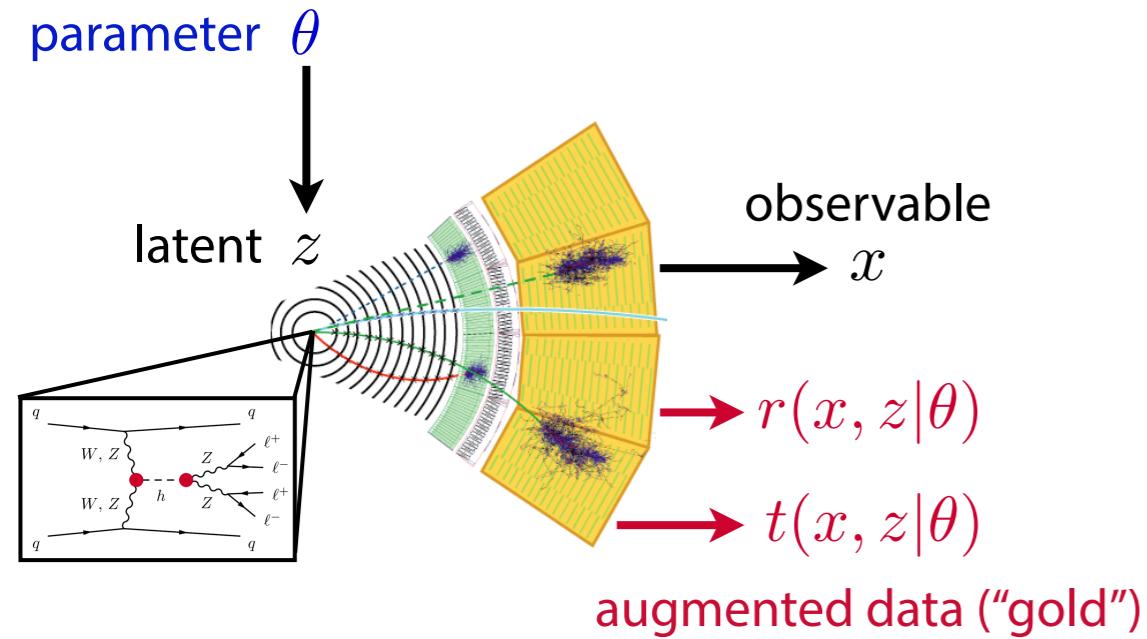
$$p(x, z|\theta) = \prod_{\text{nails } i} p_i(x, z|\theta)$$

Turns out that that can be useful...

A new approach

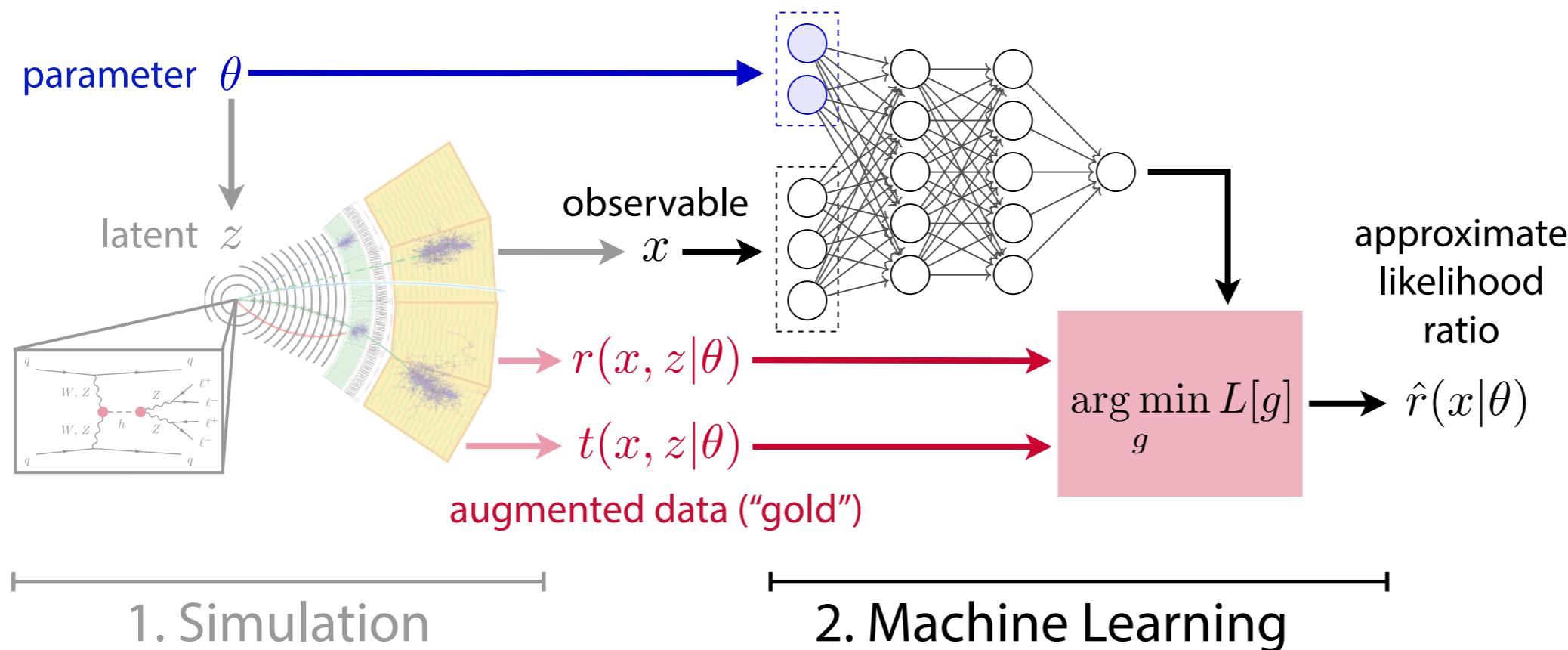
- New algorithms for simulator-based inference
 - Key idea: Extract more information from simulator, use it with ML
 - Applicable to many (but not all) likelihood-free inference problems
- Sales pitch:
 - Higher-fidelity measurements with less training data
 - Scales well to many observables (no need to pick summary statistics) and parameters
 - No approximations in simulator needed
(LHC: supports NLO, parton shower, full detector simulation)
 - Evaluation in microseconds
- Motivated by LHC legacy measurements, but applications far beyond physics

Bird's-eye view



"Mining gold": Extract additional information from simulator

Bird's-eye view



"Mining gold": Extract additional information from simulator

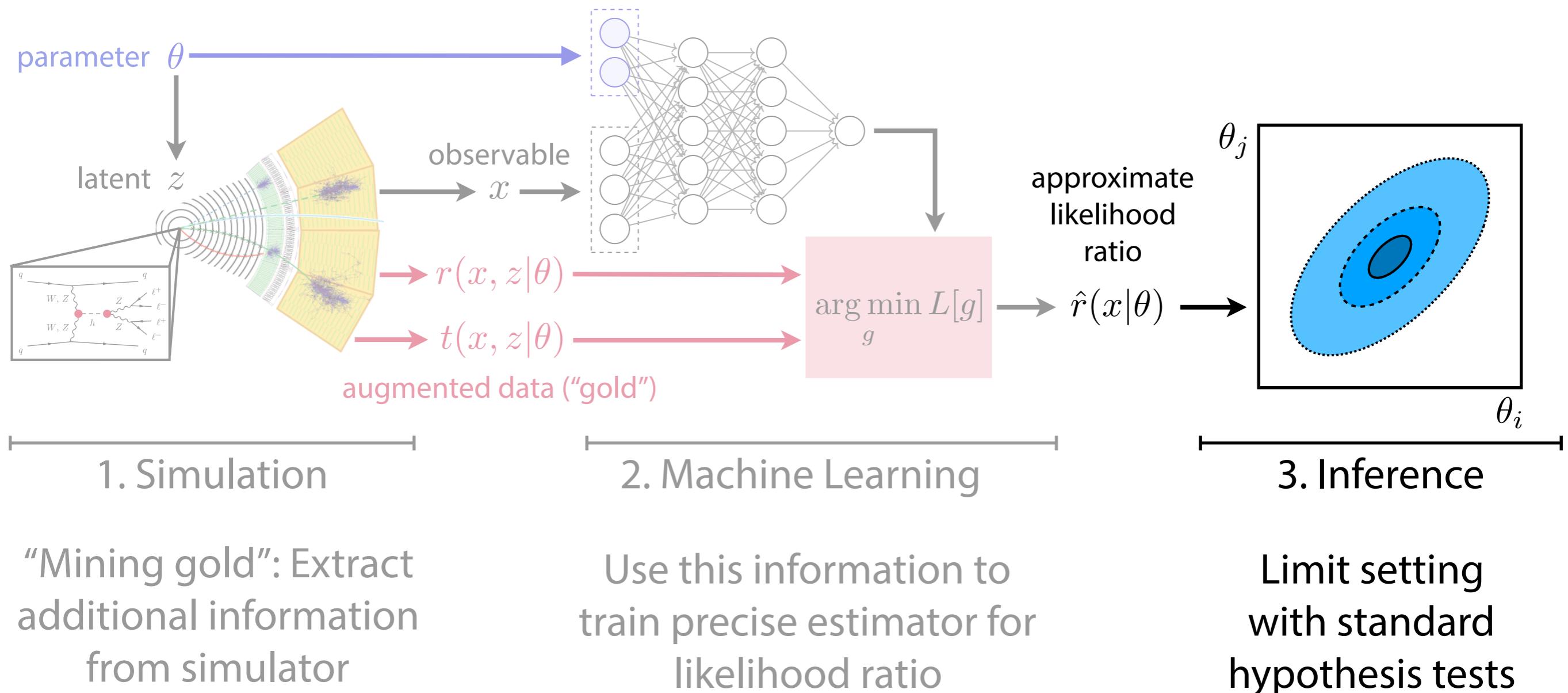
Use this information to train precise estimator for likelihood ratio

$$r(x|\theta_0, \theta_1) \equiv \frac{p(x|\theta_0)}{p(x|\theta_1)}$$

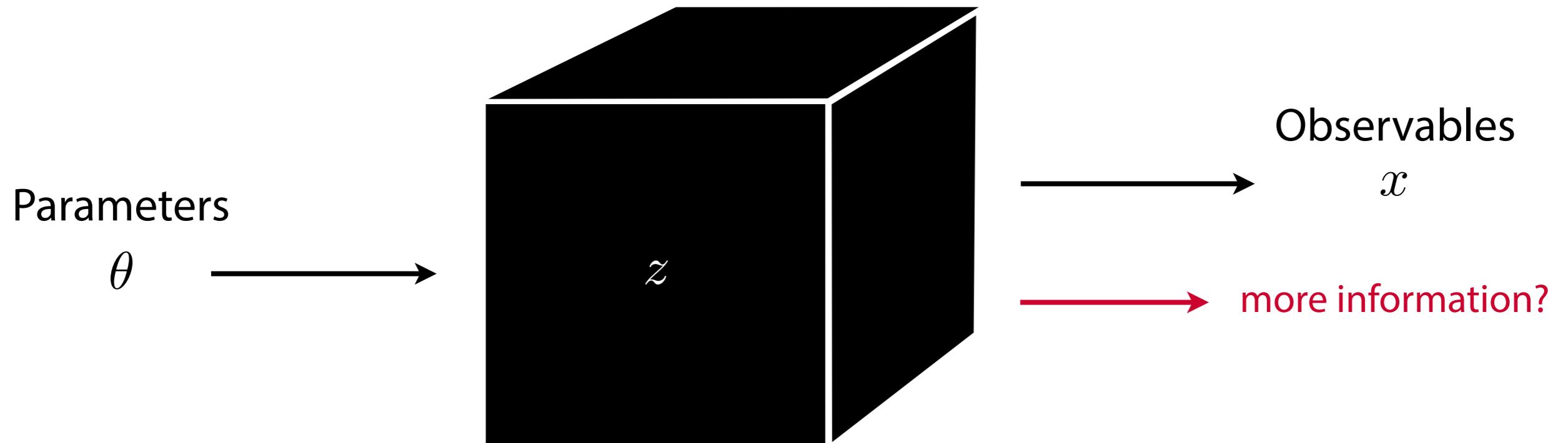
observables

model parameters

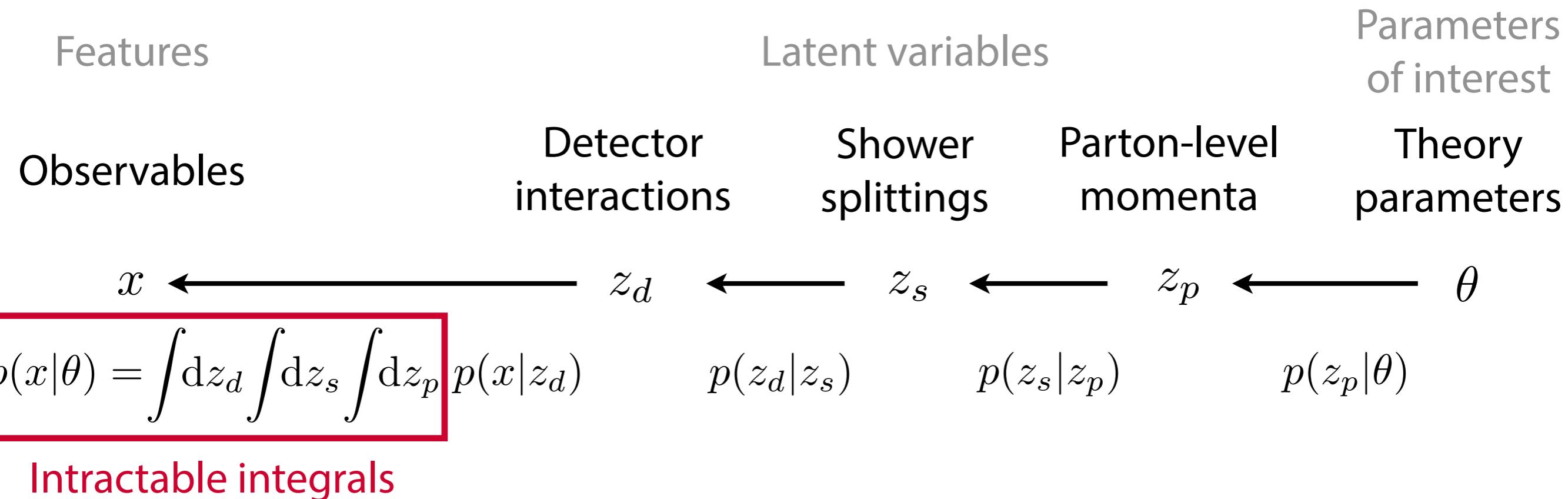
Bird's-eye view



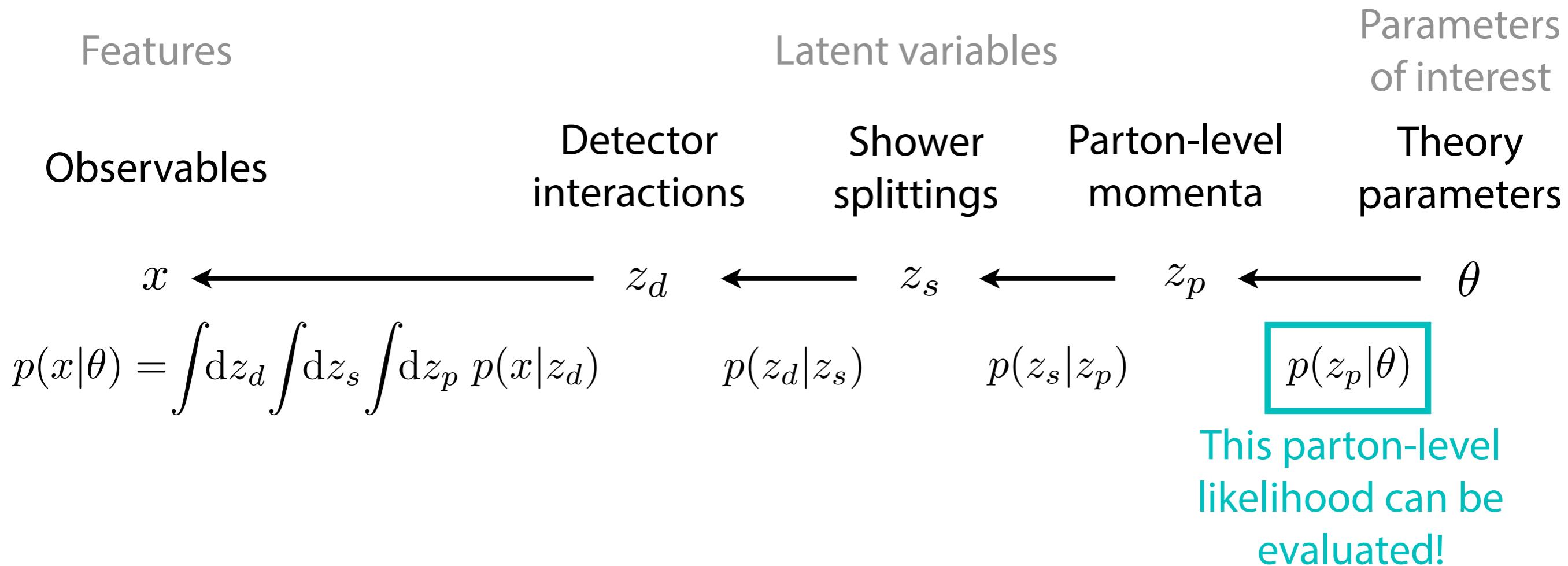
Mining gold from particle physics simulators



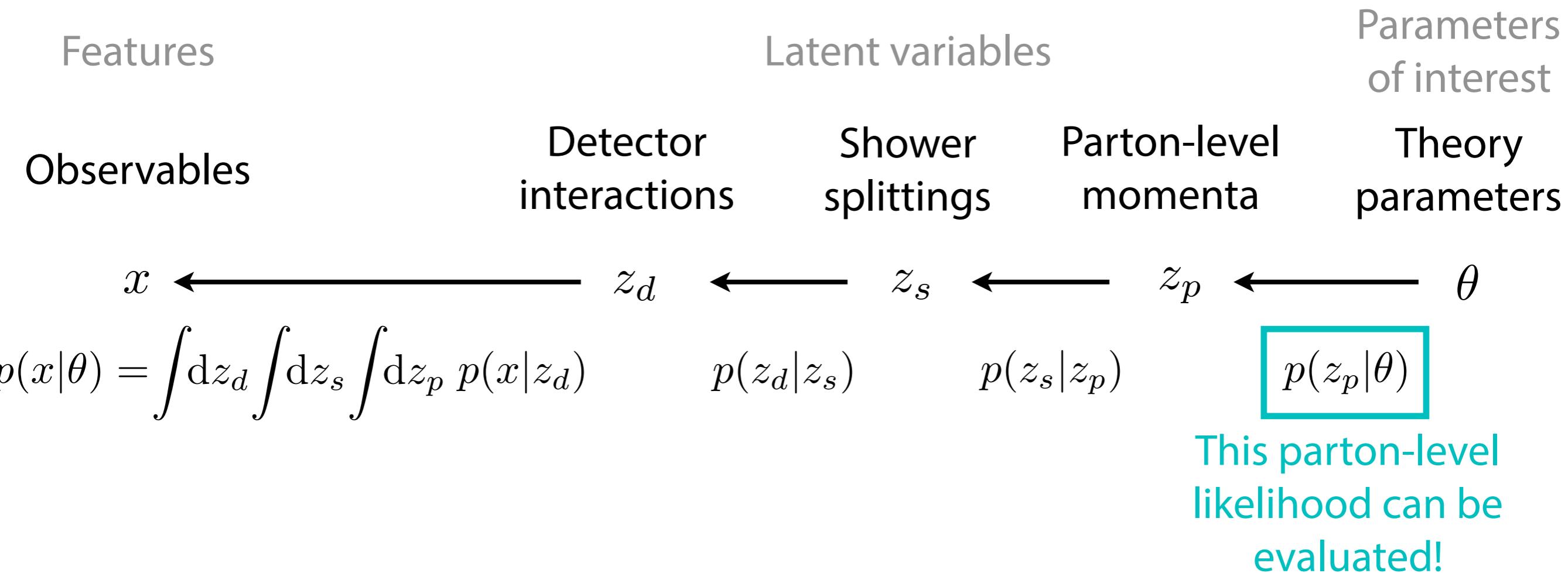
Mining gold from particle physics simulators



Mining gold from particle physics simulators



Mining gold from particle physics simulators



⇒ We can calculate the “joint” likelihood ratio conditional on a specific evolution:

$$r(x, z | \theta_0, \theta_1) \equiv \frac{p(x, z_d, z_s, z_p | \theta_0)}{p(x, z_d, z_s, z_p | \theta_1)} = \frac{p(x|z_d)}{p(x|z_d)} \frac{p(z_d|z_s)}{p(z_d|z_s)} \frac{p(z_s|z_p)}{p(z_s|z_p)}$$

$$\frac{p(z_p|\theta_0)}{p(z_p|\theta_1)}$$

The value of gold

We have **joint likelihood ratio**

$$r(x, z | \theta_0, \theta_1) \equiv \frac{p(x, z_d, z_s, z_p | \theta_0)}{p(x, z_d, z_s, z_p | \theta_1)}$$



We want **likelihood ratio**

$$r(x | \theta_0, \theta_1) \equiv \frac{p(x | \theta_0)}{p(x | \theta_1)}$$

The value of gold

We have joint likelihood ratio

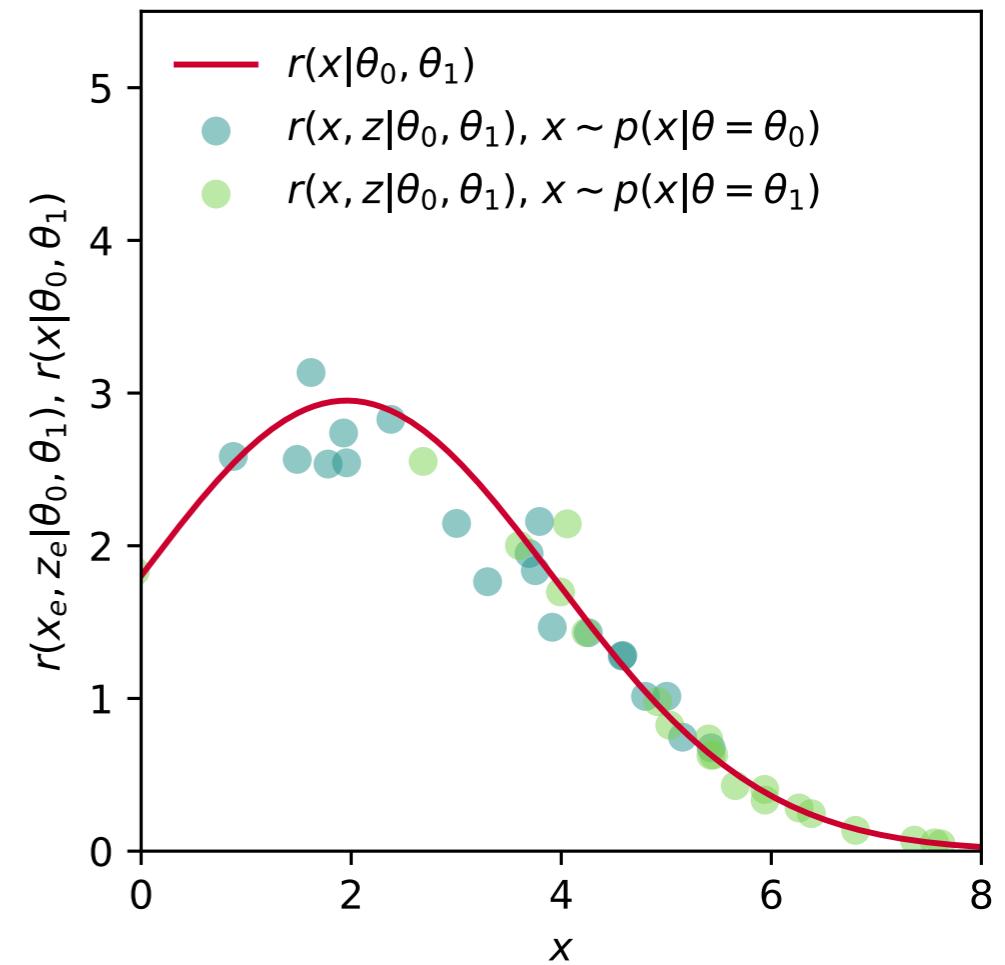
$$r(x, z|\theta_0, \theta_1) \equiv \frac{p(x, z_d, z_s, z_p|\theta_0)}{p(x, z_d, z_s, z_p|\theta_1)}$$



$r(x, z|\theta_0, \theta_1)$ are
scattered around
 $r(x|\theta_0, \theta_1)$

We want likelihood ratio

$$r(x|\theta_0, \theta_1) \equiv \frac{p(x|\theta_0)}{p(x|\theta_1)}$$



The value of gold

We have joint likelihood ratio

$$r(x, z|\theta_0, \theta_1) \equiv \frac{p(x, z_d, z_s, z_p|\theta_0)}{p(x, z_d, z_s, z_p|\theta_1)}$$



With $r(x, z|\theta_0, \theta_1)$,
we define the functional

$$L_r[\hat{r}(x)] = \mathbb{E}_{p(x, z|\theta_1)} \left[\left(\hat{r}(x) - r(x, z|\theta_0, \theta_1) \right)^2 \right].$$

One can show it is minimized by

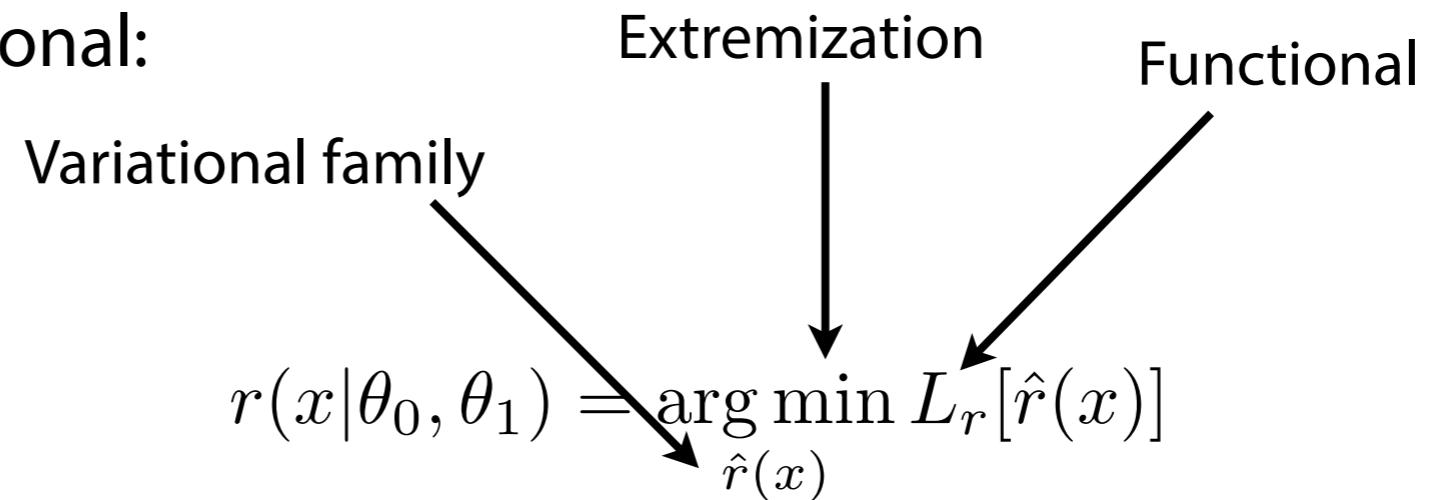
$$\arg \min_{\hat{r}(x)} L_r[\hat{r}(x)] = r(x|\theta_0, \theta_1) !$$

We want likelihood ratio

$$r(x|\theta_0, \theta_1) \equiv \frac{p(x|\theta_0)}{p(x|\theta_1)}$$

How to do this in practice?

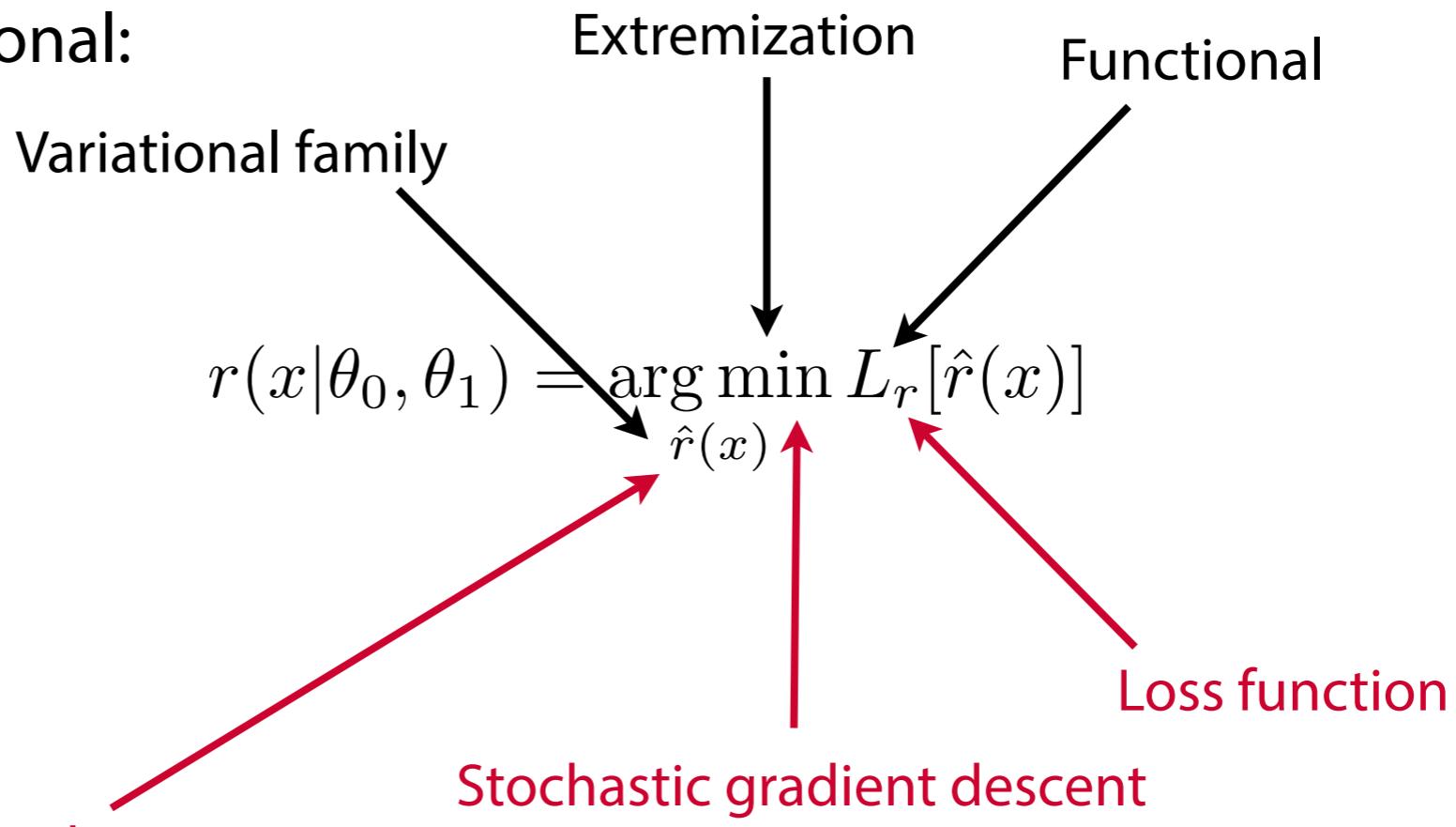
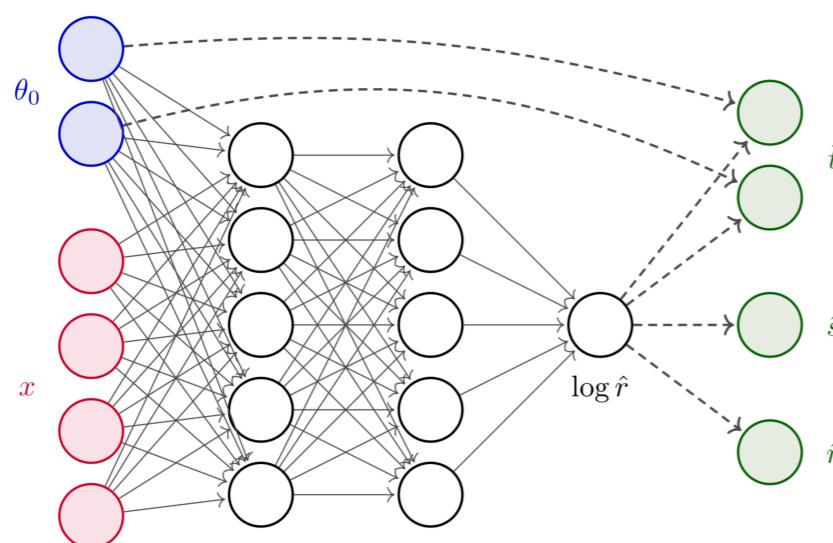
Need to minimize a functional:



How to do this in practice?

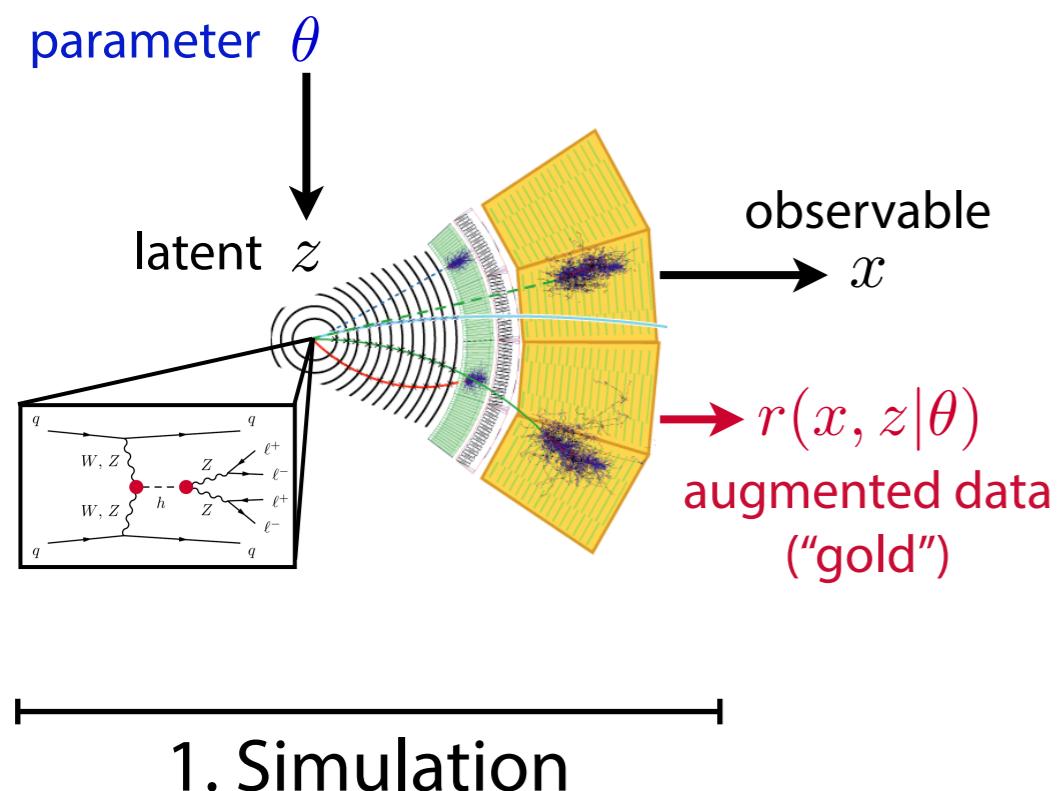
Need to minimize a functional:

This is exactly what
machine learning does!

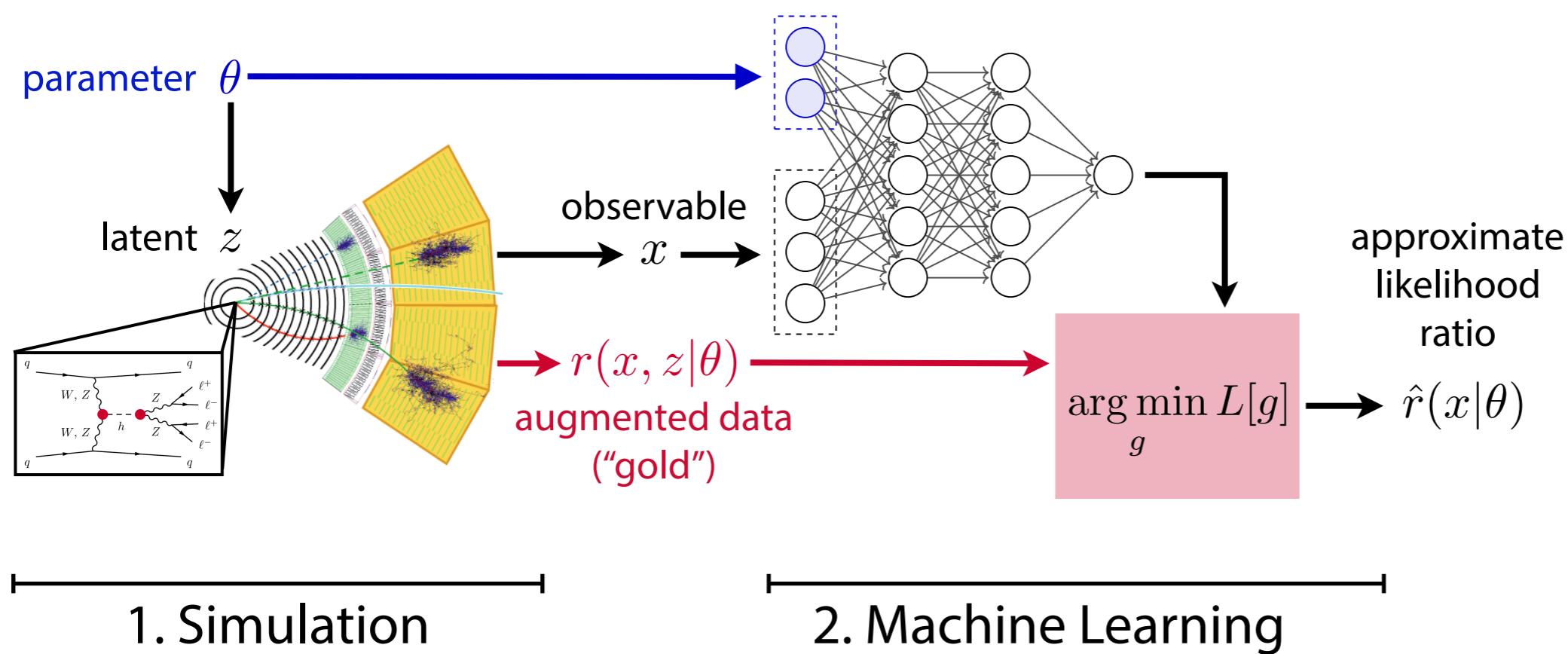


\Rightarrow We implement $\hat{r}(x|\theta_0, \theta_1)$
as a neural network trained on the
data available from the simulator

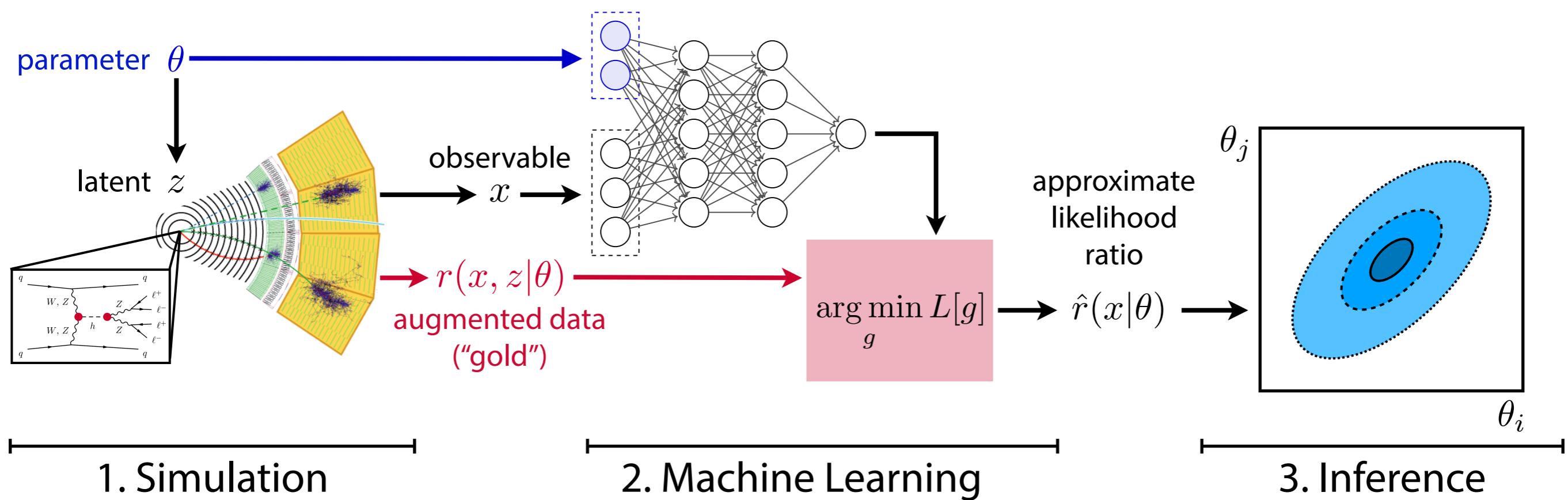
What we have so far



What we have so far



What we have so far



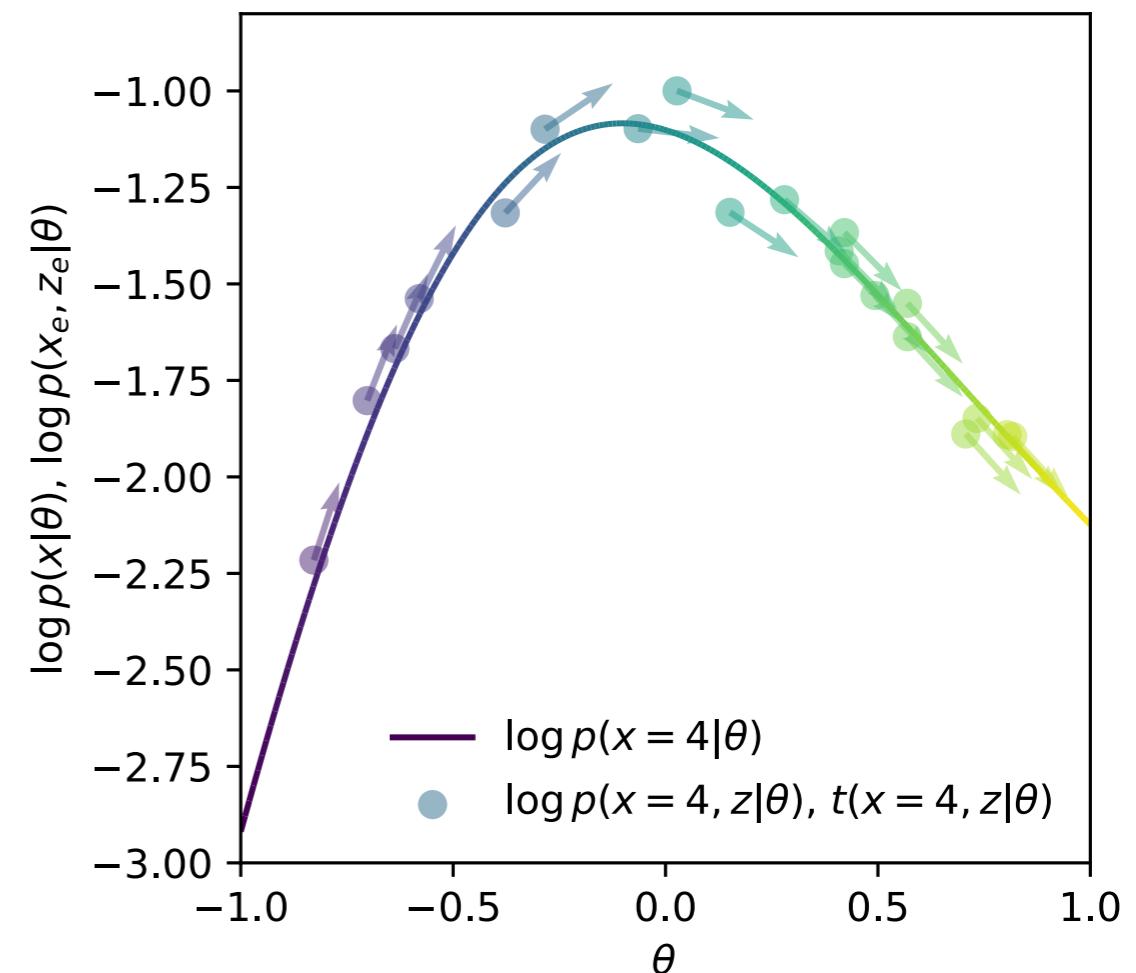
The score

- Inference just based on the joint likelihood ratio works well, but there is another powerful piece of information
- The **score**

$$t(x|\theta_0) \equiv \nabla_{\theta} \log p(x|\theta) \Big|_{\theta_0}$$

fully characterizes the likelihood function in the neighborhood of θ_0

- The score itself is intractable. But...



Learning the score

Similar to the joint likelihood ratio,
we can calculate the **joint score**

$$t(x, z|\theta_0) \equiv \nabla_{\theta} \log p(x, z_d, z_s, z|\theta) \Big|_{\theta_0}$$



We want **score**

$$t(x|\theta_0) \equiv \nabla_{\theta} \log p(x|\theta) \Big|_{\theta_0}$$

Learning the score

Similar to the joint likelihood ratio,
we can calculate the **joint score**

$$t(x, z|\theta_0) \equiv \nabla_{\theta} \log p(x, z_d, z_s, z|\theta) \Big|_{\theta_0}$$



We want **score**

$$t(x|\theta_0) \equiv \nabla_{\theta} \log p(x|\theta) \Big|_{\theta_0}$$

Given $t(x, z|\theta_0)$,
we define the functional

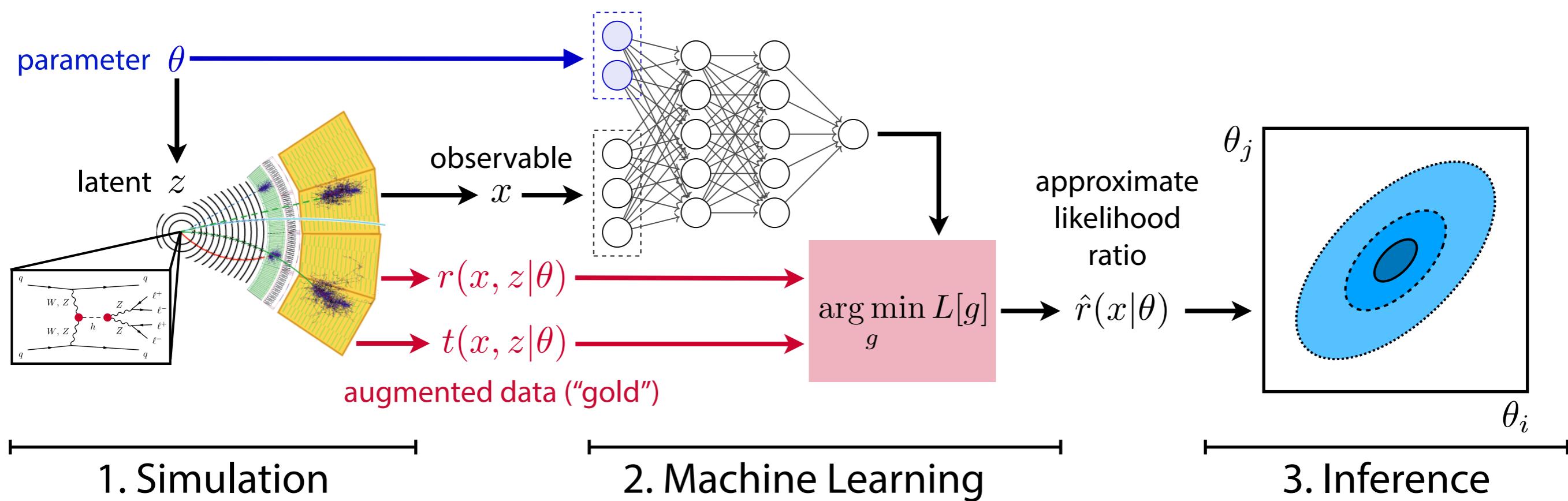
$$L_t[\hat{t}(x)] = \mathbb{E}_{p(x, z|\theta_0)} \left[\left(\hat{t}(x) - t(x, z|\theta_0) \right)^2 \right].$$

One can show it is minimized by

$$\arg \min_{\hat{t}(x)} L_t[\hat{t}(x)] = t(x|\theta_0).$$

Again, we implement this
with machine learning

Putting the pieces together



A family of likelihood-free inference strategies

Different strategies to combine the different pieces of information:

Method	L_{XE}	L_{MLE}	L_r	L_t	θ sampling
ABC (Approximate Bayesian Computation)					$\theta \sim \pi(\theta)$
NDE (Neural density estimation)			✓		$\theta \sim \pi(\theta)$
LRT / CARL (Likelihood ratio trick / calibrated approximate ratios of likelihoods)	✓				$\theta \sim \pi(\theta)$
ROLR (Regression on likelihood ratio)				✓	$\theta \sim \pi(\theta)$
SCANDAL (Score augmented neural density approximates likelihood)		✓		✓	$\theta \sim \pi(\theta)$
CASCAL (CARL and score approximate likelihood ratio).	✓			✓	$\theta \sim \pi(\theta)$
RASCAL (Ratio and score approximate likelihood ratio)			✓	✓	$\theta \sim \pi(\theta)$
SALLY (Score approximates likelihood locally)				✓	$\theta = \theta_0$
SALLINO (Score approximates likelihood locally in one dimension)				✓	$\theta = \theta_0$

A family of likelihood-free inference strategies

Different strategies to combine the different pieces of information:

Method	L_{XE}	L_{MLE}	L_r	L_t	θ sampling
ABC (Approximate Bayesian Computation)					$\theta \sim \pi(\theta)$
NDE (Neural density estimation)			✓		$\theta \sim \pi(\theta)$
LRT / CARL (Likelihood ratio trick / calibrated approximate ratios of likelihoods)	✓				$\theta \sim \pi(\theta)$
ROLR (Regression on likelihood ratio)				✓	$\theta \sim \pi(\theta)$
SCANDAL (Score augmented neural density approximates likelihood)		✓		✓	$\theta \sim \pi(\theta)$
CASCAL (CARL and score approximate likelihood ratio).	✓			✓	$\theta \sim \pi(\theta)$
RASCAL (Ratio and score approximate likelihood ratio)			✓	✓	$\theta \sim \pi(\theta)$
SALLY (Score approximates likelihood locally)				✓	$\theta = \theta_0$
SALLINO (Score approximates likelihood locally in one dimension)				✓	$\theta = \theta_0$

RASCAL loss function:

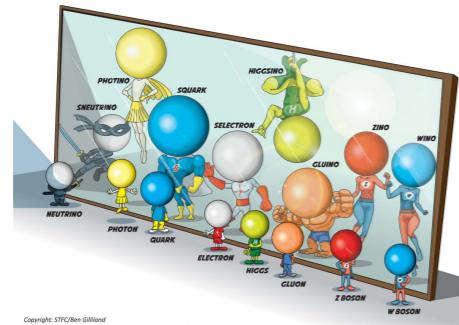
$$L_{\text{RASCAL}}[\hat{r}(x|\theta_0, \theta_1)] = L_r[\hat{r}(x|\theta_0, \theta_1)] + \alpha L_t[\nabla_{\theta_0} \log \hat{r}(x|\theta_0, \theta_1)]$$

Proof of concept: Constraining EFTs at the LHC

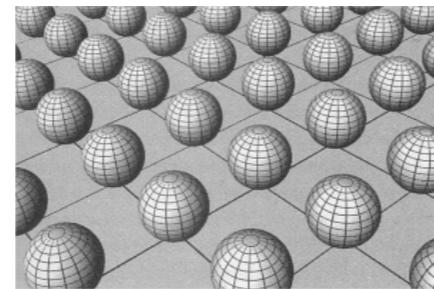
Effective field theory

Energy

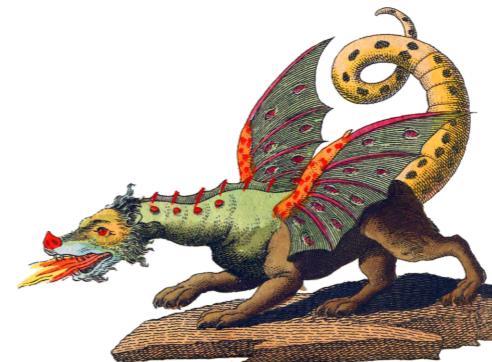
Λ



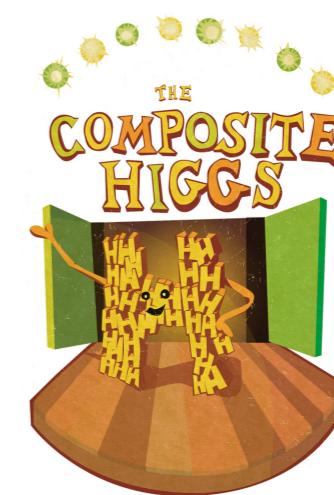
Supersymmetry?



Extra dimensions?



Dragons?



Composite Higgs?

...

$E \ll \Lambda$

LHC Higgs

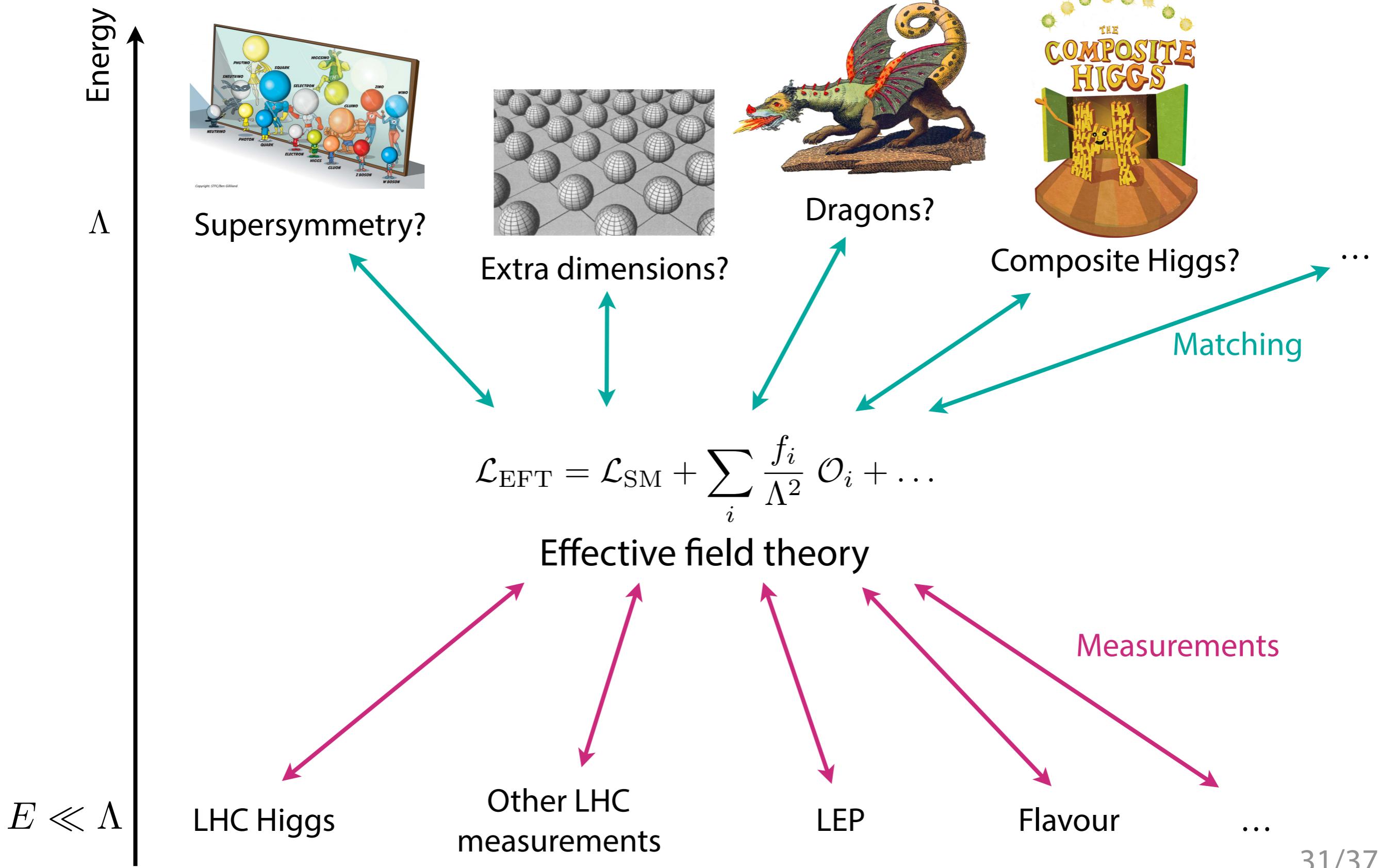
Other LHC
measurements

LEP

Flavour

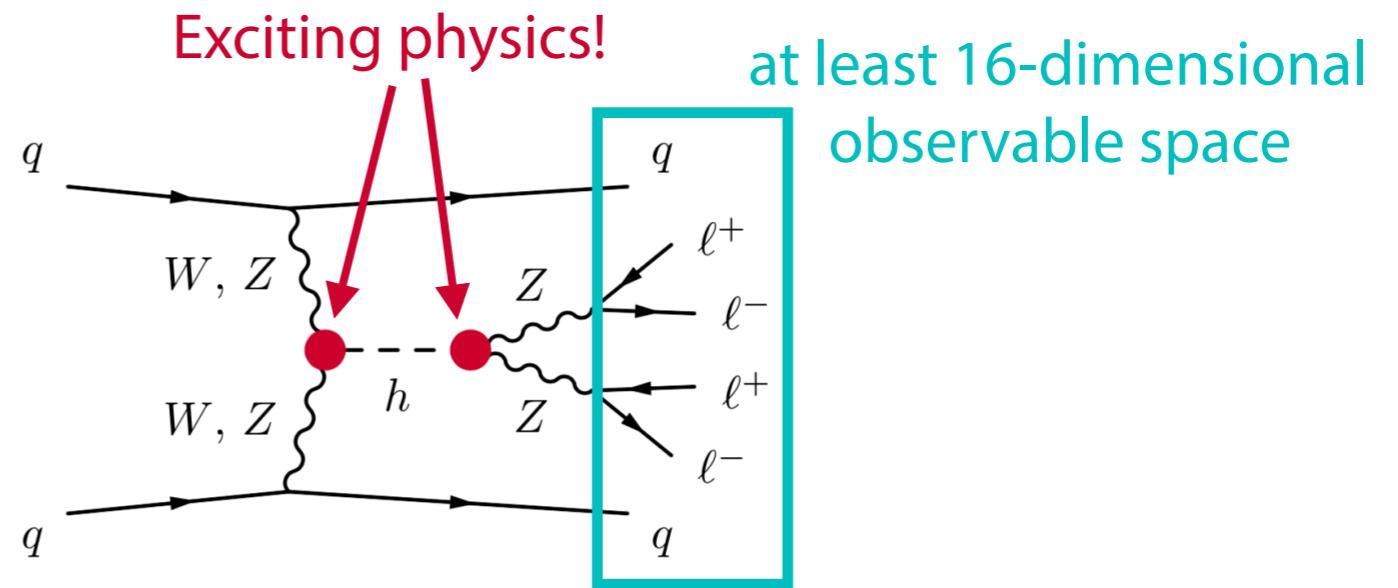
...

Effective field theory



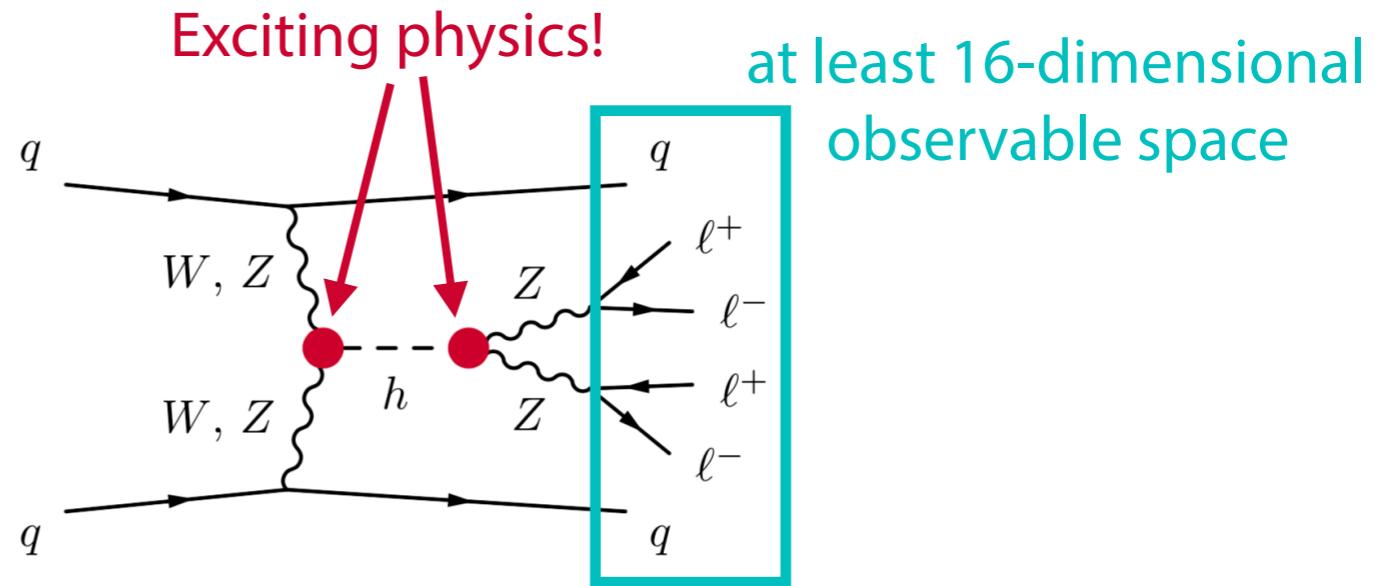
Proof of concept

- Higgs production in weak boson fusion:



Proof of concept

- Higgs production in weak boson fusion:

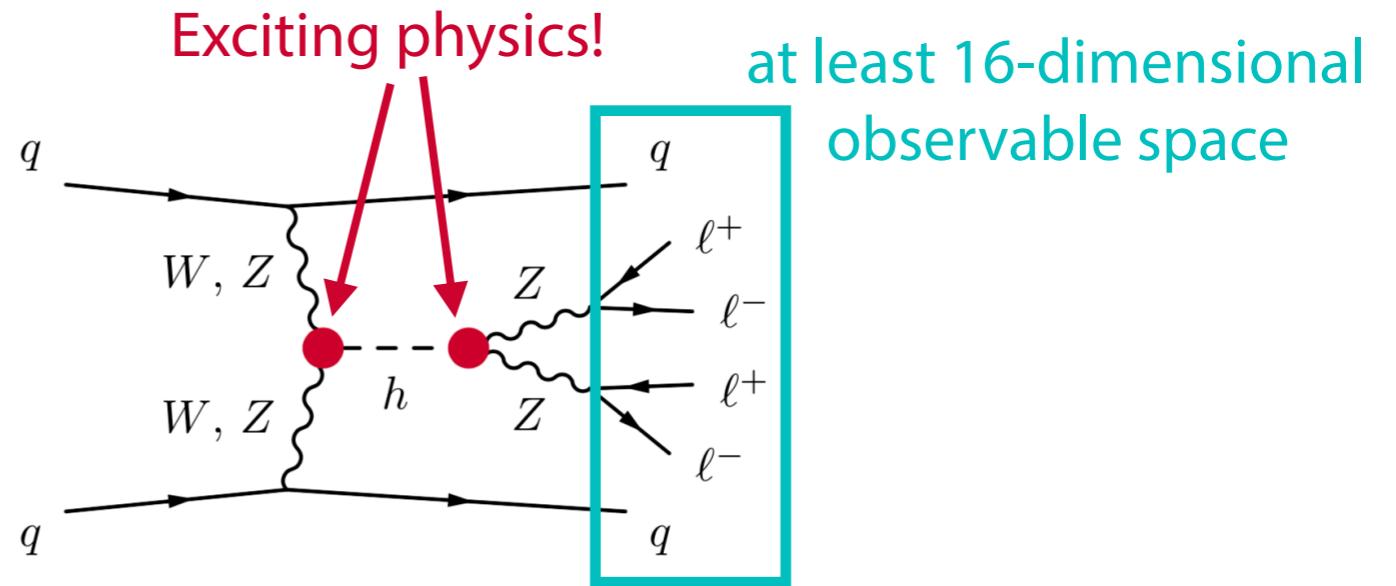


- Goal: constraints on two EFT parameters

$$\mathcal{L} = \mathcal{L}_{\text{SM}} + \boxed{\frac{f_W}{\Lambda^2}} \underbrace{\frac{ig}{2} (D^\mu \phi)^\dagger \sigma^a D^\nu \phi W_{\mu\nu}^a}_{\mathcal{O}_W} - \boxed{\frac{f_{WW}}{\Lambda^2}} \underbrace{\frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a W^{\mu\nu a}}_{\mathcal{O}_{WW}}$$

Proof of concept

- Higgs production in weak boson fusion:



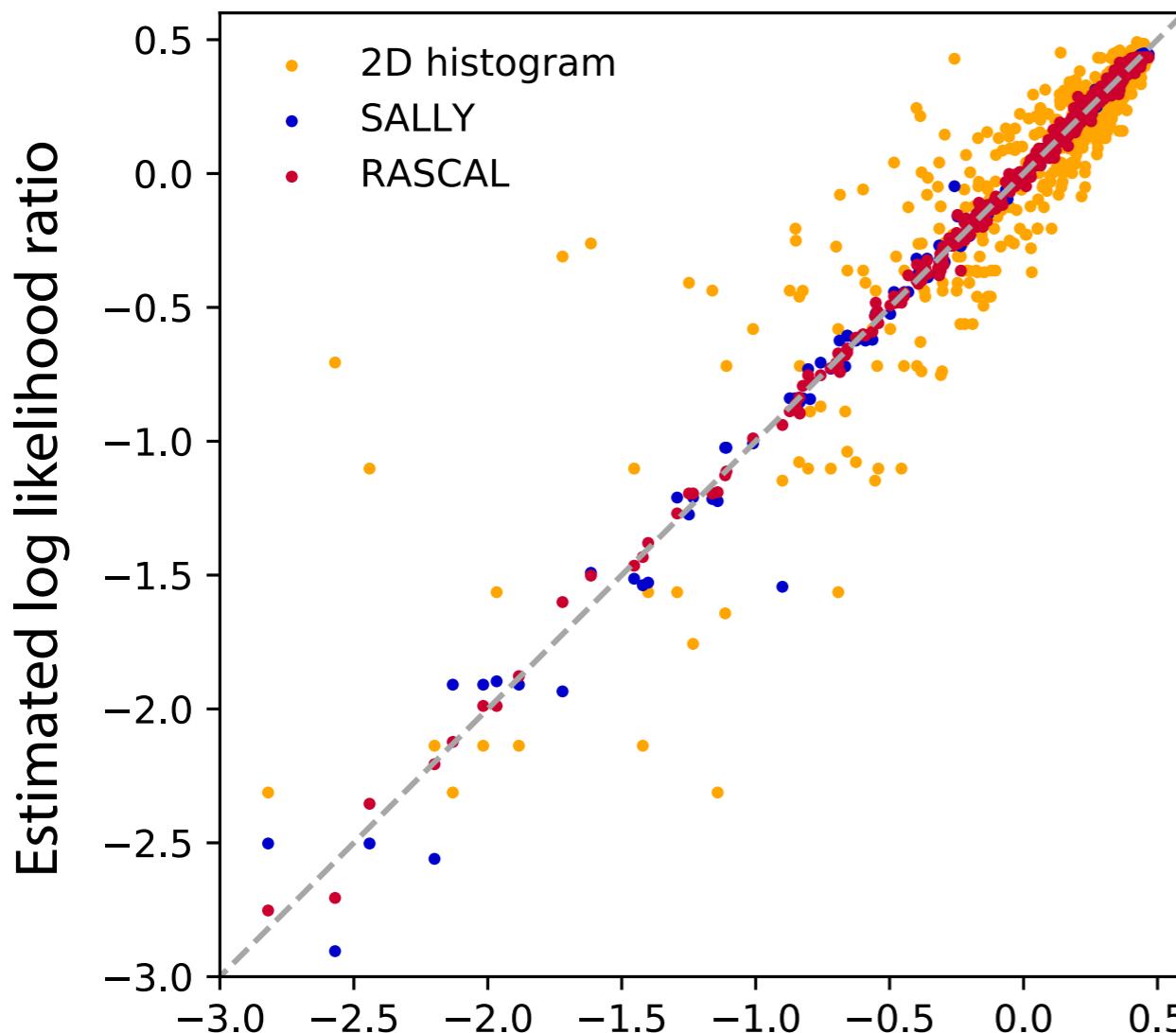
- Goal: constraints on two EFT parameters

$$\mathcal{L} = \mathcal{L}_{\text{SM}} + \boxed{\frac{f_W}{\Lambda^2}} \underbrace{\frac{ig}{2} (D^\mu \phi)^\dagger \sigma^a D^\nu \phi W_{\mu\nu}^a}_{\mathcal{O}_W} - \boxed{\frac{f_{WW}}{\Lambda^2}} \underbrace{\frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a W^{\mu\nu a}}_{\mathcal{O}_{WW}}$$

- Two setups:
 - Simplified setup in which we can compare to true likelihood
 - “Realistic” simulation with approximate detector effects
- Simulation: MadGraph [J. Alwall et al. 1405.0301] + MadMax [K. Cranmer, T. Plehn hep-ph/0605268; T. Plehn, P. Schichtel, D. Wiegand 1311.2591]

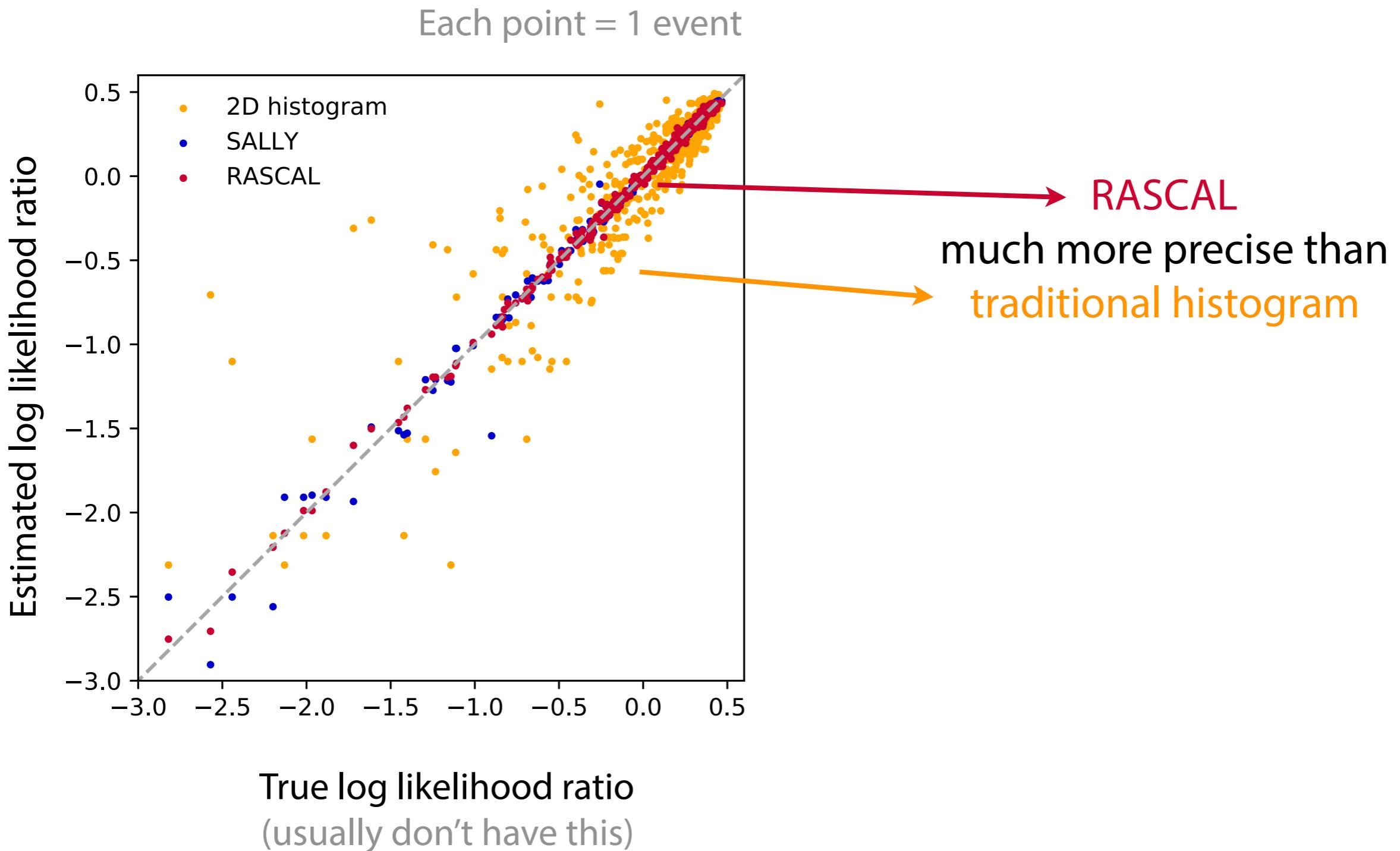
Precise likelihood ratio estimates

Each point = 1 event

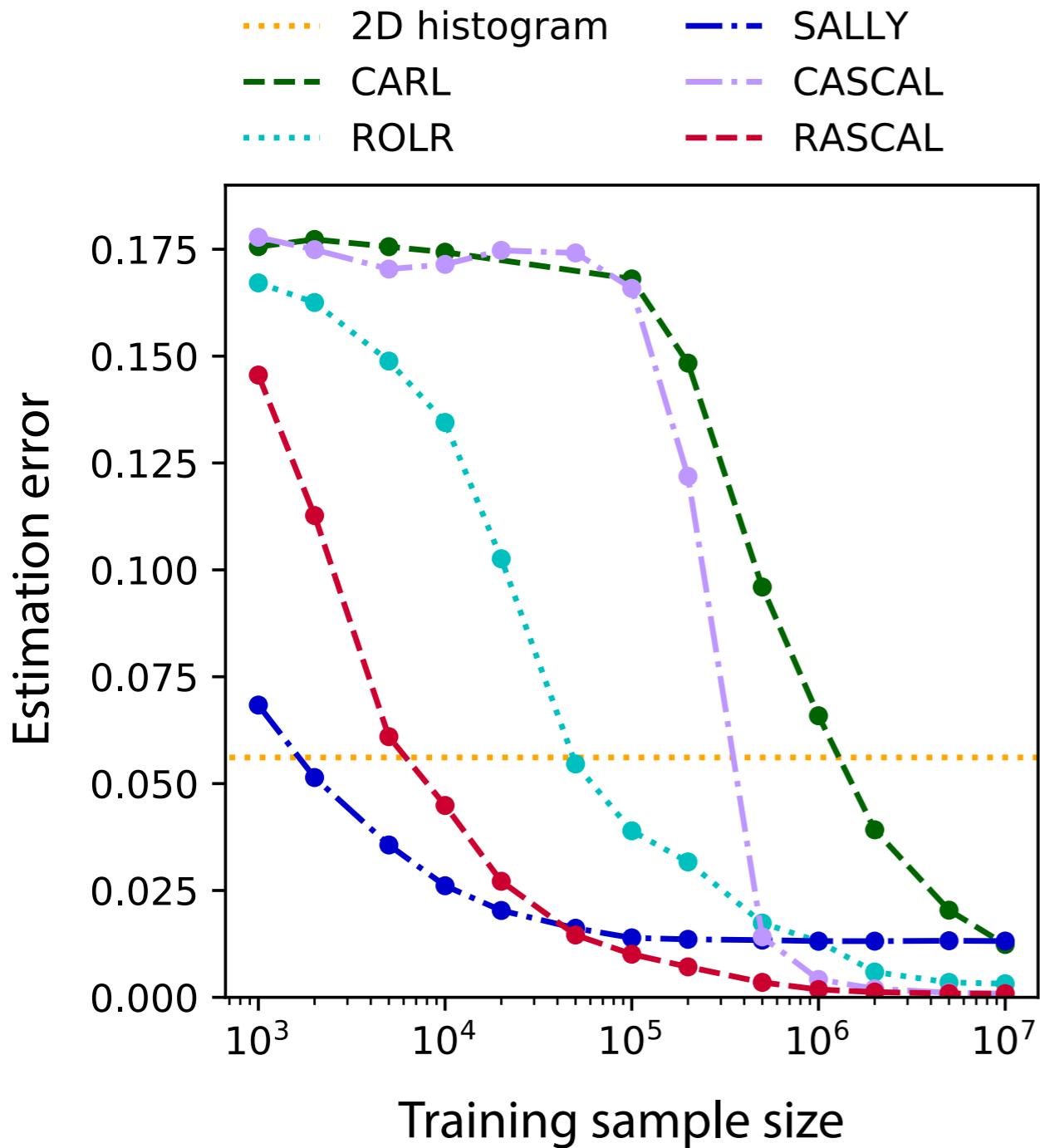


True log likelihood ratio
(usually don't have this)

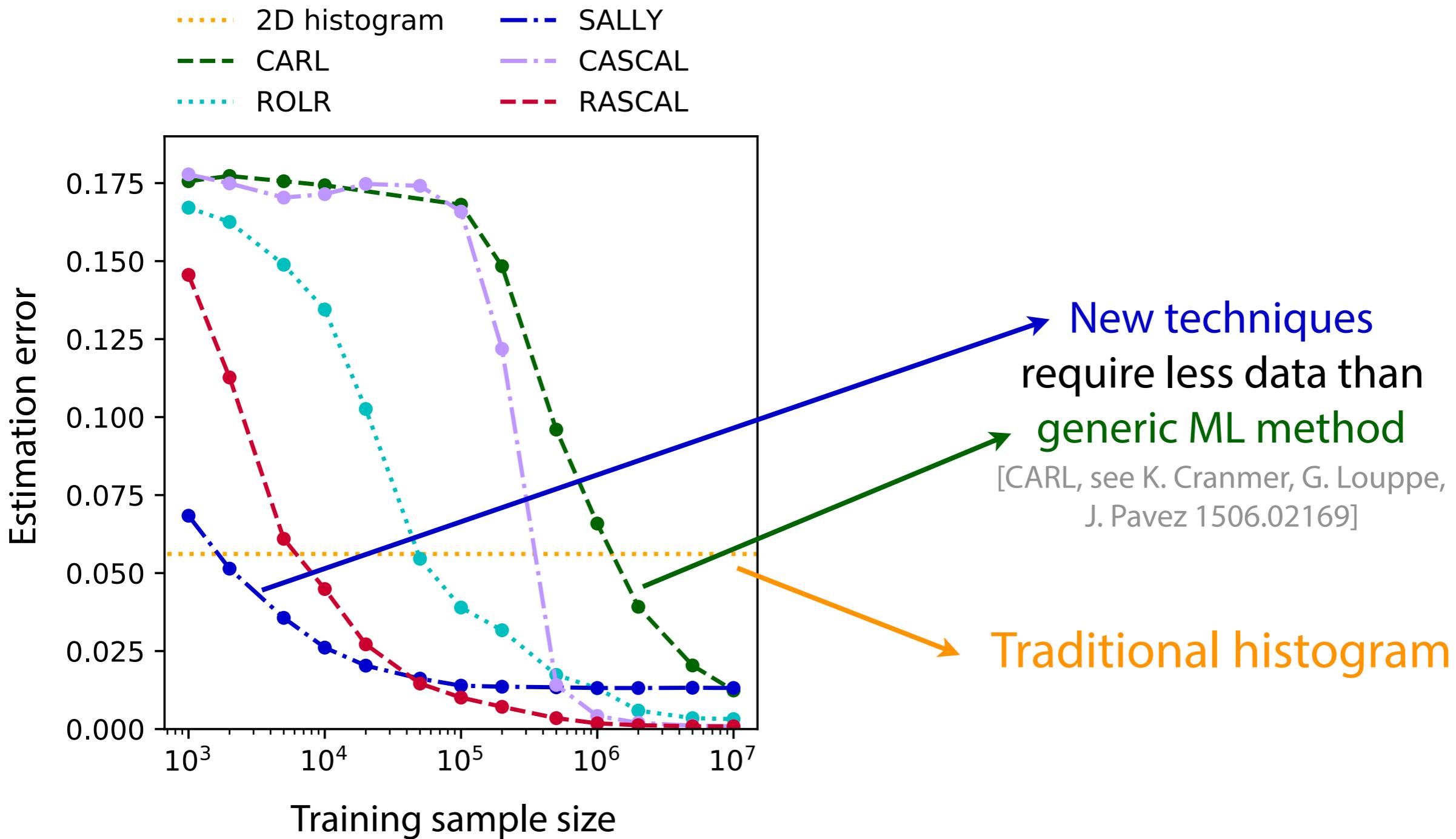
Precise likelihood ratio estimates



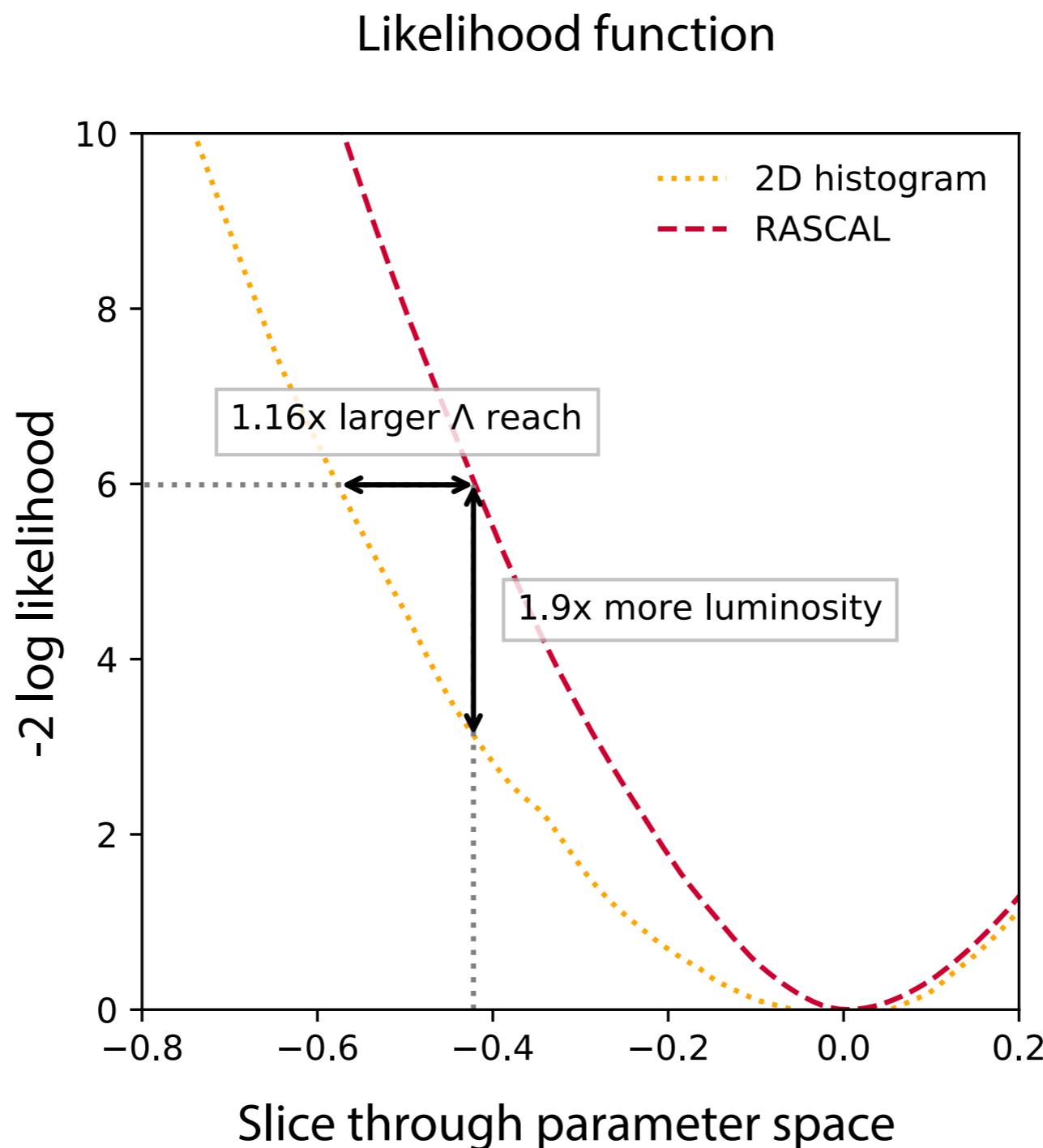
Less training data needed



Less training data needed



Better sensitivity

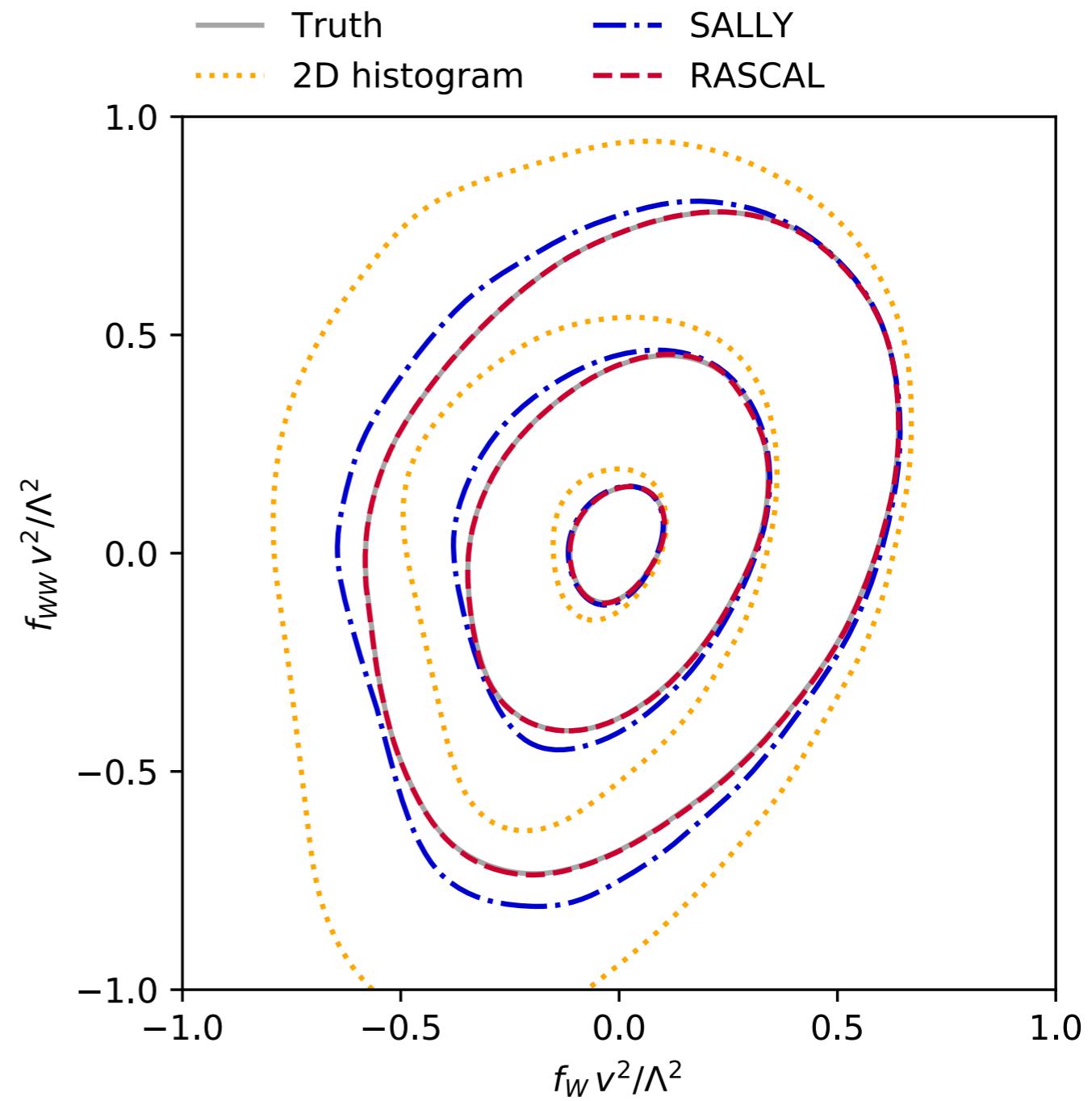


36 events, assuming SM

35/37

Stronger bounds

Expected exclusion limits at 68%, 95%, 99.7% CL



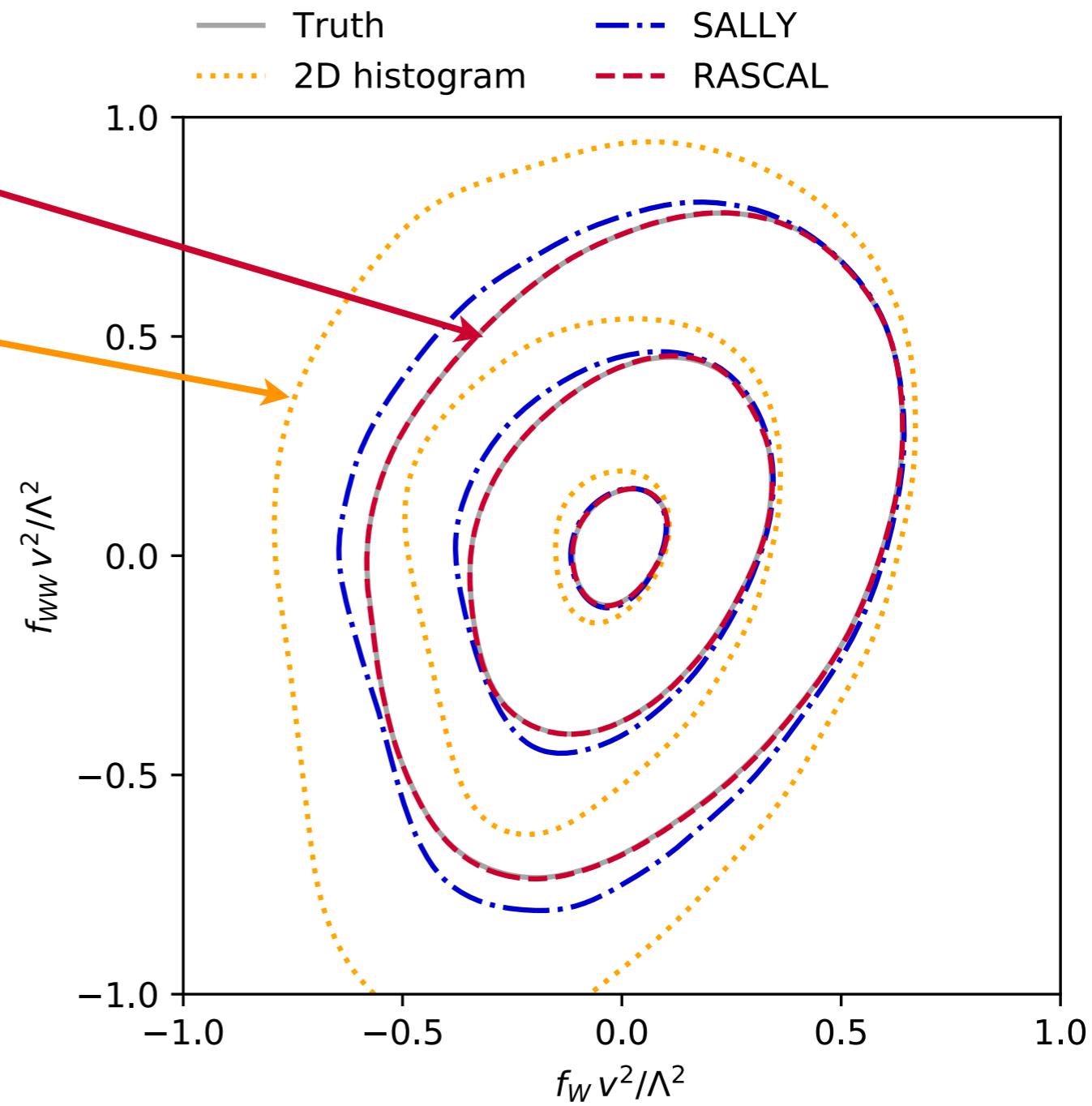
36 events, assuming SM

36/37

Stronger bounds

RASCAL
enables stronger
limits than
traditional histogram

Expected exclusion limits at 68%, 95%, 99.7% CL



36 events, assuming SM

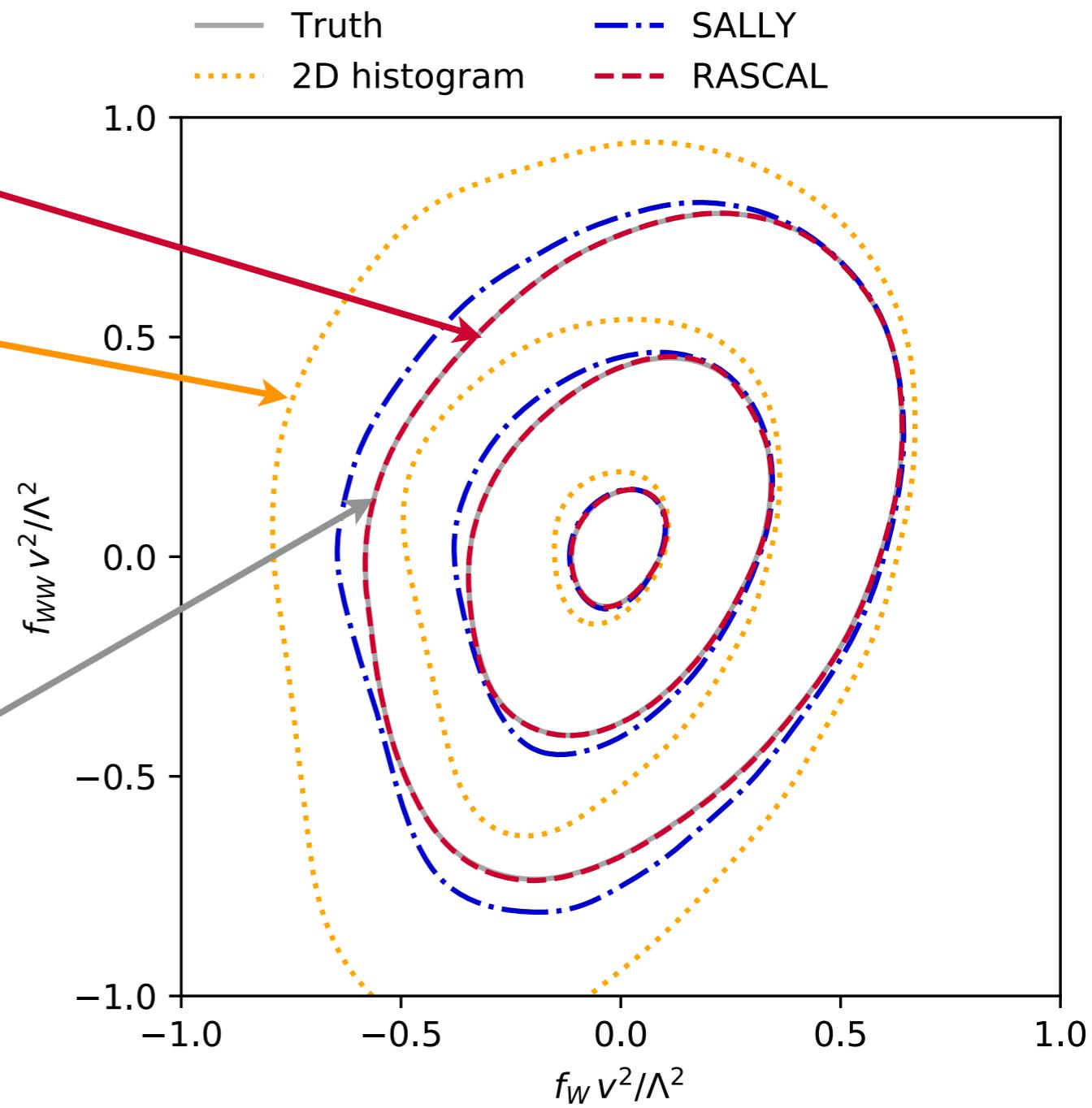
36/37

Stronger bounds

RASCAL
enables stronger
limits than
traditional histogram

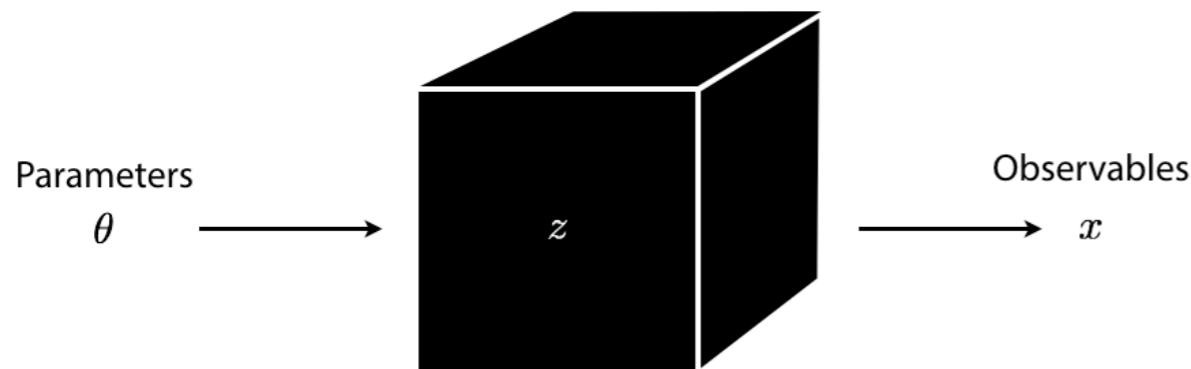
Limits from RASCAL
virtually indistinguishable
from true likelihood
(usually we don't have that)

Expected exclusion limits at 68%, 95%, 99.7% CL



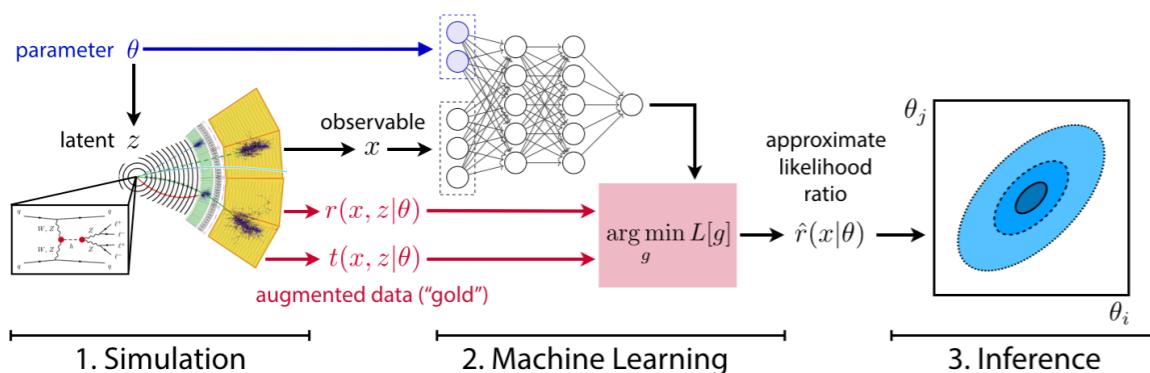
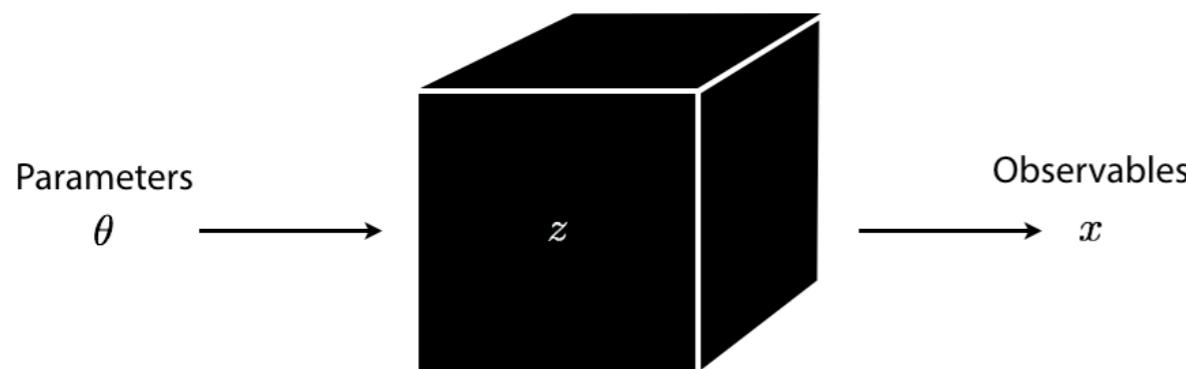
36 events, assuming SM

A new approach to simulator-based inference



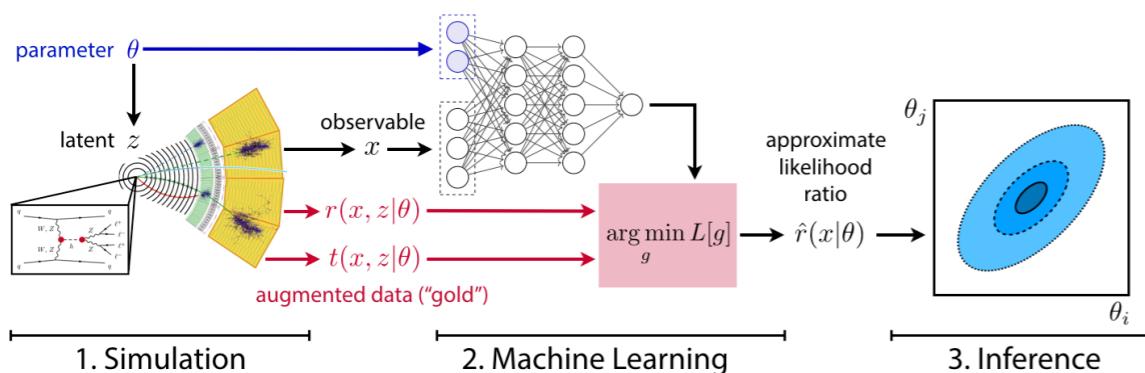
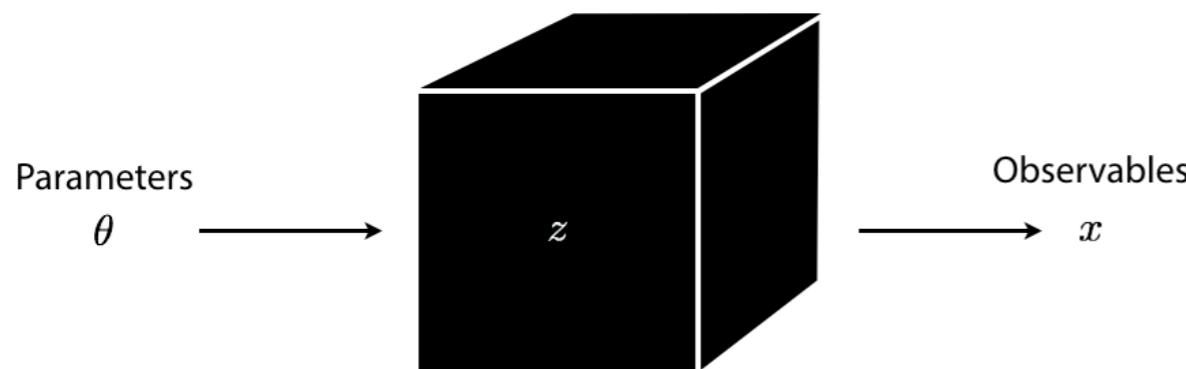
- Much of modern science is based on simulations, “likelihood-free”
- Established inference methods treat simulator as black box

A new approach to simulator-based inference



- Much of modern science is based on simulations, “likelihood-free”
- Established inference methods treat simulator as black box
- New inference techniques:
Leverage more information from simulator + power of machine learning

A new approach to simulator-based inference



- Much of modern science is based on simulations, “likelihood-free”
- Established inference methods treat simulator as black box
- New inference techniques:
Leverage more information from simulator + power of machine learning
- First application to LHC physics:
New methods allow for stronger EFT constraints with less data

