

# INFERNO: Inference-Aware Neural Optimisation

---

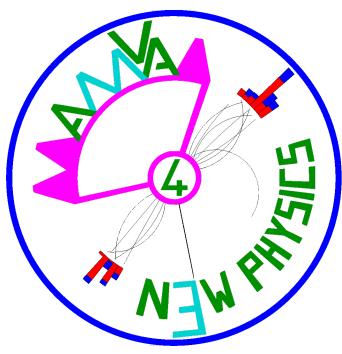
*Pablo de Castro* ( [@pablodecm](https://twitter.com/pablodecm)) and *Tommaso Dorigo* ( [@dorigo](https://twitter.com/dorigo))

19th March 2018 @ Likelihood-Free Inference Workshop (Flatiron Institute, NYC)

More details available on our preprint [arxiv:1806.04743](https://arxiv.org/abs/1806.04743)  
as well as our [GitHub code repository](https://github.com/pdecastro/inferno)

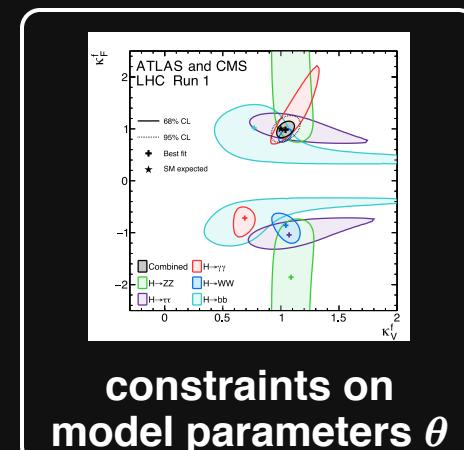
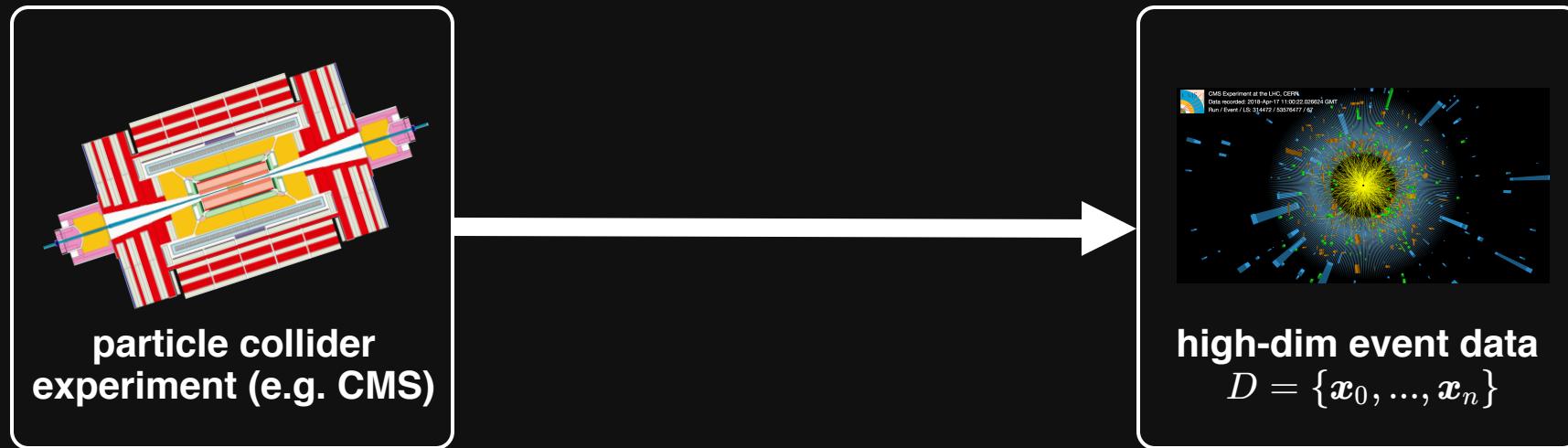
online version of this slides available at

[https://pdecastr.web.cern.ch/pdecastr/lfi\\_flatiron\\_march\\_2019](https://pdecastr.web.cern.ch/pdecastr/lfi_flatiron_march_2019)

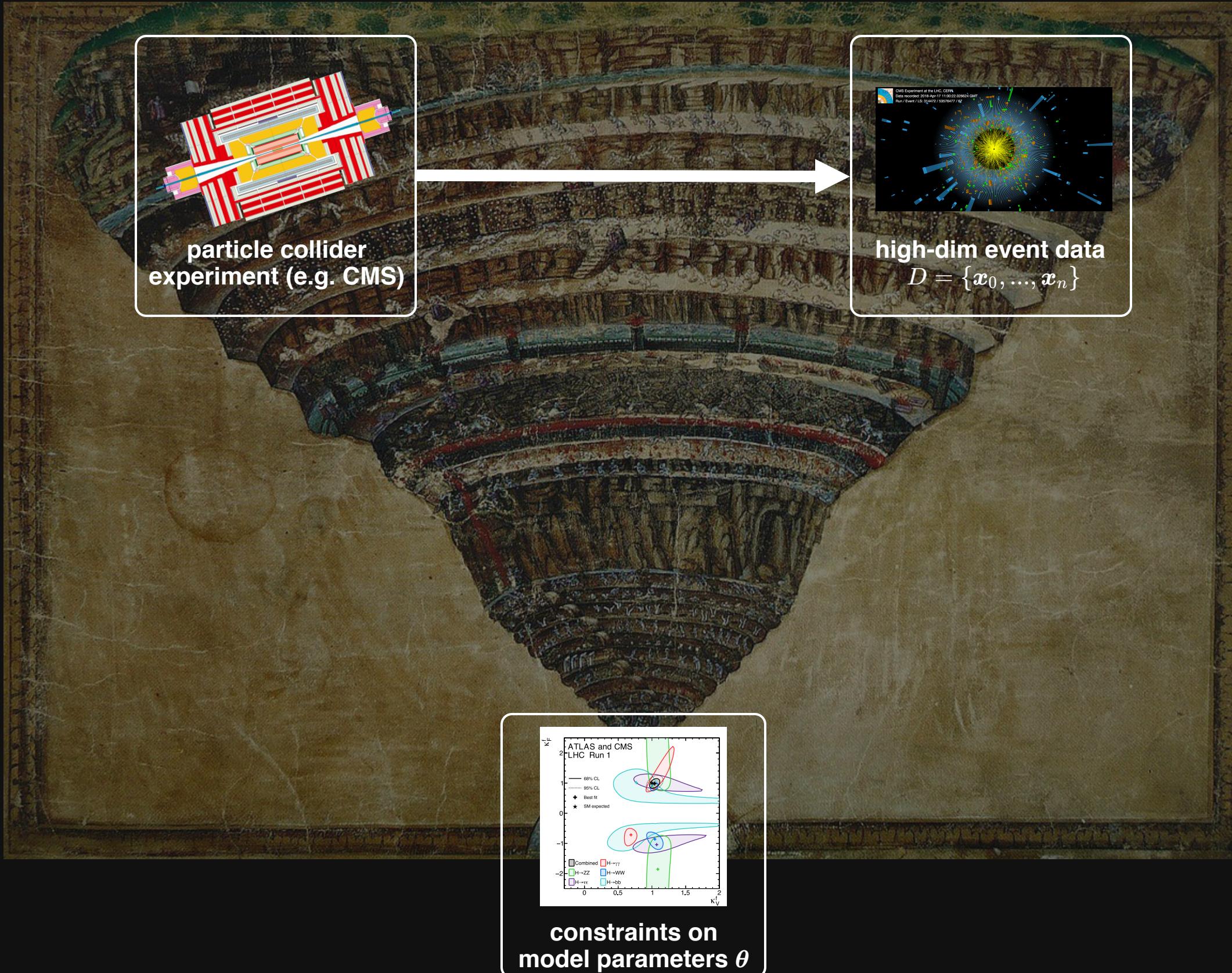


AMVA4NewPhysics has received funding from European Union's  
Horizon 2020 Programme under Grant Agreement number 675440

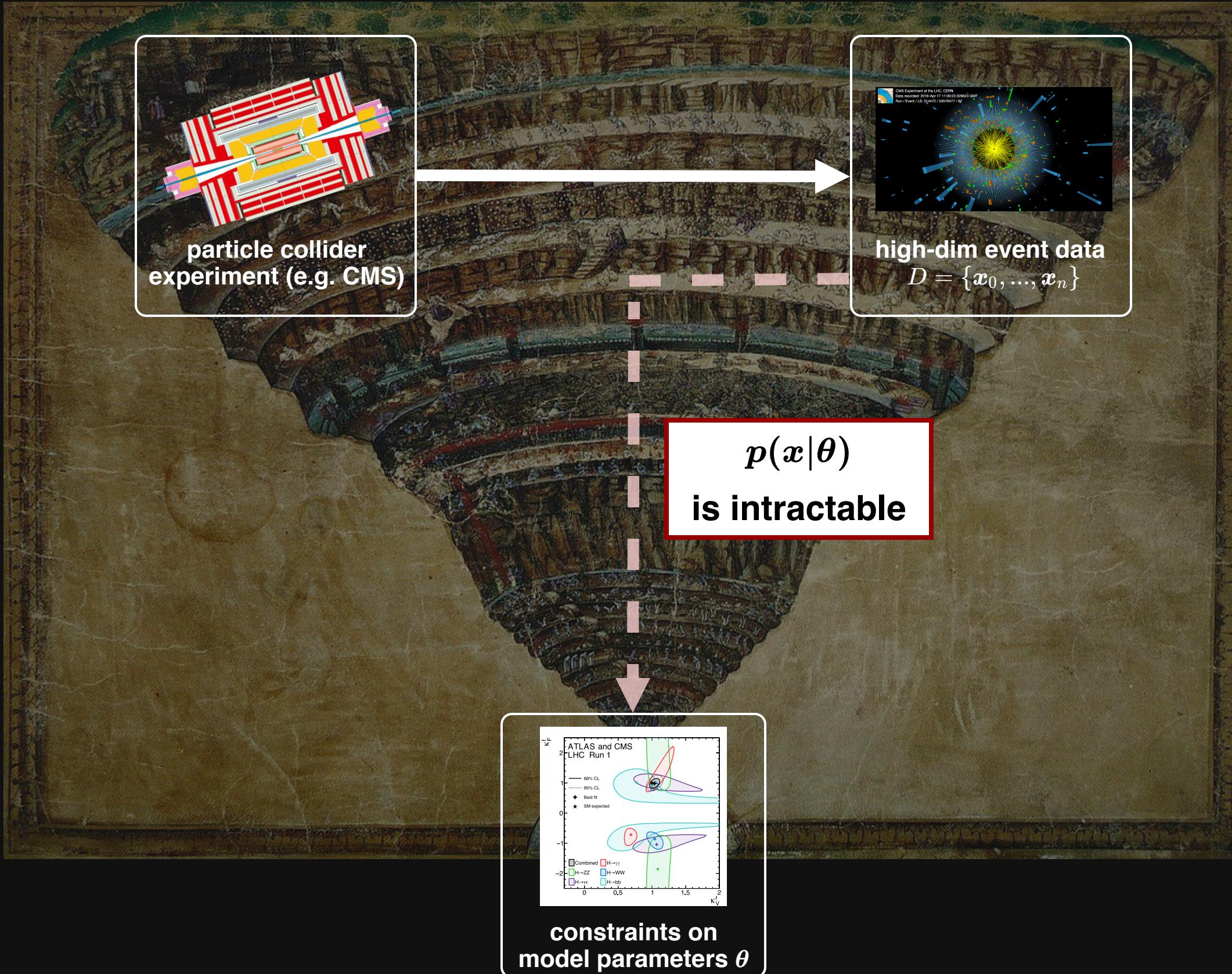
# Statistical Inference in Particle Colliders



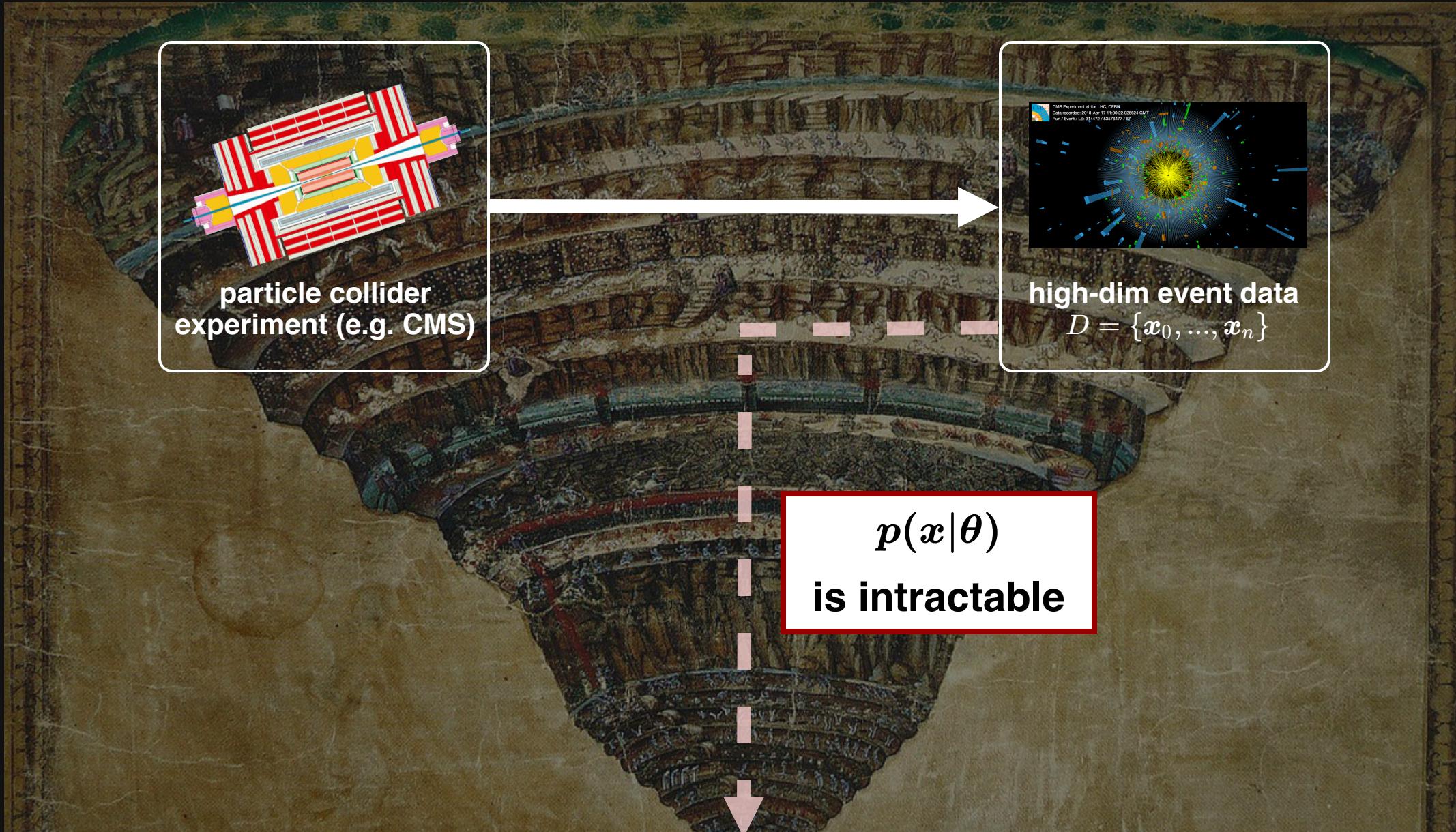
# Statistical Inference in Particle Colliders



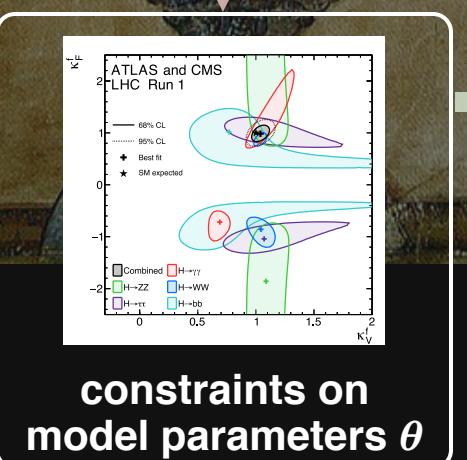
# Statistical Inference in Particle Colliders



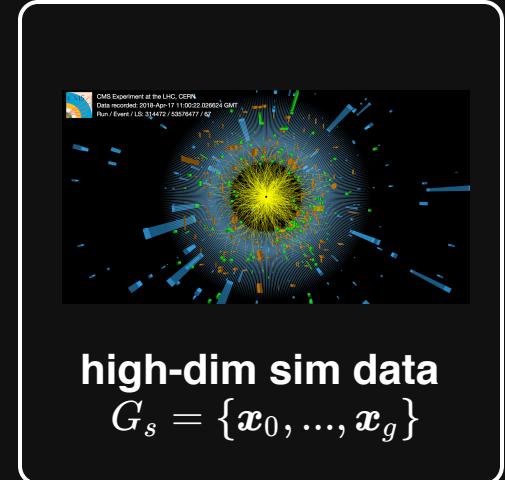
# Statistical Inference in Particle Colliders



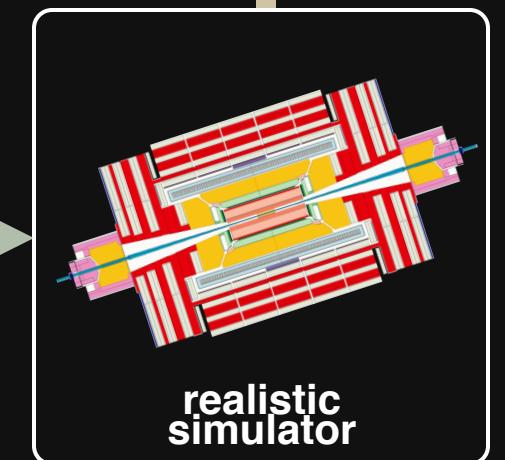
$p(x|\theta)$   
is intractable



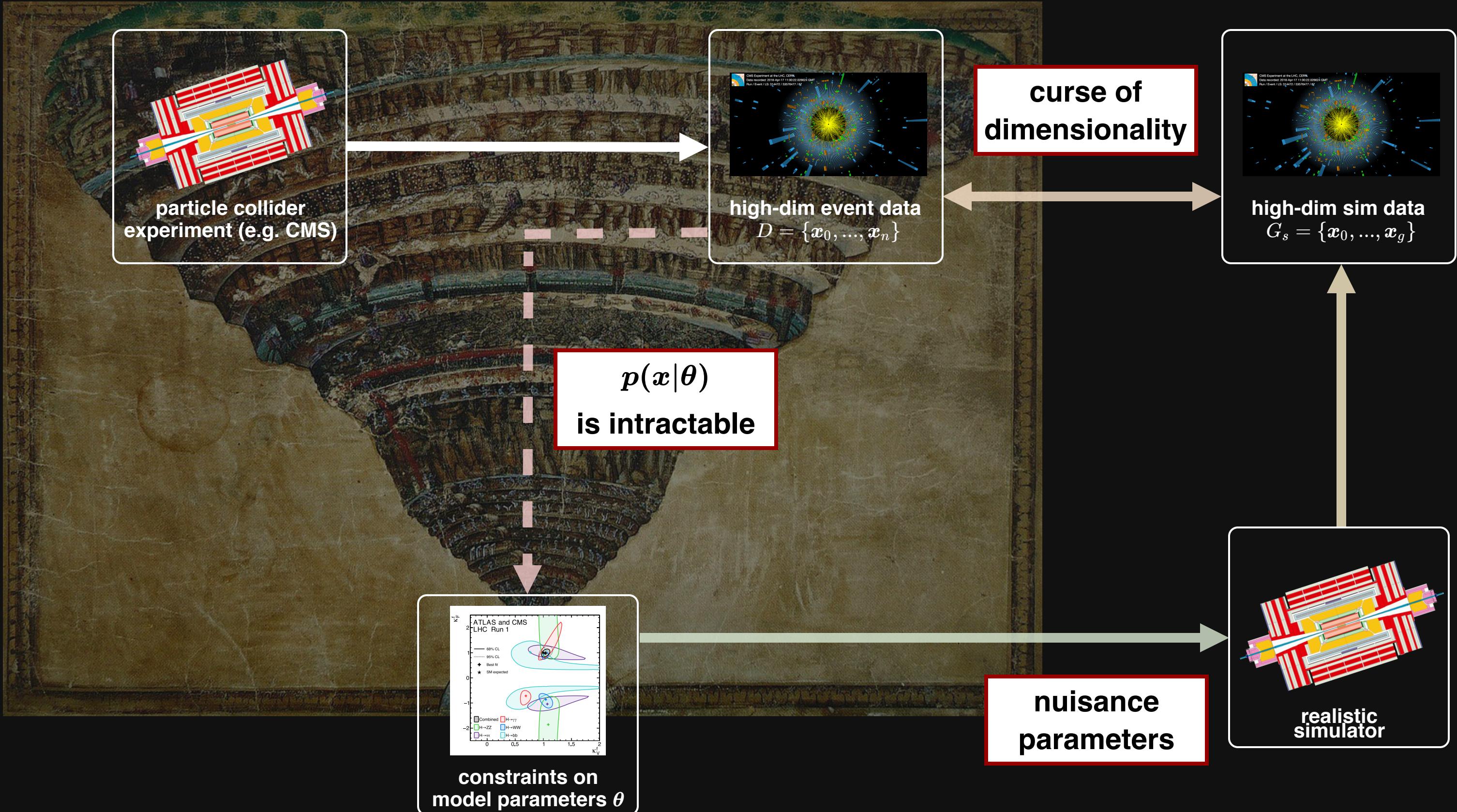
constraints on  
model parameters  $\theta$



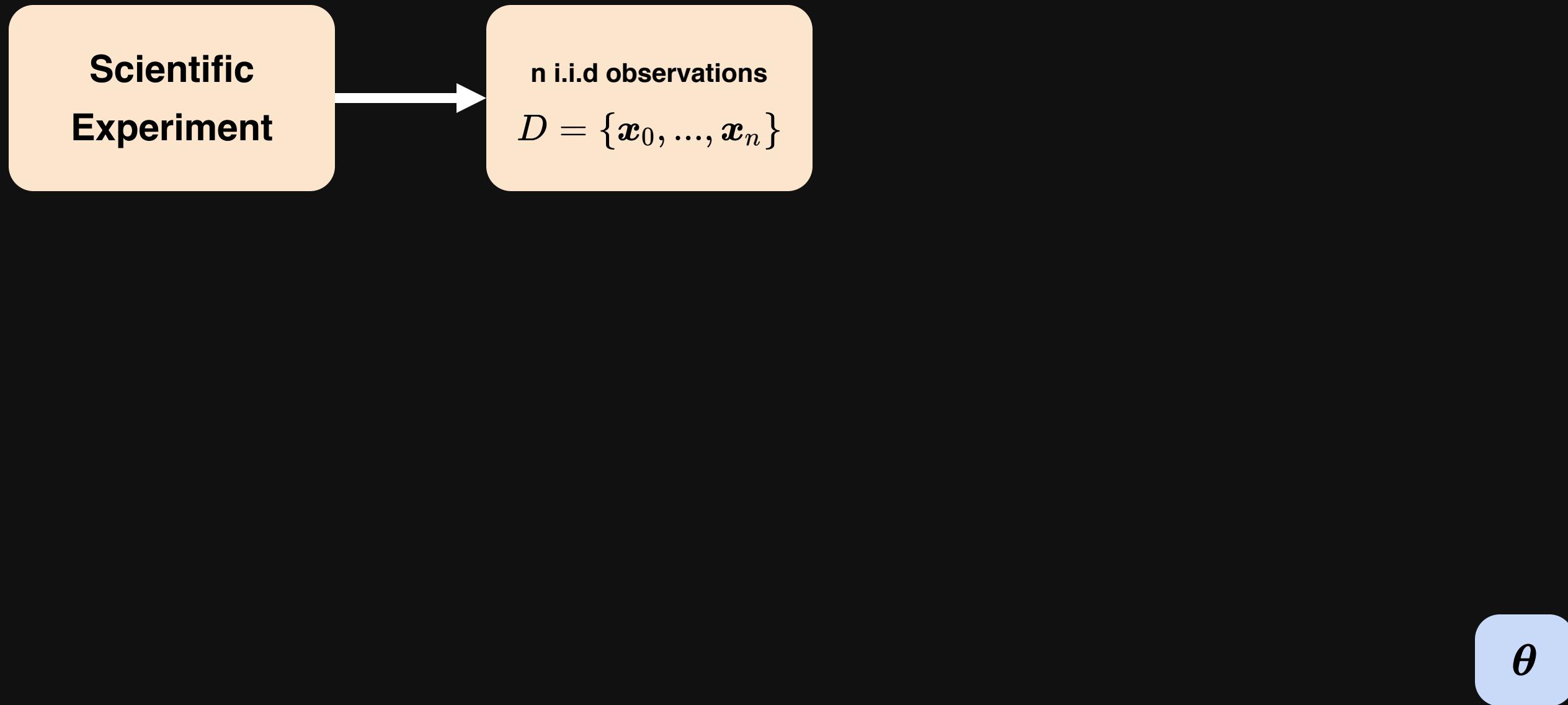
nuisance  
parameters



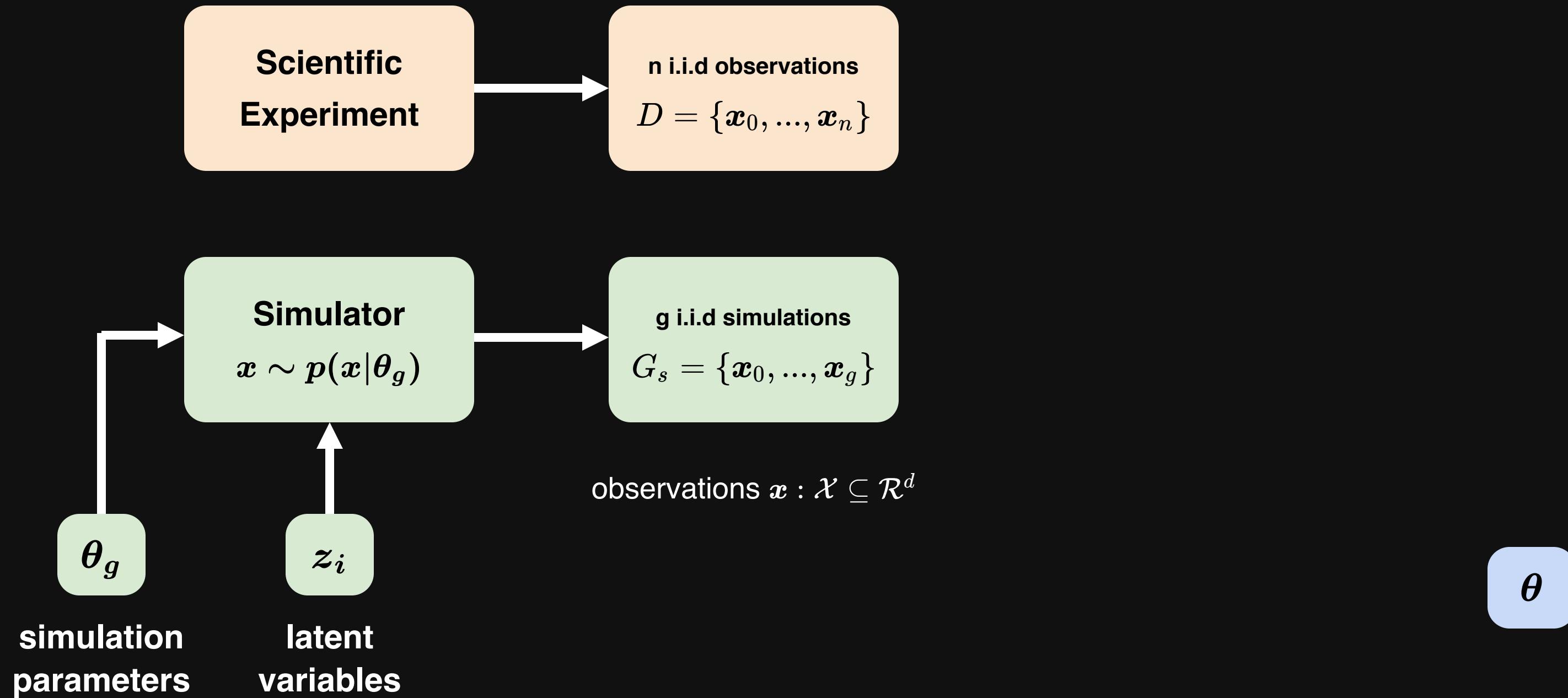
# Statistical Inference in Particle Colliders



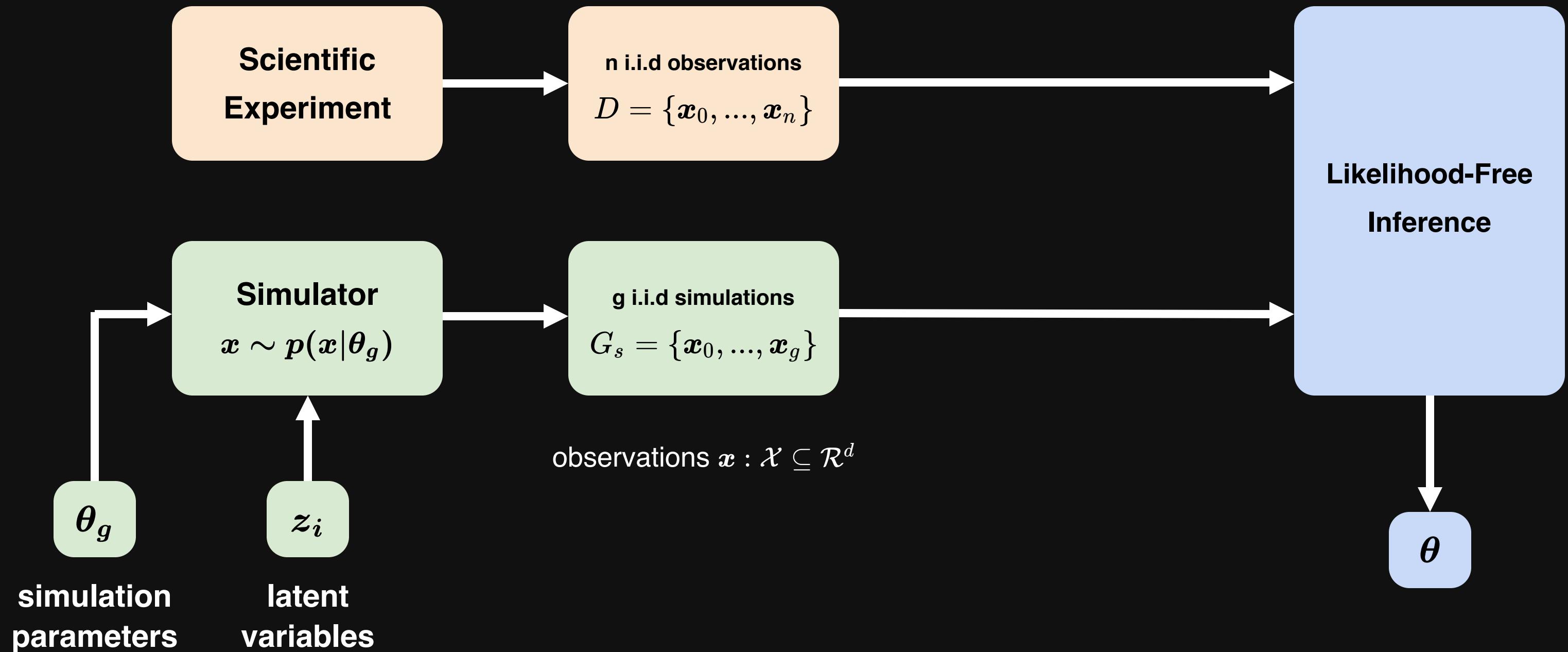
# Simulation-based Inference



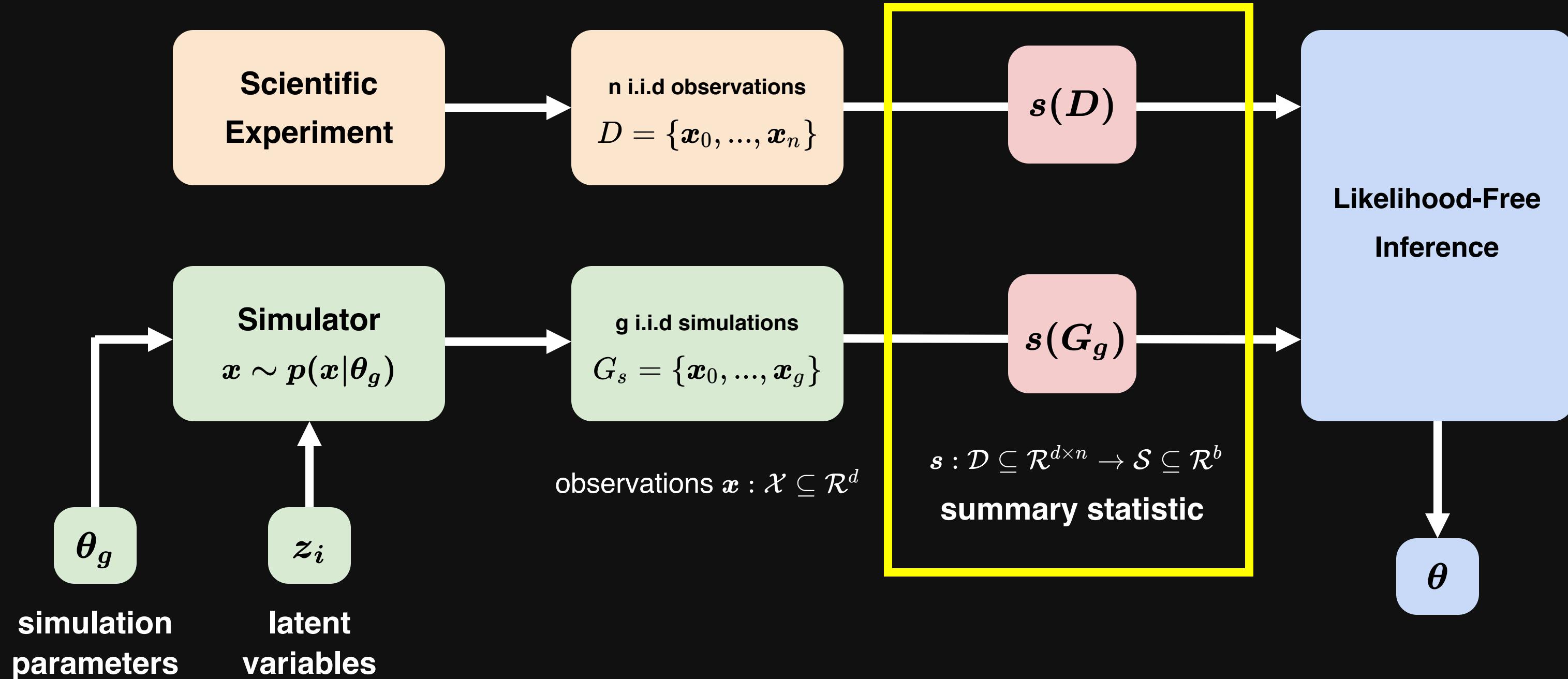
# Simulation-based Inference



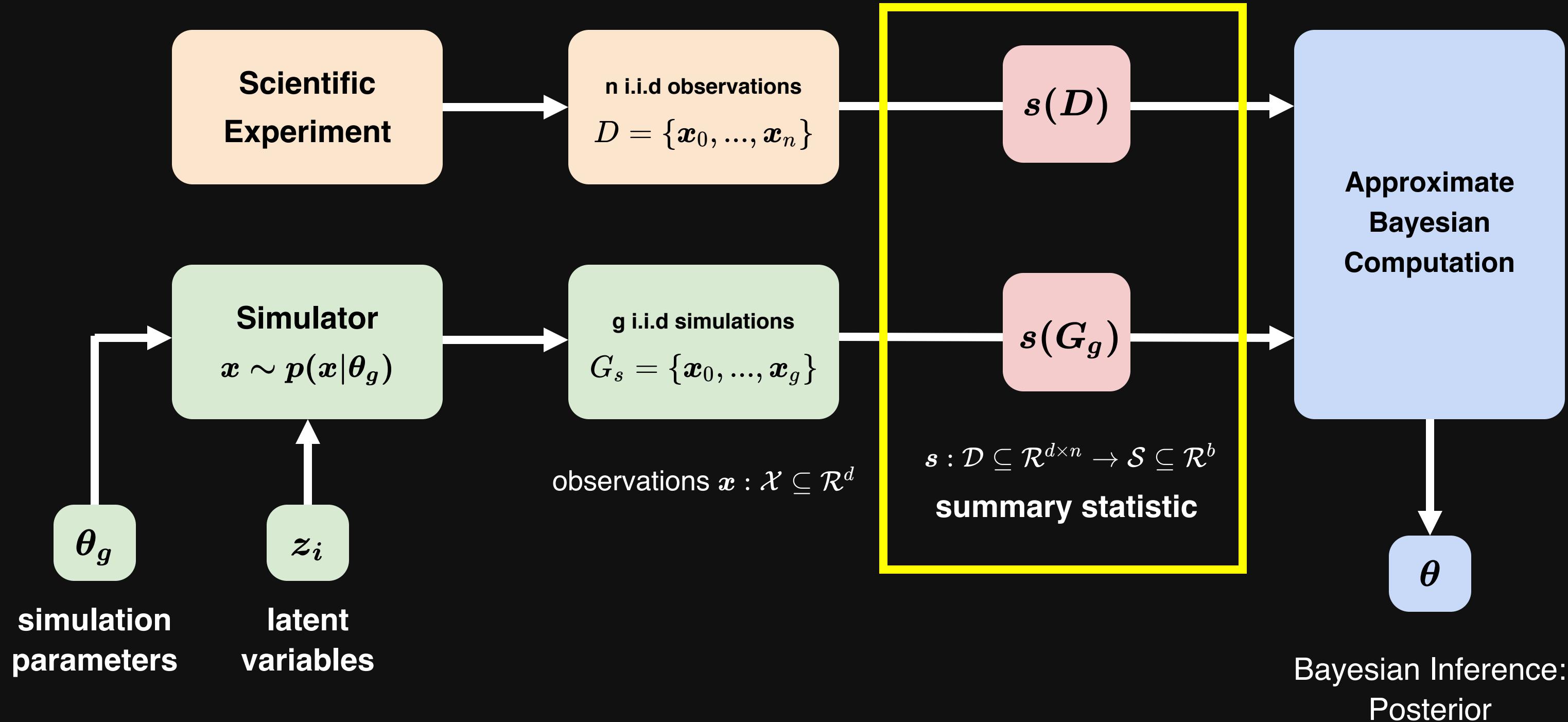
# Simulation-based Inference



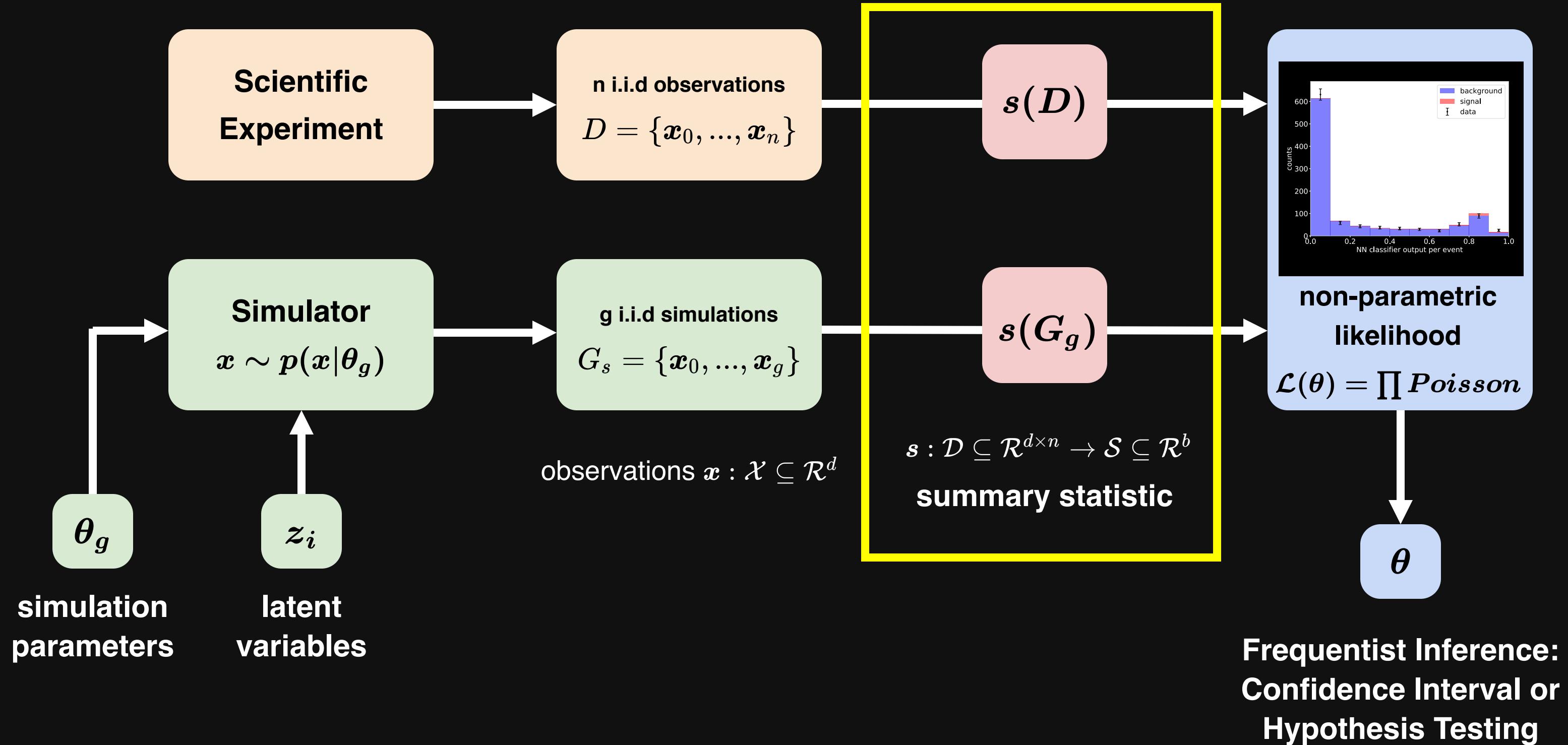
# Simulation-based Inference



# Simulation-based Inference



# Simulation-based Inference



# $p(x|\text{model})$ structure at particle colliders

simulation process  $\sim$  sampling



observation

$x$

theory  
parameters  
 $\theta$

statistical inference



# $p(\mathbf{x}|\text{model})$ structure at particle colliders

simulation process  $\sim$  sampling



observation

$\mathbf{x}$

parton  
momenta

$\mathbf{z}_p$

theory  
parameters

$\theta$

$$p(\mathbf{z}_p|\theta)$$

statistical inference



# $p(\mathbf{x}|\text{model})$ structure at particle colliders

simulation process  $\sim$  sampling



observation

$\mathbf{x}$

shower  
splittings

$\mathbf{z}_s$

parton  
momenta

$\mathbf{z}_p$

theory  
parameters

$\theta$

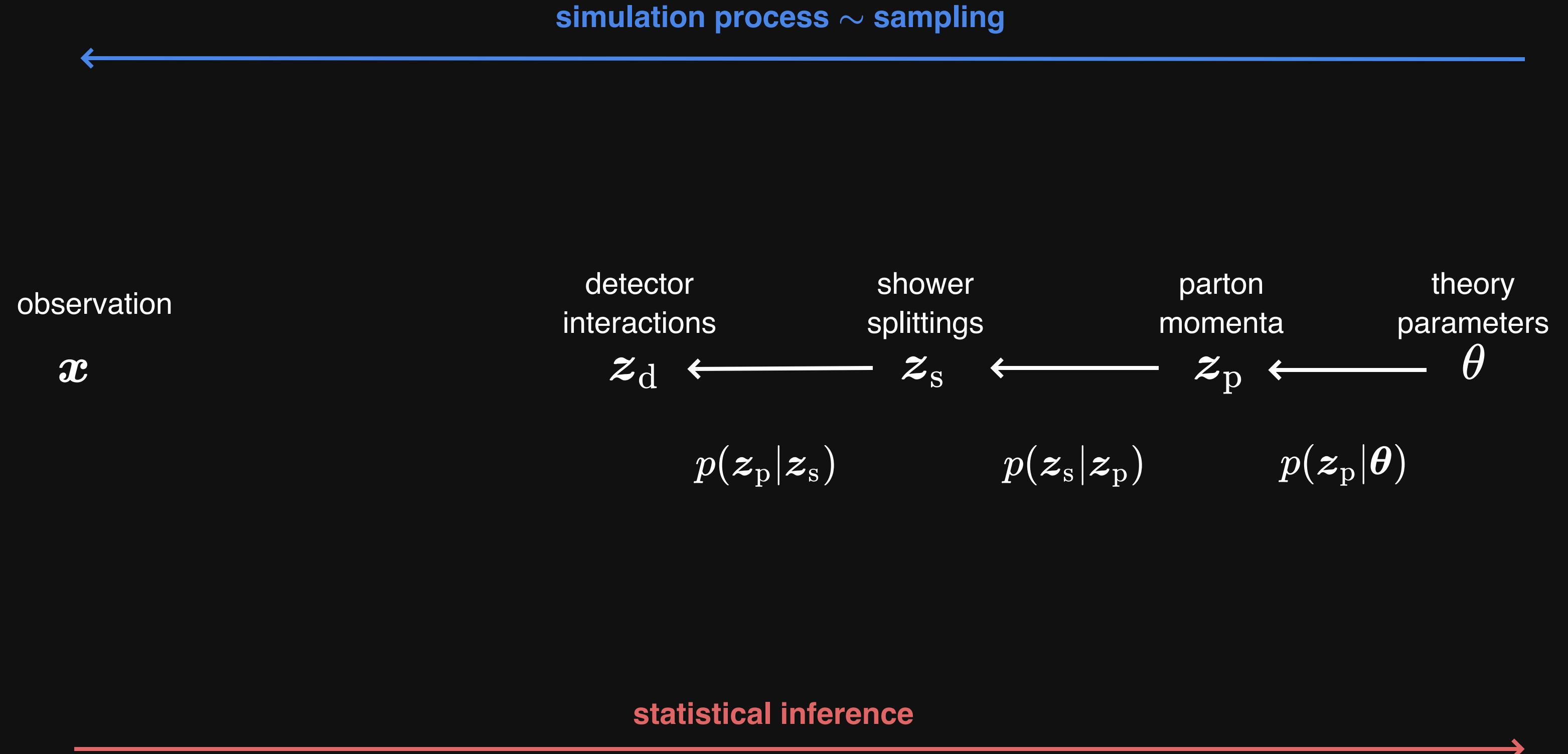
$$p(\mathbf{z}_s | \mathbf{z}_p)$$

$$p(\mathbf{z}_p | \theta)$$

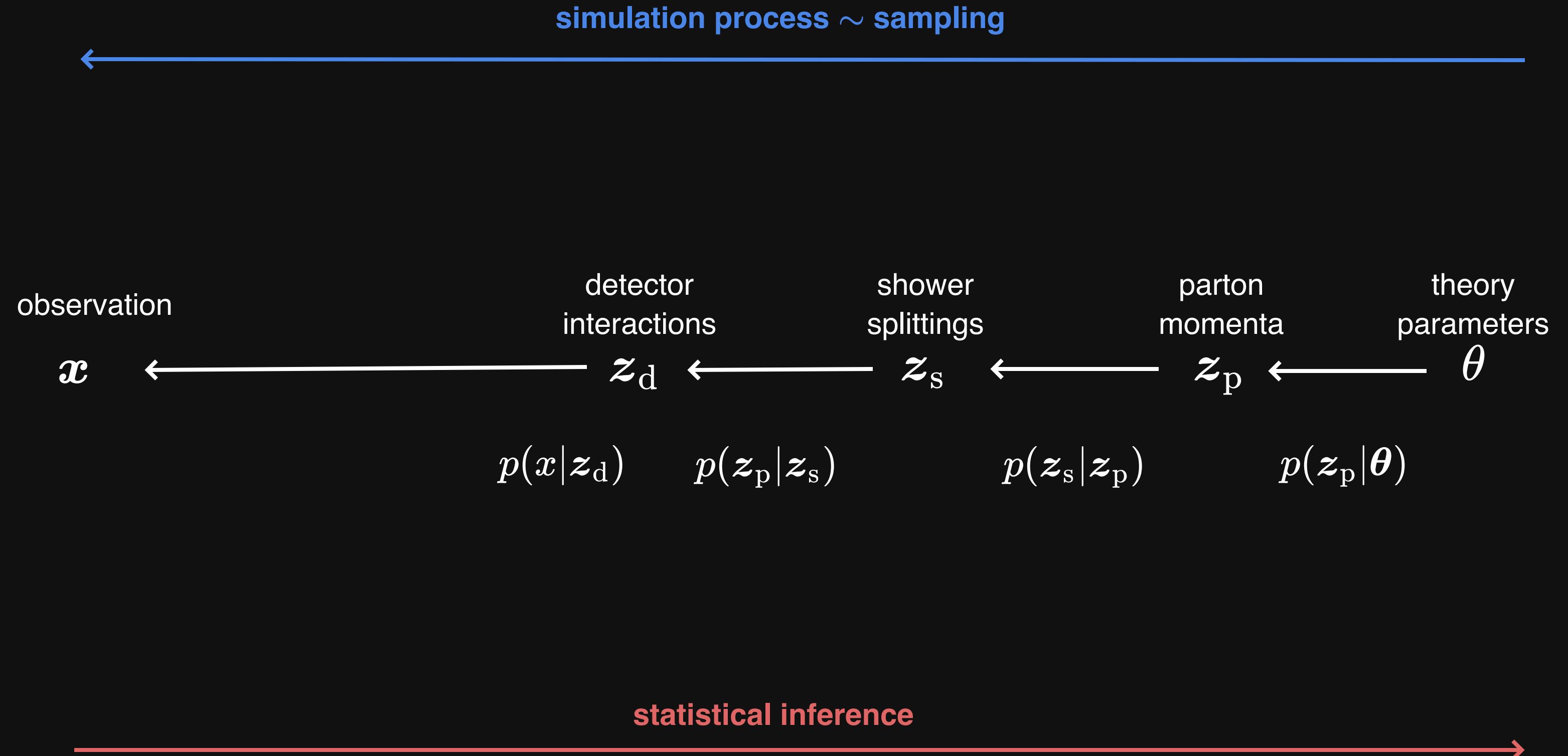
statistical inference



# $p(\mathbf{x}|\text{model})$ structure at particle colliders

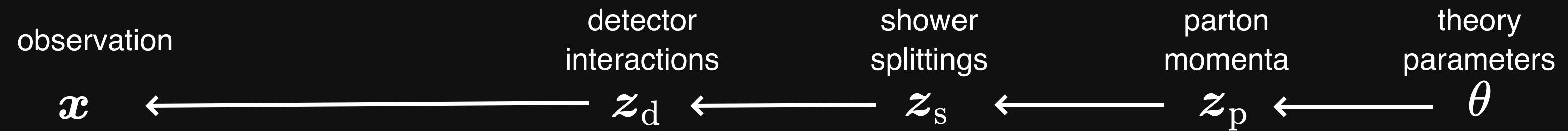


# $p(\mathbf{x}|\text{model})$ structure at particle colliders



# $p(x|\text{model})$ structure at particle colliders

simulation process  $\sim$  sampling



$$p(x|\theta) = \int dz_p \int dz_s \int dz_d \ p(x|z_d) \ p(z_p|z_s) \ p(z_s|z_p) \ p(z_p|\theta)$$

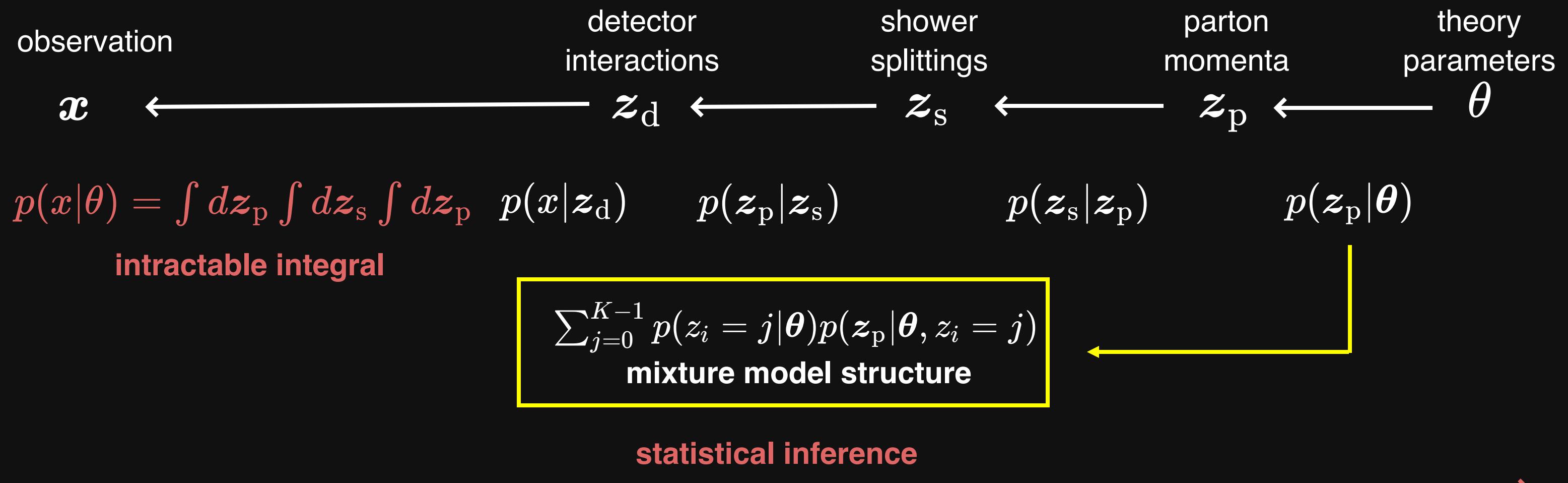
intractable integral

statistical inference

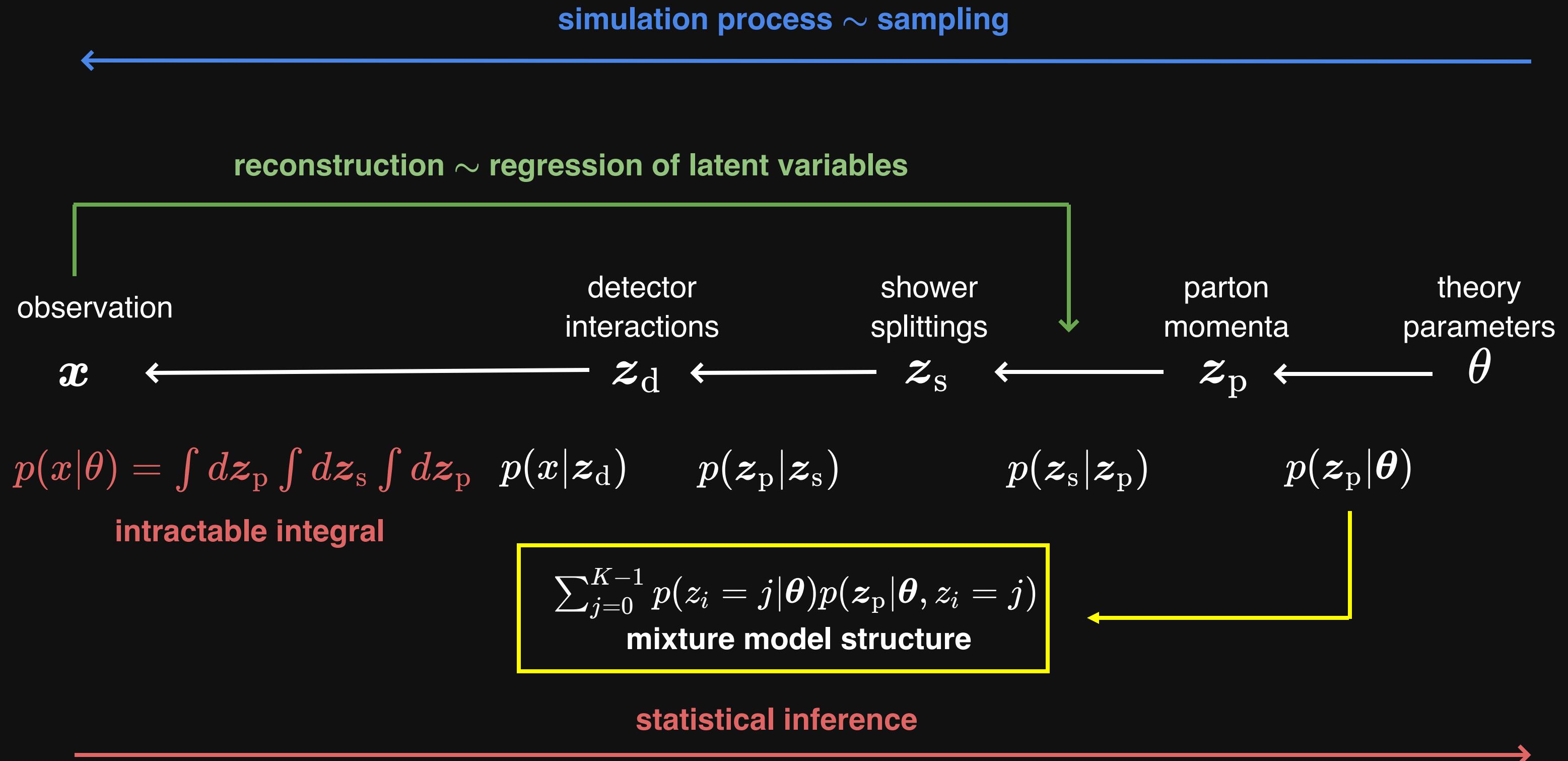


# $p(\mathbf{x}|\text{model})$ structure at particle colliders

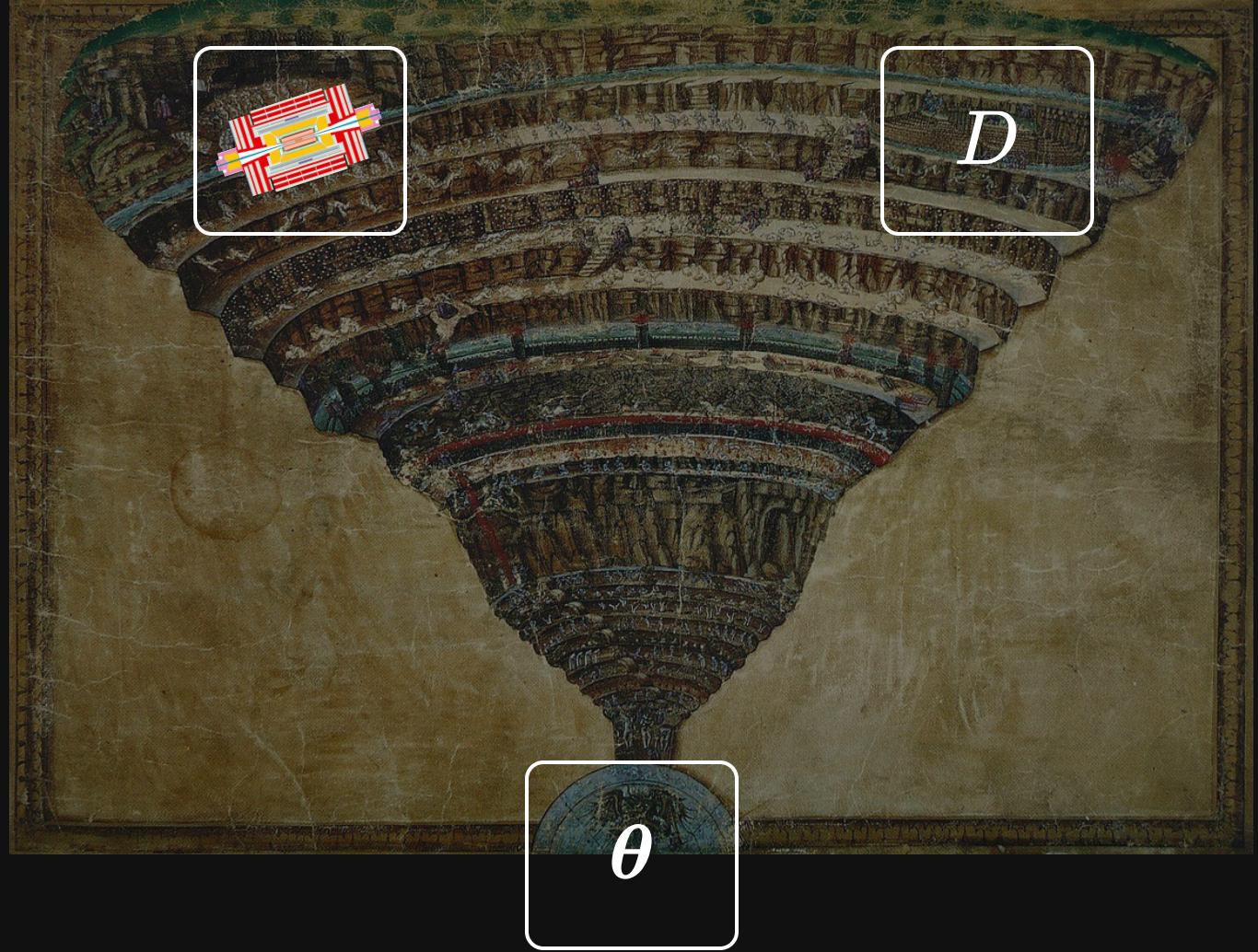
simulation process  $\sim$  sampling



# $p(\mathbf{x}|\text{model})$ structure at particle colliders

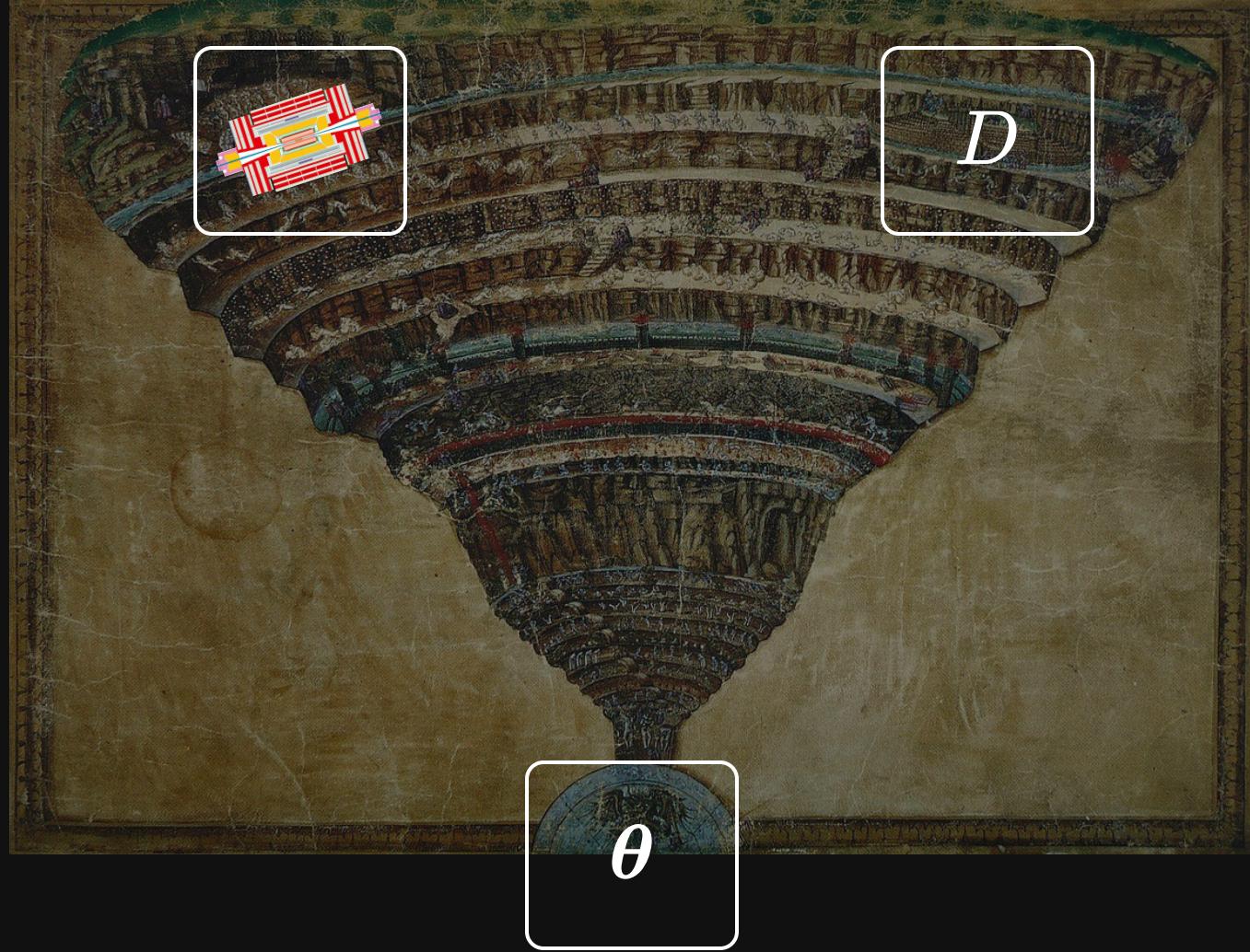


# The *hard problem* of choosing a $s(D)$



A low-dimensional summary statistics is required for inference, yet a naive choice will result on an important information loss

# The *hard problem* of choosing a $s(D)$



A low-dimensional summary statistics is required for inference, yet a naive choice will result on an important information loss

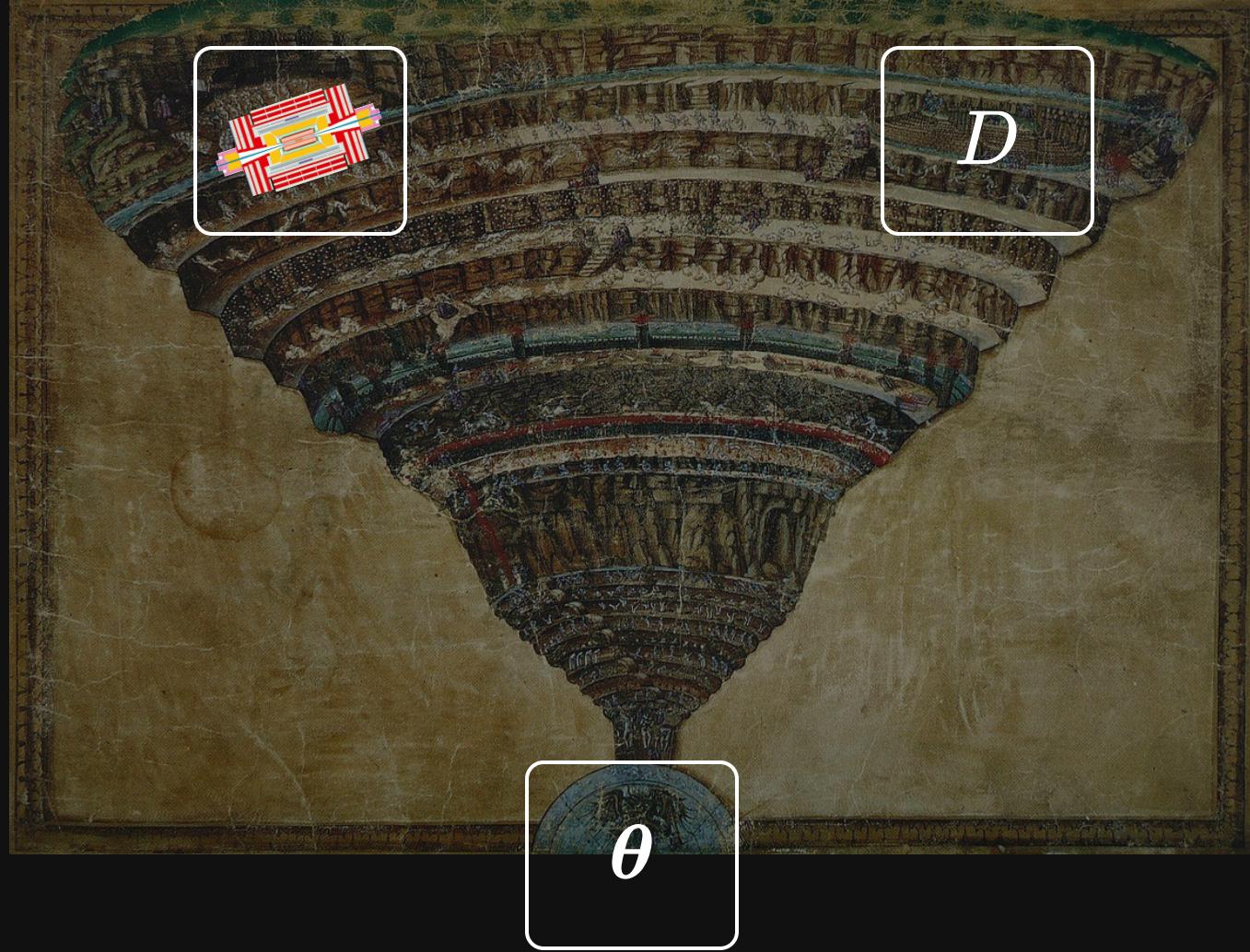
Ideally we want a sufficient summary statistic:

$$p(D|\theta) = h(D)g(s(D)|\theta)$$

classical sufficiency

cannot be obtained in general  
and might not exist

# The *hard problem* of choosing a $s(D)$



A low-dimensional summary statistics is required for inference, yet a naive choice will result on an important information loss

Ideally we want a sufficient summary statistic:

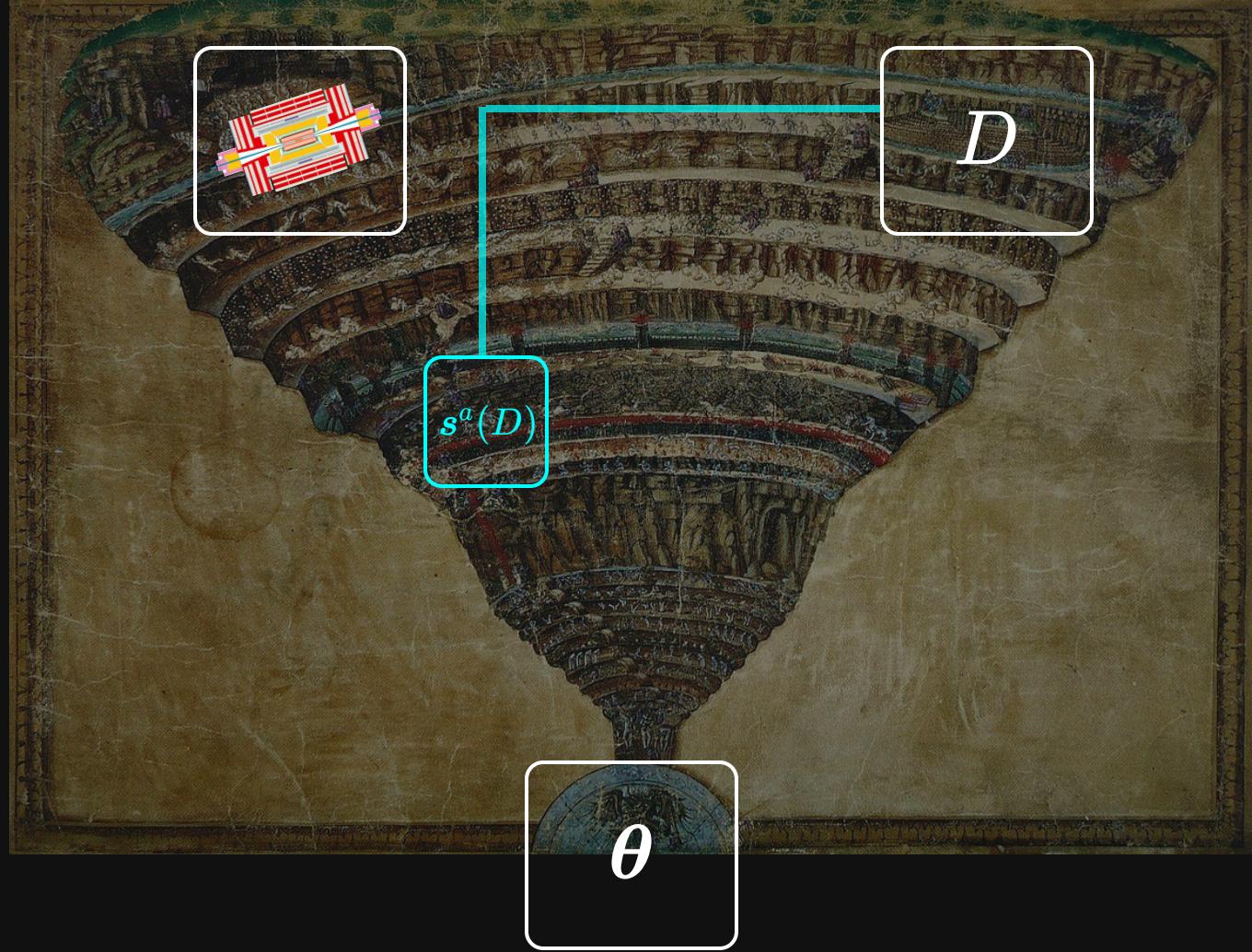
$$p(D|\theta) = h(D)g(s(D)|\theta)$$

classical sufficiency

cannot be obtained in general  
and might not exist

How to obtain good a summary statistic  $s(D)$ ?

# The *hard problem* of choosing a $s(D)$



A low-dimensional summary statistics is required for inference, yet a naive choice will result on an important information loss

Ideally we want a **sufficient summary statistic**:

$$p(D|\theta) = h(D)g(s(D)|\theta)$$

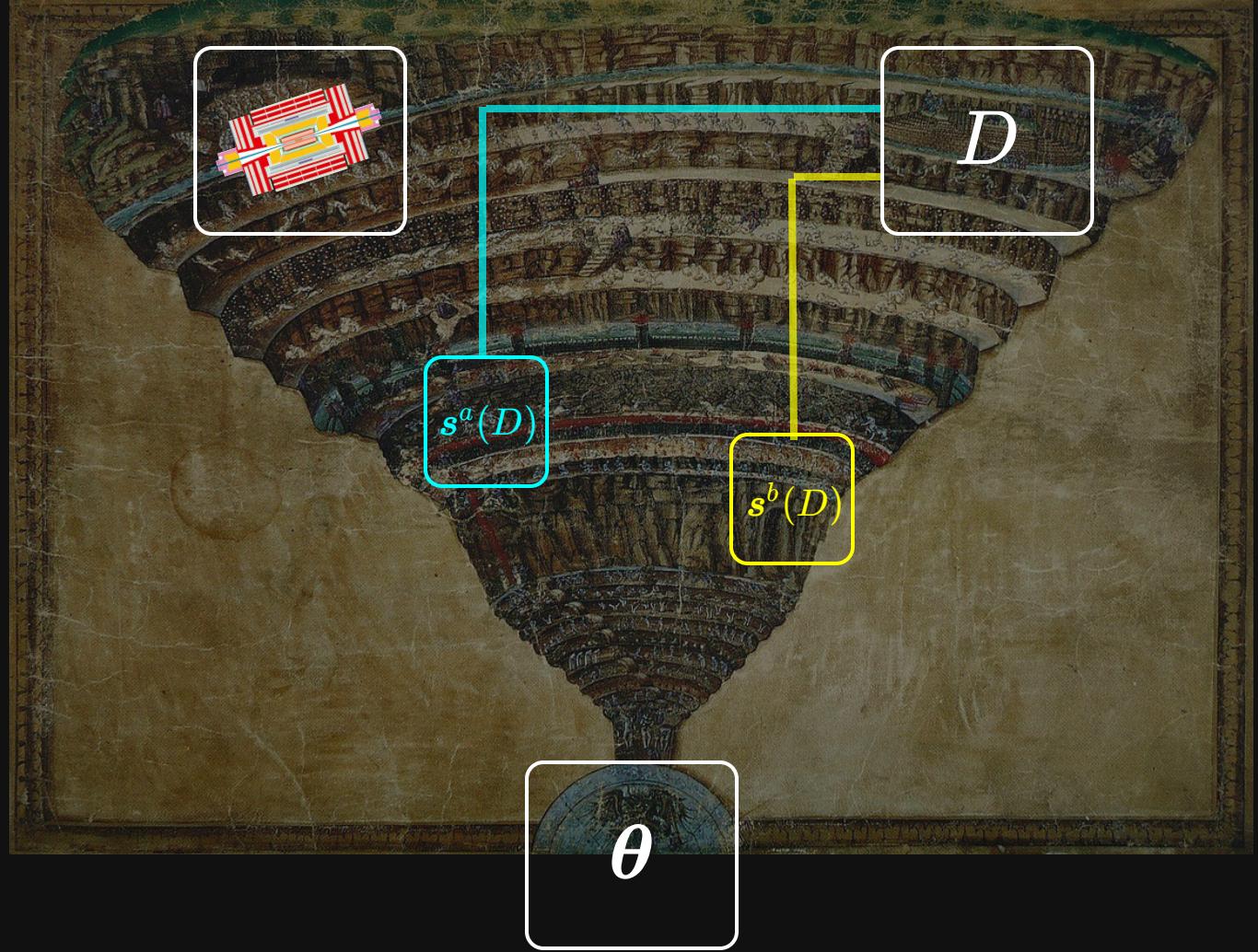
classical sufficiency

cannot be obtained in general  
and might not exist

How to obtain good a summary statistic  $s(D)$ ?

- a) Use domain inspired (i.e. physics) variables based on a combination of reconstructed variables

# The *hard problem* of choosing a $s(D)$



A low-dimensional summary statistics is required for inference, yet a naive choice will result on an important information loss

Ideally we want a sufficient summary statistic:

$$p(D|\theta) = h(D)g(s(D)|\theta)$$

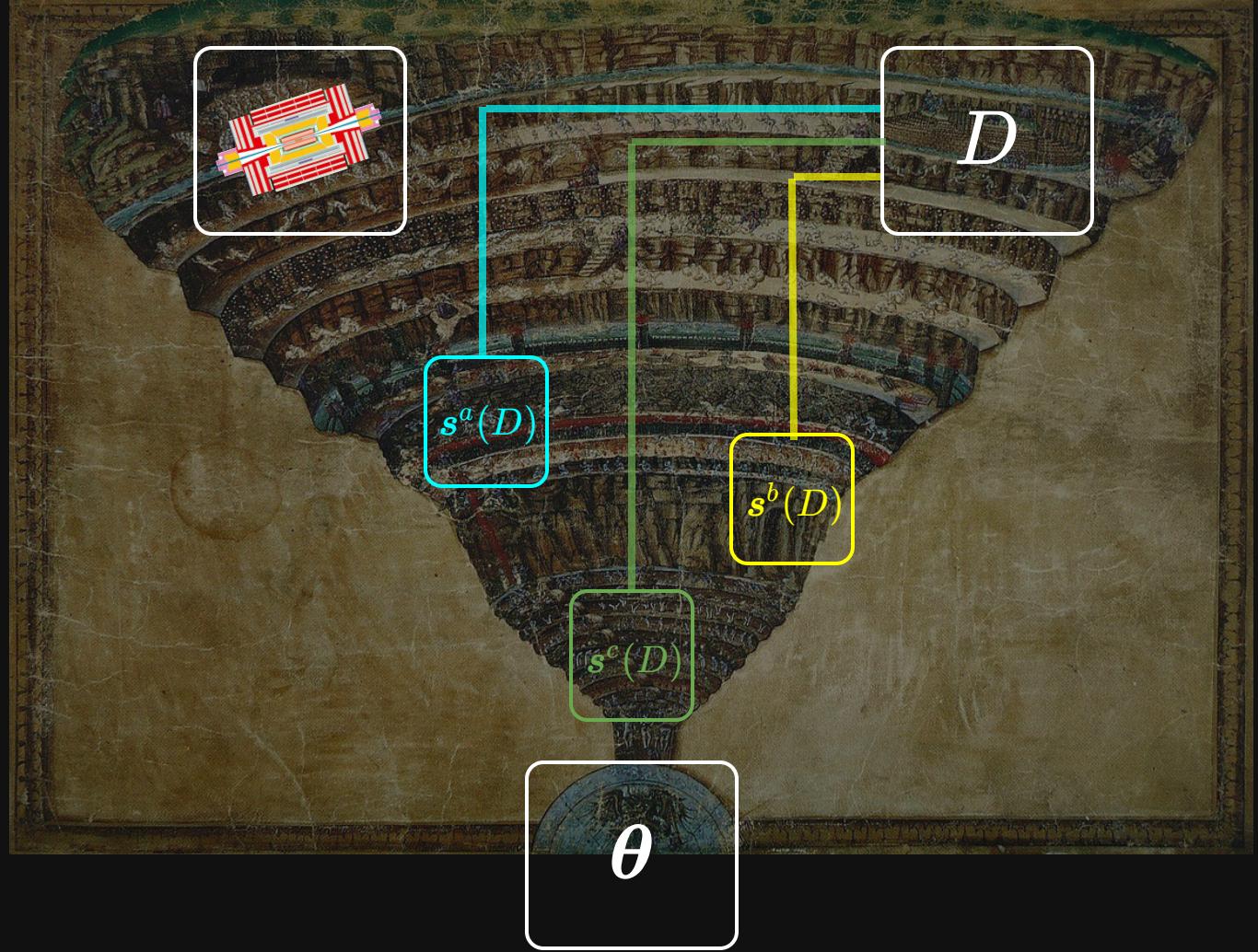
classical sufficiency

cannot be obtained in general  
and might not exist

How to obtain good a summary statistic  $s(D)$ ?

- a) Use domain inspired (i.e. physics) variables based on a combination of reconstructed variables
- b) Use the prediction of a machine learning classification or regression model trained on simulated observations

# The *hard problem* of choosing a $s(D)$



A low-dimensional summary statistics is required for inference, yet a naive choice will result on an important information loss

Ideally we want a sufficient summary statistic:

$$p(D|\theta) = h(D)g(s(D)|\theta)$$

classical sufficiency

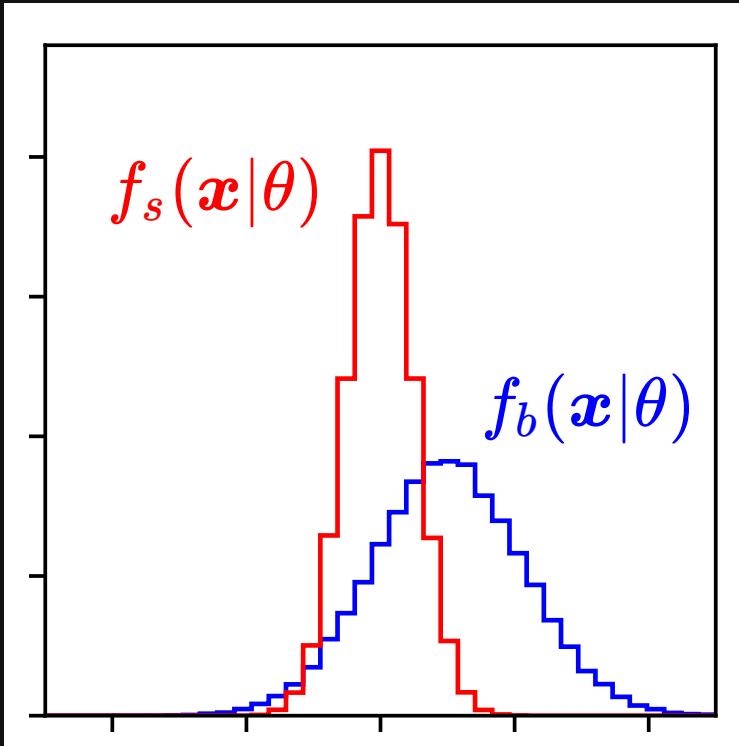
cannot be obtained in general  
and might not exist

How to obtain good a summary statistic  $s(D)$ ?

- a) Use domain inspired (i.e. physics) variables based on a combination of reconstructed variables
- b) Use the prediction of a machine learning classification or regression model trained on simulated observations
- c) Use the technique presented in this talk or one of the alternatives for obtaining powerful  $s(x)$  proposed in this workshop

# An overused trick to obtain a good $s(\mathbf{x})$

IF



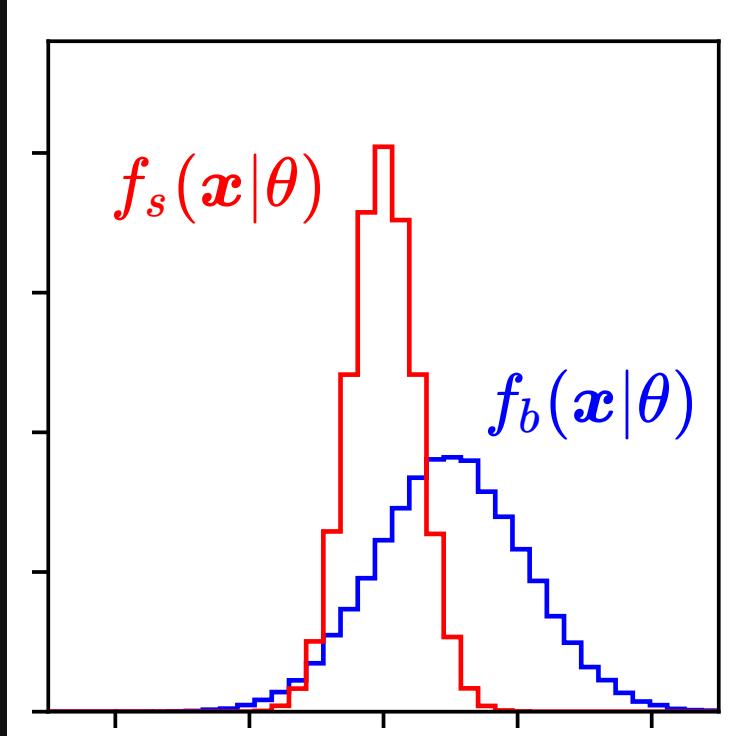
$$p(\mathbf{x}|\mu, \boldsymbol{\theta}) = \mu f_s(\mathbf{x}|\boldsymbol{\theta}) + (1 - \mu) f_b(\mathbf{x}|\boldsymbol{\theta})$$

2-component mixture

and we only want to infer about the mixture fraction  $\mu$ , **other  $\boldsymbol{\theta}$  are fixed and known**

# An overused trick to obtain a good $s(\mathbf{x})$

IF



$$p(\mathbf{x}|\mu, \boldsymbol{\theta}) = \mu f_s(\mathbf{x}|\boldsymbol{\theta}) + (1 - \mu) f_b(\mathbf{x}|\boldsymbol{\theta})$$

2-component mixture

and we only want to infer about the mixture fraction  $\mu$ , **other  $\boldsymbol{\theta}$  are fixed and known**

THEN

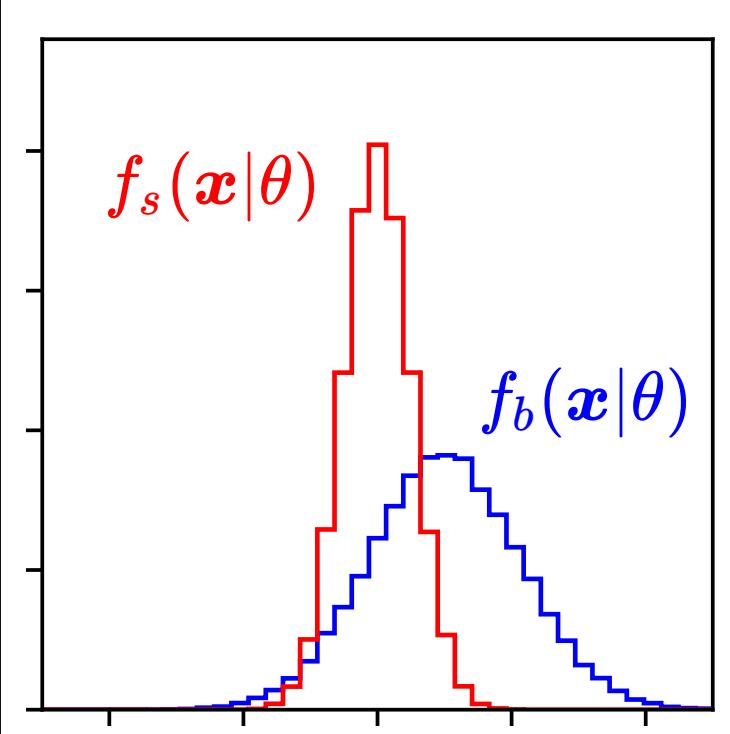
$$s_{\text{clf}}(\mathbf{x}|\boldsymbol{\theta}) = \frac{f_s(\mathbf{x}|\boldsymbol{\theta})}{f_s(\mathbf{x}|\boldsymbol{\theta}) + f_b(\mathbf{x}|\boldsymbol{\theta})}$$

bayes optimal classifier

is a **one-dimensional summary statistic**  
that is **sufficient for inference about  $\mu$**

# An overused trick to obtain a good $s(\mathbf{x})$

IF



$$p(\mathbf{x}|\mu, \boldsymbol{\theta}) = \mu f_s(\mathbf{x}|\boldsymbol{\theta}) + (1 - \mu) f_b(\mathbf{x}|\boldsymbol{\theta})$$

2-component mixture

and we only want to infer about the mixture fraction  $\mu$ , **other  $\boldsymbol{\theta}$  are fixed and known**

THEN

$$s_{\text{clf}}(\mathbf{x}|\boldsymbol{\theta}) = \frac{f_s(\mathbf{x}|\boldsymbol{\theta})}{f_s(\mathbf{x}|\boldsymbol{\theta}) + f_b(\mathbf{x}|\boldsymbol{\theta})}$$

bayes optimal classifier

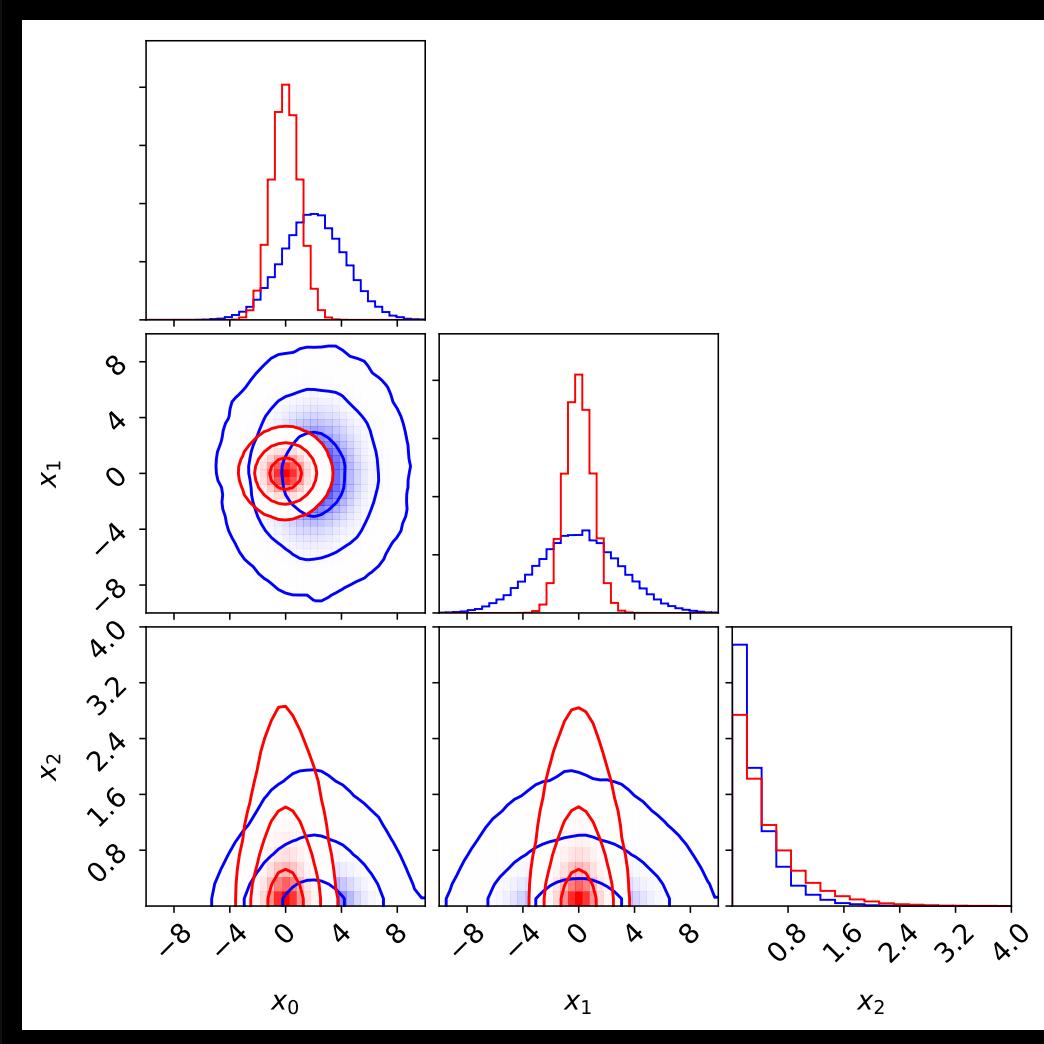
is a **one-dimensional summary statistic**  
that is **sufficient for inference about  $\mu$**

*can be approximated very efficiently by training a probabilistic classifier to distinguish signal and background observations*

$$L_{\text{BCE}} = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

*e.g. neural net minimising cross entropy*

# An example: 3D Synthetic Problem



A 3D (i.e.  $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^3$ ) mixture component model:

$$p(\mathbf{x}|\mu, r, \lambda) = \mu f_s(\mathbf{x}) + (1 - \mu) f_b(\mathbf{x}|r, \lambda)$$

where **signal** and **background** distributions are

$$f_s(\mathbf{x}) = \mathcal{N}\left((x_0, x_1) \mid (1, 1), \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right) \text{Exp}(x_2|2)$$

$$f_b(\mathbf{x}|r, \lambda) = \mathcal{N}\left((x_0, x_1) \mid (2+r, 0), \begin{bmatrix} 5 & 0 \\ 0 & 9 \end{bmatrix}\right) \text{Exp}(x_2|\lambda)$$

$$s = 50 \quad b = 1000 \quad \lambda = 3.0 \quad r = 0.0$$

**alternative parametrisation that will be used**

$$p(\mathbf{x}|s, r, \lambda, b) = \frac{s}{s+b} f_s(\mathbf{x}) + \frac{b}{s+b} f_b(\mathbf{x}|r, \lambda)$$

parameter of interest

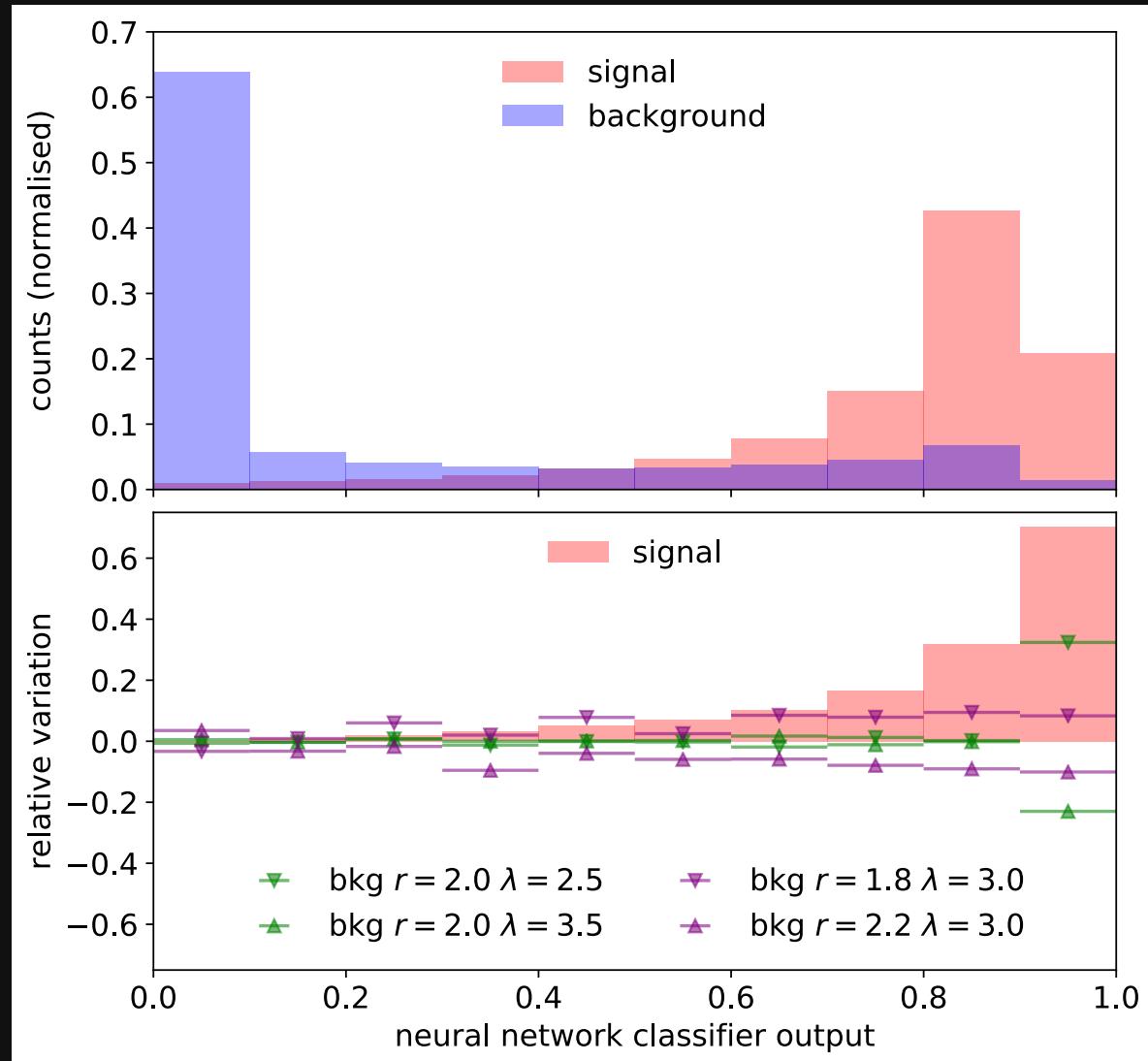
nuisance parameters

# Classification as a Surrogate Task

Probabilistic classification has been used for many years in particle colliders to obtain powerful summary statistics by discriminating signal and background simulated observations:

*neural net minimising cross entropy*

$$L_{\text{BCE}} = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

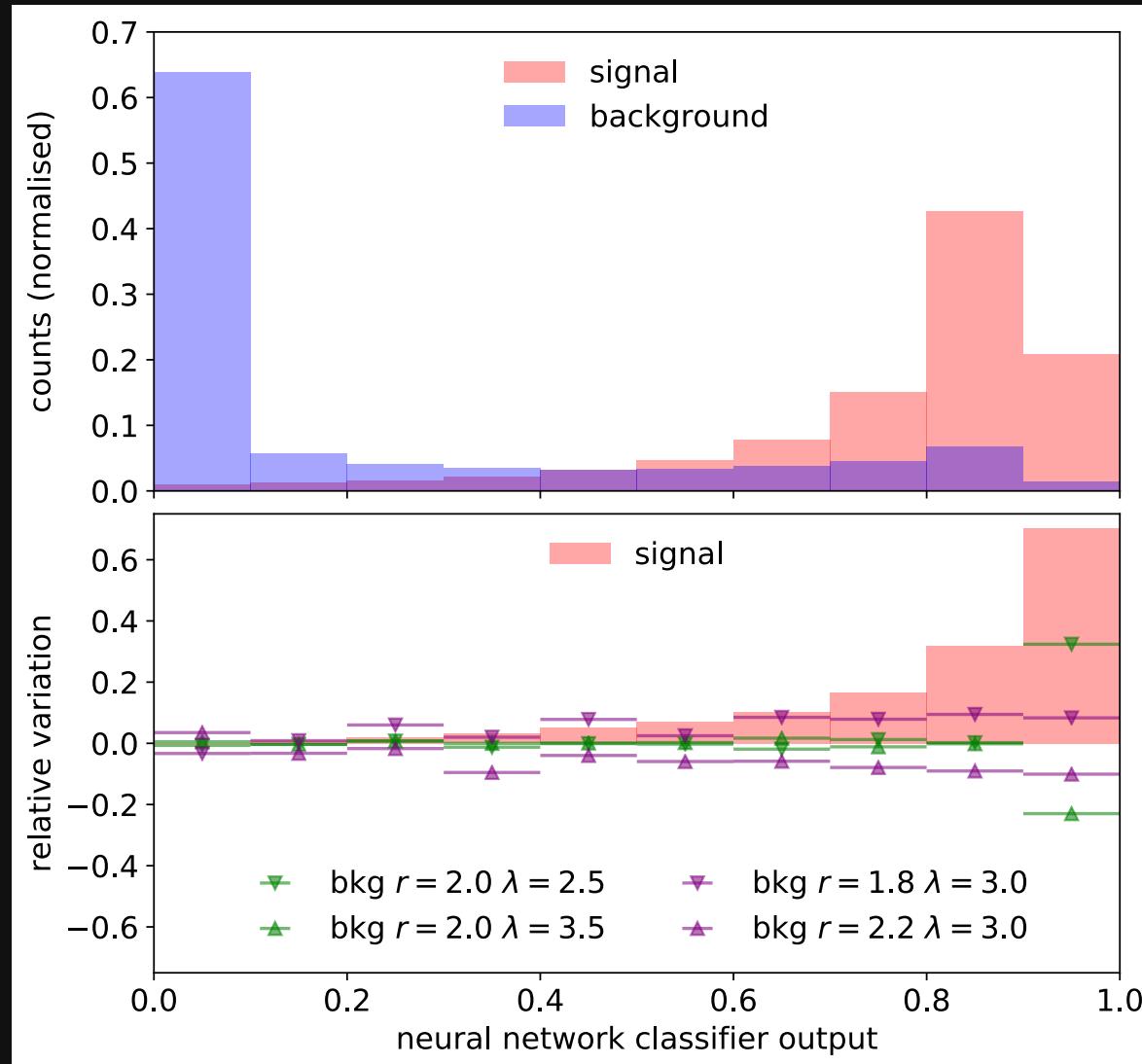


# Classification as a Surrogate Task

Probabilistic classification has been used for many years in particle colliders to obtain powerful summary statistics by discriminating signal and background simulated observations:

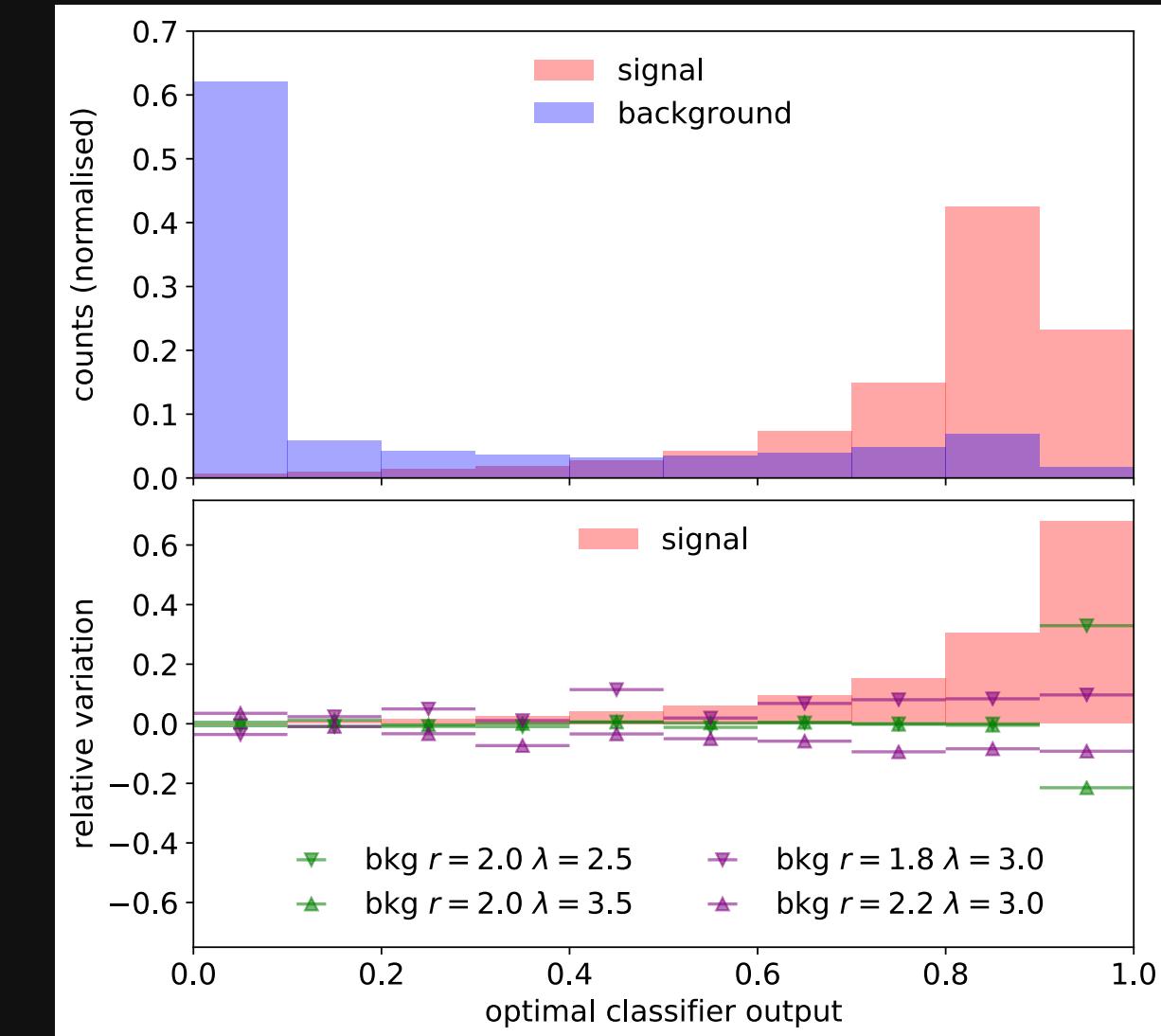
*neural net minimising cross entropy*

$$L_{\text{BCE}} = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$



*analytical*  $s_{\text{clf}}(\mathbf{x}|\boldsymbol{\theta}) = \frac{f_s(\mathbf{x}|\boldsymbol{\theta})}{f_s(\mathbf{x}|\boldsymbol{\theta}) + f_b(\mathbf{x}|\boldsymbol{\theta})}$

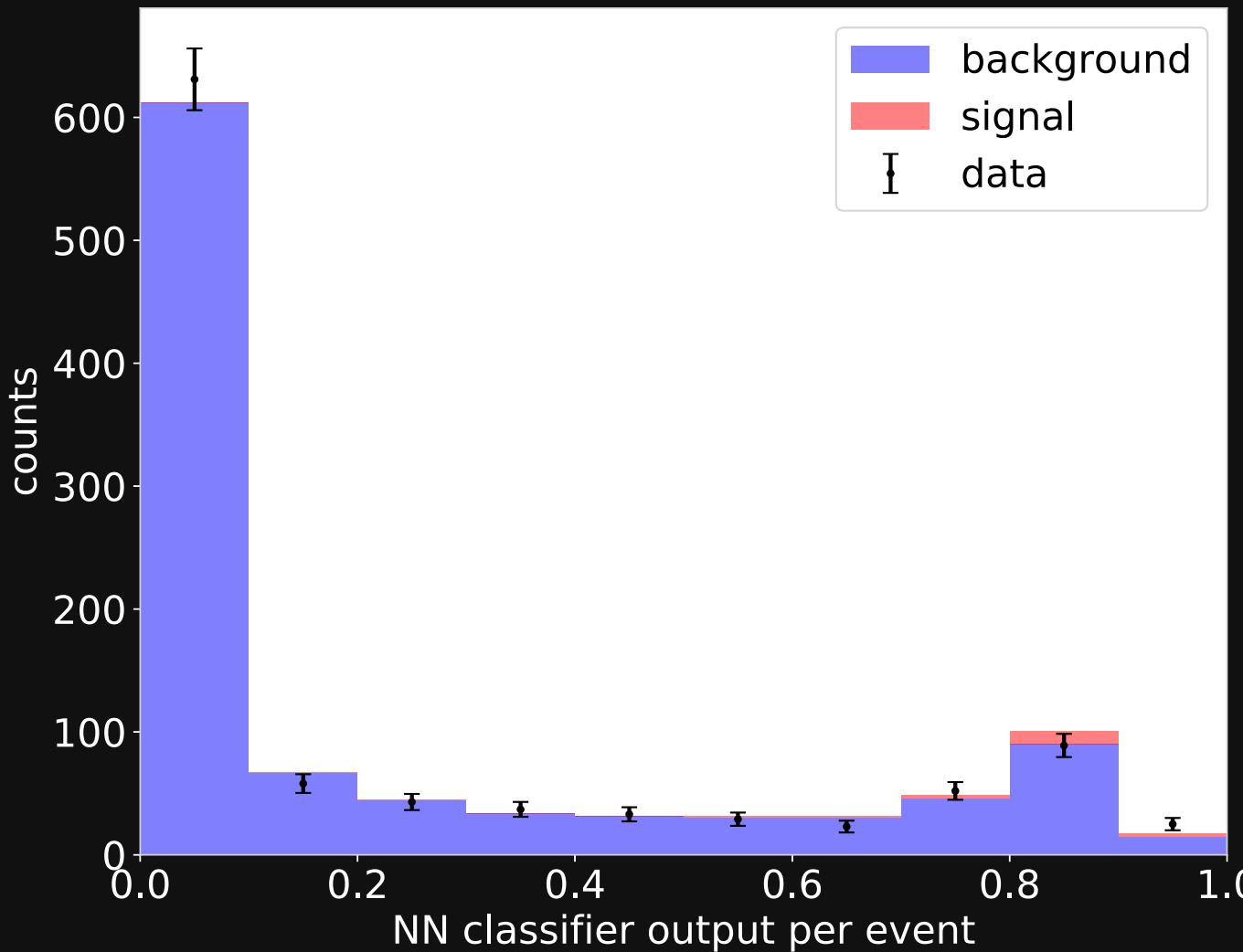
very good approximation



Given enough data and a powerful enough model and learning rule, the probabilisitic classifier output can **approximate** very closely the optimal classifier  $s_{\text{clf}}(\mathbf{x})$ , which is a sufficient statistic for  $s$

# Inference based on a low-dim $s(\mathbf{x})$

In most particle collider analyses, statistical inference is based on a single one-dimensional summary statistic  $s(\mathbf{x})$  such the output of probabilistic classifier or a clever combination of reconstructed variables based on physics



$$\mathbf{x} \longrightarrow s(\mathbf{x})$$

$\dim \mathcal{O}(10^8)$                                      $\dim \mathcal{O}(1)$

Inference is then carried out by building a non-parametric likelihood, typically a histogram associated with a Poisson count likelihood:

$$\mathcal{L}(\boldsymbol{\theta}|s(D)) = \prod_{i \in \text{bins}} \text{Pois}(n_i | n_i^s(\boldsymbol{\theta}) + n_i^b(\boldsymbol{\theta}))$$

which can in turn be used for arbitrary statistical inference given some data, such as obtaining CI on  $\mu$  or arbitrary hypothesis testing

*In the case of classifier-based  $s(\mathbf{x})$ , this approach should be near optimal when  $\mu$  (or  $s$ ) are the only unknown parameters and the histogram binning is small enough*

# Real world problems → nuisance pars

Simulations are often imperfect and often depend on additional modelling parameters:

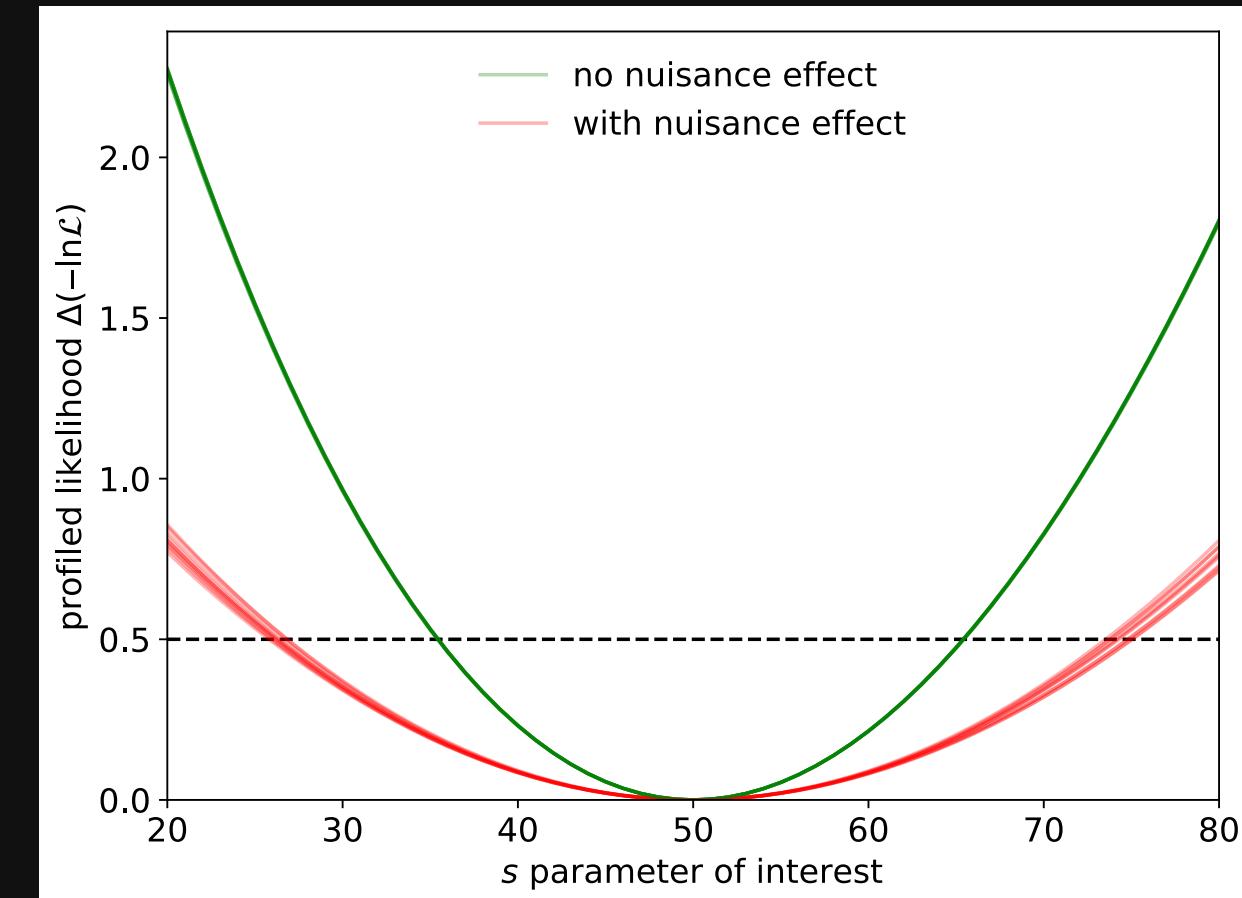
$$\theta = (\eta, \nu)$$

nuisance parameters  
(e.g.  $\lambda$  and  $r$  in 3D problem)

parameters of interest (e.g.  $s$  in 3D problem)

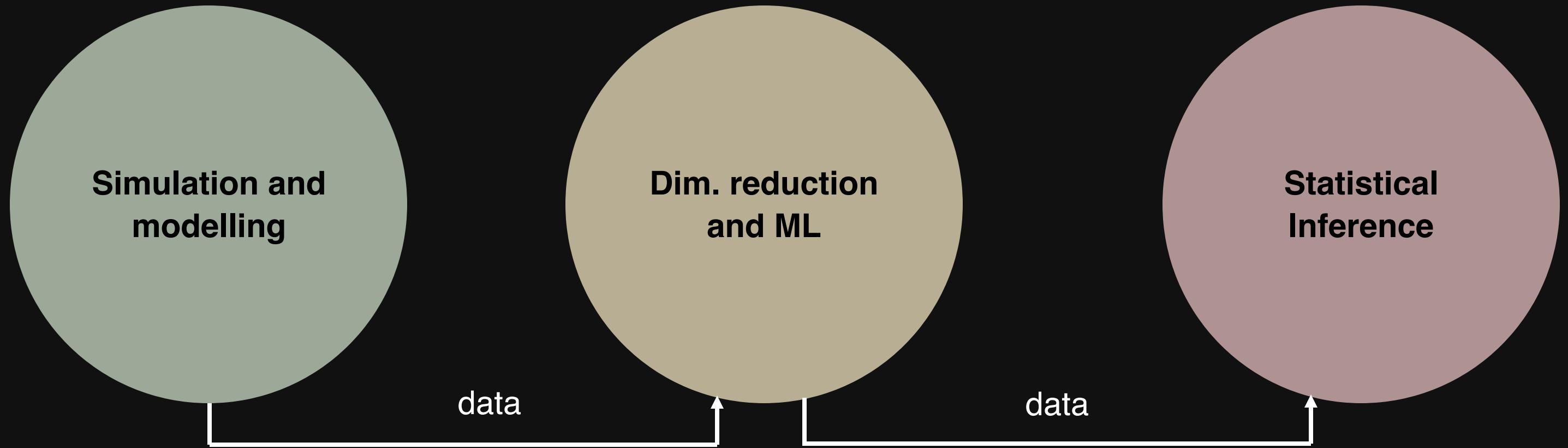
Cause a degradation of classifier-based inference, leading to larger uncertainties in the p.o.i.

**Profile likelihood of  $s$  when:**  
 $s$  unknown -  $\lambda$  and  $r$  fixed  
 $s$  unknown -  $\lambda$  and  $r$  unknown

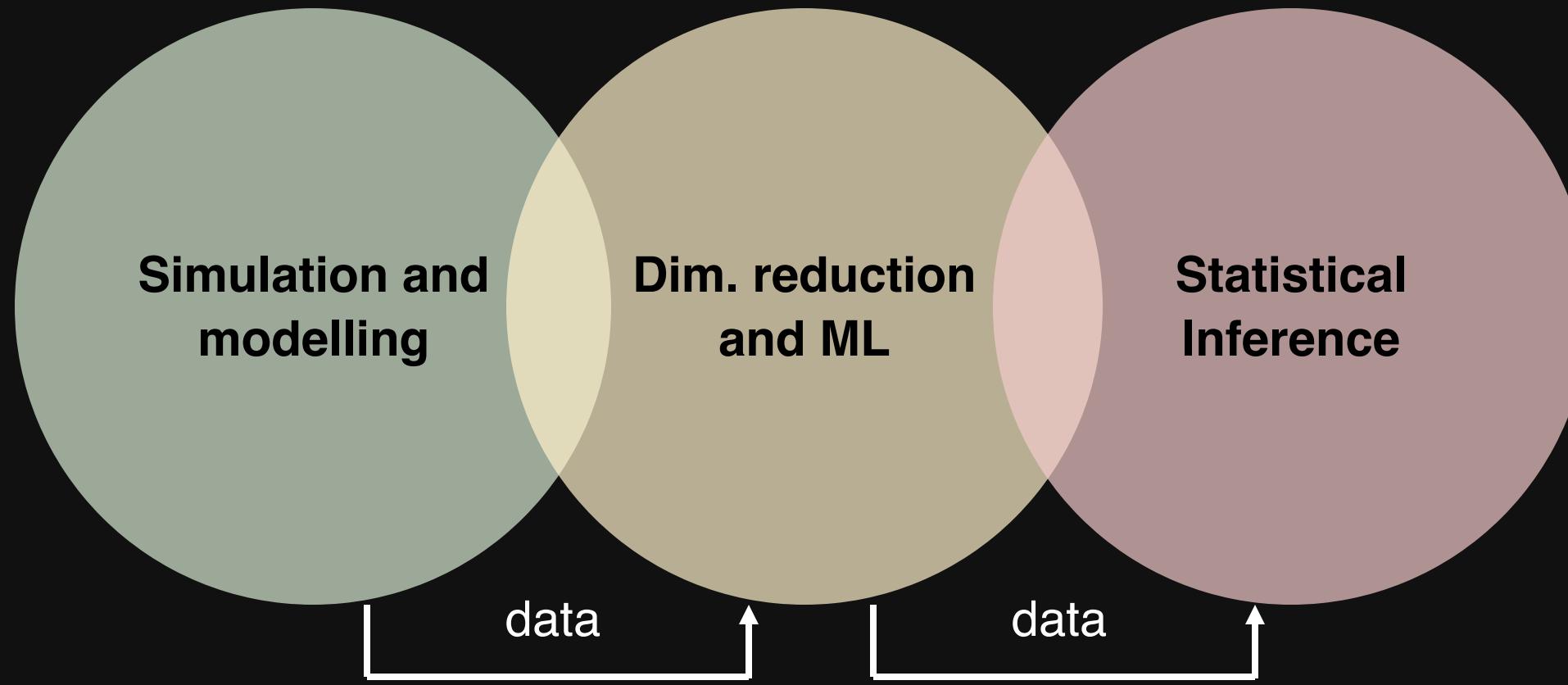


*Upper limit of the usefulness of probabilistic classifiers as  $s(x)$*

# A more holistic approach?

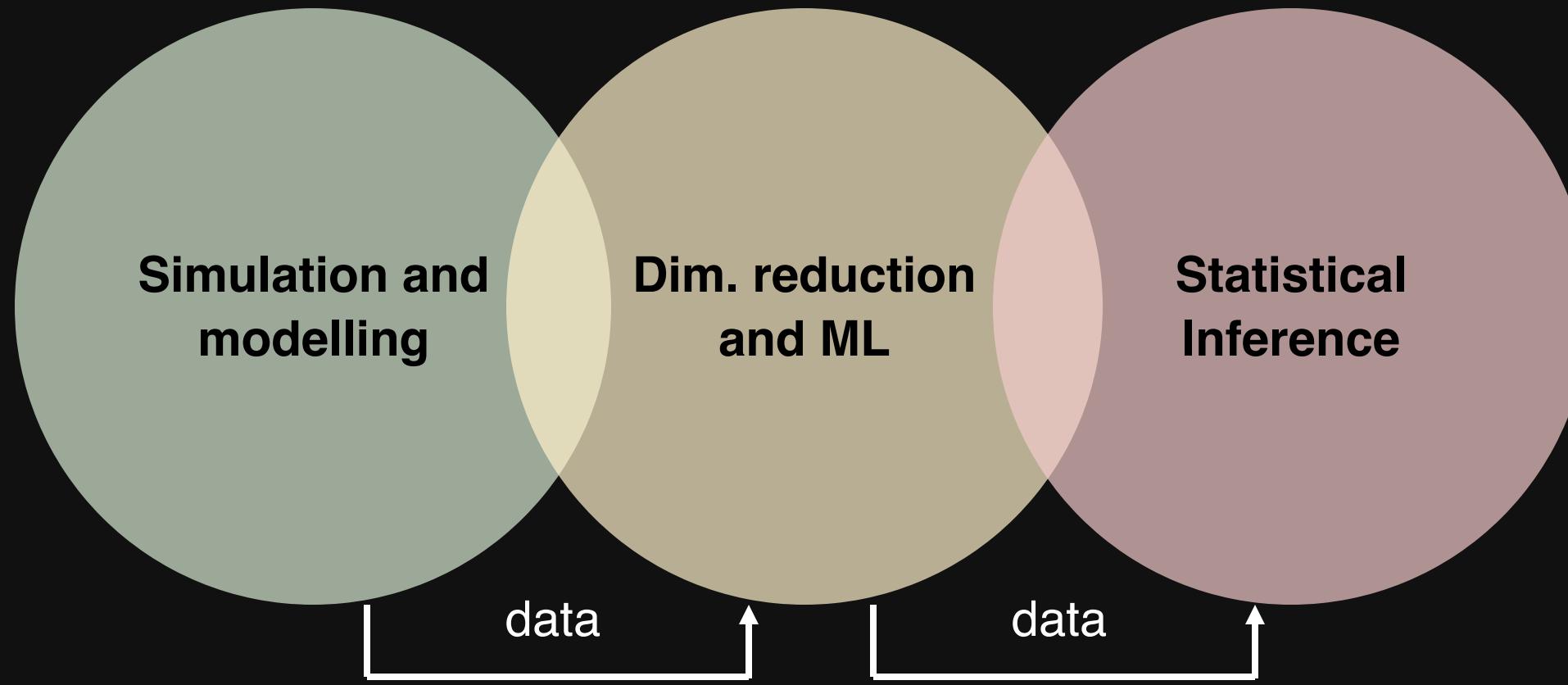


# A more holistic approach?



Optimise directly some **approximation of the expected interval width or significance**, accounting for the effect of nuisance parameters, taking advantage of modern auto-diff (e.g. TensorFlow or PyTorch)

# A more holistic approach?



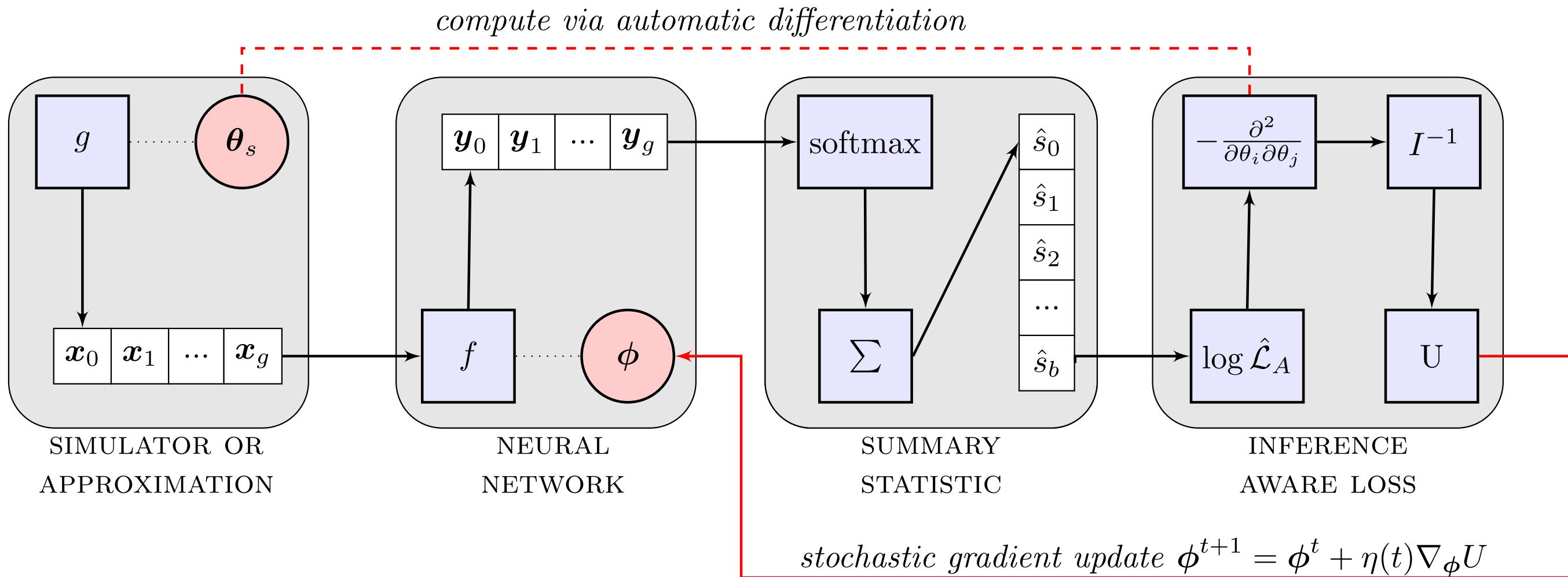
Optimise directly some **approximation of the expected interval width or significance**, accounting for the effect of nuisance parameters, taking advantage of modern auto-diff (e.g. TensorFlow or PyTorch)

**Related alternative techniques (not an exhaustive list):**

- *Learning to Pivot with Adversarial Neural Networks* [G. Louppe et al [arxiv:1611.01046](https://arxiv.org/abs/1611.01046)]
- *Mining gold from implicit models to improve likelihood-free inference* [J. Brehmer et al [arxiv:1805.12244](https://arxiv.org/abs/1805.12244)]
- *Nuisance hardened data compression for fast likelihood-free inference* [J. Alsing et al [arxiv:1903.01473](https://arxiv.org/abs/1903.01473)]
- *Automatic physical inference with information maximising neural networks* [T. Charnock et al [arxiv:1802.03537](https://arxiv.org/abs/1802.03537)]

# INFERence-Aware Neural Optimisation

An approach to learn non-linear summary statistics  $s(D)$  by **directly minimising an approximation the expected profiled (or marginalised) interval width** accounting for the effect of nuisance parameters



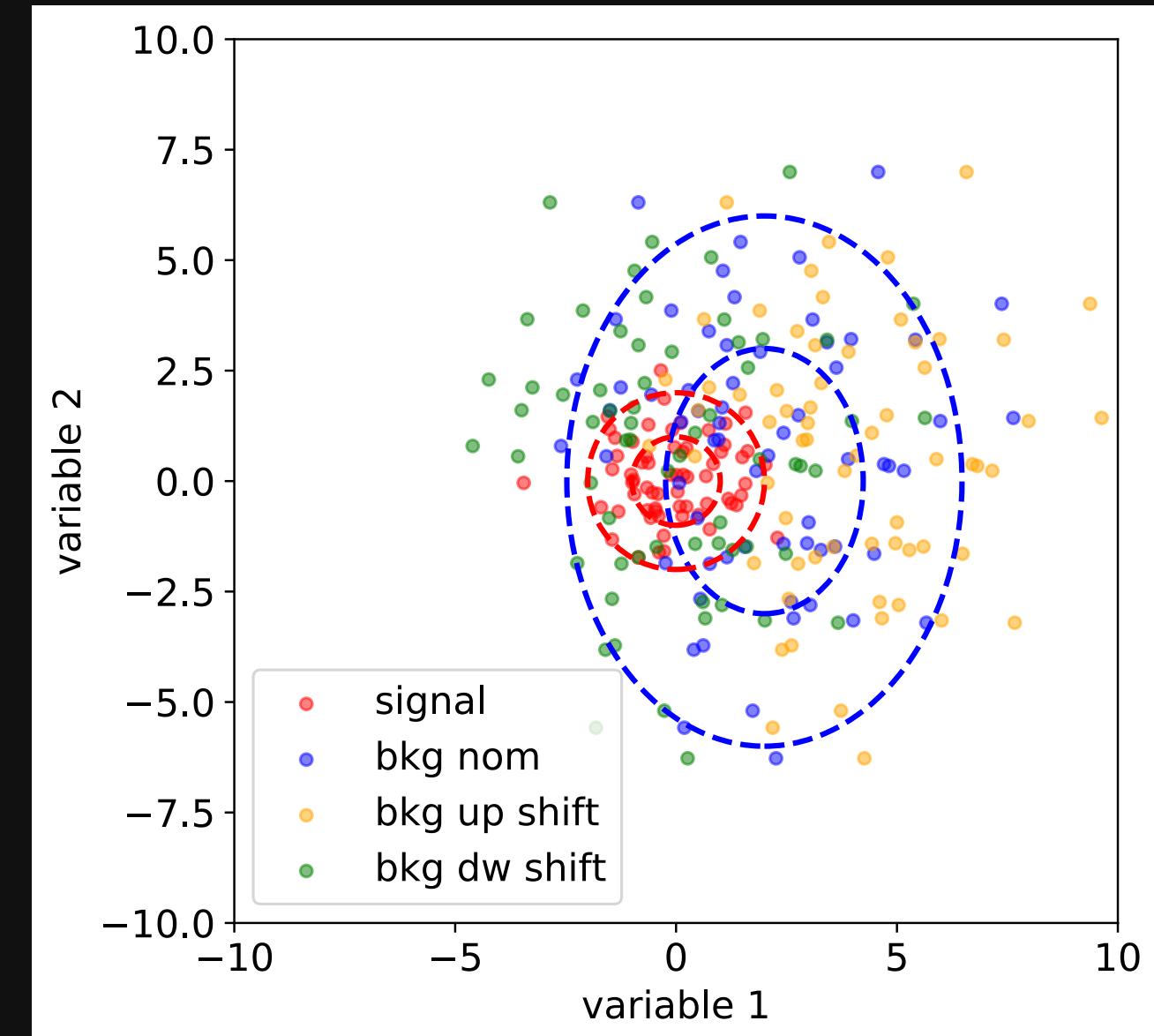
check [arxiv.org/abs/1806.04743](https://arxiv.org/abs/1806.04743) for a more detailed description

# Variational Effect of Parameters $\theta = (\eta, \nu)$

Differentiable approximation of the effect of parameters  $\theta = (\eta, \nu)$  over a given simulated observation  $(\mathbf{x}_i, \mathbf{z}_i, w_i)$

$$(\mathbf{x}_i, \mathbf{z}_i, w_i) \xrightarrow{g(\theta_s)} (\mathbf{x}'_i, \mathbf{z}'_i, w'_i)$$

a non-linear function that depends on the particularities of the problem



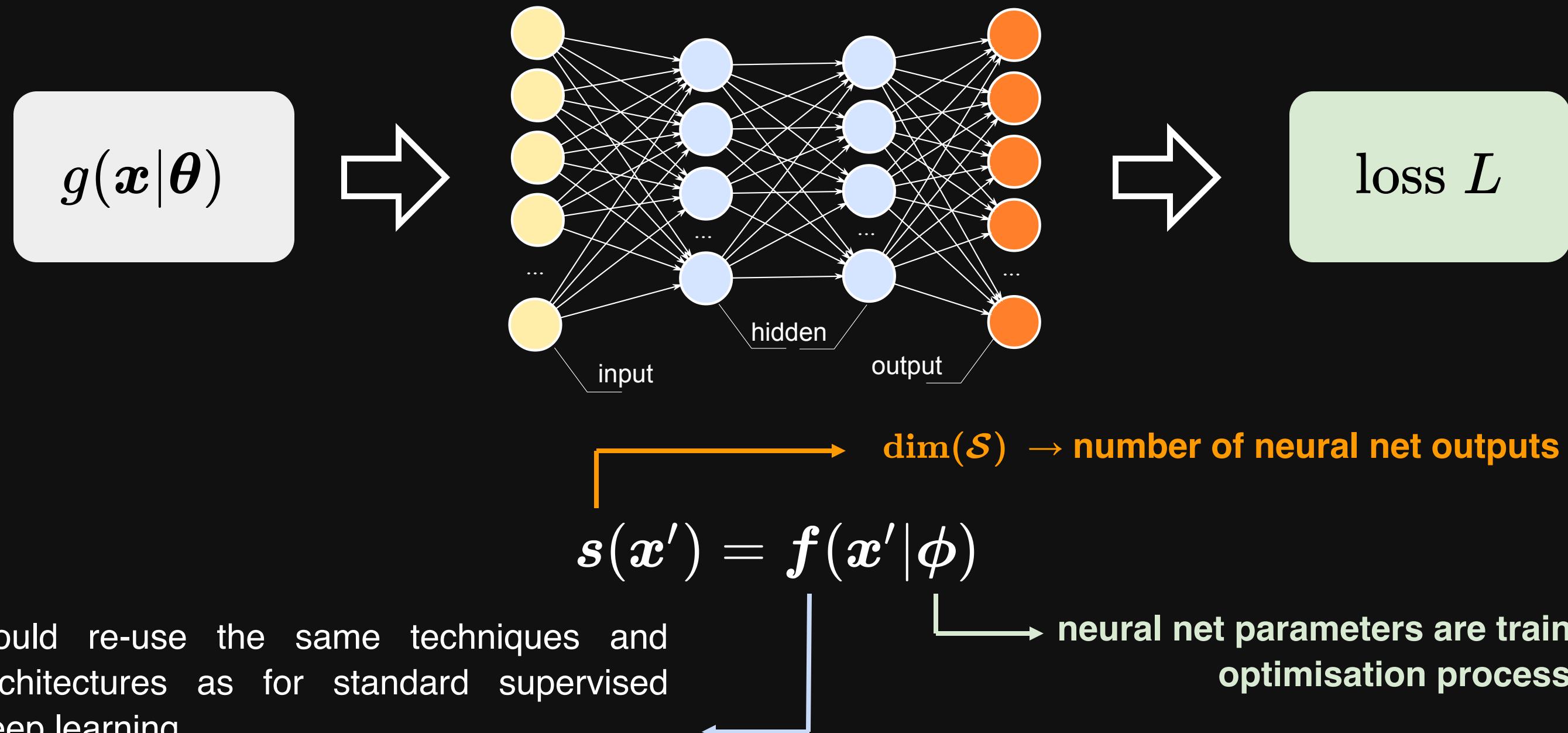
Implemented in autodiff framework (e.g. TensorFlow) to easily obtain derivatives

Alternatively could use a generative surrogate conditioned on the parameters  $\theta$

*Illustrative example based on 3D example*

$$(x'_0, x'_1) = (x_0, x_1) + r(1, 0)$$

# Neural net: a learnable non-linear $s(\mathbf{x})$

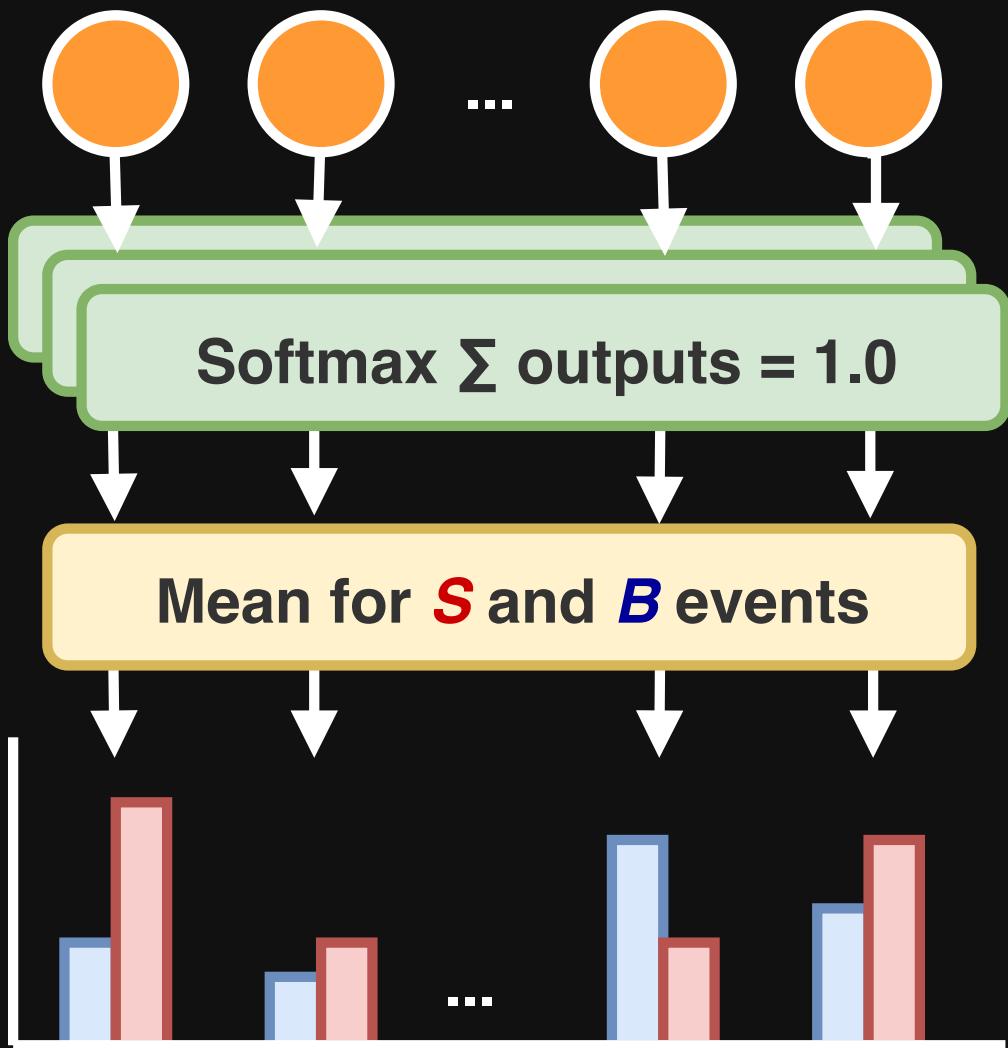


Could re-use the same techniques and architectures as for standard supervised deep learning

A two-hidden layer MLP (100 units each, ReLU activation, He normal init) used for synthetic examples in this work

# Building a non-parametric likelihood

$$\boldsymbol{x}'_i \longrightarrow s(\boldsymbol{x}'_i) \longrightarrow s(\mathbf{D}) \longrightarrow \mathcal{L}(\theta, \eta; \phi)$$



**current choice** → approximate a histogram-like summary statistic  $s(\mathbf{D})$   
*applying softmax to each observation output from the neural net and summing over each dataset*

$$\mathcal{L}(\theta; s(\mathbf{G}_s | \phi)) = \prod_{i \in \text{bins}} \text{Pois}(n_i | n_i^s + n_i^b)$$

The likelihood depends both on the neural network parameters  $\phi$  and the statistical model parameters  $(\eta, \nu)$

Alternative approaches to build  $\mathcal{L}(\theta, \eta; \phi)$  are possible (e.g. non-parametric likelihood of  $s(\mathbf{x})$  using kernel density estimation)

# Inference-motivated Loss Function

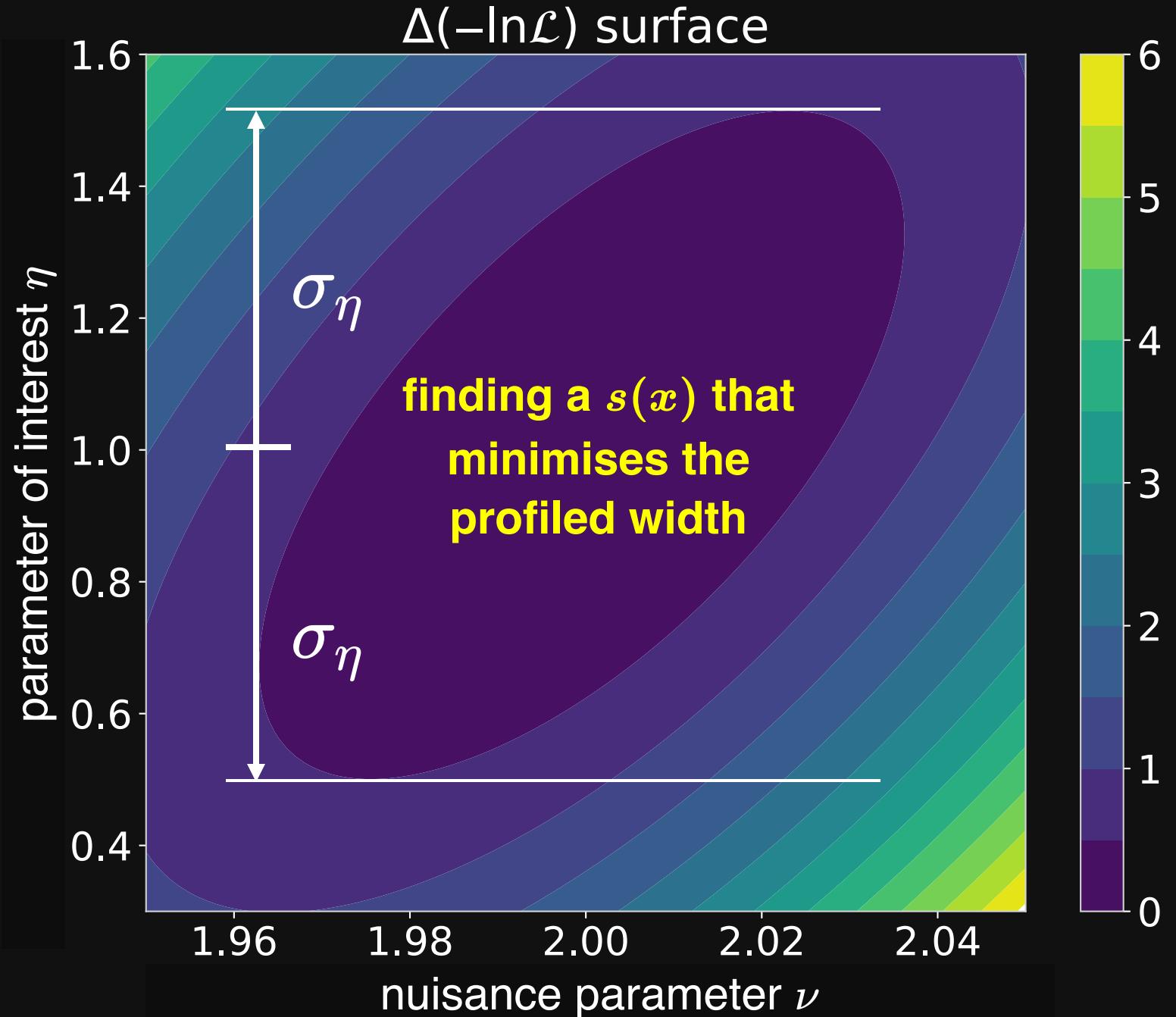
If we expand negative log-likelihood around known minimum (i.e. assuming data equals expectation  $n_i = n_i^s + n_i^b$ ) :

$$\text{covariance} \approx \mathbf{H}^{-1}(-\ln \mathcal{L})$$

can use as loss function directly the approximate variance estimator on the pars of interest:

$$\text{loss} \approx \text{Var}(\eta) \quad (\text{expected})$$

that accounts for the effect of unknowns nuisance parameters. Could be Fisher sub-determinant for multiple p.o.i.



*Equivalent to Laplace approx. in Bayesian Inference (i.e. loss  $\sim$  marginalised posterior width)*

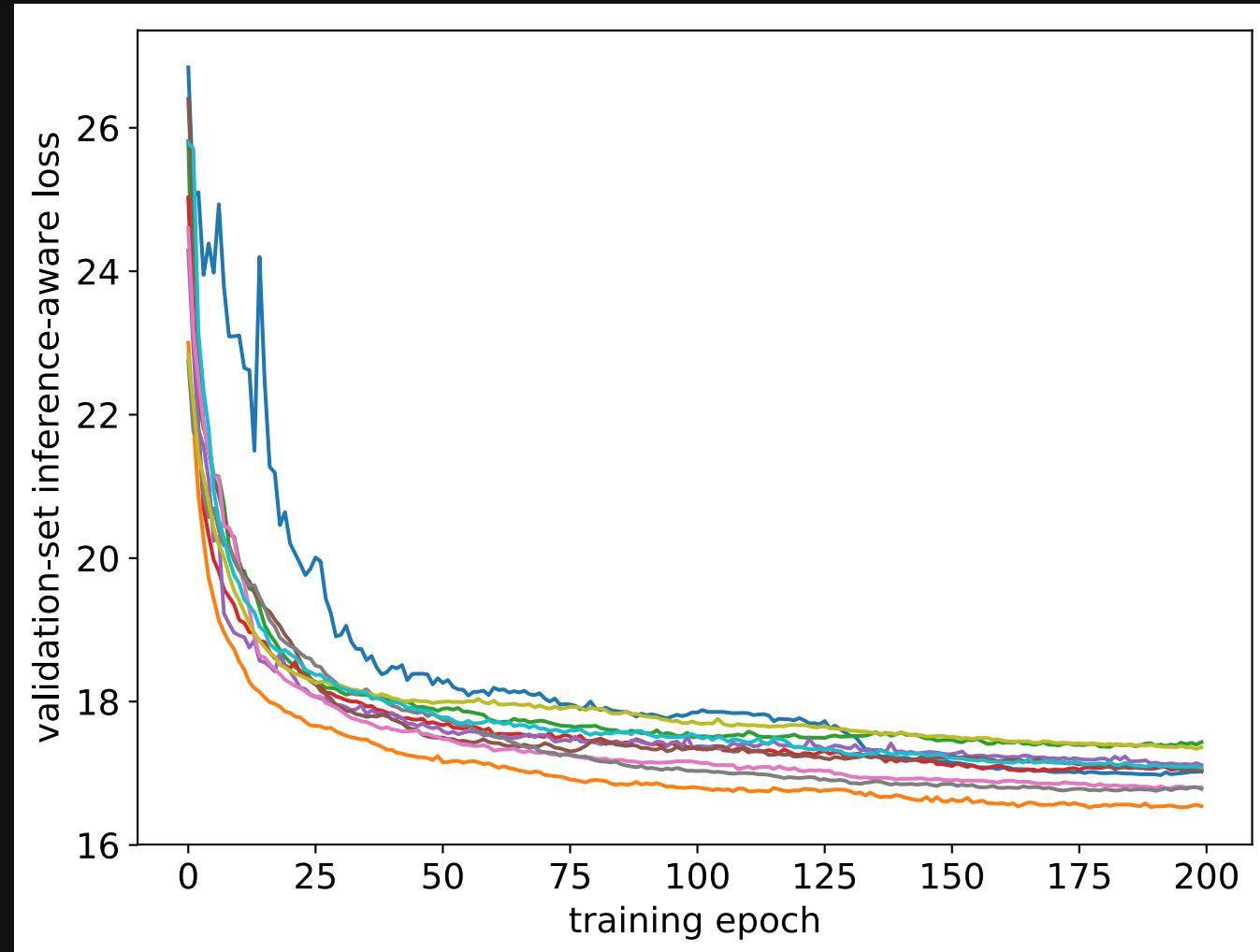
# Synthetic Inference Benchmarks

Several inference problems about  $s = \mu b / (1 - \mu)$ , with different parameter constraints, are considered based on the 3D problem mentioned before

	Benchmark 0	Benchmark 1	Benchmark 2	Benchmark 3	Benchmark 4
interest pars	1 ( $s$ )	1 ( $s$ )	1 ( $s$ )	1 ( $s$ )	1 ( $s$ )
nuisance pars	0 (all fixed)	1 ( $r$ )	2 ( $r$ and $\lambda$ )	2 ( $r$ and $\lambda$ )	3 ( $r$ , $\lambda$ and $b$ )
$r$ (bkg shift)	0.0 (fixed)	free (init 0.0)	free (init 0.0)	$\mathcal{N}(\lambda 0.0, 0.4)$	$\mathcal{N}(\lambda 0.0, 4.0)$
$\lambda$ (bkg exp rate)	3.0 (fixed)	3.0 (fixed)	free (init 3.0)	$\mathcal{N}(\lambda 3.0, 1.0)$	$\mathcal{N}(\lambda 3.0, 1.0)$
$b$ (bkg normalisation)	1000 (fixed)	1000 (fixed)	1000 (fixed)	1000 (fixed)	$\mathcal{N}(b 1000, 100)$

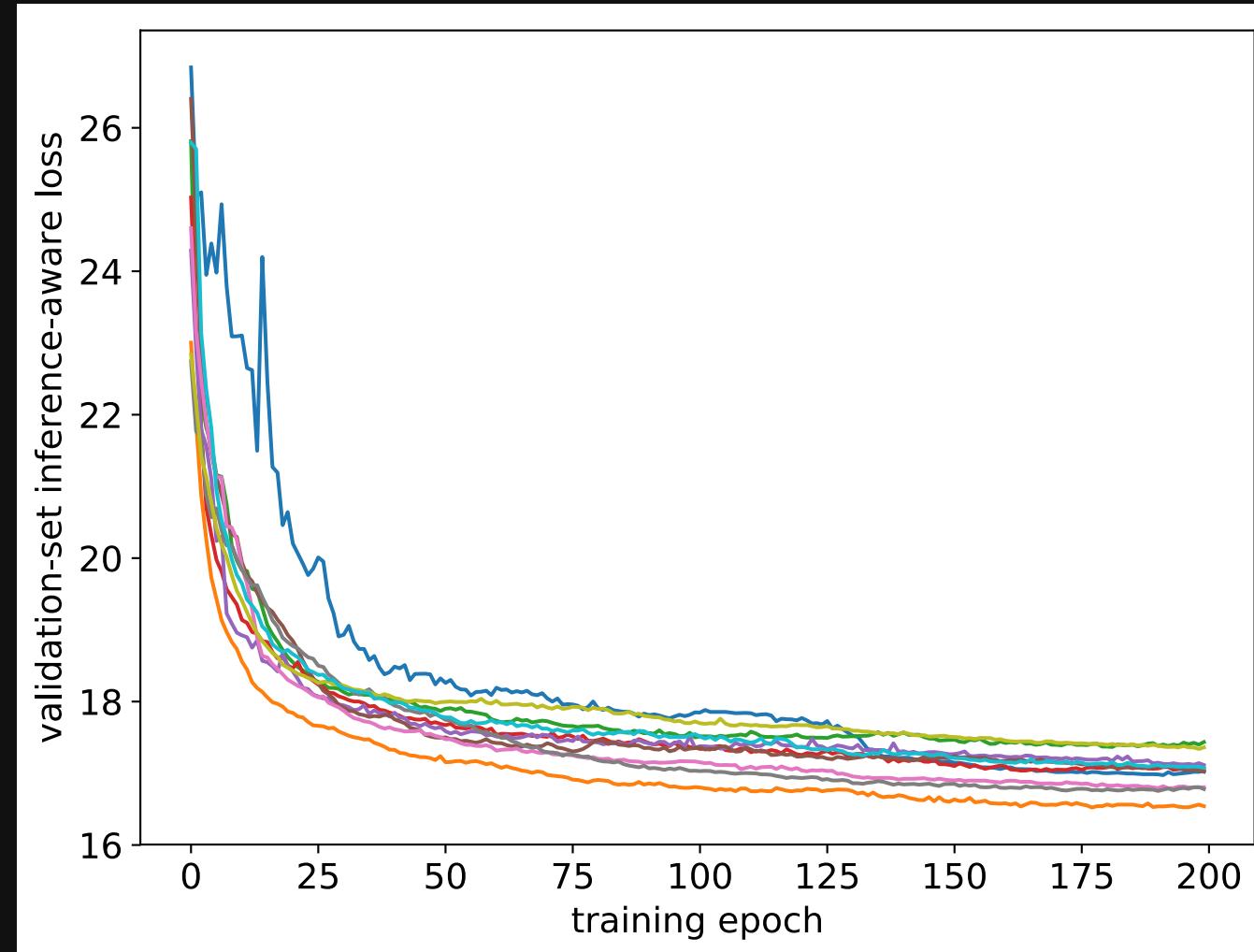
**Information about the inference problem at hand can be used within INFERNO but not with standard probabilistic classifiers**

# Technique works! (results for Benchmark 2)

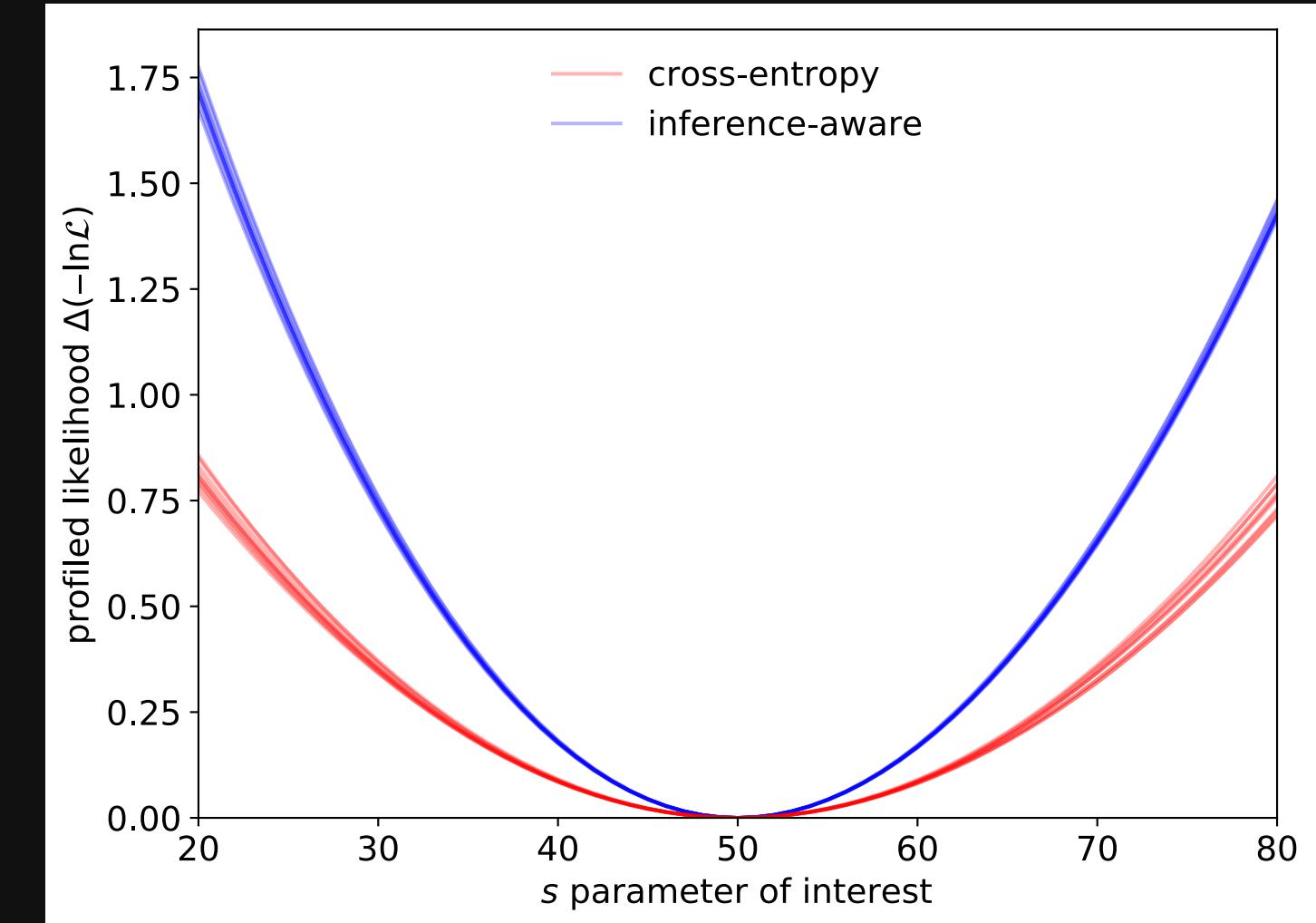


During training, INFERNO consistently converges to low-variance summary statistics

# Technique works! (results for Benchmark 2)



During training, INFERNO consistently converges to low-variance summary statistics



Clearly outperforms classifiers when nuisance parameters are relevant

# Comparison w/ Classification Approach

A more systematic comparison, shows that INFERNO clearly outperforms any classifier (even optimal Bayes) when nuisance parameters are relevant

(less is better)

Table 1: Expected uncertainty on the parameter of interest  $s$  for each of the inference benchmarks considered using a cross-entropy trained neural network model, INFERNO customised for each problem and the optimal classifier and likelihood based results.

	Benchmark 0	Benchmark 1	Benchmark 2	Benchmark 3	Benchmark 4
NN classifier	$14.99^{+0.02}_{-0.00}$	$18.94^{+0.11}_{-0.05}$	$23.94^{+0.52}_{-0.17}$	$21.54^{+0.27}_{-0.05}$	$26.71^{+0.56}_{-0.11}$
INFERNO 0	<b><math>15.51^{+0.09}_{-0.02}</math></b>	$18.34^{+5.17}_{-0.51}$	$23.24^{+6.54}_{-1.22}$	$21.38^{+3.15}_{-0.69}$	$26.38^{+7.63}_{-1.36}$
INFERNO 1	$15.80^{+0.14}_{-0.04}$	<b><math>16.79^{+0.17}_{-0.05}</math></b>	$21.41^{+2.00}_{-0.53}$	$20.29^{+1.20}_{-0.39}$	$24.26^{+2.35}_{-0.71}$
INFERNO 2	$15.71^{+0.15}_{-0.04}$	$16.87^{+0.19}_{-0.06}$	<b><math>16.95^{+0.18}_{-0.04}</math></b>	$16.88^{+0.17}_{-0.03}$	$18.67^{+0.25}_{-0.05}$
INFERNO 3	$15.70^{+0.21}_{-0.04}$	$16.91^{+0.20}_{-0.05}$	$16.97^{+0.21}_{-0.04}$	<b><math>16.89^{+0.18}_{-0.03}</math></b>	$18.69^{+0.27}_{-0.04}$
INFERNO 4	$15.71^{+0.32}_{-0.06}$	$16.89^{+0.30}_{-0.07}$	$16.95^{+0.38}_{-0.05}$	$16.88^{+0.40}_{-0.05}$	<b><math>18.68^{+0.58}_{-0.07}</math></b>
Optimal classifier	14.97	19.12	24.93	22.13	27.98
Analytical likelihood	14.71	15.52	15.65	15.62	16.89

$s_{\text{clf}}(x)$  ——————→

upper bound if we knew  $p(x|\theta)$

classifier is expected to be near optimal sufficient  $s(x)$

INFERNO adapted for each inference problem is considerable better than classifiers

# Final Remarks

Alternative ways to construct summary statistics in cases where nuisance parameters are important could greatly increase the discovery reach of scientific experiments based on simulation-based inference

The proposed INFERNO technique obtains non-linear summary statistics by minimising the expected uncertainty accounting for the effect of nuisance parameters

Early results are really promising, yet benchmarking in more complex problems and comparisons with alternative techniques are needed to shed more light on its real-world usefulness

# More details in our ArXiv preprint

The screenshot shows the arXiv.org preprint page for the paper "INFERNO: Inference-Aware Neural Optimisation" by Pablo de Castro and Tommaso Dorigo. The page includes the Cornell University logo, a search bar, and sections for download, references, and citations.

**Statistics > Machine Learning**

**INFERNO: Inference-Aware Neural Optimisation**

Pablo de Castro, Tommaso Dorigo

(Submitted on 12 Jun 2018 (v1), last revised 11 Oct 2018 (this version, v2))

Complex computer simulations are commonly required for accurate data modelling in many scientific disciplines, making statistical inference challenging due to the intractability of the likelihood evaluation for the observed data. Furthermore, sometimes one is interested on inference drawn over a subset of the generative model parameters while taking into account model uncertainty or misspecification on the remaining nuisance parameters. In this work, we show how non-linear summary statistics can be constructed by minimising inference-motivated losses via stochastic gradient descent such they provided the smallest uncertainty for the parameters of interest. As a use case, the problem of confidence interval estimation for the mixture coefficient in a multi-dimensional two-component mixture model (i.e. signal vs background) is considered, where the proposed technique clearly outperforms summary statistics based on probabilistic classification, which are a commonly used alternative but do not account for the presence of nuisance parameters.

Comments: Code available at [this https URL](https://). Version updates: – v2: fixed typos, improve text, link to code and a better synthetic experiment

Subjects: Machine Learning (stat.ML); Machine Learning (cs.LG); High Energy Physics – Experiment (hep-ex); Data Analysis, Statistics and Probability (physics.data-an); Methodology (stat.ME)

Cite as: [arXiv:1806.04743 \[stat.ML\]](https://arxiv.org/abs/1806.04743)  
(or [arXiv:1806.04743v2 \[stat.ML\]](https://arxiv.org/abs/1806.04743v2) for this version)

We gratefully acknowledge support from the Simons Foundation and member institutions.

Search or Article ID All fields

(Help | Advanced search)

**Download:**

- PDF
- Other formats

(license)

**Current browse context:**  
stat.ML  
< prev | next >  
new | recent | 1806

**Change to browse by:**

cs  
  cs.LG  
hep-ex  
physics  
  physics.data-an  
stat  
  stat.ME

**References & Citations**

- INSPIRE HEP  
(refers to | cited by )
- NASA ADS

**3 blog links** (what is this?)

**Google Scholar**

**Bookmark** (what is this?)

many more details on our stat.ML preprint [arxiv.org/abs/1806.04743](https://arxiv.org/abs/1806.04743)

feedback, comments, discussions are welcomed via

Twitter DM @pablodecm, [pablo.decastro@cern.ch](mailto:pablo.decastro@cern.ch) or workshop Gitter