

# Normalizing flows and the likelihood ratio trick in particle physics

Johann Brehmer

New York University

Deep Learning Seminar, Uni Bremen

January 6, 2020

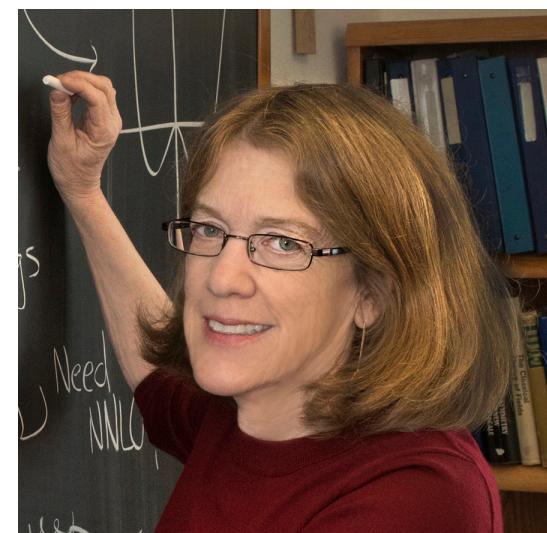
# Collaborators



Kyle Cranmer



Gilles Louppe



The SCAILFIN Project  
[scailfin.github.io](https://scailfin.github.io)



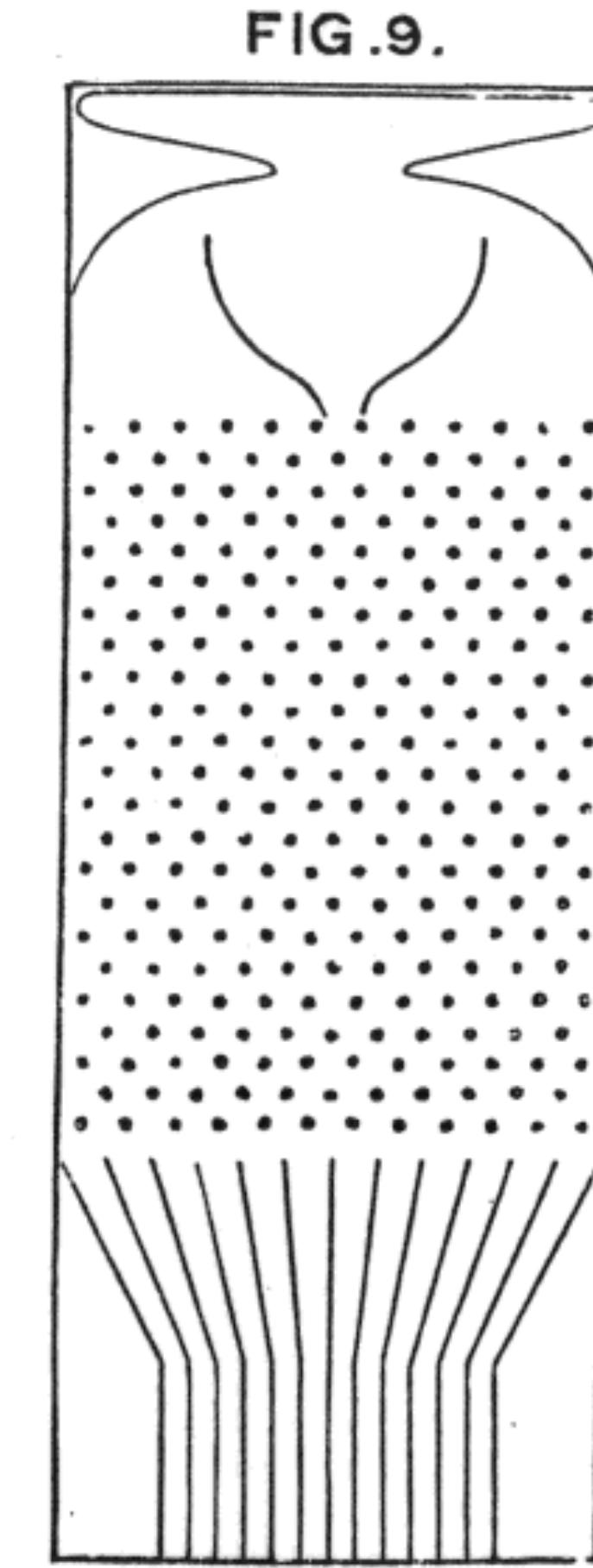
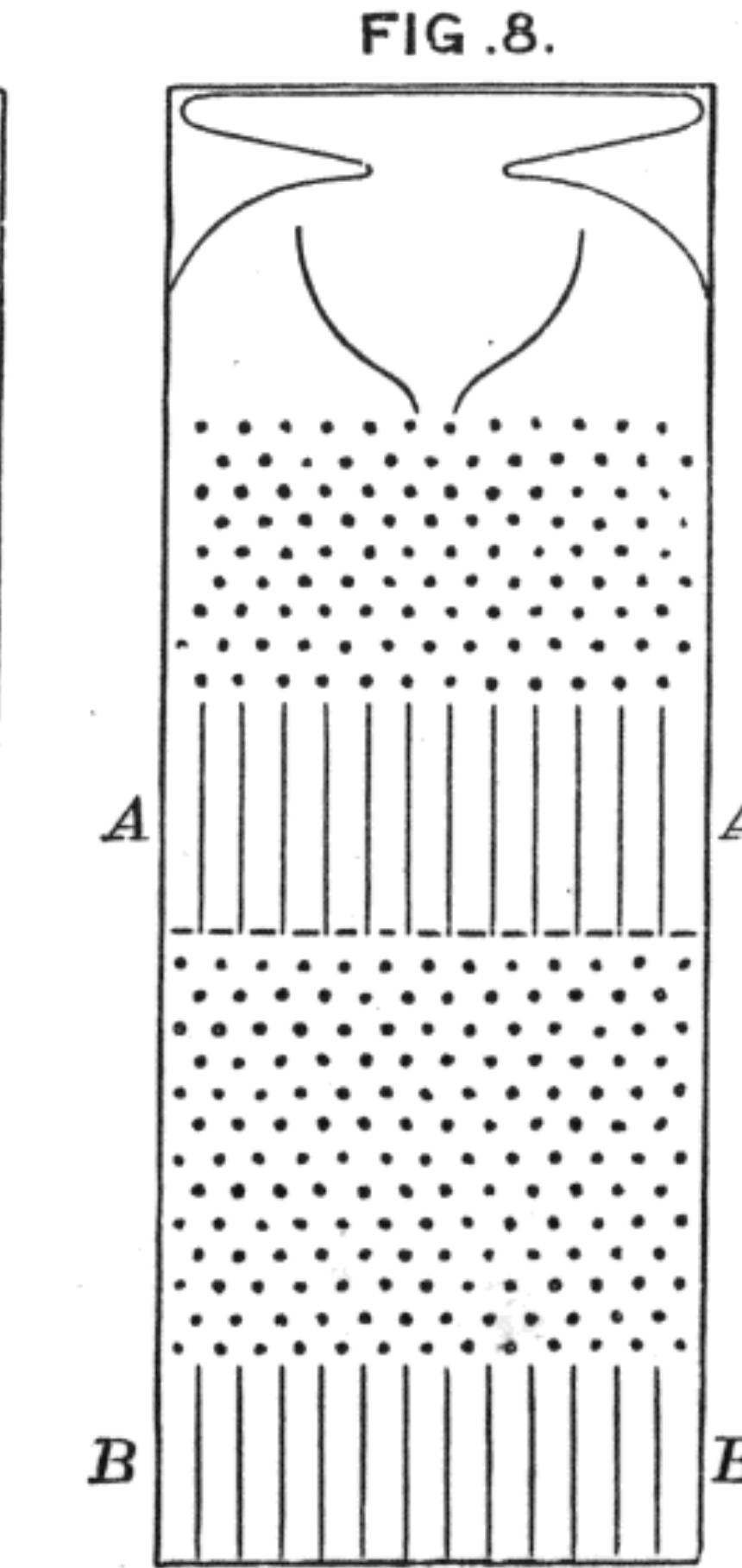
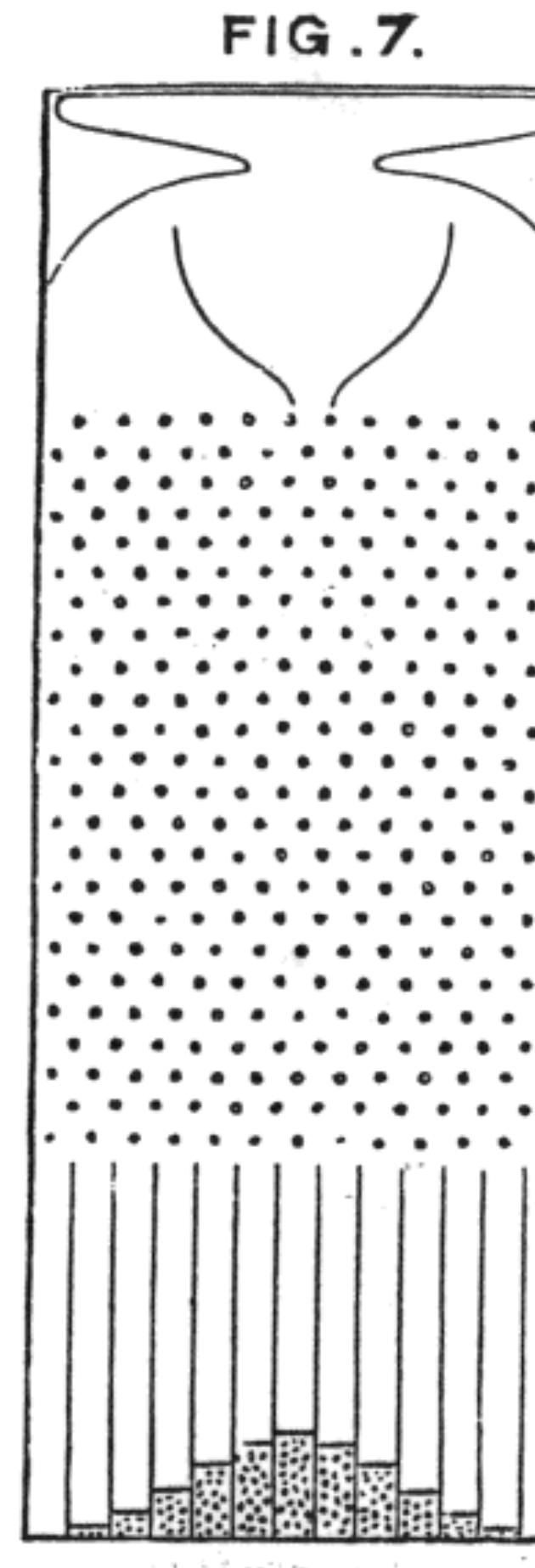
Scientists use simulations to model Nature,  
but that leads to a problem for inference.

# The Galton board



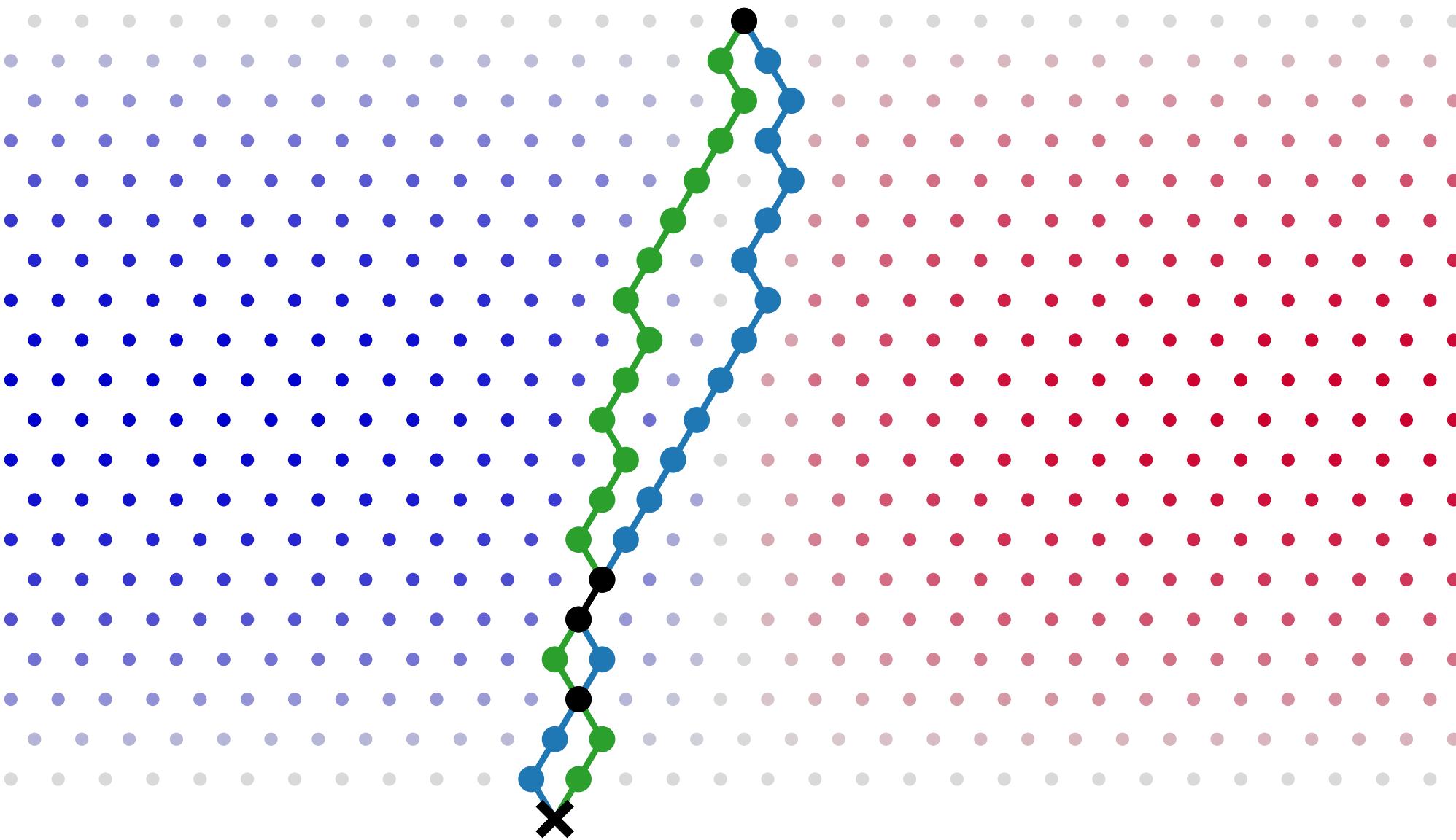
[galtonboard.com]

# The Galton board



[F. Galton 1889]

# Probabilities from integrating trajectories



$$\text{Probability of ending in bin } x : p(x) = \int dz \ p(x, z)$$

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  $\theta$  ?

- **Prediction:** given  $\theta$ , generate samples of observations  $\{x_i\}$ .

Simple: just drop balls!

# 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  $\theta$  ?

- **Prediction:** given  $\theta$ , generate samples of observations  $\{x_i\}$ .

Simple: just drop balls!

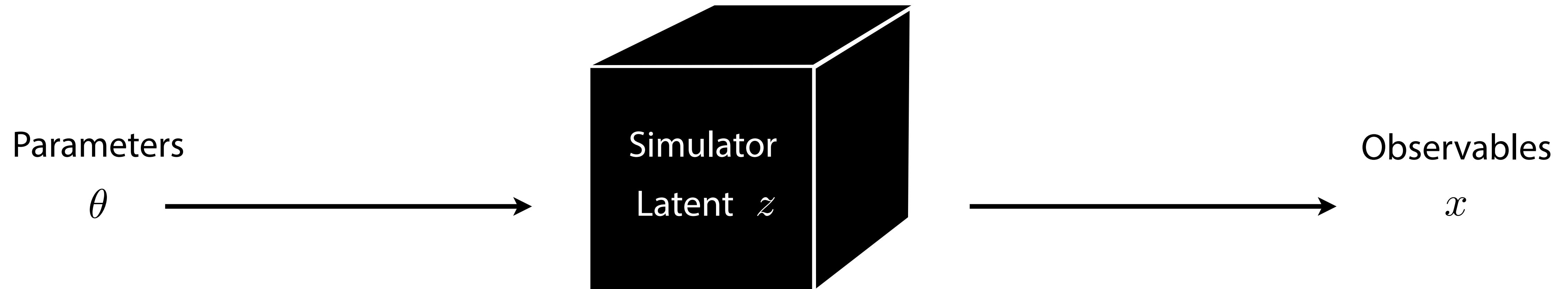
- **Inference:** given observations  $\{x_i\}$ , what are the most likely values for  $\theta$  ?

Usually we solve this with the likelihood

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

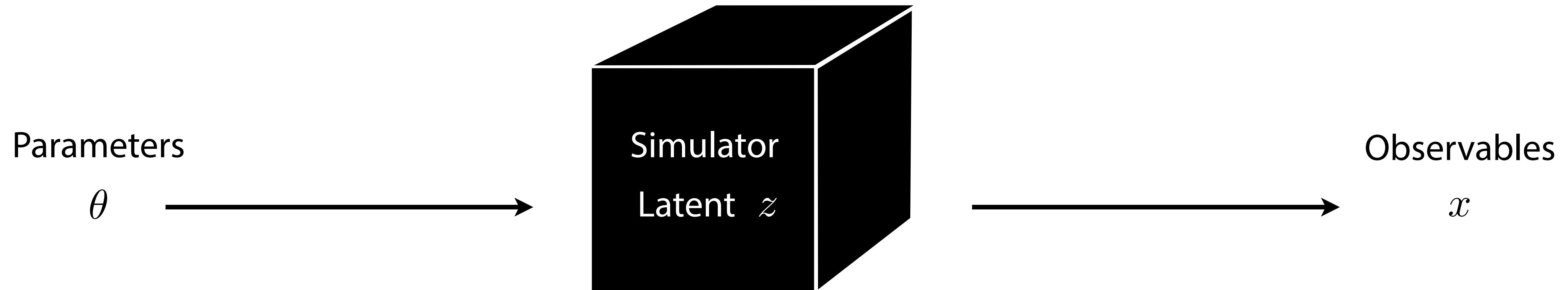
But the number of possible **paths**  $z$  can be huge, and it becomes impossible to calculate the integral!

# Simulation-based (“likelihood-free”) inference



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

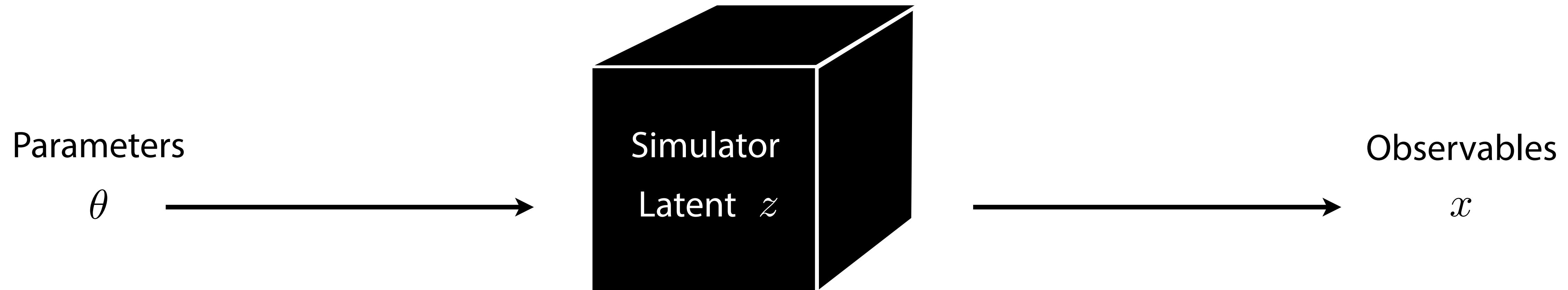
# Simulation-based (“likelihood-free”) inference



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

- 
- Prediction:
- Well-understood mechanistic model
  - Simulator can sample  $x \sim p(x|\theta)$

# Simulation-based (“likelihood-free”) inference

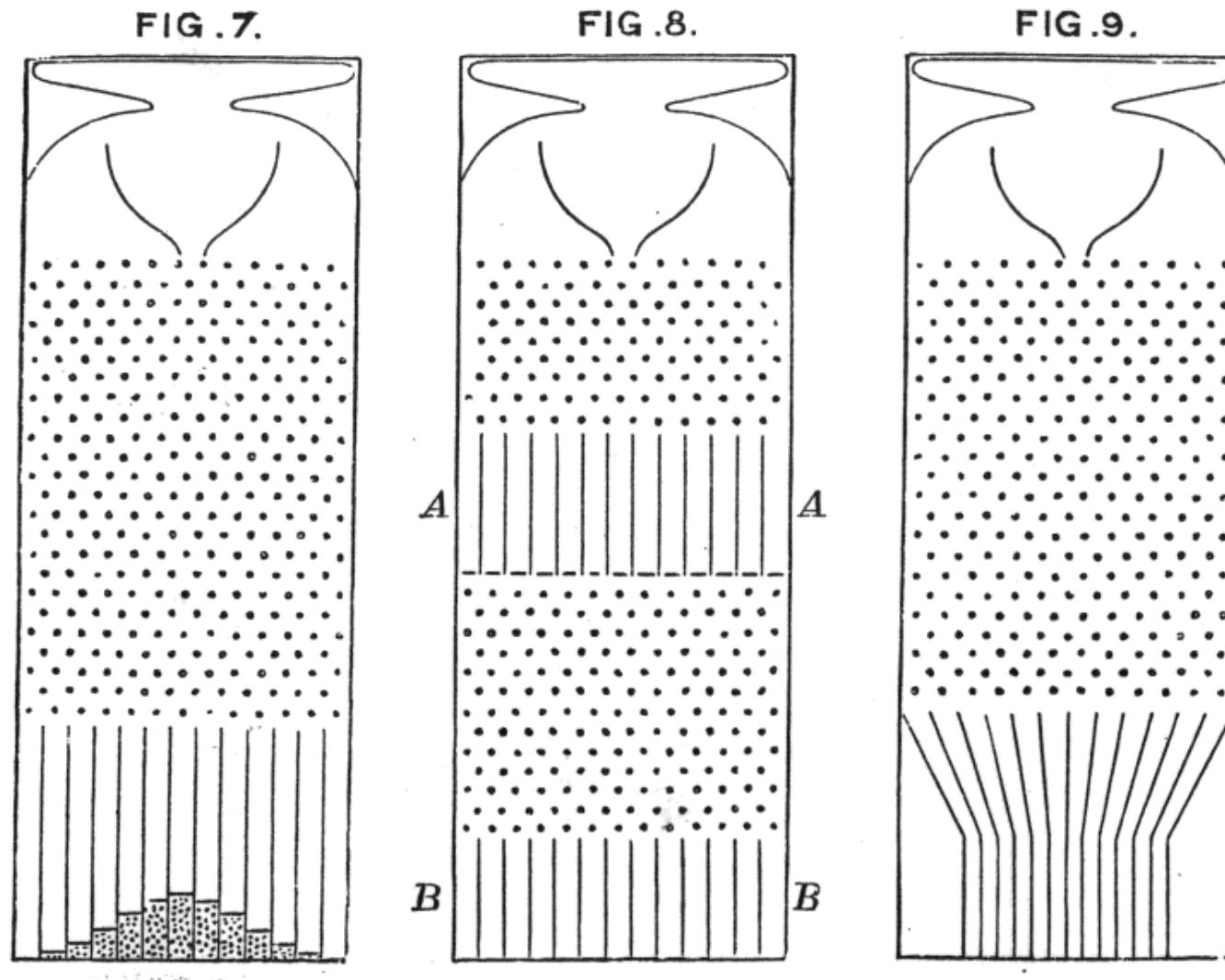


$$p(x|\theta) = \int dz p(x, z|\theta)$$

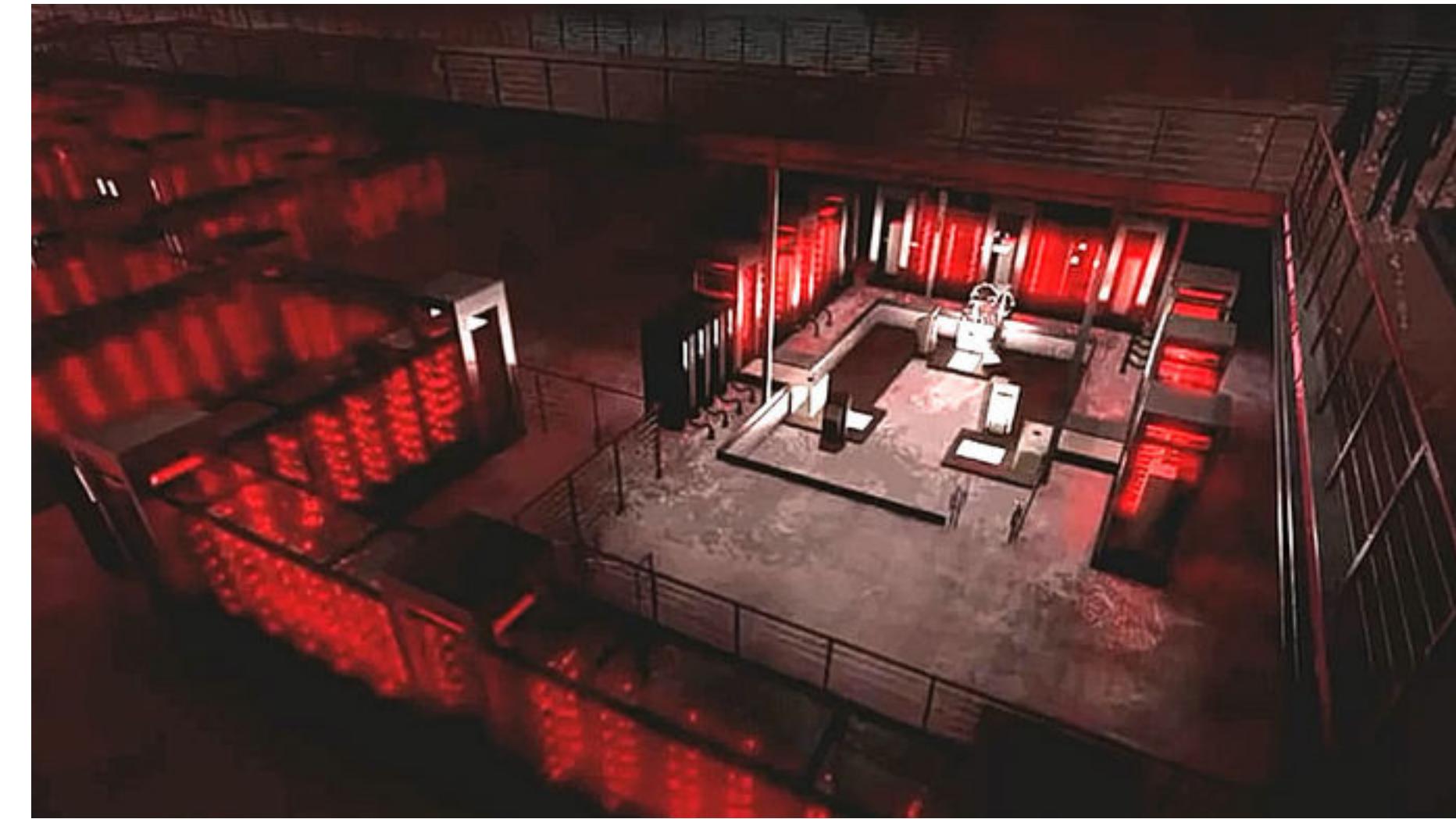
- Prediction:
- Well-understood mechanistic model
  - Simulator can sample  $x \sim p(x|\theta)$

- Inference:
- Likelihood  $p(x|\theta)$  is intractable
  - Inference challenging, needs e.g. estimator  $\hat{p}(x|\theta)$

# Galton board: metaphor for simulation-based science



[F. Galton 1889]



[HBO 2018]

Galton board device



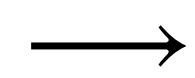
Computer simulation

Parameters  $\theta$



Model parameters  $\theta$

Bins  $x$



Observables  $x$

Path  $z$



Latent variables  $z$

(stochastic execution trace through simulator)

This inverse problem appears for instance in  
particle physics.

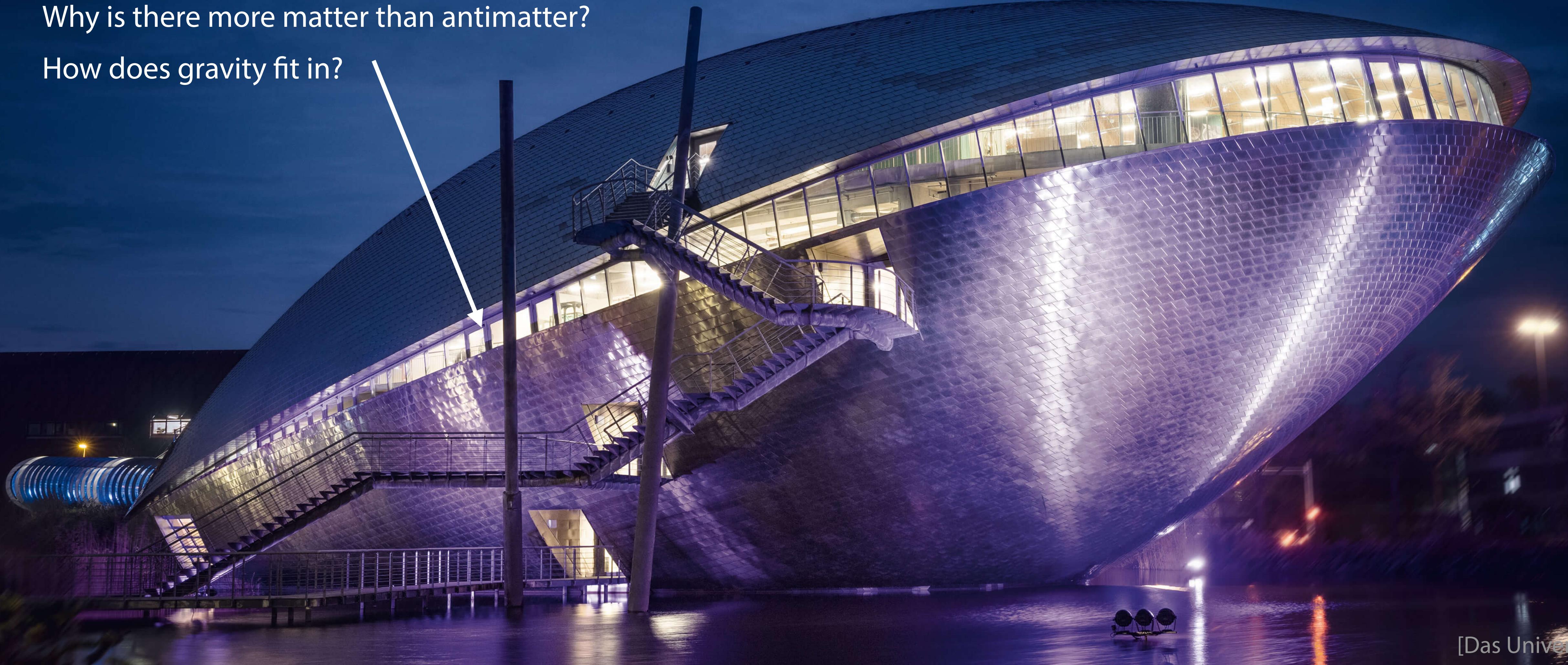
5% visible matter:

The Standard Model passes virtually every test,  
up to  $10^{-10}$  precision.

But it seems fine-tuned.

Why is there more matter than antimatter?

How does gravity fit in?



**5% visible matter:**

The Standard Model passes virtually every test,  
up to  $10^{-10}$  precision.

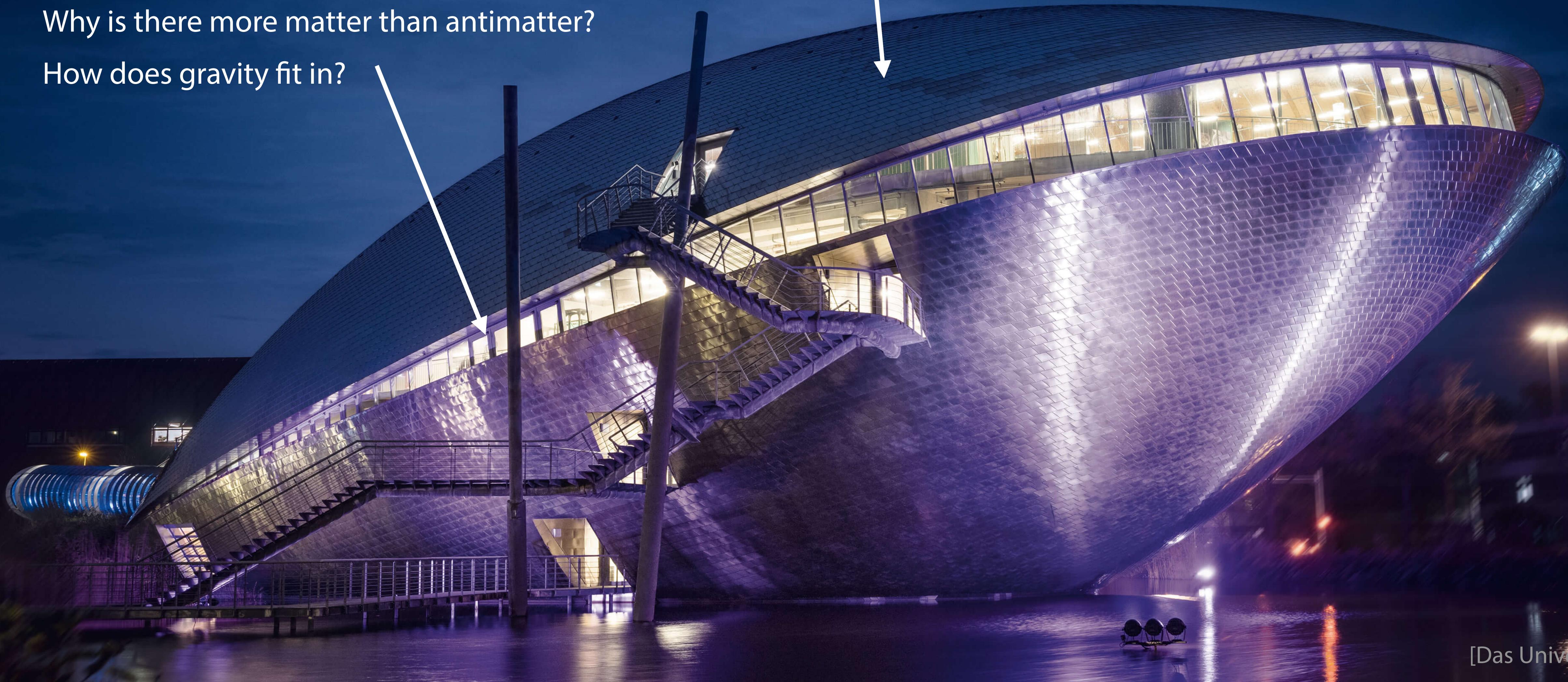
But it seems fine-tuned.

Why is there more matter than antimatter?

How does gravity fit in?

**27% Dark Matter**

What is its (particle) nature?



**5% visible matter:**

The Standard Model passes virtually every test,  
up to  $10^{-10}$  precision.

But it seems fine-tuned.

Why is there more matter than antimatter?

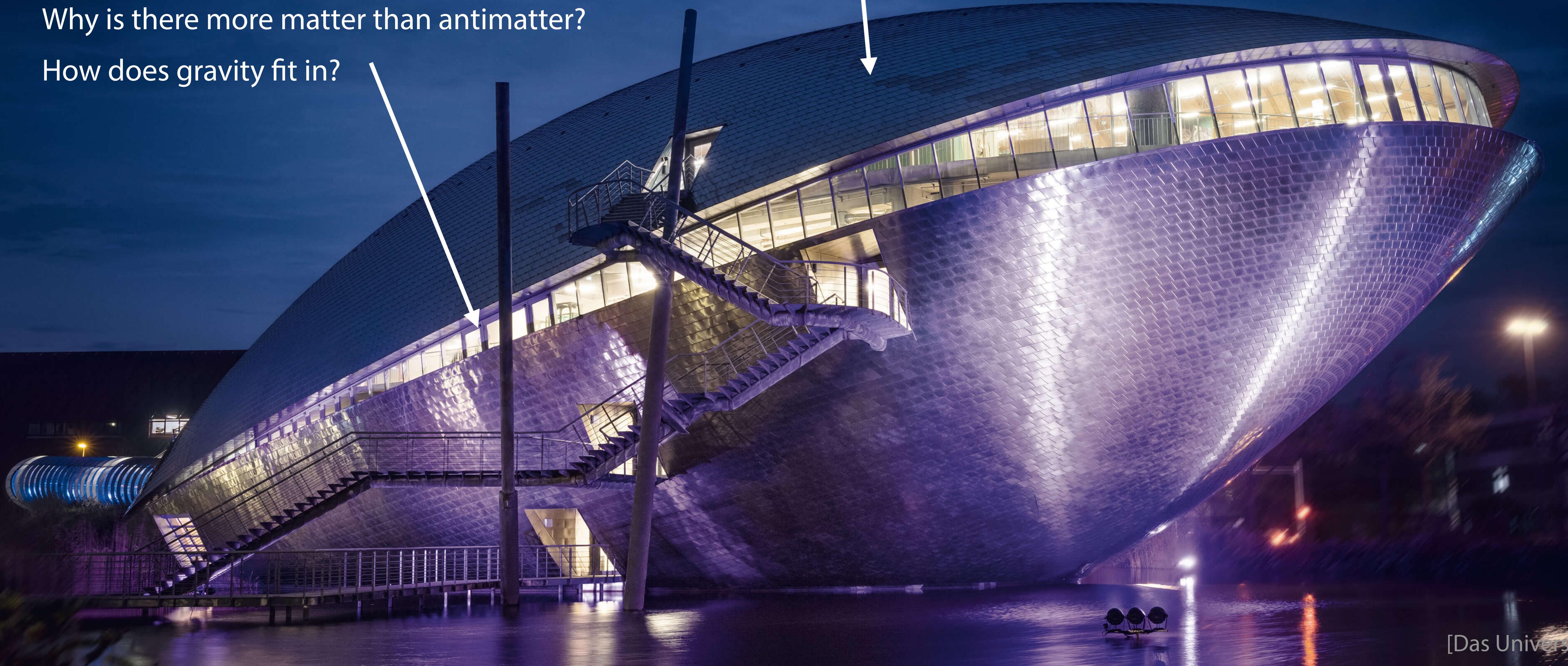
How does gravity fit in?

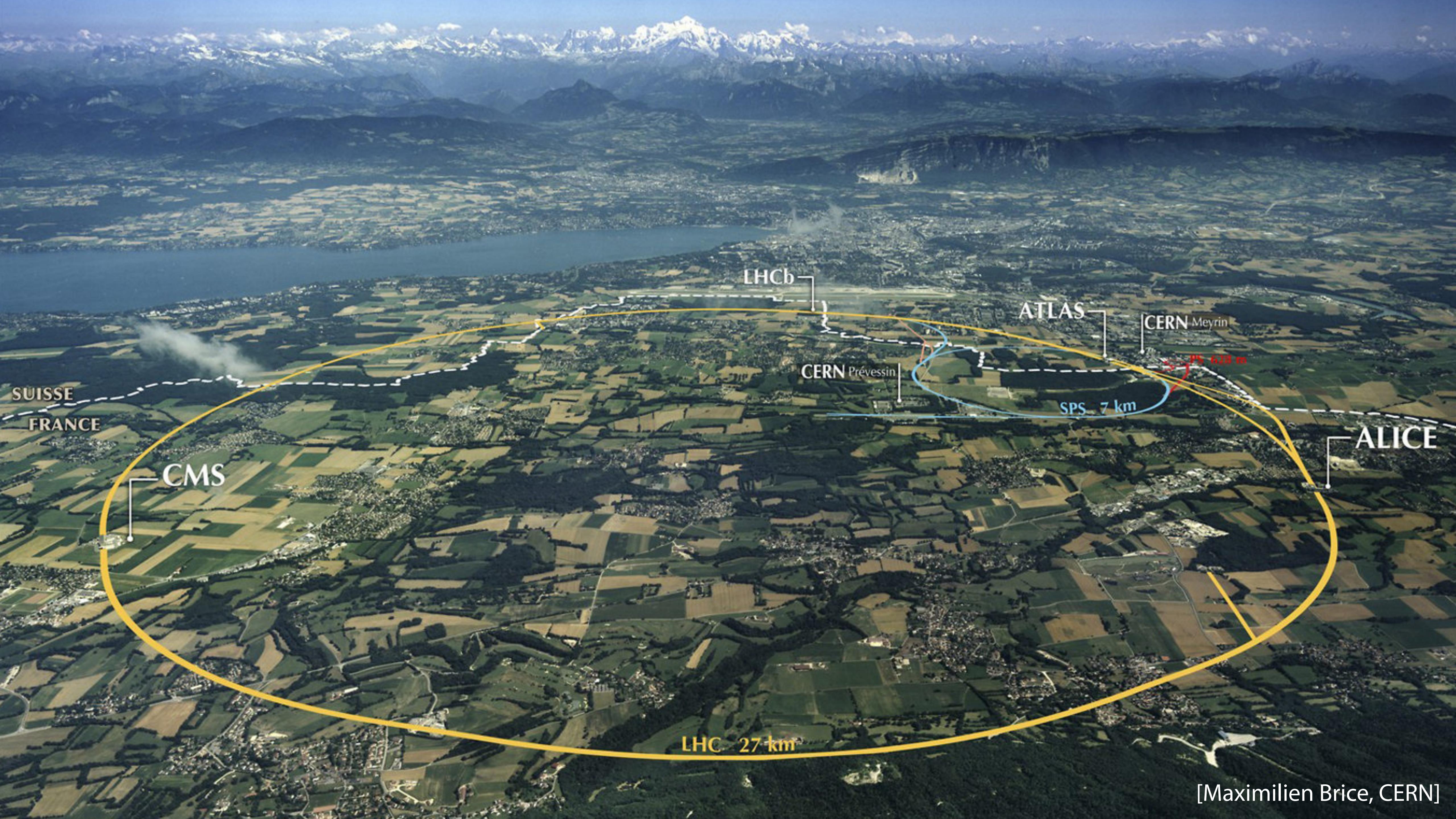
**27% Dark Matter**

What is its (particle) nature?

**68% Dark Energy:**

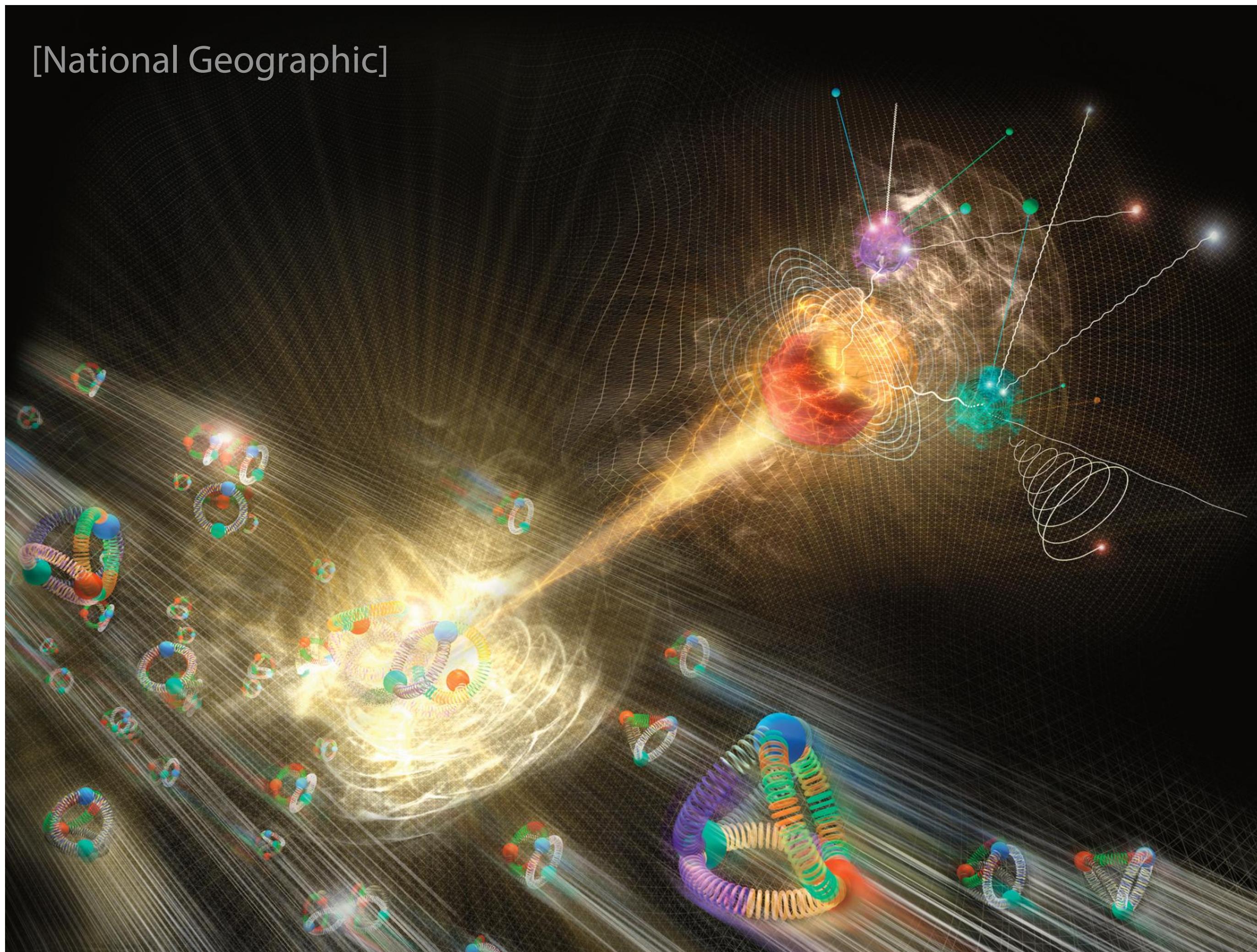
What?



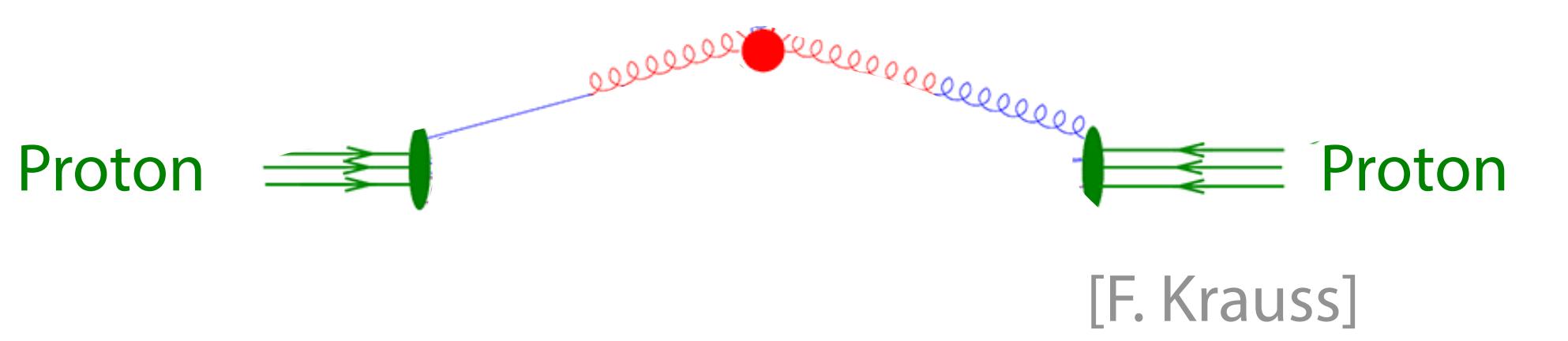


[Maximilien Brice, CERN]

# When protons meet

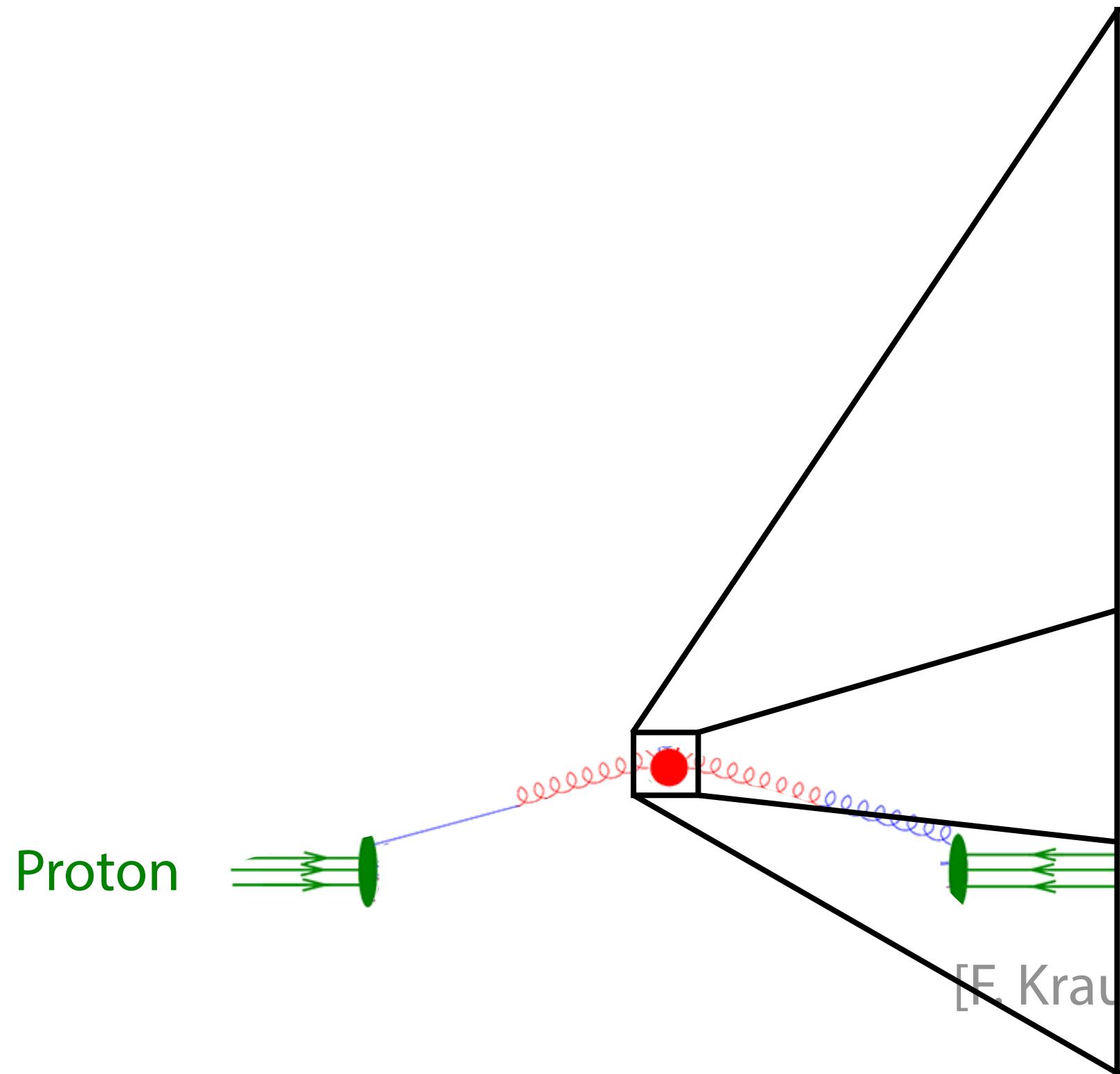


# When protons meet



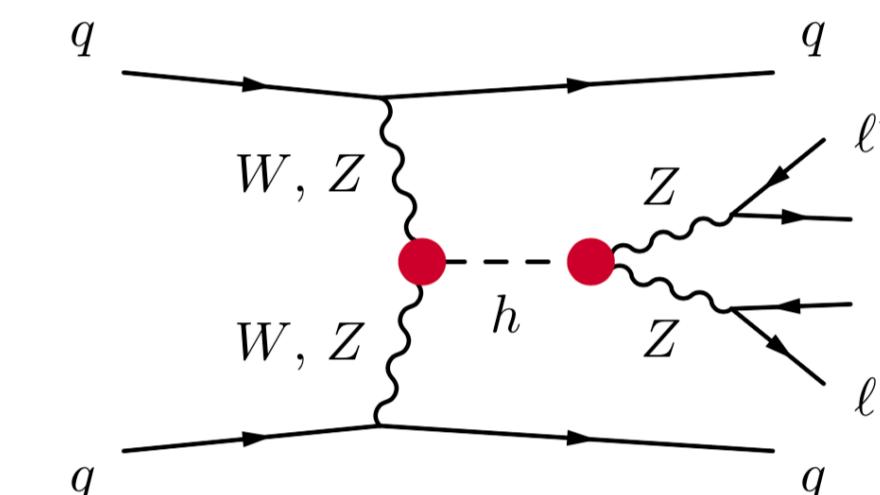
[F. Krauss]

# When protons meet



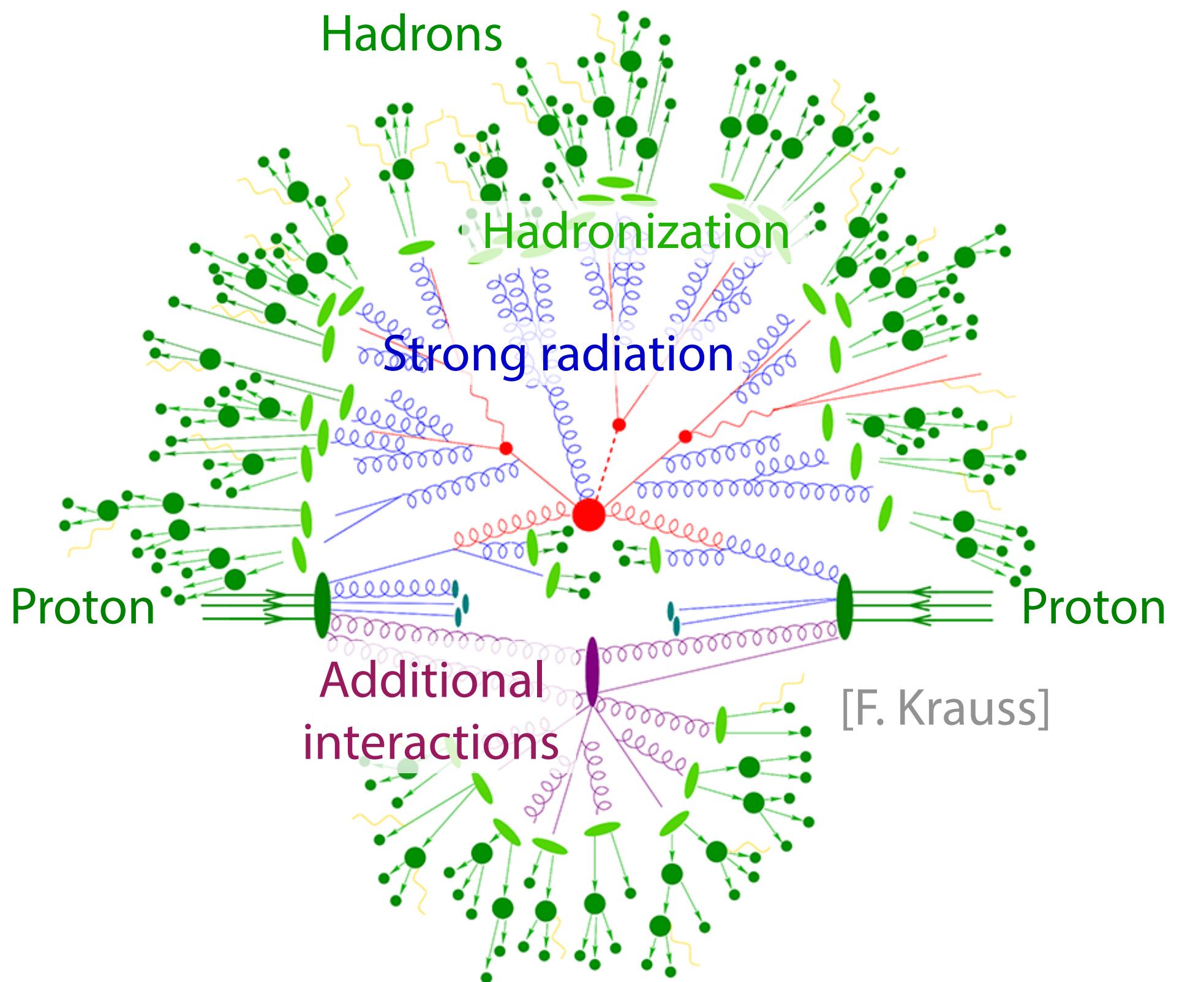
## High-energy interactions:

- Rigorously described by quantum field theories
- Sensitive to new physics

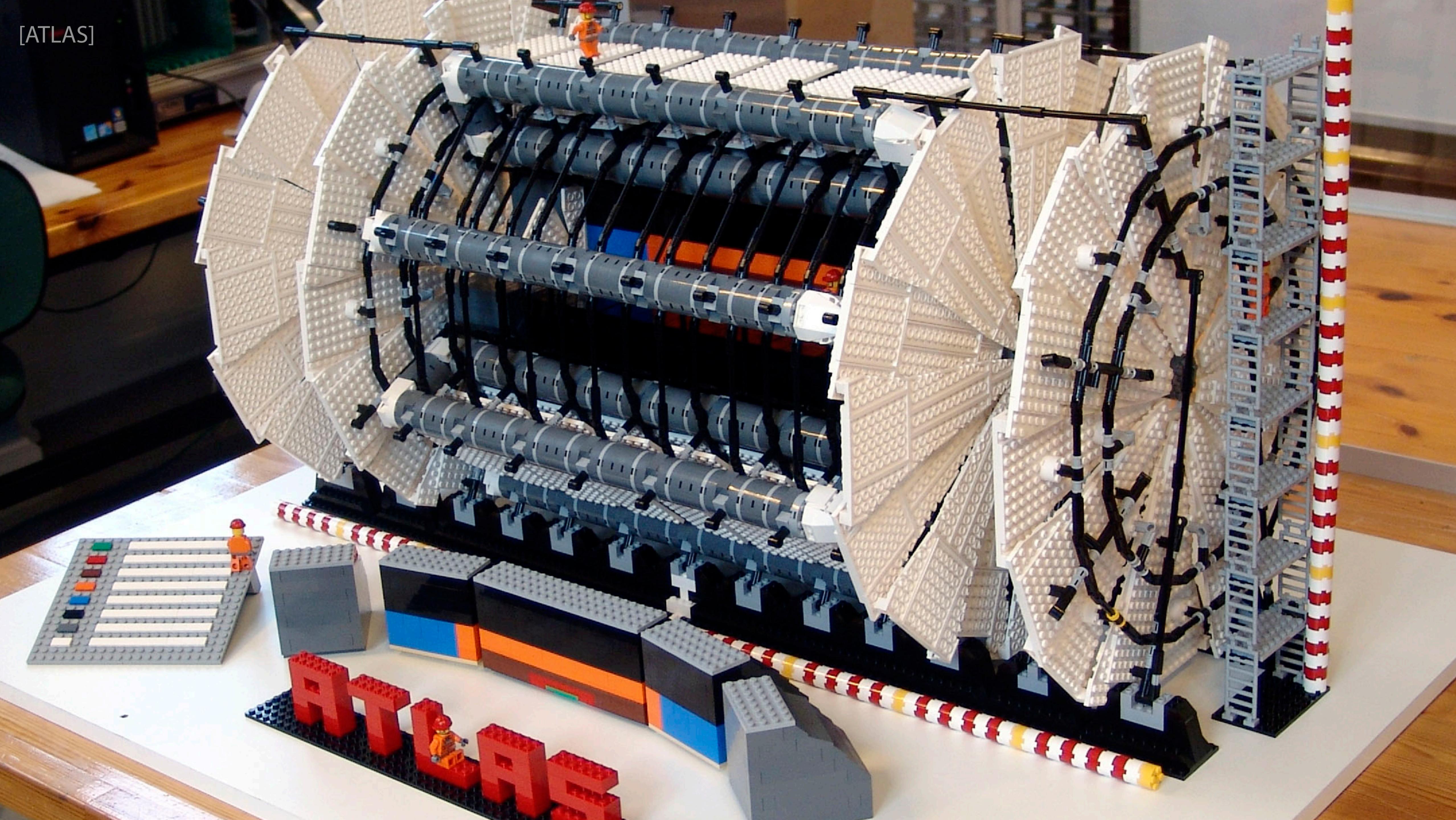


$$\begin{aligned} S = \int d^4x \left[ & \mathcal{L}_{\text{SM}} + \frac{f_{\phi,2}}{\Lambda^2} \frac{1}{2} \partial_\mu (\phi^\dagger \phi) \partial^\mu (\phi^\dagger \phi) + \frac{f_{\phi,3}}{\Lambda^2} \frac{1}{3} (\phi^\dagger \phi)^3 \right. \\ & + \frac{f_{GG}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a G^{\mu\nu a} - \frac{f_{BB}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} B^{\mu\nu} - \frac{f_{WW}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a W^{\mu\nu a} \\ & + \frac{f_B}{\Lambda^2} \frac{ig'}{2} (D^\mu \phi)^\dagger D^\nu \phi B_{\mu\nu} + \frac{f_W}{\Lambda^2} \frac{ig}{2} (D^\mu \phi)^\dagger \sigma^a D^\nu \phi W_{\mu\nu}^a \\ & + \frac{f_\ell}{\Lambda^2} (\phi^\dagger \phi) \bar{L}_L \phi \ell_R + \frac{f_u}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \tilde{\phi} u_R + \frac{f_d}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \phi d_R \\ & \left. + \frac{f_{G\widetilde{G}}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a \widetilde{G}^{\mu\nu a} - \frac{f_{B\widetilde{B}}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} \widetilde{B}^{\mu\nu} - \frac{f_{W\widetilde{W}}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a \widetilde{W}^{\mu\nu a} \right] \end{aligned}$$

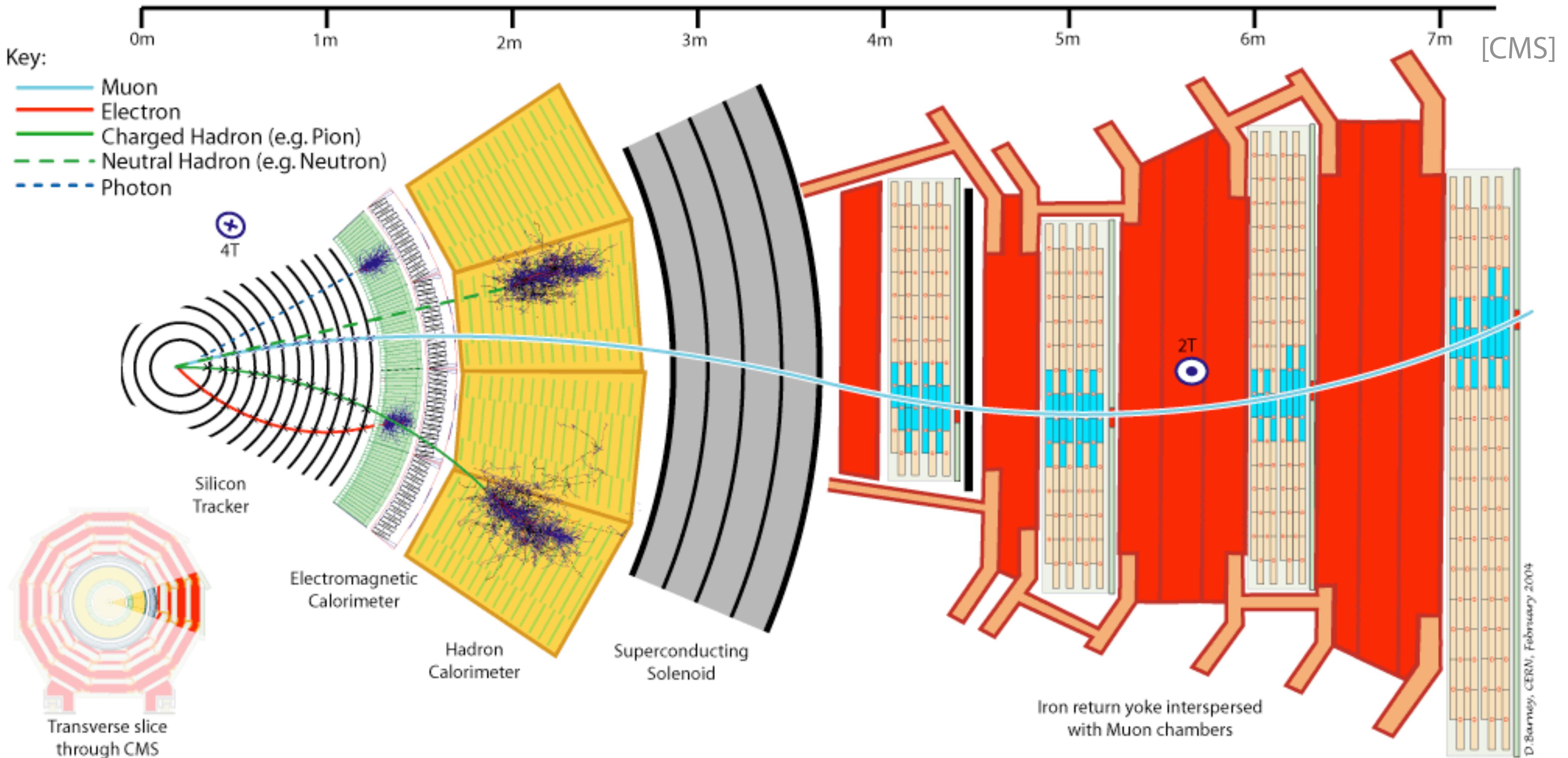
# When protons meet

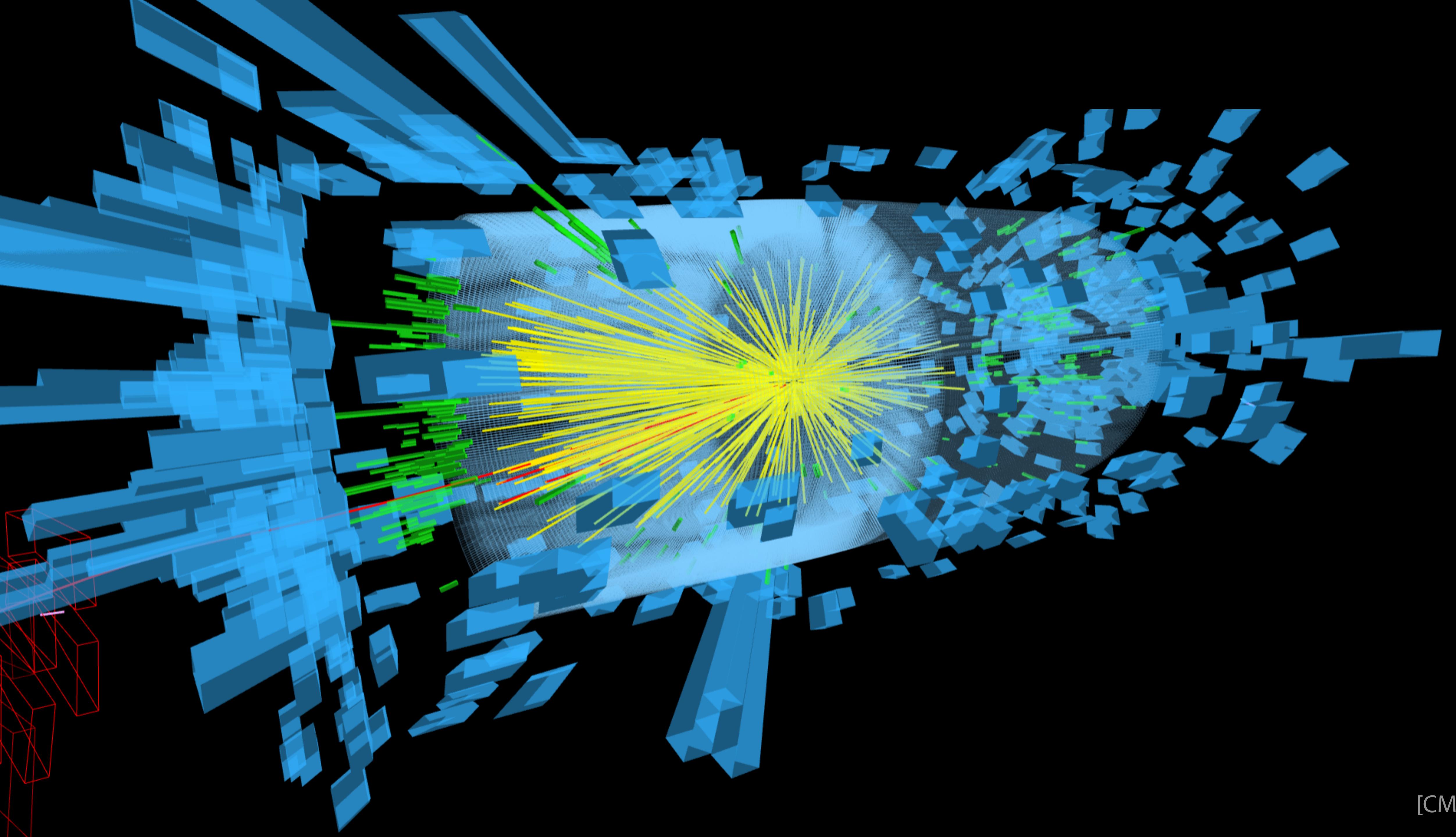


[ATLAS]



# Particle, meet detector





[CMS]

# Our model of particle physics is “likelihood-free”

$\mathcal{O}(10)$  parameters  $\theta$

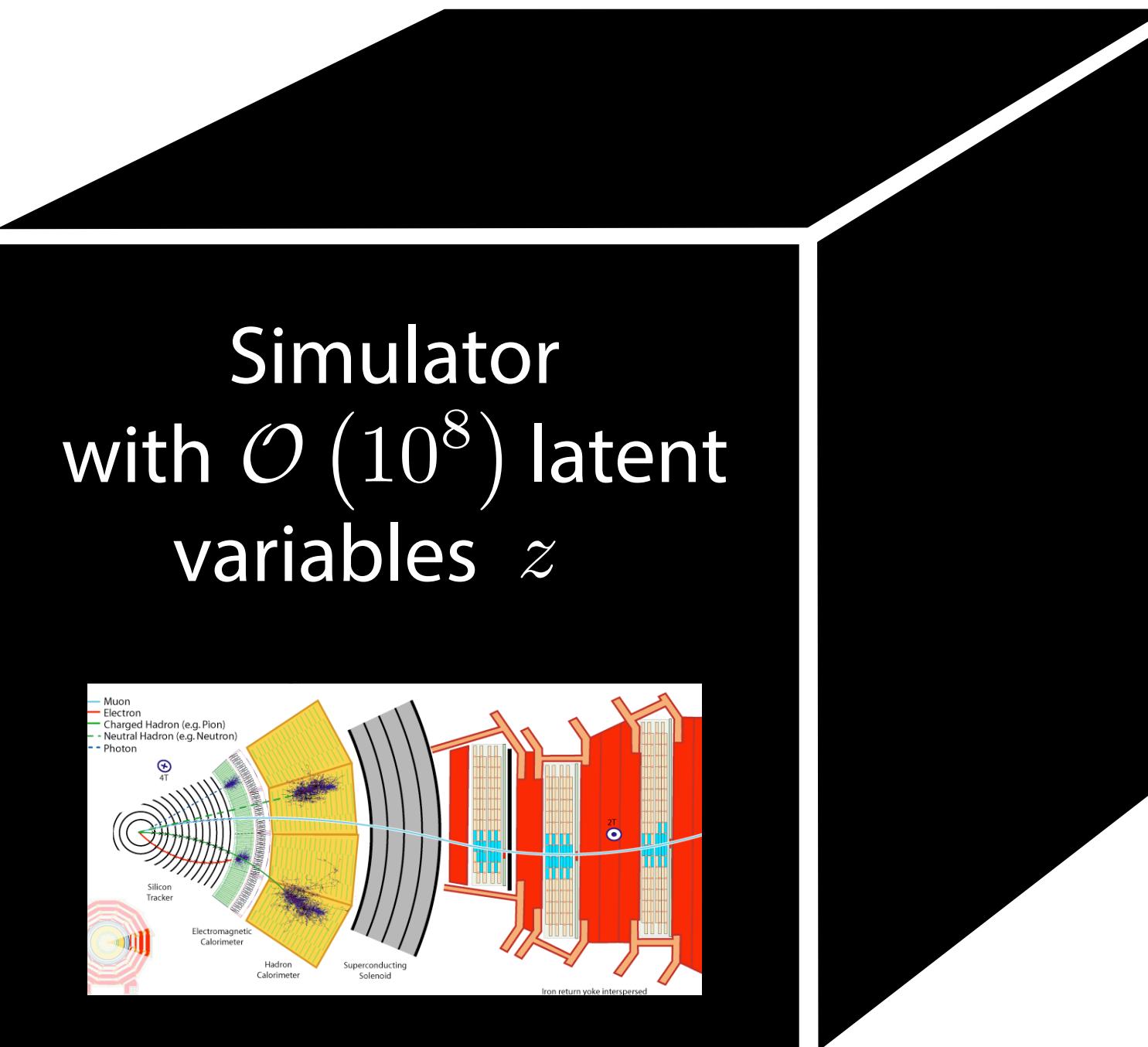
$$\begin{aligned} S = \int d^4x \left[ & \mathcal{L}_{\text{SM}} + \frac{f_{\phi,2}}{\Lambda^2} \frac{1}{2} \partial_\mu (\phi^\dagger \phi) \partial^\mu (\phi^\dagger \phi) + \frac{f_{\phi,3}}{\Lambda^2} \frac{1}{3} (\phi^\dagger \phi)^3 \right. \\ & + \frac{f_{GG}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a G^{\mu\nu a} - \frac{f_{BB}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} B^{\mu\nu} - \frac{f_{WW}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a W^{\mu\nu a} \\ & + \frac{f_B}{\Lambda^2} \frac{ig'}{2} (D^\mu \phi)^\dagger D^\nu \phi B_{\mu\nu} + \frac{f_W}{\Lambda^2} \frac{ig}{2} (D^\mu \phi)^\dagger \sigma^a D^\nu \phi W_{\mu\nu}^a \\ & + \frac{f_e}{\Lambda^2} (\phi^\dagger \phi) \bar{L}_L \phi \ell_R + \frac{f_u}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \tilde{\phi} u_R + \frac{f_d}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \phi d_R \\ & \left. + \frac{f_{G\widetilde{G}}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a \widetilde{G}^{\mu\nu a} - \frac{f_{B\widetilde{B}}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} \widetilde{B}^{\mu\nu} - \frac{f_{W\widetilde{W}}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a \widetilde{W}^{\mu\nu a} \right] \end{aligned}$$

# Our model of particle physics is “likelihood-free”

$\mathcal{O}(10)$  parameters  $\theta$



$$\begin{aligned} S = \int d^4x \left[ \mathcal{L}_{\text{SM}} + \frac{f_{\phi,2}}{\Lambda^2} \frac{1}{2} \partial_\mu (\phi^\dagger \phi) \partial^\mu (\phi^\dagger \phi) + \frac{f_{\phi,3}}{\Lambda^2} \frac{1}{3} (\phi^\dagger \phi)^3 \right. \\ + \frac{f_{GG}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a G^{\mu\nu a} - \frac{f_{BB}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} B^{\mu\nu} - \frac{f_{WW}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a W^{\mu\nu a} \\ + \frac{f_B}{\Lambda^2} \frac{ig'}{2} (D^\mu \phi)^\dagger D^\nu \phi B_{\mu\nu} + \frac{f_W}{\Lambda^2} \frac{ig}{2} (D^\mu \phi)^\dagger \sigma^a D^\nu \phi W_{\mu\nu}^a \\ + \frac{f_e}{\Lambda^2} (\phi^\dagger \phi) \bar{L}_L \phi \ell_R + \frac{f_u}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \tilde{\phi} u_R + \frac{f_d}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \phi d_R \\ \left. + \frac{f_{G\widetilde{G}}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a \widetilde{G}^{\mu\nu a} - \frac{f_{B\widetilde{B}}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} \widetilde{B}^{\mu\nu} - \frac{f_{W\widetilde{W}}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a \widetilde{W}^{\mu\nu a} \right] \end{aligned}$$



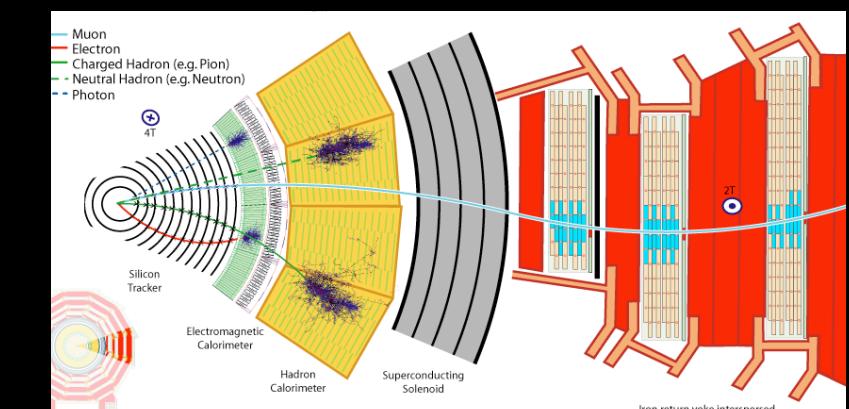
# Our model of particle physics is “likelihood-free”

$\mathcal{O}(10)$  parameters  $\theta$

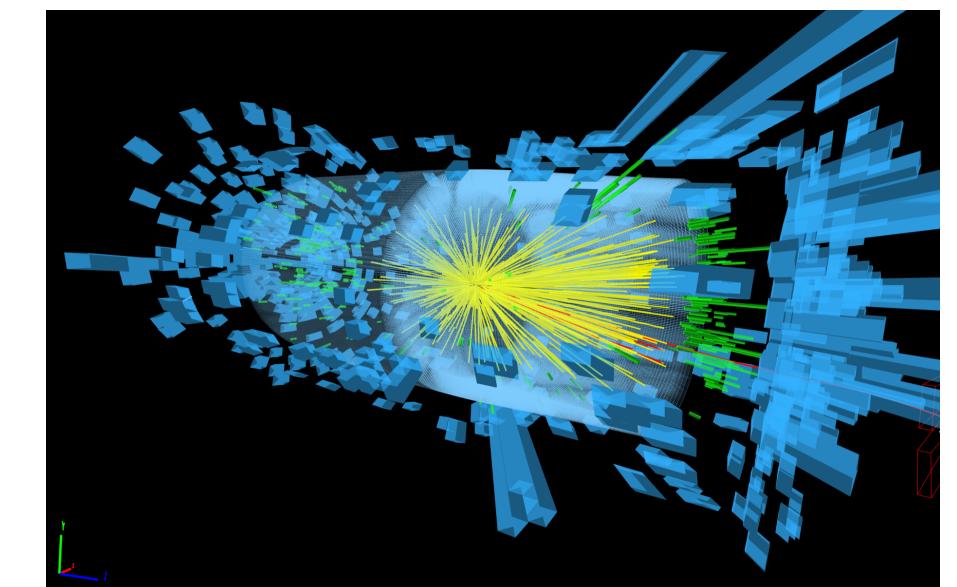


$$S = \int d^4x \left[ \mathcal{L}_{\text{SM}} + \frac{f_{\phi,2}}{\Lambda^2} \frac{1}{2} \partial_\mu (\phi^\dagger \phi) \partial^\mu (\phi^\dagger \phi) + \frac{f_{\phi,3}}{\Lambda^2} \frac{1}{3} (\phi^\dagger \phi)^3 \right. \\ + \frac{f_{GG}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a G^{\mu\nu a} - \frac{f_{BB}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} B^{\mu\nu} - \frac{f_{WW}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a W^{\mu\nu a} \\ + \frac{f_B}{\Lambda^2} \frac{ig'}{2} (D^\mu \phi)^\dagger D^\nu \phi B_{\mu\nu} + \frac{f_W}{\Lambda^2} \frac{ig}{2} (D^\mu \phi)^\dagger \sigma^a D^\nu \phi W_{\mu\nu}^a \\ + \frac{f_e}{\Lambda^2} (\phi^\dagger \phi) \bar{L}_L \phi \ell_R + \frac{f_u}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \tilde{\phi} u_R + \frac{f_d}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \phi d_R \\ \left. + \frac{f_{G\widetilde{G}}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a \widetilde{G}^{\mu\nu a} - \frac{f_{B\widetilde{B}}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} \widetilde{B}^{\mu\nu} - \frac{f_{W\widetilde{W}}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a \widetilde{W}^{\mu\nu a} \right]$$

Simulator  
with  $\mathcal{O}(10^8)$  latent  
variables  $z$



$\mathcal{O}(10 \dots 1000)$   
observables  $x$



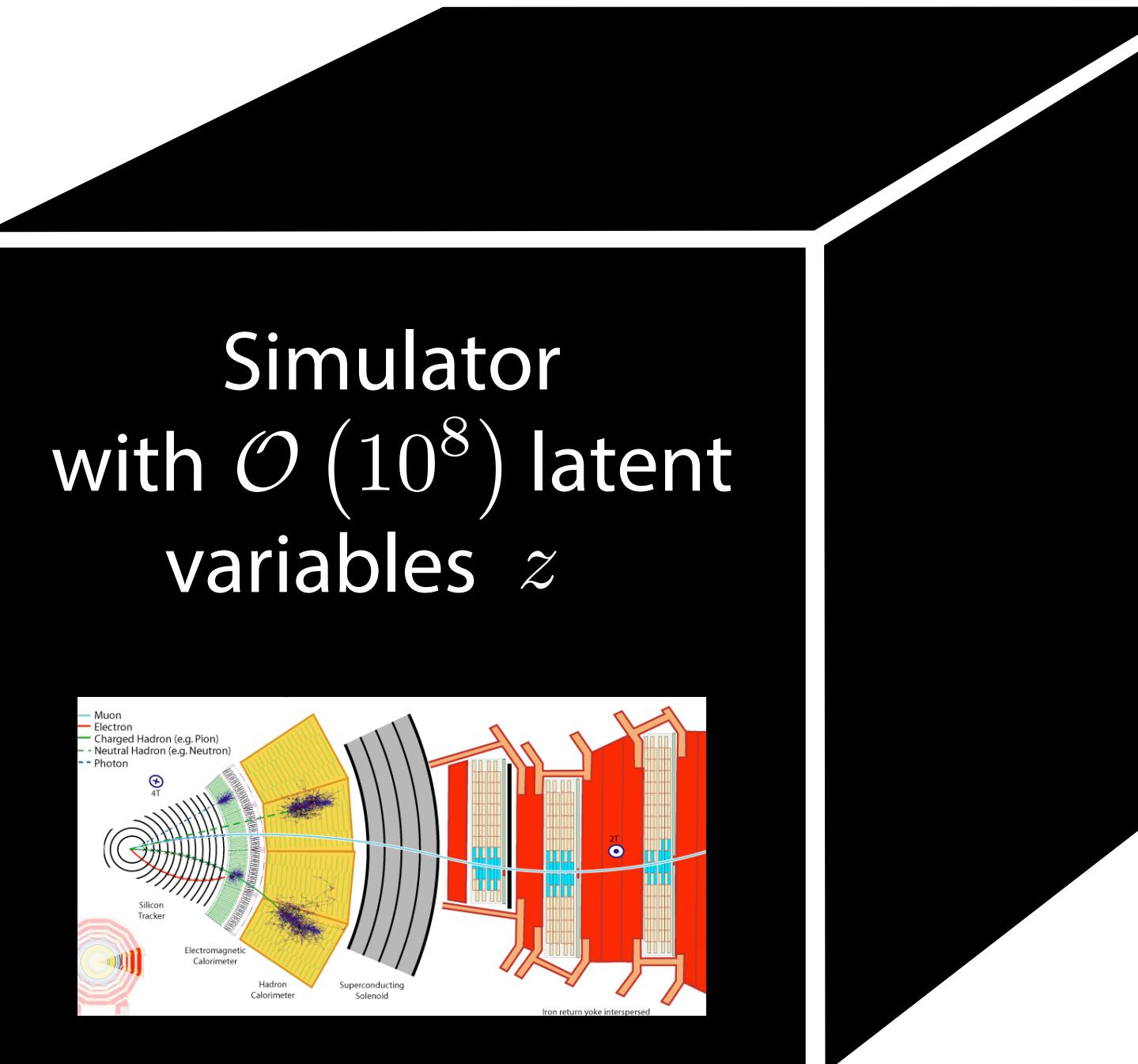
# Our model of particle physics is “likelihood-free”

$\mathcal{O}(10)$  parameters  $\theta$

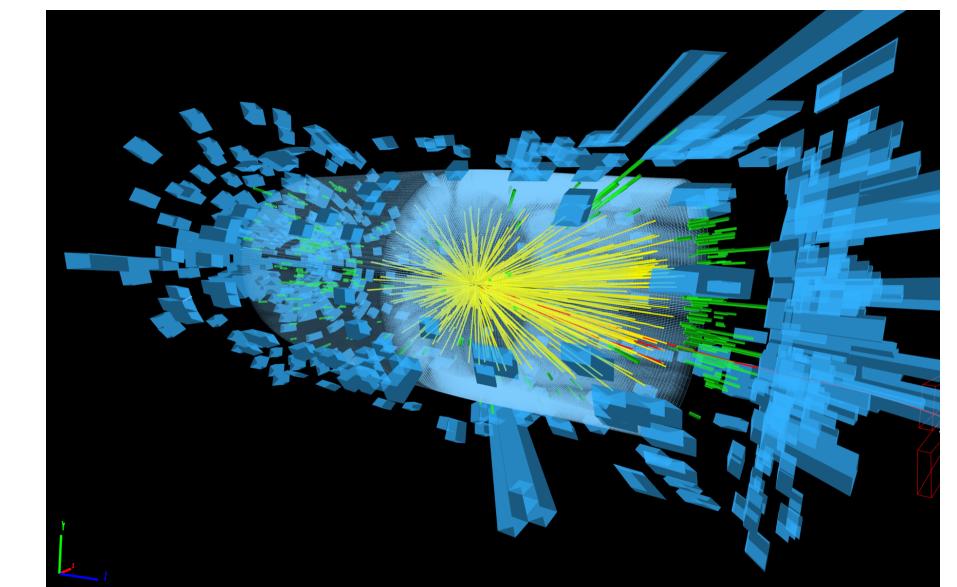
$$S = \int d^4x \left[ \mathcal{L}_{\text{SM}} + \frac{f_{\phi,2}}{\Lambda^2} \frac{1}{2} \partial_\mu (\phi^\dagger \phi) \partial^\mu (\phi^\dagger \phi) + \frac{f_{\phi,3}}{\Lambda^2} \frac{1}{3} (\phi^\dagger \phi)^3 \right. \\ + \frac{f_{GG}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a G^{\mu\nu a} - \frac{f_{BB}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} B^{\mu\nu} - \frac{f_{WW}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a W^{\mu\nu a} \\ + \frac{f_B}{\Lambda^2} \frac{ig'}{2} (D^\mu \phi)^\dagger D^\nu \phi B_{\mu\nu} + \frac{f_W}{\Lambda^2} \frac{ig}{2} (D^\mu \phi)^\dagger \sigma^a D^\nu \phi W_{\mu\nu}^a \\ + \frac{f_e}{\Lambda^2} (\phi^\dagger \phi) \bar{L}_L \phi \ell_R + \frac{f_u}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \tilde{\phi} u_R + \frac{f_d}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \phi d_R \\ \left. + \frac{f_{G\widetilde{G}}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a \widetilde{G}^{\mu\nu a} - \frac{f_{B\widetilde{B}}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} \widetilde{B}^{\mu\nu} - \frac{f_{W\widetilde{W}}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a \widetilde{W}^{\mu\nu a} \right]$$



Simulator  
with  $\mathcal{O}(10^8)$  latent  
variables  $z$



$\mathcal{O}(10 \dots 1000)$   
observables  $x$



Prediction: Simulator can sample  $x \sim p(x|\theta)$

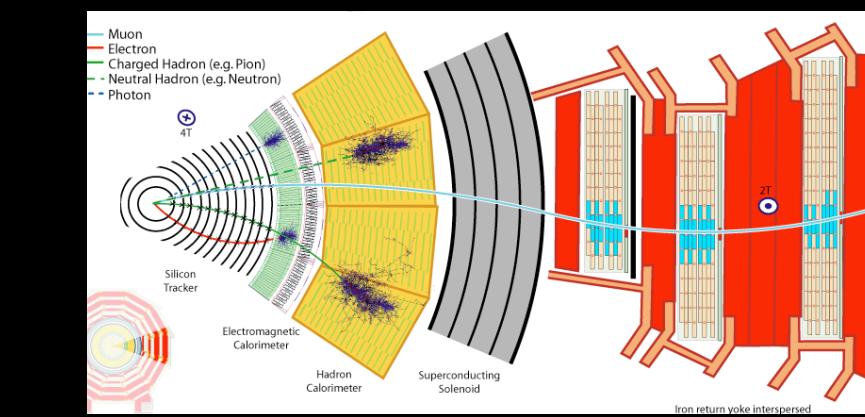
# Our model of particle physics is “likelihood-free”

$\mathcal{O}(10)$  parameters  $\theta$

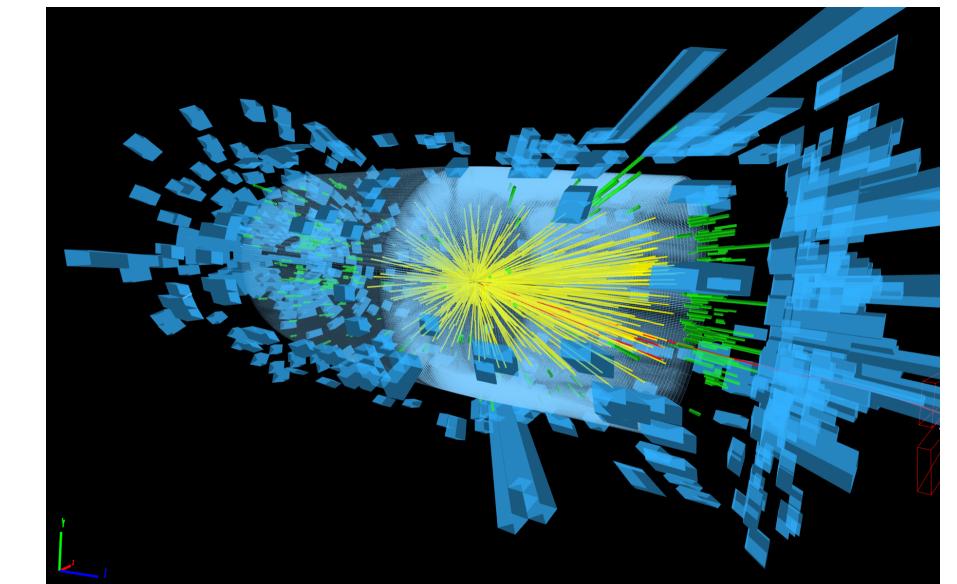
$$S = \int d^4x \left[ \mathcal{L}_{\text{SM}} + \frac{f_{\phi,2}}{\Lambda^2} \frac{1}{2} \partial_\mu (\phi^\dagger \phi) \partial^\mu (\phi^\dagger \phi) + \frac{f_{\phi,3}}{\Lambda^2} \frac{1}{3} (\phi^\dagger \phi)^3 \right. \\ + \frac{f_{GG}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a G^{\mu\nu a} - \frac{f_{BB}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} B^{\mu\nu} - \frac{f_{WW}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a W^{\mu\nu a} \\ + \frac{f_B}{\Lambda^2} \frac{ig'}{2} (D^\mu \phi)^\dagger D^\nu \phi B_{\mu\nu} + \frac{f_W}{\Lambda^2} \frac{ig}{2} (D^\mu \phi)^\dagger \sigma^a D^\nu \phi W_{\mu\nu}^a \\ + \frac{f_e}{\Lambda^2} (\phi^\dagger \phi) \bar{L}_L \phi \ell_R + \frac{f_u}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \tilde{\phi} u_R + \frac{f_d}{\Lambda^2} (\phi^\dagger \phi) \bar{Q}_L \phi d_R \\ \left. + \frac{f_{G\widetilde{G}}}{\Lambda^2} (\phi^\dagger \phi) G_{\mu\nu}^a \widetilde{G}^{\mu\nu a} - \frac{f_{B\widetilde{B}}}{\Lambda^2} \frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} \widetilde{B}^{\mu\nu} - \frac{f_{W\widetilde{W}}}{\Lambda^2} \frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a \widetilde{W}^{\mu\nu a} \right]$$



Simulator  
with  $\mathcal{O}(10^8)$  latent  
variables  $z$



$\mathcal{O}(10 \dots 1000)$   
observables  $x$



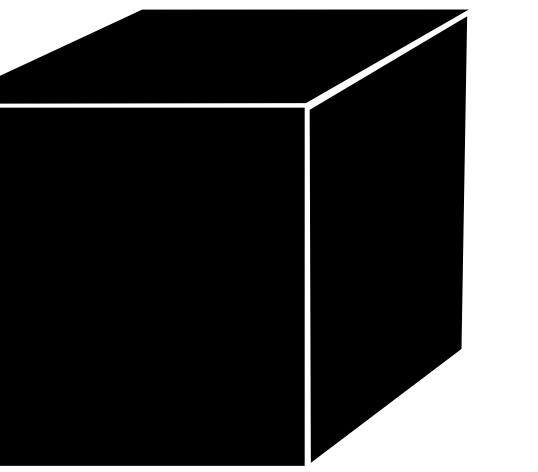
Prediction: Simulator can sample  $x \sim p(x|\theta)$

Inference: Likelihood  $p(x|\theta) = \int dz p(x, z|\theta)$  is intractable

# Problem statement(s)

You are given

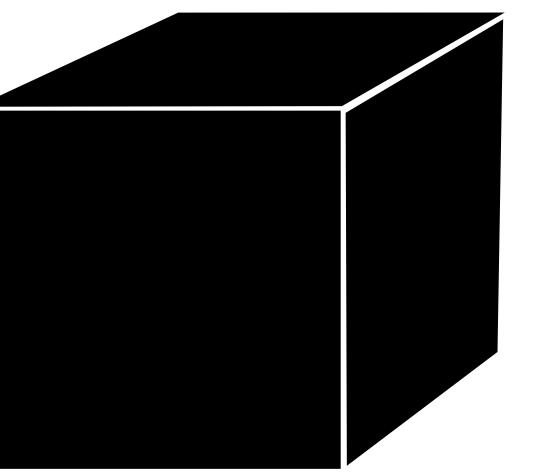
- a simulator that lets you generate  $N$  samples  $x_i \sim p(x_i|\theta_i)$  for parameter points  $\theta_i$  of your choice
- observed data  $x_{\text{obs}} \sim p(x_{\text{obs}}|\theta_{\text{true}})$
- prior belief  $p(\theta)$



# Problem statement(s)

You are given

- a simulator that lets you generate  $N$  samples  $x_i \sim p(x_i|\theta_i)$  for parameter points  $\theta_i$  of your choice
- observed data  $x_{\text{obs}} \sim p(x_{\text{obs}}|\theta_{\text{true}})$
- prior belief  $p(\theta)$



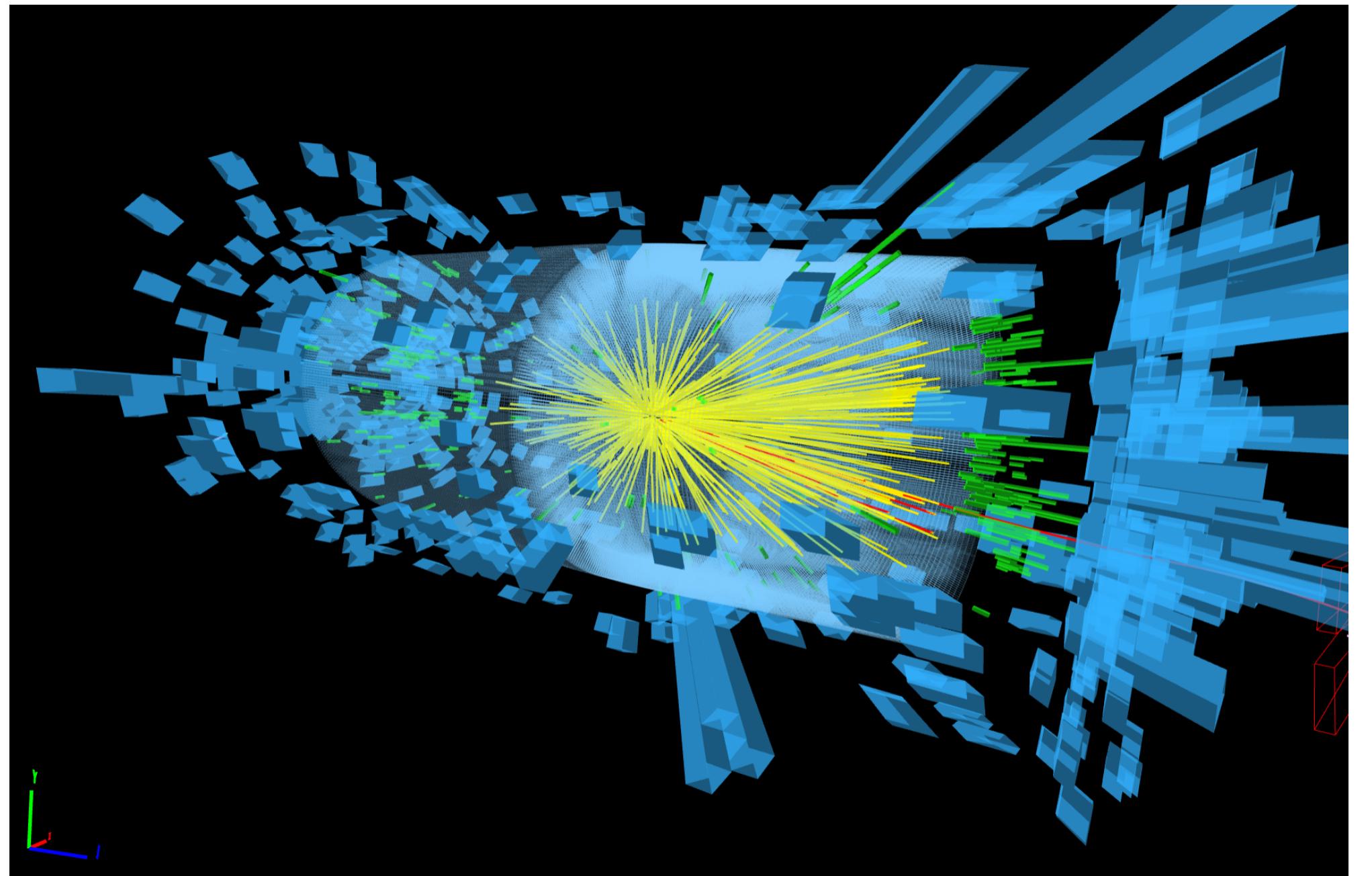
Goals: estimate either...

- true parameters  $\hat{\theta}_{\text{true}}$
- confidence sets based on likelihood  $\hat{p}(x_{\text{obs}}|\theta)$
- posterior  $\hat{p}(\theta|x_{\text{obs}}) = \frac{\hat{p}(x_{\text{obs}}|\theta) p(\theta)}{\int d\theta' \hat{p}(x_{\text{obs}}|\theta') p(\theta')}$   
or samples from posterior  $\theta \sim \hat{p}(\theta|x_{\text{obs}})$

... depending on domain conventions

We can do inference  
by compressing the data to summary statistics.

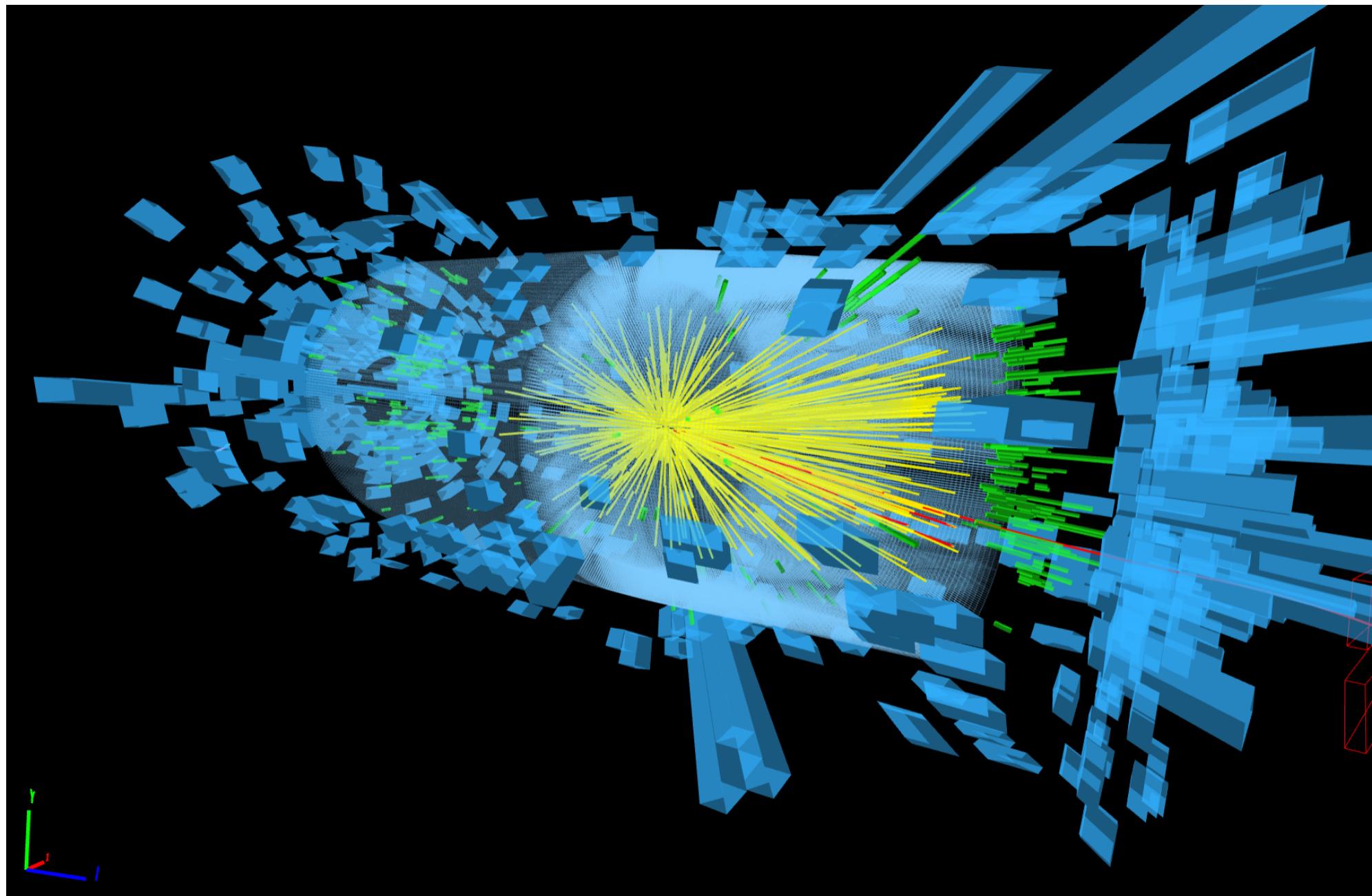
# Summary statistics



High-dimensional observable data  $x$

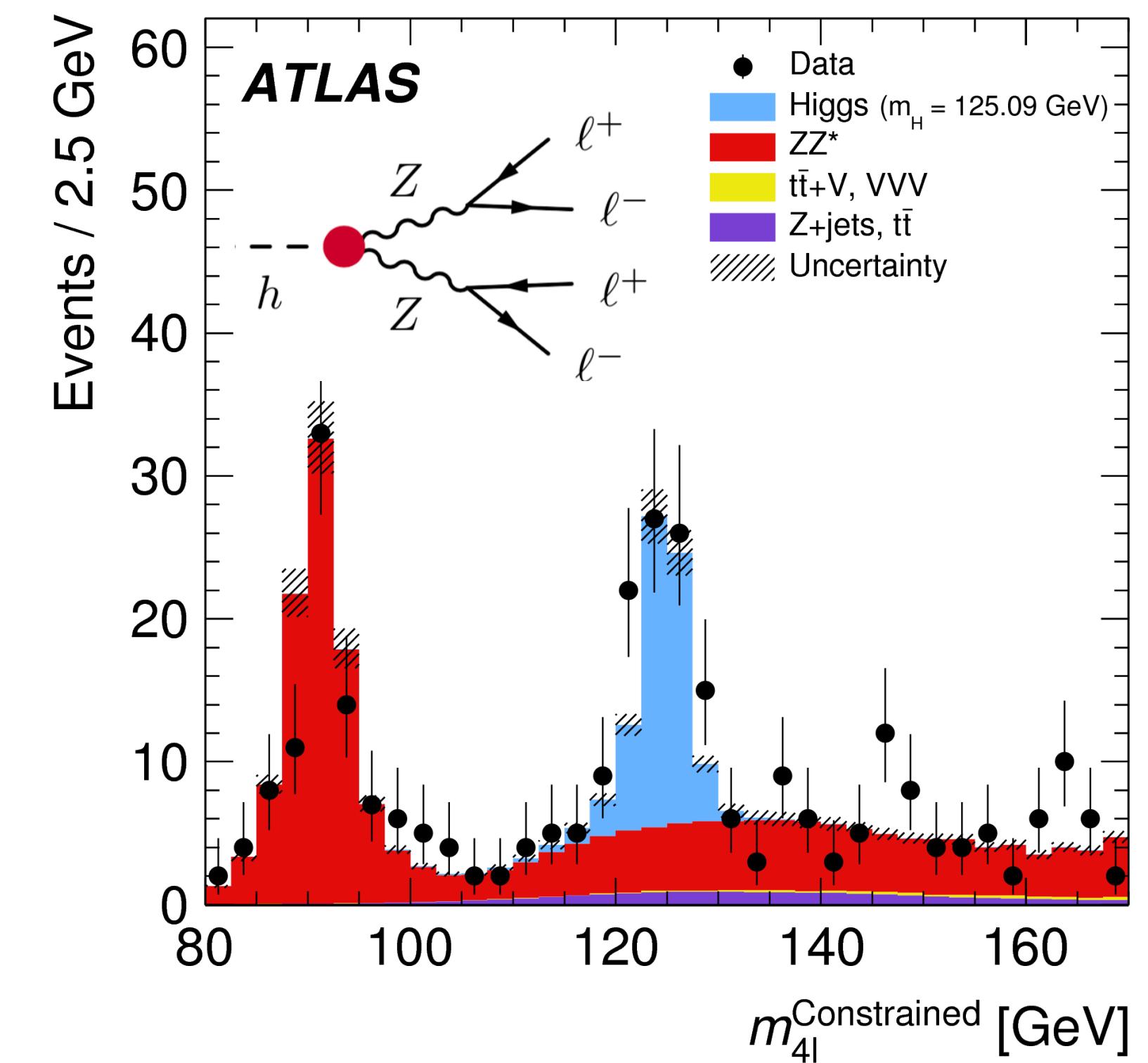
$p(x|\theta)$  cannot be calculated

# Summary statistics



High-dimensional observable data  $x$

$p(x|\theta)$  cannot be calculated



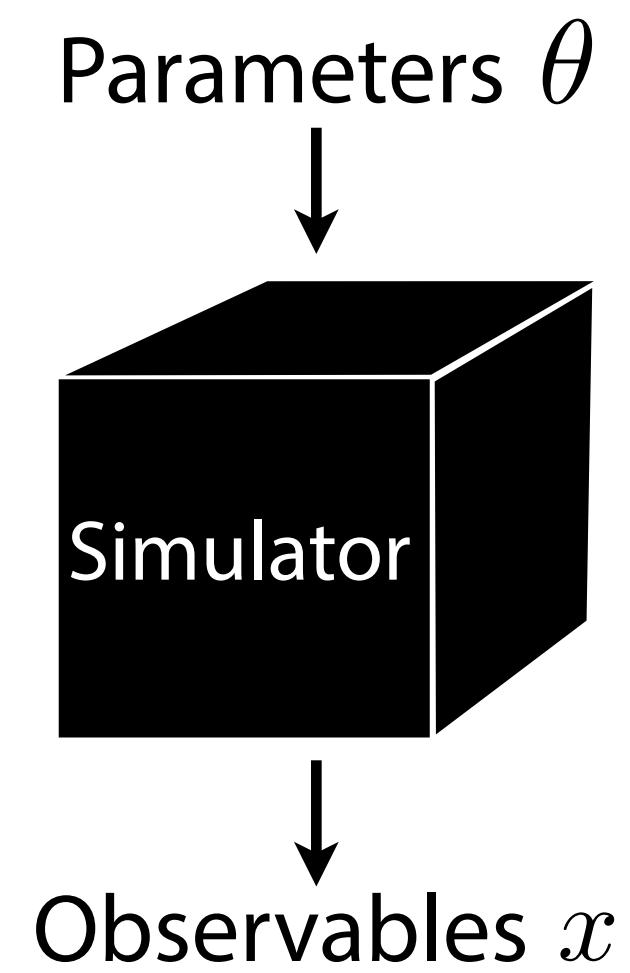
One or a few summary statistics  $x'$

$p(x'|\theta)$  can be estimated  
with histograms, KDE, ...

[ATLAS 1712.02304]

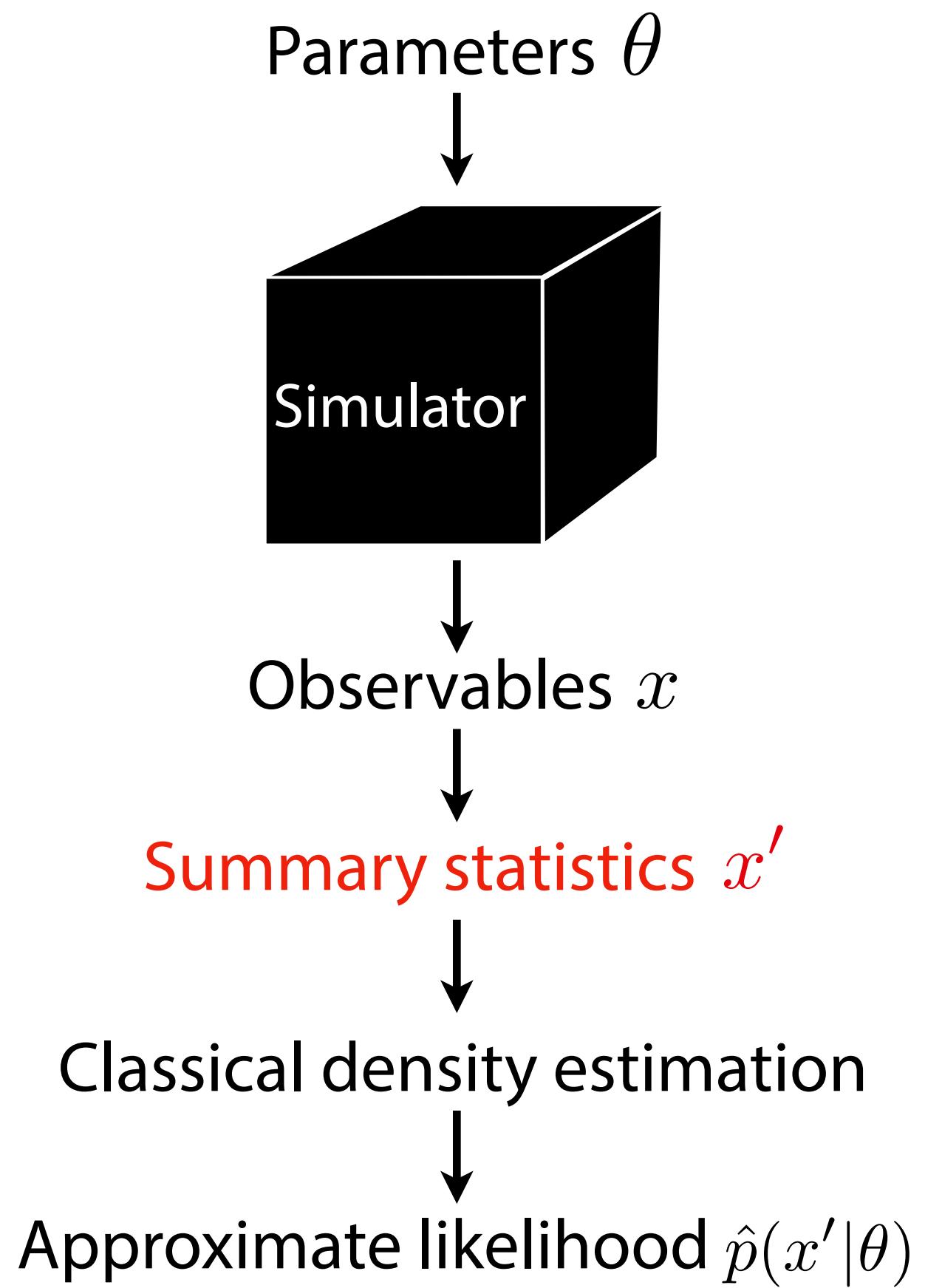
# Inference by estimating the likelihood

[e.g. P. Diggle, R. Gratton 1984]



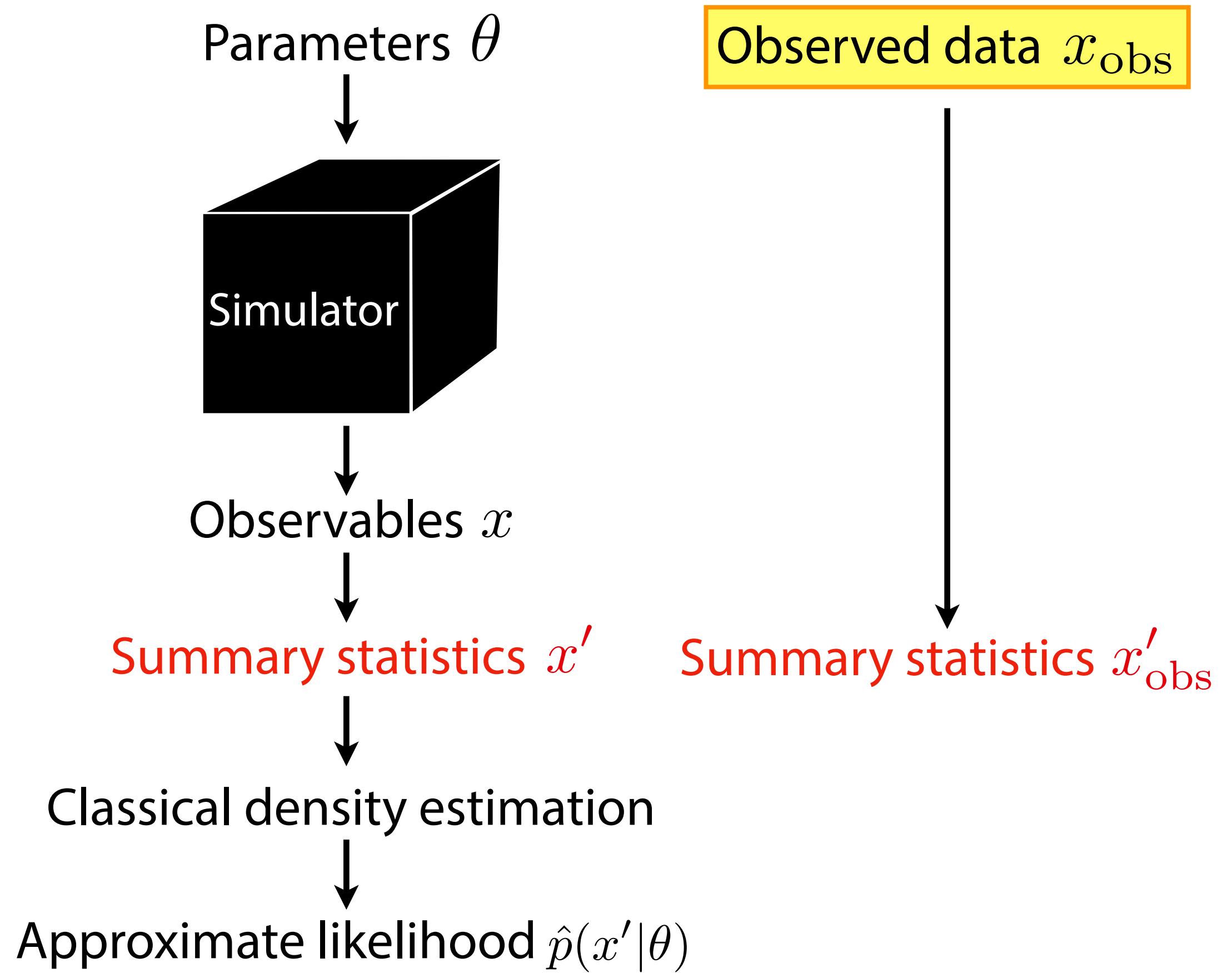
# Inference by estimating the likelihood

[e.g. P. Diggle, R. Gratton 1984]



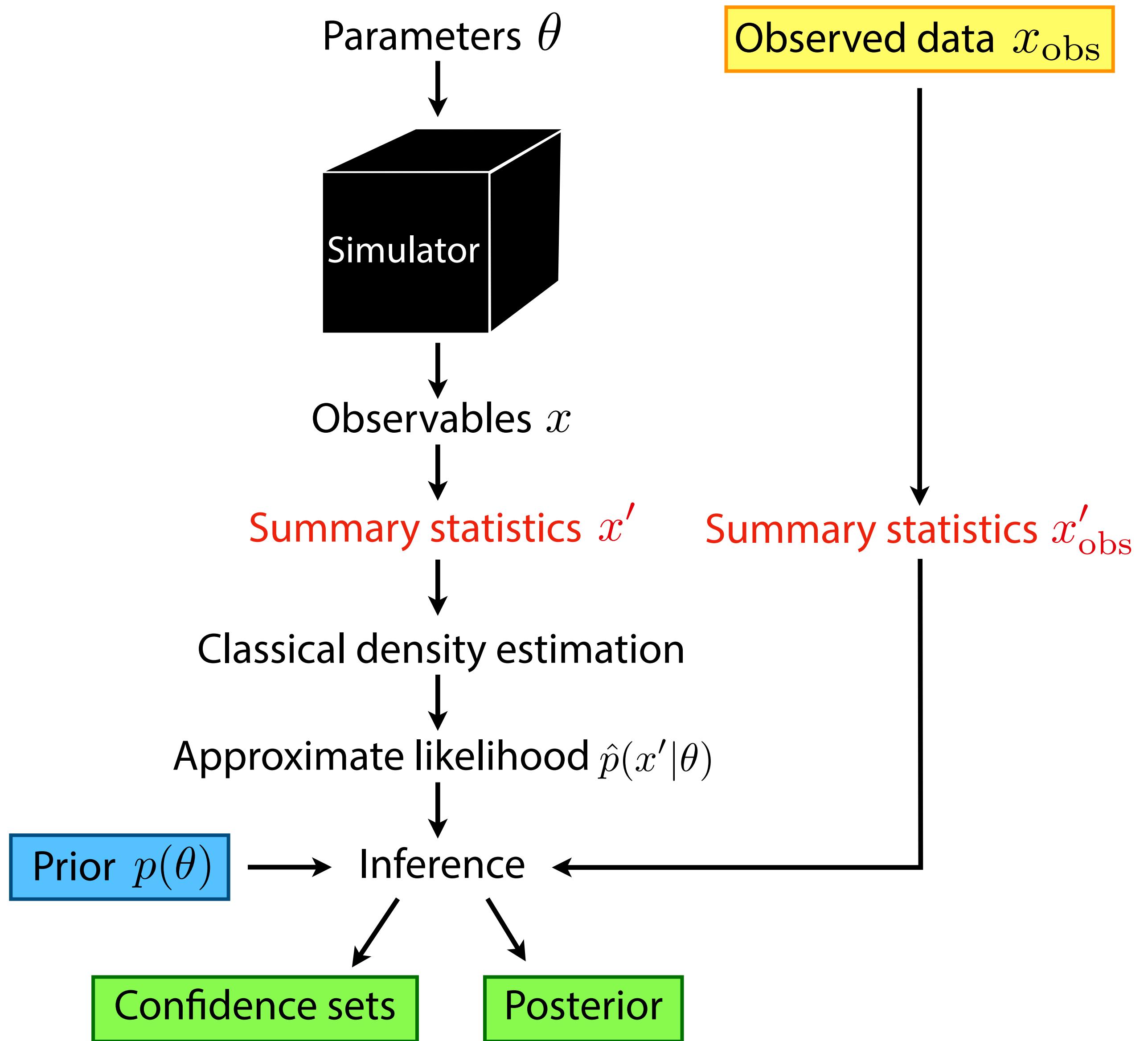
# Inference by estimating the likelihood

[e.g. P. Diggle, R. Gratton 1984]



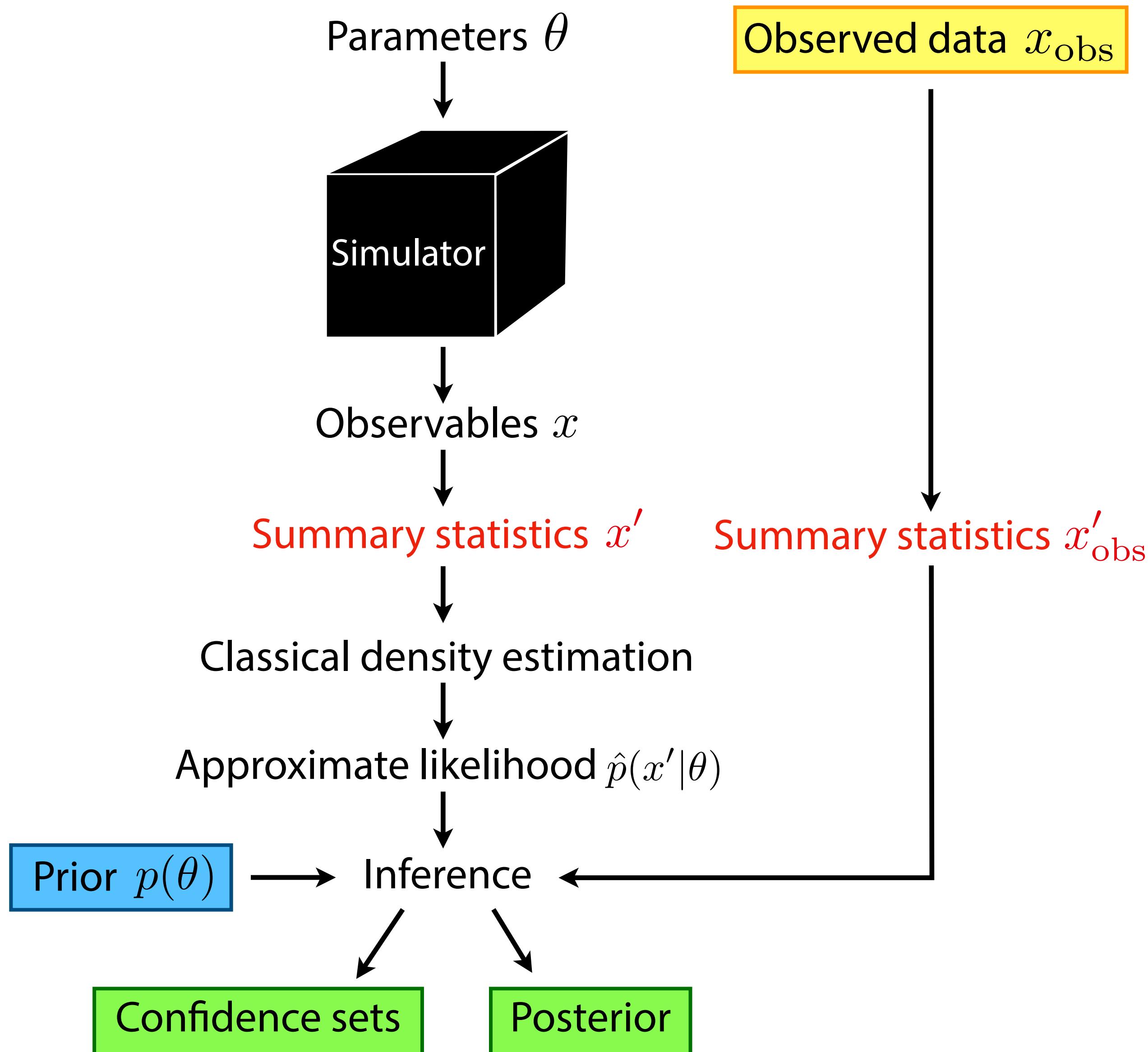
# Inference by estimating the likelihood

[e.g. P. Diggle, R. Gratton 1984]



# Inference by estimating the likelihood

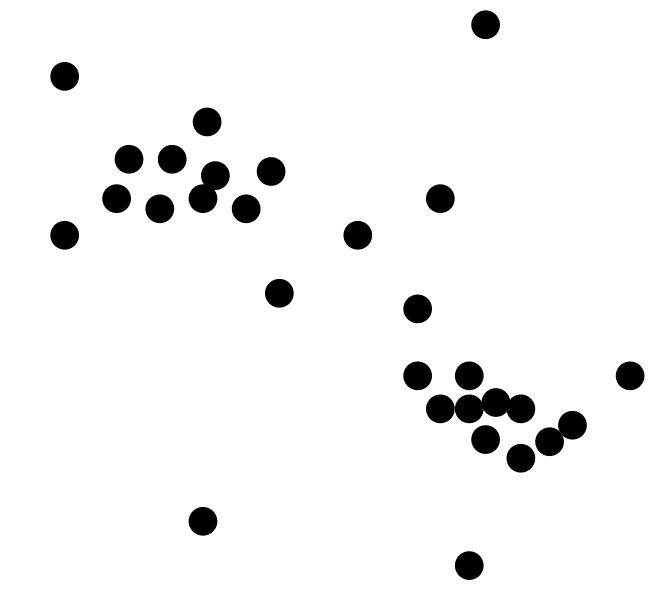
[e.g. P. Diggle, R. Gratton 1984]



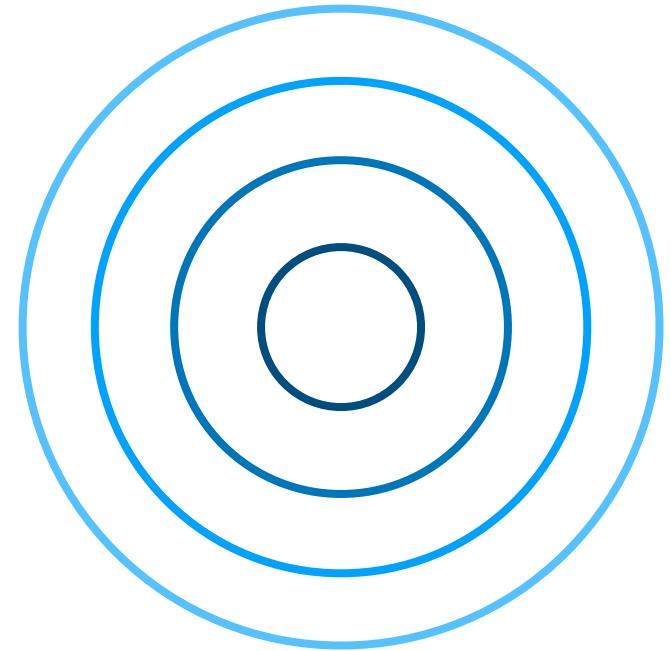
- Compression to summary statistics reduces quality of inference
- Curse of dimensionality: classical density estimation methods do not scale to more than a few summary statistics
- Popular, similar alternative: Approximate Bayesian Computation (ABC) [D. Rubin 1984]

We can do inference  
by using machine learning.

# High-dimensional density estimation with normalizing flows

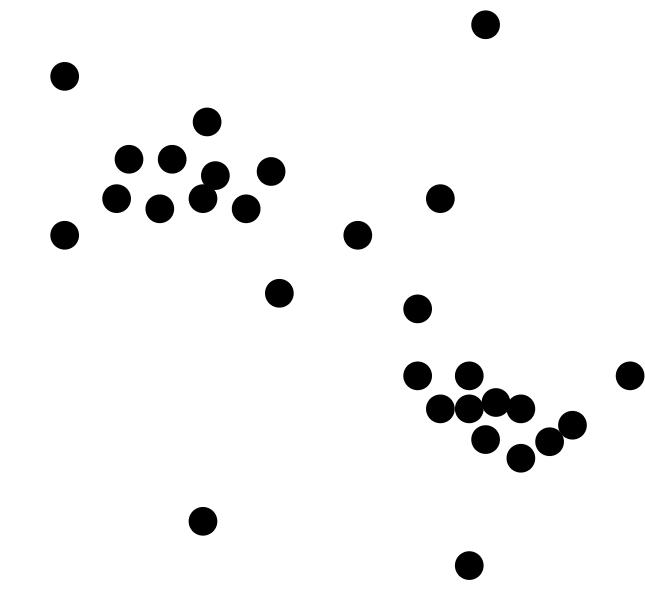


# High-dimensional density estimation with normalizing flows

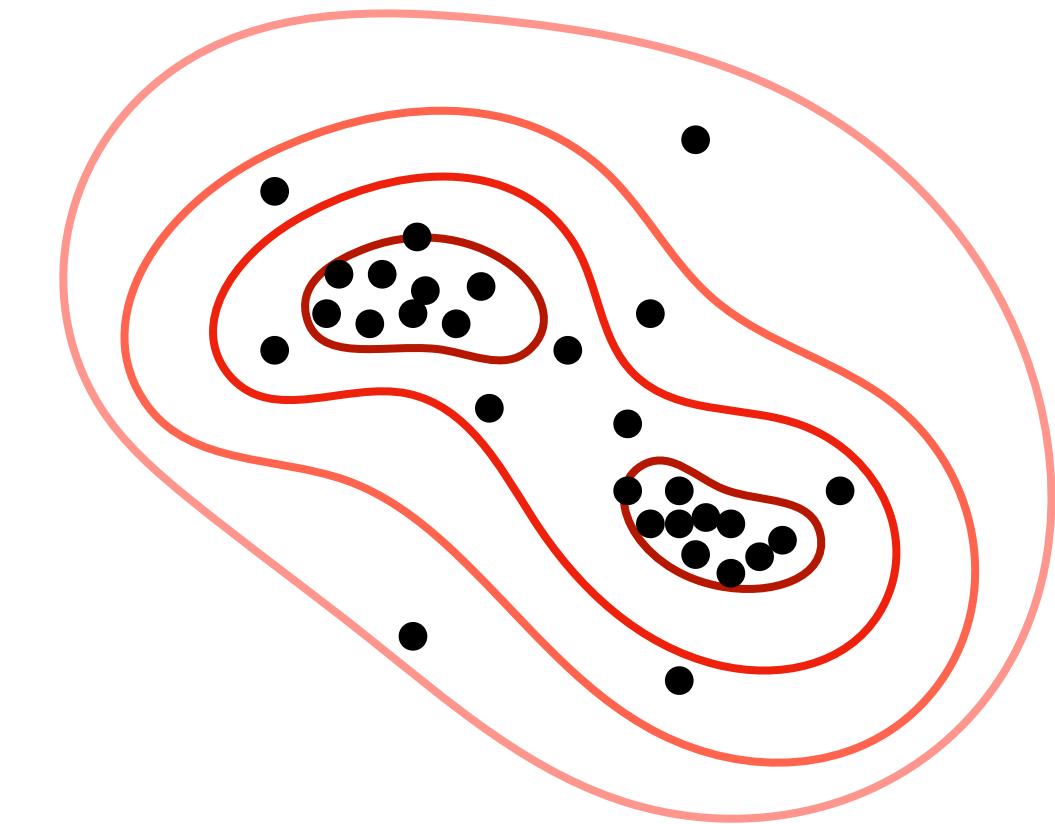
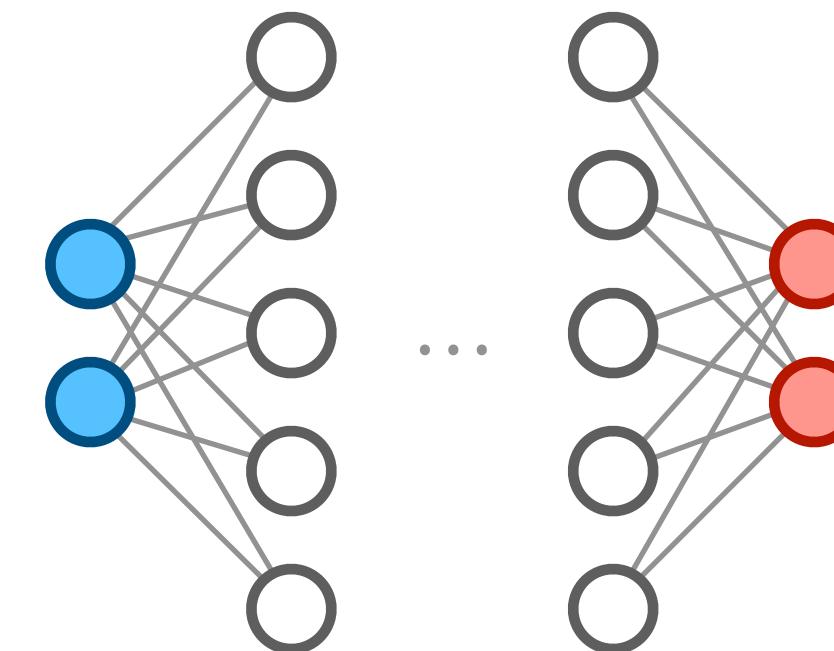
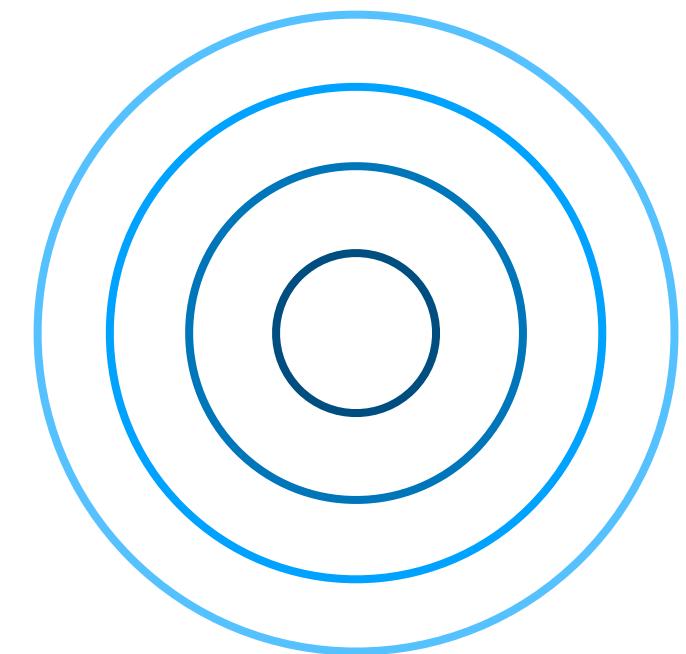


Simple base density

$$u \sim \pi(u)$$



# High-dimensional density estimation with normalizing flows



Simple base density

$$u \sim \pi(u)$$

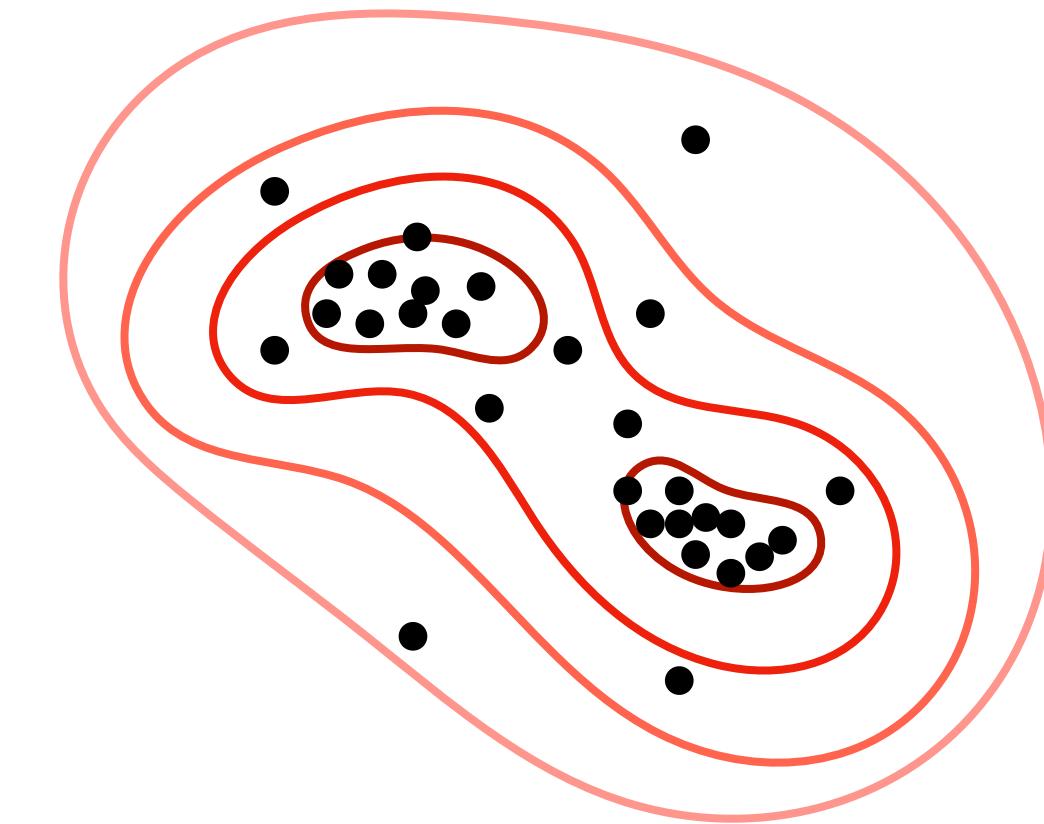
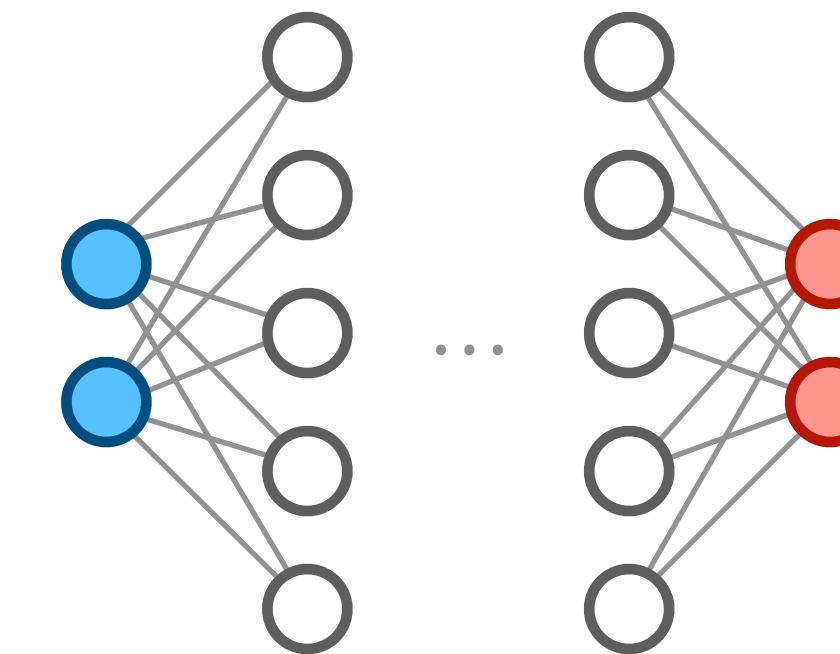
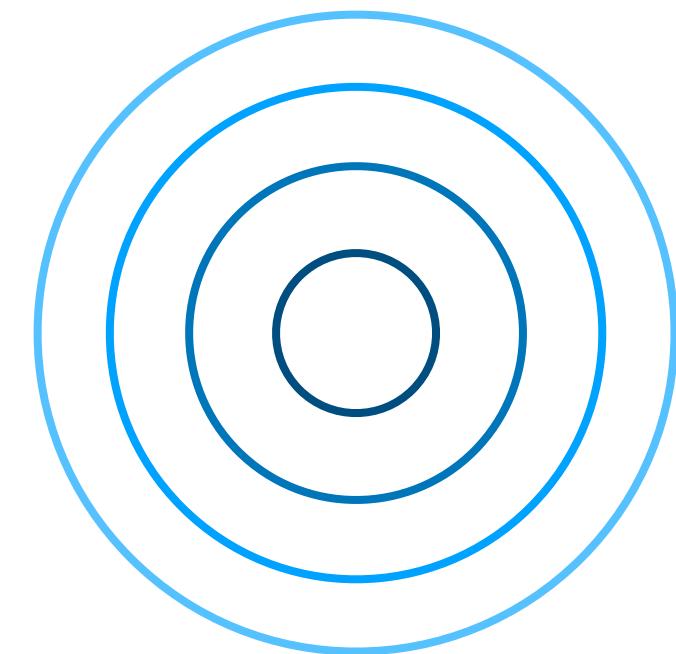
NN: transformation  $x = f(u)$

- one-to-one and invertible
- differentiable
- $f^{-1}$  and  $\det \nabla f$  are tractable

Target density is given by

$$\hat{p}(x) = \pi(f^{-1}(x)) |\det \nabla f|^{-1}$$

# High-dimensional density estimation with normalizing flows



Simple base density

$$u \sim \pi(u)$$

NN: transformation  $x = f(u)$

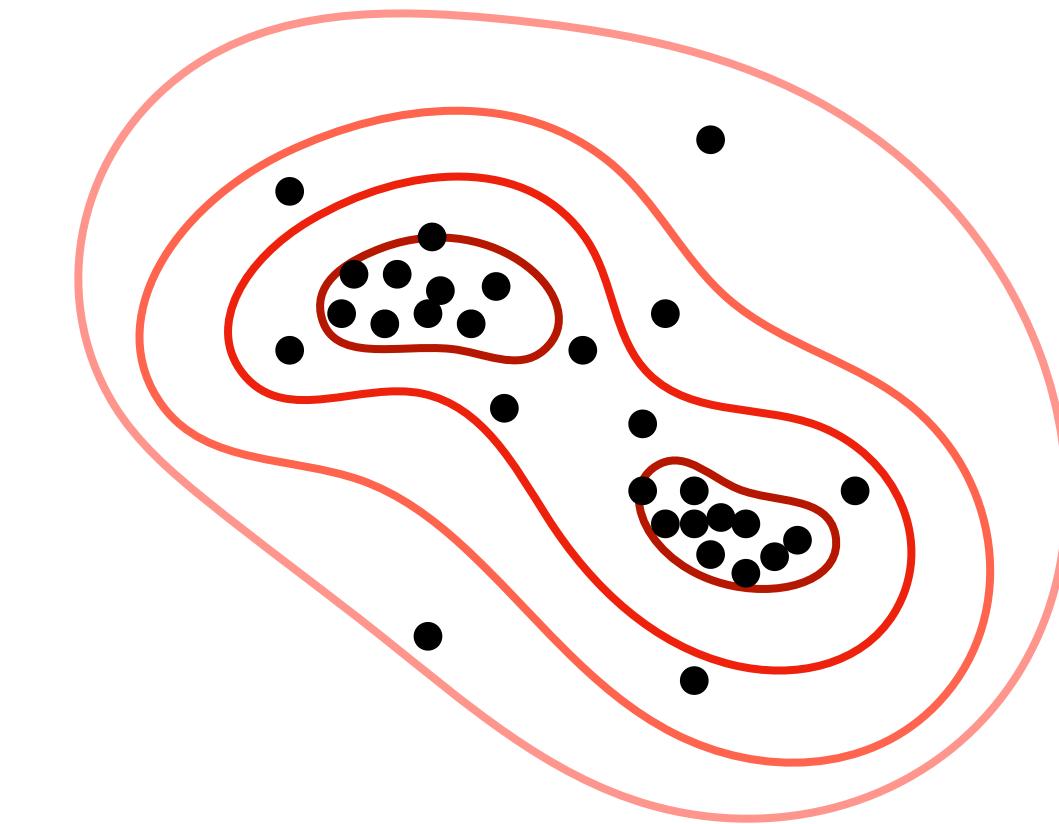
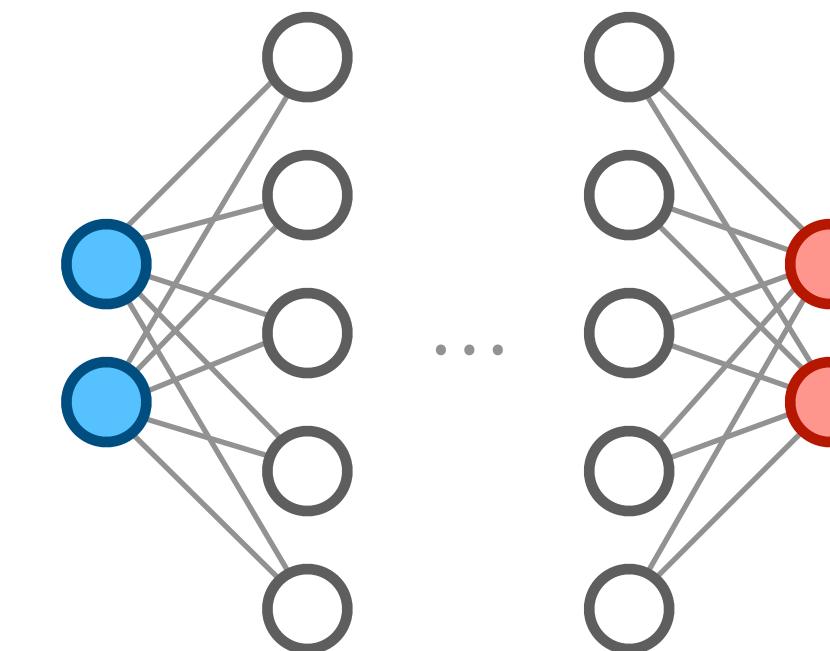
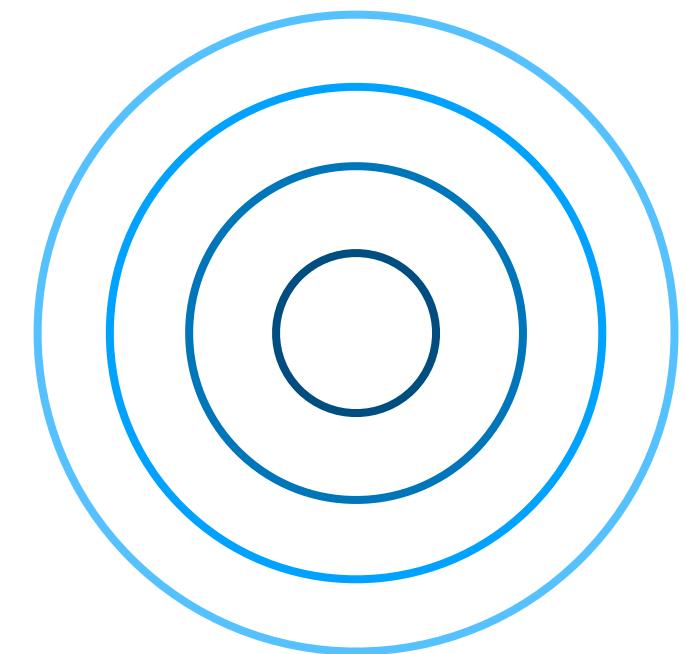
- one-to-one and invertible
- differentiable
- $f^{-1}$  and  $\det \nabla f$  are tractable

Target density is given by

$$\hat{p}(x) = \pi(f^{-1}(x)) |\det \nabla f|^{-1}$$

Train transformation by  
maximizing  $\log \hat{p}(x)$

# High-dimensional density estimation with normalizing flows



Simple base density

$$u \sim \pi(u)$$

NN: transformation  $x = f(u)$

- one-to-one and invertible
- differentiable
- $f^{-1}$  and  $\det \nabla f$  are tractable

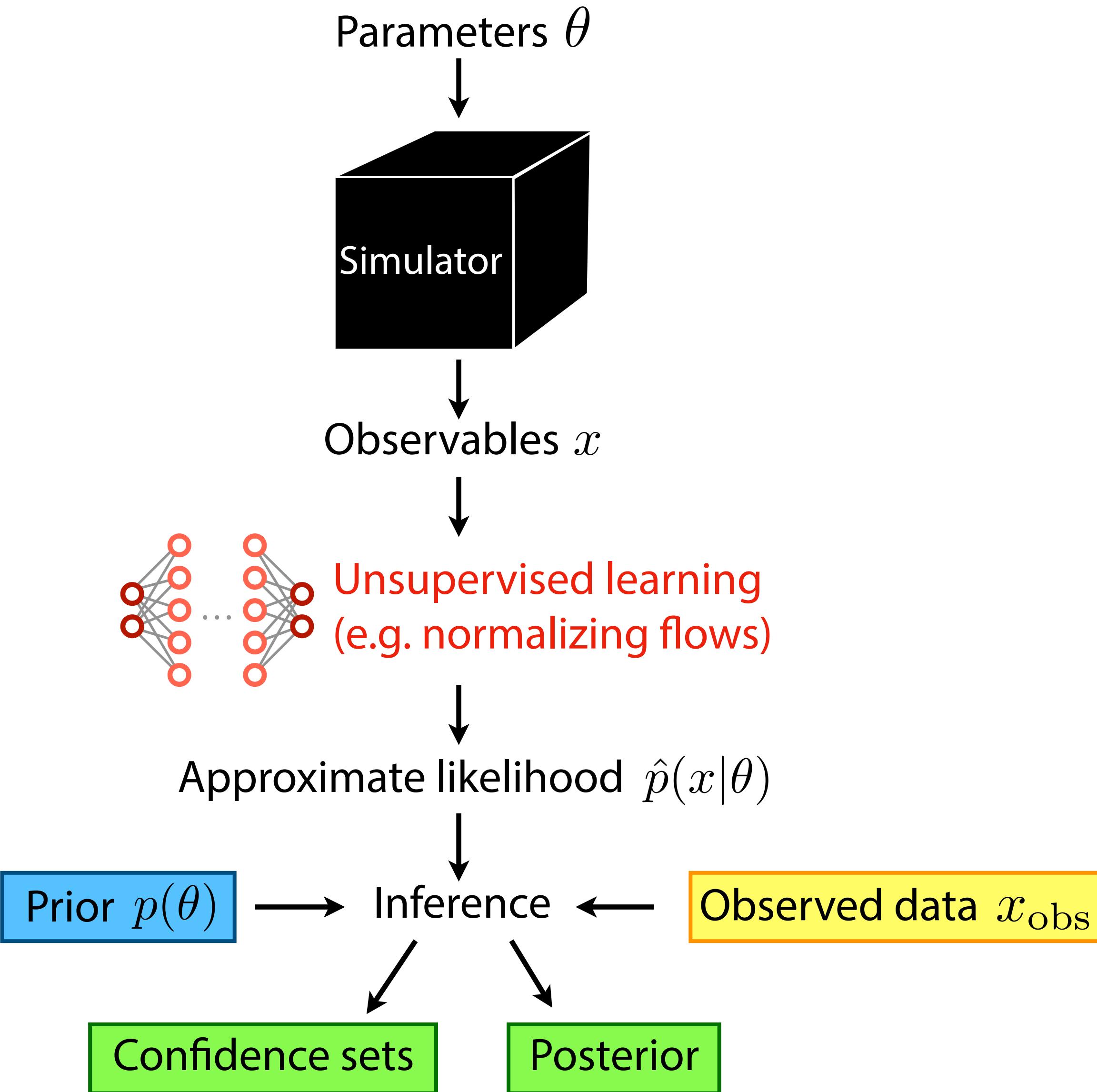
Target density is given by

$$\hat{p}(x) = \pi(f^{-1}(x)) |\det \nabla f|^{-1}$$

Train transformation by  
maximizing  $\log \hat{p}(x)$

Transformation can depend on  $\theta$   
to model conditional density  $\log \hat{p}(x|\theta)$

# Inference with neural likelihood estimation



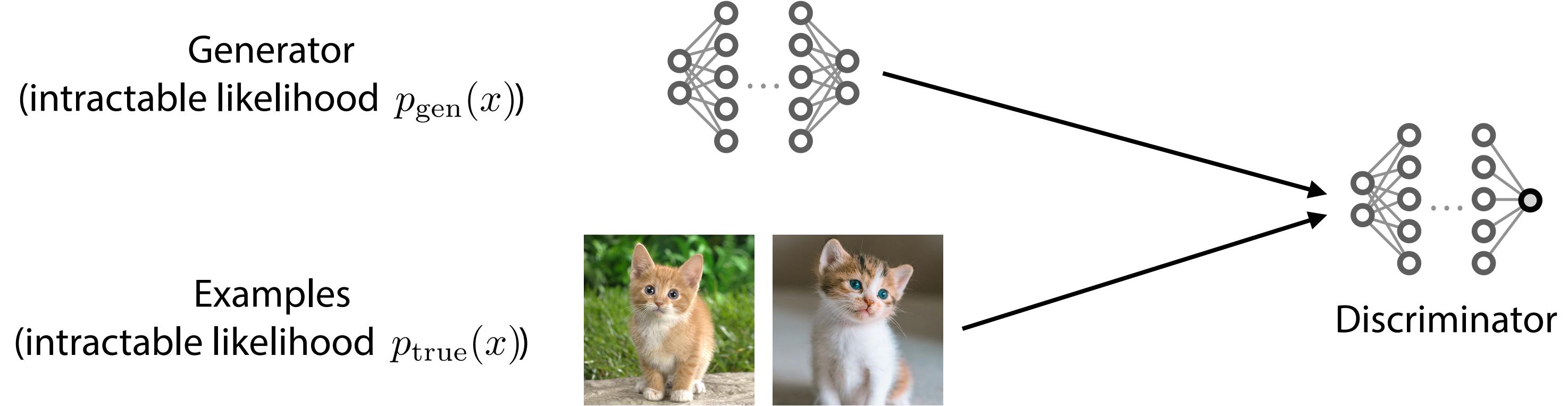
[G. Papamakarios, D. Sterratt, I. Murray 1805.07226;  
J.-M. Lueckmann, G. Bassetto, T. Karaletsos, J. Macke 1805.09294]

- Train conditional neural density estimator (e.g. normalizing flow) as tractable surrogate for simulator
- Scales well to high-dimensional data (no compression to summary stats necessary)
- Amortized: After upfront simulation + training phase, inference is efficient for new data or prior
- Alternative: learn posterior  $\hat{p}(\theta|x_{\text{obs}})$

[G. Papamakarios, I. Murray 1605.06376;  
J.-M. Lueckmann et al. 1711.01861]

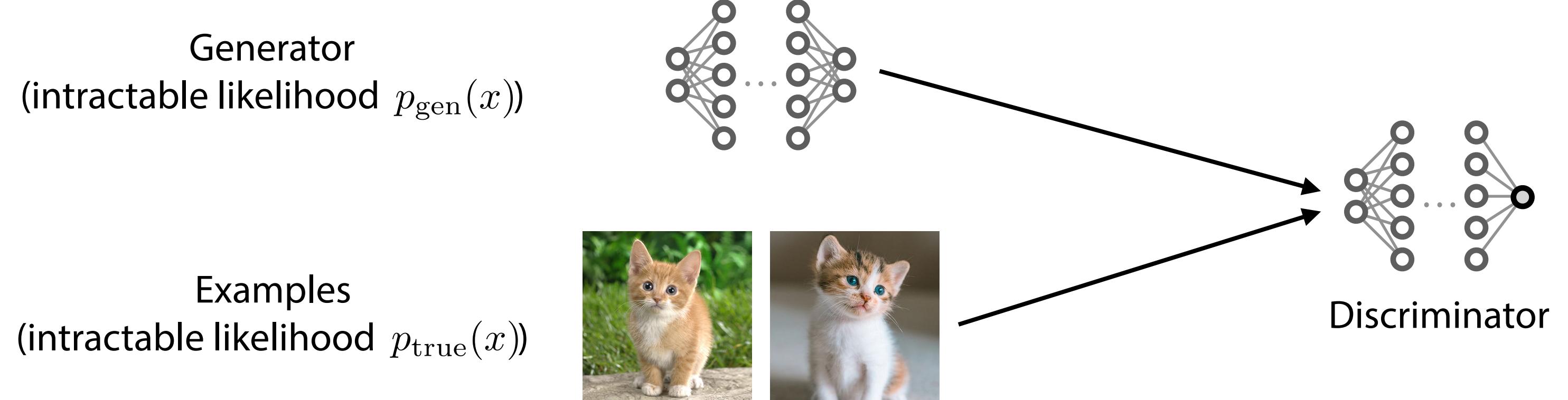
# The likelihood ratio trick

- Remember GANs?



# The likelihood ratio trick

- Remember GANs?



Discriminator learns decision function

$$s(x) \rightarrow \frac{p_{\text{true}}(x)}{p_{\text{gen}}(x) + p_{\text{true}}(x)}$$

# The likelihood ratio trick

- Remember GANs?

Generator  
(intractable likelihood  $p_g(x)$ )



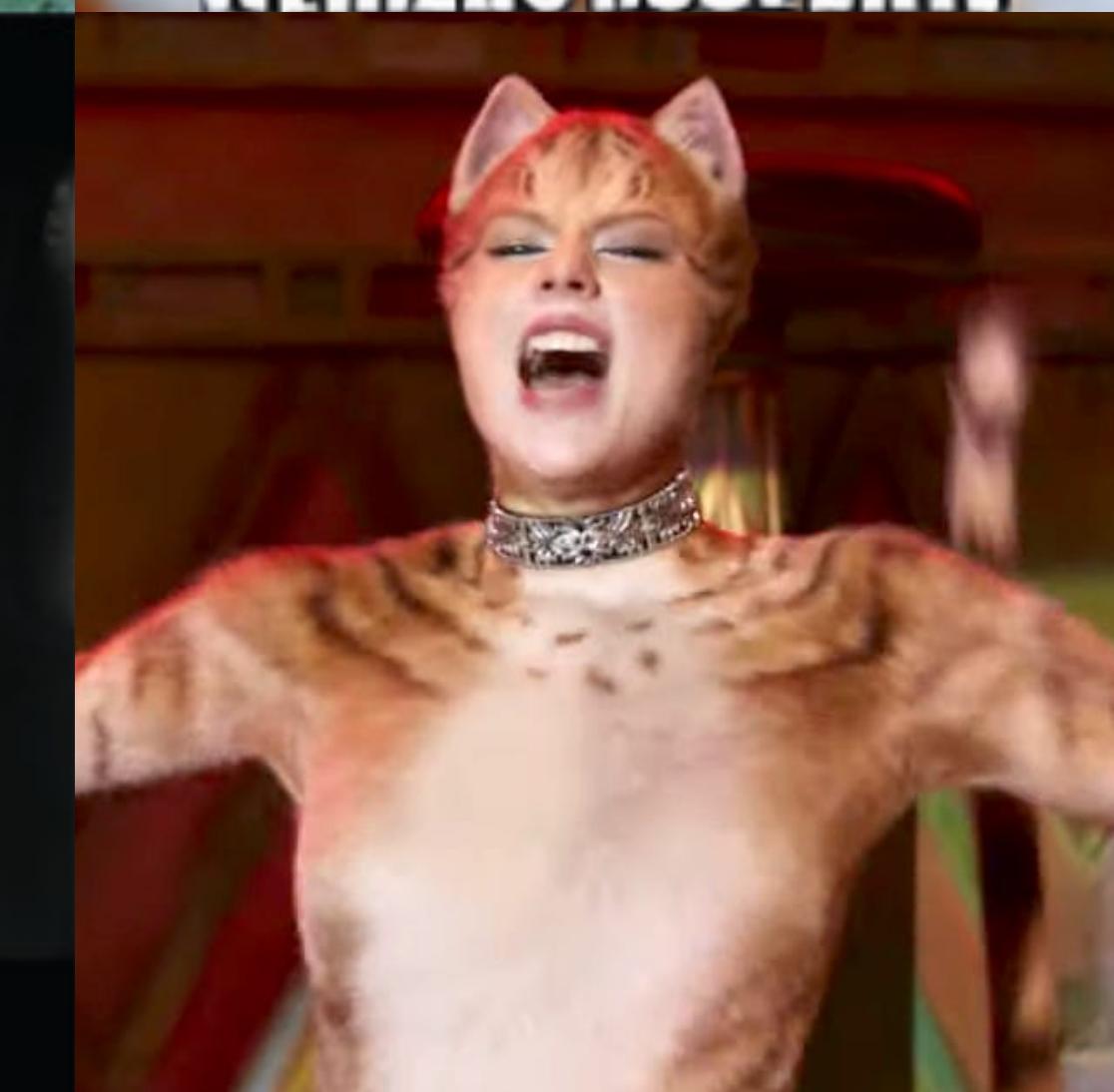
Examples  
(intractable likelihood  $p_t(x)$ )



$$\text{decision function} = \frac{p_{\text{true}}(x)}{p_{\text{true}}(x) + p_{\text{gen}}(x)}$$

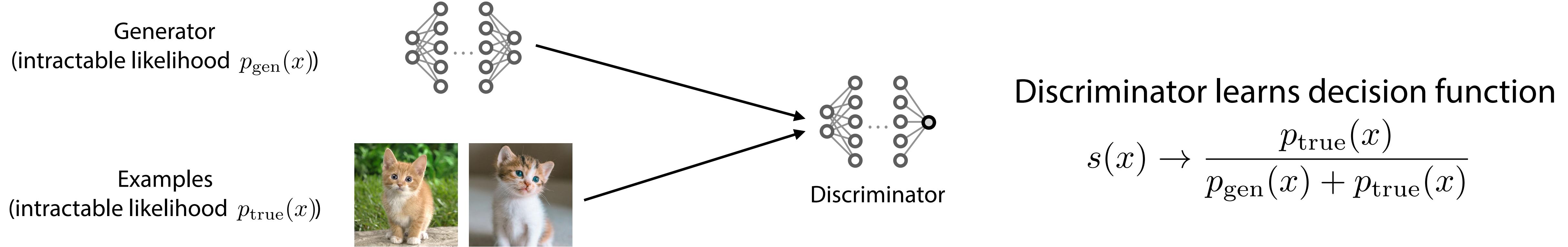


[Nvidia, Universal Pictures]

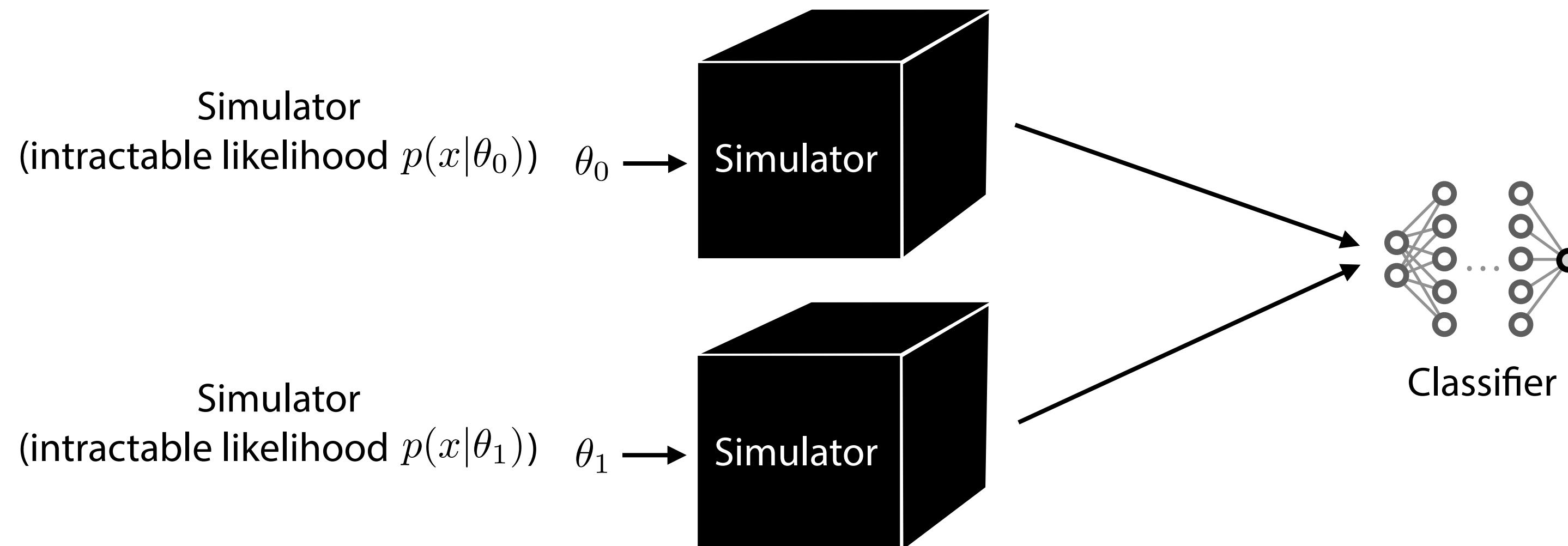


# The likelihood ratio trick

- Remember GANs?

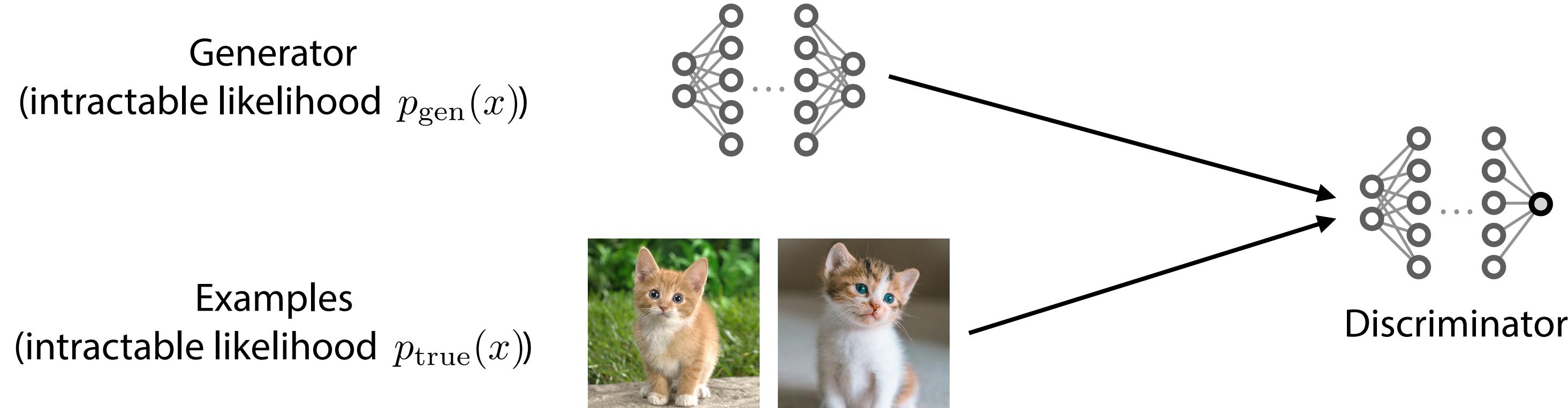


- Similarly, we can train a classifier between two sets of simulated samples



# The likelihood ratio trick

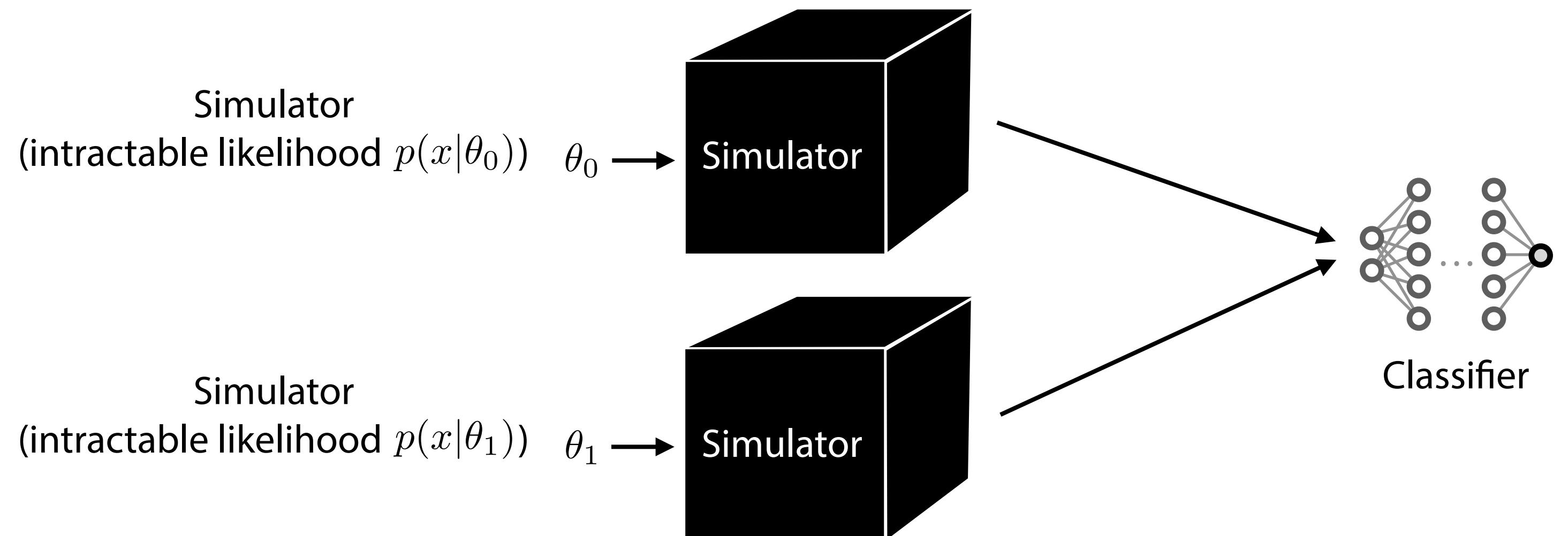
- Remember GANs?



Discriminator learns decision function

$$s(x) \rightarrow \frac{p_{\text{true}}(x)}{p_{\text{gen}}(x) + p_{\text{true}}(x)}$$

- Similarly, we can train a classifier between two sets of simulated samples



Classifier learns decision function

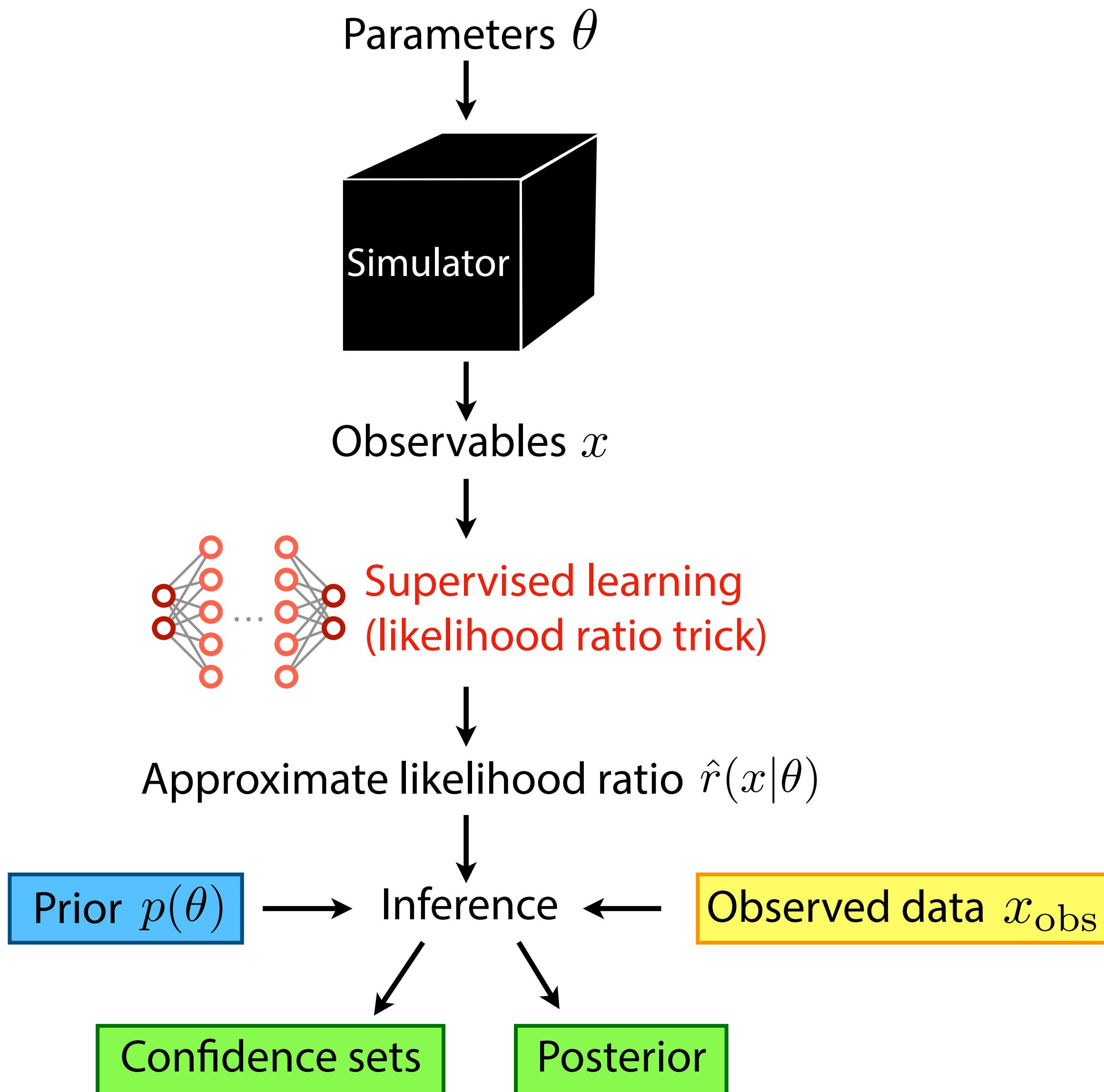
$$s(x) \rightarrow \frac{p(x|\theta_1)}{p(x|\theta_0) + p(x|\theta_1)}$$

⇒ Estimator for likelihood ratio

$$\hat{r}(x) = \frac{1 - s(x)}{s(x)} \rightarrow \frac{p(x|\theta_0)}{p(x|\theta_1)}$$

# Inference by likelihood ratio trick

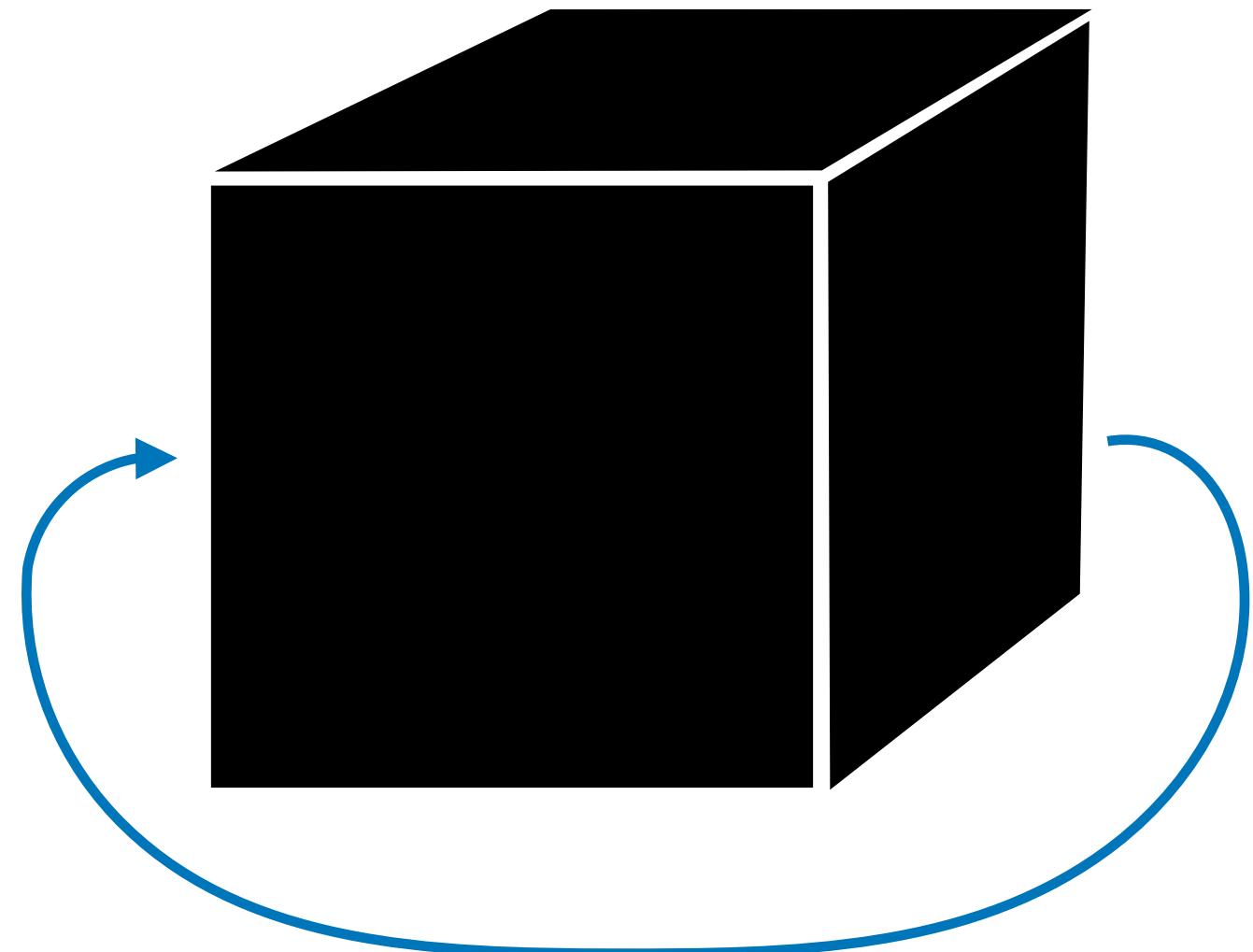
[K. Cranmer J. Pavez, G. Louppe 1506.02169]



- For inference, likelihood and likelihood ratio are interchangeable
- Advantage: Learning the likelihood ratio is often a simpler task than learning the likelihood
- Disadvantage: Cannot sample from likelihood ratio

# Frontiers of simulation-based inference

[K. Cranmer, JB, G. Louppe 1911.01429]

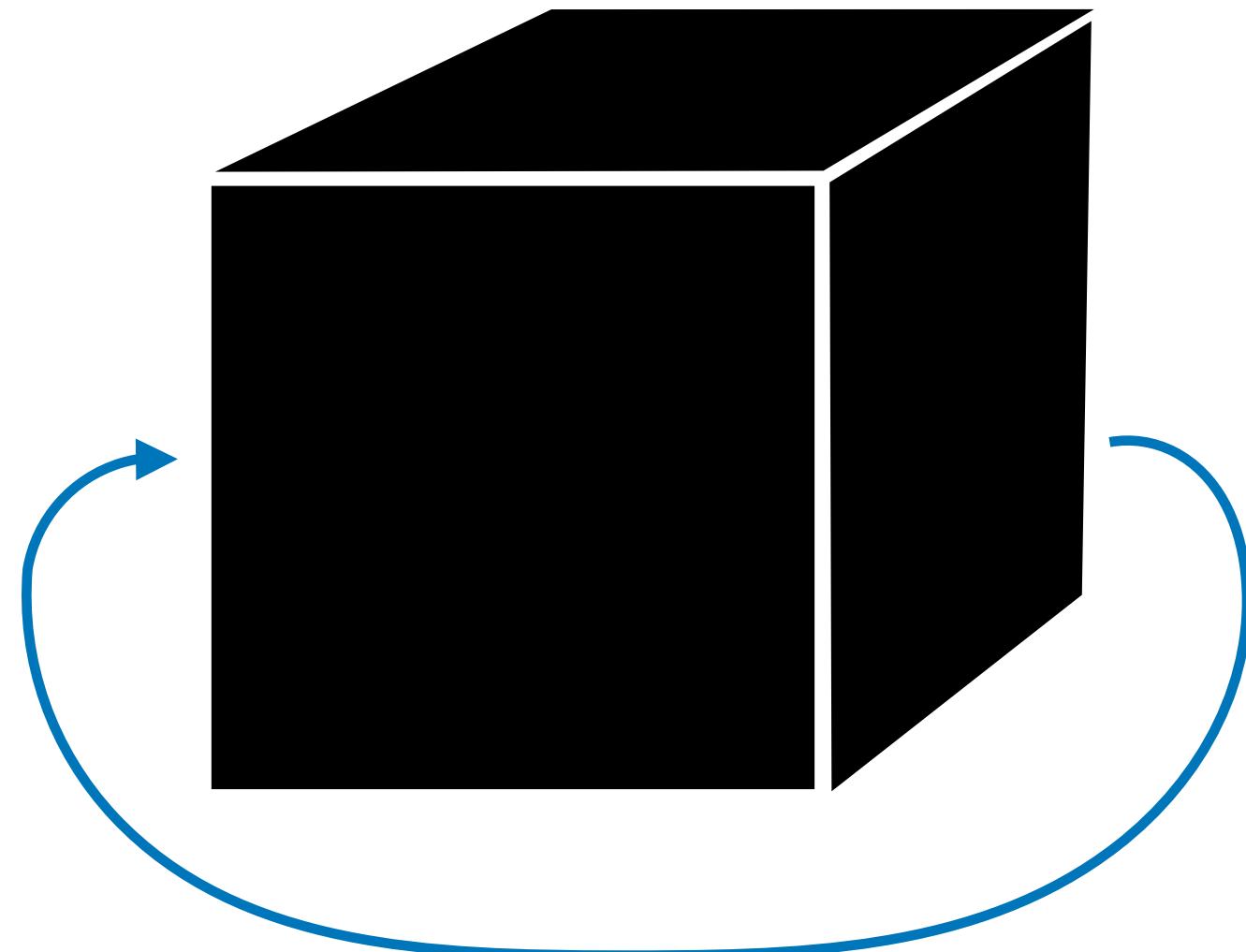


## Active learning:

Iteratively guide simulator to  
important parameter regions  
based on past results

# Frontiers of simulation-based inference

[K. Cranmer, JB, G. Louppe 1911.01429]



## Active learning:

Iteratively guide simulator to important parameter regions based on past results

## Mining gold:

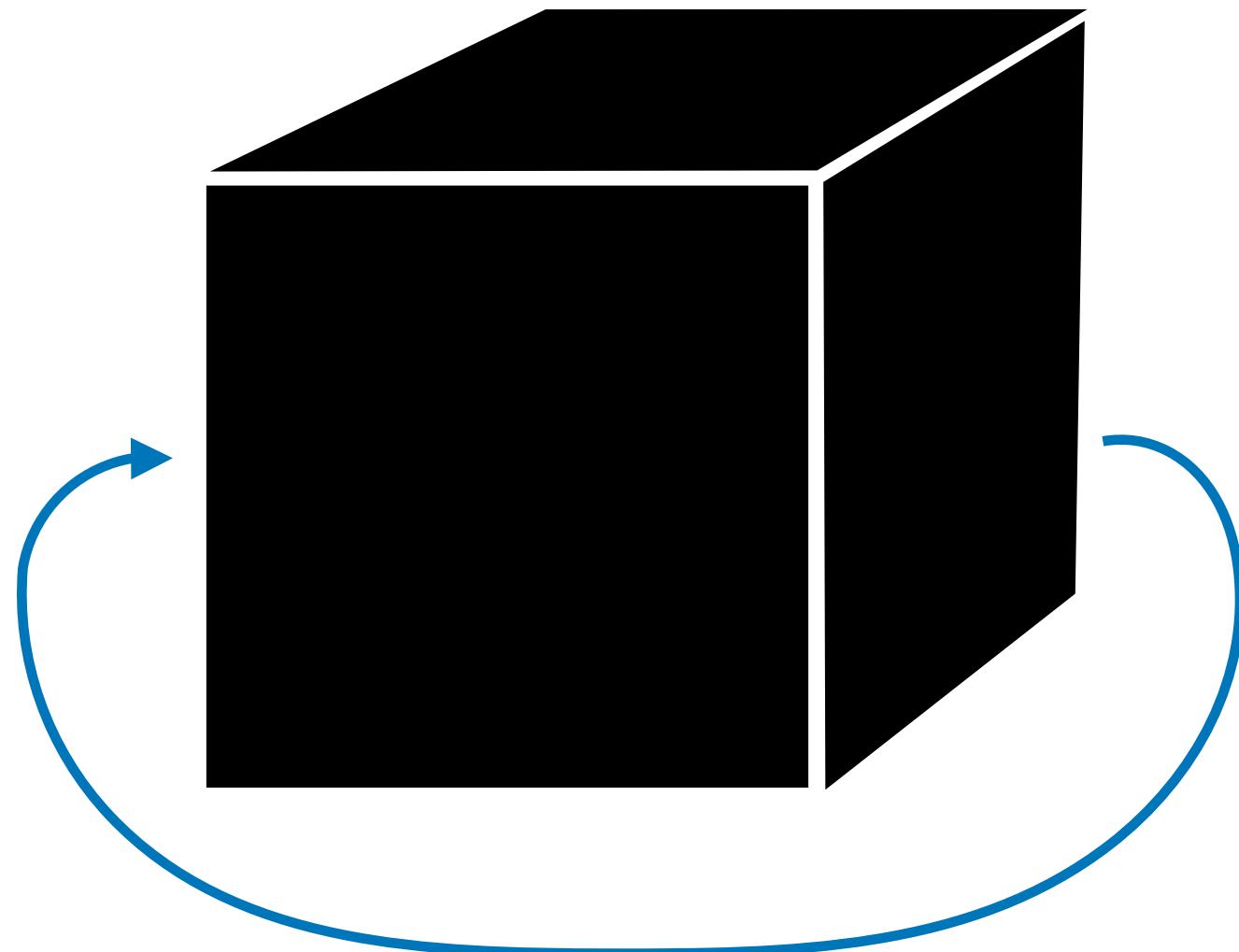
Extract and leverage information from the simulator that characterizes the latent process

[JB, G. Louppe, J. Pavez, K. Cranmer 1805.12244]



# Frontiers of simulation-based inference

[K. Cranmer, JB, G. Louppe 1911.01429]



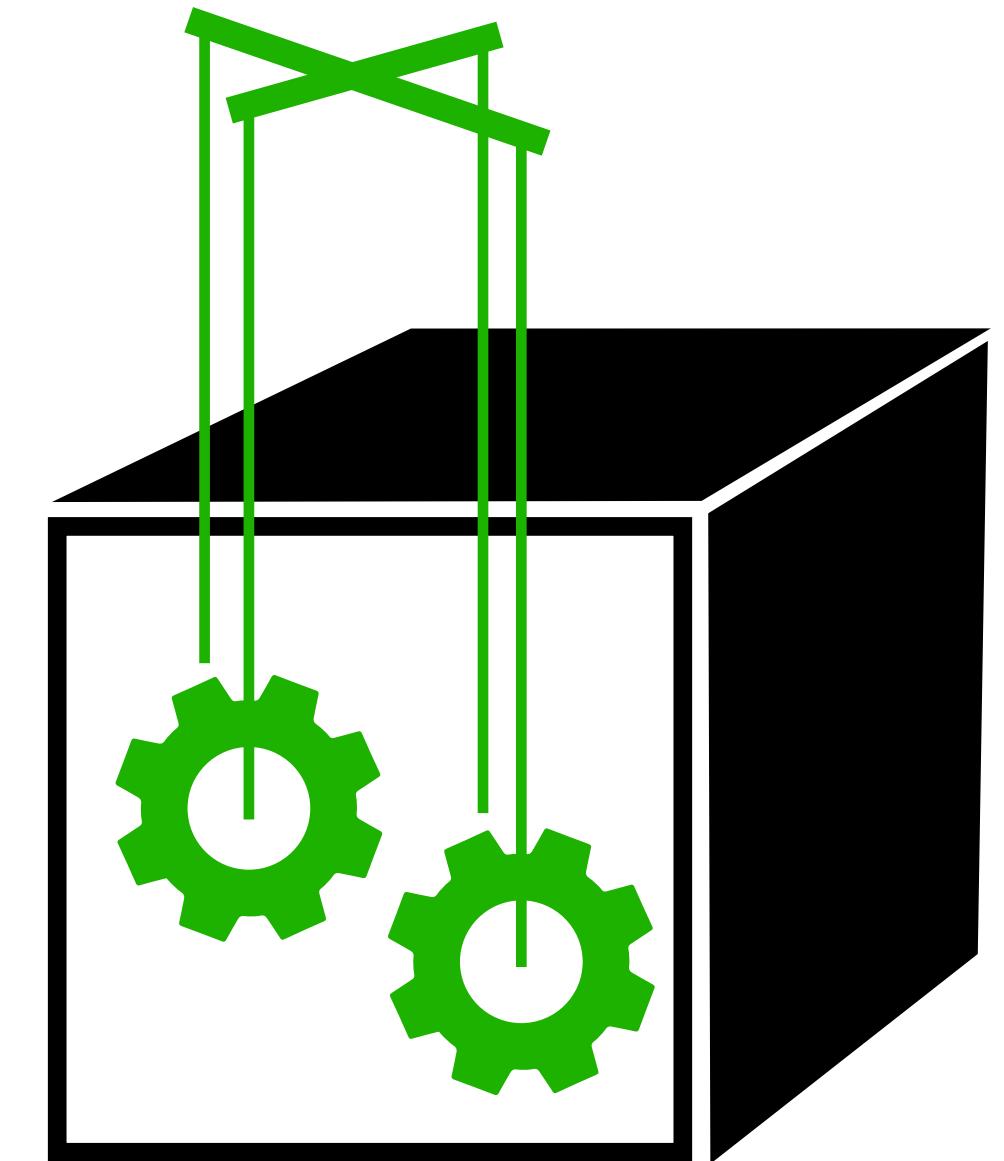
## Active learning:

Iteratively guide simulator to important parameter regions based on past results

## Mining gold:

Extract and leverage information from the simulator that characterizes the latent process

[JB, G. Louppe, J. Pavez, K. Cranmer 1805.12244]



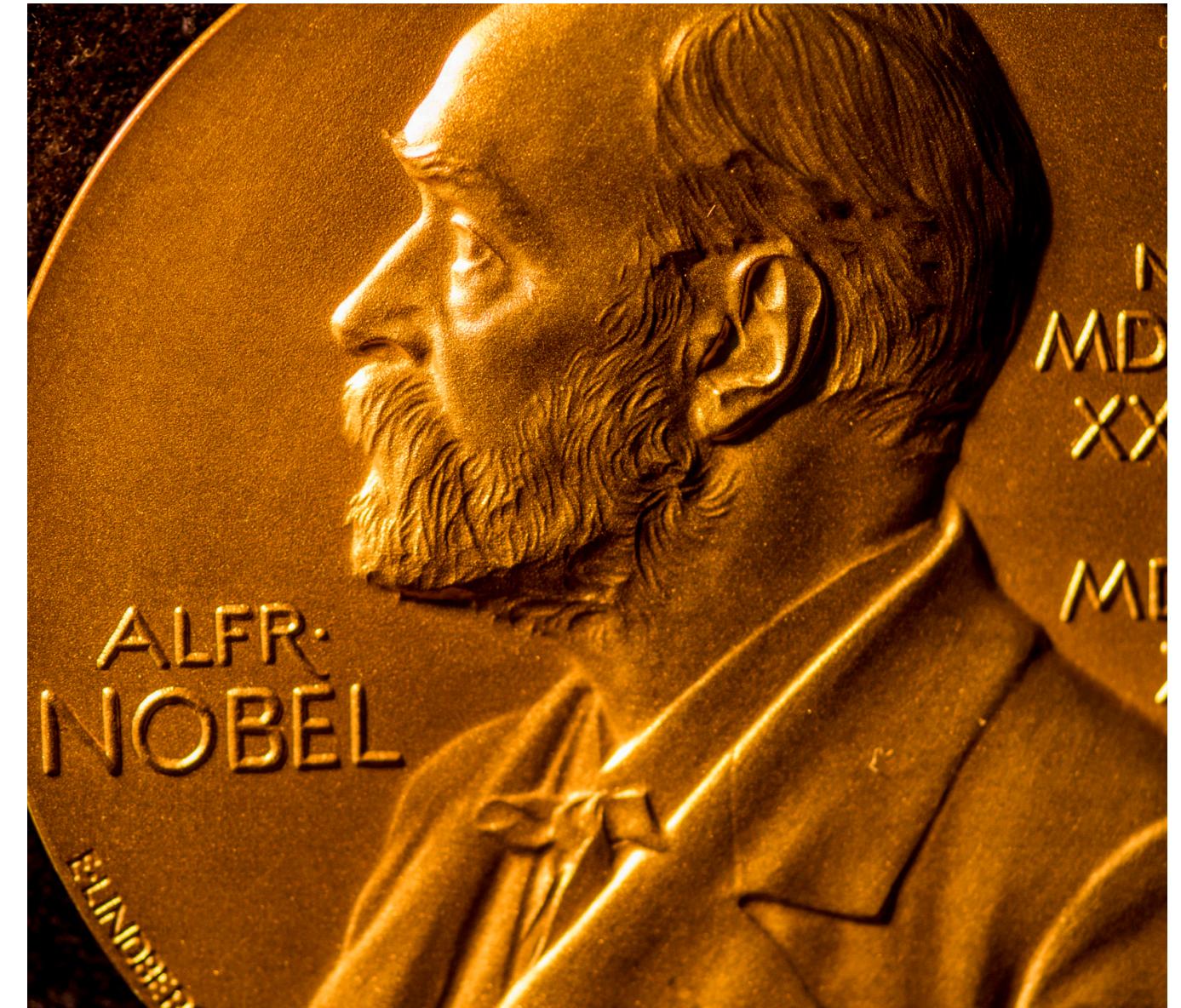
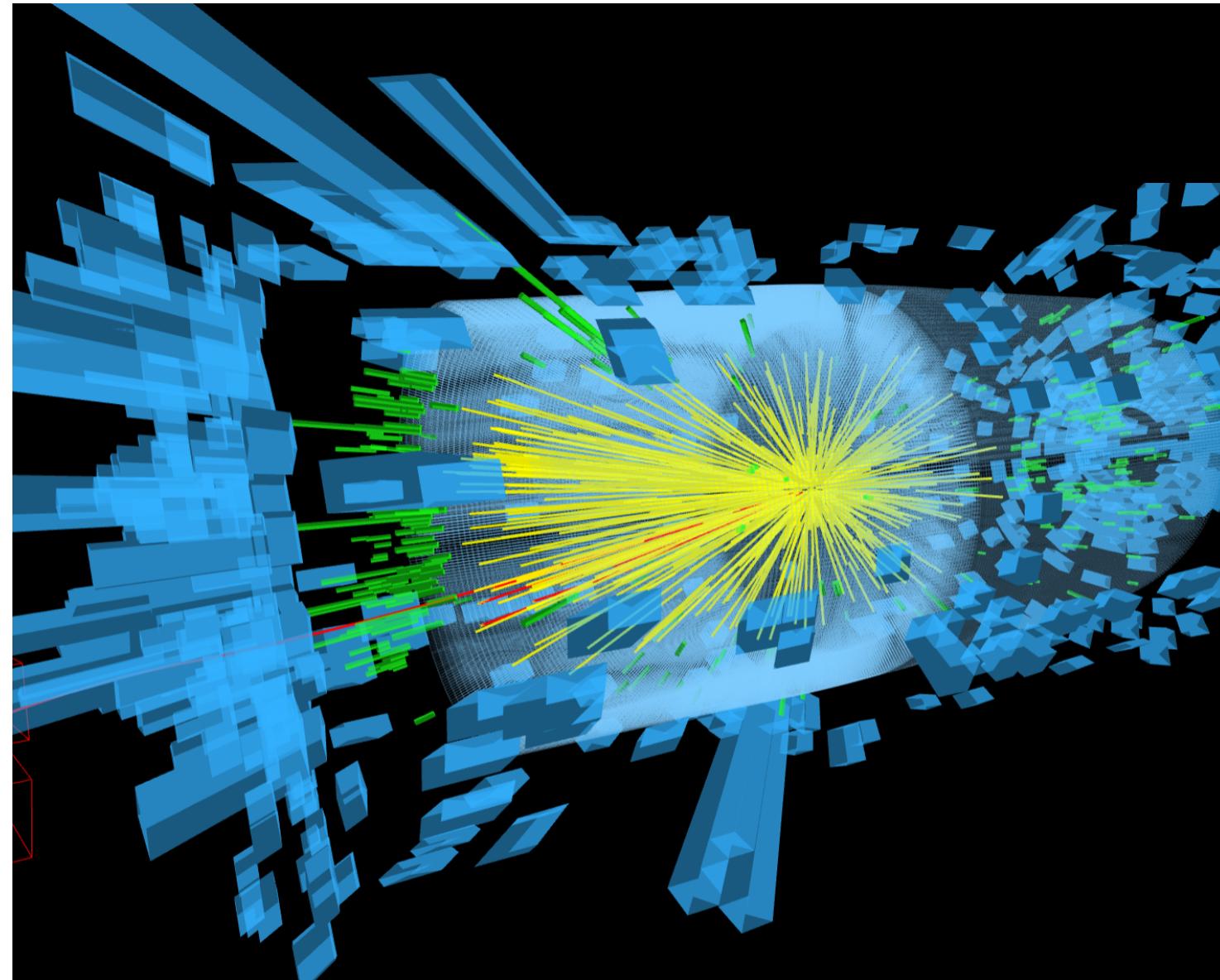
## Probabilistic programming:

Explicitly write simulator as probabilistic program, with ability to condition execution trace on observations

These new methods let us extract more knowledge from particle collisions.

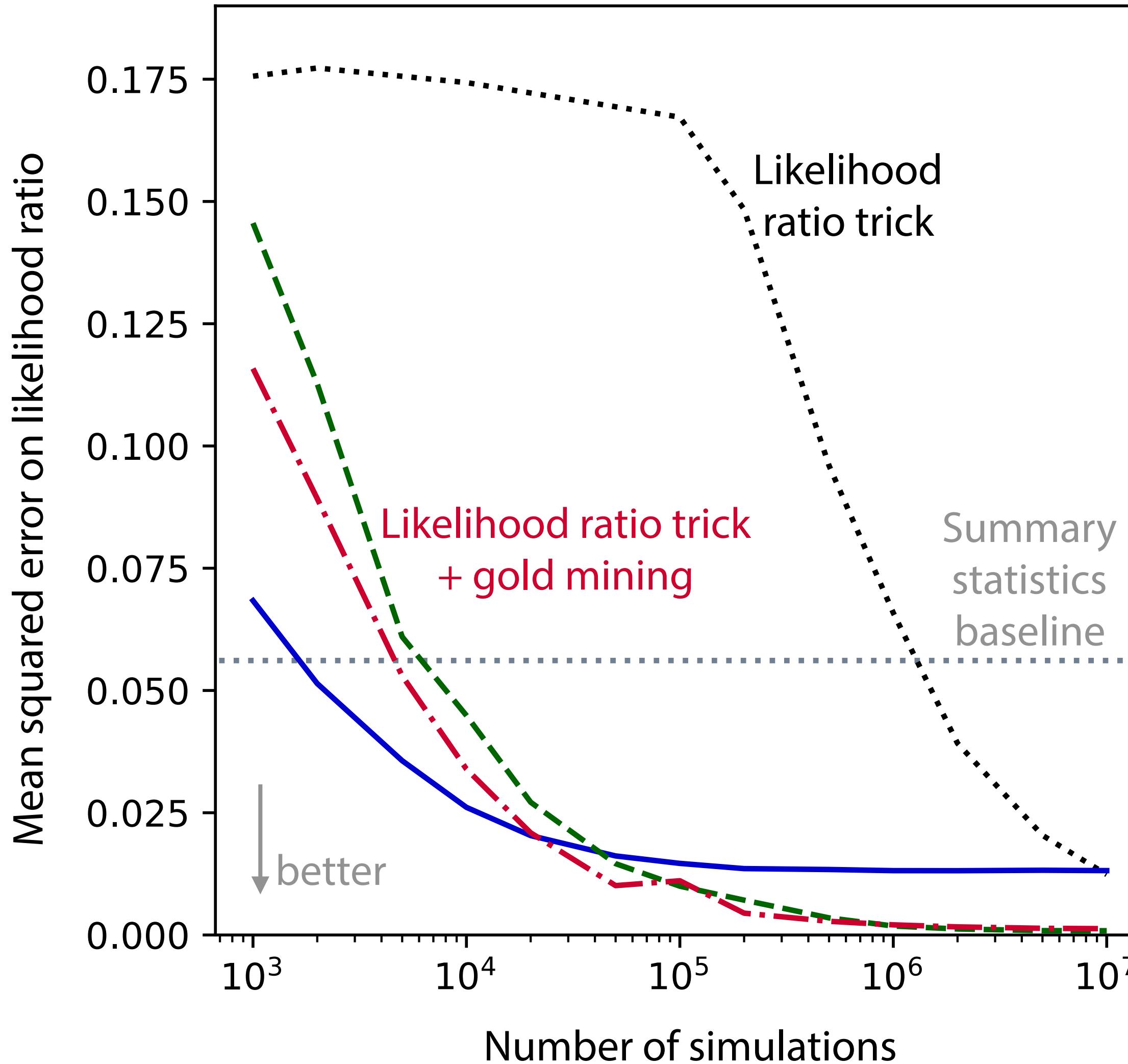
# The Higgs boson

- Discovered 2012 as latest elementary particle, last one missing from the Standard Model



- Measuring its properties might point us to answers to big questions
  - Why is gravity so weak?
  - Why are some particles so light and others so heavy?
  - Why is there more matter than antimatter?
  - Is the vacuum stable?
- We analyzed Higgs properties on synthetic data sets with old and new inference techniques

# Improving sample efficiency

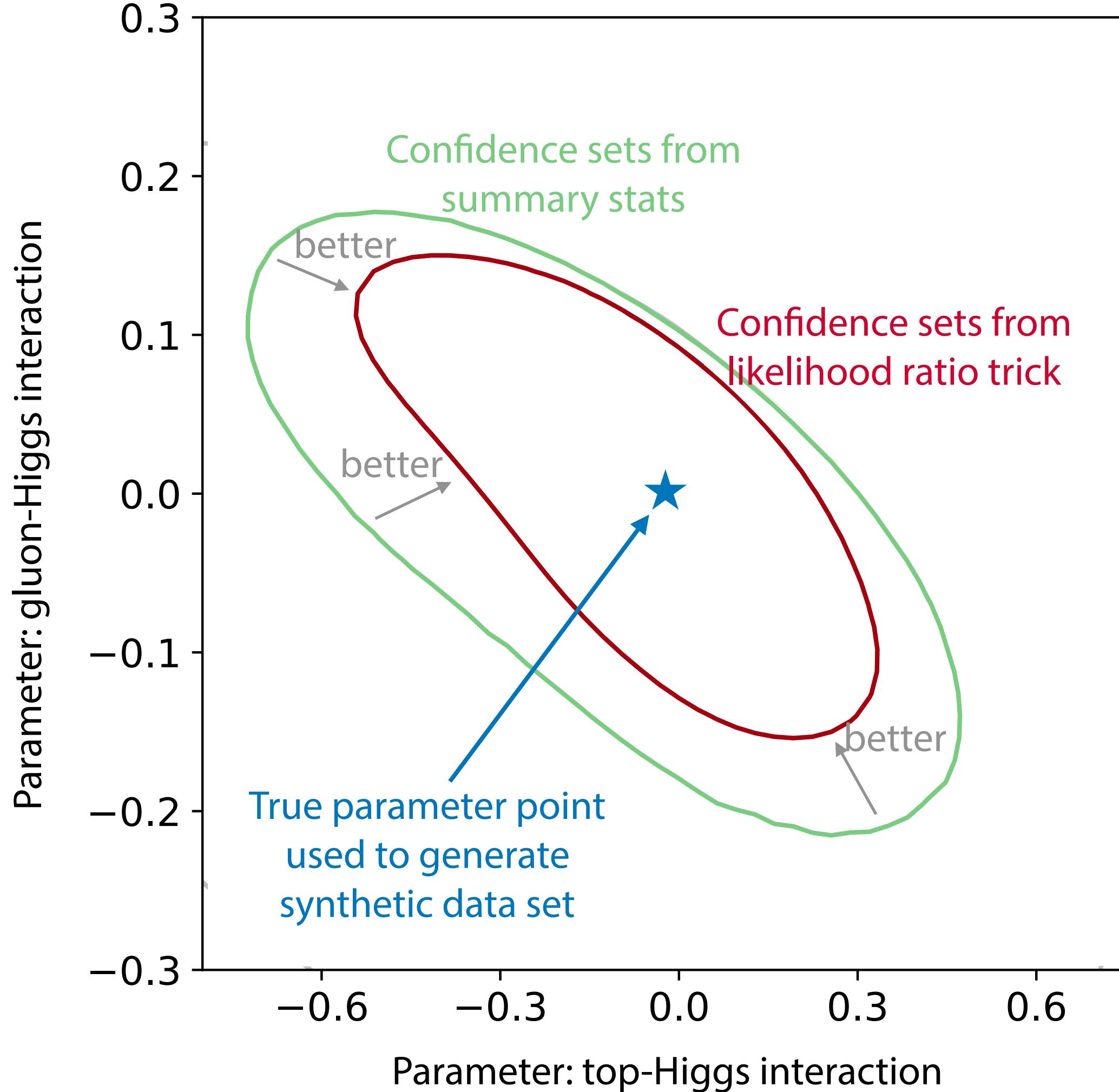


With enough training data, the ML algorithms get the likelihood function right.

Using more information from the simulator improves sample efficiency substantially.

[JB, K. Cranmer, G. Louppe, J. Pavez 1805.00013; 1805.00020;  
M. Stoye, JB, K. Cranmer, G. Louppe, J. Pavez 1808.00973]

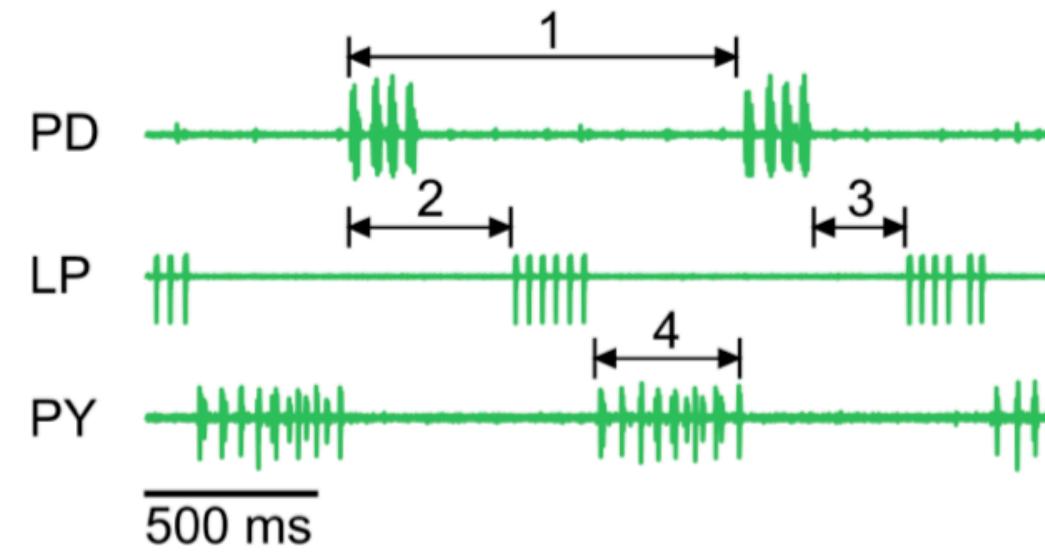
# Improving quality of inference



In some processes, the ML-based inference techniques improve the precision as much as taking 90% more data would!

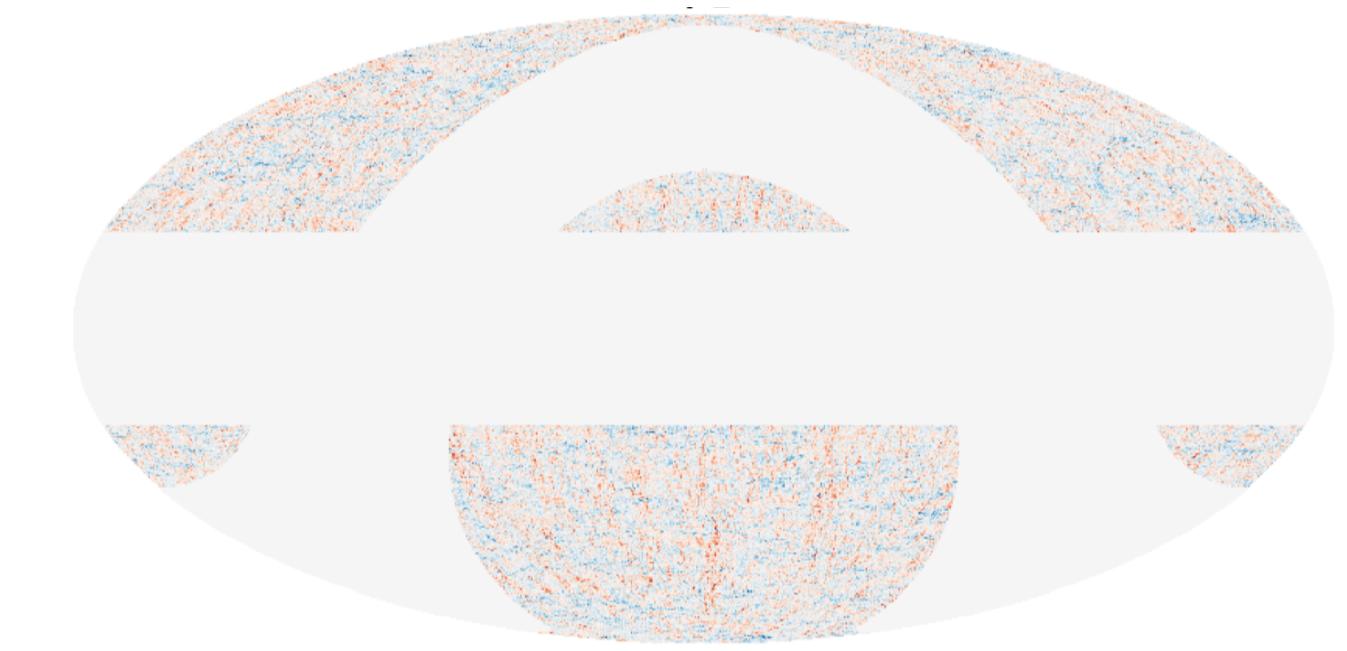
[JB, K. Cranmer, G. Louppe, J. Pavez 1805.00013; 1805.00020;  
JB, F. Kling, I. Espejo, K. Cranmer 1907.10621]

# Simulation-based inference is thriving



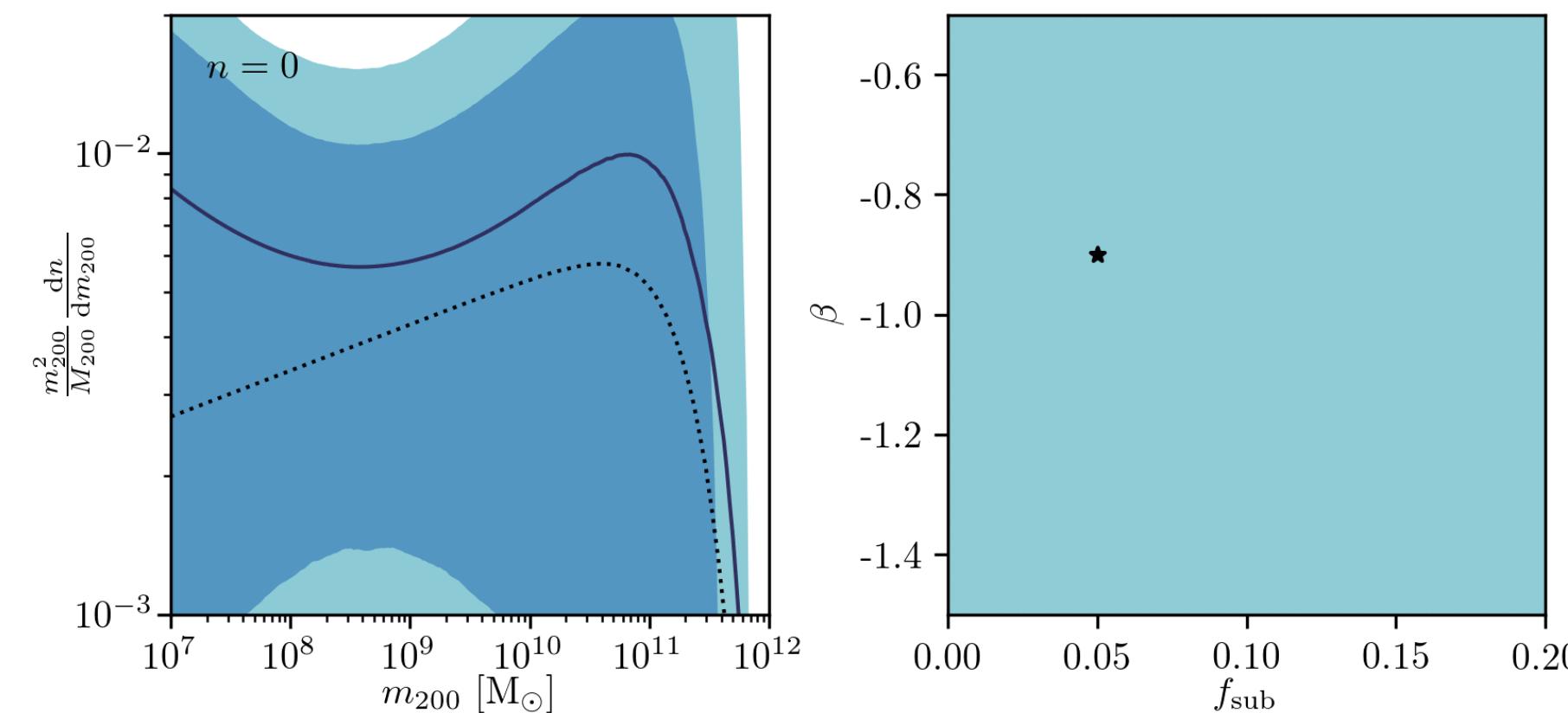
Neuroscience

[Pedro J. Goncalves et al. biorxiv:10.1101/838383]



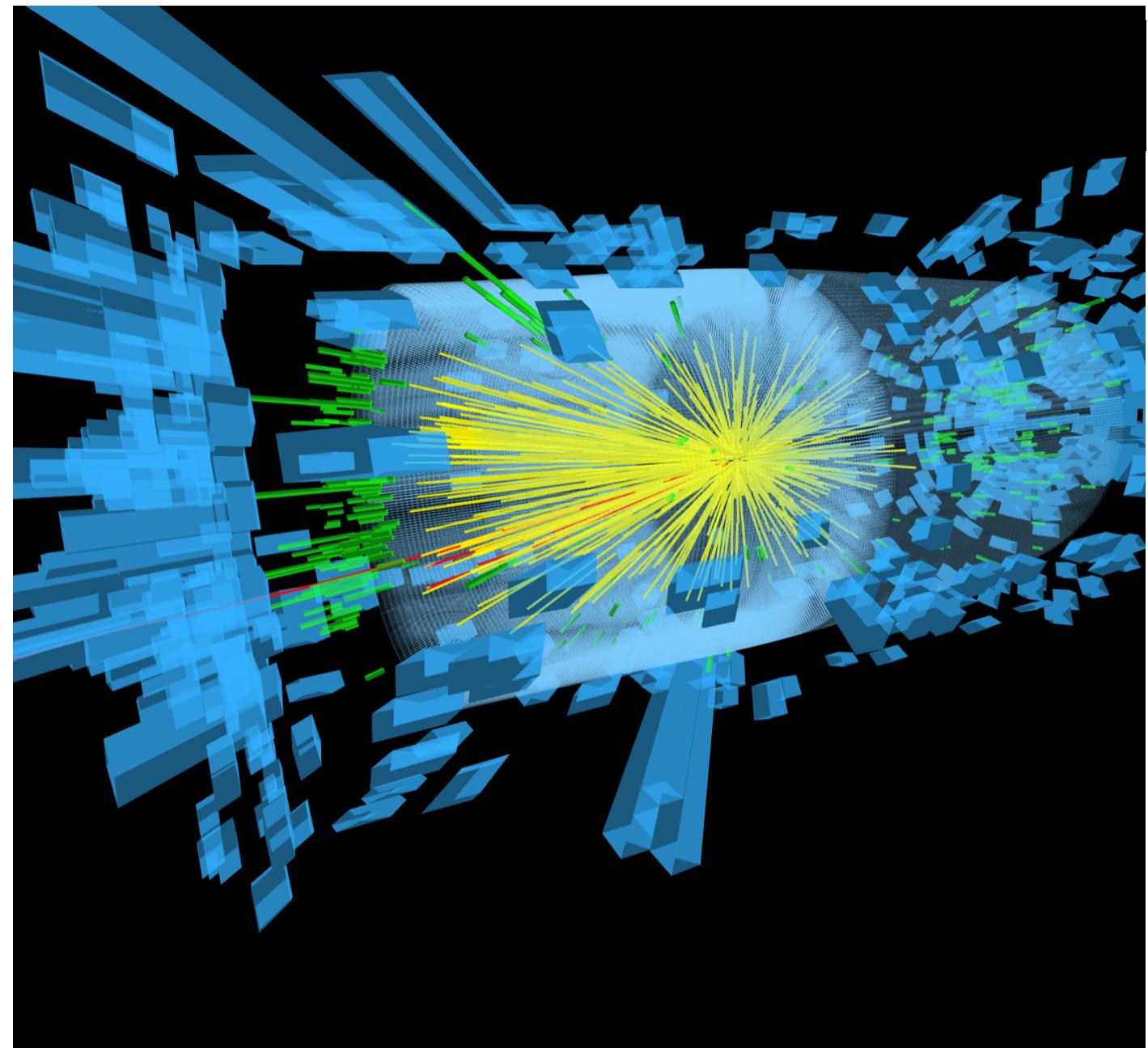
Cosmological large-scale structure

[P. Taylor et al. 1904.05364]

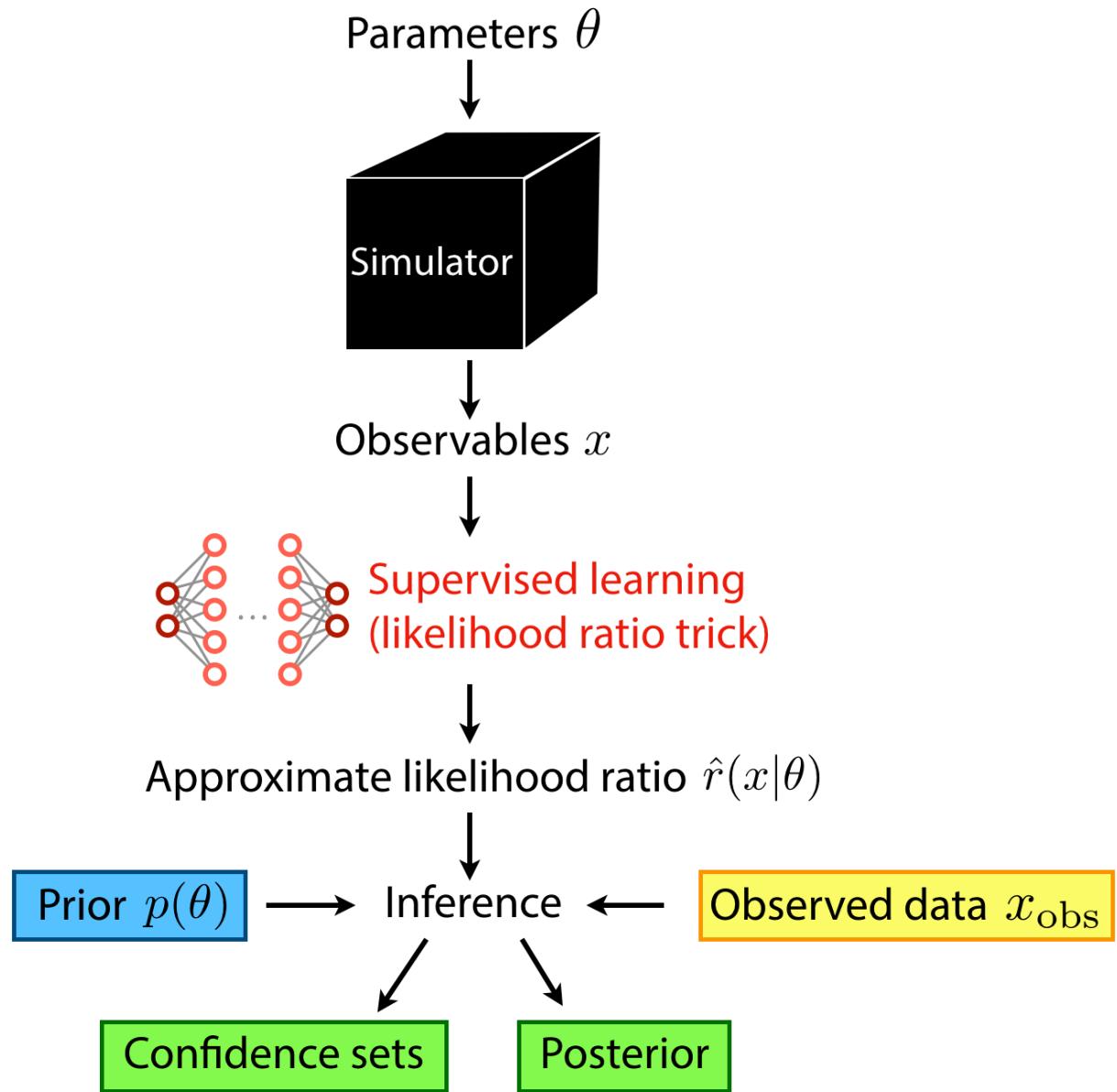


Strong gravitational lensing

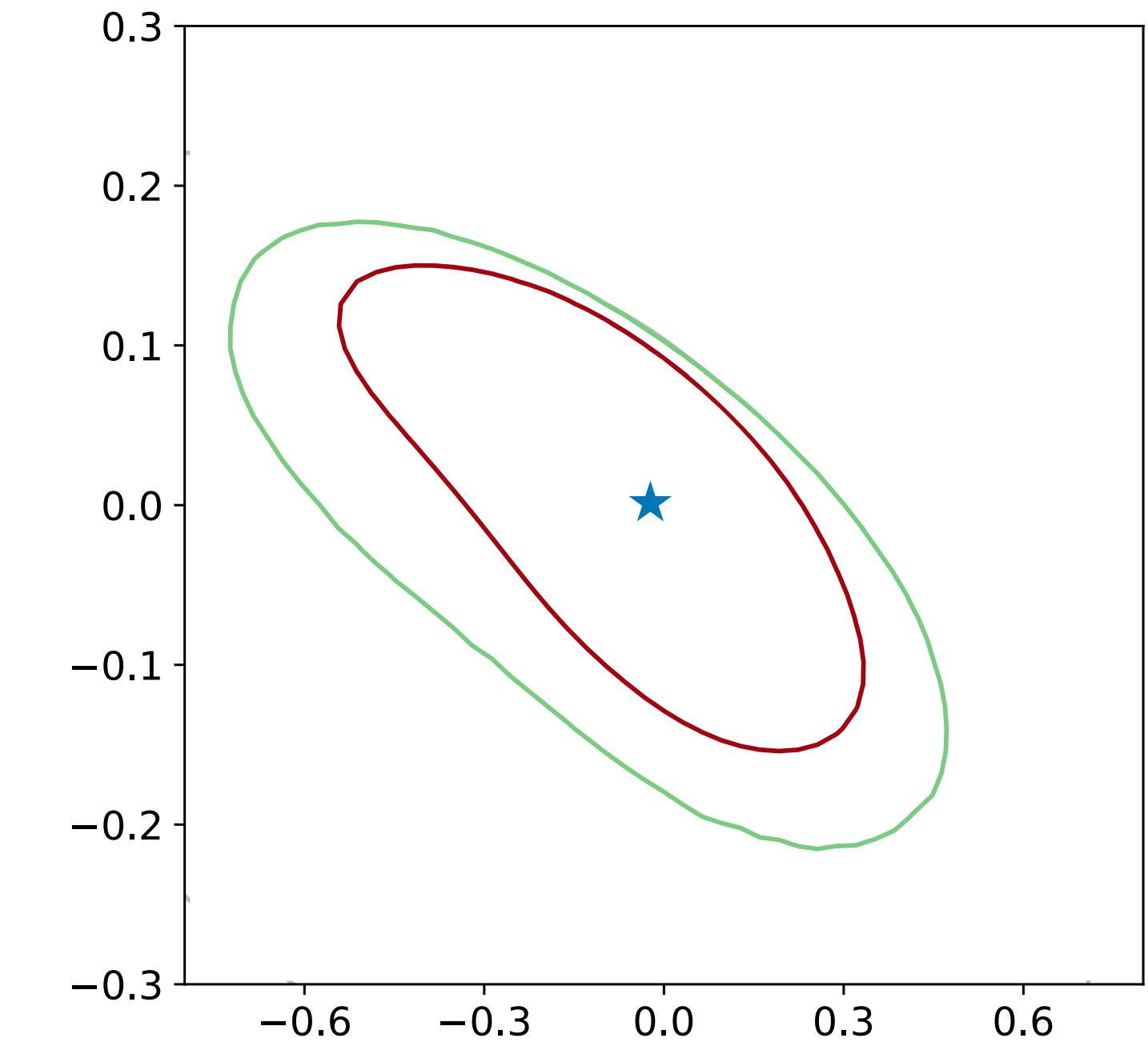
[JB, S. Mishra-Sharma et al. 1909.02005]



Scientific inference is challenging when phenomena are modeled with computer simulations



New inference algorithms are based on normalizing flows or the likelihood ratio trick



They let us extract more knowledge e.g. from particle physics experiments

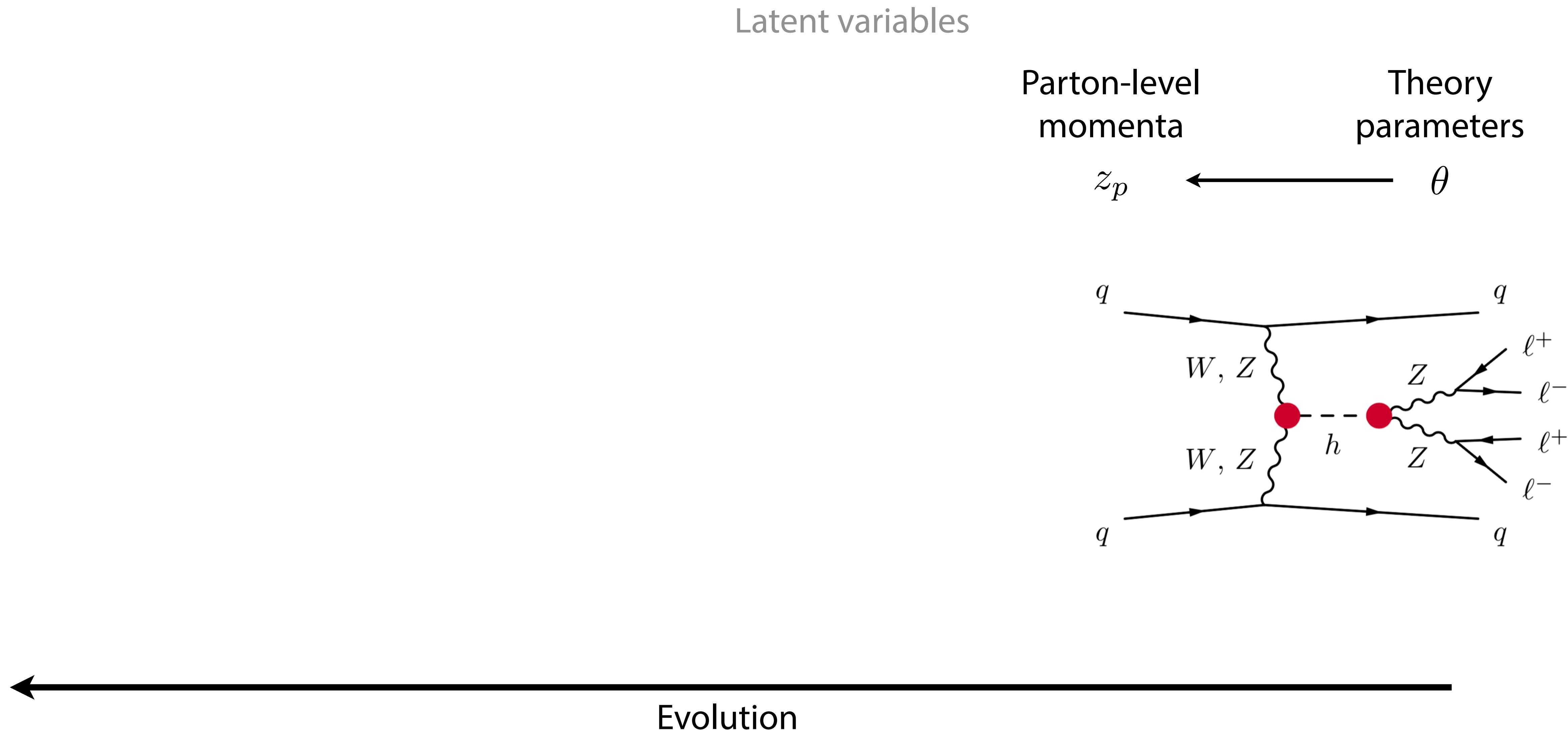
# Bonus material

# Modelling particle physics processes

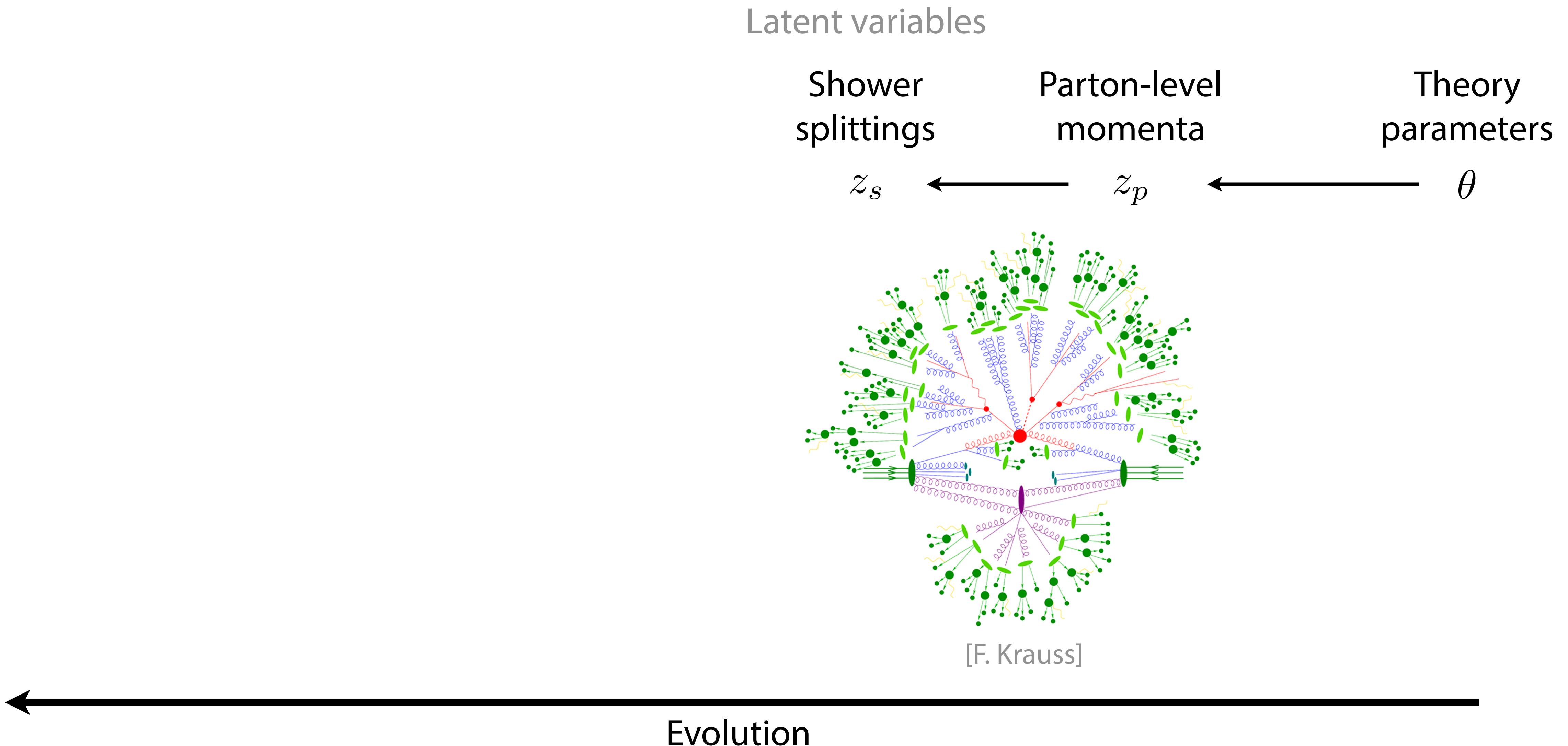
Theory  
parameters  
 $\theta$



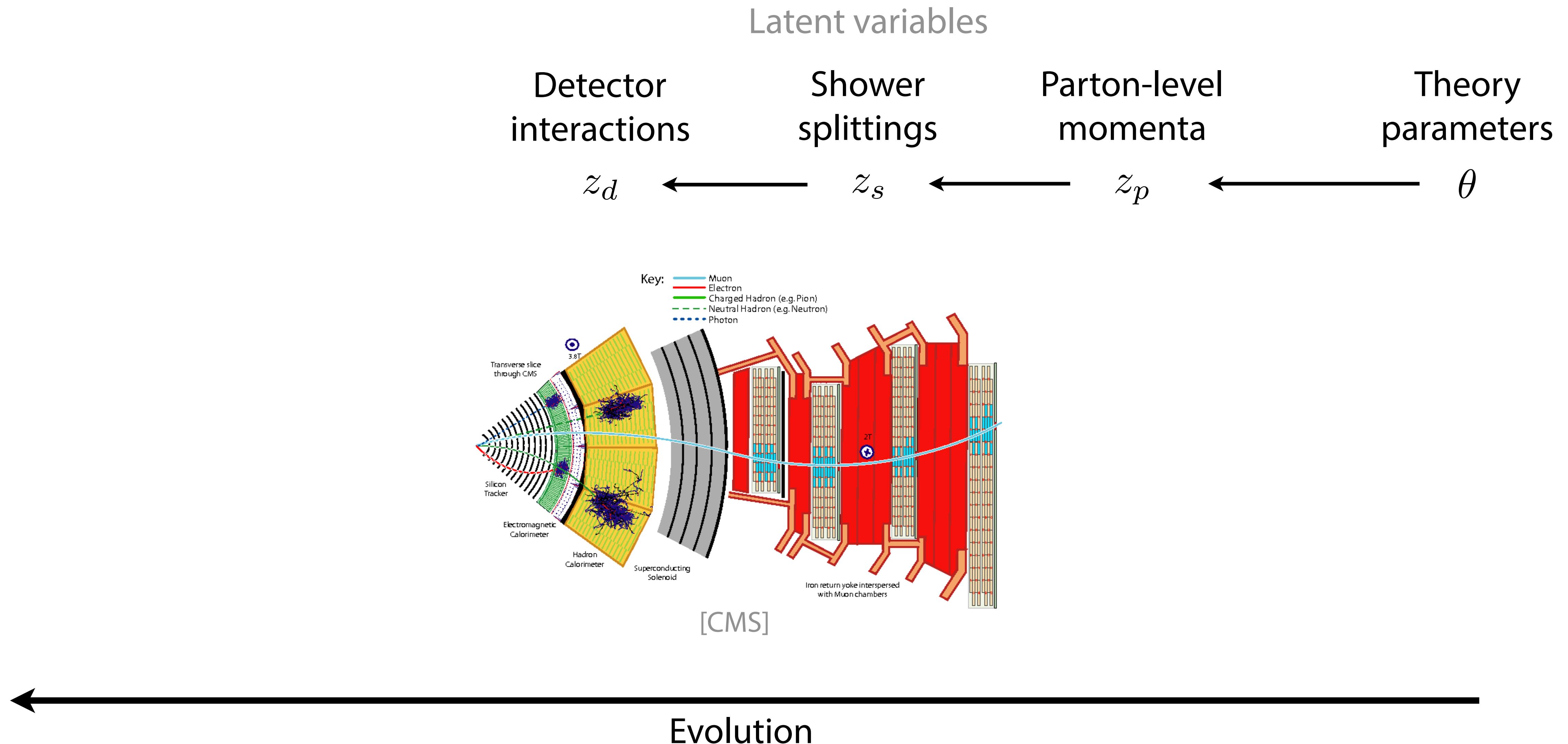
# Modelling particle physics processes



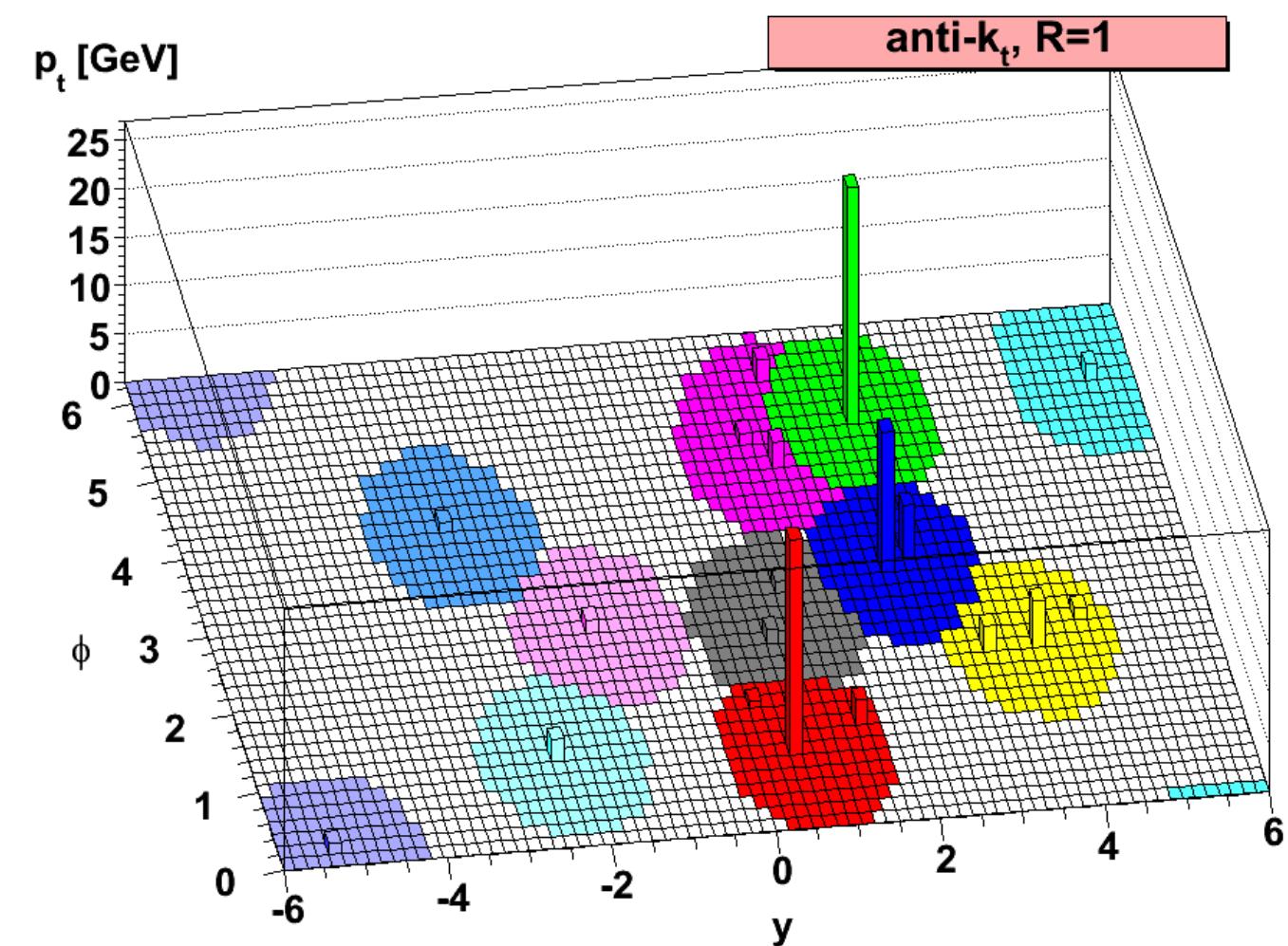
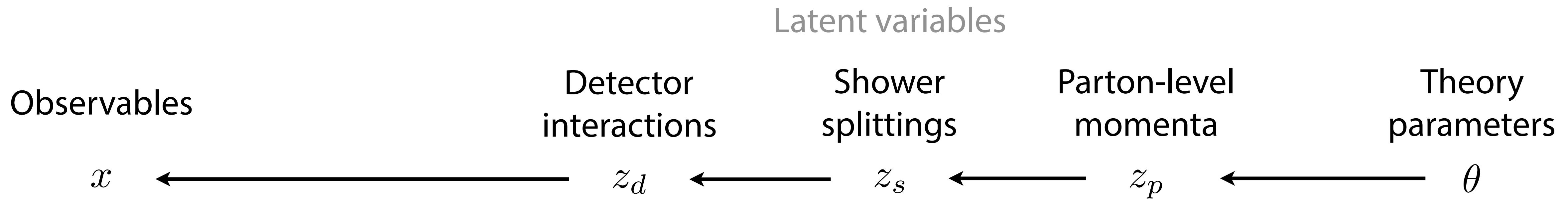
# Modelling particle physics processes



# Modelling particle physics processes

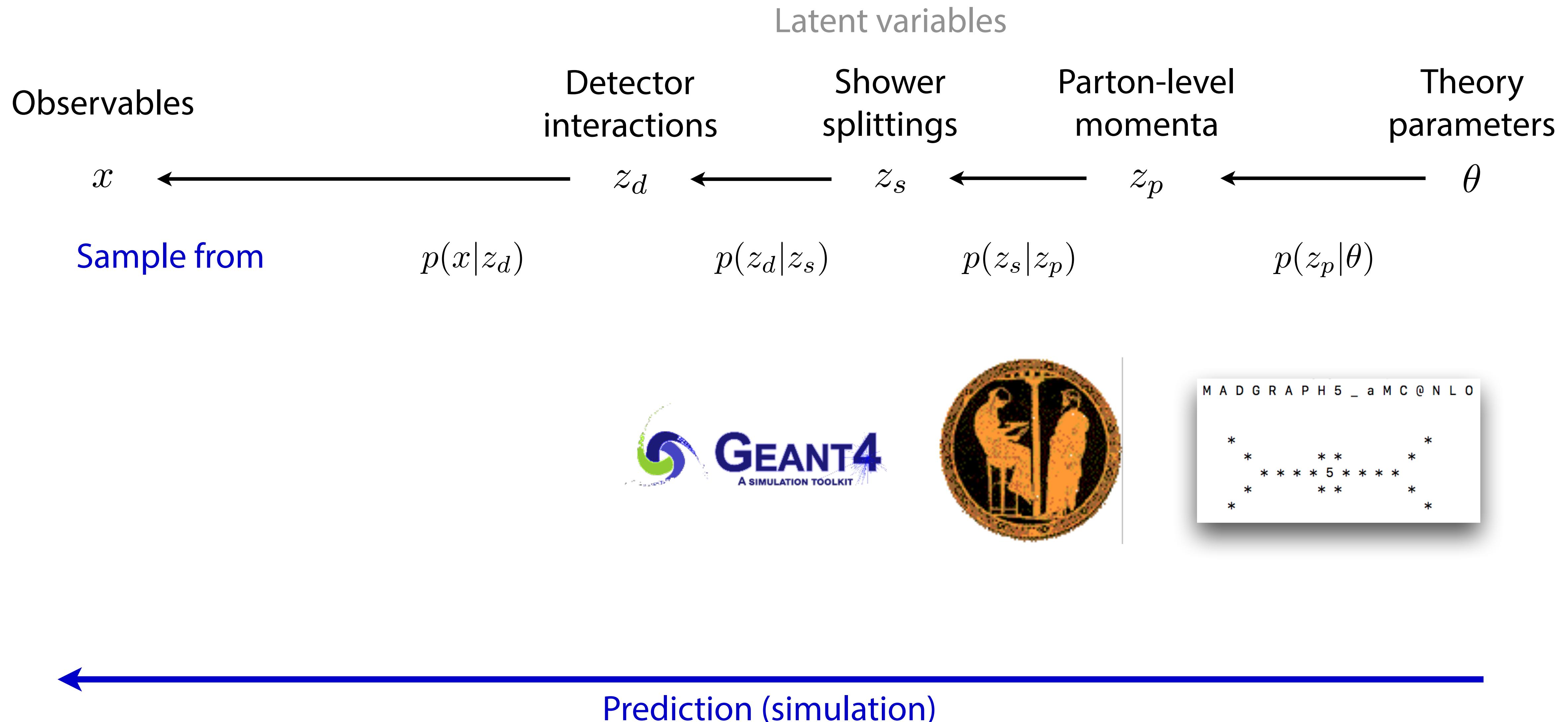


# Modelling particle physics processes

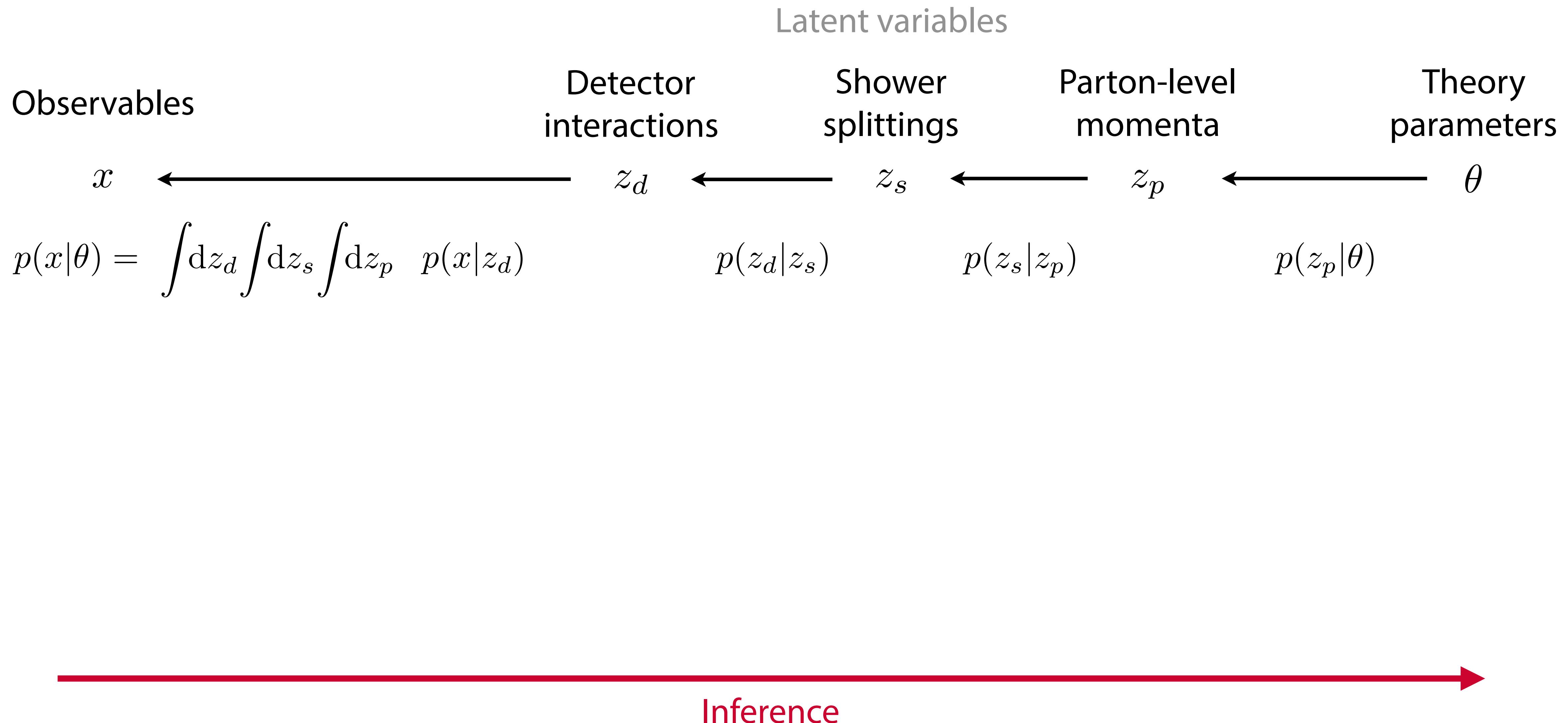


[M. Cacciari, G. Salam, G. Soyez 0802.1189]

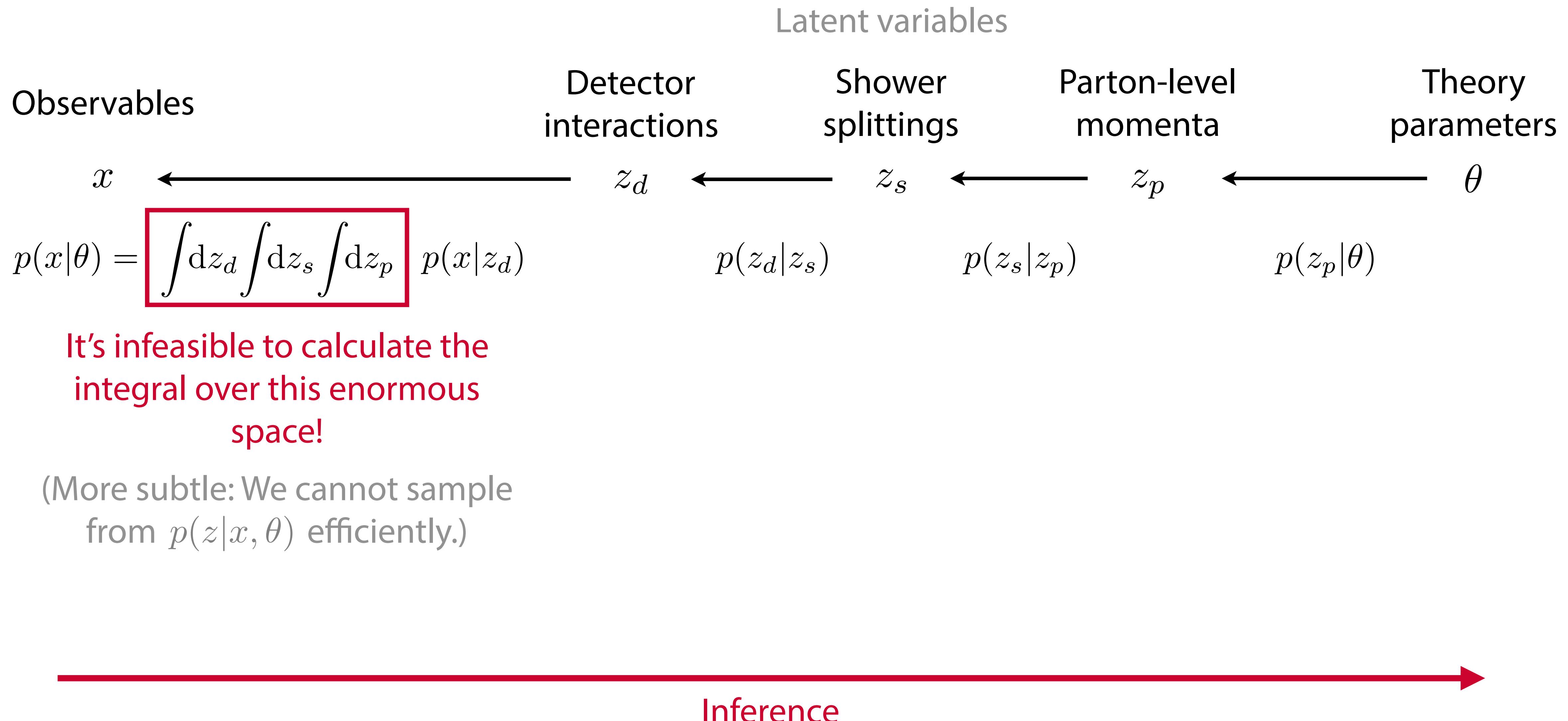
# Modelling particle physics processes



# Modelling particle physics processes



# Modelling particle physics processes



# Simulation-based inference methods

Any simulator

Use simulator directly during inference

Use surrogate model during inference

Approximate Bayesian Computation

Likelihood with summary stats

Neural likelihood

Neural posterior

Neural likelihood ratio

Additional requirements

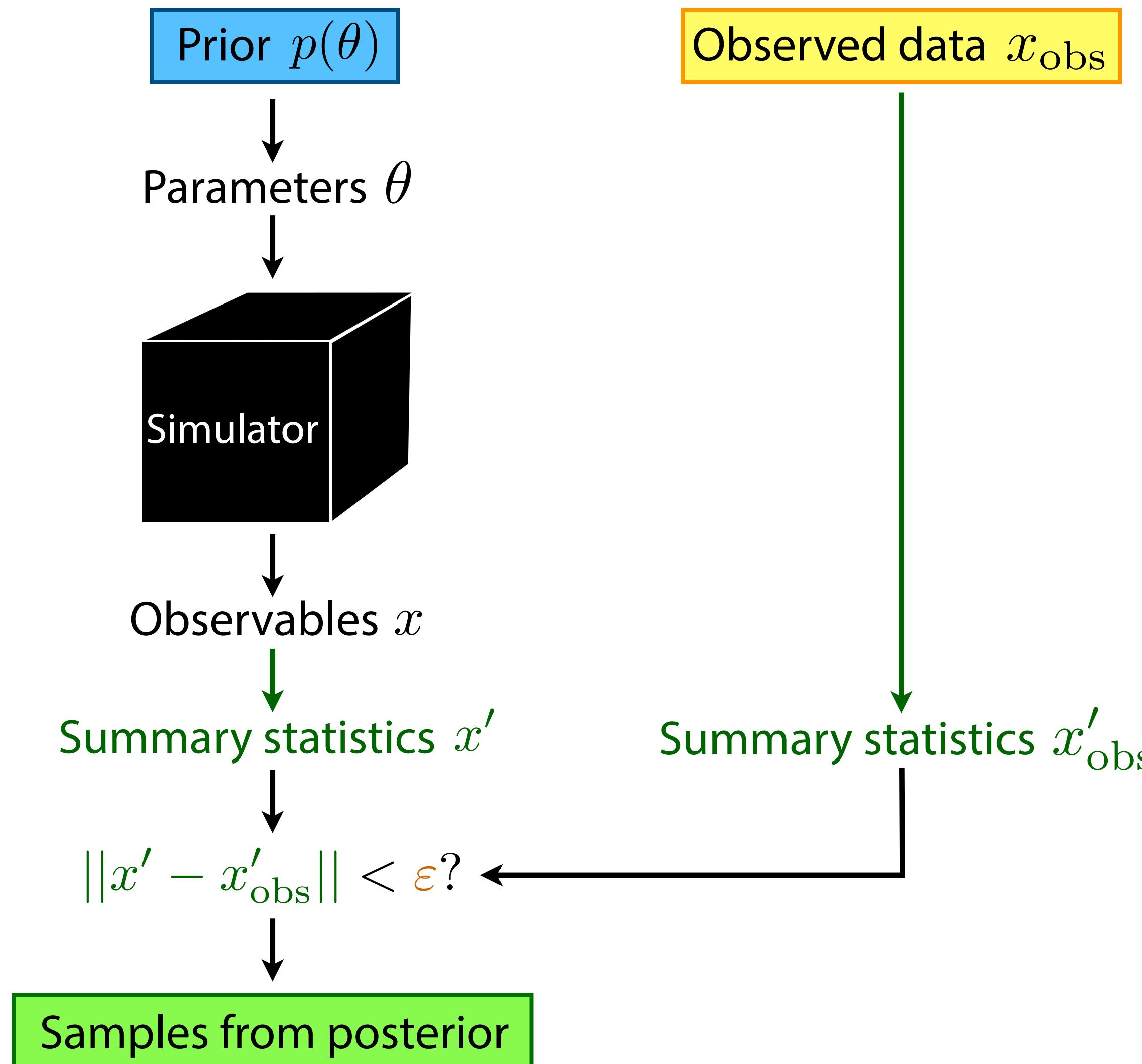
Probabilistic programming

Exact Bayesian inference for diff. sim.

Gold mining

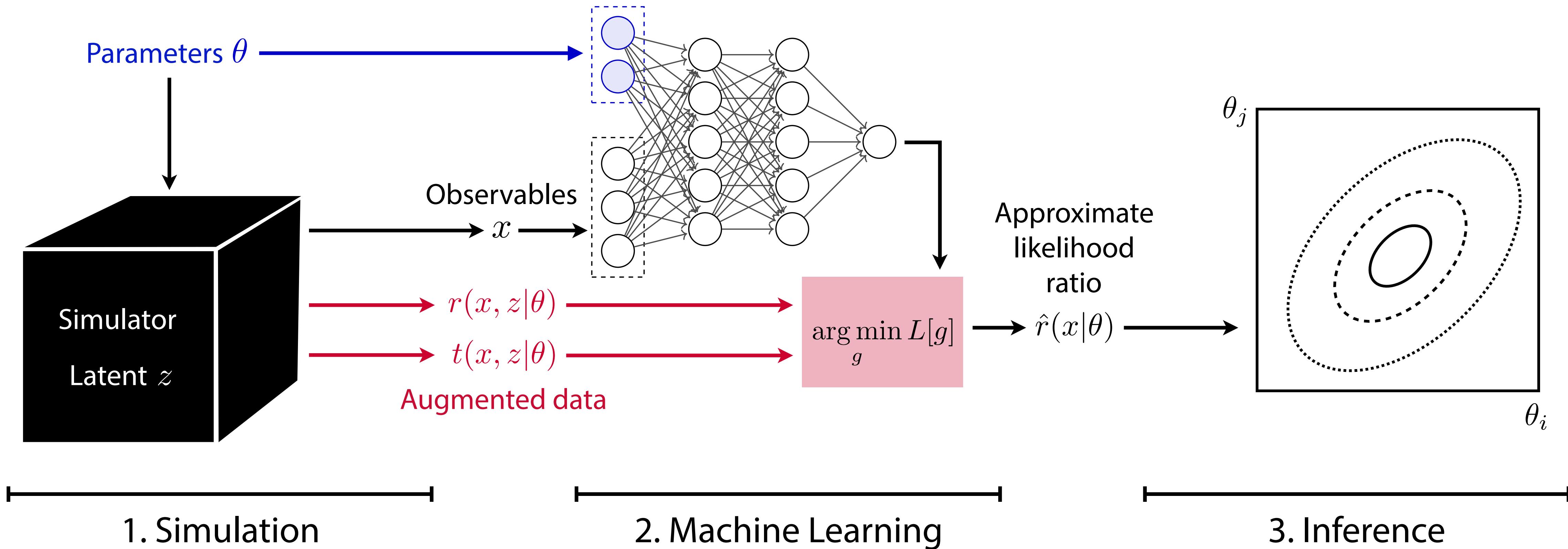
# Approximate Bayesian Computation (ABC)

[D. Rubin 1984]



- Quality of inference:  
low-dimensional **summary statistics**  $x'$  and large **bandwidths**  $\varepsilon$  reduce statistical power
- Sample efficiency:  
high-dimensional **summary statistics**  $x'$  and small **bandwidths**  $\varepsilon$  require unfeasibly many simulations

# Mining gold

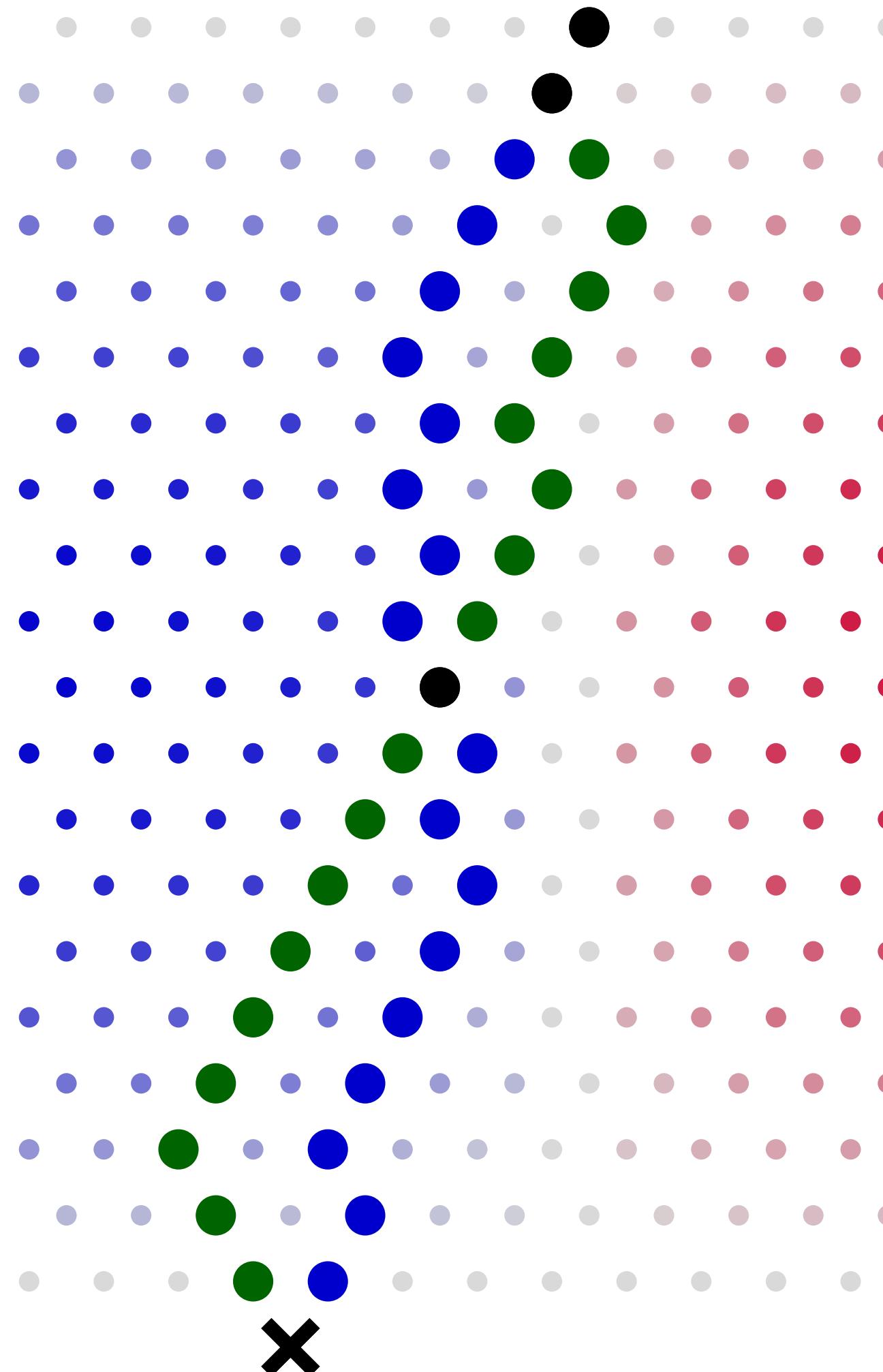


“Mining gold”: Extract additional information from simulator

Use this information to train estimator for likelihood ratio

Limit setting with standard hypothesis tests

# Mining gold on the Galton board



- Remember: the likelihood

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

is intractable because of the integral over all possible paths  $z$

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

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

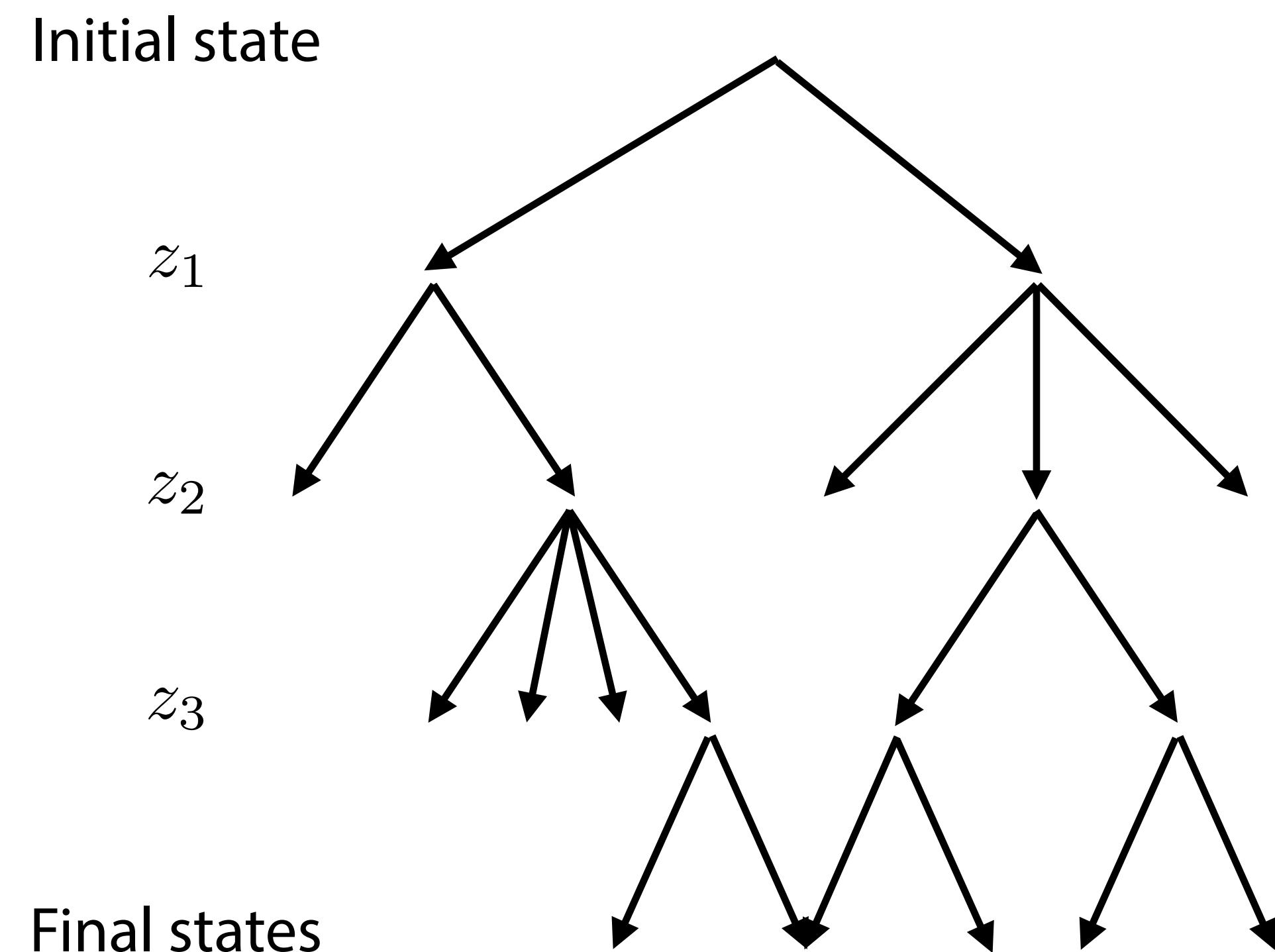
and the "joint likelihood ratio" conditional on a particular path

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

# Generalizing this idea

- Computer simulation typically evolve along a tree-like structure of successive random branchings
- The probabilities of each branching  $p_i(z_i|z_{i-1}, \theta)$  are often clearly defined in the code:

```
if random() > 0.1 + 2.5 * model_parameter:  
    do_one_thing()  
else:  
    do_another_thing()
```



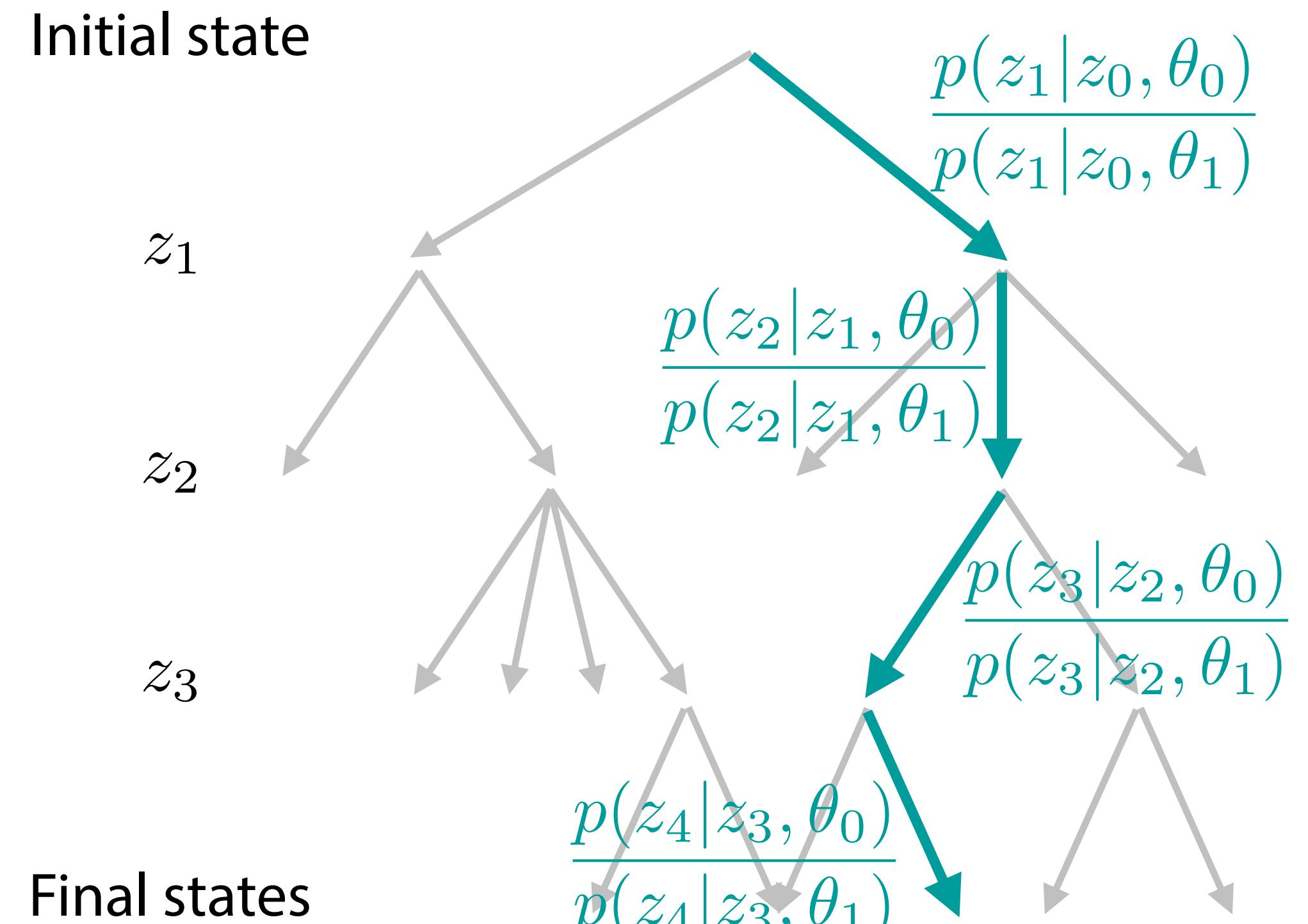
# Generalizing this idea

- Computer simulation typically evolve along a tree-like structure of successive random branchings
- The probabilities of each branching  $p_i(z_i|z_{i-1}, \theta)$  are often clearly defined in the code:

```
if random() > 0.1 + 2.5 * model_parameter:  
    do_one_thing()  
else:  
    do_another_thing()
```

- For each run of the simulator, we can calculate the probability **of the chosen path** for different values of the parameters, and the “**joint likelihood ratio**”:

$$r(x, z|\theta_0, \theta_1) = \frac{p(x, z|\theta_0)}{p(x, z|\theta_1)} = \prod_i \frac{p(z_i|z_{i-1}, \theta_0)}{p(z_i|z_{i-1}, \theta_1)}$$



# The value of gold

We can calculate the **joint likelihood ratio**

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

(“How much more likely is this simulation, including all latent variables, for  $\theta_0$  compared to  $\theta_1$ ? ”)



We want the **likelihood ratio function**

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

(“How much more likely is the observation  $x$  for  $\theta_0$  compared to  $\theta_1$ ? ”)

# The value of gold

We can calculate the joint likelihood ratio

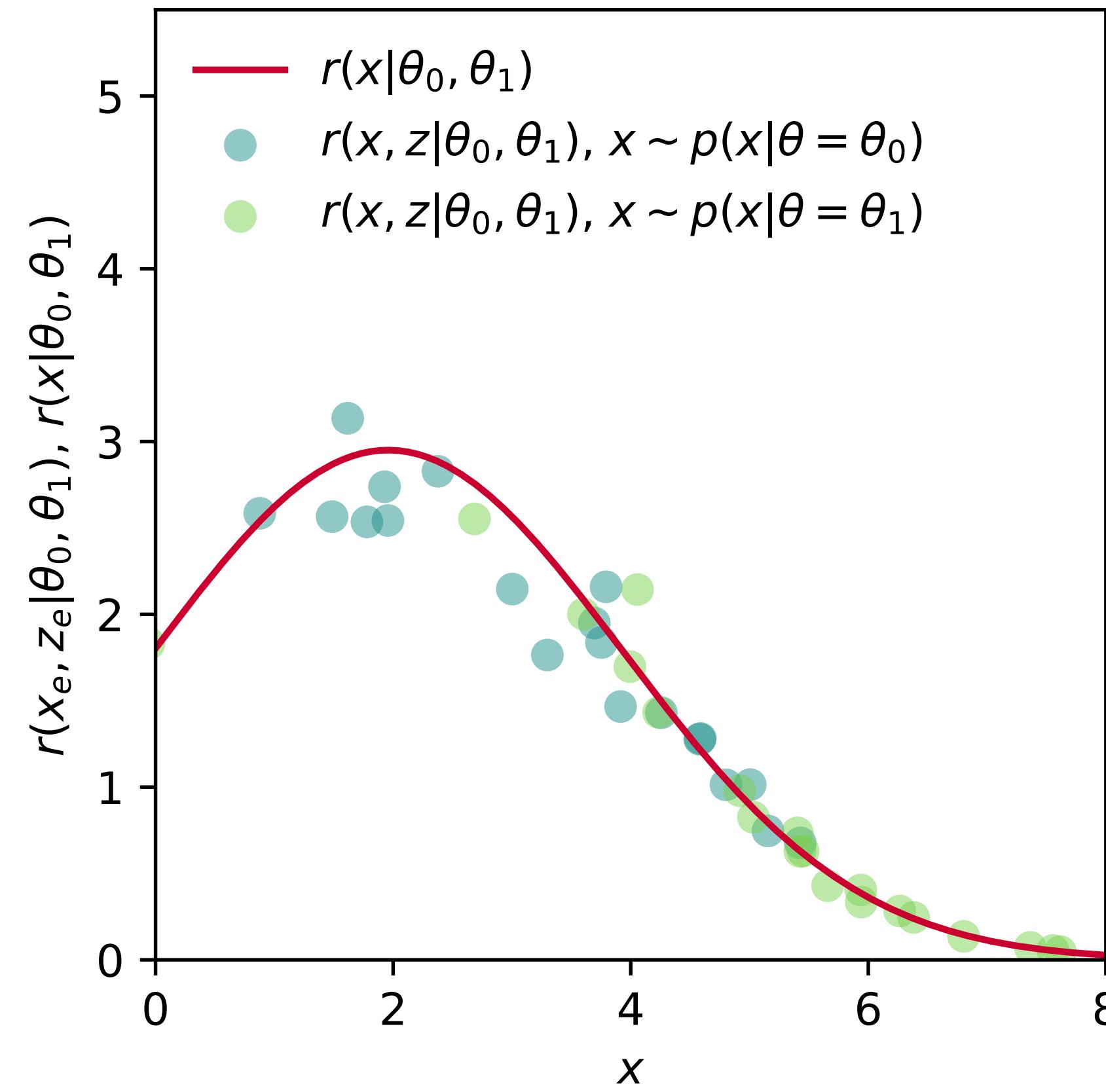
$$r(x, z | \theta_0, \theta_1) \equiv \frac{p(x, z_d, z_s, z_p | \theta_0)}{p(x, z_d, z_s, z_p | \theta_1)}$$



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

We want the likelihood ratio function

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



# The value of gold

We can calculate the 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 a functional like

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

It is minimized by

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

(And we can sample from  $p(x, z|\theta)$  by running the simulator.)

We want the likelihood ratio function

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

# Machine learning = applied calculus of variations

So to get a good estimator of the likelihood ratio, we need to minimize a functional numerically:

Variational family

Extremization

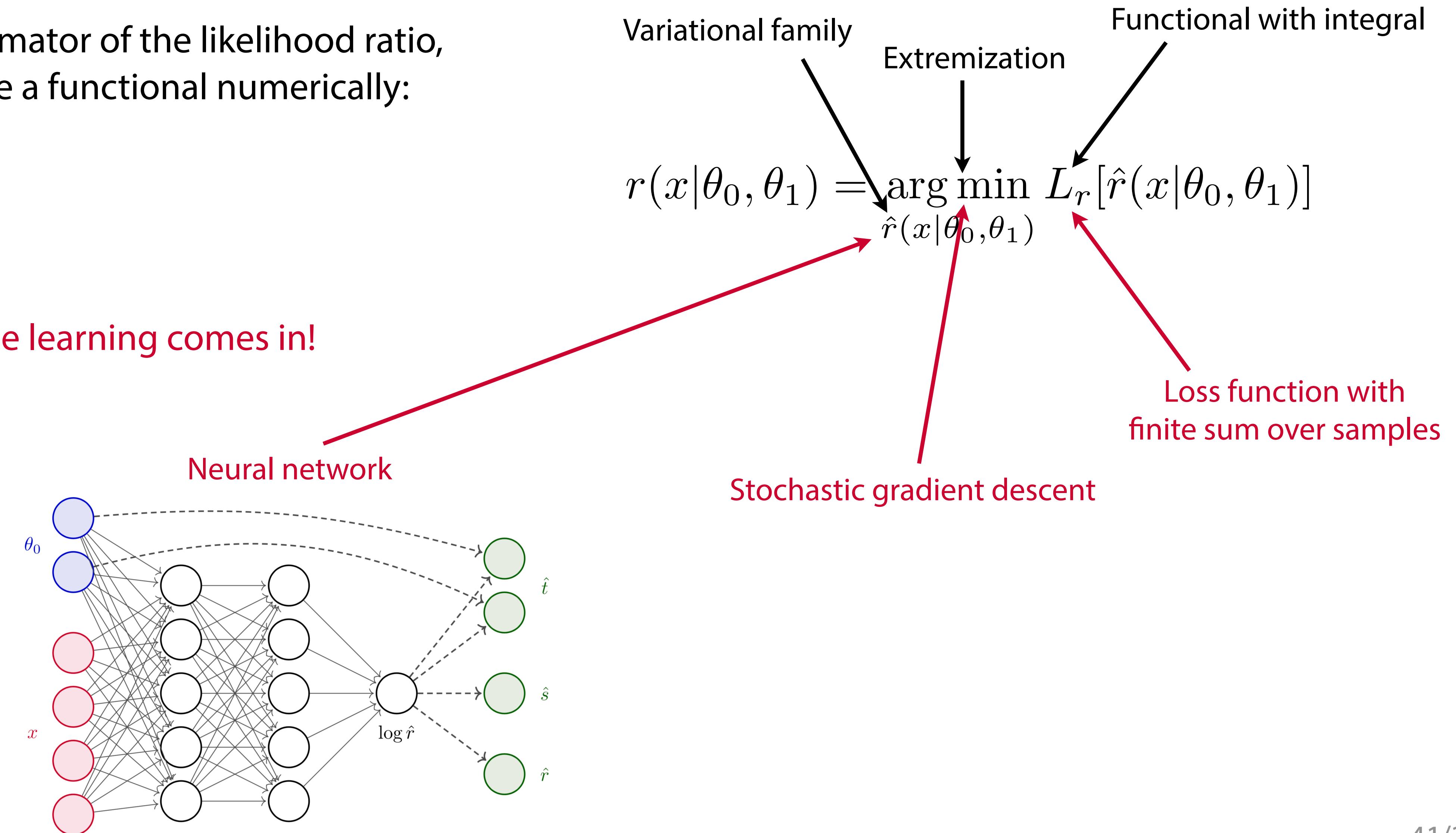
Functional with integral

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

# Machine learning = applied calculus of variations

So to get a good estimator of the likelihood ratio, we need to minimize a functional numerically:

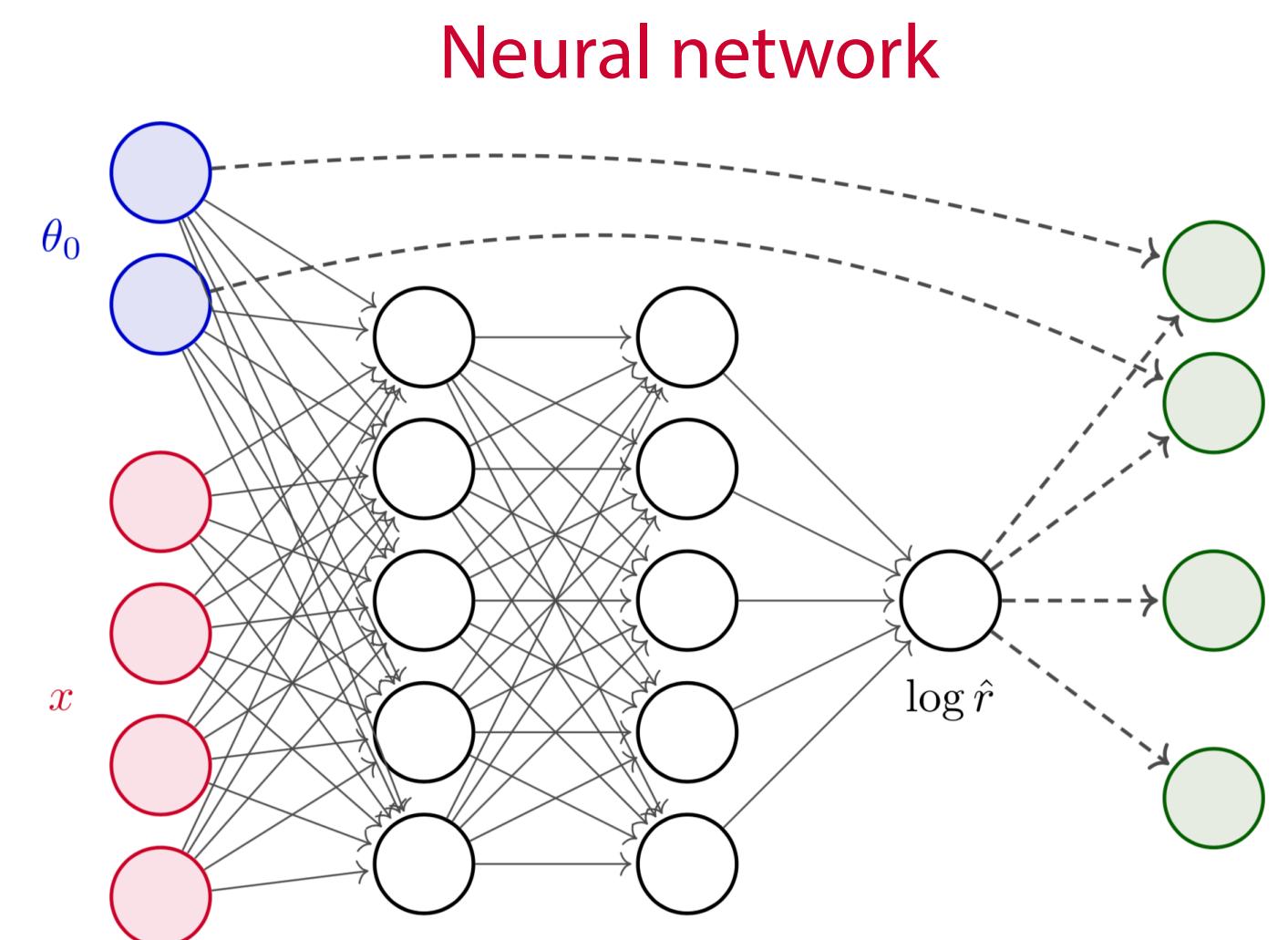
This is where machine learning comes in!



# Machine learning = applied calculus of variations

So to get a good estimator of the likelihood ratio, we need to minimize a functional numerically:

This is where machine learning comes in!



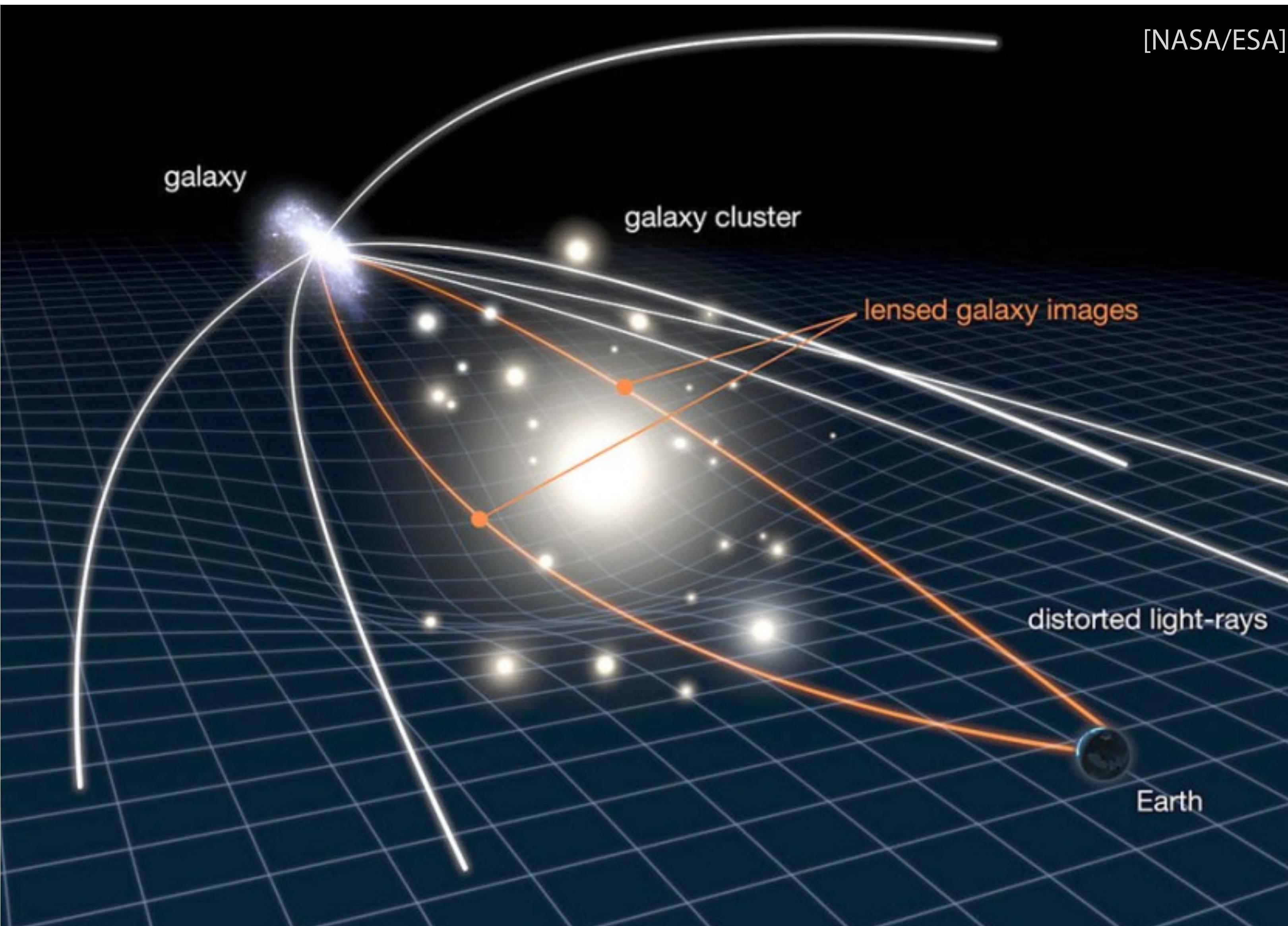
Neural network

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

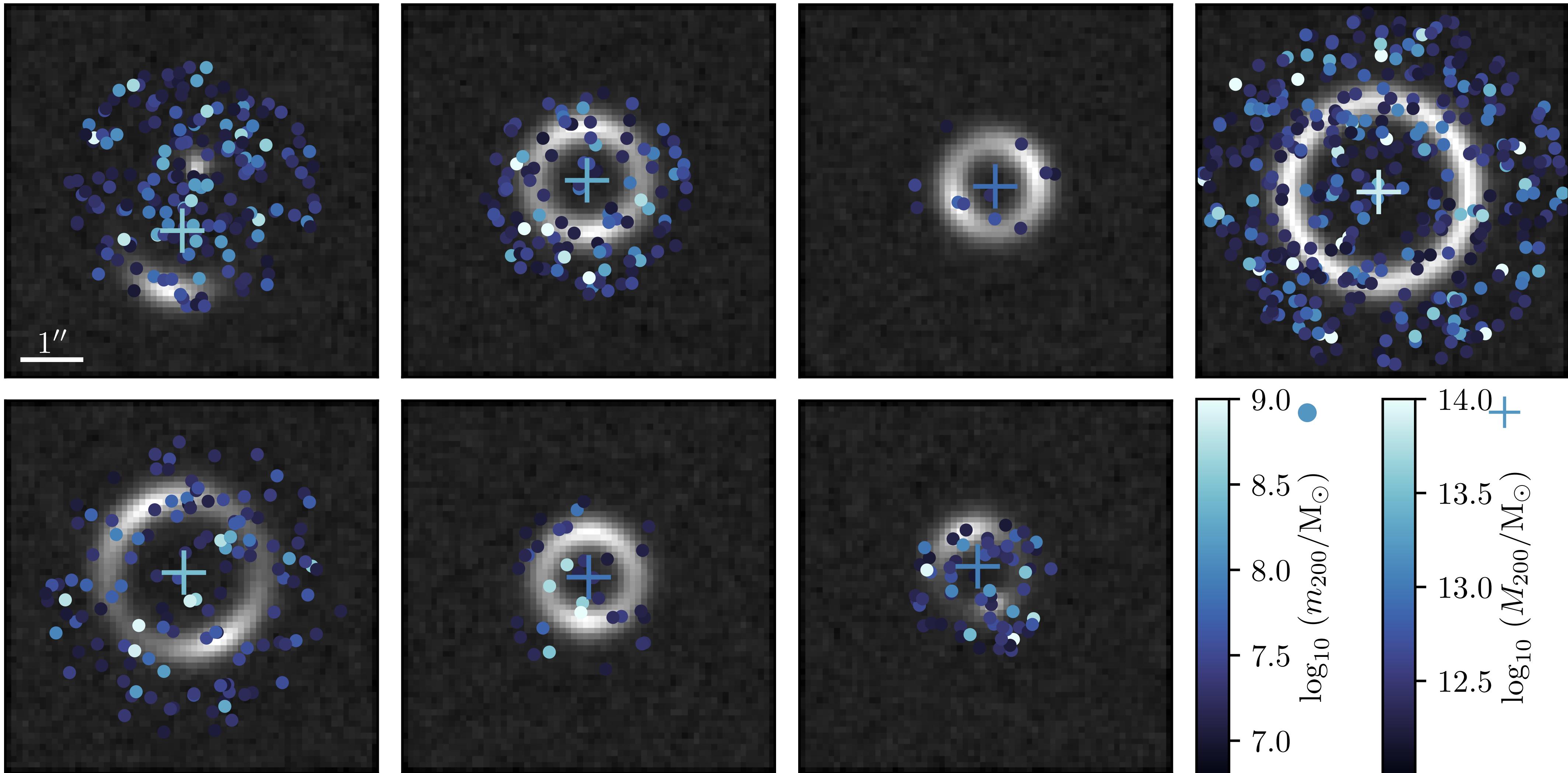
Variational family  
Extremization  
Functional with integral  
Loss function with finite sum over samples  
Stochastic gradient descent

A sufficiently expressive neural network efficiently trained in this way with enough data will learn the likelihood ratio function  $r(x|\theta_0, \theta_1)$ !

# Strong gravitational lensing



# Lensing simulations



# Mining for Dark Matter substructure

