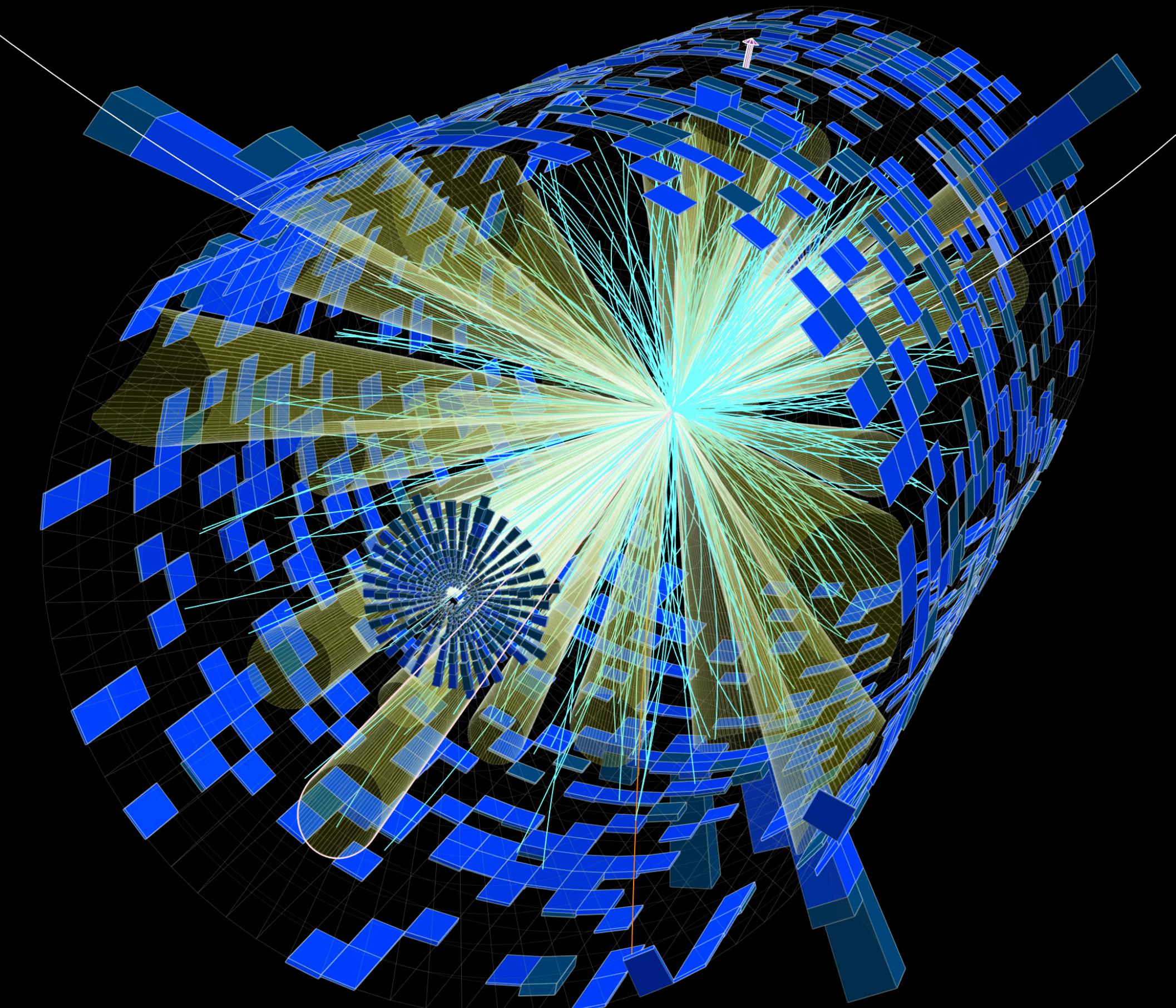




CONSTRAINING EFTS WITH MACHINE LEARNING

LECTURE 2

@KyleCranmer
New York University
Department of Physics
Center for Data Science
CILVR Lab



Many of these slides are borrowed from this talk by Johann Brehmer
(who in turn has adapted many of my slides... so it's really hybrid 😊)



The frontier of simulation-based inference

Johann Brehmer

New York University

Machine Learning at LHC workshop, Nagoya University

February 5, 2020

Particle physics processes do not have a tractable likelihood function.

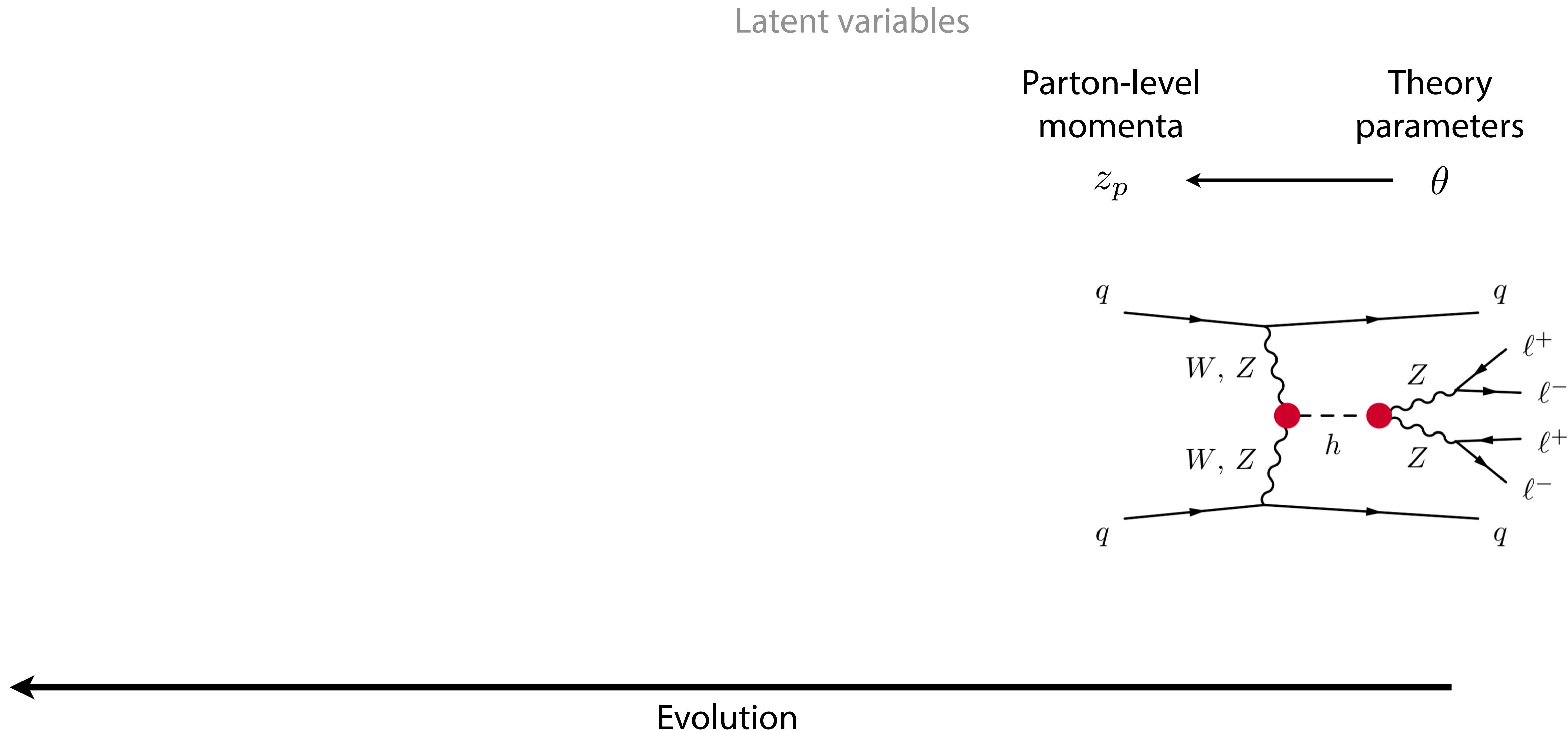
Modelling particle physics processes

Theory
parameters
 θ

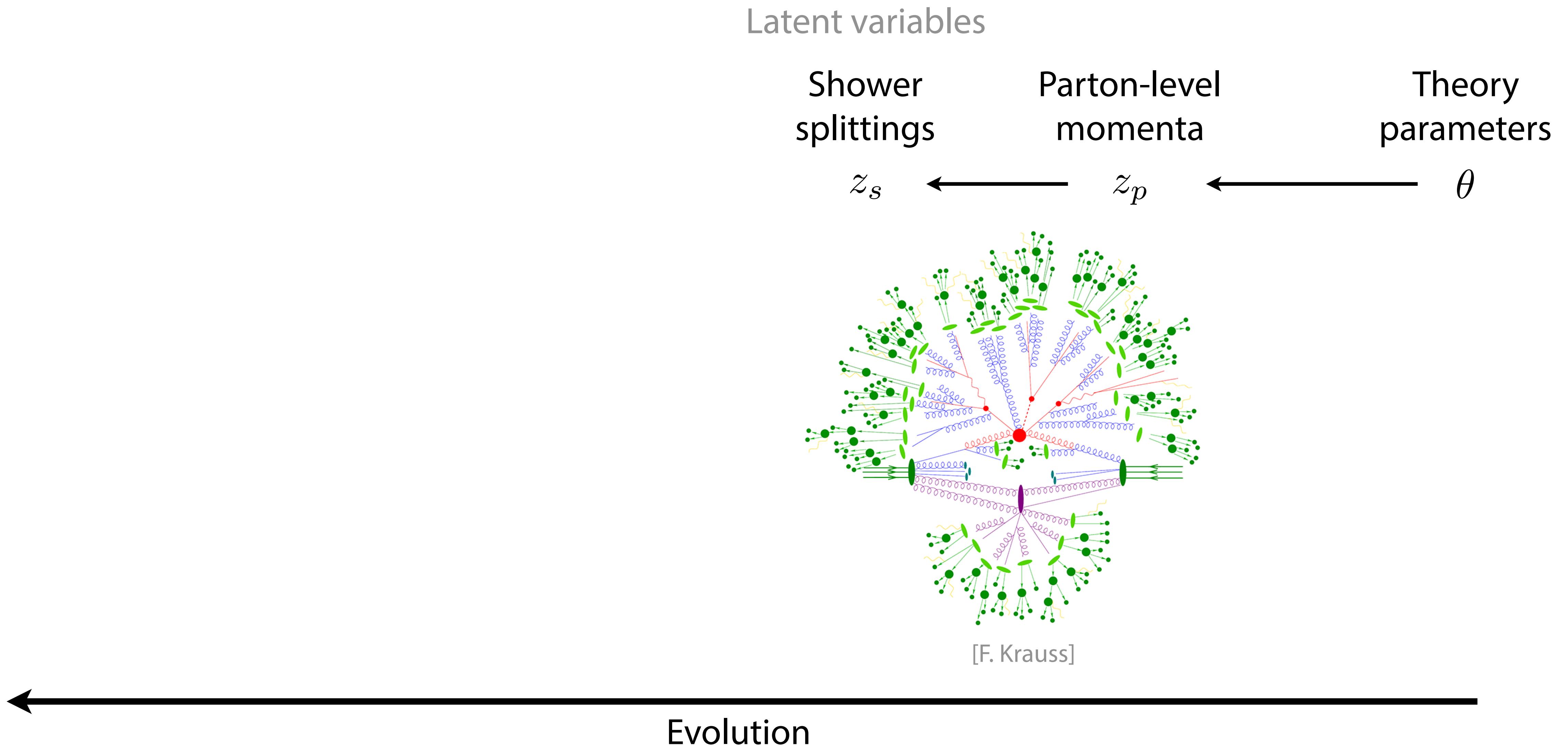


Evolution

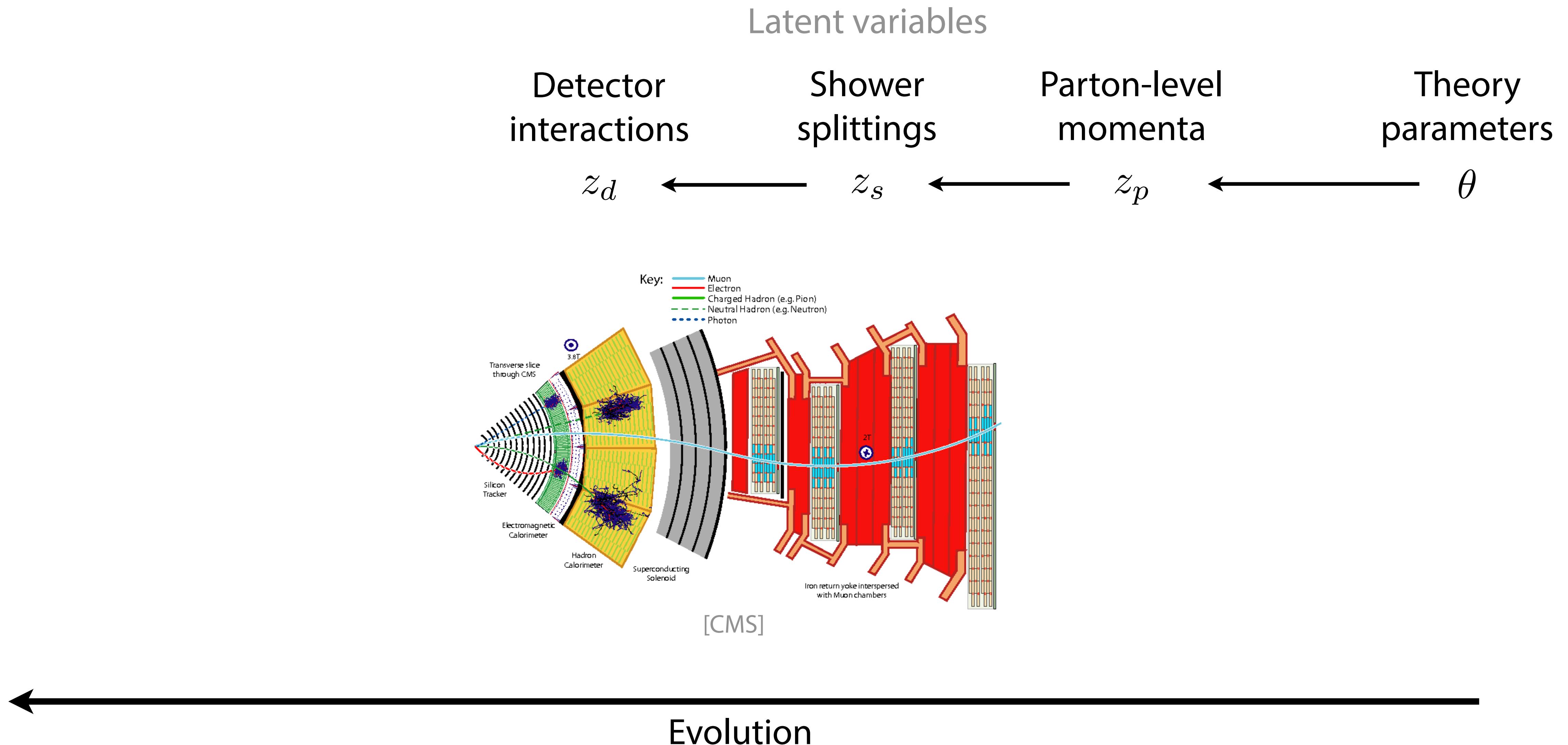
Modelling particle physics processes



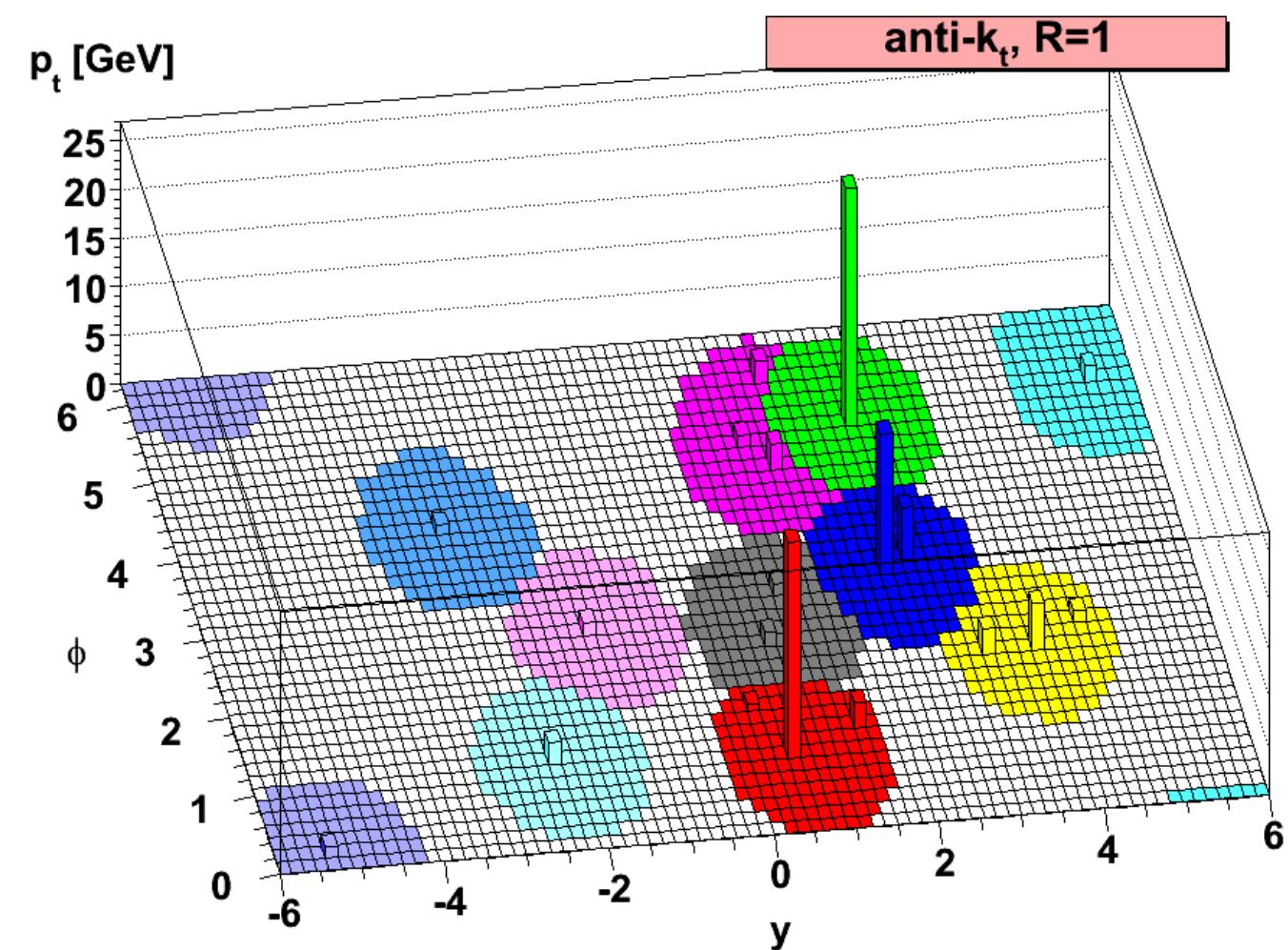
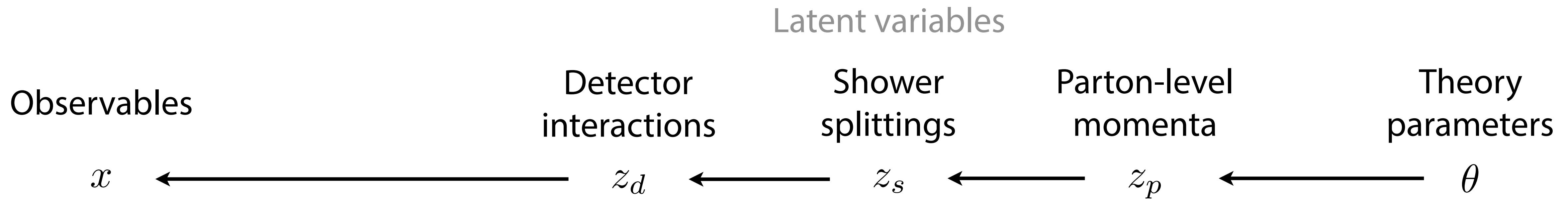
Modelling particle physics processes



Modelling particle physics processes



Modelling particle physics processes

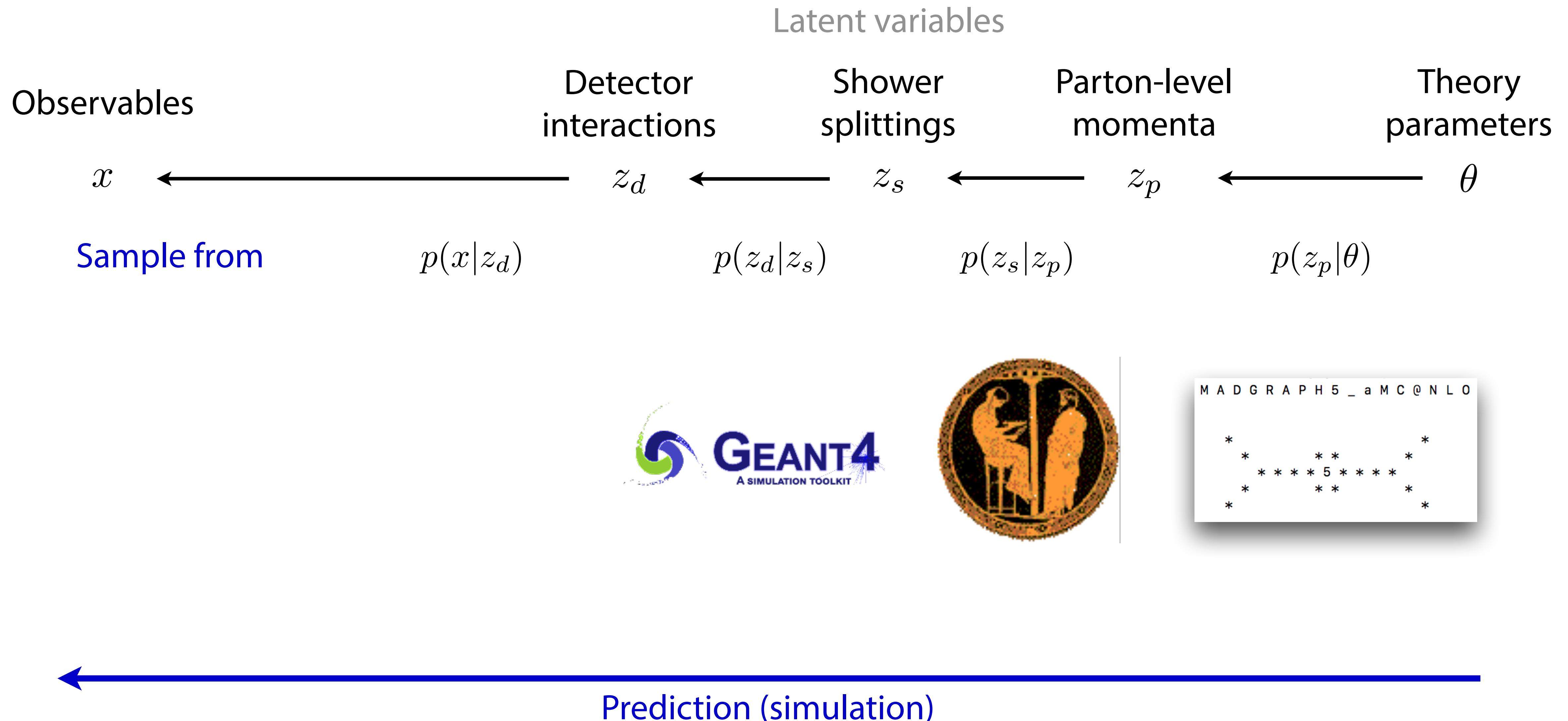


[M. Cacciari, G. Salam, G. Soyez 0802.1189]

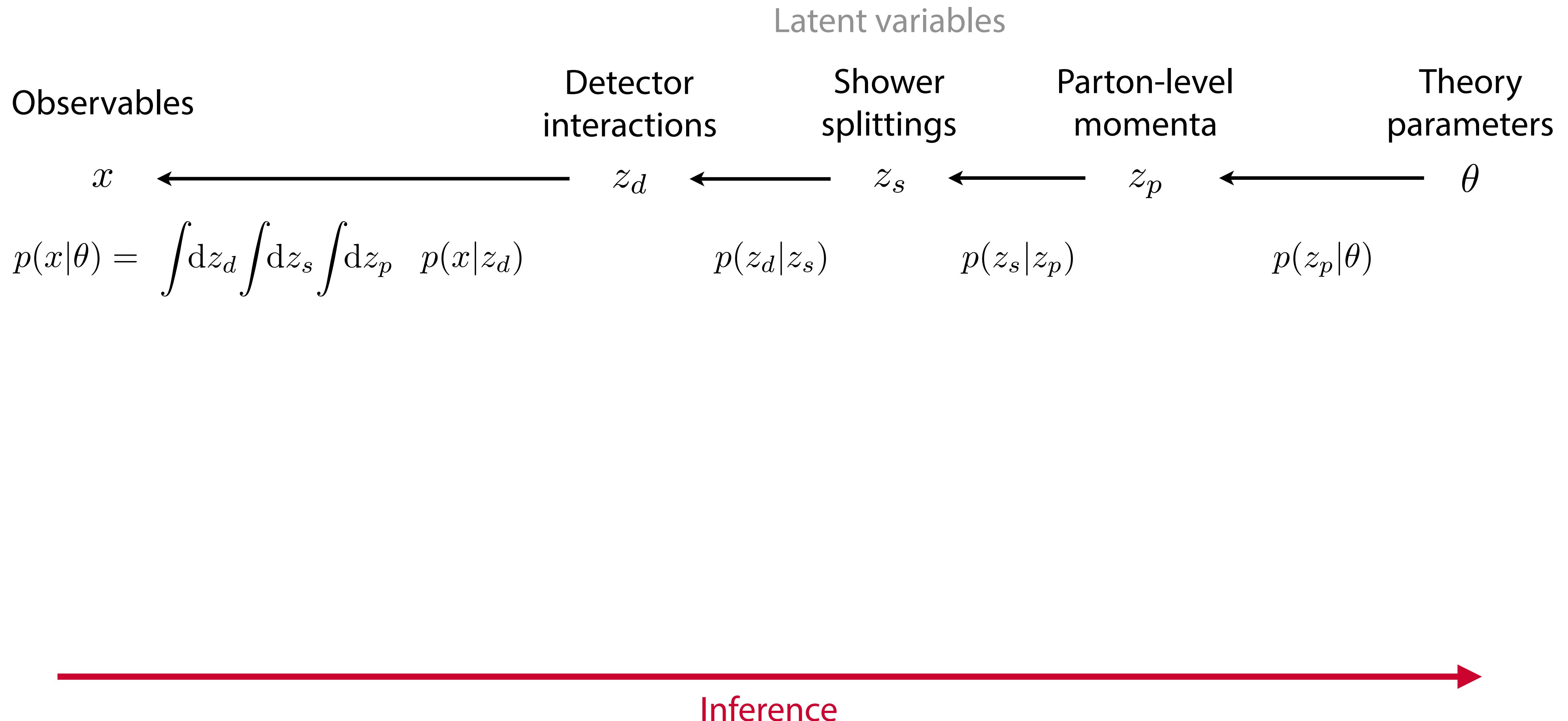
Evolution

```
graph LR; Evolution[Evolution] --> Left[ ]
```

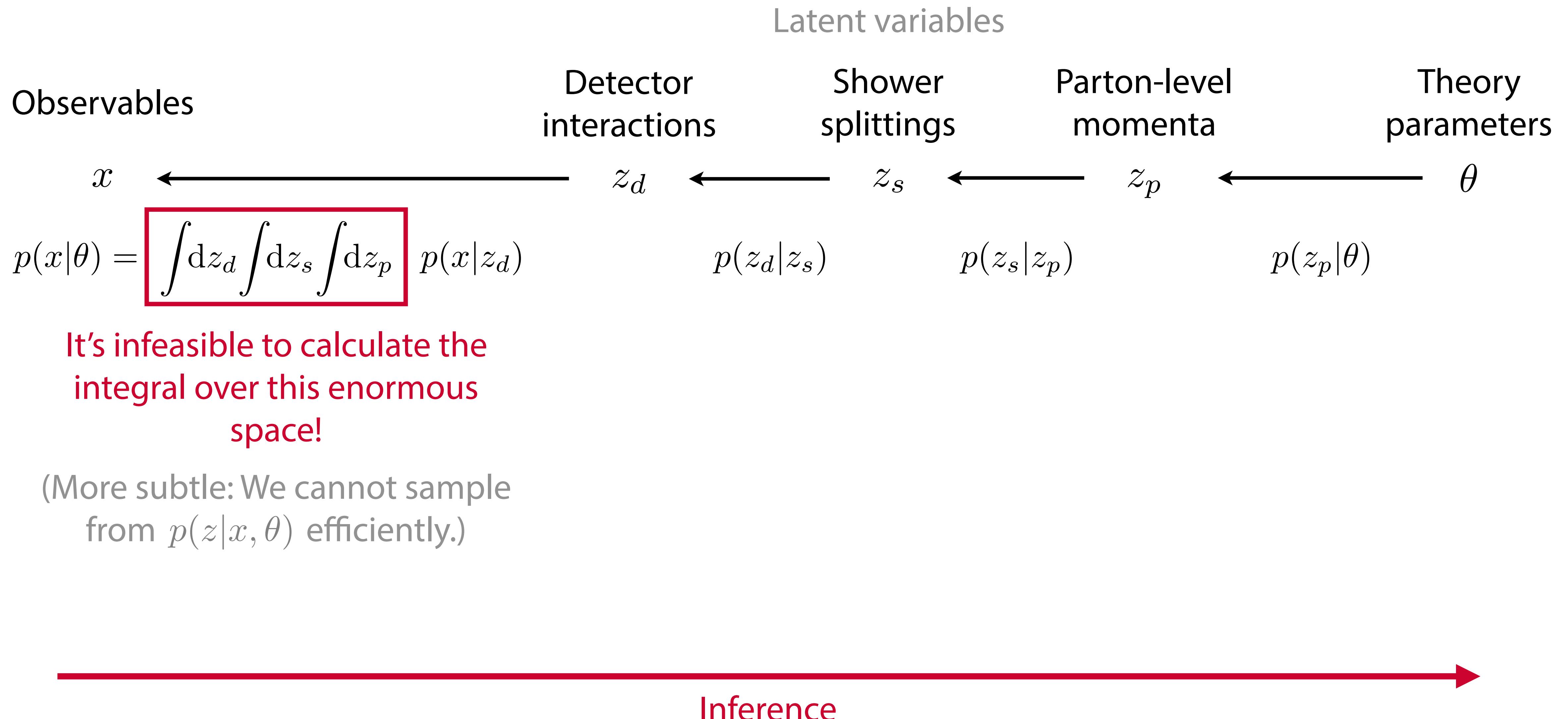
Modelling particle physics processes



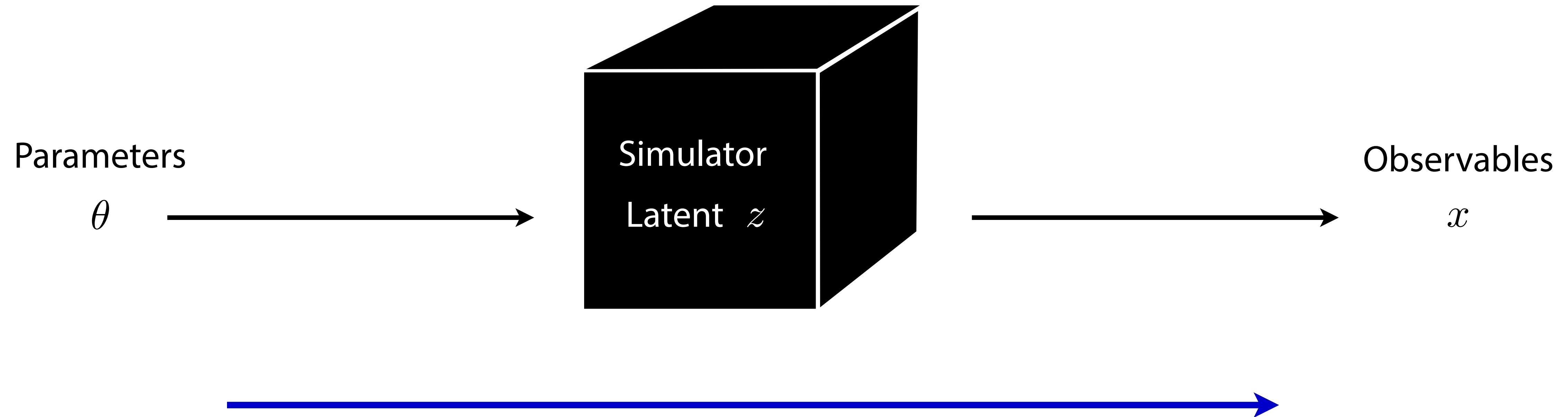
Modelling particle physics processes



Modelling particle physics processes

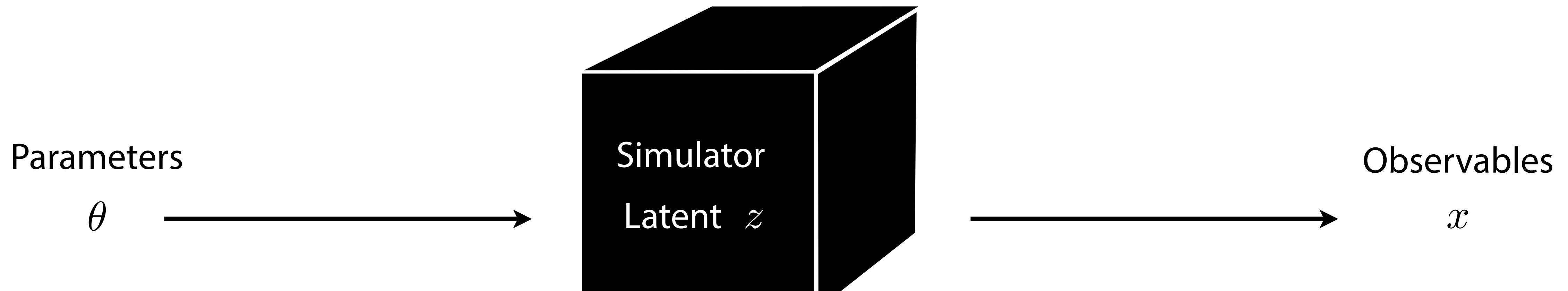


Simulation-based (“likelihood-free”) inference



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

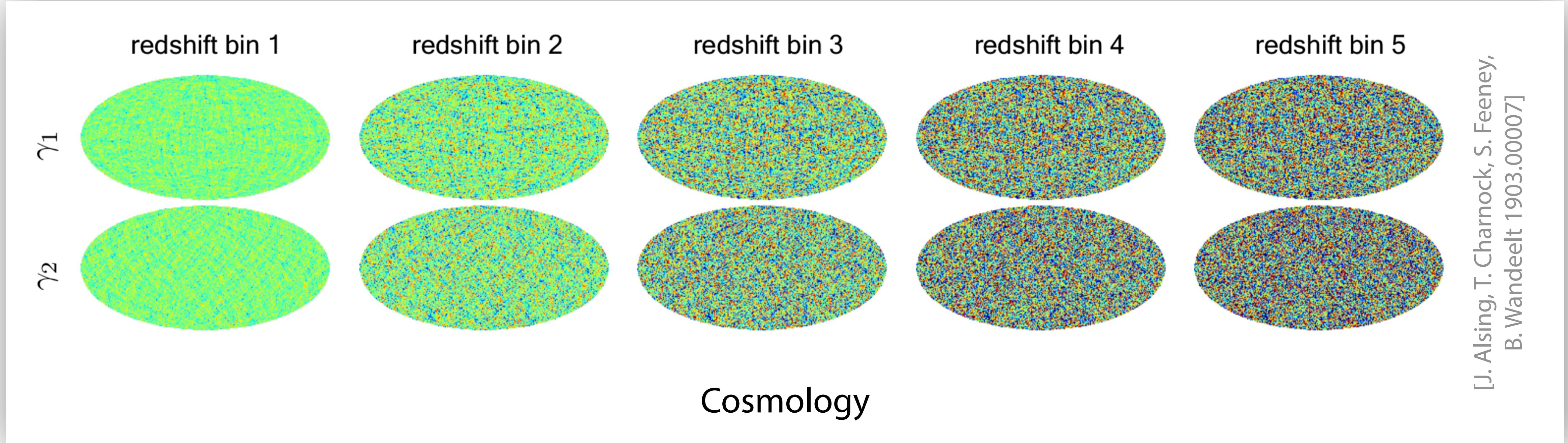
Simulation-based (“likelihood-free”) inference



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

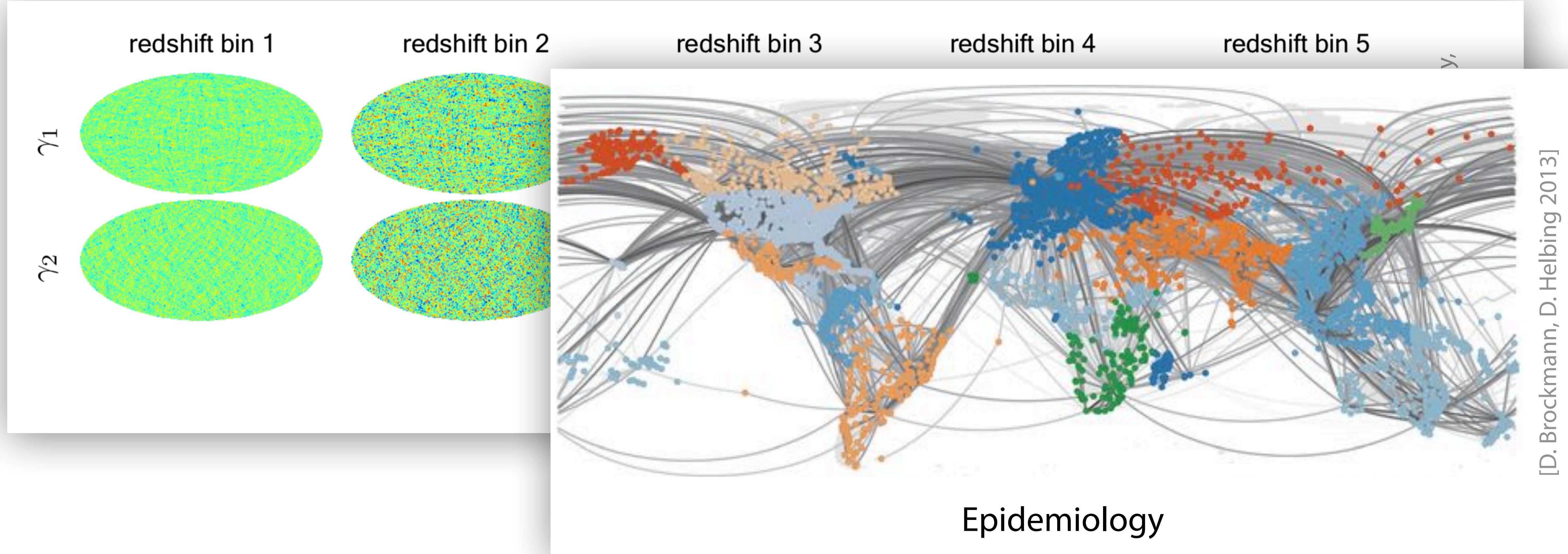
- Inference:**
- Likelihood $p(x|\theta) = \int dz p(x, z|\theta)$ is intractable
 - Inference is challenging

It's not just particle physics!

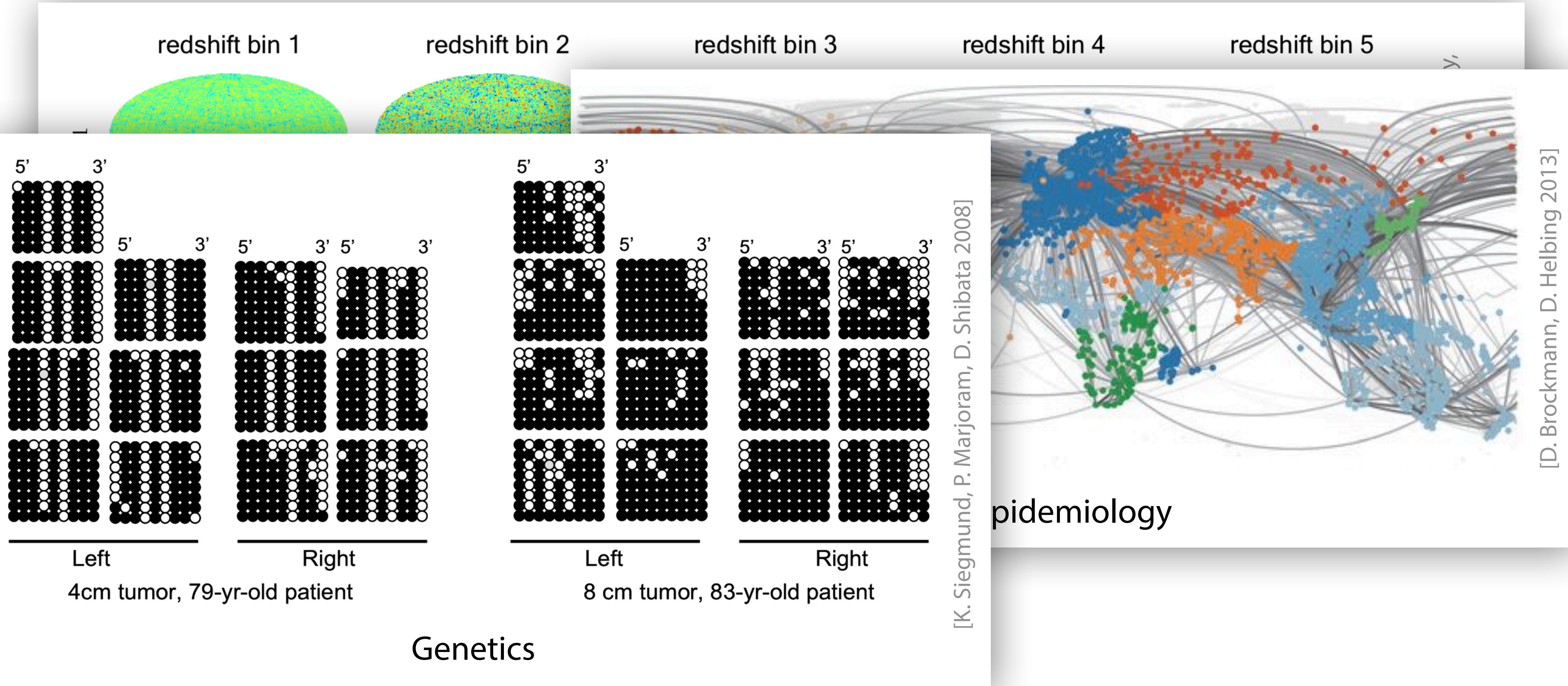


[J. Alsing, T. Charnock, S. Feeney,
B. Wandelt 1903.00007]

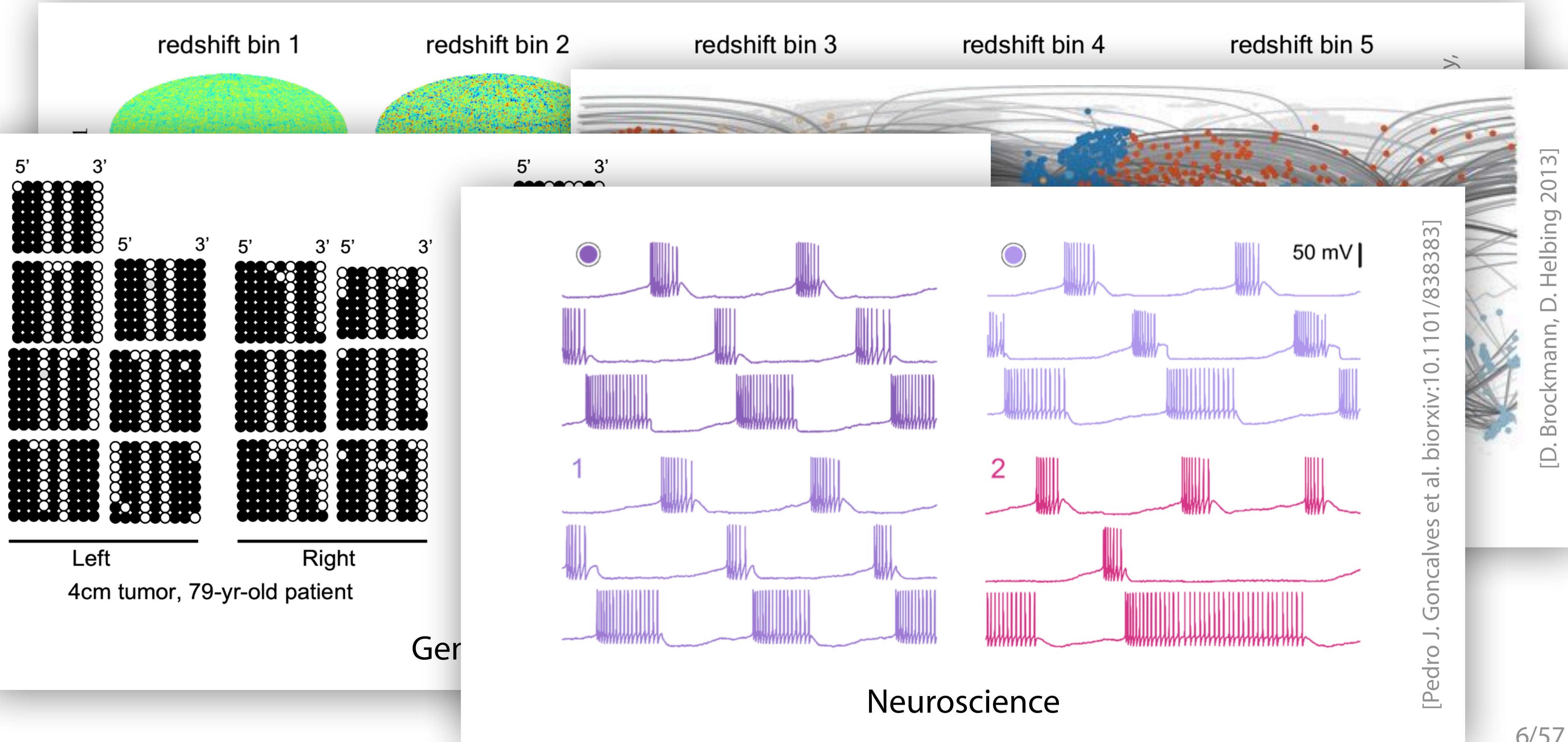
It's not just particle physics!



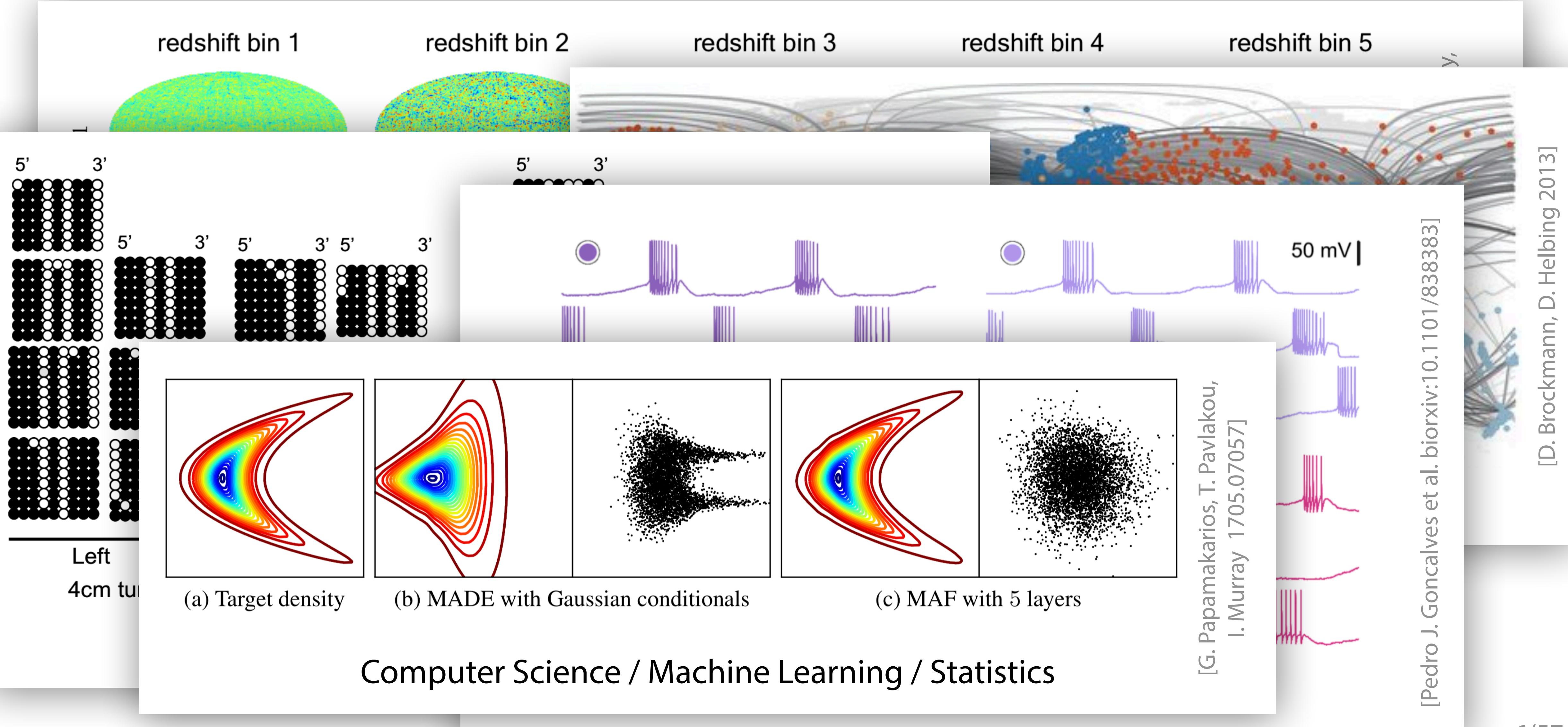
It's not just particle physics!



It's not just particle physics!



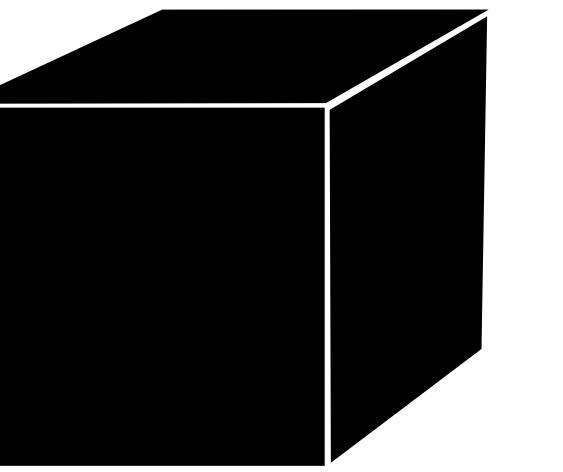
It's not just particle physics!



Formalizing the problem

You are given

- a simulator that lets you generate N samples $x_i \sim p(x_i|\theta_i)$ for parameter points θ_i of your choice

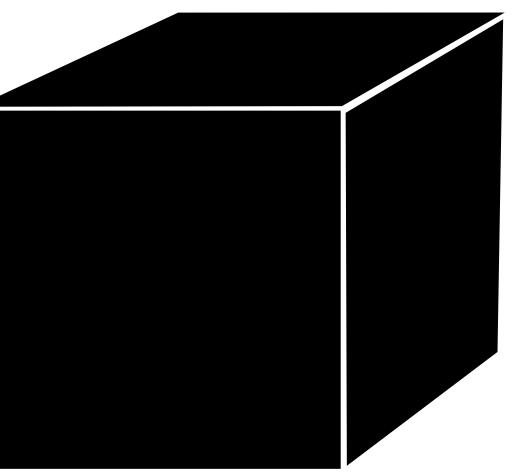


- observed data $x_{\text{obs}} \sim p(x_{\text{obs}}|\theta_{\text{true}})$
- prior belief $p(\theta)$

Formalizing the problem

You are given

- a simulator that lets you generate N samples $x_i \sim p(x_i|\theta_i)$ for parameter points θ_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

LHC footnotes

- Full LHC likelihood: $p_{\text{full}}(\{x\}|\theta) = \text{Pois}(n|L\sigma(\theta)) \prod_{\text{events } x} p(x|\theta)$

LHC footnotes

- Full LHC likelihood:

$$p_{\text{full}}(\{x\}|\theta) = \text{Pois}(n|L\sigma(\theta)) \prod_{\text{events } x} p(x|\theta)$$

Total rate term:

- How likely is it to observe n events after cuts?
- For simplicity, we ignore this part in this talk

LHC footnotes

- Full LHC likelihood:

$$p_{\text{full}}(\{x\}|\theta) = \text{Pois}(n|L\sigma(\theta)) \prod_{\text{events } x} p(x|\theta)$$

Total rate term:

- How likely is it to observe n events after cuts?
- For simplicity, we ignore this part in this talk

Kinematic term for each event:

- How likely is it that an event looks like it does?
- \sim normalized differential xsec
- Focus of this talk

LHC footnotes

- Full LHC likelihood:

$$p_{\text{full}}(\{x\}|\theta) = \text{Pois}(n|L\sigma(\theta)) \prod_{\text{events } x} p(x|\theta)$$

Total rate term:

- How likely is it to observe n events after cuts?
- For simplicity, we ignore this part in this talk

Kinematic term for each event:

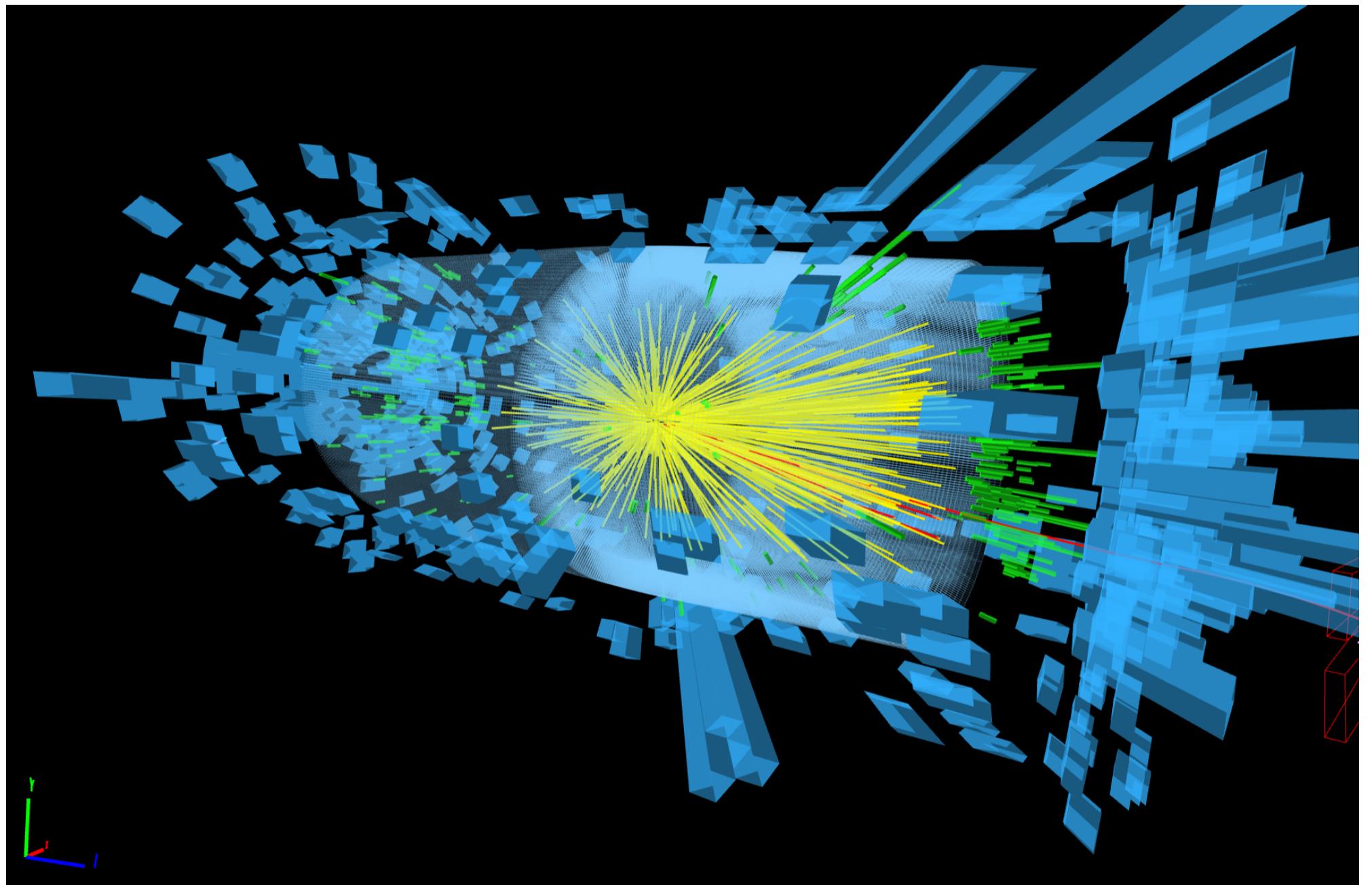
- How likely is it that an event looks like it does?
- \sim normalized differential xsec
- Focus of this talk

- Event selection:

- Choice of cuts shifts information between rate and kinematic part
- “Good” cuts depend on inference strategy
- This talk: assume fixed event selection

Traditionally, inference is made possible
with summary statistics.

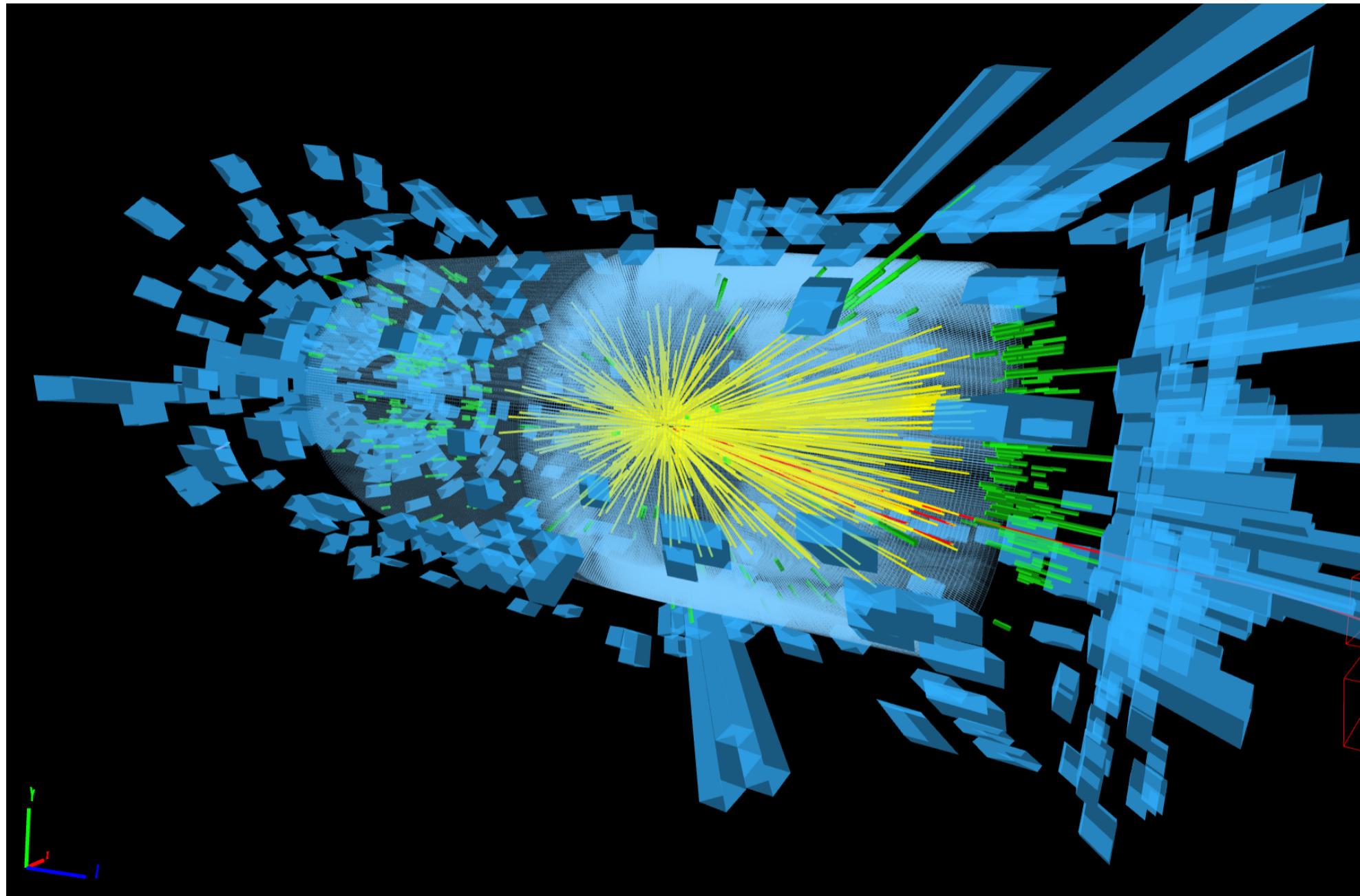
Solve it with summary statistics



High-dimensional event data x

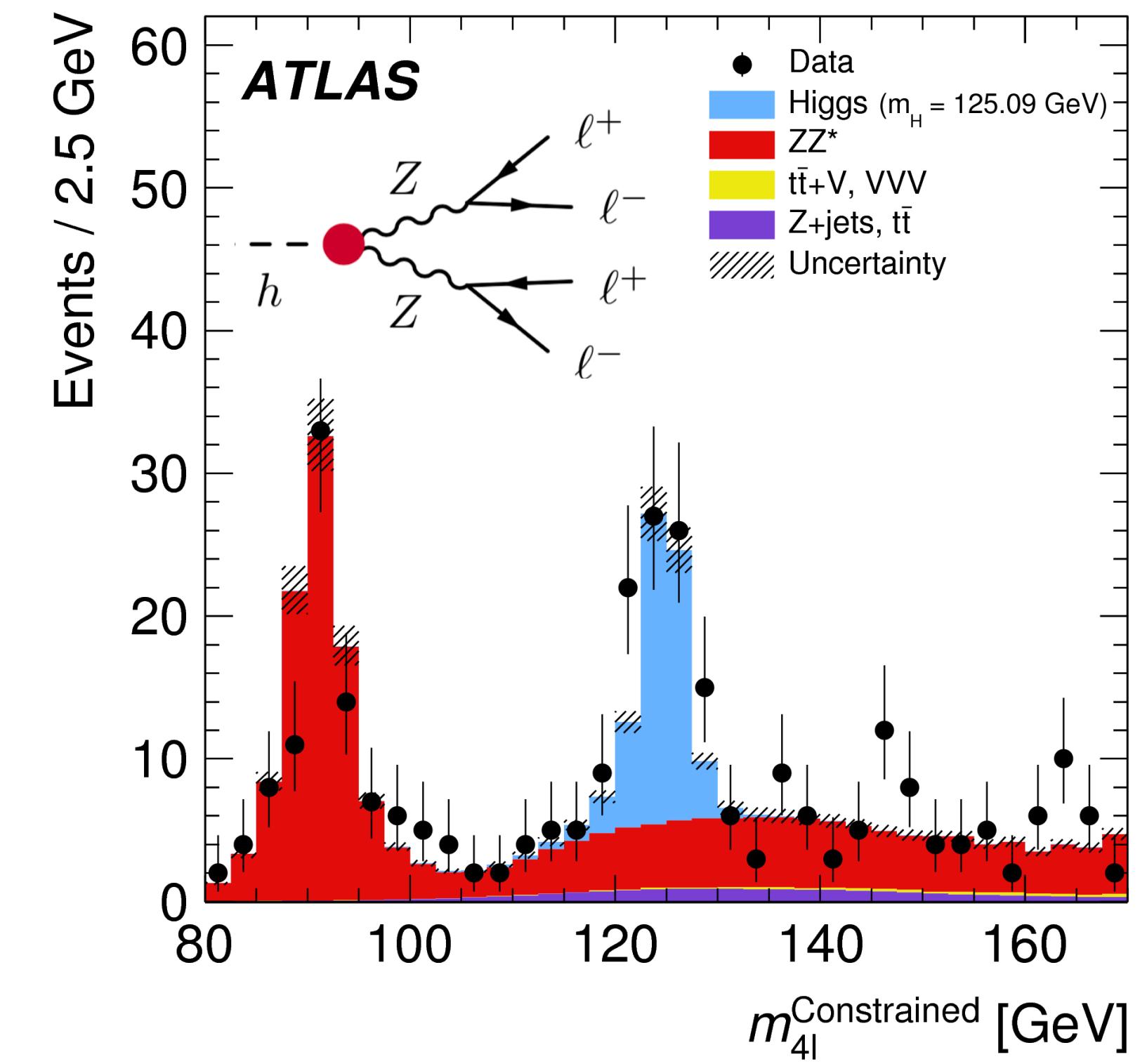
$p(x|\theta)$ cannot be calculated

Solve it with summary statistics



High-dimensional event data x

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

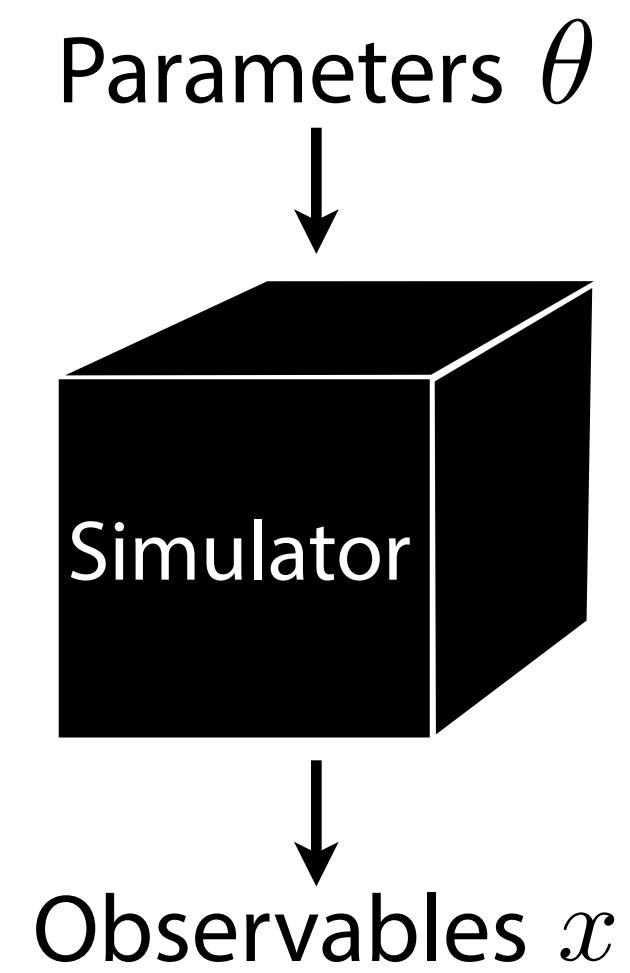


One or two summary statistics x'

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

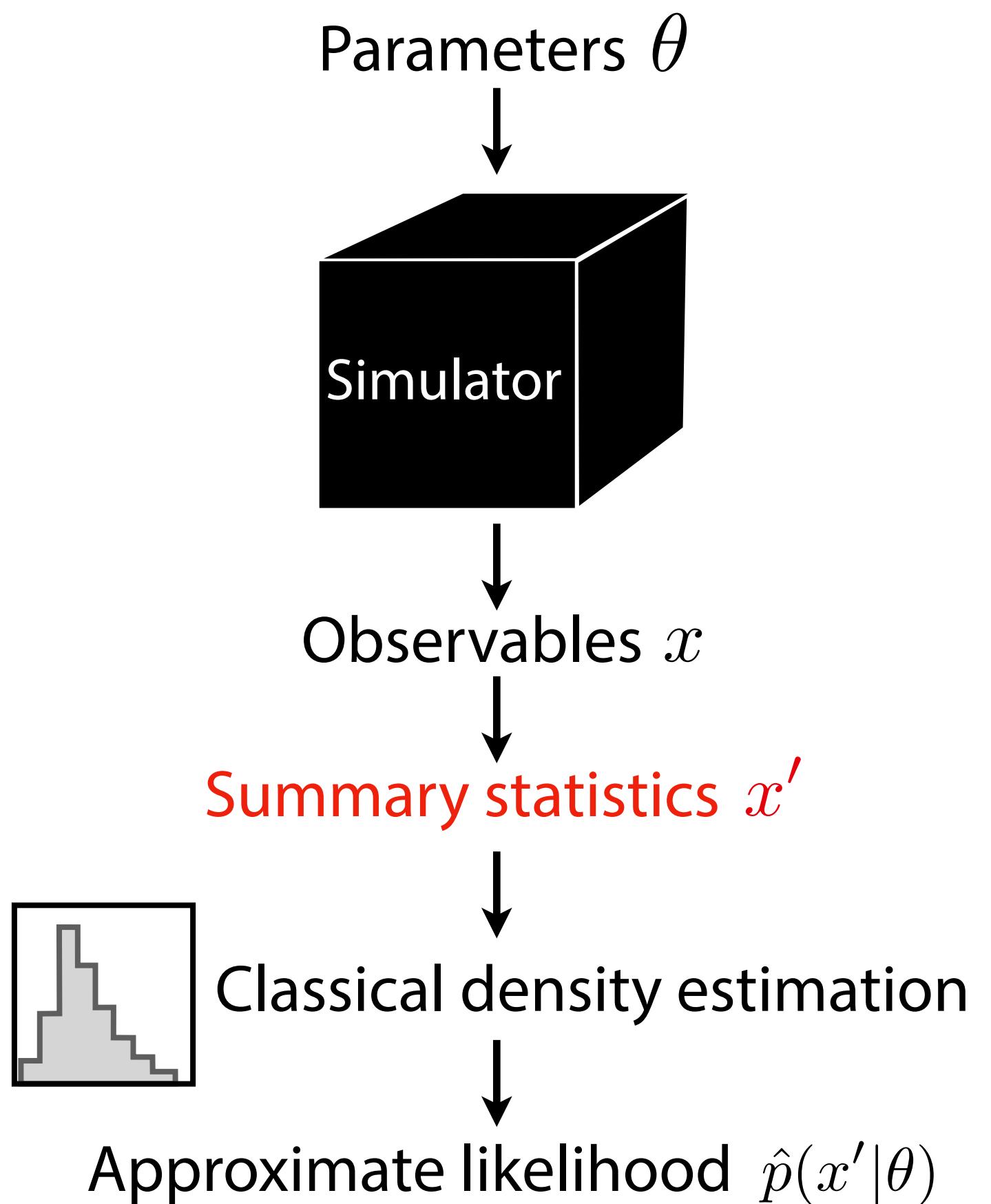
Inference by estimating the likelihood

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



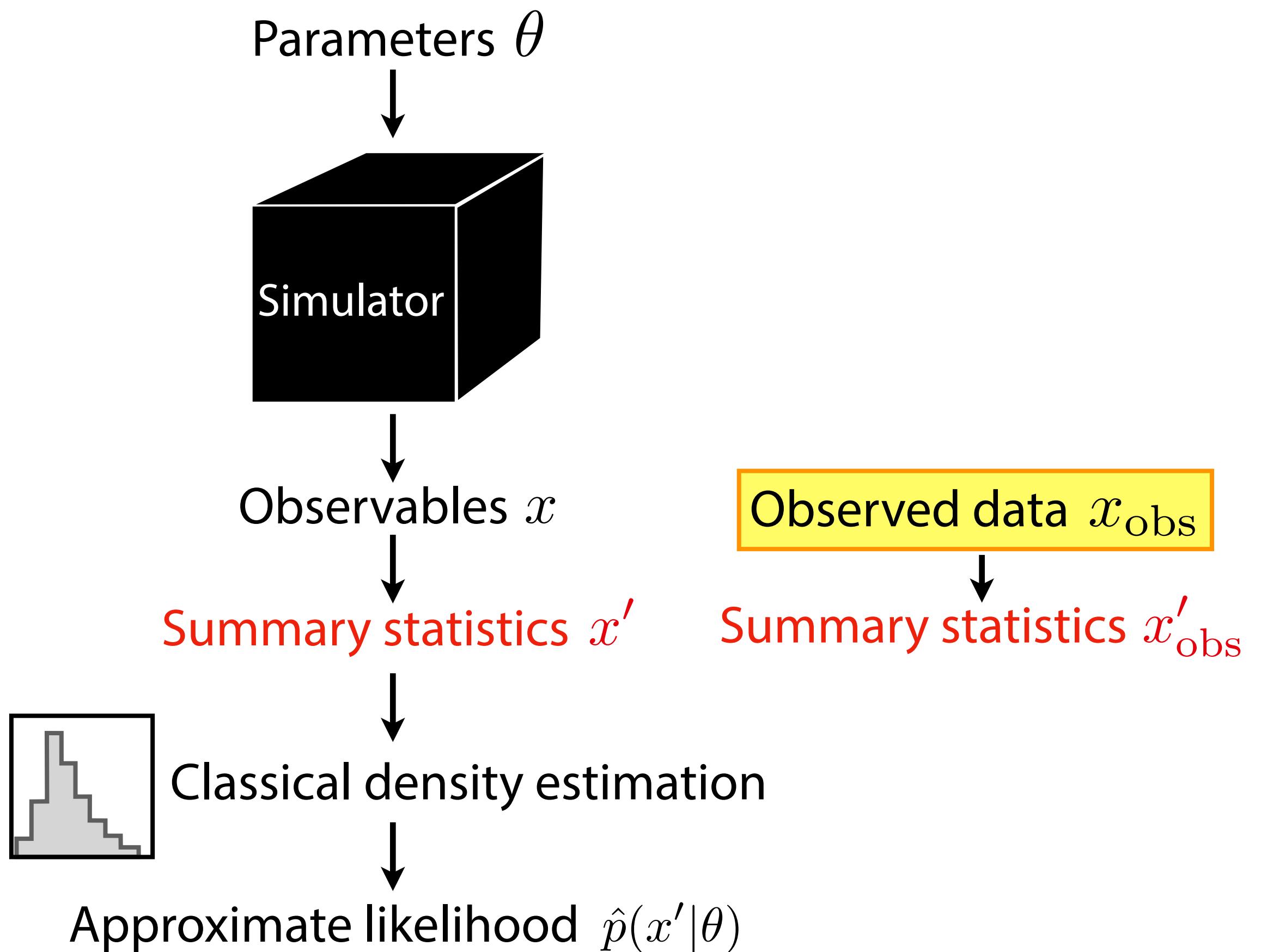
Inference by estimating the likelihood

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



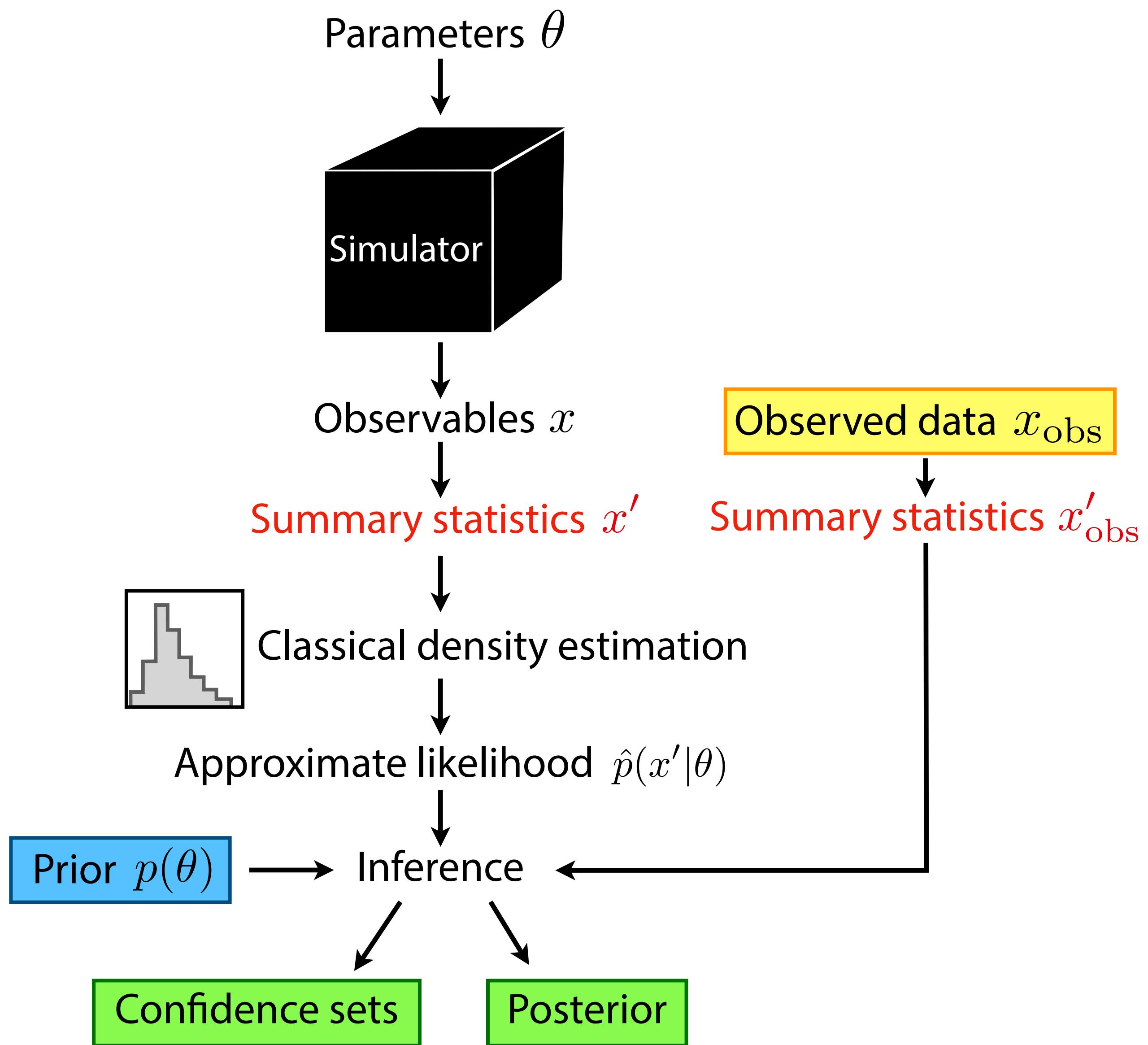
Inference by estimating the likelihood

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



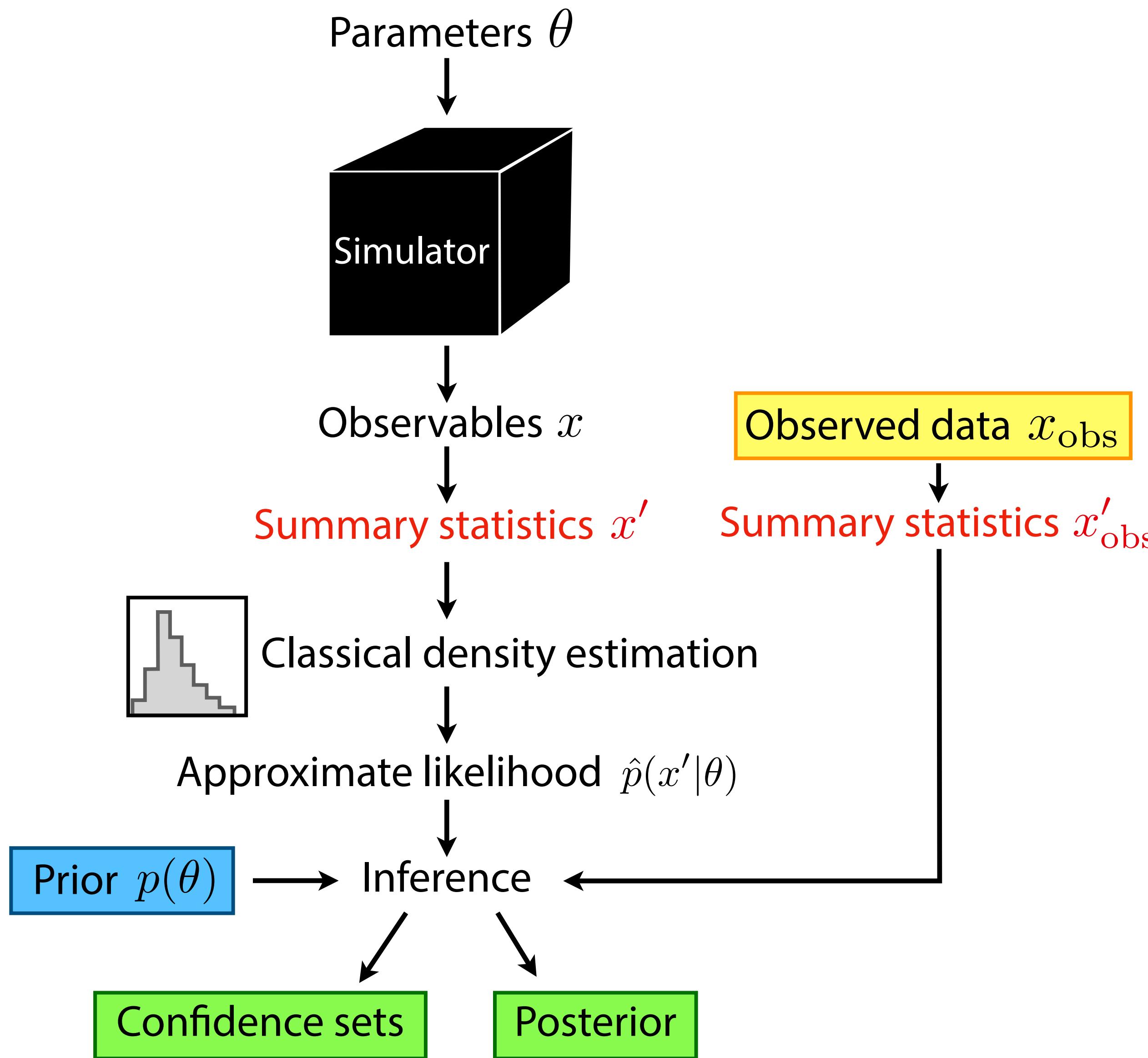
Inference by estimating the likelihood

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



Inference by estimating the likelihood

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



- Compression to summary statistics loses information & reduces quality of inference
- Curse of dimensionality: does not scale to more than a few summary statistics
- Related alternative: Approximate Bayesian Computation (ABC) [D. Rubin 1984]

Summary statistics for LHC measurements?

- In many LHC problems there is no single good summary statistic: 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 all trustworthy high-level features (reconstructed four-momenta...), or some form of low-level features, including correlations

("fully differential cross section")

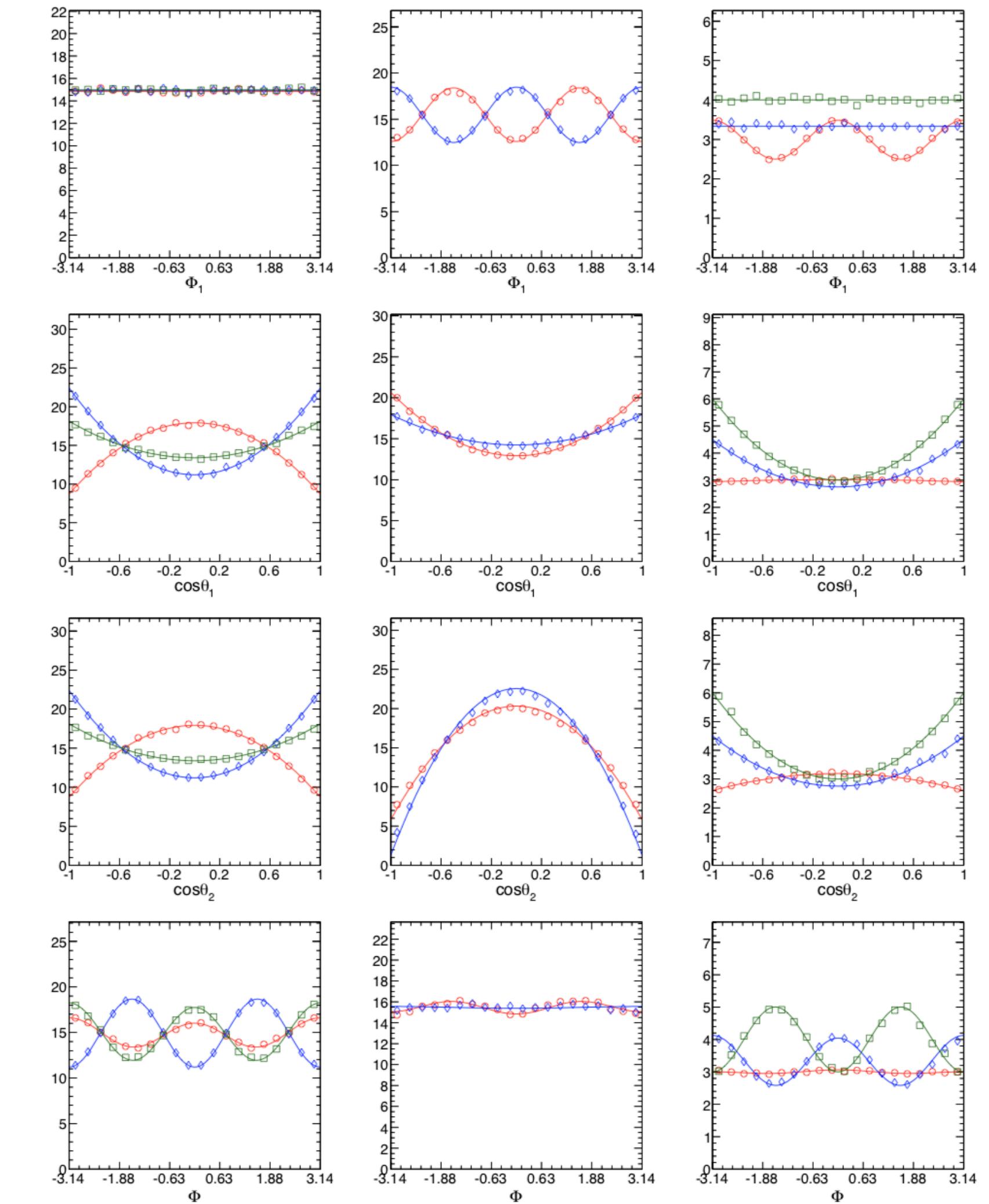
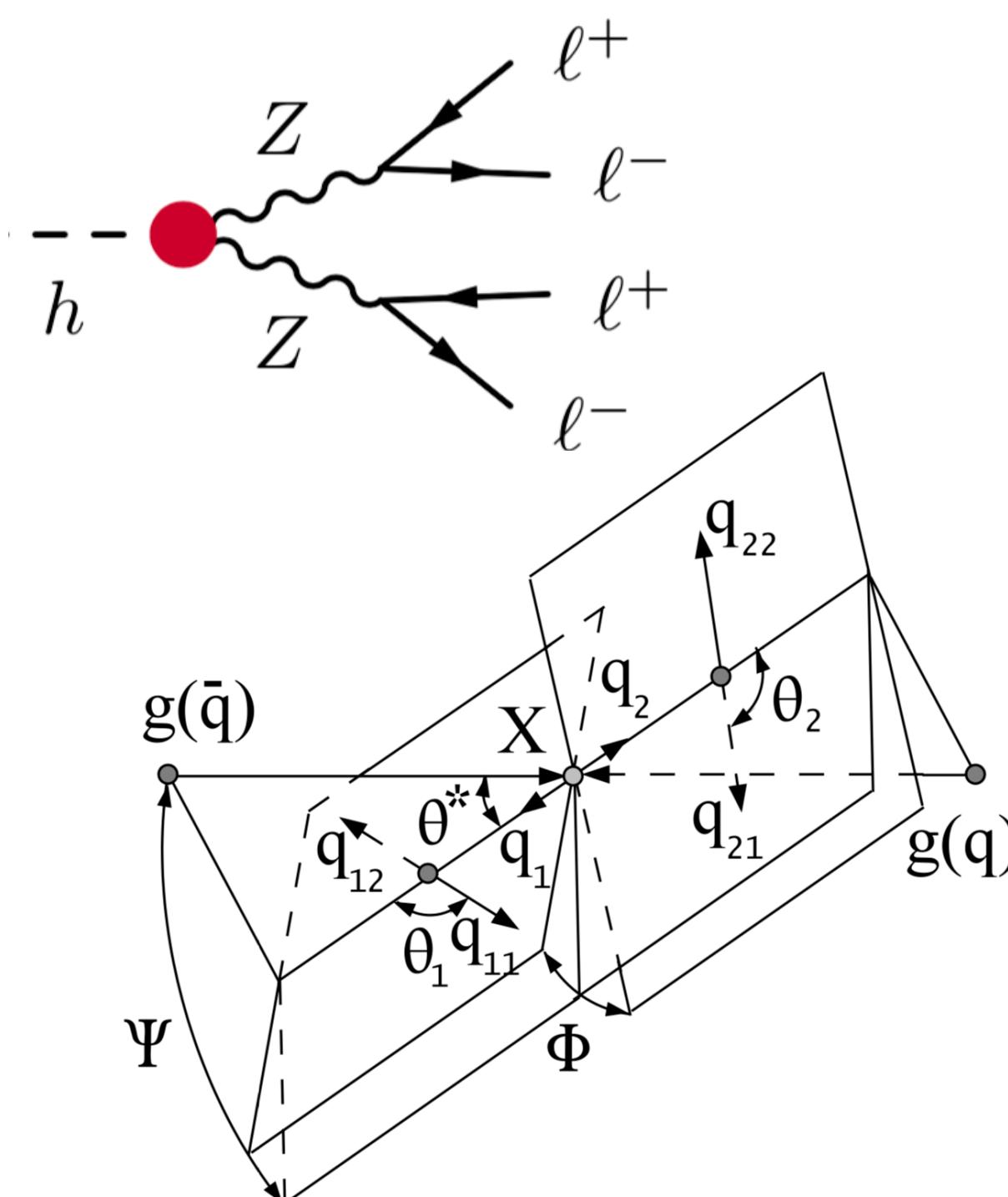
Summary statistics for LHC measurements?

- In many LHC problems there is no single good summary statistic: 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 all trustworthy high-level features (reconstructed four-momenta...), or some form of low-level features, including correlations

("fully differential cross section")

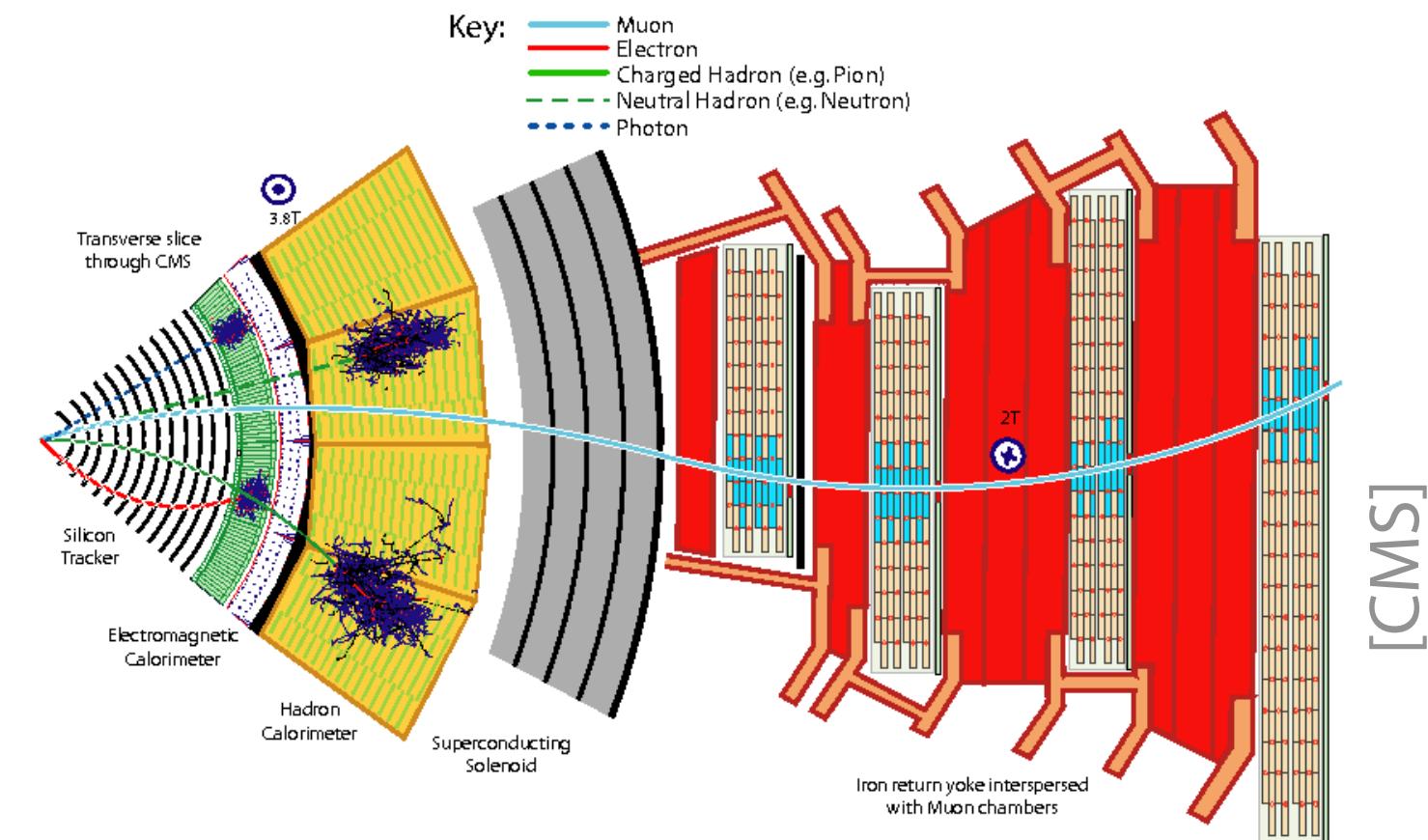


[Bolognesi et al. 1208.4018]

Solve it by approximating the integral

- Problem: high-dimensional integral over shower / detector trajectories

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



Solve it by approximating the integral

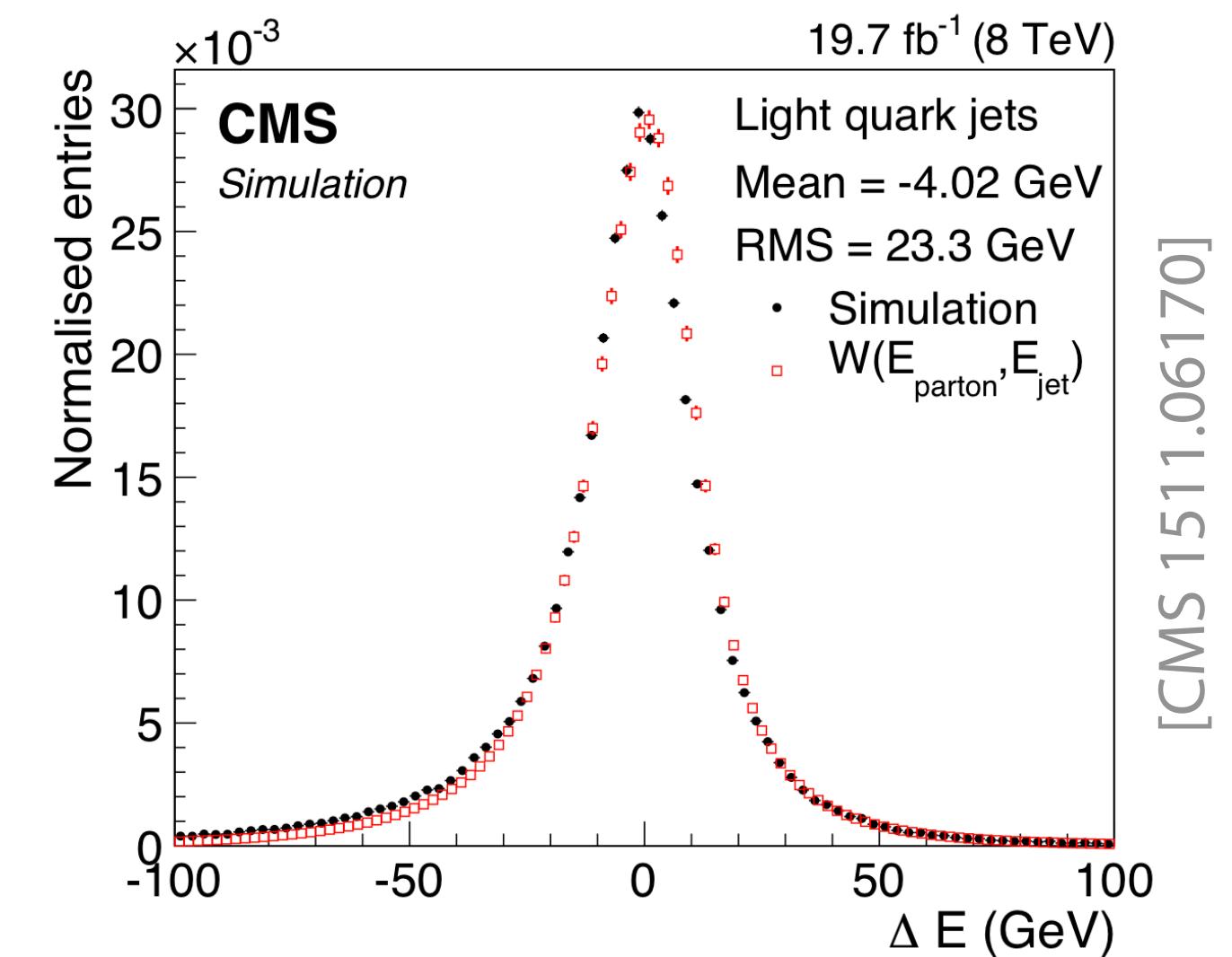
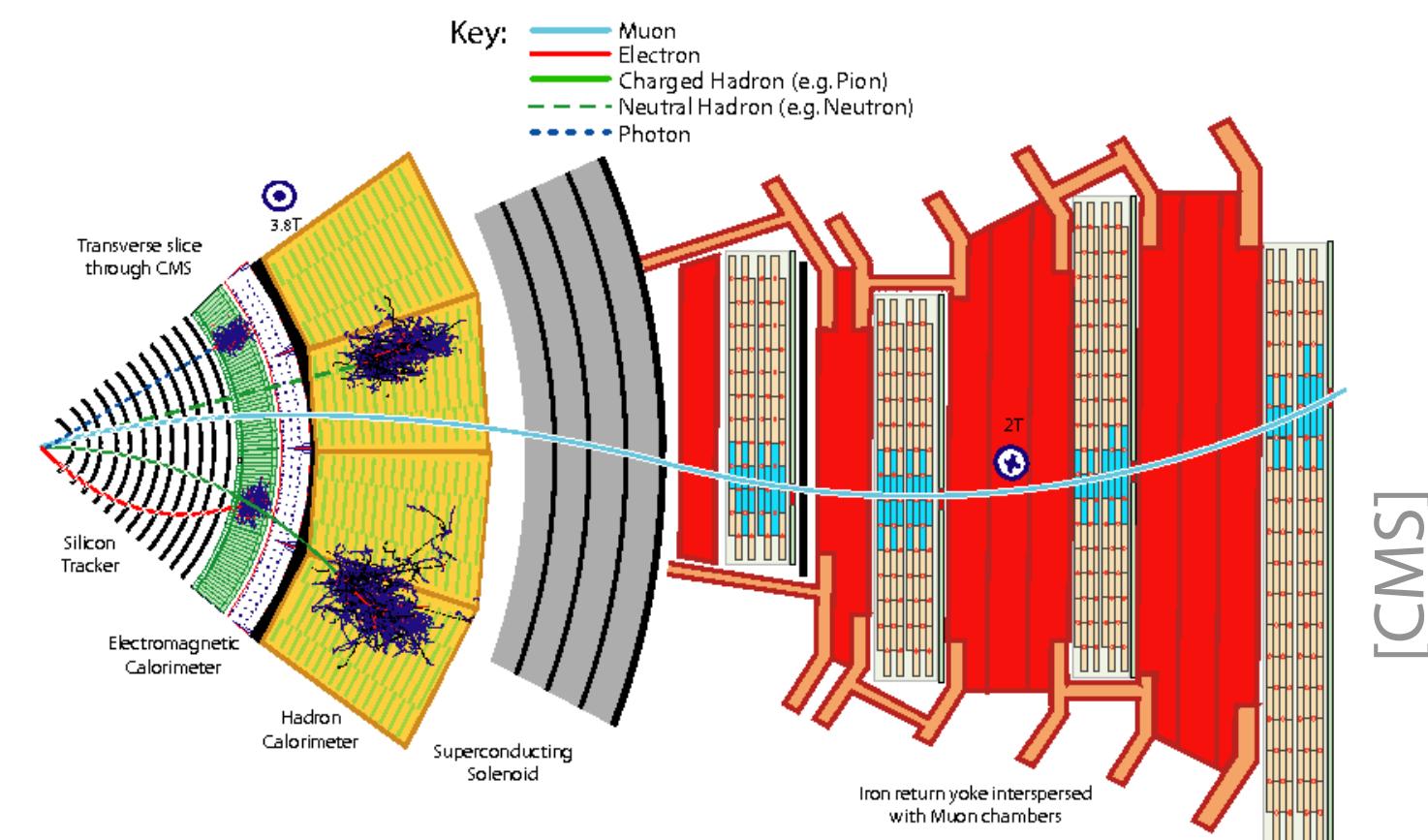
- Problem: high-dimensional integral over **shower / detector trajectories**

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

- Matrix Element Method (and similarly Optimal Observables): [K. Kondo 1988]

- approximate **shower + detector effects** into **transfer function** $\hat{p}(x|z_p)$
- explicitly calculate remaining integral

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



Solve it by approximating the integral

- Problem: high-dimensional integral over **shower / detector trajectories**

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

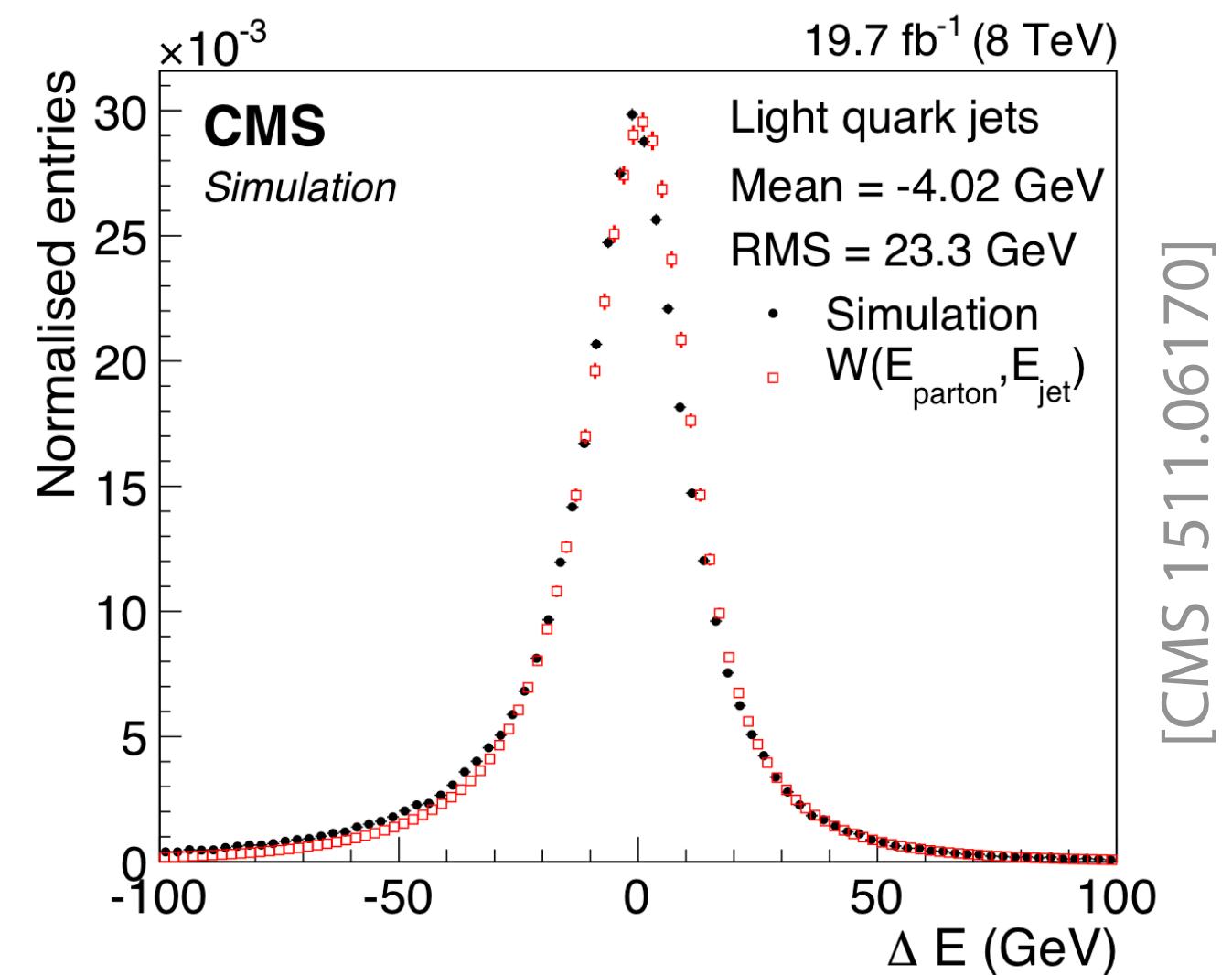
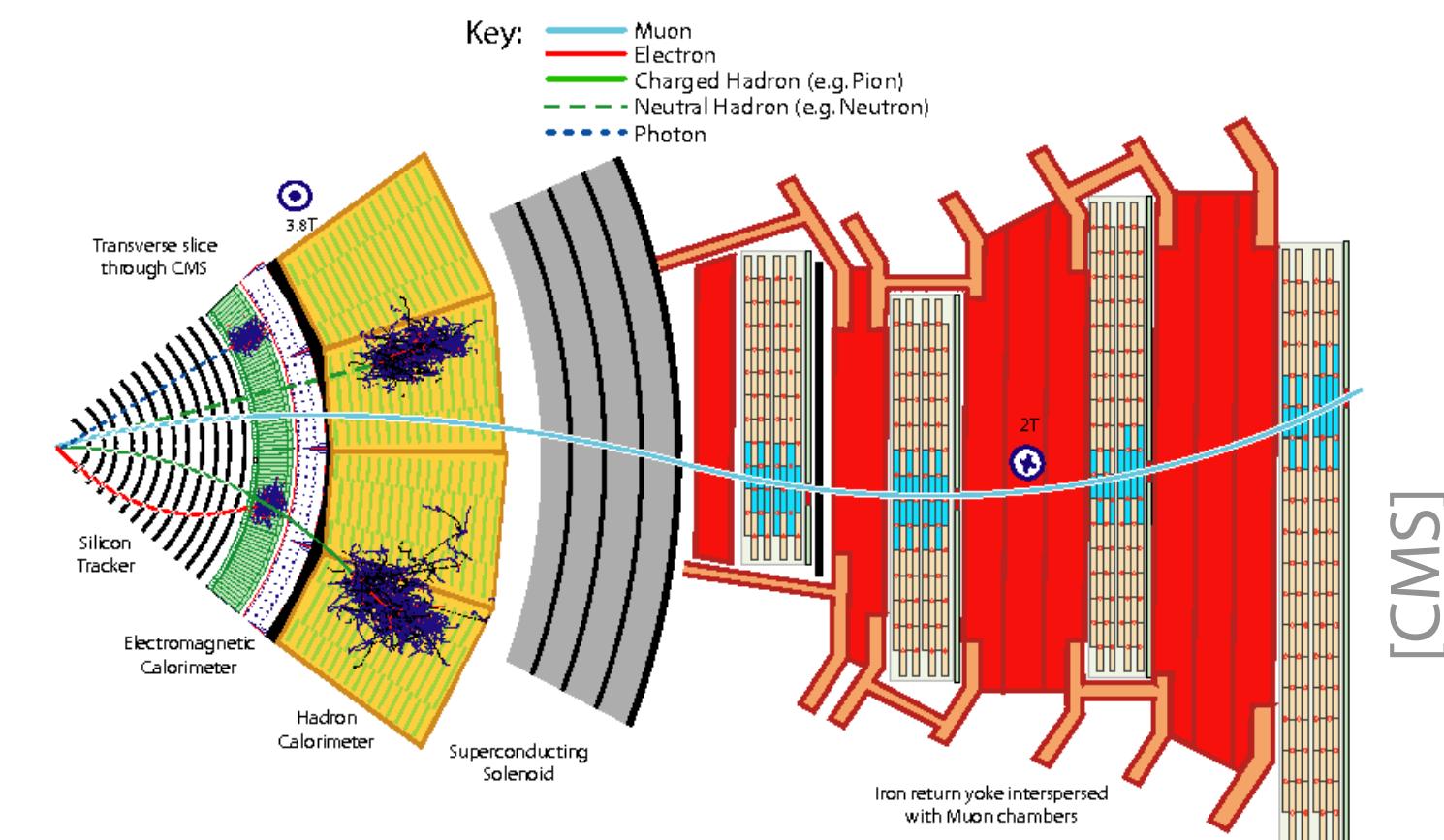
- Matrix Element Method (and similarly Optimal Observables): [K. Kondo 1988]

- approximate **shower + detector effects** into **transfer function** $\hat{p}(x|z_p)$
- explicitly calculate remaining integral

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

⇒ Uses matrix-element information, no summary statistics necessary, but:

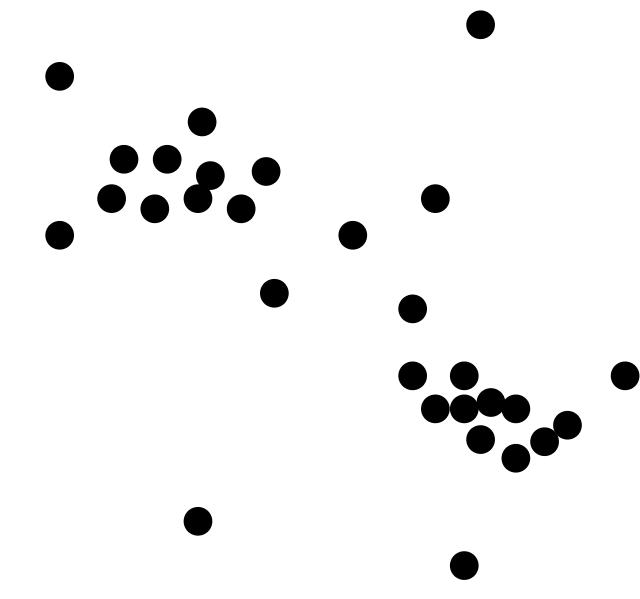
- ad-hoc transfer functions (what about extra radiation?)
- evaluation still requires calculating an expensive integral



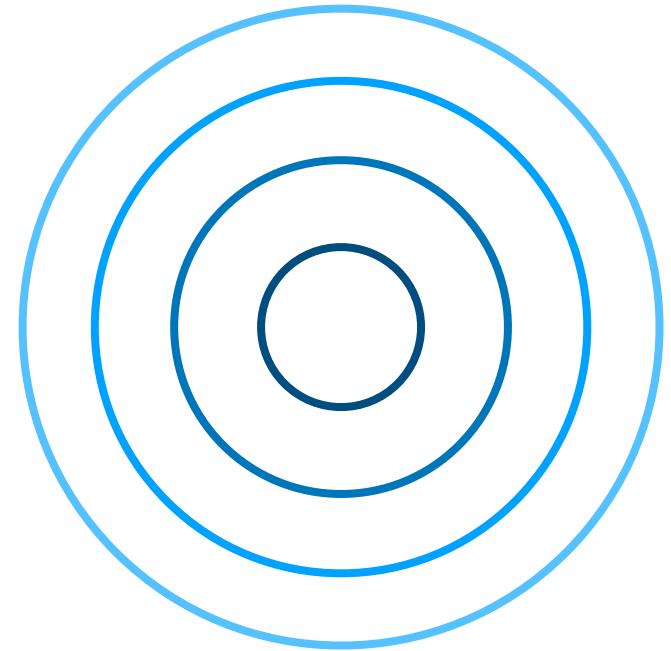
But we can use machine learning
to make inference possible.

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

High-dimensional density estimation with normalizing flows

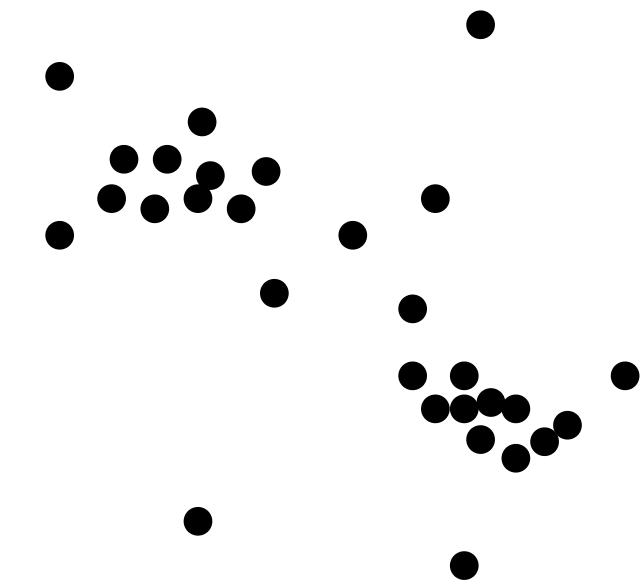


High-dimensional density estimation with normalizing flows

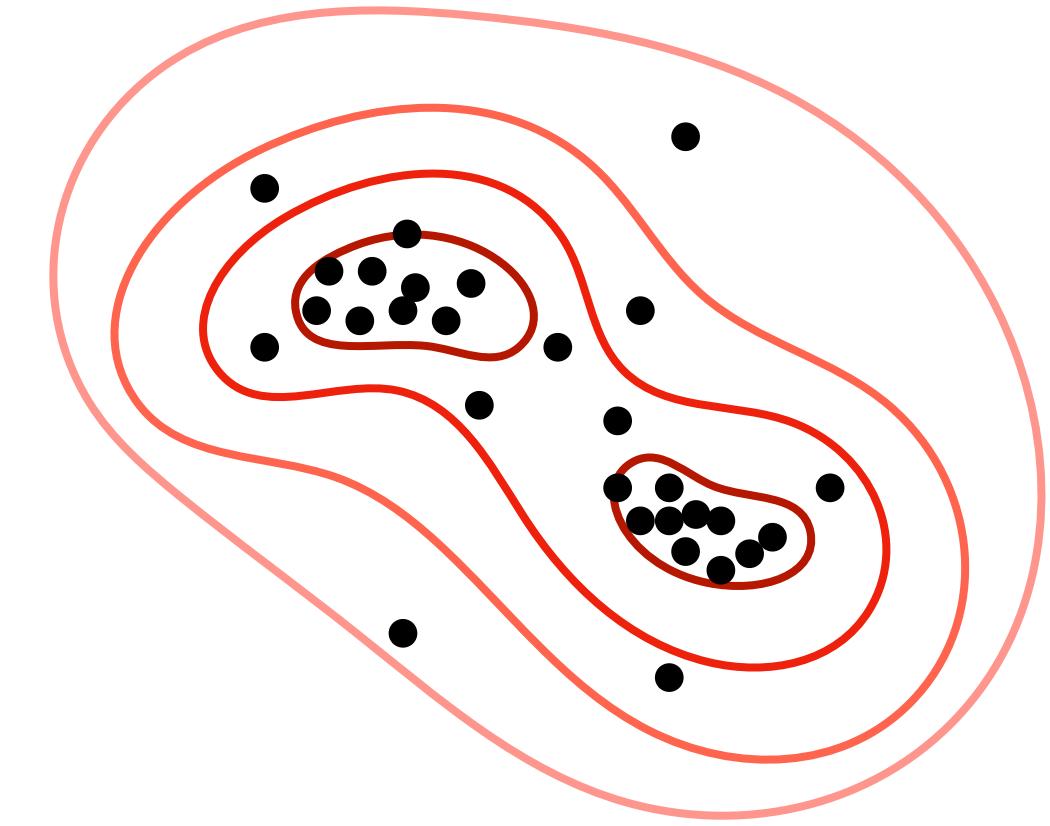
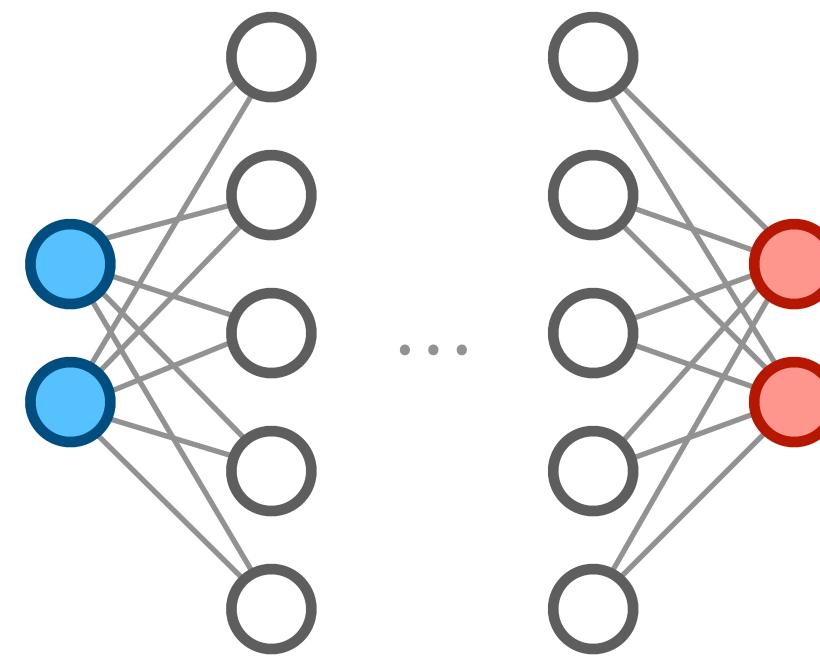
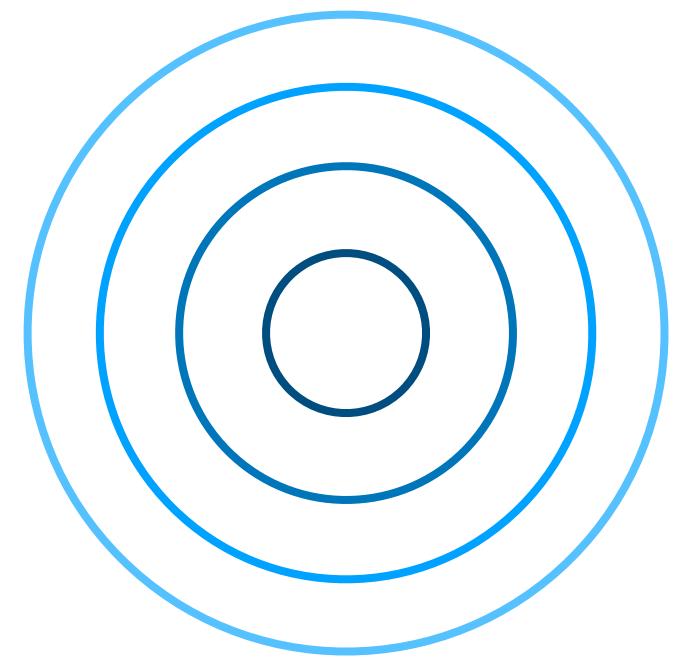


Simple base density

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



High-dimensional density estimation with normalizing flows



Simple base density

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

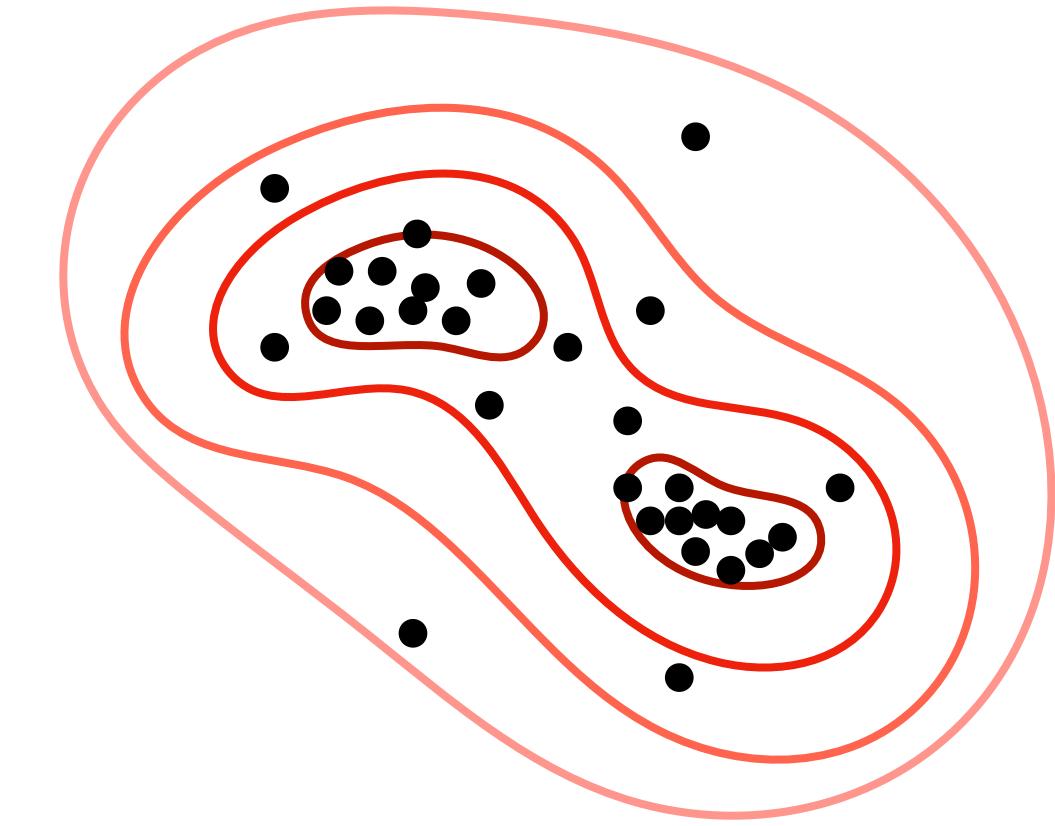
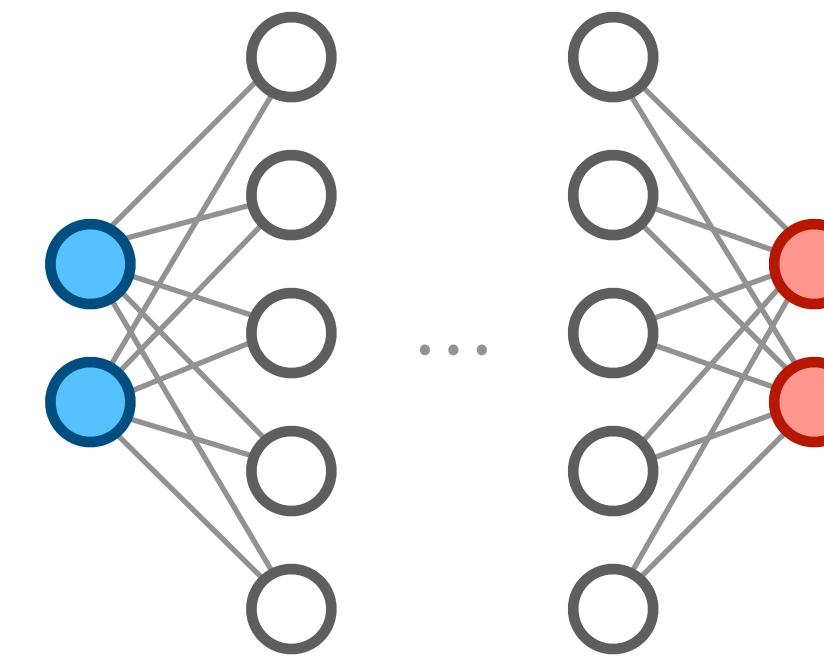
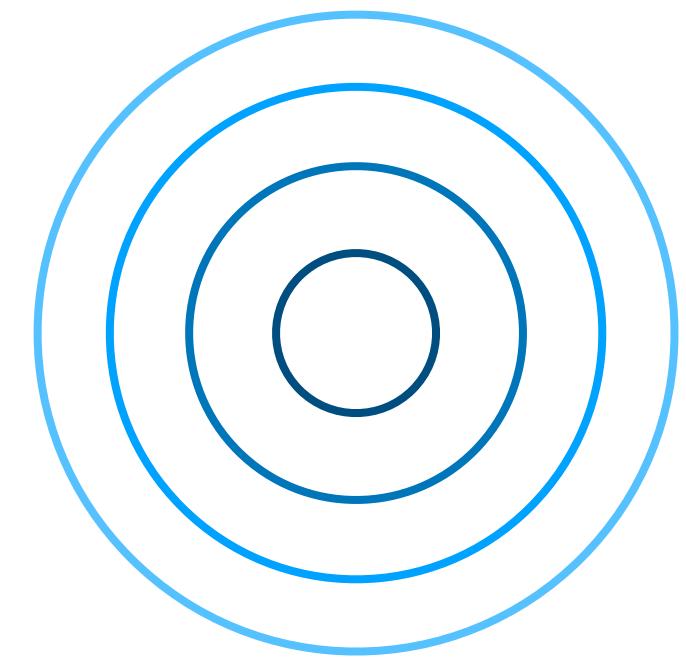
NN: transformation $x = f(u)$

- 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 $x = f(u)$

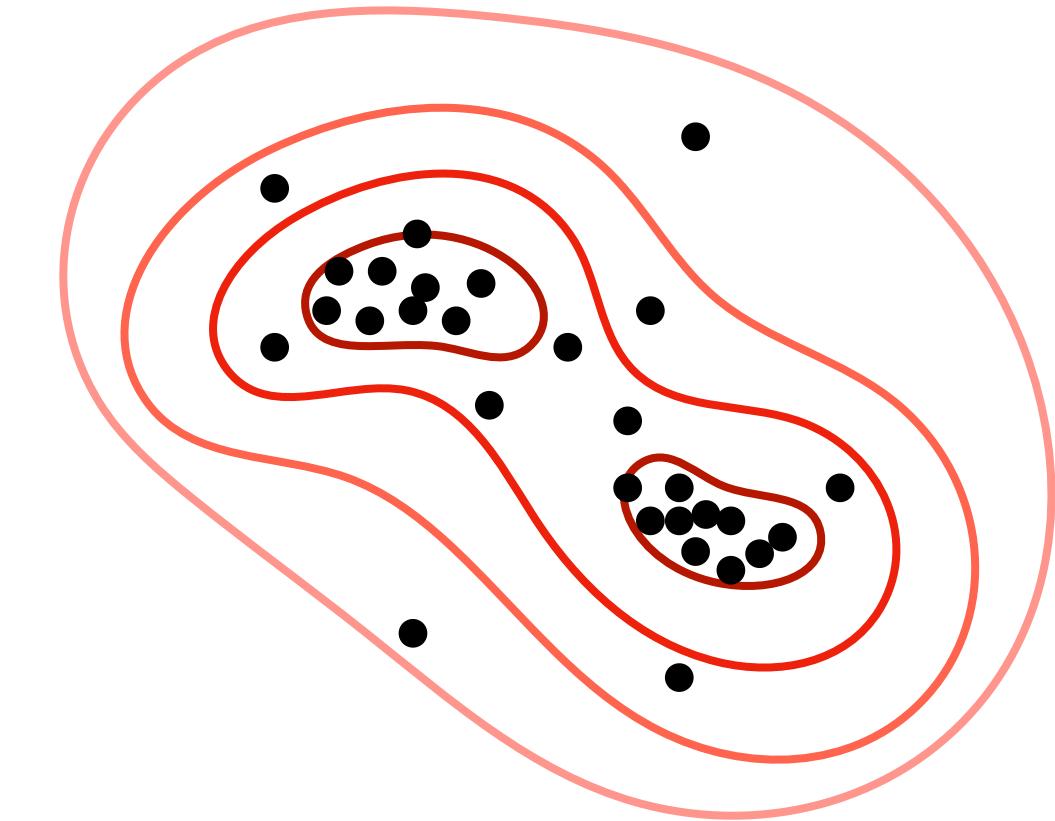
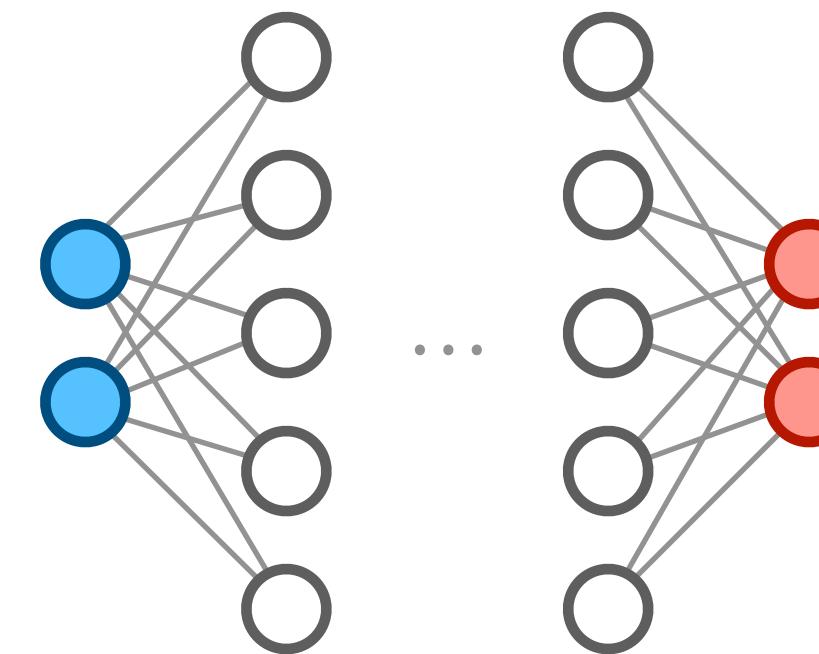
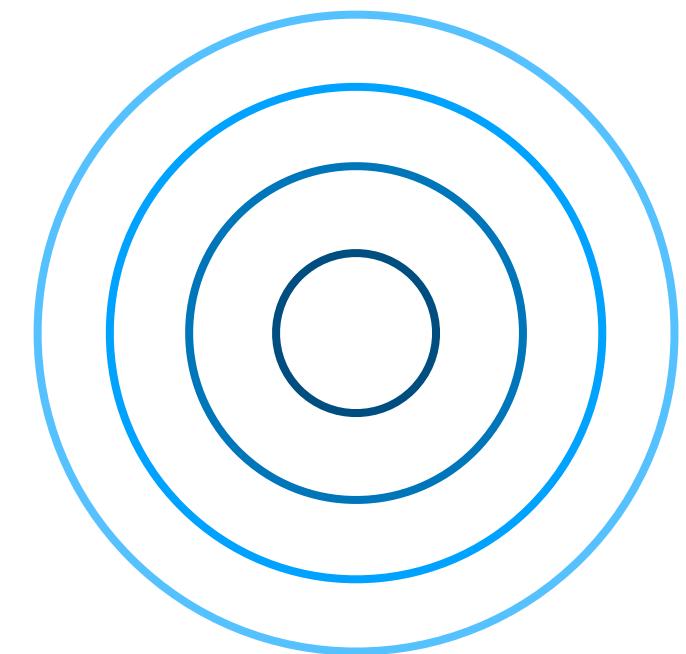
- one-to-one and invertible
- differentiable
- f^{-1} and $\det \nabla f$ are tractable

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

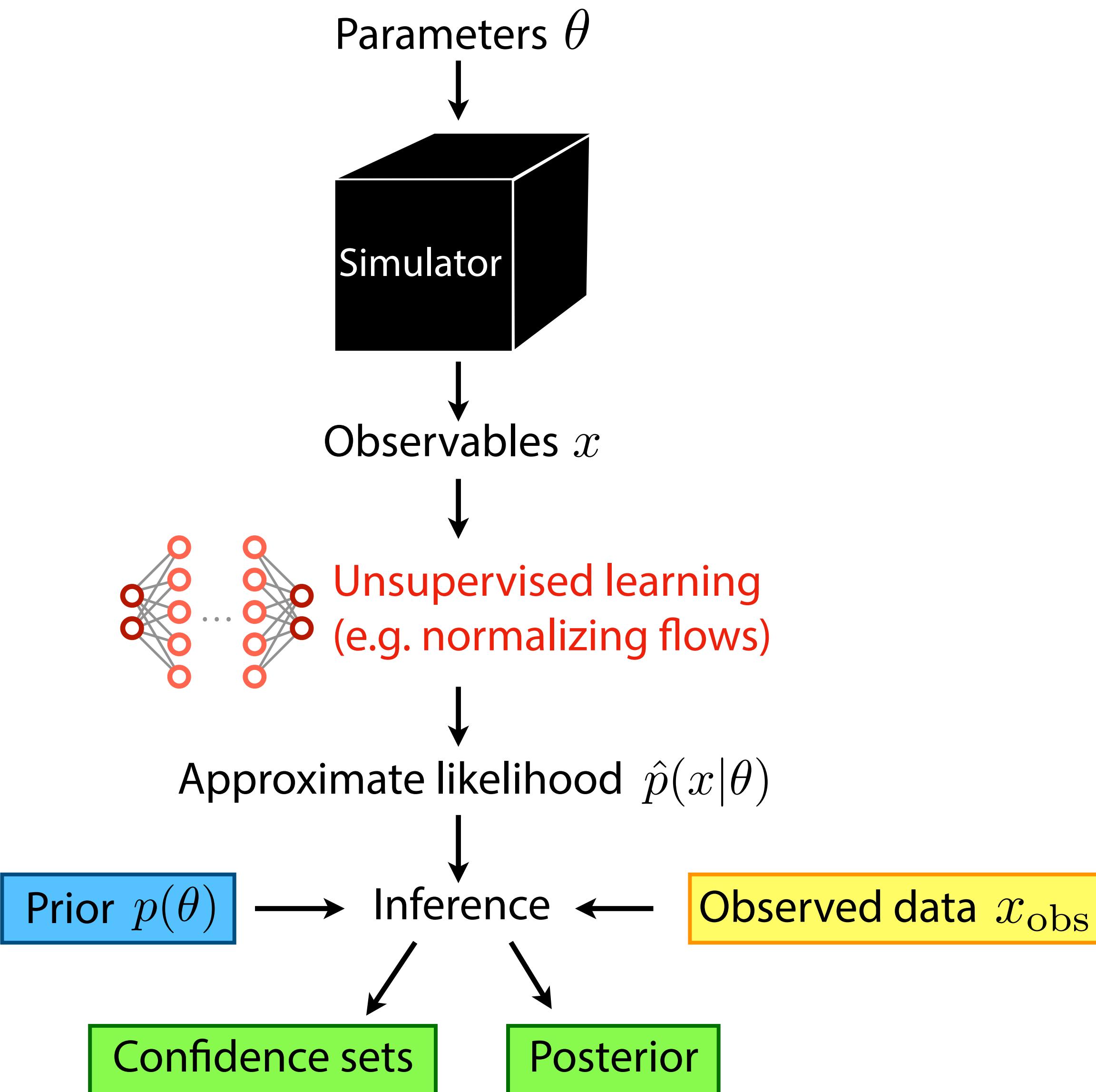
Target density is given by

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

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

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

Inference with neural likelihood estimation



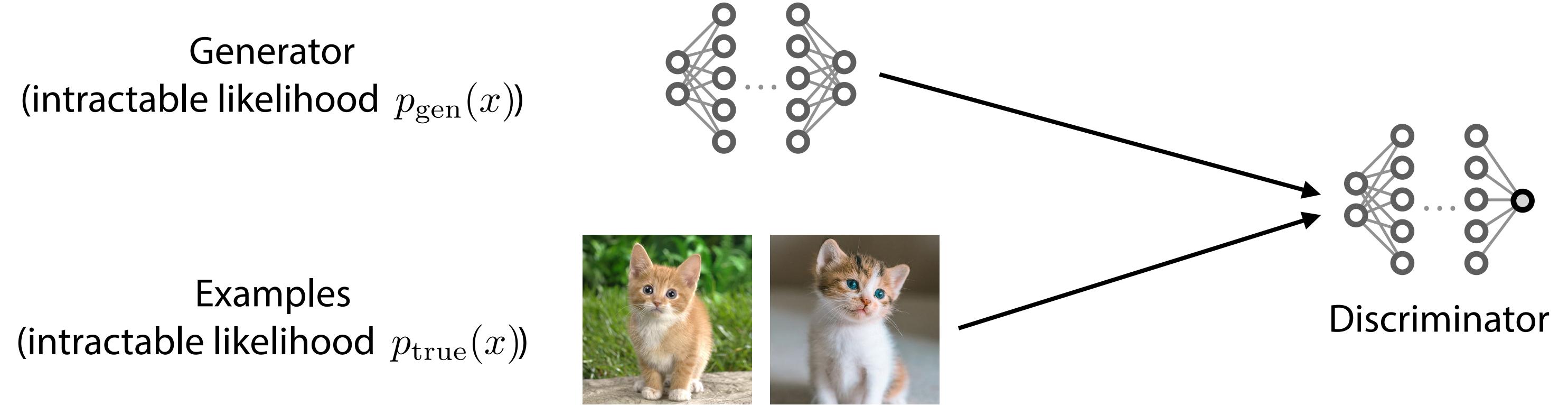
[G. Papamakarios, D. Sterratt, I. Murray 1805.07226;
J.-M. Lueckmann, G. Bassetto, T. Karaletsos, J. Macke 1805.09294]

- Conditional neural density estimator (e.g. normalizing flow) as tractable surrogate for simulator likelihood
- 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
- Related alternative: learn posterior $\hat{p}(\theta|x_{\text{obs}})$

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

The likelihood ratio trick

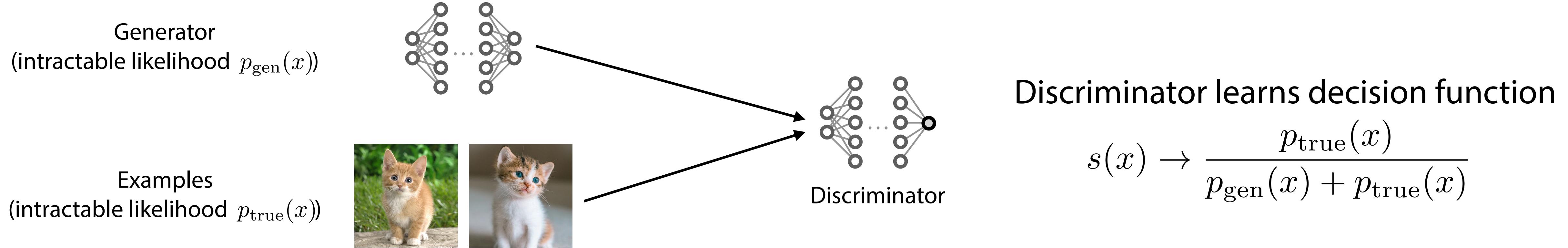
- Generative Adversarial Networks (GANs):



[I. Goodfellow et al. 1406.2661]

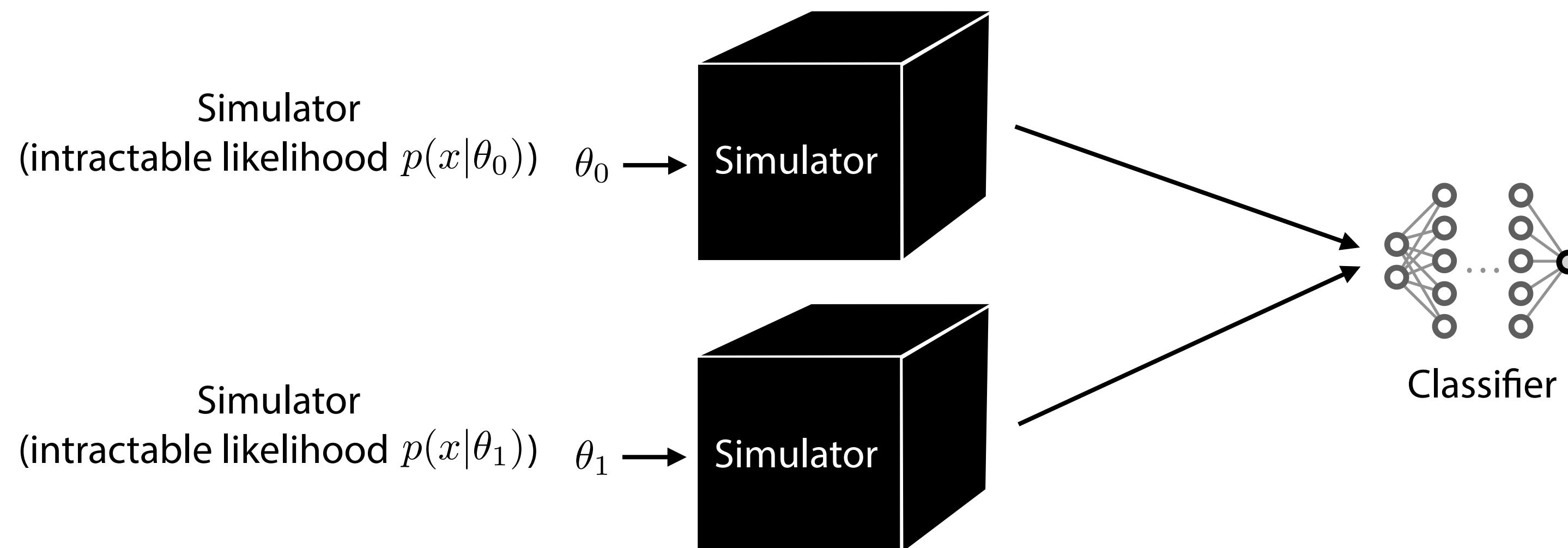
The likelihood ratio trick

- Generative Adversarial Networks (GANs):



[I. Goodfellow et al. 1406.2661]

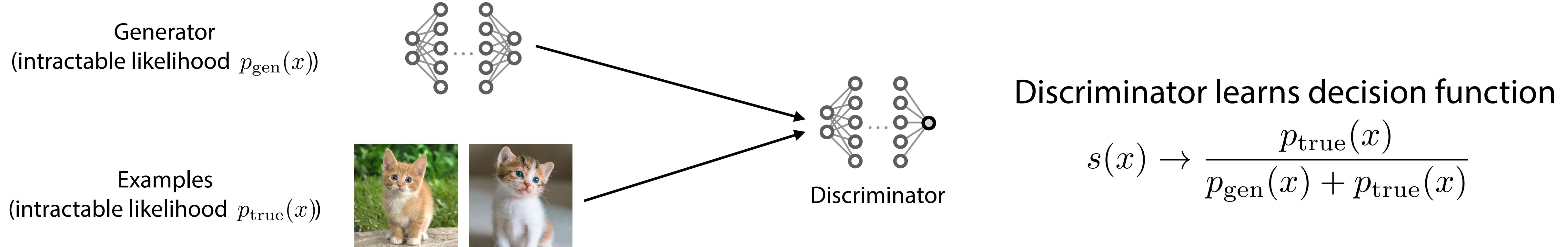
- Similarly, we can train a classifier between two sets of simulated samples



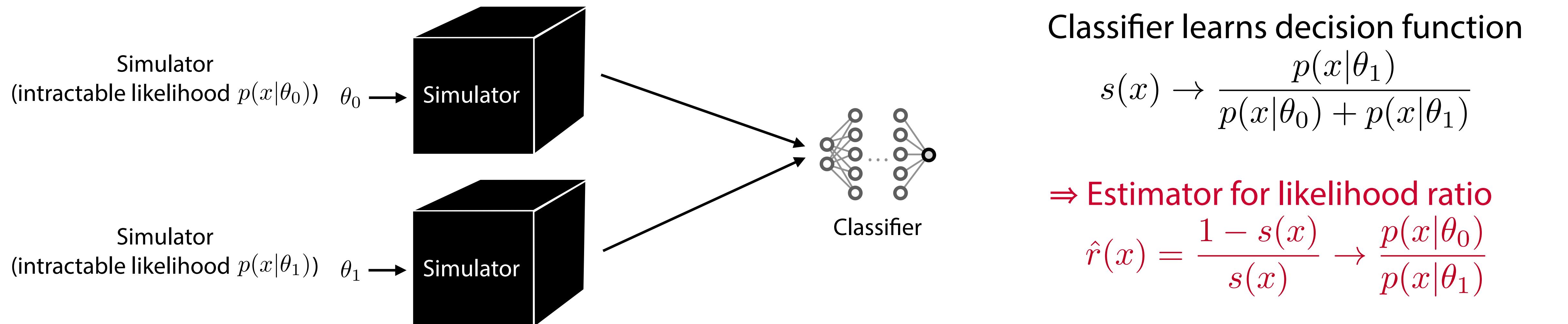
[K. Cranmer, J. Pavez, G. Louppe 1506.02169]

The likelihood ratio trick

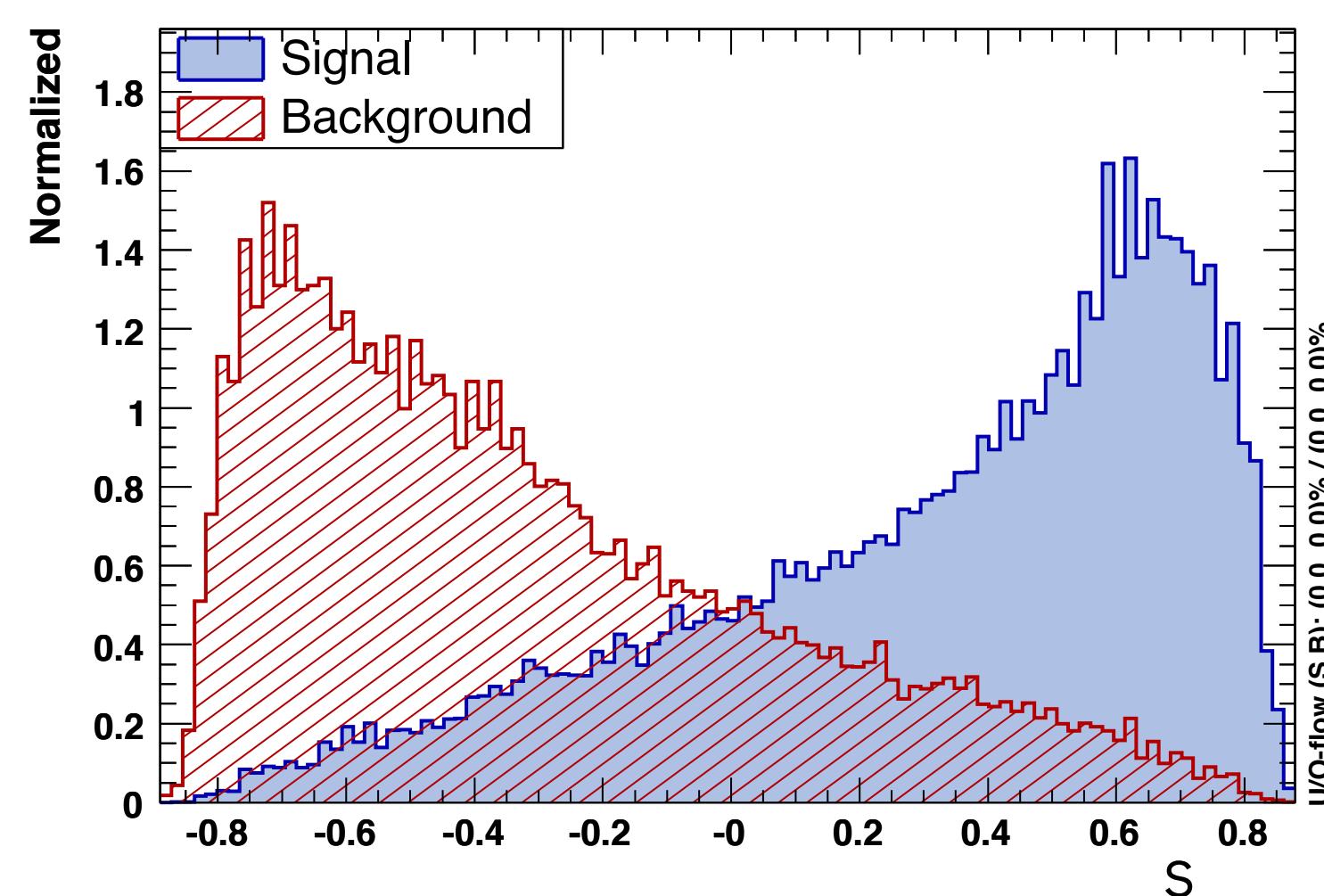
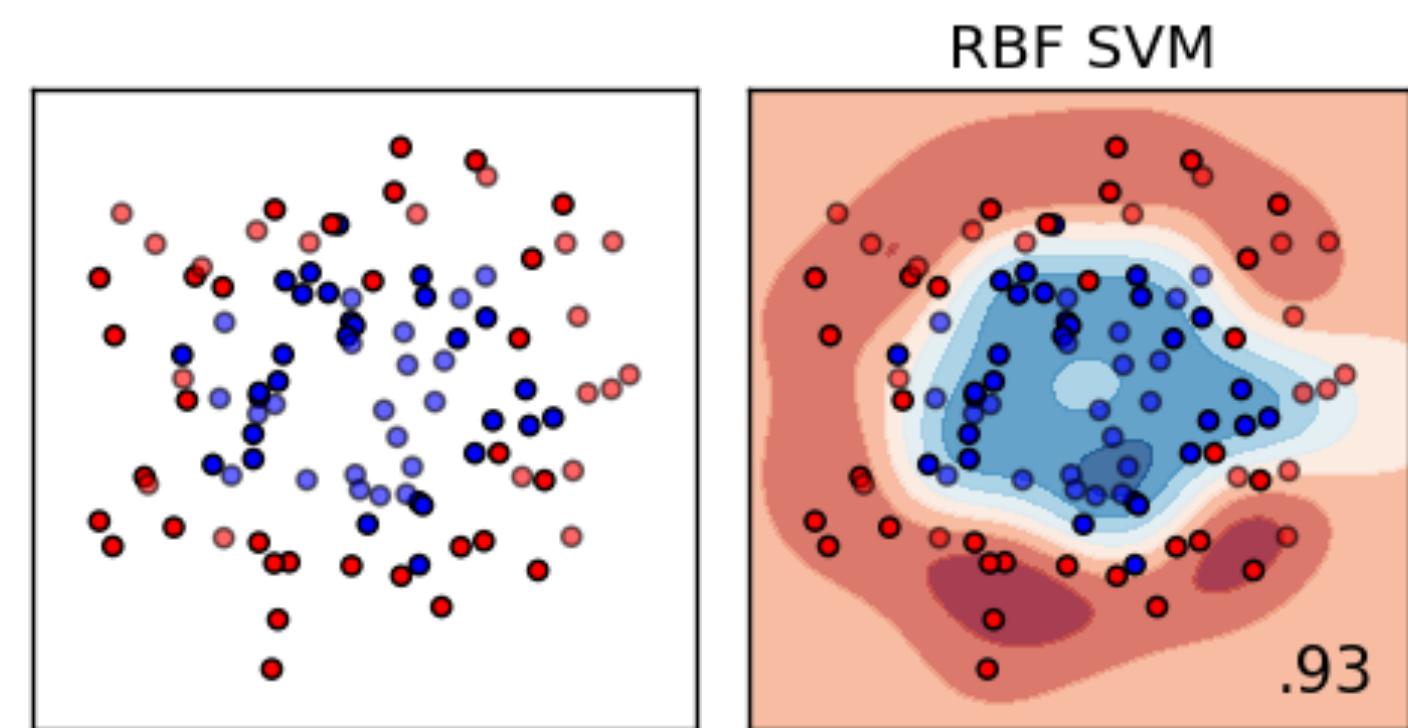
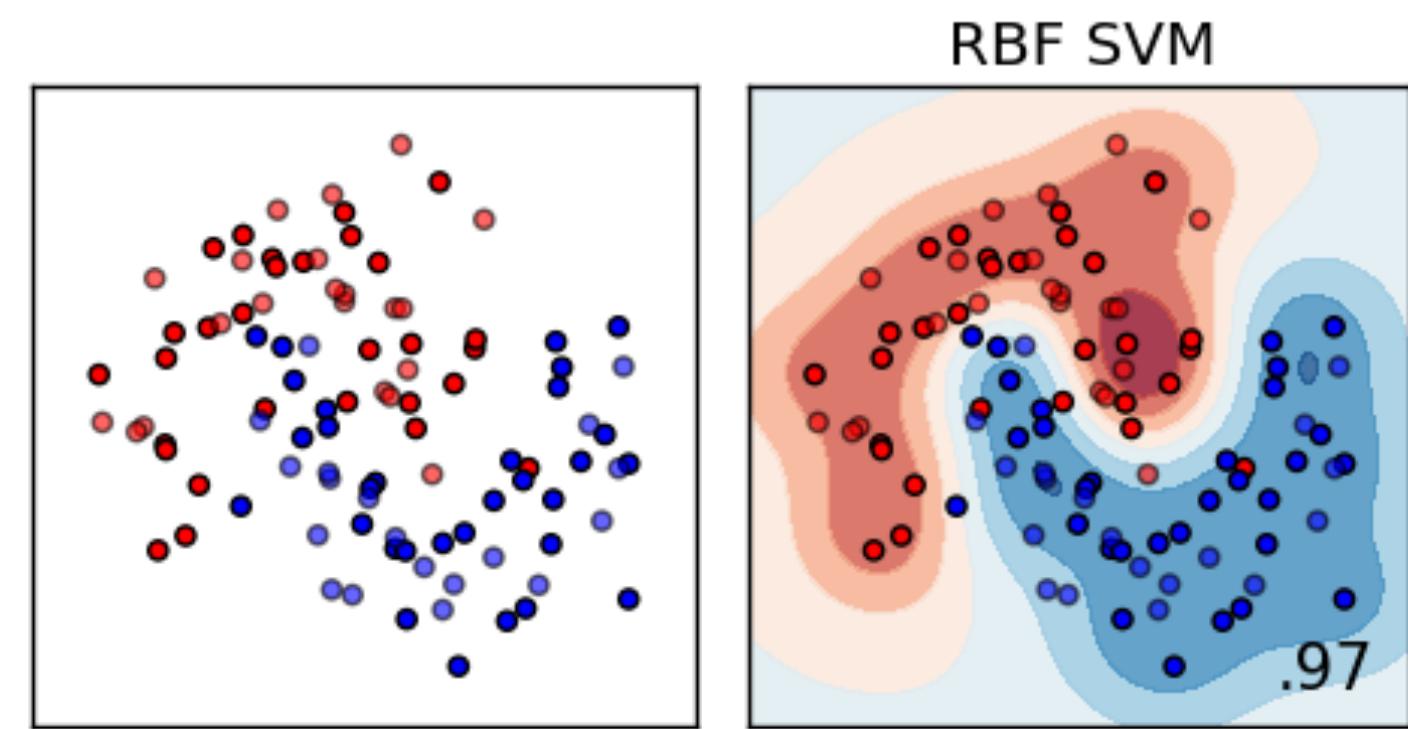
- Generative Adversarial Networks (GANs):



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



LIKELIHOOD RATIO TRICK



- **binary classifier:** find function $s(x)$ that minimizes **loss**:

$$L[s] = \int p(x|H_0) (0 - s(x))^2 dx + \int p(x|H_1) (1 - s(x))^2 dx$$

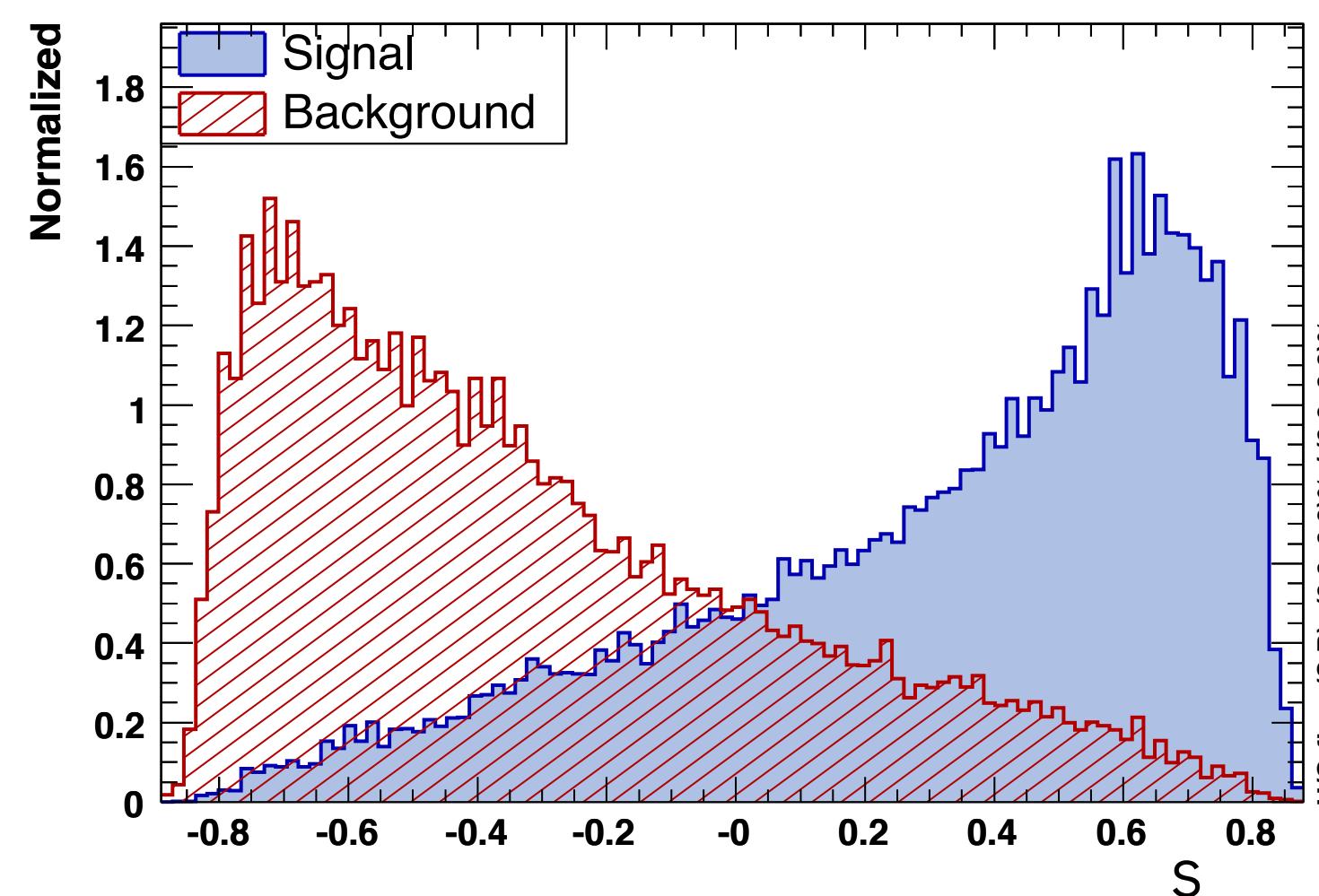
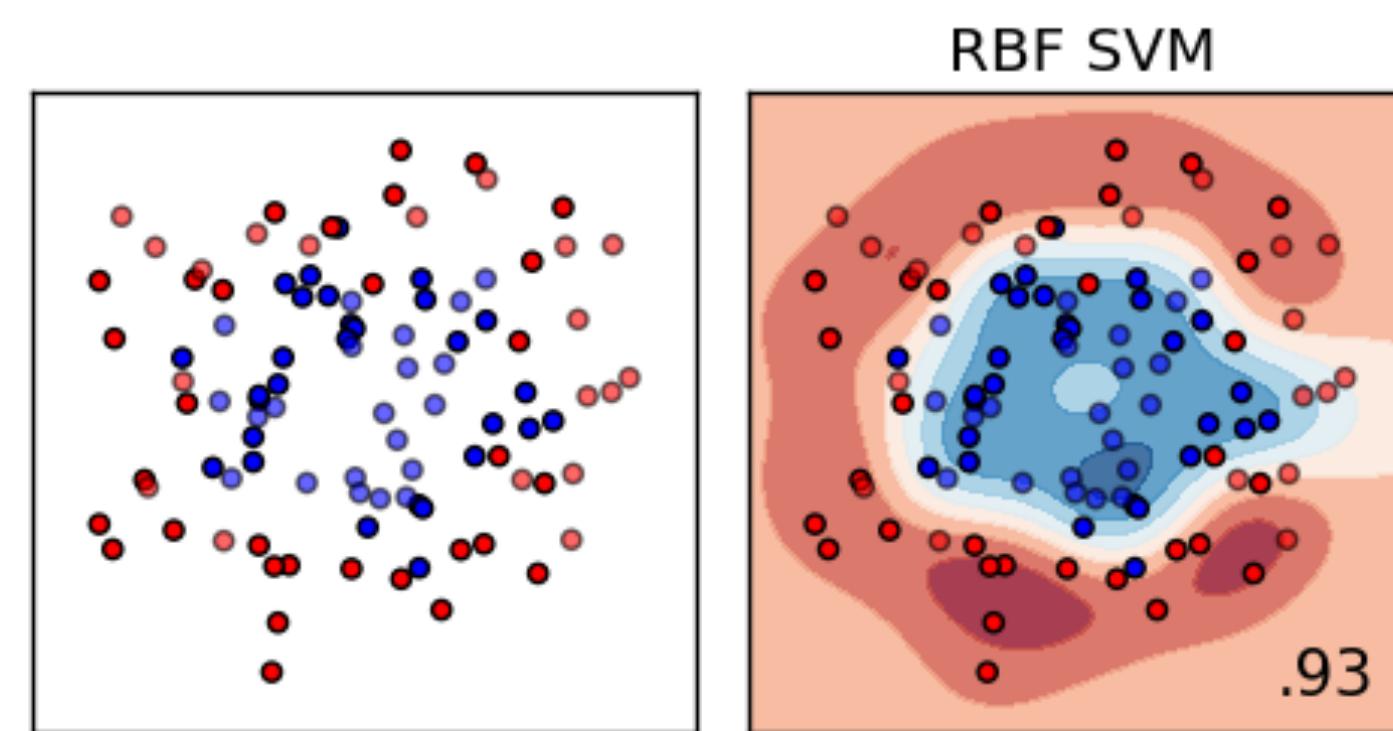
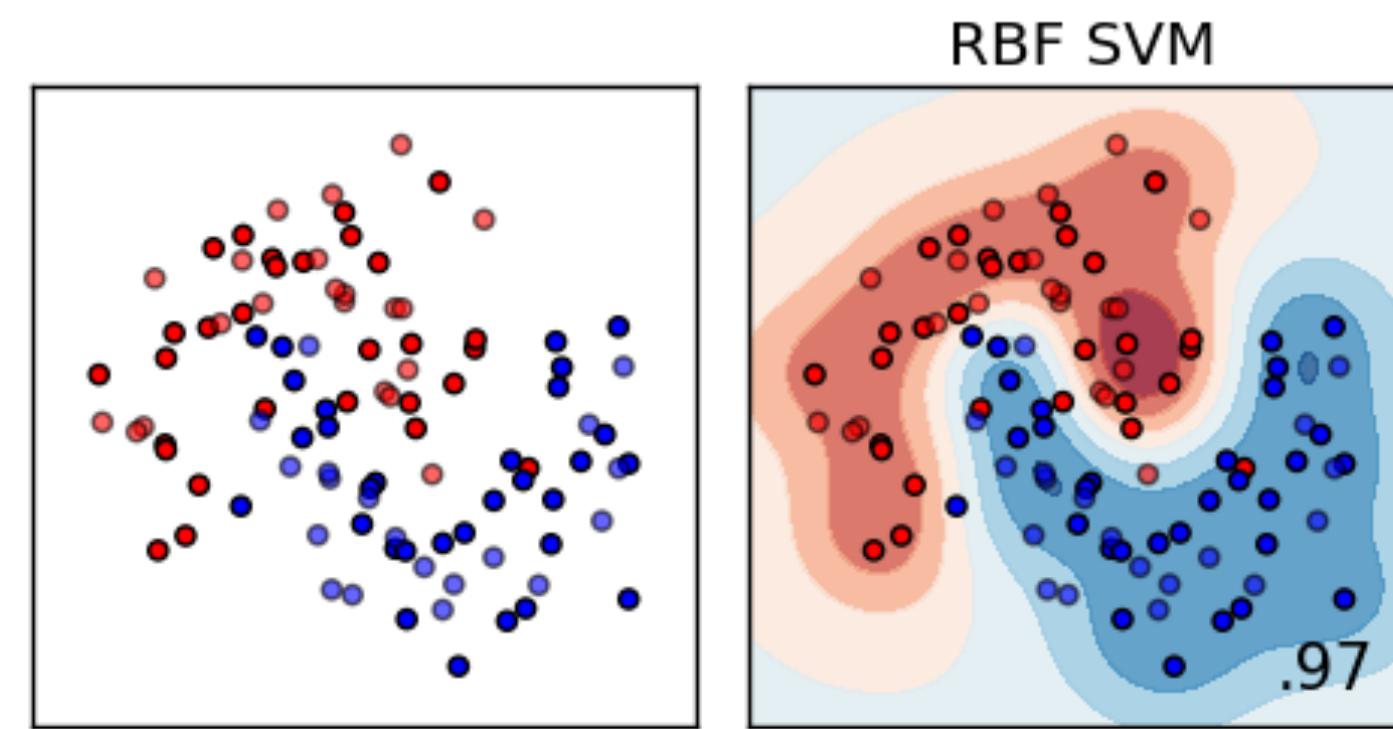
- i.e. approximate the optimal classifier

$$s(x) = \frac{p(x|H_1)}{p(x|H_0) + p(x|H_1)}$$

- which is 1-to-1 with the likelihood ratio

$$\frac{p(x|H_1)}{p(x|H_0)}$$

LIKELIHOOD RATIO TRICK



- **binary classifier:** find function $s(x)$ that minimizes **loss**:

$$L[s] = \int p(x|H_0) (0 - s(x))^2 dx + \int p(x|H_1) (1 - s(x))^2 dx$$

$$\approx \frac{1}{N} \sum_{i=1}^N (y_i - s(x_i))^2$$

- i.e. approximate the optimal classifier

$$s(x) = \frac{p(x|H_1)}{p(x|H_0) + p(x|H_1)}$$

- which is 1-to-1 with the likelihood ratio

$$\frac{p(x|H_1)}{p(x|H_0)}$$

EXTENDING THE LIKELIHOOD RATIO TRICK

A binary classifier approximates

$$s(x) = \frac{p(x|H_1)}{p(x|H_0) + p(x|H_1)}$$

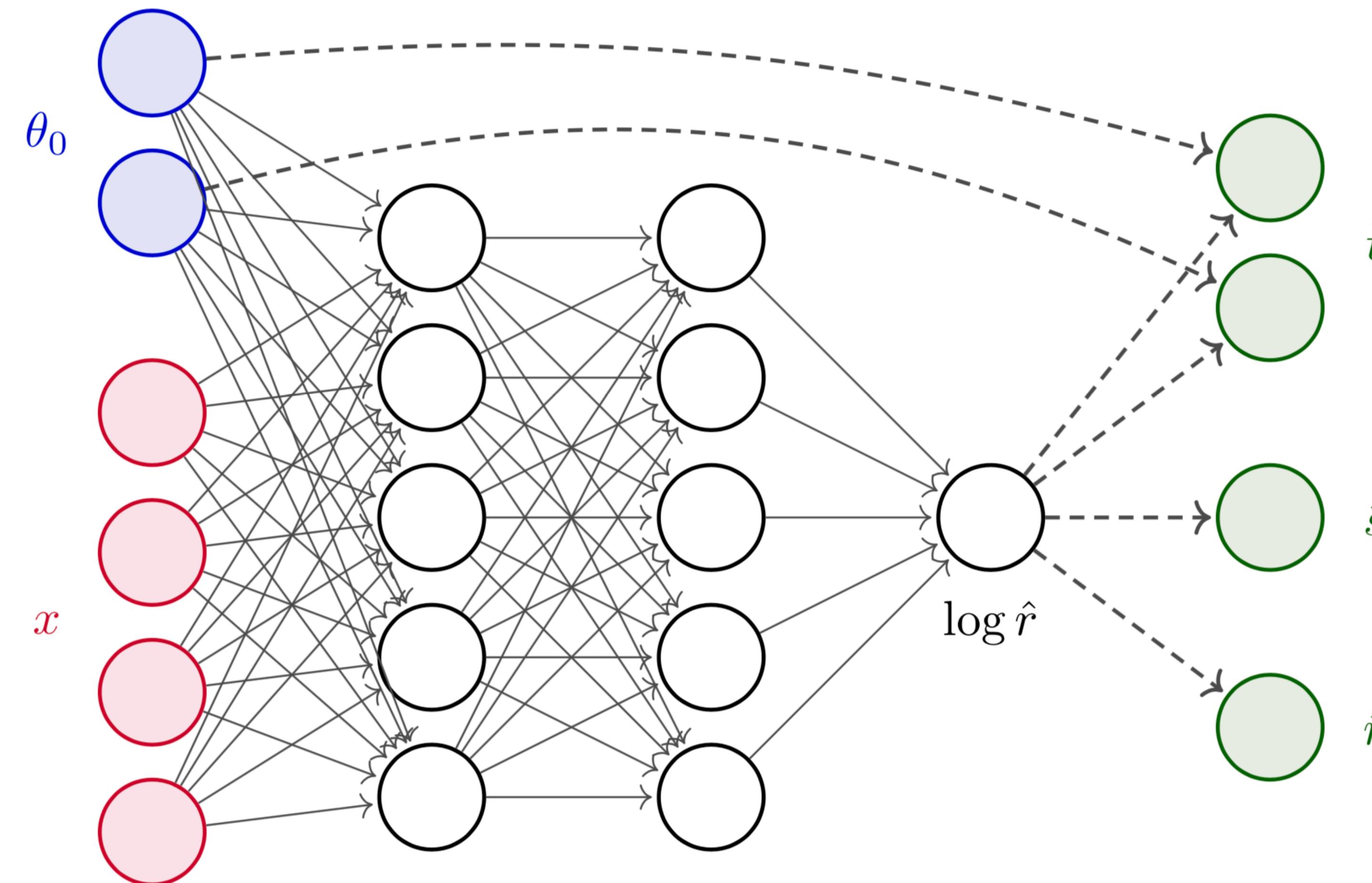
Which is one-to-one with the likelihood ratio

$$\frac{p(x|H_1)}{p(x|H_0)} = 1 - \frac{1}{s(x)}$$

Can do the same thing for any two points θ_0 & θ_1 in parameter space Θ . I call this a **parametrized classifier**

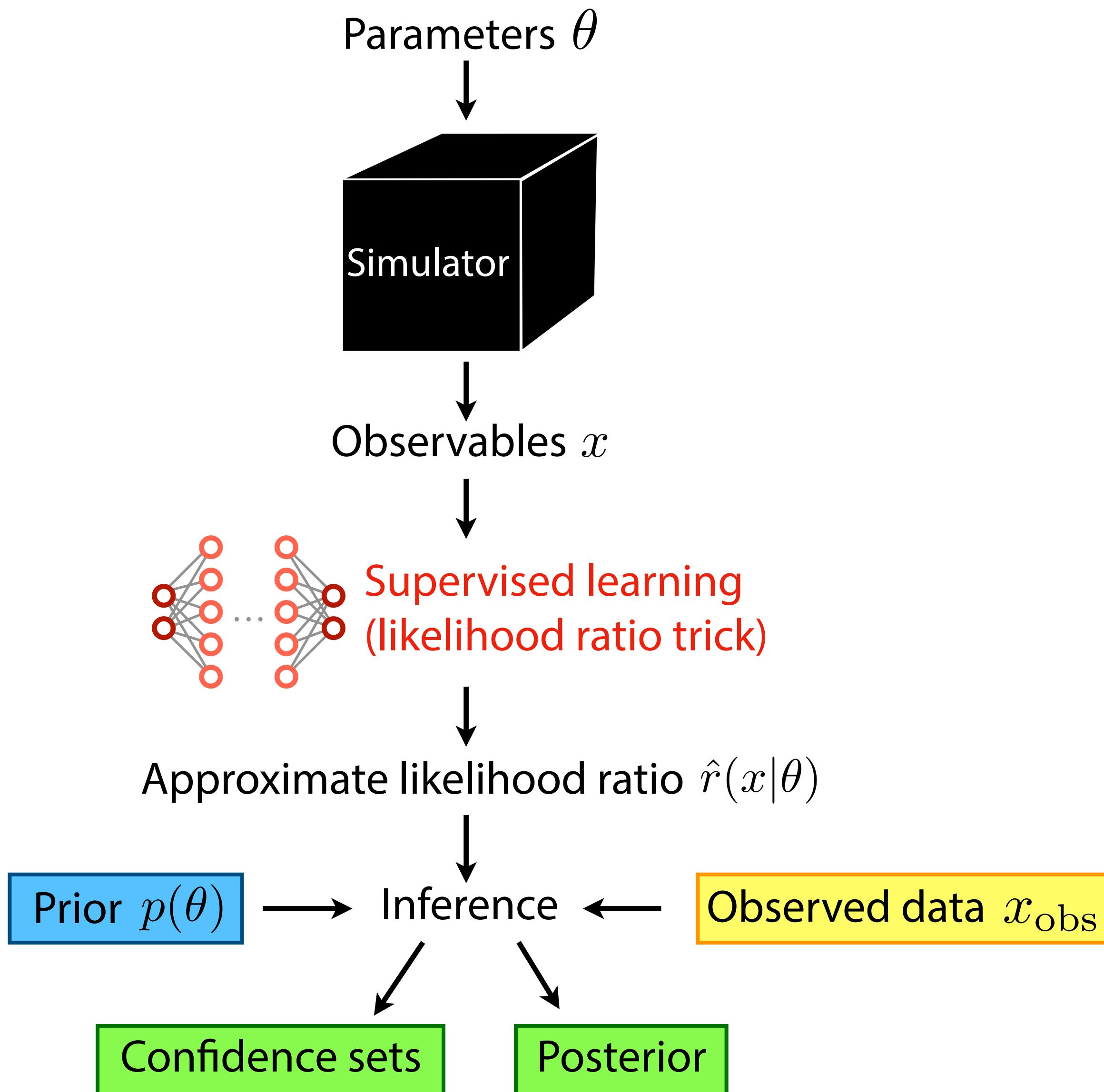
$$s(x; \theta_0, \theta_1) = \frac{p(x|\theta_1)}{p(x|\theta_0) + p(x|\theta_1)}$$

PARAMETERIZED CLASSIFIERS



Inference by likelihood ratio trick

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

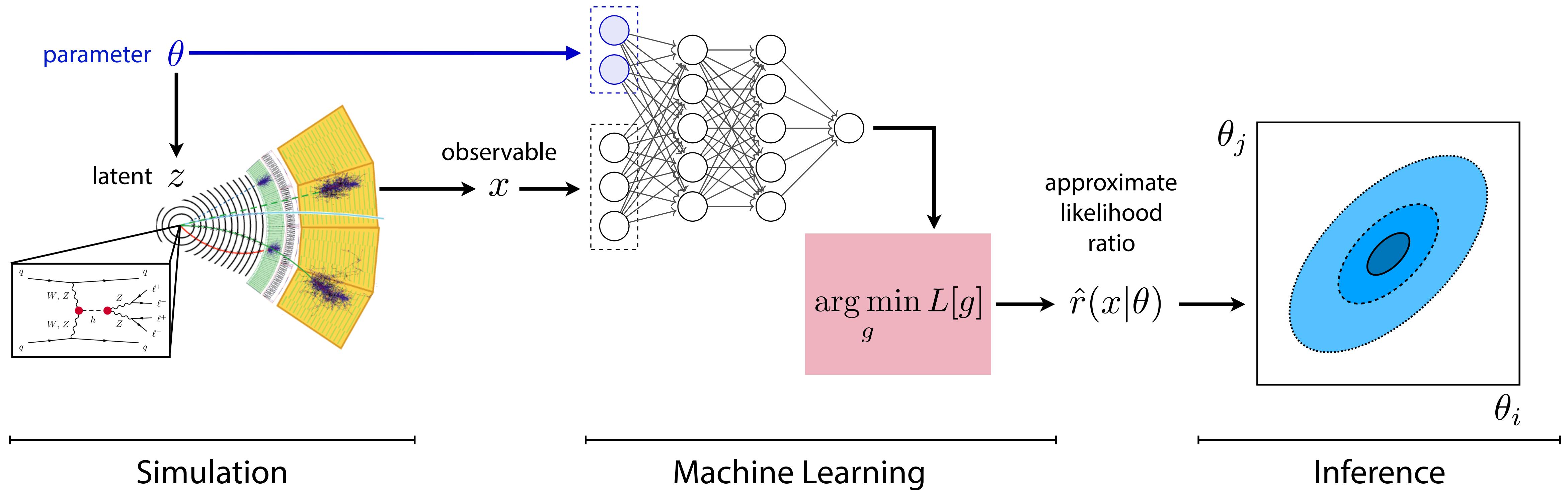


- For inference, likelihood and likelihood ratio are interchangeable
- Advantage: Learning the likelihood ratio can be a simpler task than learning the likelihood
- Disadvantage: Cannot sample from likelihood ratio

“Mining gold”:
Physics insights can make these inference methods
more efficient.

[JB, K. Cranmer, G. Louppe, J. Pavez 1805.00013, 1805.00020, 1805.12244]

Learning with Simulated Data

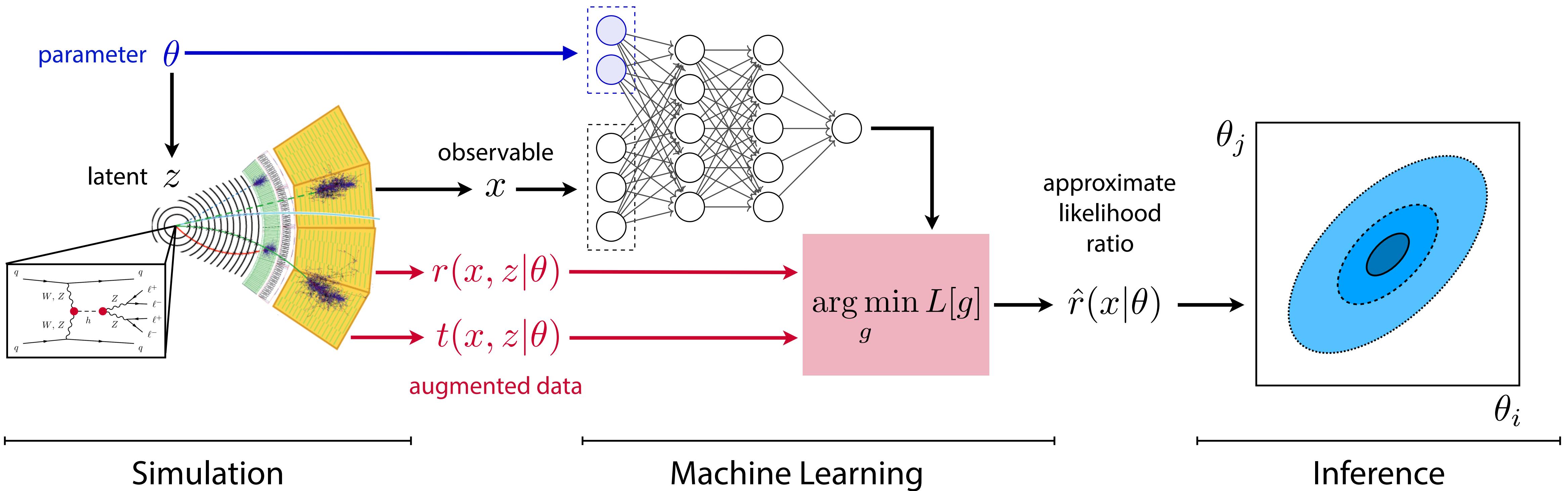


“Mining gold”: Extract additional information from simulator

Use this information to train estimator for likelihood ratio

Limit setting with standard hypothesis tests

Learning with Augmented Data

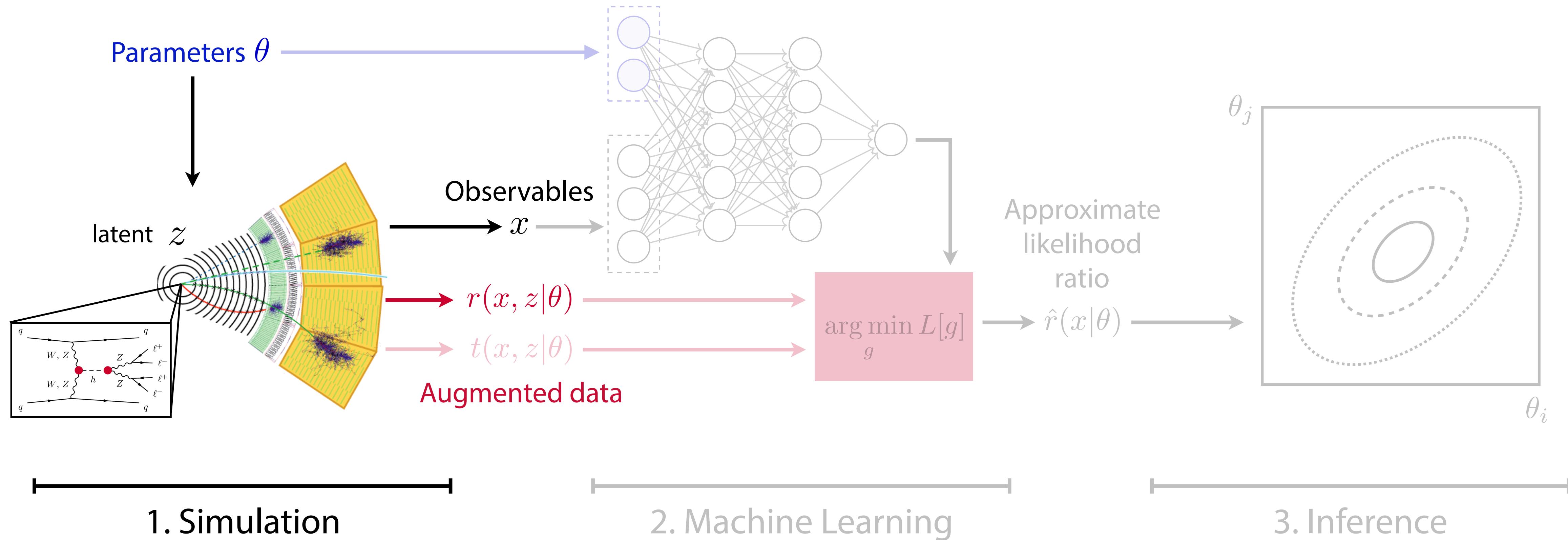


“Mining gold”: Extract additional information from simulator

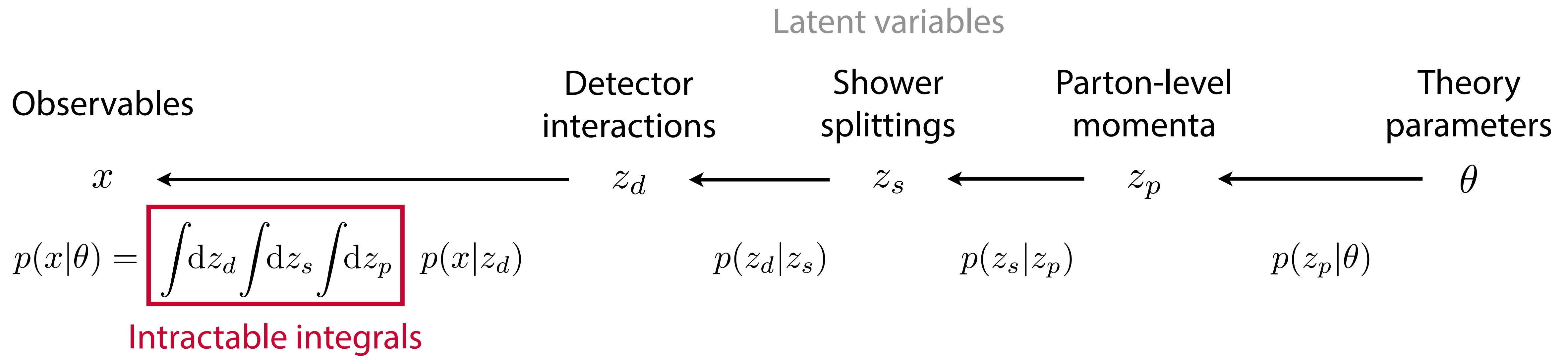
Use this information to train estimator for likelihood ratio

Limit setting with standard hypothesis tests

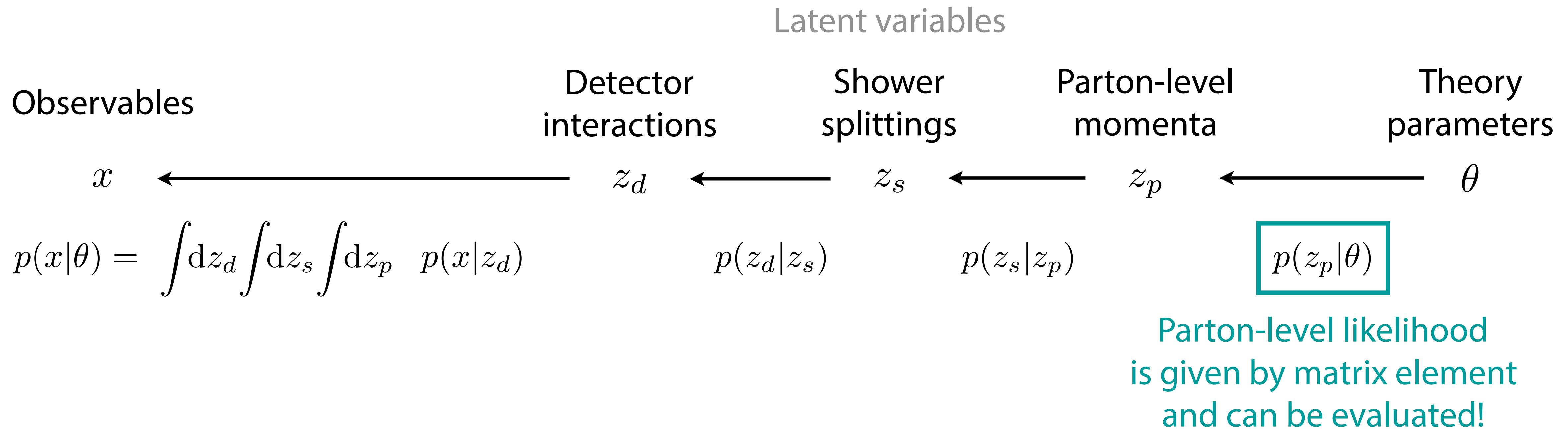
Learning with Augmented Data



Mining gold from the simulator



Mining gold from the simulator



⇒ For each simulated event, we can calculate the **joint likelihood ratio** which depends on the specific evolution of the simulation:

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

$$\frac{p(z_p|\theta_0)}{p(z_p|\theta_1)} \sim \frac{|\mathcal{M}(z_p|\theta_0)|^2}{|\mathcal{M}(z_p|\theta_1)|^2}$$

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 simulated event, including all intermediate states, for θ_0 compared to θ_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 θ_0 compared to θ_1 ?)

The value of gold

We can calculate the **joint likelihood ratio**

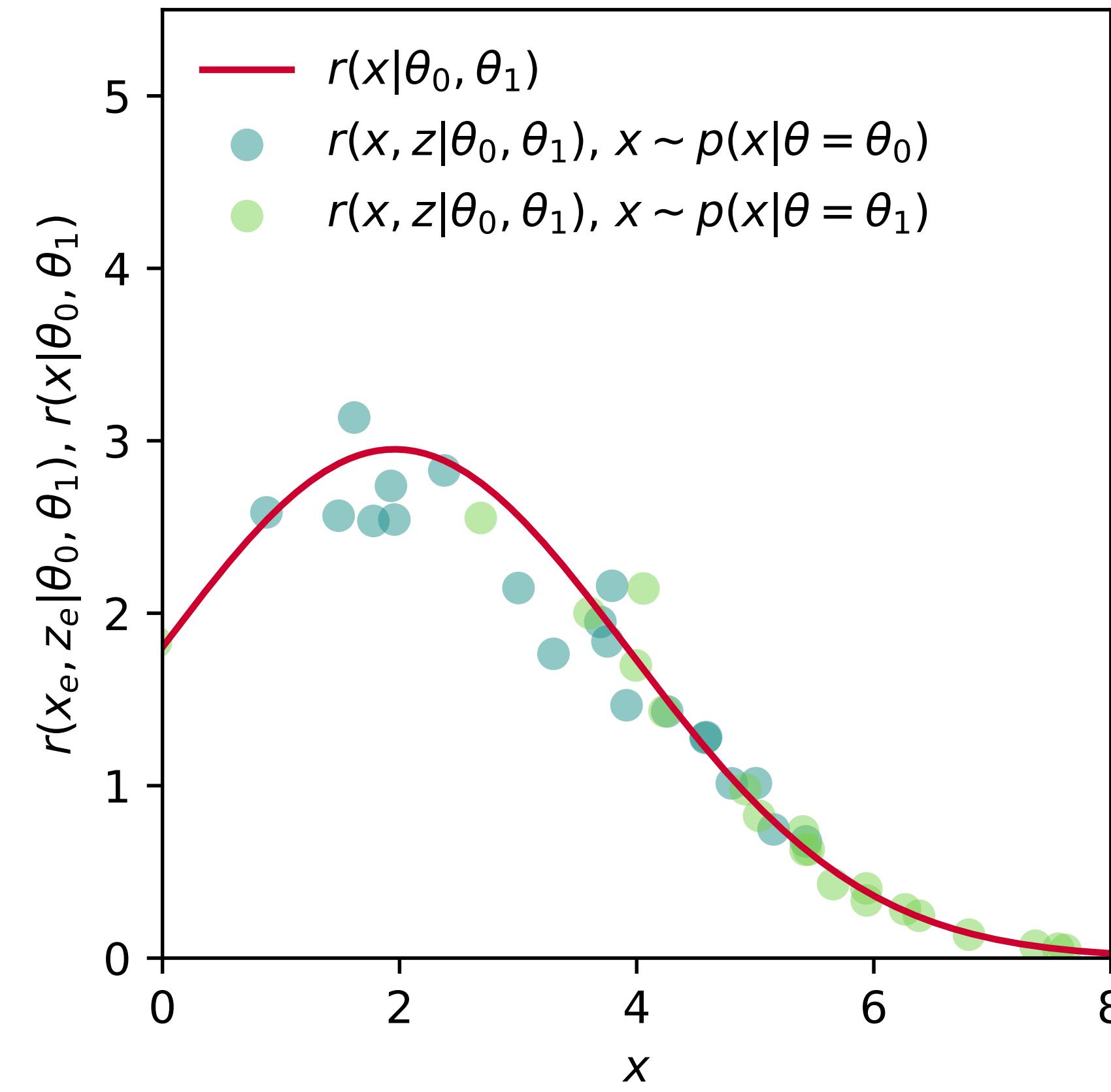
$$r(x, z | \theta_0, \theta_1) \equiv \frac{p(x, z_d, z_s, z_p | \theta_0)}{p(x, z_d, z_s, z_p | \theta_1)}$$



$r(x, z | \theta_0, \theta_1)$ 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

$$\mathbb{E}_{z \sim p(z|x, \theta_1)} [r(x, z|\theta_0, \theta_1)] = \arg \min_{\hat{r}(x|\theta_0, \theta_1)} L_r[\hat{r}(x|\theta_0, \theta_1)]$$

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



We want the **likelihood ratio function**

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

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

$$\mathbb{E}_{z \sim p(z|x, \theta_1)} [r(x, z|\theta_0, \theta_1)] = \arg \min_{\hat{r}(x|\theta_0, \theta_1)} L_r[\hat{r}(x|\theta_0, \theta_1)]!$$

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

.... and then magic ...

$$\begin{aligned} \mathbb{E}_{z \sim p(z|x, \theta_1)} [r(x, z|\theta_0, \theta_1)] &= \int dz p(z|x, \theta_1) \frac{p(x, z|\theta_0)}{p(x, z|\theta_1)} \\ &= \int dz \frac{p(x, z|\theta_1)}{p(x|\theta_1)} \frac{p(x, z|\theta_0)}{p(x, z|\theta_1)} \\ &= r(x|\theta_0, \theta_1) ! \end{aligned}$$

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:

Extremization

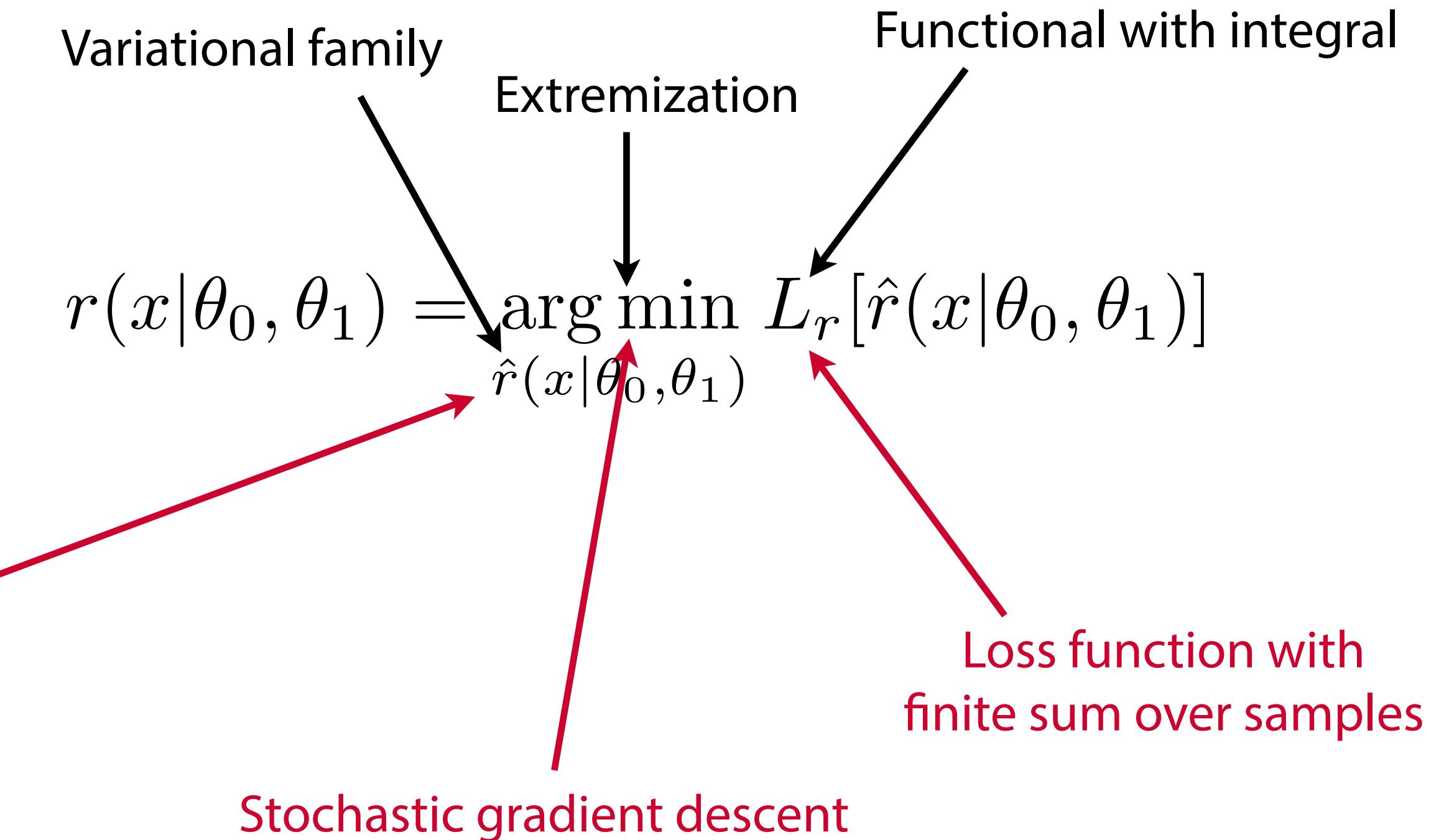
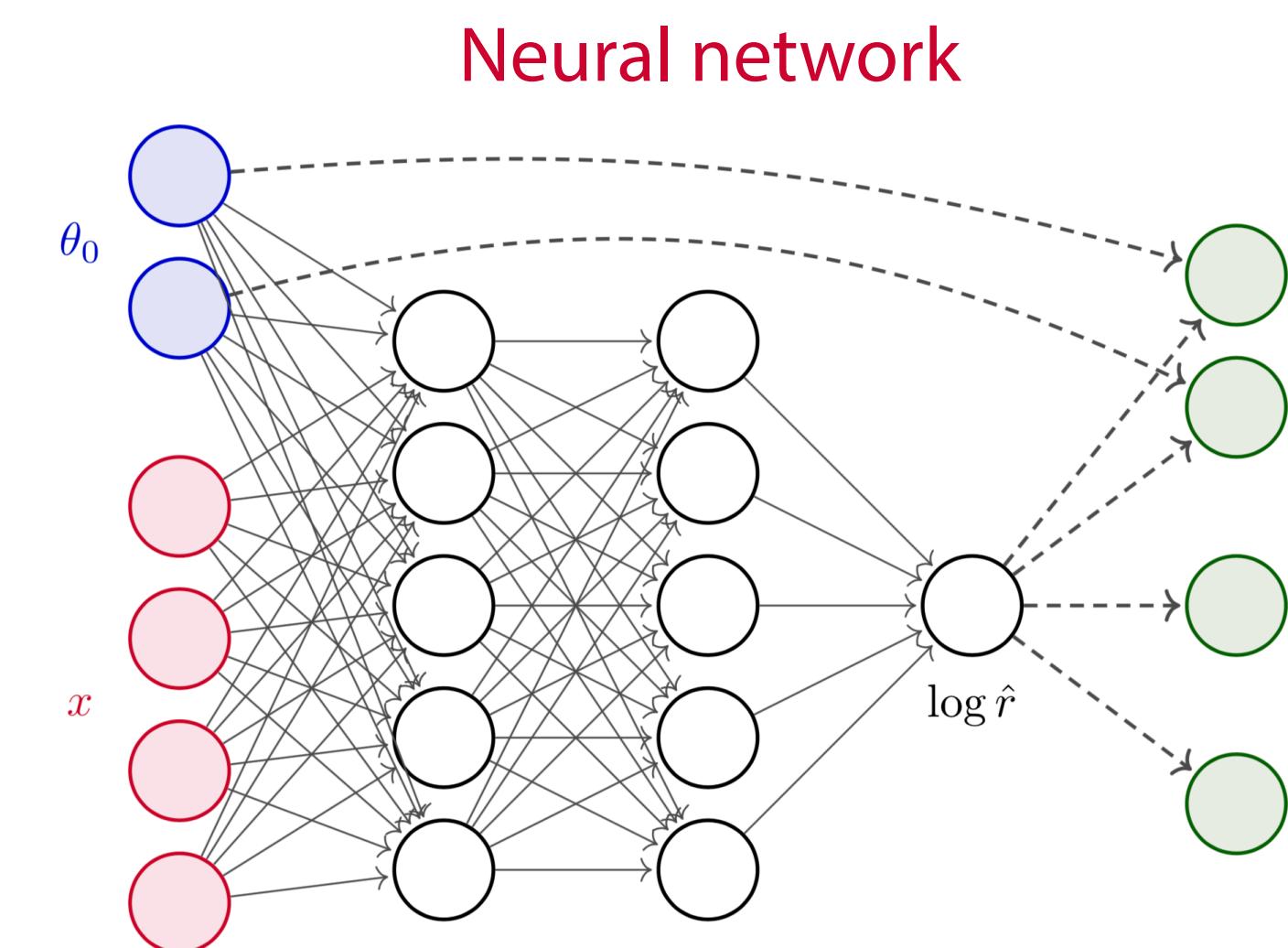
$$r(x|\theta_0, \theta_1) = \arg \min_{\hat{r}(x|\theta_0, \theta_1)} L_r[\hat{r}(x|\theta_0, \theta_1)]$$

Functional with integral

Machine learning = applied calculus of variations

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

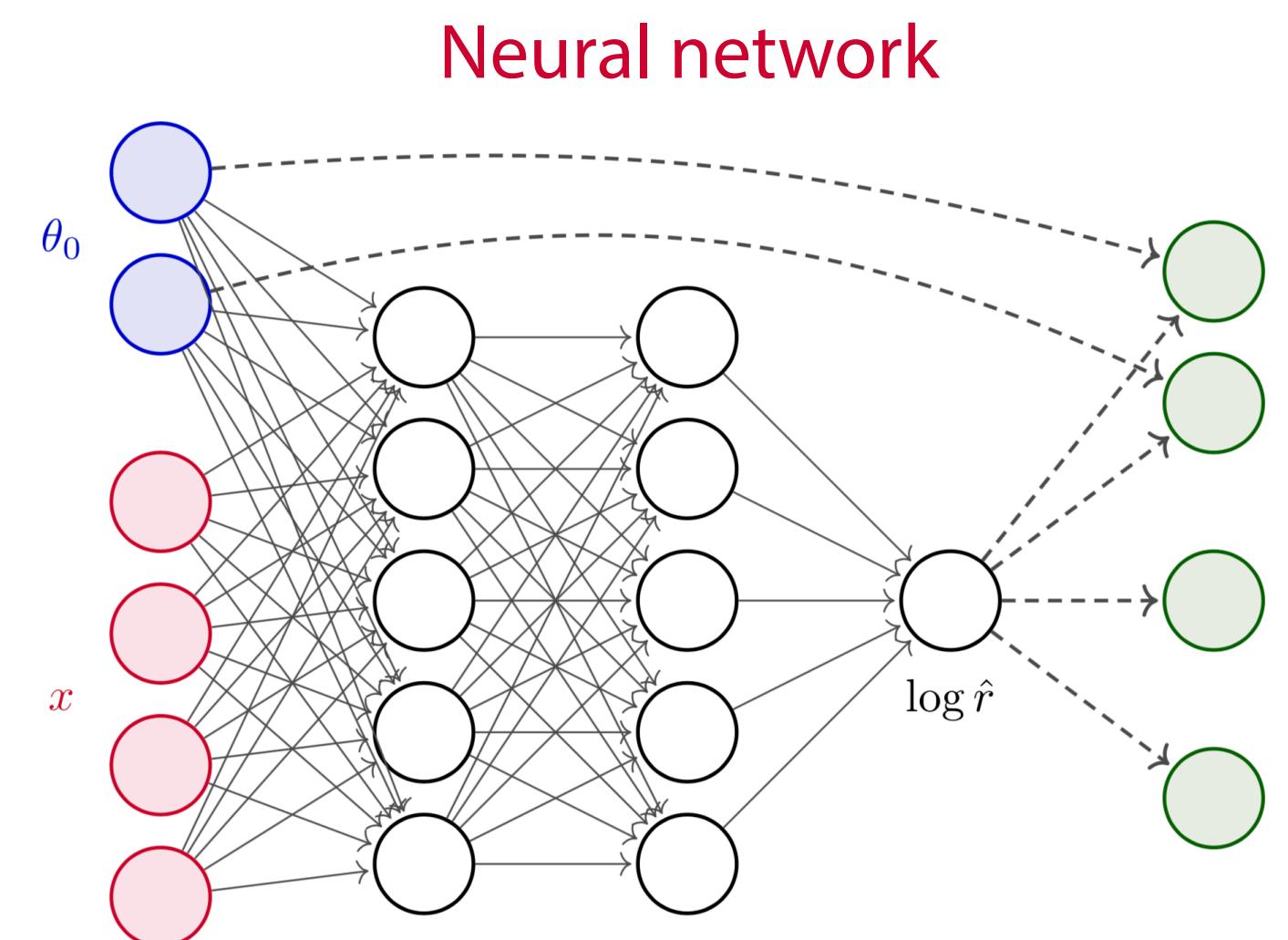
This is where machine learning comes in!



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!



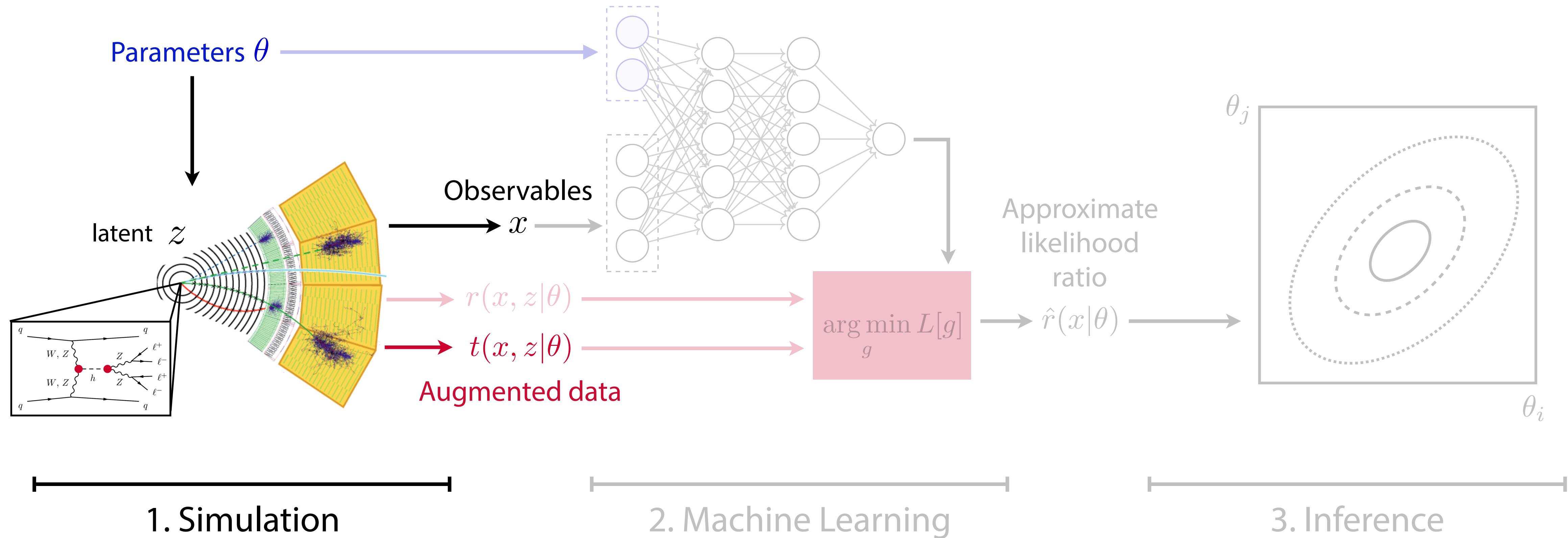
$$r(x|\theta_0, \theta_1) = \arg \min_{\hat{r}(x|\theta_0, \theta_1)} L_r[\hat{r}(x|\theta_0, \theta_1)]$$

Diagram annotations:

- Variational family: Points to the term $\arg \min_{\hat{r}(x|\theta_0, \theta_1)}$.
- Extremization: Points to the term $L_r[\hat{r}(x|\theta_0, \theta_1)]$.
- Functional with integral: Points to the term $L_r[\hat{r}(x|\theta_0, \theta_1)]$.
- Loss function with finite sum over samples: Points to the term $L_r[\hat{r}(x|\theta_0, \theta_1)]$.
- Stochastic gradient descent: Points to the term $\hat{r}(x|\theta_0, \theta_1)$.

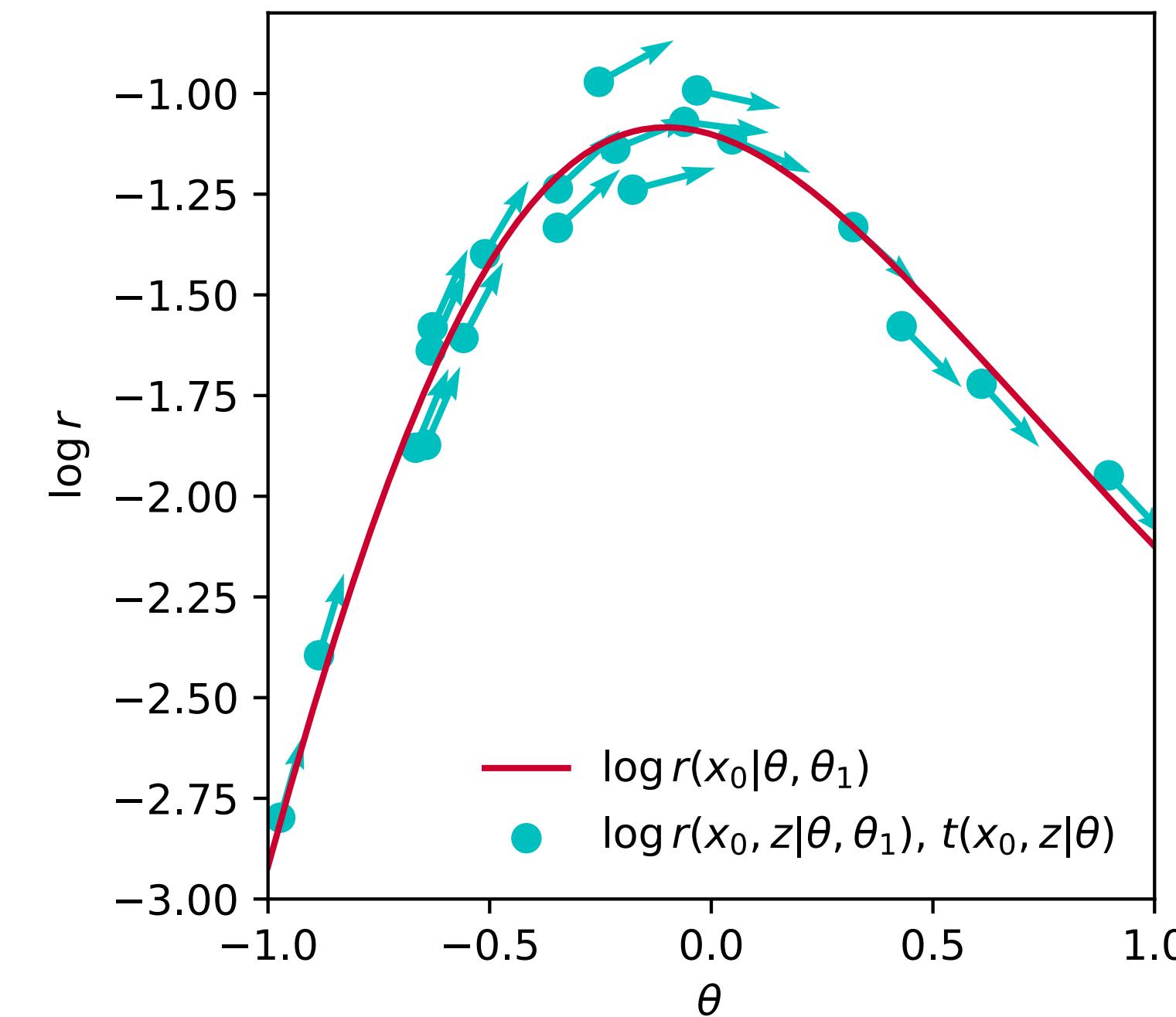
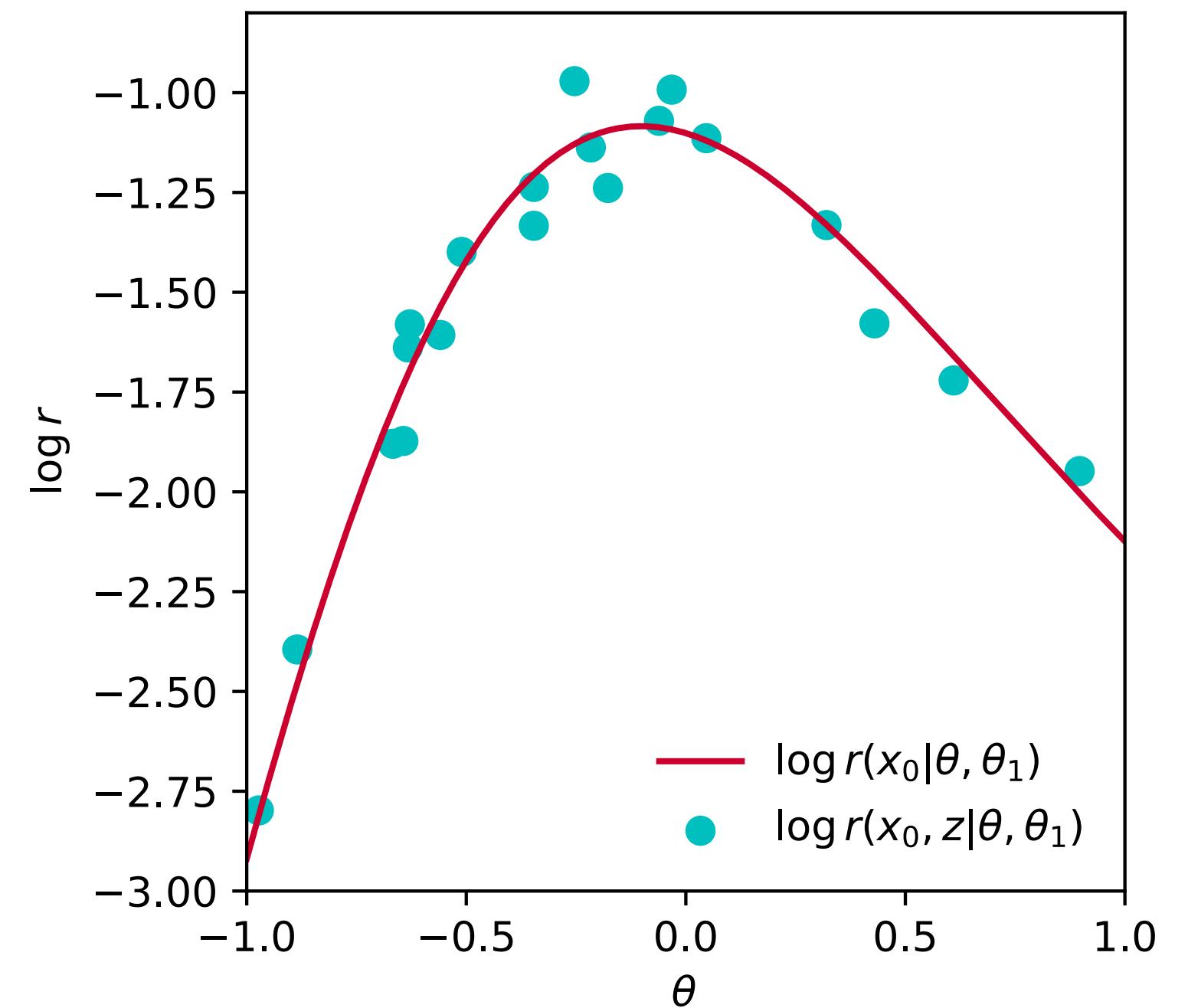
A sufficiently expressive neural network efficiently trained in this way with enough data will learn the likelihood ratio function $r(x|\theta_0, \theta_1)$!

Learning with Augmented Data



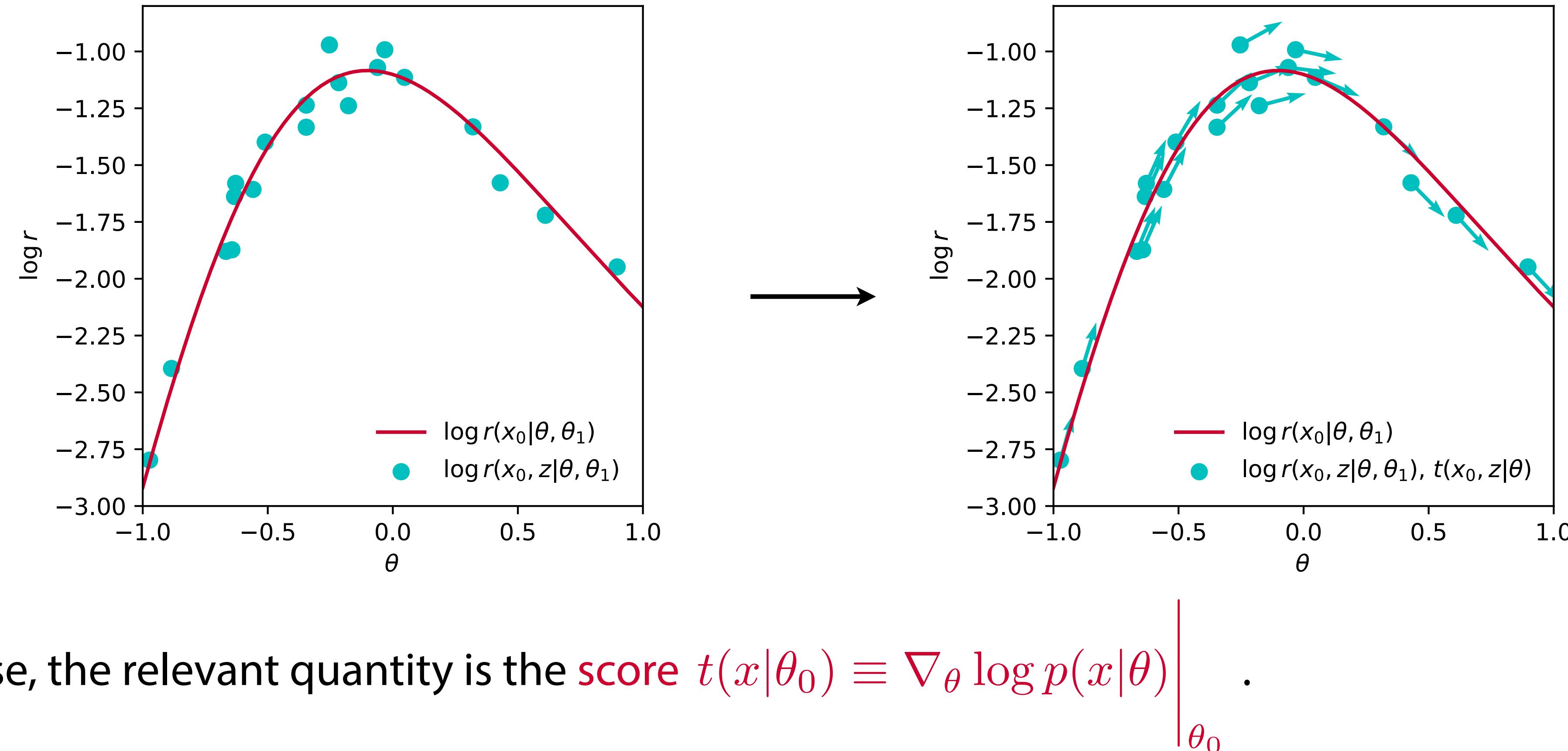
One more piece: the score

- Knowing derivative often helps fitting:



One more piece: the score

- Knowing derivative often helps fitting:



- In our case, the relevant quantity is the **score** $t(x|\theta_0) \equiv \nabla_{\theta} \log p(x|\theta) \bigg|_{\theta_0}$.
- The score itself is intractable. But...

Learning the score

Similar to the joint likelihood ratio, from the simulator we can extract the **joint score**

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



We want the **score**

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

Learning the score

Similar to the joint likelihood ratio, from the simulator we can extract the **joint score**

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



We want the **score**

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

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

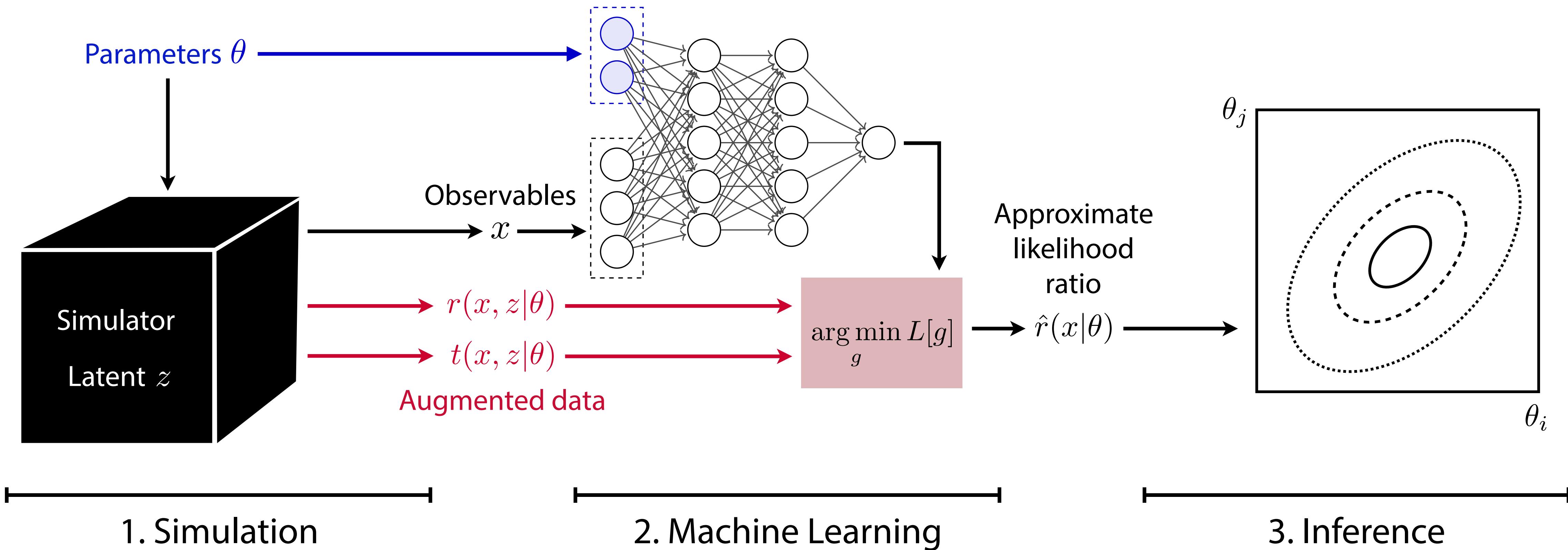
$$L_t[\hat{t}(x|\theta_0)] = \int dx \int dz \ p(x, z|\theta_0) \left[(\hat{t}(x|\theta_0) - t(x, z|\theta_0))^2 \right].$$

One can show it is minimized by

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

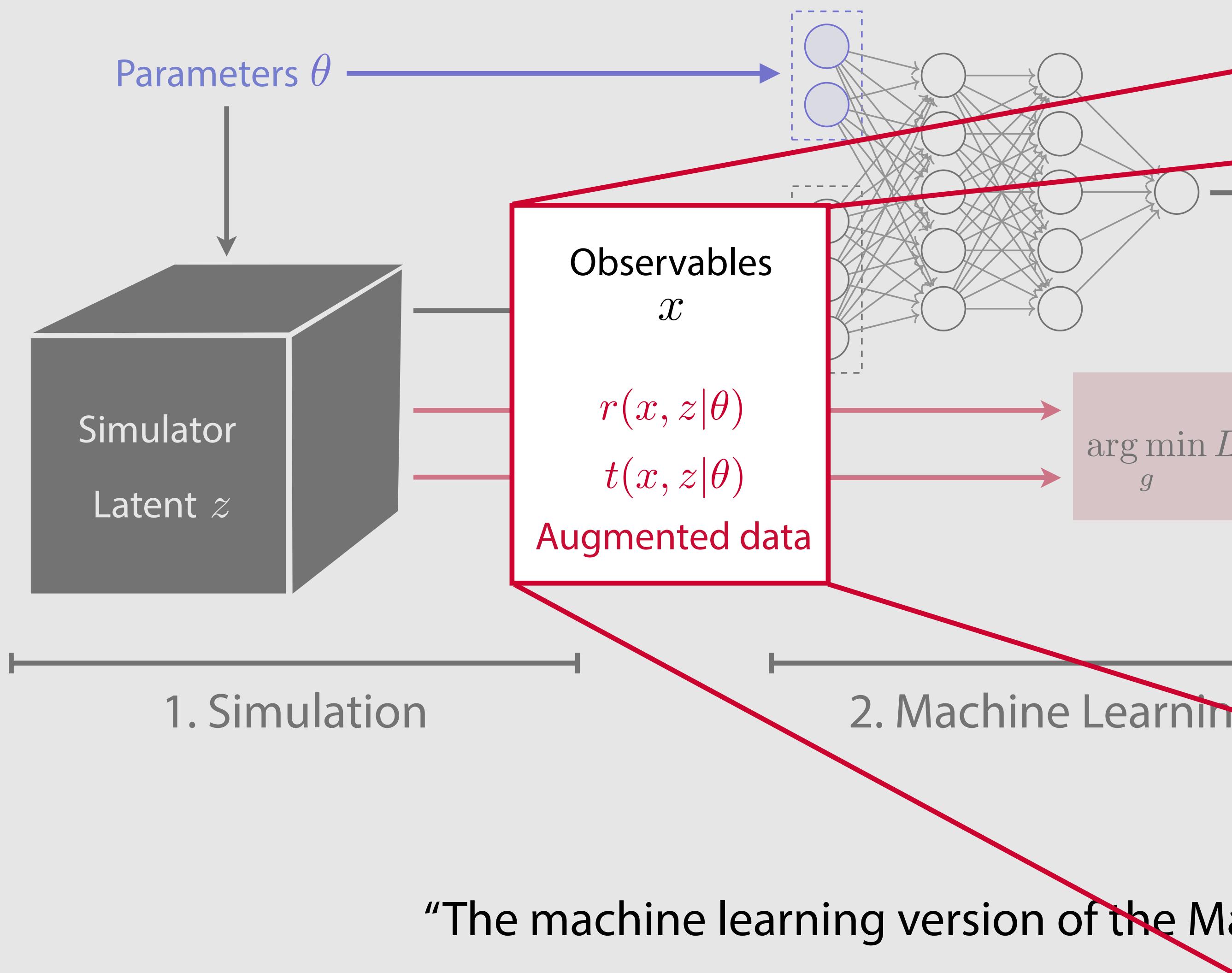
Again, we implement this minimization through machine learning.

Putting the pieces together: RASCAL & ALICES

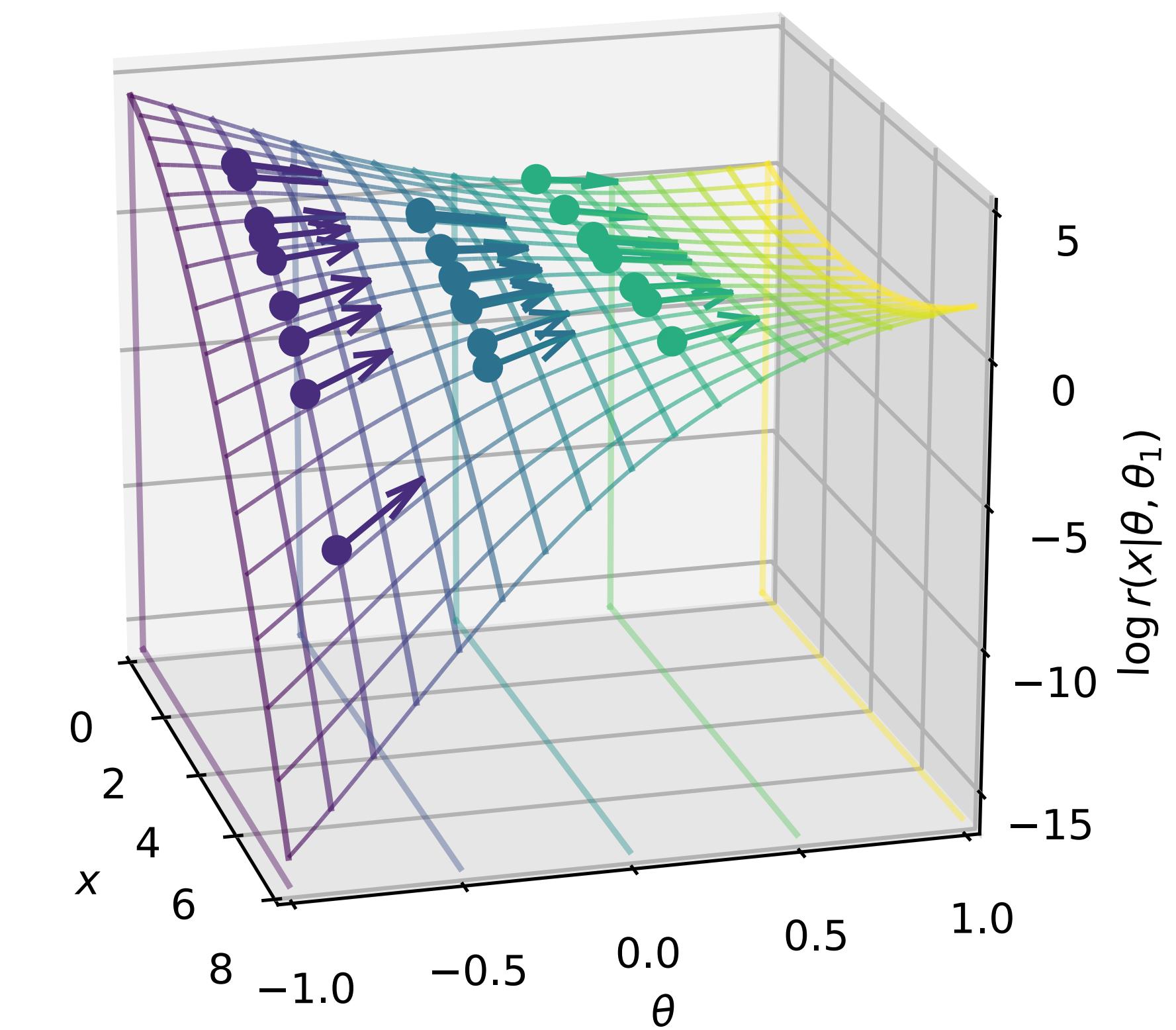


“The machine learning version of the Matrix Element Method”

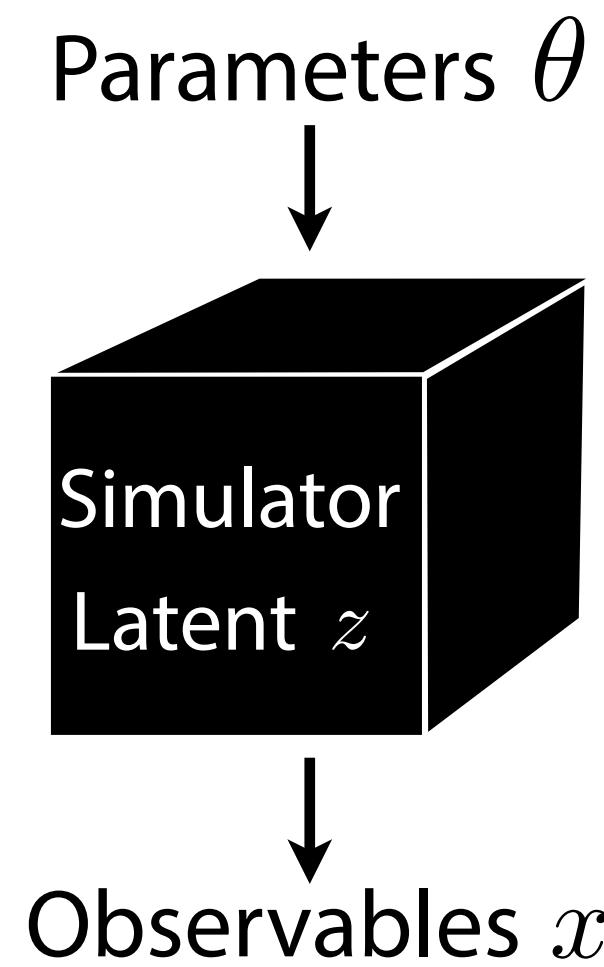
Putting the pieces together: RASCAL & ALICES



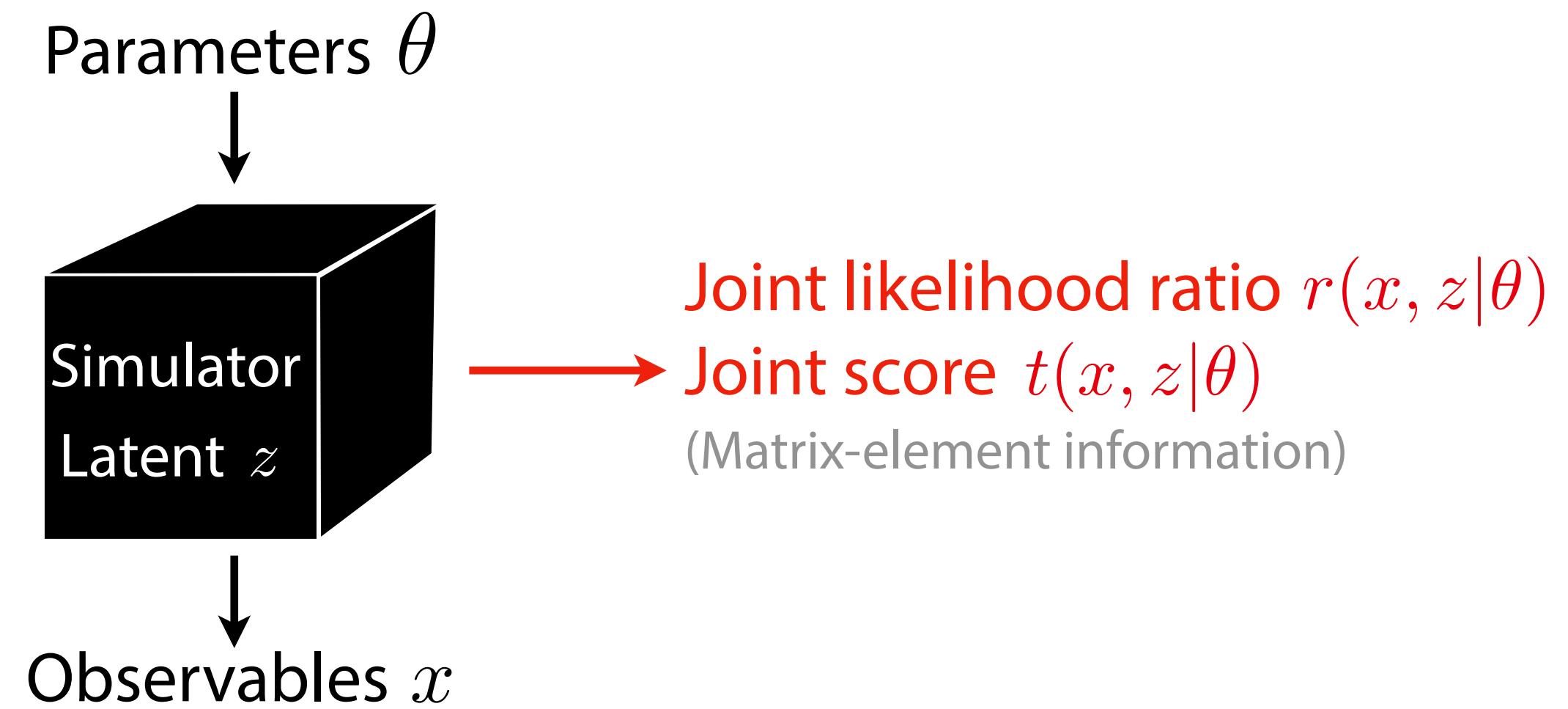
RASCAL & ALICES combines three orthogonal pieces of information



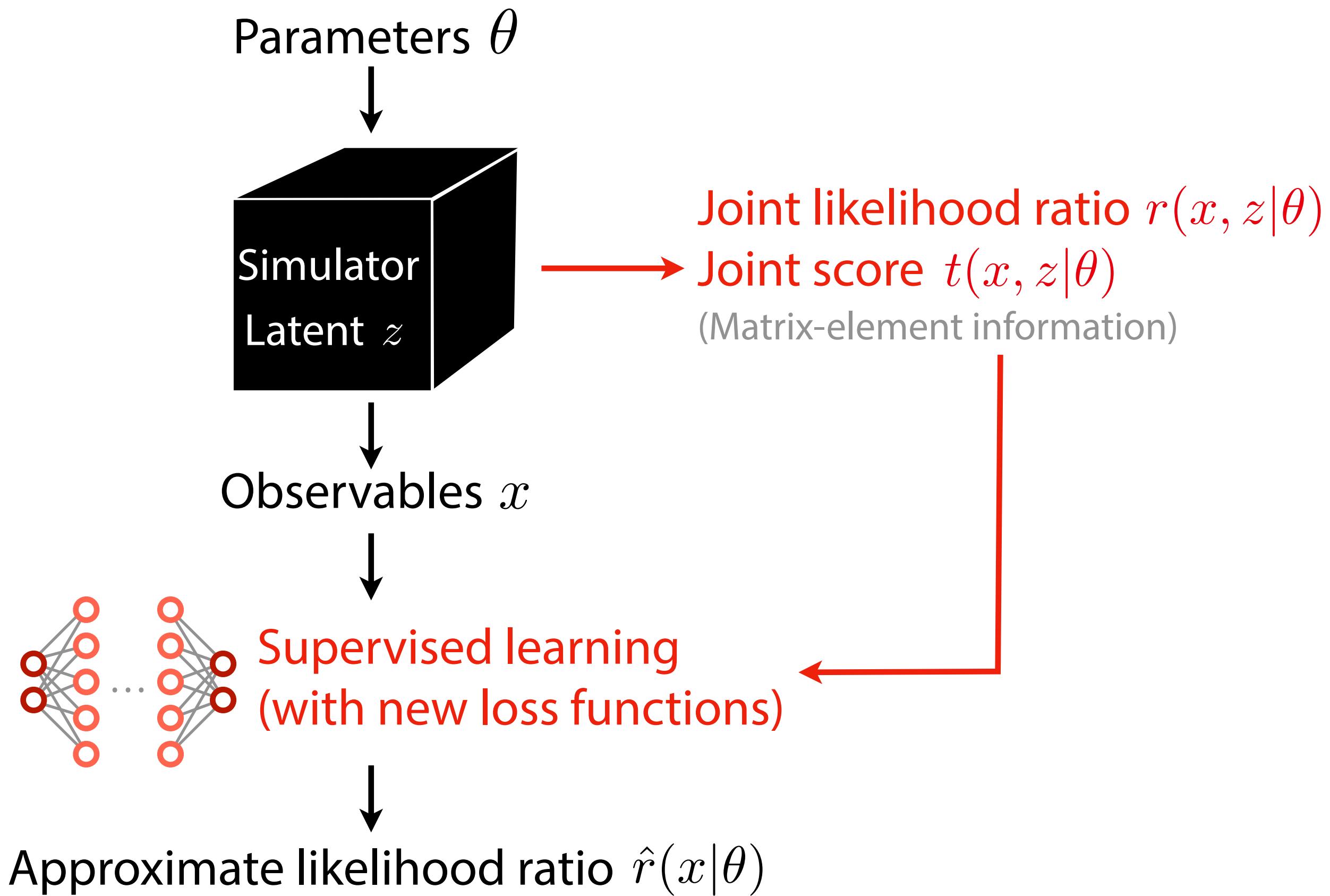
Mining gold: Recap



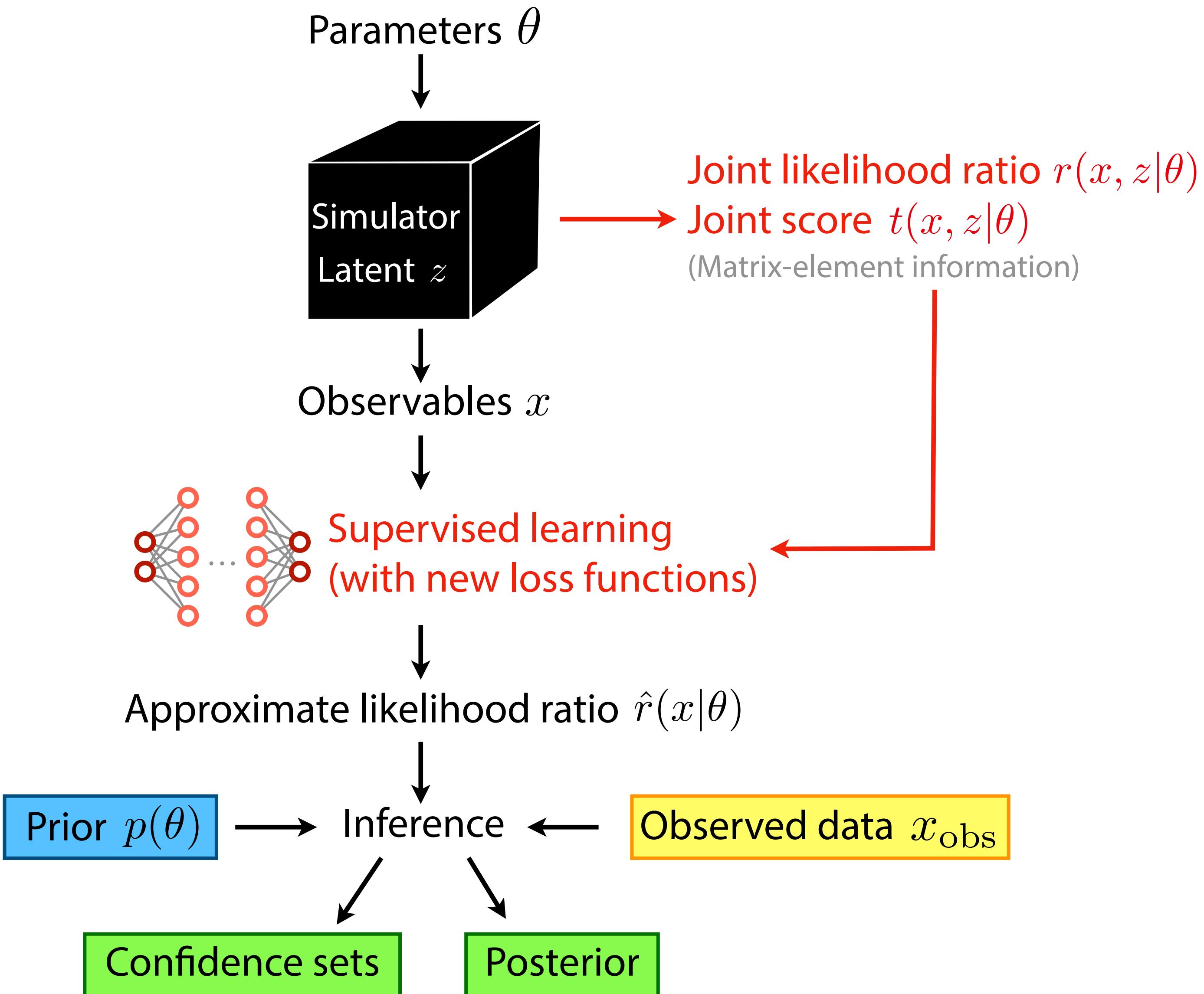
Mining gold: Recap



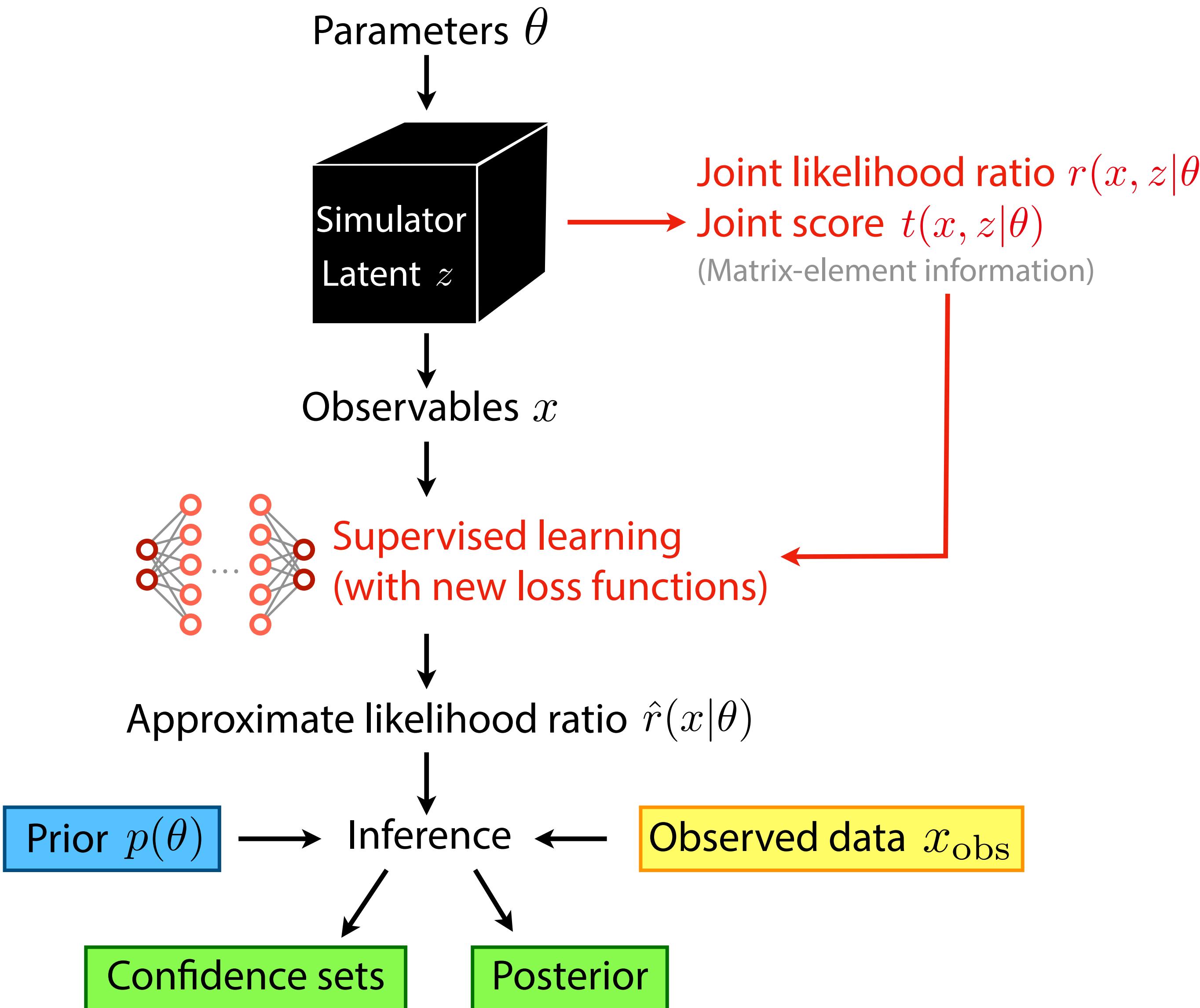
Mining gold: Recap



Mining gold: Recap



Mining gold: Summary



“The machine learning version of the Matrix Element Method”

- Scales to high-dimensional data (no summary statistics necessary)
- Uses matrix-element information to improve sample efficiency
- Neural networks learn effect of shower + detector (no transfer functions & works with ME+PS matching)
- Amortized (evaluation in μs)

We can also machine learn optimal observables.

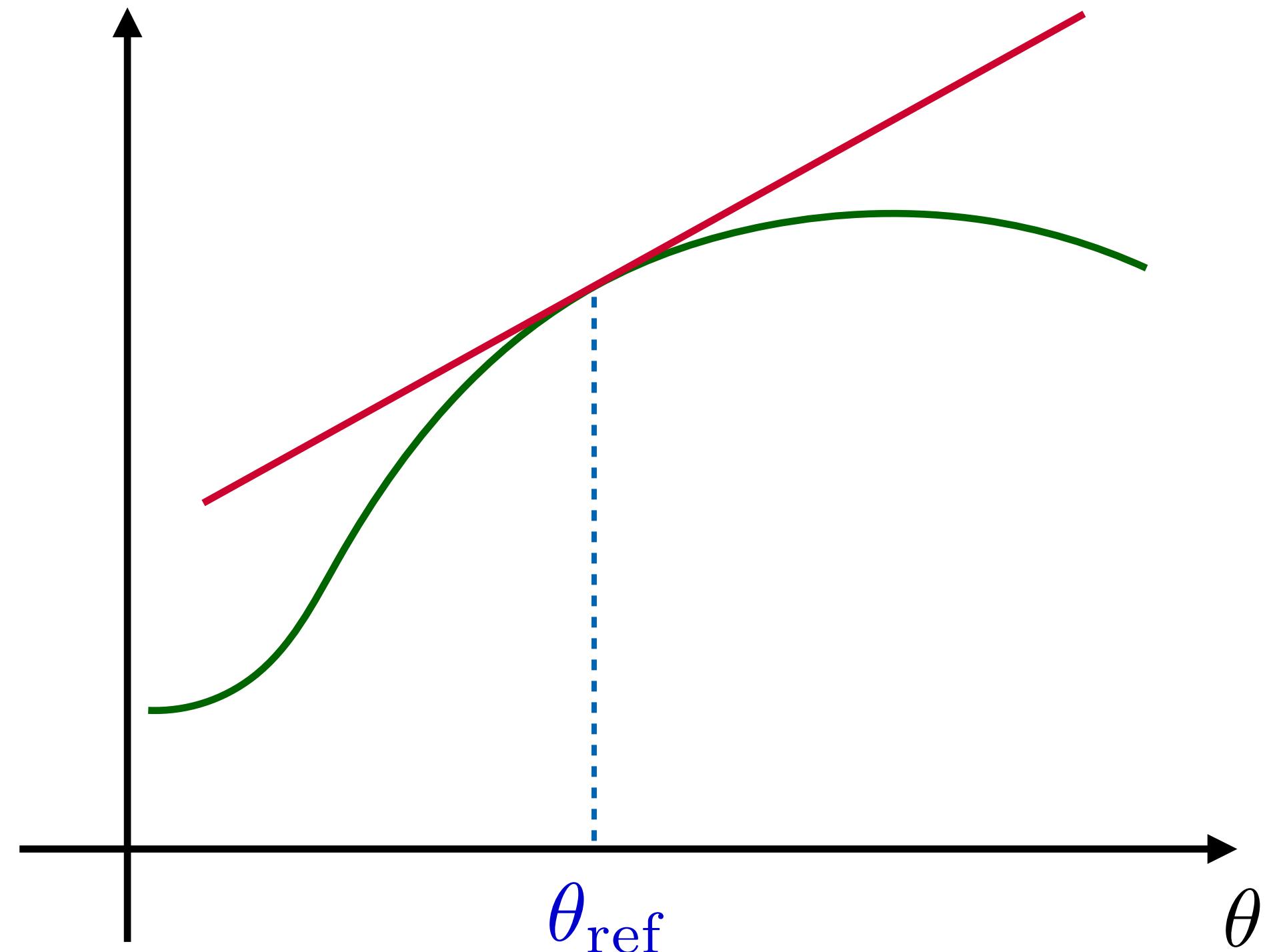
[JB, K. Cranmer, G. Louppe, J. Pavez 1805.00013, 1805.00020, 1805.12244]

The local model

[see also J. Alsing, B. Wandelt 1712.00012; J. Alsing, B. Wandelt, S. Freeney 1801.01497; P. de Castro, T. Dorigo 1806.04743; J. Alsing, B. Wandelt 1903.01473]

Taylor expansion of $\log p(x|\theta)$ around θ_{ref} :

$$\begin{aligned}\log p(x|\theta) &= \log p(x|\theta_{\text{ref}}) \\ &+ \underbrace{\nabla_{\theta} \log p(x|\theta) \Big|_{\theta_{\text{ref}}} \cdot (\theta - \theta_{\text{ref}})}_{\equiv t(x|\theta_{\text{ref}})} \\ &+ \mathcal{O}((\theta - \theta_{\text{ref}})^2)\end{aligned}$$



The local model

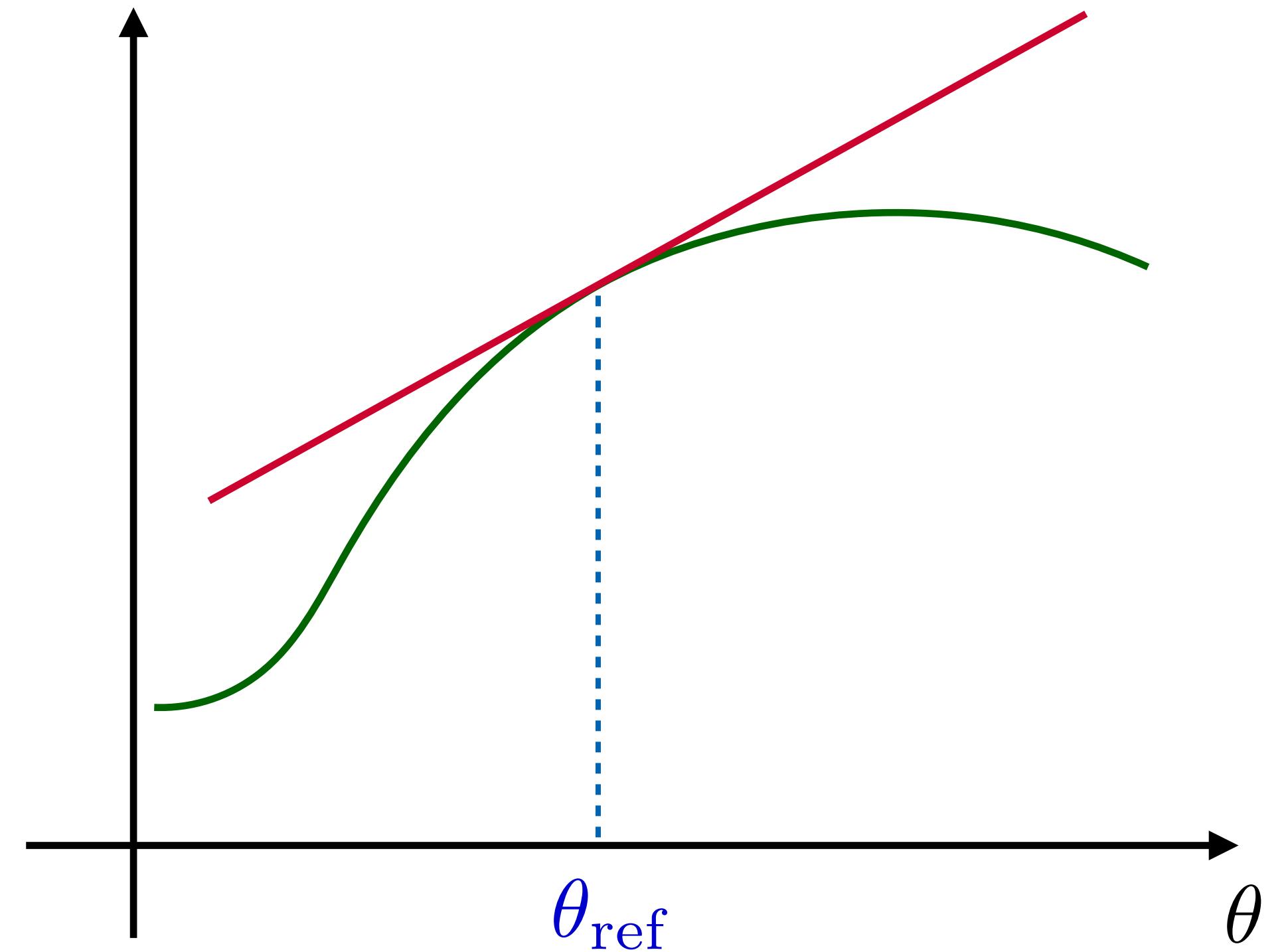
[see also J. Alsing, B. Wandelt 1712.00012; J. Alsing, B. Wandelt, S. Freeney 1801.01497; P. de Castro, T. Dorigo 1806.04743; J. Alsing, B. Wandelt 1903.01473]

Taylor expansion of $\log p(x|\theta)$ around θ_{ref} :

$$\begin{aligned}\log p(x|\theta) &= \log p(x|\theta_{\text{ref}}) \\ &+ \underbrace{\nabla_{\theta} \log p(x|\theta) \Big|_{\theta_{\text{ref}}} \cdot (\theta - \theta_{\text{ref}})}_{\equiv t(x|\theta_{\text{ref}})} \\ &+ \mathcal{O}((\theta - \theta_{\text{ref}})^2)\end{aligned}$$

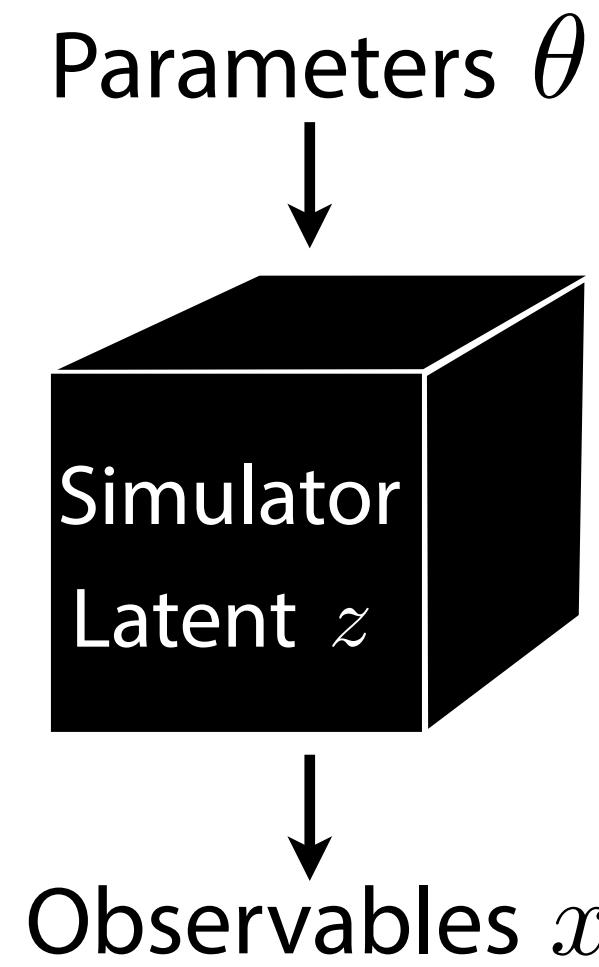
In the neighborhood of θ_{ref} (e.g. close to the SM):

- the **score vector** $t(x|\theta_{\text{ref}})$ components are sufficient statistics
- knowing $t(x|\theta_{\text{ref}})$ is just as powerful as knowing the full function $\log p(x|\theta)$
- $t(x|\theta_{\text{ref}})$ are the most powerful observables

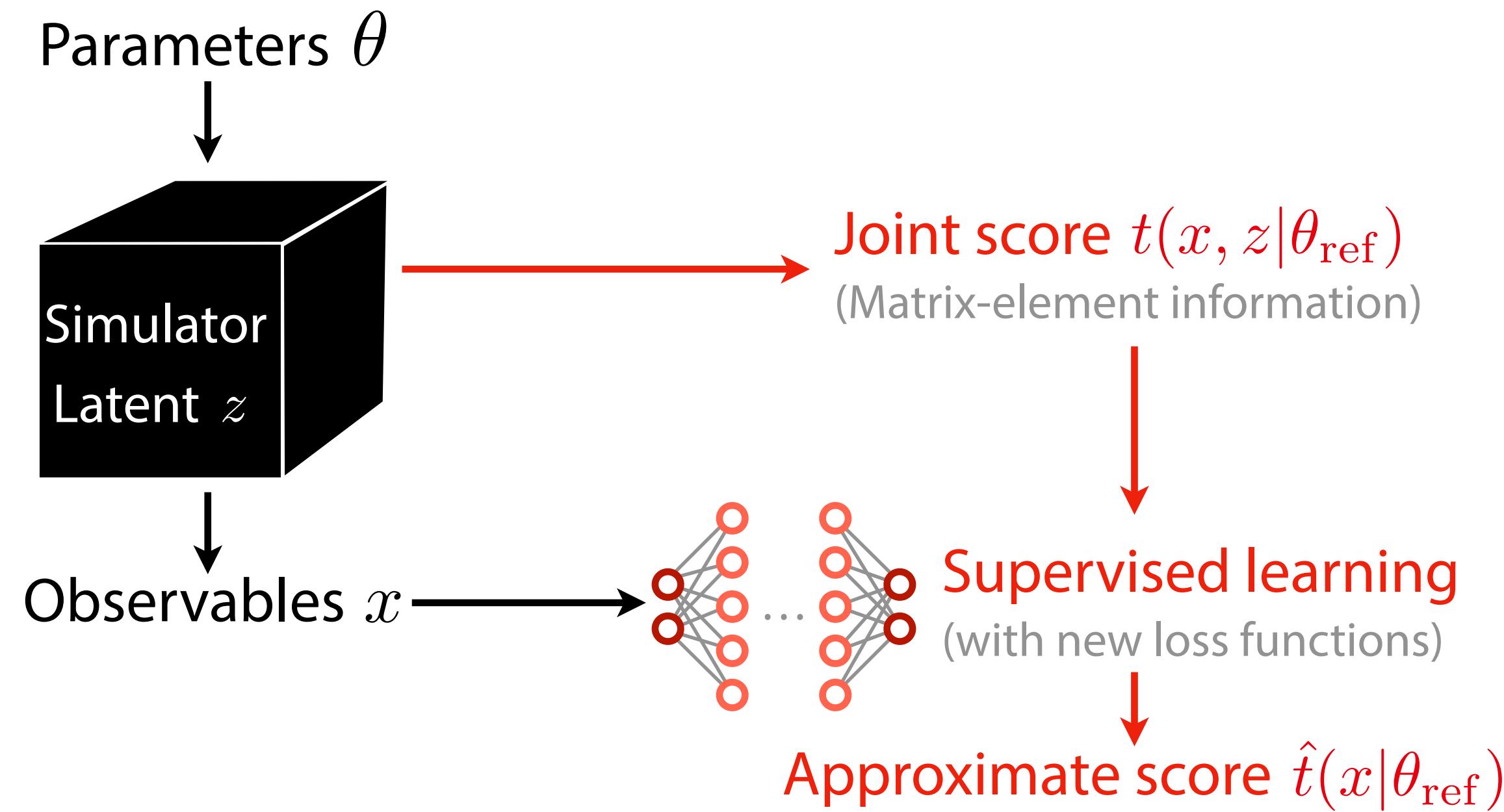


The score itself is intractable. But we can use the same trick as for the likelihood ratio!

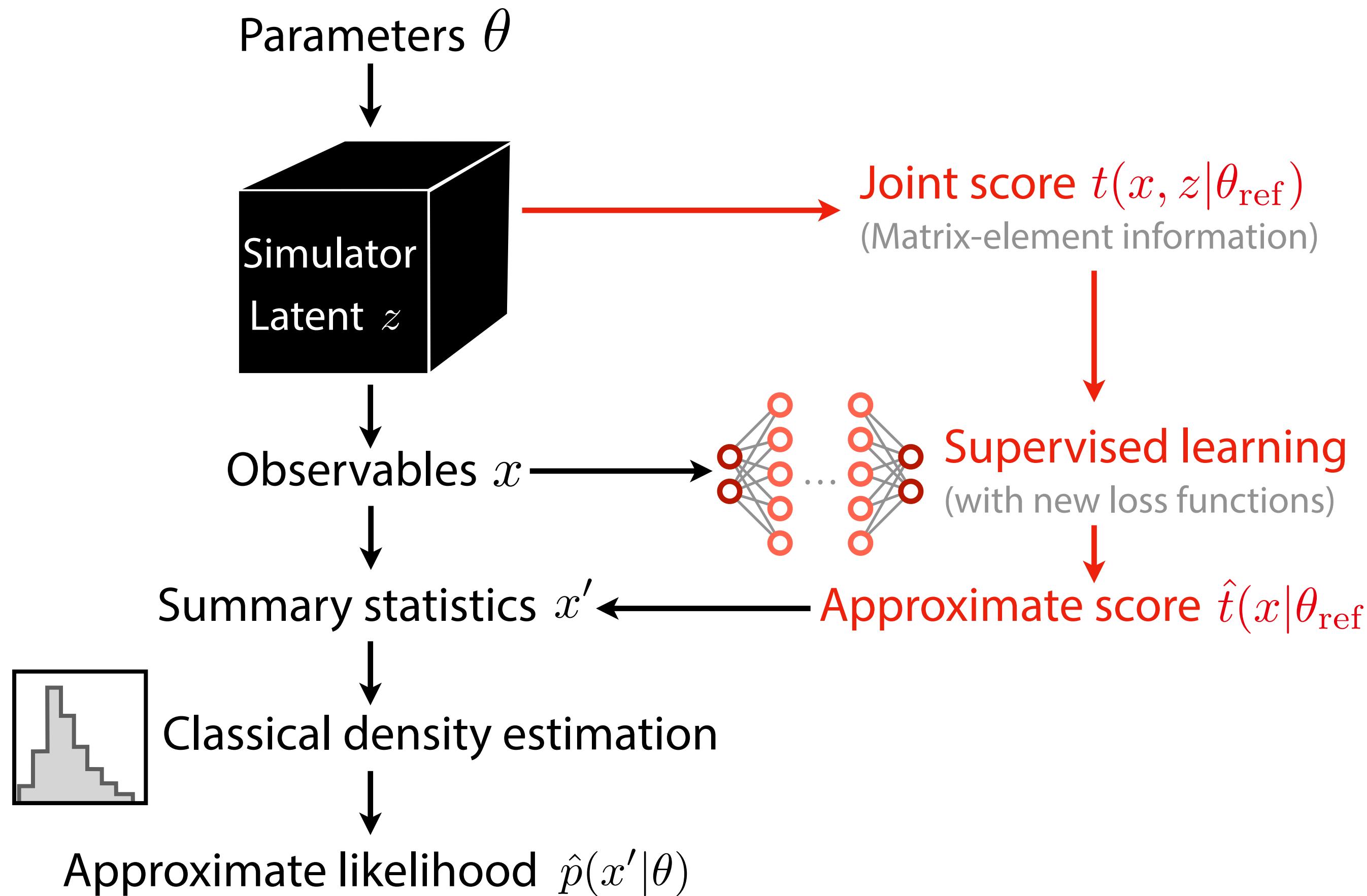
Neural optimal observables (SALLY)



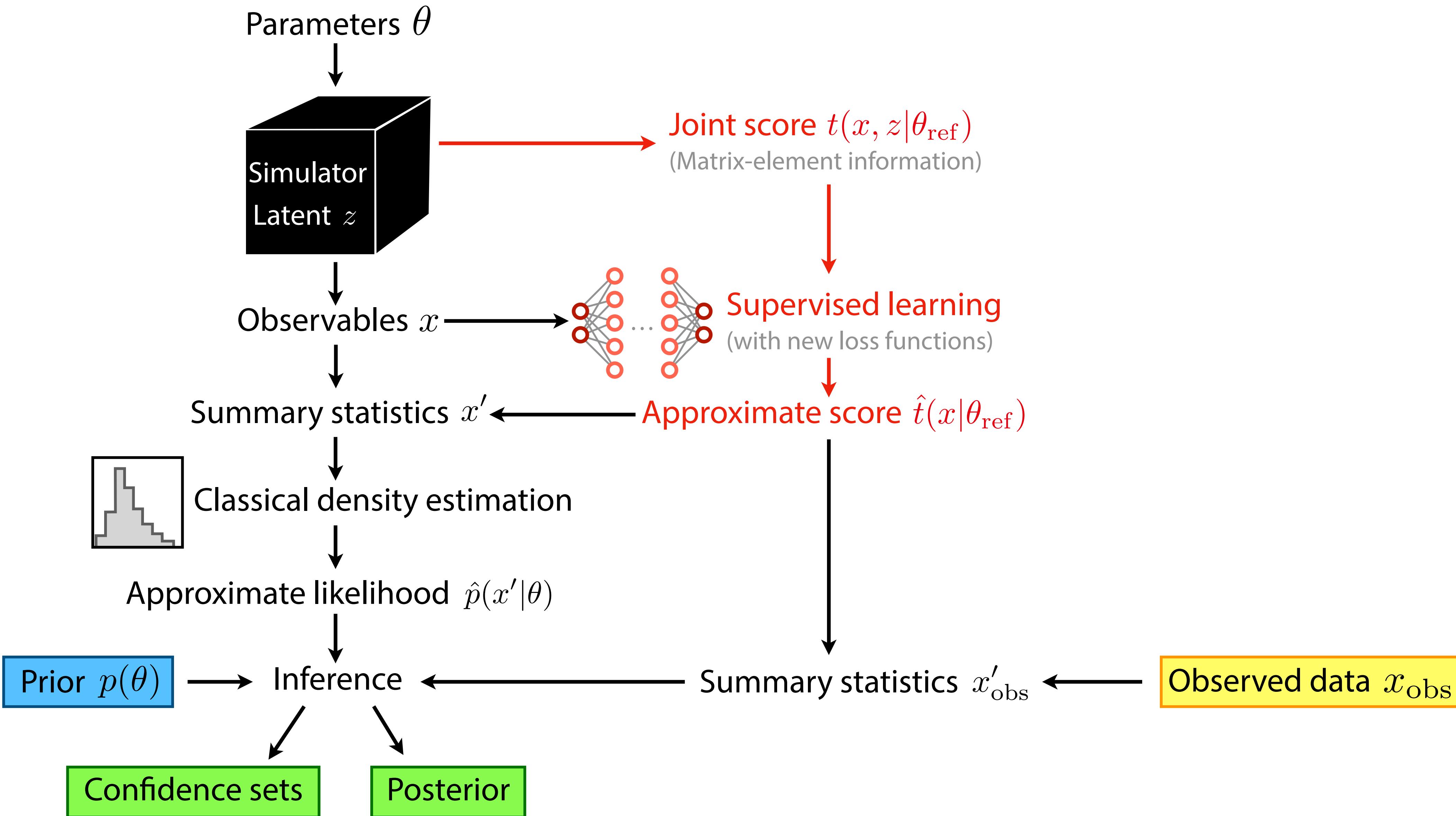
Neural optimal observables (SALLY)



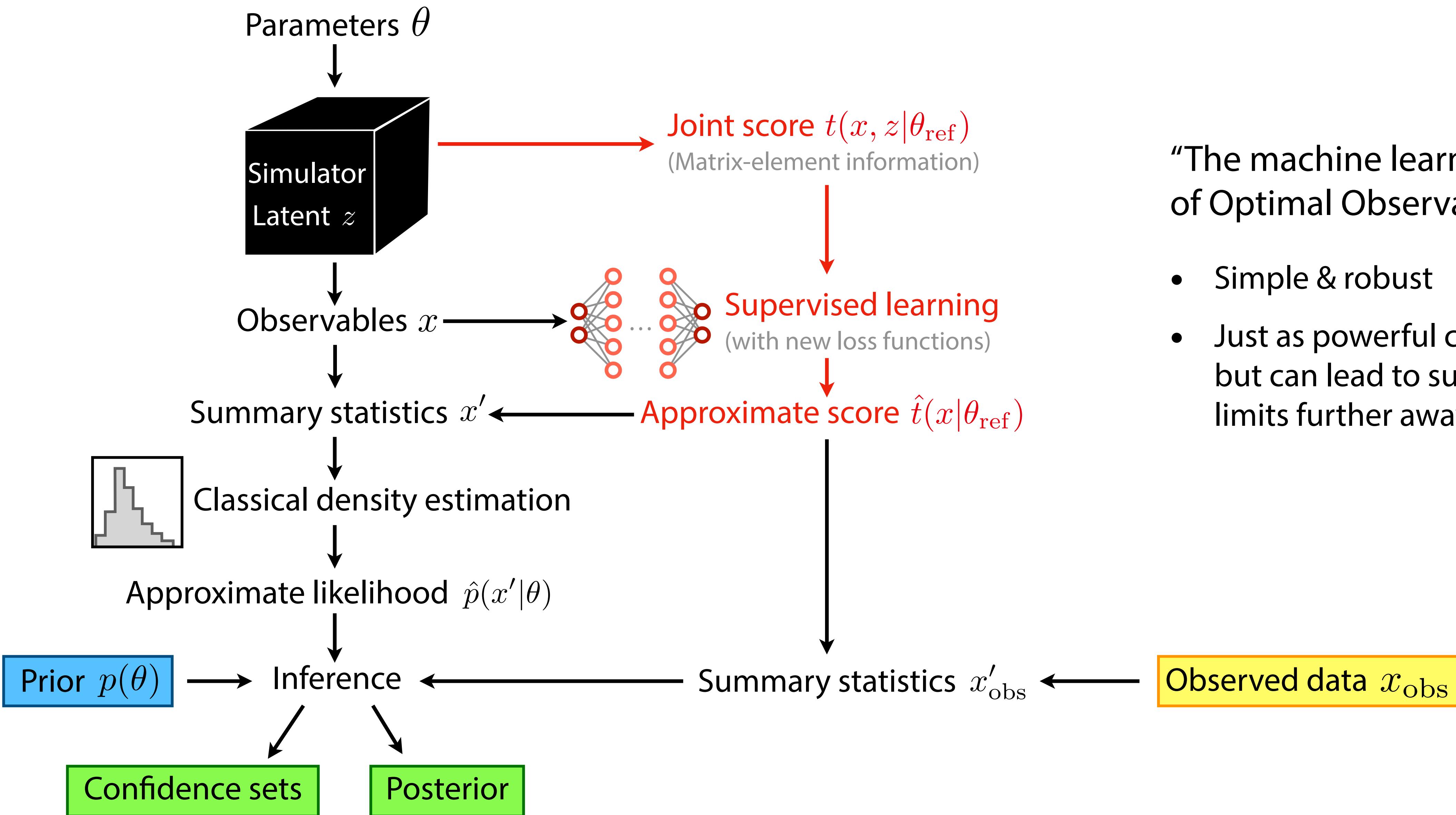
Neural optimal observables (SALLY)



Neural optimal observables (SALLY)



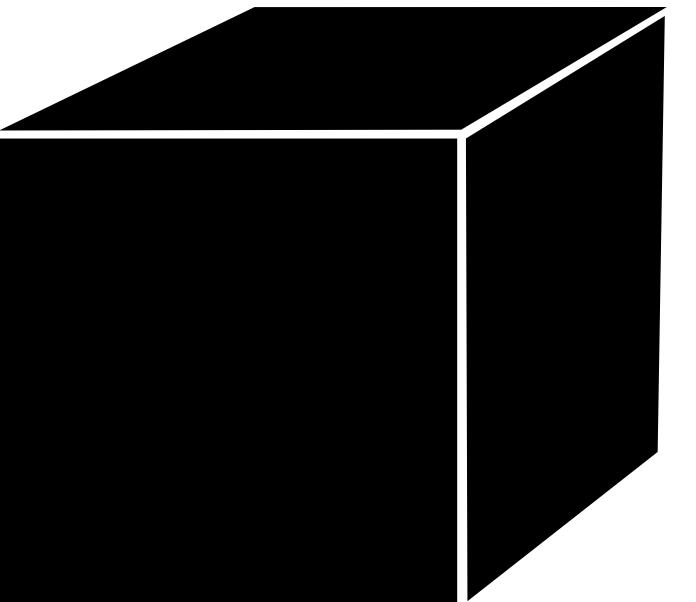
Neural optimal observables (SALLY)



An incomplete wrap-up of simulation-based inference methods

Method	Approximations	Upfront cost	Eval
Summary statistics:			
Likelihood for summary stats (standard histograms)	Reduction to summary stats	Fast	Fast
Approximate Bayesian Computation	Reduction to summary stats	Depends	Depends
Matrix elements:			
Matrix Element Method	Transfer fns	Fast	Slow
Optimal Observables	Transfer fns, optimal only locally	Fast	Slow
Neural networks:			
Neural likelihood	NN	Needs many samples	Fast
Neural posterior	NN	Needs many samples	Fast
Neural likelihood ratio	NN	Needs many samples	Fast
Neural networks + matrix elements:			
Neural likelihood (ratio) + gold mining (RASCAL etc)	NN	Needs less samples	Fast
Neural optimal observables (SALLY)	NN, optimal only locally	Needs less samples	Fast

Systematics



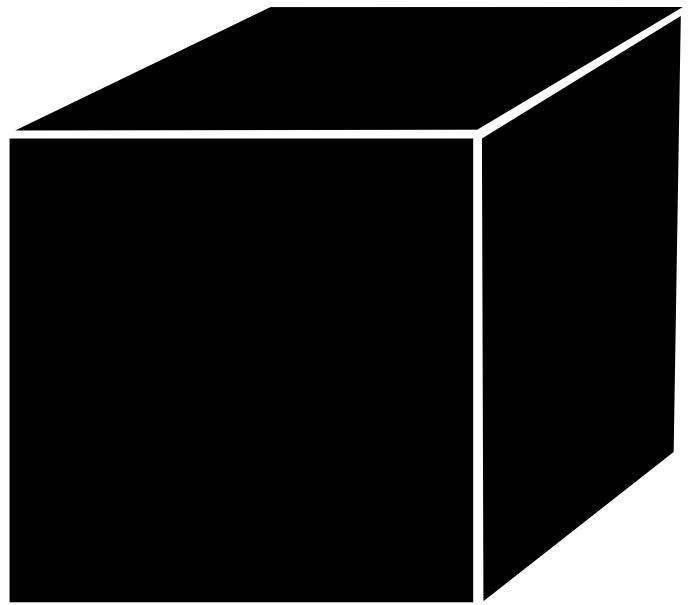
Don't fully trust the simulator?

- Nuisance parameters to model systematic uncertainties
- Methods learn dependence both on parameters of interest and nuisance parameters. Then we can construct profile likelihood and “nuisance-hardened” score

[J. Alsing, B. Wandelt 1903.01473;
see also P. de Castro, T. Dorigo 1806.04743]

- Alternatively: Robustness to nuisance with adversarial training
[G. Louppe, M. Kagan, K. Cranmer 1611.01046]

Systematics



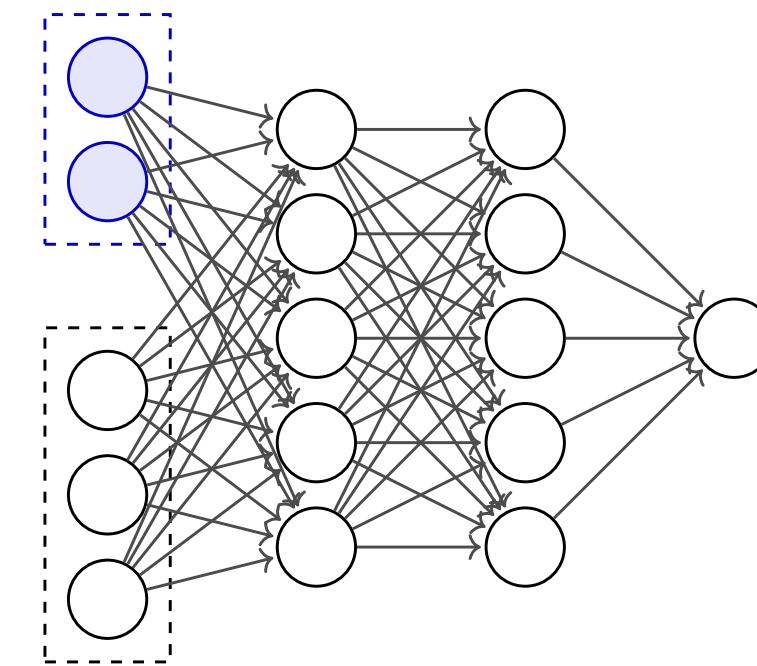
Don't fully trust the simulator?

- Nuisance parameters to model systematic uncertainties
- Methods learn dependence both on parameters of interest and nuisance parameters. Then we can construct profile likelihood and “nuisance-hardened” score

[J. Alsing, B. Wandelt 1903.01473;
see also P. de Castro, T. Dorigo 1806.04743]

- Alternatively: Robustness to nuisance with adversarial training

[G. Louppe, M. Kagan, K. Cranmer 1611.01046]



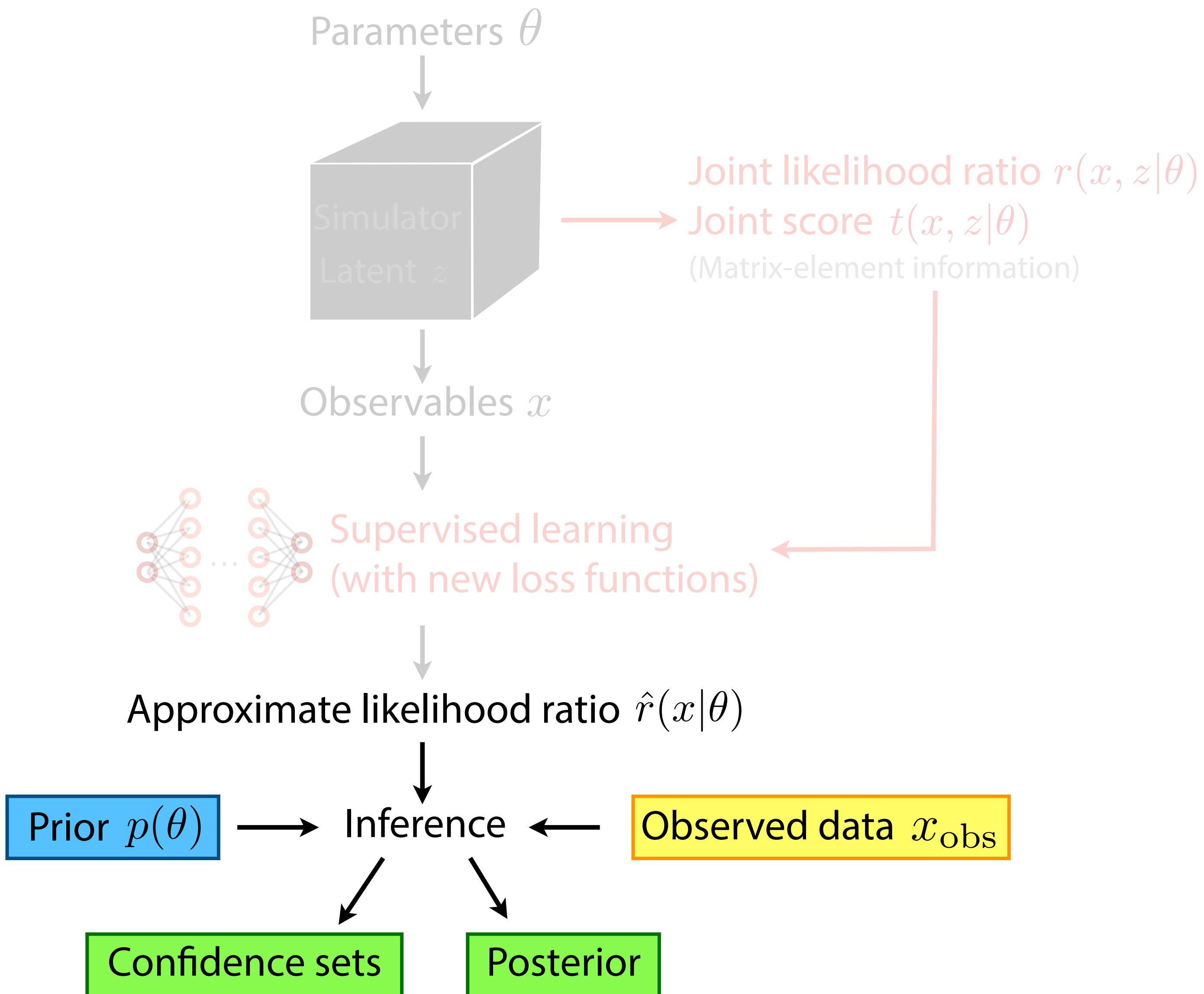
Don't blindly trust the neural network?

- Sanity checks: expectation values, “critic” tests
- Calibration / Neyman construction with toys
(badly trained network can lead to suboptimal limits, but not to wrong limits)

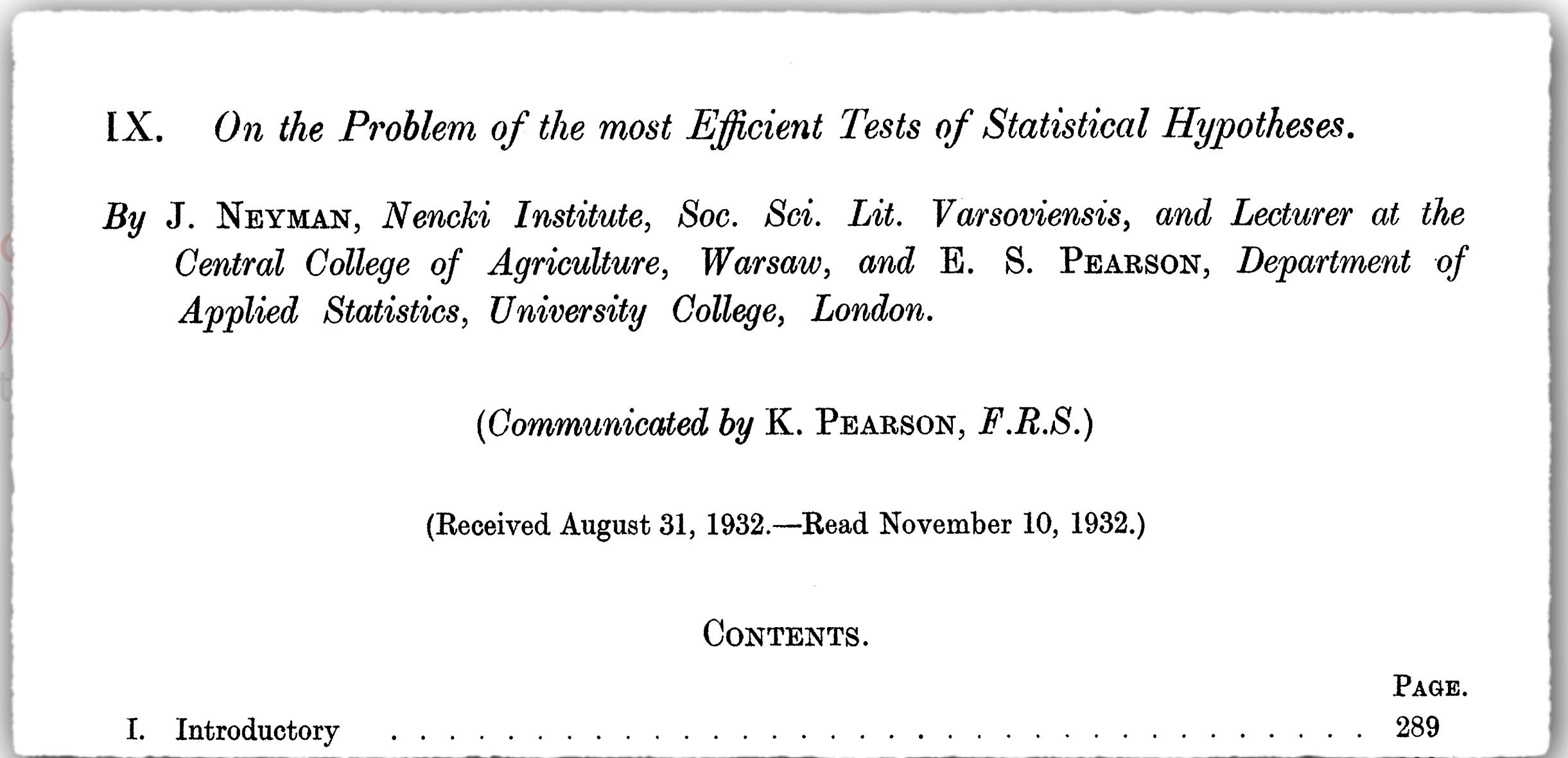
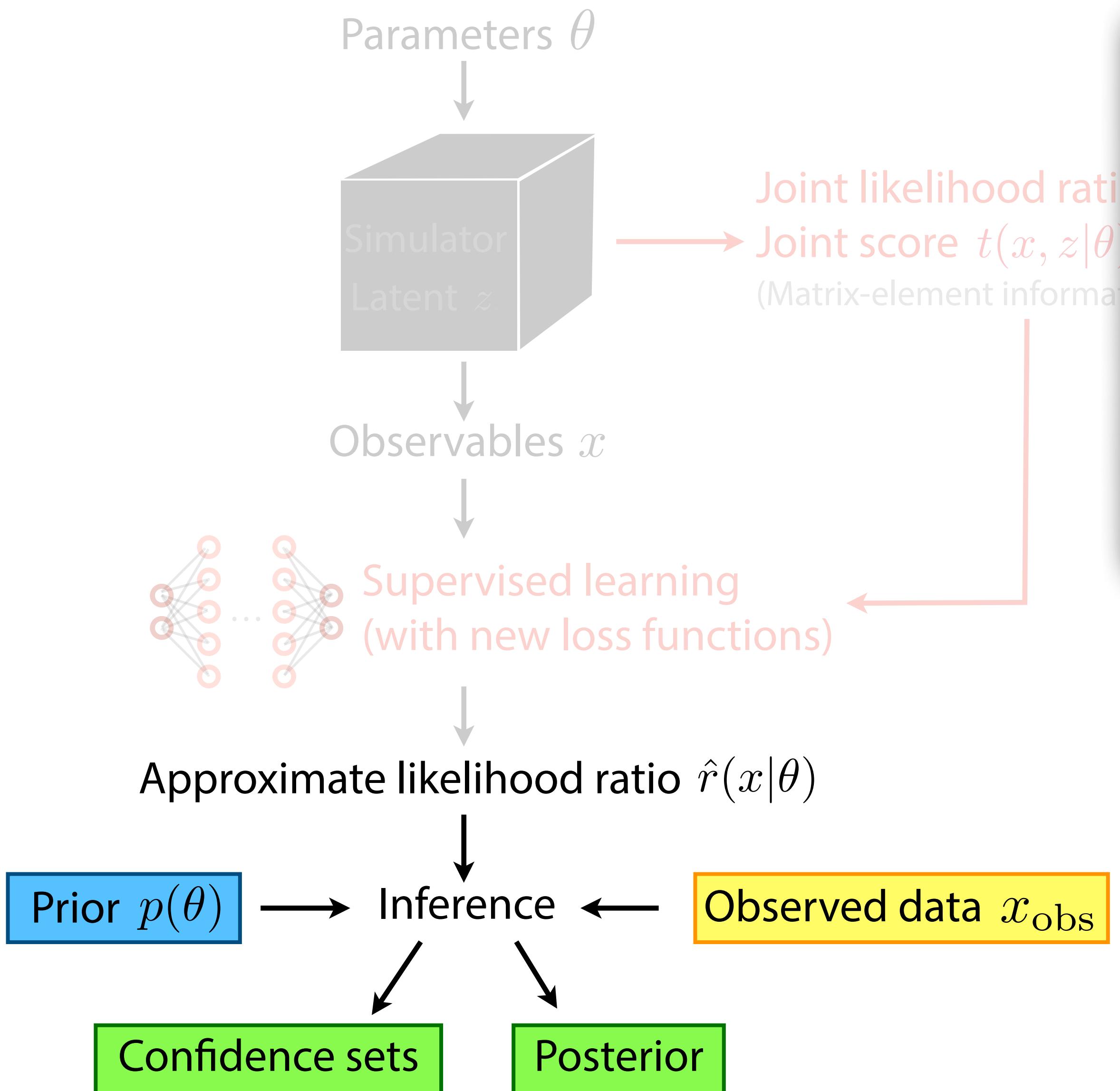
Inference, Limits, & Fisher Information Geometry

[JB, K. Cranmer, G. Louppe, J. Pavez 1805.00013, 1805.00020, 1805.12244]

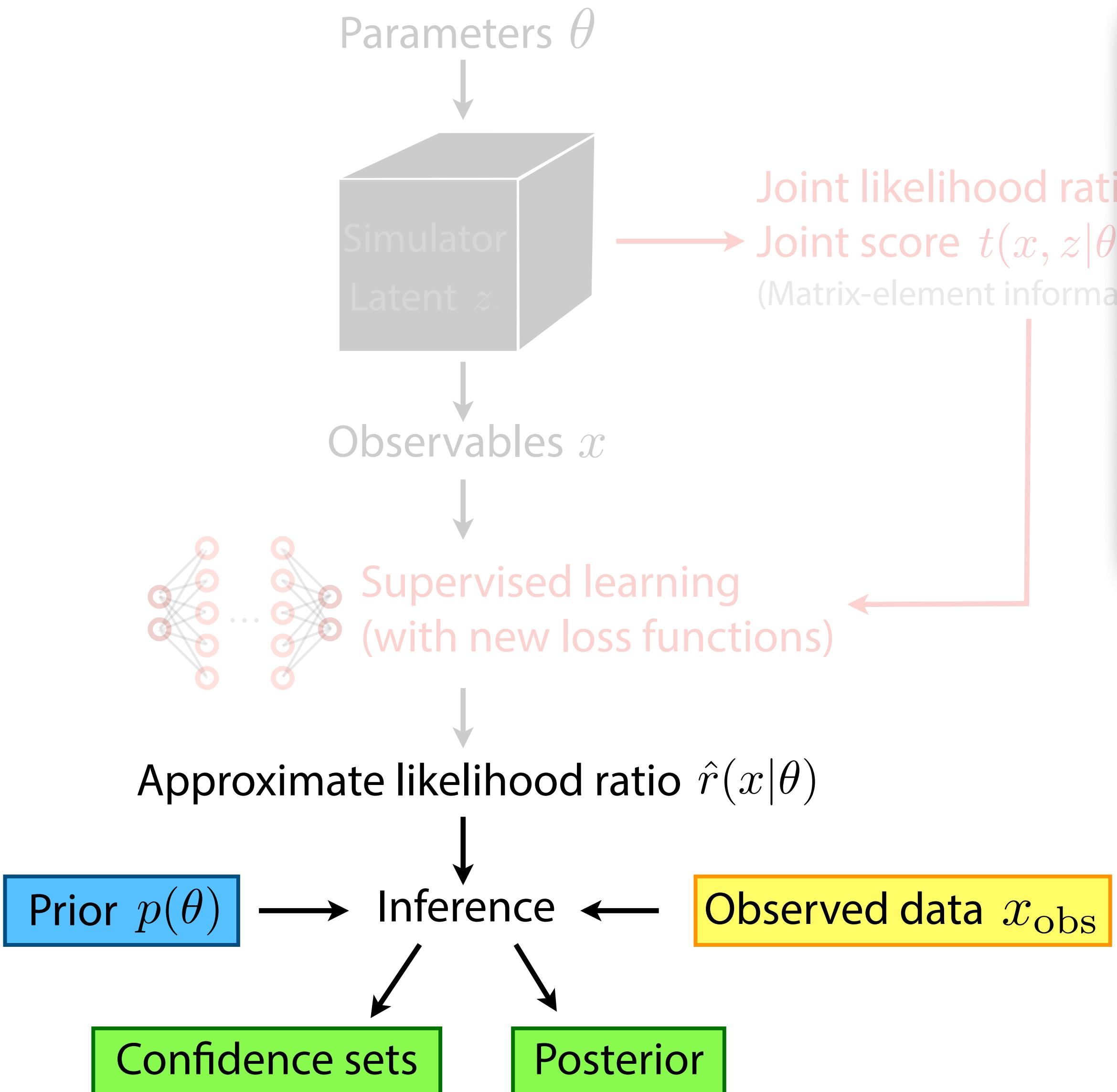
Step 3: Inference



Step 3: Inference



Step 3: Inference



IX. On the Problem of the most Efficient Tests of Statistical Hypotheses.

By J. NEYMAN, *Nencki Institute, Soc. Sci. Lit. Varsoviensis, and Lecturer at the Central College of Agriculture, Warsaw*, and E. S. PEARSON, *Department of Applied Statistics, University College, London*.

(Communicated by K. PEARSON, F.R.S.)

(Received August 31, 1932.—Read November 10, 1932.)

Eur. Phys. J. C (2011) 71: 1554
DOI 10.1140/epjc/s10052-011-1554-0

THE EUROPEAN
PHYSICAL JOURNAL C

Special Article - Tools for Experiment and Theory

Asymptotic formulae for likelihood-based tests of new physics

Glen Cowan¹, Kyle Cranmer², Eilam Gross³, Ofer Vitells^{3,a}

¹Physics Department, Royal Holloway, University of London, Egham TW20 0EX, UK

²Physics Department, New York University, New York, NY 10003, USA

³Weizmann Institute of Science, Rehovot 76100, Israel

Received: 15 October 2010 / Revised: 6 January 2011 / Published online: 9 February 2011

© The Author(s) 2011. This article is published with open access at Springerlink.com

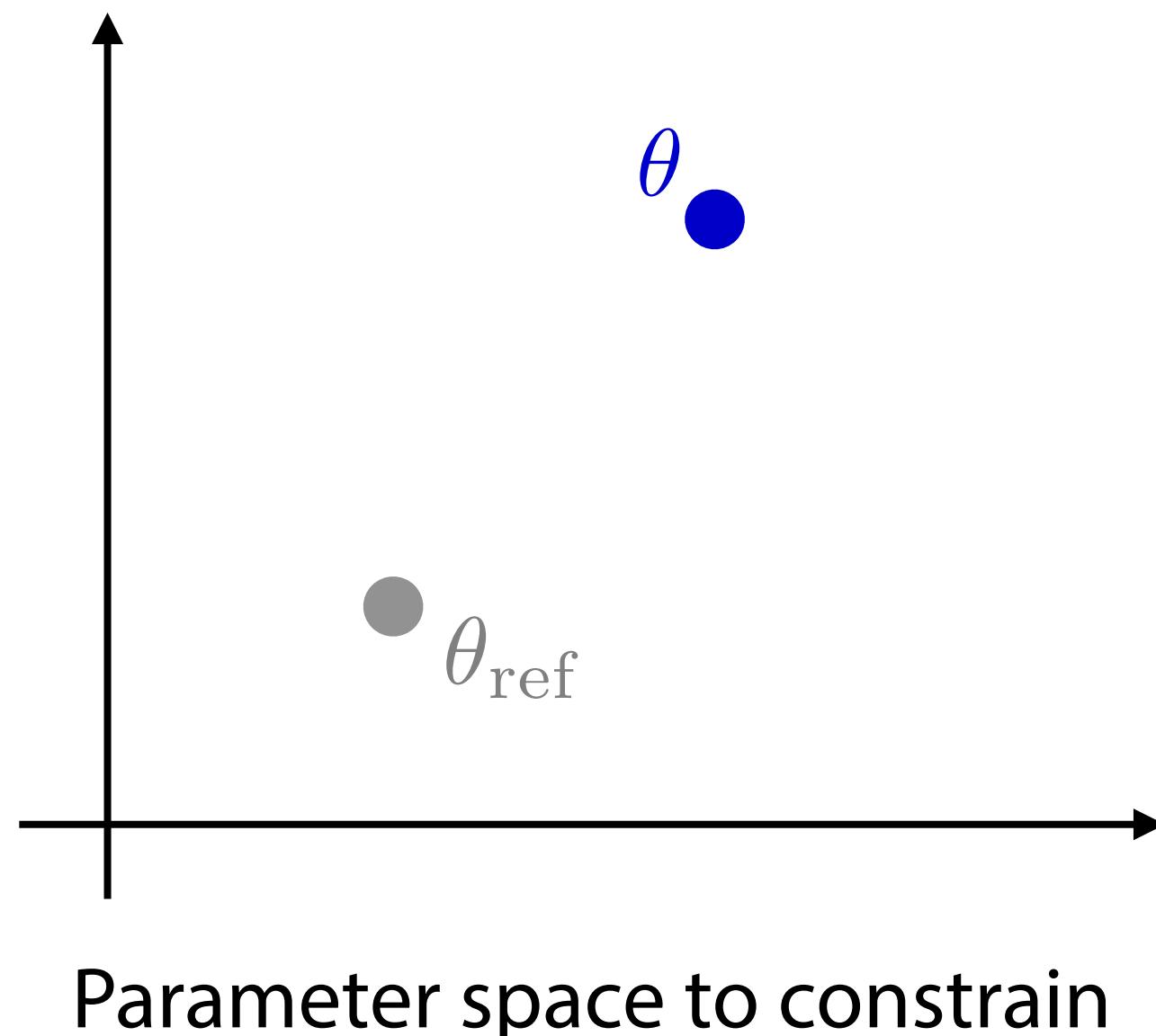
Abstract We describe likelihood-based statistical tests for use in high energy physics for the discovery of new phenomena and for construction of confidence intervals on model

data sets by a single representative one, referred to here as the “Asimov” data set.¹ In the past, this method has been used and justified intuitively (e.g., [4, 5]). Here we provide

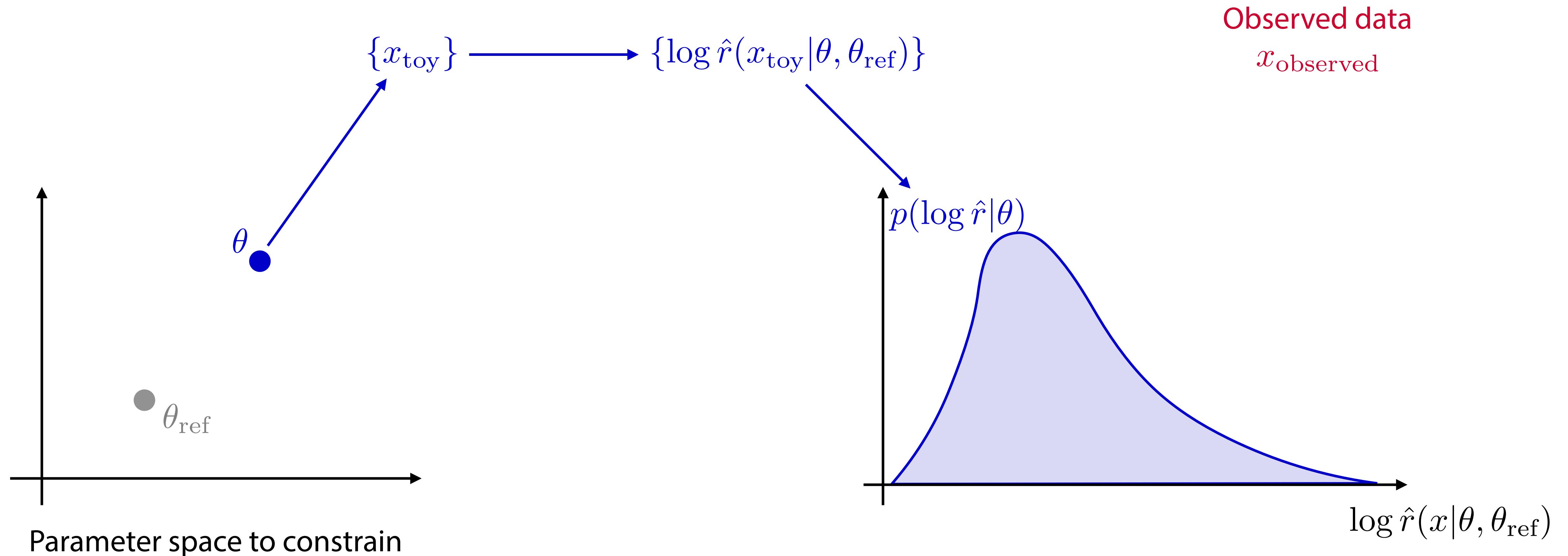
Limit setting (frequentist, standard ATLAS / CMS practice)

Observed data

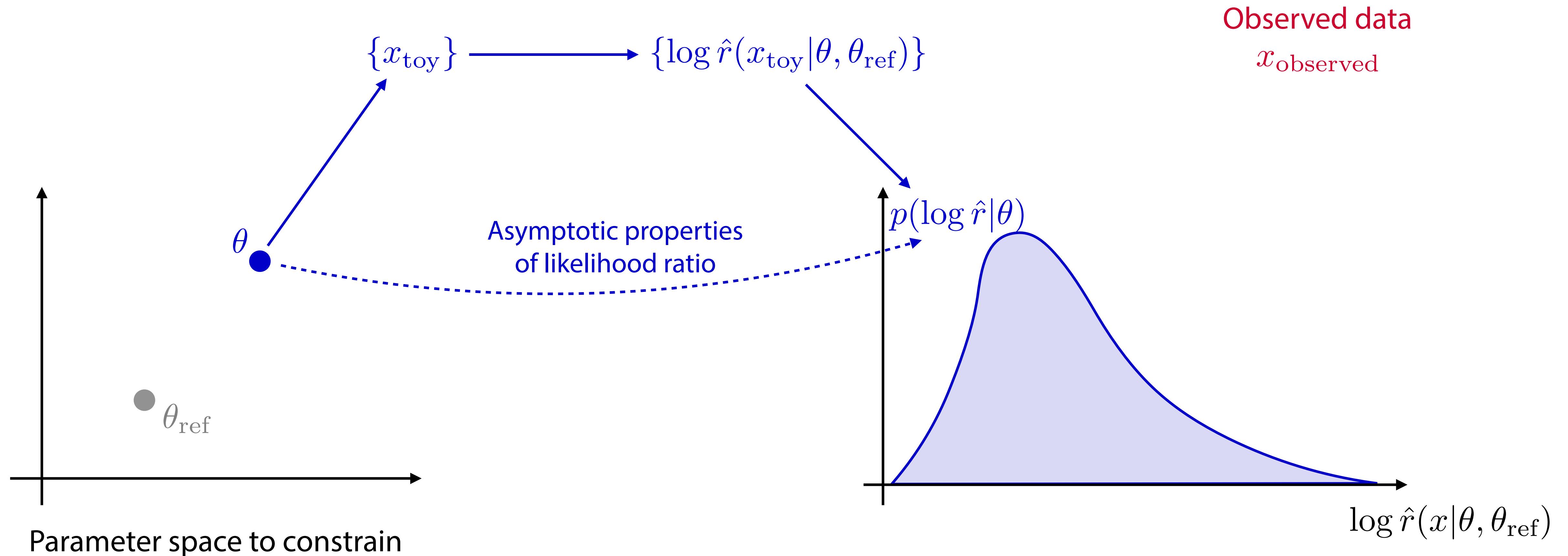
x_{observed}



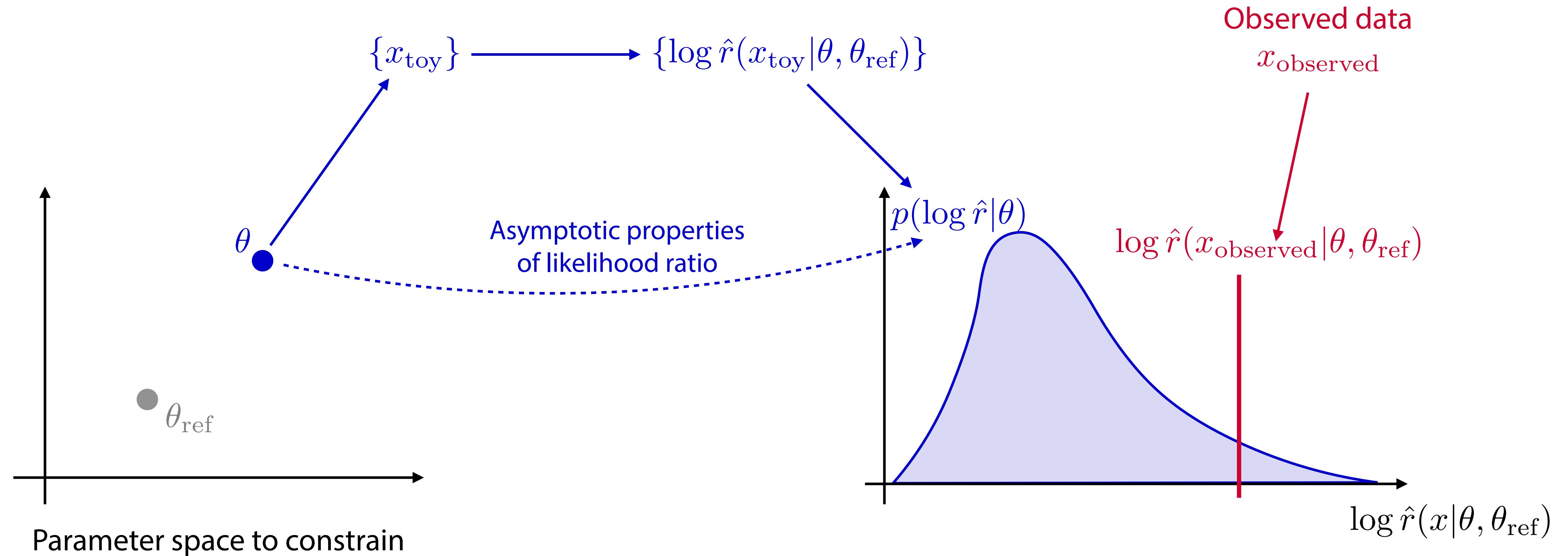
Limit setting (frequentist, standard ATLAS / CMS practice)



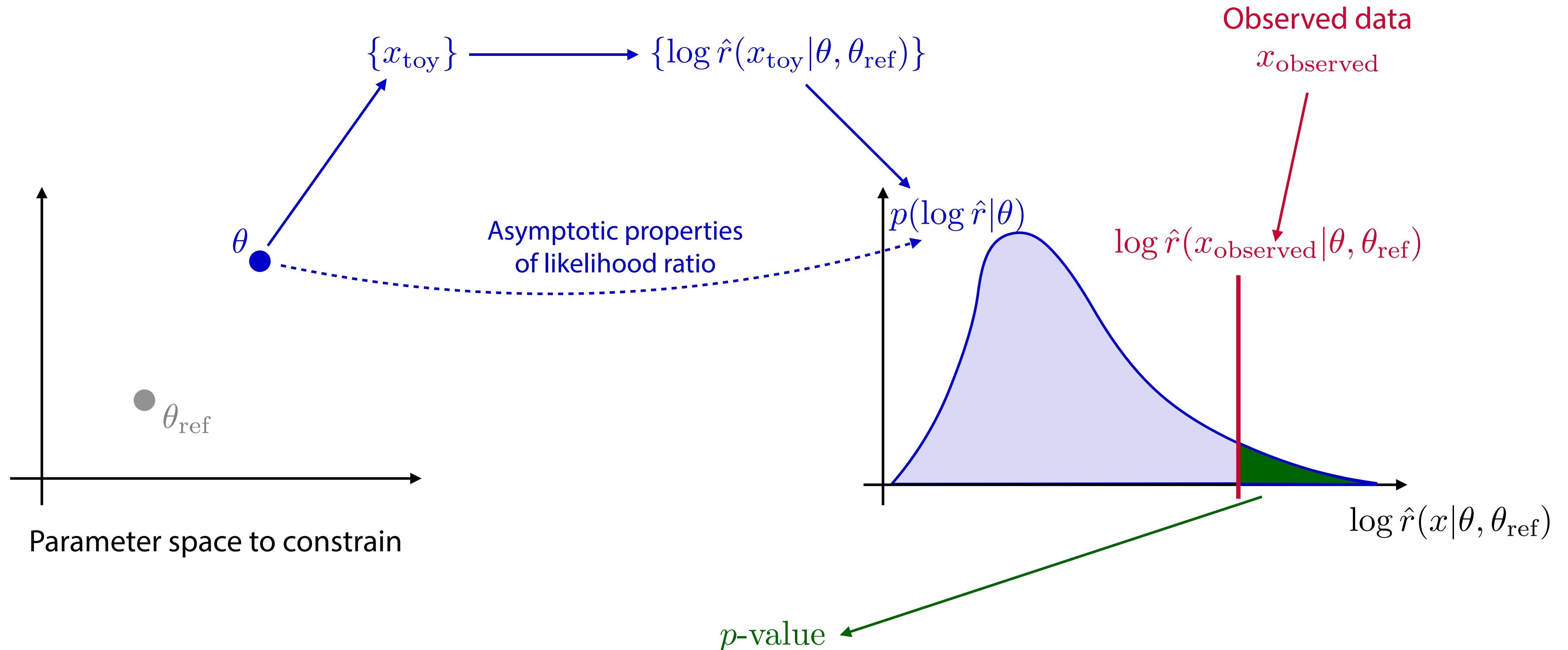
Limit setting (frequentist, standard ATLAS / CMS practice)



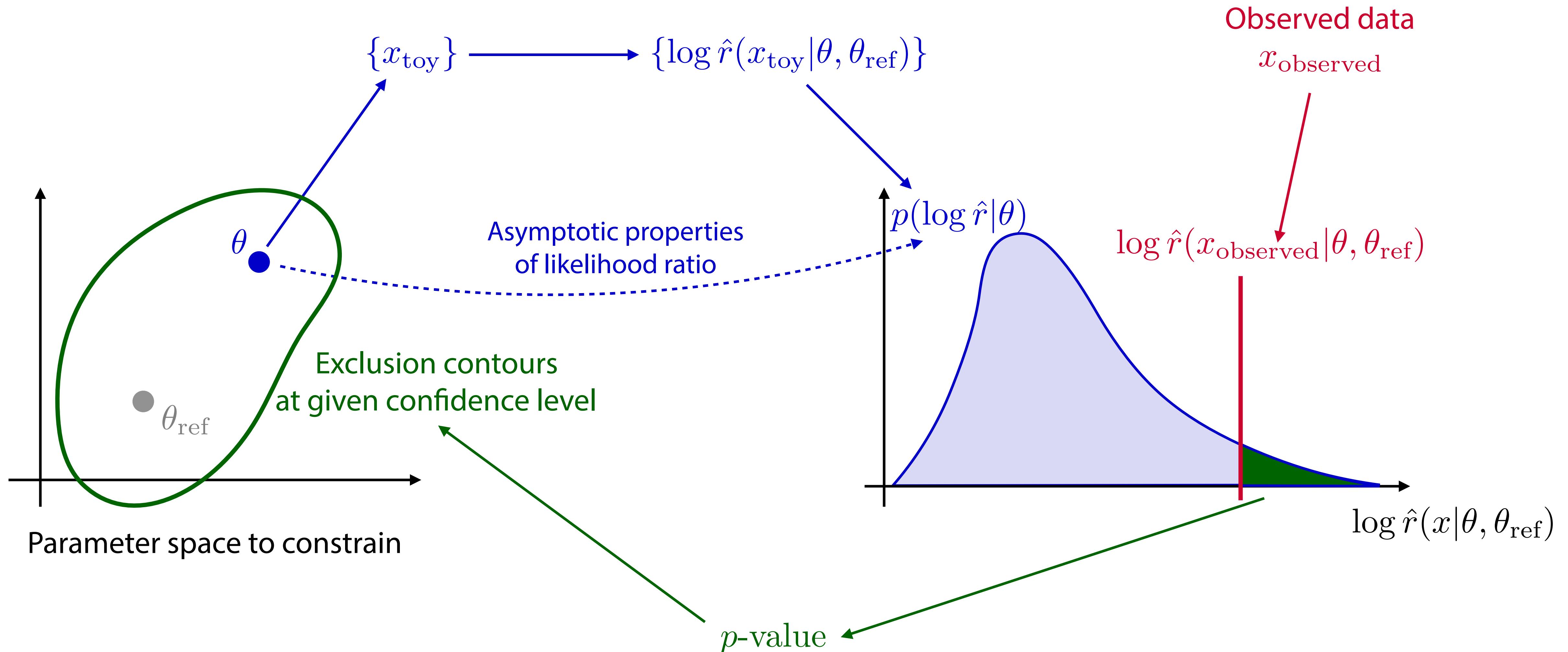
Limit setting (frequentist, standard ATLAS / CMS practice)



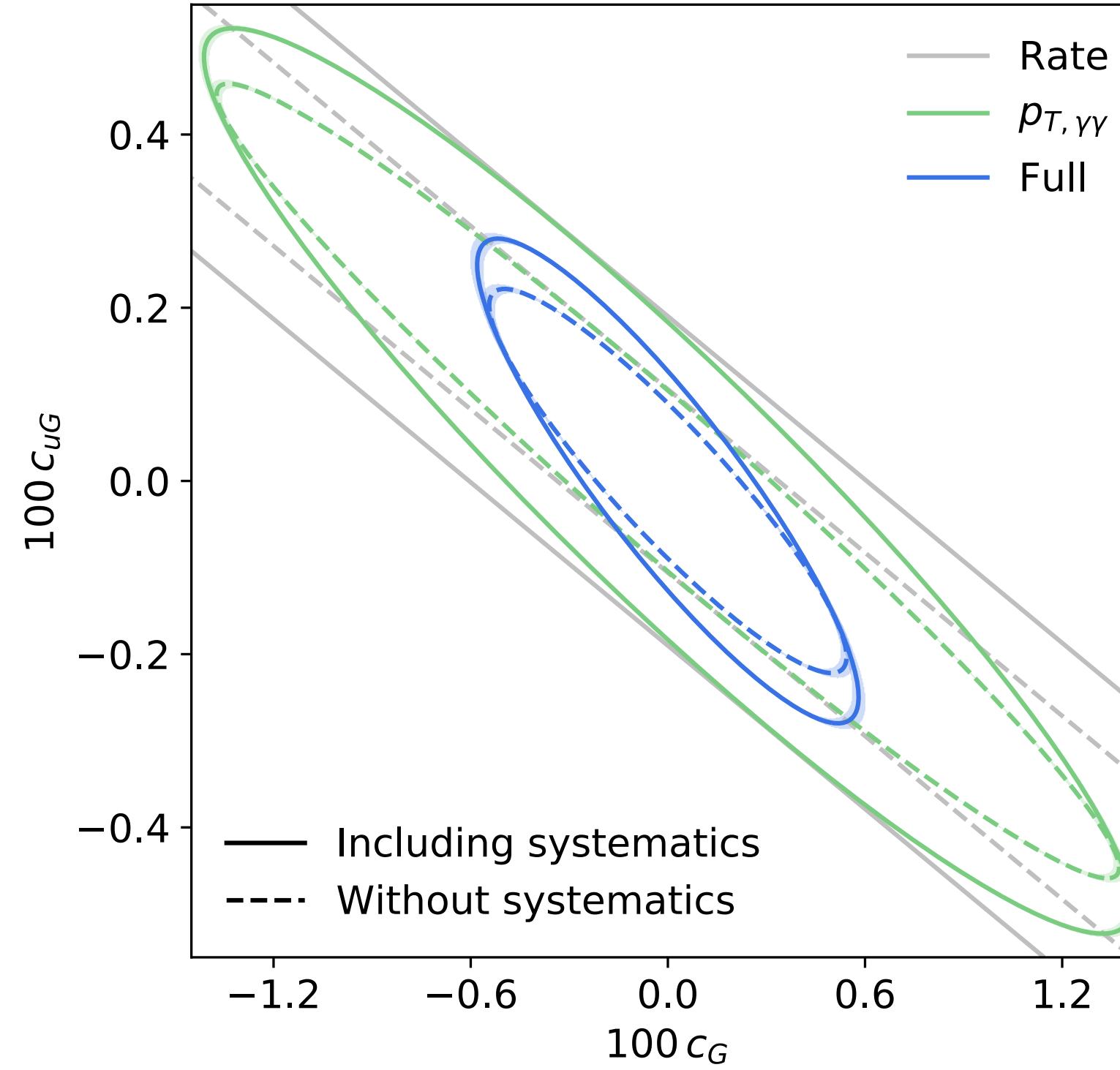
Limit setting (frequentist, standard ATLAS / CMS practice)



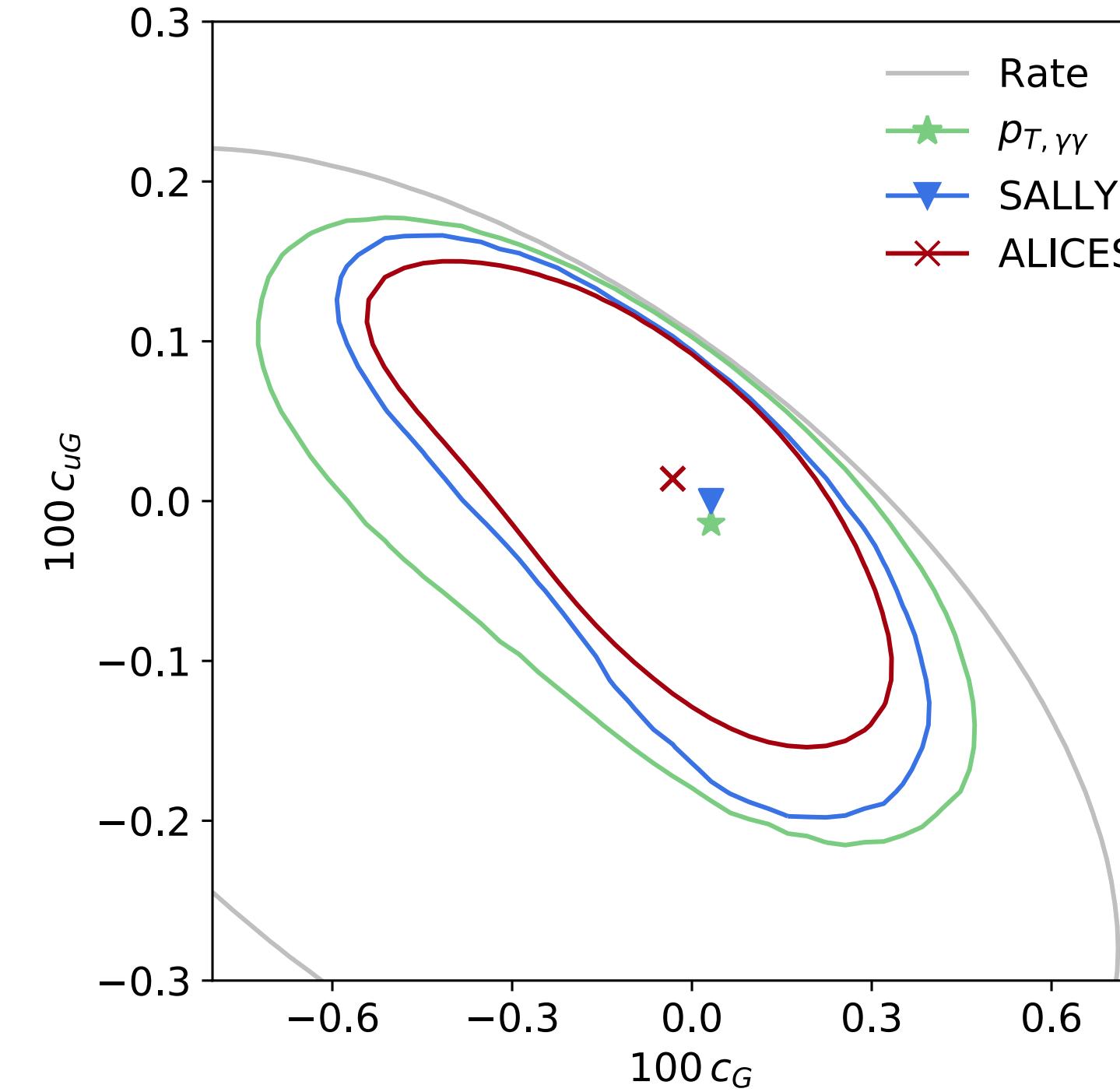
Limit setting (frequentist, standard ATLAS / CMS practice)



Types of Inference Results & Sensitivity Summaries



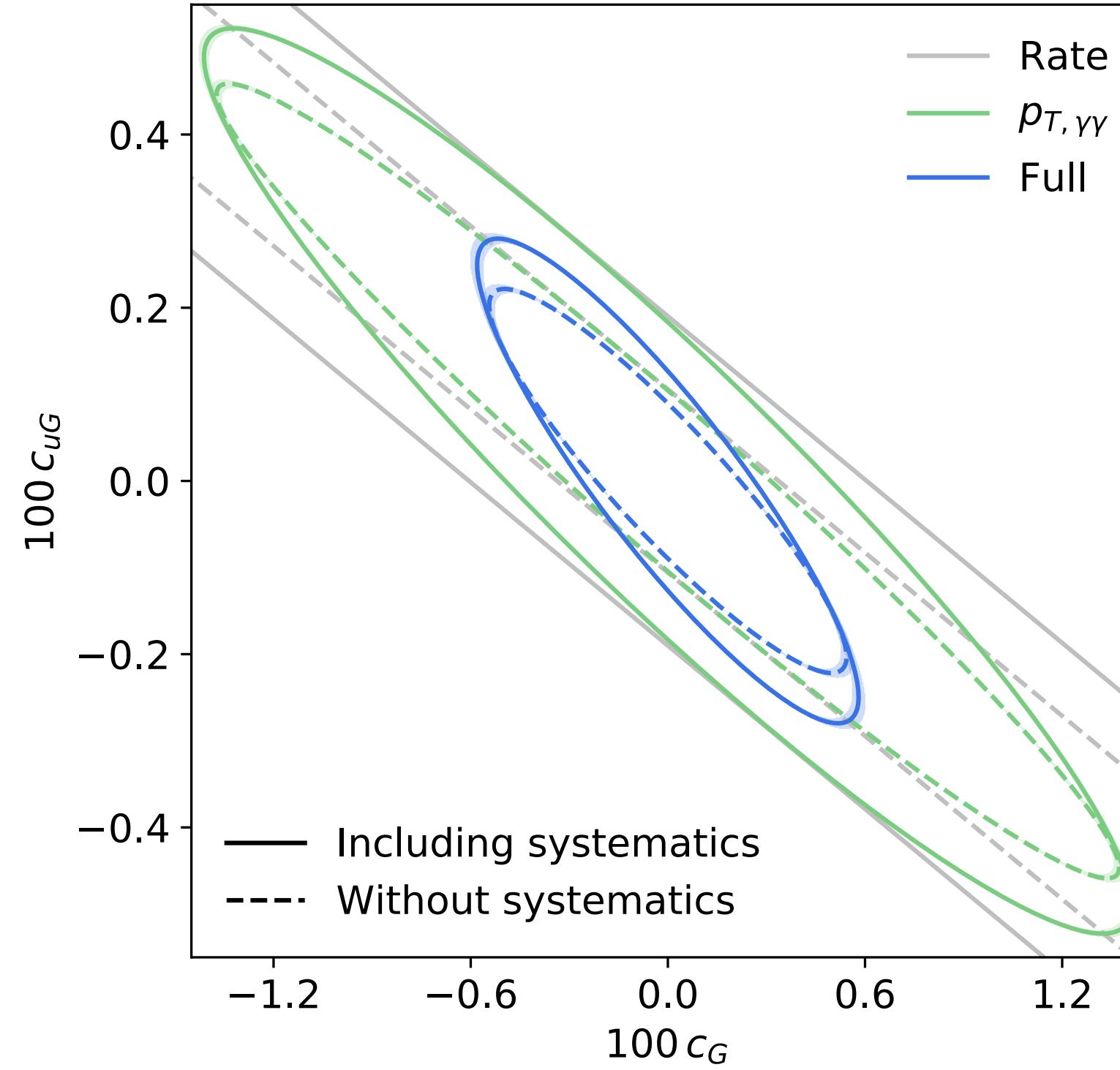
Elliptical Contours
(with / without systematics)



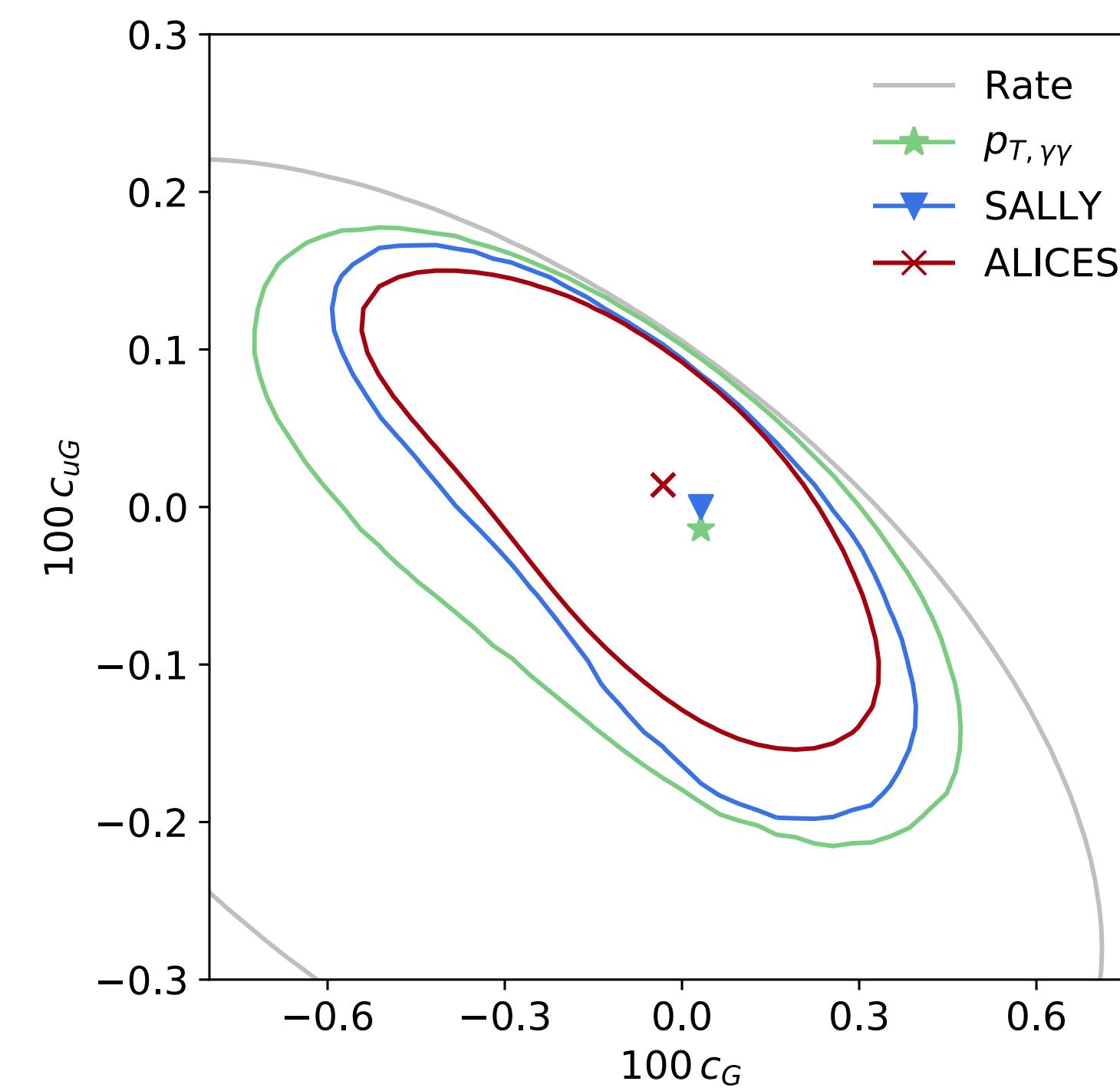
Non-Elliptical Contours

Fisher Information
&
Information Geometry

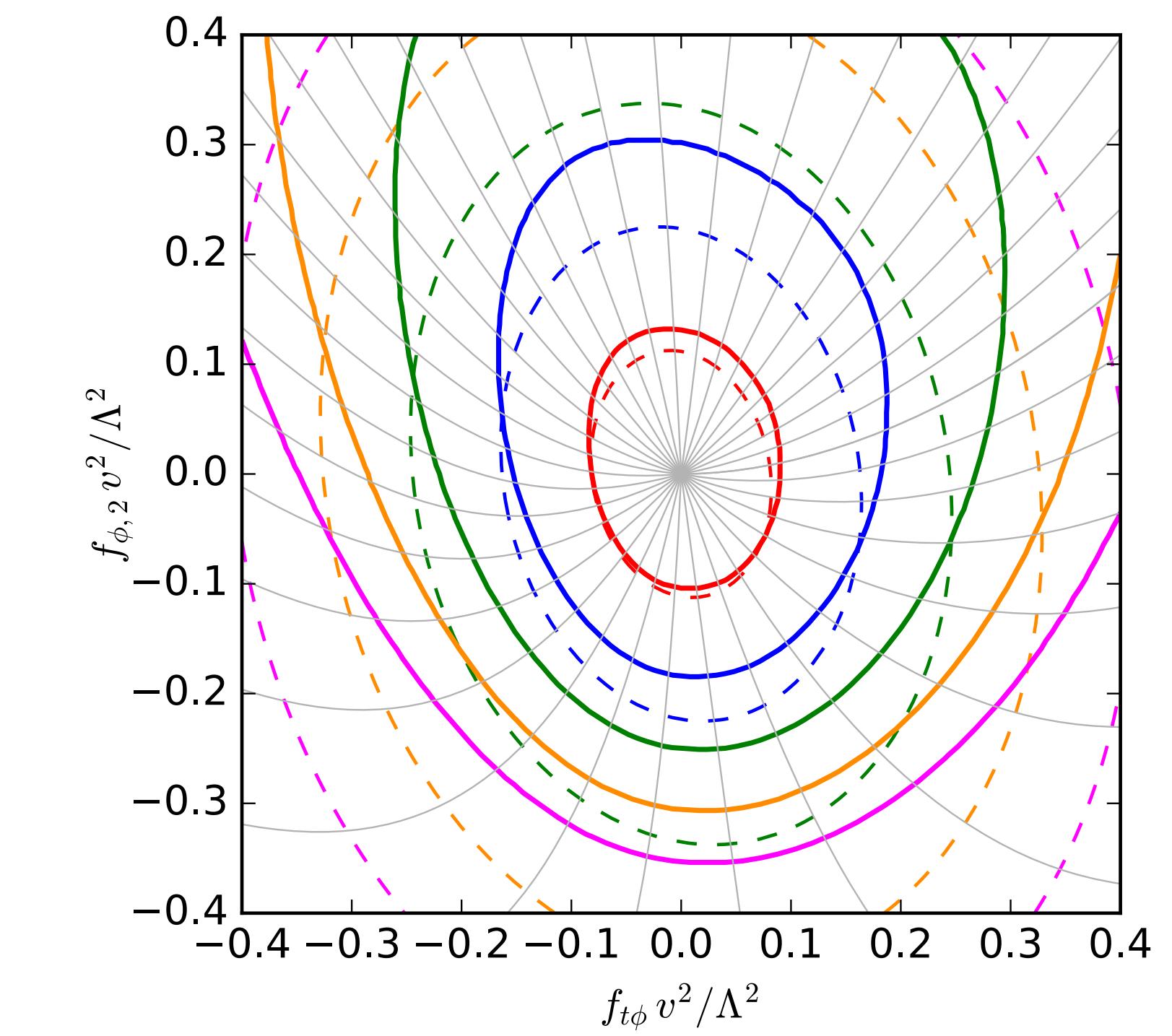
Types of Inference Results & Sensitivity Summaries



Elliptical Contours
(with / without systematics)



Non-Elliptical Contours



Fisher Information
&
Information Geometry

MadMiner automates all of these methods.

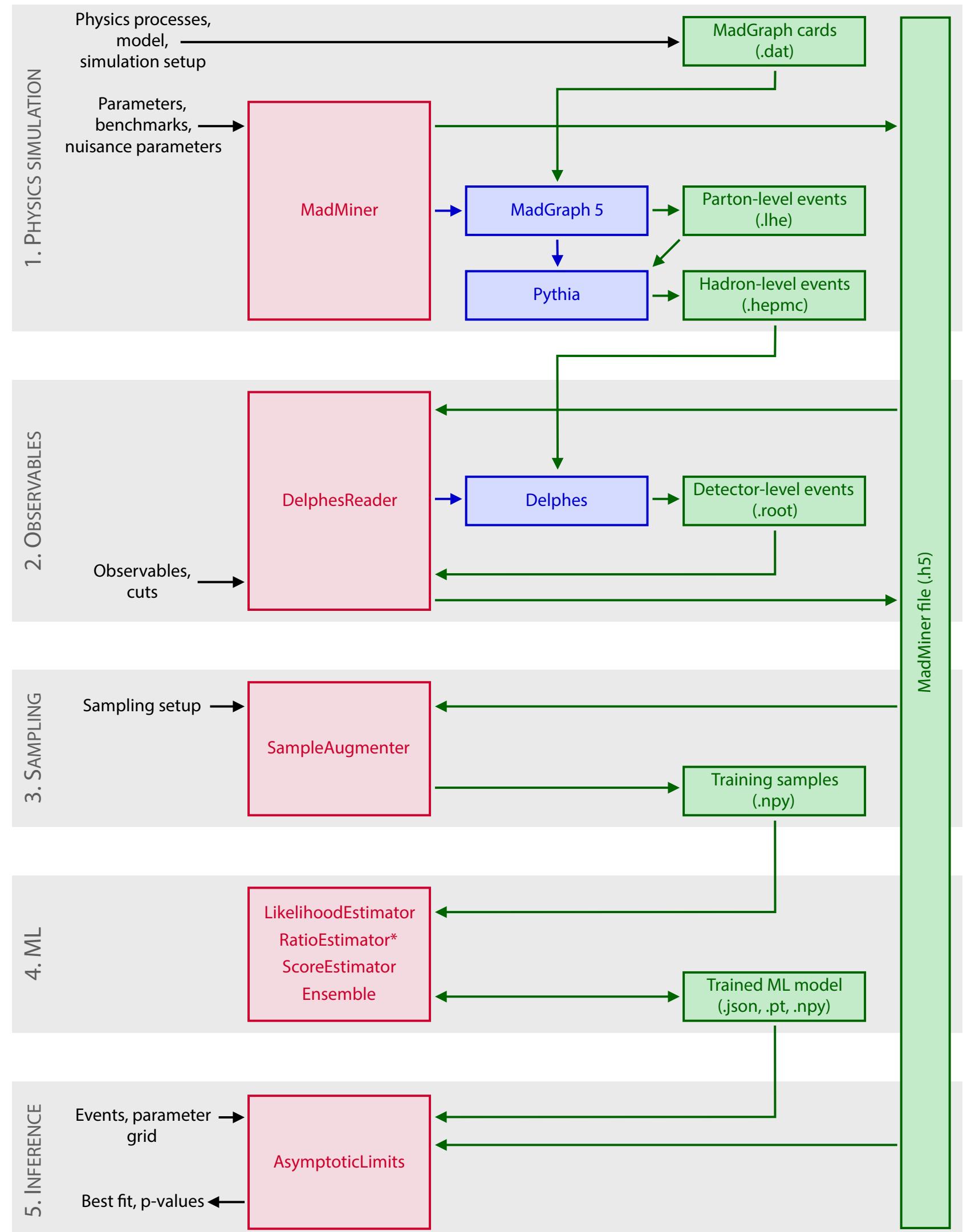
[JB, F. Kling, I. Espejo, K. Cranmer 1907.10621]

Automation

[JB, F. Kling, I. Espejo, K. Cranmer 1907.10621]

New Python package **MadMiner** makes it straightforward to apply the new techniques to LHC problems

- Out of the box: Pheno-level analyses
 - MadGraph, Pythia, Delphes
 - Systematic uncertainties from PDF / scale variation
- Scalable to state-of-the-art experimental tools
 - Mostly requires bookkeeping of fully differential cross sections
- Modular interface
 - Extensive documentation
 - Embedded into Python / ML ecosystem



MadMiner: Machine learning–based inference for particle physics

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

[pypi package 0.6.3](#) [build passing](#) [docs failing](#) [chat on gitter](#) [code style black](#) [License MIT](#) [DOI 10.5281/zenodo.1489147](#)
[arXiv 1907.10621](#)

Introduction

Particle physics processes are usually modeled with complex Monte-Carlo simulations of the hard process, parton shower, and detector interactions. These simulators typically do not admit a tractable likelihood function: given a (potentially high-dimensional) set of observables, it is usually not possible to calculate the probability of these observables for some model parameters. Particle physicists usually tackle this problem of "likelihood-free inference" by hand-picking a few "good" observables or summary statistics and filling histograms of them. But this conventional

UCI-TR-2019-16, SLAC-PUB-17461

MadMiner: Machine learning–based inference for particle physics

Johann Brehmer,^{1,*} Felix Kling,^{2,3,†} Irina Espejo,^{1,‡} and Kyle Cranmer^{1,§}

¹ *Center for Data Science and Center for Cosmology and Particle Physics,
New York University, New York, NY 10003, USA*

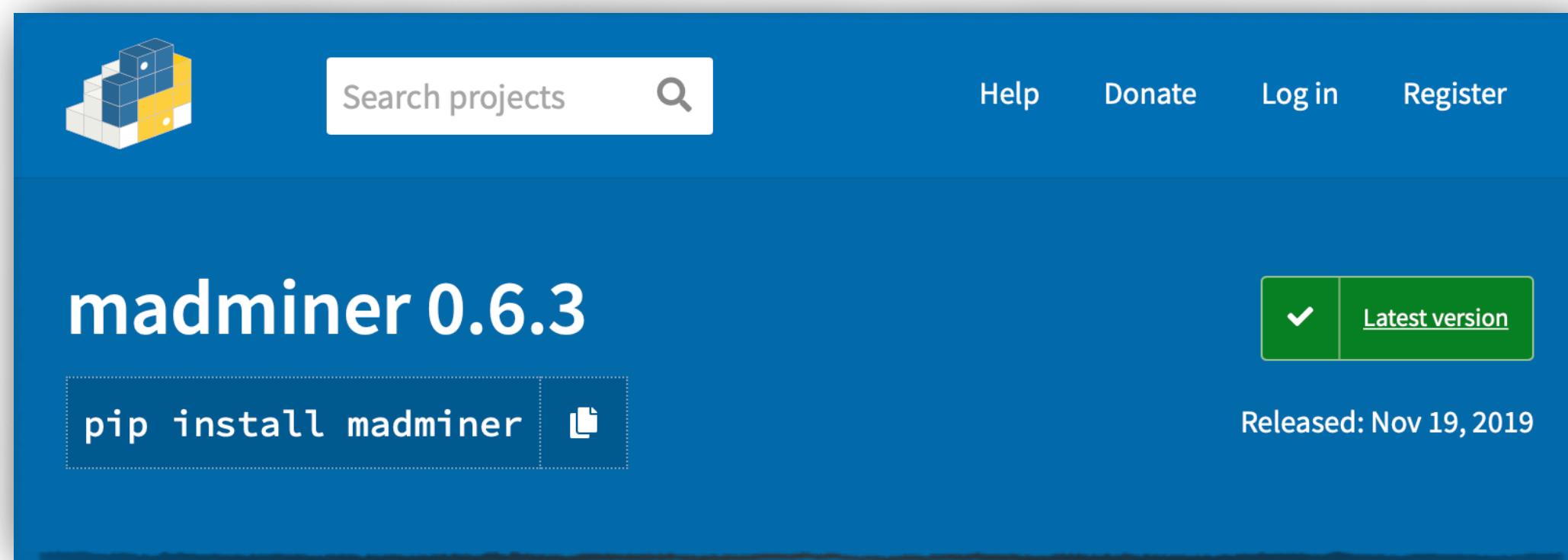
² *Department of Physics and Astronomy, University of California, Irvine, CA 92697, USA*

³ *SLAC National Accelerator Laboratory, 2575 Sand Hill Road, Menlo Park, CA 94025, USA*

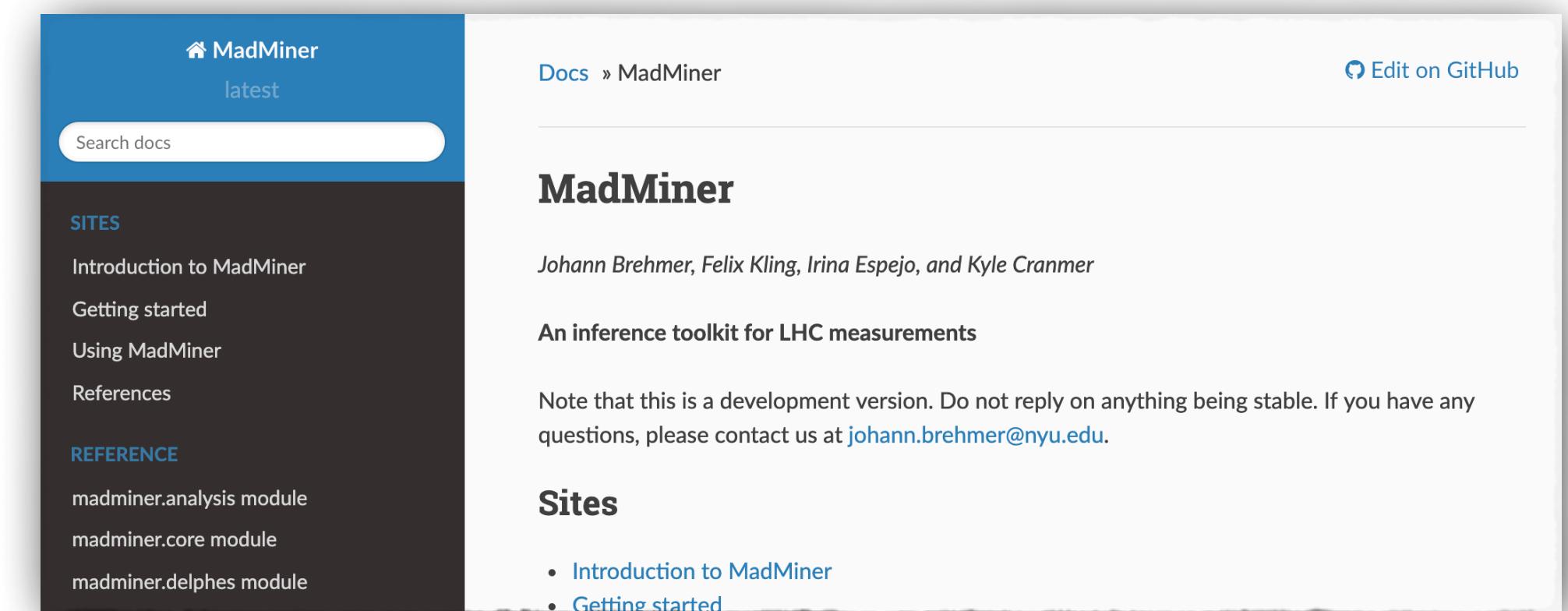
Precision measurements at the LHC often require analyzing high-dimensional event data for subtle kinematic signatures, which is challenging for established analysis methods. Recently, a powerful family of multivariate inference techniques that leverage both matrix element information and machine learning has been developed. This approach neither requires the reduction of high-dimensional data to summary statistics nor any simplifications to the underlying physics or detector response. In this paper we introduce **MadMiner**, a Python module

Repository and tutorials:
github.com/johannbrehmer/madminer

Paper with detailed explanations:
[1907.10621](https://arxiv.org/abs/1907.10621)



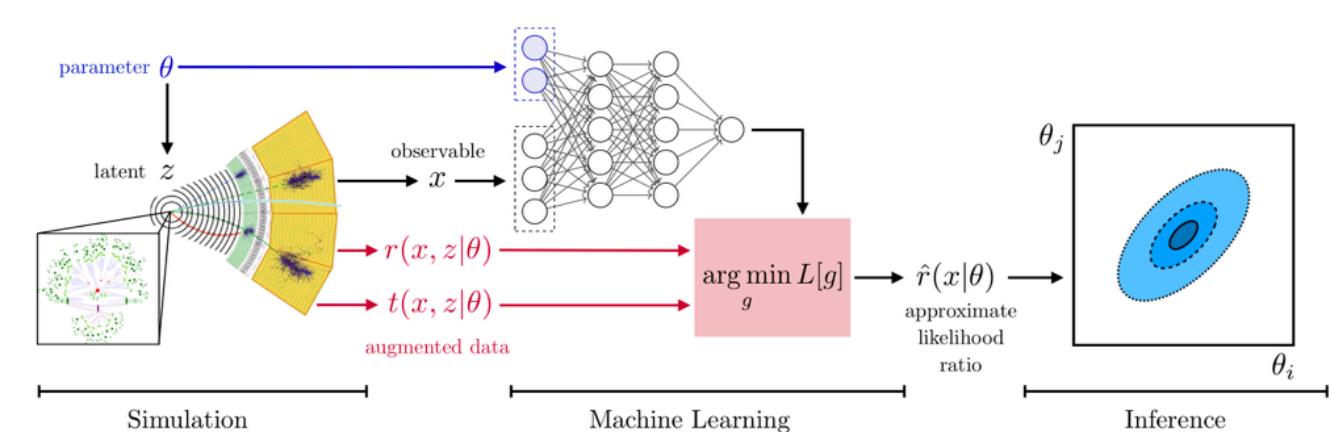
Installation:
`pip install madminer`



API documentation:
madminer.readthedocs.io

Hands-on Tutorial

- Step-by-step instructions
- Do Preliminaries
 - Need to install Docker
 - And then start Jupyter
- Step 1 is fast
- Step 2 takes ~25 min
- Step 3 takes ~20 min
- While they are running I will lecture
- Then we will finish with results



← 

Introduction

MadMiner tutorial

This is a tutorial on [MadMiner](#) developed by Johann Brehmer, Felix Kling, Irina Espejo, and Kyle Cranmer. It is built using [Jupyter Book](#).

MadMiner Tutorial

- Introduction
- MadMiner Tutorial
- Preliminaries
- Overview
- Define process to study *
- Morphing
- Interactive Morphing Demo
- Create training data
- Set MadGraph Directory
- Parton Level *
- With Delphes
- Train model
- Likelihood Ratio *
- Score *
- Likelihood
- Statistical Analysis
- Limits on EFT parameters *
- Fisher Information
- Information Geometry
- Congratulations

Introduction to MadMiner

Particle physics processes are usually modelled with complex Monte-Carlo simulations of the hard process, parton shower, and detector interactions. These simulators typically do not admit a tractable likelihood function: given a (potentially high-dimensional) set of observables, it is usually not possible to calculate the probability of these observables for some model parameters. Particle physicists usually tackle this problem of “likelihood-free inference” by hand-picking a few “good” observables or summary statistics and filling histograms of them. But this conventional approach discards the information in all other observables and often does not scale well to high-dimensional problems.

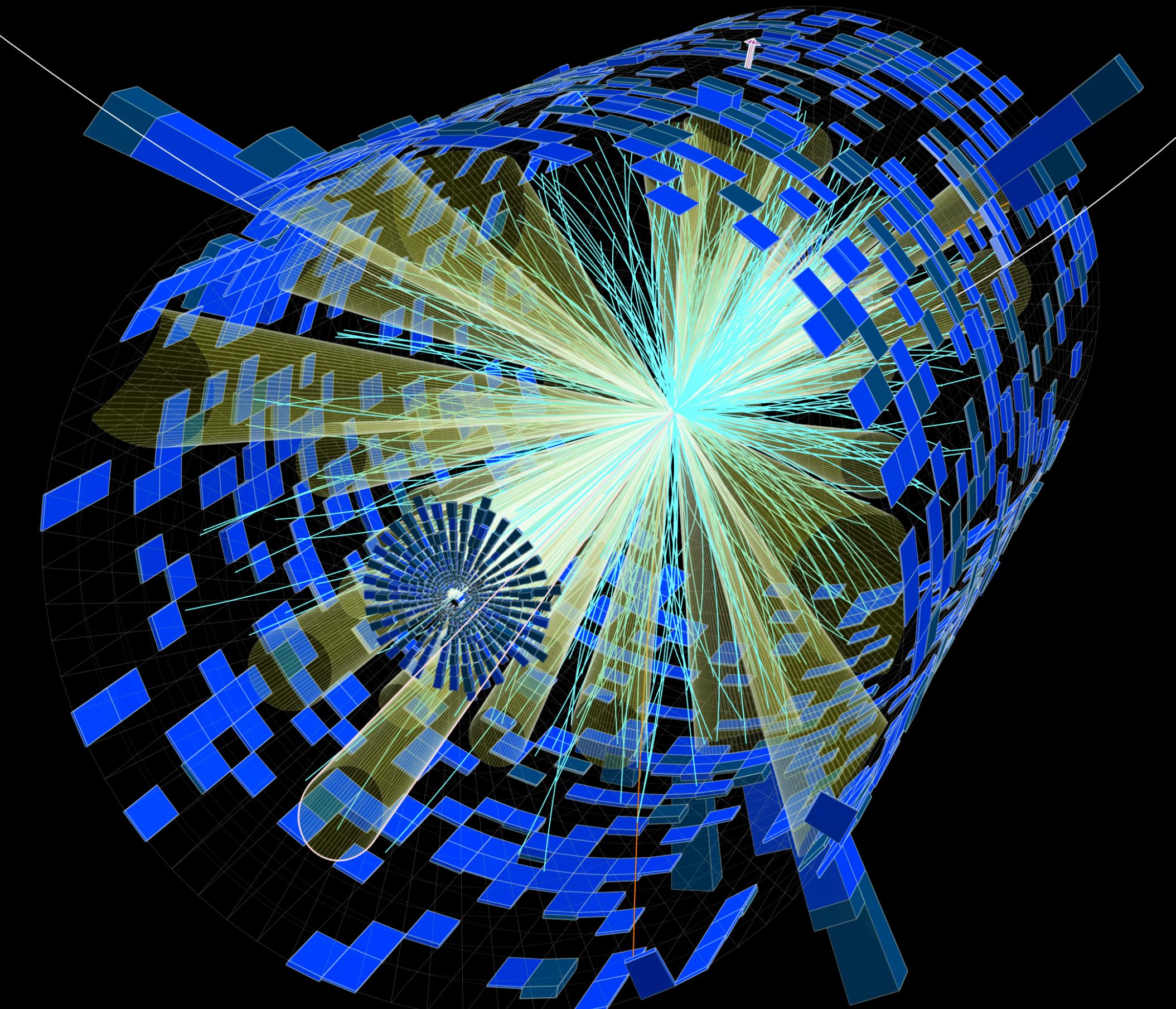
In the three publications [“Constraining Effective Field Theories With Machine Learning”](#), [“A Guide to Constraining Effective Field Theories With Machine Learning”](#), and [“Mining gold from implicit models to improve likelihood-free inference”](#), a new approach has been developed. In a nut shell, additional information is extracted from the simulations that is closely related to the matrix elements that determine the hard process. This “augmented data” can be used to train neural networks to efficiently approximate arbitrary likelihood ratios. We playfully call this process “mining gold” from the simulator, since this information may be hard to get, but turns out to be very valuable for inference.



MACHINE LEARNING

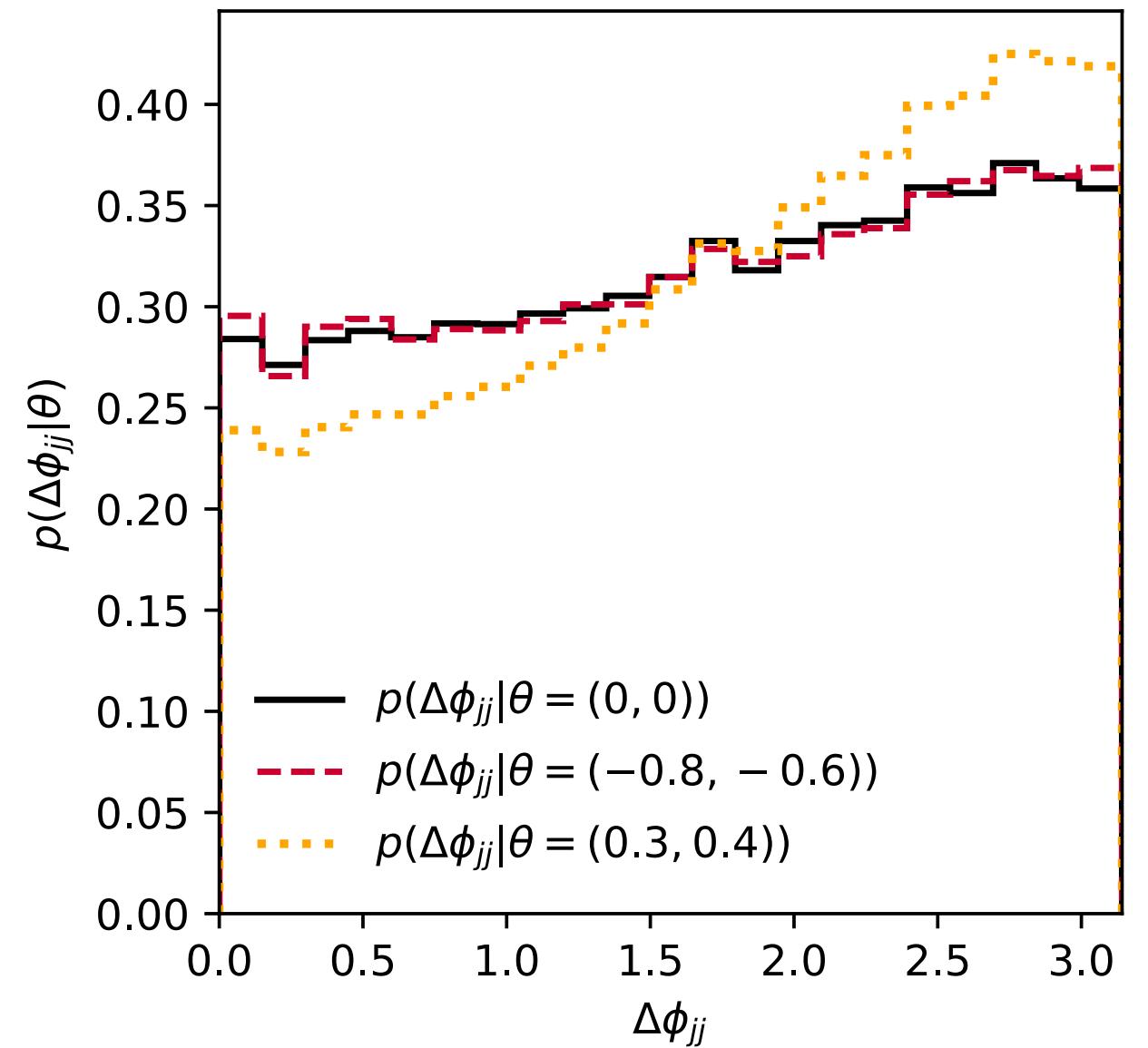
LECTURE 4

@KyleCranmer
New York University
Department of Physics
Center for Data Science
CILVR Lab



These techniques let us constrain
effective theories more effectively.

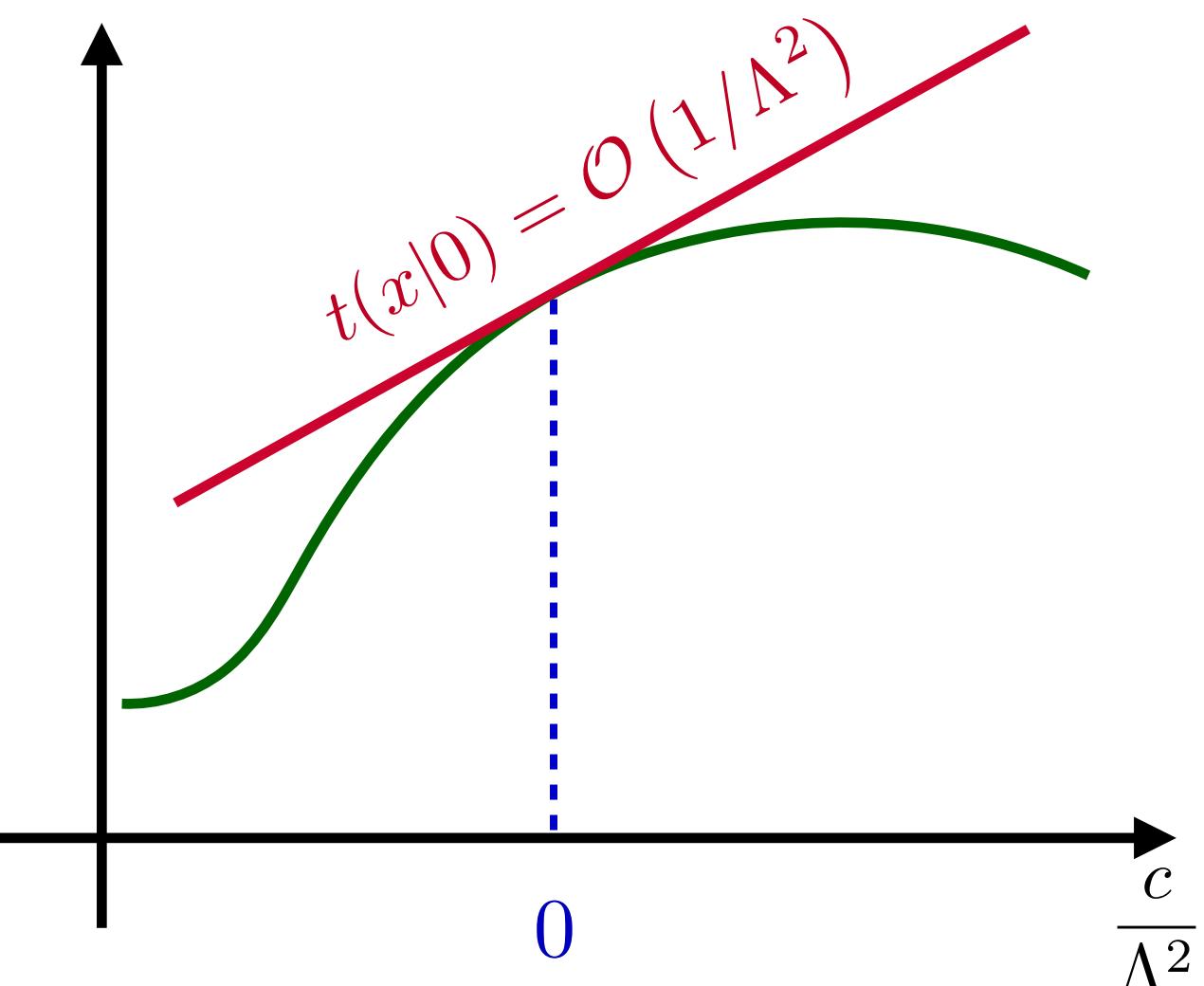
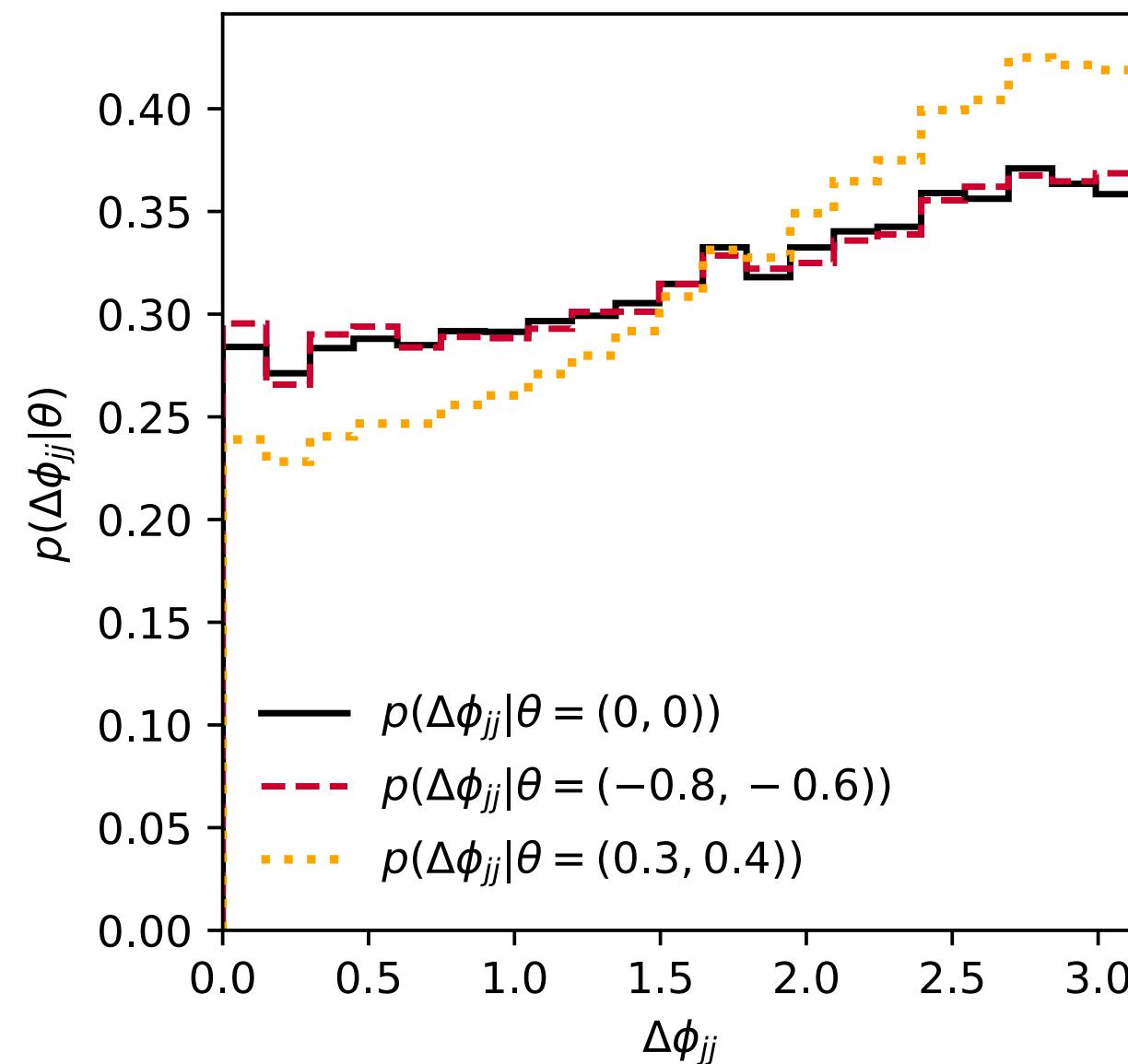
Perfect match for EFT measurements



- Good for subtle kinematic effects

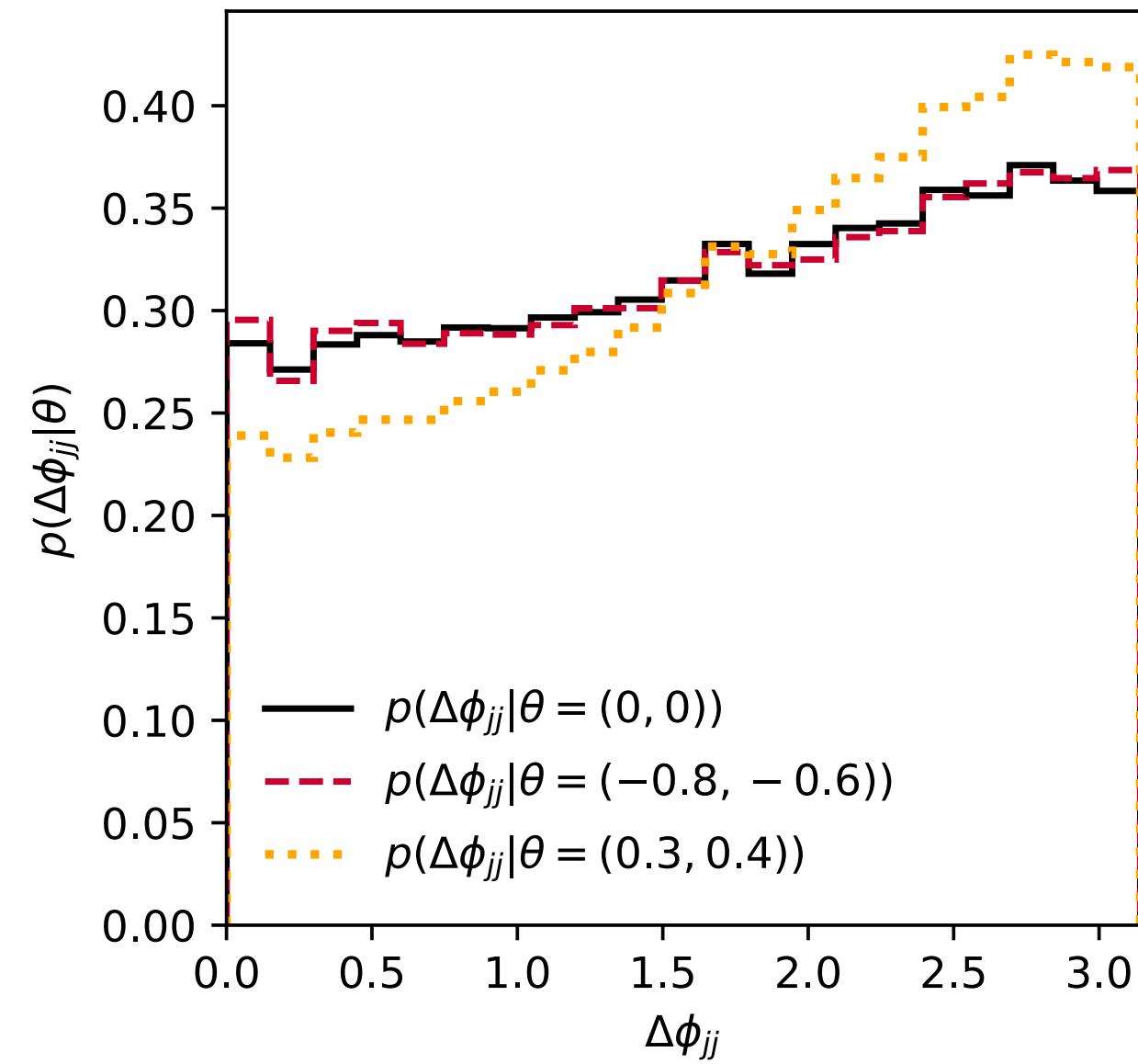
(Subtle point: Large overlap of kinematic distributions reduces variance of joint likelihood ratio / joint score)

Perfect match for EFT measurements

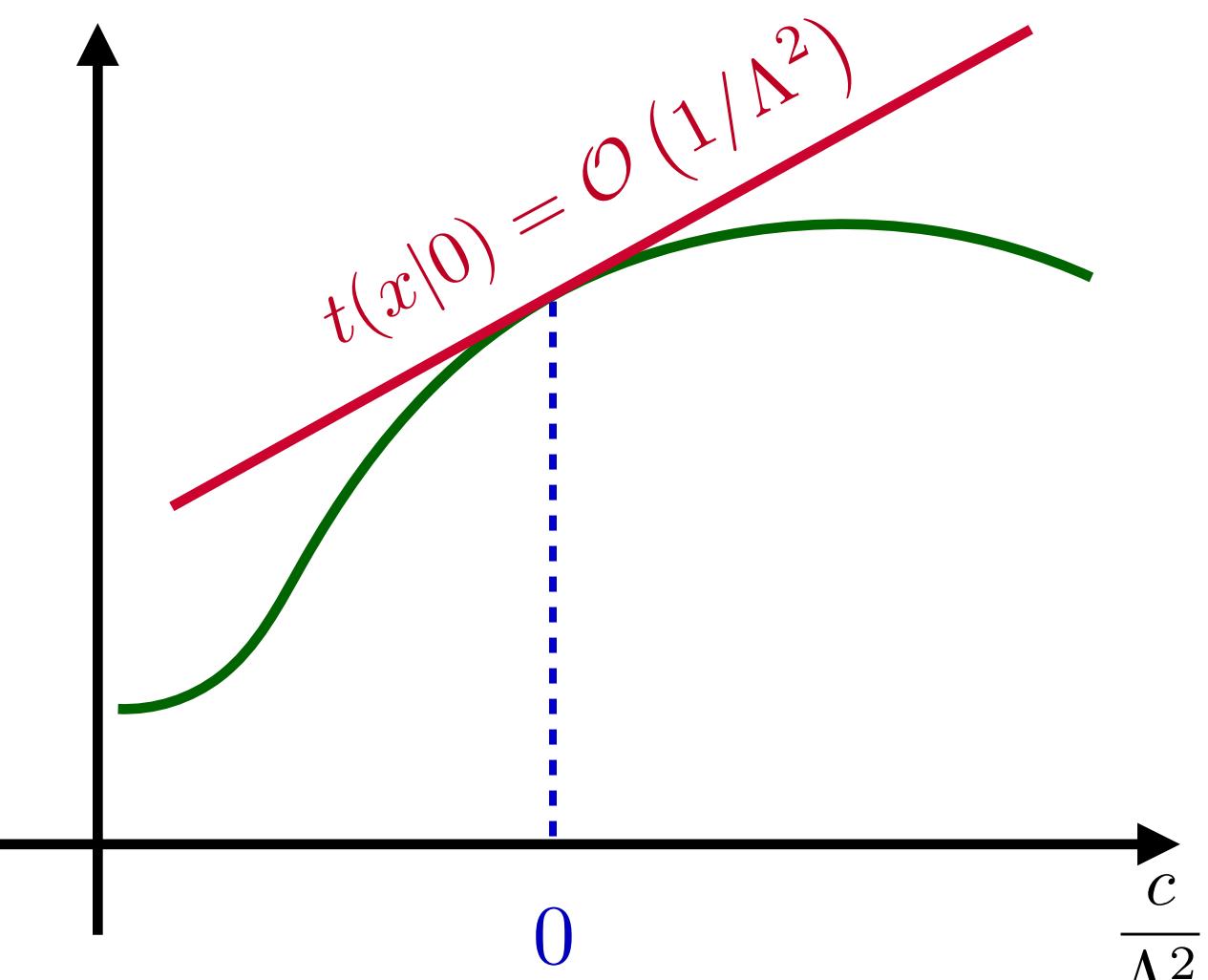


- Good for subtle kinematic effects
(Subtle point: Large overlap of kinematic distributions reduces variance of joint likelihood ratio / joint score)
- Interference effects can be isolated using SALLY at the SM (SMALLY?)

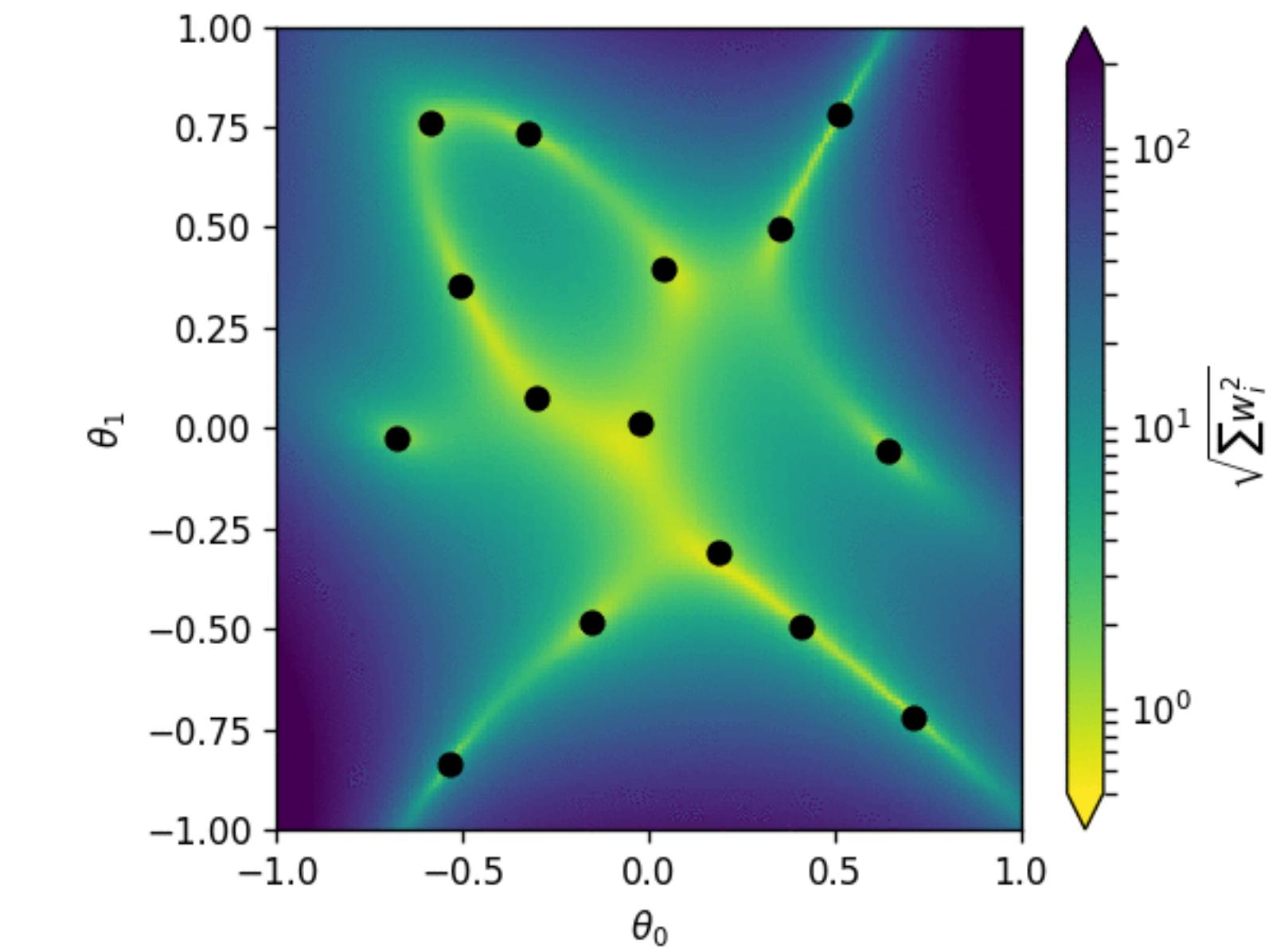
Perfect match for EFT measurements



- Good for subtle kinematic effects
(Subtle point: Large overlap of kinematic distributions reduces variance of joint likelihood ratio / joint score)



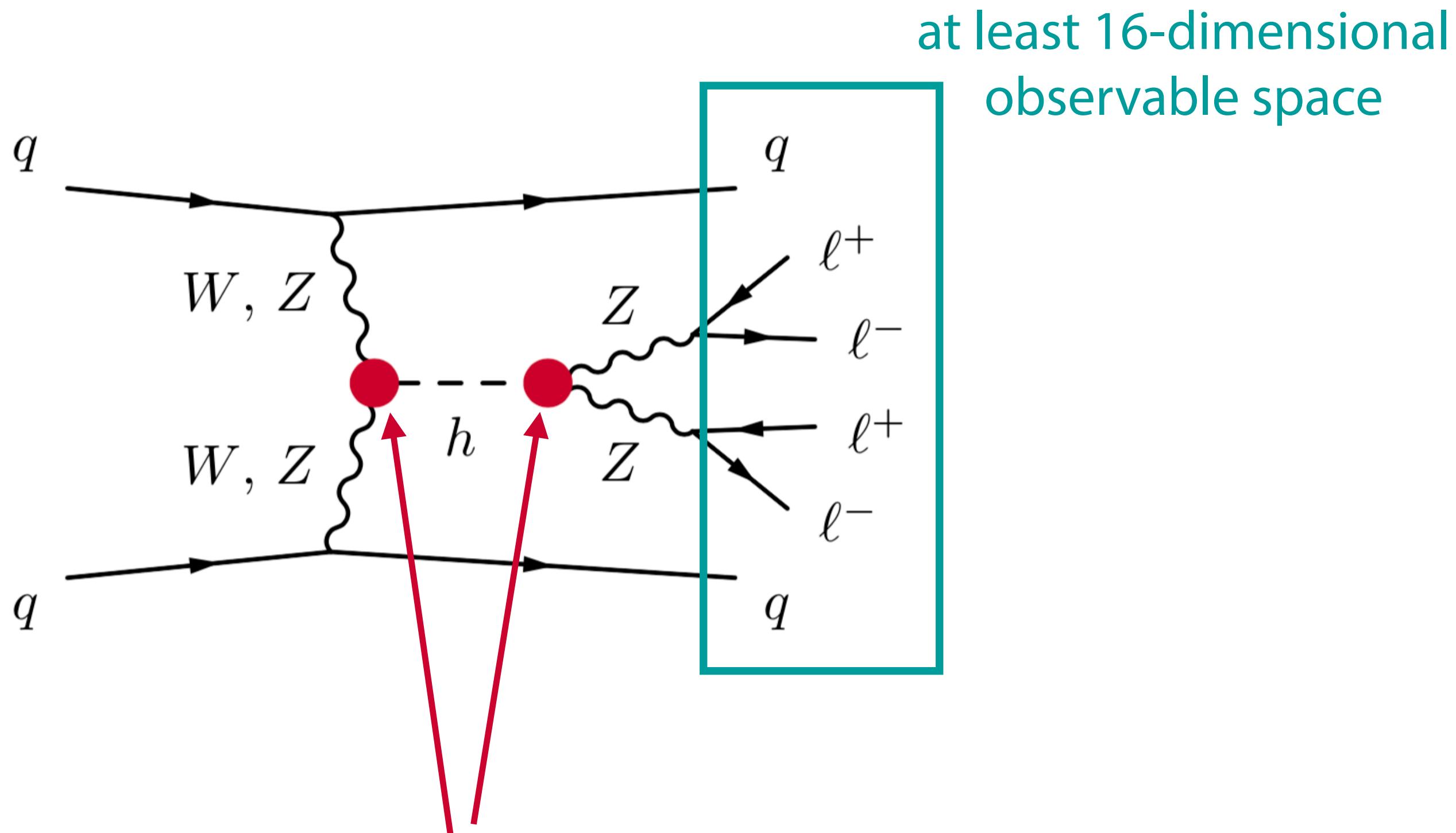
- Interference effects can be isolated using SALLY at the SM (SMALLY?)



- Morphing techniques allow fast reweighting to any parameter points
[e.g. ATL-PHYS-PUB-2015-047]

Proof of concept: Higgs production in weak boson fusion

[JB, K. Cranmer, G. Louppe, J. Pavez
1805.00013, 1805.00020]



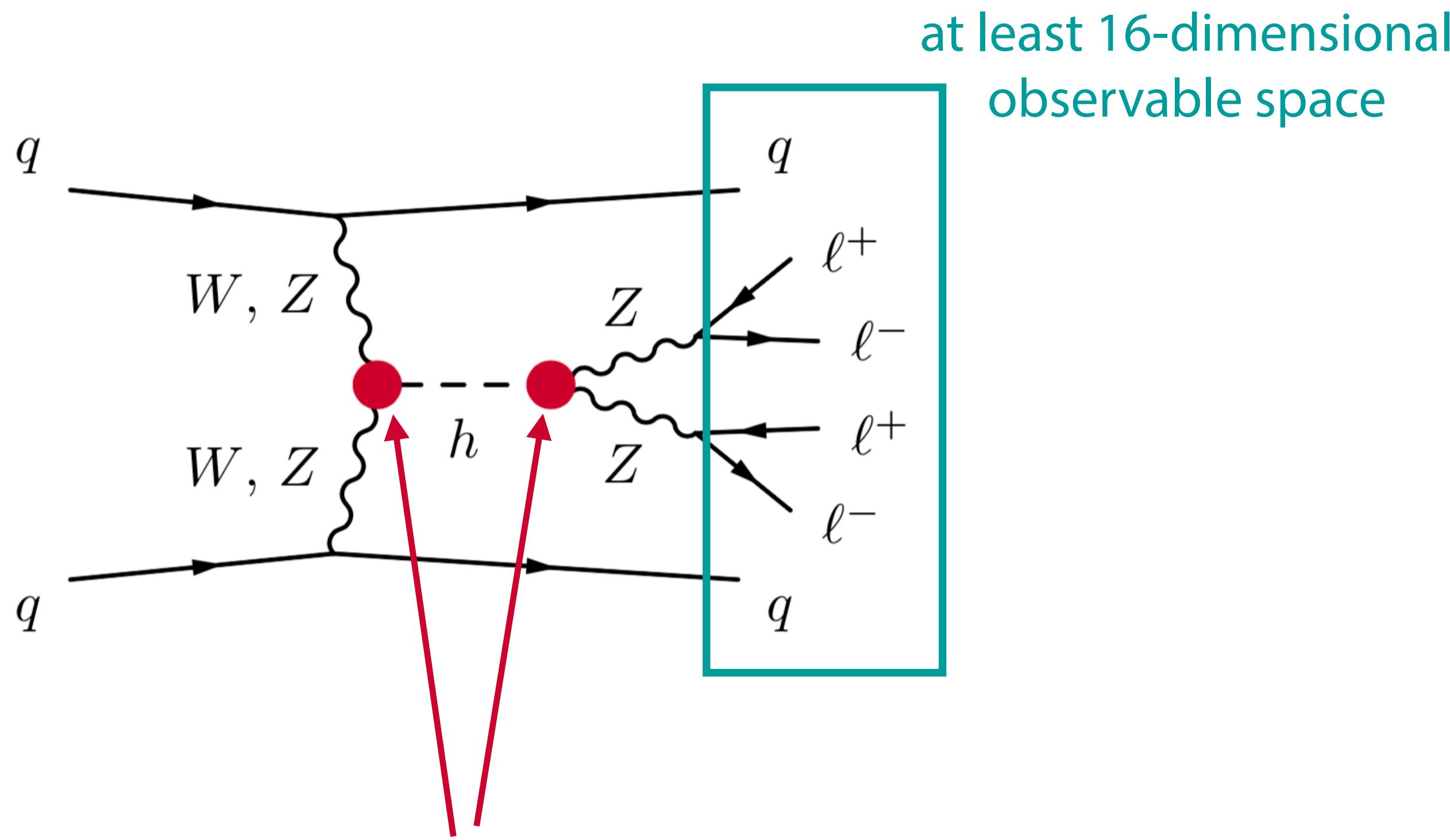
Exciting new physics might hide here!

We parameterize it with two EFT coefficients:

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

Proof of concept: Higgs production in weak boson fusion

[JB, K. Cranmer, G. Louppe, J. Pavez
1805.00013, 1805.00020]



We parameterize it with two EFT coefficients:

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

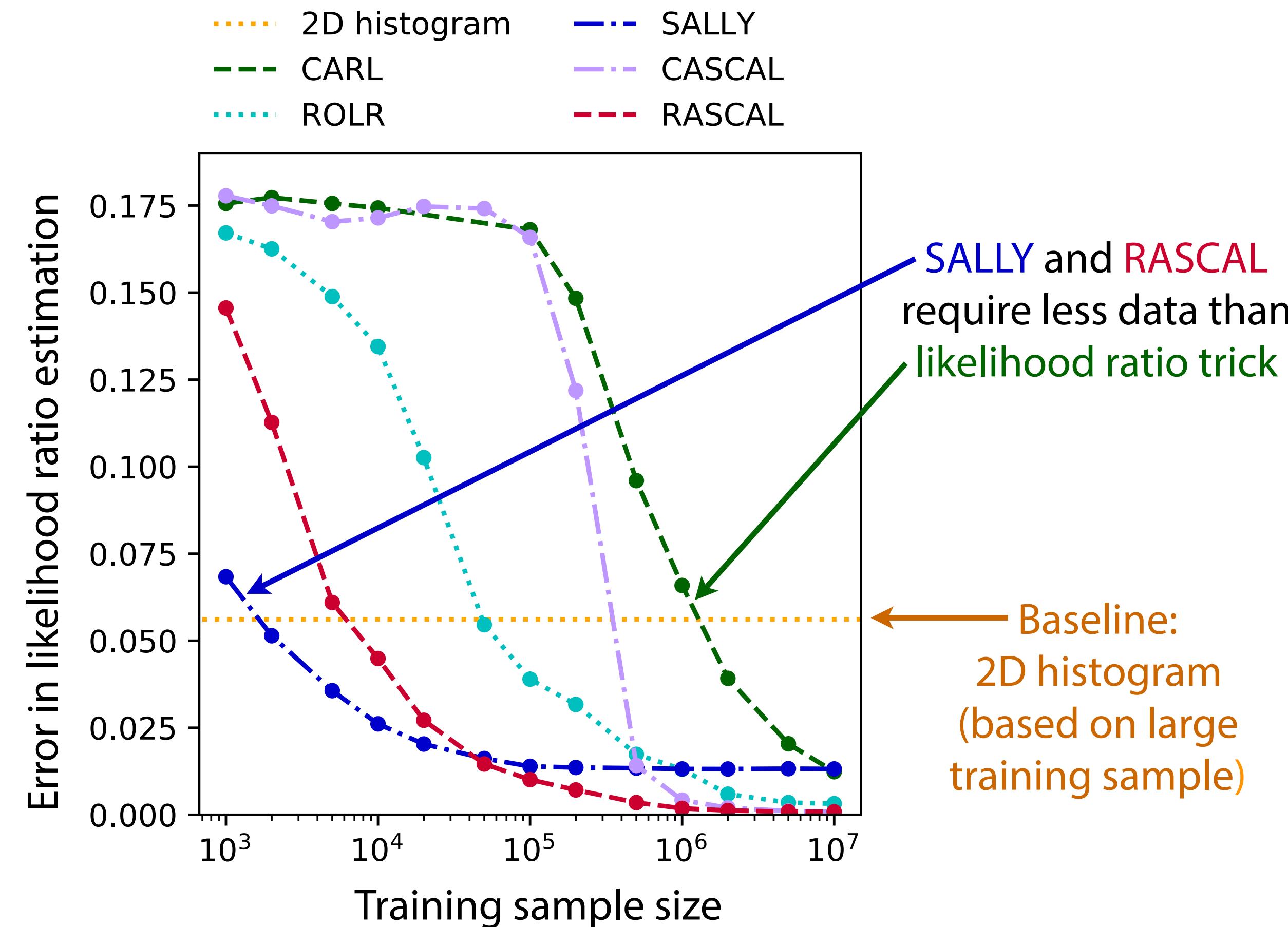
Goal: constrain the two EFT parameters

- new inference methods
- baseline: 2d histogram analysis of jet momenta & angular correlations

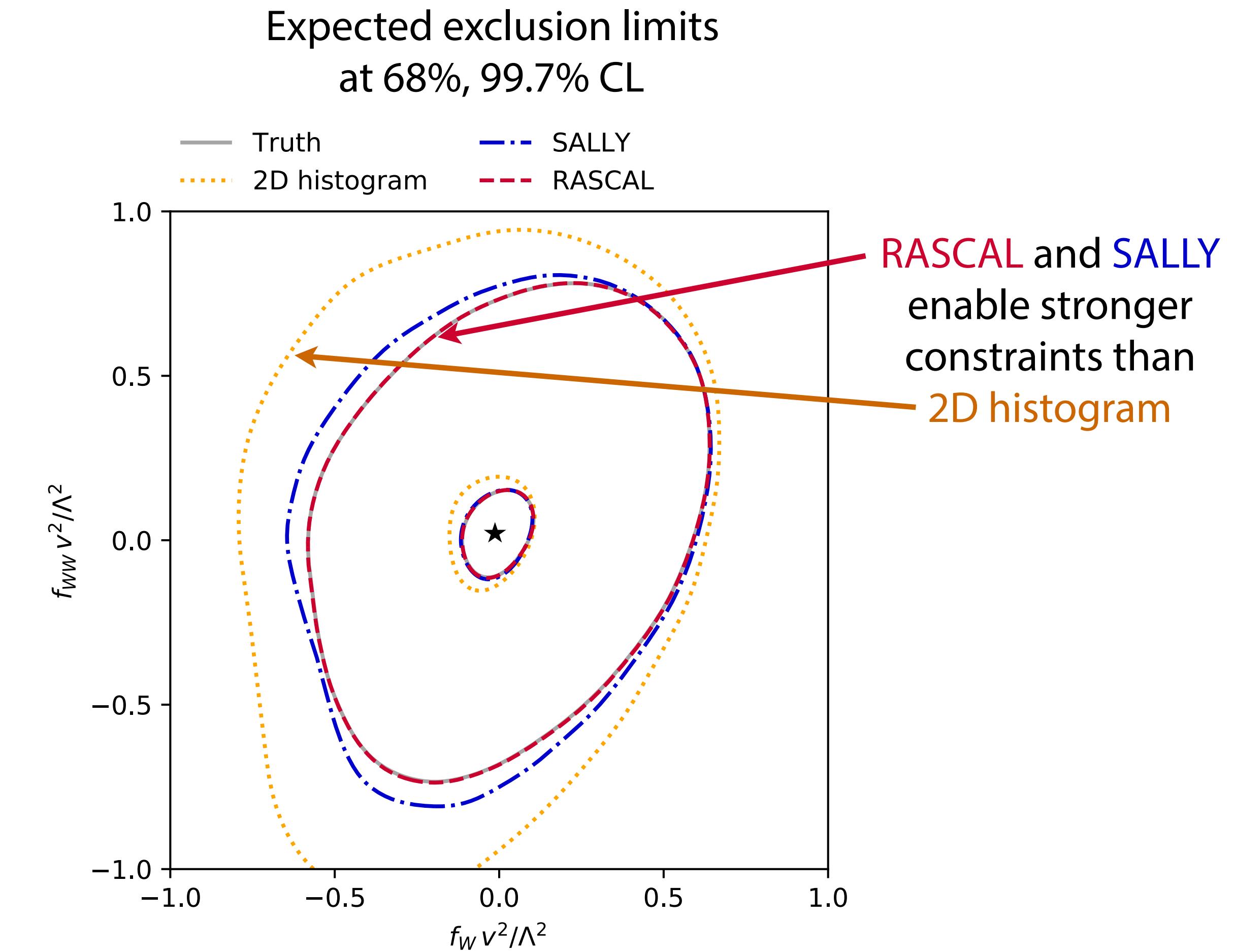
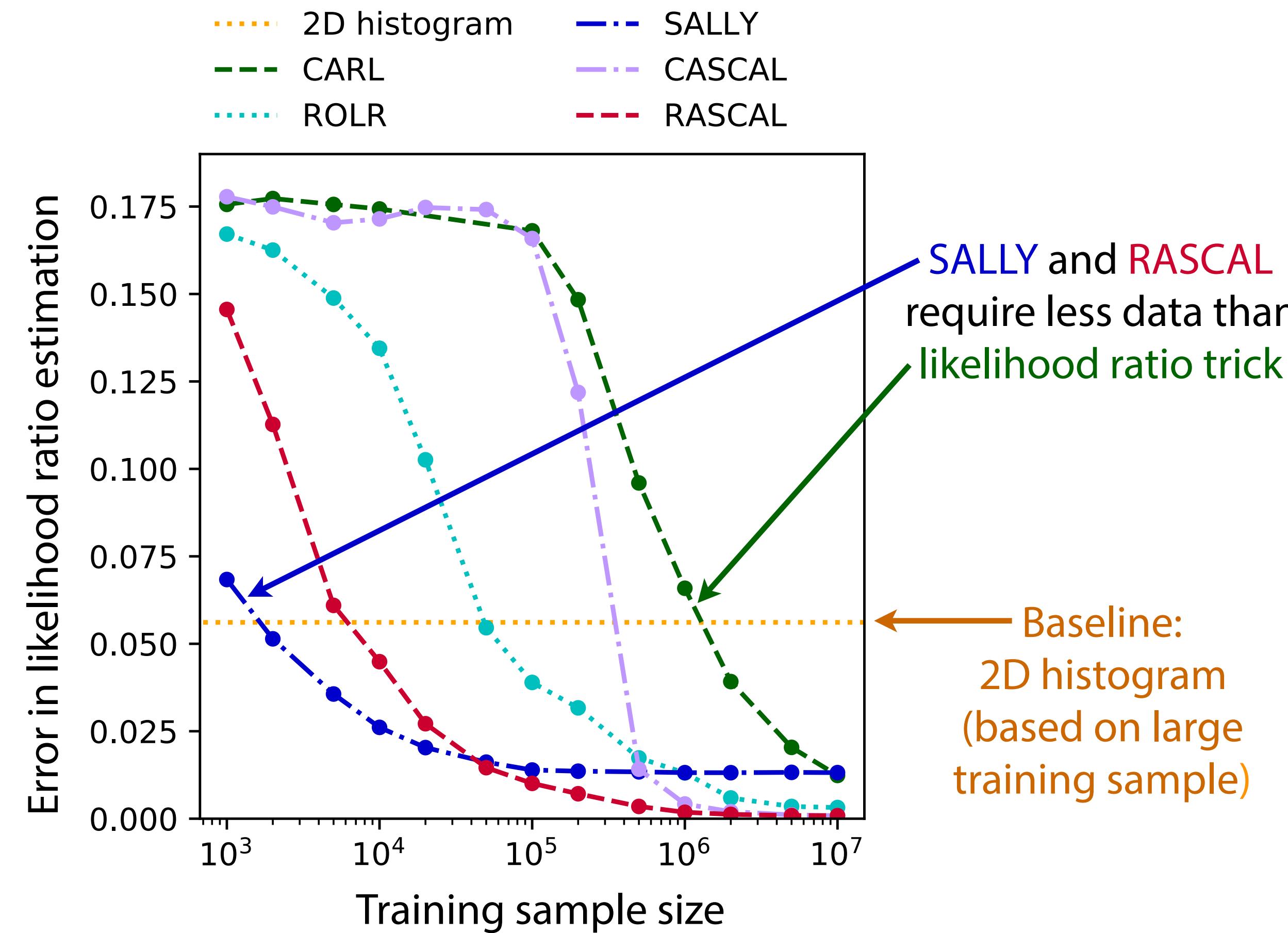
Two scenarios:

- Simplified setup in which we can compare to true likelihood
- “Realistic” simulation with approximate detector effects

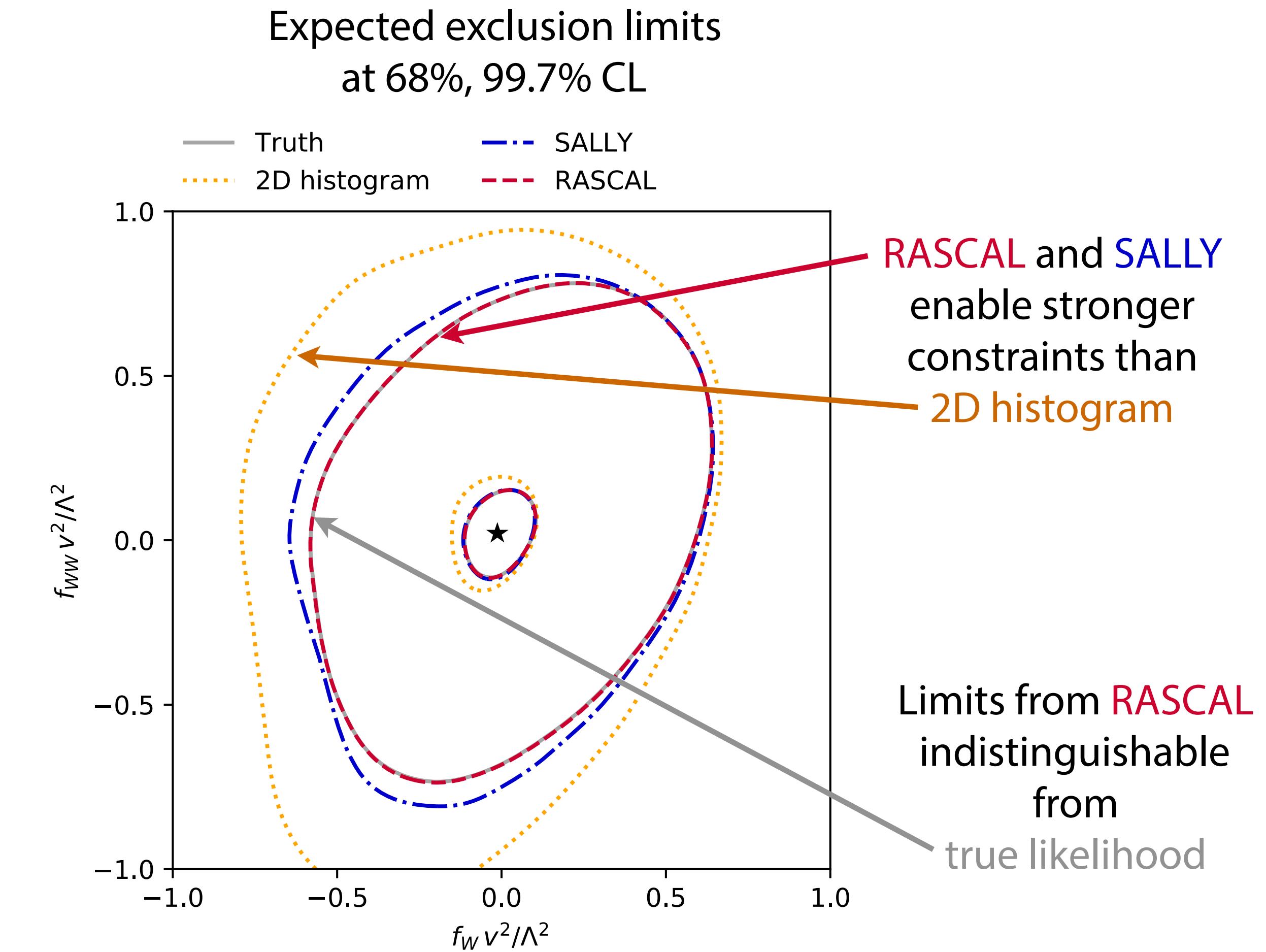
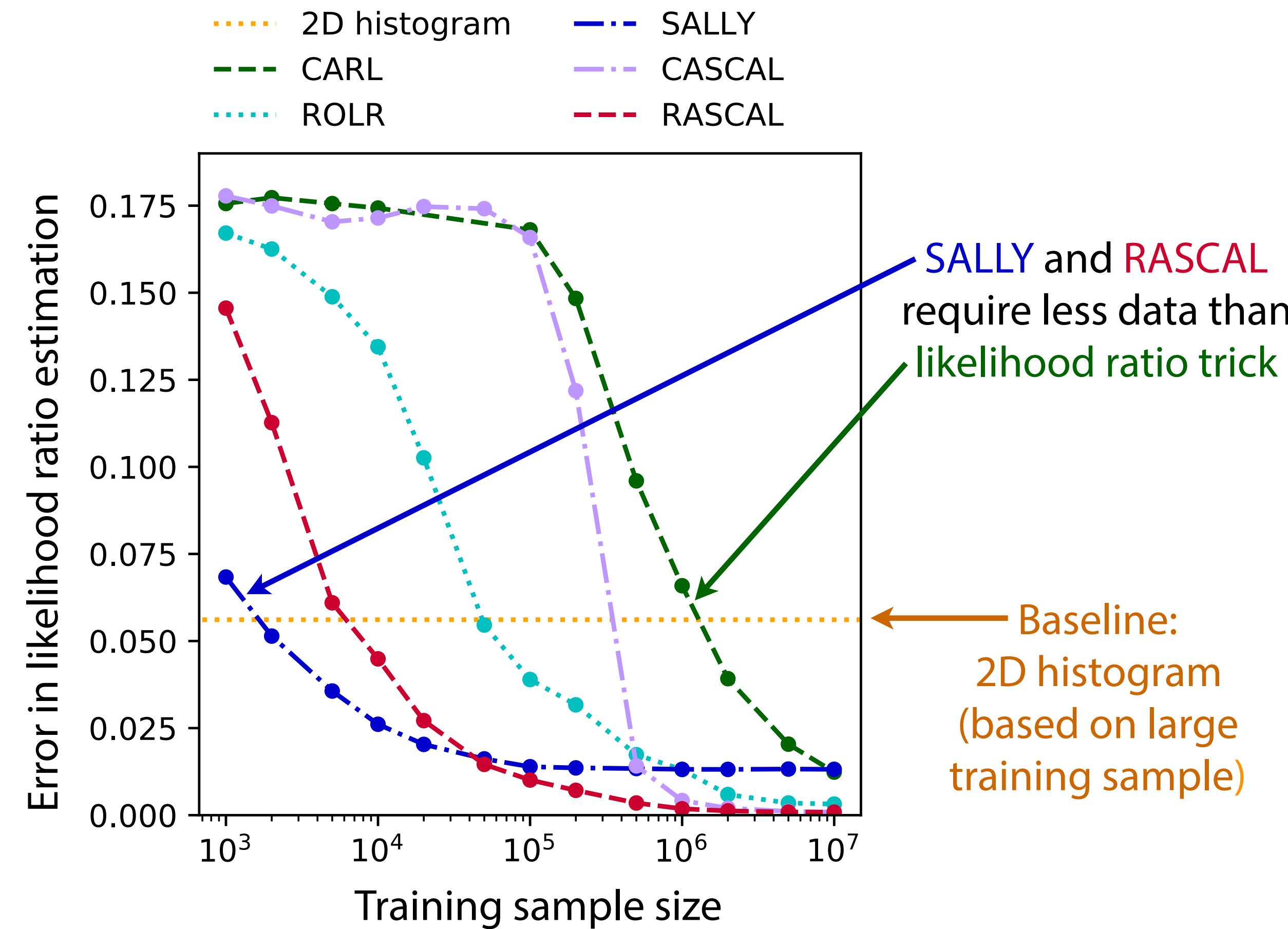
Proof of concept: Stronger constraints with less training data



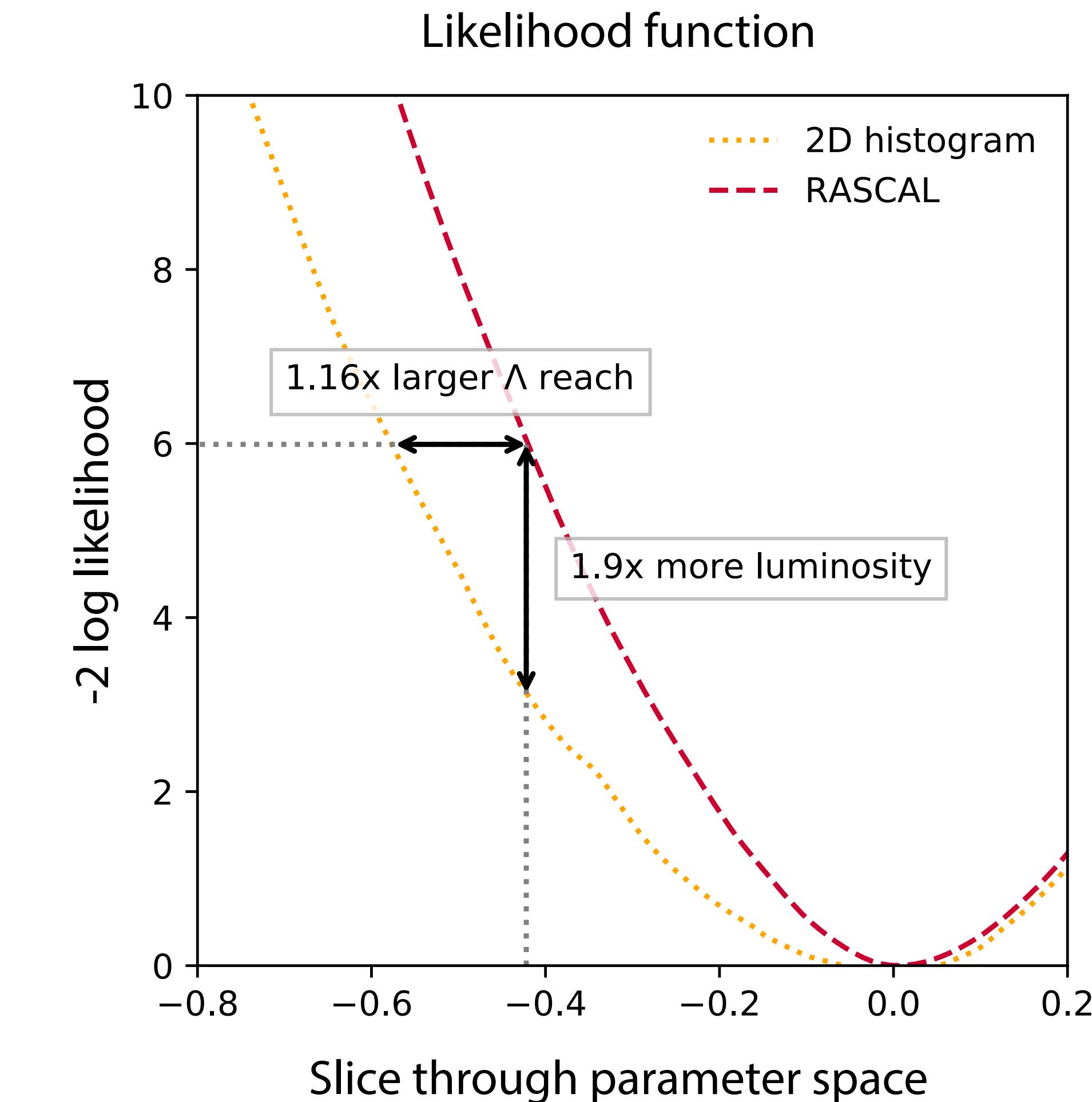
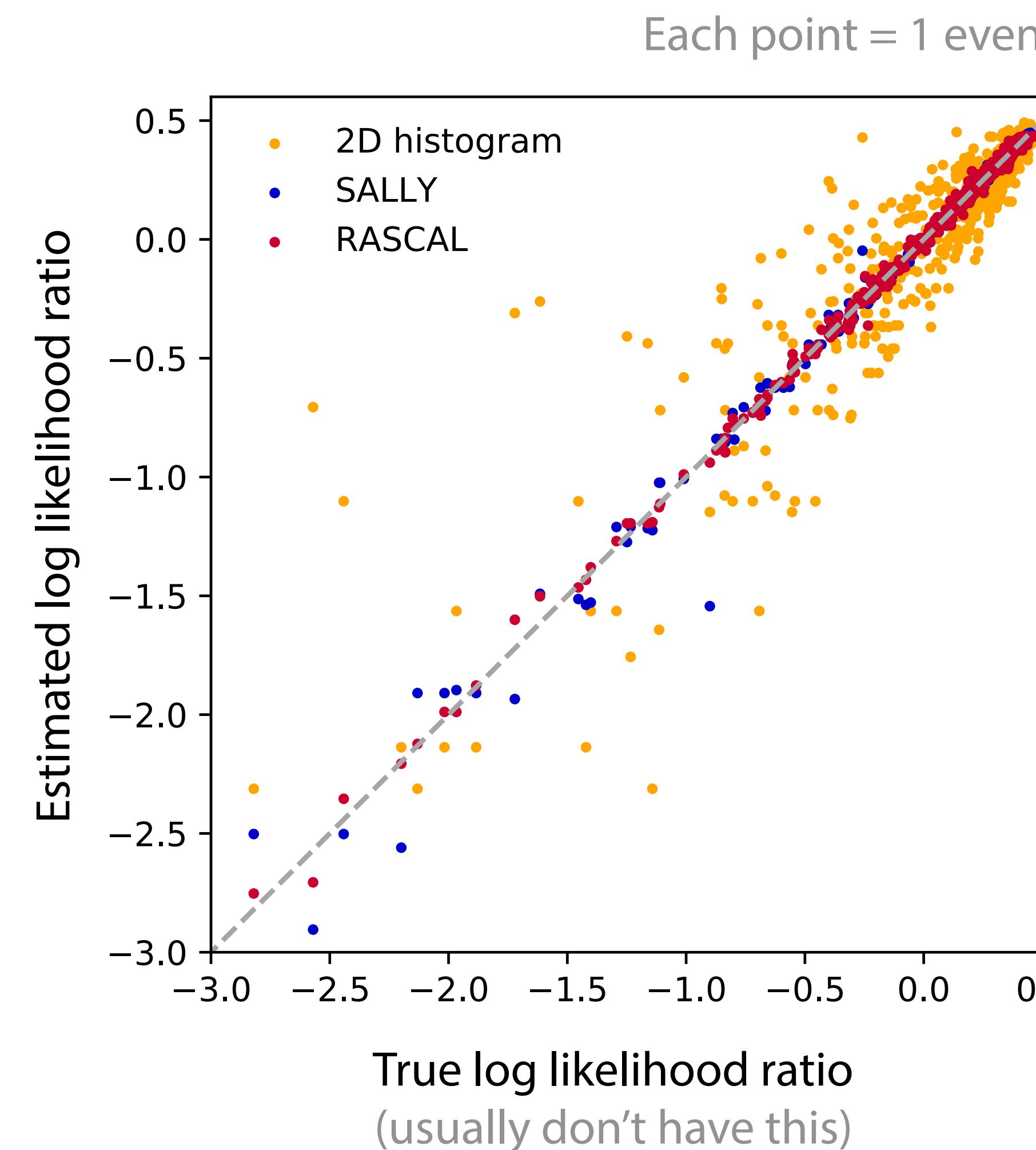
Proof of concept: Stronger constraints with less training data



Proof of concept: Stronger constraints with less training data



Proof of concept: 16% greater reach (~90% more data)



Constraining operators in ttH effectively

[JB, F. Kling, I. Espejo, K. Cranmer 1907.10621]

- Pheno-level analysis of

$$pp \rightarrow t\bar{t} h \rightarrow (b\ell^+) (\bar{b}\ell^-) (\gamma\gamma) E_T^{\text{miss}}$$

with MadGraph + Pythia + Delphes

- Inference on three EFT operators:

$$\mathcal{O}_u = -\frac{1}{v^2} (H^\dagger H) (H^\dagger \bar{Q}_L) u_R, \quad \mathcal{O}_G = \frac{g_s^2}{m_W^2} (H^\dagger H) G_{\mu\nu}^a G_a^{\mu\nu},$$

$$\mathcal{O}_{uG} = -\frac{4g_s}{m_W^2} y_u (H^\dagger \bar{Q}_L) \gamma^{\mu\nu} T_a u_R G_{\mu\nu}^a$$

Constraining operators in ttH effectively

[JB, F. Kling, I. Espejo, K. Cranmer 1907.10621]

- Pheno-level analysis of

$$pp \rightarrow t\bar{t} h \rightarrow (b\ell^+) (\bar{b}\ell^-) (\gamma\gamma) E_T^{\text{miss}}$$

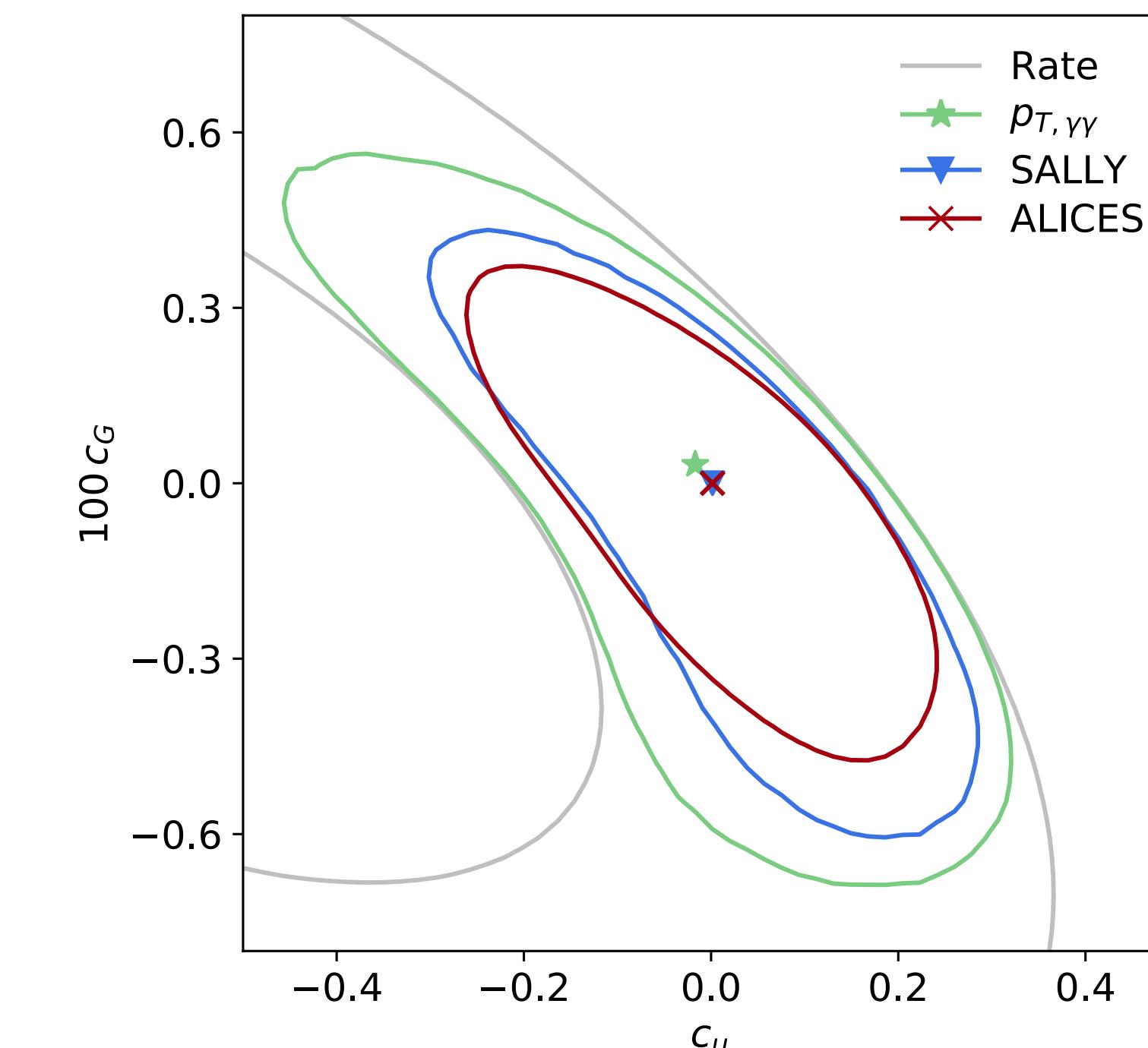
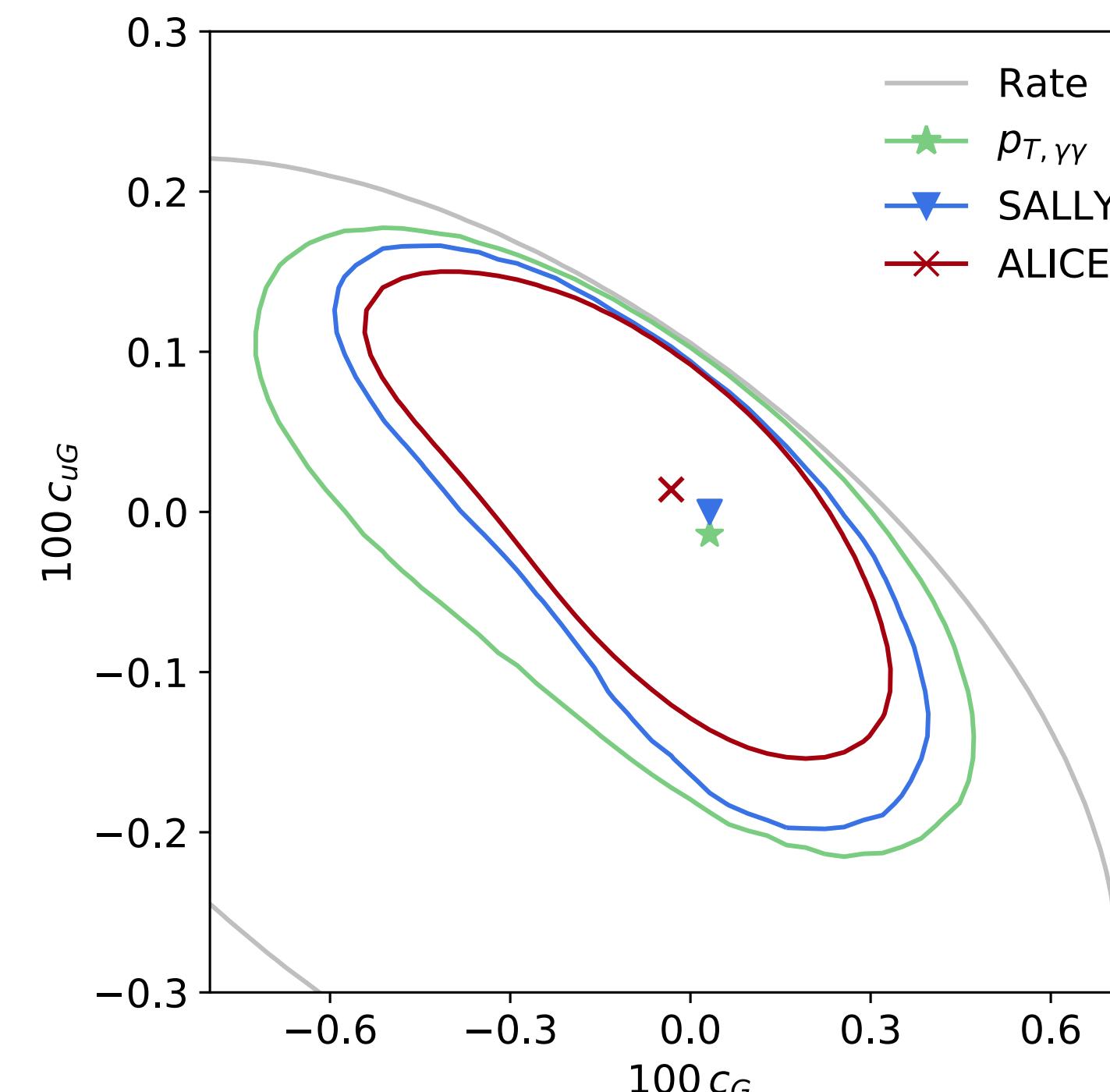
with MadGraph + Pythia + Delphes

- Inference on three EFT operators:

$$\mathcal{O}_u = -\frac{1}{v^2} (H^\dagger H) (H^\dagger \bar{Q}_L) u_R, \quad \mathcal{O}_G = \frac{g_s^2}{m_W^2} (H^\dagger H) G_{\mu\nu}^a G_a^{\mu\nu},$$

$$\mathcal{O}_{uG} = -\frac{4g_s}{m_W^2} y_u (H^\dagger \bar{Q}_L) \gamma^{\mu\nu} T_a u_R G_{\mu\nu}^a$$

- New **inference techniques** improve expected HL-LHC limits compared to **histogram baseline**:

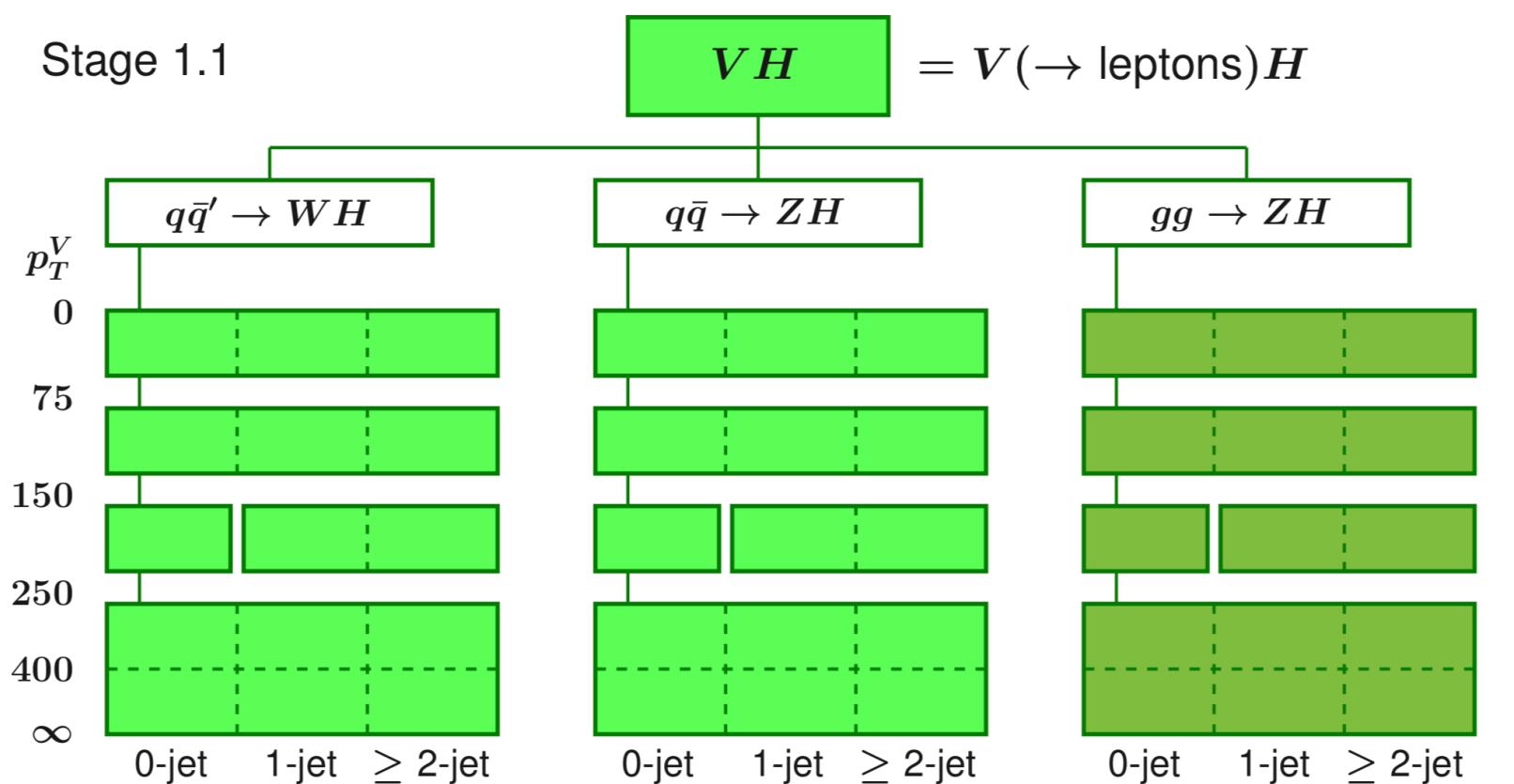


Benchmarking STXS in WH

[JB, S. Dawson, S. Homiller, F. Kling, T. Plehn 1908.06980]

- Simplified Template Cross-Sections (STXS) define observable bins that are supposed to capture as much information on NP as possible

[N. Berger et al. 1906.02754; HXSWG YR4]



- Let's check! How much information on

$$\tilde{\mathcal{O}}_{HD} = \mathcal{O}_{H\square} - \frac{\mathcal{O}_{HD}}{4} = (\phi^\dagger \phi) \square (\phi^\dagger \phi) - \frac{1}{4} (\phi^\dagger D^\mu \phi)^* (\phi^\dagger D_\mu \phi)$$

$$\mathcal{O}_{HW} = \phi^\dagger \phi W_{\mu\nu}^a W^{\mu\nu a}$$

$$\mathcal{O}_{Hq}^{(3)} = (\phi^\dagger i \overleftrightarrow{D}_\mu^a \phi) (\bar{Q}_L \sigma^a \gamma^\mu Q_L) ,$$

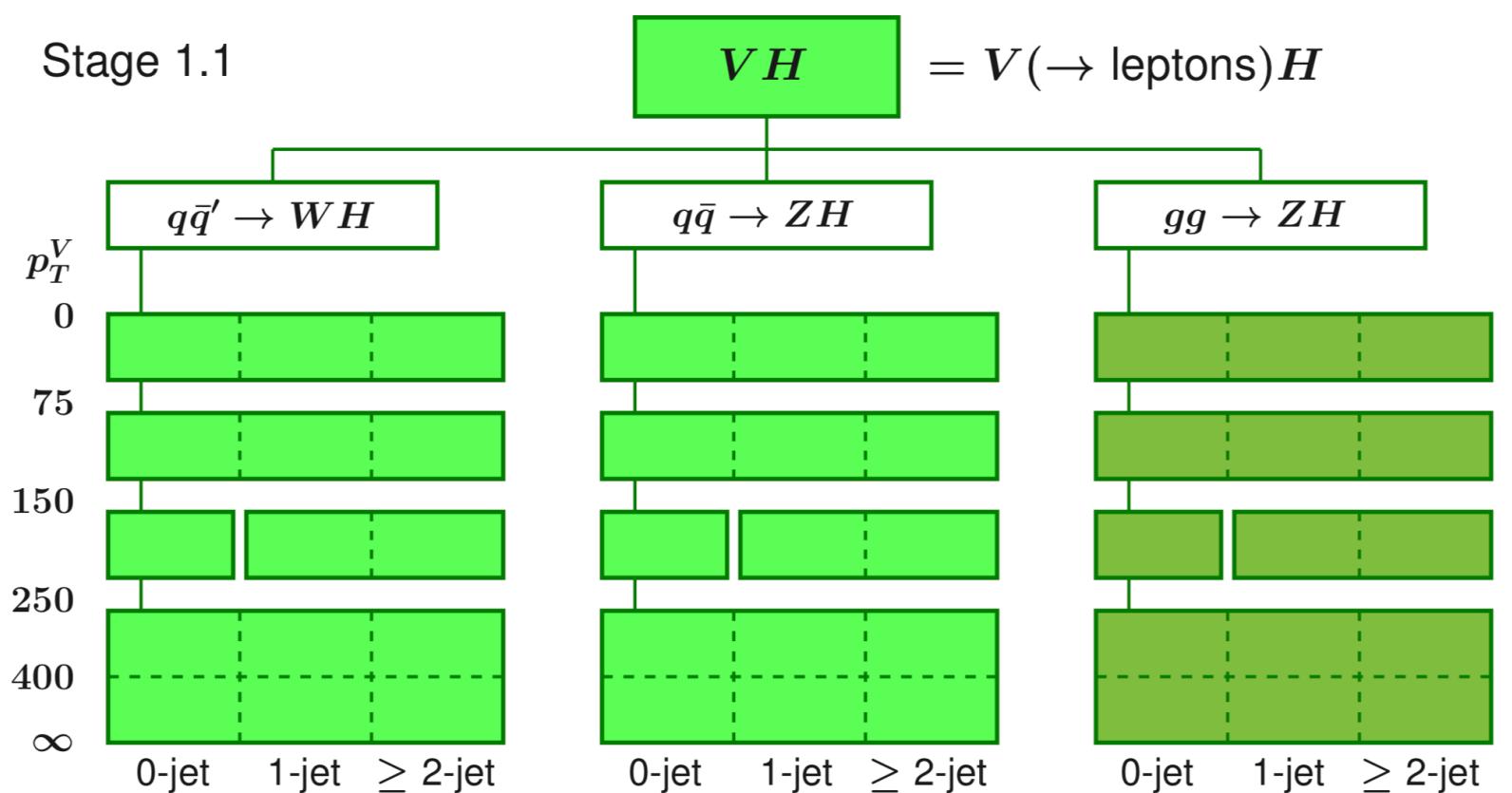
can we extract from $pp \rightarrow WH \rightarrow \ell\nu b\bar{b}$?

Benchmarking STXS in WH

[JB, S. Dawson, S. Homiller, F. Kling, T. Plehn 1908.06980]

- Simplified Template Cross-Sections (STXS) define observable bins that are supposed to capture as much information on NP as possible

[N. Berger et al. 1906.02754; HXSWG YR4]



- Let's check! How much information on

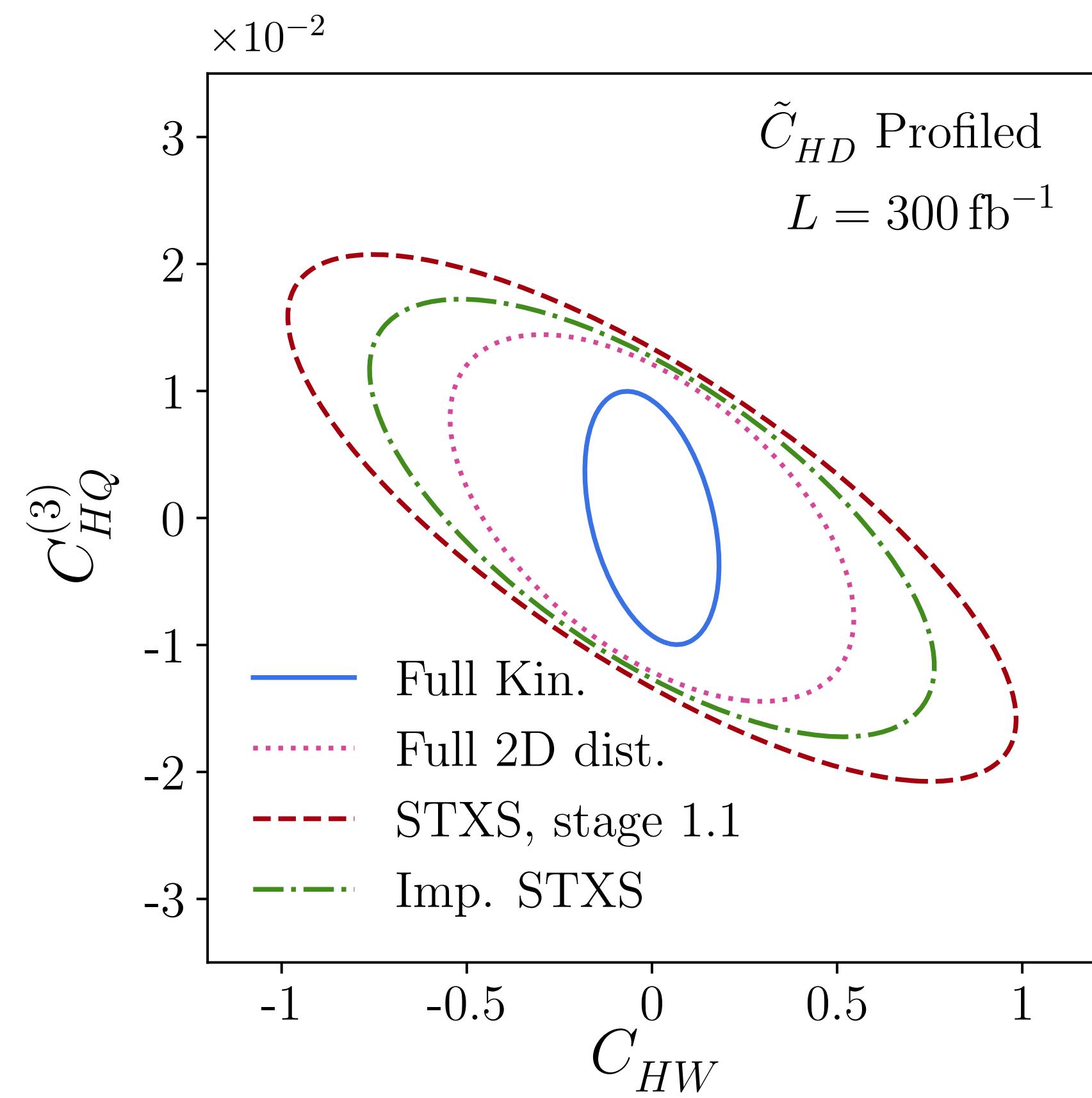
$$\tilde{\mathcal{O}}_{HD} = \mathcal{O}_{H\Box} - \frac{\mathcal{O}_{HD}}{4} = (\phi^\dagger \phi) \Box (\phi^\dagger \phi) - \frac{1}{4} (\phi^\dagger D^\mu \phi)^* (\phi^\dagger D_\mu \phi)$$

$$\mathcal{O}_{HW} = \phi^\dagger \phi W_{\mu\nu}^a W^{\mu\nu a}$$

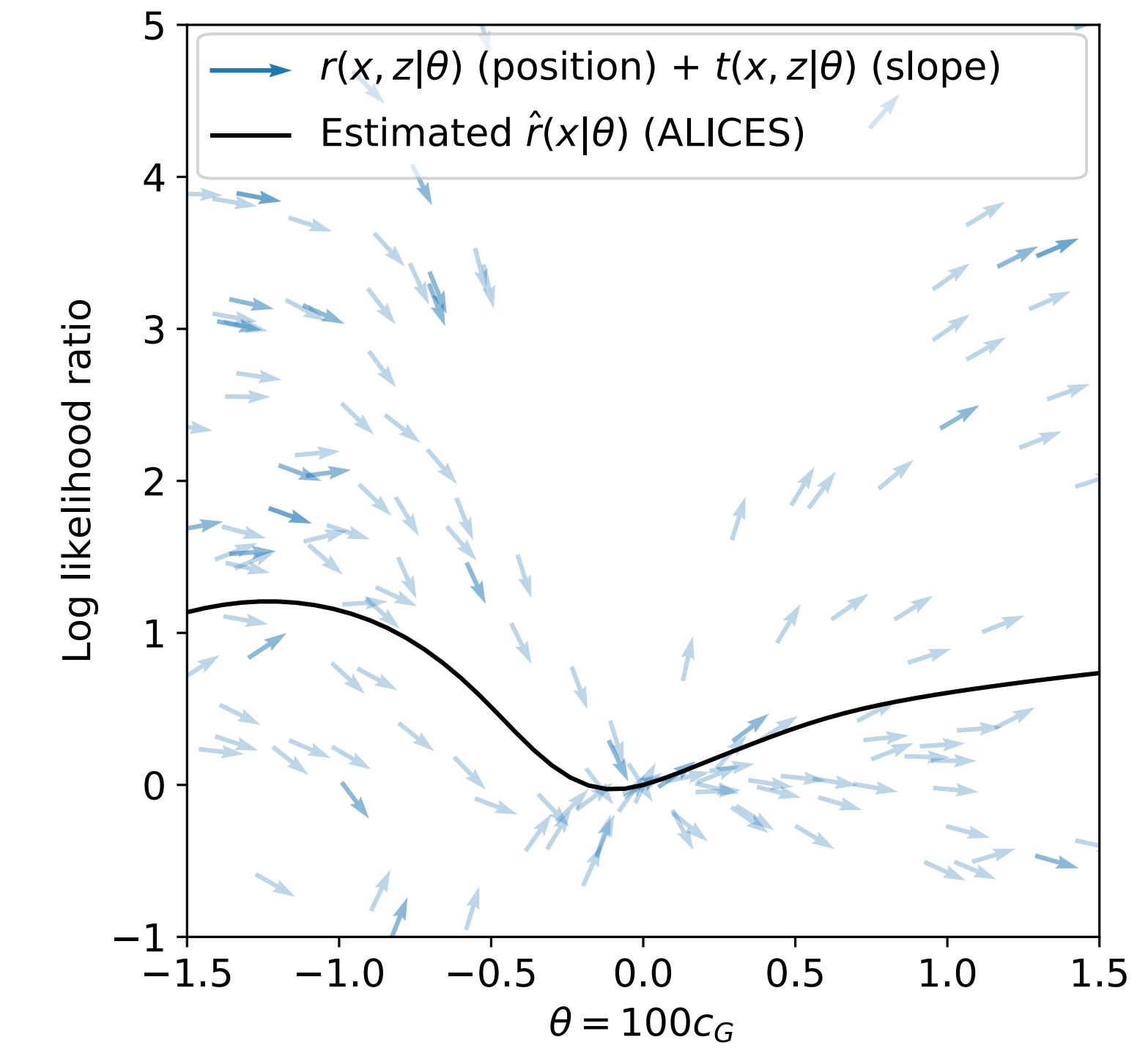
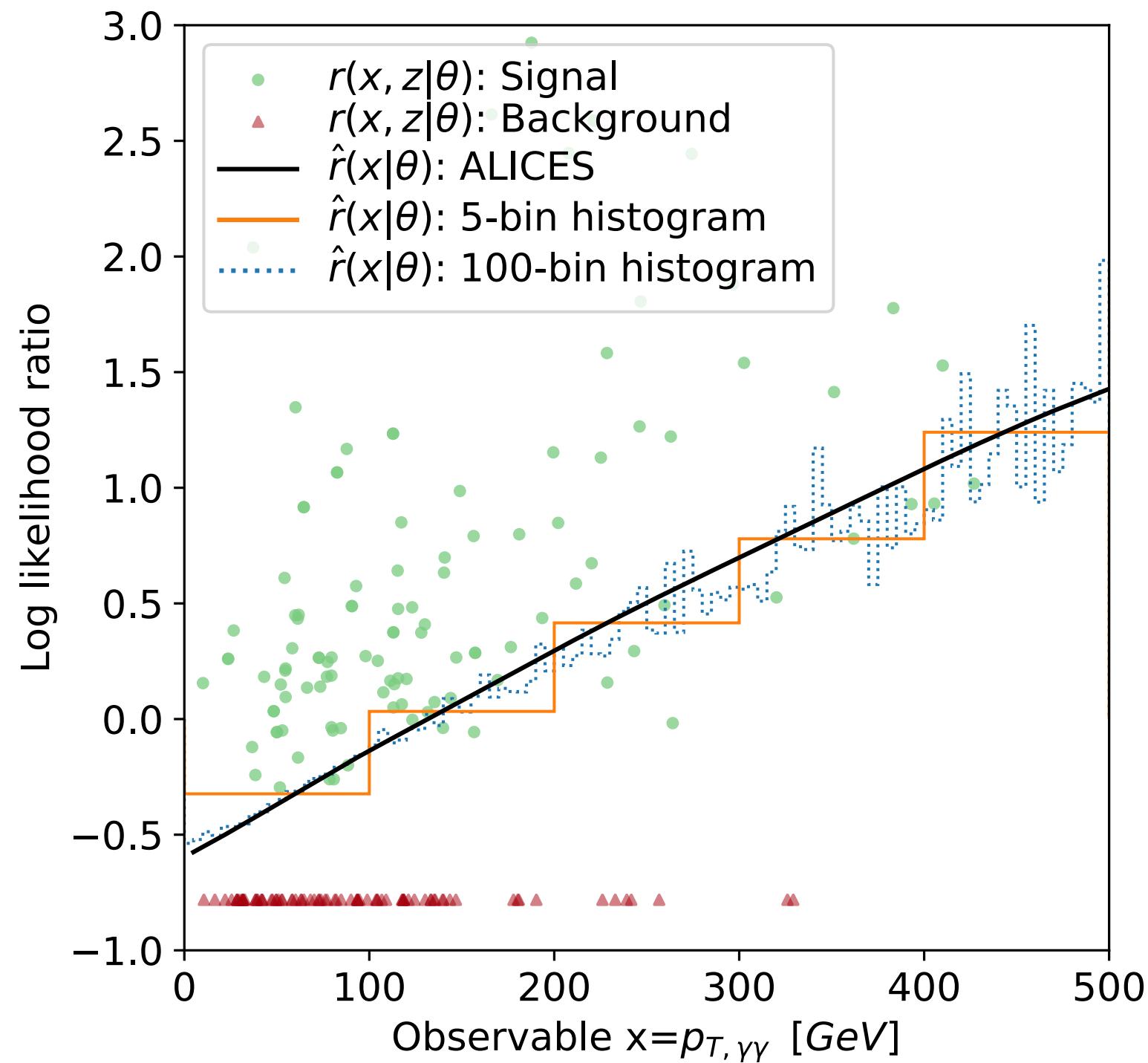
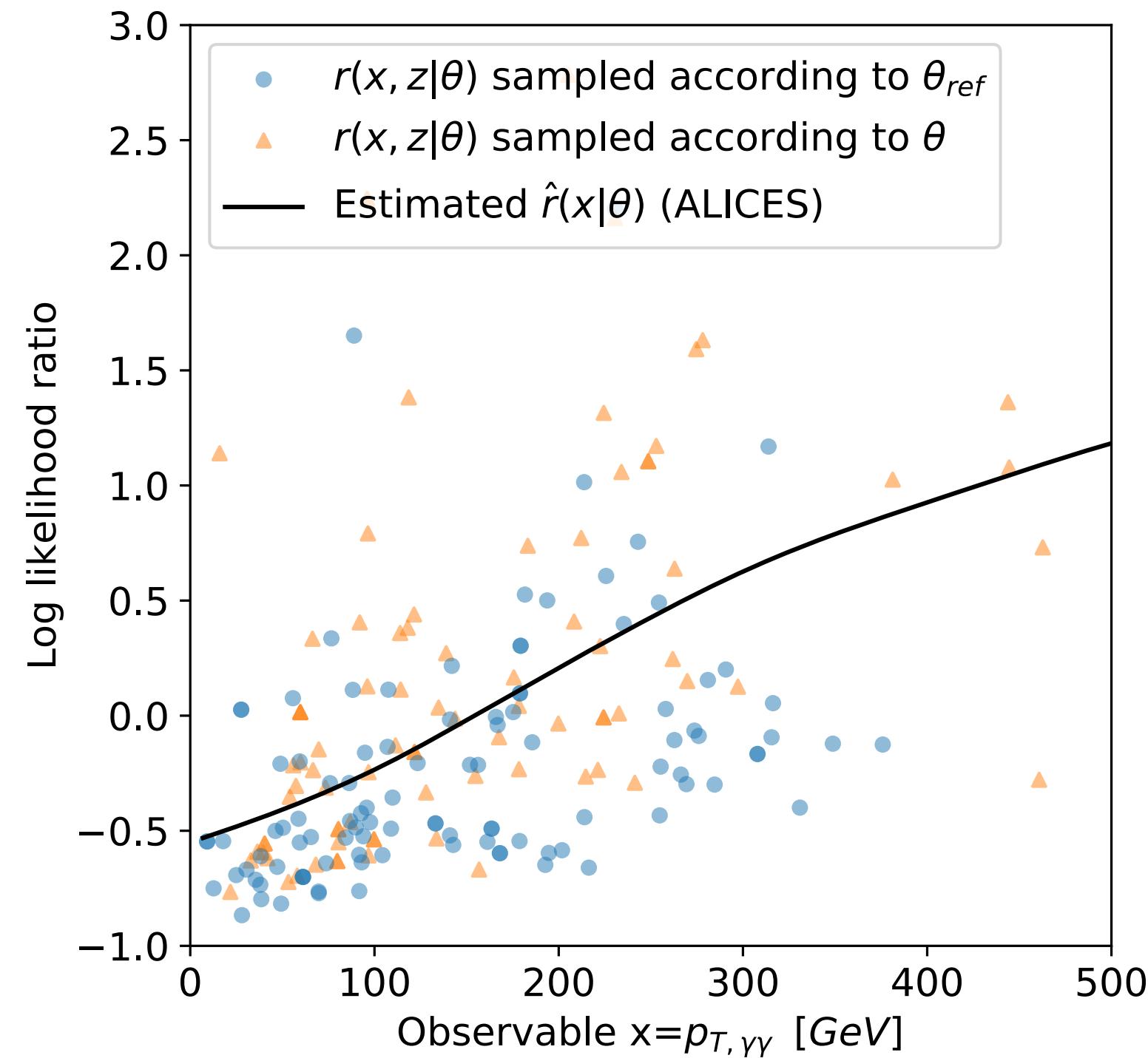
$$\mathcal{O}_{Hq}^{(3)} = (\phi^\dagger i \overleftrightarrow{D}_\mu^a \phi) (\bar{Q}_L \sigma^a \gamma^\mu Q_L) ,$$

can we extract from $pp \rightarrow WH \rightarrow \ell\nu b\bar{b}$?

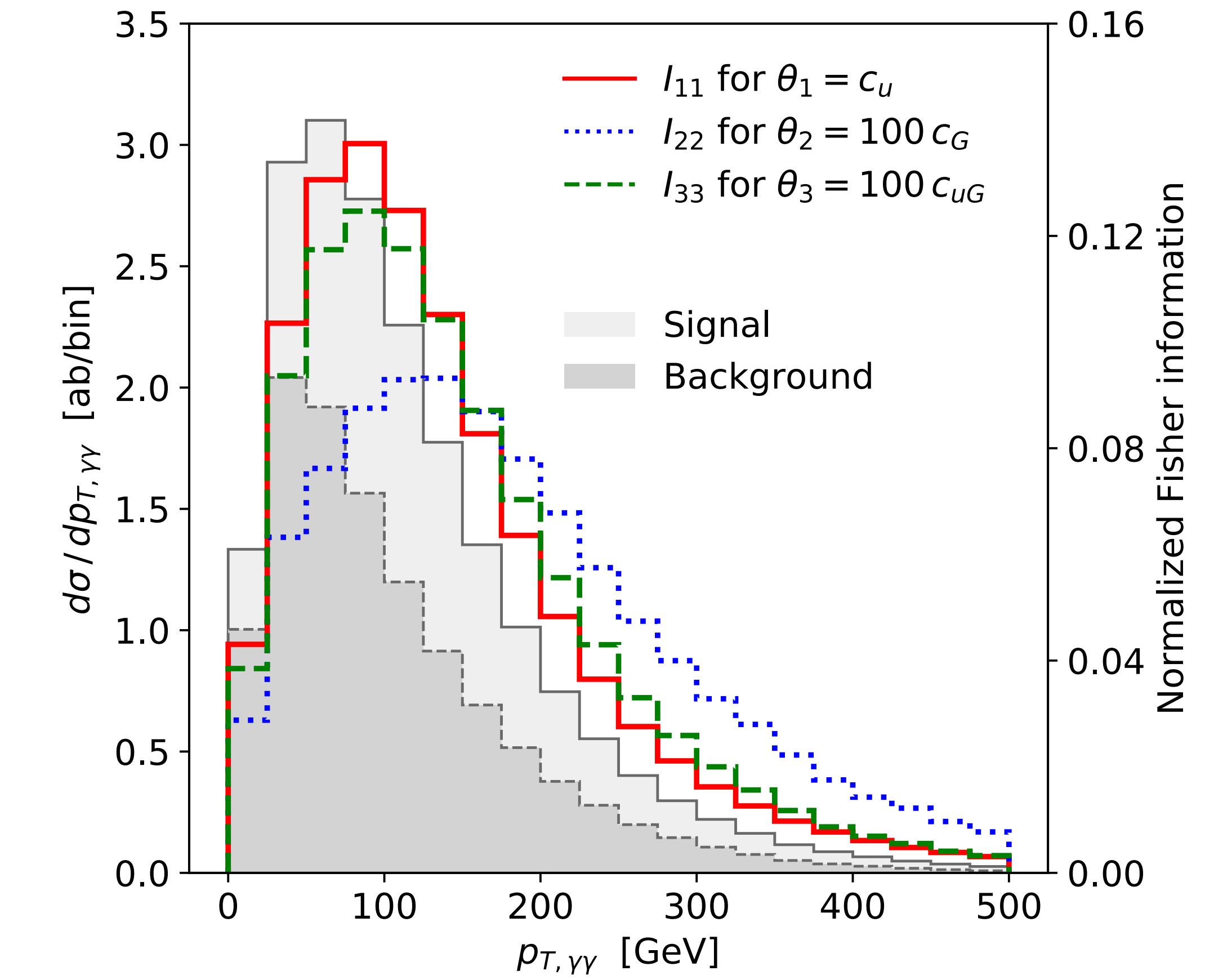
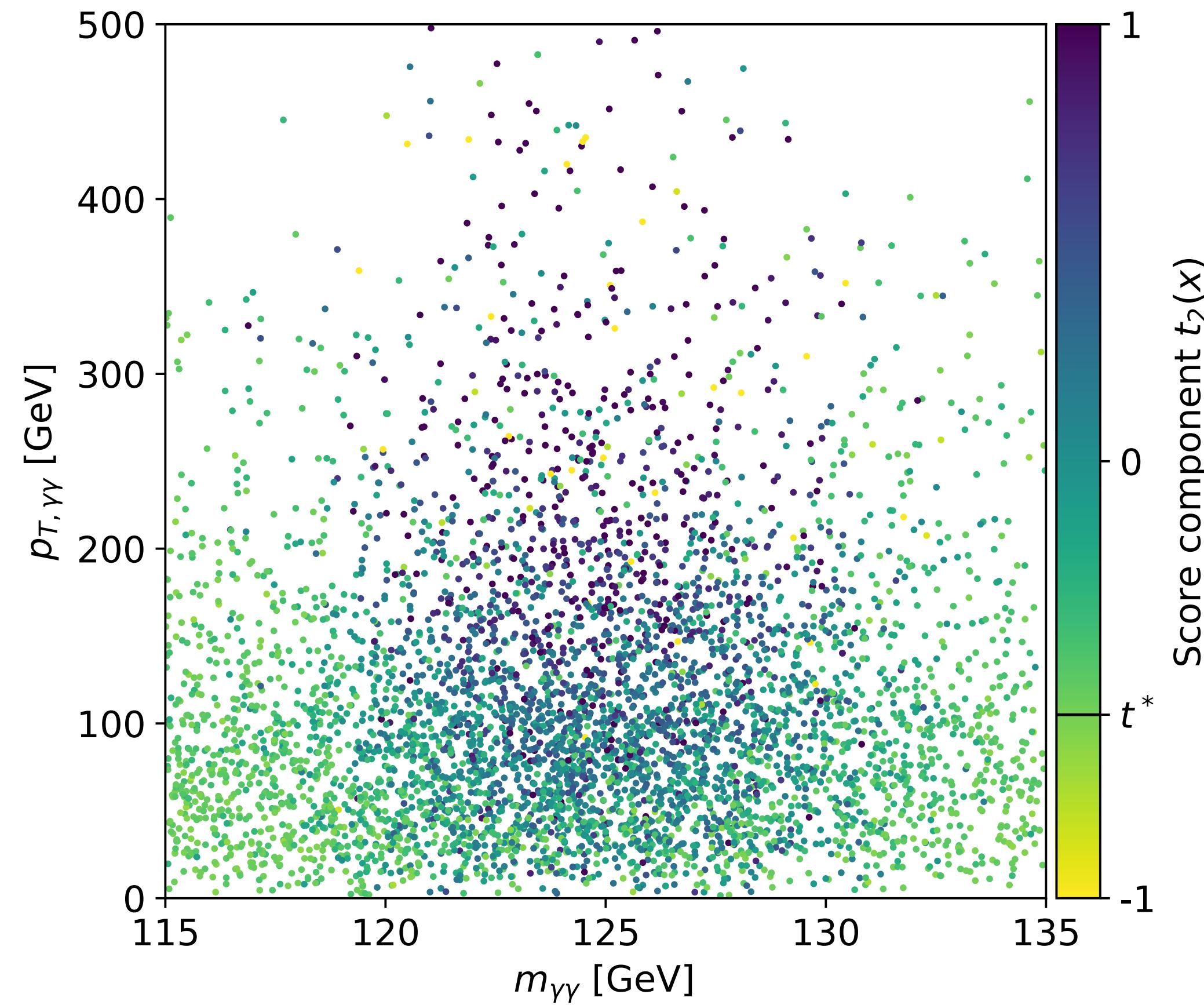
- Results: STXS are indeed sensitive to operators, adding a few more bins improve them, but a multivariate analysis is still stronger



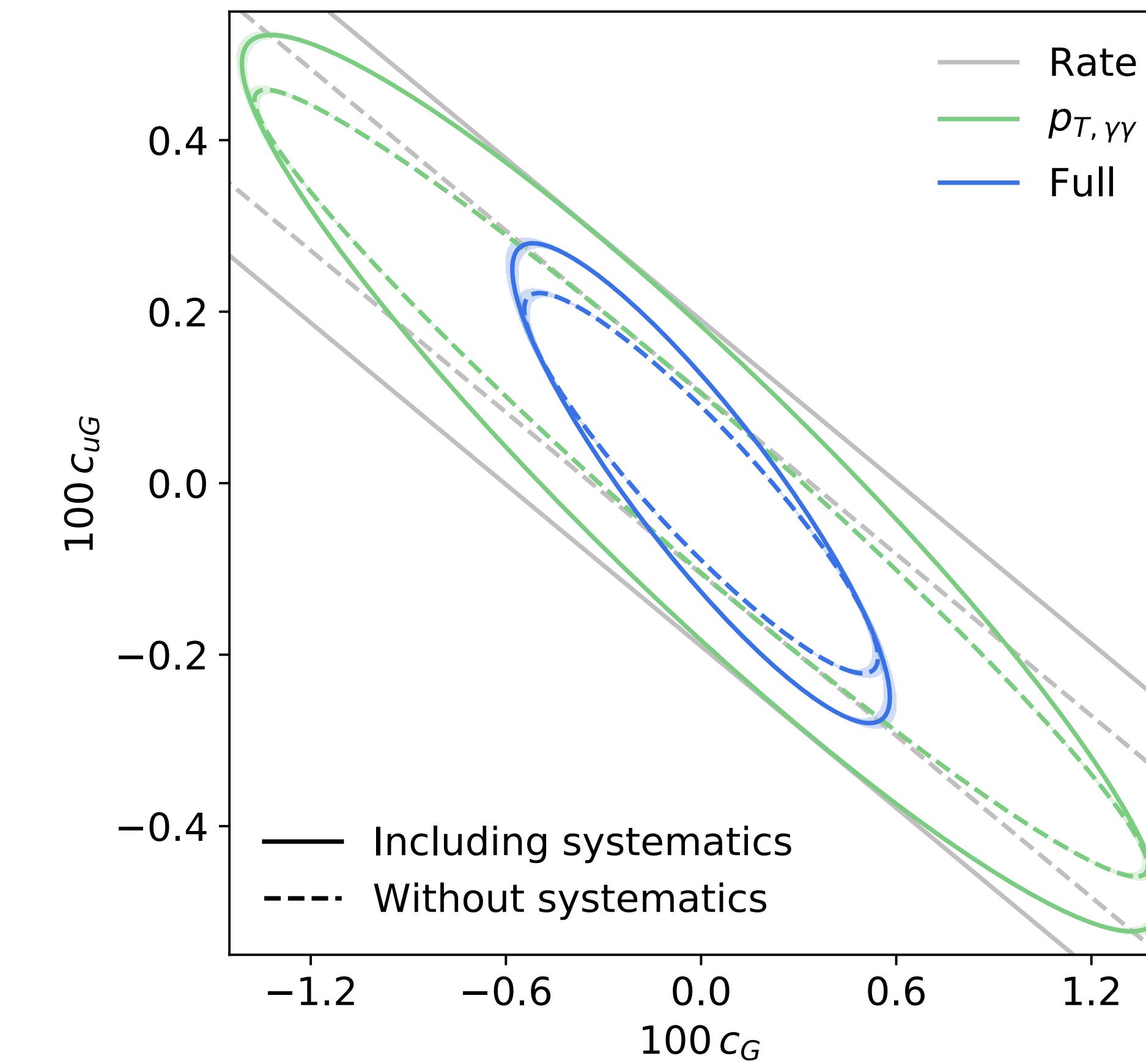
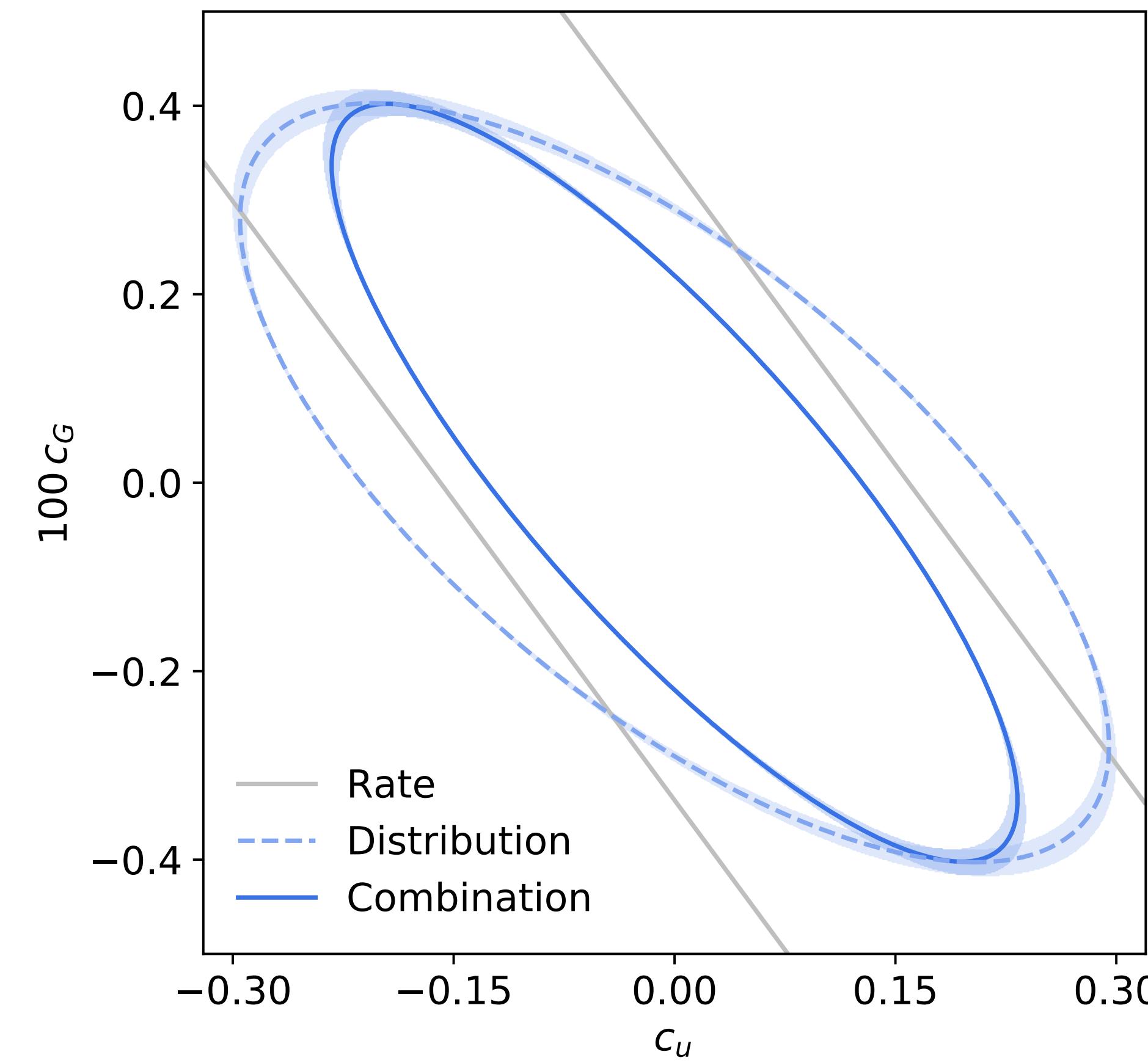
ttH: 1D illustration



ttH: score and information vs pT and m

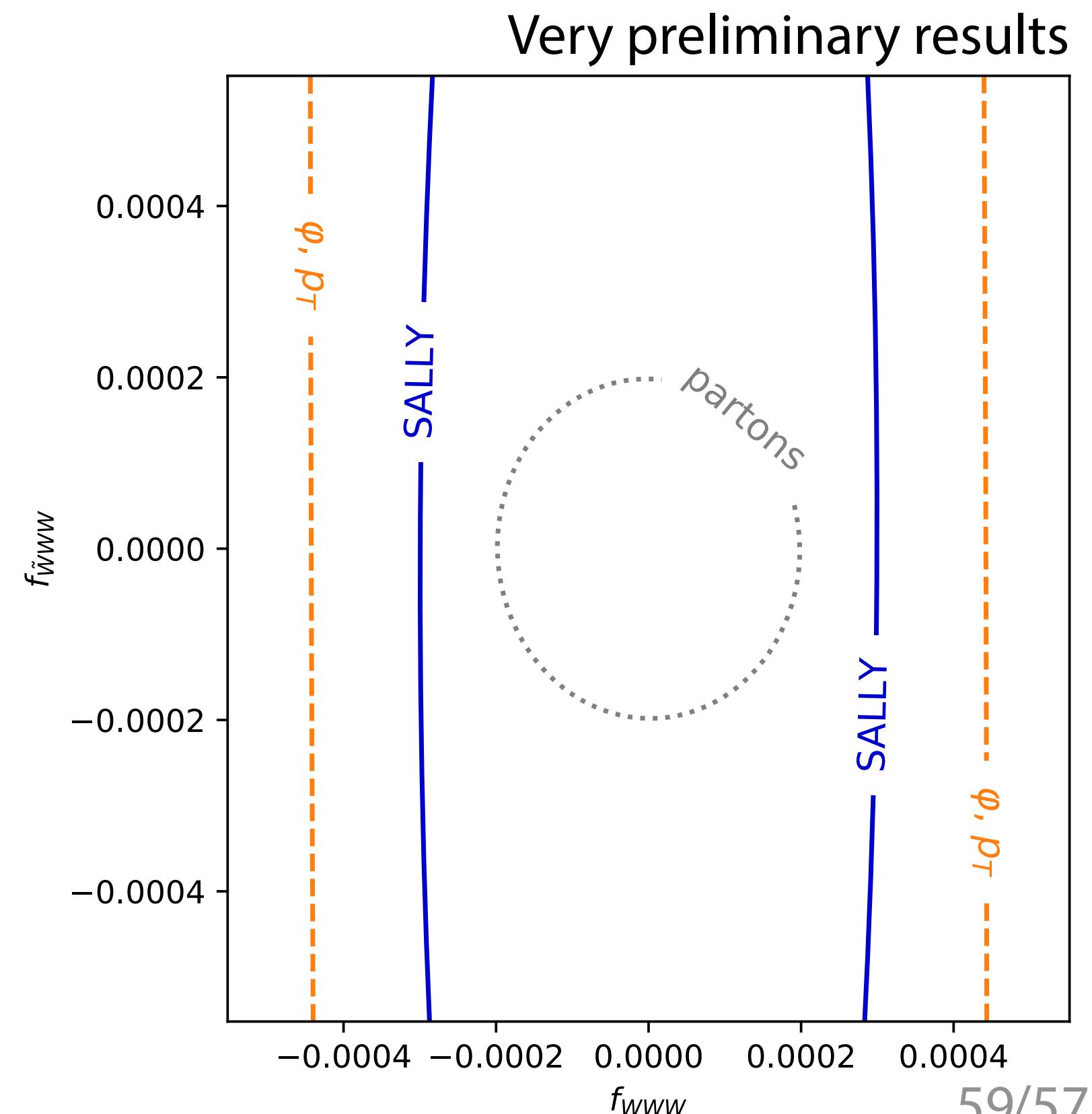
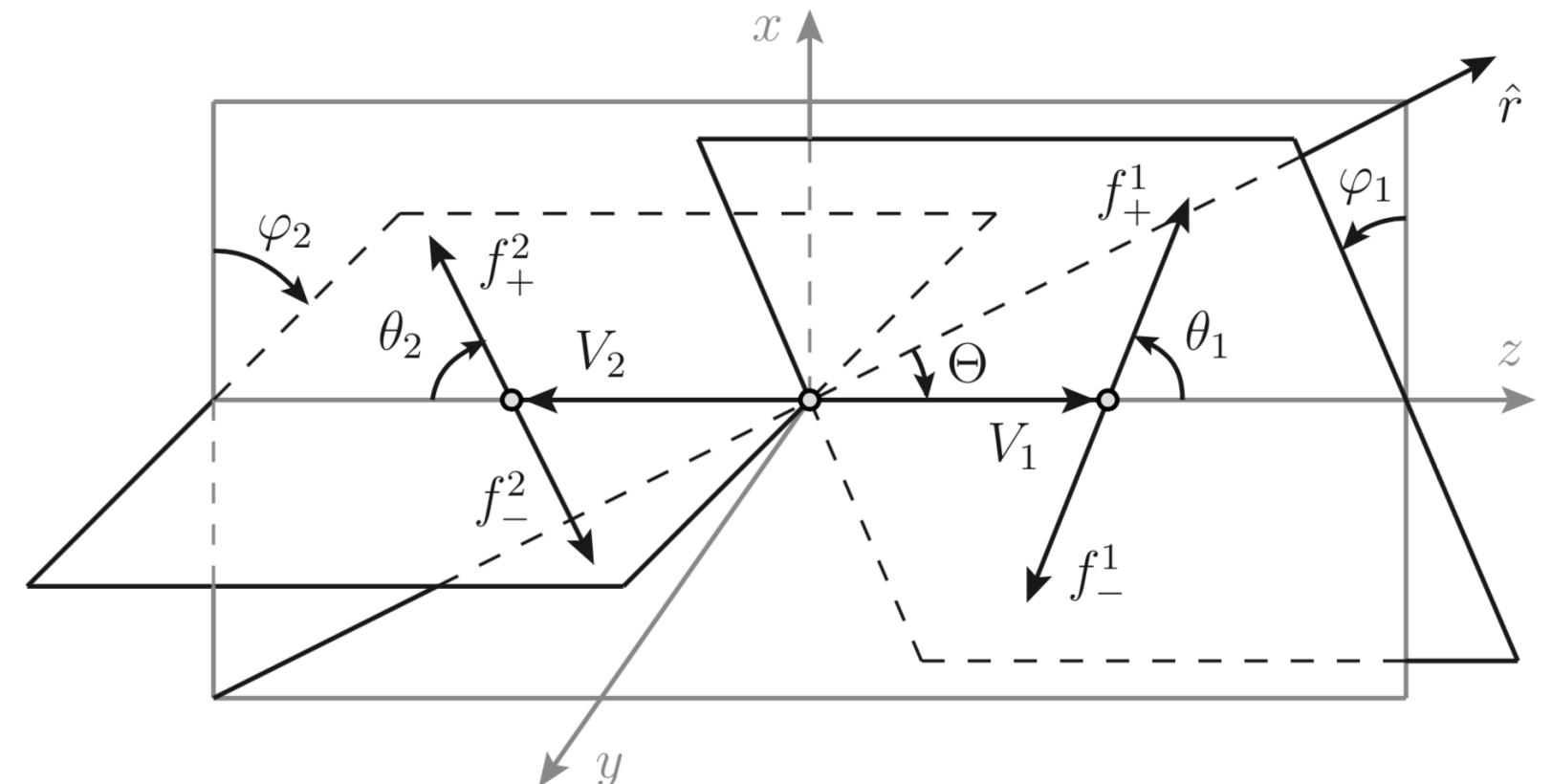


ttH: more information results



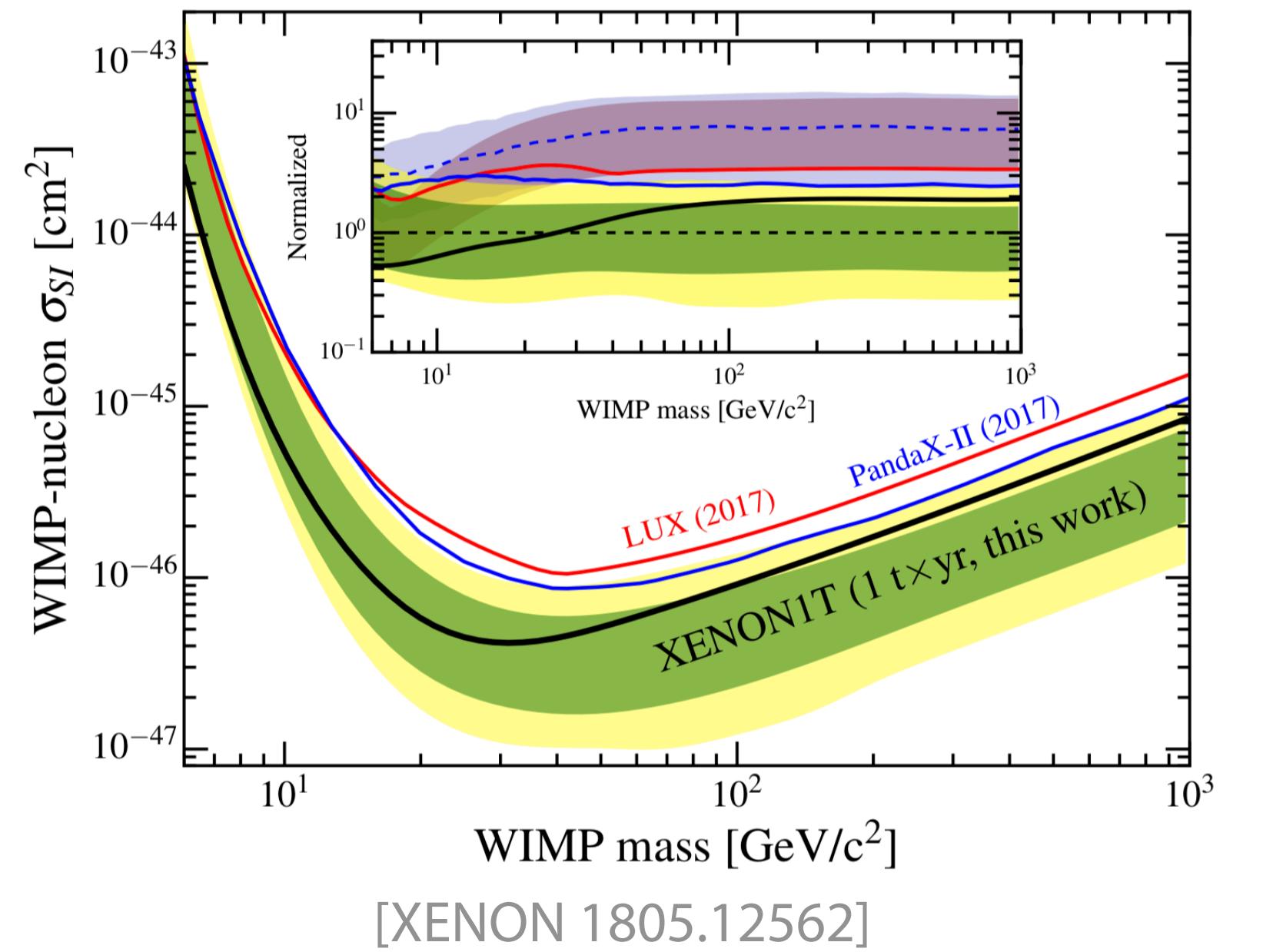
Diboson production

- In inclusive observables, the interference between SM and new physics amplitudes vanishes
⇒ Reduced sensitivity to new physics
- “Diboson interference resurrection”: an **angular variable** φ can be constructed to be sensitive to this interference
[G. Panico, F. Riva, A. Wulzer 1708.07823;
A. Azatov, D. Barducci, E. Venturini 1901.04821]
- We test the ML approach in EFT measurements in $W\gamma \rightarrow \ell\nu \gamma$
[JB, K. Cranmer, M. Farina, F. Kling, D. Pappadopulo, J. Ruderman in progress]
- Preliminary results: we can extract more information when we **analyze events with SALLY** than with **histograms of φ** and standard observables



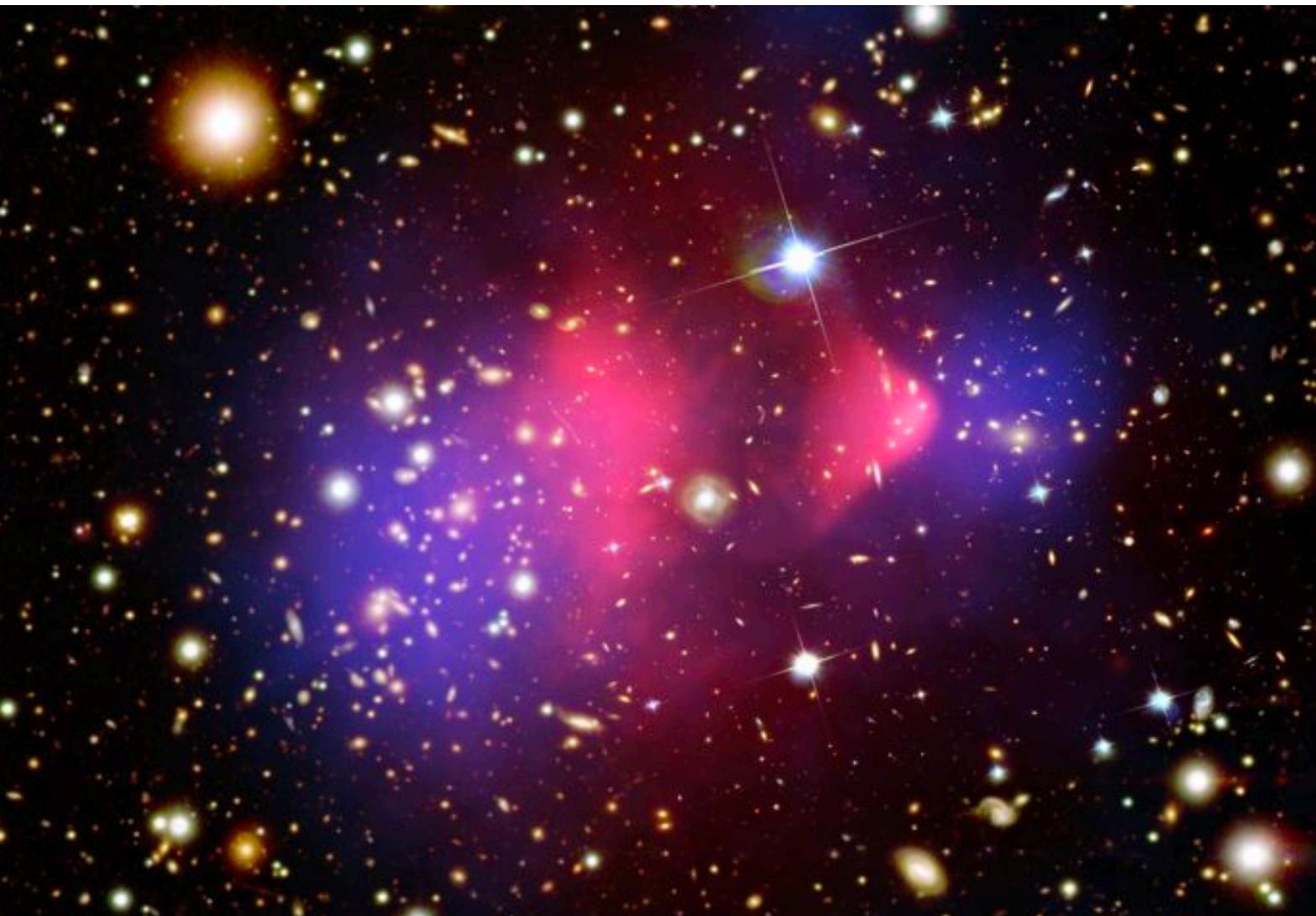
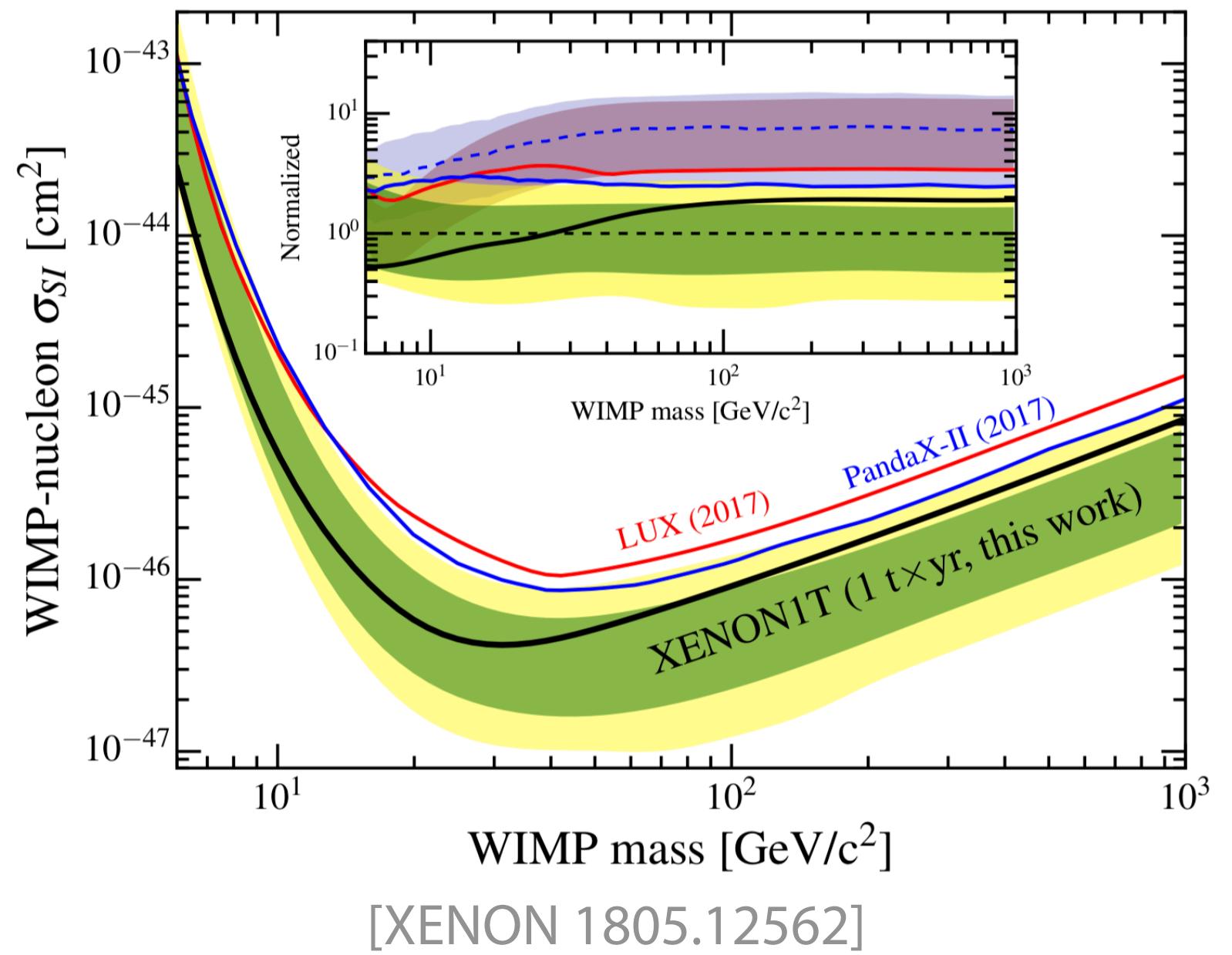
These methods work beyond the LHC.

Dark matter (DM) status



Non-gravitational interactions:
no evidence so far

Dark matter (DM) status

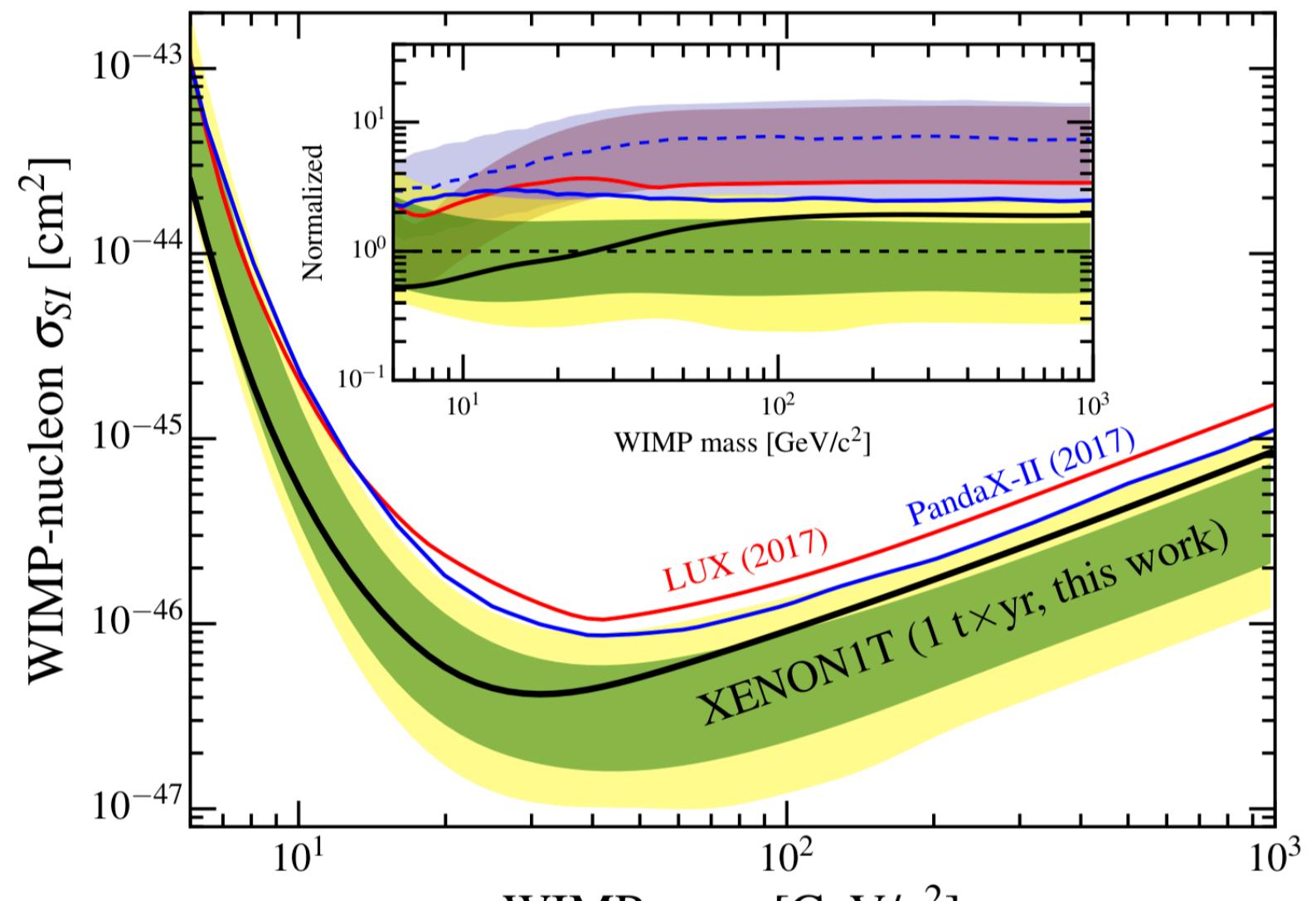


[NASA]

Non-gravitational interactions:
no evidence so far

Gravitational interactions at **large length scales**: plenty of evidence, consistent with Λ CDM

Dark matter (DM) status



[XENON 1805.12562]



[NASA]



[T. Brown, J.Tumlinson]

Non-gravitational interactions:
no evidence so far

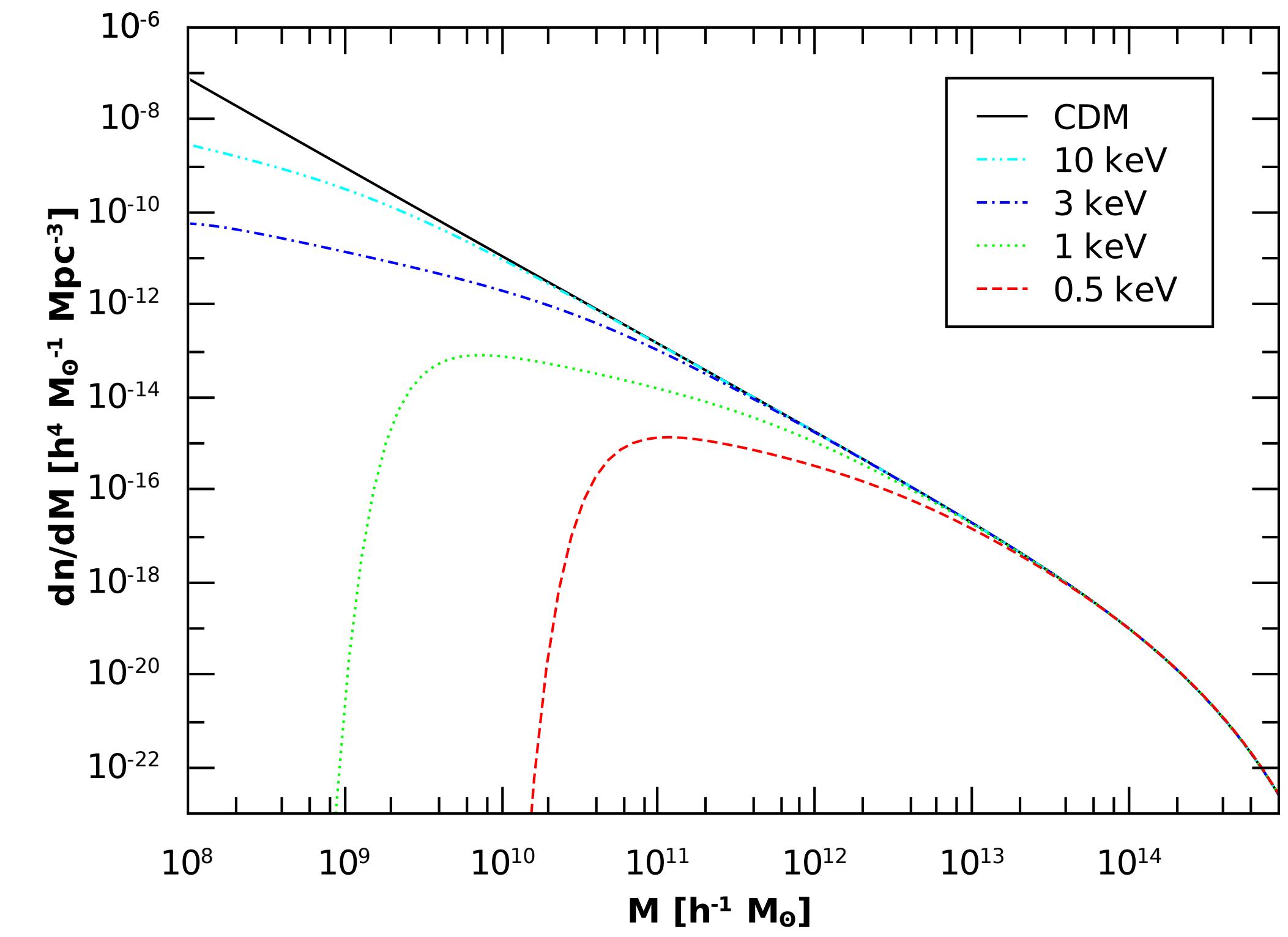
Gravitational interactions at **large length scales**: plenty of evidence, consistent with Λ CDM

Gravitational interactions at **small scales (DM substructure)**: beginning to be probed, many models predict deviations from Λ CDM

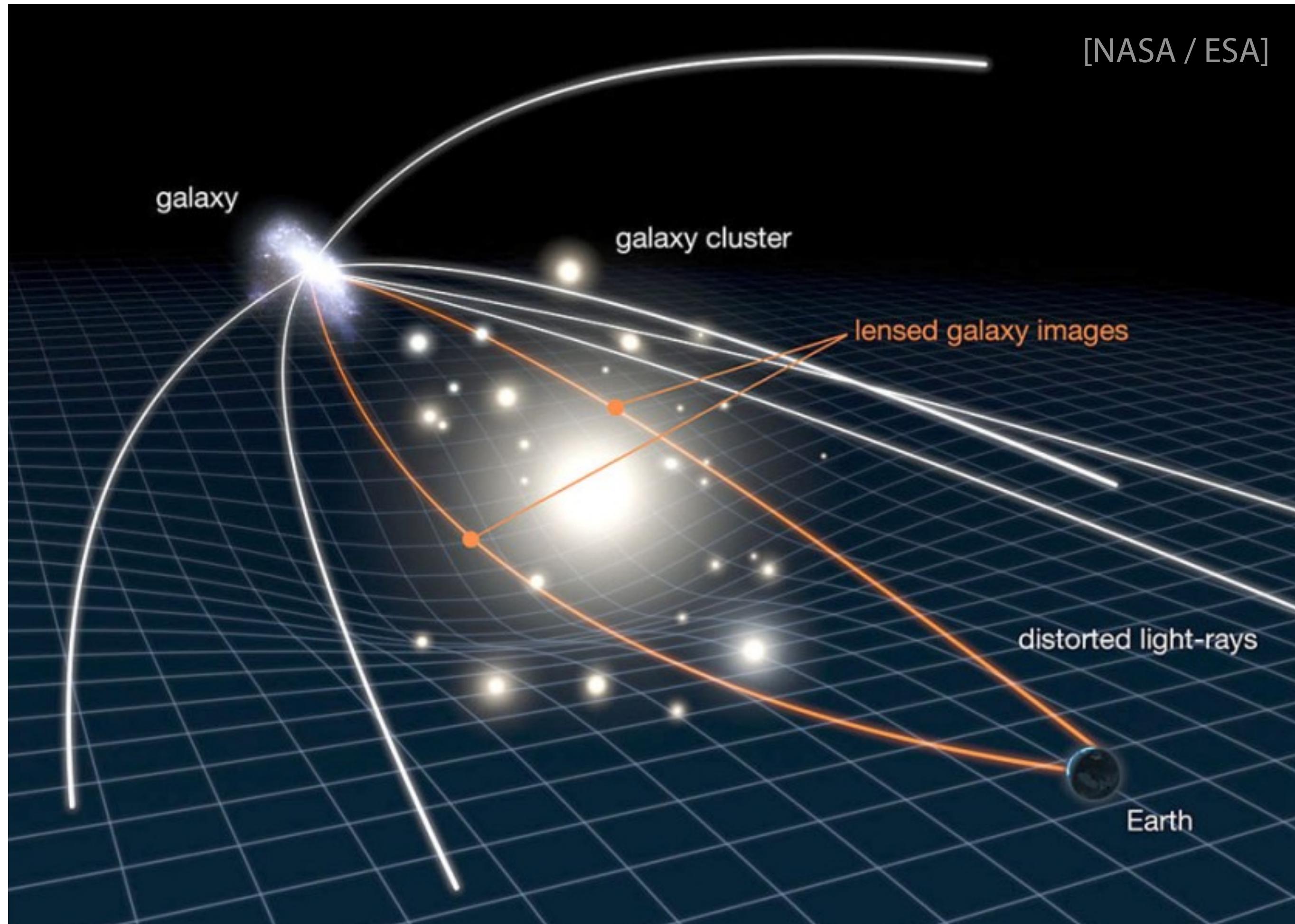
DM small-scale structure as a probe of DM particle properties



Abundance of DM subhalos vs mass:



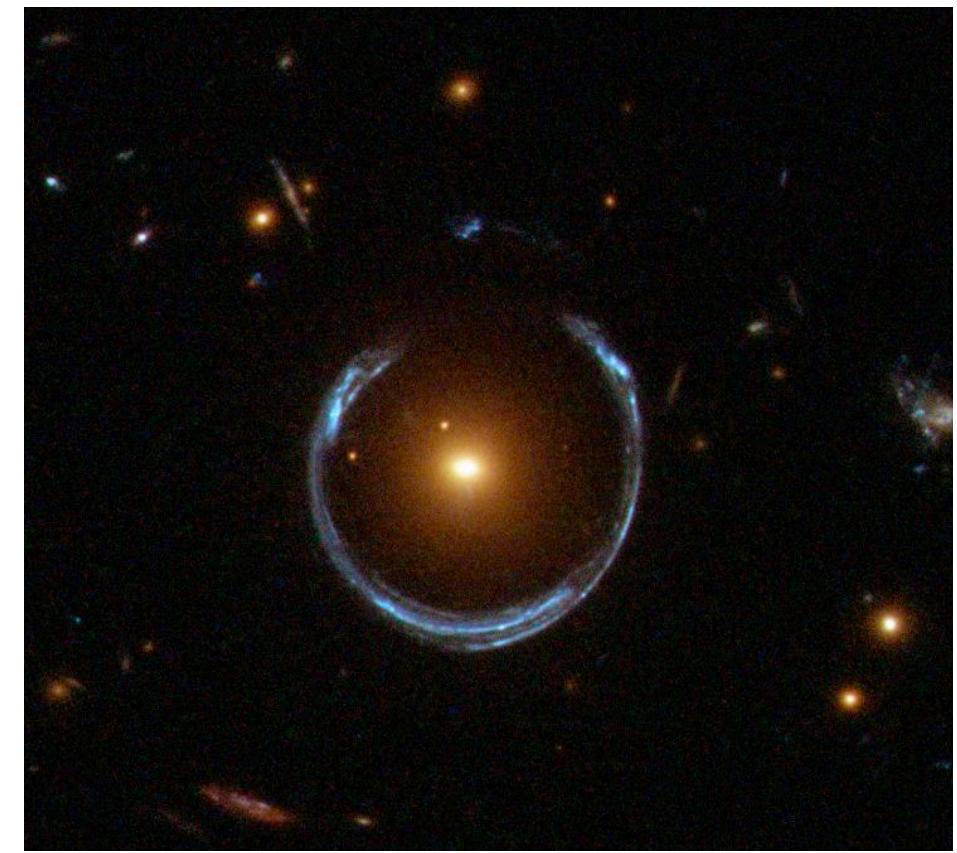
Strong gravitational lensing



Multiple images
of quasars

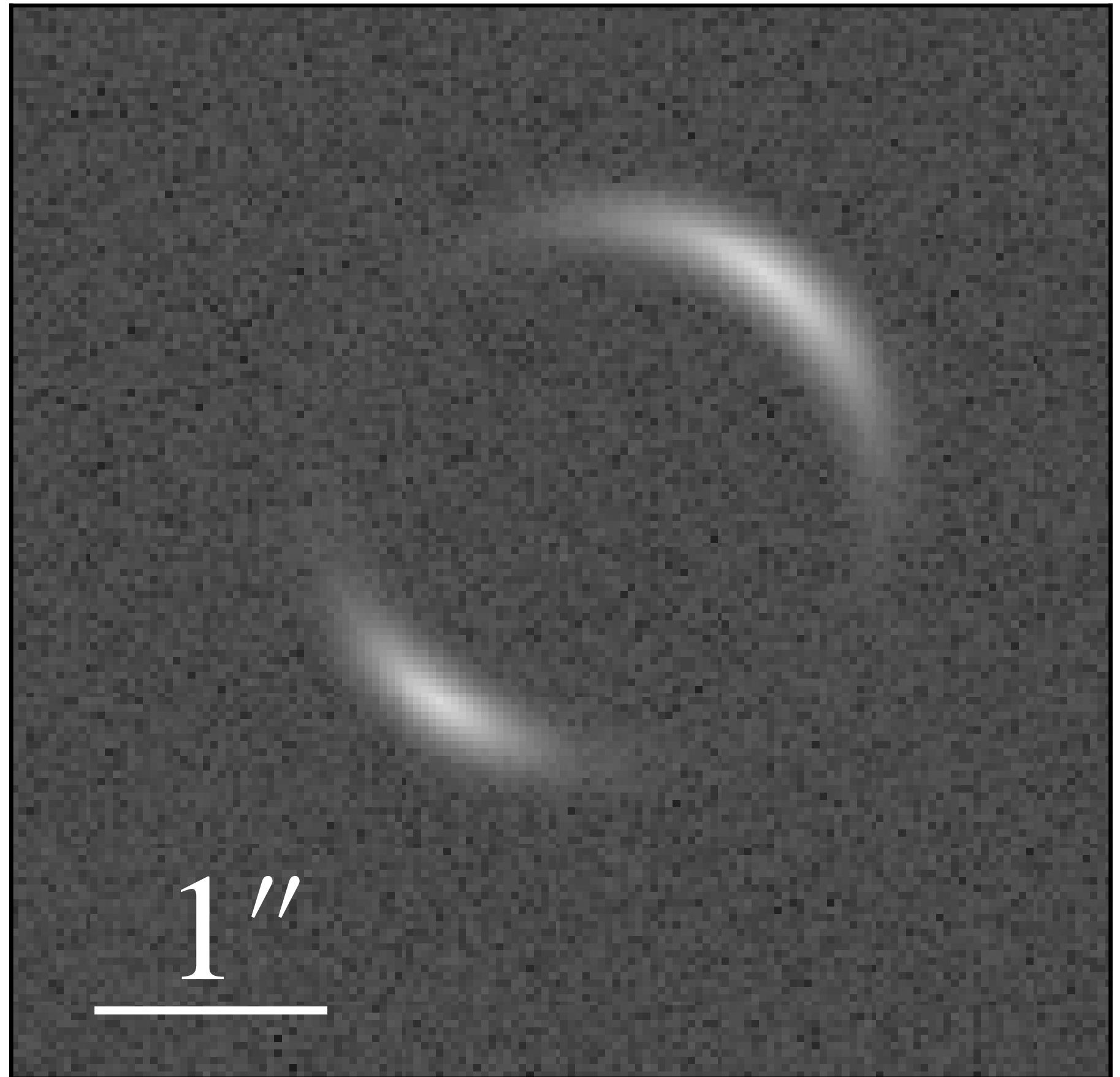


Extended arcs
from galaxies



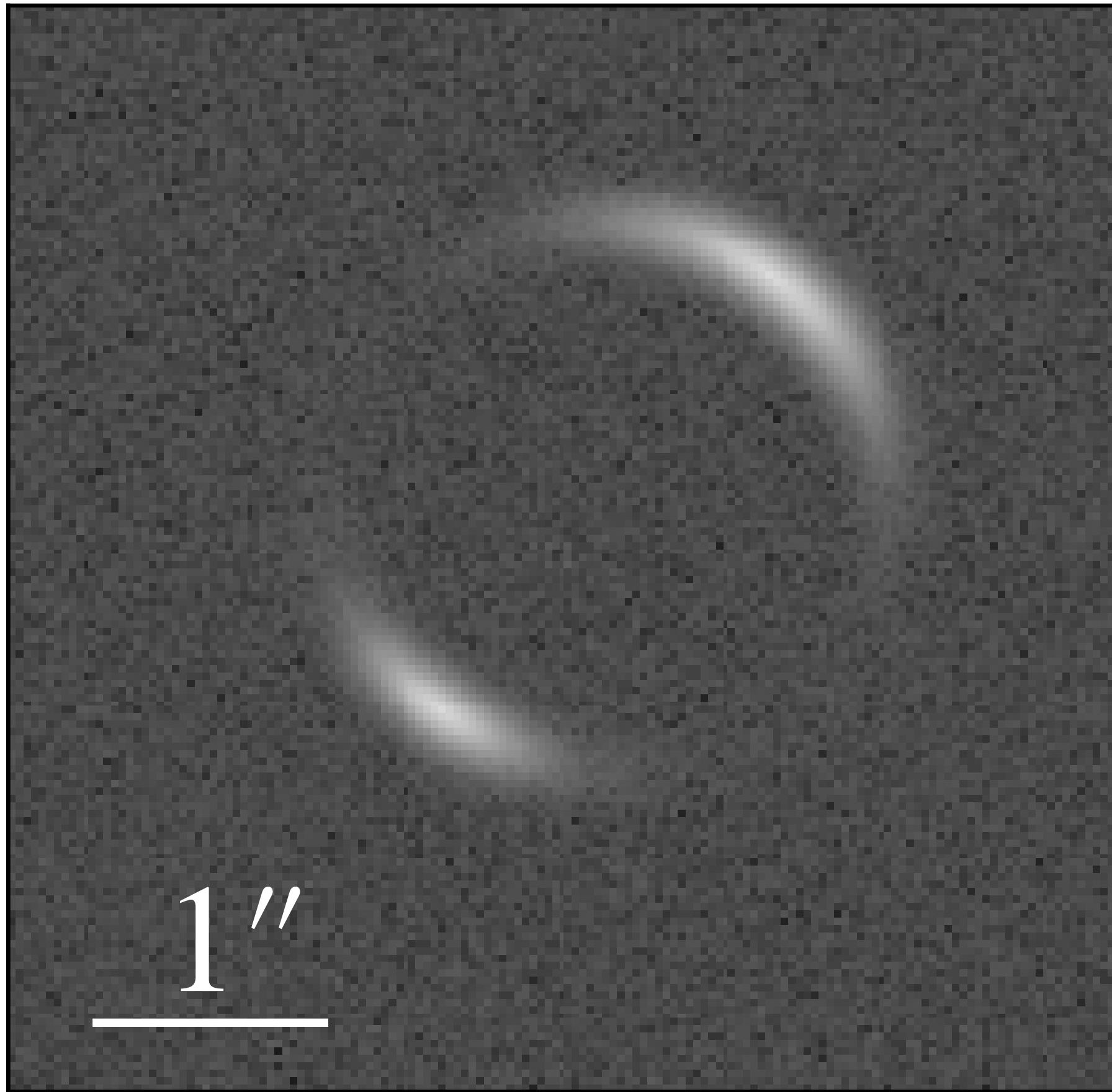
Subhalos affect strong lensing

Smooth halo only

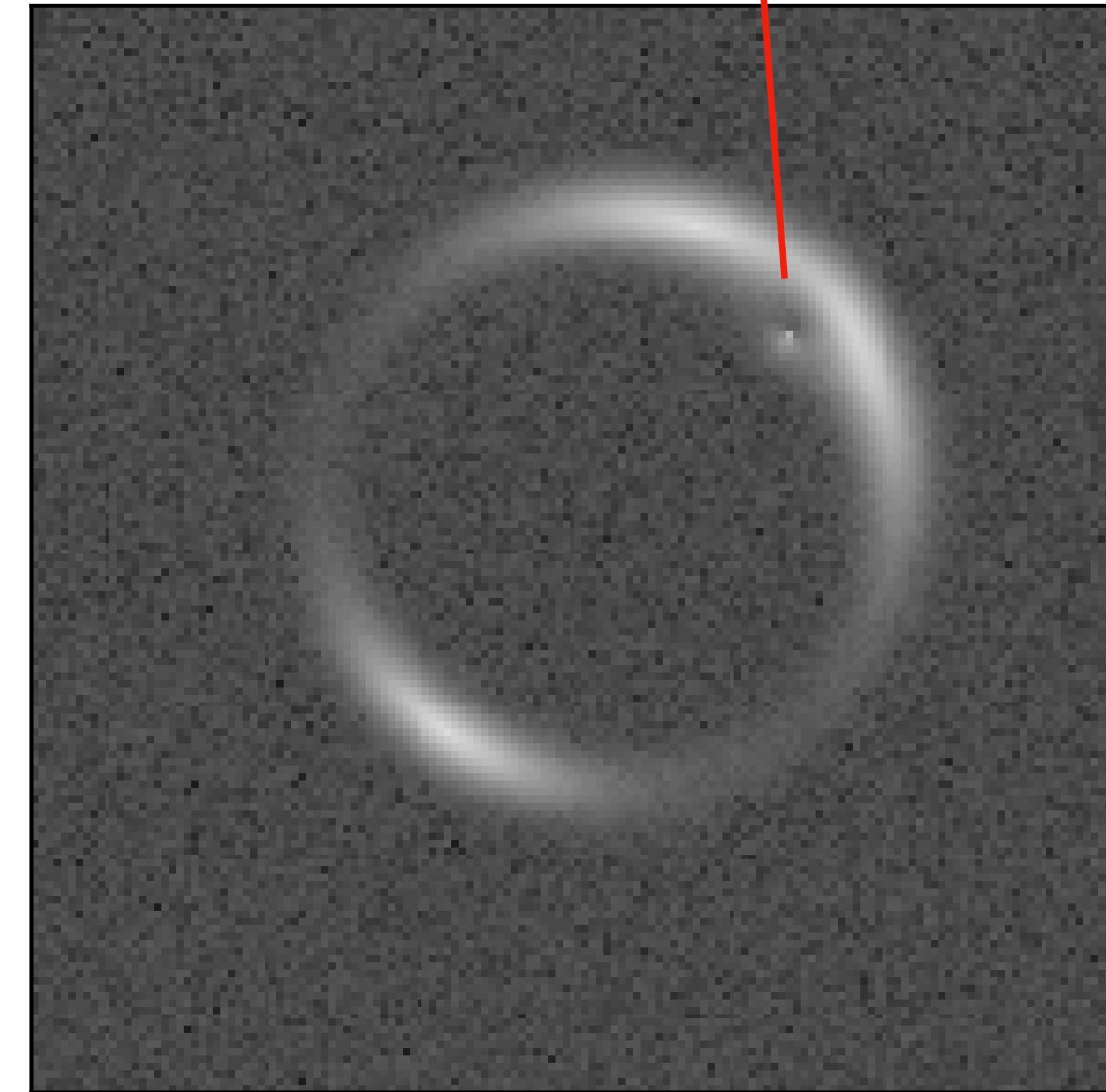


Subhalos affect strong lensing

Smooth halo only

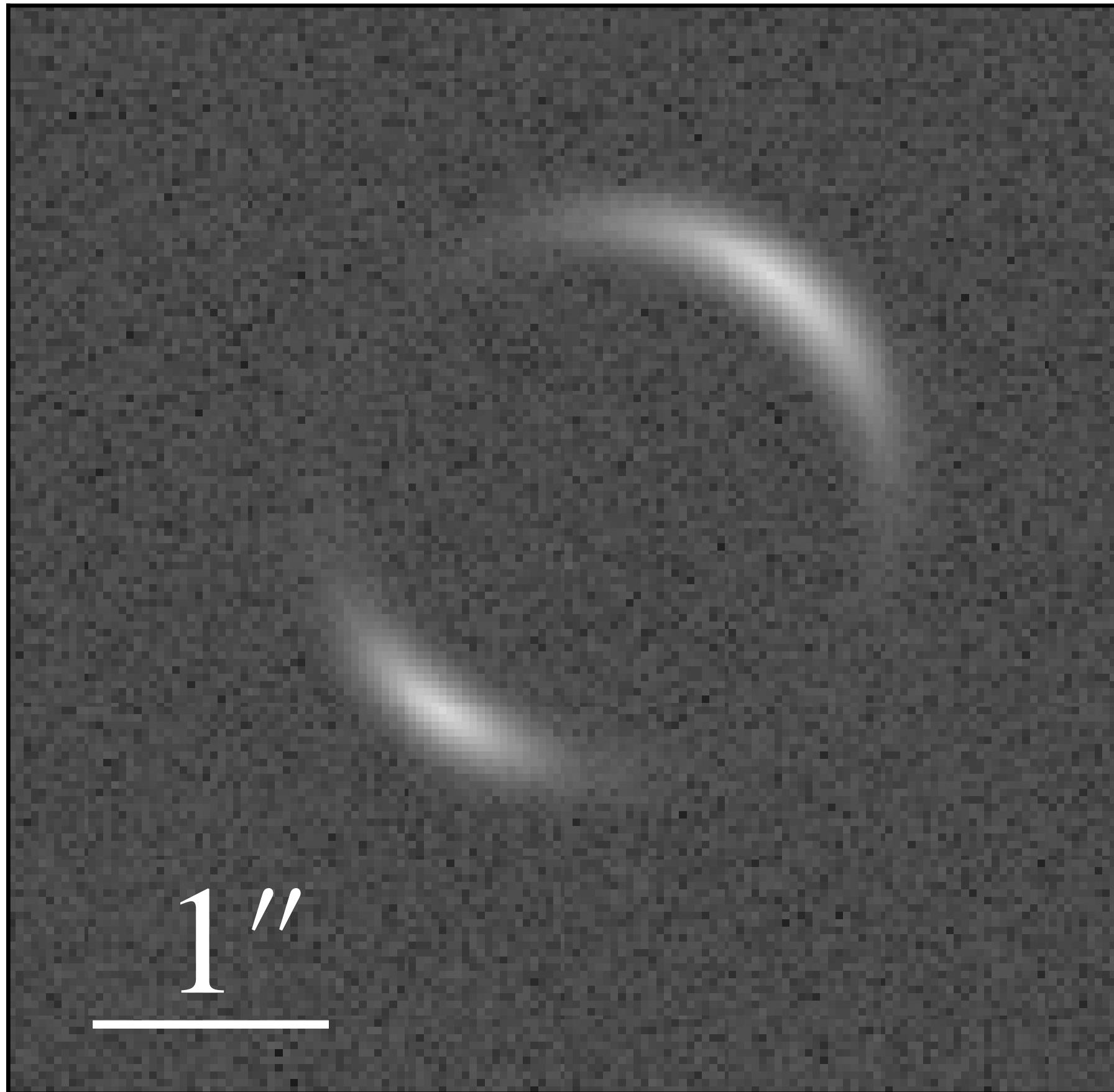


Smooth halo + **subhalo**

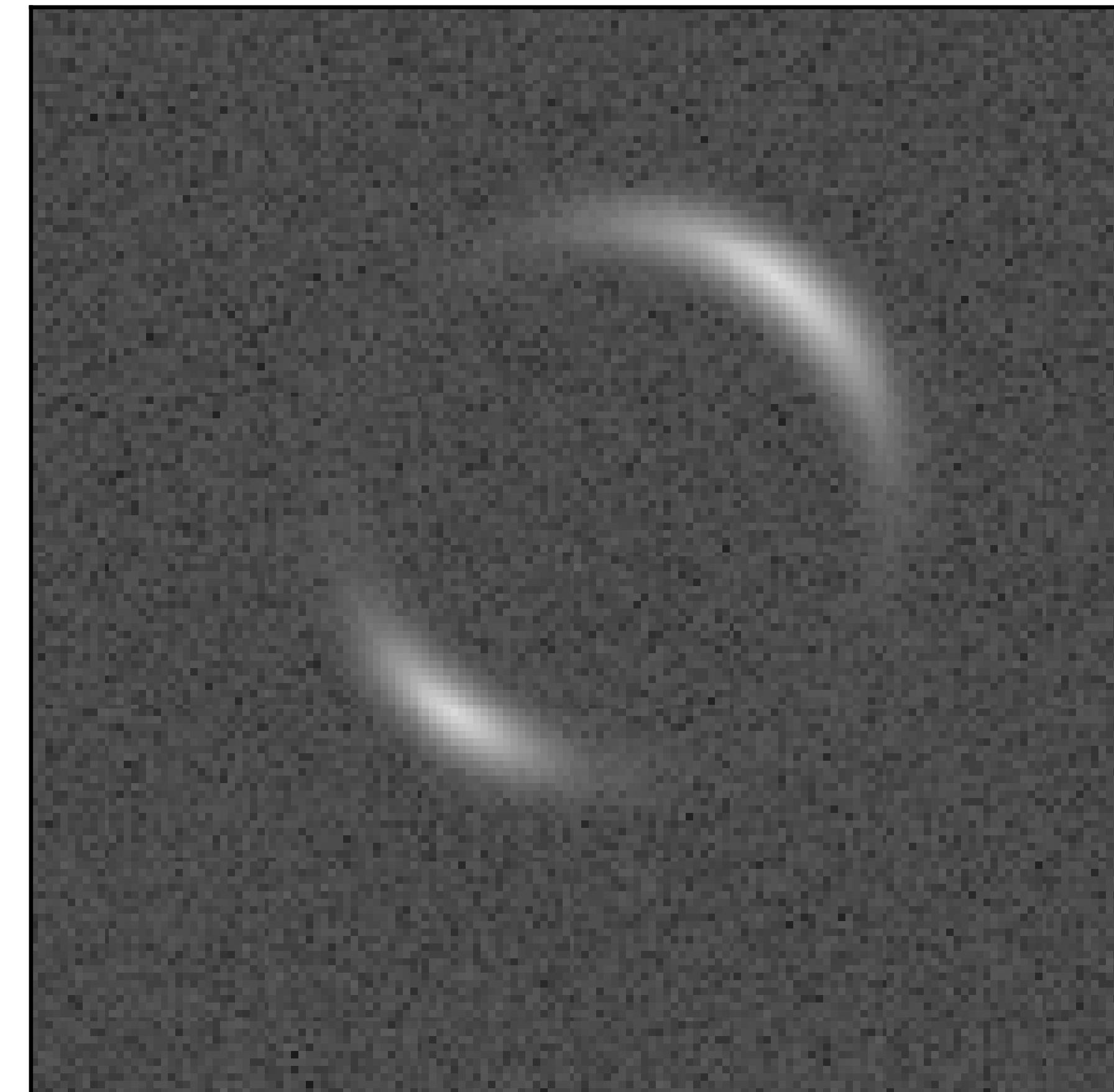


Subhalos affect strong lensing... realistically

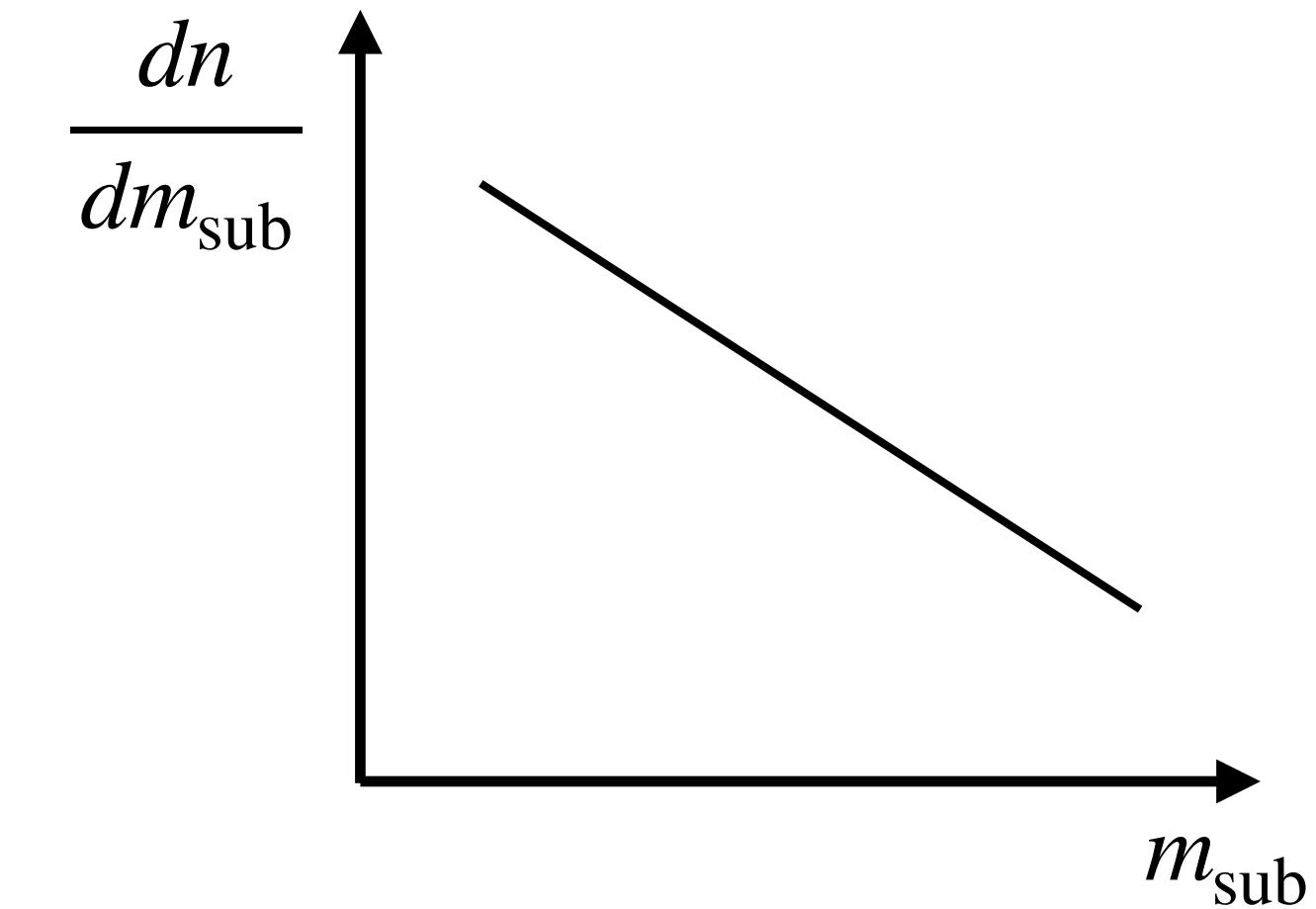
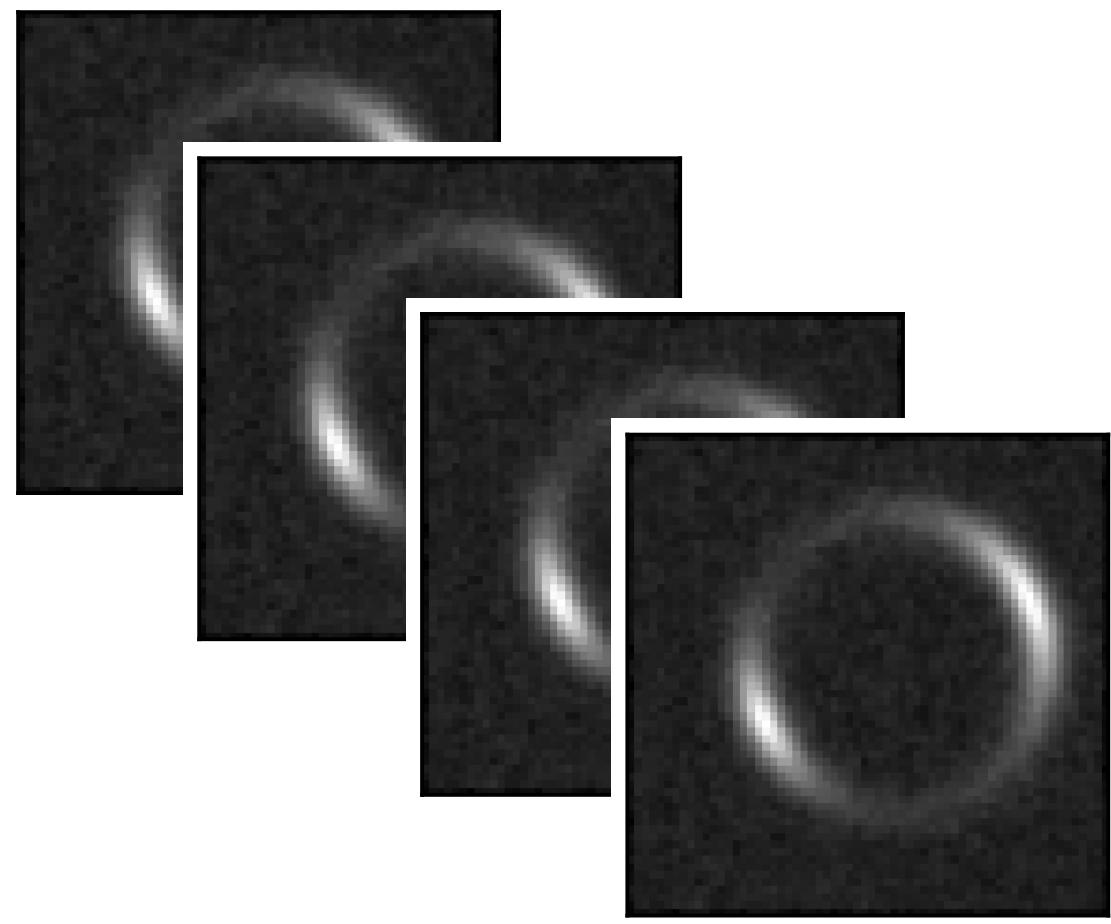
Smooth halo only



Smooth halo + subhalos



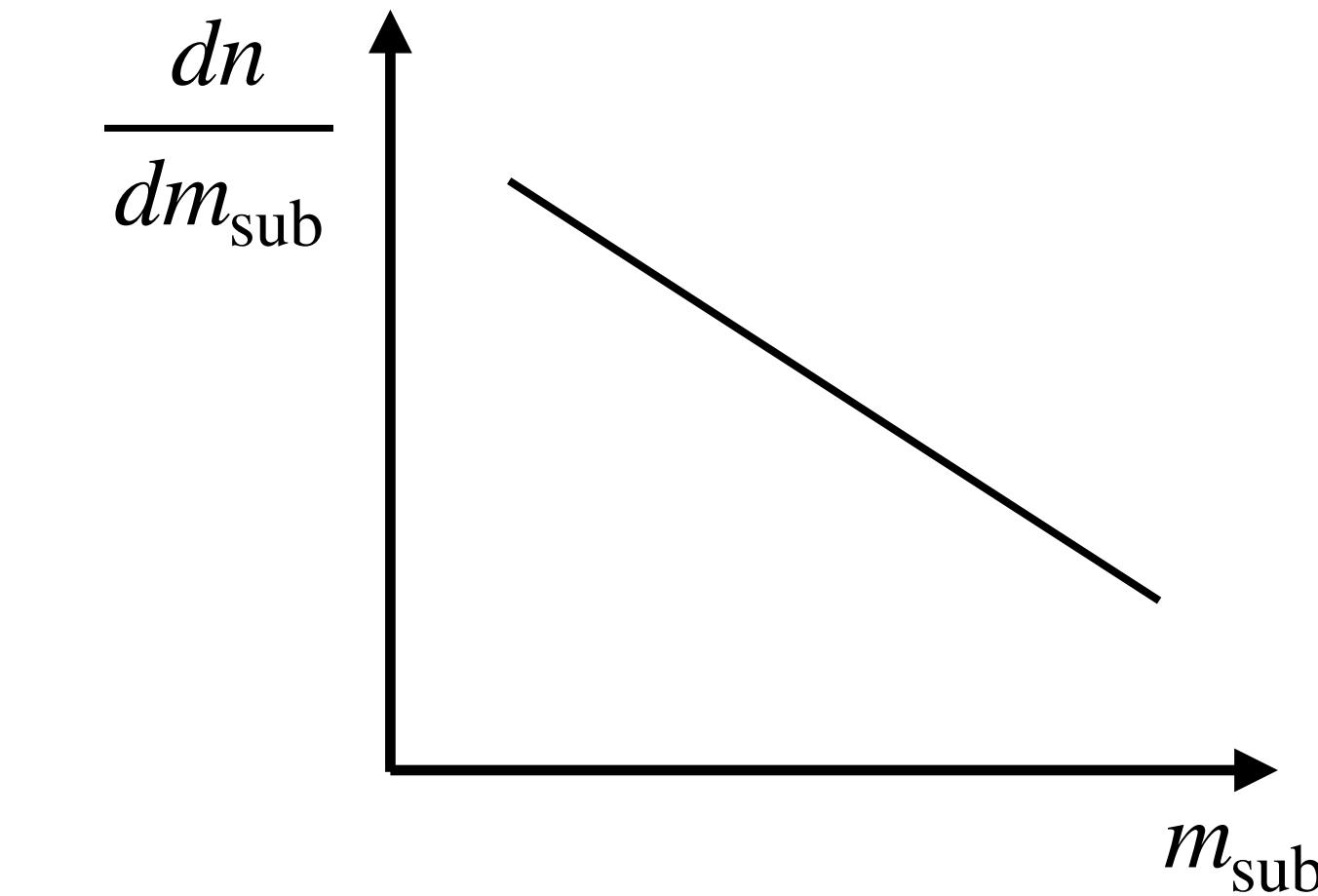
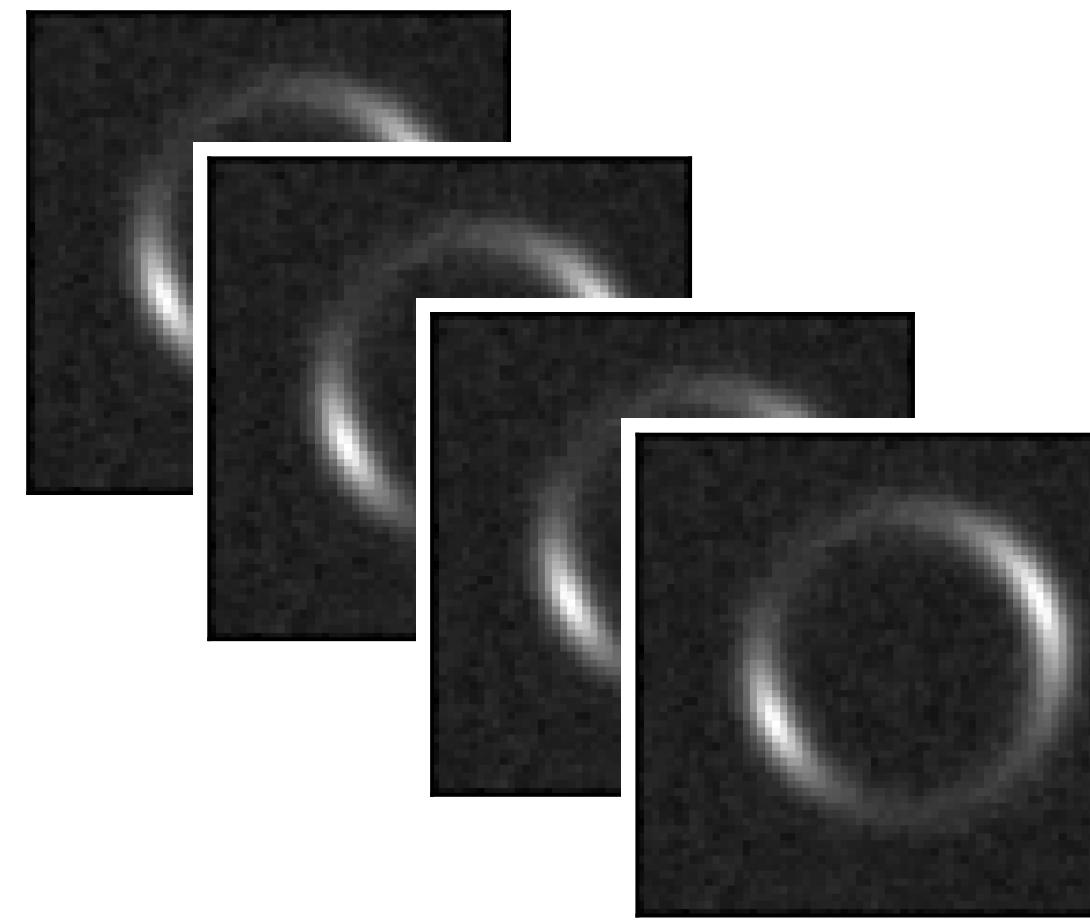
Scalable inference for small subhalos



Future surveys (LSST, Euclid) are expected to deliver large samples of galaxy-galaxy strong lenses [Collett et al 1507.02657]

Goal: infer subhalo mass distribution through collective effects of many light subhalos

Scalable inference for small subhalos



Future surveys (LSST, Euclid) are expected to deliver large samples of galaxy-galaxy strong lenses [Collett et al 1507.02657]

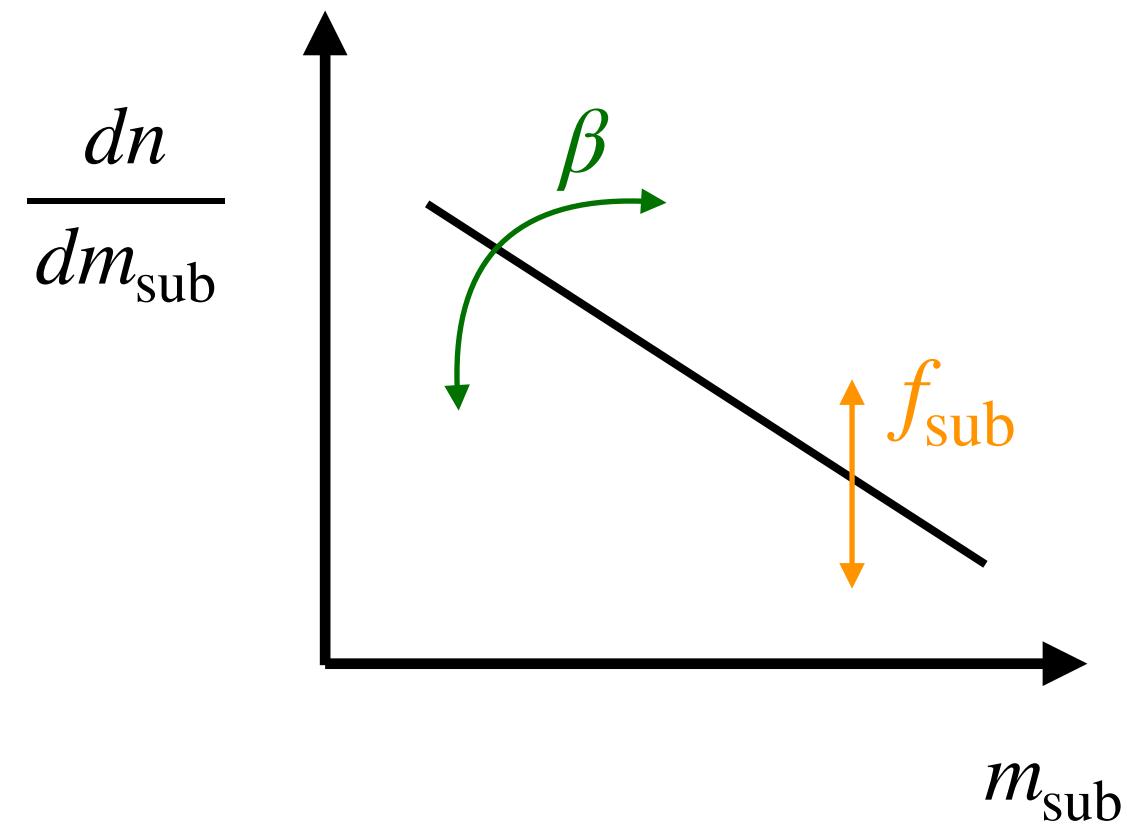
Goal: infer subhalo mass distribution through collective effects of many light subhalos

⇒ Need inference technique that

- scales to many lenses (fast evaluation)
- captures subtle effects in high-dimensional image data
- can deal with a large number of subhalos (latent variables)

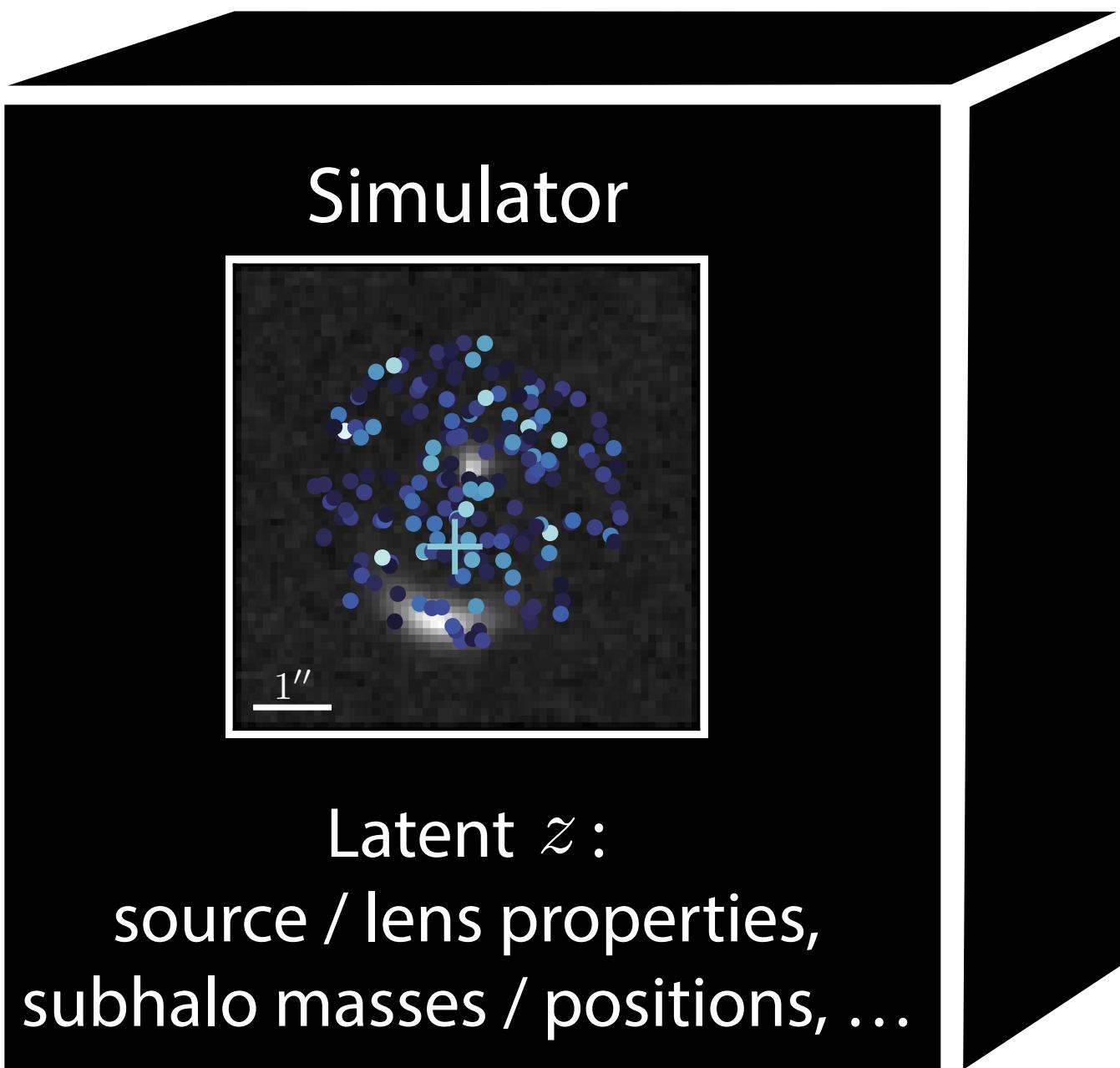
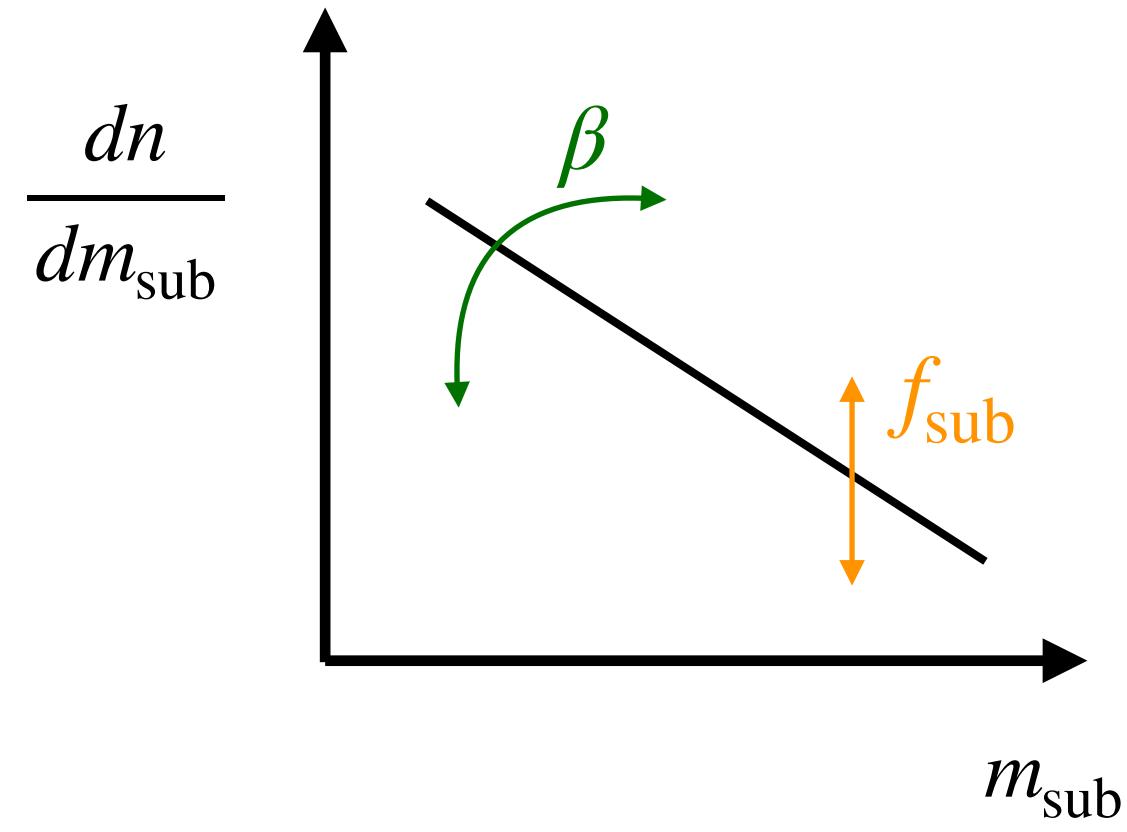
Simulation-based inference for strong lensing

2 parameters $\theta = (\beta, f_{\text{sub}})$



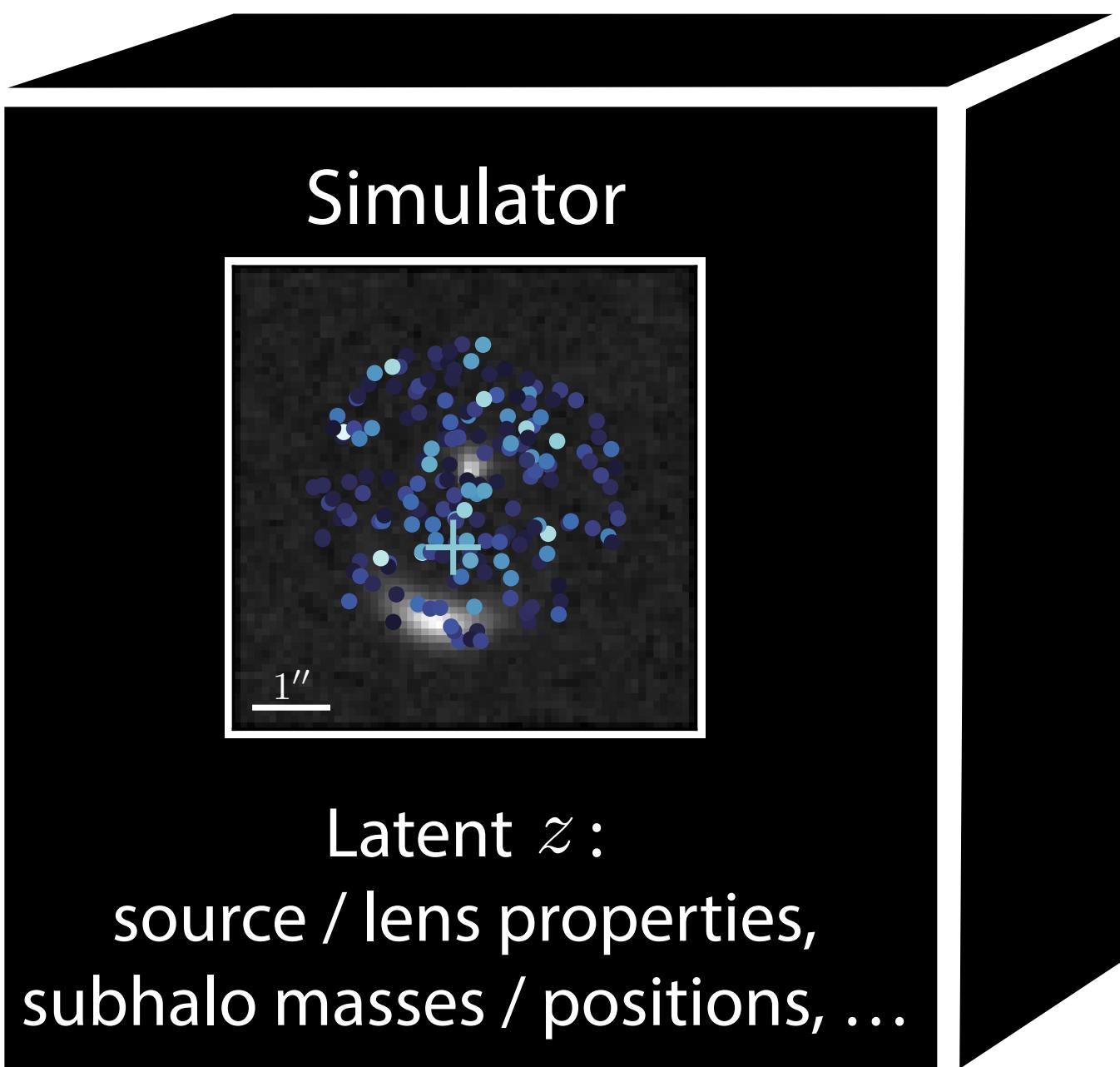
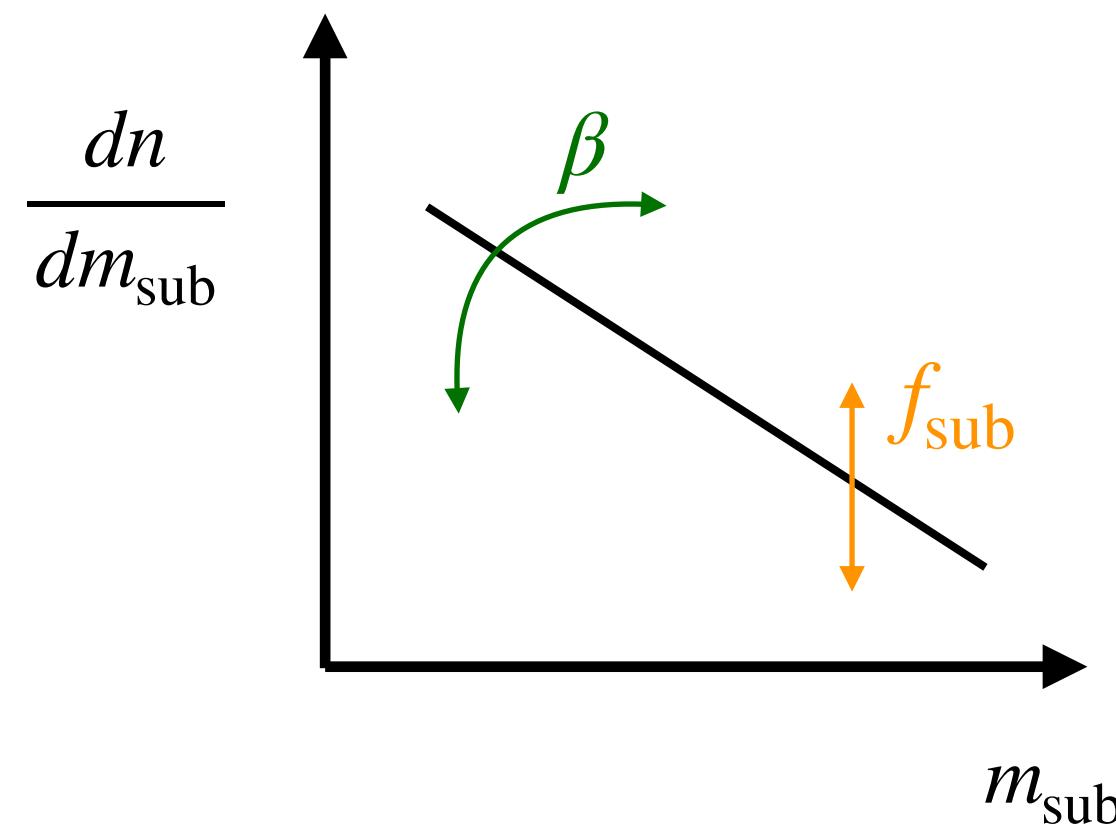
Simulation-based inference for strong lensing

2 parameters $\theta = (\beta, f_{\text{sub}})$

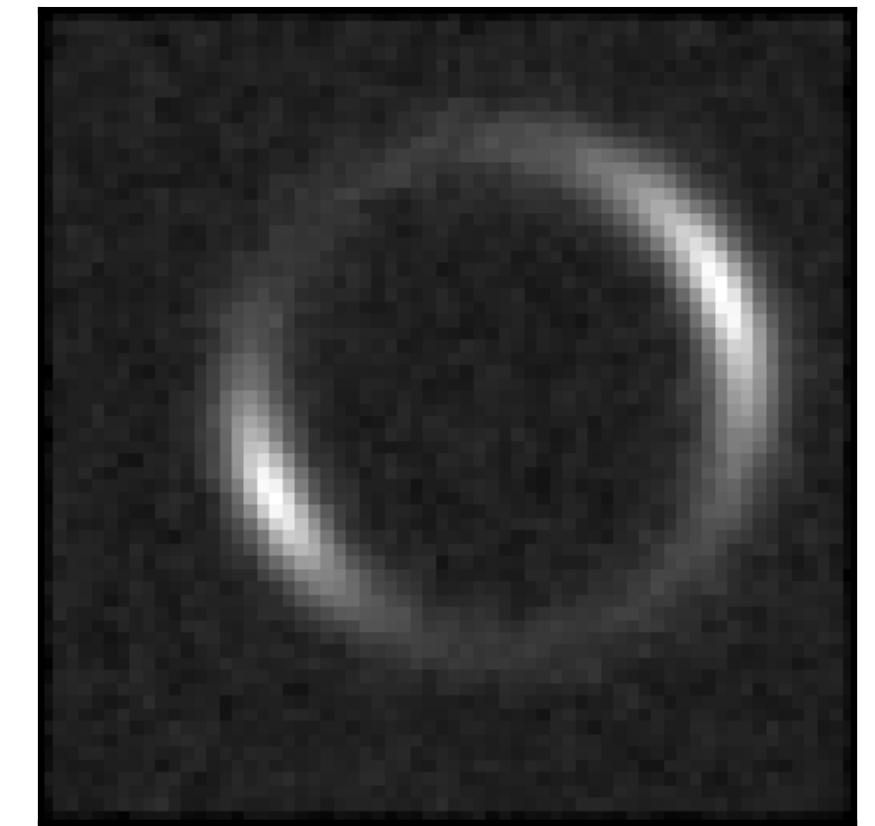


Simulation-based inference for strong lensing

2 parameters $\theta = (\beta, f_{\text{sub}})$

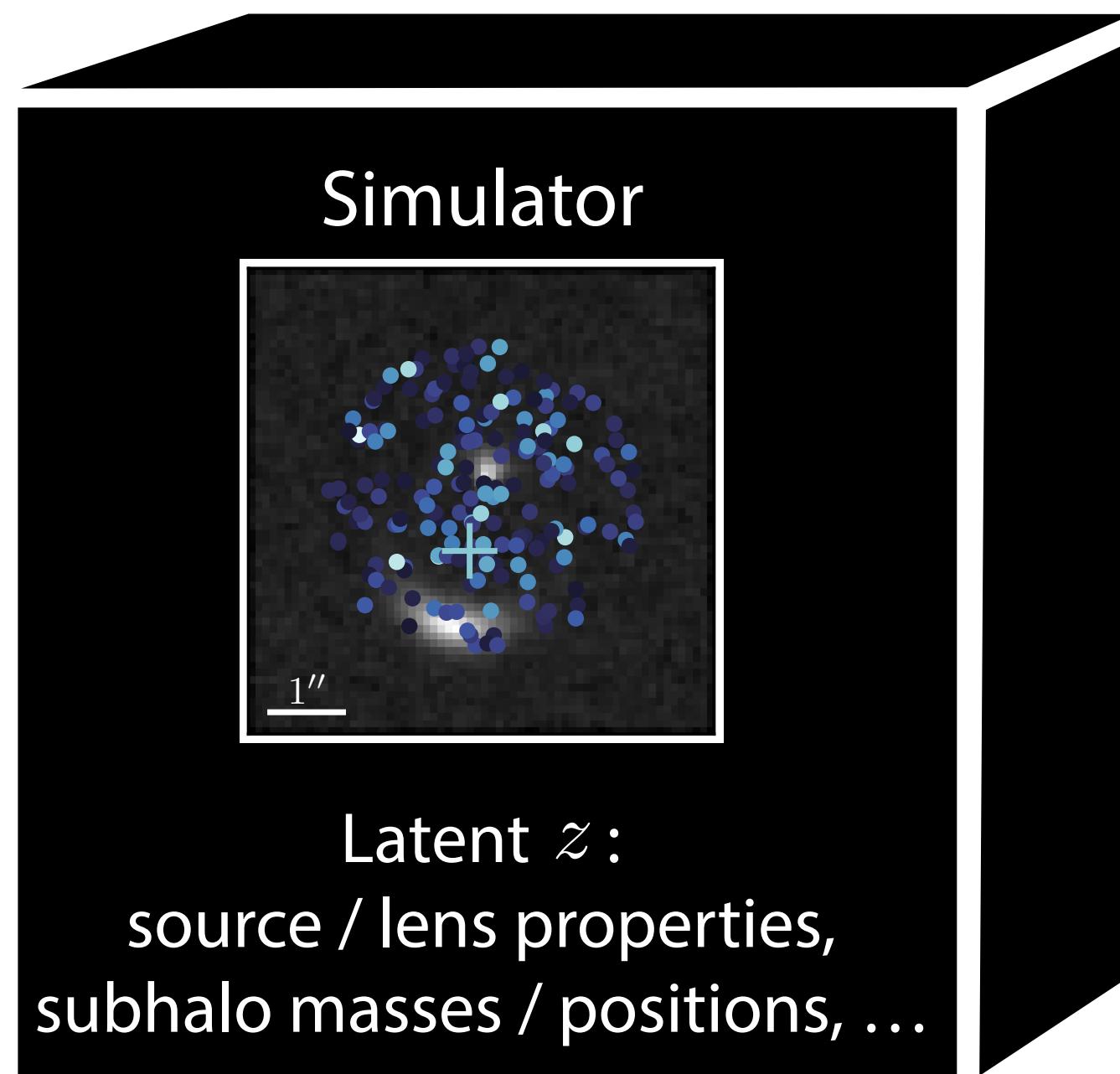
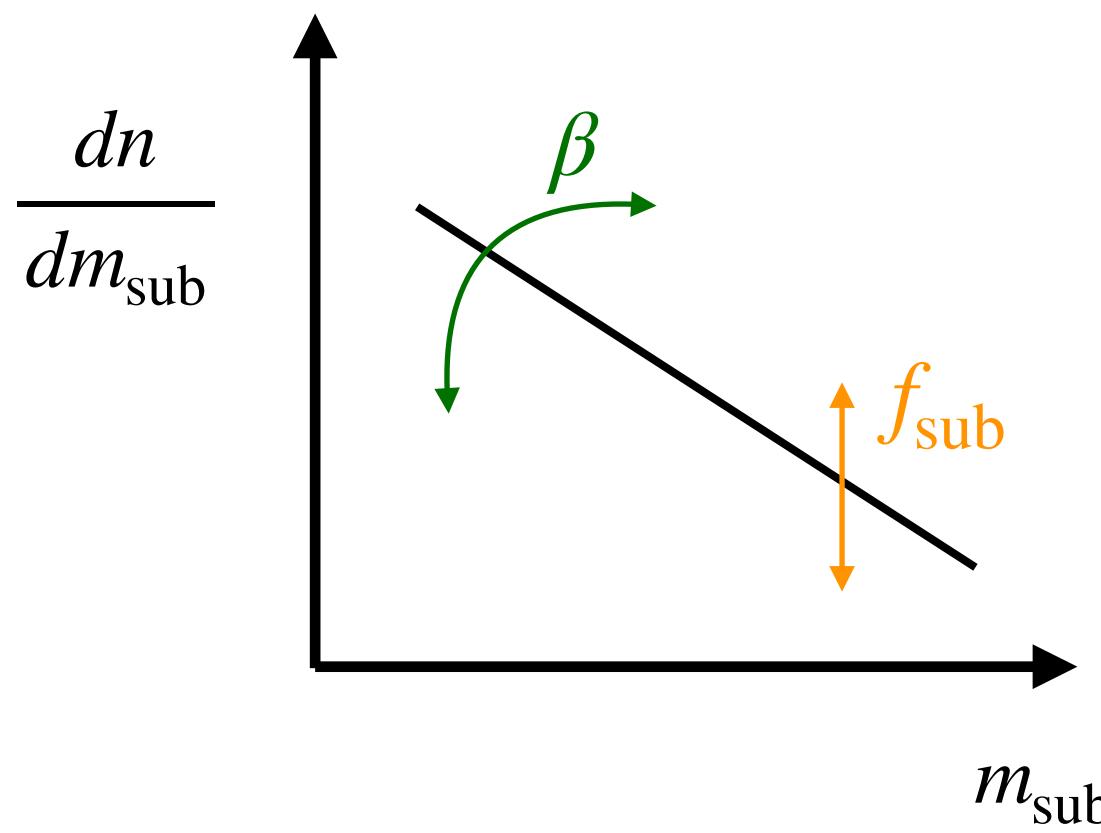


64² observables x

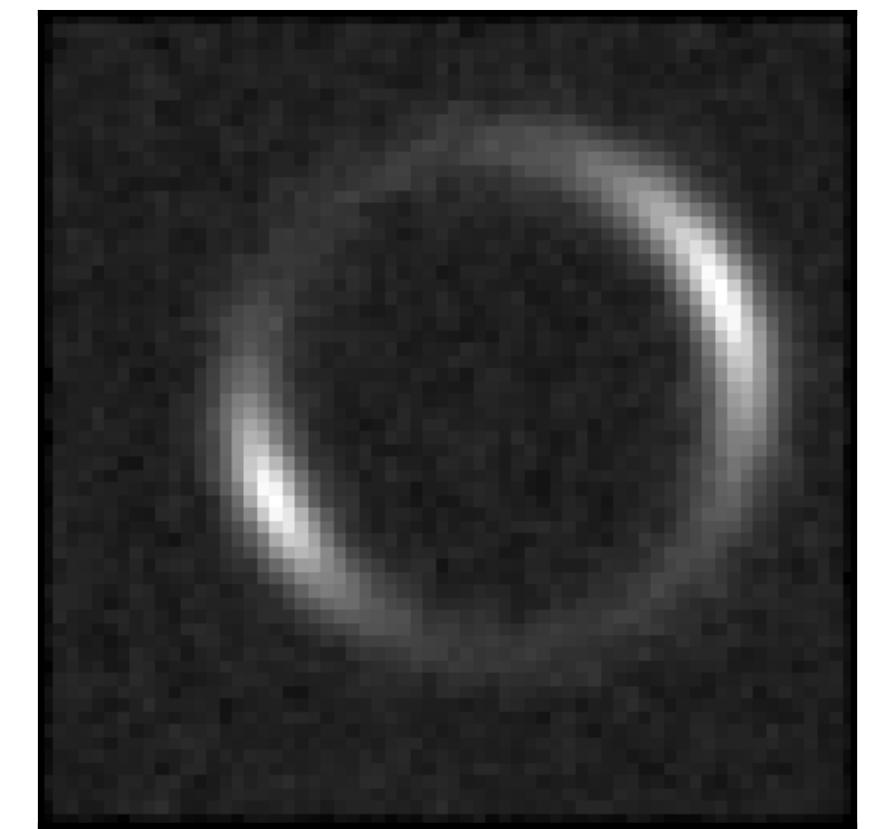


Simulation-based inference for strong lensing

2 parameters $\theta = (\beta, f_{\text{sub}})$



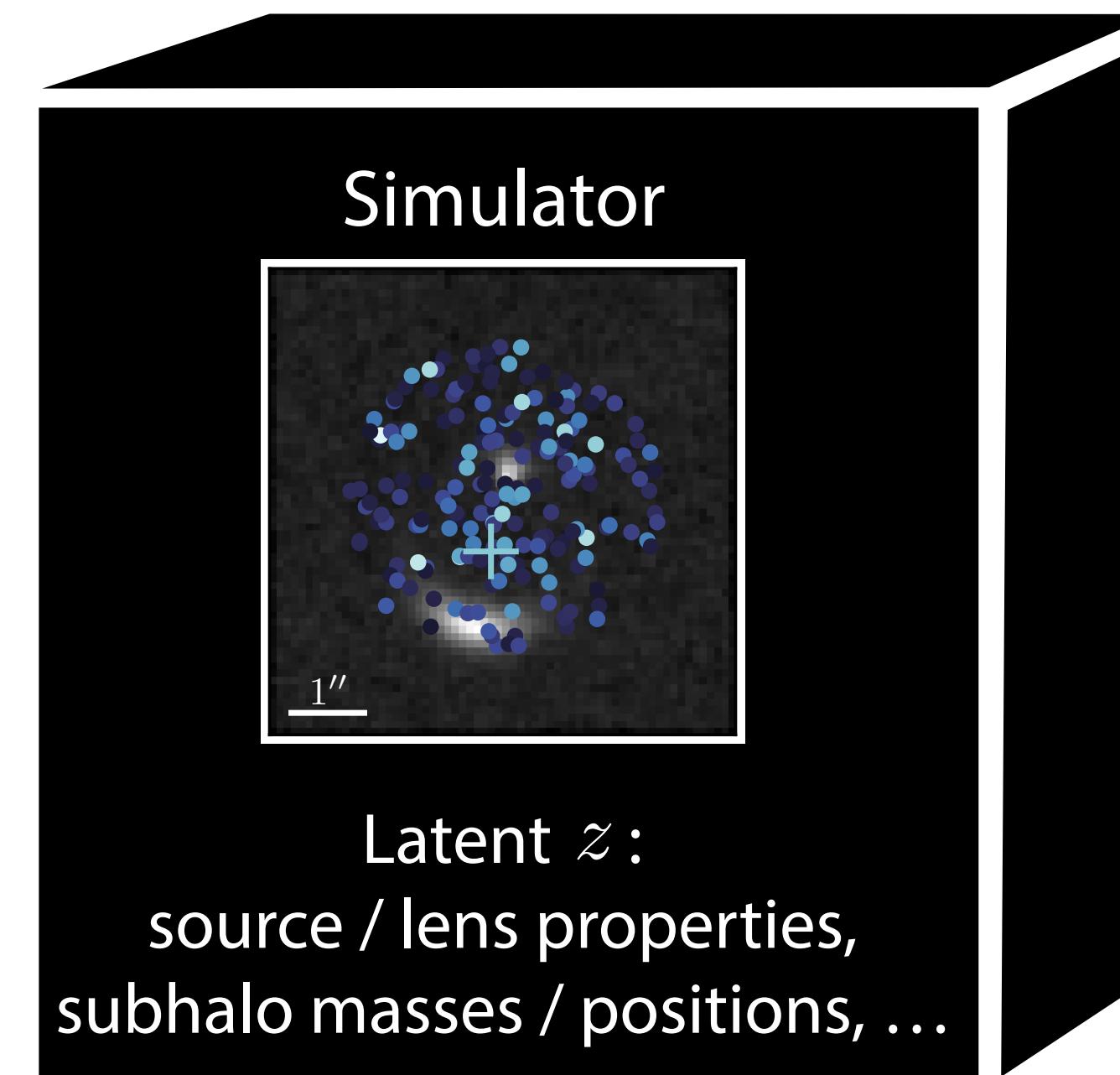
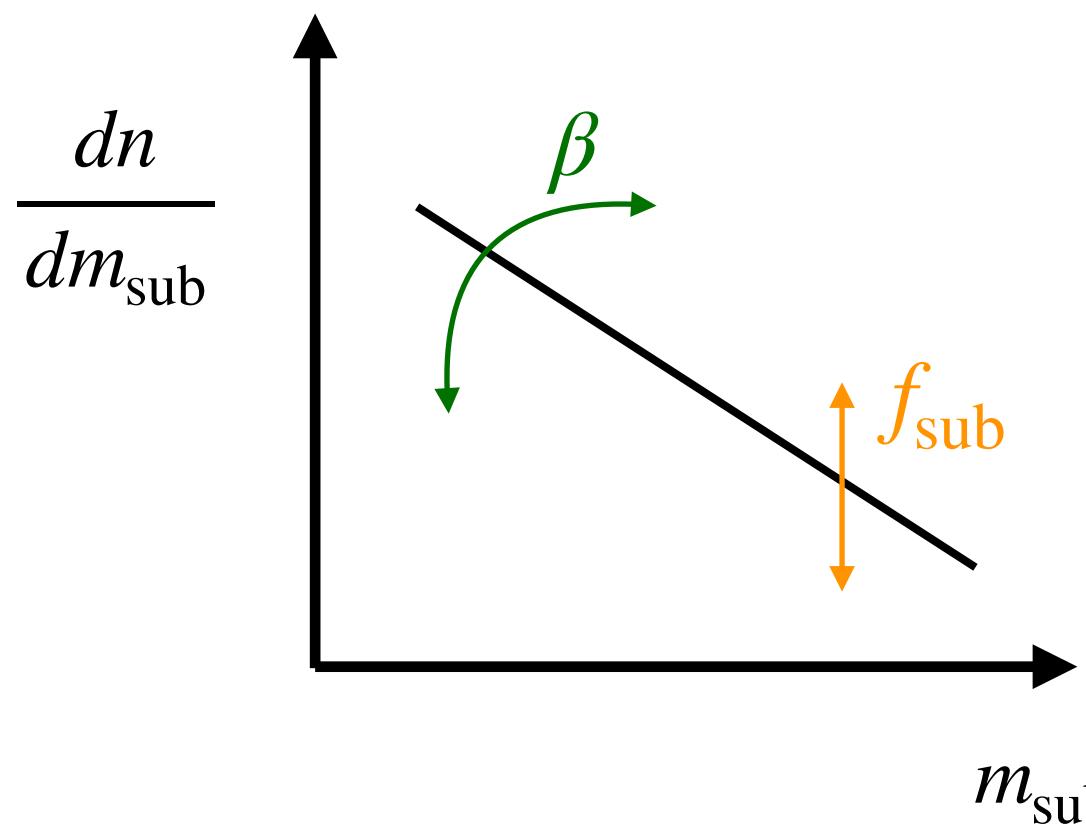
64² observables x



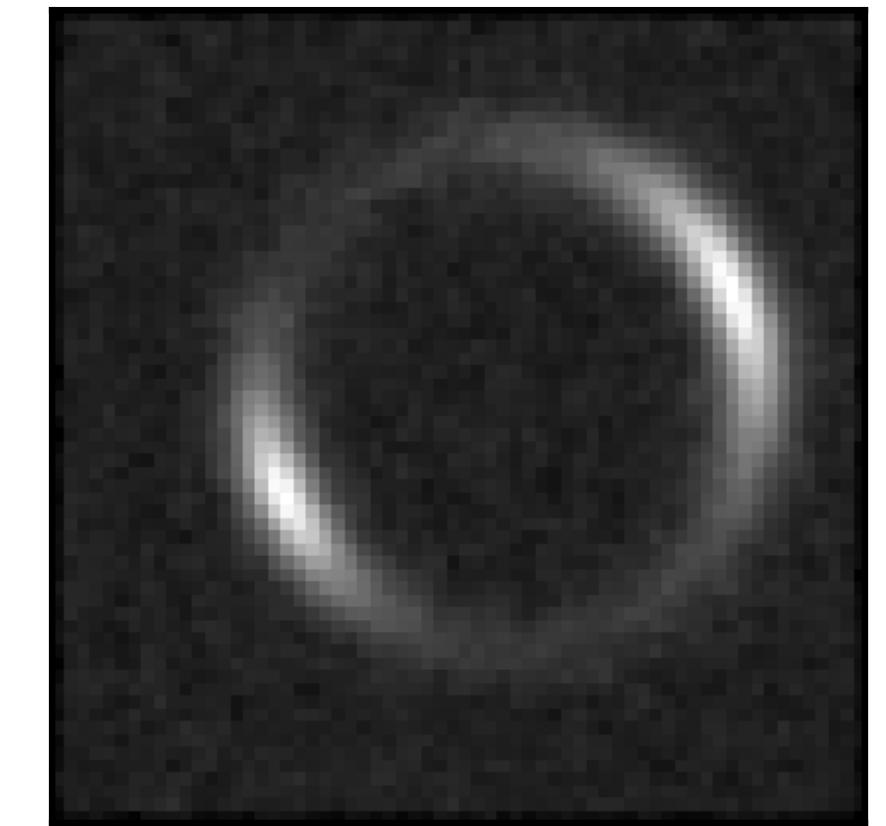
Prediction: We construct a simulator that can sample $x \sim p(x|\theta)$

Simulation-based inference for strong lensing

2 parameters $\theta = (\beta, f_{\text{sub}})$



64² observables x

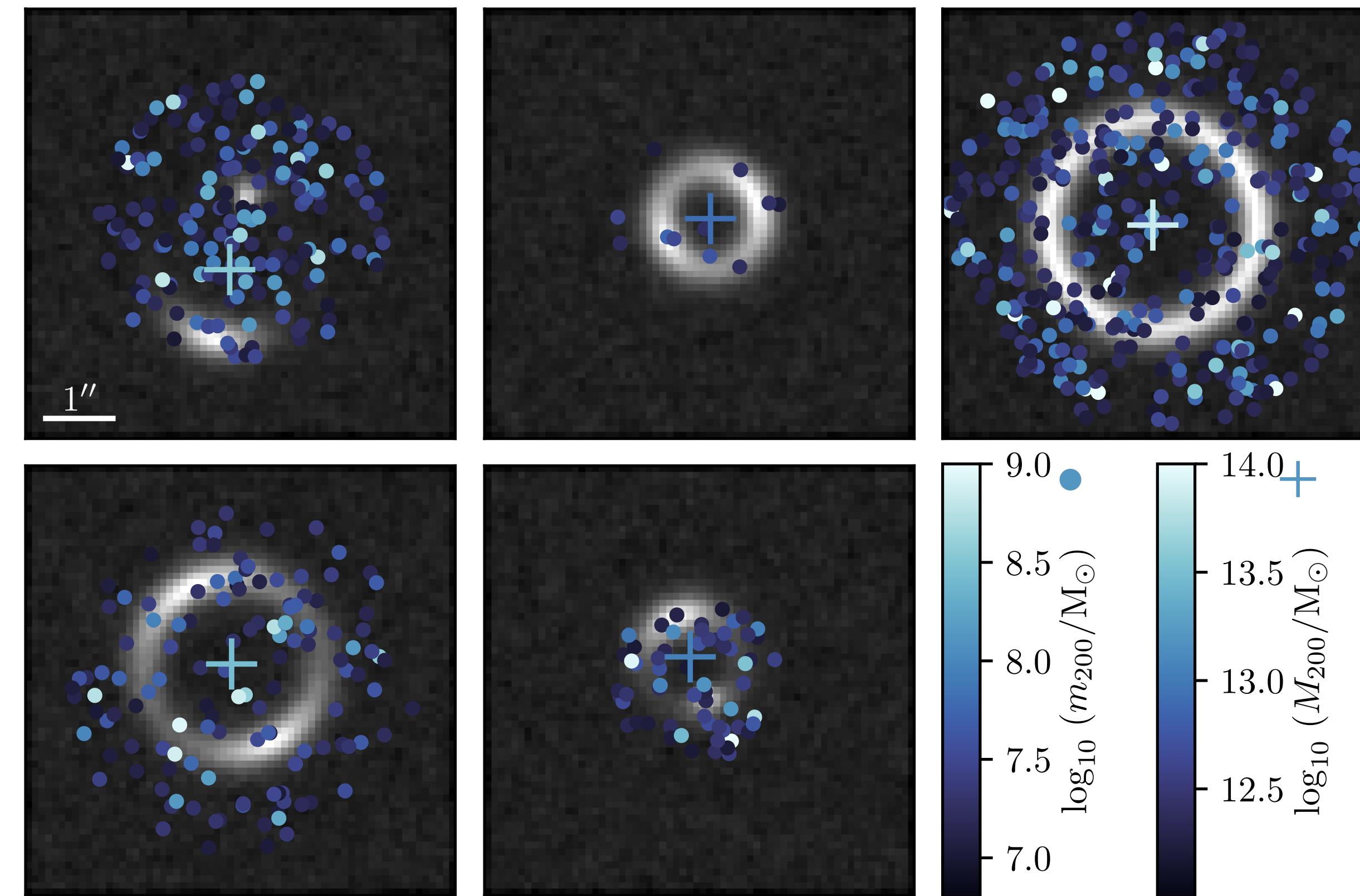


Prediction: We construct a simulator that can sample $x \sim p(x|\theta)$

Inference: We train neural likelihood ratio estimators $\hat{r}(x|\theta)$

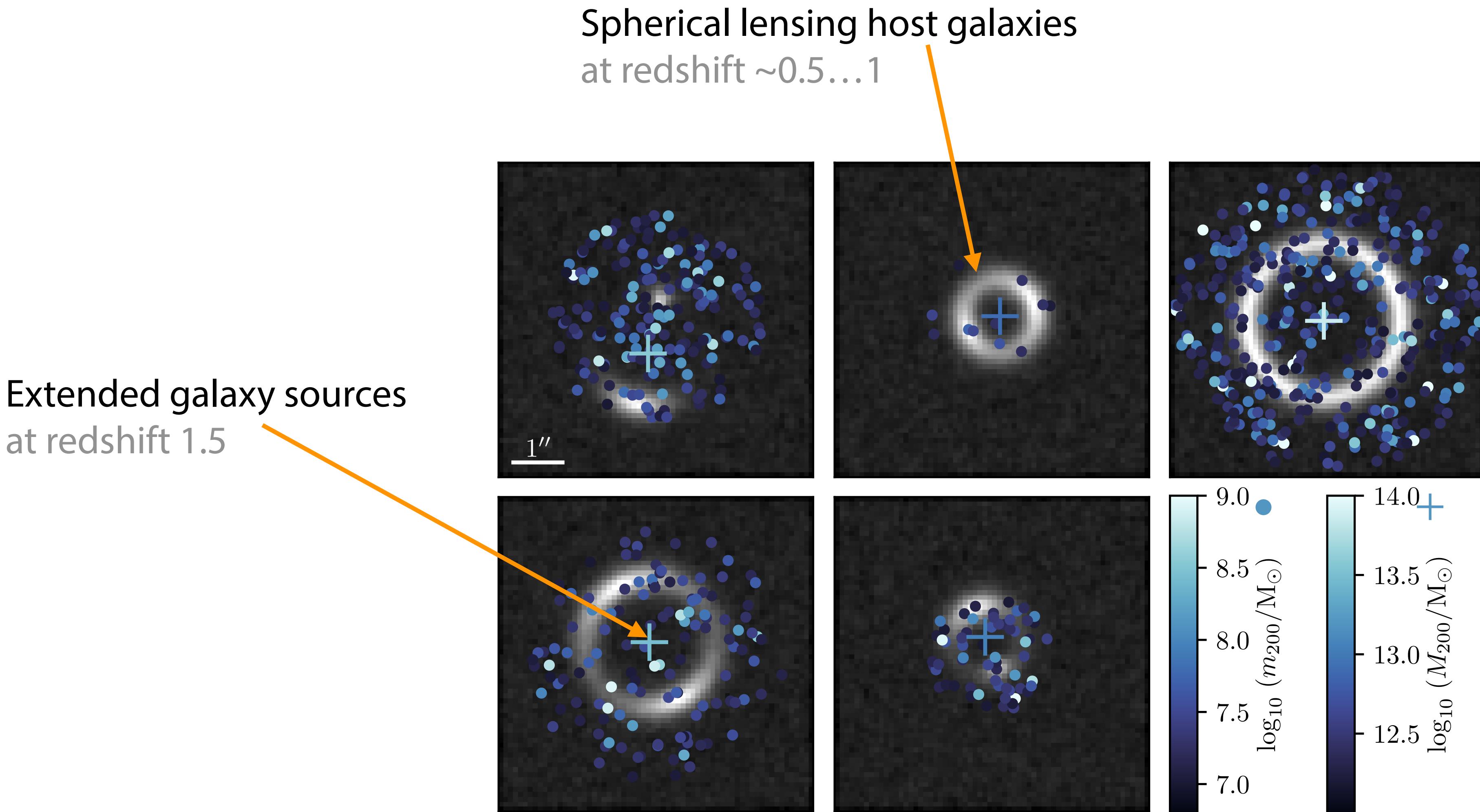
Proof-of-principle simulator

[following T. Collett 1507.02657]



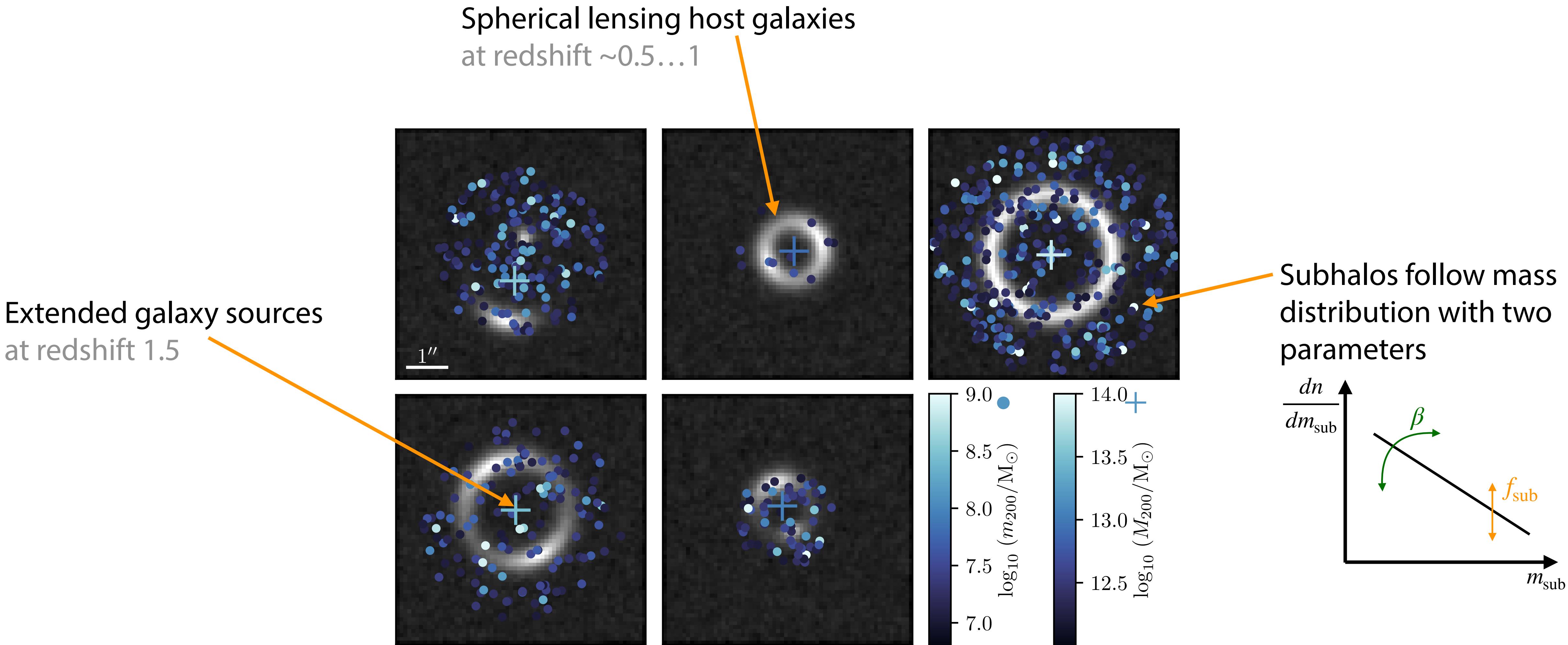
Proof-of-principle simulator

[following T. Collett 1507.02657]



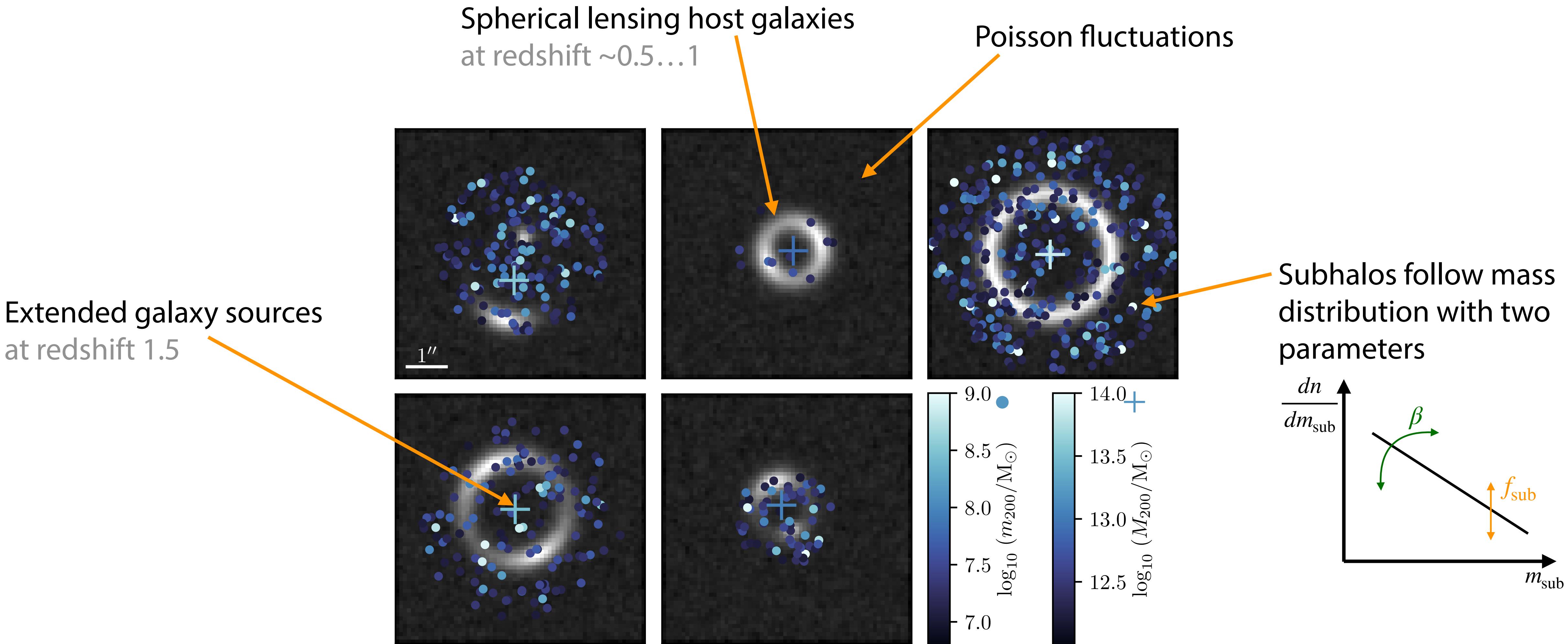
Proof-of-principle simulator

[following T. Collett 1507.02657]



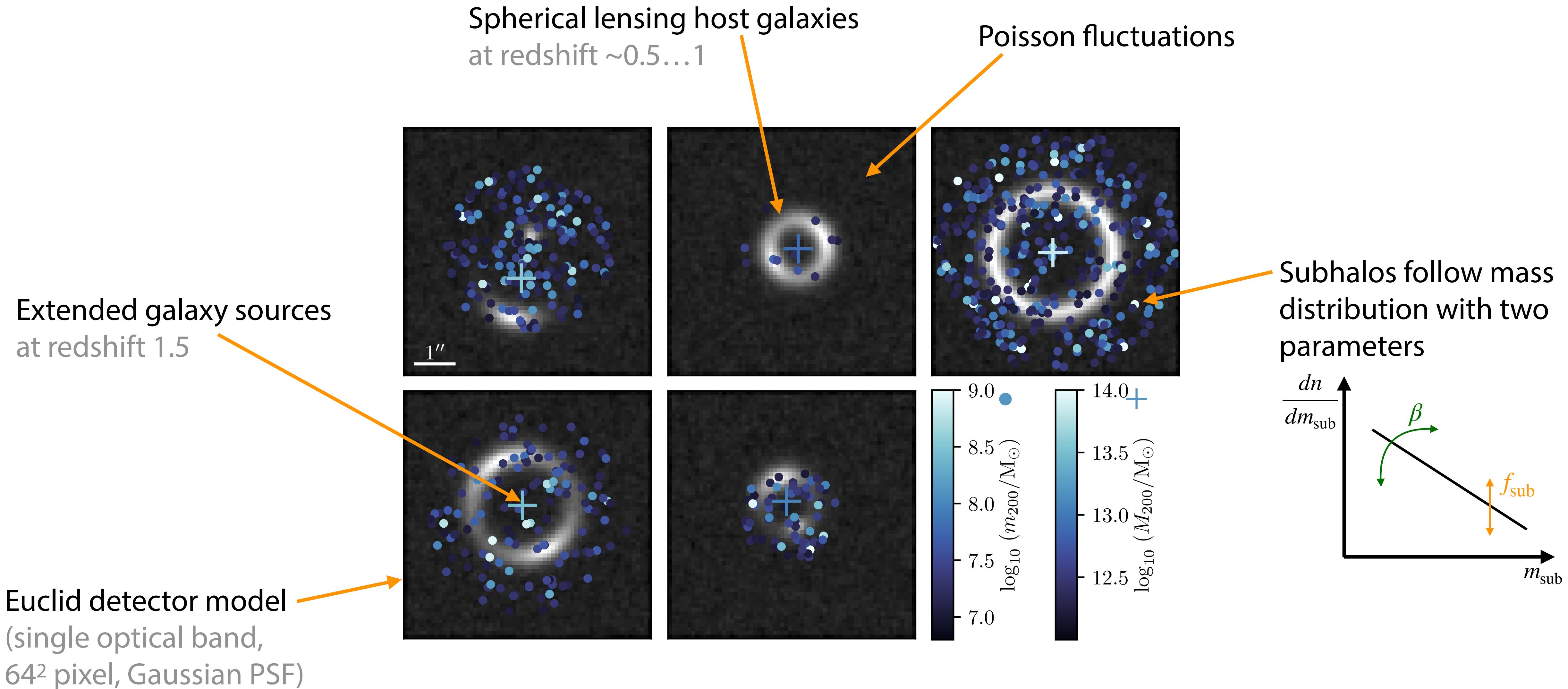
Proof-of-principle simulator

[following T. Collett 1507.02657]

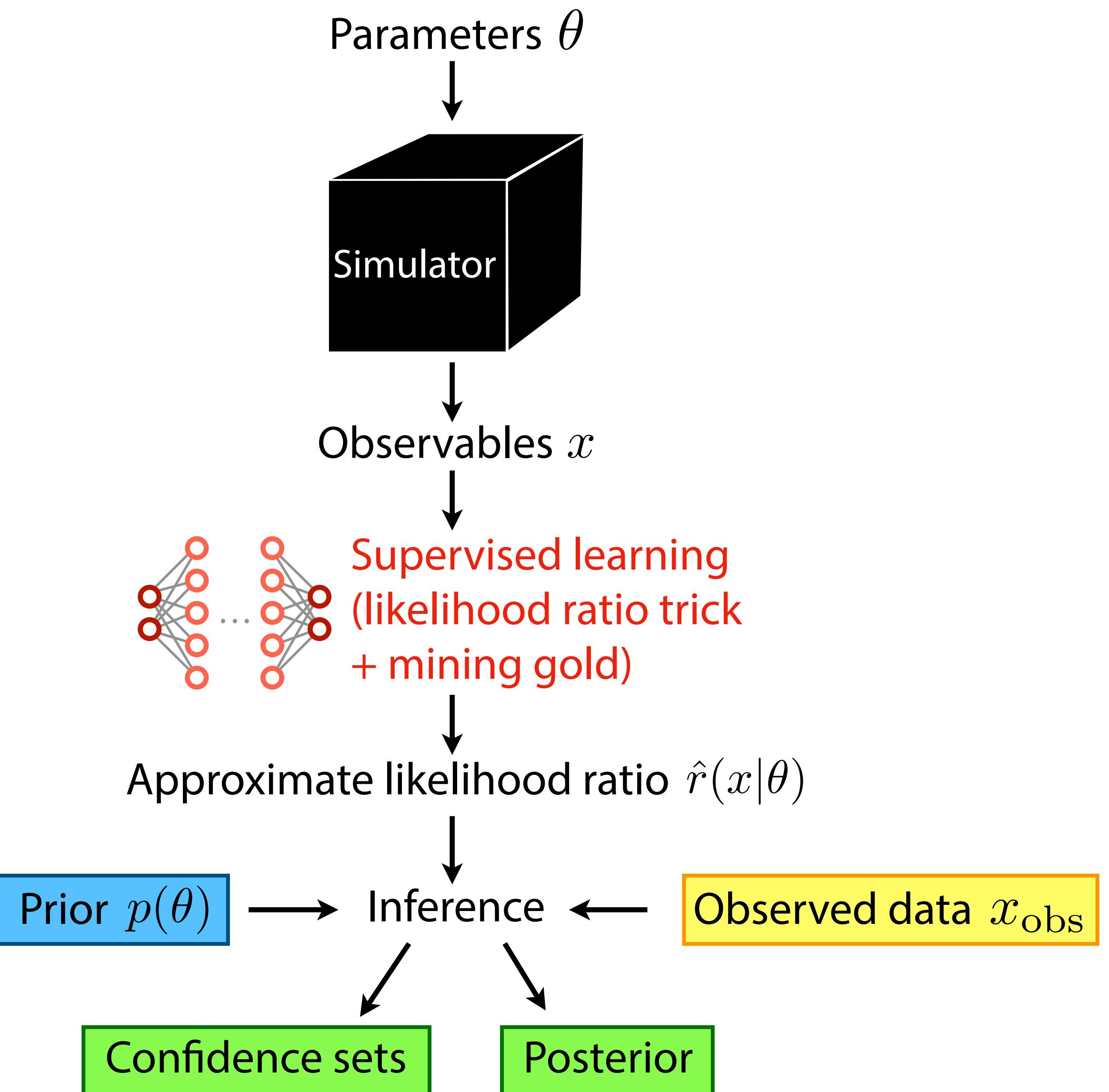


Proof-of-principle simulator

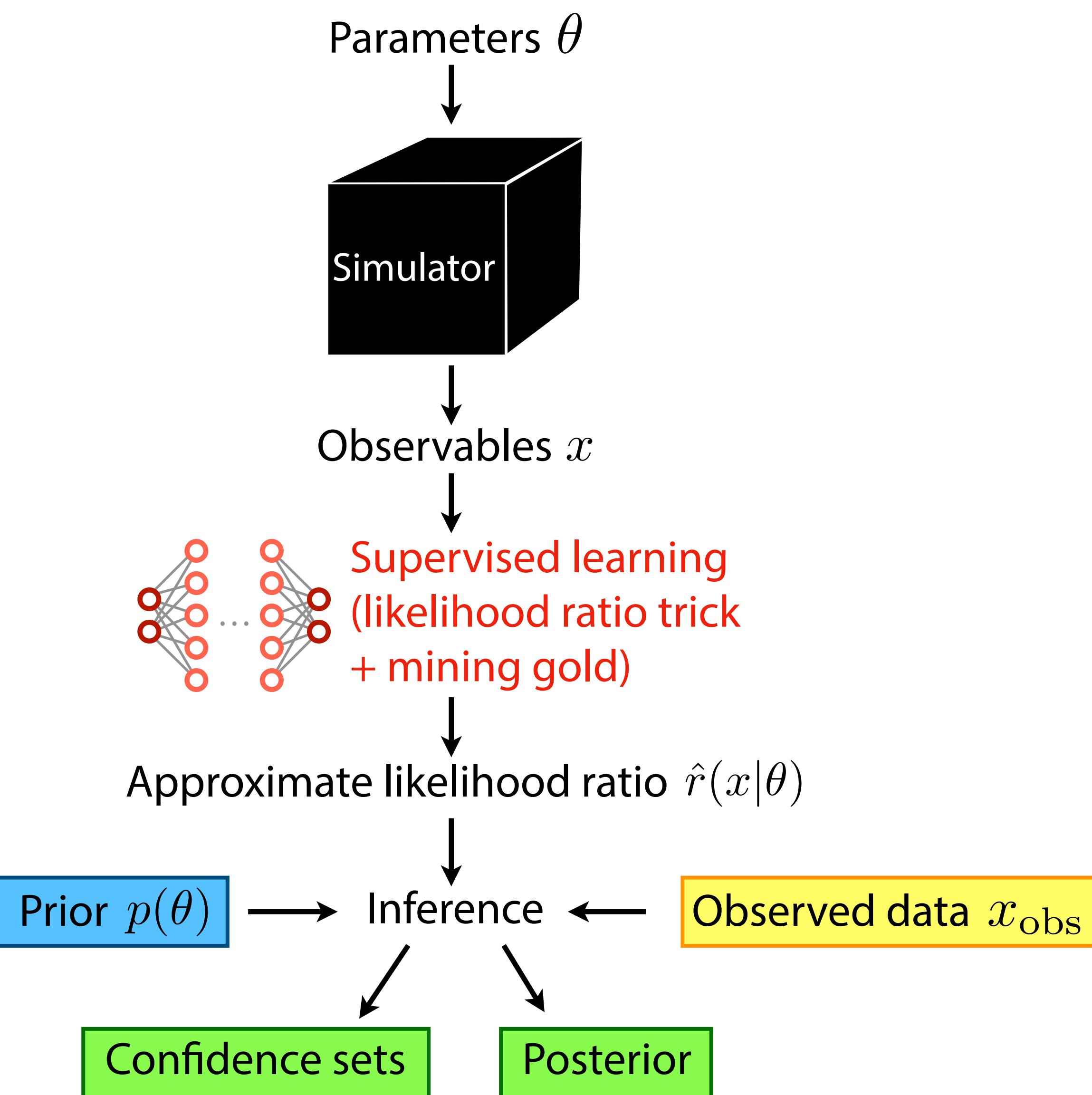
[following T. Collett 1507.02657]



Inference setup

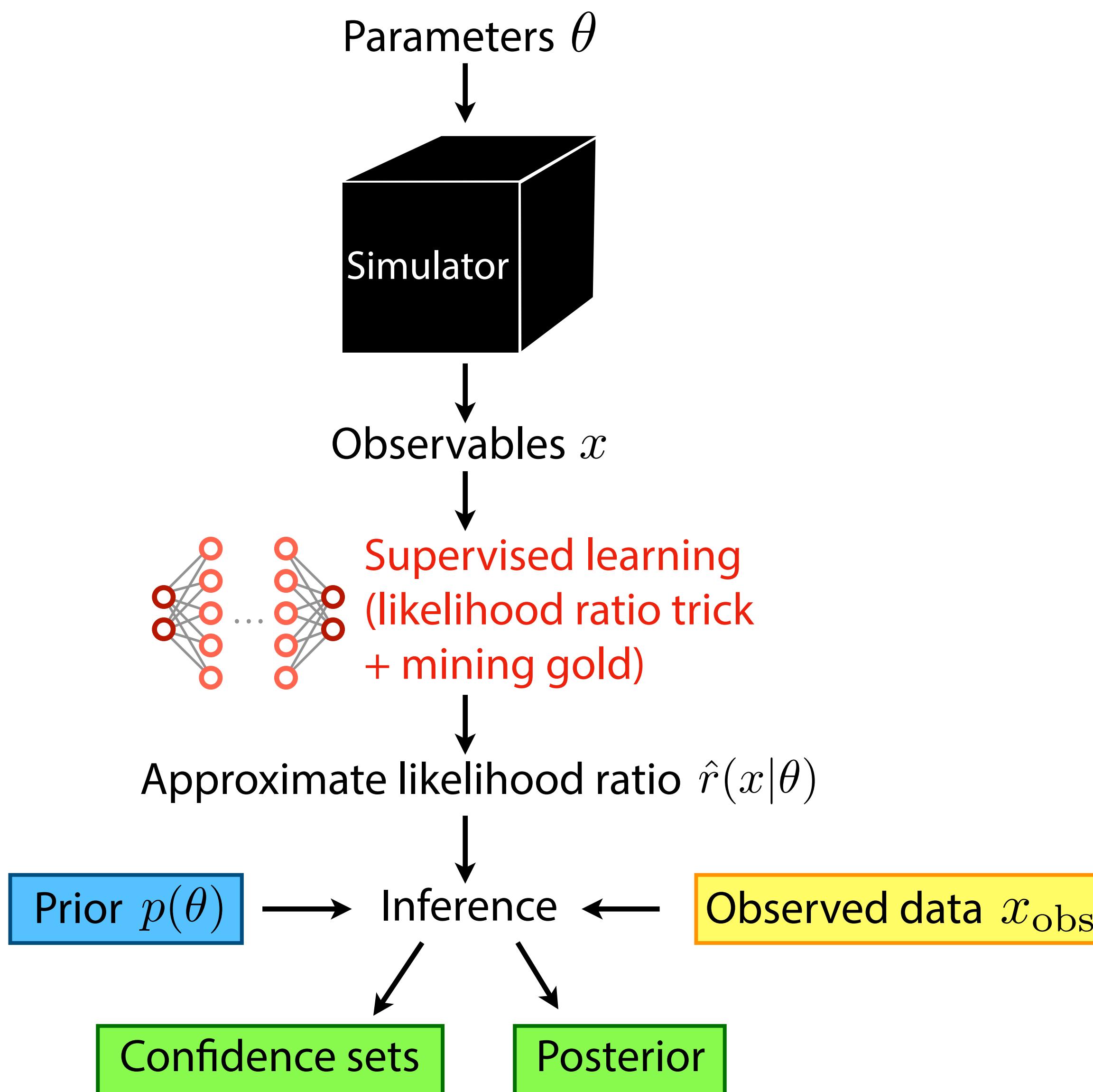


Inference setup



Training data: 10^6 lensed images with
 $0 \leq f_{\text{sub}} \leq 0.2, -1.5 \leq \beta \leq -0.5$

Inference setup

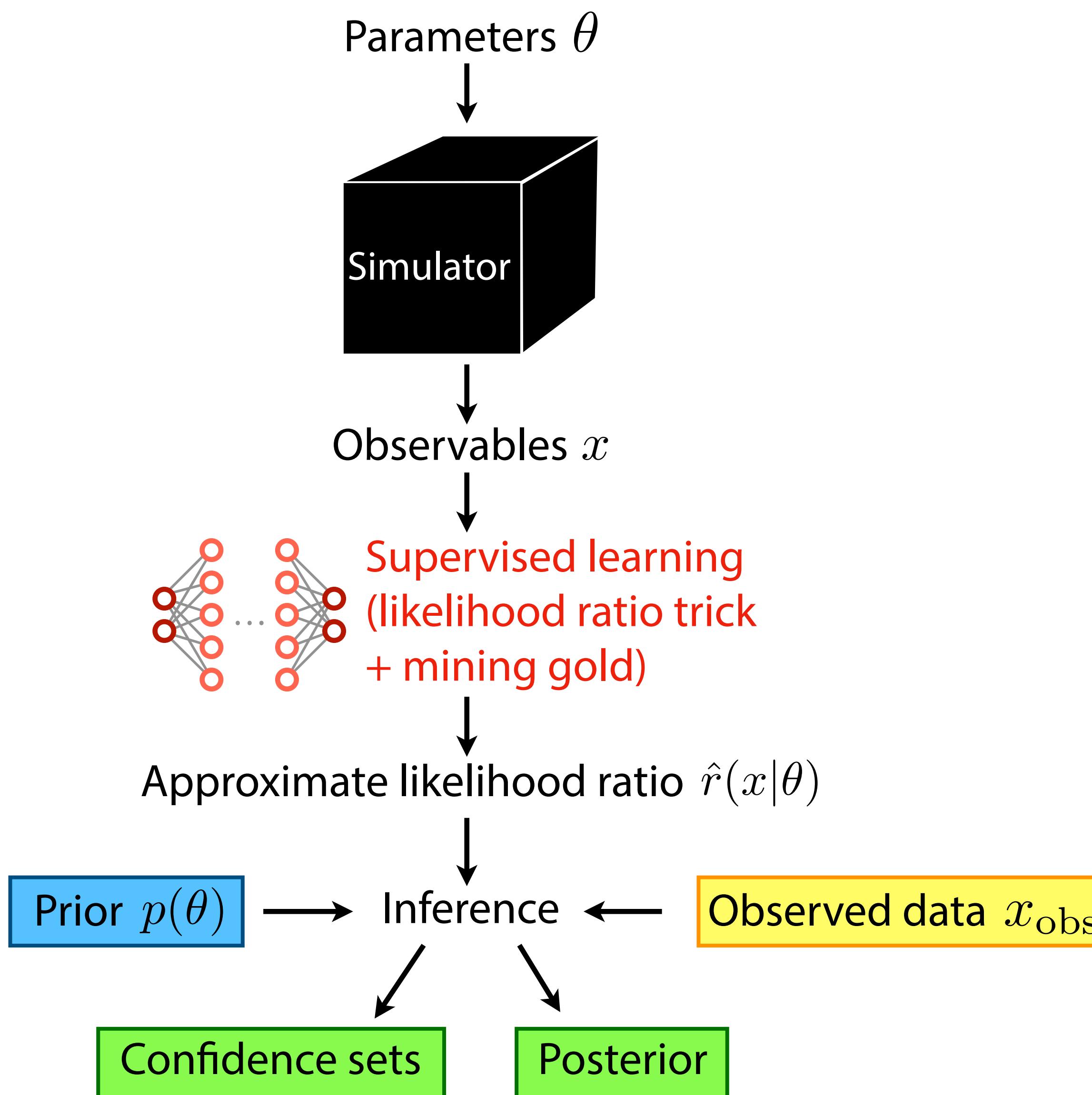


Training data: 10^6 lensed images with
 $0 \leq f_{\text{sub}} \leq 0.2, -1.5 \leq \beta \leq -0.5$

Convolutional neural network (modified ResNet-18)
trained on ALICES loss
[M. Stoye, JB, J. Pavez, G, Louppe, K. Cranmer 1808.00973]

Probabilistic calibration of network output

Inference setup



Training data: 10^6 lensed images with
 $0 \leq f_{\text{sub}} \leq 0.2, -1.5 \leq \beta \leq -0.5$

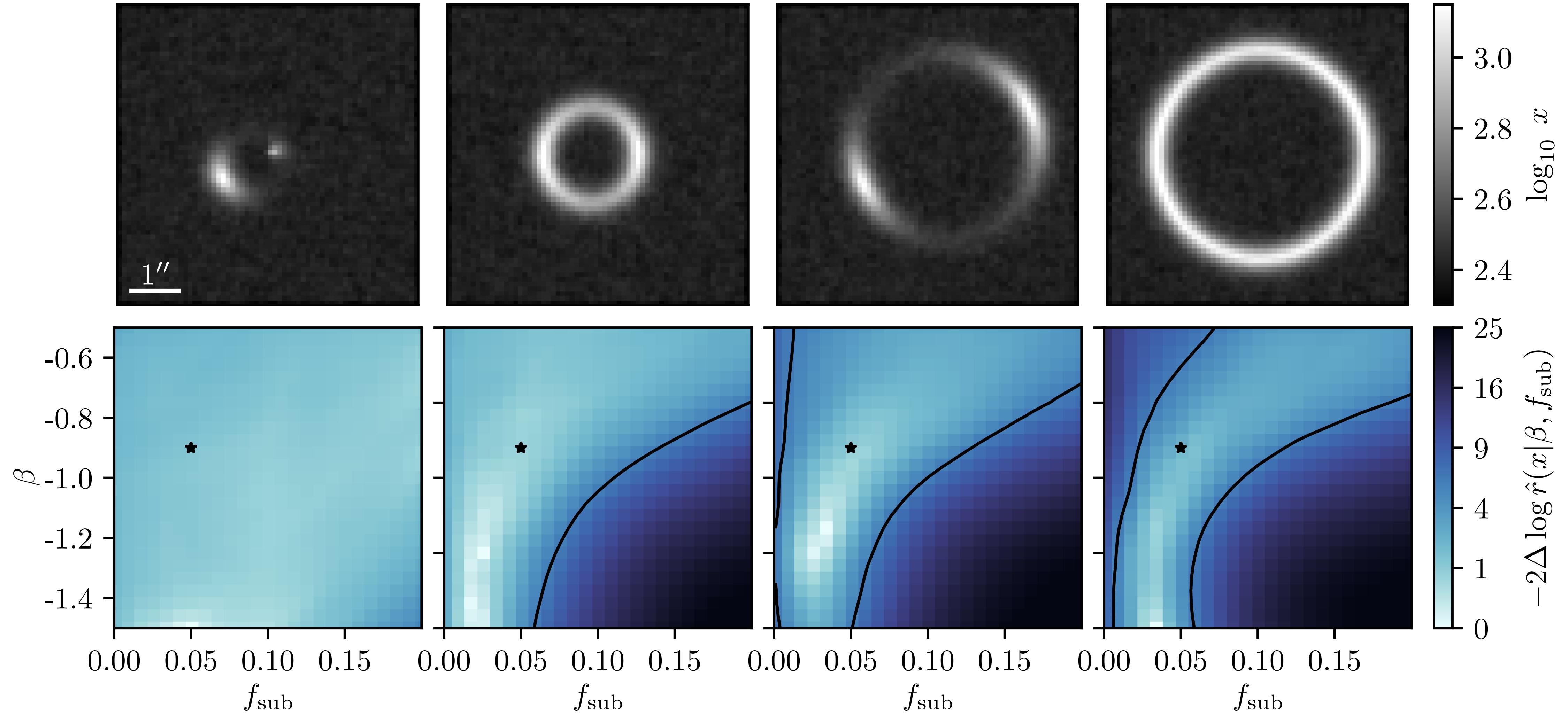
Convolutional neural network (modified ResNet-18)
trained on ALICES loss
[M. Stoye, JB, J. Pavez, G, Louppe, K. Cranmer 1808.00973]

Probabilistic calibration of network output

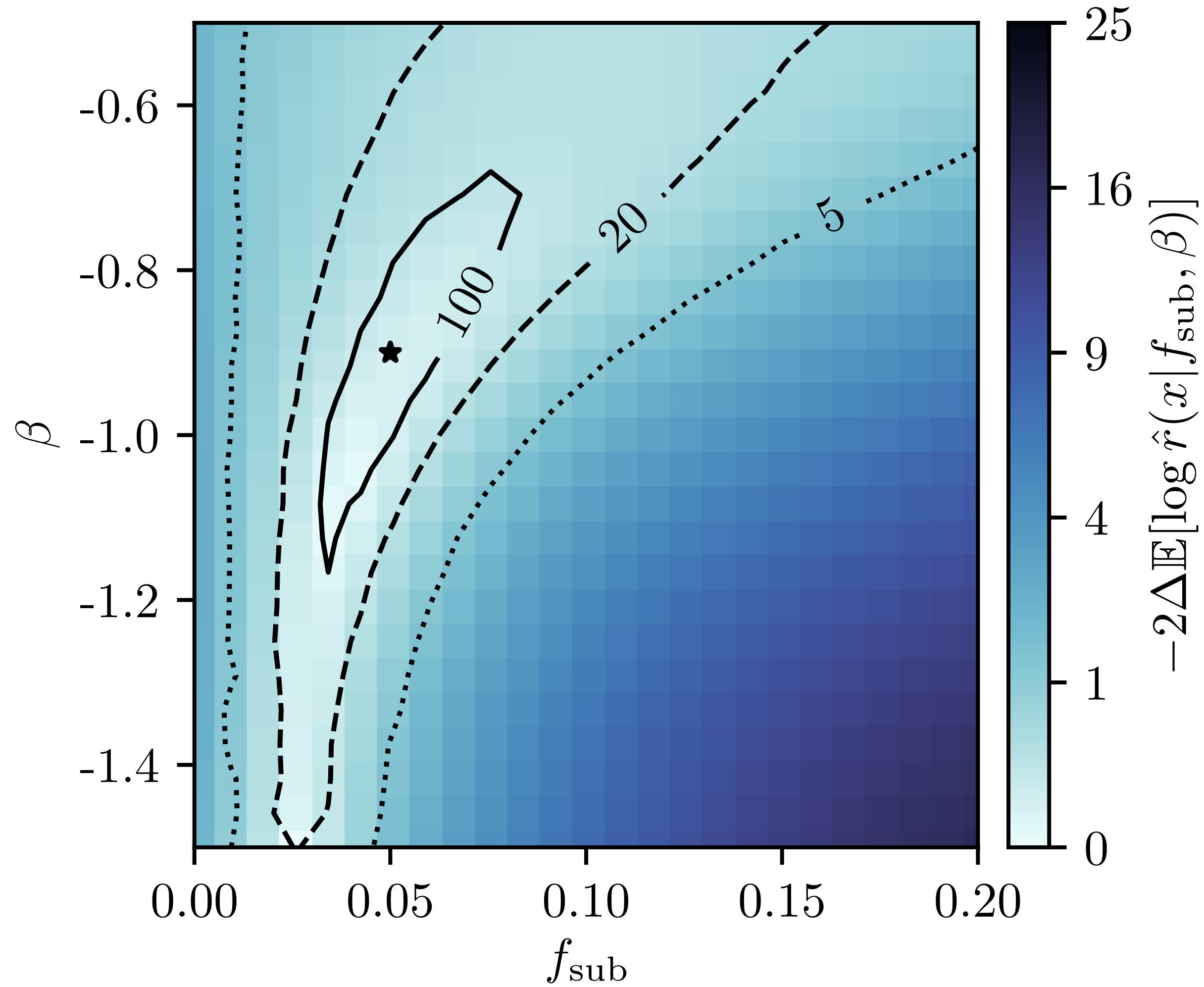
Synthetic “observed” data set simulated for
 $f_{\text{sub}} = 0.05, \beta = -0.9$

Bayesian & frequentist inference

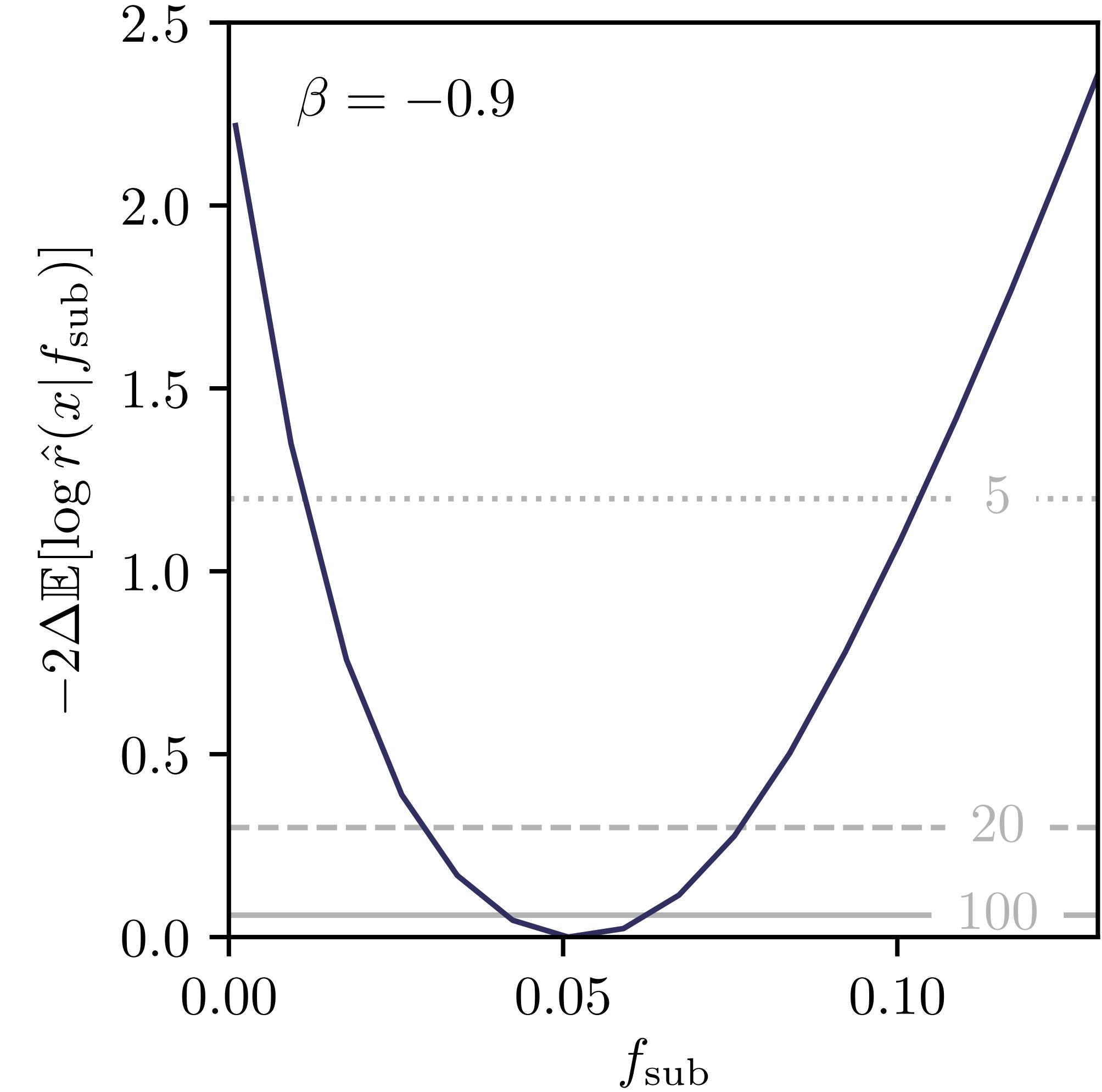
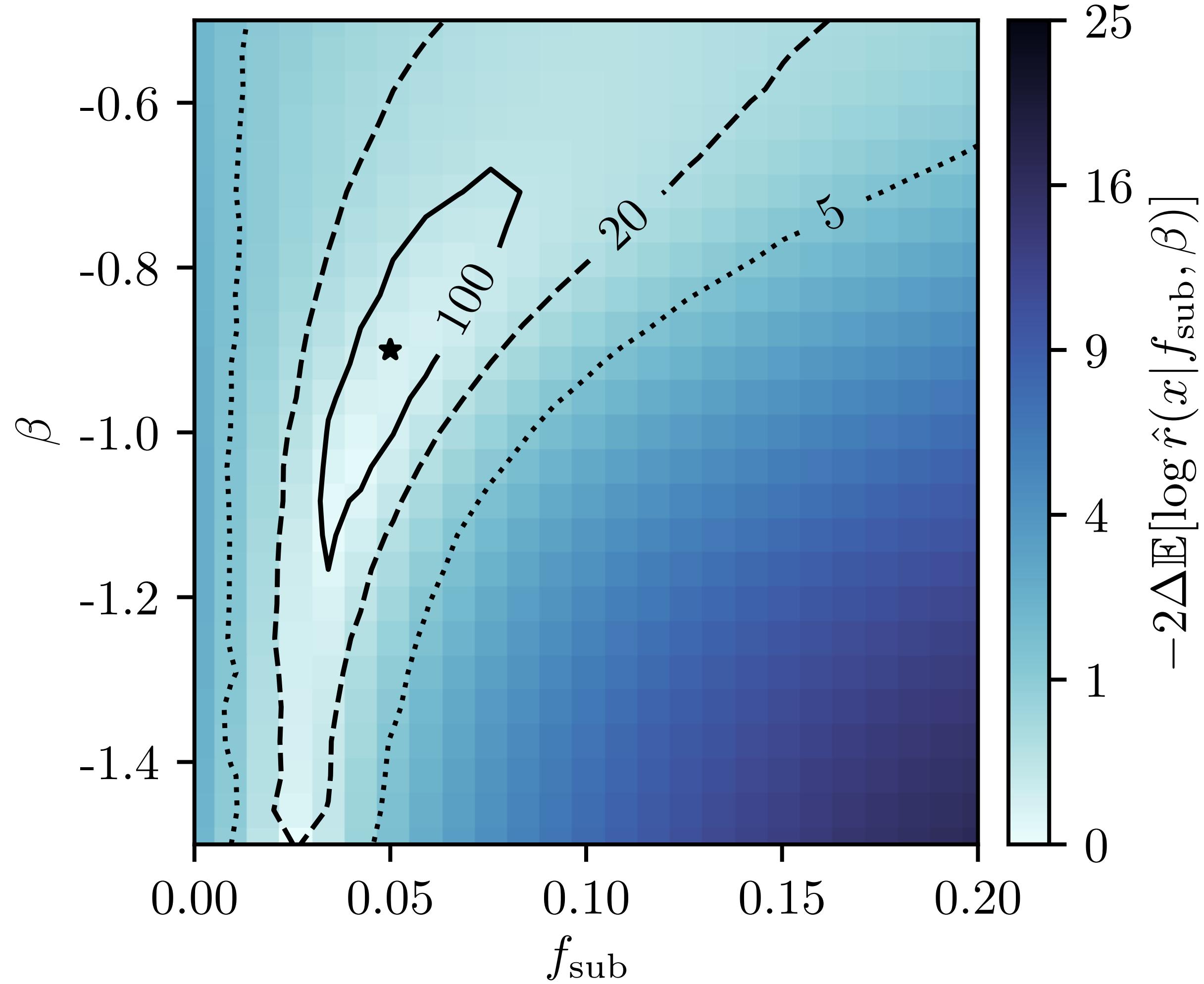
Inferring parameters from individual images



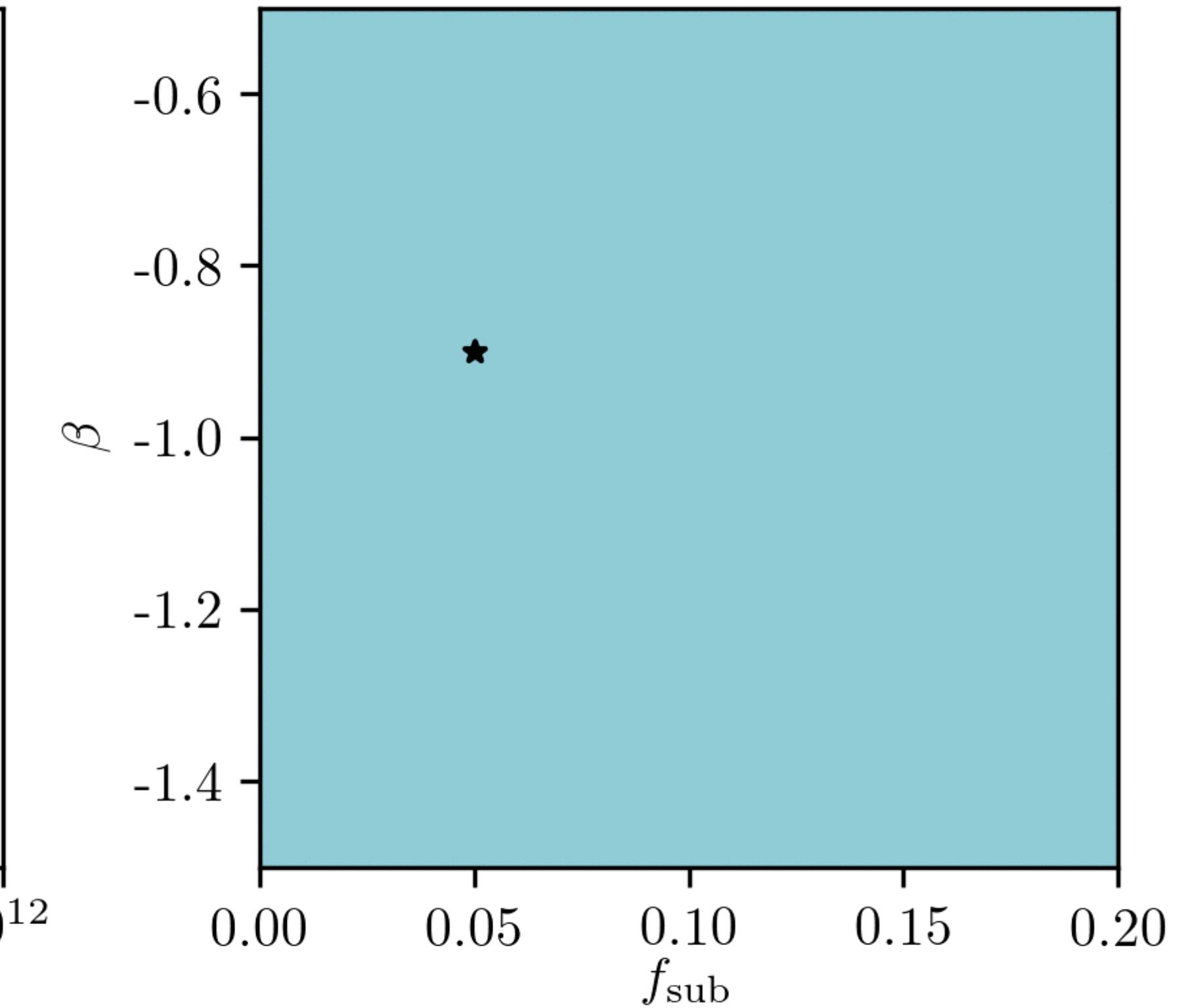
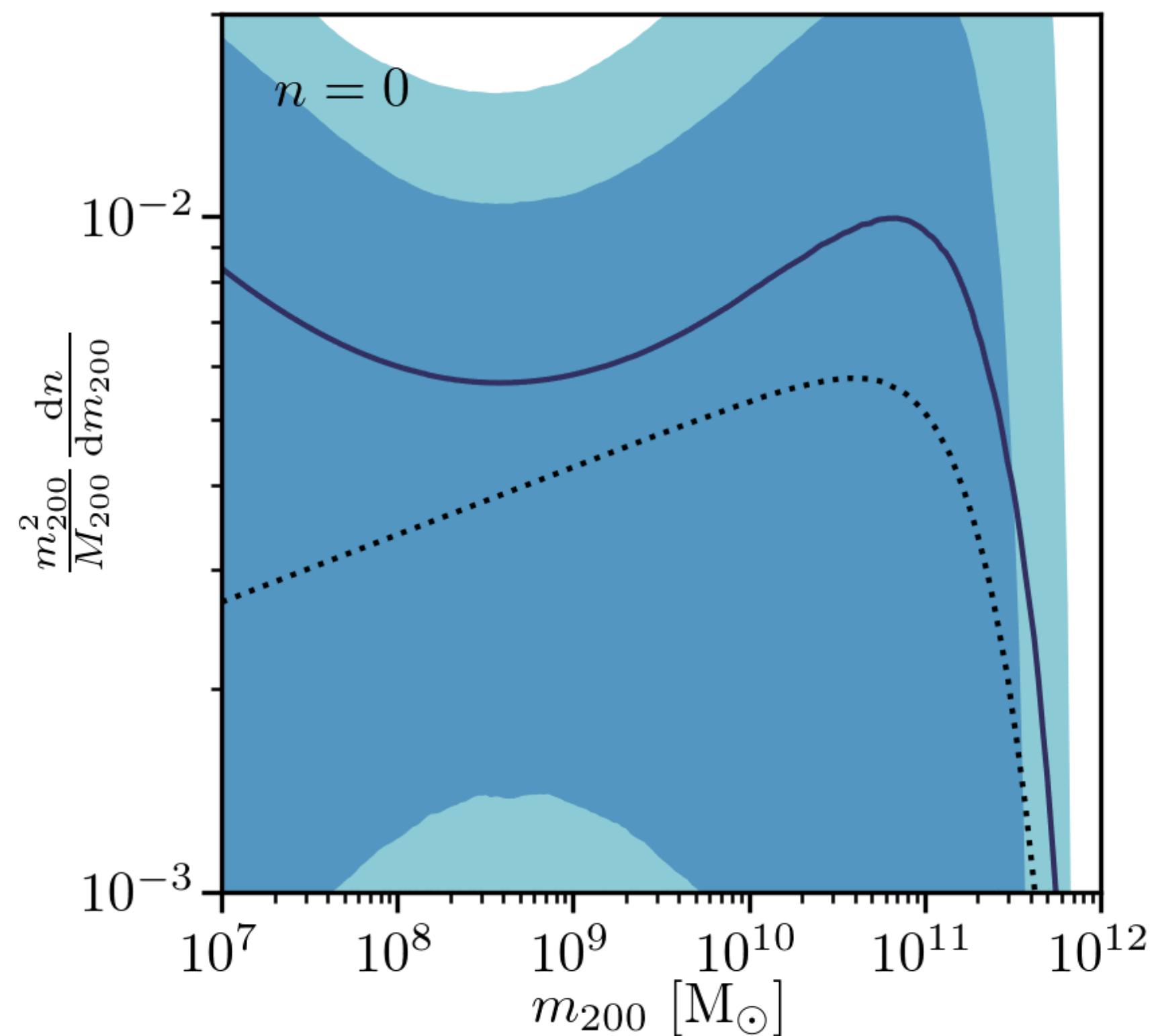
Expected likelihood ratio map



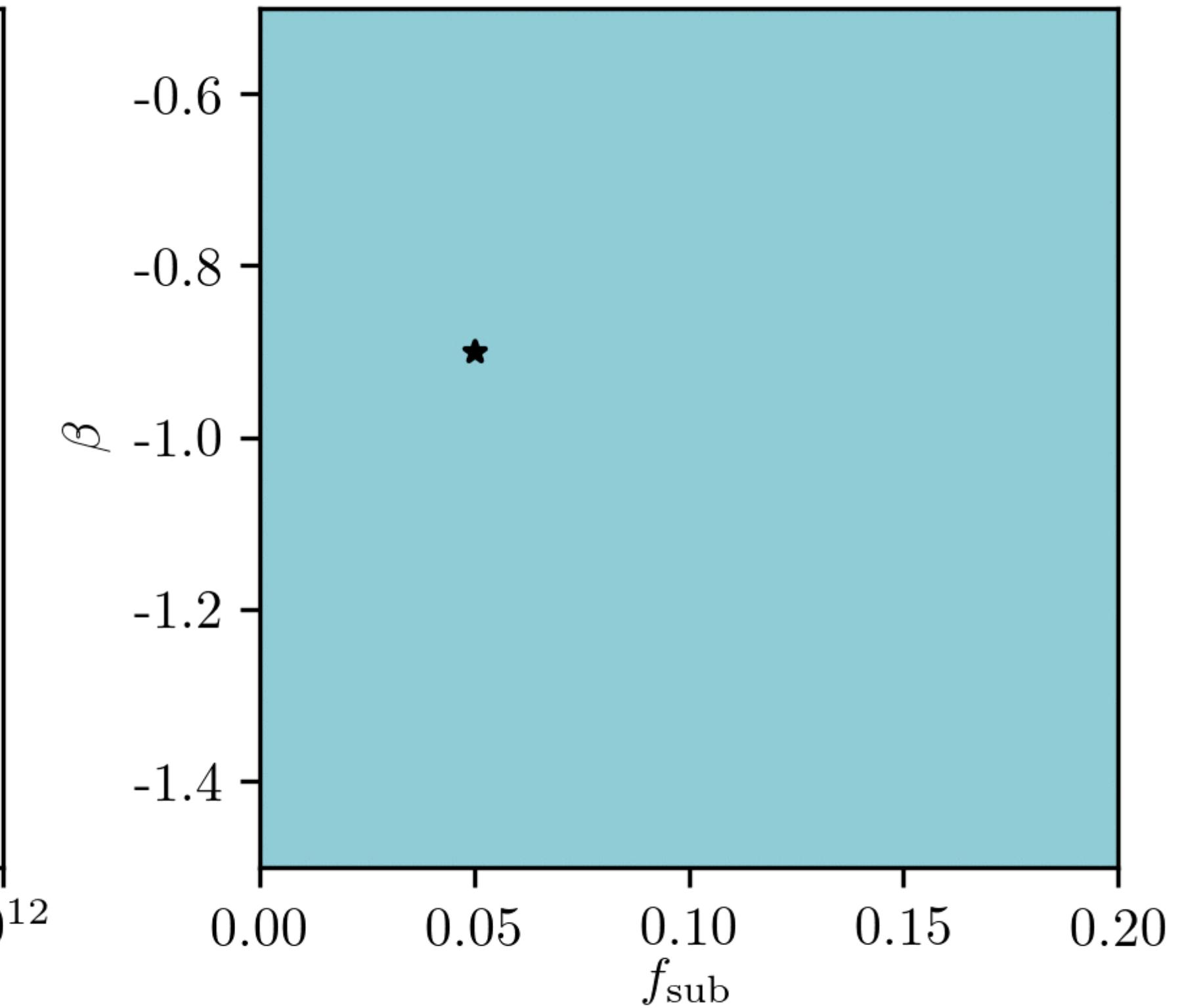
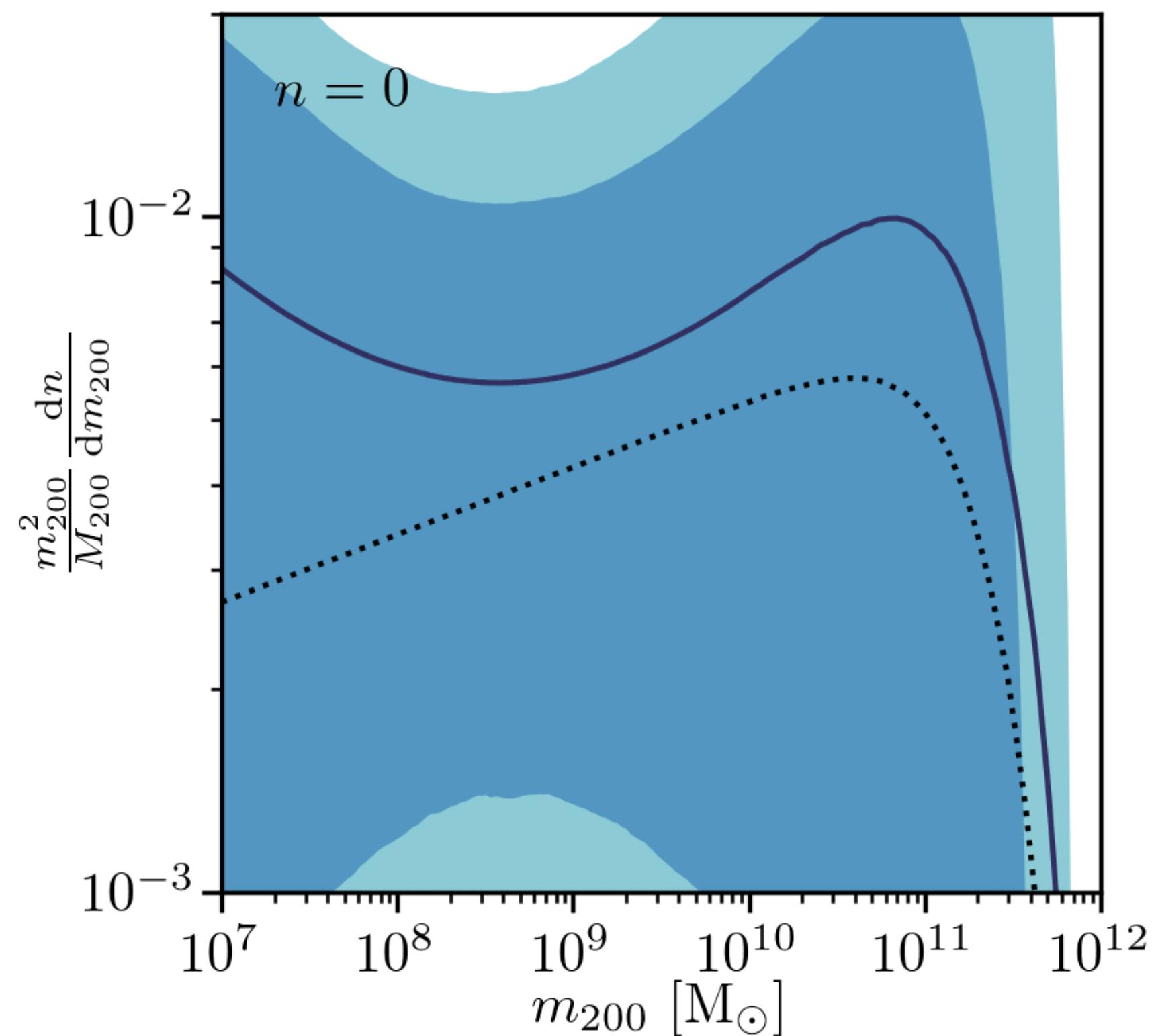
Expected likelihood ratio map

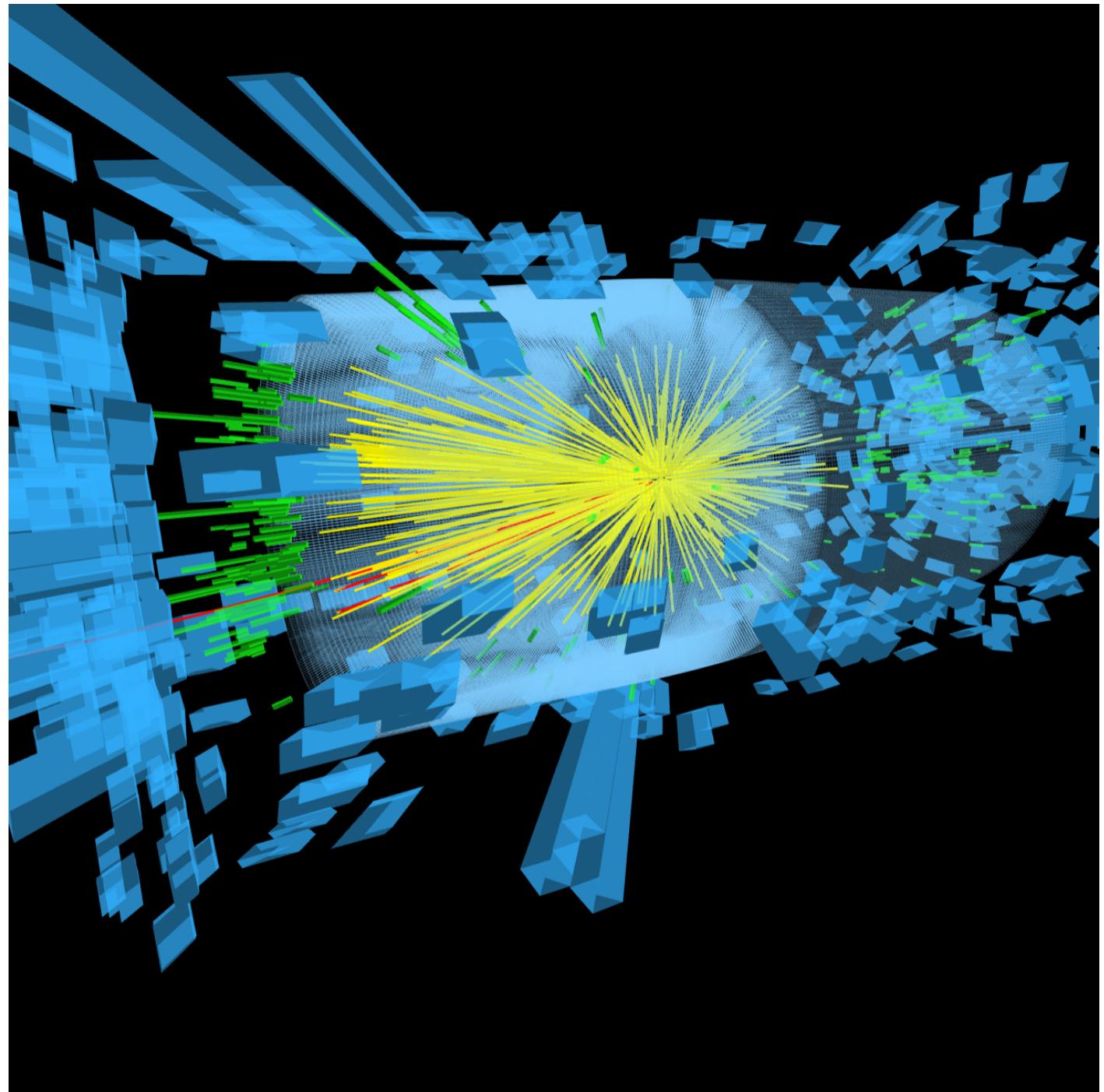


Bayesian inference

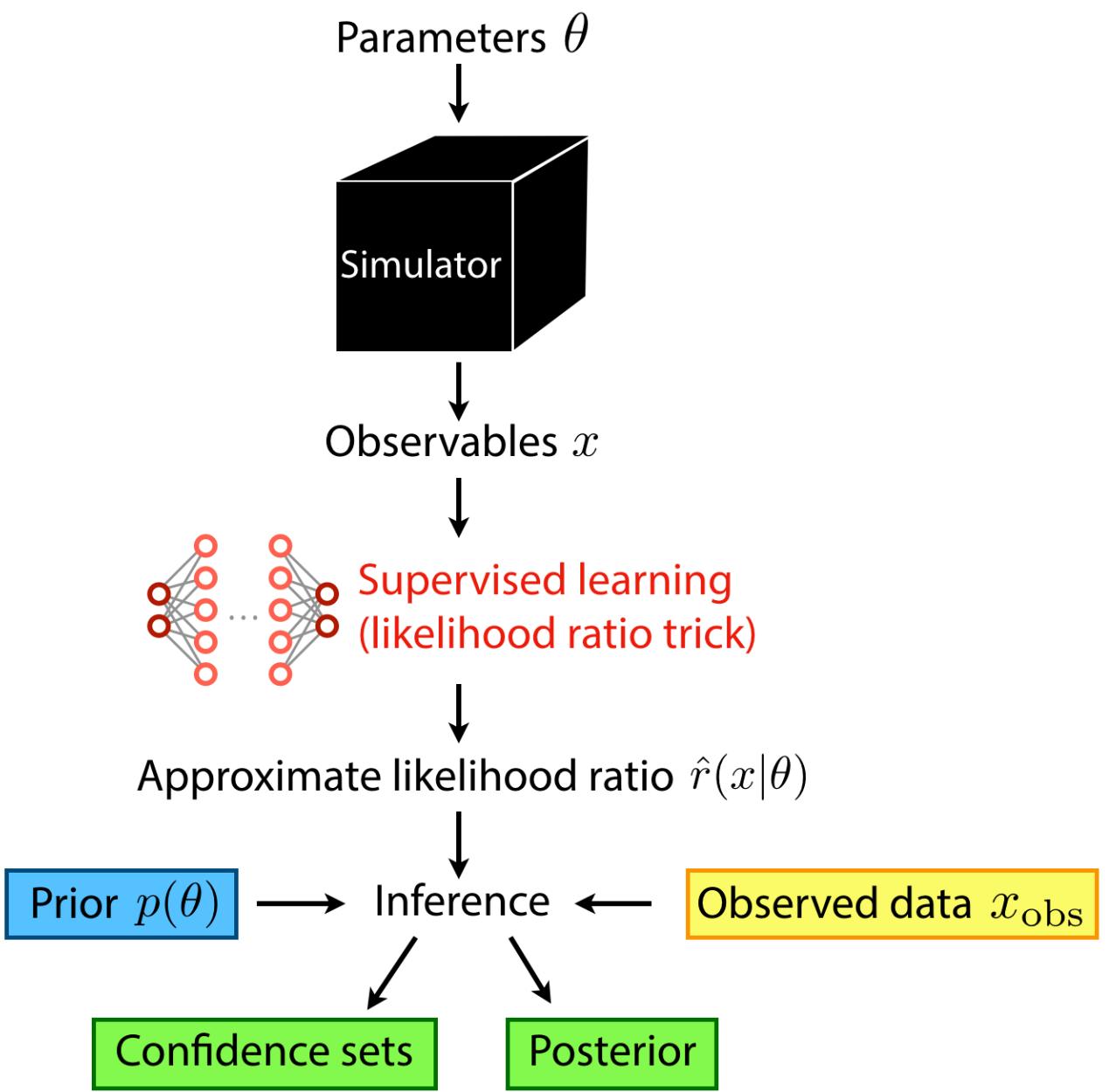


Bayesian inference

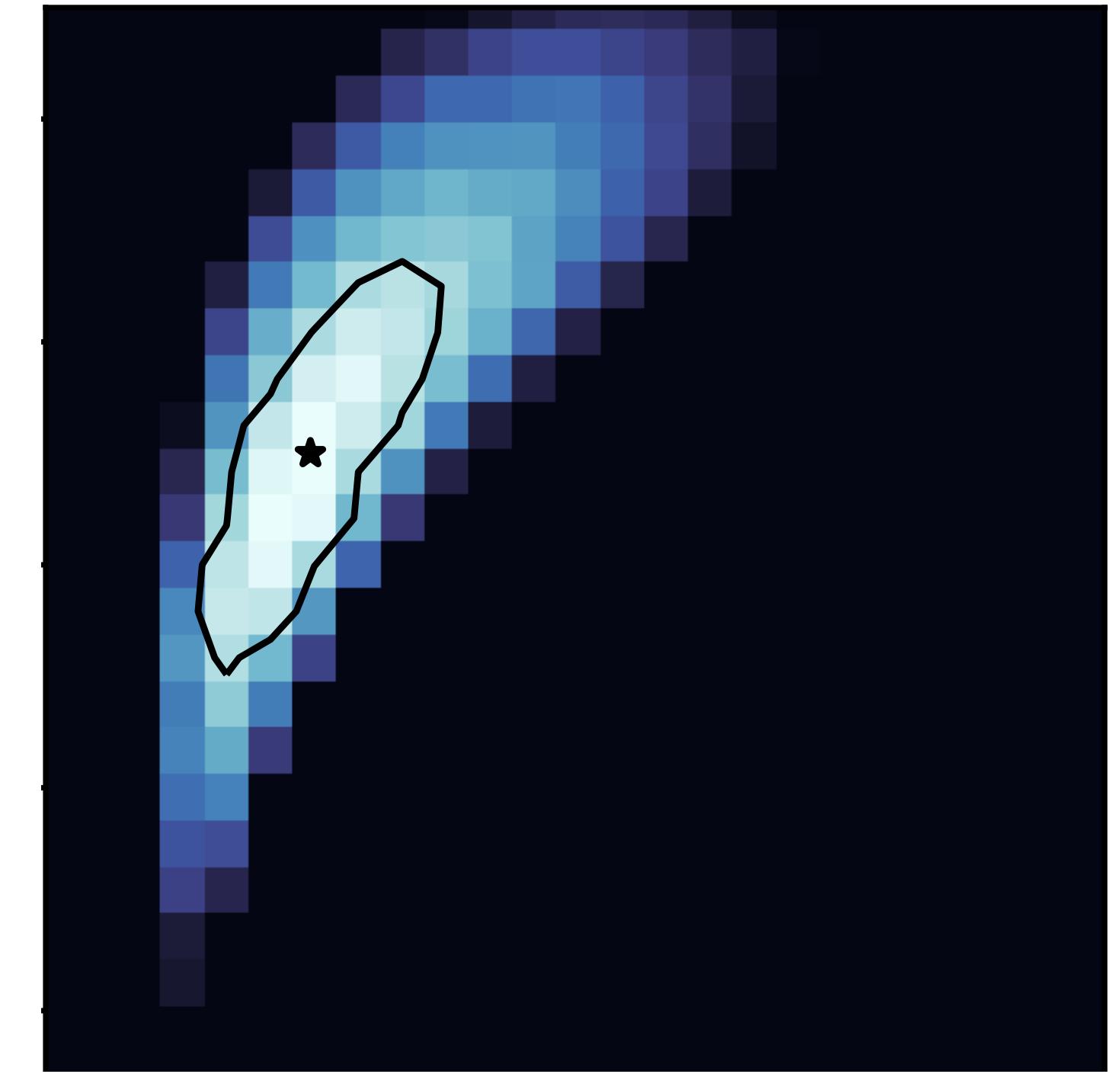




Inference is challenging when phenomena are modeled with computer simulations and the data are high-dimensional



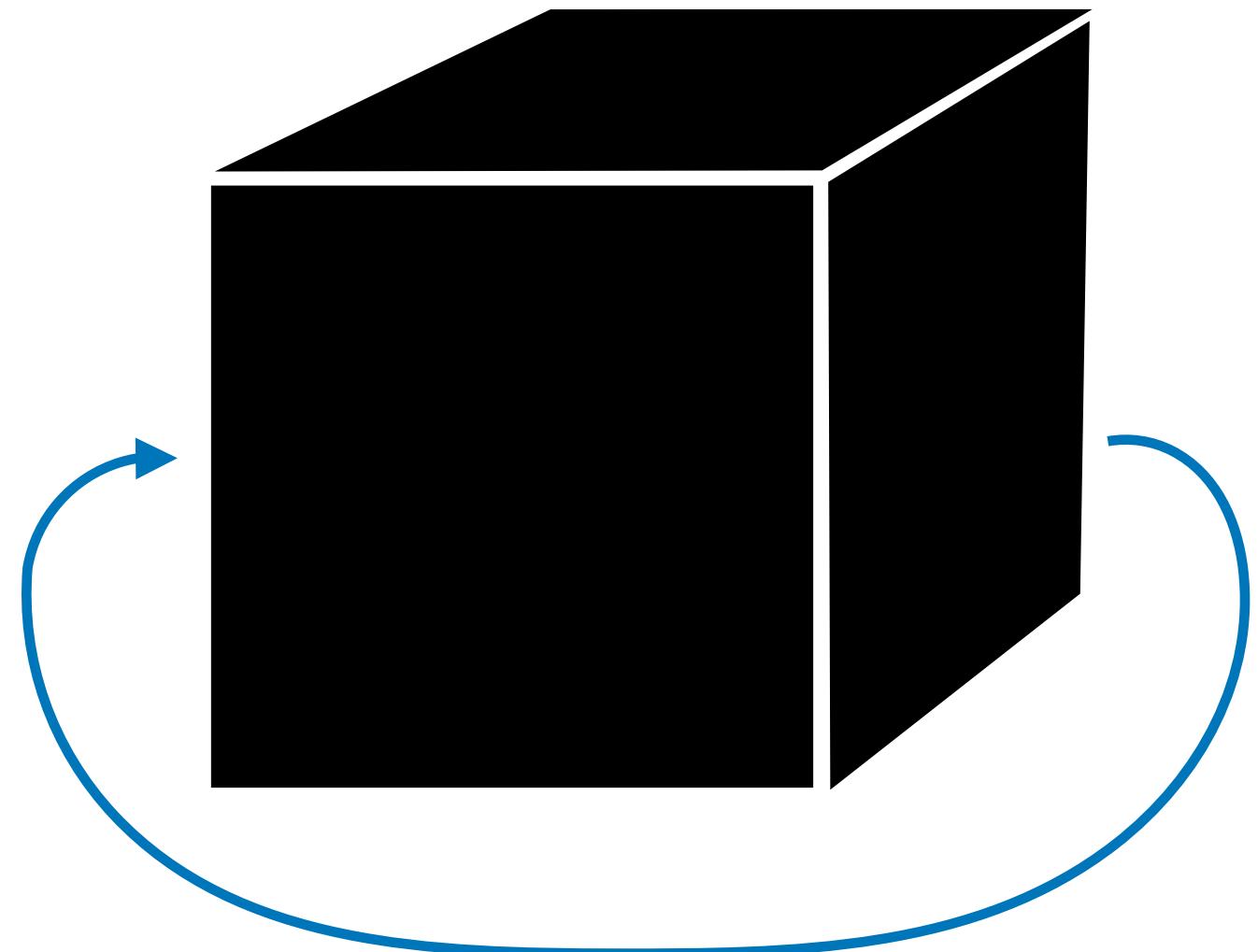
New inference algorithms are based on normalizing flows or the likelihood ratio trick; we can improve them with physics understanding



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

Frontiers of simulation-based inference

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

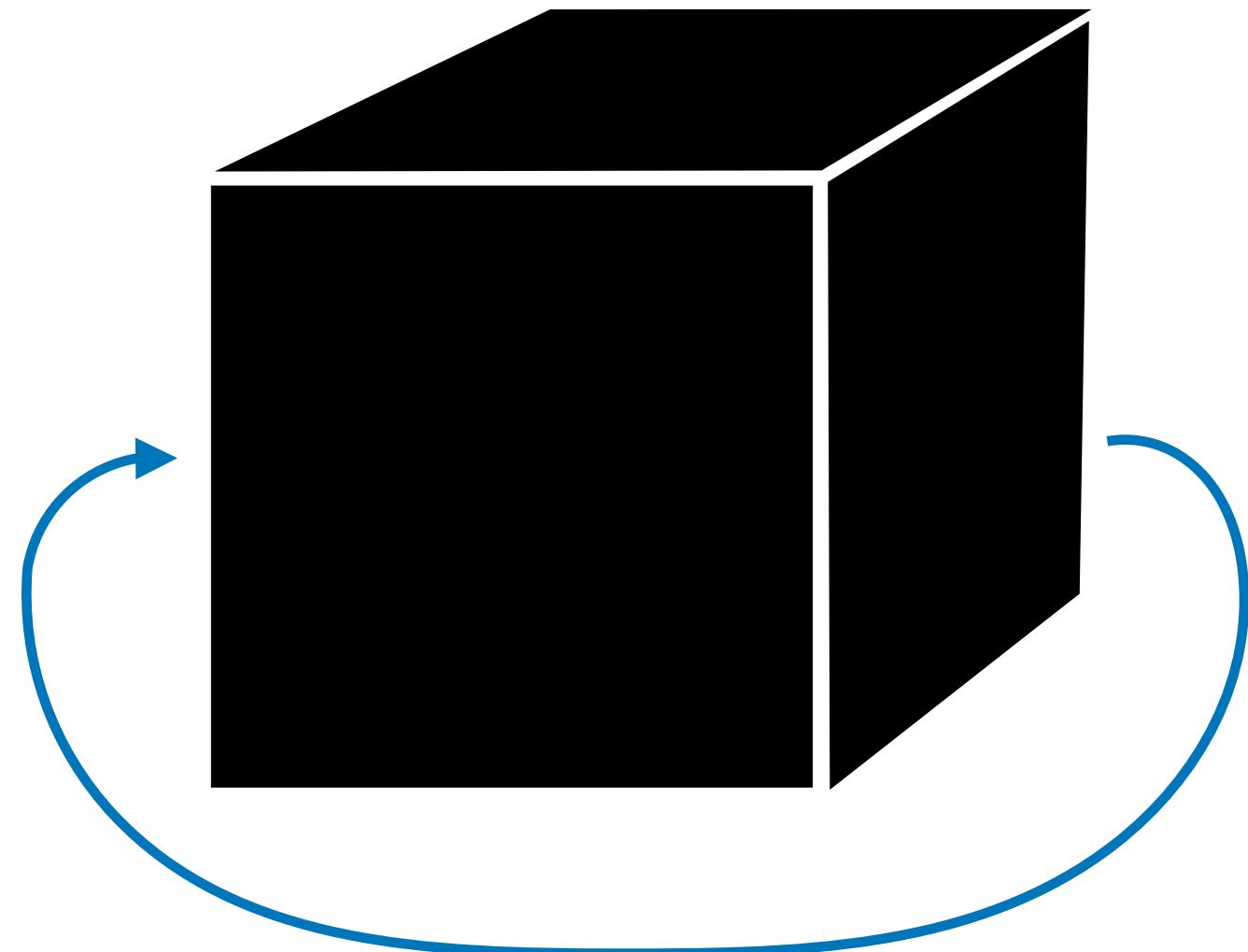


Active learning:

Iteratively guide simulator to
important parameter regions
based on past results

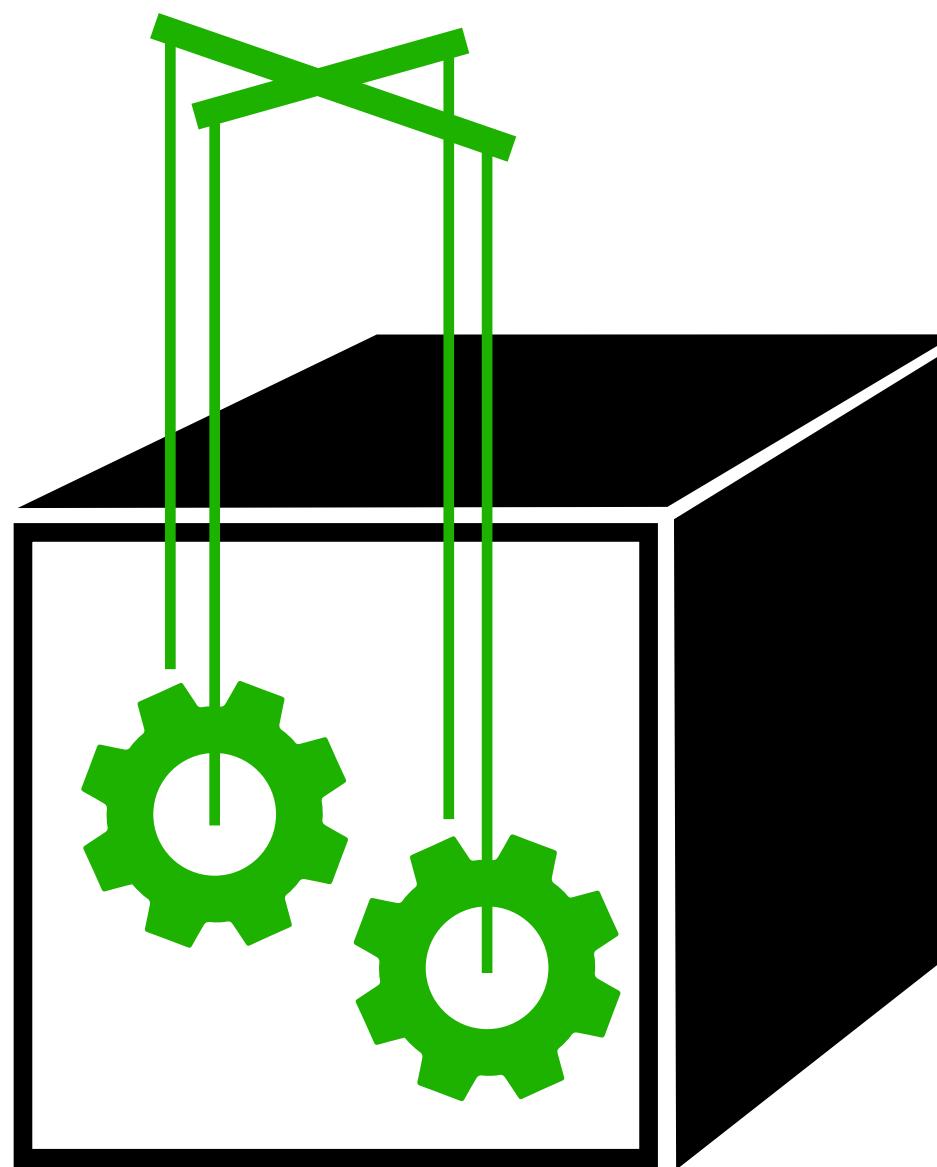
Frontiers of simulation-based inference

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



Active learning:

Iteratively guide simulator to important parameter regions based on past results

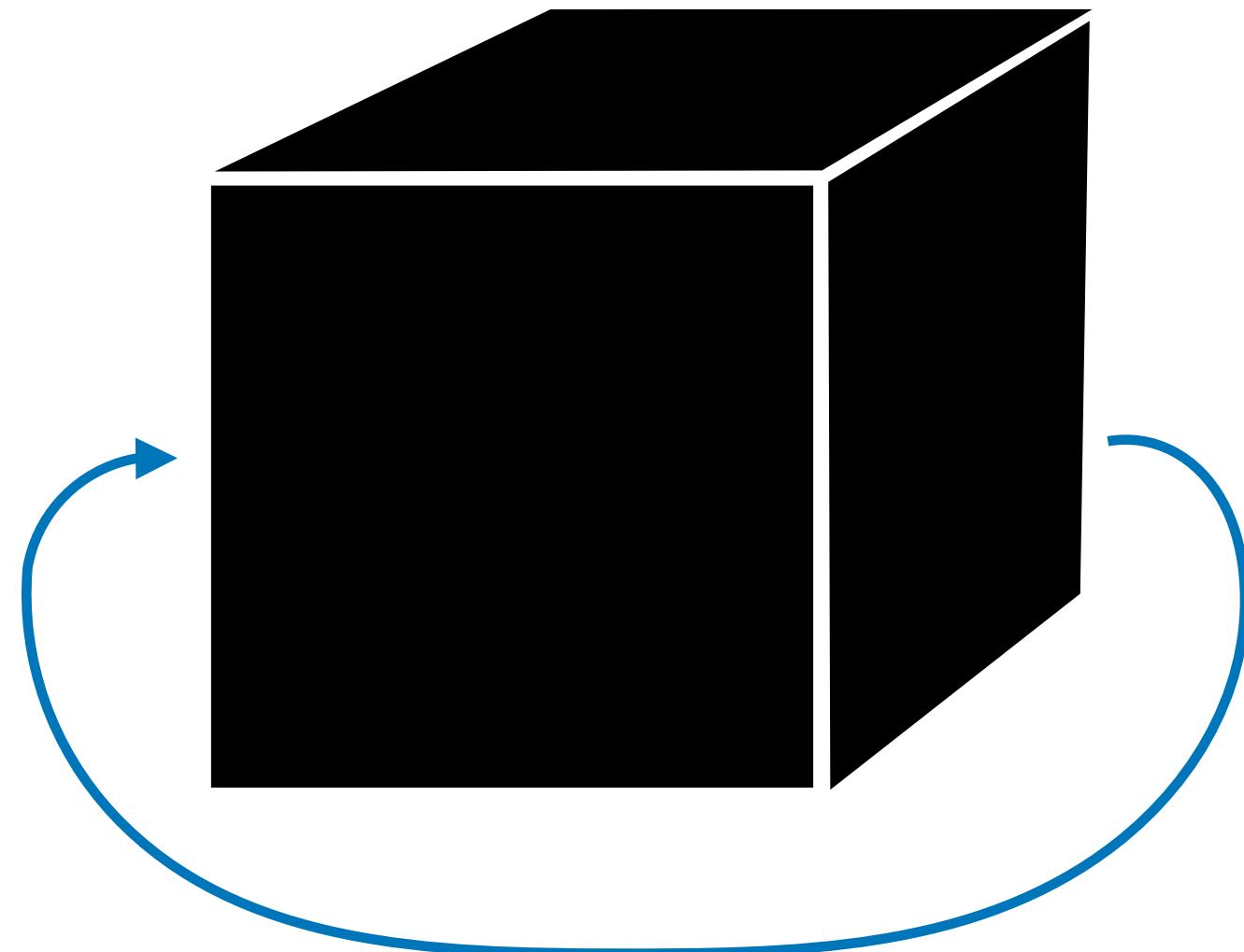


Probabilistic programming:

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

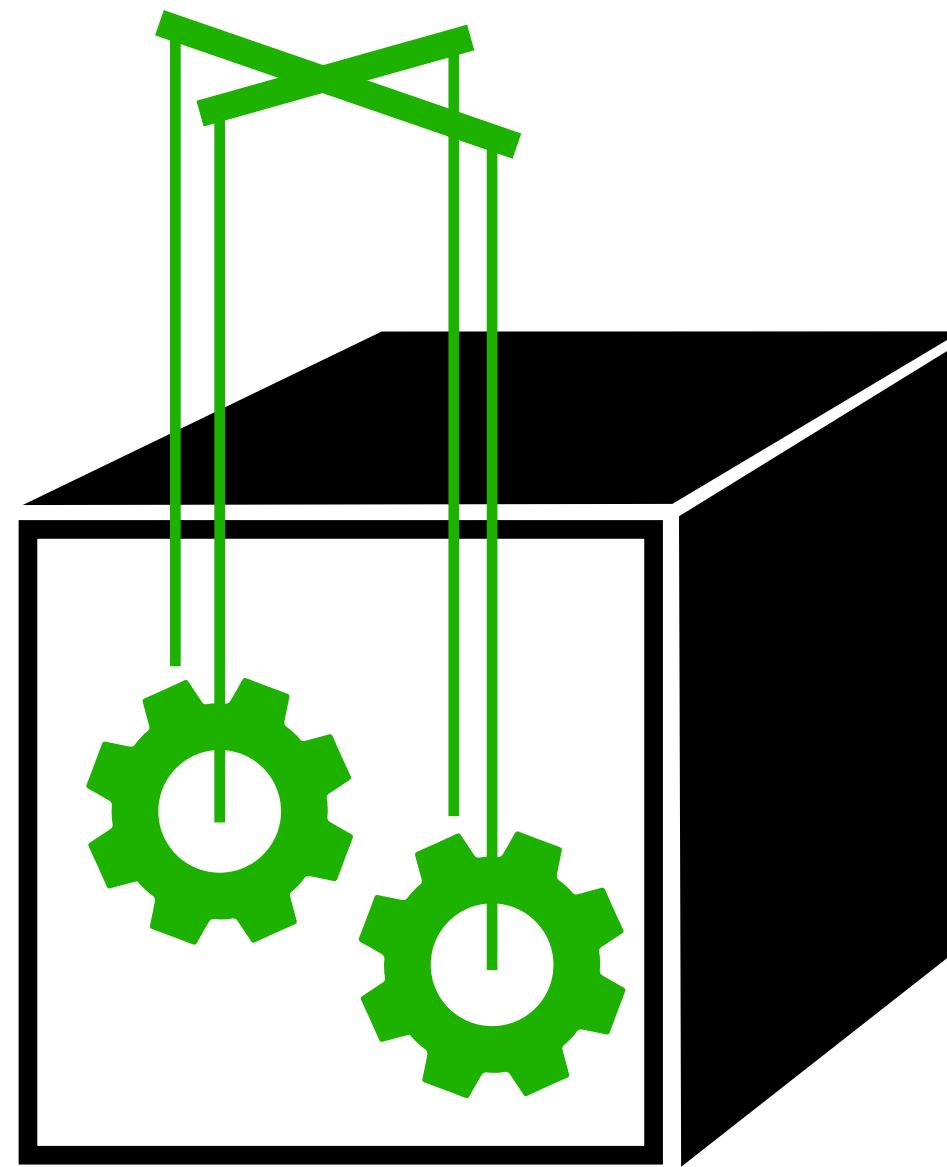
Frontiers of simulation-based inference

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



Active learning:

Iteratively guide simulator to important parameter regions based on past results



Probabilistic programming:

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



Mining gold:

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

References

Opinionated review

K. Cranmer, JB, G. Louppe:
“The frontier of simulation-based inference”
[1911.01429]

Do It Yourself (for LHC physics)

JB, F. Kling, I. Espejo, K. Cranmer:
“MadMiner: Machine learning—based inference for particle physics”
[CSBS, 1907.10621, <https://github.com/diana-hep/madminer>]

Strong lensing

JB, S. Mishra-Sharma, J. Hermans, G. Louppe, K. Cranmer
“Mining for Dark Matter Substructure: Inferring subhalo population properties
from strong lenses with machine learning”
[ApJ, 1909.02005]

LHC HXSWG YR4 STXS

JB, S. Dawson, S. Homiller, F. Kling, T. Plehn:
“Benchmarking simplified template cross sections in WH production”
[JHEP, 1908.06980]

Meh, the original trilogy was better

JB, K. Cranmer, G. Louppe, J. Pavez:
“Constraining Effective Field Theories with
machine learning”
[PRL, 1805.00013]

JB, K. Cranmer, G. Louppe, J. Pavez:
“A guide to constraining Effective Field Theories
with machine learning”
[PRD, 1805.00020]

JB, G. Louppe, J. Pavez, K. Cranmer:
“Mining gold from implicit models to improve
likelihood-free inference”
[PNAS, 1805.12244]

Follow-up with incremental improvements

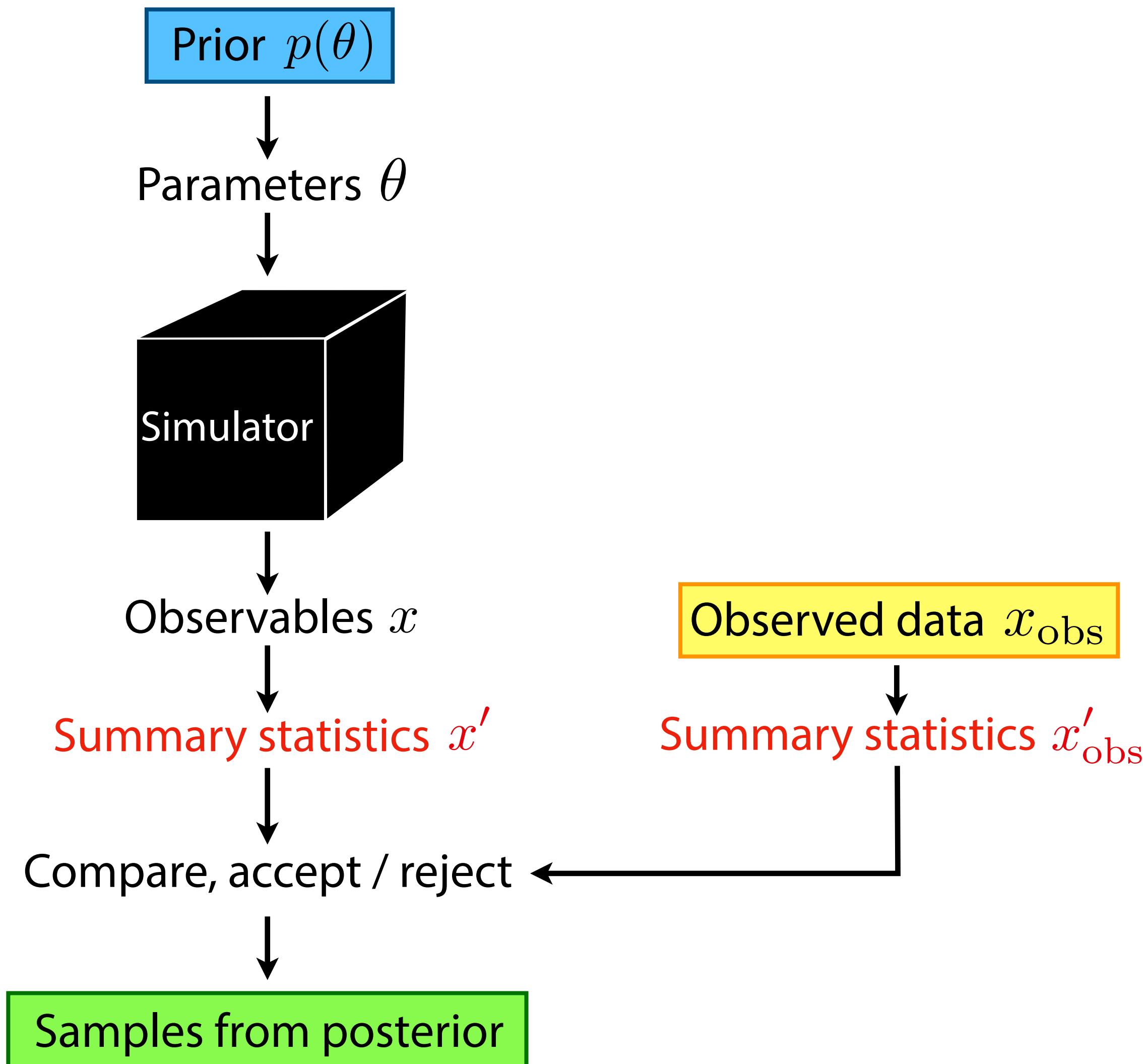
M. Stoye, JB, K. Cranmer, G. Louppe, J. Pavez:
“Likelihood-free inference with an improved
cross-entropy estimator”
[NeurIPS workshop, 1808.00973]

“Hooray! Question mark?”

— Todd Chavez

Approximate Bayesian Computation (ABC)

[D. Rubin 1984]

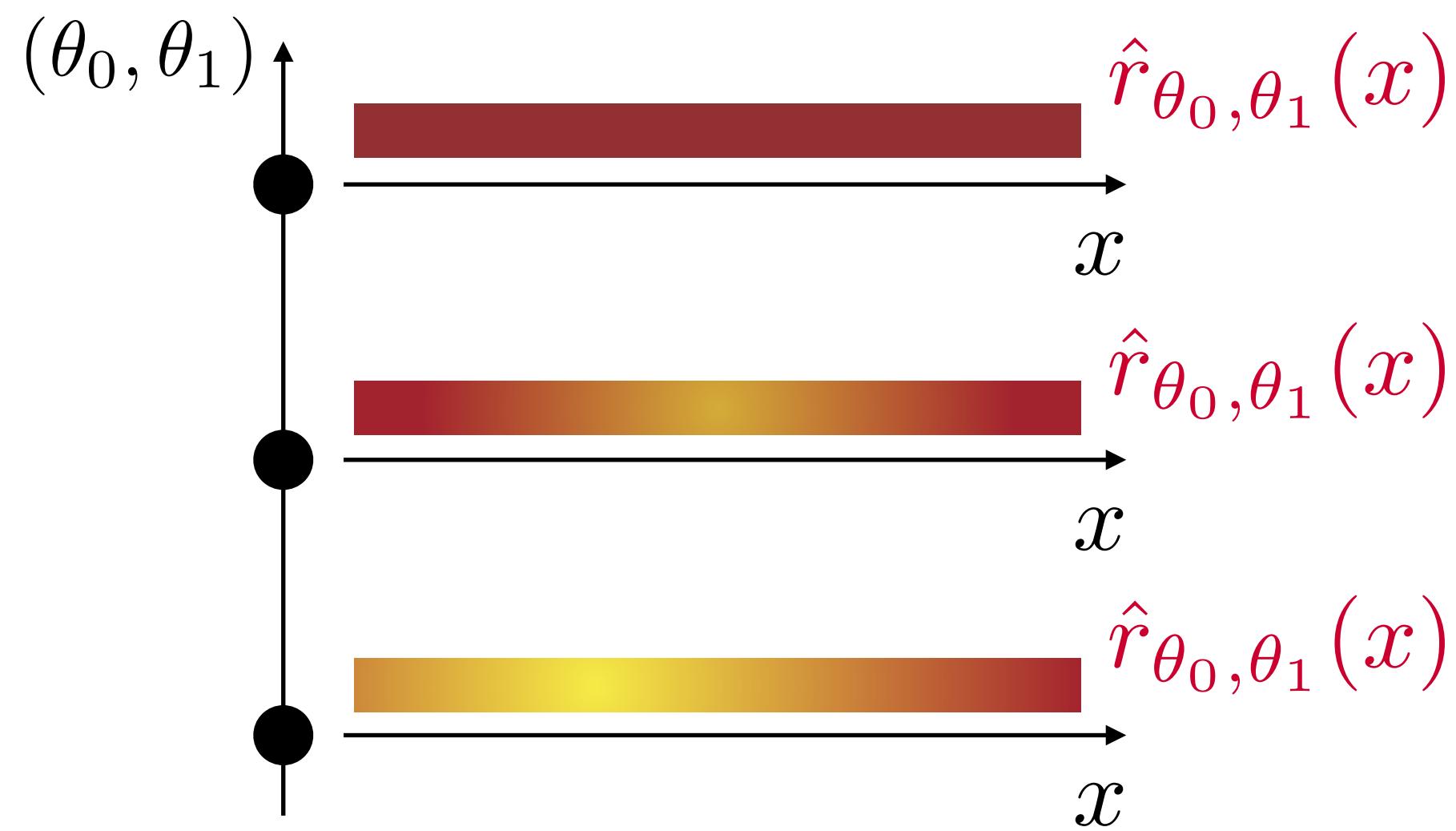


- Compression to summary statistics and acceptance threshold reduce quality of inference
- Rejection algorithm can be very sample inefficient

Two types of likelihood ratio estimators

A) Point by point:

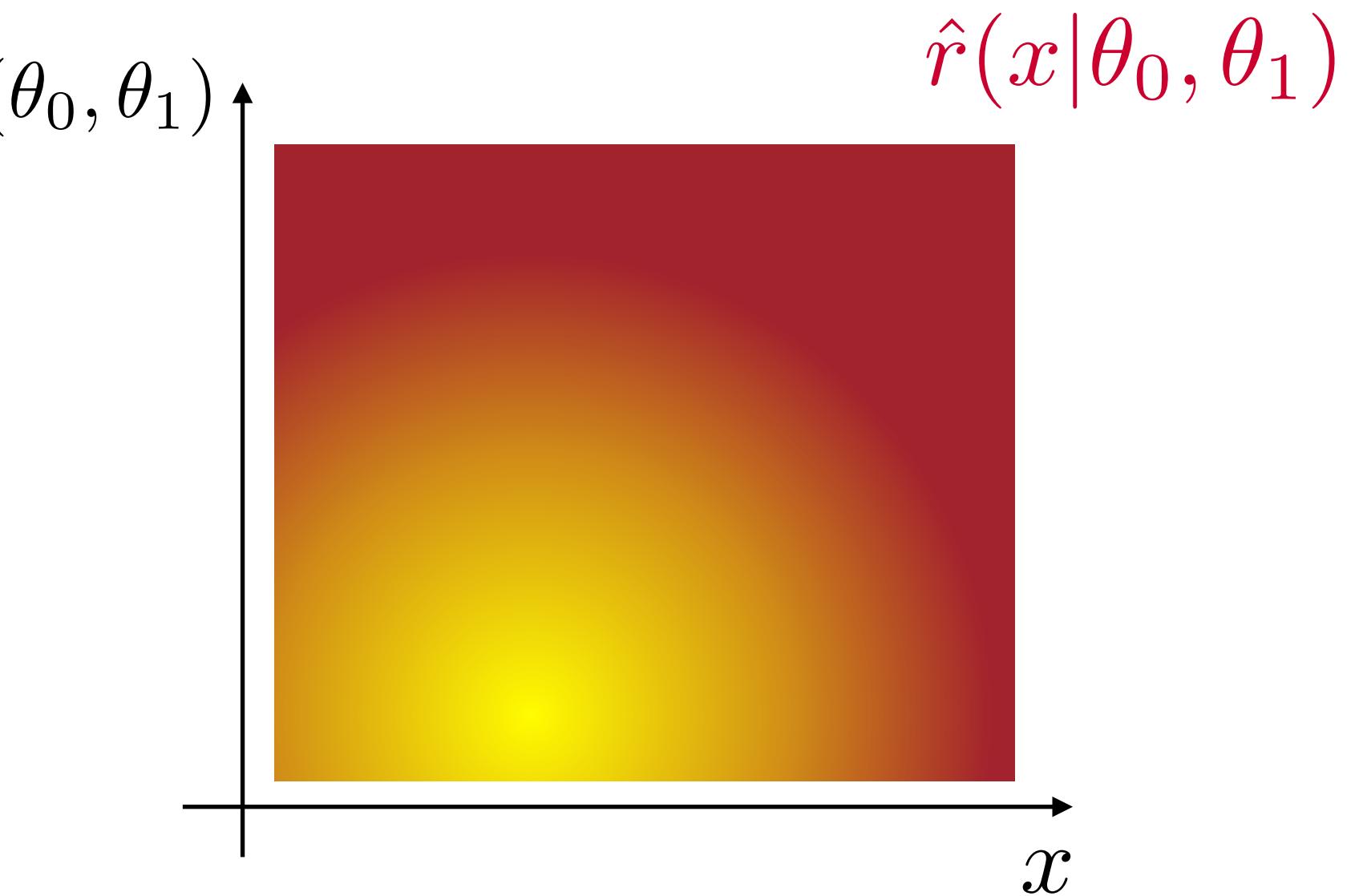
- first, define grid of parameter points $\{(\theta_0, \theta_1)\}$
- for each combination (θ_0, θ_1) , create separate estimator $\hat{r}_{\theta_0, \theta_1}(x)$
- final results can be interpolated between grid points



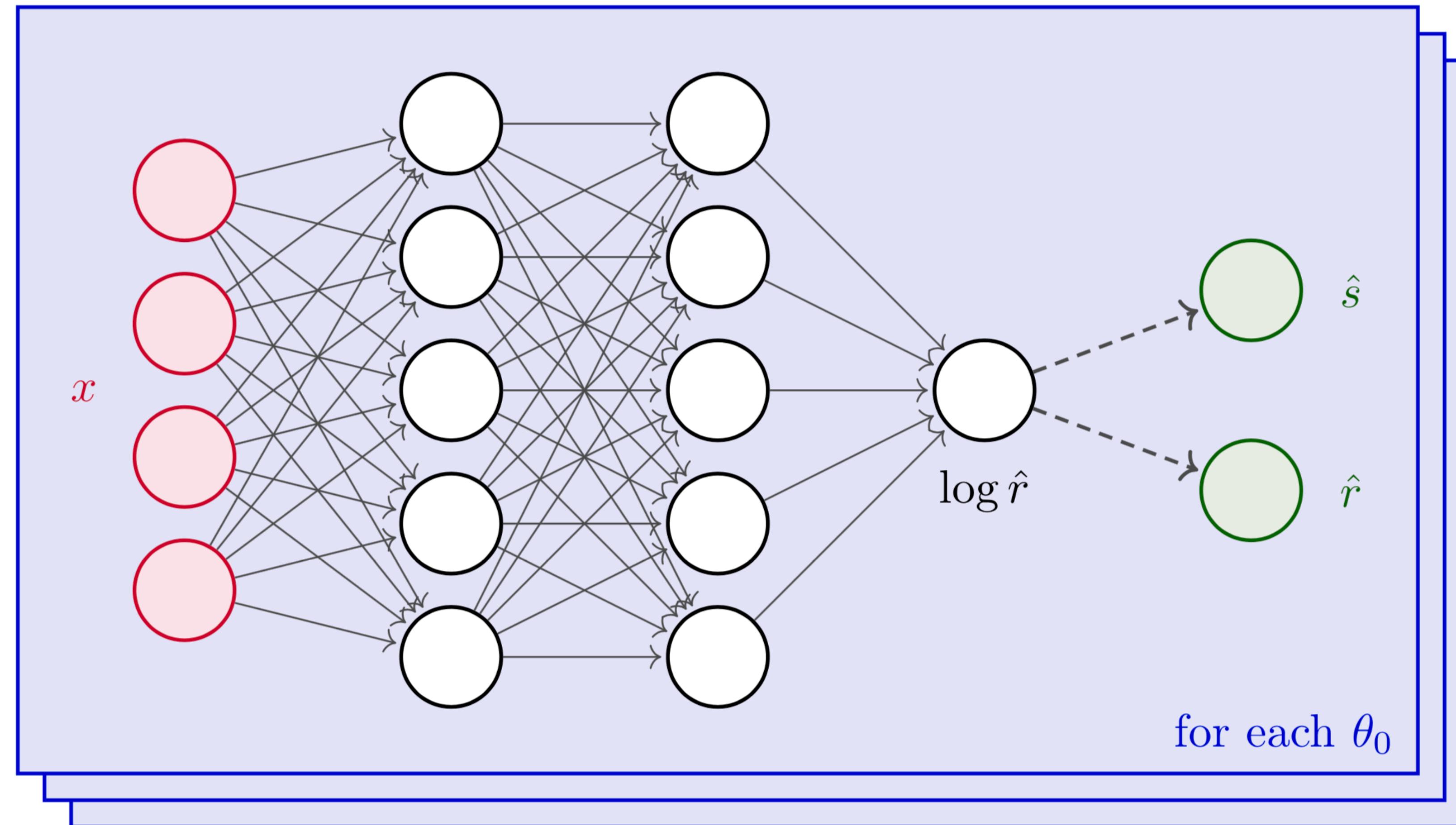
B) Parameterized:

[K. Cranmer, J. Pavez, G. Louppe 1506.02169;
P. Baldi et al. 1601.07913]

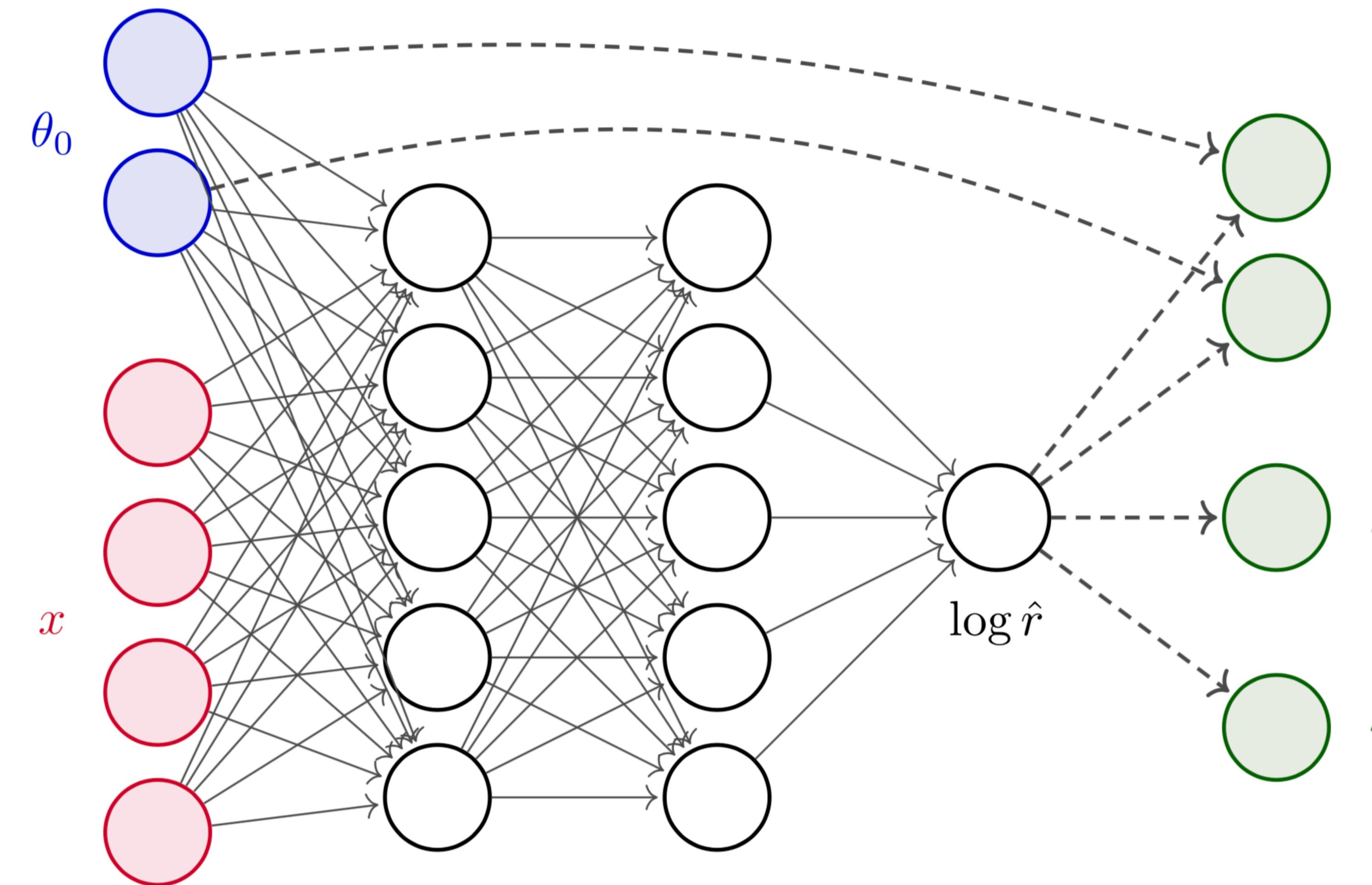
- create one estimator $\hat{r}(x|\theta_0, \theta_1)$ that is a function of θ_0 and θ_1
- no further interpolation necessary
- “borrows information” from close points



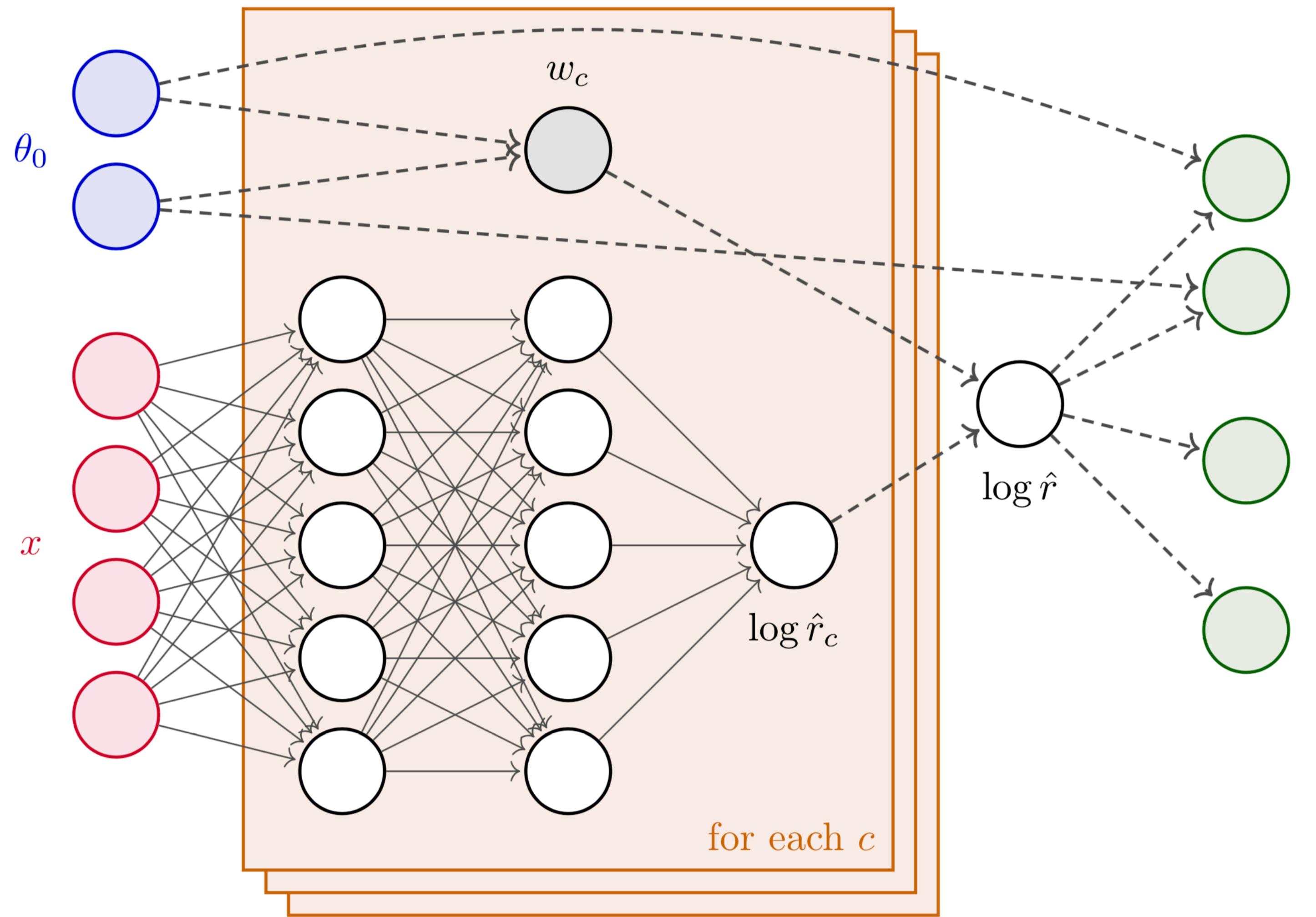
Point by point



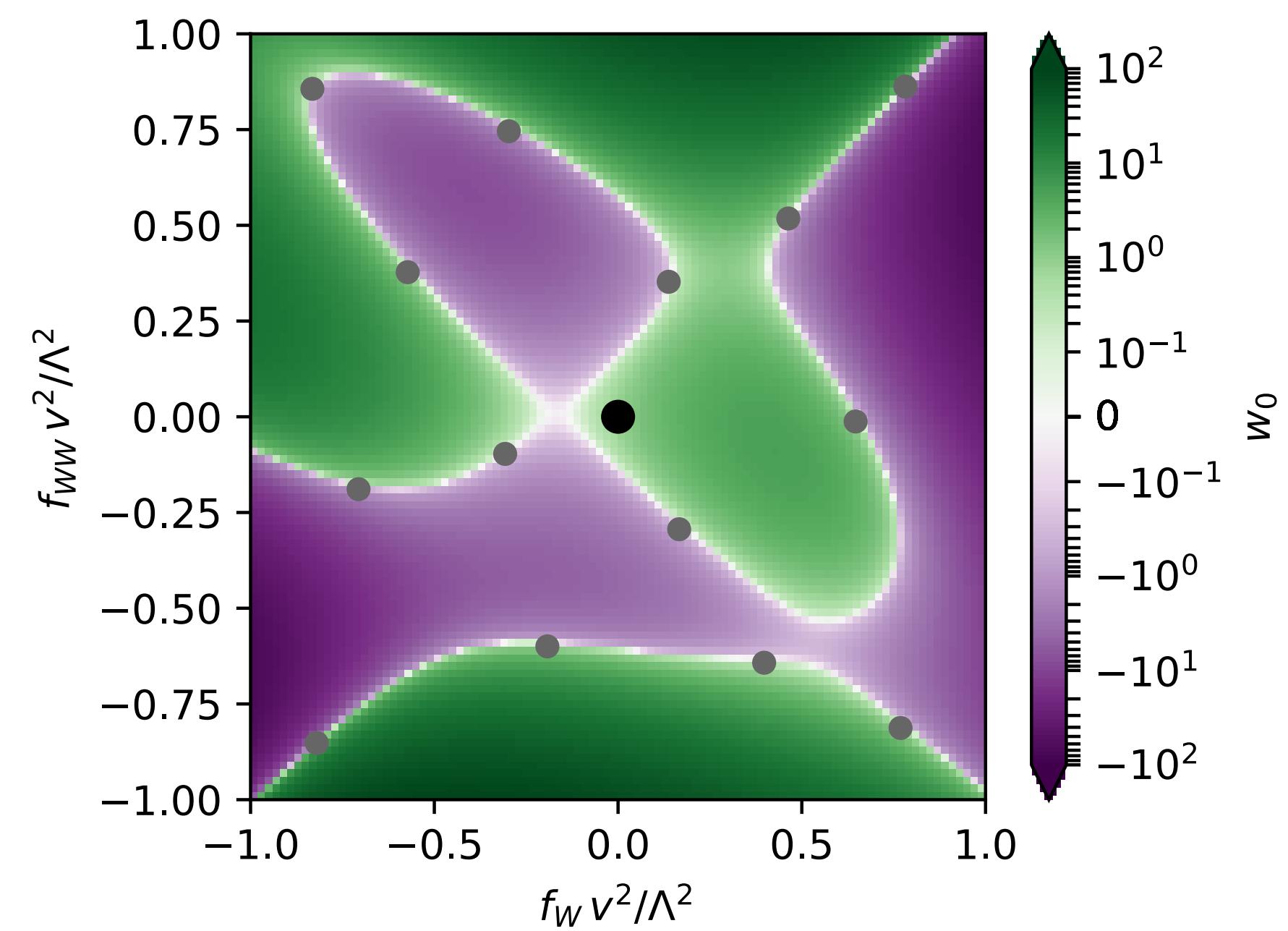
(Physics-agnostic) parameterized estimators



Morphing-aware parameterized estimators



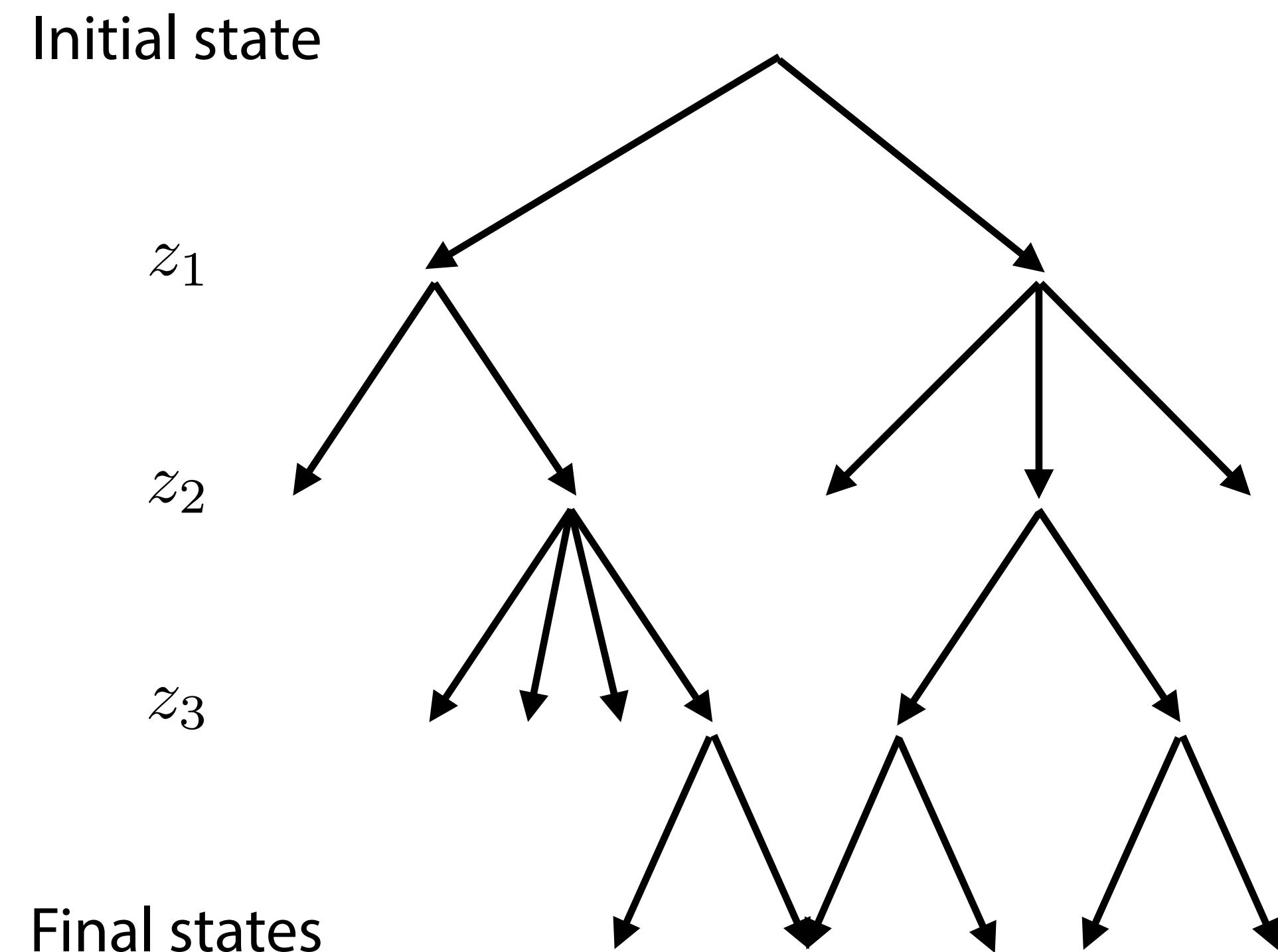
$$\hat{r}(x|\theta_0, \theta_1) = \sum_c w_c(\theta_0) \hat{r}_c(x)$$



Mining gold from any simulation

- 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()
```



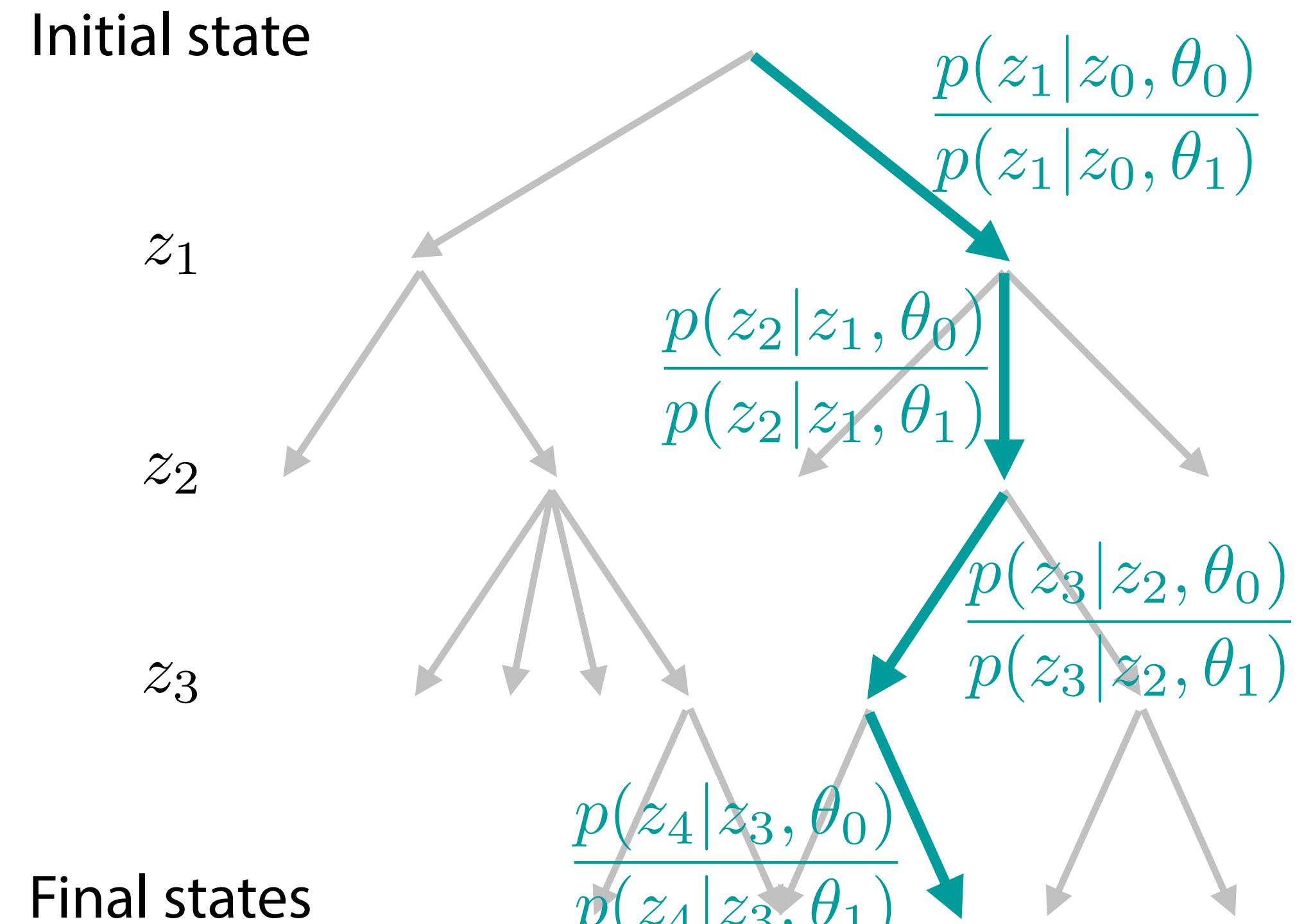
Mining gold from any simulation

- 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**”:

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



Mining gold: A family of new inference techniques

Method	Simulate	Extract		NN estimates	Asympt. exact	Generative
		$r(x, z)$	$t(x, z)$			
ROLR	$\theta_0 \sim \pi(\theta), \theta_1$	✓		$\hat{r}(x \theta_0, \theta_1)$	✓	
CASCAL	$\theta_0 \sim \pi(\theta), \theta_1$		✓	$\hat{r}(x \theta_0, \theta_1)$	✓	
ALICE	$\theta_0 \sim \pi(\theta), \theta_1$		✓	$\hat{r}(x \theta_0, \theta_1)$	✓	
RASCAL	$\theta_0 \sim \pi(\theta), \theta_1$	✓	✓	$\hat{r}(x \theta_0, \theta_1)$	✓	
ALICES	$\theta_0 \sim \pi(\theta), \theta_1$	✓	✓	$\hat{r}(x \theta_0, \theta_1)$	✓	
SCANDAL	$\theta \sim \pi(\theta)$		✓	$\hat{p}(x \theta)$	✓	✓
SALLY	θ_{ref}		✓	$\hat{t}(x \theta_{\text{ref}})$	in local approx.	
SALLINO	θ_{ref}		✓	$\hat{t}(x \theta_{\text{ref}})$	in local approx.	

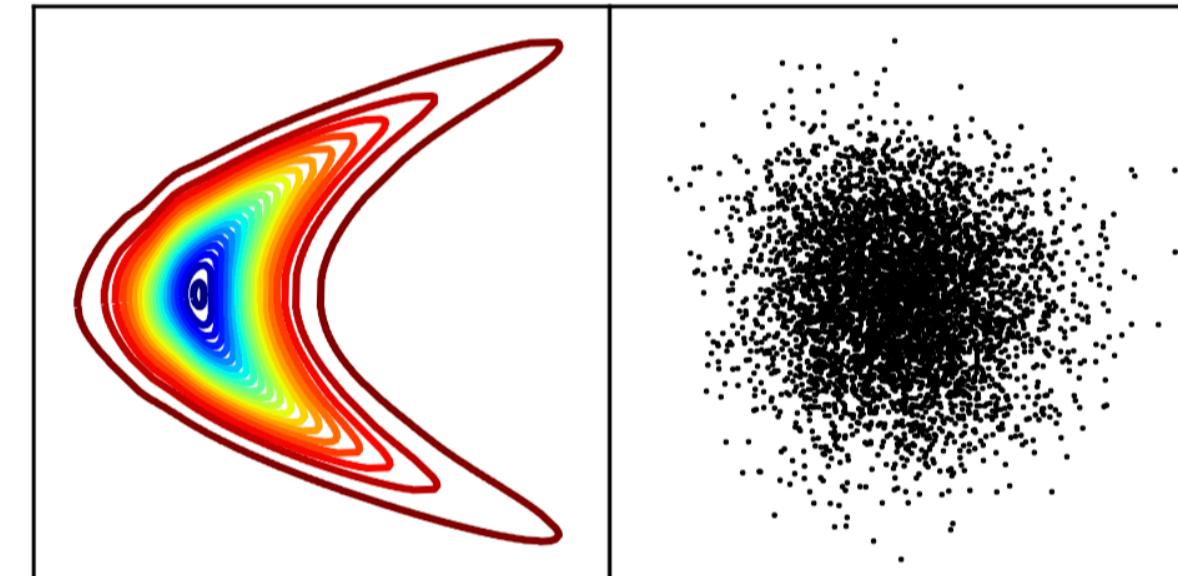
Mining gold: A family of new inference techniques

Method	Simulate	Extract $r(x, z)$ $t(x, z)$	NN estimates	Asympt. exact	Generative
ROLR	$\theta_0 \sim \pi(\theta), \theta_1$	✓	$\hat{r}(x \theta_0, \theta_1)$	✓	
CASCAL	$\theta_0 \sim \pi(\theta), \theta_1$	✓	$\hat{r}(x \theta_0, \theta_1)$	✓	
ALICE	$\theta_0 \sim \pi(\theta), \theta_1$	✓	$\hat{r}(x \theta_0, \theta_1)$	✓	
RASCAL	$\theta_0 \sim \pi(\theta), \theta_1$	✓	✓	$\hat{r}(x \theta_0, \theta_1)$	✓
ALICES	$\theta_0 \sim \pi(\theta), \theta_1$	✓	✓	$\hat{r}(x \theta_0, \theta_1)$	✓
SCANDAL	$\theta \sim \pi(\theta)$		✓	$\hat{p}(x \theta)$	✓
SALLY	θ_{ref}		✓	$\hat{t}(x \theta_{\text{ref}})$	in local approx.
SALLINO	θ_{ref}		✓	$\hat{t}(x \theta_{\text{ref}})$	in local approx.

Performance gains with
cross-entropy-based loss
[M. Stoye, JB, K. Cranmer, G.
Louppe, J. Pavez 1808.00973]

Mining gold: A family of new inference techniques

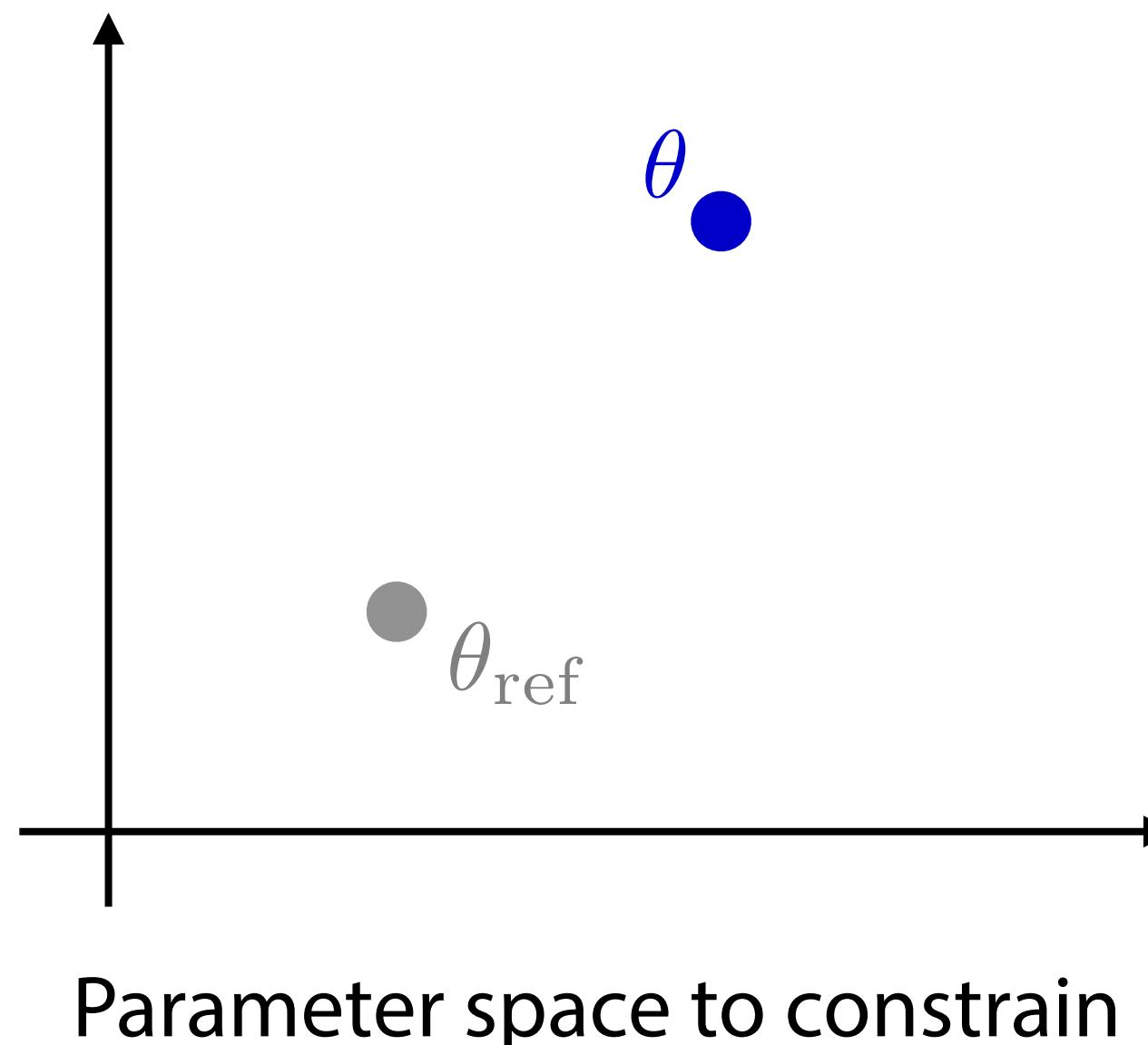
Method	Simulate	Extract $r(x, z)$ $t(x, z)$	NN estimates	Asympt. exact	Generative
ROLR	$\theta_0 \sim \pi(\theta), \theta_1$	✓	$\hat{r}(x \theta_0, \theta_1)$	✓	
CASCAL	$\theta_0 \sim \pi(\theta), \theta_1$		✓	$\hat{r}(x \theta_0, \theta_1)$	✓
ALICE	$\theta_0 \sim \pi(\theta), \theta_1$		✓	$\hat{r}(x \theta_0, \theta_1)$	✓
RASCAL	$\theta_0 \sim \pi(\theta), \theta_1$	✓	✓	$\hat{r}(x \theta_0, \theta_1)$	✓
ALICES	$\theta_0 \sim \pi(\theta), \theta_1$	✓	✓	$\hat{r}(x \theta_0, \theta_1)$	✓
SCANDAL	$\theta \sim \pi(\theta)$		✓	$\hat{p}(x \theta)$	✓
SALLY	θ_{ref}		✓	$\hat{t}(x \theta_{\text{ref}})$	in local approx.
SALLINO	θ_{ref}		✓	$\hat{t}(x \theta_{\text{ref}})$	in local approx.



Combination with state-of-the-art conditional neural density estimators, e.g. normalizing flows

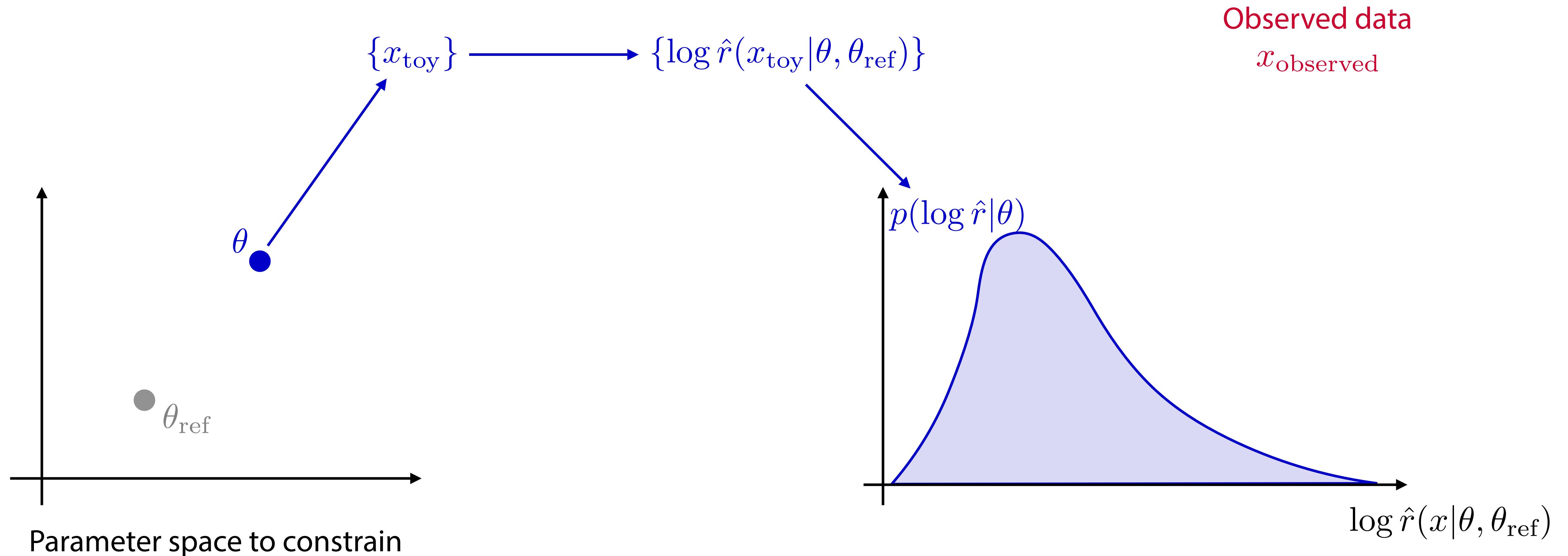
[everything by G. Papamakarios: G. Papamakarios, T. Pavlakou, I. Murray 1705.07057; G. Papamakarios, D. Sterratt, I. Murray 1805.07226; ...]

Limit setting (frequentist, standard ATLAS / CMS practice)

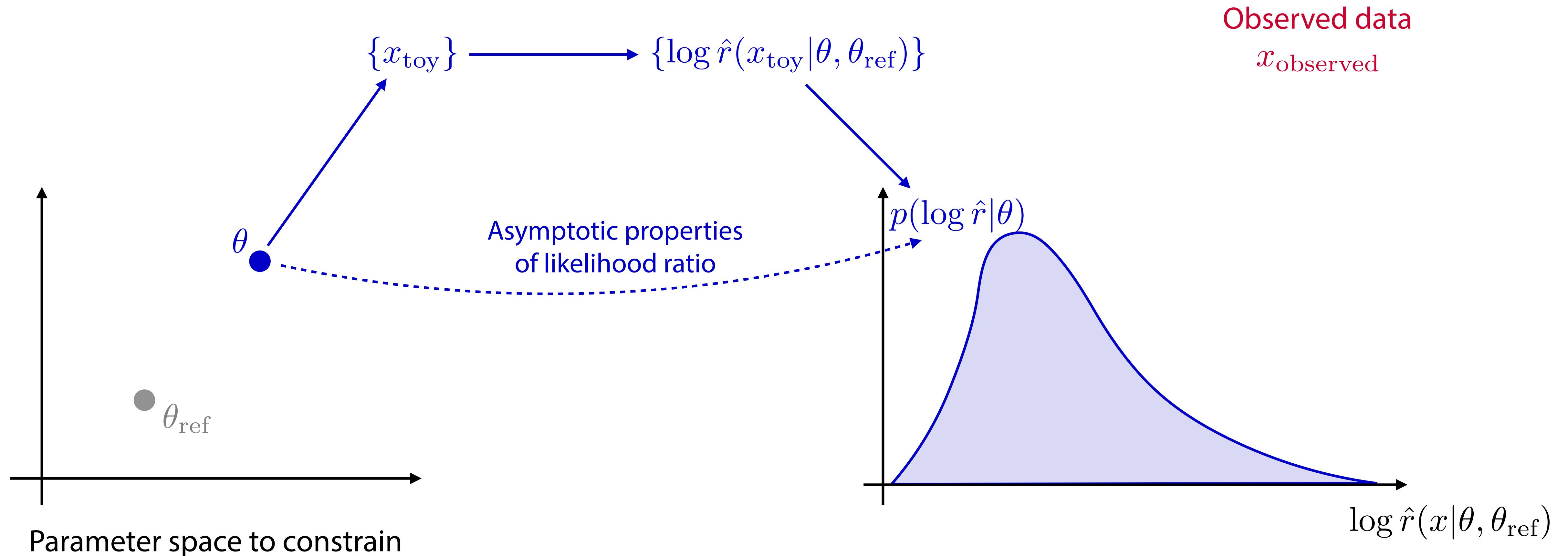


Observed data
 x_{observed}

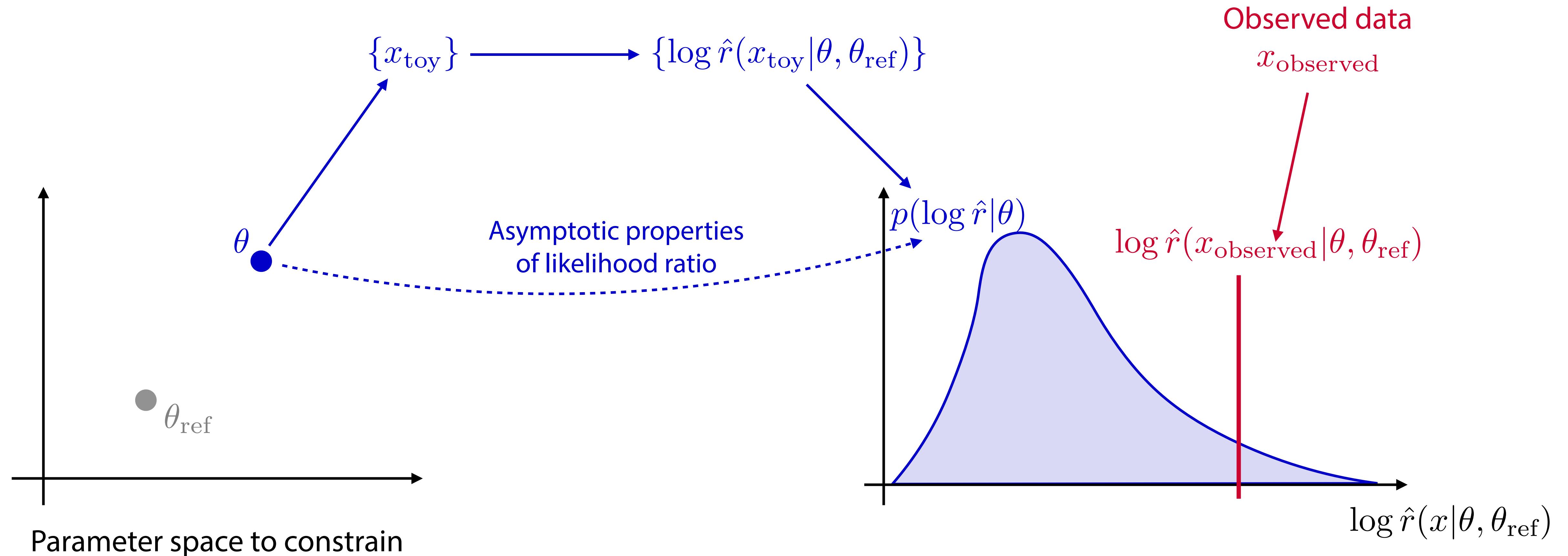
Limit setting (frequentist, standard ATLAS / CMS practice)



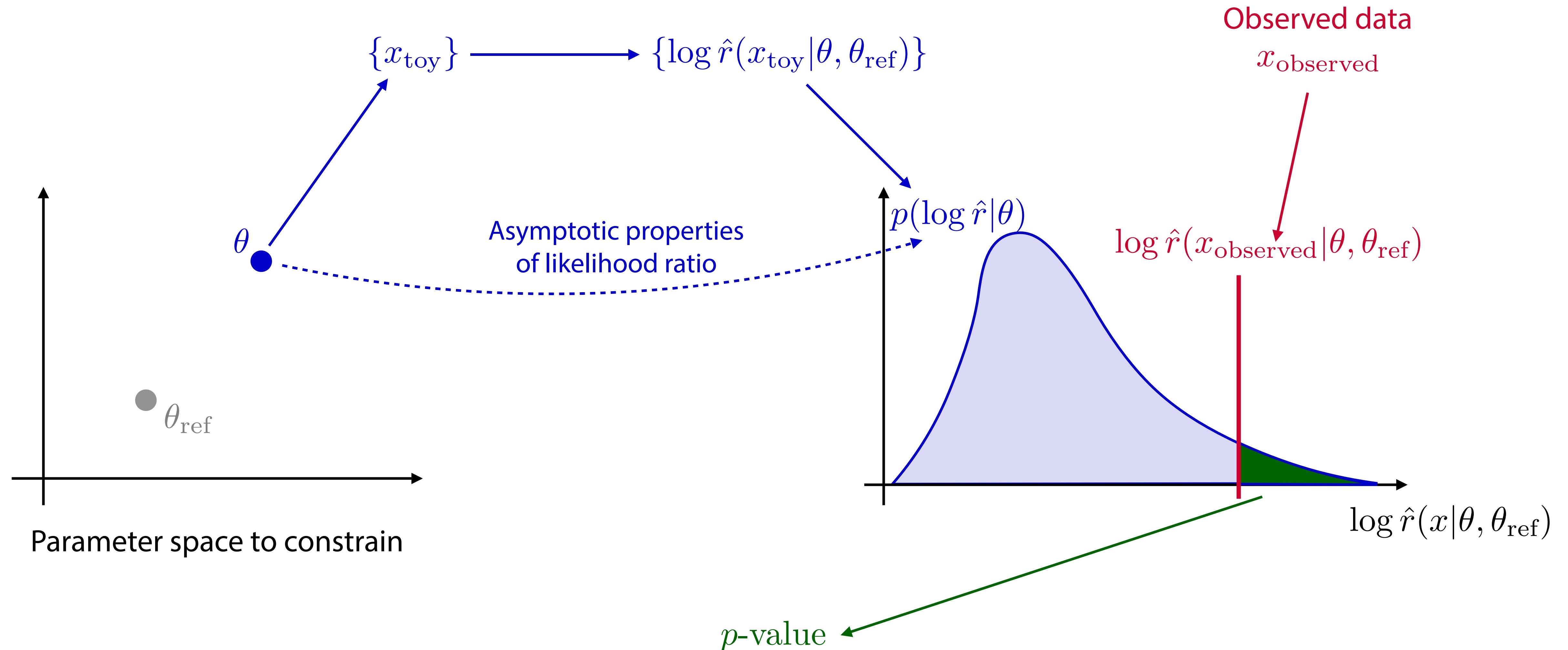
Limit setting (frequentist, standard ATLAS / CMS practice)



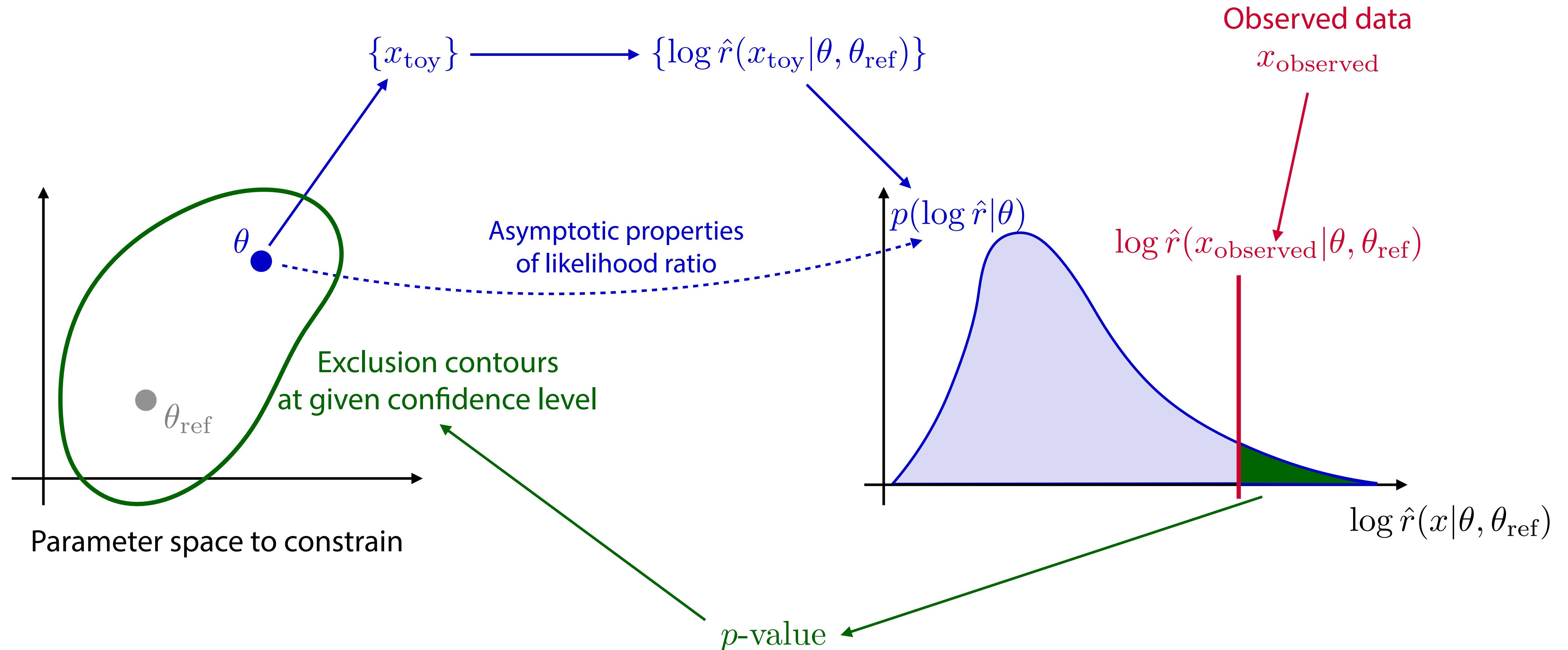
Limit setting (frequentist, standard ATLAS / CMS practice)



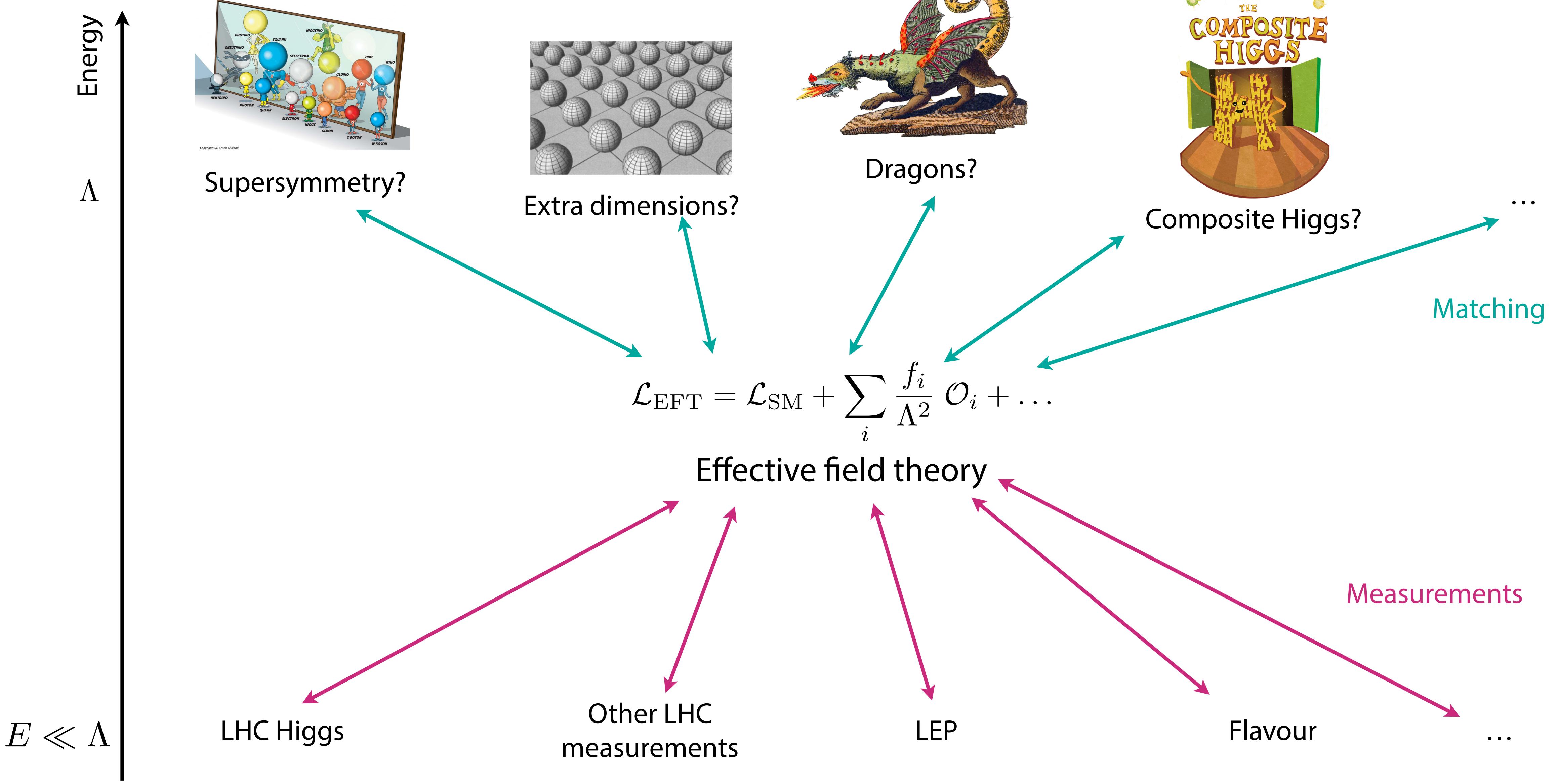
Limit setting (frequentist, standard ATLAS / CMS practice)



Limit setting (frequentist, standard ATLAS / CMS practice)

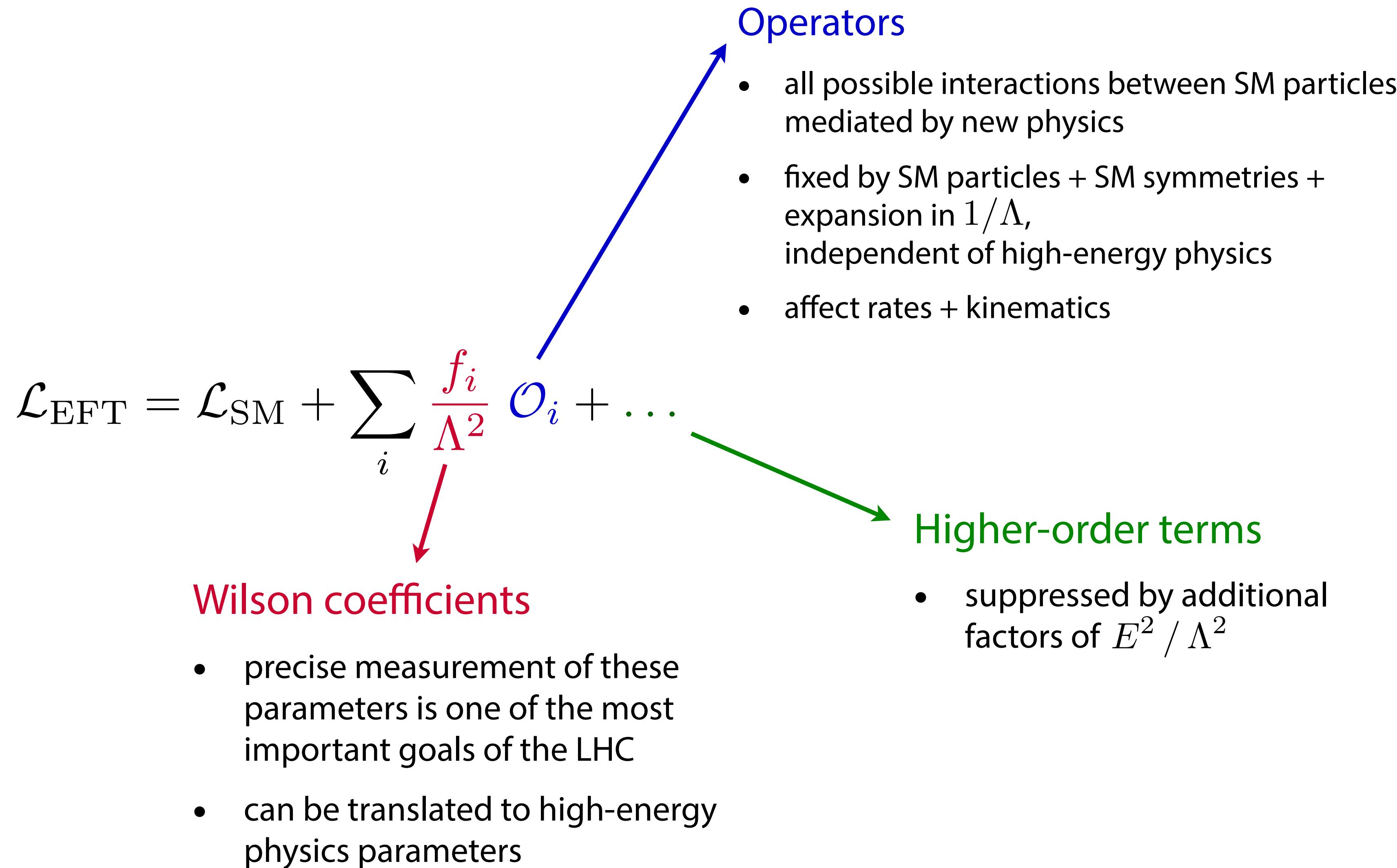


Effective field theory



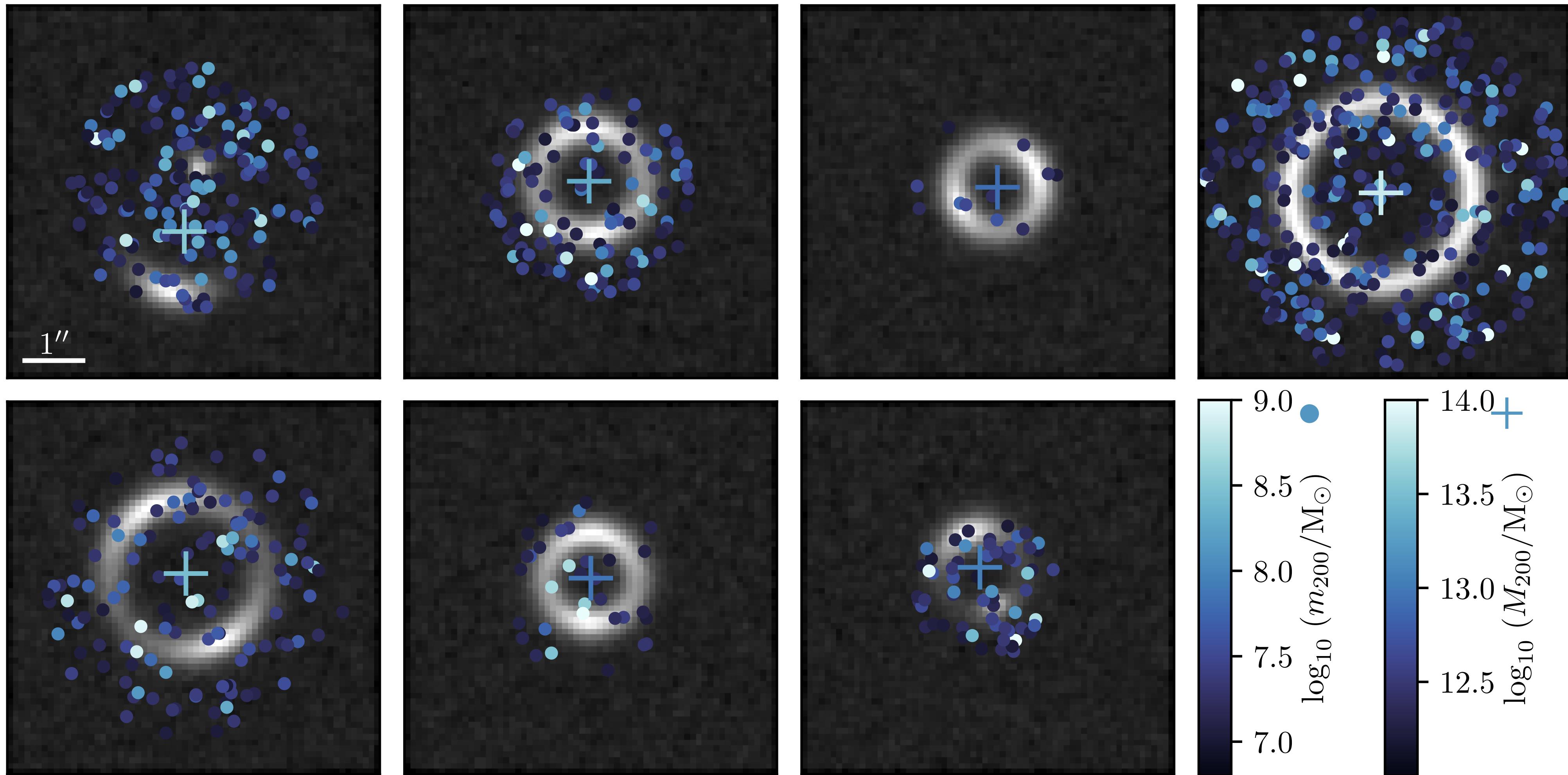
SMEFT (Standard Model Effective Field Theory)

[W. Buchmuller, D. Wyler 85;
B. Grzadkowski, M. Iskrzynski, M. Misiak, J. Rosiek 1008.4884; ...]



$\mathcal{O}_{\phi,1} = (D_\mu \phi)^\dagger \phi \phi^\dagger D^\mu \phi$	$\mathcal{O}_{GG} = (\phi^\dagger \phi) G_{\mu\nu}^a G^{\mu\nu a}$
$\mathcal{O}_{\phi,2} = \frac{1}{2} \partial_\mu (\phi^\dagger \phi) \partial^\mu (\phi^\dagger \phi)$	$\mathcal{O}_{BB} = -\frac{g'^2}{4} (\phi^\dagger \phi) B_{\mu\nu} B^{\mu\nu}$
$\mathcal{O}_{\phi,3} = \frac{1}{3} (\phi^\dagger \phi)^3$	$\mathcal{O}_{WW} = -\frac{g^2}{4} (\phi^\dagger \phi) W_{\mu\nu}^a W^{\mu\nu a}$
$\mathcal{O}_{\phi,4} = (\phi^\dagger \phi) (D_\mu \phi)^\dagger D^\mu \phi$	$\mathcal{O}_{BW} = -\frac{g g'}{4} (\phi^\dagger \sigma^a \phi) B_{\mu\nu} W^{\mu\nu a}$
	$\mathcal{O}_B = \frac{ig'}{2} (D^\mu \phi)^\dagger D^\nu \phi B_{\mu\nu}$
	$\mathcal{O}_W = \frac{ig}{2} (D^\mu \phi)^\dagger \sigma^a D^\nu \phi W_{\mu\nu}^a$

Strong lensing: simulated images



Strong lensing: expected posterior

