

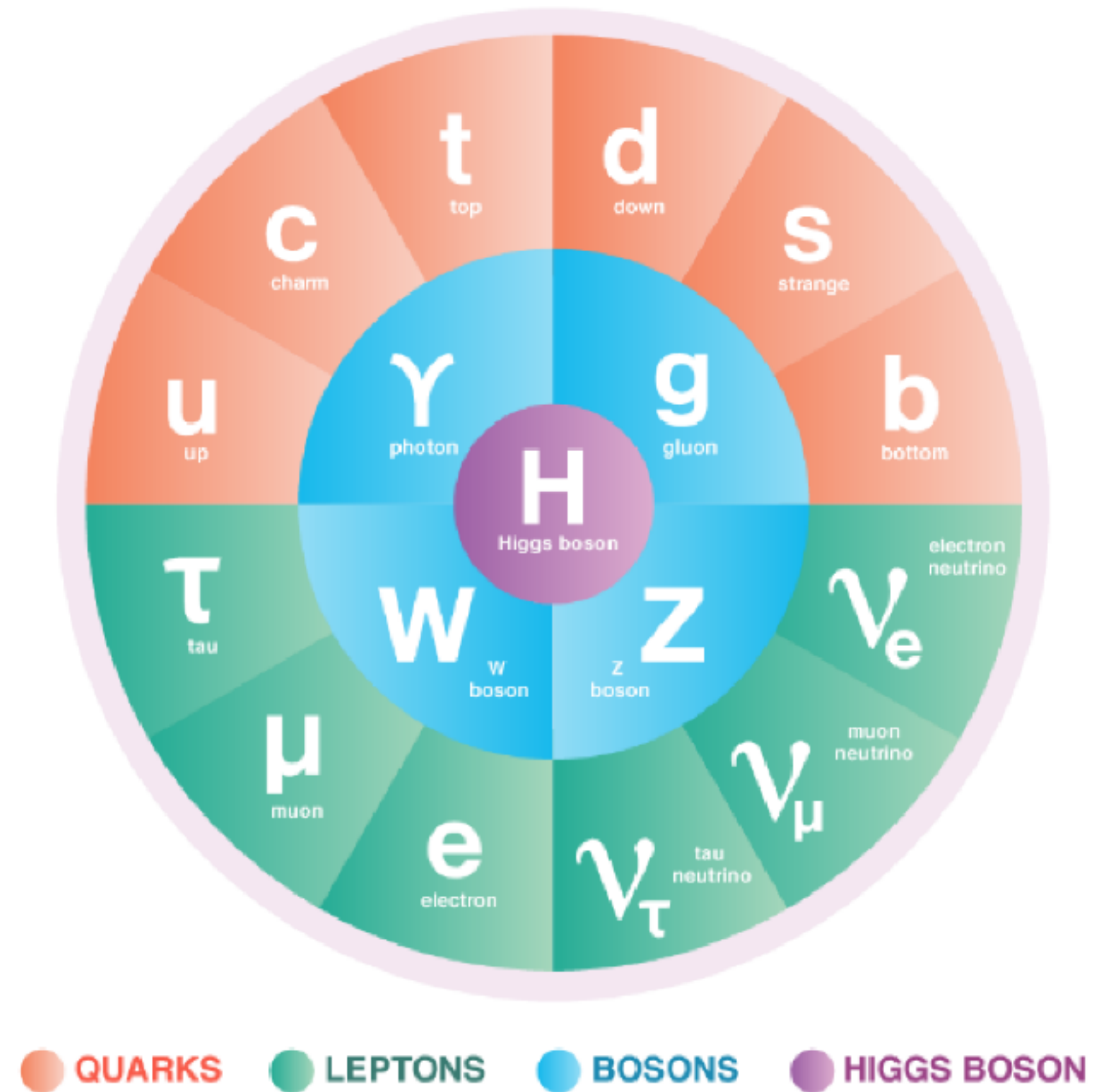
Aspects of Likelihood Free / Simulation-based Inference for Particle Physics

Lukas Heinrich, CMU STAMPS 2026



MDSI **TUM**

My Physics Context



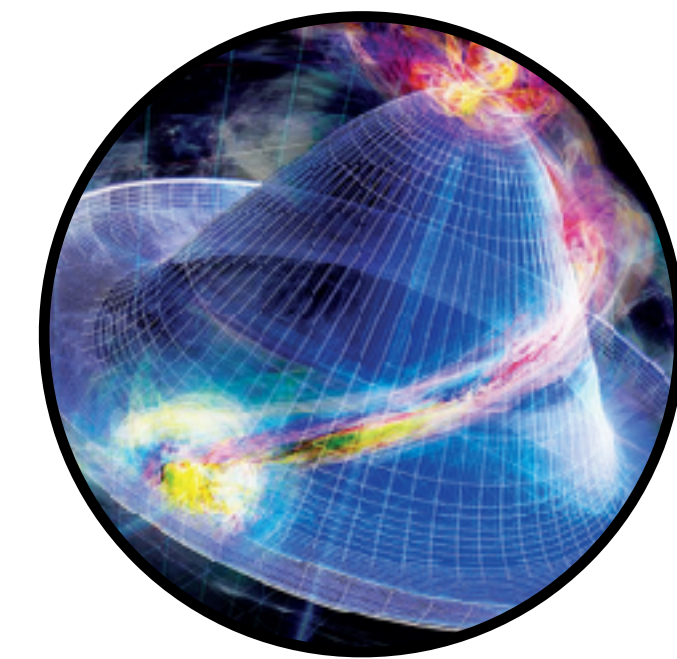
Dark Matter



Nature of Neutrinos



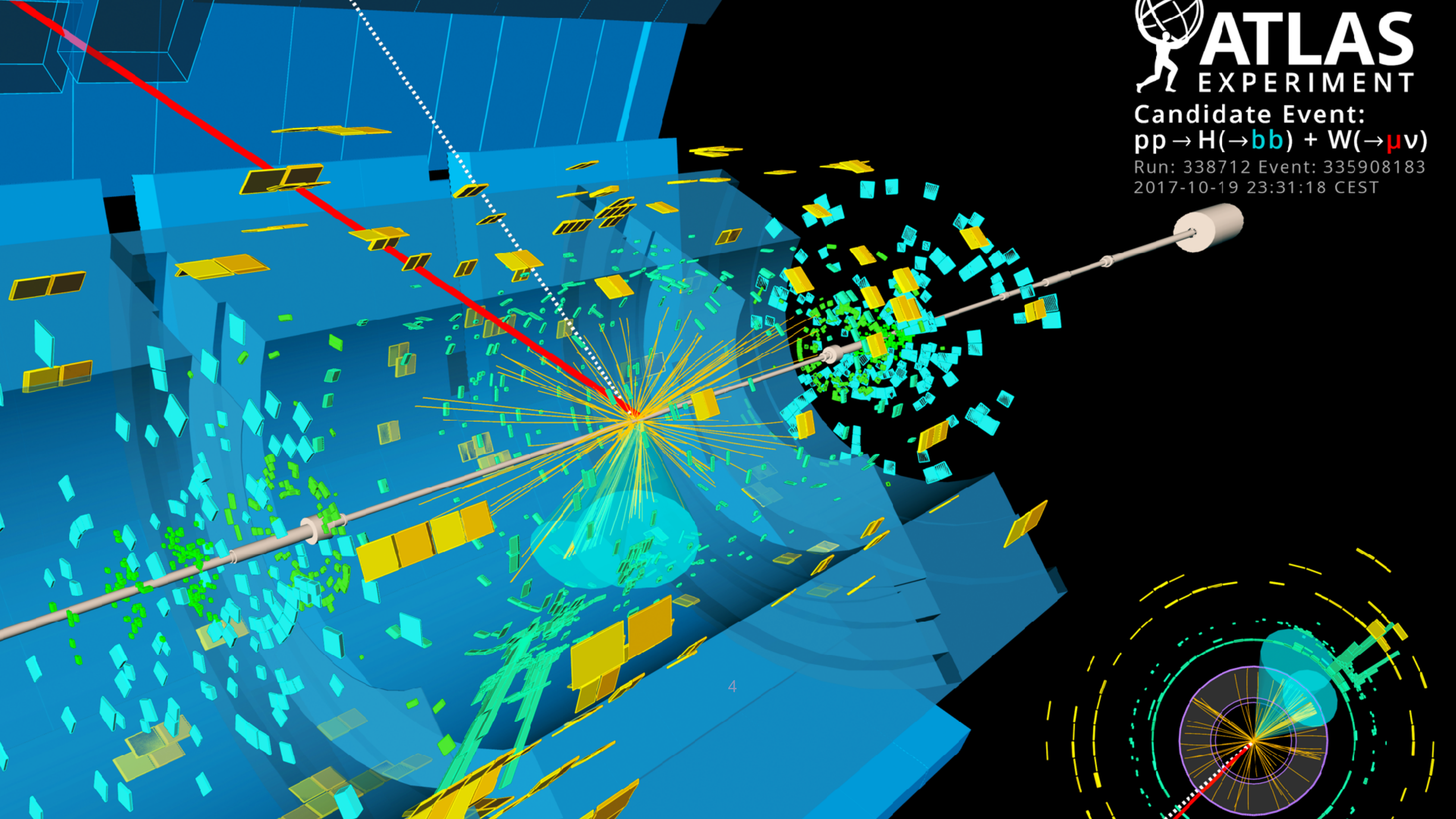
Matter & Antimatter



Origin of Mass

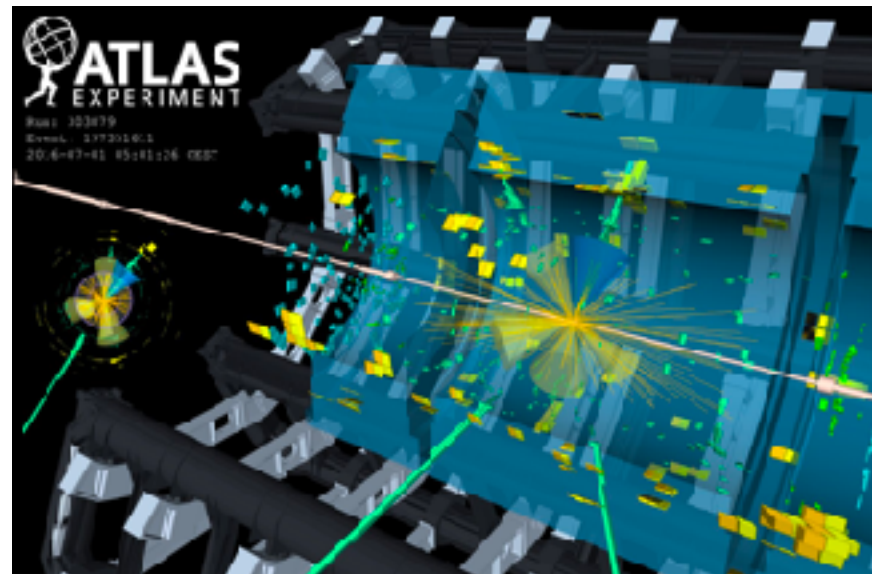
Candidate Event:
 $pp \rightarrow H(\rightarrow bb) + W(\rightarrow \mu\nu)$

Run: 338712 Event: 335908183
2017-10-19 23:31:18 CEST

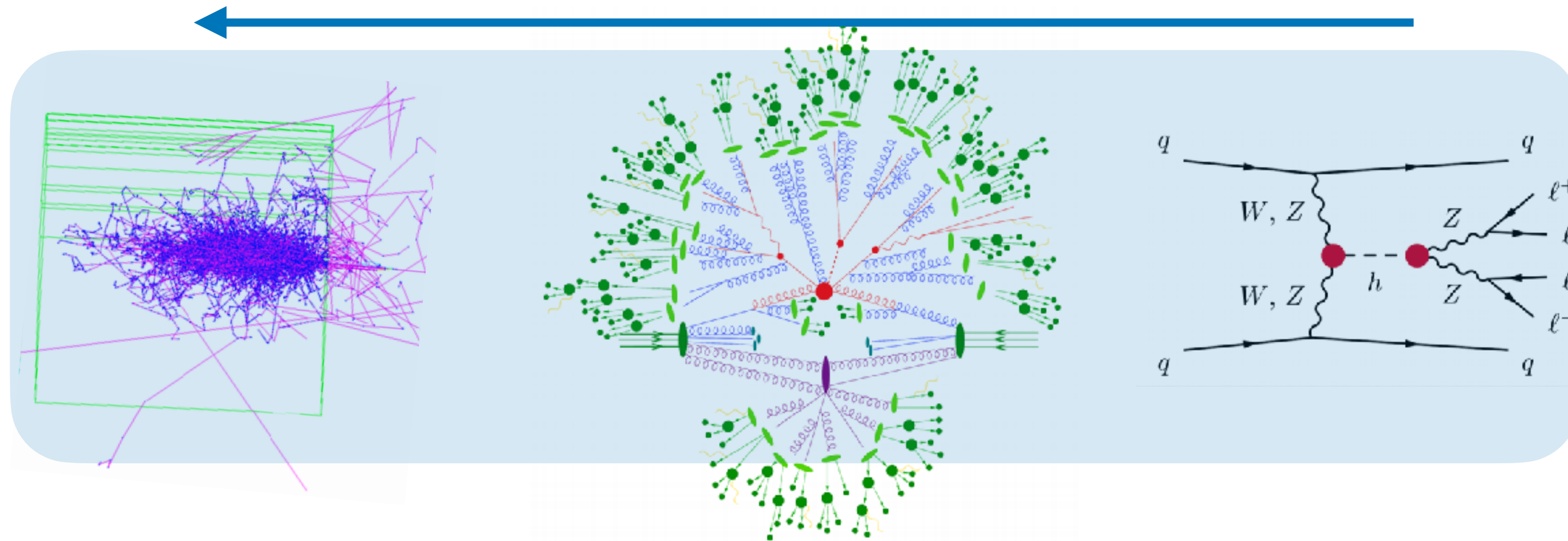


HEP is Likelihood-Free

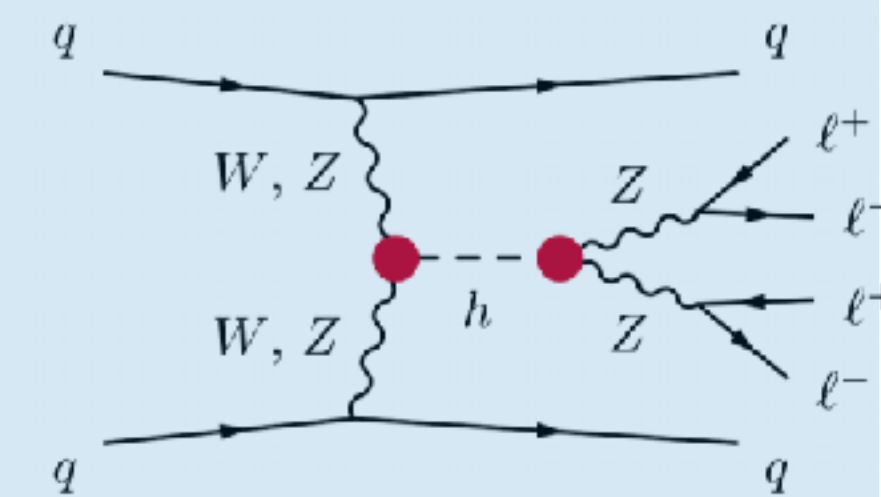
The data-generating process in particle physics is defined by a long sequence of complicated stochastic transitions



x



θ

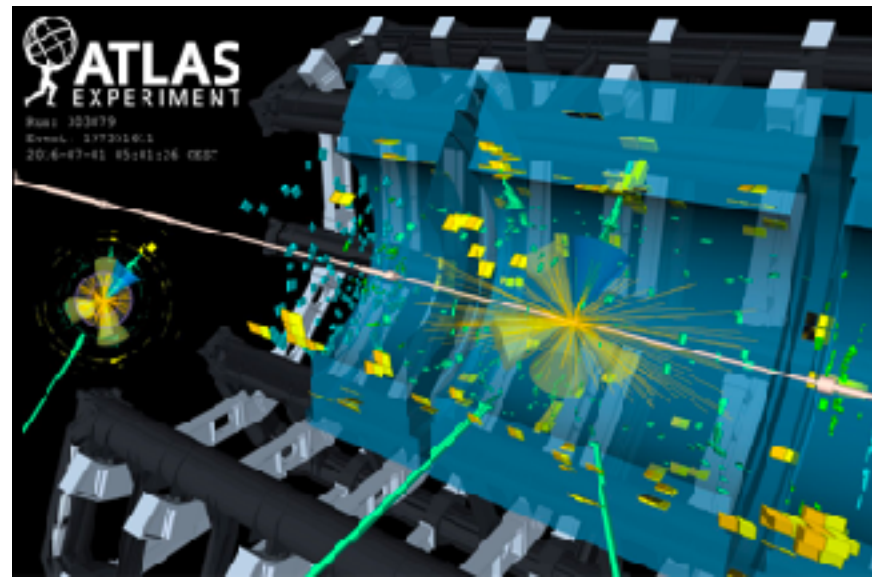


$$\mathcal{L} = -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} + i\bar{\psi}\not{D}\psi + h.c. + \bar{\psi}_i y_{ij} \psi_j \phi + h.c. + |\mathcal{D}_\mu \phi|^2 - V(\phi)$$

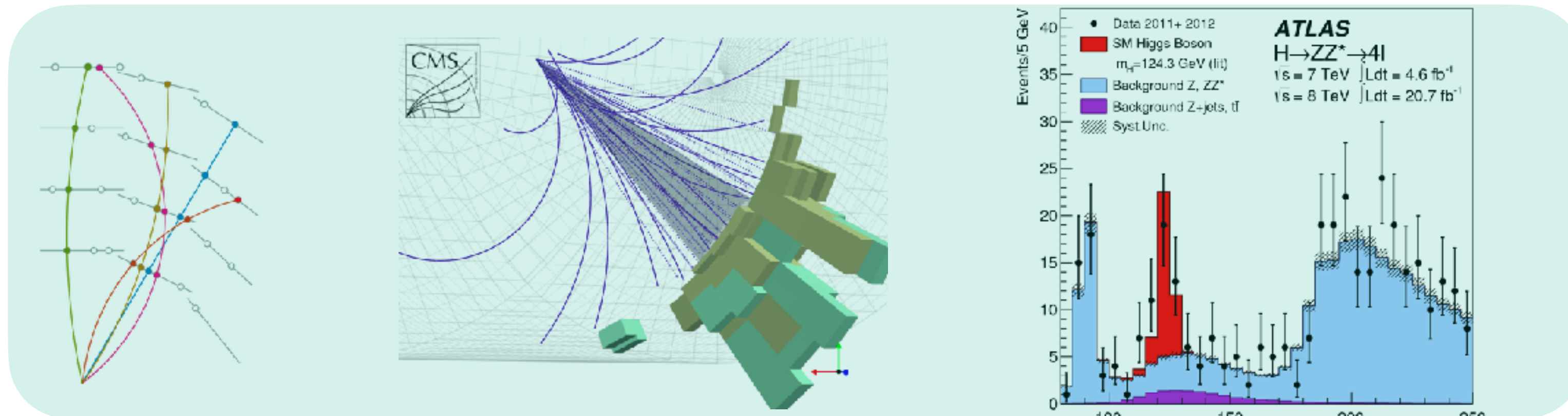
$$p(x | \theta) = \int dz \ p(x | z_h) \ p(z_h | z_p) \ p(z_p | \theta)$$

HEP is Likelihood-Free

Traditionally, inference is done by a series of per-event manually derived summary statistics



x



$$\mathcal{L} = -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} + i\bar{\psi}\not{D}\psi + h.c. + \bar{\psi}_i y_{ij} \psi_j \phi + h.c. + \frac{1}{2} \partial_\mu \phi^\dagger \partial^\mu \phi - V(\phi)$$

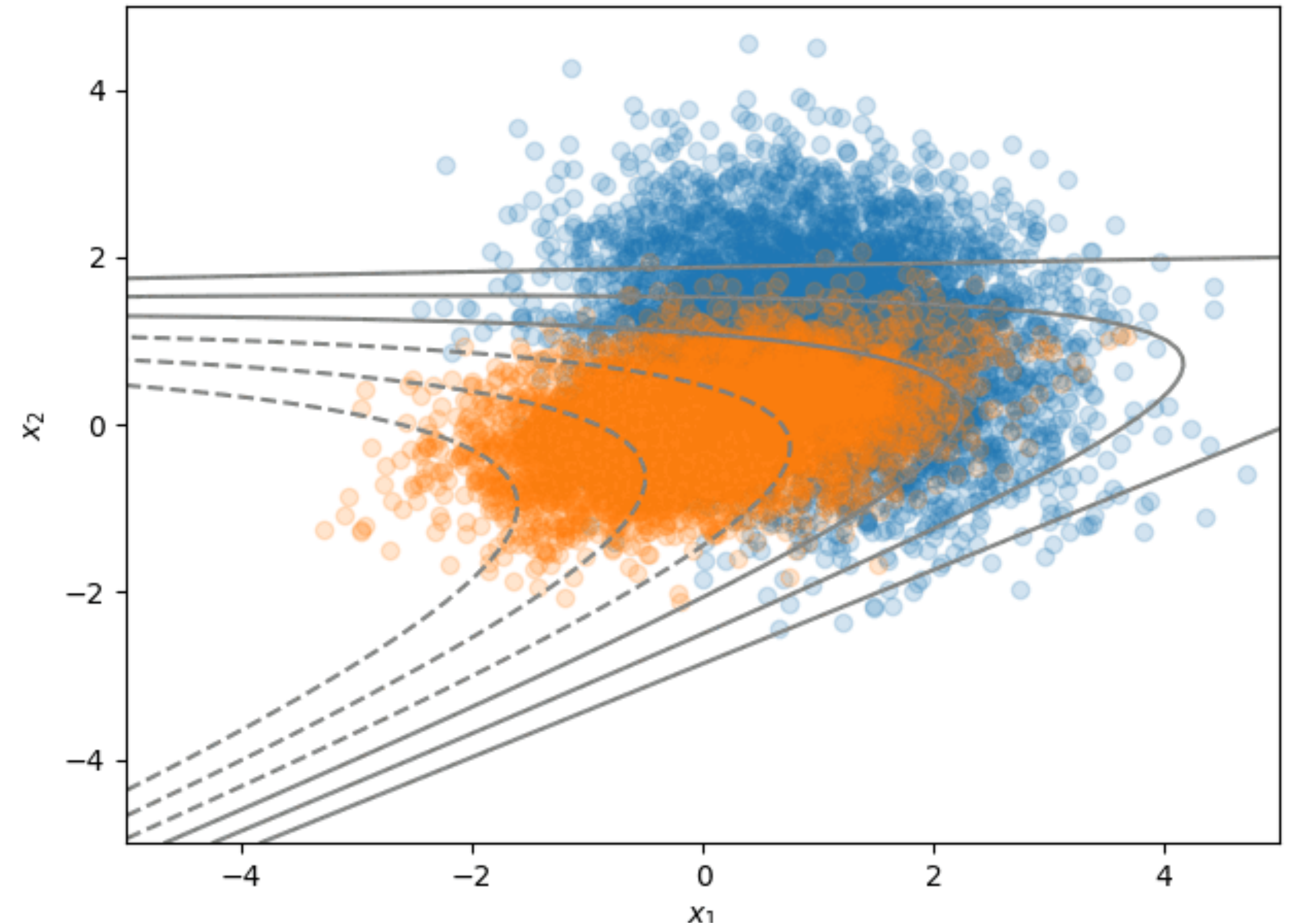
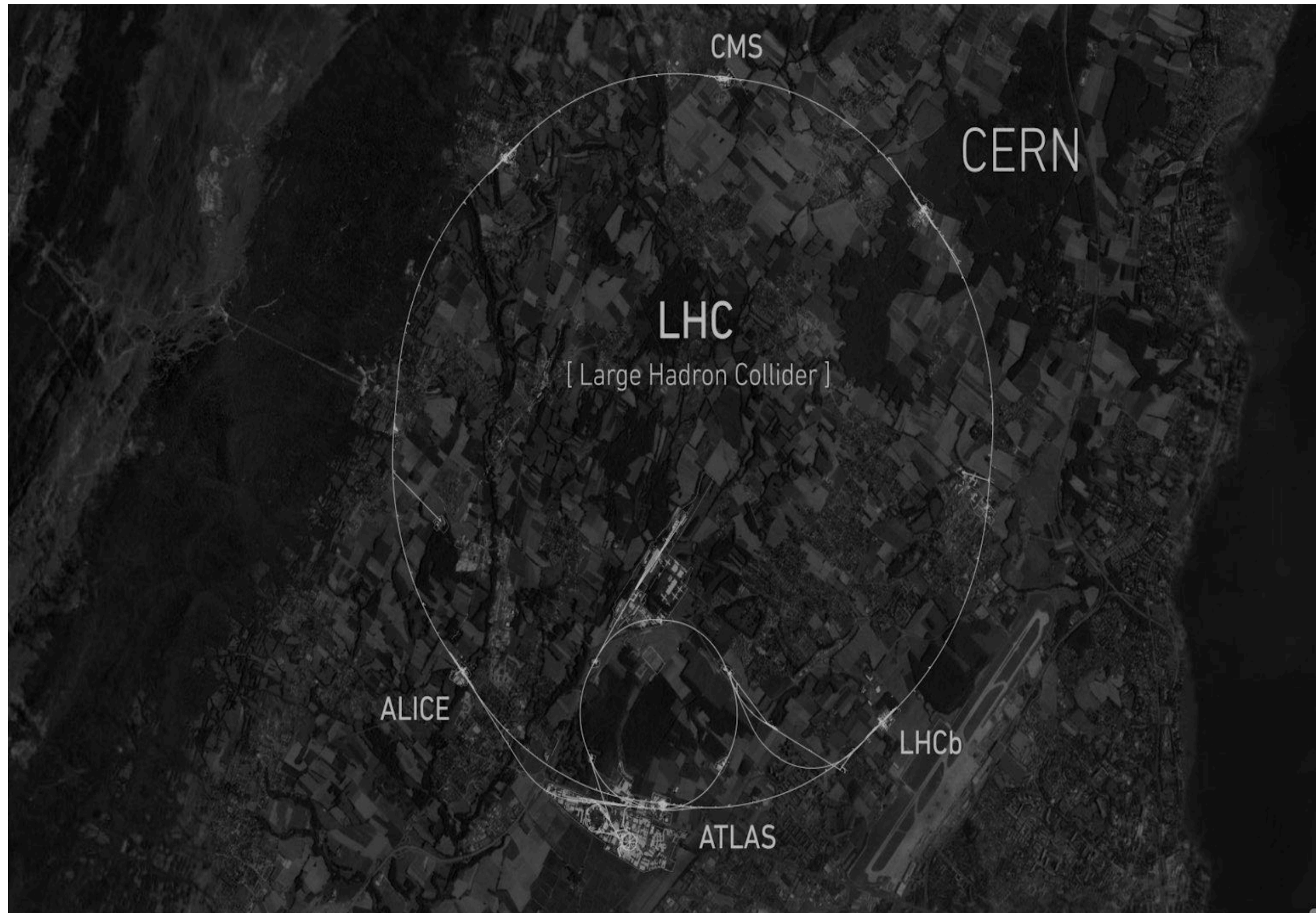
θ

approximate inference: $p(\theta | x) \approx p(\theta | \lambda_\phi(x))$

what's the best way to compress the data?

ML-driven SBI

More recently, ML techniques have changed how we think about this simulation-based inference workflow



Likelihood-Ratio Trick

Three Types of High Dimension

In particle physics we want to do inference on the full set of particle events we record (potentially billions!)

$$p(\theta | \mathcal{D}) \quad \frac{p(\mathcal{D} | \theta)}{p(\mathcal{D} | \hat{\theta})} \quad \mathcal{D} = \{x_1, x_2, x_3 \dots x_n\}$$

Bayes *Frequentist*

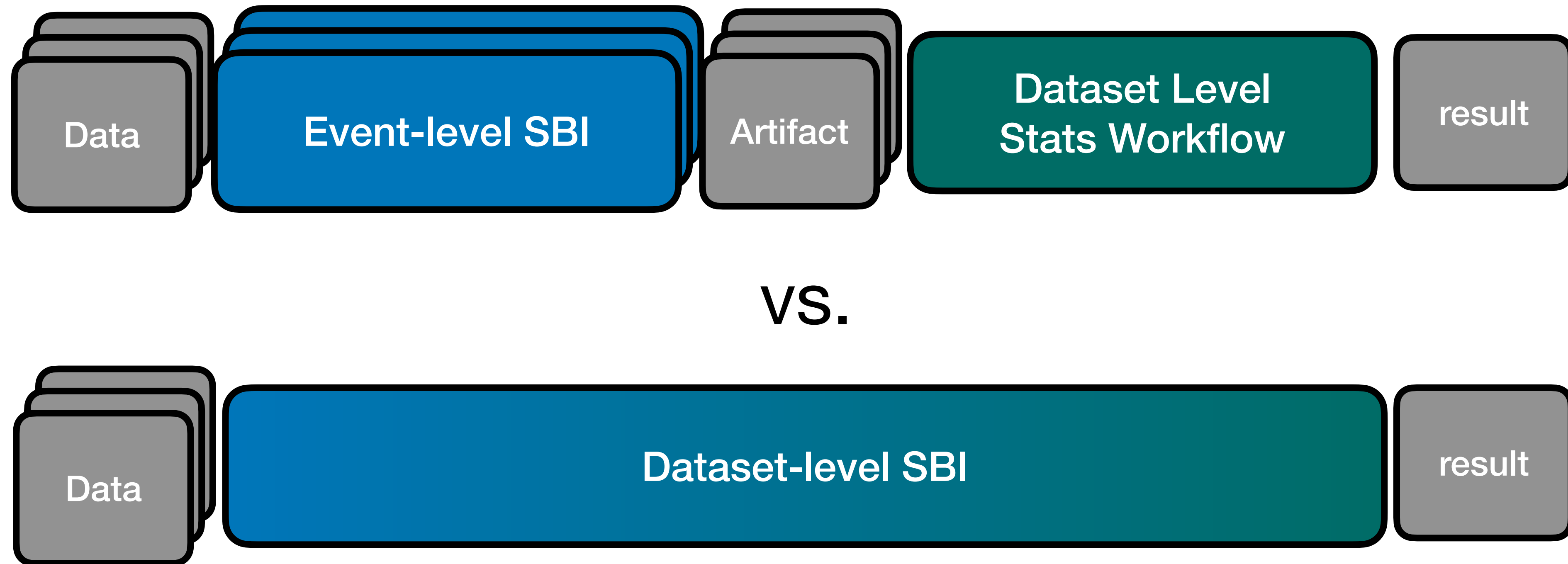
where we assume the data emerges from an i.i.d process

$$x_i \in \mathcal{D} \quad x_i \sim p(x | \theta)$$

*potentially set cardinality depends on θ
(in HEP-speak “extended likelihood”)*

A stymied situation

The SBI workflow in HEP was therefore never really fully “end-to-end” but always had a “Traditional Stats” component



Example

1) Train an event-level L'hood-ratio estimator $r_\phi(x, \theta) = \frac{p(x | \theta)}{p(x | \theta_0)}$

2) Full-dataset likelihood ratio

Important:
*exploiting the fact that
Lhood ratios compose*

$$R(\theta) = \frac{p(\mathcal{D} | \theta)}{p(\mathcal{D} | \theta_0)} = \prod_i \frac{p(x_i | \theta)}{p(x_i | \theta_0)}$$

3) Run inference (MCMC, MLE, ...) on the dataset-level $R(\theta)$

Can we get to a fully end-to-end setup?

Another Problem

A lot of inference targets the subset of the full parameter space: parameters of interest vs nuisance parameters

$$p(\mu | \mathcal{D}) = \int_{\theta \neq \mu} p(\theta | \mathcal{D})$$

Bayes

$$\frac{p(\mathcal{D} | \mu, \hat{\nu})}{p(\mathcal{D} | \hat{\mu}, \hat{\nu})} = \frac{\sup_{\theta \neq \mu} p(\dots)}{\sup_{\theta} p(\dots)}$$

Frequentist

But to compose full-dataset likelihood-(ratio) we need to train SBI methods on the full parameter space → **expensive!**

profile likelihoods / marginal posteriors do not compose

Dataset-level SBI

With data-set level SBI you can target the **low-dimensional** inference result directly & quickly without intermediate steps

[stat.ME] 24 Mar 2022

Learning Optimal Test Statistics in the Presence of Nuisance Parameters

Lukas Heinrich
Technical University of Munich
E-mail: lukas.heinrich@cern.ch

Abstract. The design of optimal number of scenarios optimal test s turning this argument around we ca where only samples from a simulat scenarios. We propose a likelihood-equivalent to the profile likelihood

1. Introduction

2203.13079

Published in Transactions on Machine Learning Research (02/2024)

Hierarchical Neural Simulation-Based Inference Over Event Ensembles

Lukas Heinrich
Technical University of Munich

Siddharth Mishra-Sharma
MIT, Harvard University, IAIFI

Chris Pollard
University of Warwick

Philipp Windischhofer
University of Chicago

Reviewed on OpenReview: <https://openreview.net/>

2306.12584

Abstr

When analyzing real-world data it is common to sets of observations that collectively constrain interest. Such models often have a hierarchical individual events and "global" parameters influer

[stat.ML] 21 Feb 2024

[stat.ME] 8 May 2026

**It Just Takes Two:
Scaling Amortized Inference to Large Sets**

Antoine Wehenkel*
Apple
awehenkel@apple.com

Michael Kagan
SLAC National Accelerator Laboratory
makagan@slac.stanford.edu

Lukas Heinrich
TU München
Lukas.Heinrich@cern.ch

Chris Pollard*
University of Warwick
christopher.pollard@warwick.ac.uk

2605.07972

Abstract

Neural posterior estimation has emerged as a powerful tool for amortized inference, with growing adoption across scientific and applied domains. In many of these applications the conditioning variable is a set of observations whose elements

Example: Learning to Profile

The likelihood ratio trick is a simple argument

The likelihood ratio is the optimal statistic (a function) for hypothesis testing, so function optimization (ML) for discrimination power must converge to it

The low-dimensional equivalent is the profile likelihood ratio

If we can find an argument why / how the profile likelihood statistic (a function) is optimal, we can find it with machine learning!

Example: Learning to Profile

Wald (1943) showed that there is a sense in which the profile likelihood is optimal (and thus findable with optimization)

TESTS OF STATISTICAL HYPOTHESES CONCERNING SEVERAL PARAMETERS WHEN THE NUMBER OF OBSERVATIONS IS LARGE⁽¹⁾

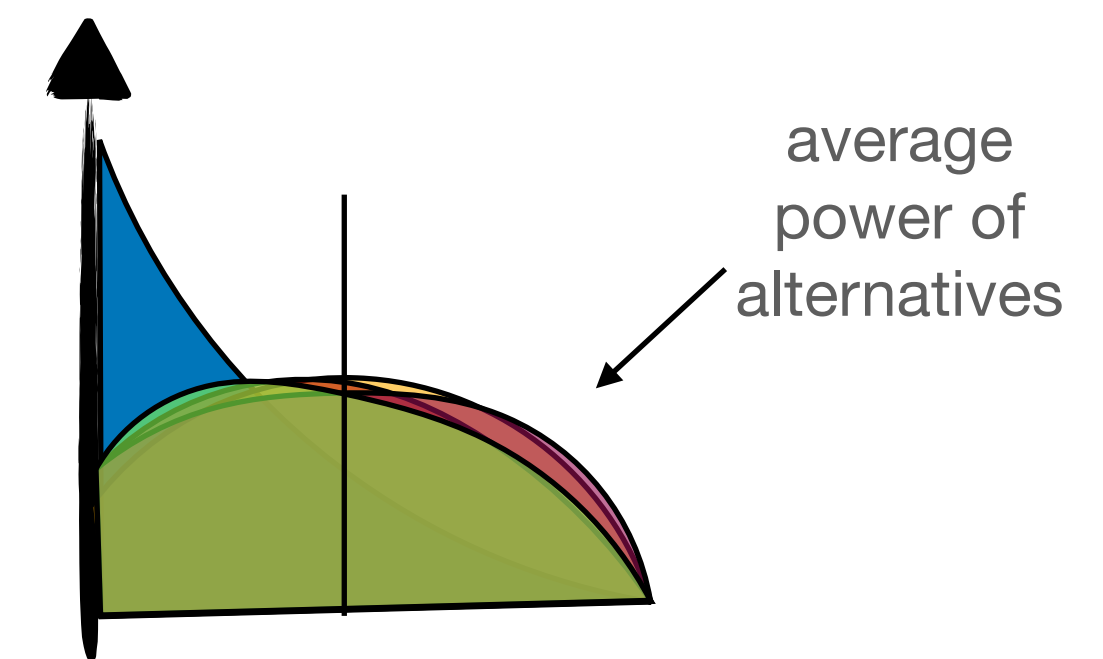
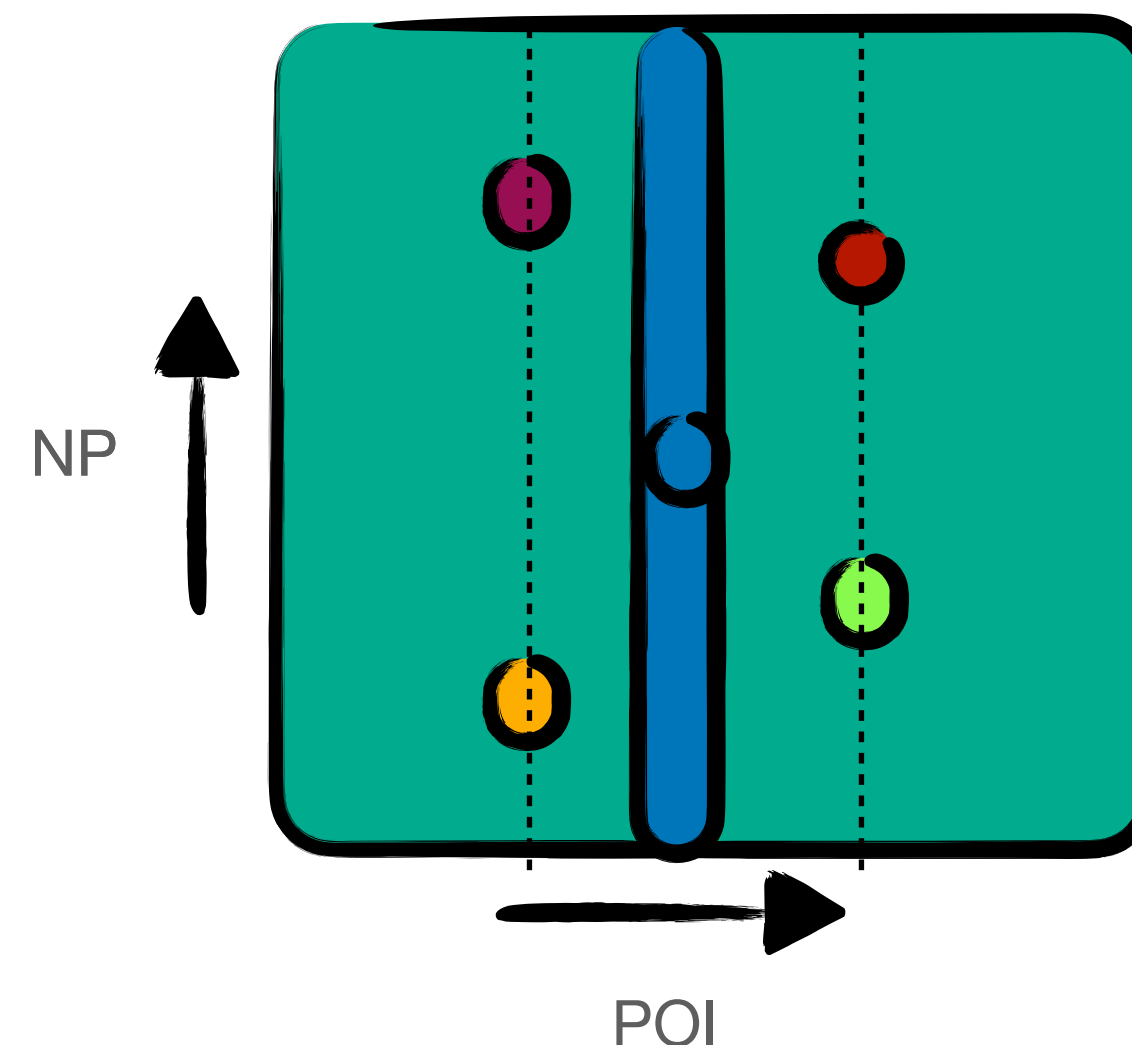
BY
ABRAHAM WALD

TABLE OF CONTENTS

1. Introduction.....	426
2. Assumptions on the density function $f(x, \theta)$	428
3. The joint limit distribution of $\hat{\theta}_n$	429
4. Reduction of the general problem to the case of a multivariate normal distribution..	433
5. Tests of simple hypotheses which have uniformly best average power over a family of surfaces.....	445
6. Tests of simple hypotheses which have best constant power on a family of surfaces...	450
7. Most stringent tests of simple hypotheses.....	451
8. Definitions of "best" tests of composite hypotheses.....	453
9. Tests of linear composite hypotheses which have uniformly best average power over a family of surfaces.....	455
10. Tests of linear composite hypotheses which have best constant power on a family of surfaces.....	461
11. Most stringent tests of linear composite hypotheses.....	461
12. The general composite hypothesis.....	463
13. Optimum properties of the likelihood ratio test.....	470
14. Large sample distribution of the likelihood ratio.....	478
15. Summary.....	481

1. **Introduction.** In this paper we shall deal with the following general problem: Let $f(x^1, x^2, \dots, x^r, \theta^1, \dots, \theta^k)$ be the joint probability density func-

pLR has best average power for hypos equidistant to a given POI value

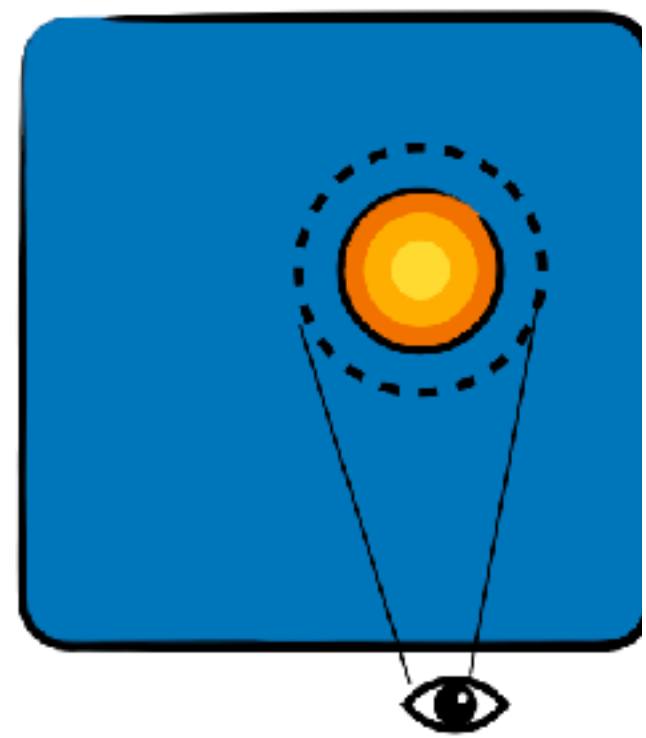


$$t(x|\mu) = -2 \log \frac{p(x|\mu, \hat{\nu})}{p(x|\hat{\mu}, \hat{\nu})}$$

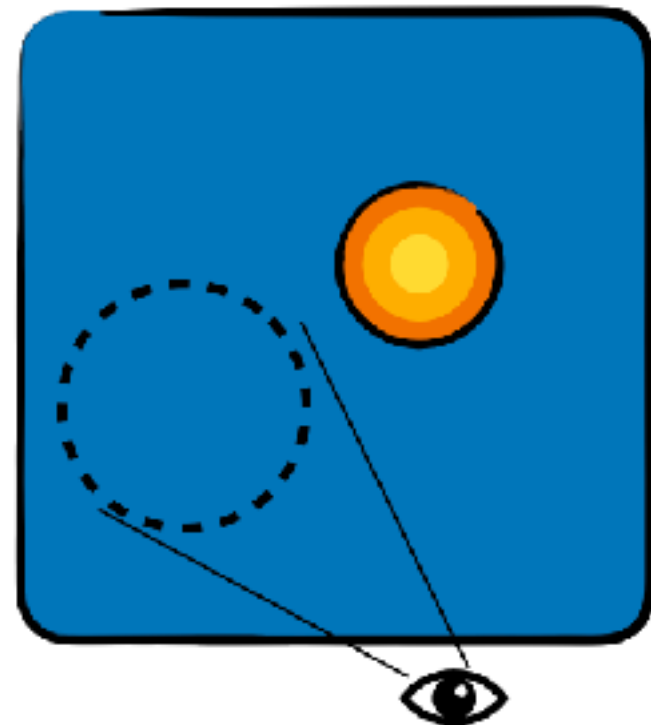
Example: Learning to Profile

For counting experiments, this was easy to show, i.e. sets where we ~only consider the cardinality!

$$f(\mathcal{D}) \rightarrow c_1, c_2, \dots, c_n$$

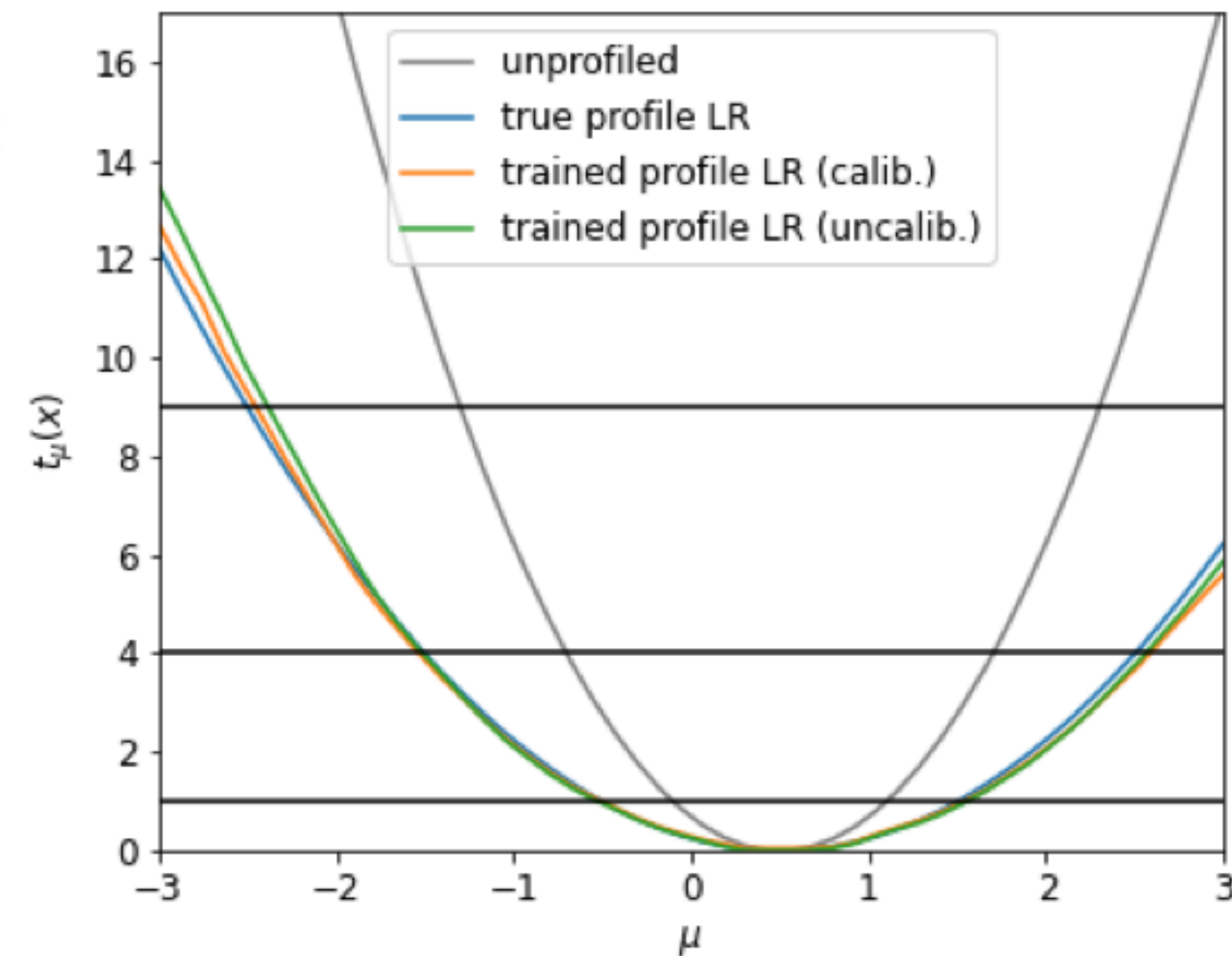
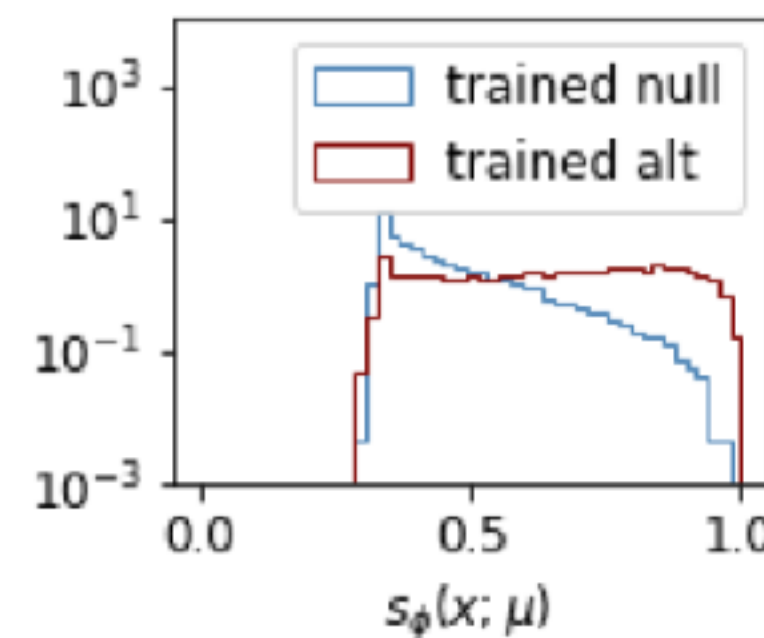
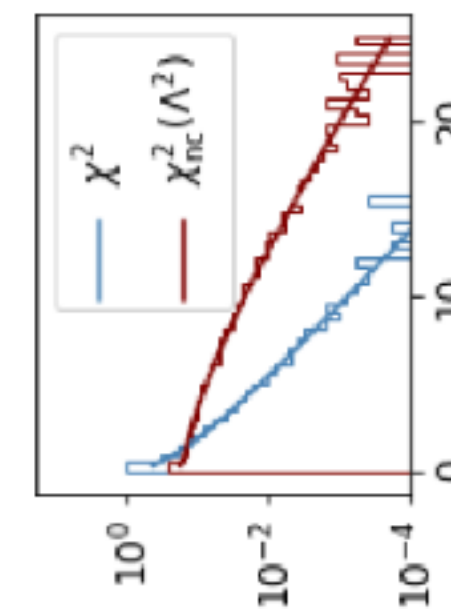
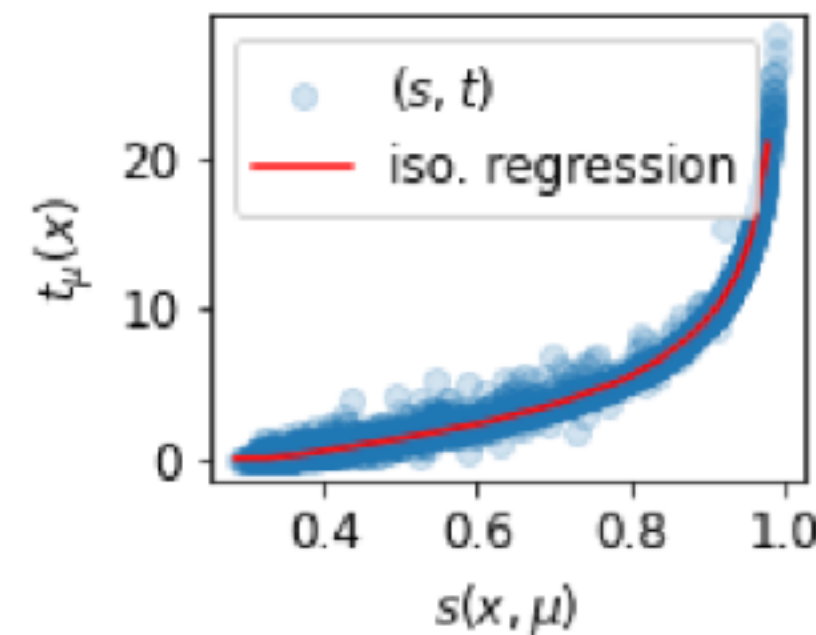


main measurement



side-band

$$p(c_1, c_2 | \mu, \nu) = \text{Pois}(c_1 | \mu s + \nu b) \text{Pois}(c_2 | \nu)$$



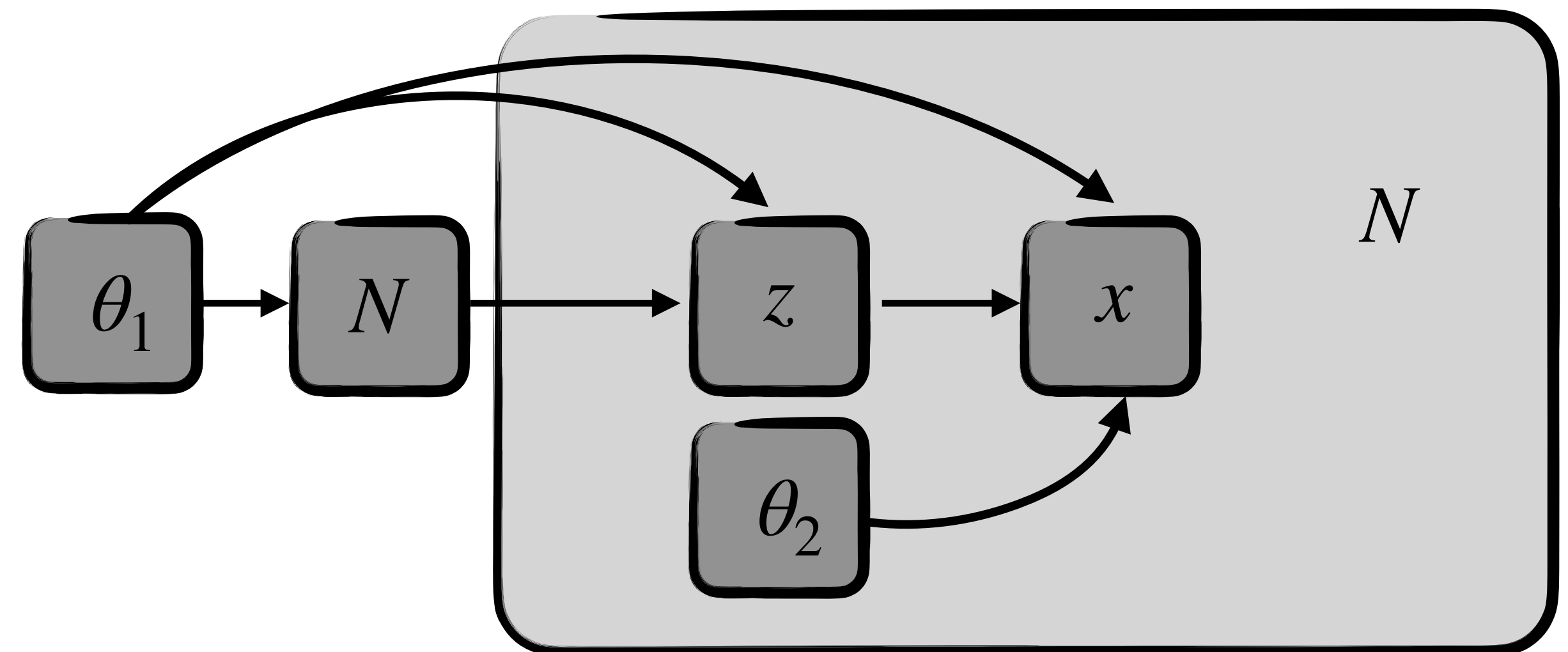
Beyond Cardinalities

If we want to go beyond simple examples, the SBI workflow at the dataset level becomes more complicated

The simulator is now a **simulator of variable-sized sets**

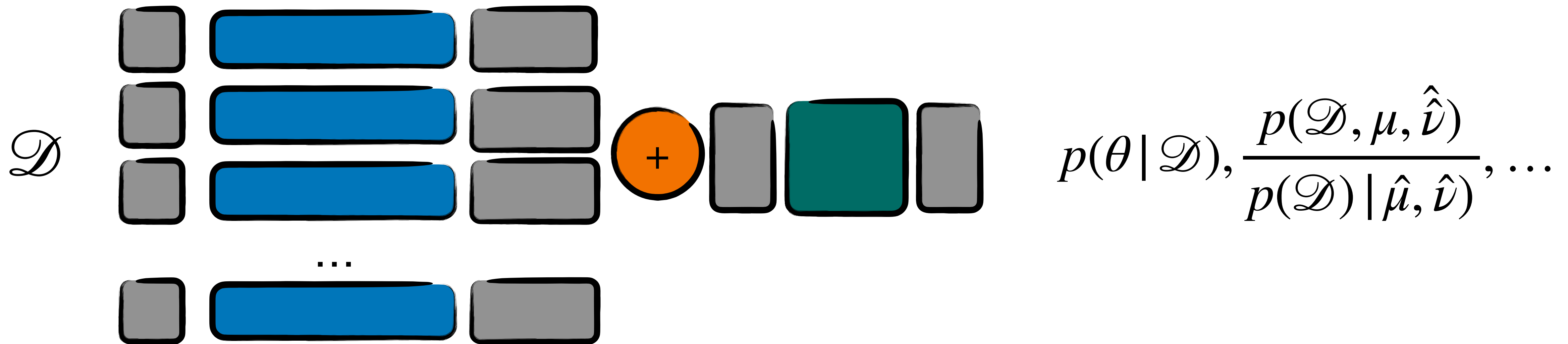
→ some θ impact the set itself, some only set members

$$(\mathcal{D}, \theta) \sim p(\mathcal{D} | \theta)p(\theta)$$



The naive approach

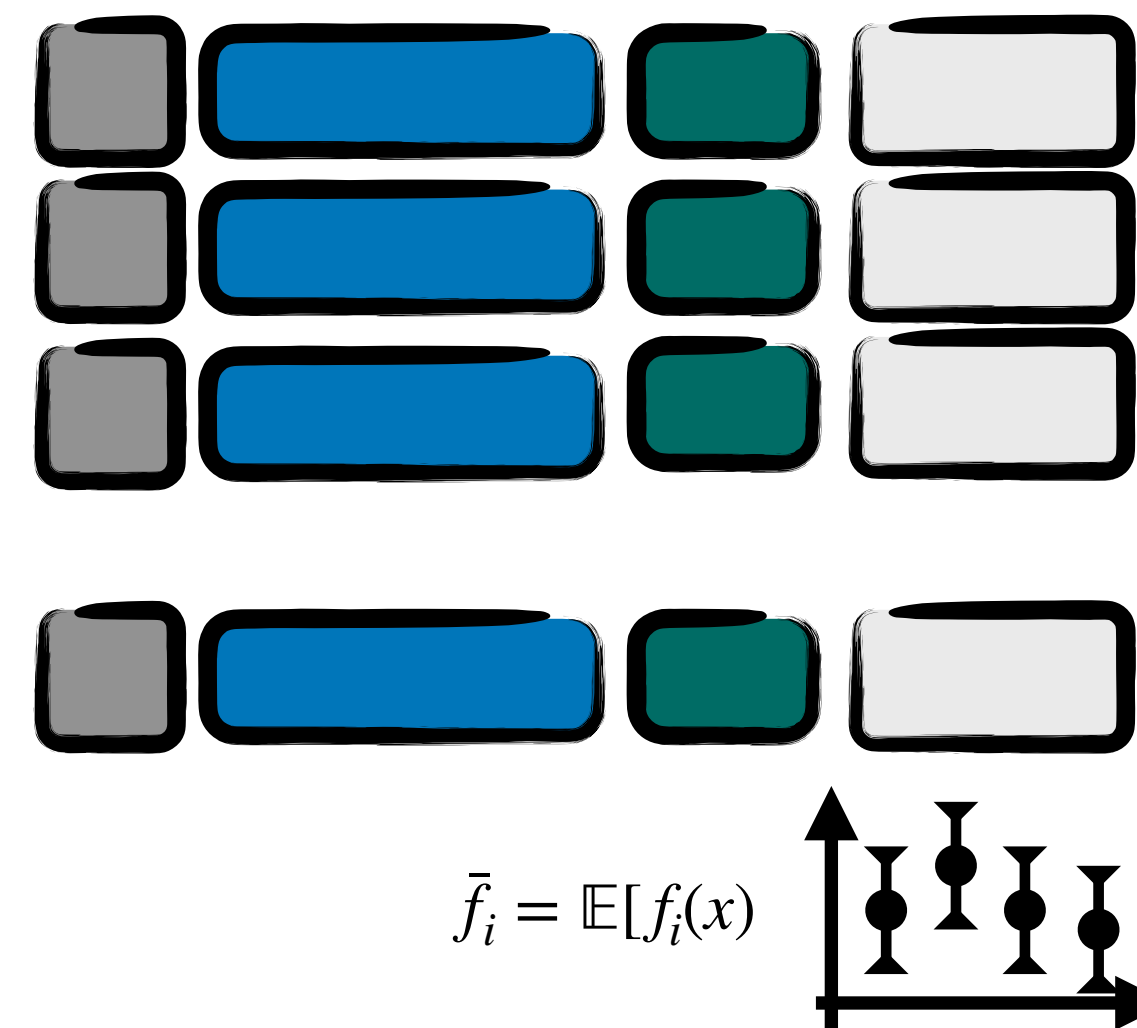
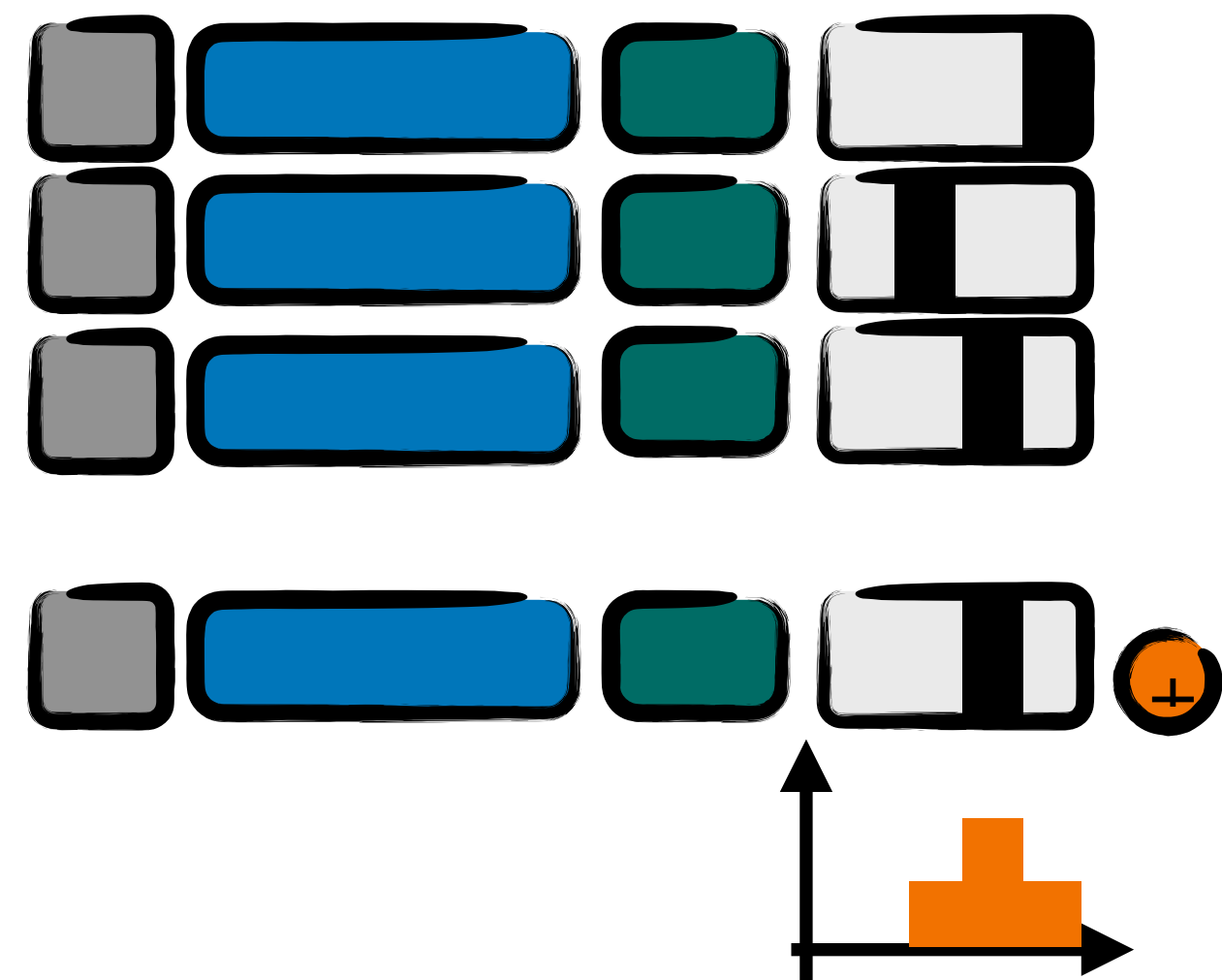
There is a rich literature for set embeddings in ML literature
→ foundational method Deep Sets, then Transformers ...



This structure is sufficient to model any functions of a set

Small Diversion

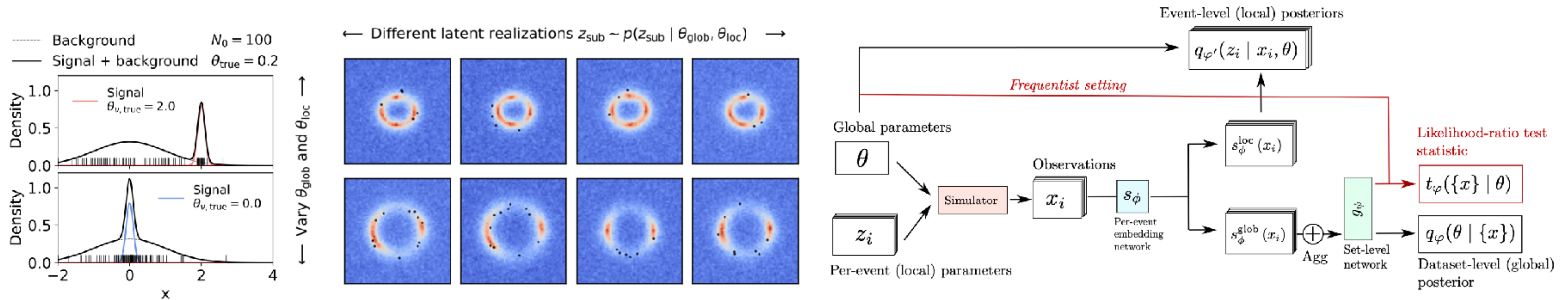
We can see that Deep Sets are sufficient, because they are either collections of moments or essentially histograms
→ Riesz-Markov-Kakutani Theorem and related



During learning, Deep Sets choose what projection to histogram what expectations to measure in order to fully characterize $p(x | \theta)$

Unsurprisingly it works

The standard SBI workflow on sets (of course) works



.. but it doesn't scale to very large sets (thousands, millions)

→ every forward pass processes the full set, scales with $|\mathcal{D}|$

Scaling to Large Sets

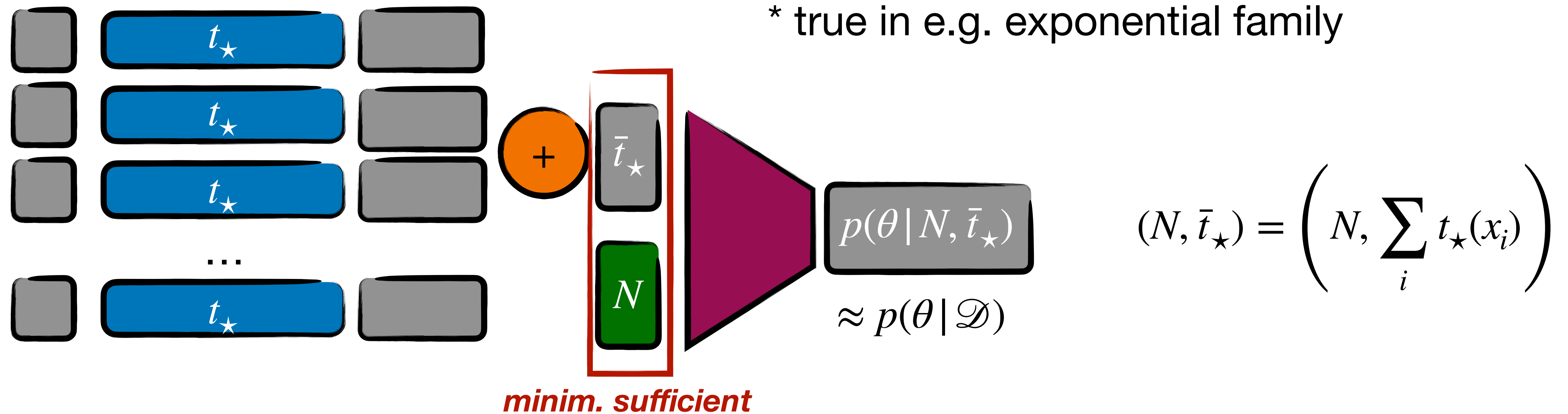
To scale to large sets, we're looking for embeddings of \mathcal{D} that are approximately (Bayesian) sufficient for

$$p(\theta | \mathcal{D}) \approx p(\theta | E_\phi(\mathcal{D}))$$

Question: Can we find a training procedure that produces a sufficient embedding of datasets without training on full \mathcal{D}

Scaling to Large Sets

Assumption*: for inferring θ there is a per-sample embedder $t_{\star}(x)$ whose pooling + cardinality $|\mathcal{D}|$ is minim. sufficient



How can we efficiently find the embedder $t_{\star}(x)$?

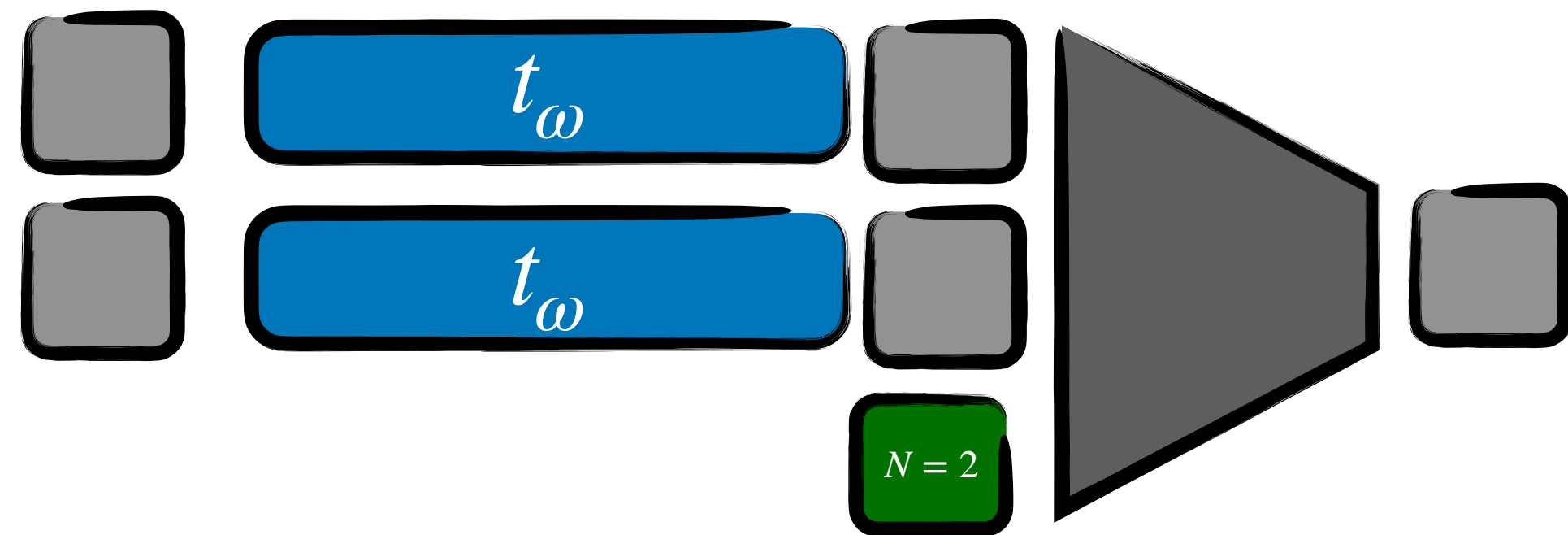
PAIRS*: Small-Cardinality Pretraining

Simple result: we can find a sufficient embedder by pre-training on small-cardinalities ($N=2$)

* *Pretraining Aggregators for Inference at Arbitrary Set-Sizes*

PAIRS*: Small-Cardinality Pretraining

Simple result: we can find a sufficient embedder by pre-training on small-cardinalities ($N=2$)

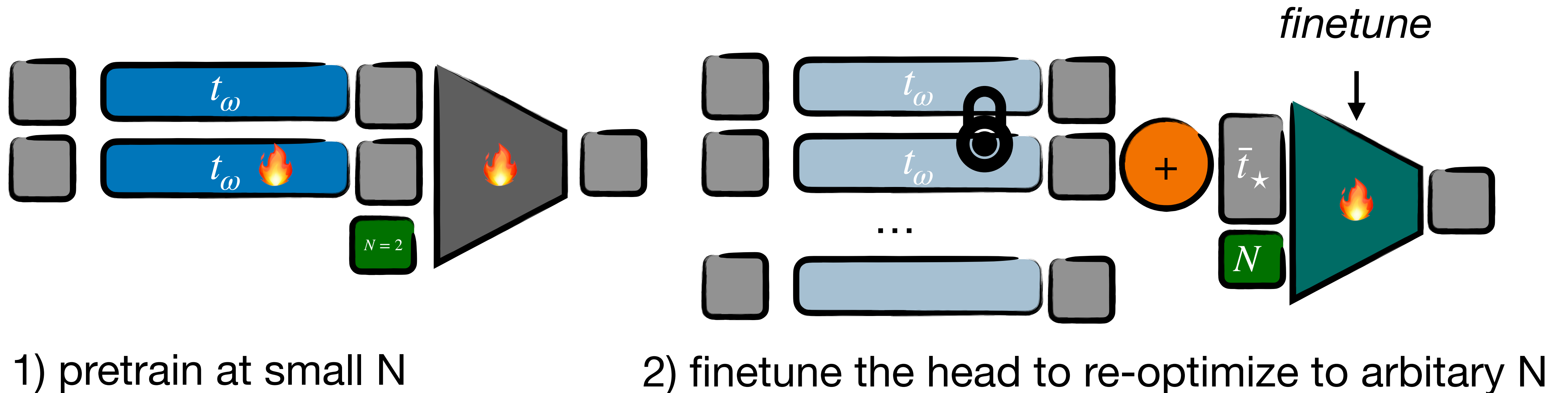


pretrain at small N
($N = 1, 2$)

* *Pretraining Aggregators for Inference at Arbitrary Set-Sizes*

PAIRS*: Small-Cardinality Pretraining

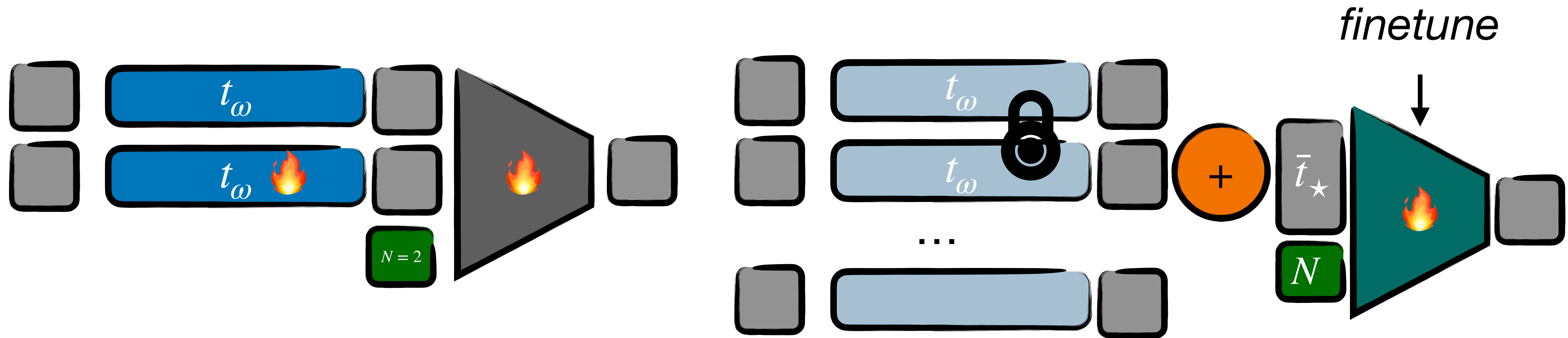
Simple result: we can find a sufficient embedder by pre-training on small-cardinalities ($N=2$)



* *Pretraining Aggregators for Inference at Arbitrary Set-Sizes*

PAIRS*: Small-Cardinality Pretraining

Simple result: we can find a sufficient embedder by pre-training on small-cardinalities ($N=2$)



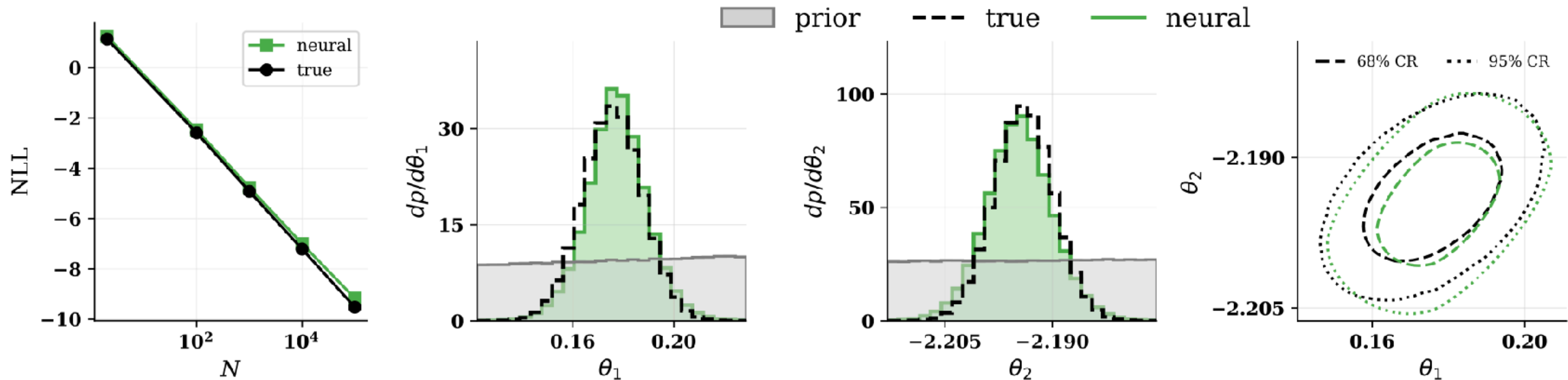
1) pretrain at small N

2) finetune the head to re-optimize to arbitrary N

Large- N Set Embed can be cached! Not part of forward pass

PAIRS*: Small-Cardinality Pretraining

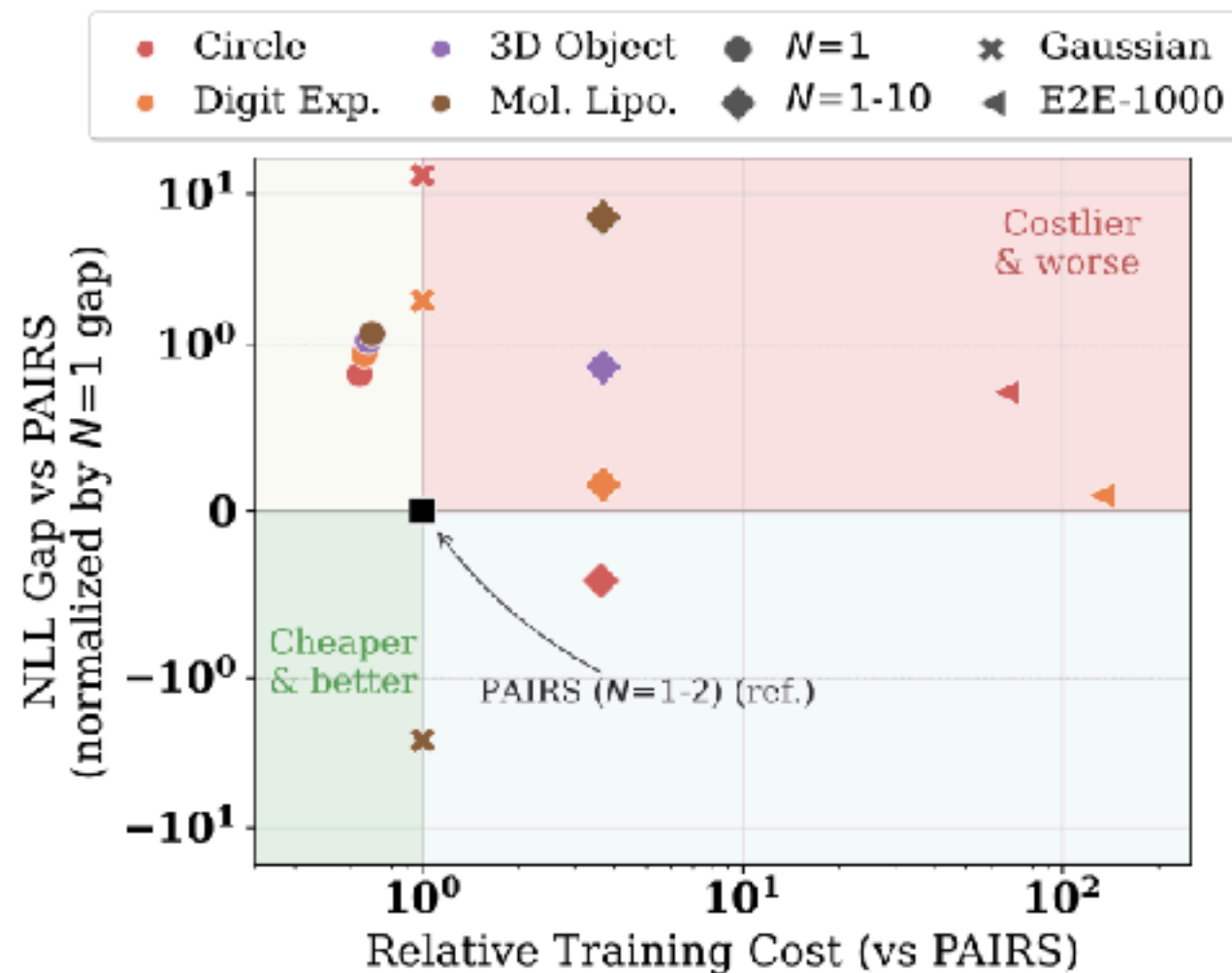
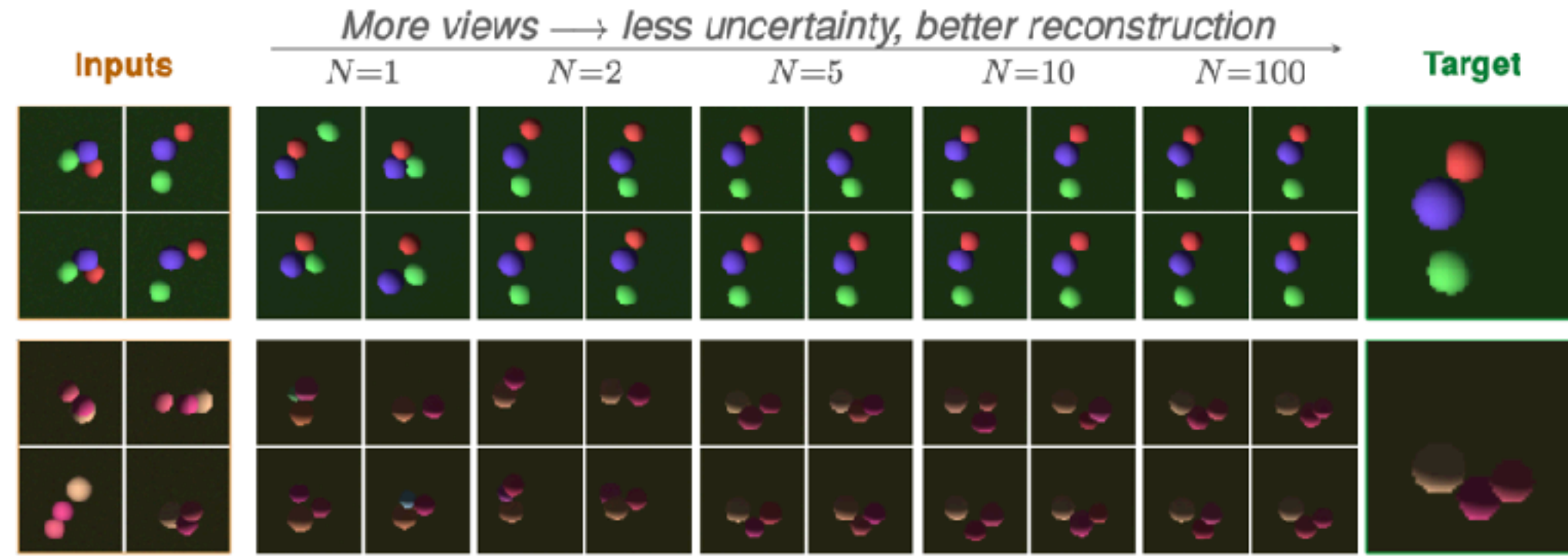
Consequence: we can do a fully neural end-to-end inference pipeline with no explicit MCMC / Profiling, etc
→ without ever training on full-scale sets



* *Pretraining Aggregators for Inference at Arbitrary Set-Sizes*

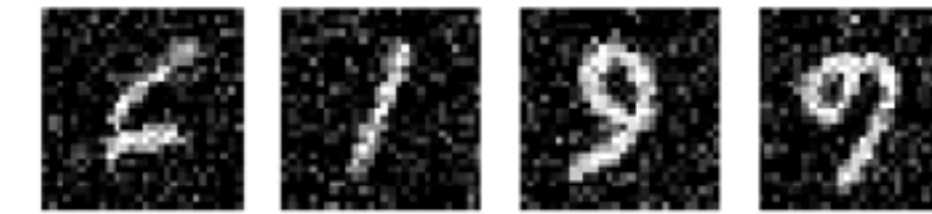
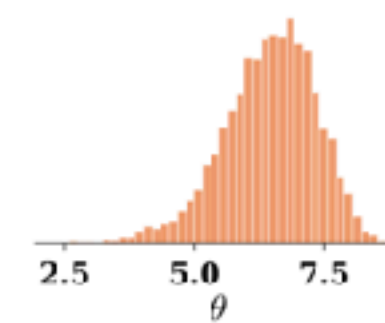
Applications

Very motivated by HEP, but set-level inference appears in many applications in biology, astronomy, vision,

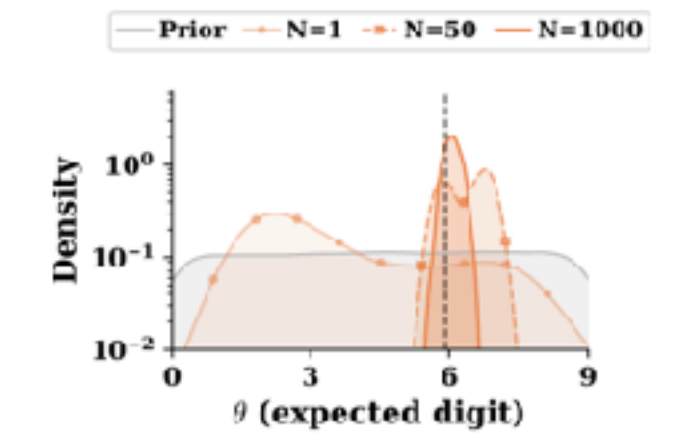


Digit Expectation

$\theta = \mathbb{E}[\text{digit}] \in [0, 9]$, $x_i \in \mathbb{R}^{28 \times 28}$ (rotated MNIST image)

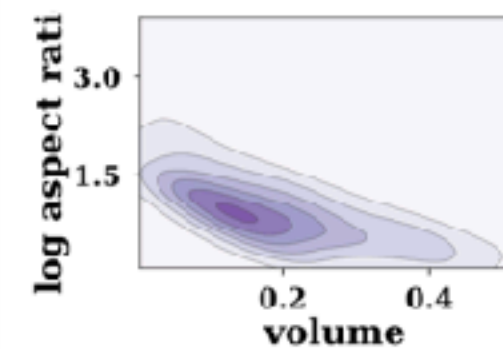


shared rotation, additive noise ($\sigma=0.2$)

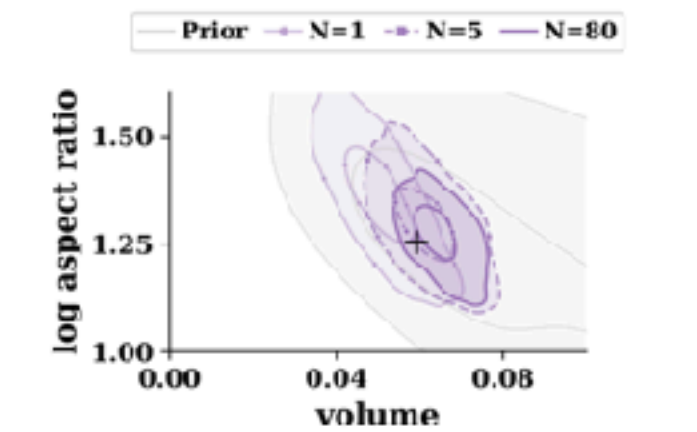


3D Object Properties

$\theta \in \mathbb{R}^2$ (volume, log aspect ratio), $x_i \in \mathbb{R}^{224 \times 224}$ (rendered view)

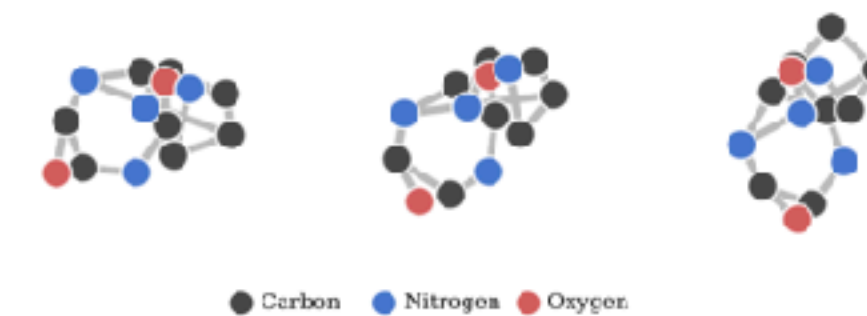
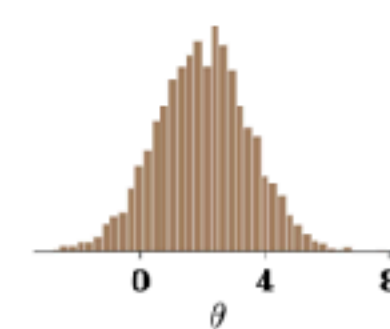


clean reference \rightarrow partial view (35-50% visible)

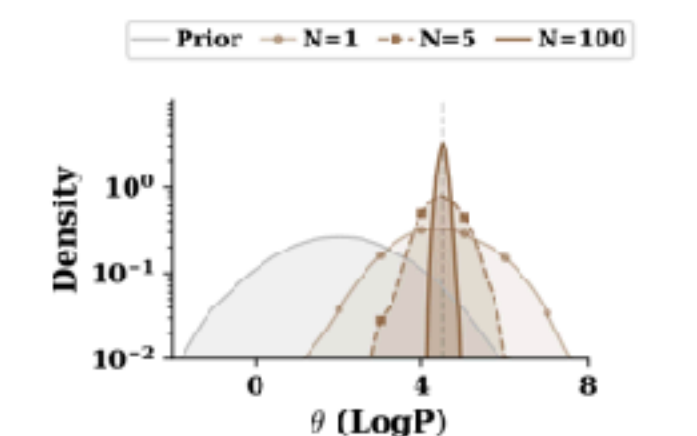


Molecular Lipophilicity

$\theta = \text{LogP} \in \mathbb{R}$ (partition coefficient), $x_i = 3\text{D molecular graph (conformer)}$



Carbon Nitrogen Oxygen



Summary & Outlook

SBI in Particle Physics is tough because of three “high-D”:

→ instance size, dataset cardinality, parameter space

Problem can be generalized / formalized as “end-to-end SBI on large-scale i.i.d. sets”

For a large-class of problems pretrain-finetune is an efficient strategy to produces an embedder that is still ~optimal

Can do dataset-level SBI without $O(N)$ forward pass and without a secondary stats. step

Proof Sketch

To know a iid dataset \mathcal{D} “sufficiently” you “just” need to know

i) know $p(x)$

ii) know N

ad i) if you know sufficiently many moments you know $p(x)$ (Riesz)

$$\bar{f}_k = \int f_k(x)p(x) \quad \text{with } f_k \in \mathcal{F} \quad \bar{f}_k \approx \frac{1}{N} \sum_k f_k(x_i)$$

ergo: knowing $\left(\frac{1}{N} \sum_k f_k(x_i), N \right) \approx$ sufficient stat $f(\mathcal{D})$