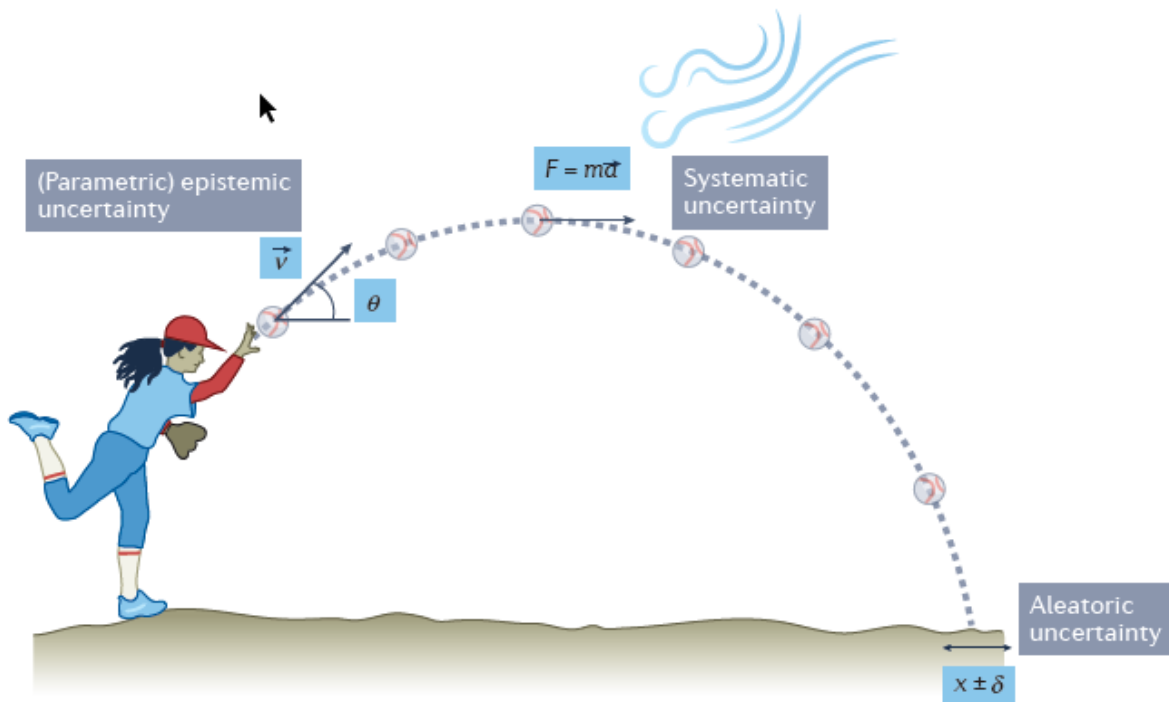# Simulation-based inference for the physical sciences

Grenoble Artificial Intelligence for Physical Sciences
May 29, 2024

Gilles Louppe
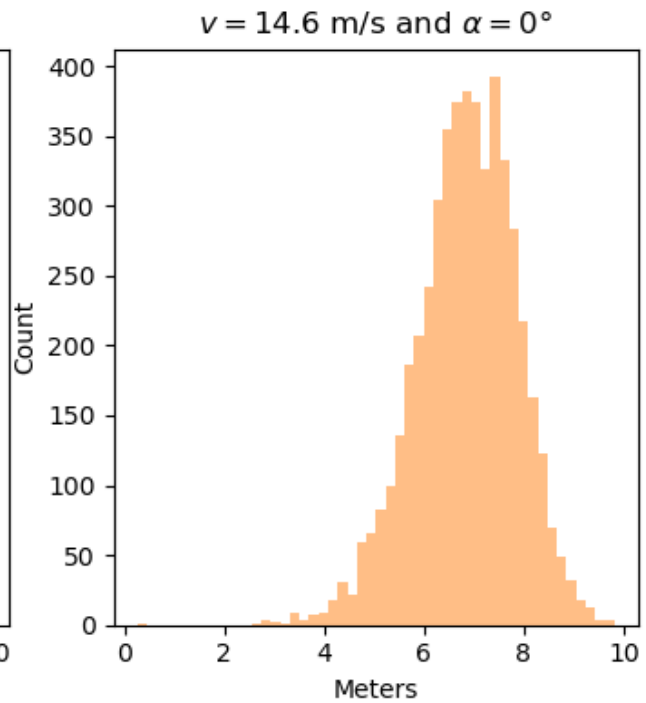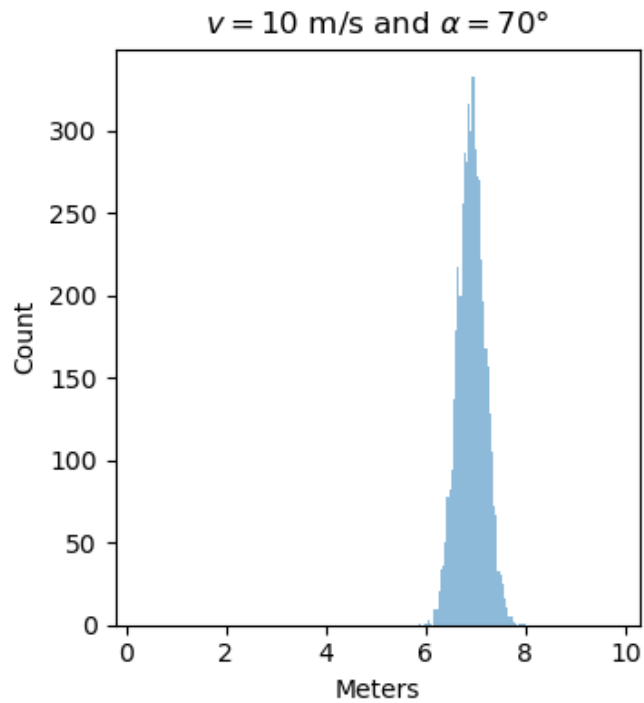g.louppe@uliege.be

(Parametric) epistemic uncertainty

$\vec{v}$

$\theta$

$F = m\vec{a}$

Systematic uncertainty

Aleatoric uncertainty

$x \pm \delta$
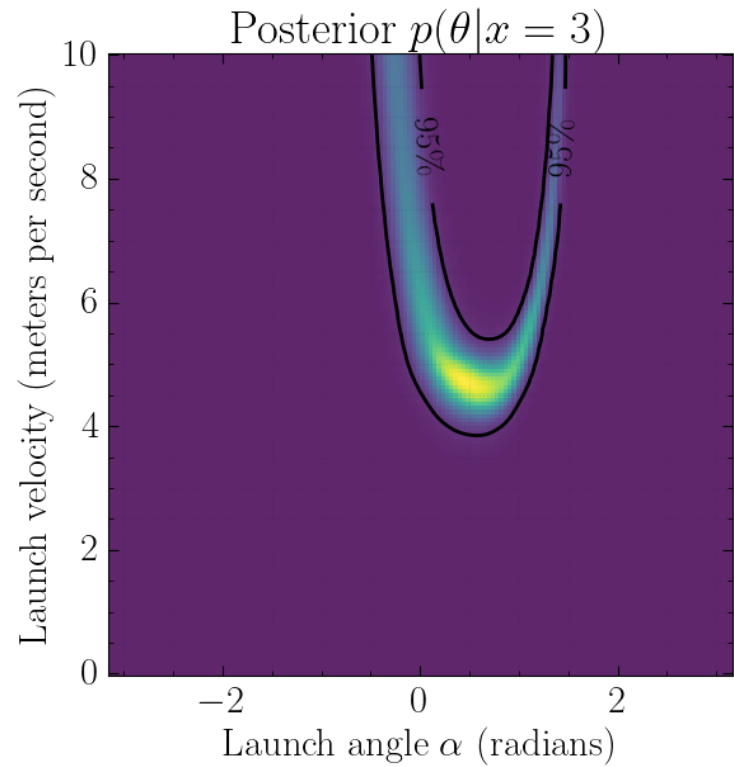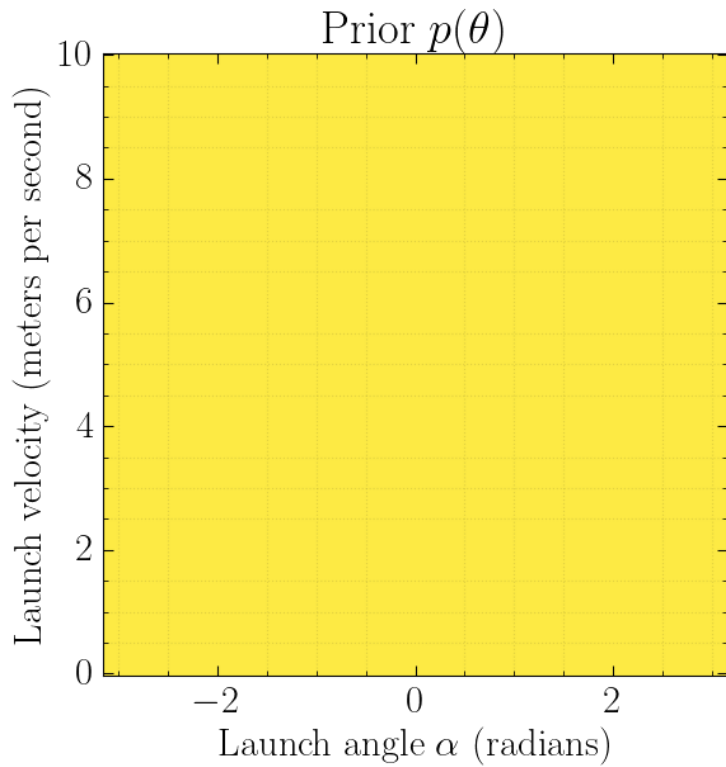
$$v_x = v\cos(\alpha), \quad v_y = v\sin(\alpha),$$

$$\frac{dx}{dt} = v_x, \quad \frac{dy}{dt} = v_y, \frac{dv_y}{dt} = -G.$$

```python
def simulate(v, alpha, dt=0.001):
    v_x = v * np.cos(alpha)  # x velocity m/s
    v_y = v * np.sin(alpha)  # y velocity m/s
    y = 1.1 + 0.3 * random.normal()
    x = 0.0

    while y > 0: # simulate until ball hits floor
        v_y += dt * -G  # acceleration due to gravity
        x += dt * v_x
        y += dt * v_y

    return x + 0.25 * random.normal()
```

What parameter values $\theta$ are the most plausible?

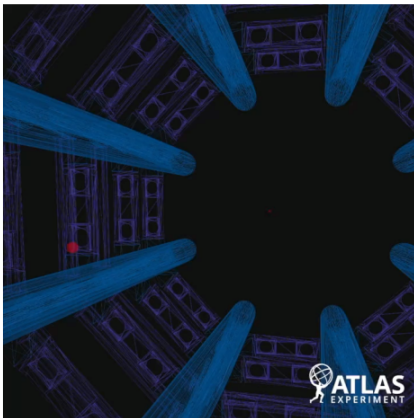Prior $p(\theta)$     Posterior $p(\theta|x=3)$

# Simulation-based inference

# Simulators as generative models

A simulator prescribes a generative model that can be used to simulate data $\mathbf{x}$.

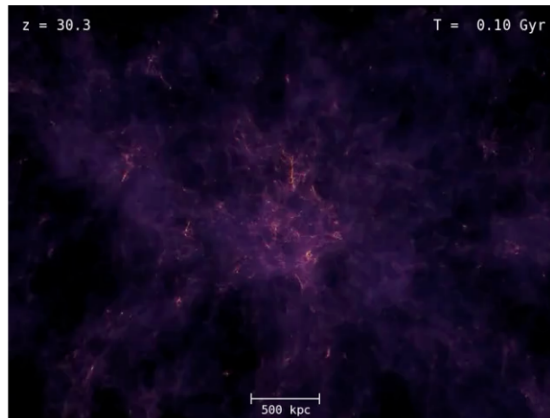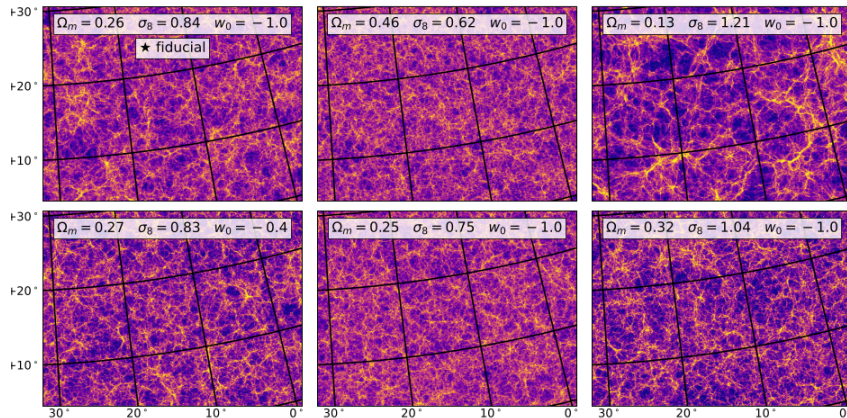| Collider data | Cosmology data | Molecular dynamics |
|:---:|:---:|:---:|
| particles $\sim p(\text{particles})$ | particles $\sim p(\text{particles})$ | configurations $\sim p(\text{configurations})$ |



[C. Cesarotti with ATLAS]

[Aquarius simulation]

[E. Cances et al]

# Conditional simulators

A conditional simulator prescribes a way to sample from the likelihood $p(\mathbf{x}|\theta)$, where $\theta$ is a set of conditioning variables or parameters.

**Cosmology data**

$$\text{map} \sim p(\text{map} \mid \{\Omega_m, \sigma_8, w_0\})$$



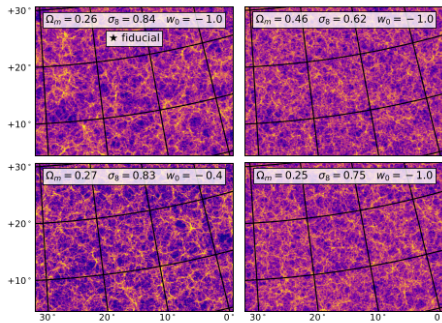[Kacprzak et al 2022]

$$x \sim p(x; \mathcal{M})$$

*Model*

*or*

$$x \sim p(x \mid \theta)$$

*Model parameters*

# What can we do with generative models?

| Produce samples and make predictions | Evaluate densities | Encode complex priors |
|---|---|---|

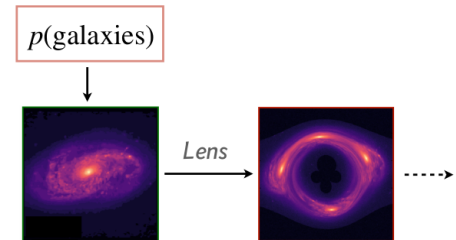$$\mathbf{x} \sim p(\mathbf{x}|\theta)$$

$$p(\mathbf{x}|\theta)$$

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})}$$

$$p(\mathbf{x})$$



[Kacprzak et al 2022]





$p(\text{galaxies})$

Lens

# Inference



- Frequentist inference: find the parameters $\hat{\theta}$ that maximizes the likelihood $p(\mathbf{x}|\theta)$ or build a confidence interval thereof.

- Bayesian inference: compute the posterior distribution

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})}$$

of the parameters $\theta$ given the data $\mathbf{x}$.
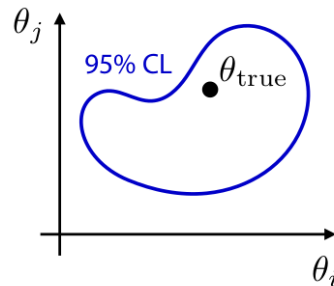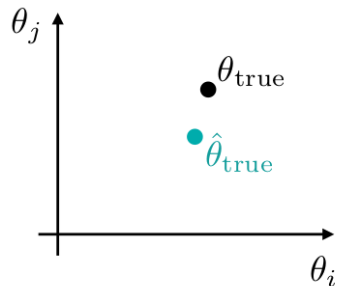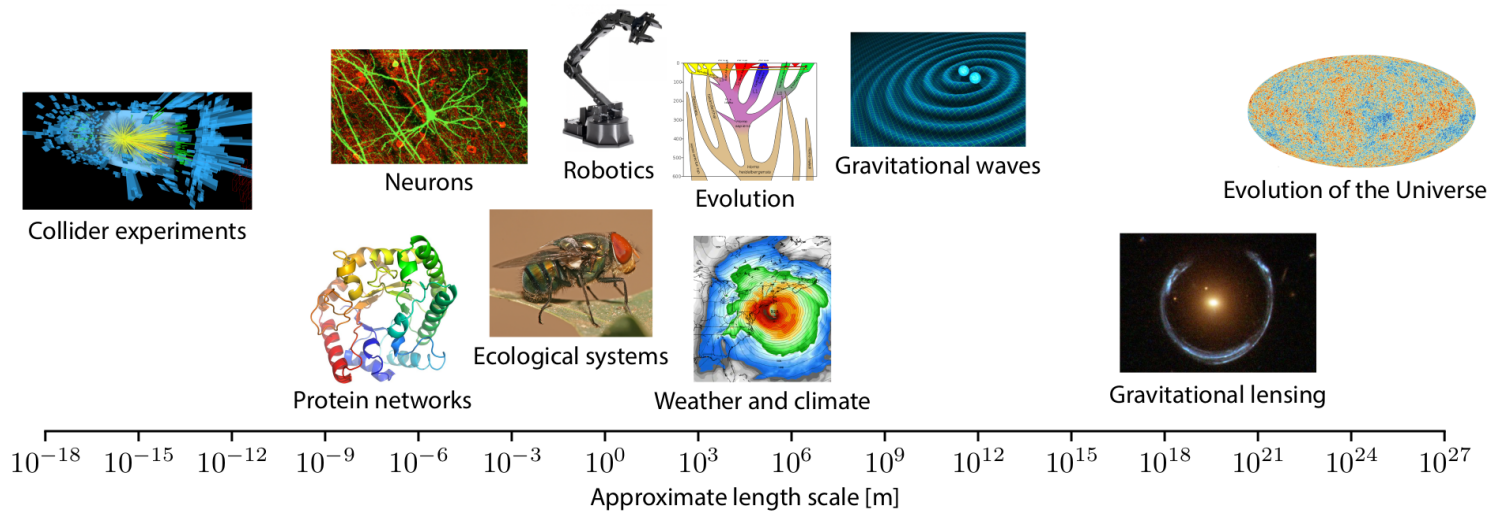
Collider experiments — Neurons — Robotics — Evolution — Gravitational waves — Evolution of the Universe — Protein networks — Ecological systems — Weather and climate — Gravitational lensing

Approximate length scale [m]

$10^{-18}$ $10^{-15}$ $10^{-12}$ $10^{-9}$ $10^{-6}$ $10^{-3}$ $10^{0}$ $10^{3}$ $10^{6}$ $10^{9}$ $10^{12}$ $10^{15}$ $10^{18}$ $10^{21}$ $10^{24}$ $10^{27}$

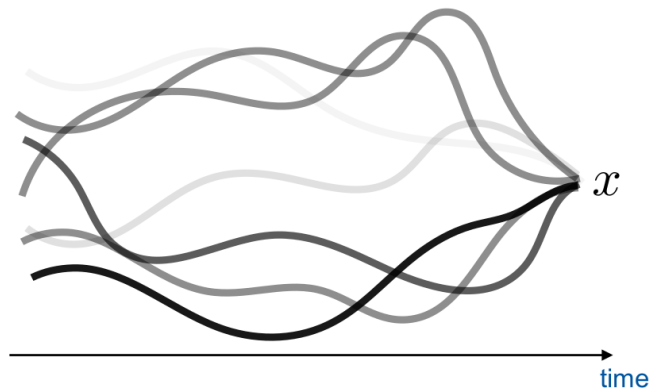Examples of inference problems across the physical sciences

- Discovering new particles in high-energy physics

- Data assimilation in weather forecasting

- Estimating gravitational wave parameters

- Retrieving atmospheric properties of exoplanets
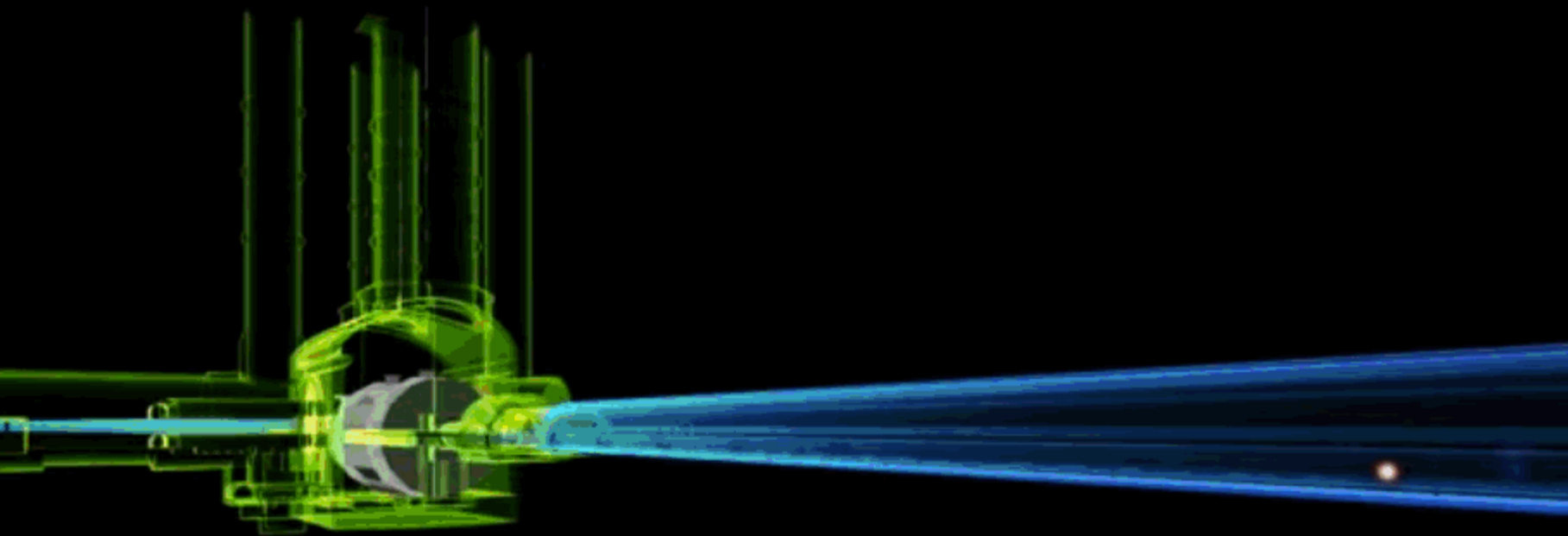
- Constraining cosmological models from galaxy surveys

# Intractable likelihoods

The (modeled) data generating process may involve additional latent variables $\mathbf{z}$ that are not observed, leading to likelihoods

$$p(\mathbf{x}|\theta) = \int p(\mathbf{x}, \mathbf{z}|\theta)d\mathbf{z}.$$

In this case, evaluating the likelihood becomes intractable.
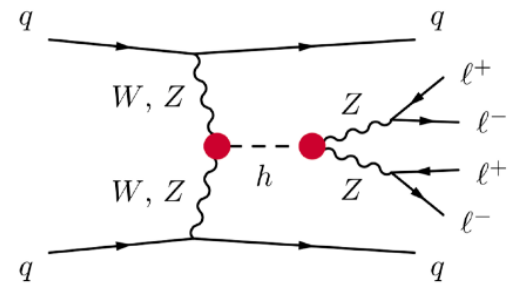
$p(\mathbf{z}_p | \theta)$

Latent variables                    Parameters
                                    of interest

Parton-level              Theory
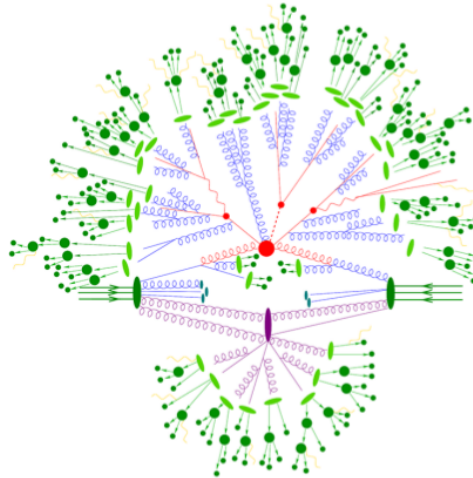momenta                   parameters

$z_p \longleftarrow \theta$

$$p(\mathbf{z}_s|\theta) = \int p(\mathbf{z}_p|\theta)p(\mathbf{z}_s|\mathbf{z}_p)d\mathbf{z}_p$$

Latent variables

Parameters
of interest
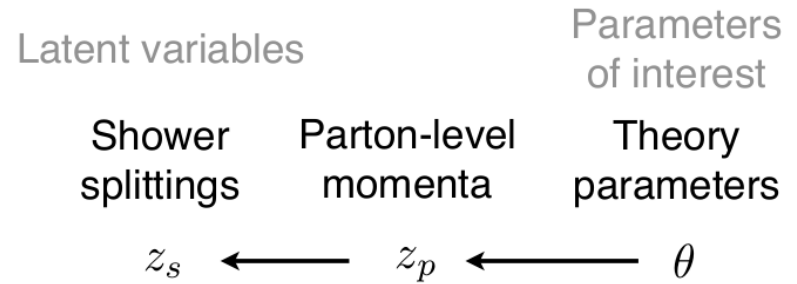
Shower
splittings

Parton-level
momenta

Theory
parameters

$z_s \longleftarrow z_p \longleftarrow \theta$

$$p(\mathbf{z}_d|\theta) = \iint p(\mathbf{z}_p|\theta)p(\mathbf{z}_s|\mathbf{z}_p)p(\mathbf{z}_d|\mathbf{z}_s)d\mathbf{z}_p d\mathbf{z}_s$$

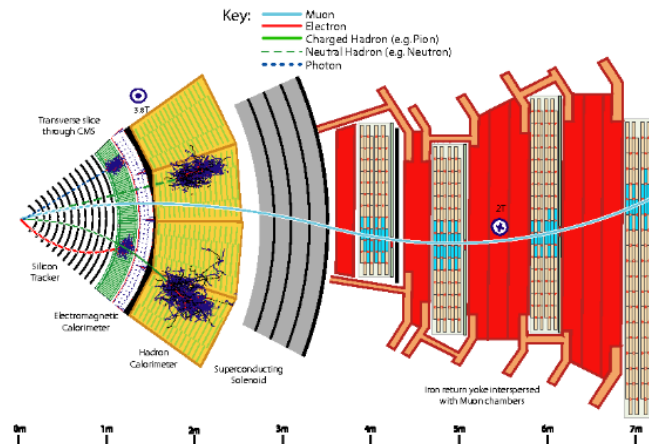$$p(\mathbf{x}|\theta) = \underbrace{\iiint}_{\text{yikes!}} p(\mathbf{z}_p|\theta)p(\mathbf{z}_s|\mathbf{z}_p)p(\mathbf{z}_d|\mathbf{z}_s)p(\mathbf{x}|\mathbf{z}_d)\,d\mathbf{z}_p\,d\mathbf{z}_s\,d\mathbf{z}_d$$

| Features | | Latent variables | | | Parameters of interest |
|---|---|---|---|---|---|
| Observables | | Detector interactions | Shower splittings | Parton-level momenta | Theory parameters |
| $x$ | $\longleftarrow$ | $z_d$ $\longleftarrow$ | $z_s$ $\longleftarrow$ | $z_p$ $\longleftarrow$ | $\theta$ |



[Image source: M. Cacciari, G. Salam, G. Soyez 0802.1189]

Parameters $\theta$ $\longrightarrow$ Simulator / Latent $z$ $\longrightarrow$ Observables $x$

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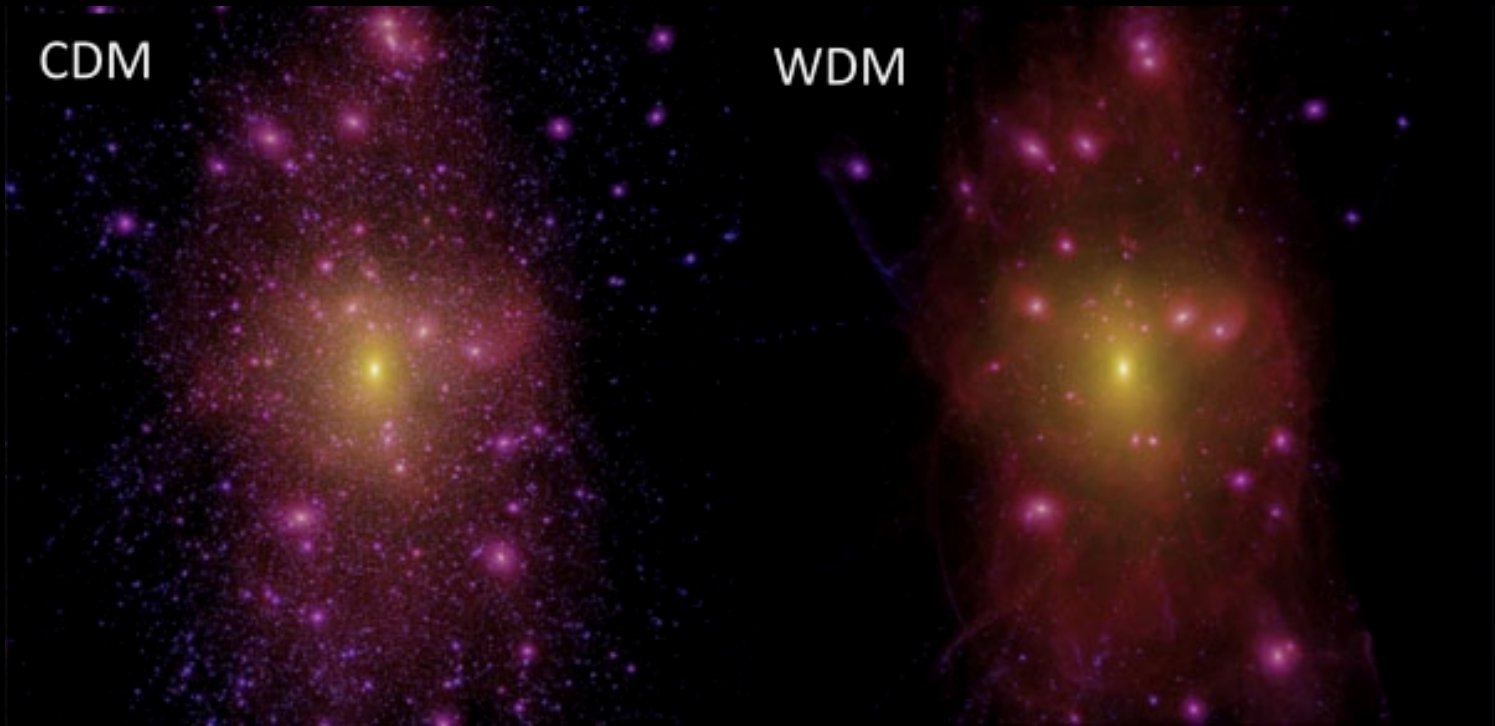- Simulator can generate samples $x \sim p(x|\theta)$

Inference:
- Interactions between low-level components lead to challenging inverse problems
- Likelihood $p(x|\theta) = \int \mathrm{d}z\ p(x,z|\theta)$ is intractable

Statistical inference becomes challenging when the likelihood $p(\mathbf{x}|\boldsymbol{\theta})$ is implicit or intractable. **Simulation-based inference algorithms are needed.**

# A case study

Hermans et al, "Constraining dark matter with stellar streams", 2021.

Can we constrain the nature of dark matter from cosmological observations?

# Constraining dark matter with stellar streams



**Palomar 5
(Pal5) stream**
Pal5 was discovered in 2001 as
the first thin stream formed from
a globular cluster. Its current orbit
takes it far over the galactic center.

**Globular clusters**
These hives typically hold
100,000 stars or fewer and give
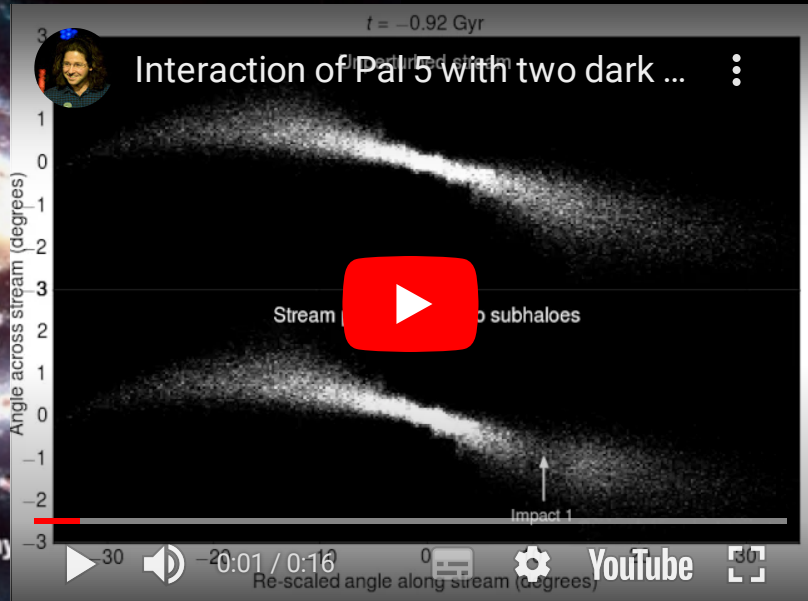rise to long, thin streams.
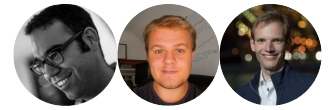
Gap

**GD1 stream**
Discovered in 2006, GD1 is
the longest known thin stream,
stretching across more than half the
northern sky. It contains a gap that could
be the scar of a dark matter collision
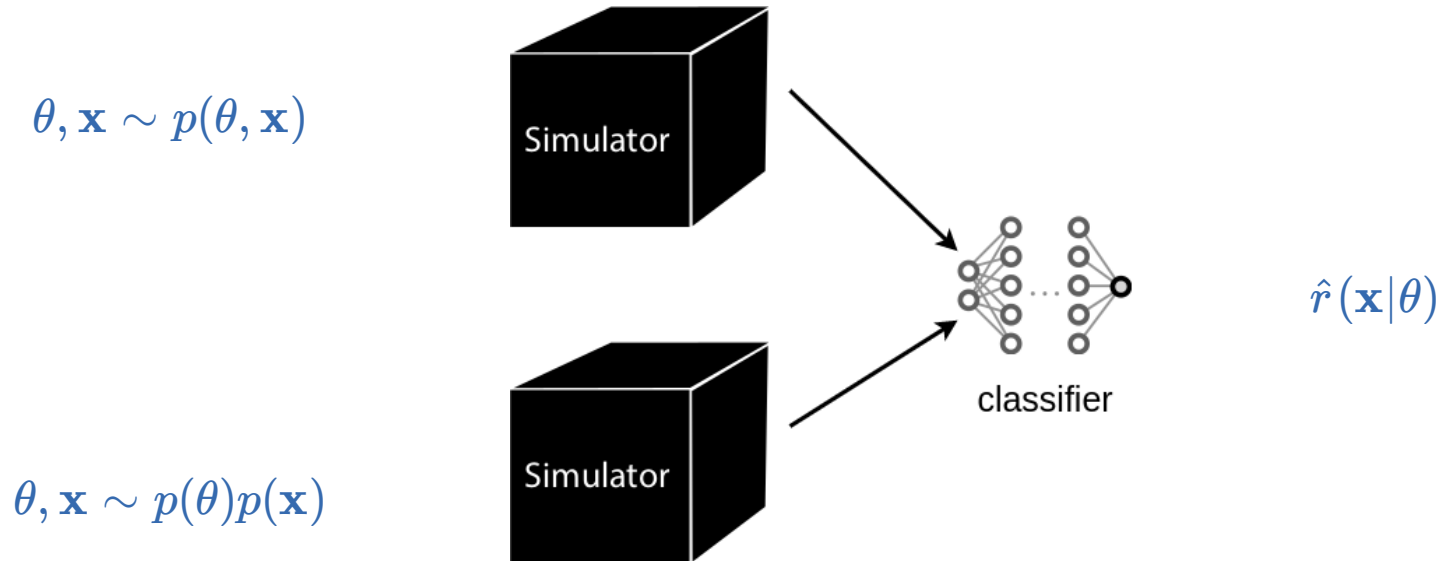500 million years ago.

Sun

Milky Way

$t = -0.92$ Gyr

Interaction of Pal 5 with two dark ...

Stream p...o subhaloes

Angle across stream (degrees)

Re-scaled angle along stream (degrees)

Impact 1

0:01 / 0:16        YouTube

$$p(m_{\mathrm{WDM}}, t_{\mathrm{age}}|\mathrm{GD1}) = \frac{p(\mathrm{GD1}|m_{\mathrm{WDM}}, t_{\mathrm{age}})p(m_{\mathrm{WDM}}, t_{\mathrm{age}})}{p(\mathrm{GD\text{-}1})}$$
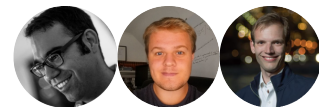
# Neural ratio estimation (NRE)

The likelihood-to-evidence $r(\mathbf{x}|\theta) = \frac{p(\mathbf{x}|\theta)}{p(\mathbf{x})} = \frac{p(\theta,\mathbf{x})}{p(\theta)p(\mathbf{x})}$ ratio can estimated from a binary classifier $d(\theta, \mathbf{x})$, even if neither the likelihood nor the evidence can be evaluated.



$$\theta, \mathbf{x} \sim p(\theta, \mathbf{x})$$

$$\theta, \mathbf{x} \sim p(\theta)p(\mathbf{x})$$
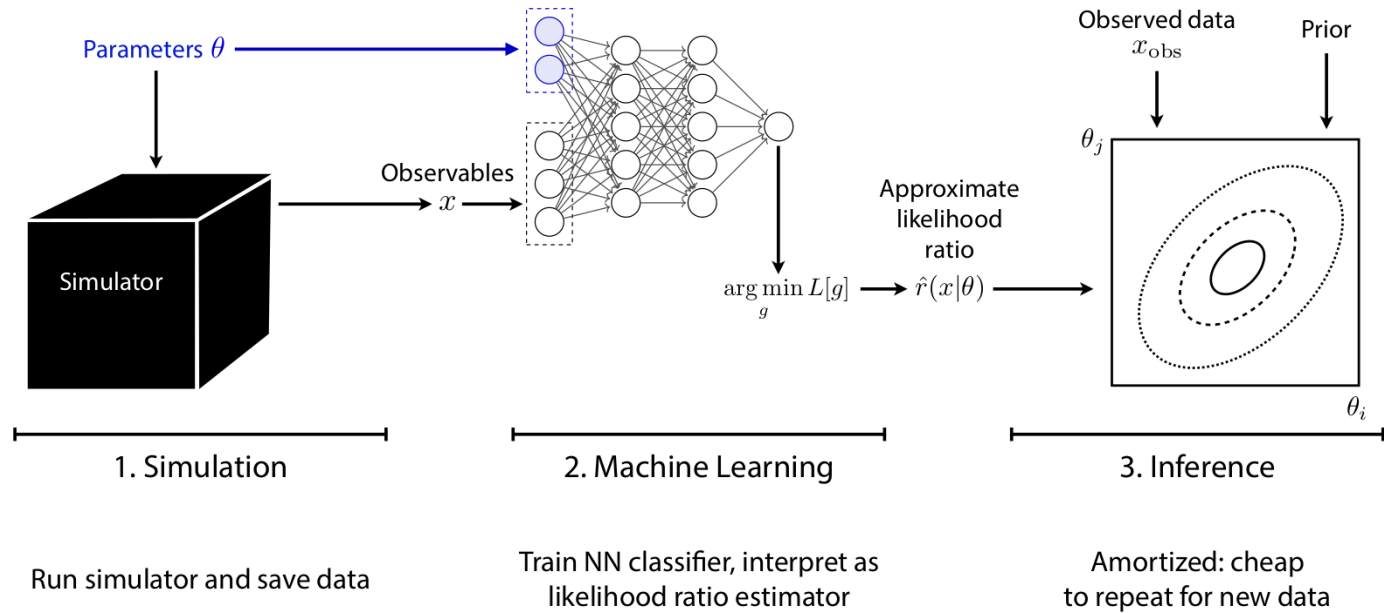
$$\hat{r}(\mathbf{x}|\theta)$$

The solution $d$ found after training approximates the optimal classifier

$$d(\theta, \mathbf{x}) \approx d^*(\theta, \mathbf{x}) = \frac{p(\theta, \mathbf{x})}{p(\theta, \mathbf{x}) + p(\theta)p(\mathbf{x})}.$$
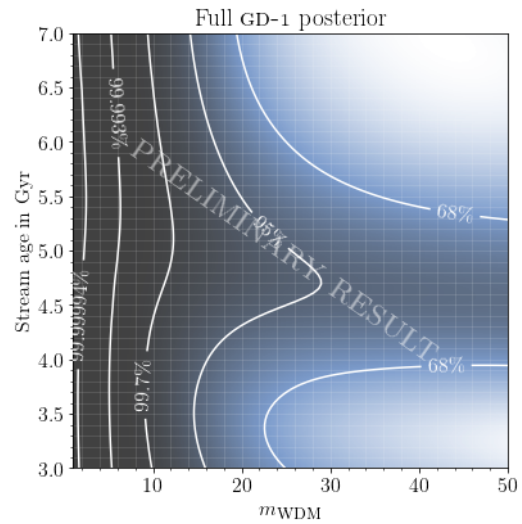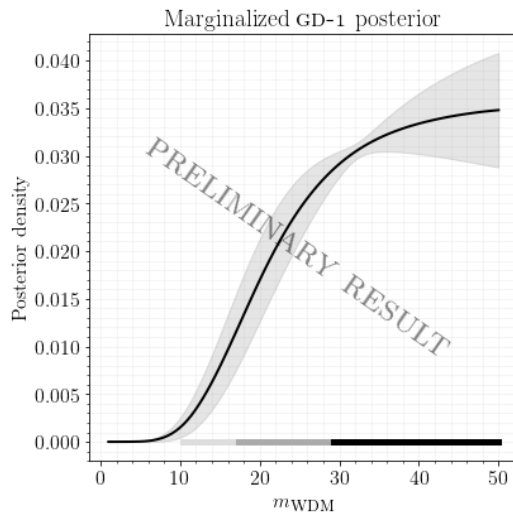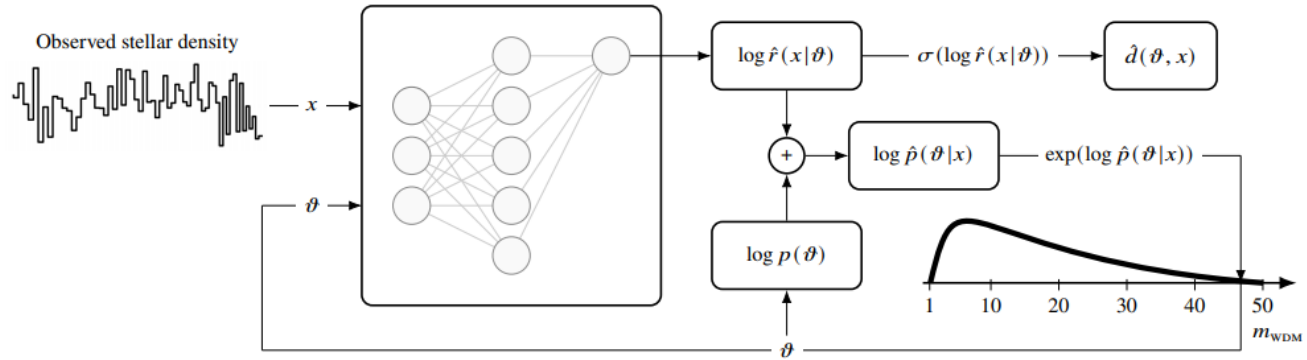
Therefore,

$$r(\mathbf{x}|\theta) = \frac{p(\mathbf{x}|\theta)}{p(\mathbf{x})} = \frac{p(\theta, \mathbf{x})}{p(\theta)p(\mathbf{x})} \approx \frac{d(\theta, \mathbf{x})}{1 - d(\theta, \mathbf{x})} = \hat{r}(\mathbf{x}|\theta).$$

1. Simulation

2. Machine Learning

3. Inference

Run simulator and save data

Train NN classifier, interpret as likelihood ratio estimator

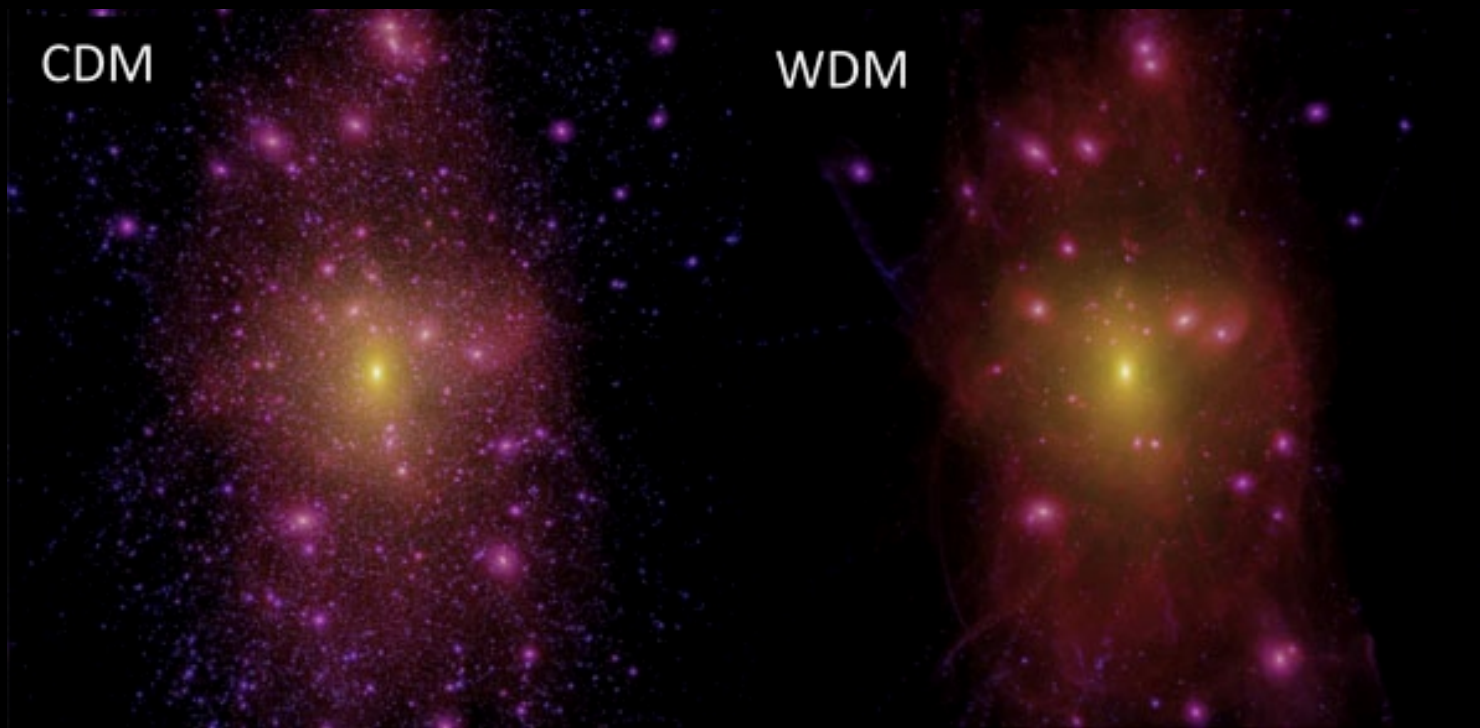Amortized: cheap to repeat for new data

$$p(\theta|\mathbf{x}) \approx \hat{r}(\mathbf{x}|\theta)p(\theta)$$

# NRE for stellar streams

Preliminary results for GD-1 suggest a preference for CDM over WDM.

Wait a minute Gilles...
I can't claim that in a paper!
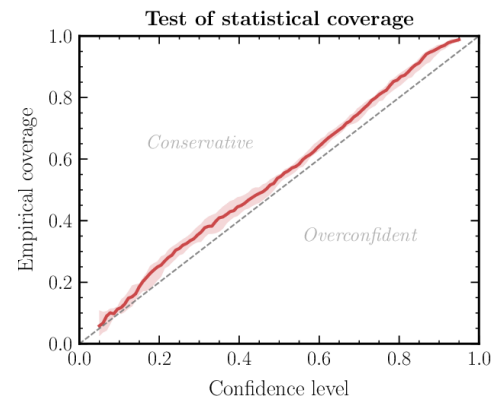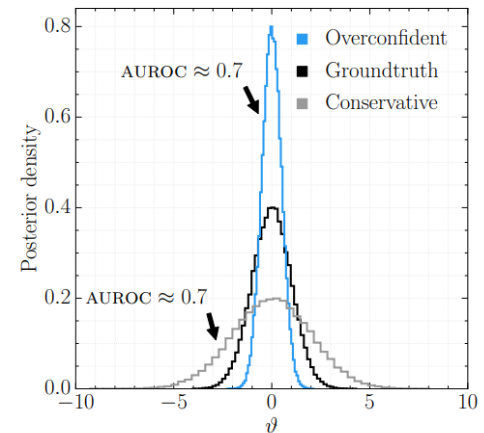Your neural network must be wrong!

# Expected coverage

$$\mathrm{EC}(\hat{p}, \alpha) = \mathbb{E}_{p(\theta, \mathbf{x})}[\theta \in \Theta_{\hat{p}(\theta|\mathbf{x})}(\alpha)]$$
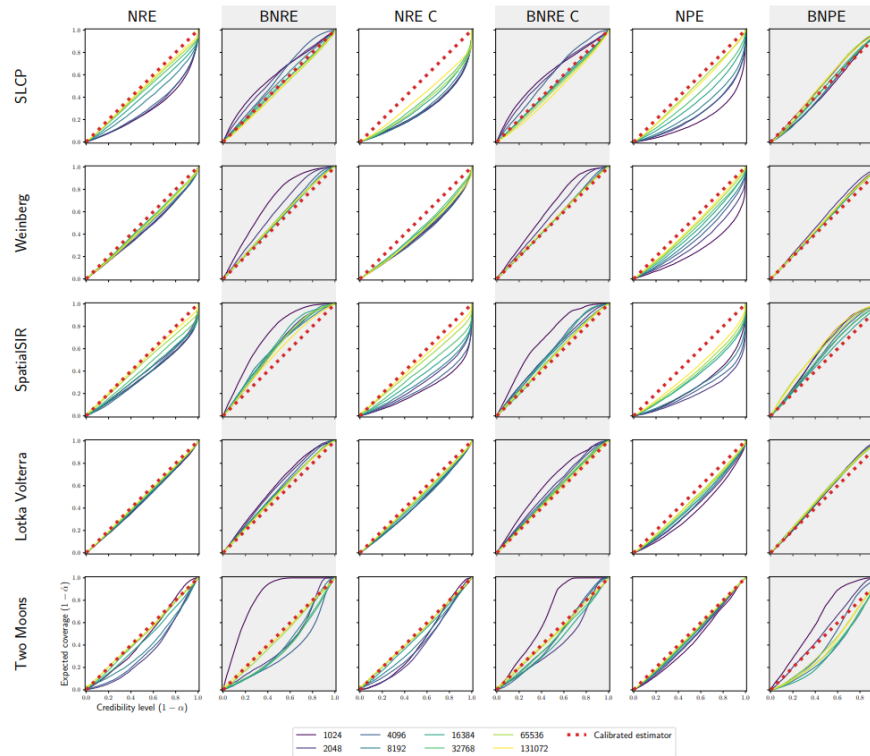
If the expected coverage is close to the nominal coverage probability $\alpha$, then the approximate posterior $\hat{p}$ is calibrated.

- If $\mathrm{EC} < \alpha$, then the posterior is underdispersed and overconfident.

- If $\mathrm{EC} > \alpha$, then the posterior is overdispersed and conservative.

# Balancing inference
# for conservative posteriors



Conservative posteriors can be obtained by enforcing $d$ to be balanced, i.e. such that $\mathbb{E}_{p(\theta,\mathbf{x})}\left[d(\theta,\mathbf{x})\right] = \mathbb{E}_{p(\theta)p(\mathbf{x})}\left[1 - d(\theta,\mathbf{x})\right]$.

# Summary

Simulation-based inference is a major evolution in the statistical capabilities for science, as it enables the analysis of complex models and data without simplifying assumptions.

Obstacles remain to be overcome, such as the curse of dimensionality, the need for large amounts of data, or the necessary robustness of the inference network.