

The frontier of simulation-based inference

DataLearning seminars, Imperial College London
February 8, 2022

Gilles Louppe
g.louppe@uliege.be

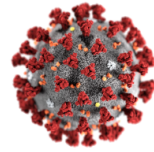
Scientific simulators



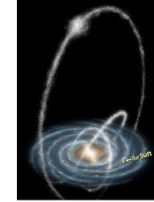
Chemical reactions



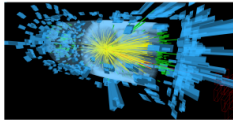
Flames



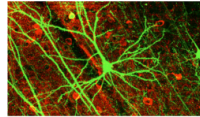
Epidemics



Stellar streams



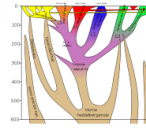
Collider experiments



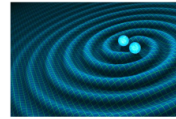
Neurons



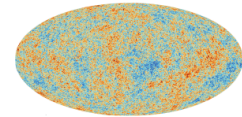
Robotics



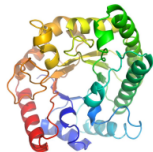
Evolution



Gravitational waves



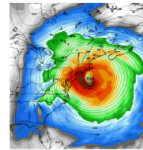
Evolution of the Universe



Protein networks



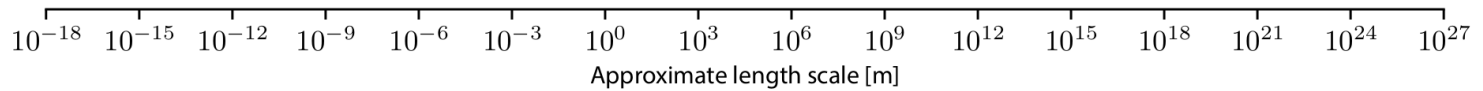
Ecological systems



Weather and climate

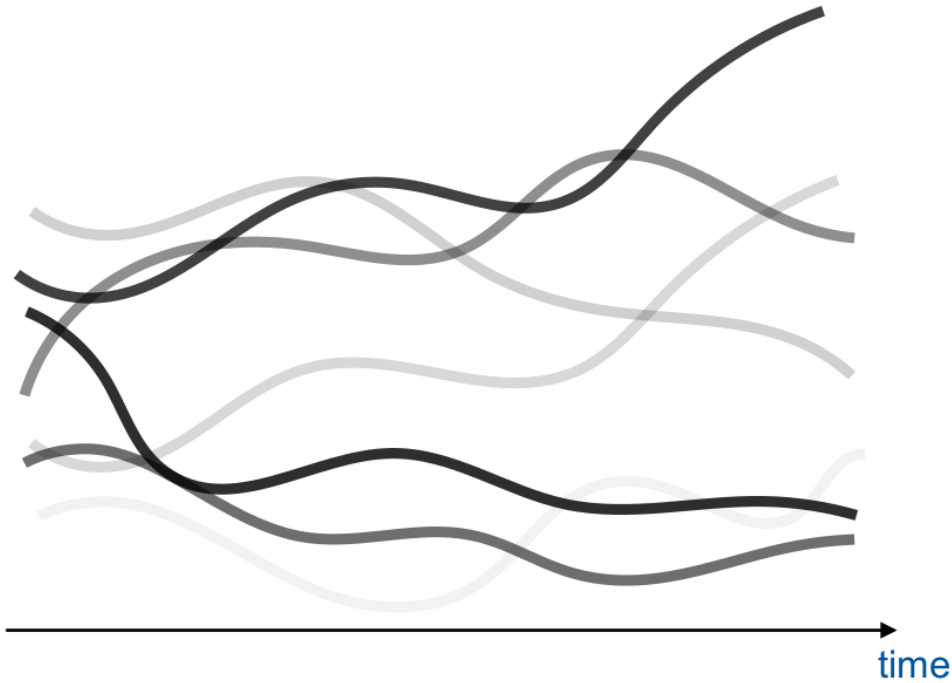


Gravitational lensing





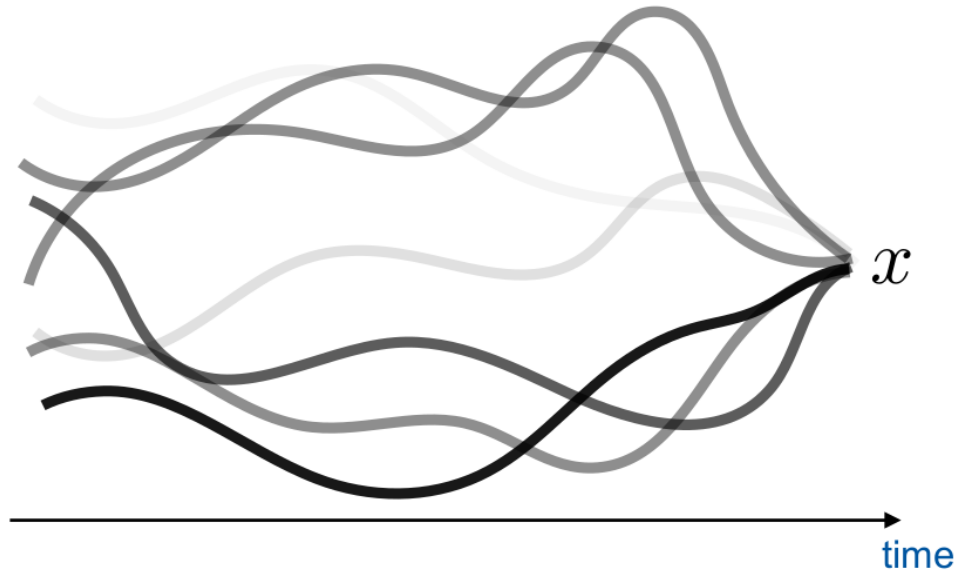
Prediction



$$\theta, z, x \sim p(\theta, z, x)$$



Inference



This results in the likelihood $p(x|\theta) = \int p(x, z|\theta)dz$ to be intractable.

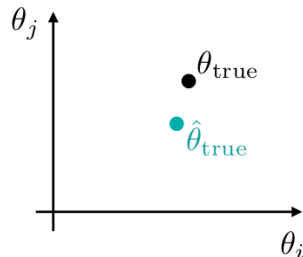
Problem statement

Start with

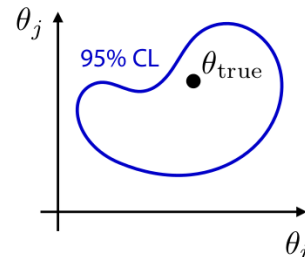
- a simulator that lets you generate N samples $x_i \sim p(x_i | \theta_i)$,
- observed data $x_{\text{obs}} \sim p(x_{\text{obs}} | \theta_{\text{true}})$,
- a prior $p(\theta)$.

Then,

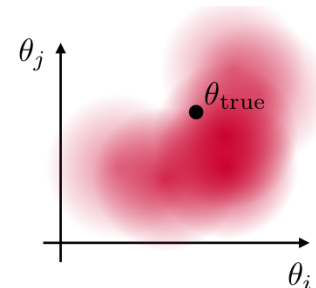
a) estimate θ_{true}
(e.g., MLE)



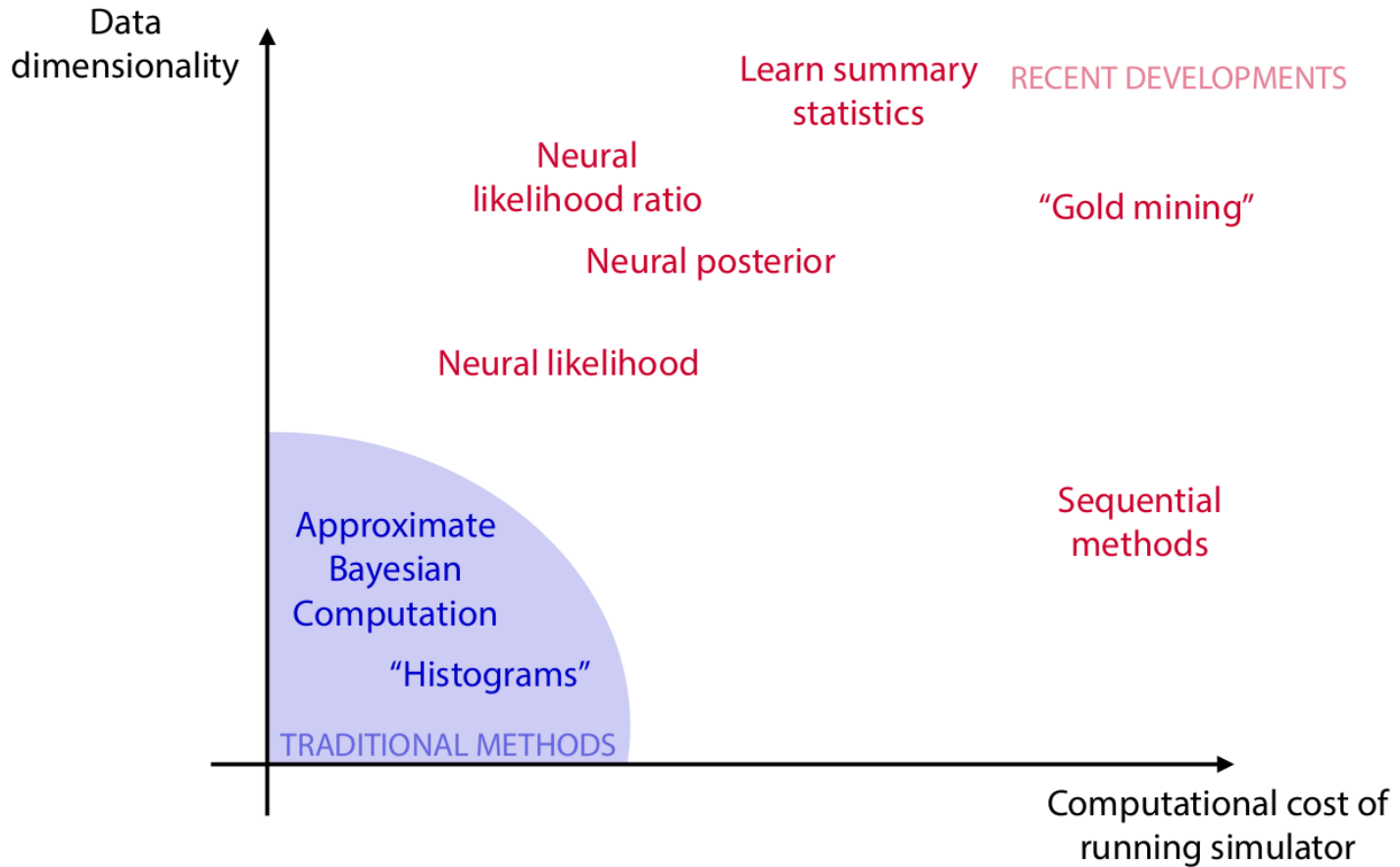
b) construct confidence sets



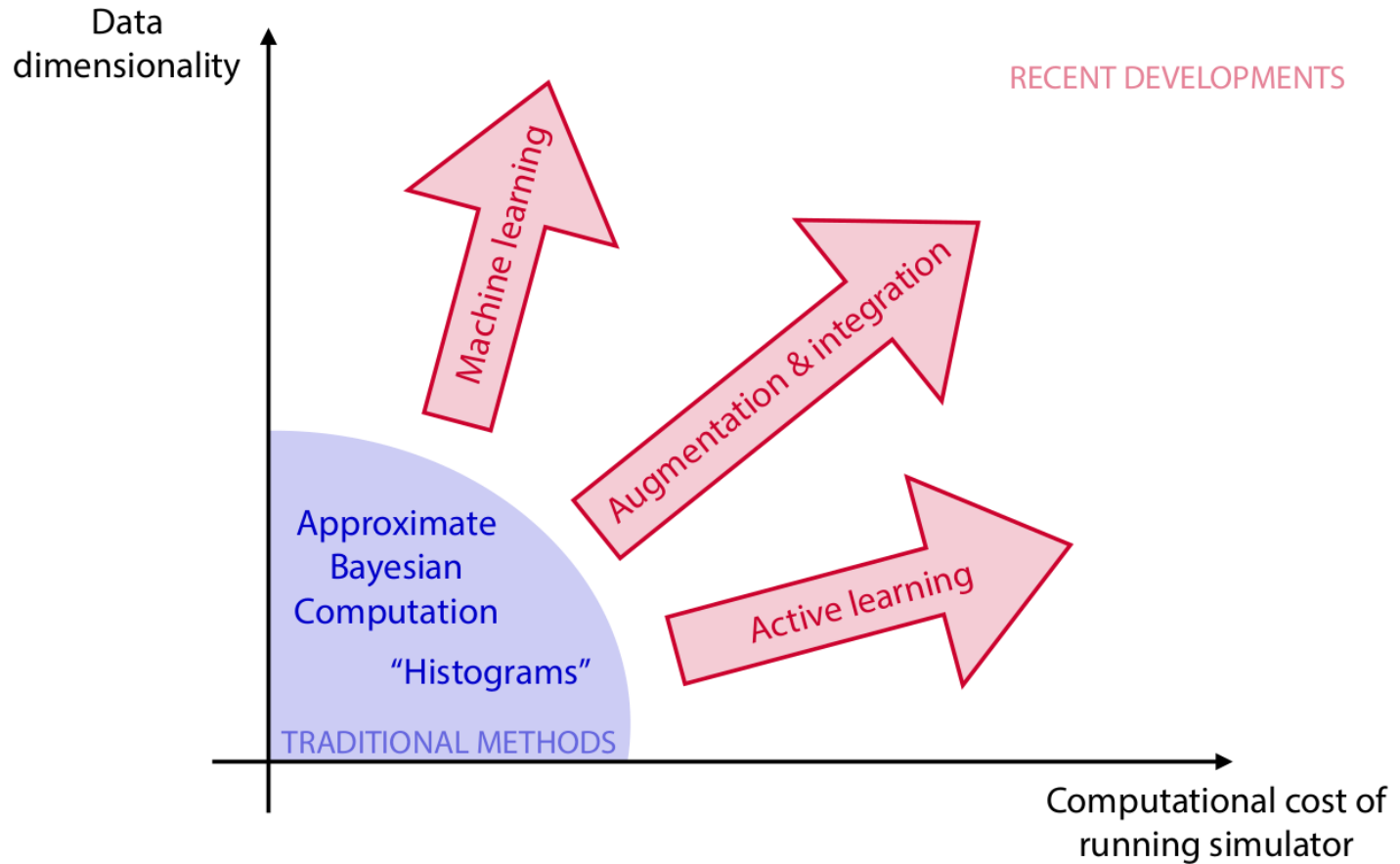
c) estimate the posterior $p(\theta | x_{\text{obs}})$



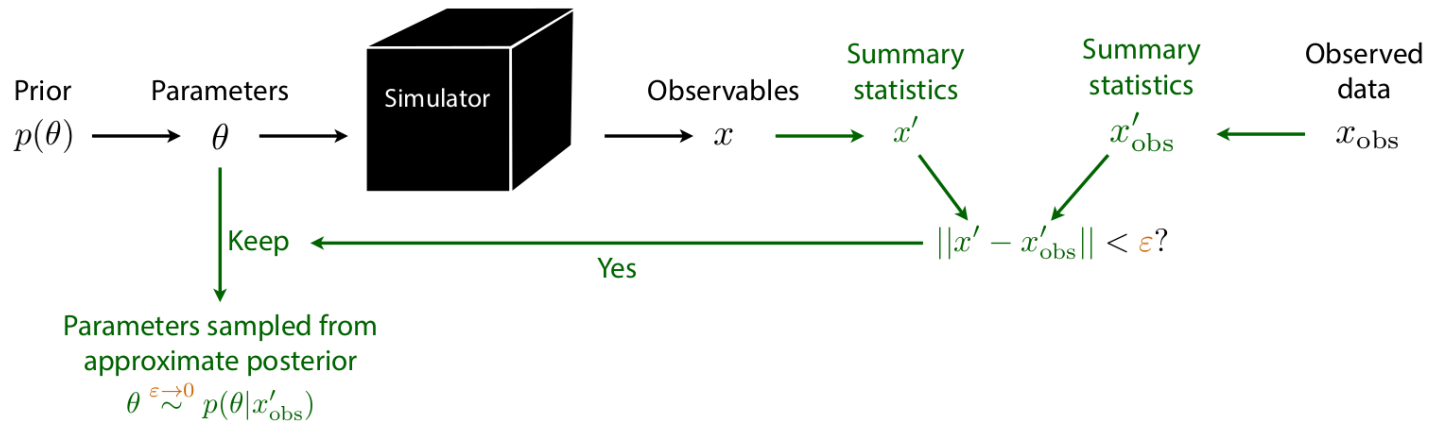
Inference algorithms



Inference algorithms



Approximate Bayesian Computation (ABC)



Issues

- How to choose x' ? ϵ ? $\| \cdot \|$?
- No tractable posterior.
- Need to run new simulations for new data or new prior.

Neural Ratio Estimation (NRE)

The Bayes rule can be rewritten as

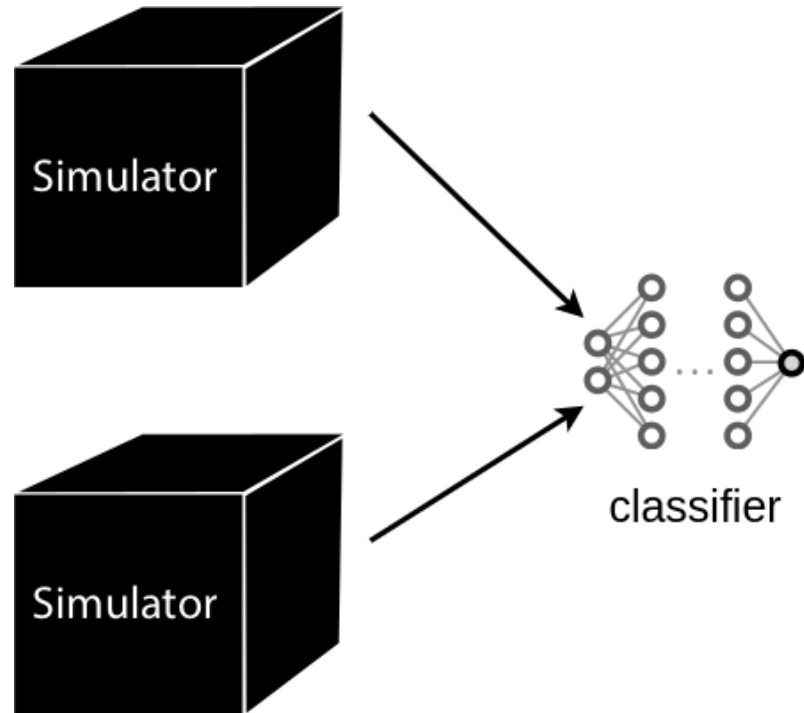
$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)} = r(x|\theta)p(\theta) \approx \hat{r}(x|\theta)p(\theta),$$

where $r(x|\theta) = \frac{p(x|\theta)}{p(x)}$ is the likelihood-to-evidence ratio.

The ratio can be learned with machine learning, even neither the likelihood nor the evidence can be evaluated!

The likelihood ratio trick

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



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

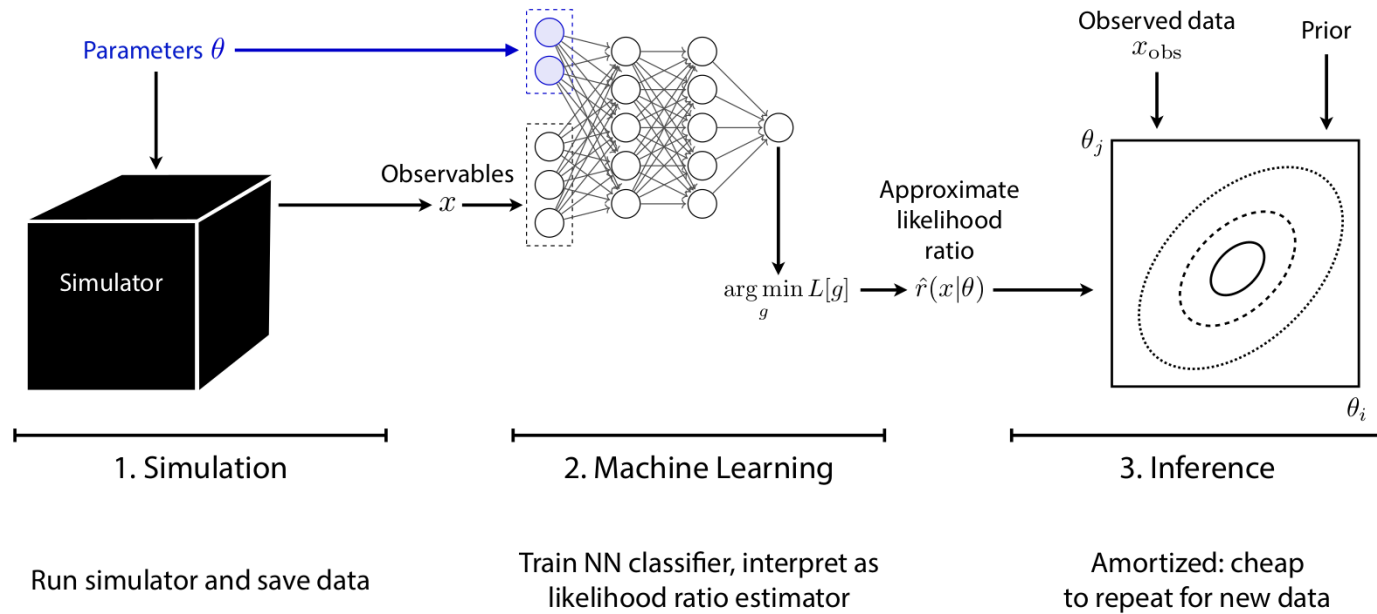
The solution d found after training approximates the optimal classifier

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

Therefore,

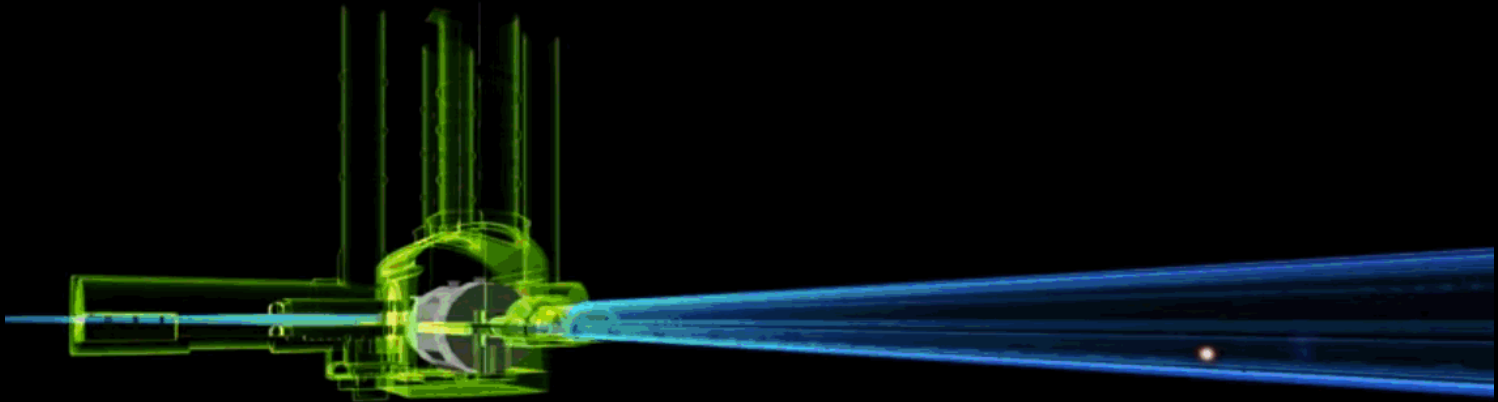
$$r(x|\theta) = \frac{p(x|\theta)}{p(x)} = \frac{p(x, \theta)}{p(x)p(\theta)} \approx \frac{d(x, \theta)}{1 - d(x, \theta)} = \hat{r}(x|\theta).$$

Inference

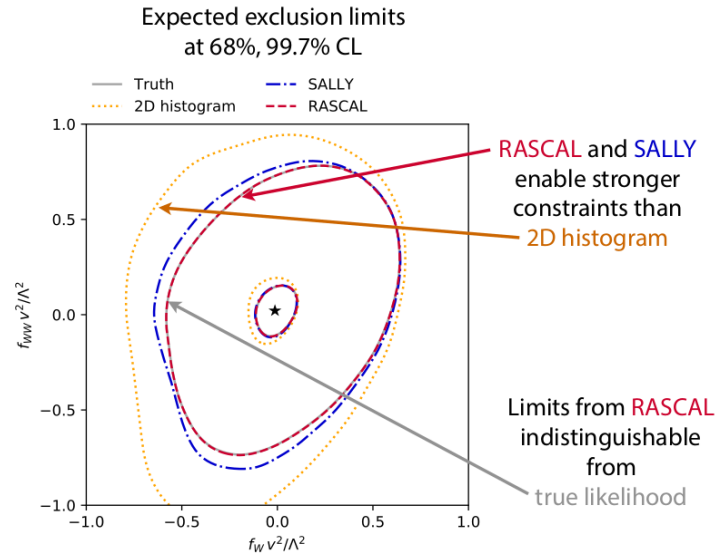
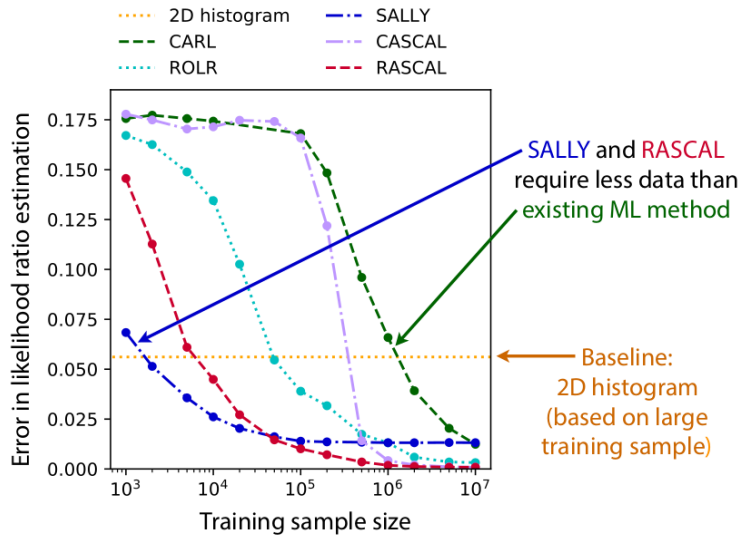


Showtime!

Some applications of physics and astrophysics.



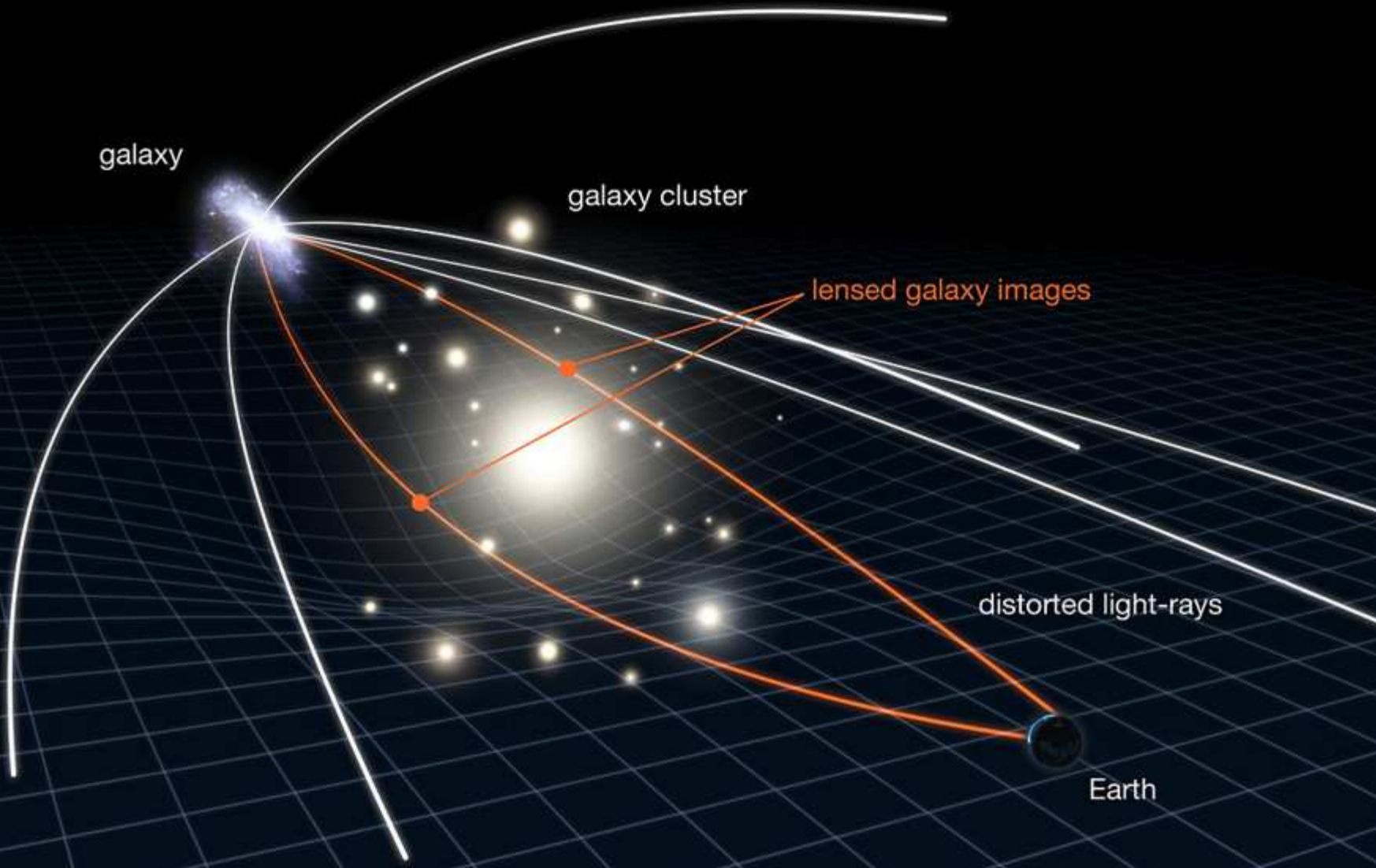
Case 1: Hunting new physics at particle colliders



With enough training data, NRE gets the likelihood-ratio statistic right.

Using more information from the simulator improves sample efficiency substantially.

Case 2: Dark matter substructure from gravitational lensing



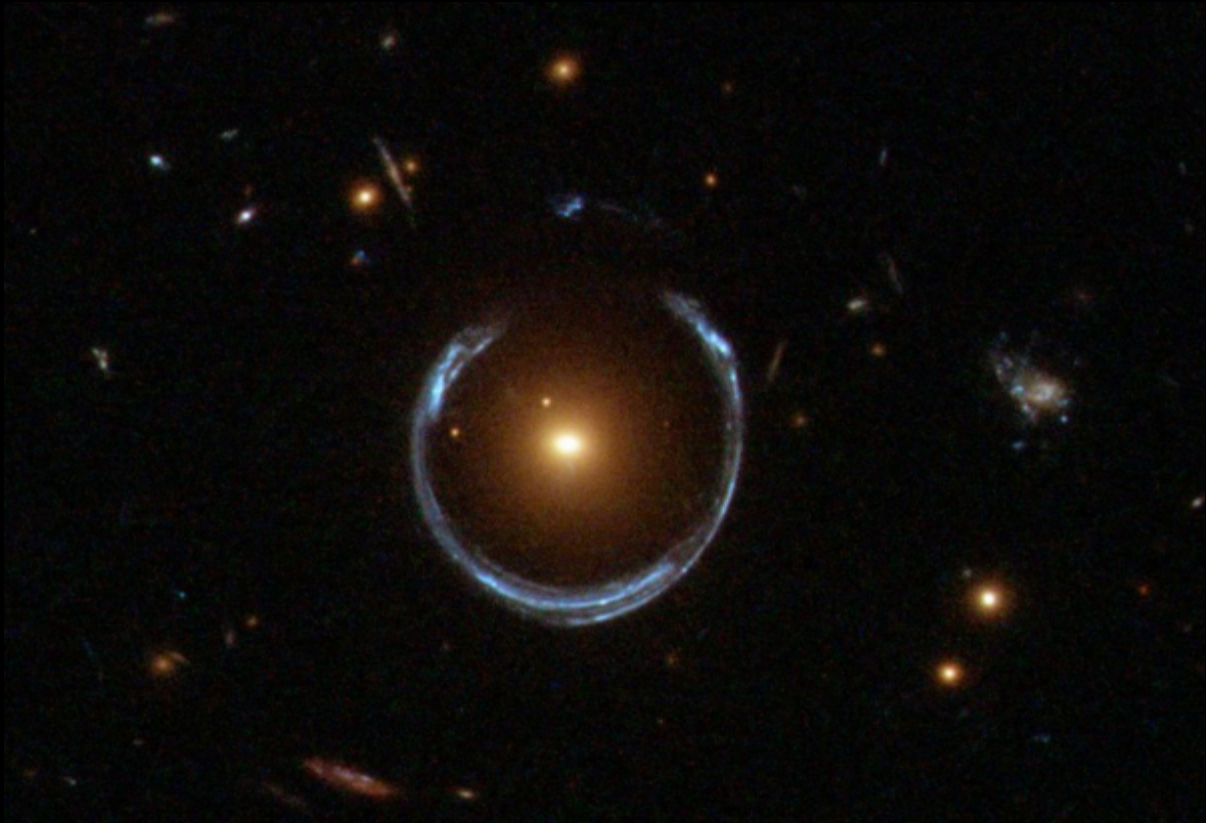
galaxy

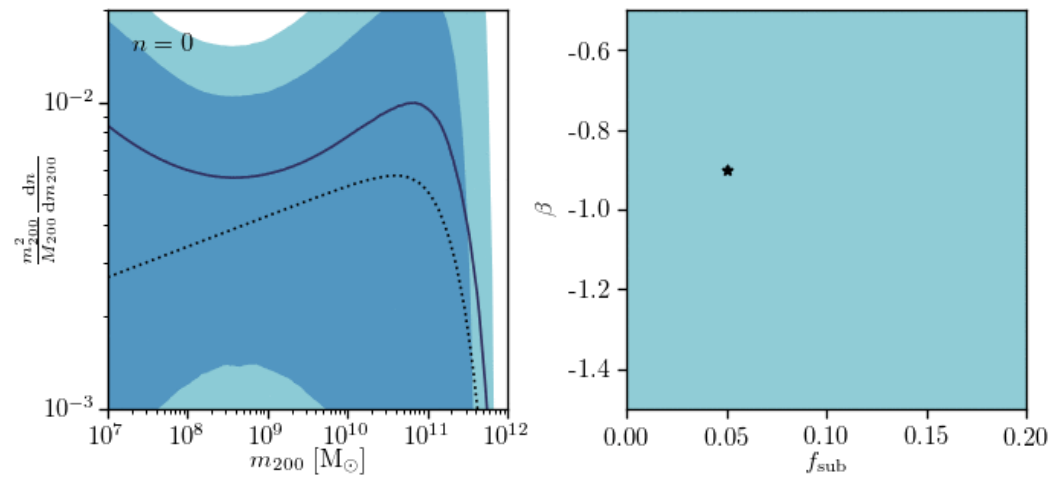
galaxy cluster

lensed galaxy images

distorted light-rays

Earth





Case 3: Constraining dark matter with stellar streams

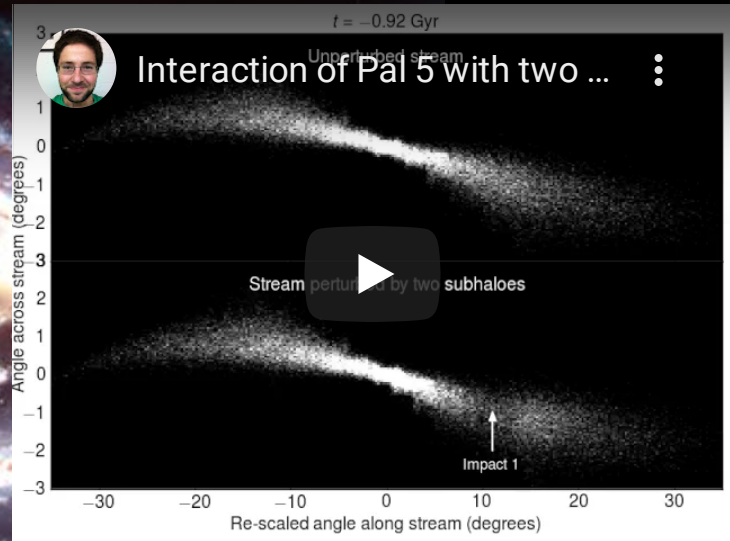
Palomar 5 (Pal5) stream
Pal5 was discovered in 2001 as the first thin stream formed from a globular cluster. Its current orbit takes it far over the galactic center.

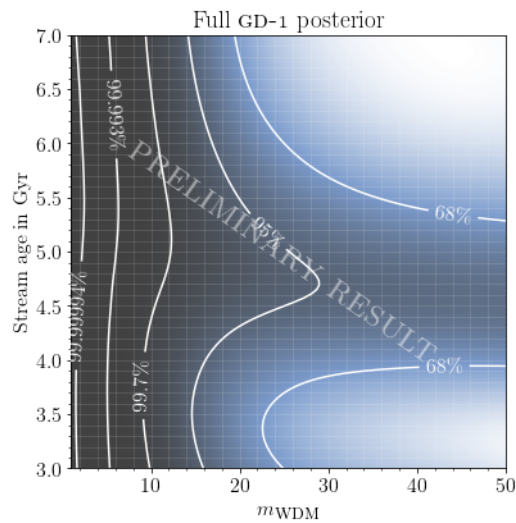
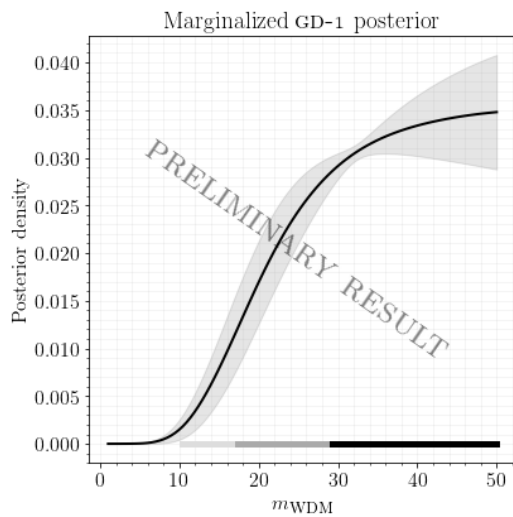
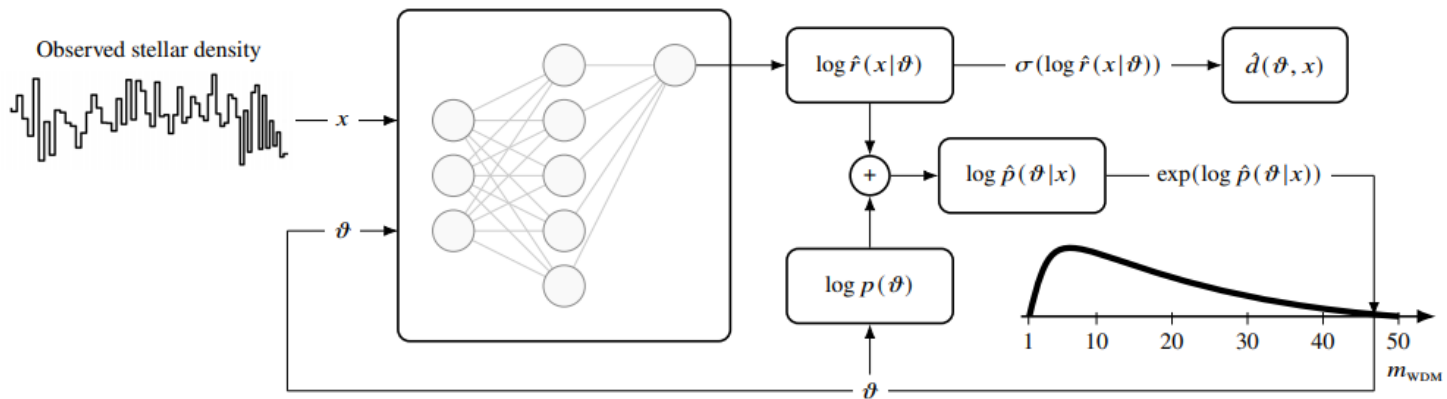
Globular clusters
These hives typically hold 100,000 stars or fewer and give rise to long, thin streams.

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.

Milky Way





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

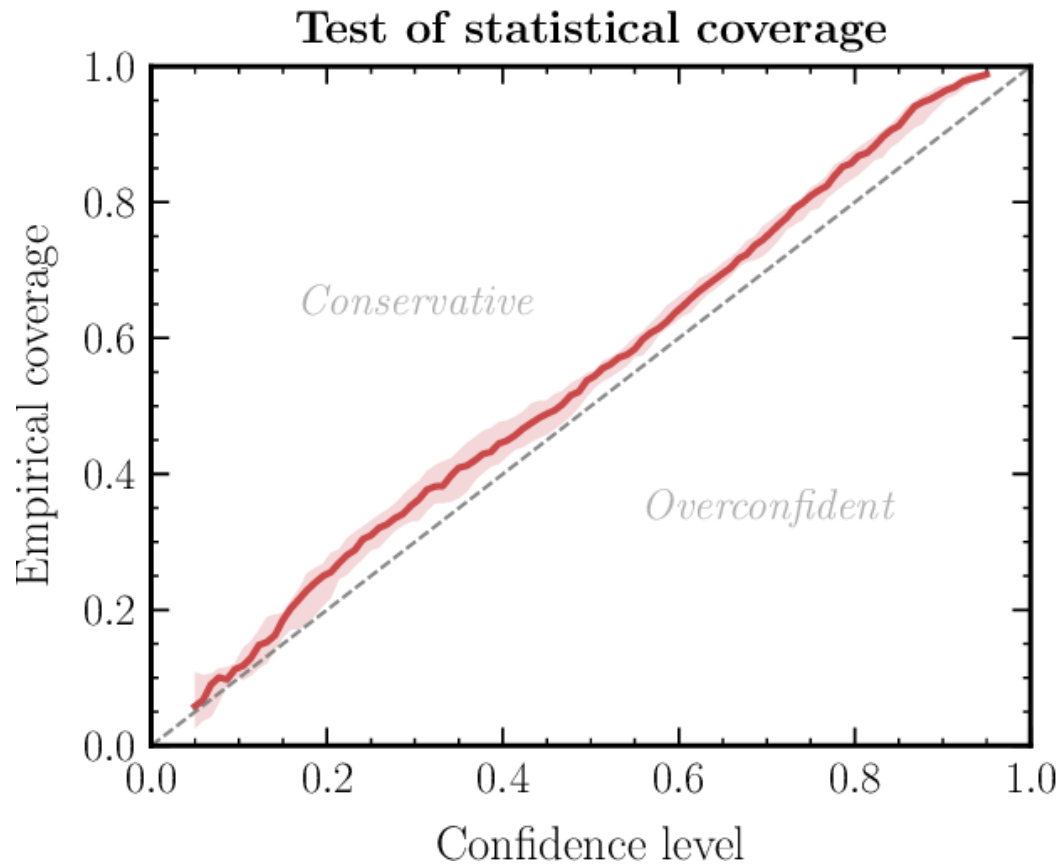
Diagnosing inference

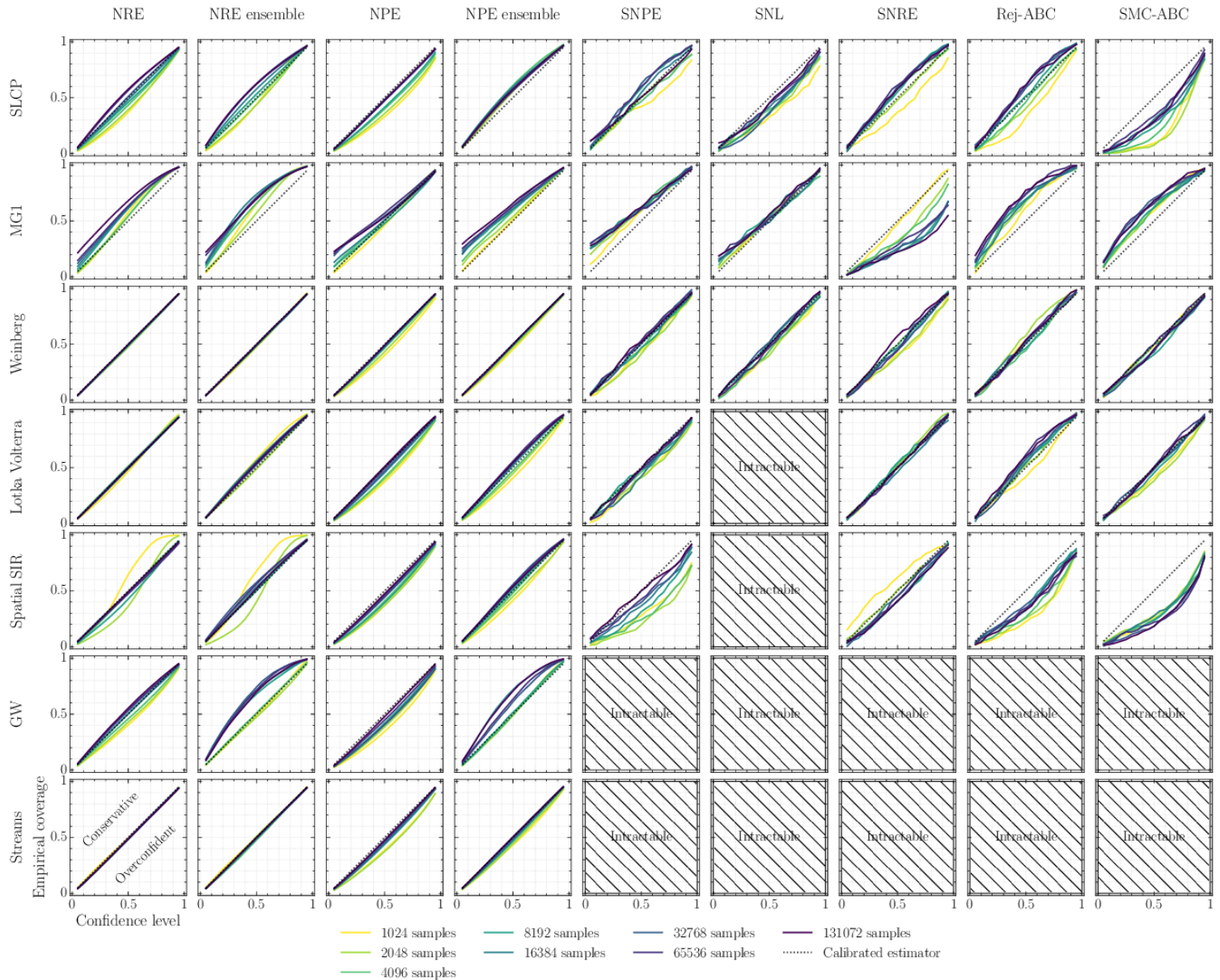
How to assess that approximate posteriors are not too wrong?

Coverage

- For every $x, \theta \sim p(x, \theta)$ in a validation set, compute the $1 - \alpha$ credible interval based on $\hat{p}(\theta|x) = \hat{r}(x|\theta)p(\theta)$.
- The fraction of samples for which θ is contained within the interval corresponds to the empirical coverage probability.

If the empirical coverage is larger than the nominal coverage probability $1 - \alpha$, then the ratio estimator \hat{r} passes the diagnostic.





All benchmarked algorithms can produce non-conservative posterior approximations.

~~The frontier of simulation-based inference~~

~~Averting a crisis in simulation- based inference?~~

DataLearning seminars, Imperial College London
February 8, 2022

Gilles Louppe
g.louppe@uliege.be

Summary

- Much of modern science is based on simulators making precise predictions, but in which inference is challenging.
- Machine learning enables powerful inference methods, which work in problems from the smallest to the largest scales.
- Advances in simulation-based inference will translate into scientific progress.
- However, further work is needed to make these methods more robust and reliable.

The end.