

Advanced Machine Learning

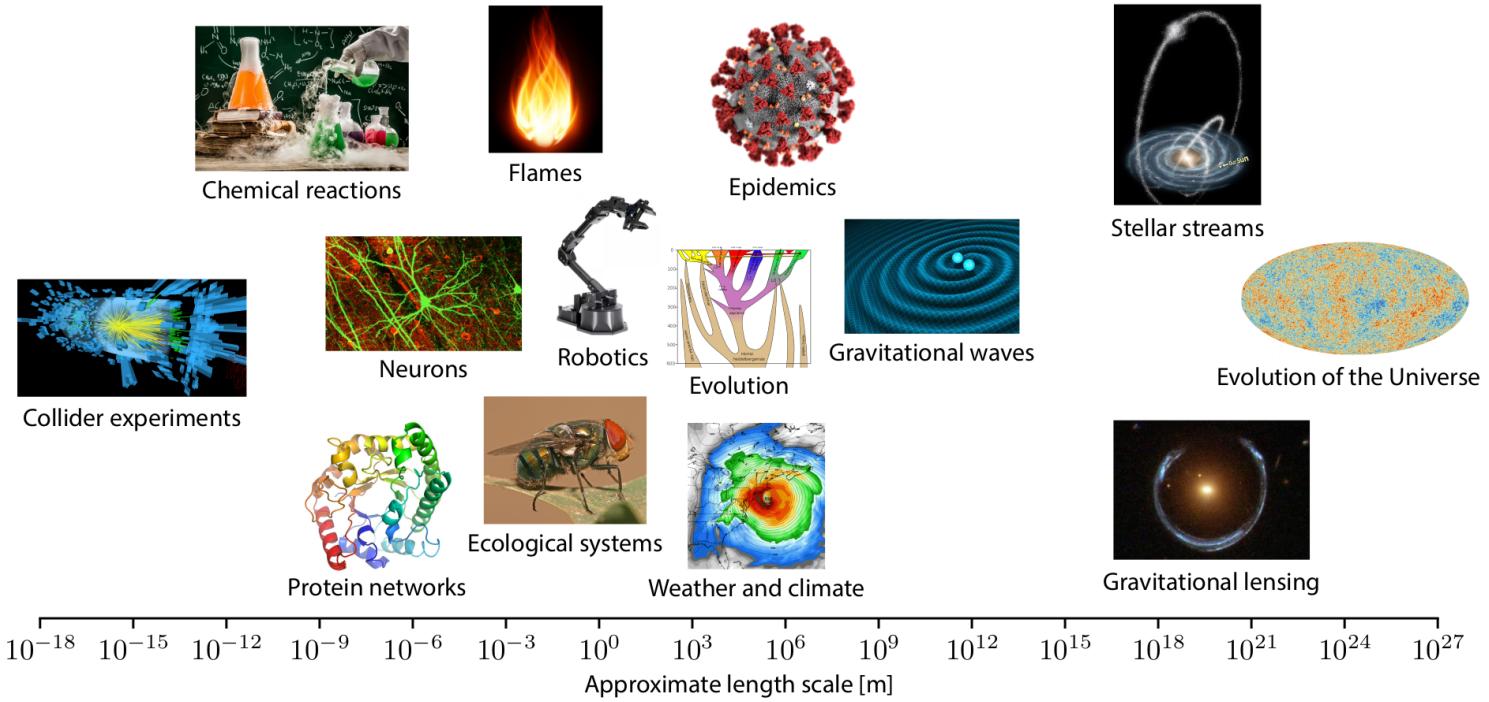
Paper: Cranmer, Brehmer, and Louppe, *The frontier of simulation-based inference*, 2020.

Gilles Louppe
g.louppe@uliege.be

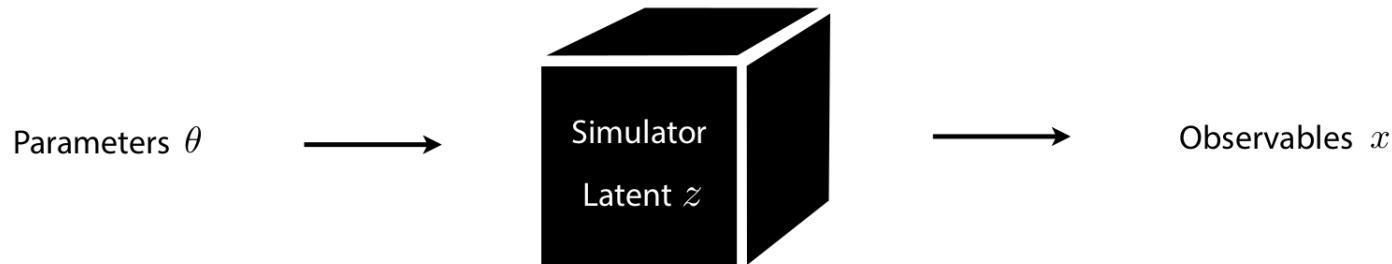


This talk is inspired and adapted from previous talks given by my wonderful co-authors [Kyle Cranmer](#) and [Johann Brehmer](#). Thanks to them!



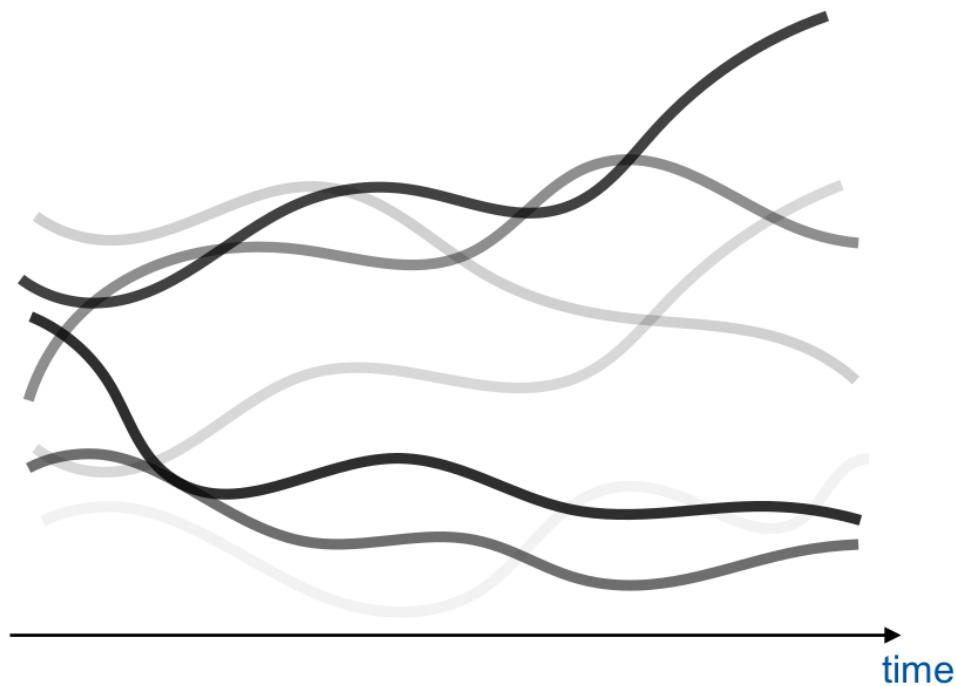


Simulation-based inference

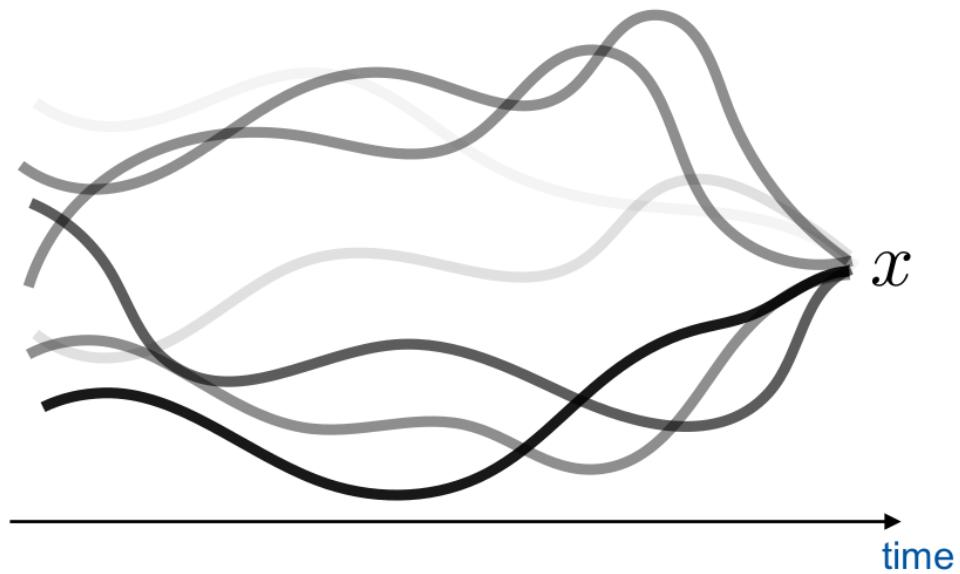


- Prediction:
- Well-motivated mechanistic, causal model
 - Simulator can generate samples $x \sim p(x|\theta)$

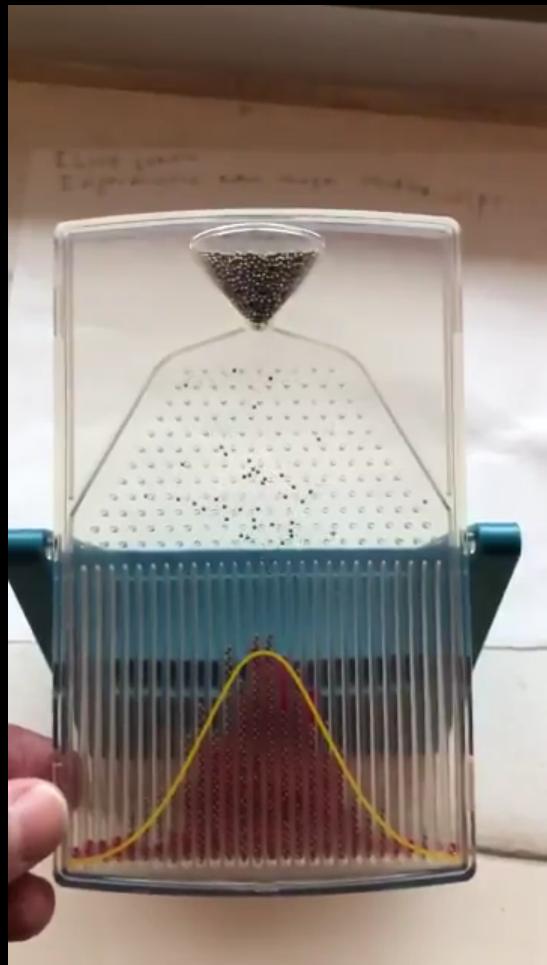
- Inference:
- Interactions between low-level components lead to challenging inverse problems
 - Likelihood $p(x|\theta) = \int dz p(x, z|\theta)$ is intractable



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



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

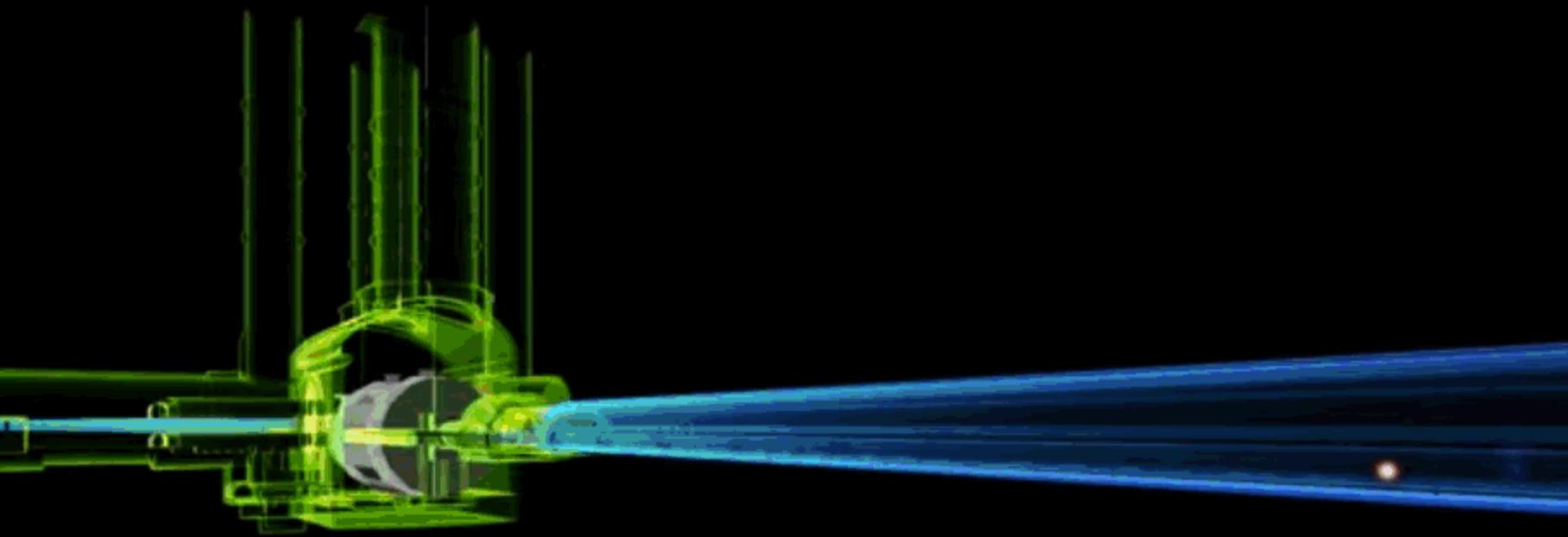


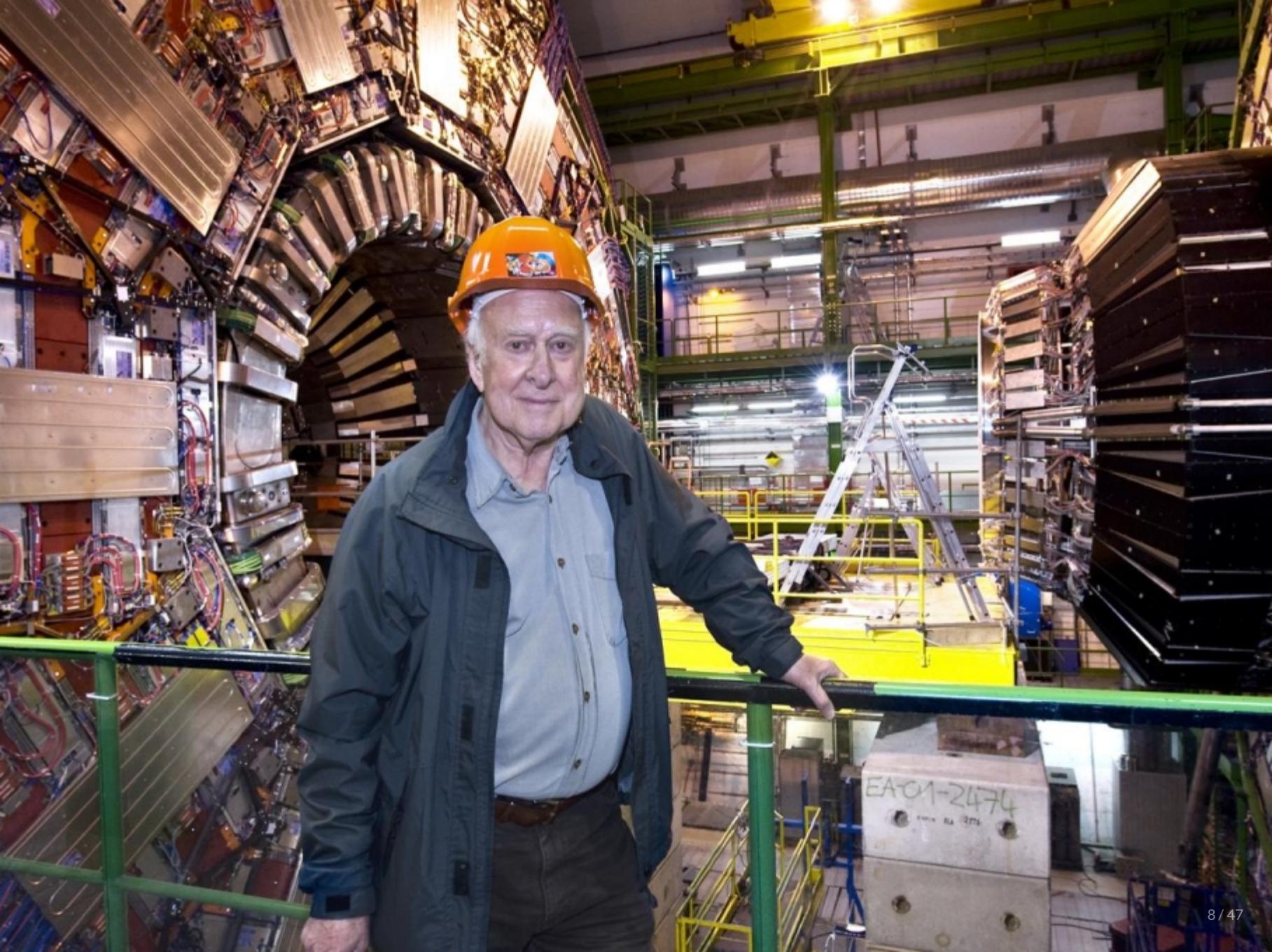
The Bean machine is a **metaphore** of **simulation-based science**:

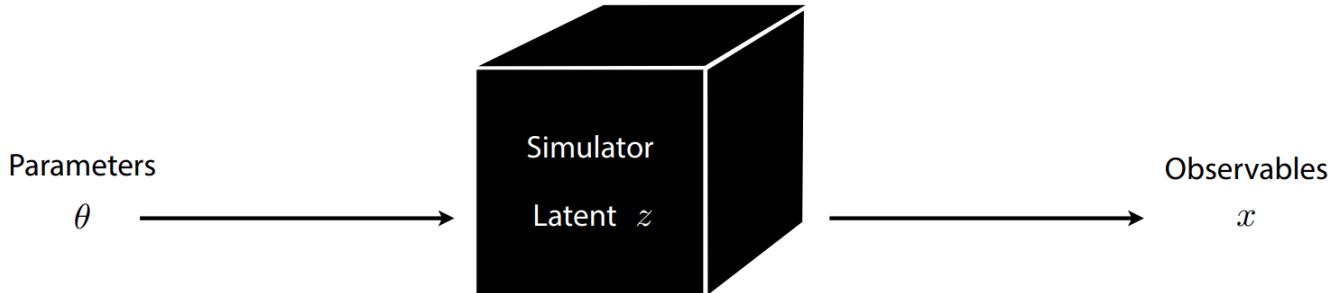
Bean machine	→	Computer simulation
Parameters θ	→	Model parameters θ
Buckets x	→	Observables x
Random paths z	→	Latent variables z (stochastic execution traces through simulator)

The case of particle physics

$$\begin{aligned}
\mathcal{L}_{SM} = & -\frac{1}{2}\partial_\mu g_\mu^a \partial_\nu g_\mu^a - g_s f^{abc} \partial_\mu g_\nu^a g_\mu^b g_\nu^c - \frac{1}{4}g_s^2 f^{abc} f^{acd} g_\mu^b g_\mu^c g_\mu^d g_\nu^e - \partial_\nu W_\mu^+ \partial_\nu W_\mu^- - \\
& M^2 W_\mu^+ W_\mu^- - \frac{1}{2}\partial_\nu Z_\mu^0 \partial_\nu Z_\mu^0 - \frac{1}{2c_w^2} M^2 Z_\mu^0 Z_\mu^0 - \frac{1}{2}\partial_\mu A_\nu \partial_\mu A_\nu - ig s_w (\partial_\nu Z_\mu^0 (W_\mu^+ W_\nu^- - \\
& W_\nu^+ W_\mu^-) - Z_\mu^0 (W_\mu^+ \partial_\nu W_\mu^- - W_\nu^+ \partial_\nu W_\mu^+) + Z_\mu^0 (W_\mu^+ \partial_\nu W_\mu^- - W_\nu^- \partial_\nu W_\mu^+)) - \\
& ig s_w (\partial_\nu A_\mu (W_\mu^+ W_\nu^- - W_\nu^+ W_\mu^-) - A_\nu (W_\mu^+ \partial_\nu W_\mu^- - W_\nu^- \partial_\nu W_\mu^+) + A_\mu (W_\nu^+ \partial_\nu W_\mu^- - \\
& W_\nu^- \partial_\nu W_\mu^+)) - \frac{1}{2}g^2 W_\mu^+ W_\mu^- W_\nu^+ W_\nu^- + \frac{1}{2}g^2 W_\mu^+ W_\nu^+ W_\mu^- W_\nu^- + g^2 c_w^2 (Z_\mu^0 W_\mu^+ Z_\nu^0 W_\nu^- - \\
& Z_\mu^0 Z_\nu^0 W_\mu^+ W_\nu^-) + g^2 s_w^2 (A_\mu W_\mu^+ A_\nu W_\nu^- - A_\mu A_\nu W_\mu^+ W_\nu^-) + g^2 s_w c_w (A_\mu Z_\nu^0 (W_\mu^+ W_\nu^- - \\
& W_\nu^+ W_\mu^-) - 2A_\mu Z_\mu^0 W_\nu^+ W_\nu^-) - \frac{1}{2}\partial_\mu H \partial_\mu H - 2M^2 \alpha_h H^2 - \partial_\mu \phi^+ \partial_\mu \phi^- - \frac{1}{2}\partial_\mu \phi^0 \partial_\mu \phi^0 - \\
& \beta_h \left(\frac{2M^2}{g^2} + \frac{2M}{g} H + \frac{1}{2}(H^2 + \phi^0 \phi^0 + 2\phi^+ \phi^-) \right) + \frac{2M^4}{g^2} \alpha_h - \\
& g \alpha_h M (H^3 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-) - \\
& \frac{1}{8}g^2 \alpha_h (H^4 + (\phi^0)^4 + 4(\phi^+ \phi^-)^2 + 4(\phi^0)^2 \phi^+ \phi^- + 4H^2 \phi^+ \phi^- + 2(\phi^0)^2 H^2) - \\
& g M W_\mu^+ W_\mu^- H - \frac{1}{2}g \frac{M}{c_w^2} Z_\mu^0 Z_\mu^0 H - \\
& \frac{1}{2}ig (W_\mu^+ (\phi^0 \partial_\mu \phi^- - \phi^- \partial_\mu \phi^0) - W_\mu^- (\phi^0 \partial_\mu \phi^+ - \phi^+ \partial_\mu \phi^0)) + \\
& \frac{1}{2}g (W_\mu^+ (H \partial_\mu \phi^- - \phi^- \partial_\mu H) + W_\mu^- (H \partial_\mu \phi^+ - \phi^+ \partial_\mu H)) + \frac{1}{2}g \frac{1}{c_w} (Z_\mu^0 (H \partial_\mu \phi^0 - \phi^0 \partial_\mu H) + \\
& M (\frac{1}{c_w} Z_\mu^0 \partial_\mu \phi^0 + W_\mu^+ \partial_\mu \phi^- + W_\mu^- \partial_\mu \phi^+) - ig \frac{s_w^2}{c_w} M Z_\mu^0 (W_\mu^+ \phi^- - W_\mu^- \phi^+) + ig s_w M A_\mu (W_\mu^+ \phi^- - \\
& W_\mu^- \phi^+) - ig \frac{1-2c_w^2}{2c_w} Z_\mu^0 (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) + ig s_w A_\mu (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) - \\
& \frac{1}{4}g^2 W_\mu^+ W_\mu^- (H^2 + (\phi^0)^2 + 2\phi^+ \phi^-) - \frac{1}{8}g^2 \frac{1}{c_w^2} Z_\mu^0 Z_\mu^0 (H^2 + (\phi^0)^2 + 2(2s_w^2 - 1)^2 \phi^+ \phi^-) - \\
& \frac{1}{2}g^2 \frac{s_w^2}{c_w} Z_\mu^0 \phi^0 (W_\mu^+ \phi^- - W_\mu^- \phi^+) - \frac{1}{2}ig \frac{s_w^2}{c_w} Z_\mu^0 H (W_\mu^+ \phi^- - W_\mu^- \phi^+) + \frac{1}{2}g^2 s_w A_\mu \phi^0 (W_\mu^+ \phi^- + \\
& W_\mu^- \phi^+) + \frac{1}{2}ig^2 s_w A_\mu H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - g^2 \frac{s_w}{c_w} (2c_w^2 - 1) Z_\mu^0 A_\mu \phi^+ \phi^- - \\
& g^2 s_w^2 A_\mu A_\mu \phi^+ \phi^- + \frac{1}{2}ig_s \lambda_{ij}^a (\bar{q}_i^a \gamma^\mu q_j^a) g_\mu^a - \bar{e}^\lambda (\gamma \partial + m_e^\lambda) e^\lambda - \bar{\nu}^\lambda (\gamma \partial + m_\nu^\lambda) \nu^\lambda - \bar{u}^\lambda (\gamma \partial + \\
& m_u^\lambda) u^\lambda - \bar{d}_j^\lambda (\gamma \partial + m_d^\lambda) d^\lambda_j + ig s_w A_\mu ((-\bar{e}^\lambda \gamma^\mu e^\lambda) + \frac{2}{3}(\bar{u}_j^\lambda \gamma^\mu u_j^\lambda) - \frac{1}{3}(\bar{d}_j^\lambda \gamma^\mu d_j^\lambda)) + \\
& \frac{ig}{4c_w} Z_\mu^0 \{(\bar{\nu}^\lambda \gamma^\mu (1 + \gamma^5) \nu^\lambda) + (\bar{e}^\lambda \gamma^\mu (4s_w^2 - 1 - \gamma^5) e^\lambda) + (\bar{d}_j^\lambda \gamma^\mu (\frac{4}{3}s_w^2 - 1 - \gamma^5) d_j^\lambda)\} + \\
& (\bar{u}_j^\lambda \gamma^\mu (1 - \frac{8}{3}s_w^2 + \gamma^5) u_j^\lambda) \} + \frac{ig}{2\sqrt{2}} W_\mu^+ ((\bar{\nu}^\lambda \gamma^\mu (1 + \gamma^5) U^{lep} \lambda_\kappa e^\kappa) + (\bar{u}_j^\lambda \gamma^\mu (1 + \gamma^5) C_{\lambda\kappa} d_j^\kappa)) + \\
& \frac{ig}{2\sqrt{2}} W_\mu^- ((\bar{e}^\kappa U^{lep\dagger} \lambda_\lambda \gamma^\mu (1 + \gamma^5) \nu^\lambda) + (\bar{d}_j^\kappa C_{\lambda\lambda}^\dagger \gamma^\mu (1 + \gamma^5) u_j^\lambda)) + \\
& \frac{ig}{2M\sqrt{2}} \phi^+ (-m_e^\kappa (\bar{\nu}^\lambda U^{lep} \lambda_\kappa (1 - \gamma^5) e^\kappa) + m_\nu^\kappa (\bar{\nu}^\lambda U^{lep} \lambda_\kappa (1 + \gamma^5) e^\kappa) + \\
& \frac{ig}{2M\sqrt{2}} \phi^- (m_e^\lambda (\bar{e}^\lambda U^{lep\dagger} \lambda_\kappa (1 + \gamma^5) \nu^\kappa) - m_\nu^\kappa (\bar{e}^\lambda U^{lep\dagger} \lambda_\kappa (1 - \gamma^5) \nu^\kappa)) - \frac{g}{2} \frac{m_e^\lambda}{M} H (\bar{\nu}^\lambda \nu^\lambda) - \\
& \frac{g}{2} \frac{m_\nu^\lambda}{M} H (\bar{e}^\lambda e^\lambda) + \frac{ig}{2} \frac{m_e^\lambda}{M} \phi^0 (\bar{\nu}^\lambda \gamma^5 \nu^\lambda) - \frac{ig}{2} \frac{m_\nu^\lambda}{M} \phi^0 (\bar{e}^\lambda \gamma^5 e^\lambda) - \frac{1}{4} \bar{\nu}_\lambda M_{\lambda\kappa}^R (1 - \gamma_5) \bar{\nu}_\kappa - \\
& \frac{1}{4} \bar{\nu}_\lambda M_{\lambda\kappa}^R (1 - \gamma_5) \bar{\nu}_\kappa + \frac{ig}{2M\sqrt{2}} \phi^+ (-m_d^\kappa (\bar{u}_j^\lambda C_{\lambda\kappa} (1 - \gamma^5) d_j^\kappa) + m_u^\lambda (\bar{u}_j^\lambda C_{\lambda\kappa} (1 + \gamma^5) d_j^\kappa) + \\
& \frac{ig}{2M\sqrt{2}} \phi^- (m_d^\lambda (\bar{d}_j^\lambda C_{\lambda\kappa}^\dagger (1 + \gamma^5) u_j^\kappa) - m_u^\kappa (\bar{d}_j^\lambda C_{\lambda\kappa}^\dagger (1 - \gamma^5) u_j^\kappa)) - \frac{g}{2} \frac{m_e^\lambda}{M} H (\bar{u}_j^\lambda u_j^\lambda) - \\
& \frac{g}{2} \frac{m_\nu^\lambda}{M} H (\bar{d}_j^\lambda d_j^\lambda) + \frac{ig}{2} \frac{m_e^\lambda}{M} \phi^0 (\bar{u}_j^\lambda \gamma^5 u_j^\lambda) - \frac{ig}{2} \frac{m_\nu^\lambda}{M} \phi^0 (\bar{d}_j^\lambda \gamma^5 d_j^\lambda) + \bar{G}^a \partial^2 G^a + g_s f^{abc} \partial_\mu \bar{G}^a G^b g_c^\mu + \\
& \bar{X}^+ (\partial^2 - M^2) X^+ + \bar{X}^- (\partial^2 - M^2) X^- + \bar{X}^0 (\partial^2 - \frac{M^2}{c_w^2}) X^0 + \bar{Y} \partial^2 Y + ig c_w W_\mu^+ (\partial_\mu \bar{X}^0 X^- - \\
& \partial_\mu \bar{X}^+ X^0) + ig s_w W_\mu^+ (\partial_\mu \bar{Y} X^- - \partial_\mu \bar{X}^+ Y) + ig c_w W_\mu^- (\partial_\mu \bar{X}^- X^0 - \\
& \partial_\mu \bar{X}^0 X^+) + ig s_w W_\mu^- (\partial_\mu \bar{X}^- Y - \partial_\mu \bar{Y} X^+) + ig c_w Z_\mu^0 (\partial_\mu \bar{X}^+ X^+ - \\
& \partial_\mu \bar{X}^- X^-) + ig s_w A_\mu (\partial_\mu \bar{X}^+ X^+ - \\
& \partial_\mu \bar{X}^- X^-) - \frac{1}{2}g M \left(\bar{X}^+ X^+ H + \bar{X}^- X^- H + \frac{1}{c_w^2} \bar{X}^0 X^0 H \right) + \frac{1-2c_w^2}{2c_w} ig M (\bar{X}^+ X^0 \phi^- - \bar{X}^- X^0 \phi^-) + \\
& \frac{1}{2c_w} ig M (\bar{X}^0 X^- \phi^+ - \bar{X}^0 X^+ \phi^-) + ig M s_w (\bar{X}^0 X^- \phi^+ - \bar{X}^0 X^+ \phi^-) + \\
& \frac{1}{2}ig M (\bar{X}^+ X^+ \phi^0 - \bar{X}^- X^- \phi^0) .
\end{aligned}$$



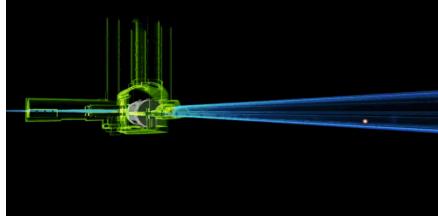




SM with parameters θ

Simulated observables \mathbf{x}

Real observations x_{obs}



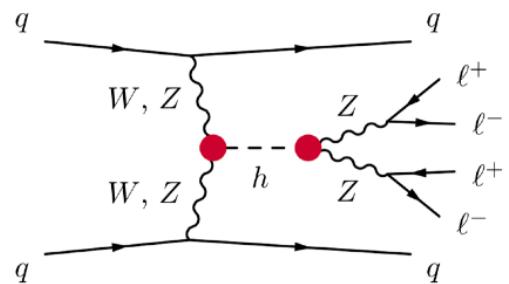
Latent variables

Parameters
of interest

Parton-level
momenta

Theory
parameters

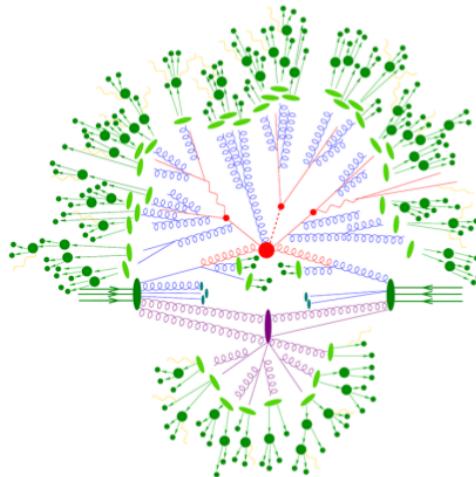
$$z_p \leftarrow \theta$$

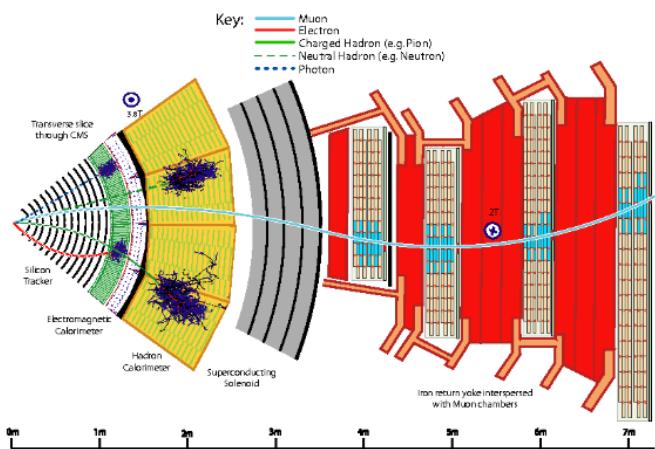
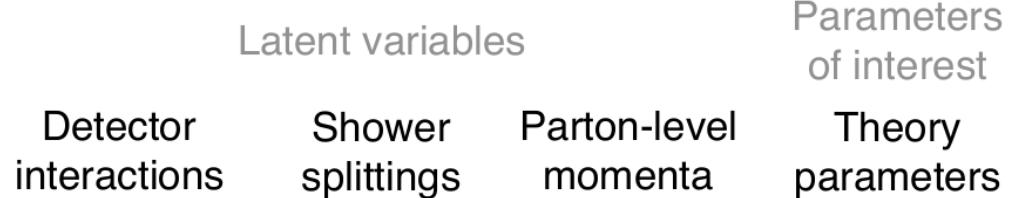


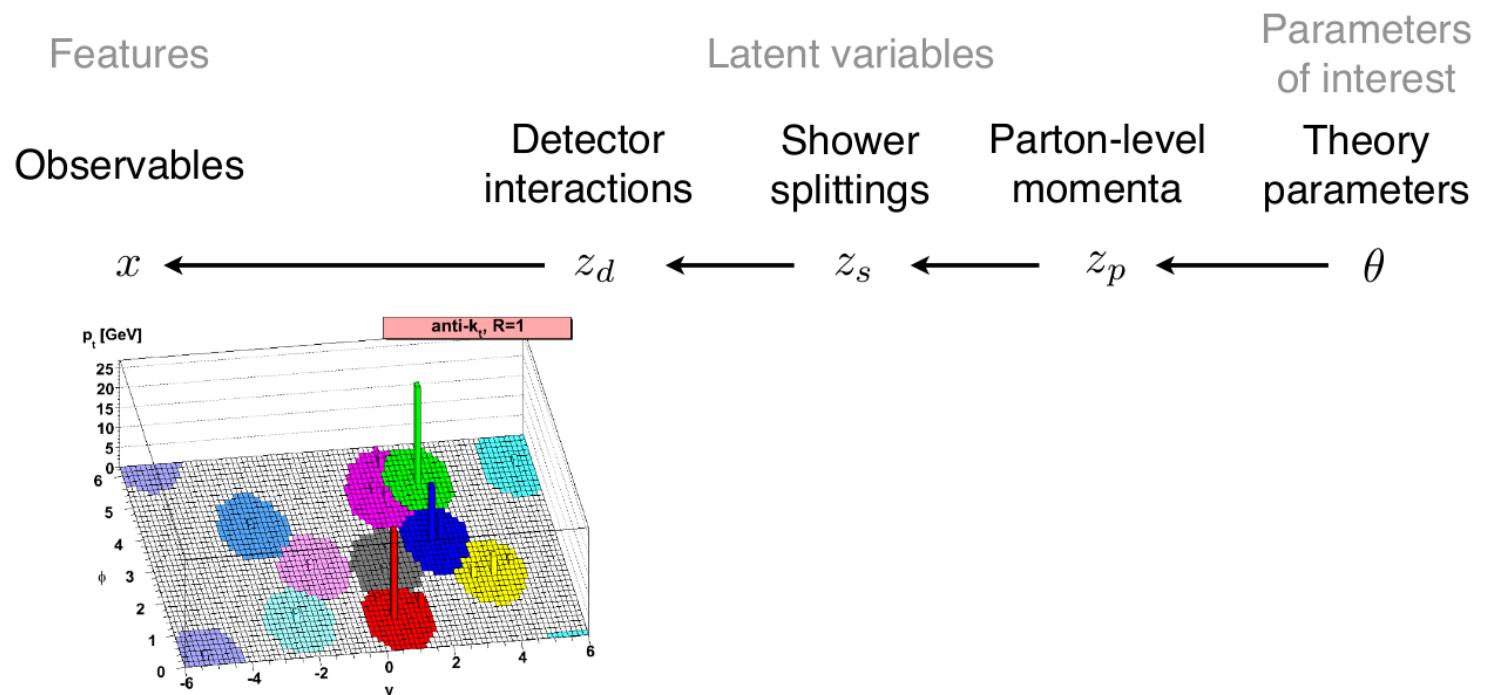
Latent variables Parameters
of interest

Shower Parton-level Theory
splittings momenta parameters

$$z_s \leftarrow z_p \leftarrow \theta$$







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

$$p(x|\theta) = \underbrace{\iiint}_{\text{yikes!}} p(z_p|\theta)p(z_s|z_p)p(z_d|z_s)p(x|z_d)dz_p dz_s dz_d$$

Inference

Problem statement(s)

Start with

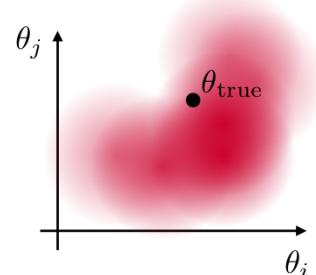
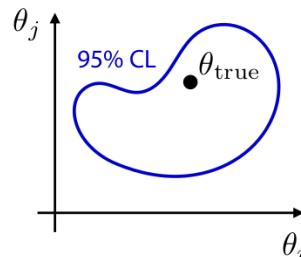
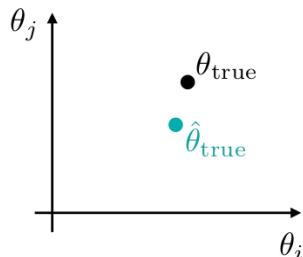
- a simulator that lets you generate N samples $x_i \sim p(x_i | \theta_i)$ (for parameters θ_i of our choice),
- observed data $x_{\text{obs}} \sim p(x_{\text{obs}} | \theta_{\text{true}})$,
- a prior $p(\theta)$.

Then,

a) estimate θ_{true}
(e.g., MLE)

b) construct confidence
sets

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



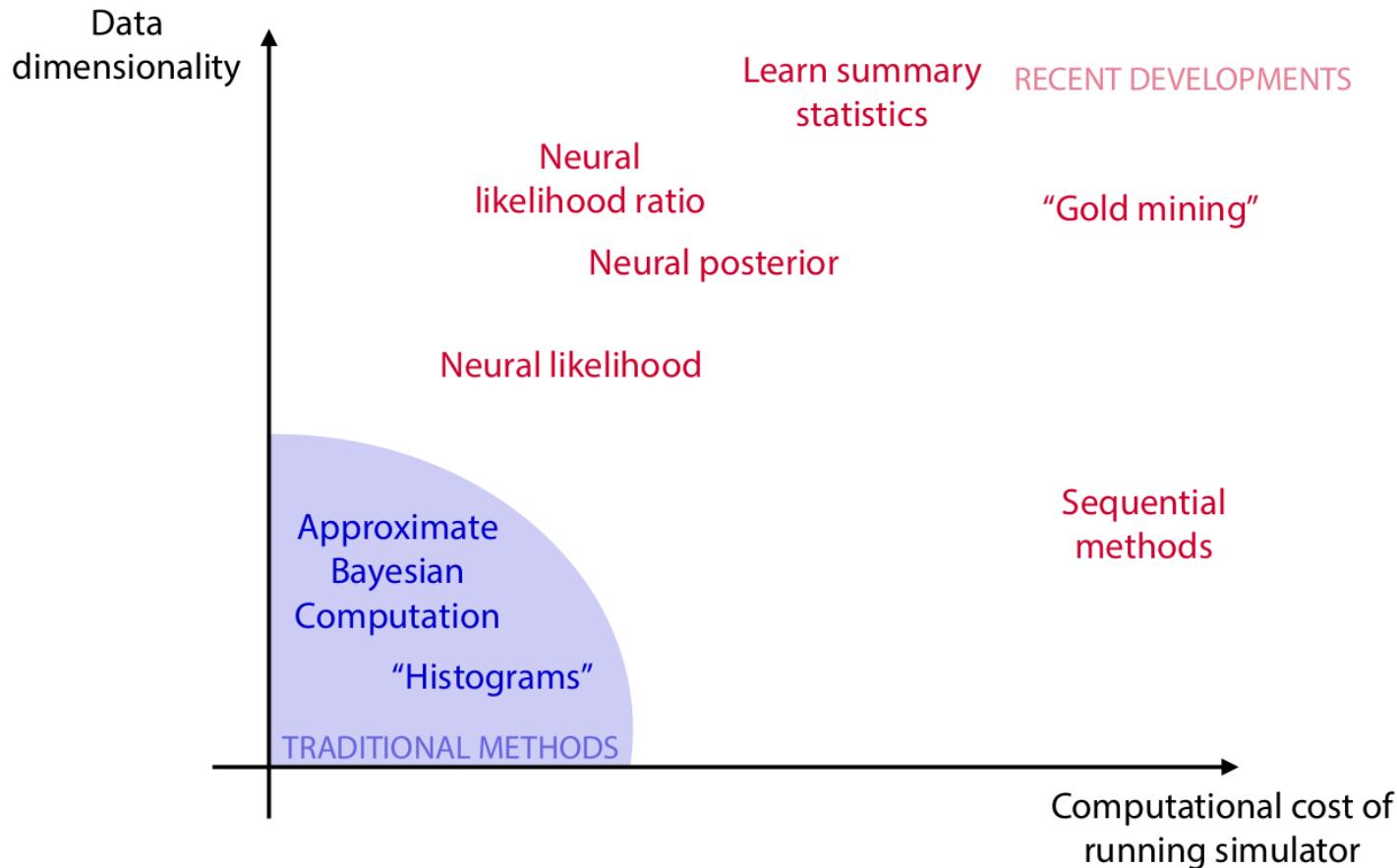
Ingredients

Statistical inference requires the computation of **key ingredients**, such as

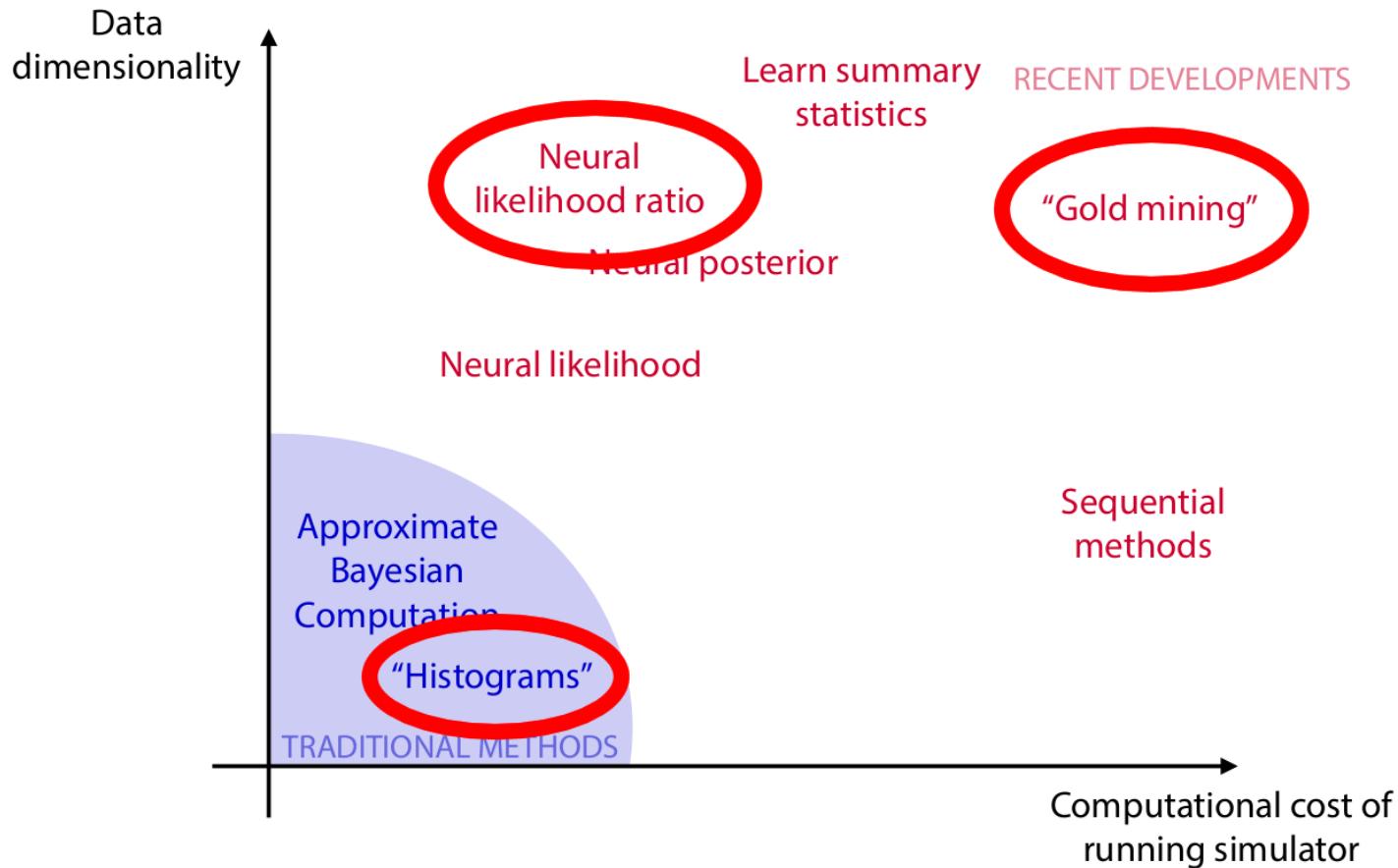
- the likelihood $p(x|\theta)$,
- the likelihood ratio $r(x|\theta_0, \theta_1) = \frac{p(x|\theta_0)}{p(x|\theta_1)}$,
- or the posterior $p(\theta|x)$,

but none are usually tractable in simulation-based science!

Inference algorithms



Inference algorithms



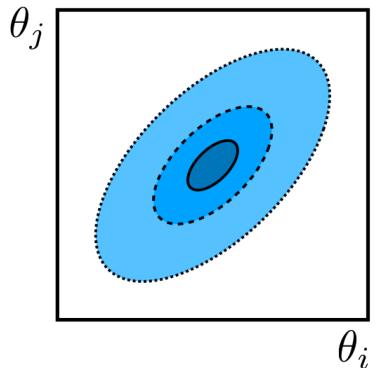
The frequentist approach

The frequentist (physicist's) way

The Neyman-Pearson lemma states that the likelihood ratio

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

is the **most powerful test statistic** to discriminate between a null hypothesis θ_0 and an alternative θ_1 .



JOURNAL OF MATHEMATICAL STATISTICS
VOLUME 11, NUMBER 1, MARCH 1933

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

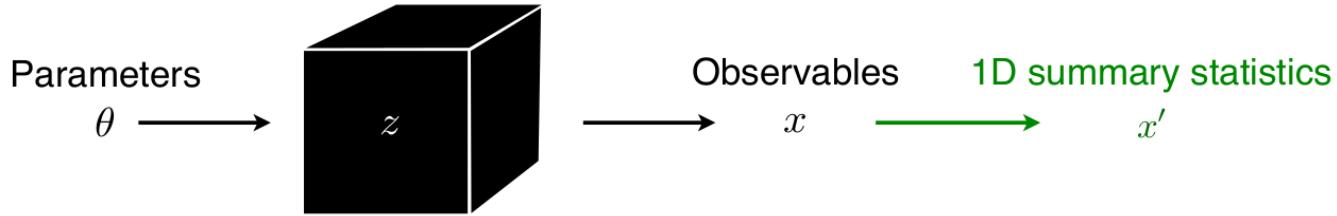
By J. NEYMAN, Nencki Institute, Soc. Sci. Lit. Varsoviensis, and Lecturer at the Central College of Agriculture, Warsaw, and E. S. PEARSON, Department of Applied Statistics, University College, London.

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

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

CONTENTS.

	PAGE.
I. Introductory	289
II. Outline of General Theory	293
III. Simple Hypotheses	293



Define a projection function $s : \mathcal{X} \rightarrow \mathbb{R}$ mapping observables x to a summary statistic $x' = s(x)$.

Then, **approximate** the likelihood $p(x|\theta)$ with the surrogate $\hat{p}(x|\theta) = p(x'|\theta)$.

From this it comes

$$\frac{p(x|\theta_0)}{p(x|\theta_1)} \approx \frac{\hat{p}(x|\theta_0)}{\hat{p}(x|\theta_1)} = \hat{r}(x|\theta_0, \theta_1).$$

Wilks theorem

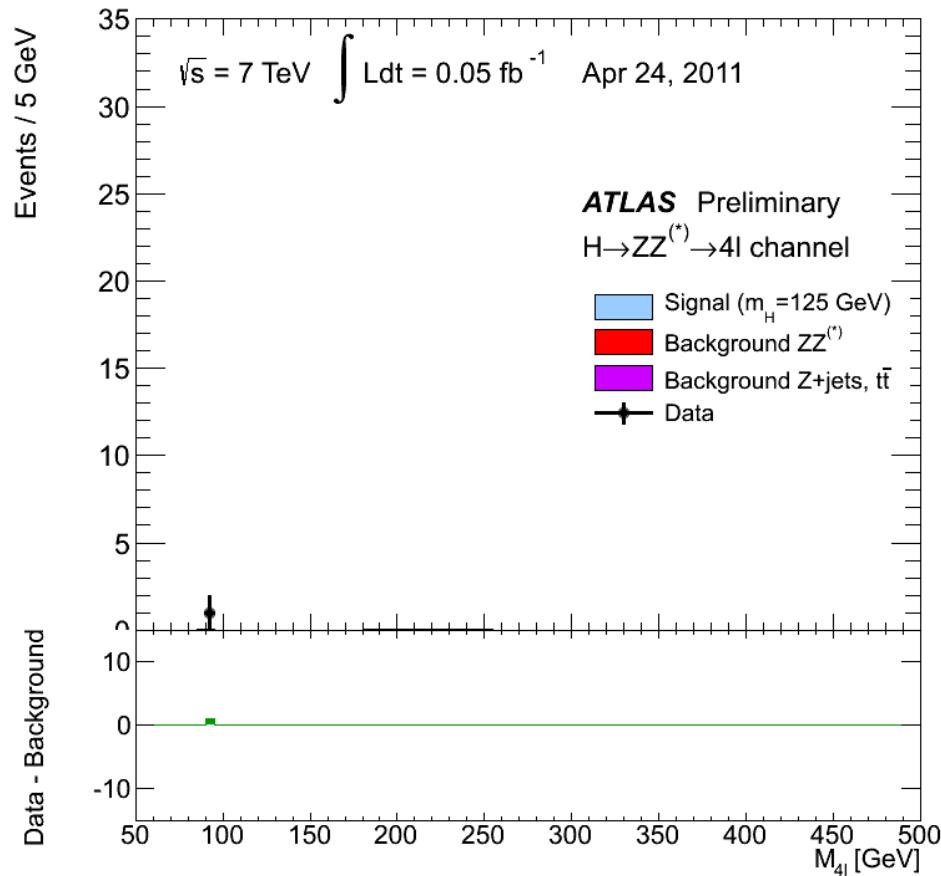
Consider the test statistic

$$q(\theta) = -2 \sum_x \log \frac{p(x|\theta)}{p(x|\hat{\theta})} = -2 \sum_x \log r(x|\theta, \hat{\theta})$$

for a fixed number N of observations $\{x\}$ and where $\hat{\theta}$ is the maximum likelihood estimator.

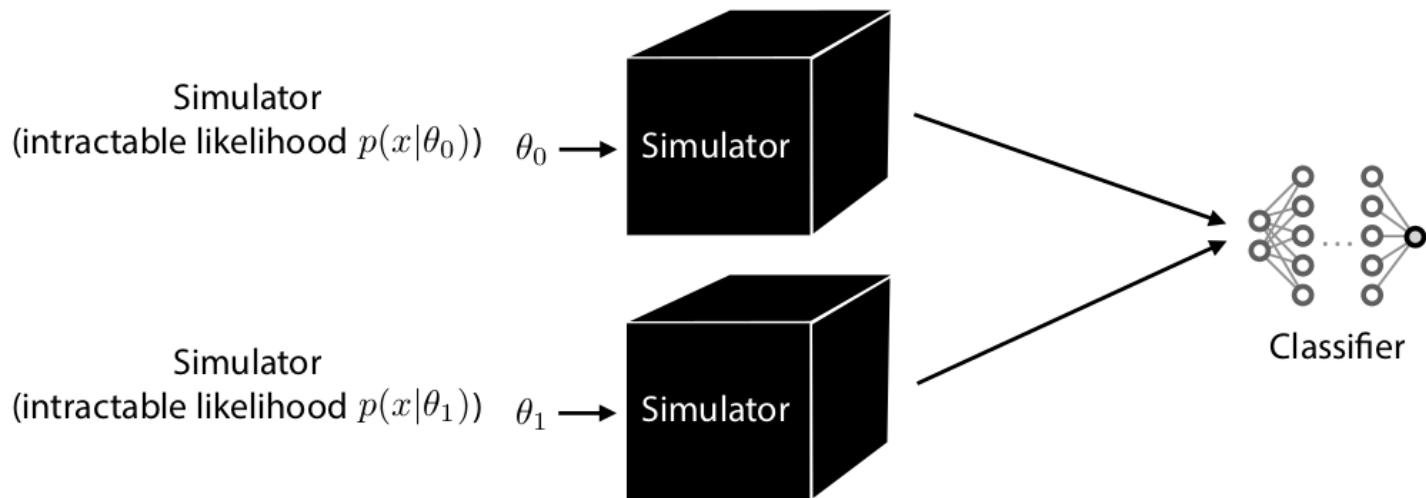
When $N \rightarrow \infty, q(\theta) \sim \chi_2$.

Therefore (and provided the assumptions apply!), an observed value $q_{\text{obs}}(\theta)$ translates directly to a p-value that measures the confidence with which θ can be excluded.



Discovery of the Higgs boson at 5σ

The likelihood ratio trick



The solution \hat{s} found after training approximates the optimal classifier

$$\hat{s}(x) \approx s^*(x) = \frac{p(x|\theta_1)}{p(x|\theta_0) + p(x|\theta_1)}.$$

Therefore,

$$r(x|\theta_0, \theta_1) \approx \hat{r}(x|\theta_0, \theta_1) = \frac{1 - \hat{s}(x)}{\hat{s}(x)}.$$

To avoid retraining a classifier \hat{s} for every (θ_0, θ_1) pair, fix θ_1 to θ_{ref} and train a single **parameterized** classifier $\hat{s}(x|\theta_0, \theta_{\text{ref}})$ where θ_0 is also given as input.

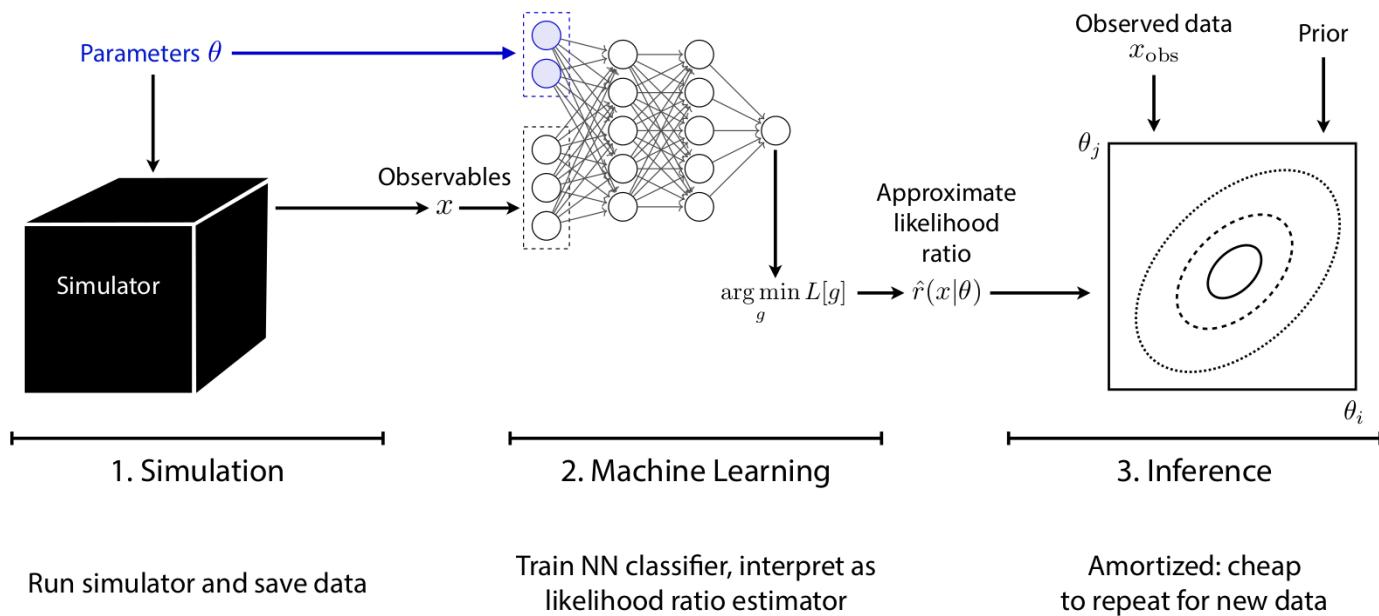
Therefore, we have

$$\hat{r}(x|\theta_0, \theta_{\text{ref}}) = \frac{1 - \hat{s}(x|\theta_0, \theta_{\text{ref}})}{\hat{s}(x|\theta_0, \theta_{\text{ref}})}$$

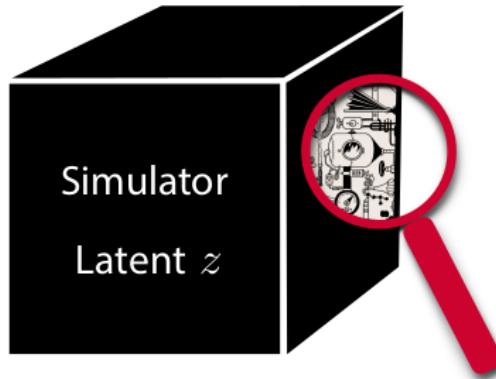
such that for any (θ_0, θ_1) ,

$$r(x|\theta_0, \theta_1) \approx \frac{\hat{r}(x|\theta_0, \theta_{\text{ref}})}{\hat{r}(x|\theta_1, \theta_{\text{ref}})}.$$

Inference



Gold mining



We cannot compute $p(x|\theta) = \int p(x, z|\theta) dz$.

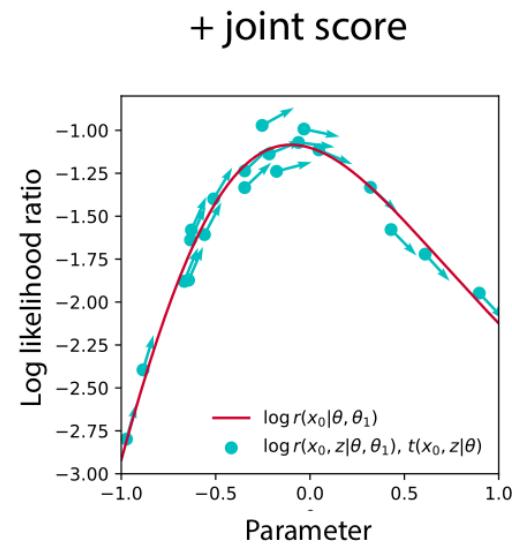
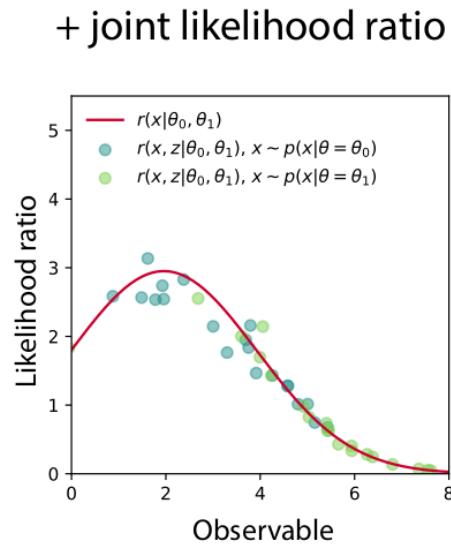
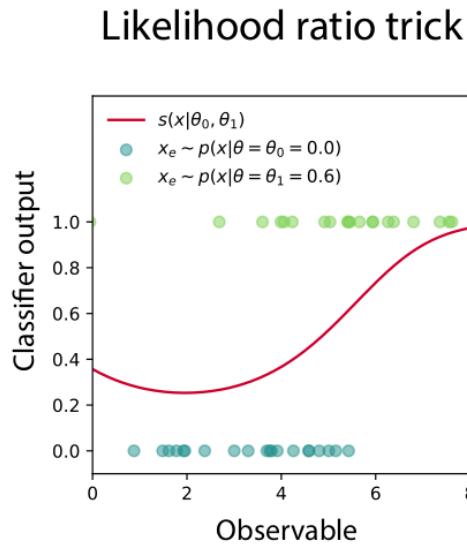
However, using techniques from probabilistic programming we can often extract

- the joint likelihood ratio $r(x, z|\theta) = \frac{p(x, z|\theta)}{p_{\text{ref}}(x, z)}$
- the joint score $t(x, z|\theta) = \nabla_{\theta} \log p(x, z|\theta)$.

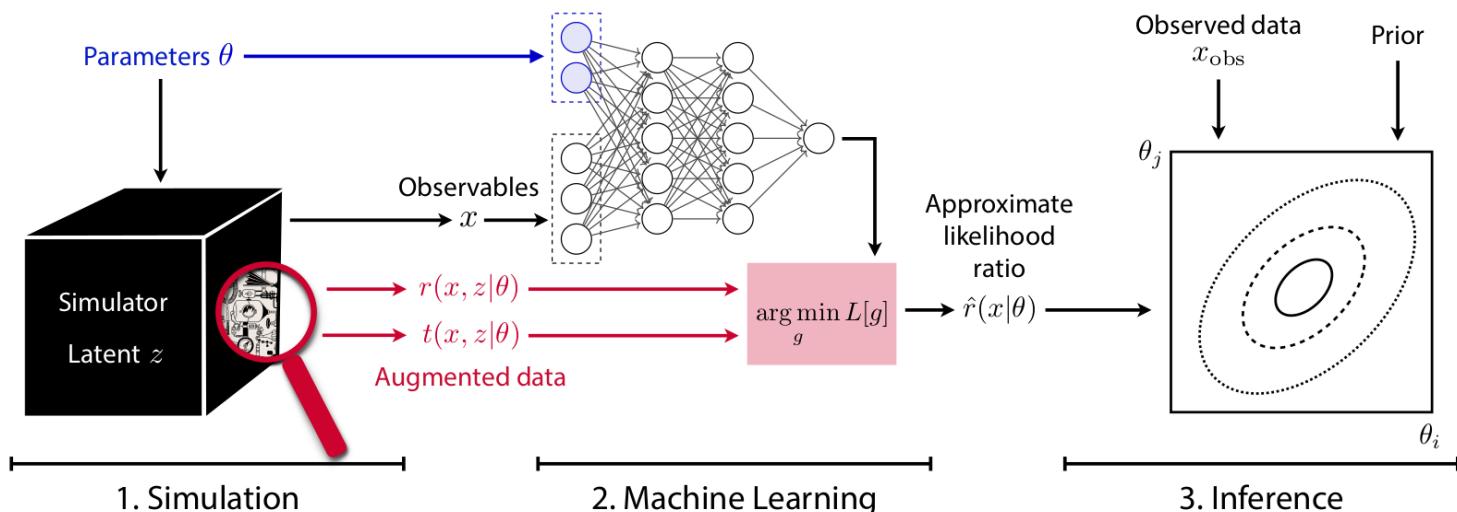
This is interesting because

- the joint likelihood ratio is an unbiased estimator of the likelihood ratio,
- the joint score provides unbiased gradient information

⇒ use them as labels in supervised NN training!



RASCAL

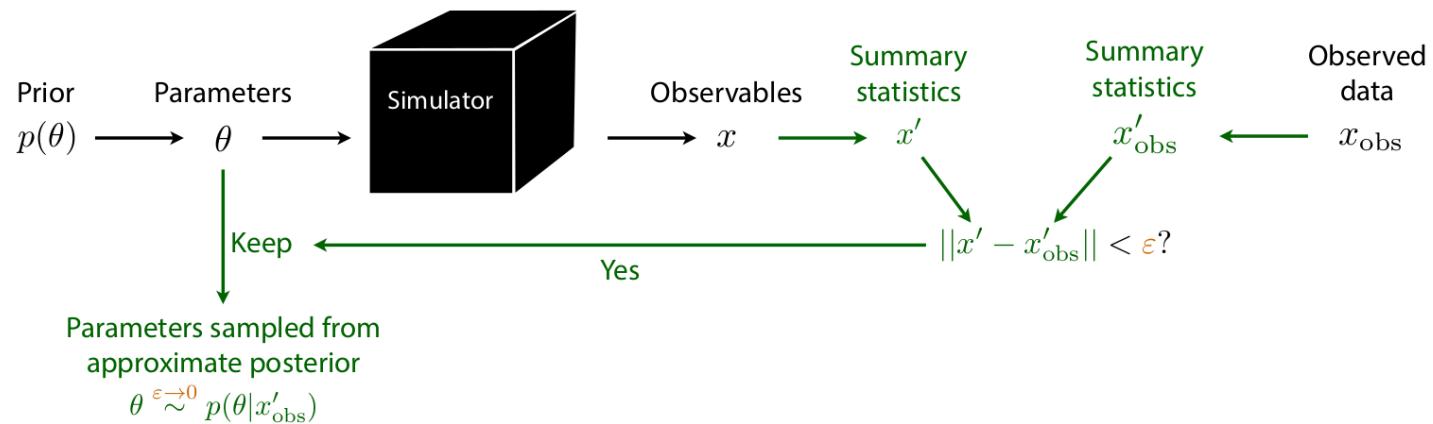


Extract joint likelihood ratio
and joint score from simulator

Augment training data &
use as labels in new loss functions
⇒ improve training efficiency

The Bayesian way

Approximate Bayesian Computation (ABC)



Issues

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

Amortizing Bayes

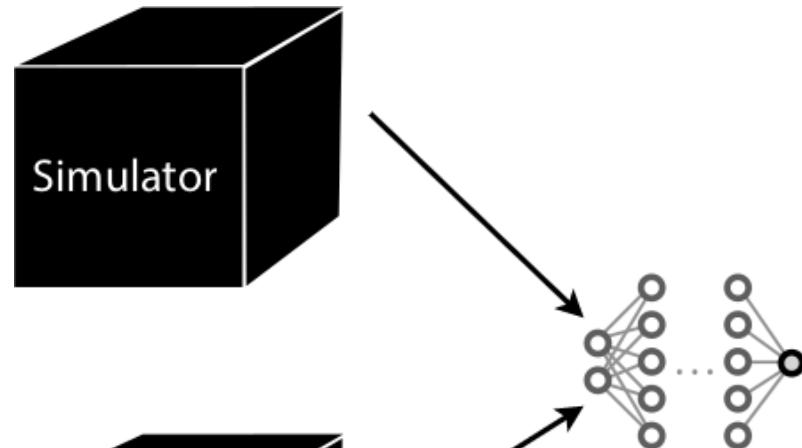
The Bayes rule can be rewritten as

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

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

As previously, the ratio can be learned with the likelihood ratio trick!

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



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

The solution \hat{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,

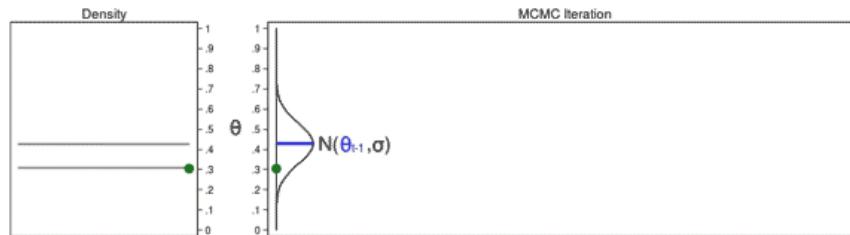
$$r(x|\theta) = \frac{p(x|\theta)}{p(x)} = \frac{p(x, \theta)}{p(x)p(\theta)} \approx \frac{d(x, \theta)}{1 - d(x, \theta)} = \hat{r}(x|\theta).$$

Likelihood-free MCMC

MCMC samplers require the evaluation of the posterior ratios, which can be obtained by evaluating the ratio of ratios:

$$\begin{aligned}\frac{p(\theta_{\text{new}} | x)}{p(\theta_{t-1} | x)} &= \frac{p(x | \theta_{\text{new}}) p(\theta_{\text{new}}) / p(x)}{p(x | \theta_{t-1}) p(\theta_{t-1}) / p(x)} \\ &= \frac{r(x | \theta_{\text{new}})}{r(x | \theta_{t-1})} \frac{p(\theta_{\text{new}})}{p(\theta_{t-1})}.\end{aligned}$$

Extensions with HMC is possible since $\nabla_{\theta} p(x | \theta) = \frac{\nabla_{\theta} r(x | \theta)}{r(x | \theta)}$.



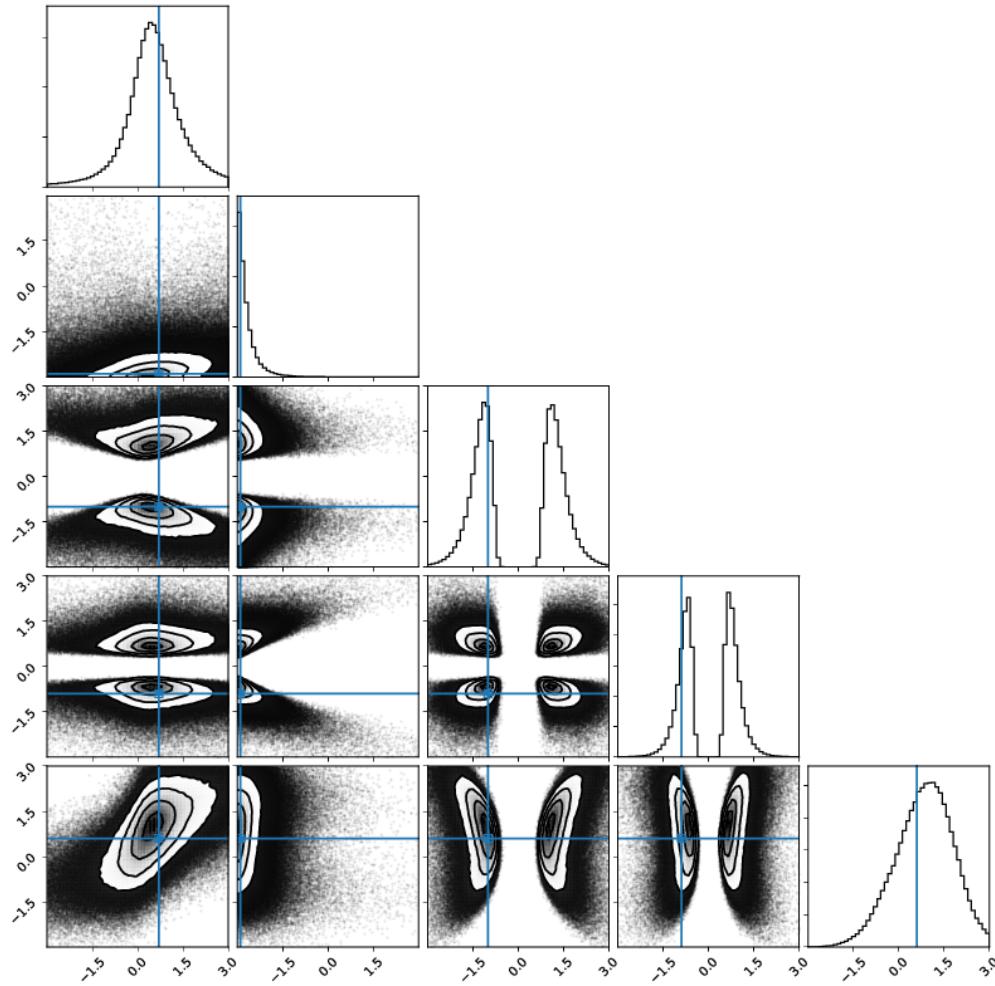
$$\text{Step 1: } r(\theta_{\text{new}}, \theta_{t-1}) = \frac{\text{Posterior}(\theta_{\text{new}})}{\text{Posterior}(\theta_{t-1})} = \frac{\text{Beta}(1,1, 0.306) \times \text{Binomial}(10,4, 0.306)}{\text{Beta}(1,1, 0.429) \times \text{Binomial}(10,4, 0.429)} = 0.834$$

$$\text{Step 2: Acceptance probability } \alpha(\theta_{\text{new}}, \theta_{t-1}) = \min\{r(\theta_{\text{new}}, \theta_{t-1}), 1\} = \min\{0.834, 1\} = 0.834$$

$$\text{Step 3: Draw } u \sim \text{Uniform}(0,1) = 0.617$$

$$\text{Step 4: If } u < \alpha(\theta_{\text{new}}, \theta_{t-1}) \rightarrow \text{If } 0.617 < 0.834 \quad \text{Then } \theta_t = \theta_{\text{new}} = 0.306 \\ \text{Otherwise } \theta_t = \theta_{t-1} = 0.429$$

Diagnostics



How to assess that the approximate posterior is not 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.

Convergence towards the nominal value θ^*

If the approximation \hat{r} is correct, then the posterior

$$\begin{aligned}\hat{p}(\theta|\mathcal{X}) &= \frac{p(\theta)p(\mathcal{X}|\theta)}{p(\mathcal{X})} = p(\theta) \left[\int p(\theta') \prod_{x \in \mathcal{X}} \frac{p(x|\theta')}{p(x|\theta)} d\theta' \right]^{-1} \\ &\approx p(\theta) \left[\int p(\theta') \prod_{x \in \mathcal{X}} \frac{\hat{r}(x|\theta')}{\hat{r}(x|\theta)} d\theta' \right]^{-1}\end{aligned}$$

should concentrate around θ^* as the number of observations

$$\mathcal{X} = \{x_1, \dots, x_n\},$$

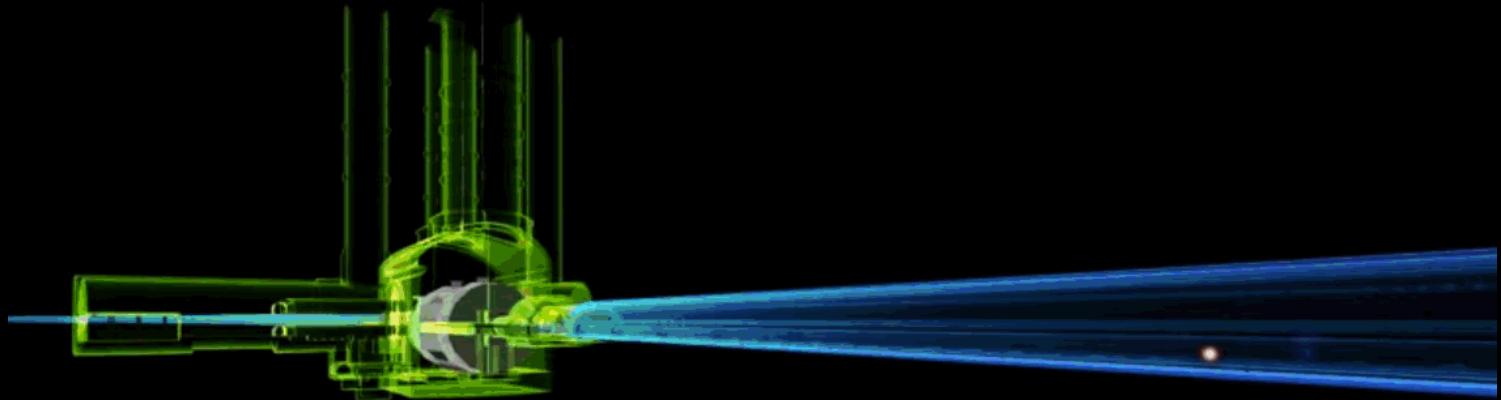
for $x_i \sim p(x|\theta^*)$, increases.

ROC AUC score

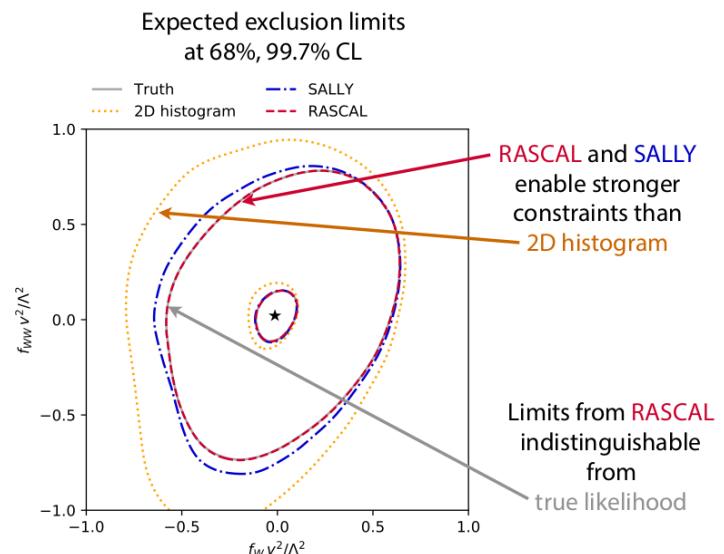
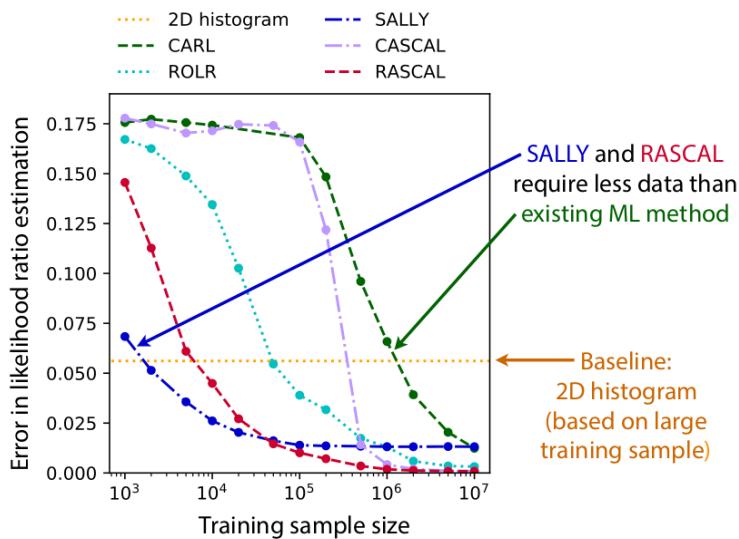
The ratio estimator $\hat{r}(x|\theta)$ is only exact when samples x from the reweighted marginal model $p(x)\hat{r}(x|\theta)$ cannot be distinguished from samples x from a specific likelihood $p(x|\theta)$.

Therefore, the predictive ROC AUC performance of a classifier should be close to **0.5** if the ratio is correct.

Showtime!

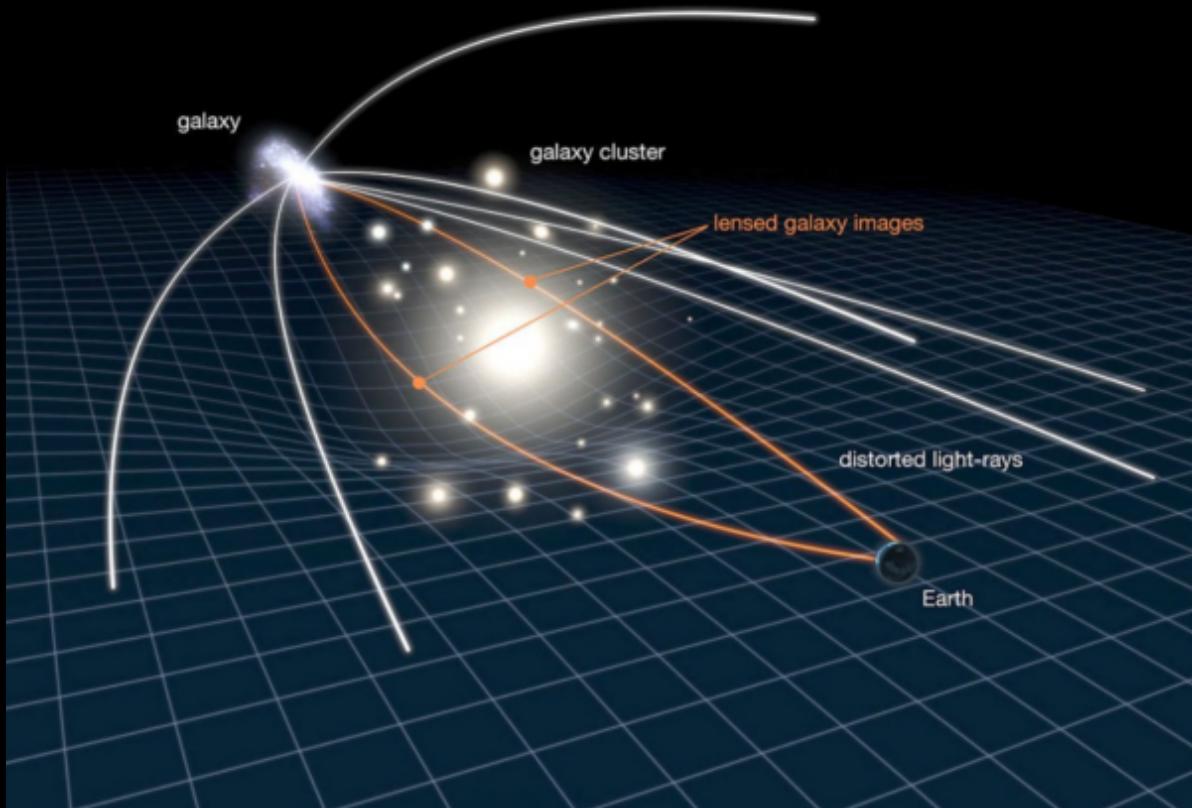


Case 1: Hunting new physics at particle colliders



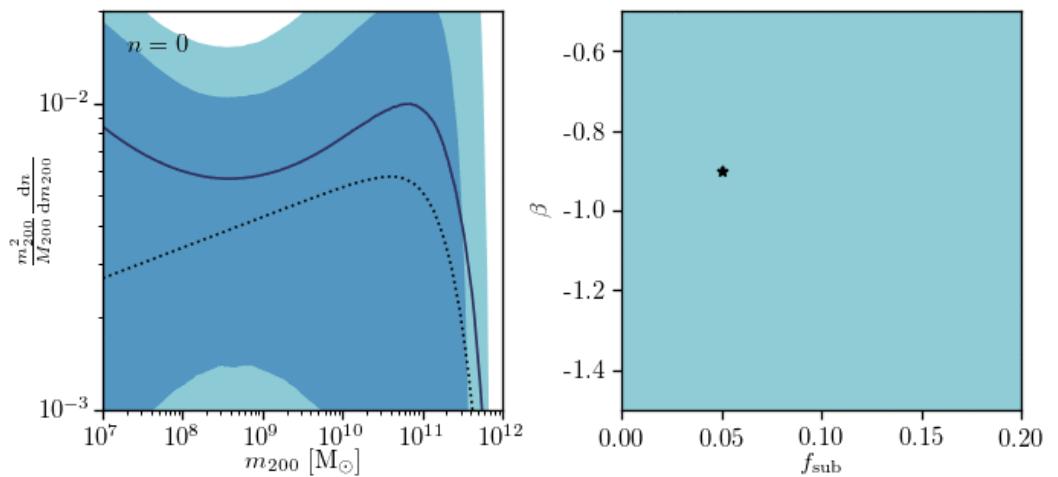
With enough training data, the ML algorithms get the likelihood function right.

Using more information from the simulator improves sample efficiency substantially.



Case 2: Dark matter substructure from gravitational lensing





Case 3: Constraining dark matter with stellar streams



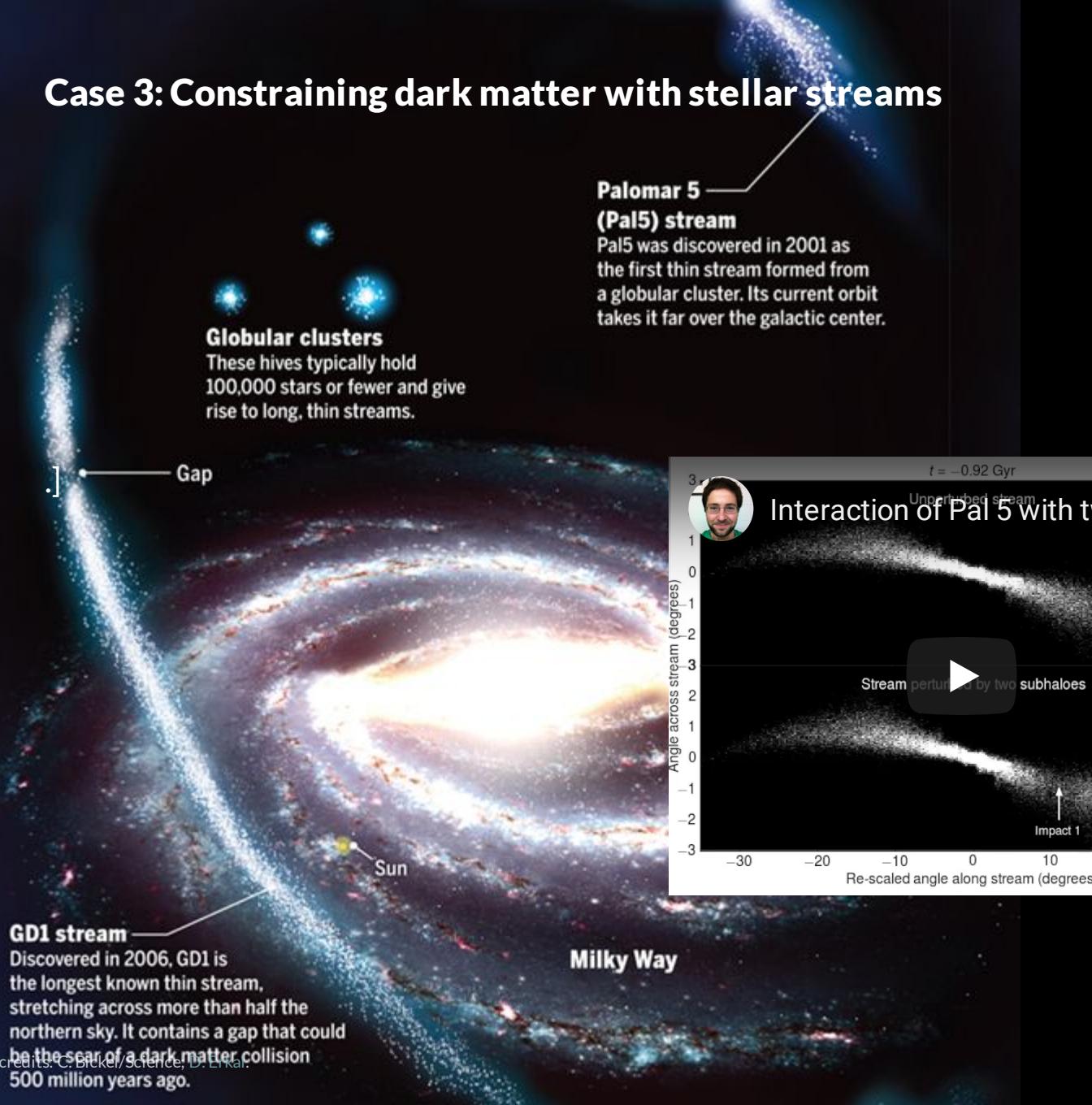
Globular clusters

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

Palomar 5

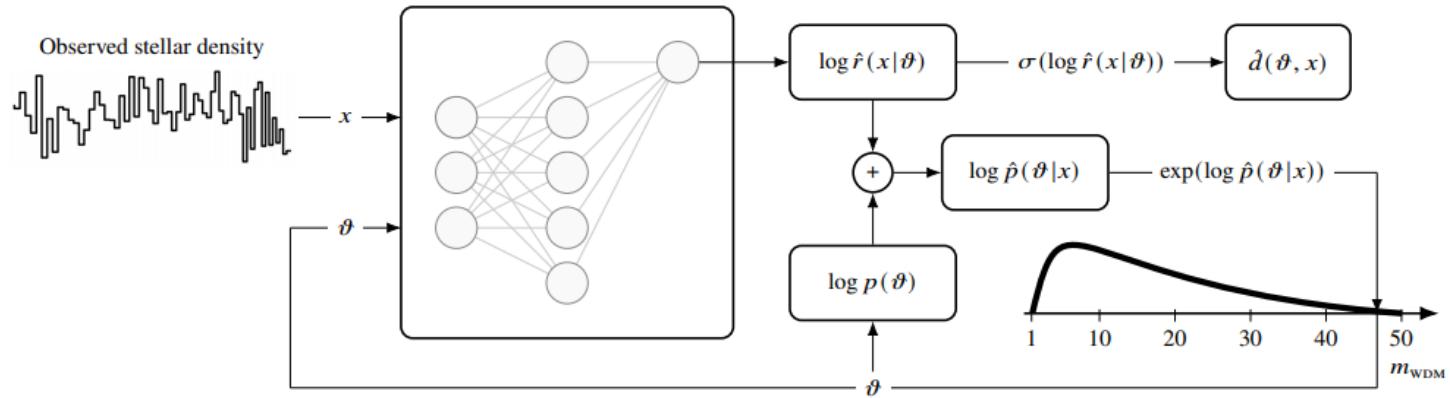
(Pal5) stream

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.

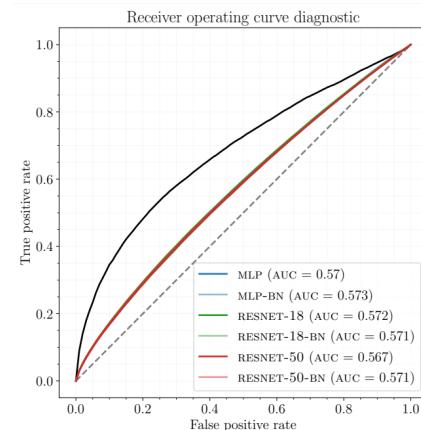
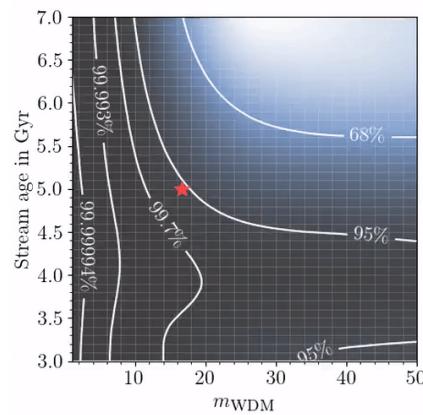


GD1 stream

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.



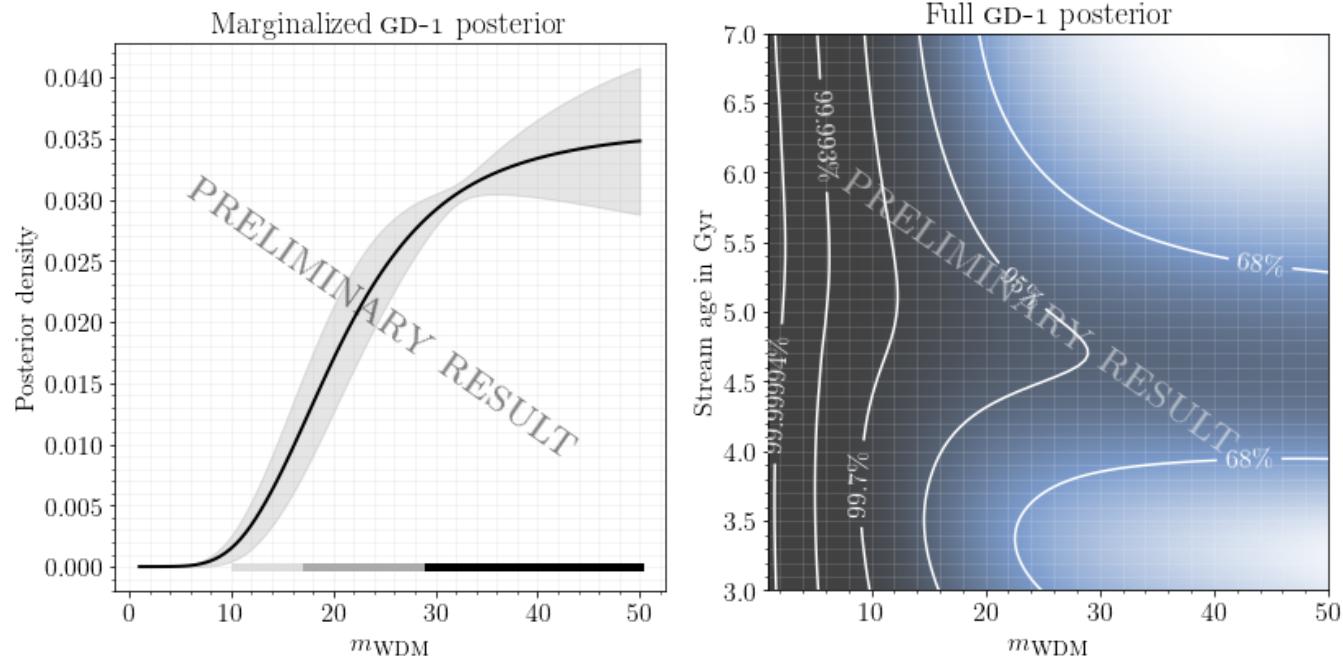
Architecture	68% CR	95% CR
$\hat{r}(x \theta)$ with $\theta \triangleq (m_{\text{WDM}}, t_{\text{age}})$		
MLP	0.685 ± 0.004	0.954 ± 0.002
MLP-BN	0.687 ± 0.006	0.951 ± 0.002
RESNET-18	0.667 ± 0.004	0.943 ± 0.002
RESNET-18-BN	0.672 ± 0.004	0.945 ± 0.001
RESNET-50	0.671 ± 0.005	0.947 ± 0.003
RESNET-50-BN	0.678 ± 0.004	0.949 ± 0.004
$\hat{r}(x \theta)$ with $\theta \triangleq (m_{\text{WDM}}, t_{\text{age}})$		
MLP	0.685 ± 0.005	0.953 ± 0.002
MLP-BN	0.685 ± 0.004	0.952 ± 0.003
RESNET-18	0.666 ± 0.005	0.945 ± 0.002
RESNET-18-BN	0.671 ± 0.003	0.945 ± 0.003
RESNET-50	0.674 ± 0.006	0.944 ± 0.002
RESNET-50-BN	0.677 ± 0.004	0.947 ± 0.003



Coverage

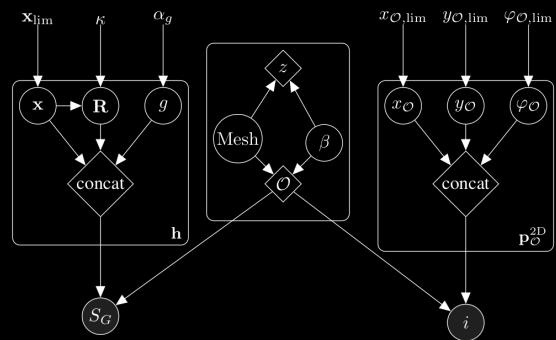
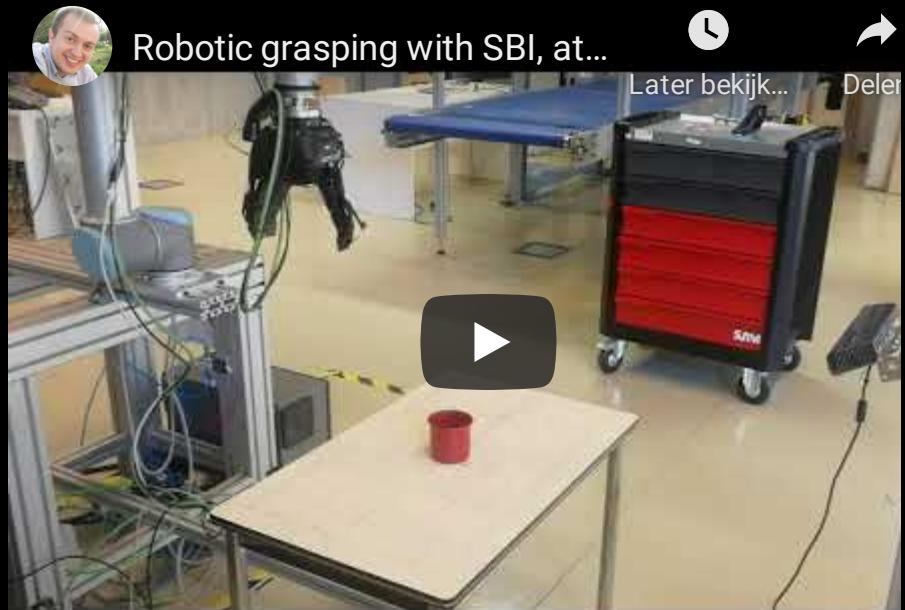
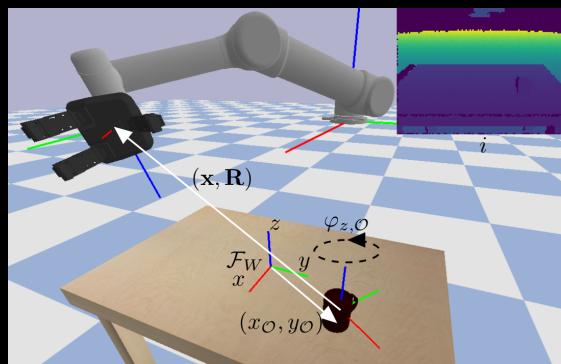
Convergence to θ^*

ROC AUC score

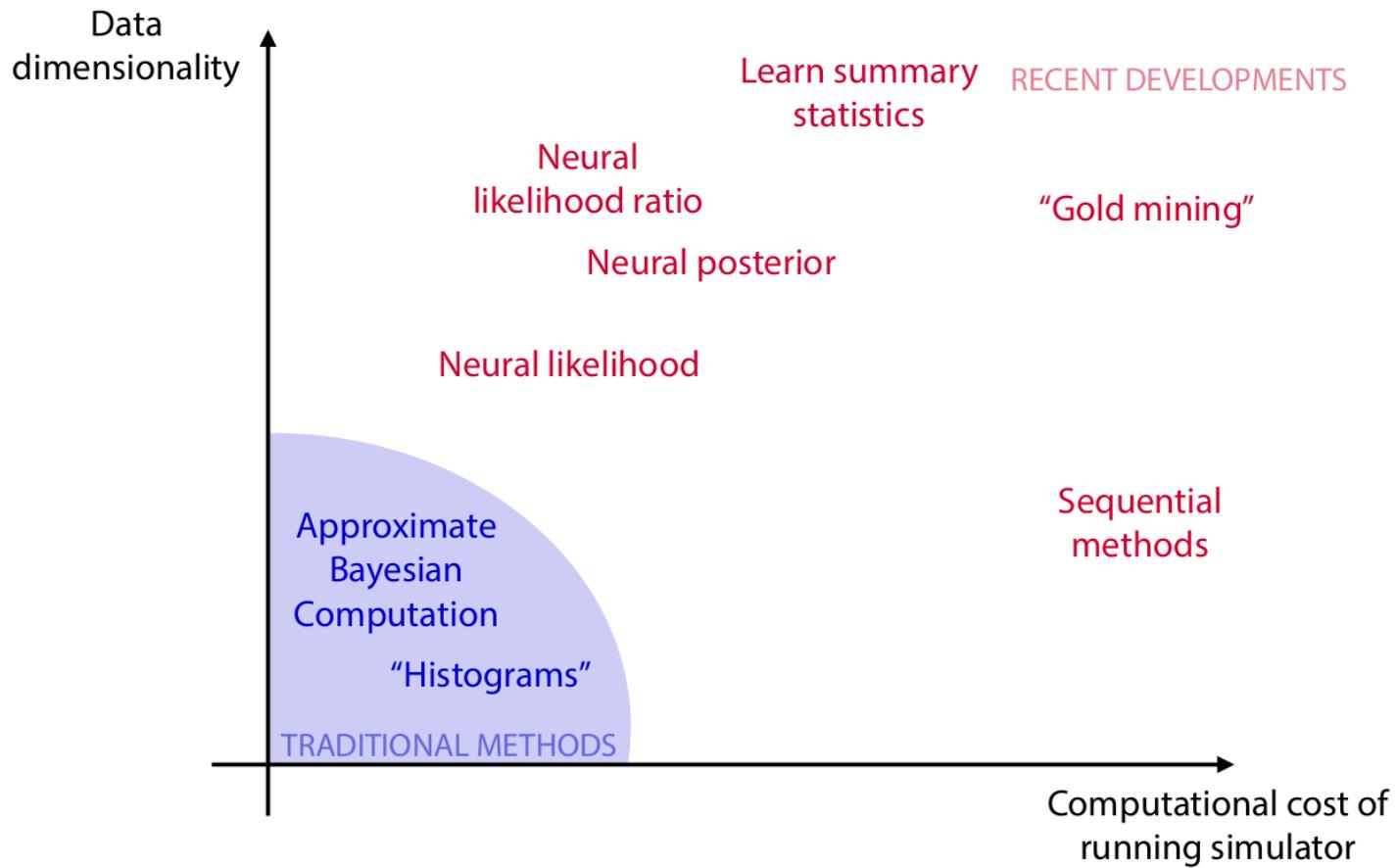


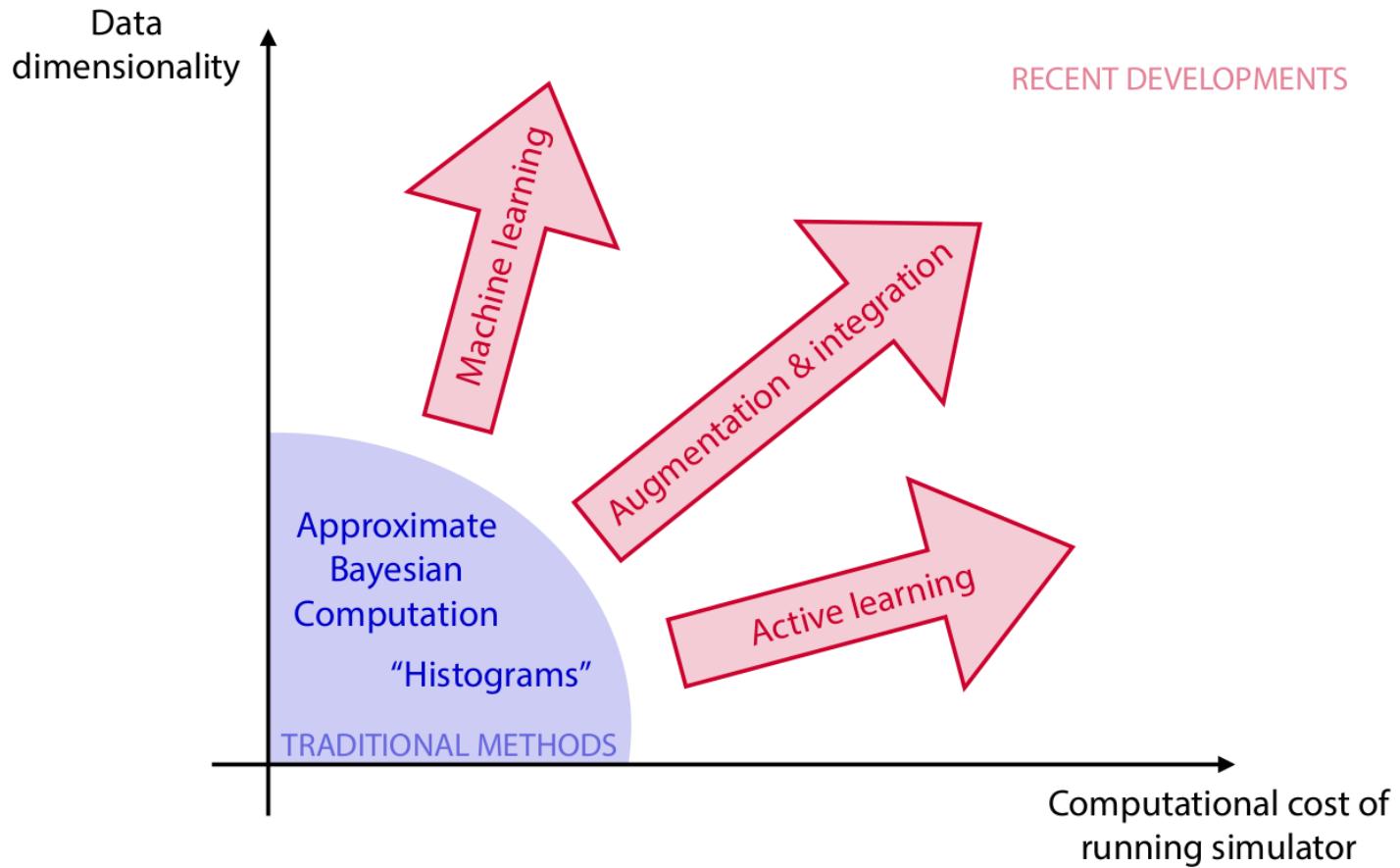
Preliminary results for GD-1 suggest a preference for CDM over WDM.

Case 4: Robotic grasping



The frontier







The frontier of simulation-based inference

Kyle Cranmer^{a,b,1}, Johann Brehmer^{a,b}, and Gilles Louppe^c

^aCenter for Cosmology and Particle Physics, New York University, New York, NY 10003; ^bCenter for Data Science, New York University, New York, NY 10011; and ^cMontefiore Institute, University of Liège, B-4000 Liège, Belgium

Edited by Jitendra Malik, University of California, Berkeley, CA, and approved April 10, 2020 (received for review November 4, 2019)

Many domains of science have developed complex simulations to describe phenomena of interest. While these simulations provide high-fidelity models, they are poorly suited for inference and lead to challenging inverse problems. We review the rapidly developing field of simulation-based inference and identify the forces giving additional momentum to the field. Finally, we describe how the frontier is expanding so that a broad audience can appreciate the profound influence these developments may have on science.

statistical inference | implicit models | likelihood-free inference | approximate Bayesian computation | neural density estimation

Mechanistic models can be used to predict how systems will behave in a variety of circumstances. These run the gamut of distance scales, with notable examples including particle physics, molecular dynamics, protein folding, population genetics, neuroscience, epidemiology, economics, ecology, climate science, astrophysics, and cosmology. The expressiveness of programming languages facilitates the development of complex, high-fidelity simulations and the power of modern computing provides the ability to generate synthetic data from them. Unfortunately, these simulators are poorly suited for statistical inference. The source of the challenge is that the probability density (or likelihood) for a given observation—an essential ingredient for both frequentist and Bayesian inference methods—is typically intractable. Such models are often referred to as implicit models and contrasted against prescribed models where the likelihood for an observation can be explicitly calculated (1). The problem setting of statistical inference under intractable likelihoods has been dubbed likelihood-free inference—although it is a bit of a misnomer as typically one attempts to estimate the intractable likelihood, so we feel the term simulation-based inference is more apt.

The intractability of the likelihood is an obstruction for scientific progress as statistical inference is a key component of the scientific method. In areas where this obstruction has appeared, scientists have had to give up on the field, if only

the simulator—is being recognized as a key idea to improve the sample efficiency of various inference methods. A third direction of research has stopped treating the simulator as a black box and focused on integrations that allow the inference engine to tap into the internal details of the simulator directly.

Amidst this ongoing revolution, the landscape of simulation-based inference is changing rapidly. In this review we aim to provide the reader with a high-level overview of the basic ideas behind both old and new inference techniques. Rather than discussing the algorithms in technical detail, we focus on the current frontiers of research and comment on some ongoing developments that we deem particularly exciting.

Simulation-Based Inference

Simulators. Statistical inference is performed within the context of a statistical model, and in simulation-based inference the simulator itself defines the statistical model. For the purpose of this paper, a simulator is a computer program that takes as input a vector of parameters θ , samples a series of internal states or latent variables $z_i \sim p_i(z_i|\theta, z_{<i})$, and finally produces a data vector $x \sim p(x|\theta, z)$ as output. Programs that involve random samplings and are interpreted as statistical models are known as probabilistic programs, and simulators are an example. Within this general formulation, real-life simulators can vary substantially:

- The parameters θ describe the underlying mechanistic model and thus affect the transition probabilities $p_i(z_i|\theta, z_{<i})$. Typically the mechanistic model is interpretable by a domain scientist and θ has relatively few components and a fixed dimensionality. Examples include coefficients found in the Hamiltonian of a physical system, the virulence and incubation rate of a pathogen, or fundamental constants of Nature.
- The latent variables z that appear in the data-generating process may directly or indirectly correspond to a physically meaningful state of a system, but typically this state is unobservable in practice. The structure of the latent space varies substantially between simulators. The latent variables may be continuous

In summary

- Much of modern science is based on simulators making precise predictions, but in which inference is challenging.
- Machine learning enables powerful inference methods.
- They work in problems from the smallest to the largest scales.
- Further advances in machine learning will translate into scientific progress.

Thanks!



References

- Hermans, J., Banik, N., Weniger, C., Bertone, G., & Louppe, G. (2020). Towards constraining warm dark matter with stellar streams through neural simulation-based inference. arXiv preprint arXiv:2011.14923.
- Cranmer, K., Brehmer, J., & Louppe, G. (2020). The frontier of simulation-based inference. Proceedings of the National Academy of Sciences, 117(48), 30055-30062.
- Brehmer, J., Mishra-Sharma, S., Hermans, J., Louppe, G., Cranmer, K. (2019). Mining for Dark Matter Substructure: Inferring subhalo population properties from strong lenses with machine learning. arXiv preprint arXiv 1909.02005.
- Hermans, J., Begy, V., & Louppe, G. (2019). Likelihood-free MCMC with Approximate Likelihood Ratios. arXiv preprint arXiv:1903.04057.
- Brehmer, J., Louppe, G., Pavez, J., & Cranmer, K. (2018). Mining gold from implicit models to improve likelihood-free inference. arXiv preprint arXiv:1805.12244.
- Brehmer, J., Cranmer, K., Louppe, G., & Pavez, J. (2018). Constraining Effective Field Theories with Machine Learning. arXiv preprint arXiv:1805.00013.
- Brehmer, J., Cranmer, K., Louppe, G., & Pavez, J. (2018). A Guide to Constraining Effective Field Theories with Machine Learning. arXiv preprint arXiv:1805.00020.
- Cranmer, K., Pavez, J., & Louppe, G. (2015). Approximating likelihood ratios with calibrated discriminative classifiers. arXiv preprint arXiv:1506.02169.

The end.