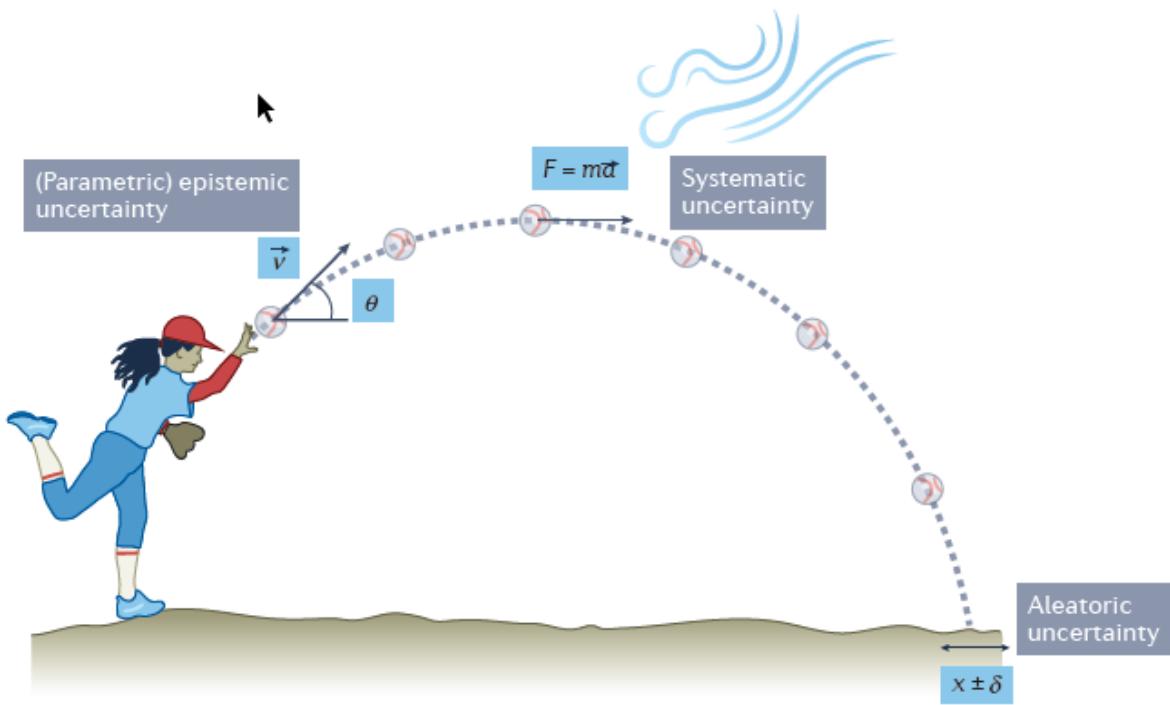


An introduction to simulation-based inference

Prof. Gilles Louppe
g.louppe@uliege.be



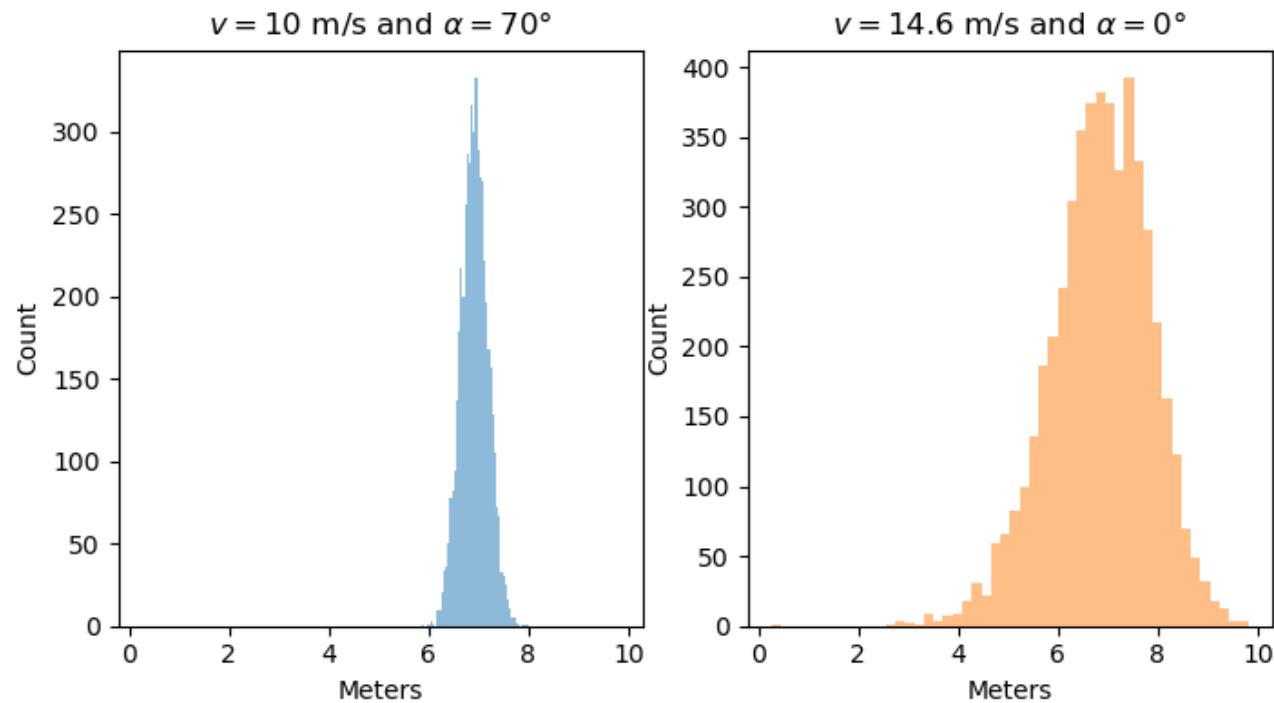
$$v_x = v \cos(\alpha), \quad v_y = v \sin(\alpha),$$

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

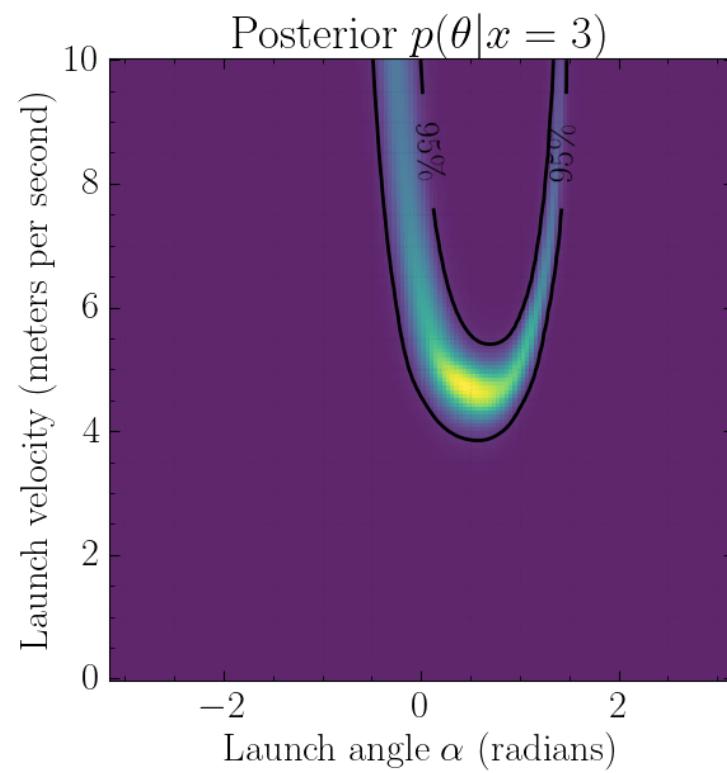
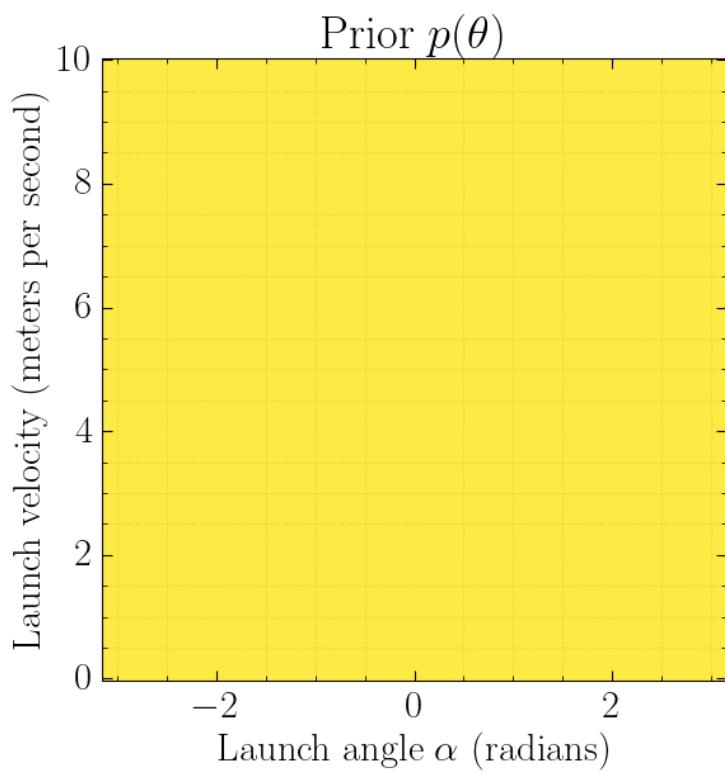
```
def simulate(v, alpha, dt=0.001):
    v_x = v * np.cos(alpha) # x velocity m/s
    v_y = v * np.sin(alpha) # y velocity m/s
    y = 1.1 + 0.3 * random.normal()
    x = 0.0

    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 θ are the most plausible?



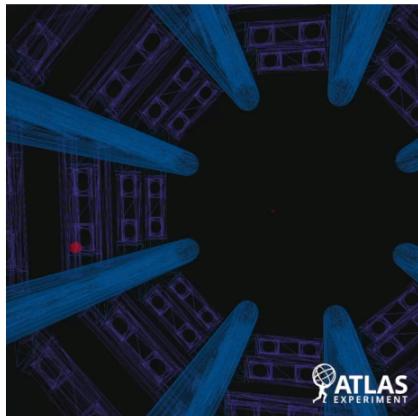
Simulation-based inference

Simulators as generative models

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

Collider data

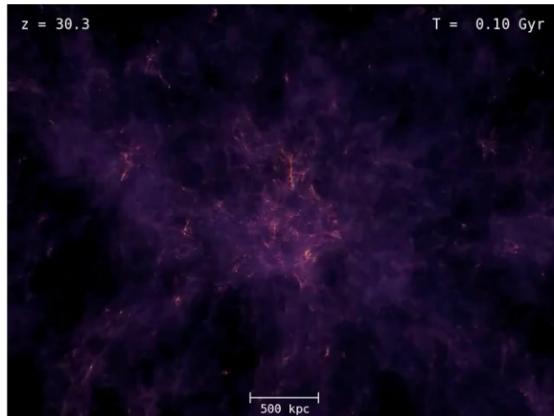
particles $\sim p(\text{particles})$



[C. Cesarotti with ATLAS]

Cosmology data

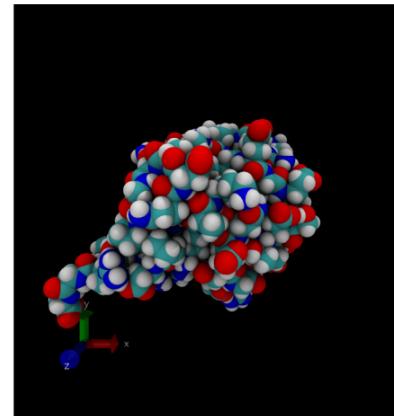
particles $\sim p(\text{particles})$



[Aquarius simulation]

Molecular dynamics

configurations $\sim p(\text{configurations})$



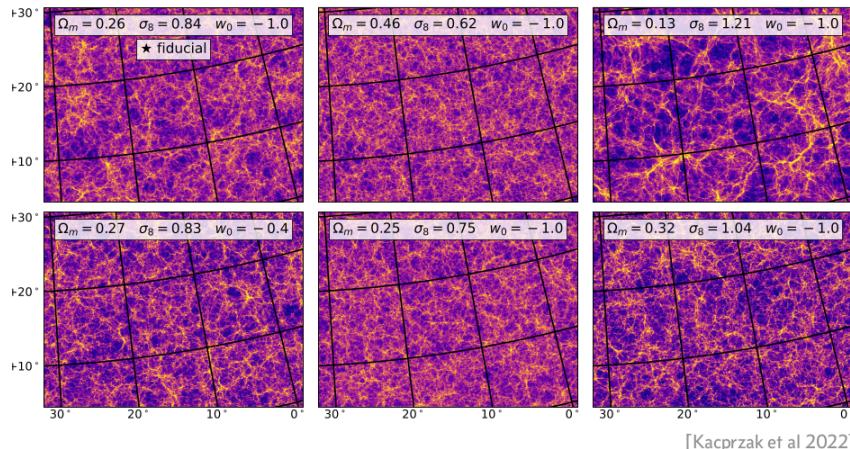
[E. Cancès et al.]

Conditional simulators

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

Cosmology data

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



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

Model

or

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

Model
parameters

What can we do with generative models?

Produce samples and
make predictions

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

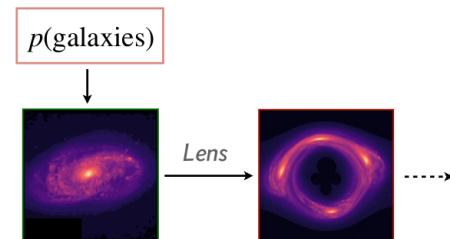
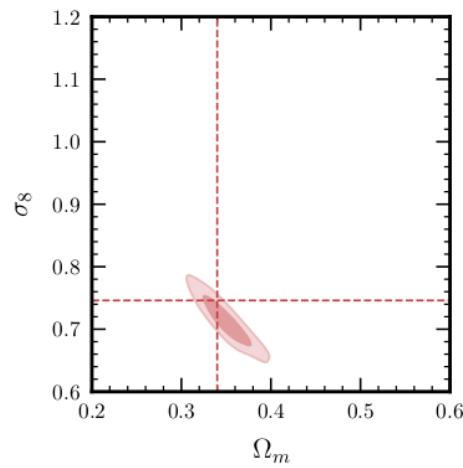
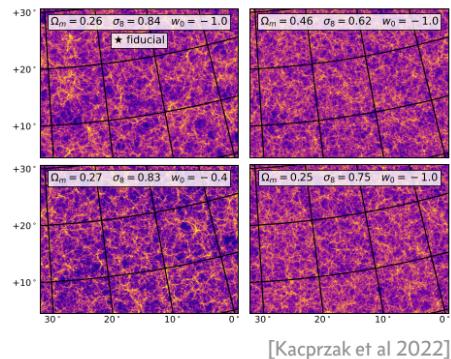
Evaluate densities

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

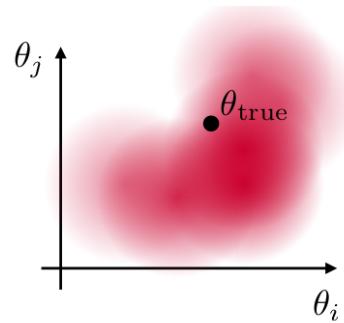
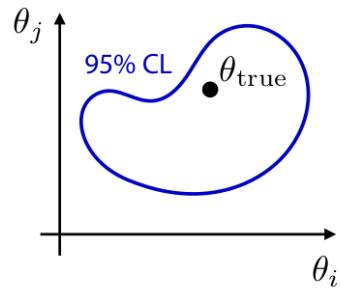
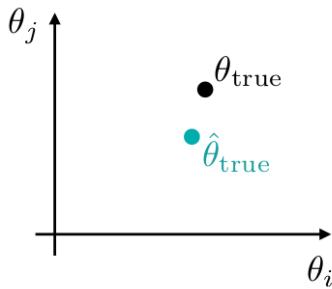
$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})}$$

Encode complex priors

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



Inference



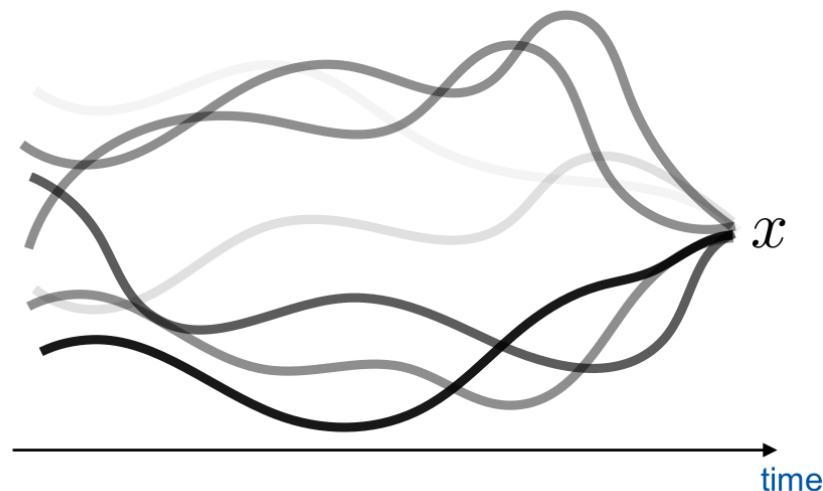
- Frequentist inference: find $\hat{\theta}$ that maximizes the likelihood $p(\mathbf{x}|\theta)$ or build a confidence interval thereof.
- Bayesian inference: compute the posterior distribution $p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})}$.

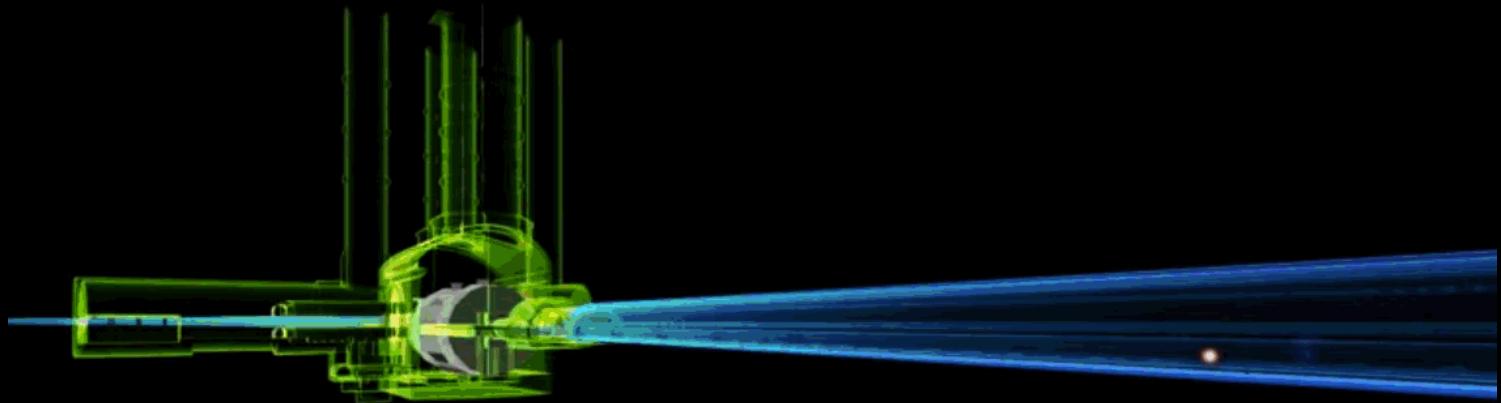
Intractable likelihoods

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.





$$p(\mathbf{z}_p | \theta)$$

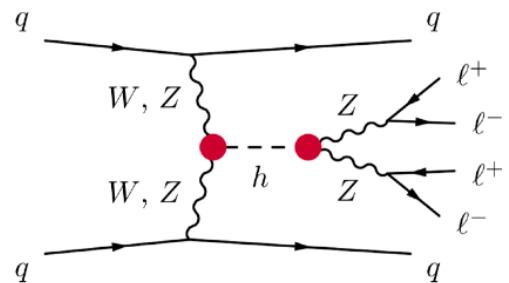
Latent variables

Parameters
of interest

Parton-level
momenta

Theory
parameters

$$z_p \xleftarrow{} \theta$$

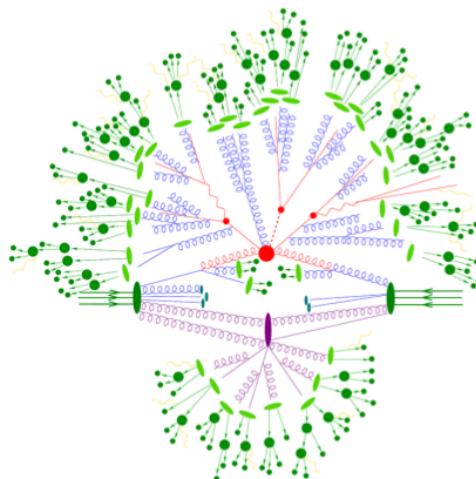


$$p(\mathbf{z}_s | \theta) = \int p(\mathbf{z}_p | \theta) p(\mathbf{z}_s | \mathbf{z}_p) d\mathbf{z}_p$$

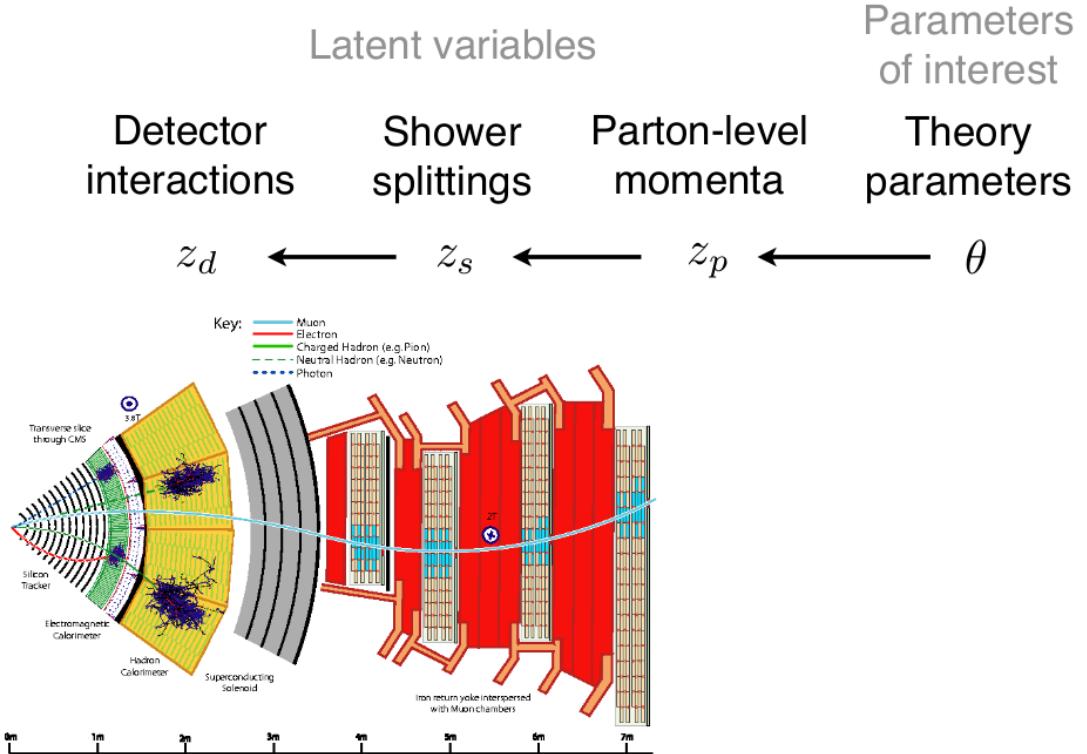
Latent variables
Parameters
of interest

Shower splittings Parton-level momenta Theory parameters

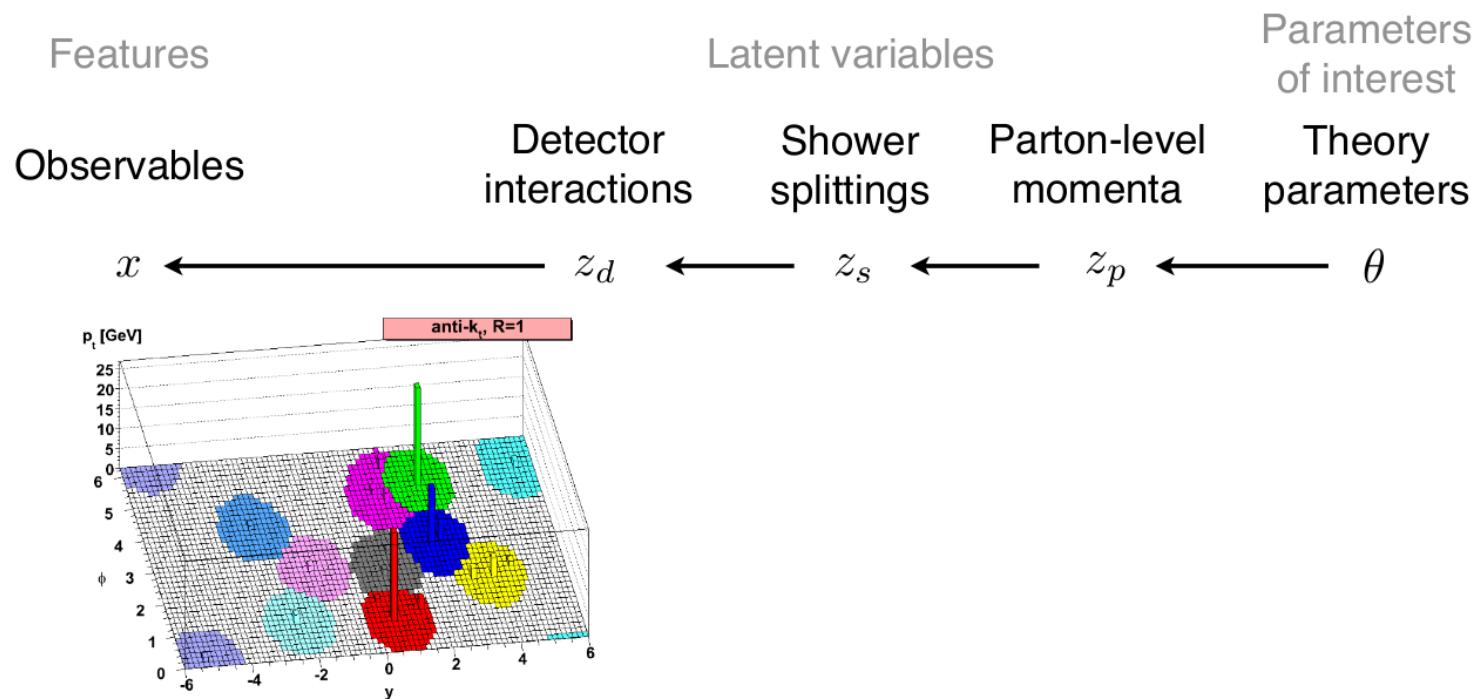
$$z_s \quad \leftarrow \quad z_p \quad \leftarrow \quad \theta$$



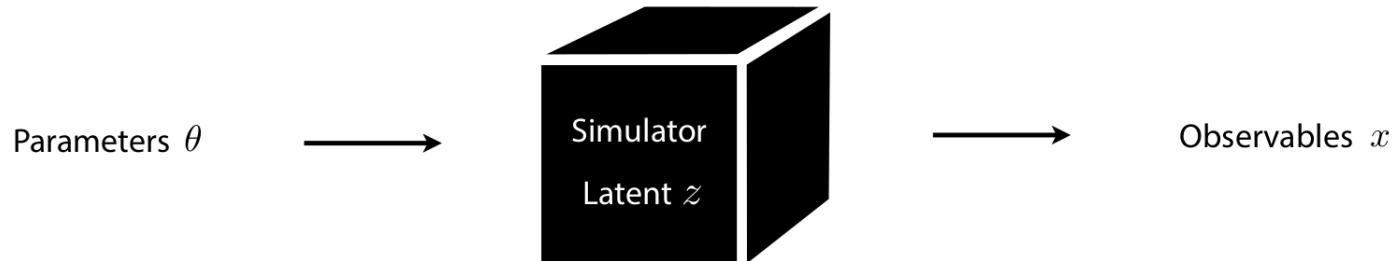
$$p(\mathbf{z}_d | \theta) = \iint p(\mathbf{z}_p | \theta) p(\mathbf{z}_s | \mathbf{z}_p) p(\mathbf{z}_d | \mathbf{z}_s) d\mathbf{z}_p d\mathbf{z}_s$$



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



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

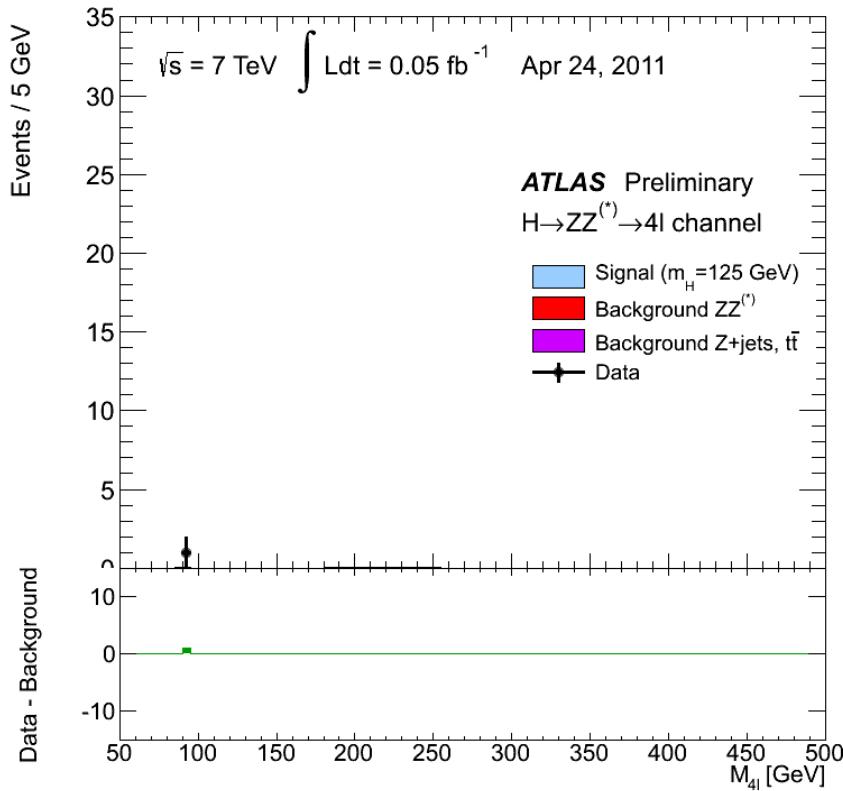


- Prediction:**
- Well-motivated mechanistic, causal model
 - Simulator can generate samples $x \sim p(x|\theta)$

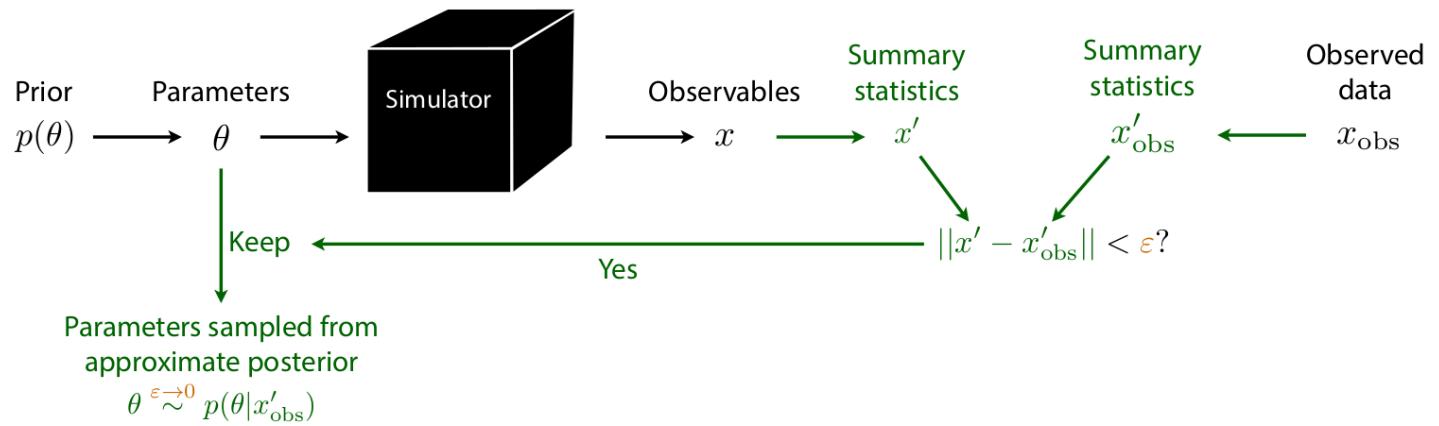
- Inference:**
- Interactions between low-level components lead to challenging inverse problems
 - Likelihood $p(x|\theta) = \int dz p(x, z|\theta)$ is intractable

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

pre-2019



(Frequentist) Approximate the likelihood $p(\mathbf{x}|\theta)$ as
 $p(\mathbf{x}|\theta) \approx \hat{p}(\mathbf{x}|\theta) = p(s(\mathbf{x})|\theta)$ for some (well-chosen) summary statistic $s(\cdot)$.



(Bayesian) Approximate the posterior $p(\theta|\mathbf{x})$ using Approximate Bayesian Computation.

Issues:

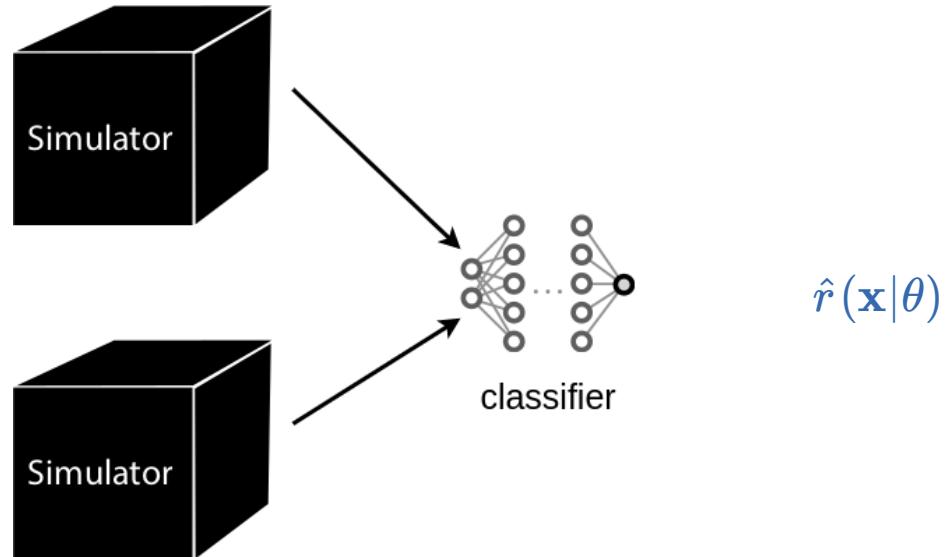
- How to choose $\mathbf{x}' = s(\mathbf{x})$? ϵ ? $\|\cdot\|$?
- No tractable posterior.
- Need to run new simulations for new data or new prior.

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 be estimated from a binary classifier $d(\theta, \mathbf{x})$, even if neither the likelihood nor the evidence can be evaluated.

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



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

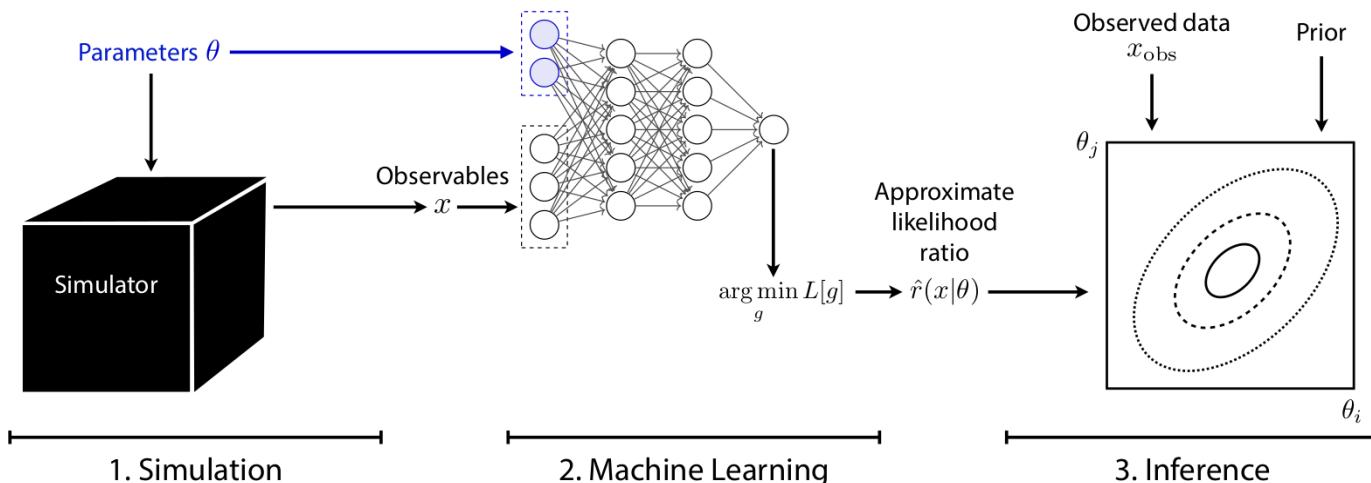


The solution \mathbf{d} found after training approximates the optimal classifier

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

Therefore,

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



Run simulator and save data

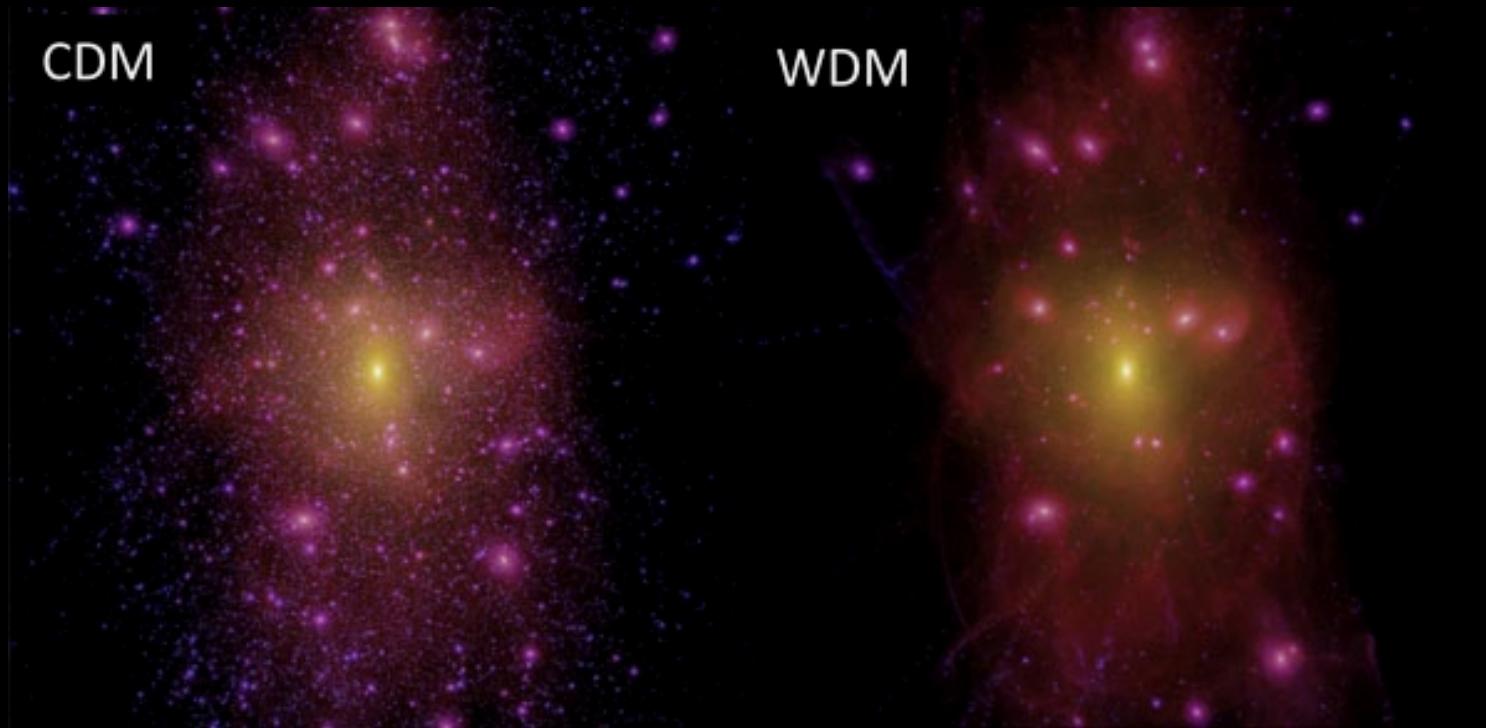
Train NN classifier, interpret as likelihood ratio estimator

Amortized: cheap to repeat for new data

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

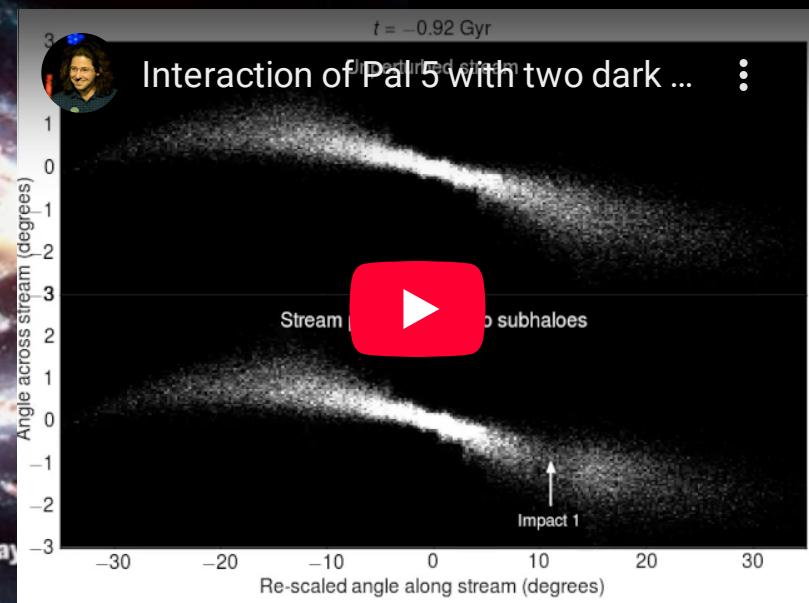
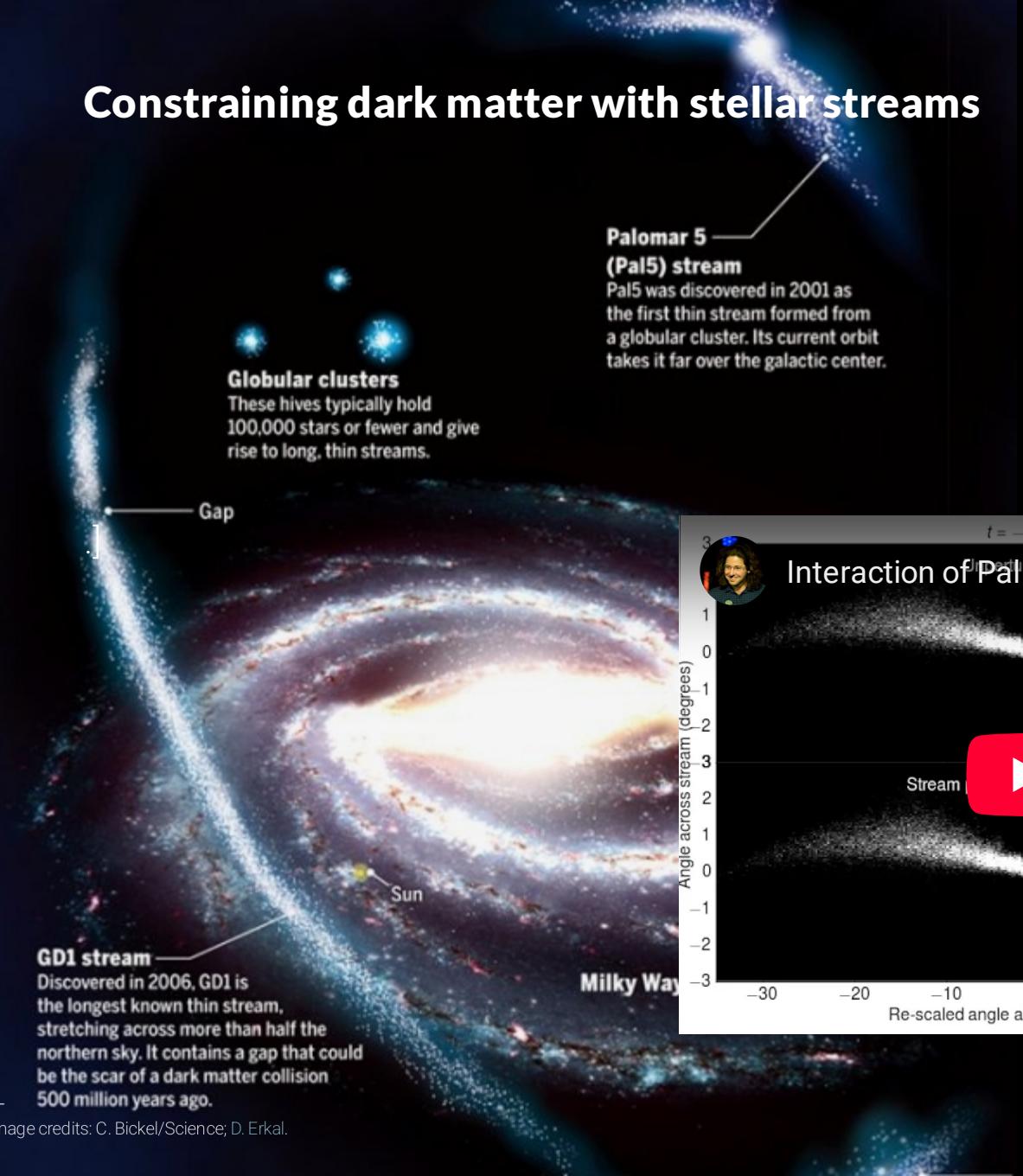


A case study (Hermans et al, 2021)



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

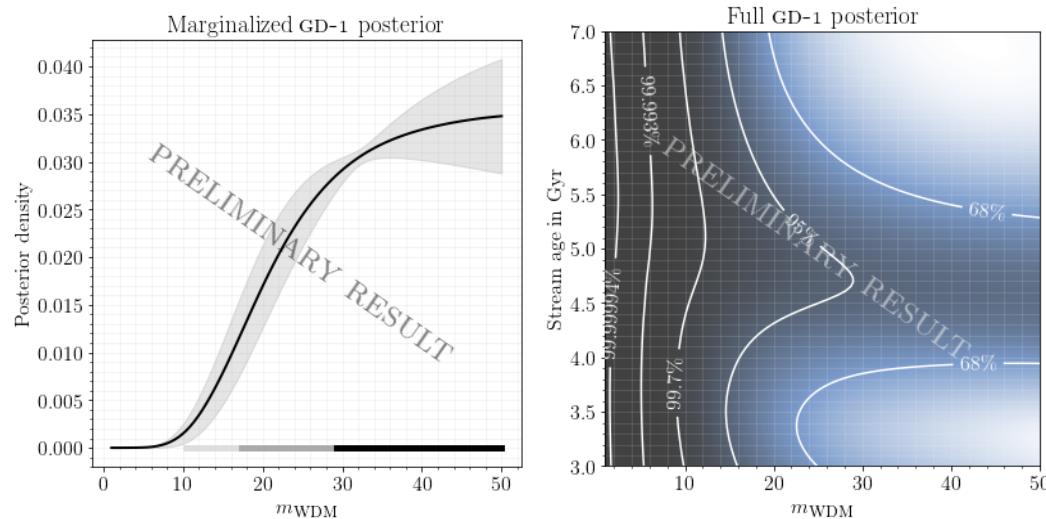
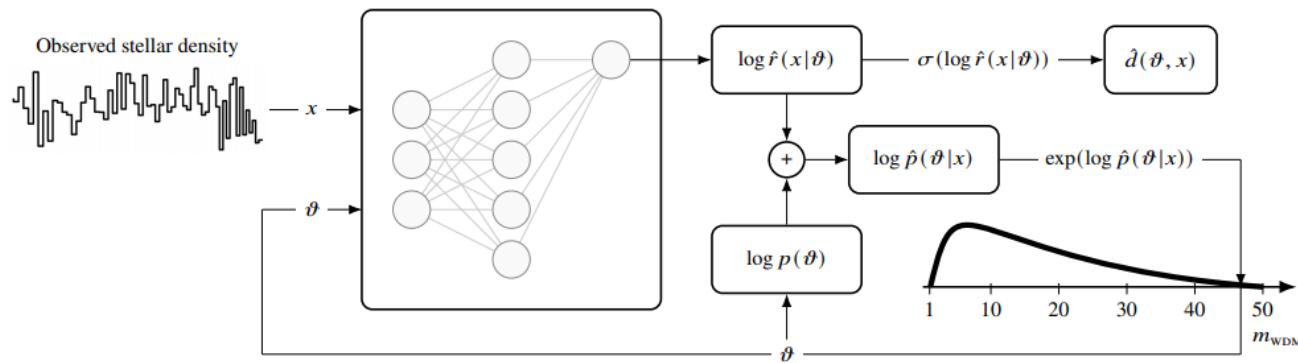
Constraining dark matter with stellar streams

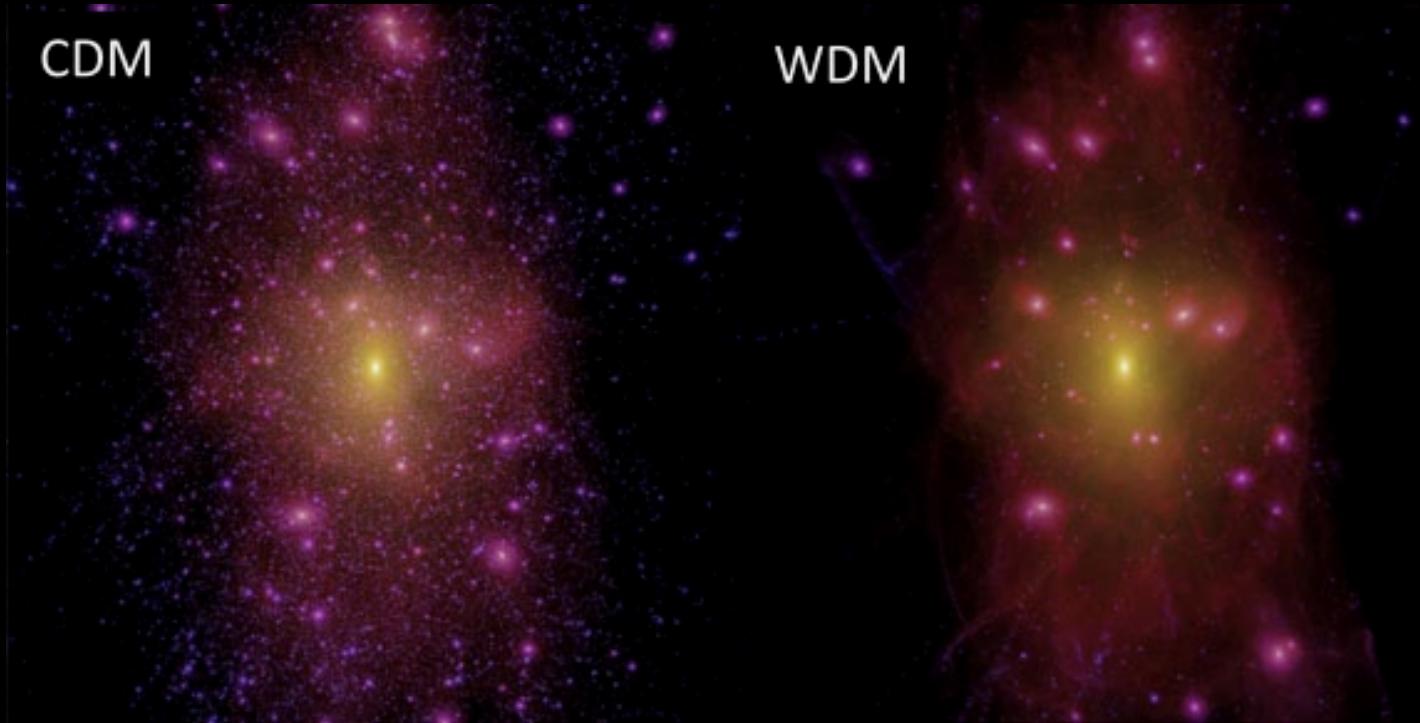


$$p(m_{\text{WDM}}, t_{\text{age}} | \text{GD1}) = \frac{p(\text{GD1} | m_{\text{WDM}}, t_{\text{age}}) p(m_{\text{WDM}}, t_{\text{age}})}{p(\text{GD-1})}$$



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!

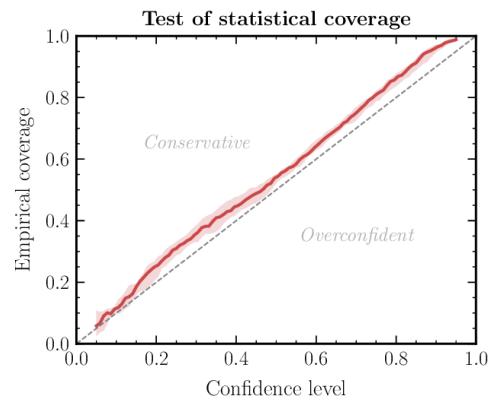
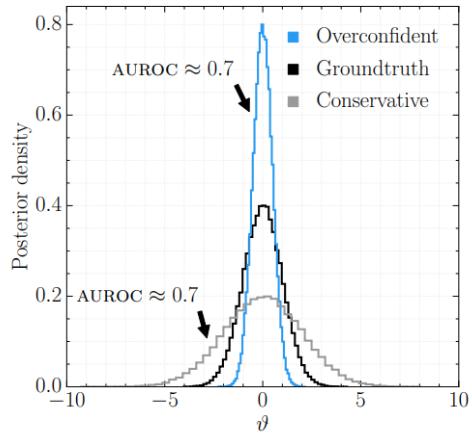


Expected coverage

$$\text{EC}(\hat{p}, \alpha) = \mathbb{E}_{p(\theta, \mathbf{x})} [\theta \in \Theta_{\hat{p}(\theta|\mathbf{x})}(\alpha)]$$

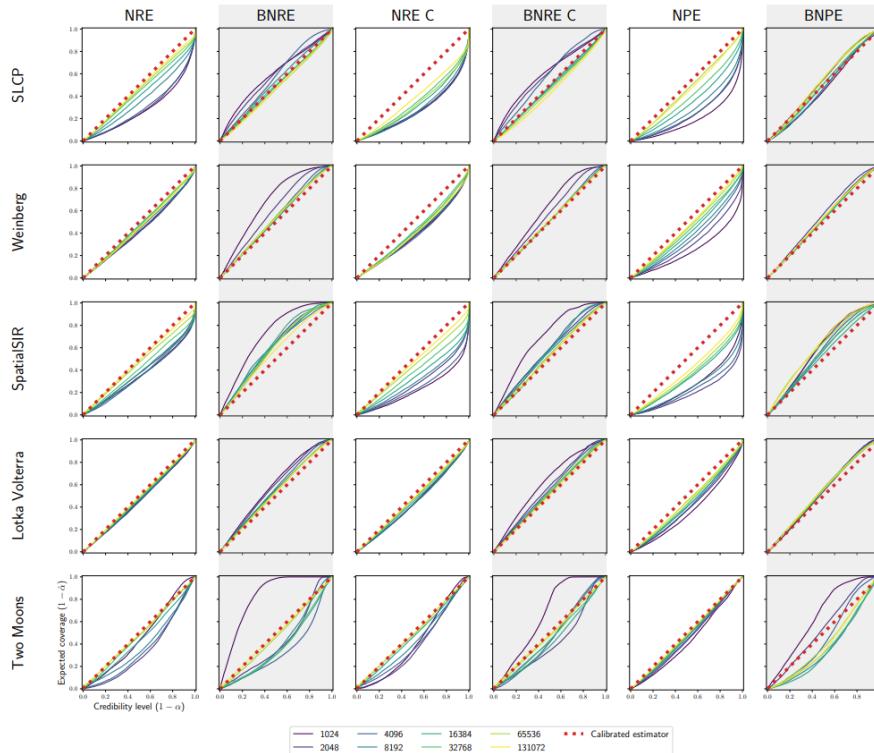
If the expected coverage is close to the nominal coverage probability α , then the approximate posterior \hat{p} is calibrated.

- If $\text{EC} < \alpha$, then the posterior is underdispersed and overconfident.
- If $\text{EC} > \alpha$, then the posterior is overdispersed and conservative.



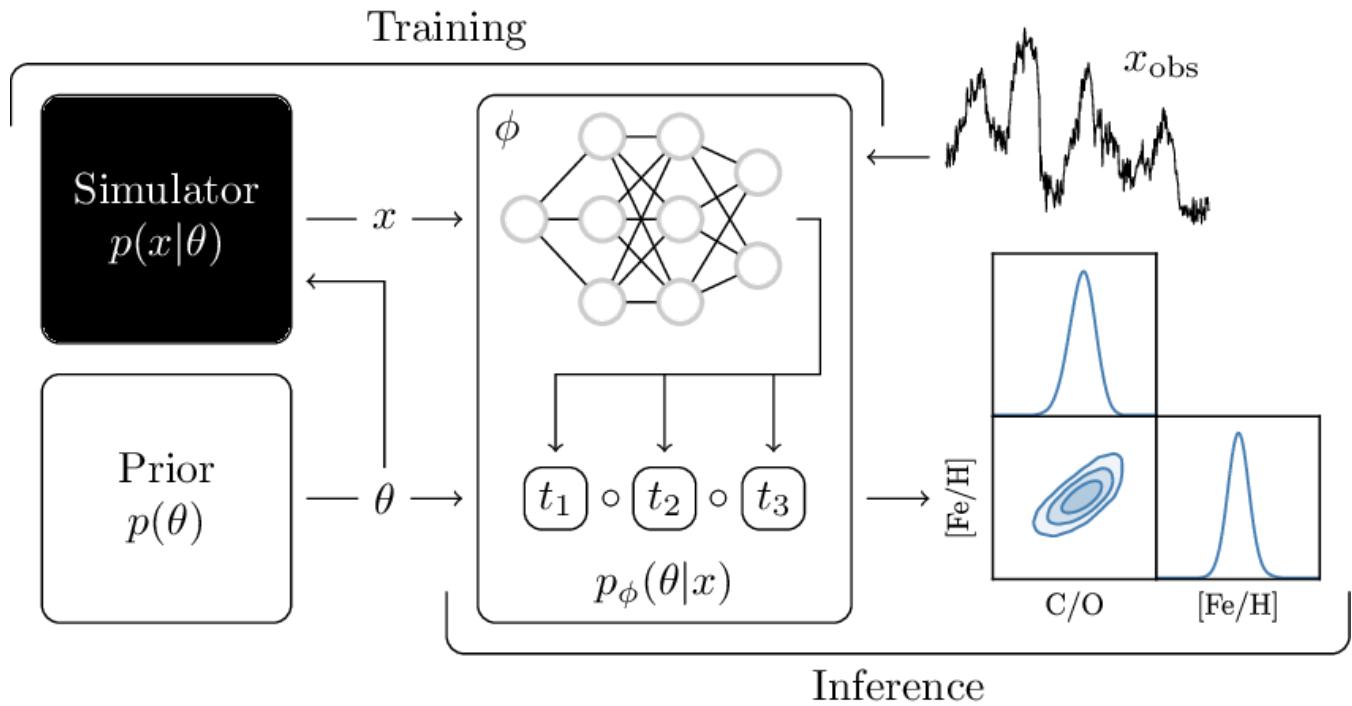


Balancing inference for conservative posteriors



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

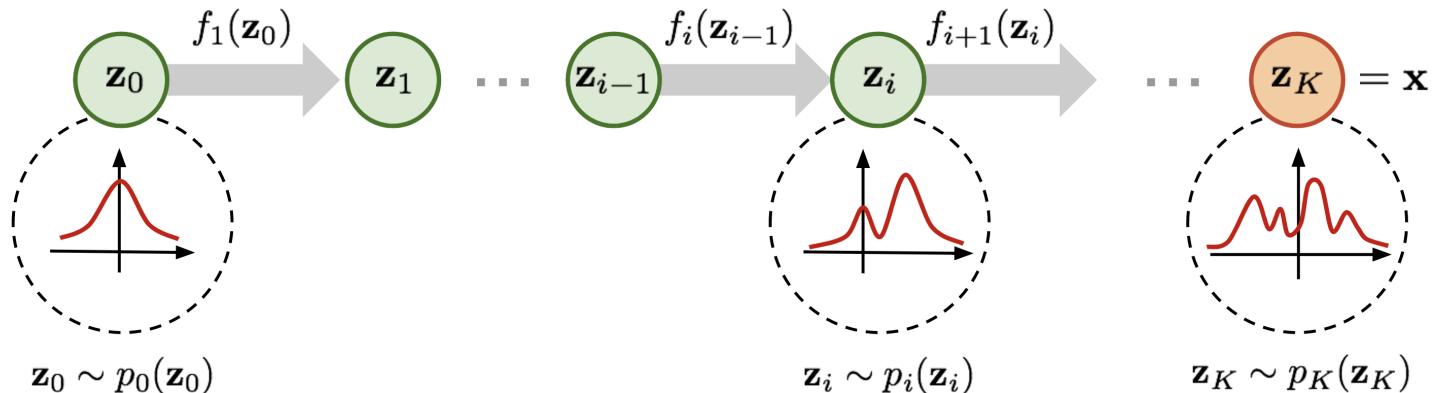
Neural Posterior Estimation (NPE)



$$\min_{q_\phi} \mathbb{E}_{p(x)} [\text{KL}(p(\theta|x) || q_\phi(\theta|x))]$$

Normalizing flows

A normalizing flow is a sequence of invertible transformations f_k that map a simple distribution p_0 to a more complex distribution p_K :

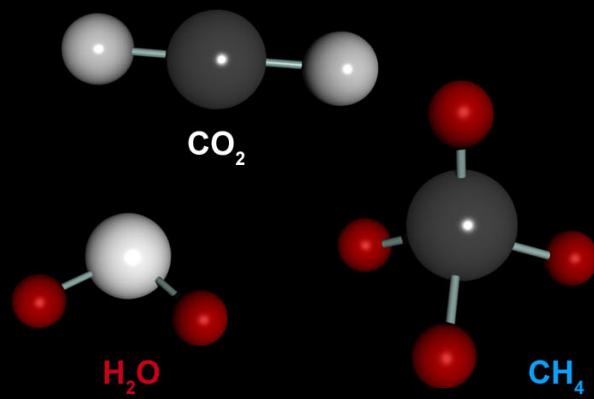
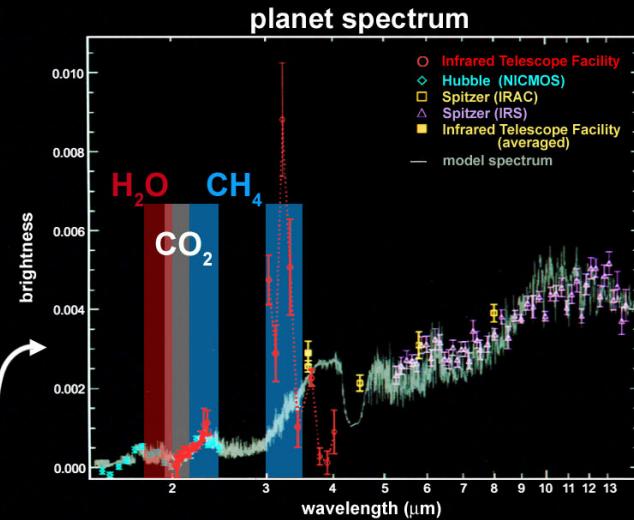
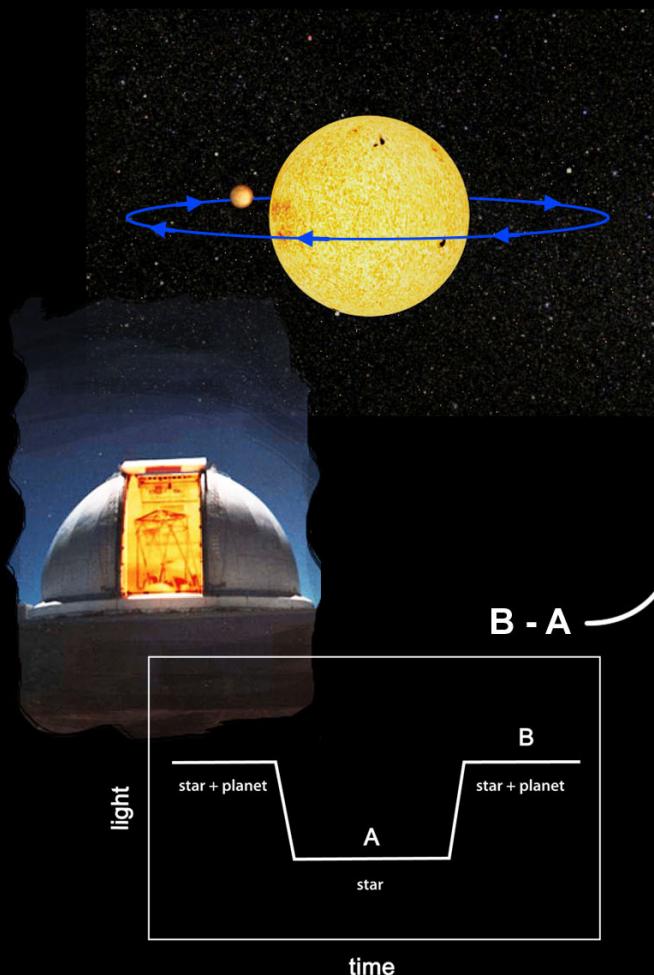


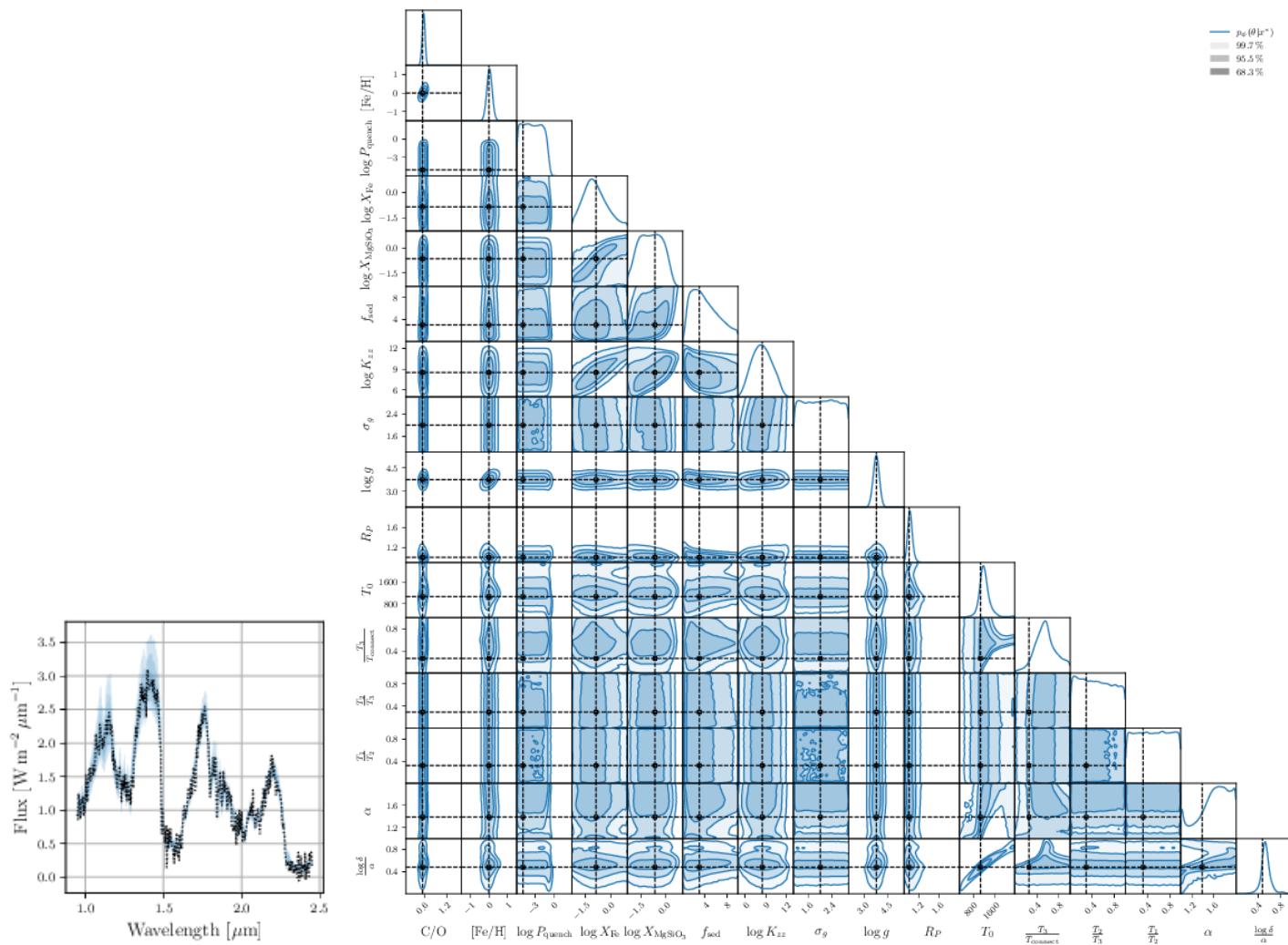
By the change of variables formula, the log-likelihood of a sample \mathbf{x} is given by

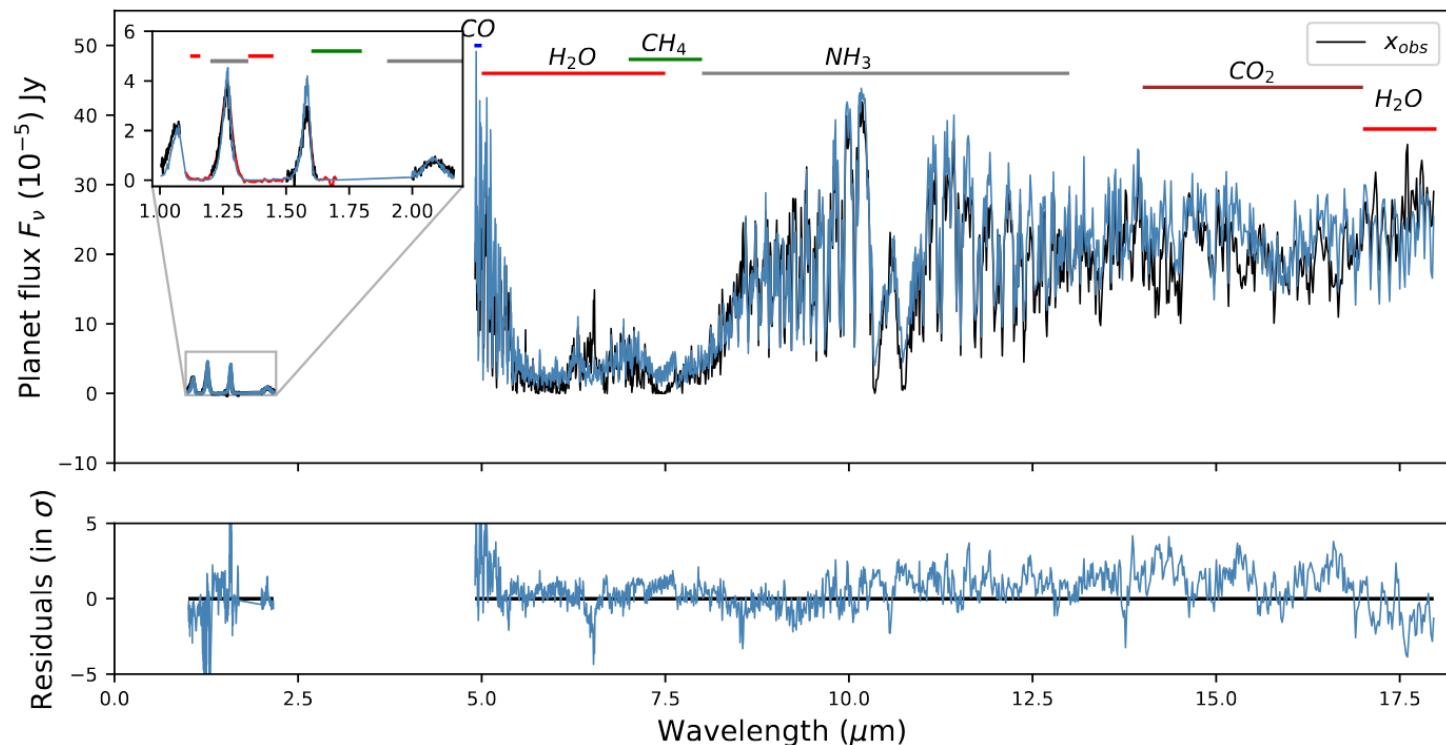
$$\log p(\mathbf{x}) = \log p(\mathbf{z}_0) - \sum_{k=1}^K \log |\det J_{f_k}(\mathbf{z}_{k-1})|.$$



Exoplanet atmosphere characterization







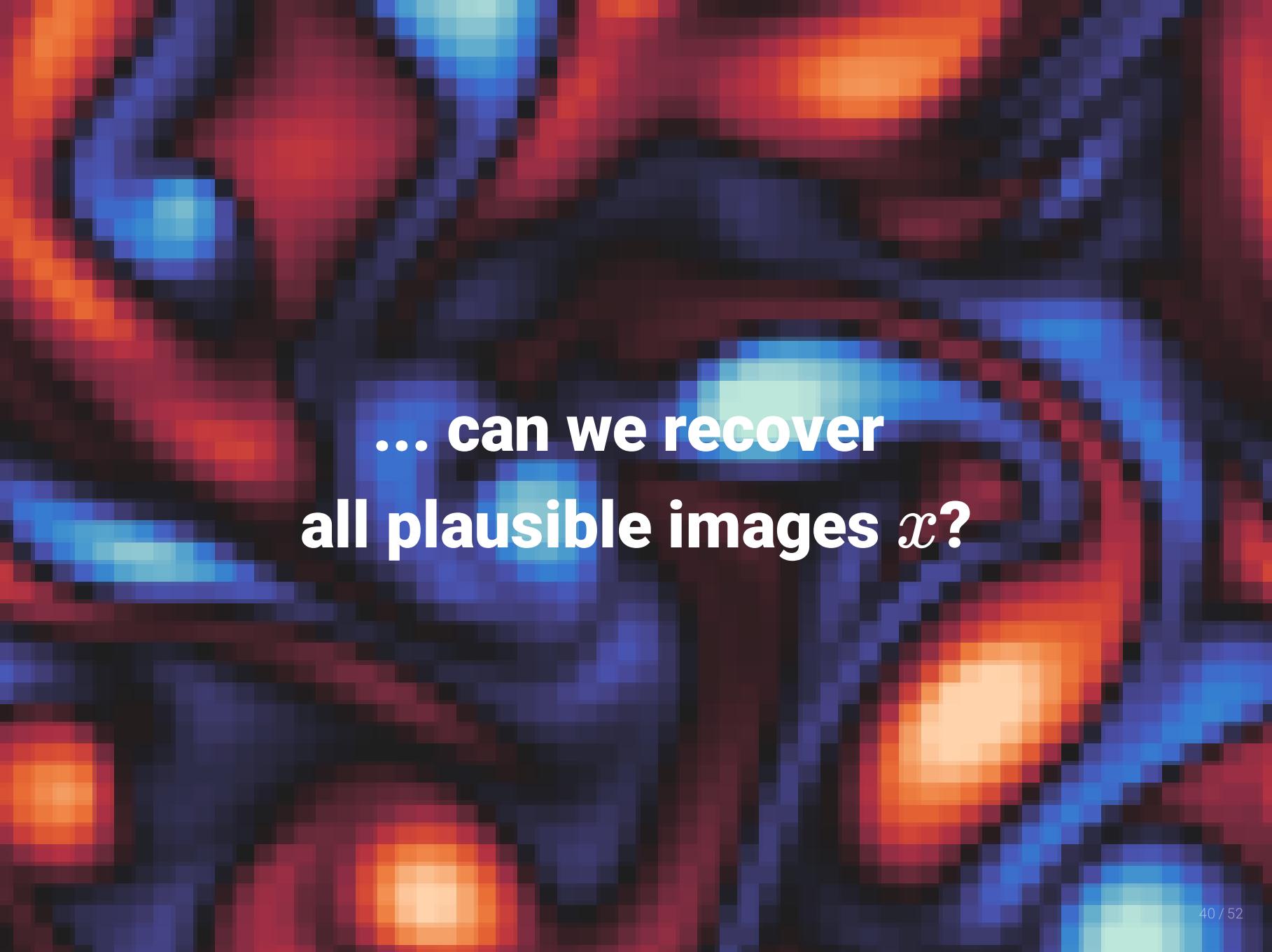
Panchromatic characterization of WISE 1738 using JWST/MIRI
(Vasist et al, submitted).

Neural score estimation (NSE)

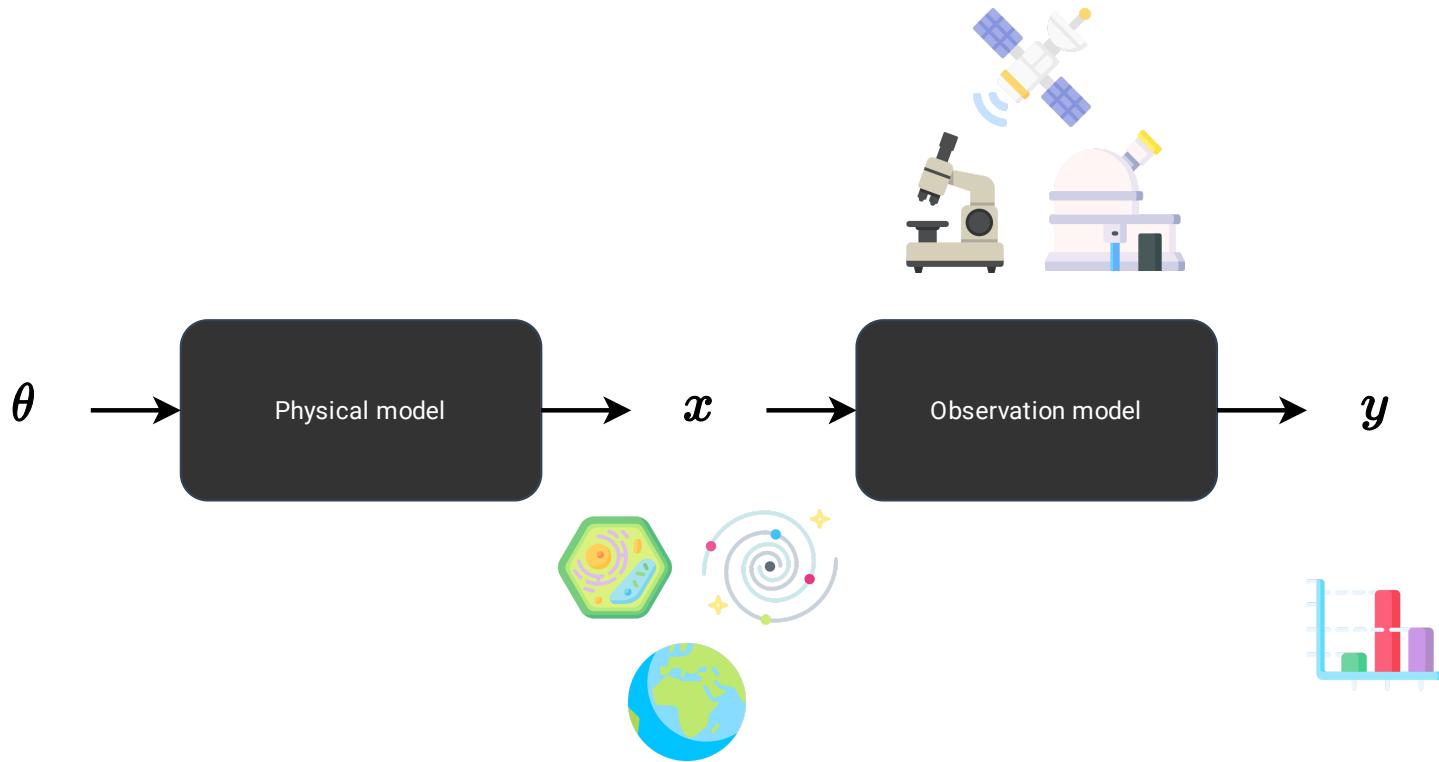


How do we estimate $p(x|y)$ when x is high-dimensional?

From a noisy observation y ...



... can we recover
all plausible images x ?



Problem statement

Given a noisy observation y , estimate the posterior distribution $p(x|y) = \frac{p(y|x)p(x)}{p(y)}$ of plausible latent states x .

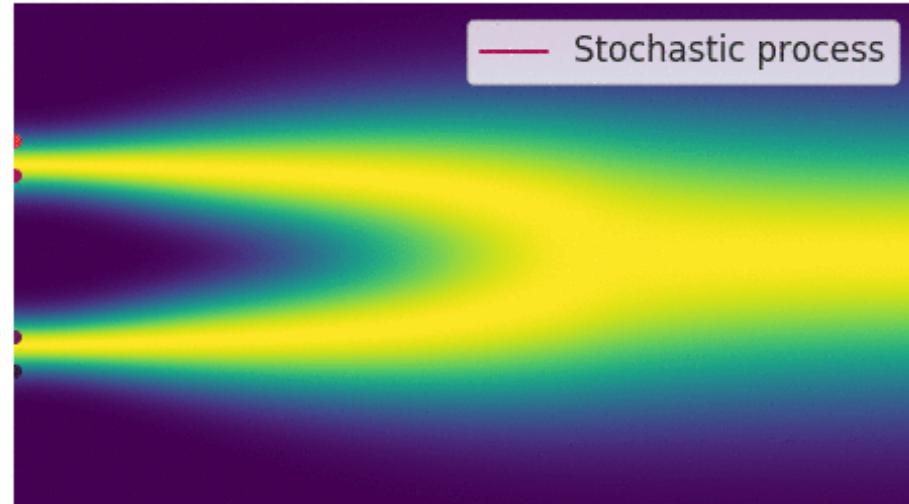
Score-based generative models 101

Samples $\mathbf{x} \sim p(\mathbf{x})$ are progressively perturbed through a diffusion process described by the forward SDE

$$d\mathbf{x}_t = f_t \mathbf{x}_t dt + g_t dw_t,$$

where \mathbf{x}_t is the perturbed sample at time t , leading to a Gaussian diffusion kernel

$$p(\mathbf{x}_t | \mathbf{x}) = \mathcal{N}(\mathbf{x}_t | \alpha_t \mathbf{x}, \Sigma_t).$$

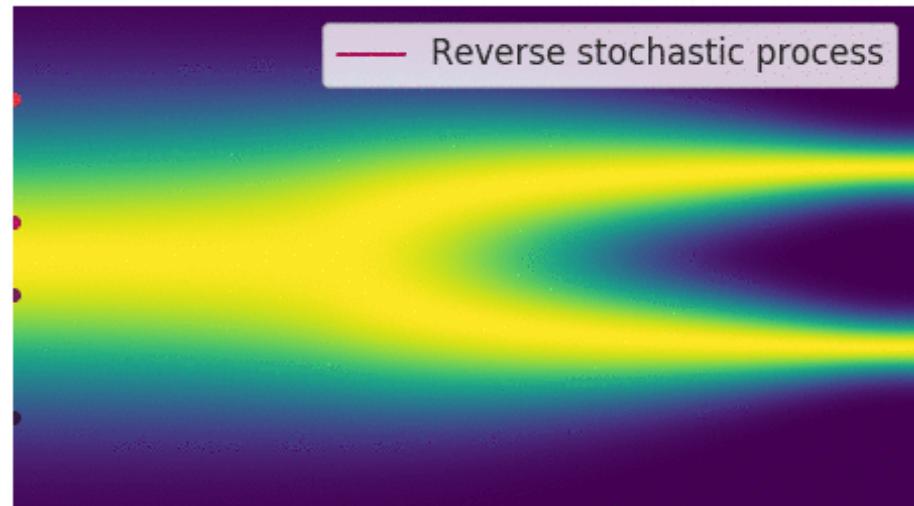


Forward diffusion process.

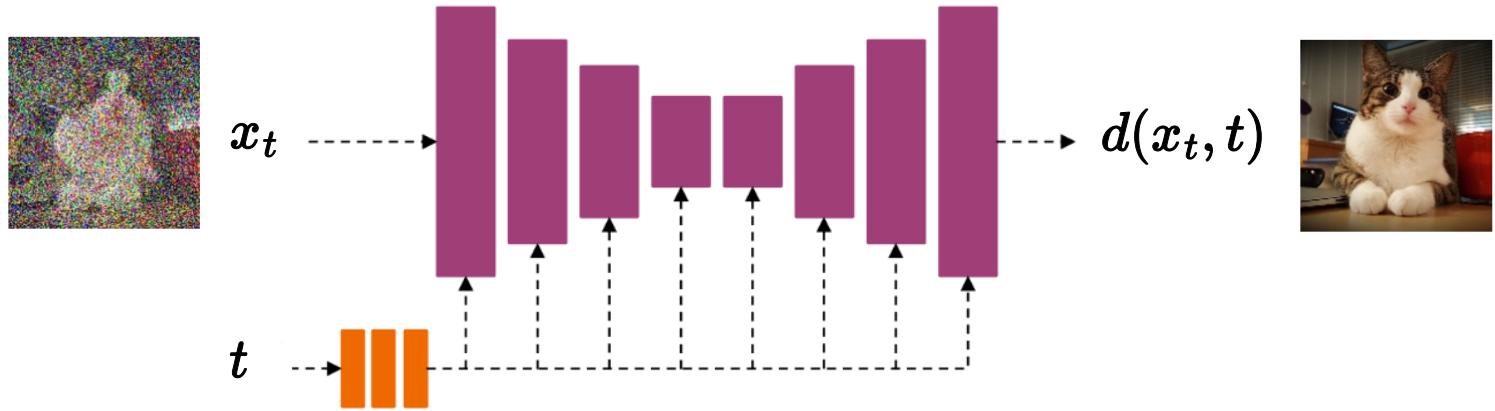
The reverse process satisfies a reverse-time SDE that can be derived analytically from the forward SDE as

$$dx_t = \left[f_t x_t - \frac{1 + \eta^2}{2} g_t^2 \nabla_{x_t} \log p(x_t) \right] dt + \eta g_t dw_t.$$

Therefore, to generate data samples $x_0 \sim p(x_0) \approx p(x)$, we can draw noise samples $x_1 \sim p(x_1) \approx \mathcal{N}(0, \Sigma_1)$ and gradually remove the noise therein by simulating the reverse SDE from $t = 1$ to 0 .



Reverse denoising process.



The score function $\nabla_{x_t} \log p(x_t)$ is unknown, but can be approximated by a neural network $d_\theta(x_t, t)$ by minimizing the denoising score matching objective

$$\mathbb{E}_{p(x)p(t)p(x_t|x)} [\lambda_t \|d_\theta(x_t, t) - x\|_2^2].$$

The optimal denoiser d_θ is the mean $\mathbb{E}[x|x_t]$ which, via Tweedie's formula, allows to use $s_\theta(x_t, t) = \Sigma_t^{-1}(d_\theta(x_t, t) - x_t)$ as a score estimate in the reverse SDE.



Inverting single observations

Because of the Bayes' rule, the posterior score $\nabla_{x_t} \log p(x_t|y)$ to inject in the reverse SDE can be decomposed as

$$\nabla_{x_t} \log p(x_t|y) = \nabla_{x_t} \log p(x_t) + \nabla_{x_t} \log p(y|x_t) - \nabla_{x_t} \log p(y).$$

This is particularly convenient as it enables **zero-shot posterior sampling** from a diffusion prior $p(x_0)$ without having to pre-wire the neural denoiser to the observation model $p(y|x)$.



Approximating $\nabla_{x_t} \log p(y|x_t)$

We want to estimate the score $\nabla_{x_t} \log p(y|x_t)$ of the noise-perturbed likelihood

$$p(y|x_t) = \int p(y|x)p(x|x_t)dx.$$

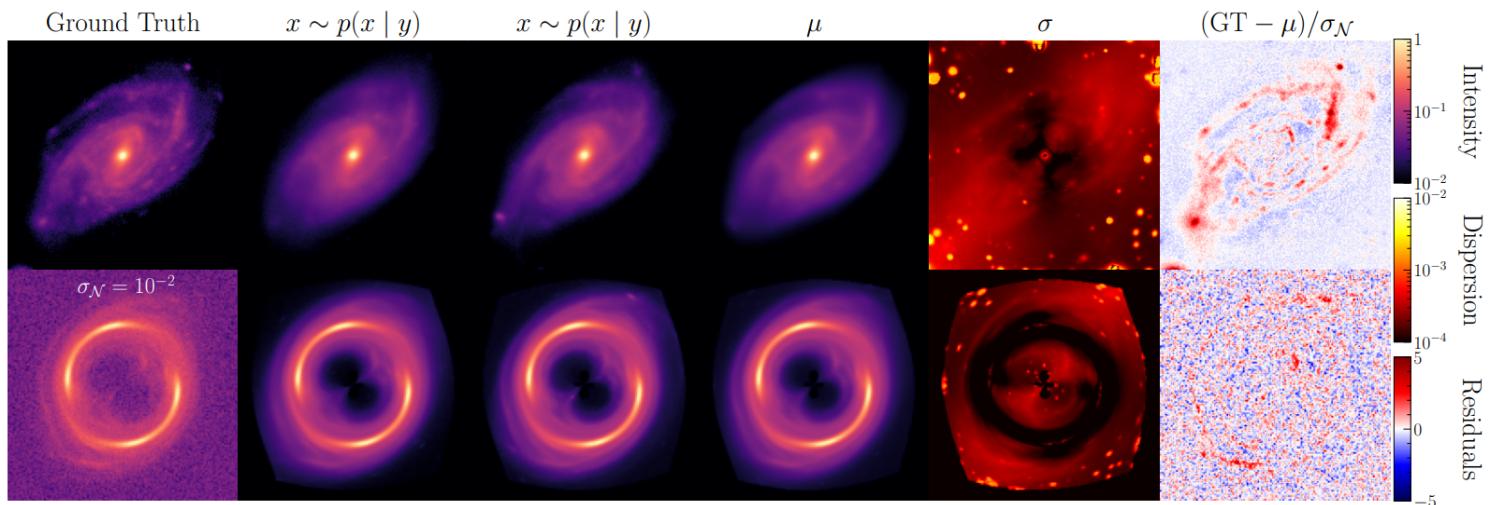
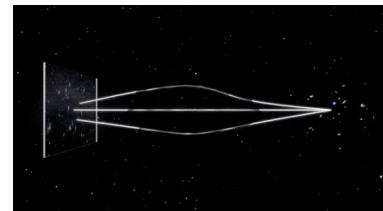
Our approach:

- Assume a linear Gaussian observation model $p(y|x) = \mathcal{N}(y|Ax, \Sigma_y)$.
- Assume the approximation $p(x|x_t) \approx \mathcal{N}(x|\mathbb{E}[x|x_t], \mathbb{V}[x|x_t])$, where $\mathbb{E}[x|x_t]$ is estimated by the denoiser and $\mathbb{V}[x|x_t]$ is estimated using Tweedie's covariance formula.
- Then $p(y|x_t) \approx \mathcal{N}(y|A\mathbb{E}[x|x_t], \Sigma_y + A\mathbb{V}[x|x_t]A^T)$.
- The score $\nabla_{x_t} \log p(y|x_t)$ then approximates to

$$\nabla_{x_t} \mathbb{E}[x|x_t]^T A^T (\Sigma_y + A\mathbb{V}[x|x_t]A^T)^{-1} (y - A\mathbb{E}[x|x_t]).$$

 x 

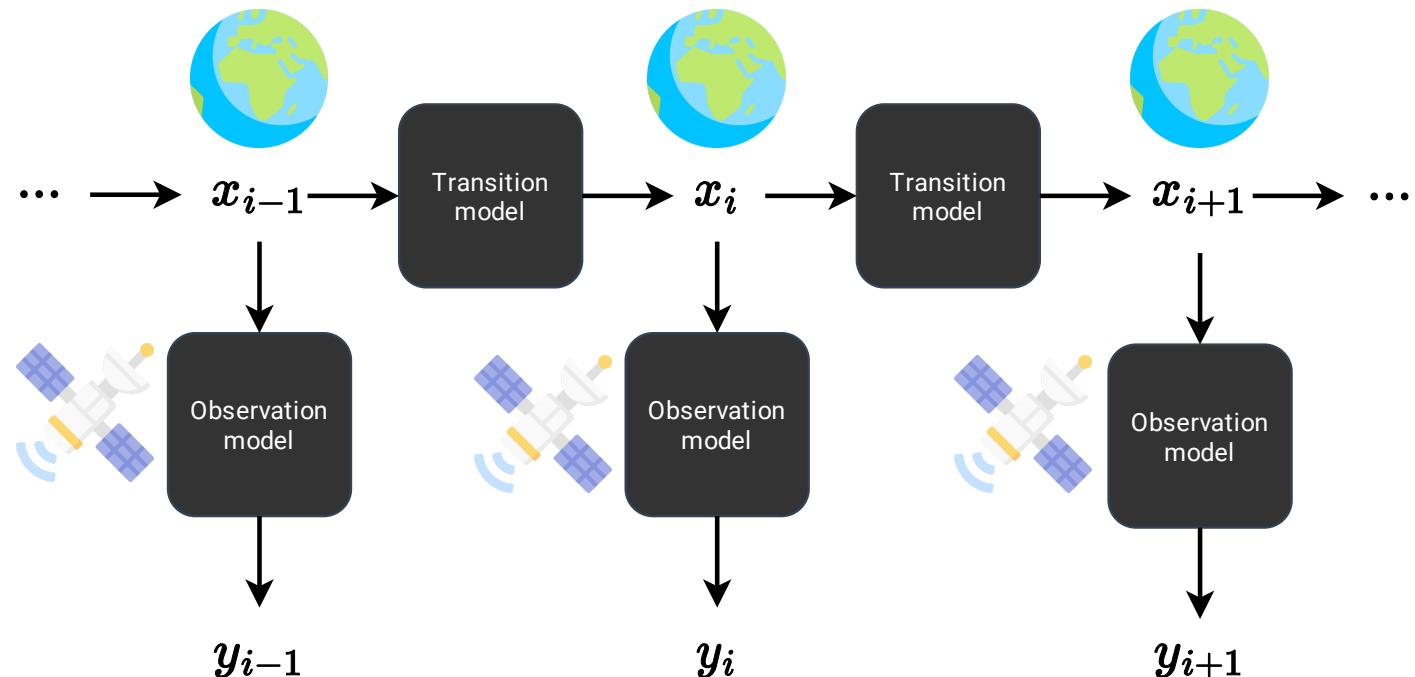
Lensing operator

 y 

Plausible galaxy images x can be recovered from lensed observations y by zero-shot posterior sampling from a diffusion prior $p(x)$.

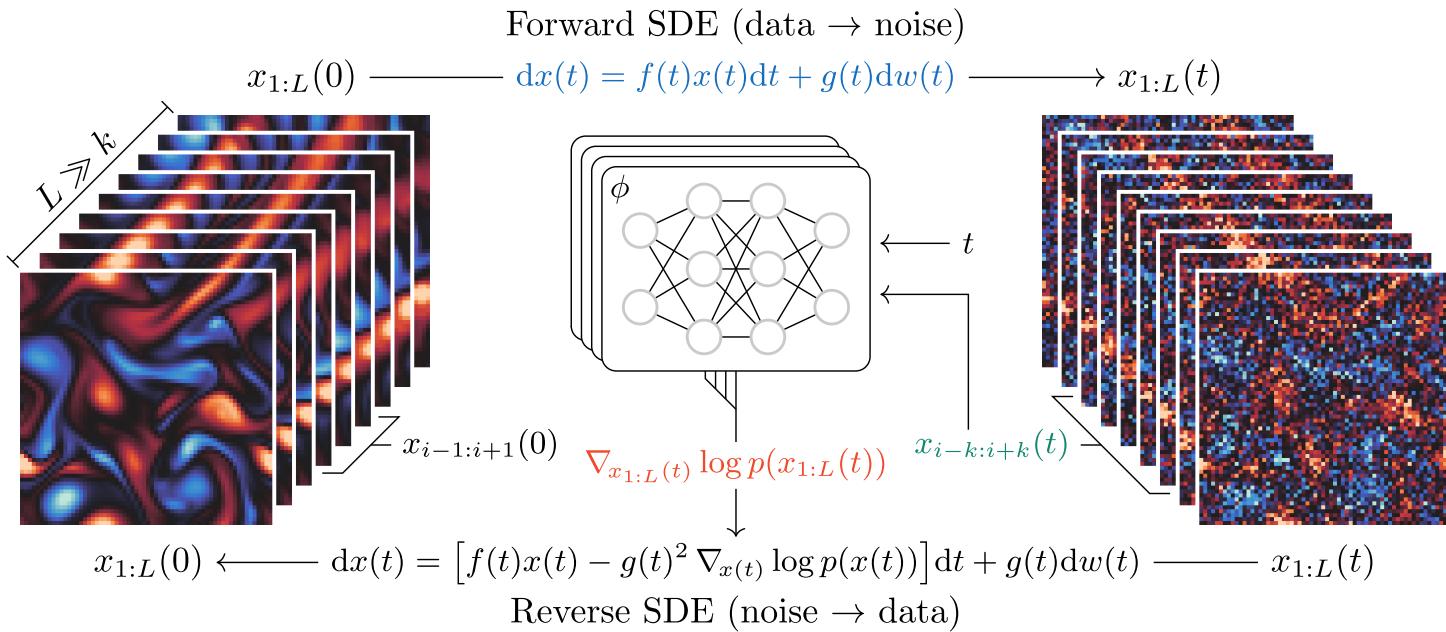


Score-based data assimilation in dynamical systems



The goal of **data assimilation** is to estimate plausible trajectories $\mathbf{x}_{1:L}$ given one or more noisy observations \mathbf{y} (or $\mathbf{y}_{1:L}$) as the posterior

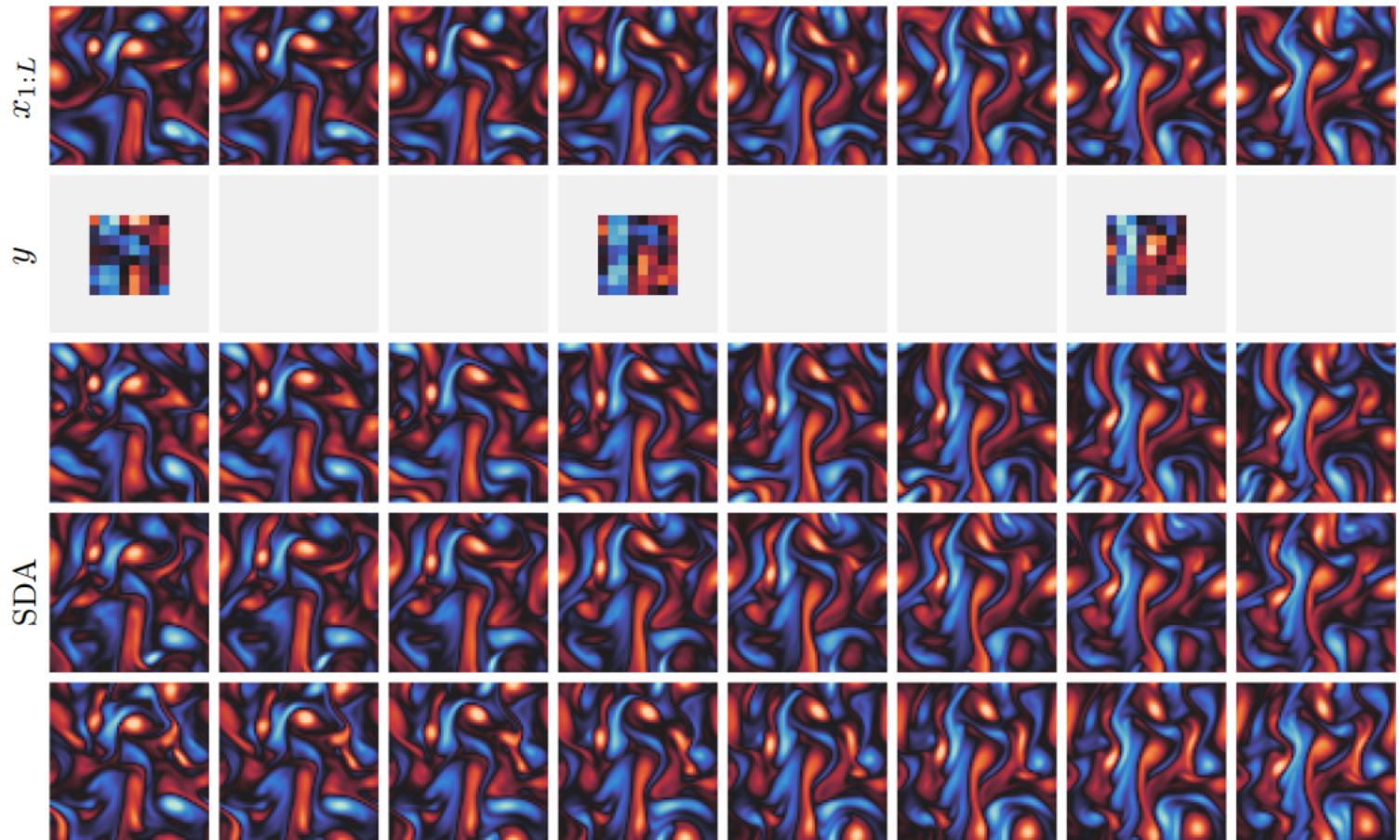
$$p(\mathbf{x}_{1:L} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{x}_{1:L})}{p(\mathbf{y})} p(\mathbf{x}_0) \prod_{i=1}^{L-1} p(\mathbf{x}_{i+1} | \mathbf{x}_i).$$



Our approach:

- Build a score-based generative model $p(\mathbf{x}_{1:L})$ of arbitrary-length trajectories*.
- Use zero-shot posterior sampling to generate plausible trajectories from noisy observations \mathbf{y} .

*:The score of a (noise perturbed) trajectory can be approximated by a sum of scores. See paper for details.



Sampling trajectories $\textcolor{blue}{x}_{1:L}$ from
noisy, incomplete and coarse-grained observations $\textcolor{blue}{y}$.

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.