

Sequential simulation-based inference for cosmological initial conditions

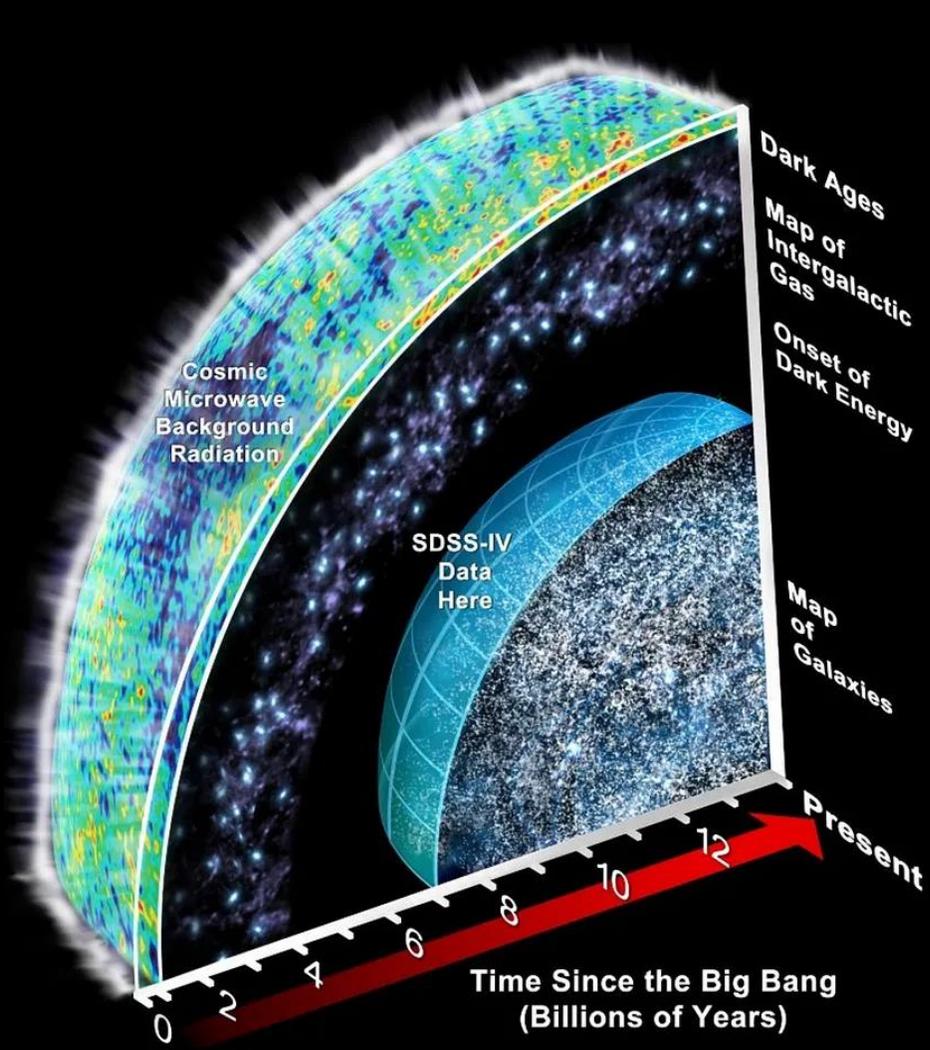
Oleg Savchenko

Work with: Guillermo Franco Abellán, Florian List, Noemi Anau Montel, Christoph Weniger

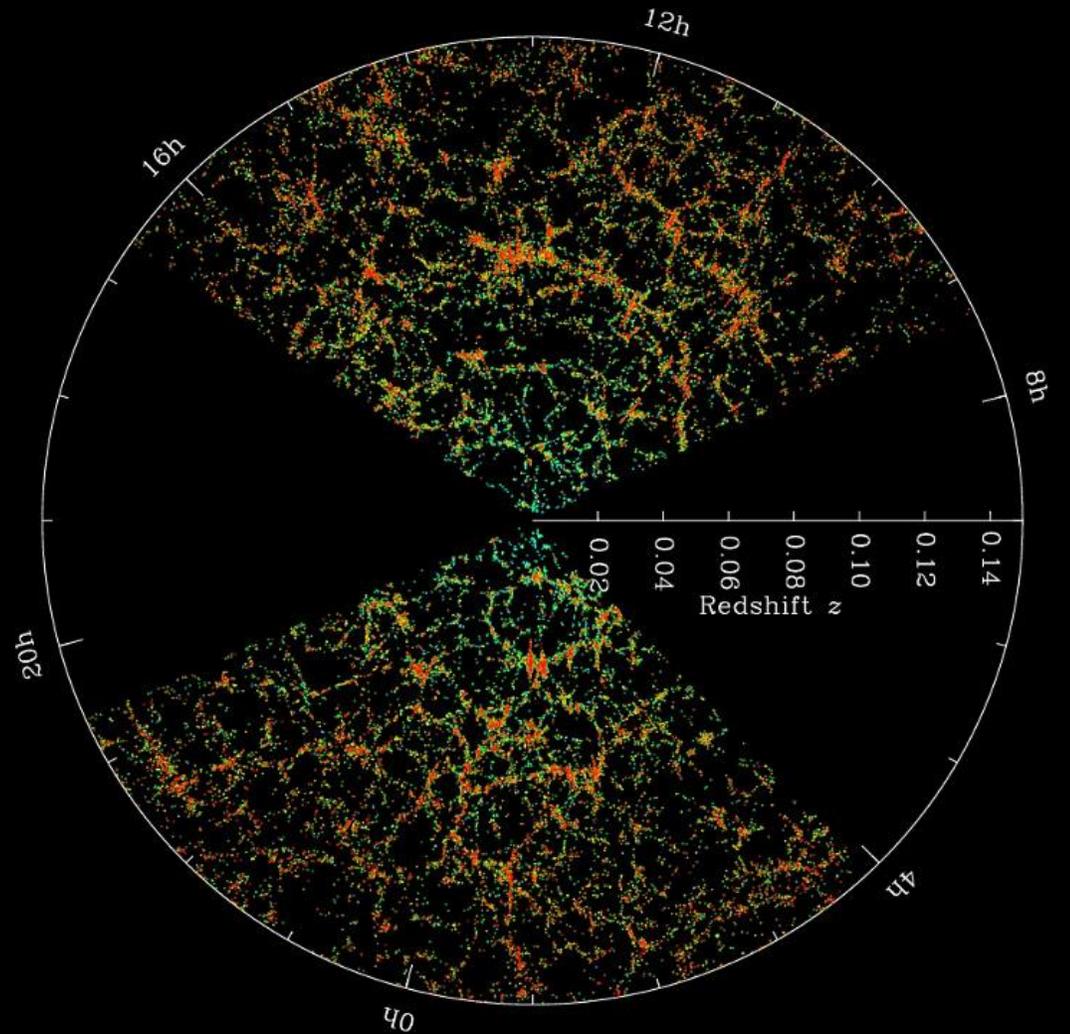


EUROPEAN AI FOR
FUNDAMENTAL PHYSICS
CONFERENCE
EuCAIFCon 2025

Large scale structure

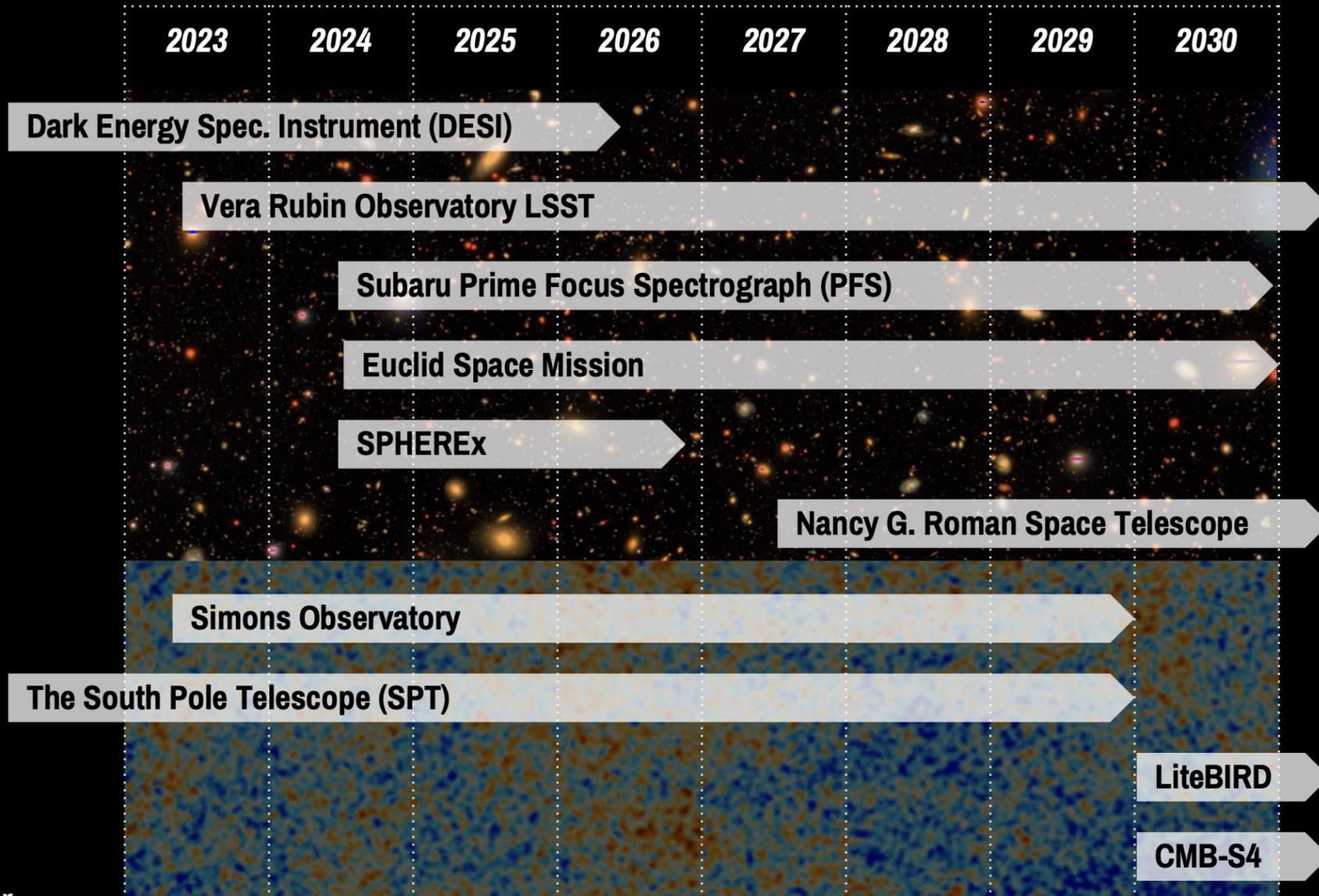


<https://mapoftheuniverse.net/>

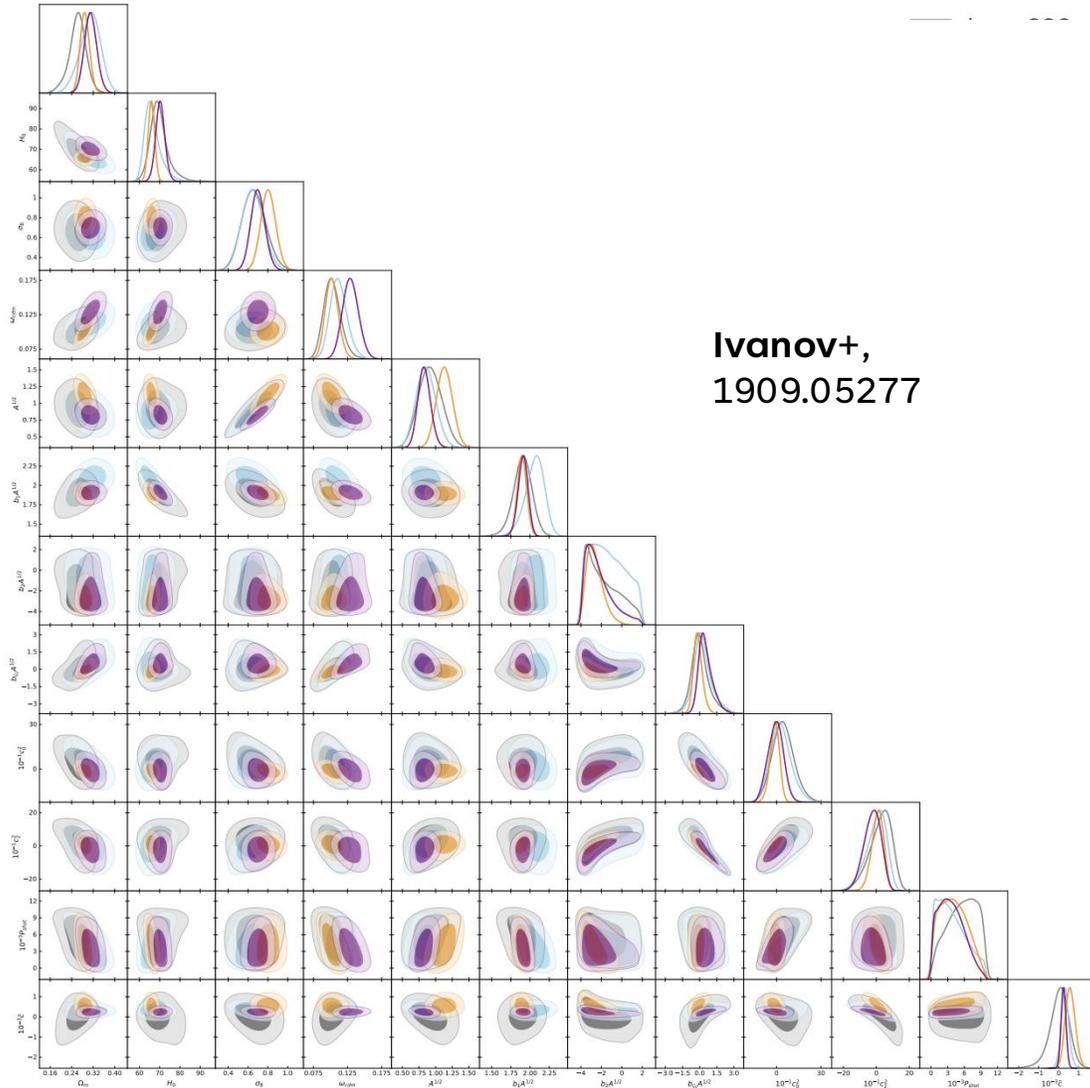


SDSS collaboration

Next decade



How to analyse LSS data?



Classic approach:

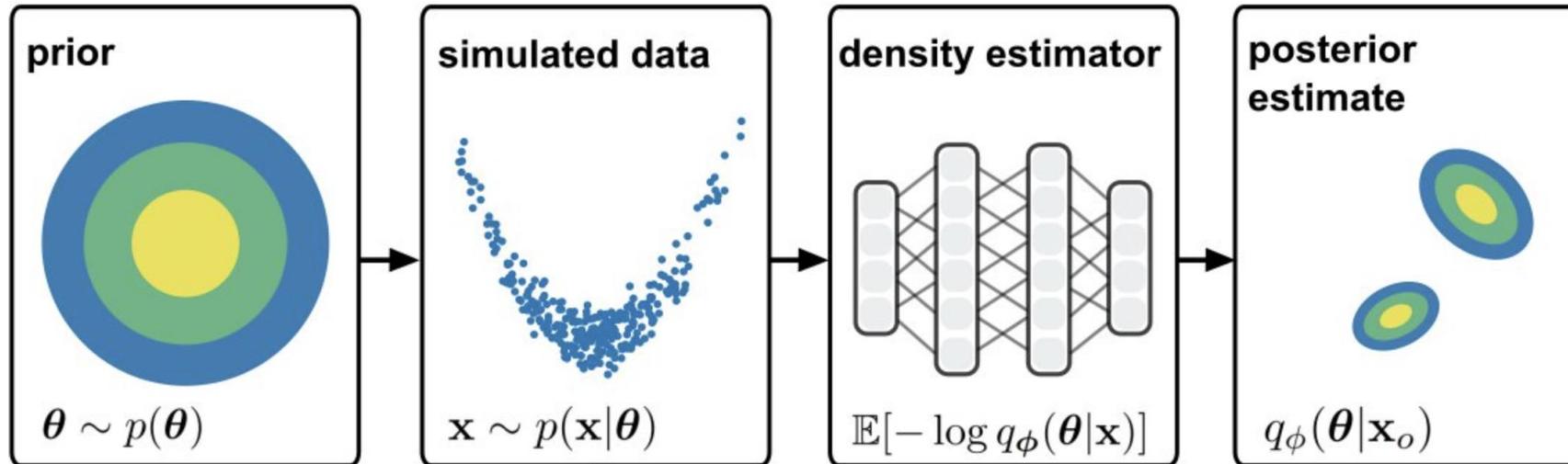
- Come up with some summary statistic $s(data)$
- Develop theory predictions for s
- Construct an analytic likelihood model (usually Gaussian)
- Use Bayes theorem and run MCMC:

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

Figure 11: The triangle plot for cosmological and nuisance parameters of four independent BOSS datasets.

Simulation-based inference

How to perform inference if all that you have is a forward simulator that can generate samples?

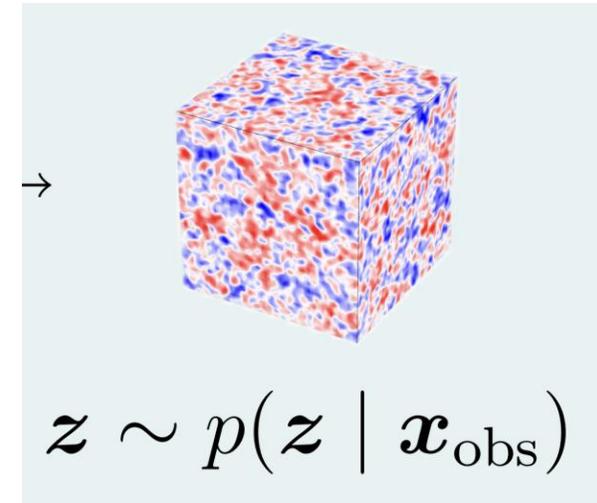
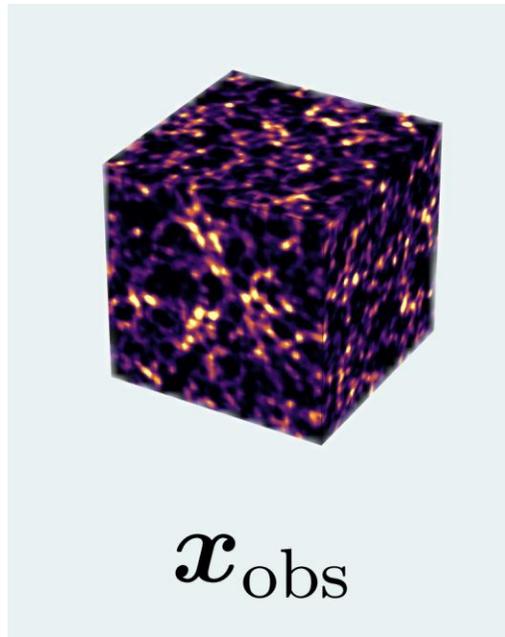


$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$

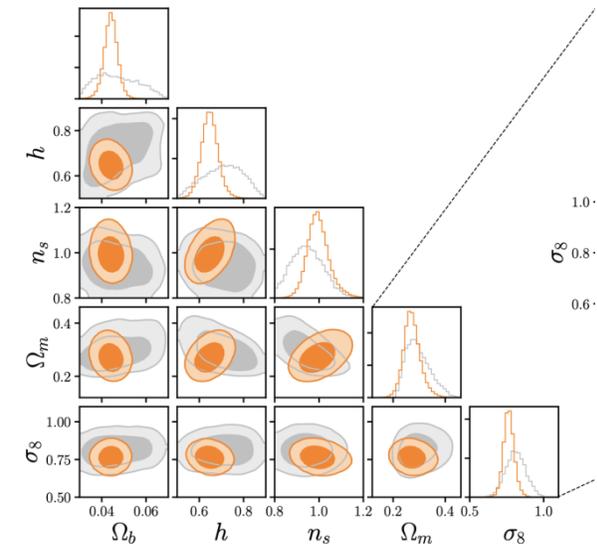
- Neural Posterior Estimation
- Neural Likelihood Estimation
- Neural Ratio Estimation

Field-level inference

The whole field contains much more information than some summary like a power spectrum!



$$p(\mathbf{x}_{\text{IC}}, \boldsymbol{\theta} | \mathbf{x}_{\text{obs}}) = ?$$



Initial conditions reconstruction

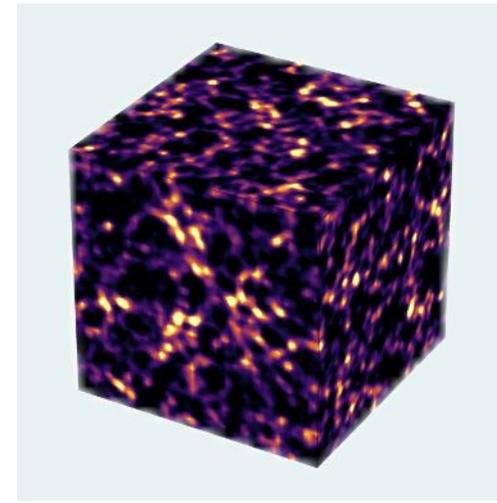
Very early universe had very simple properties!

→ feasible way to do field reconstruction is to infer these initial conditions 😊

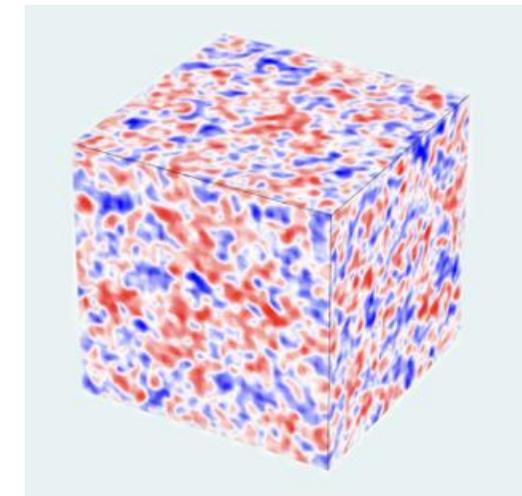
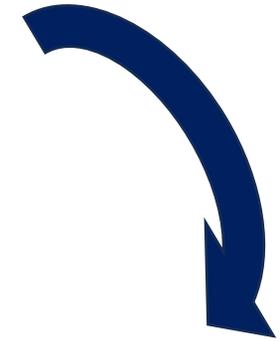
$$P(\boldsymbol{\delta}^{\text{IC}} | \{N_i^g\}) = \frac{P(\boldsymbol{\delta}^{\text{IC}})P(\{N_i^g\} | G_i(\boldsymbol{\delta}^{\text{IC}}))}{P(\{N_i^g\})}$$

$\boldsymbol{\delta}^{\text{IC}}$ - initial density field

$\{N_i^g\}$ - galaxy catalog data



$z = 0$

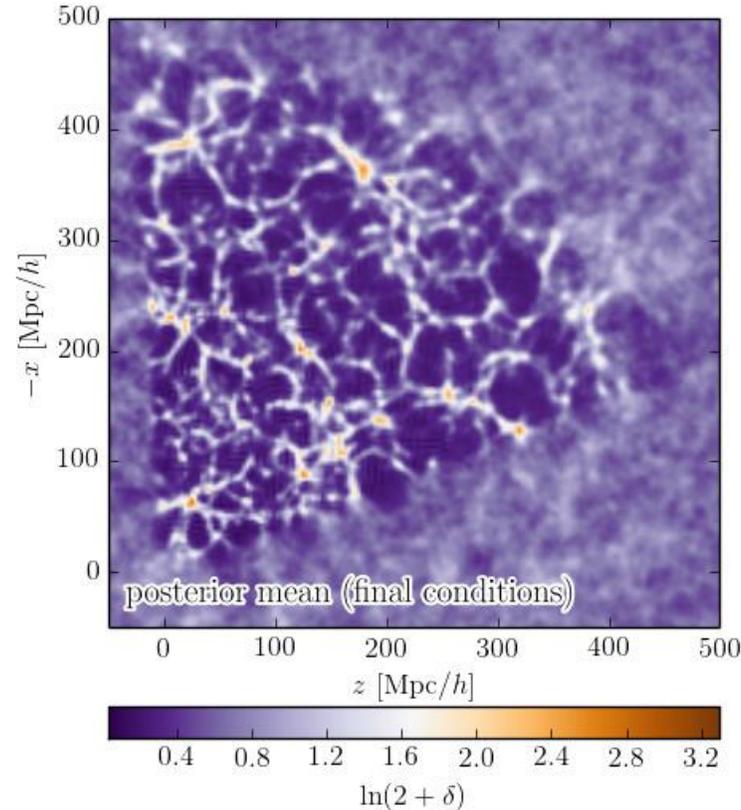
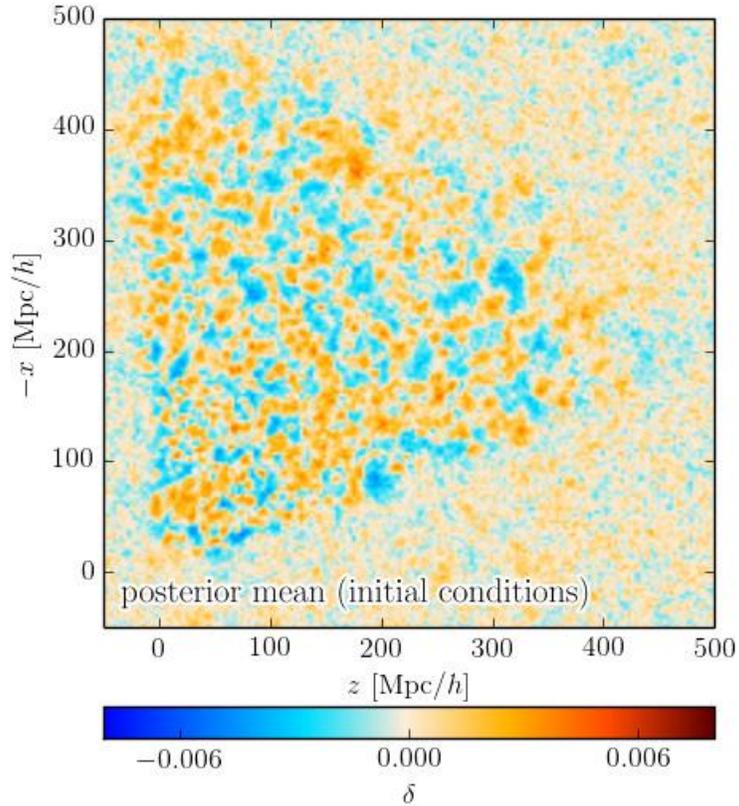


$z = 1000$

Gaussian initial conditions

Hamiltonian Monte Carlo

Jasche+, 1203.3639, 1806.11117



- Most developed method so far
- Explicit likelihood
- Requires gradients from the simulator
- Computationally expensive

Samples produced with **Bayesian Origin Reconstruction from Galaxies (BORG)** algorithm

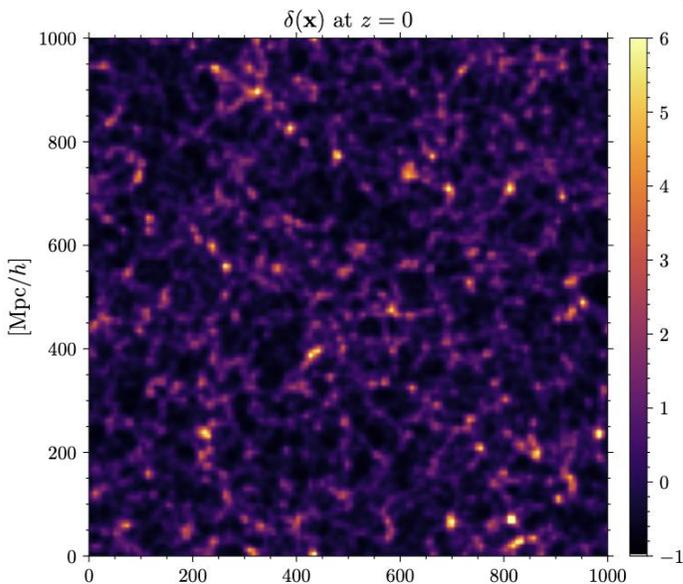
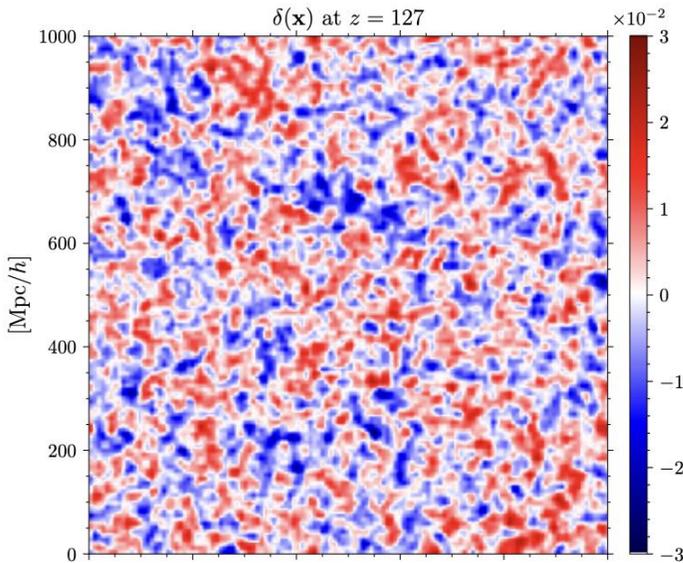
Our setting

128³ resolution: ~million-dimensional parameter space!

- Want to explore the full posterior, not only get point estimates
- Want to keep things as simple as possible: model the likelihood as

$$p(\mathbf{x}|\mathbf{z}) \propto \exp\left\{-\frac{1}{2}(\mathbf{z} - \hat{\mathbf{z}}_{\theta}(\mathbf{x}))^T \mathbf{Q}_{\theta}^L (\mathbf{z} - \hat{\mathbf{z}}_{\theta}(\mathbf{x}))\right\}$$

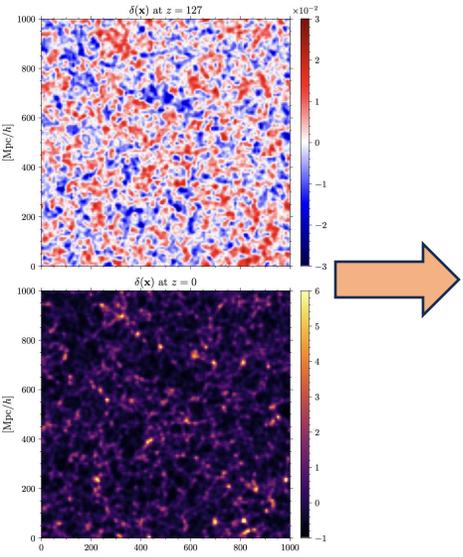
- Trainable parts of the model: estimator $\boldsymbol{\mu}_{MAP}(\mathbf{x})$ and Q matrix Fourier diagonal values



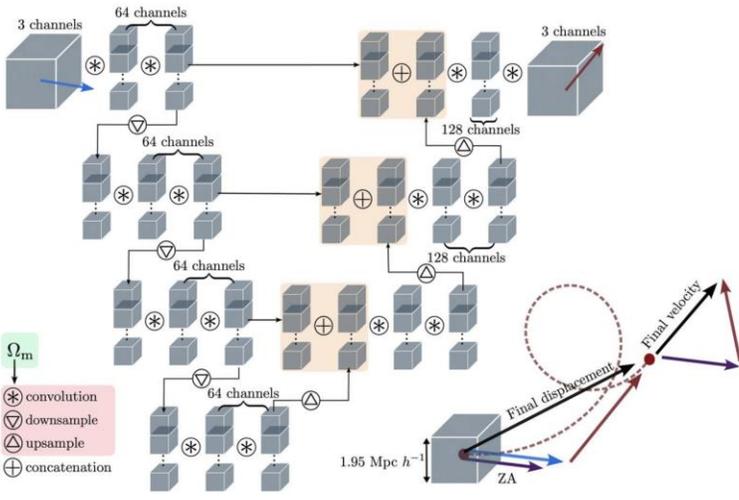
Training data: 2000
128³ (1 Gpc/h)³ Quijote
N-body simulations

Apply a U-Net to estimate $\mu_{MAP}(x)$

Our approach



Training data



map2map U-Net
Jamieson+, 2206.04594

Train the model with a simple MAP loss function

$$\mathcal{L} = \frac{1}{2} \sum_{i=1}^N \{ (z_i - \hat{\mu}_{\theta}(x_i))^T Q_{\theta} (z_i - \hat{\mu}_{\theta}(x_i)) \} - \frac{N}{2} \text{tr} \log Q_{\theta}$$

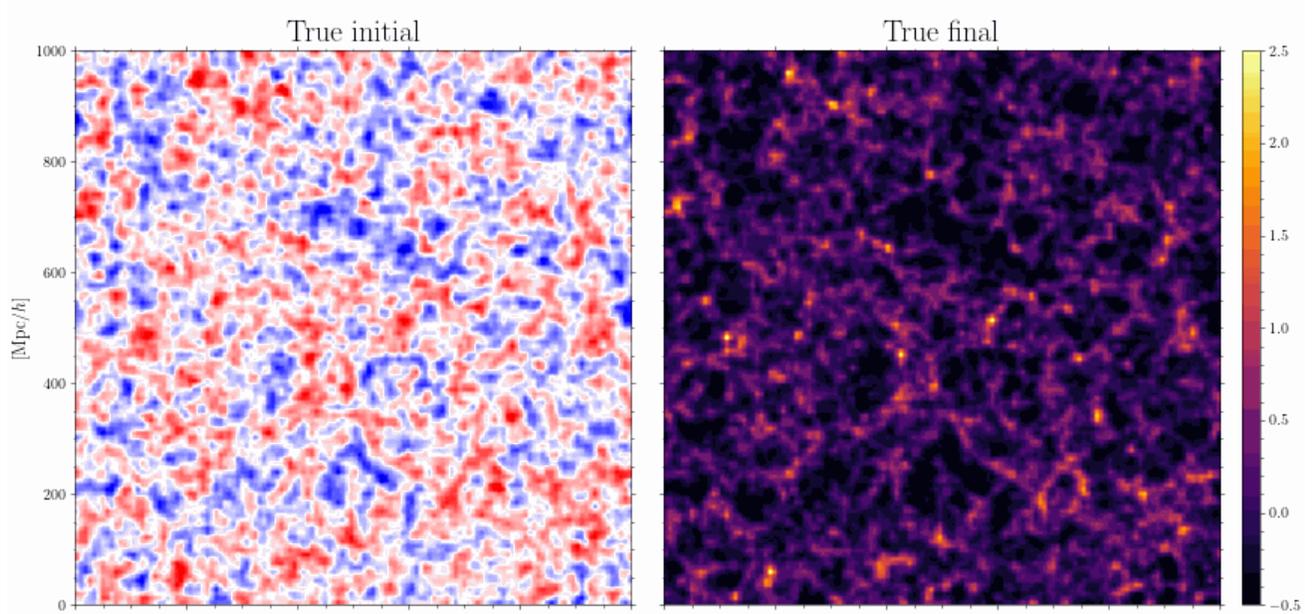
and a Fourier-diagonal Q matrix

Super-fast sampling from a Gaussian

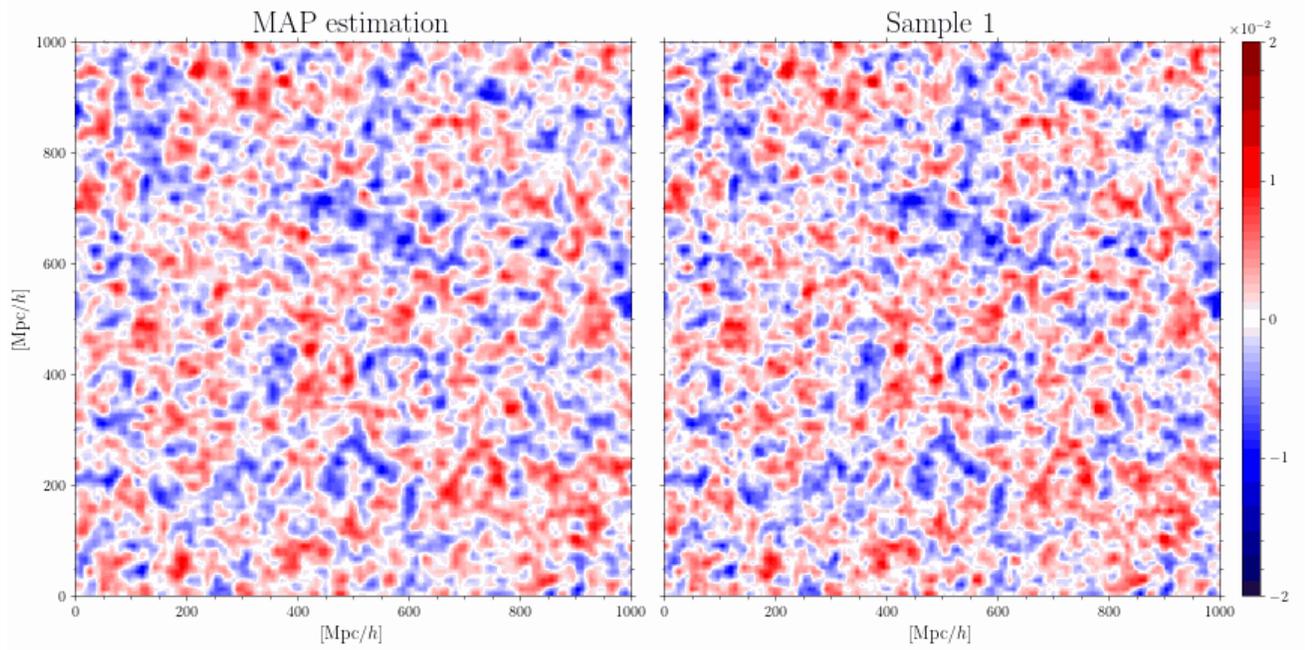
Learn precision matrix and the embedding network simultaneously

Results

True



Inferred



$V = (1 \text{ Gpc}/h)^3$

$N_{\text{grid}} = 128^3$

1.5 hrs of training

on 1 GPU

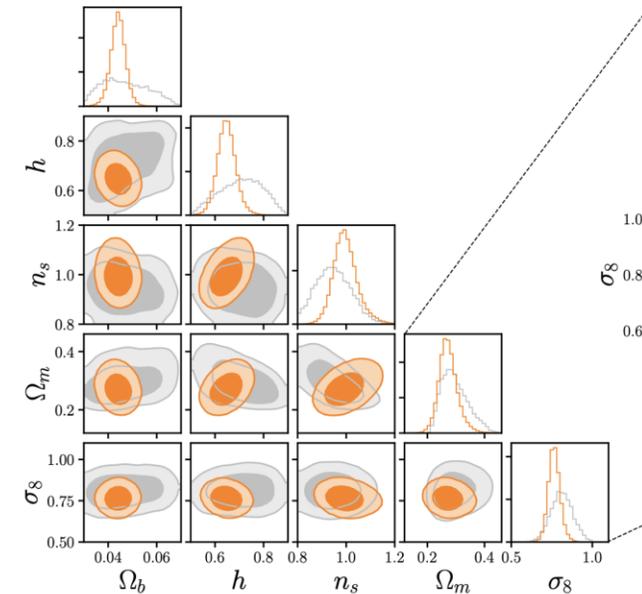
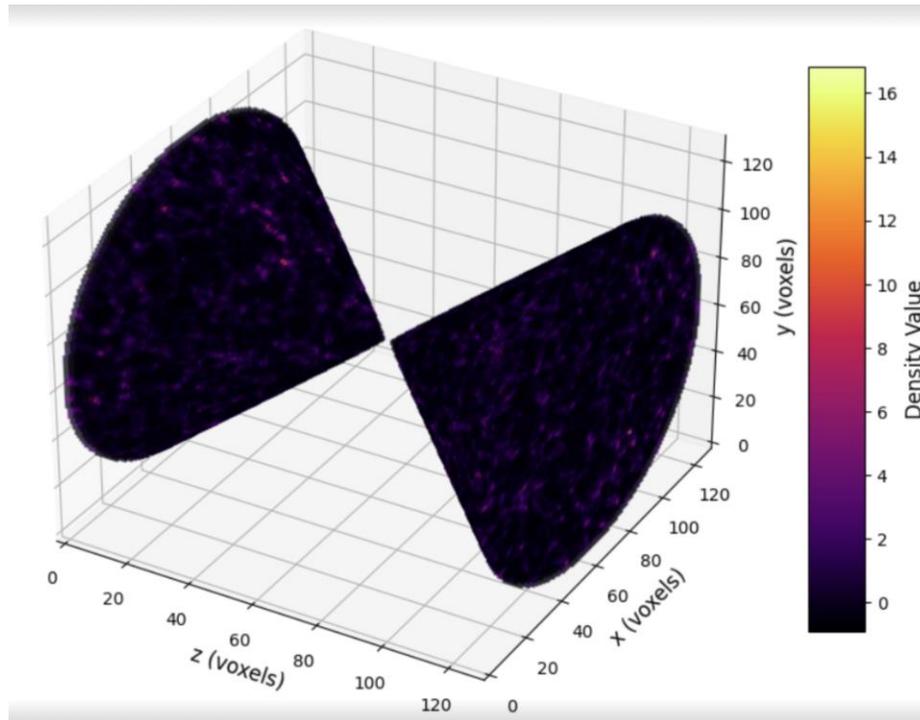
NVIDIA 40GB A100

10^3 samples in $< 3 \text{ s}$

Next steps

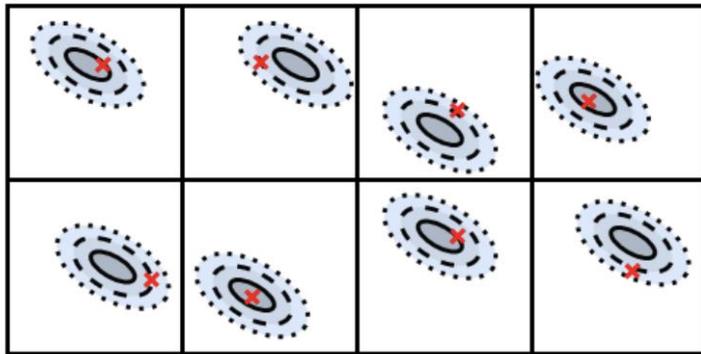
Moving to a more realistic setting:

- Survey effects: mask, spatially varying noise, etc.
- Varying cosmological parameters



Amortized vs sequential inference

- Fast inference for new \mathbf{x}_{obs}
- High simulation cost for training
- Can be not very precise

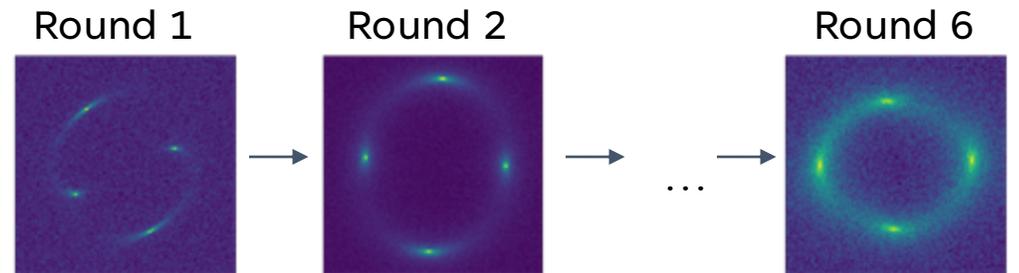


Cole+, 2111.08030

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

simultaneously estimates the posteriors
for all simulated observations

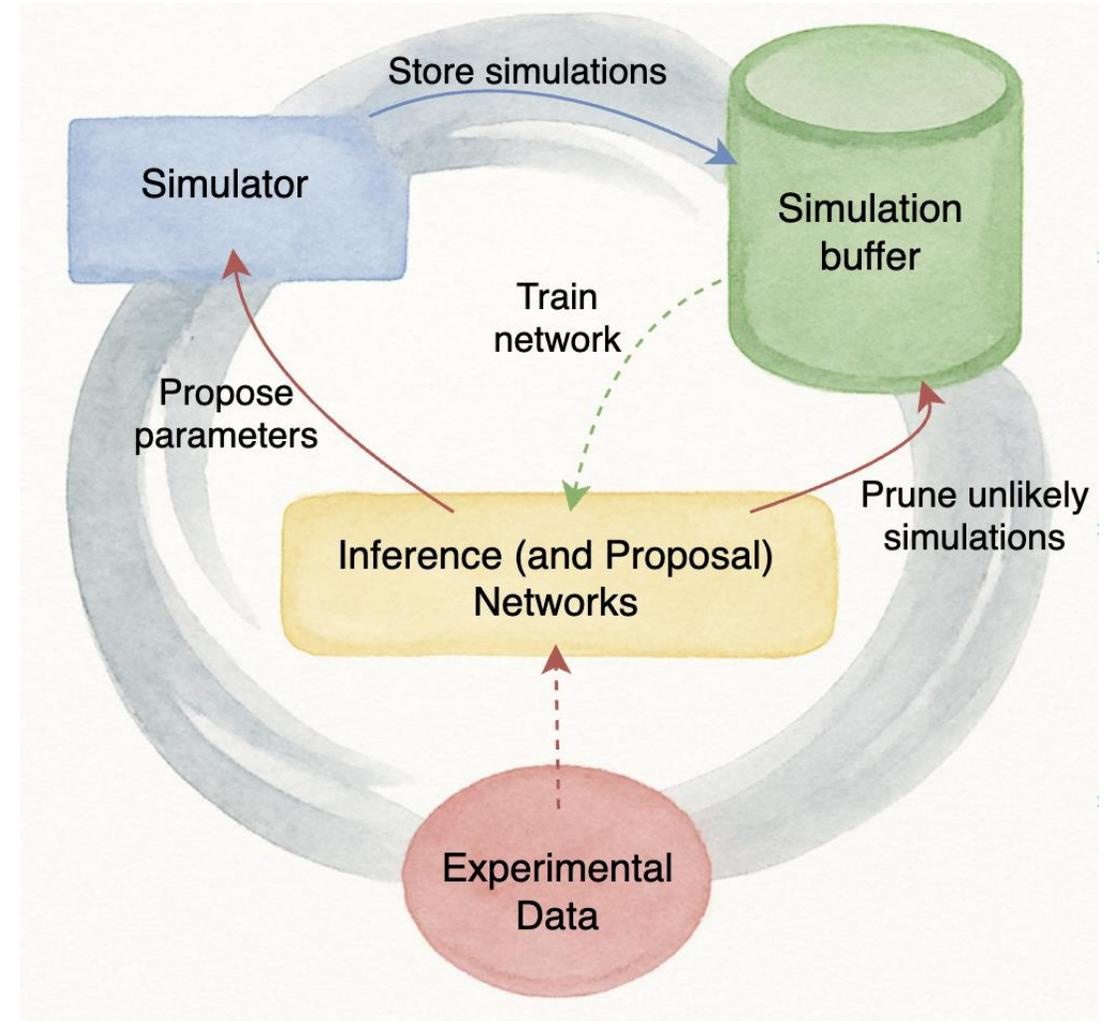
- Slow inference for new \mathbf{x}_{obs}
- Low simulation cost
- High precision



Animation credit:
N. Anau Montel

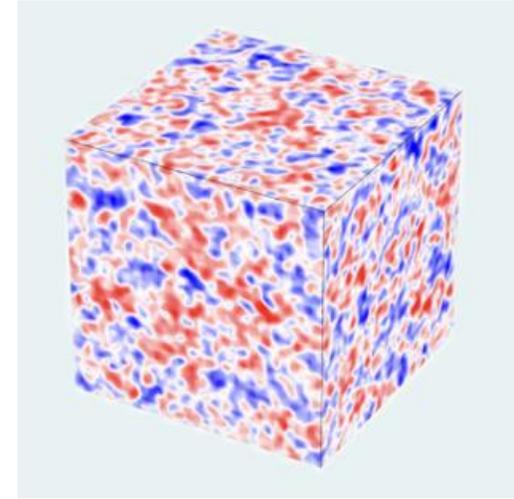
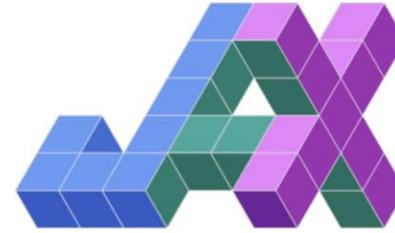
Sequential inference in Falcon

- Framework for **distributed computing** on multiple nodes
- Training set evolves during training for sequential inference: new samples are drawn from the **proposal distribution**
- Some nodes **simulate**, some nodes train to do **inference**



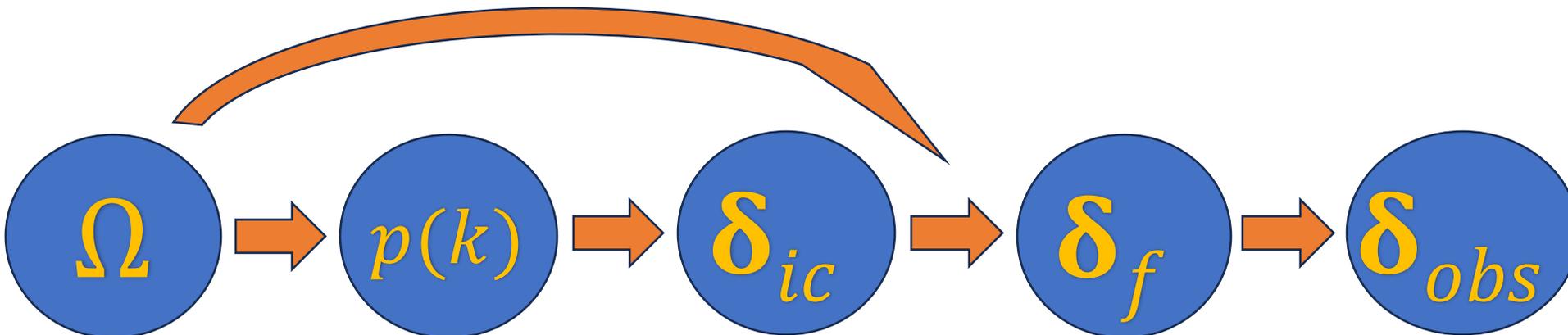
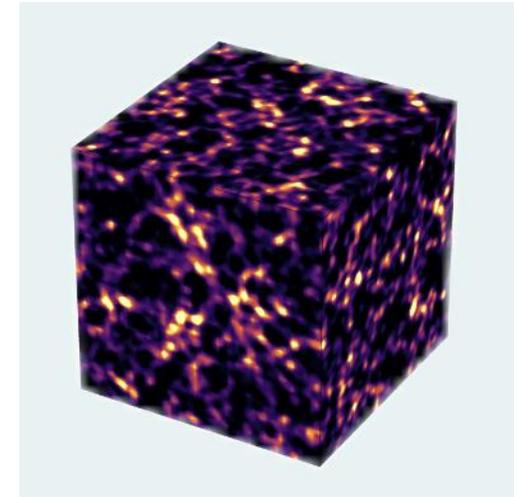
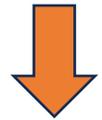
Forward model

Code developed by F. List, O. Hahn et al.



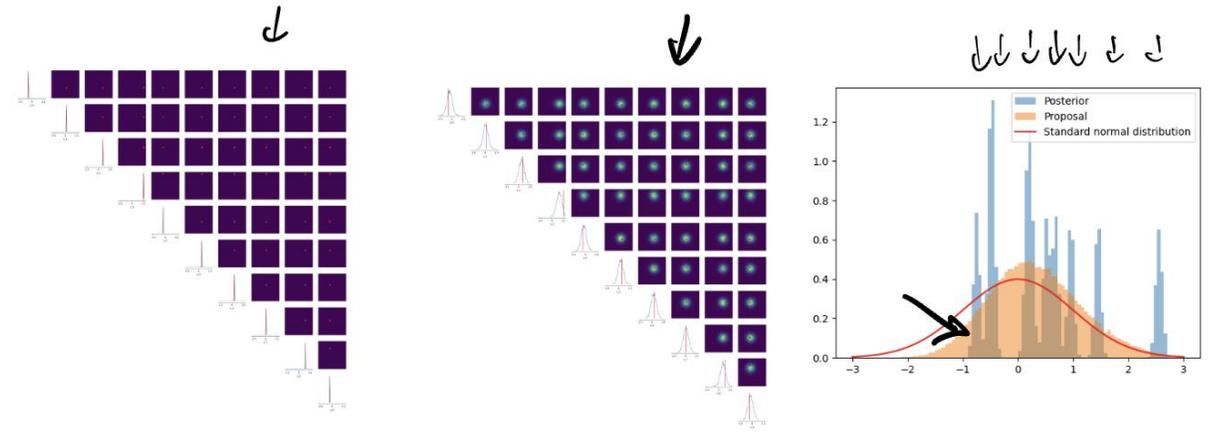
Differentiable, highly-parallelisable particle-mesh simulator written in JAX.

Hierarchical model implemented in Falcon:

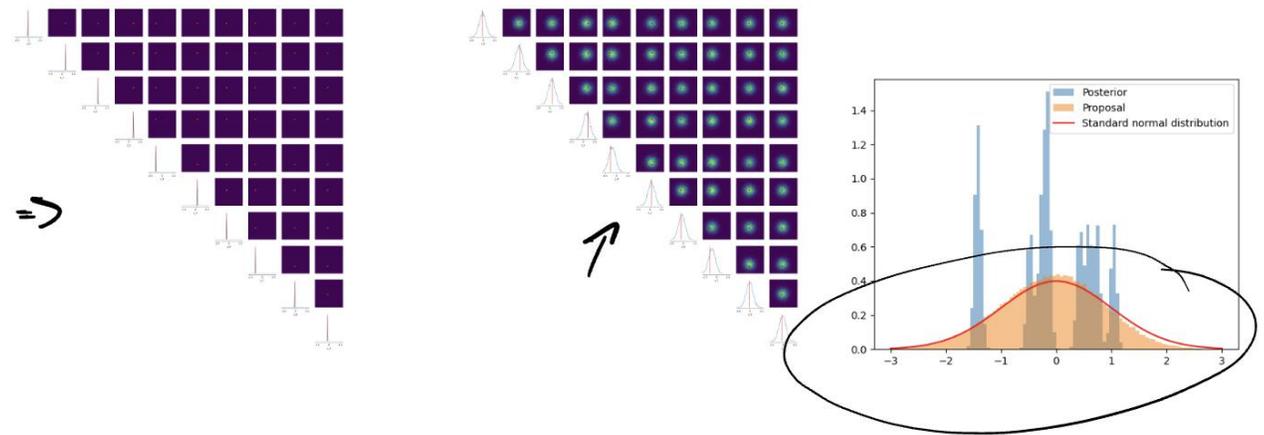


Sequential inference in Falcon

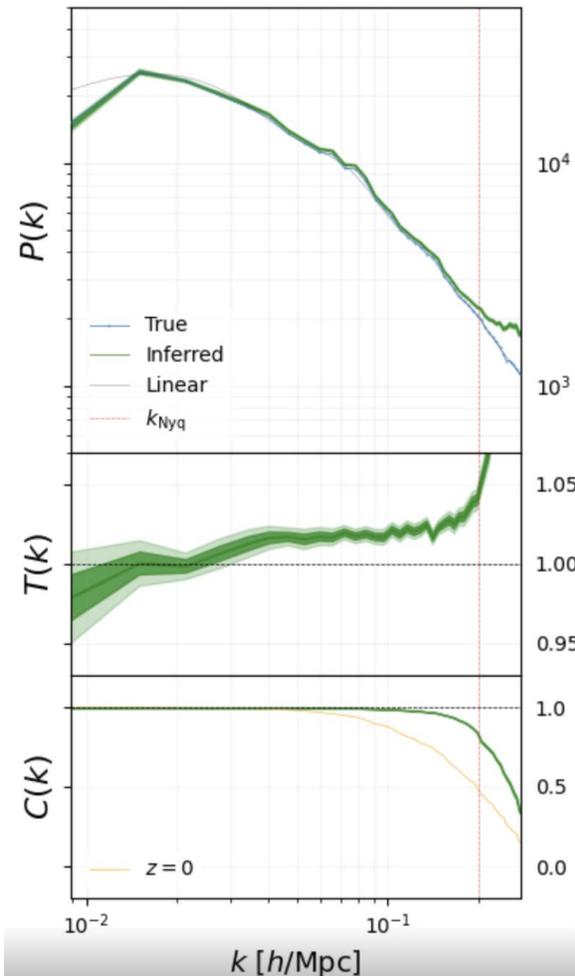
Why is this hard in high dimensions? Naïve way of tempering the likelihood fails



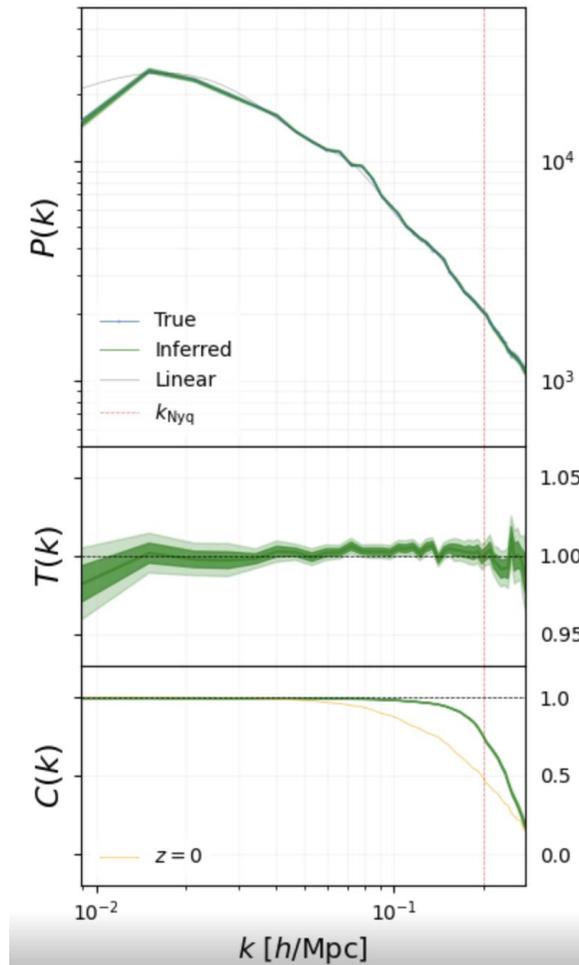
Raising likelihood to some power produces samples with incorrect summaries



Tempering the likelihood

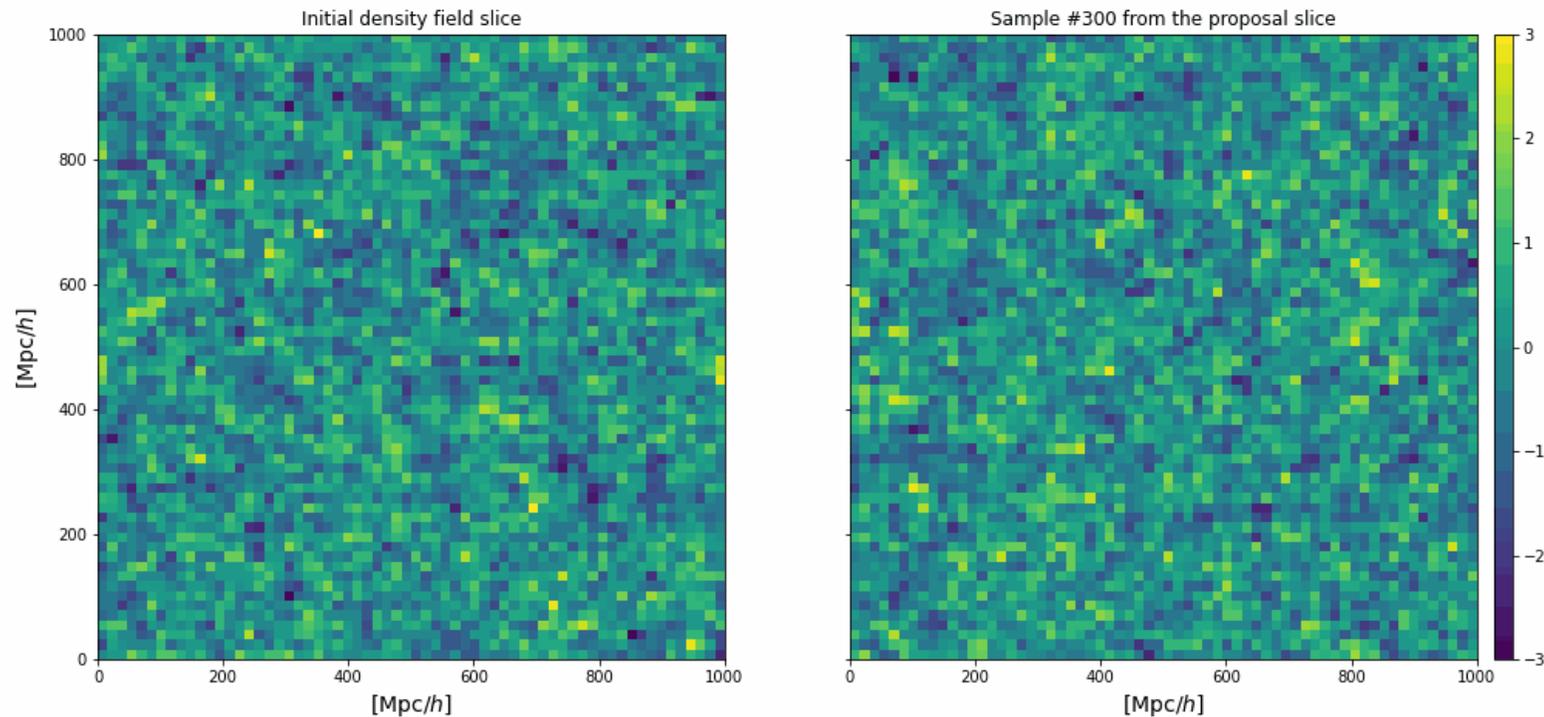


VS



- Need to have a **proposal distribution**
- Can temper the likelihood by adding more noise to the data

Tempering the likelihood

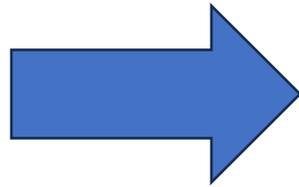
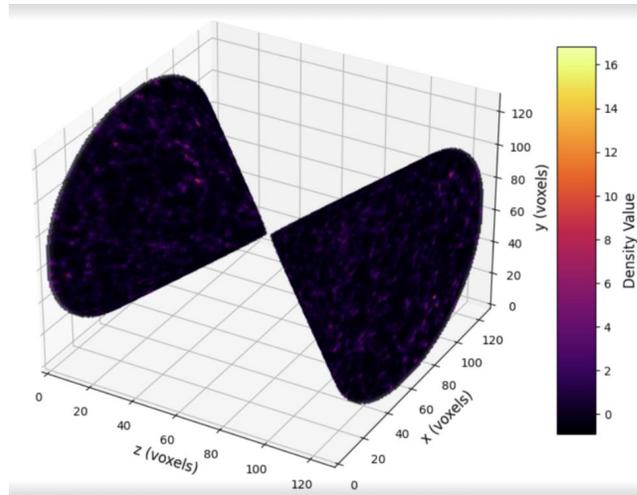


- Need to have a **proposal distribution**
- Can temper the likelihood by adding specific noise to the proposal data:

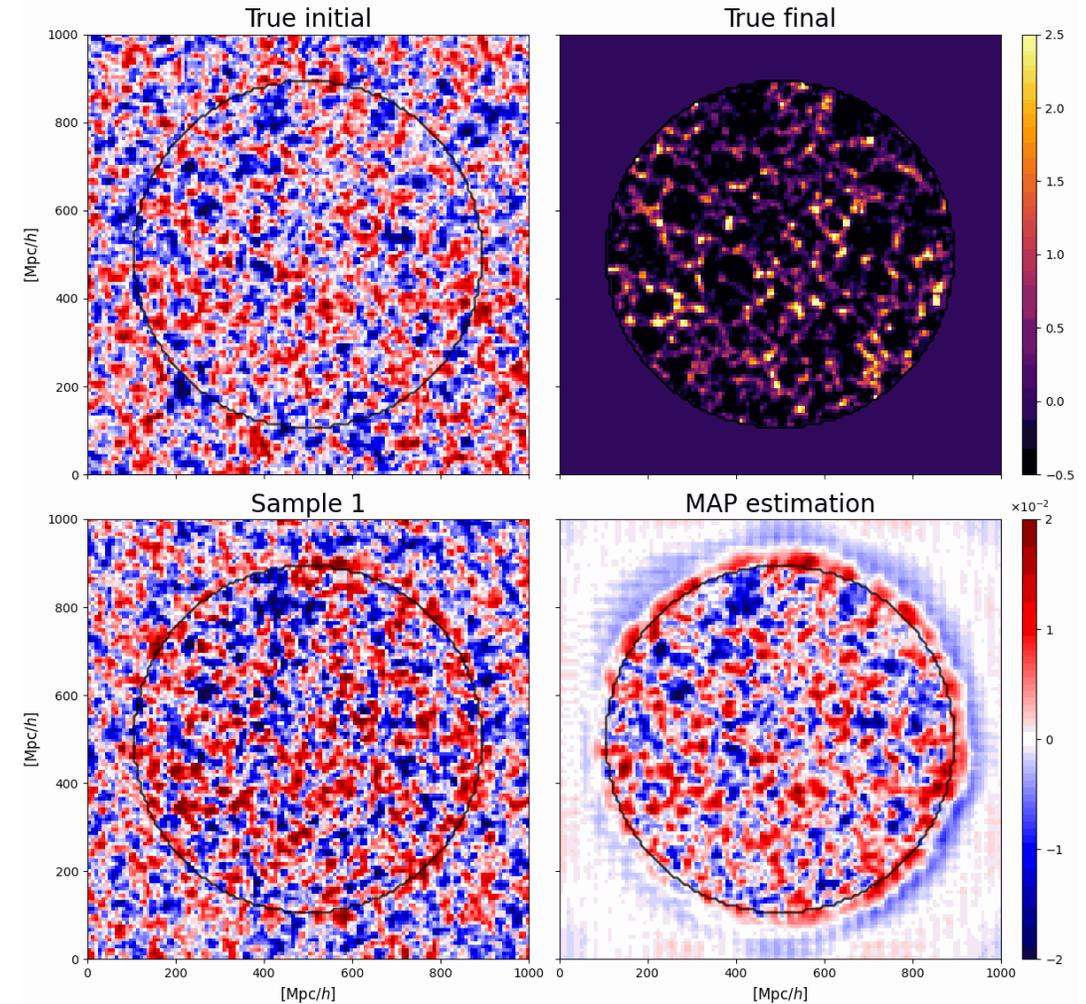
$$Q_{\text{eff}} = (Q^P + \gamma Q_{\theta}^L) [Q^P + (2\gamma - \gamma^2) Q_{\theta}^L]^{-1} (Q^P + \gamma Q_{\theta}^L)$$

Incomplete data

Field is constrained in the observed region and sampled from the prior in the unobserved region.



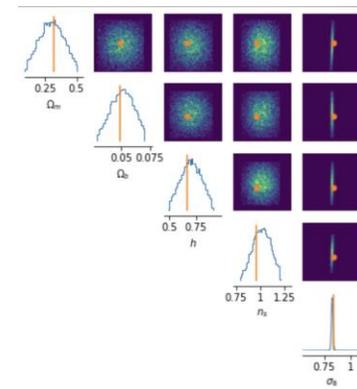
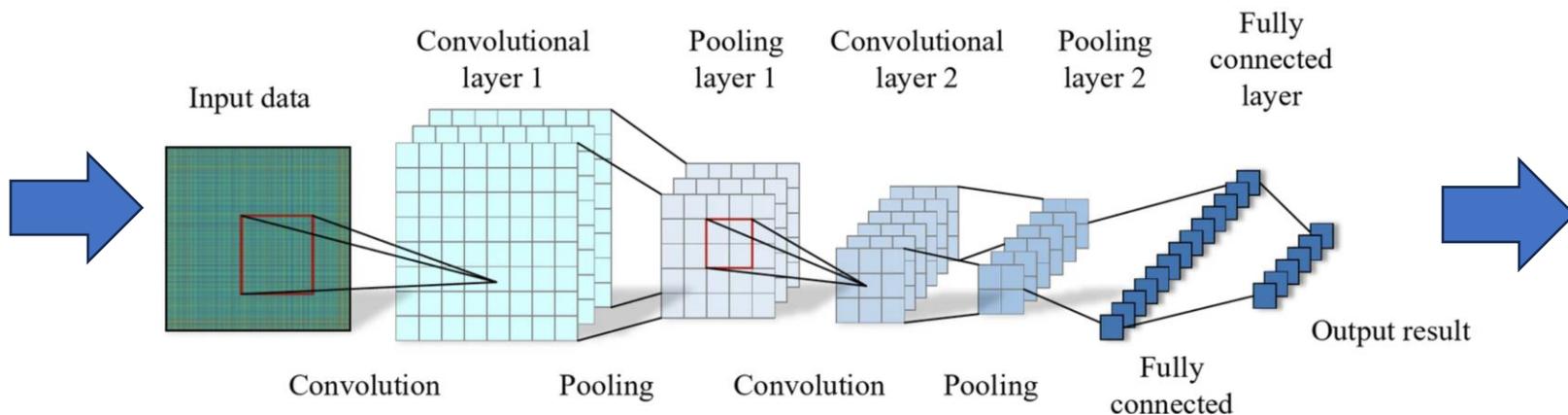
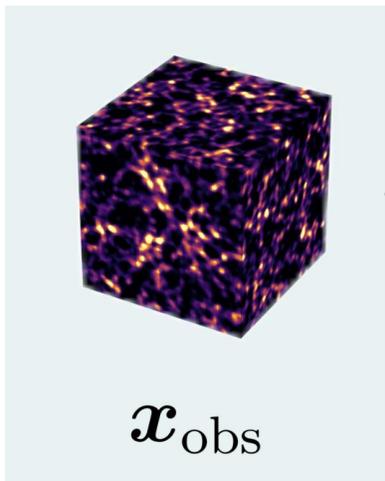
PRELIMINARY



Animation credit:
A. Orban

Parameter inference

CNN trainable summary statistic



$$p(\delta_{\text{IC}}, \theta | \delta_{\text{obs}}) = p(\delta_{\text{IC}} | \theta, \delta_{\text{obs}}) \times p(\theta | \delta_{\text{obs}})$$

Summary

- Reconstruction of cosmological **initial conditions** is an important problem that allows to **analyse LSS data in the fullest way**
- **Falcon** framework allows for powerful **distributed computing** applications for **active learning**
- **Tempering the likelihood** in high dimensions can be done by adding a specific amount of additional noise to the proposal
- Moving towards realistic settings (**incomplete data, varying cosmology**) requires **sequential SBI**

BACK UP SLIDES

Cosmological simulations

Types:

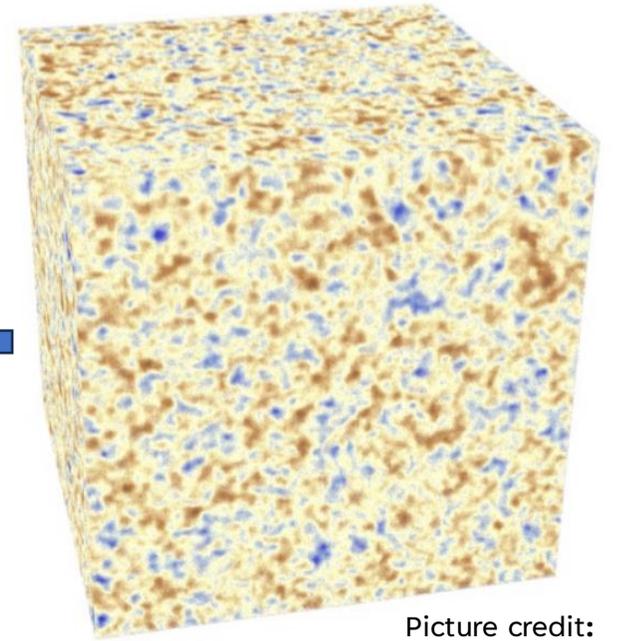
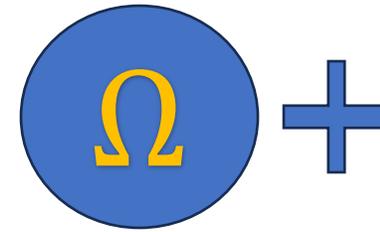
- LPT
- COLA
- Particle mesh
- N-body
- Hydrodynamical

Quijote

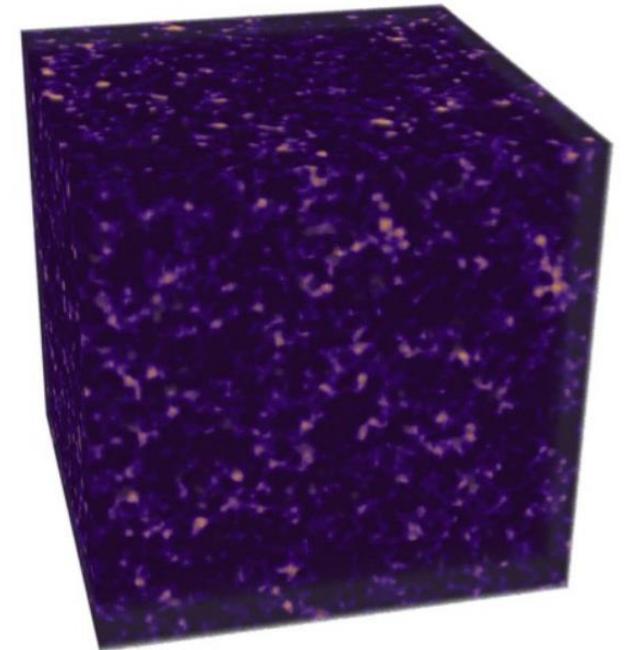
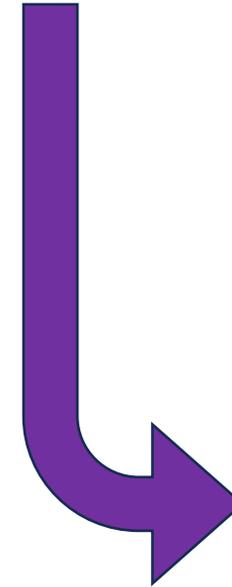
<https://quijote-simulations.readthedocs.io/>

CAMELS

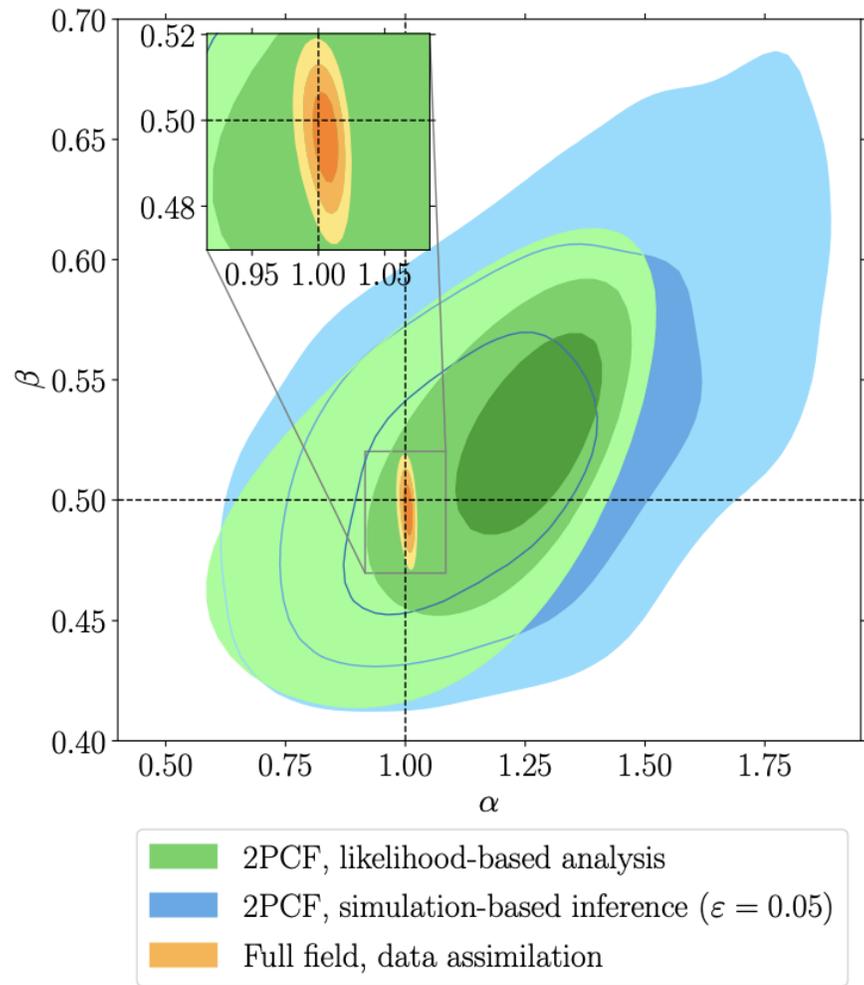
<https://camels.readthedocs.io/>



Picture credit:
Legin+, 2304.03788



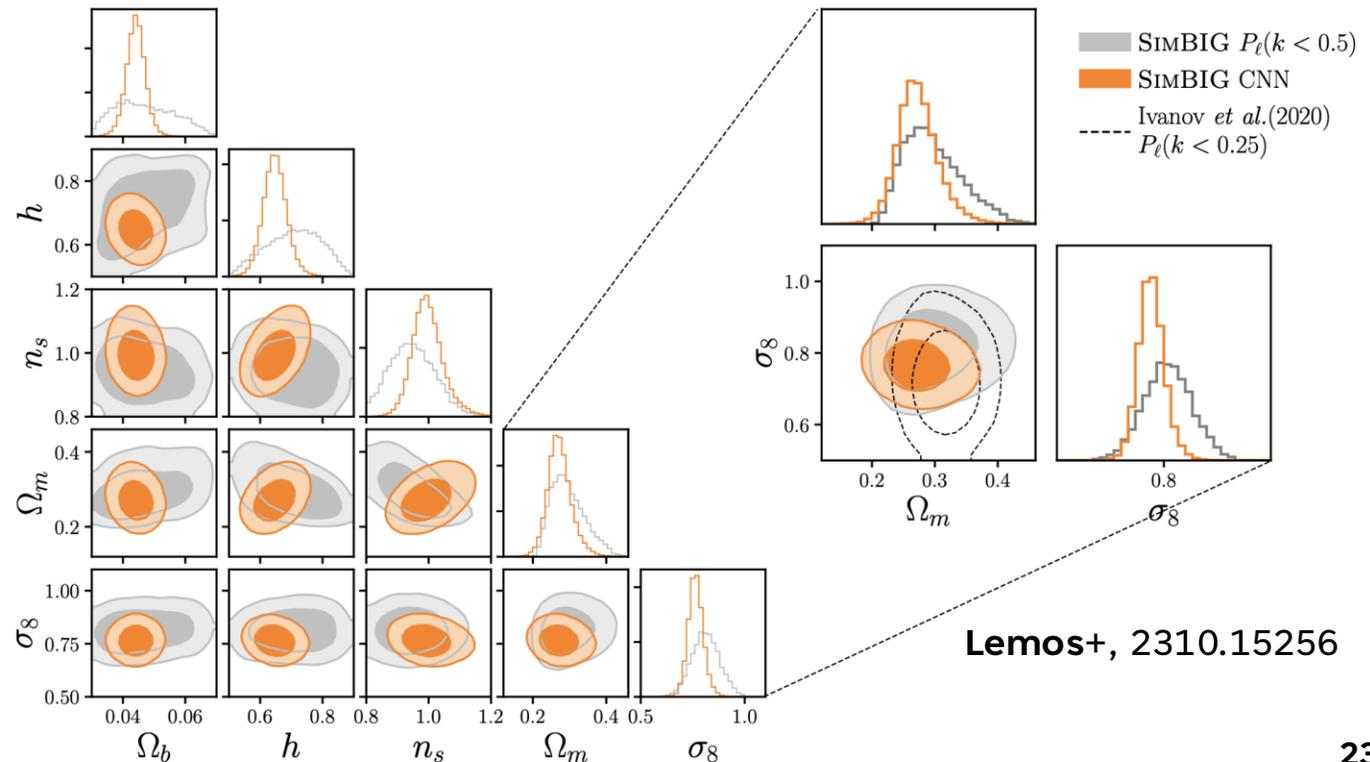
Field-level inference



Leclercq+, 2103.04158

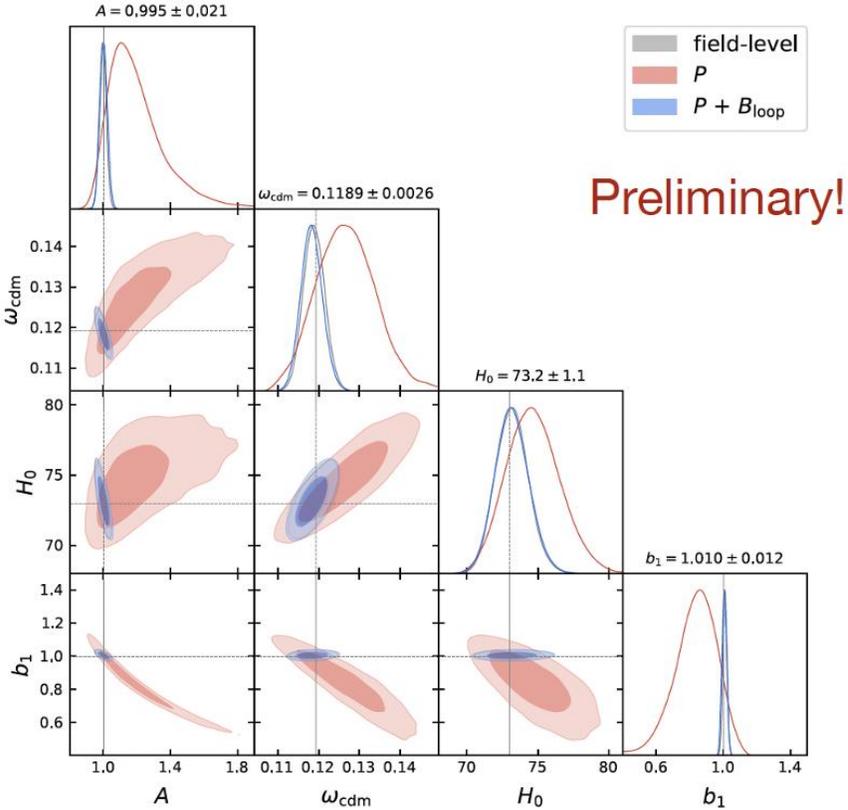
The whole field contains much more information than some summary like a power spectrum!

$$\langle \delta_m(\mathbf{k}) \delta_m(\mathbf{k}') \rangle = (2\pi)^3 \delta_D^3(\mathbf{k} + \mathbf{k}') P_{mm}(k)$$

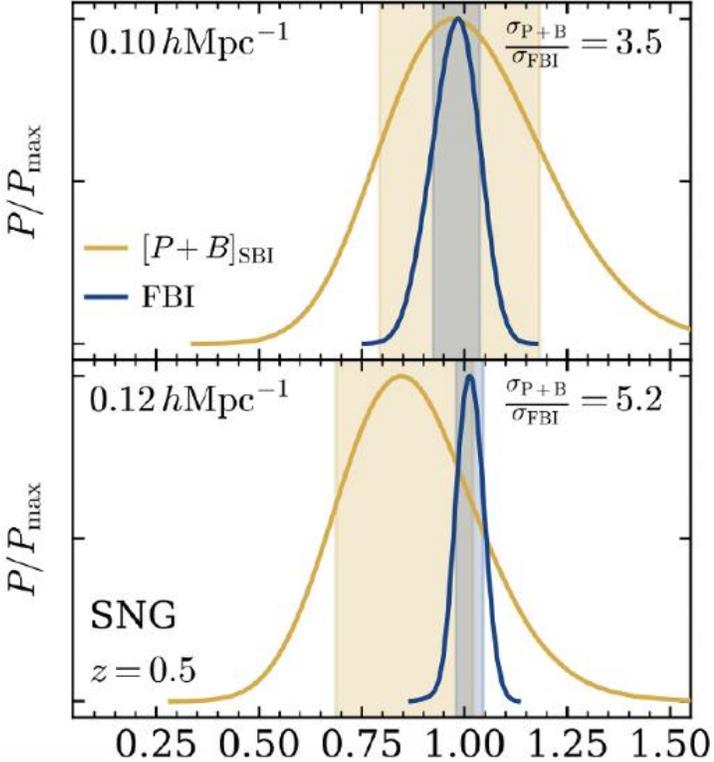


Ongoing debate

Akitsu et al. (to appear)



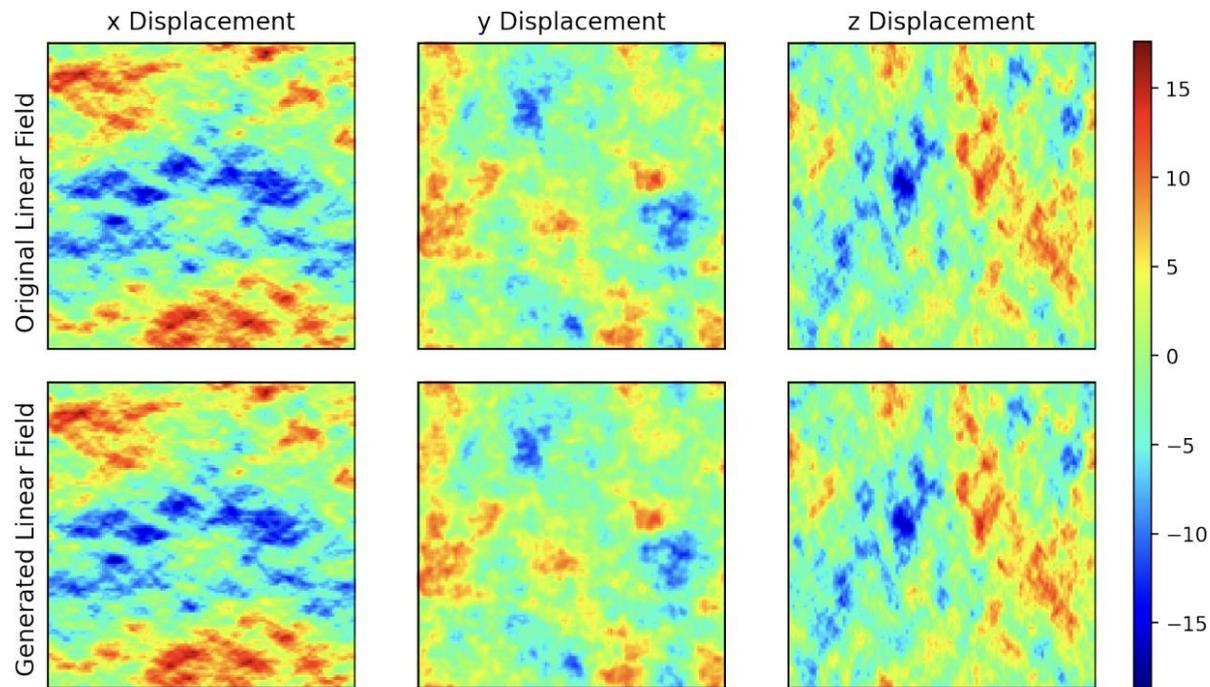
Nguyen, Schmidt, Tucci, Reinecke, Kostic (2024)



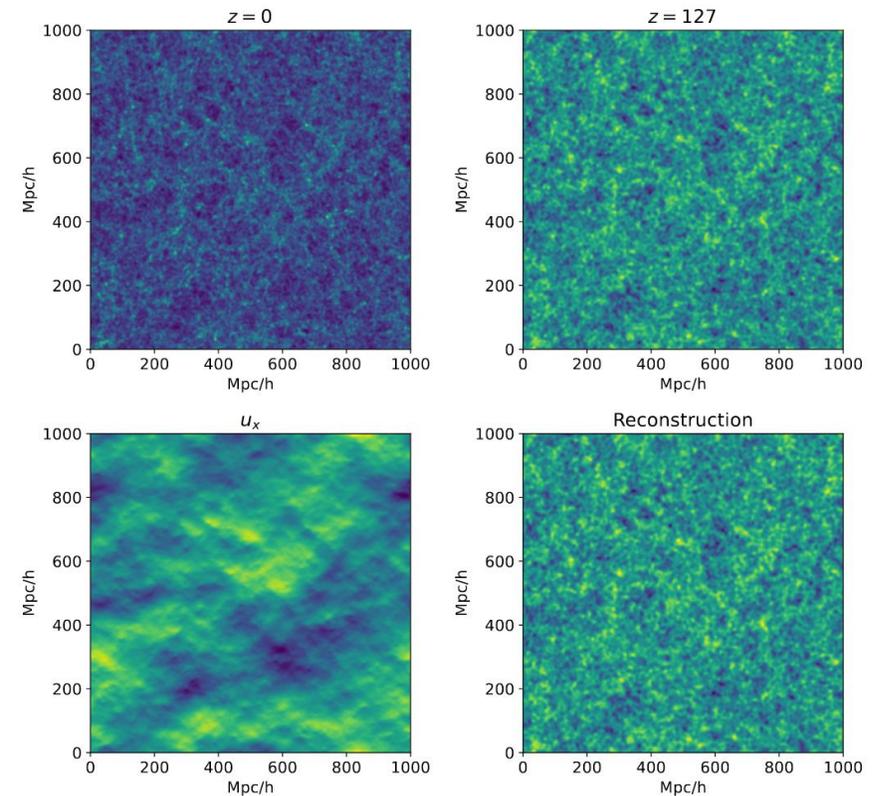
Different results from different groups!

ML approaches

- **Point estimates:** train a neural net to give single deterministic prediction (e.g. via MSE loss)



Jindal+, 2303.13056



Flöss+, 2305.07018

Diffusion models

- Train a network to approximate the **score**: $\nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{y})$
- Generate samples via reverse-diffusion process.

24 hrs of training on 4 80GB NVIDIA A100 GPU's

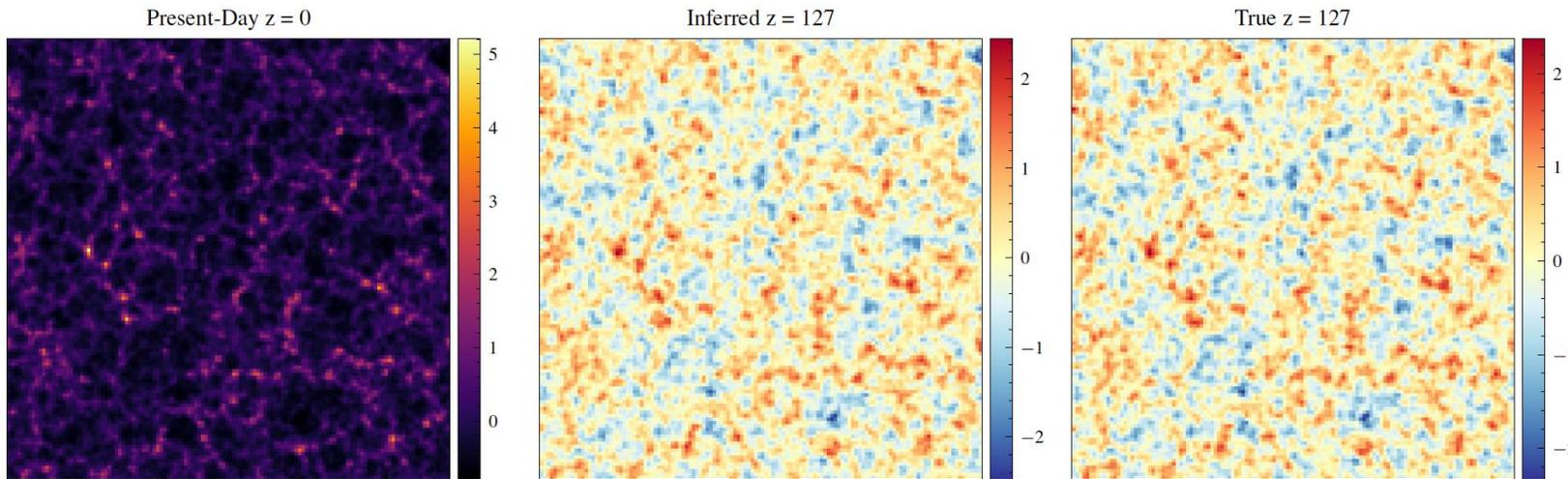
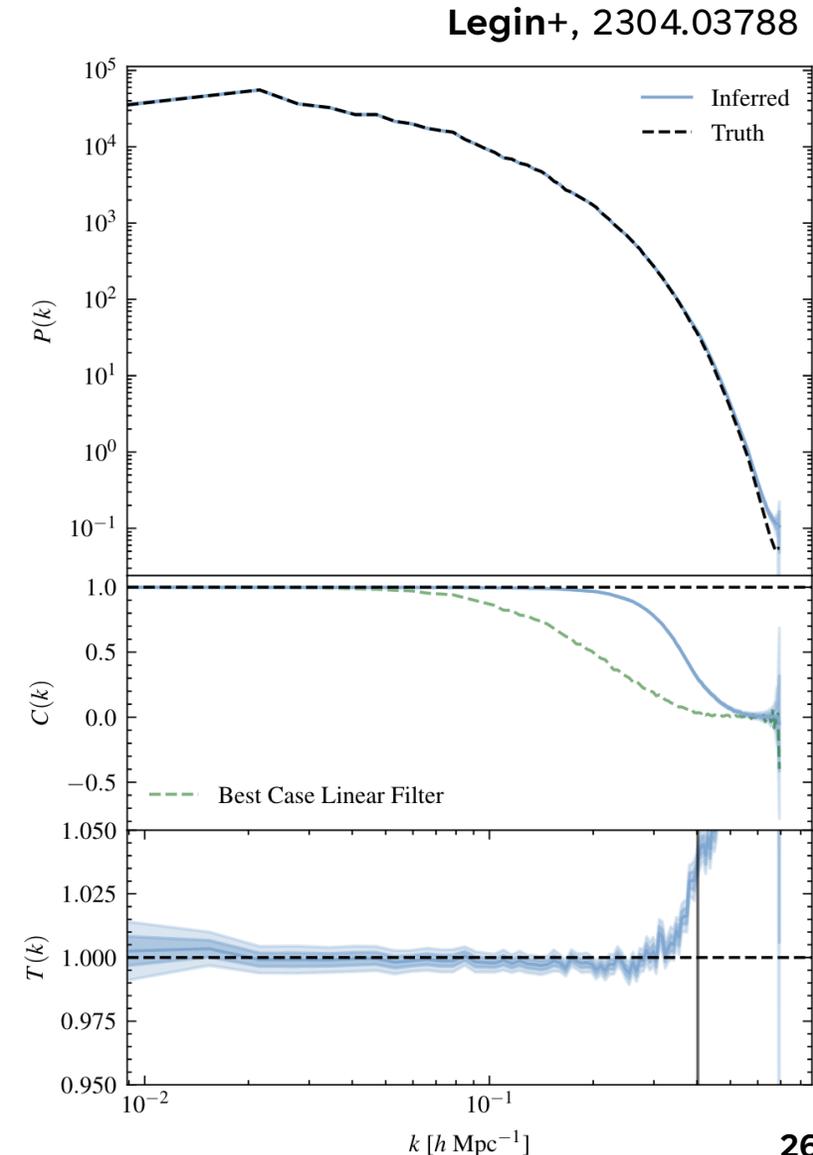
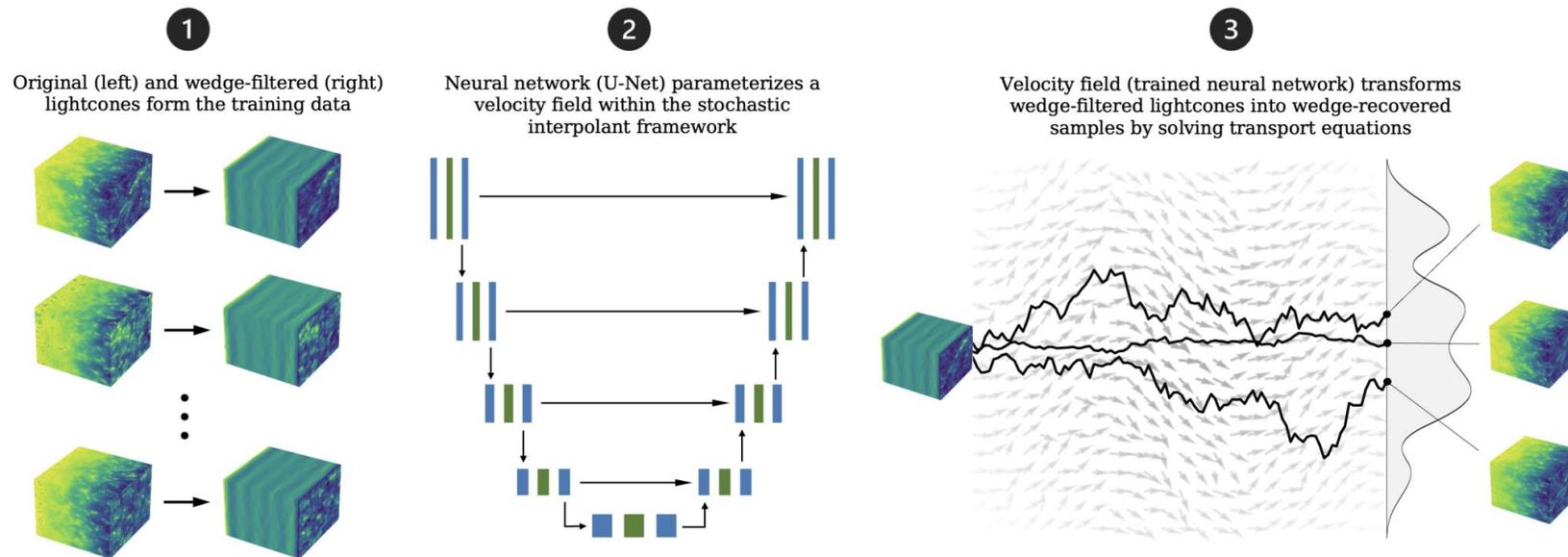


Figure 2. Left: The density field at redshift $z = 0$ for the fiducial Planck cosmology. Center: Initial conditions sampled from the posterior $p(\mathbf{x}|\mathbf{y})$. Right: The true initial conditions. All three density fields span a $1000 \times 1000 \times 125 (h^{-1} \text{Mpc})^3$ region averaged over the third axis. This example demonstrates the capability of score-based generative models to sample highly detailed initial conditions consistent with the ground truth. See Figure 3 for quantification of uncertainty.



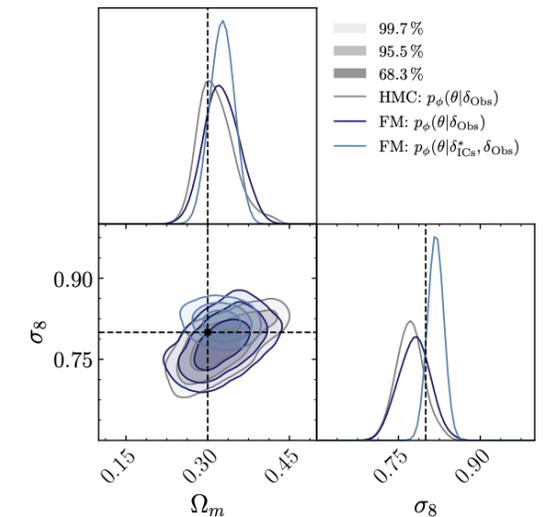
Flow-based models (stochastic interpolants)

- Train a network to approximate the **flow velocity field**
- Velocity field then stochastically transforms the samples



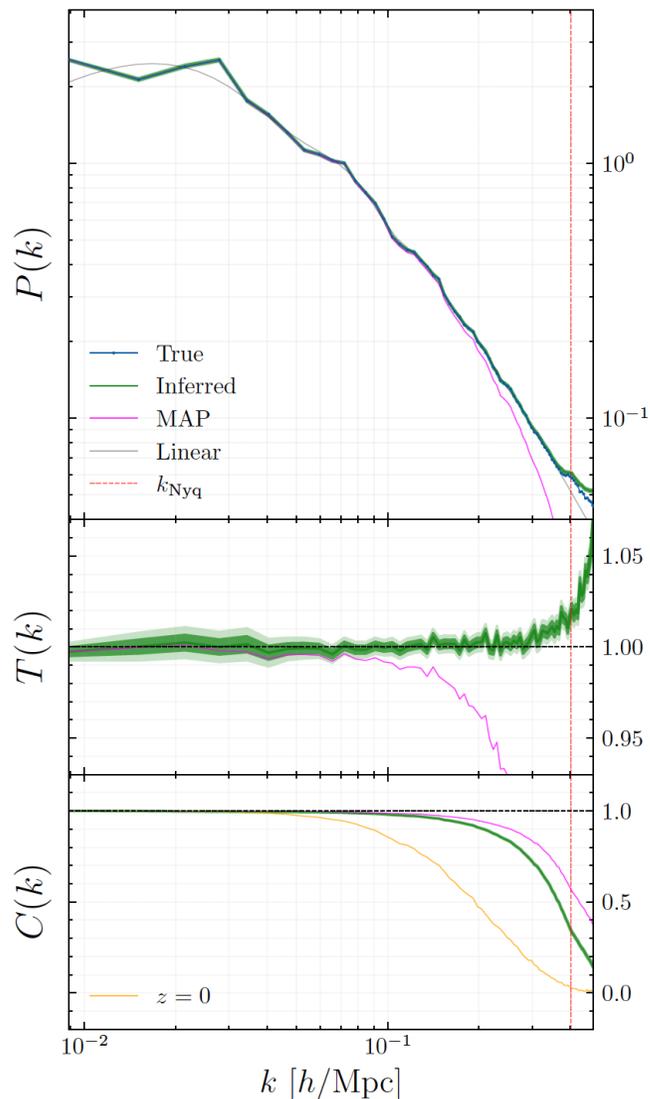
3 days of training on an 80GB NVIDIA A100 GPU,
sampling in minutes

Cuesta-Lazaro+,
NeurIPS 2024

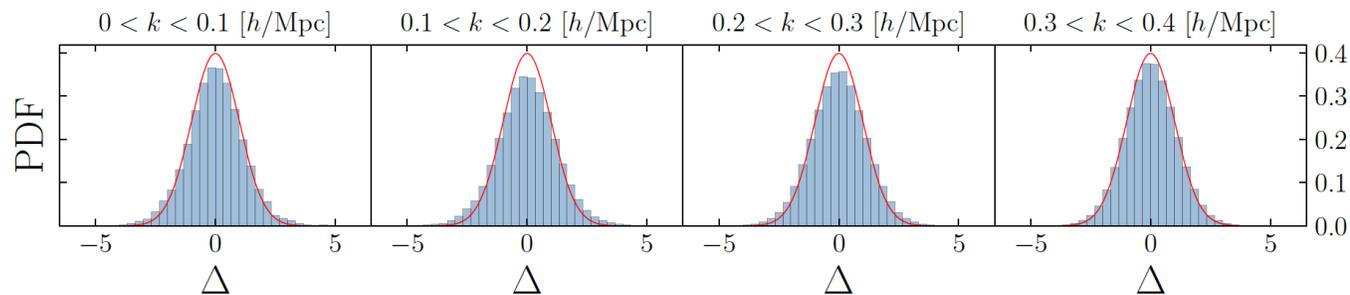


(a) Posteriors over cosmological parameters obtained through Stochastic Interpolants and HMC in a random test simulation. The true value is highlighted with dashed lines.

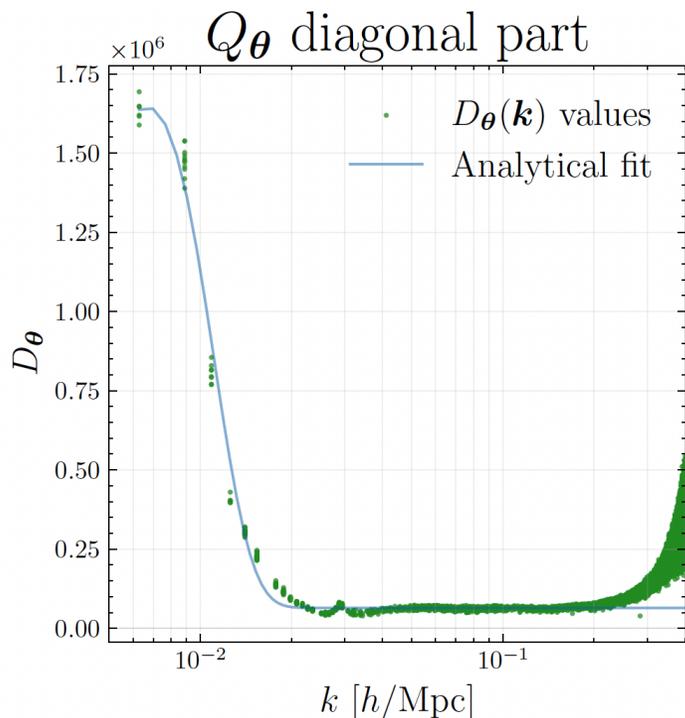
Summary statistics comparison



1-2% agreement in the power spectrum



Coverage test shows that samples follow the correct distribution



Knowledge of the $Q(|k|)$ dependence allows to turn any point estimator into a fast sampler

[OS+, 2502.03139](#)

Incomplete data

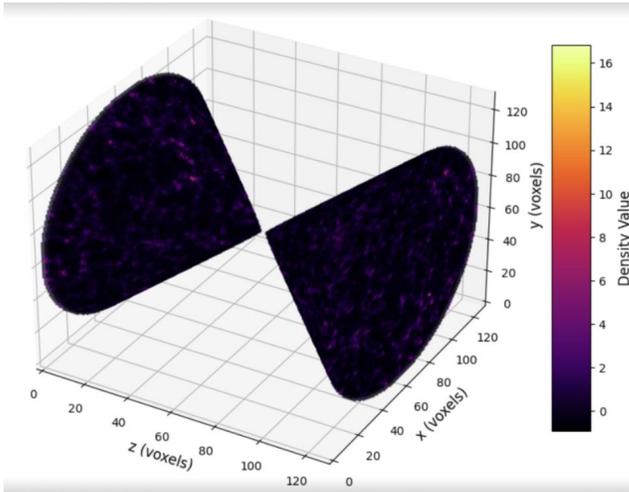
- Modify \mathbf{Q}_L matrix with real-space factors P :

$$\mathbf{Q}_\theta^L = \mathcal{F}^\dagger \mathbf{D}_\theta^L \mathcal{F} \quad \Rightarrow \quad \mathbf{Q}_\theta^L = \mathcal{P} \mathcal{F}^\dagger \mathbf{D}_\theta^L \mathcal{F} \mathcal{P}$$

- Use Conjugate Gradient to estimate $\text{tr} \log \mathbf{Q}$:

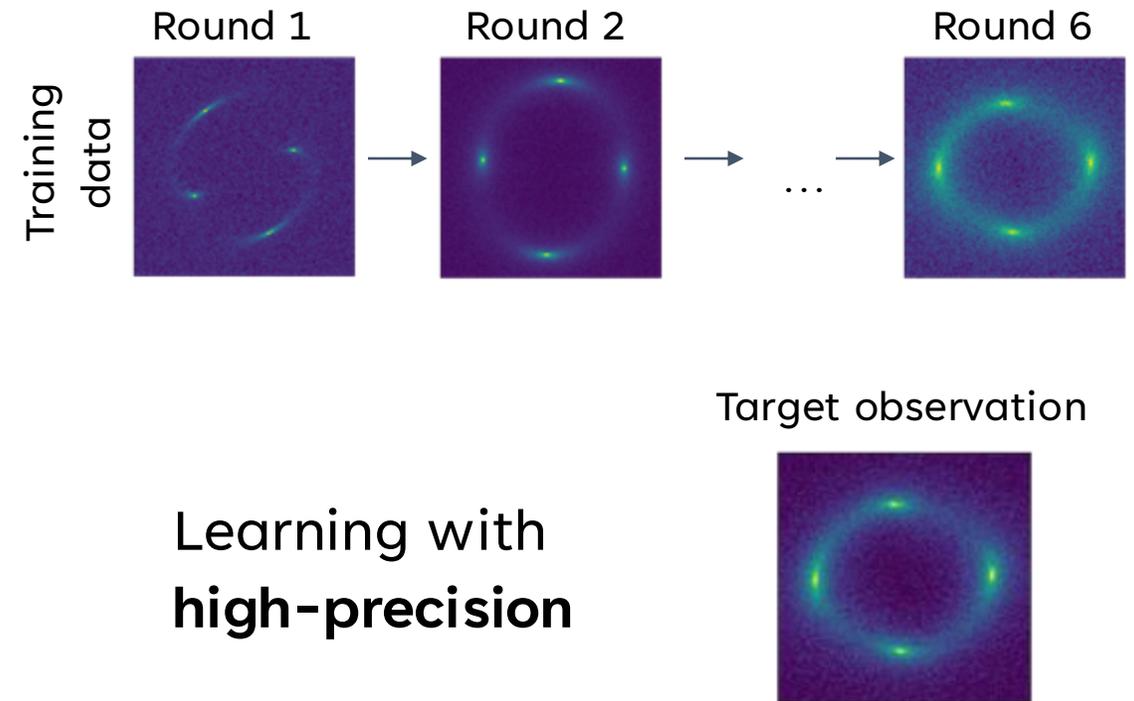
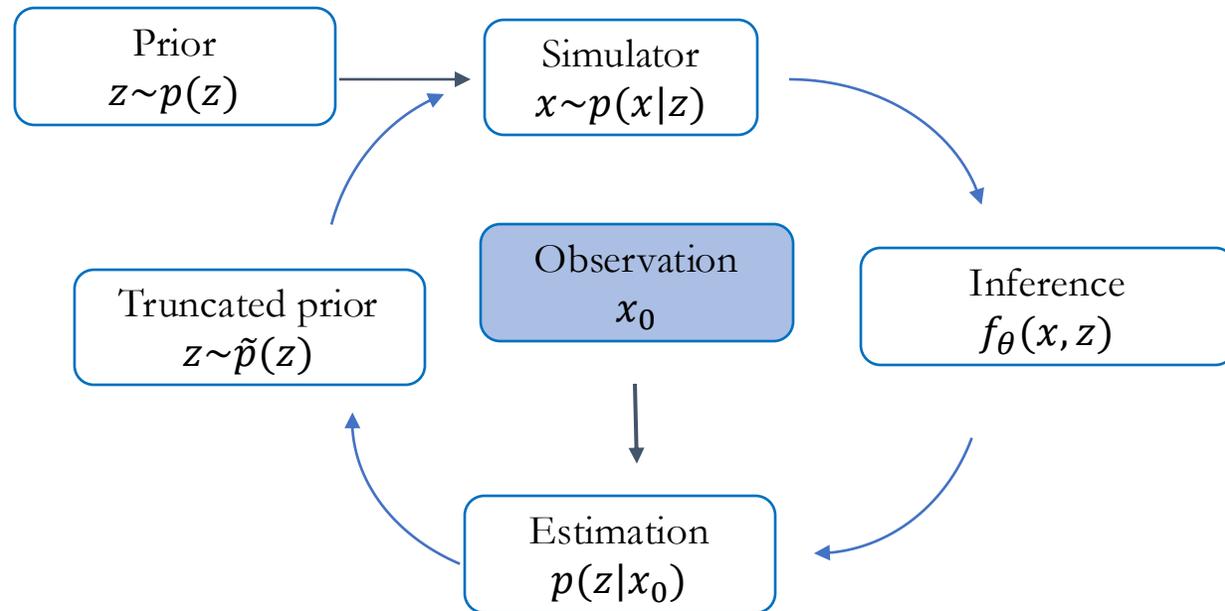
$$\mathcal{L} = \frac{1}{2} \sum_{i=1}^N \left\{ (z_i - \hat{\mu}_\theta(x_i))^T \mathbf{Q}_\theta (z_i - \hat{\mu}_\theta(x_i)) \right\} - \frac{N}{2} \text{tr} \log \mathbf{Q}_\theta$$

- Then sampling requires GEDA



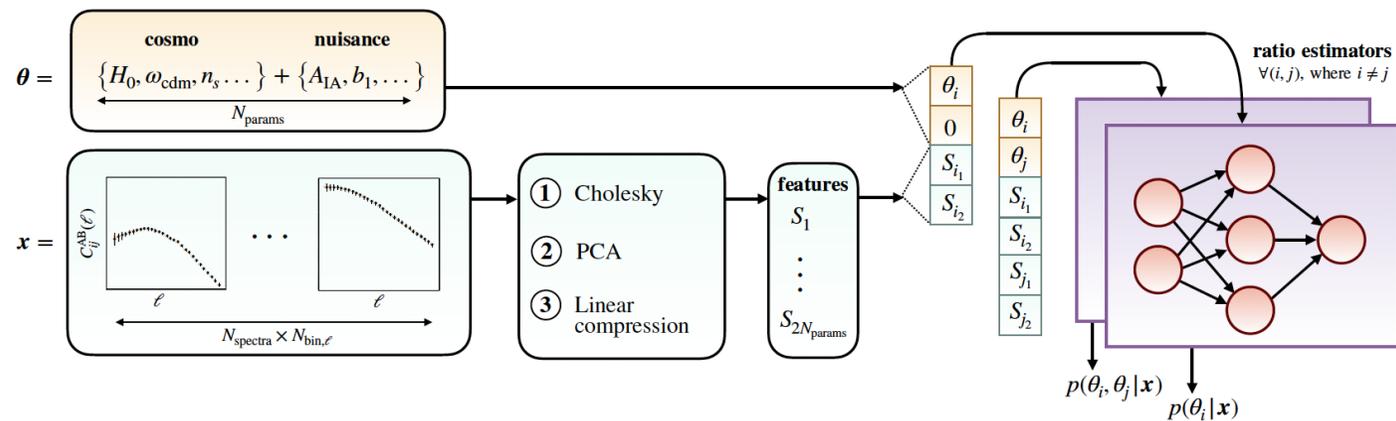
Sequential inference & adaptive learning

- Our parameter space is too vast to explore
- Want to ‘zoom in’ into it and obtain precise results with a low number of simulations

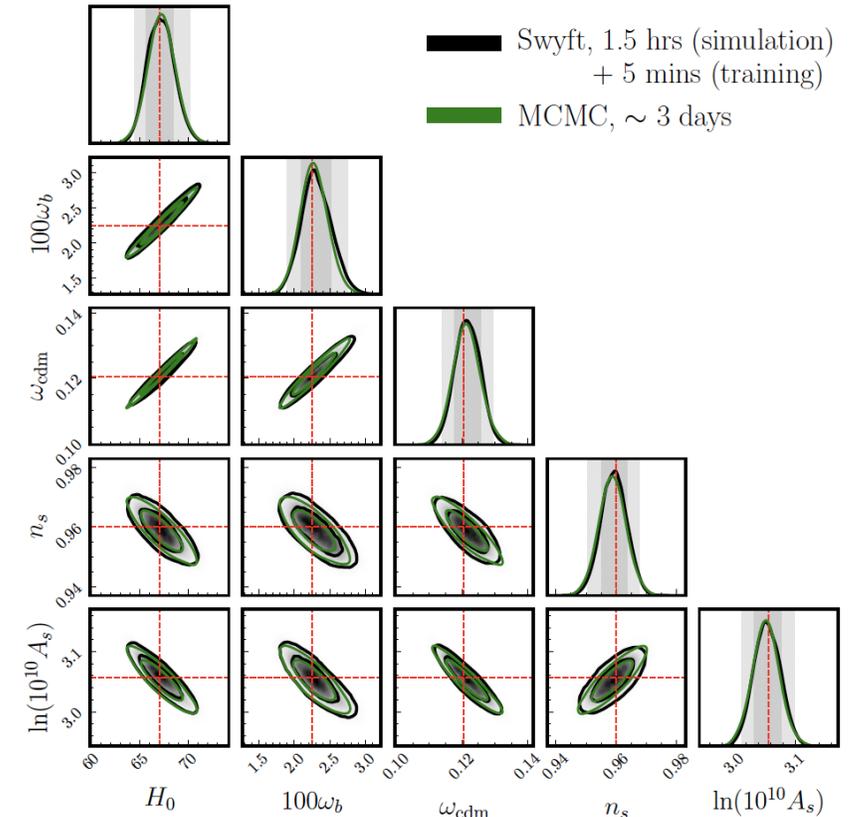


Parameter inference applied to Euclid

Apply **Marginal Neural Ratio Estimation** algorithm via **swyft** code to pre-compressed **3x2pt** statistics



G.F. Abellán+, [2403.14750](#)



Excellent agreement with MCMC and a dramatic reduction in CPU time