

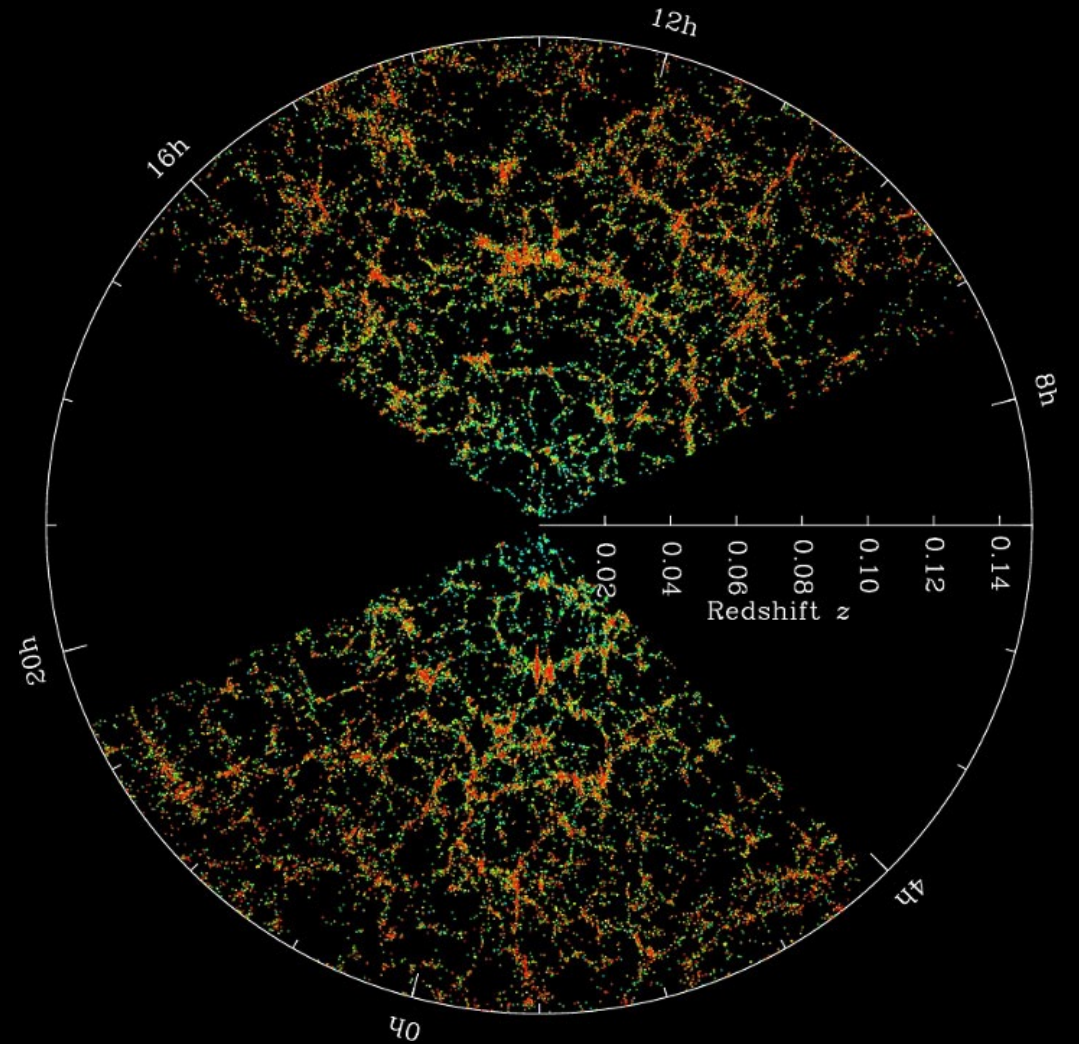
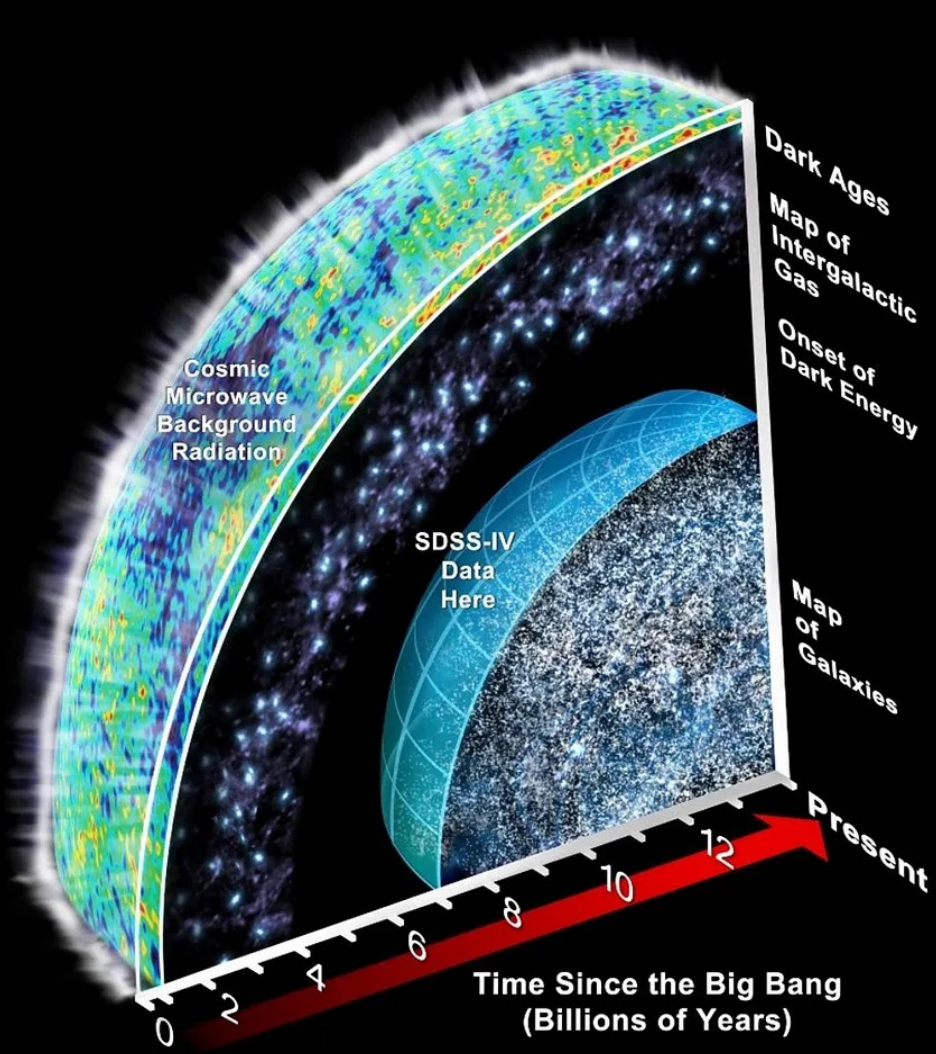
Reconstruction of cosmological initial conditions with sequential simulation-based inference

Oleg Savchenko

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



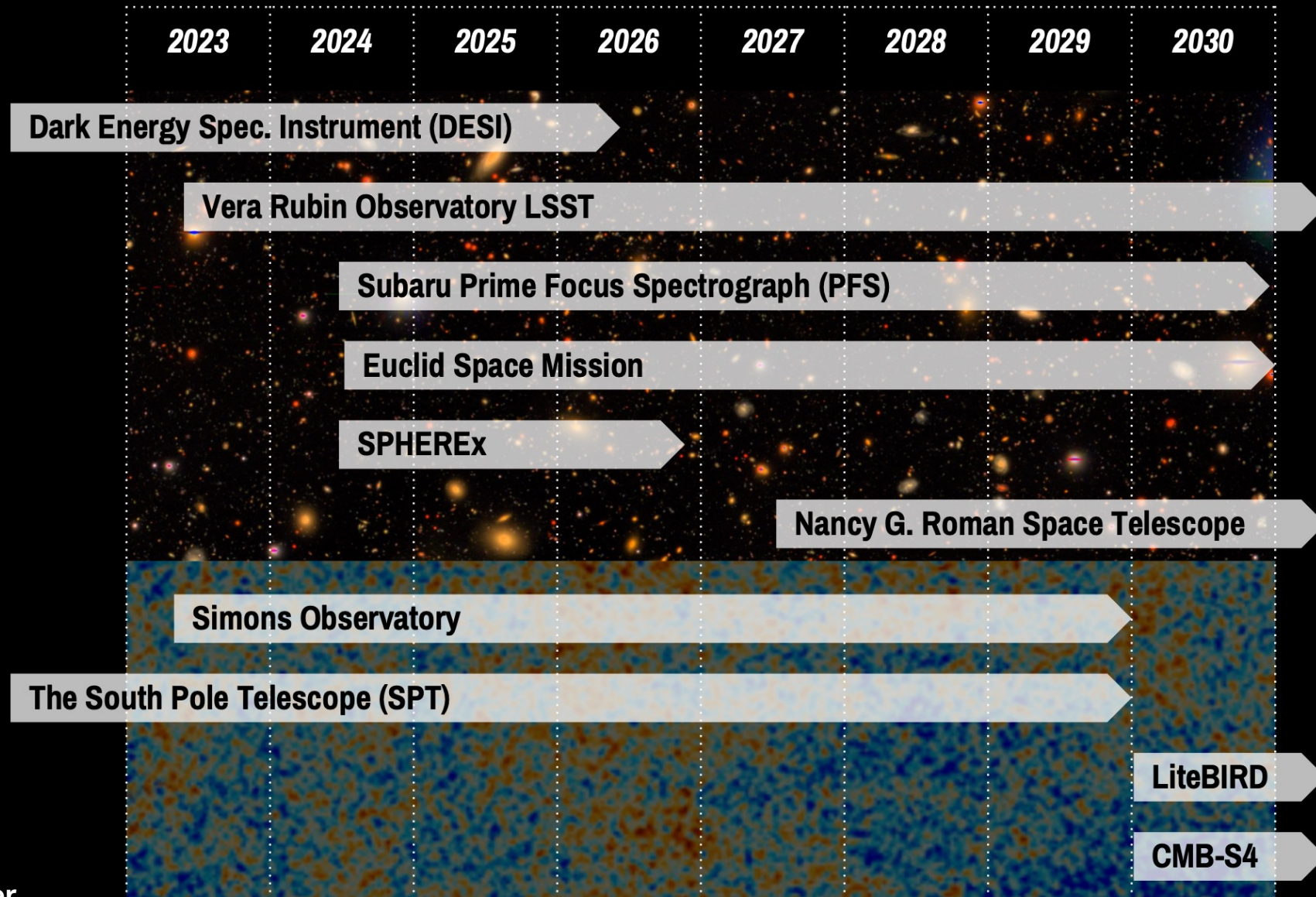
Large scale structure



<https://mapoftheuniverse.net/>

SDSS collaboration

Next decade



Cosmological simulations

Types:

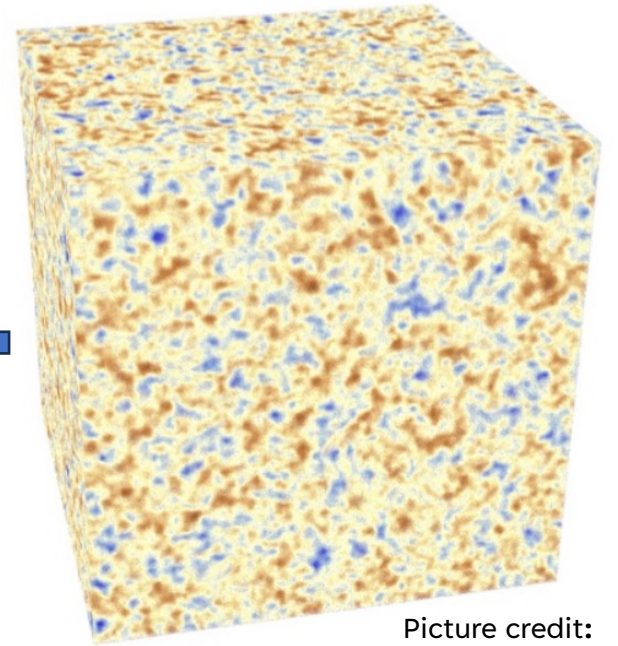
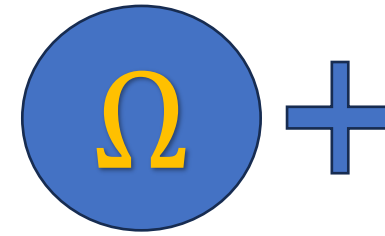
- LPT
- COLA
- Particle mesh
- N-body
- Hydrodynamical

Quijote

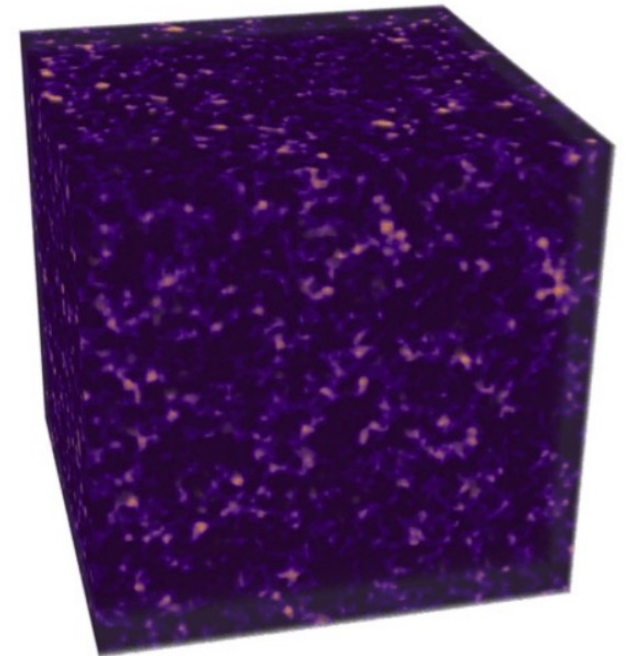
<https://quijote-simulations.readthedocs.io/>

CAMELS

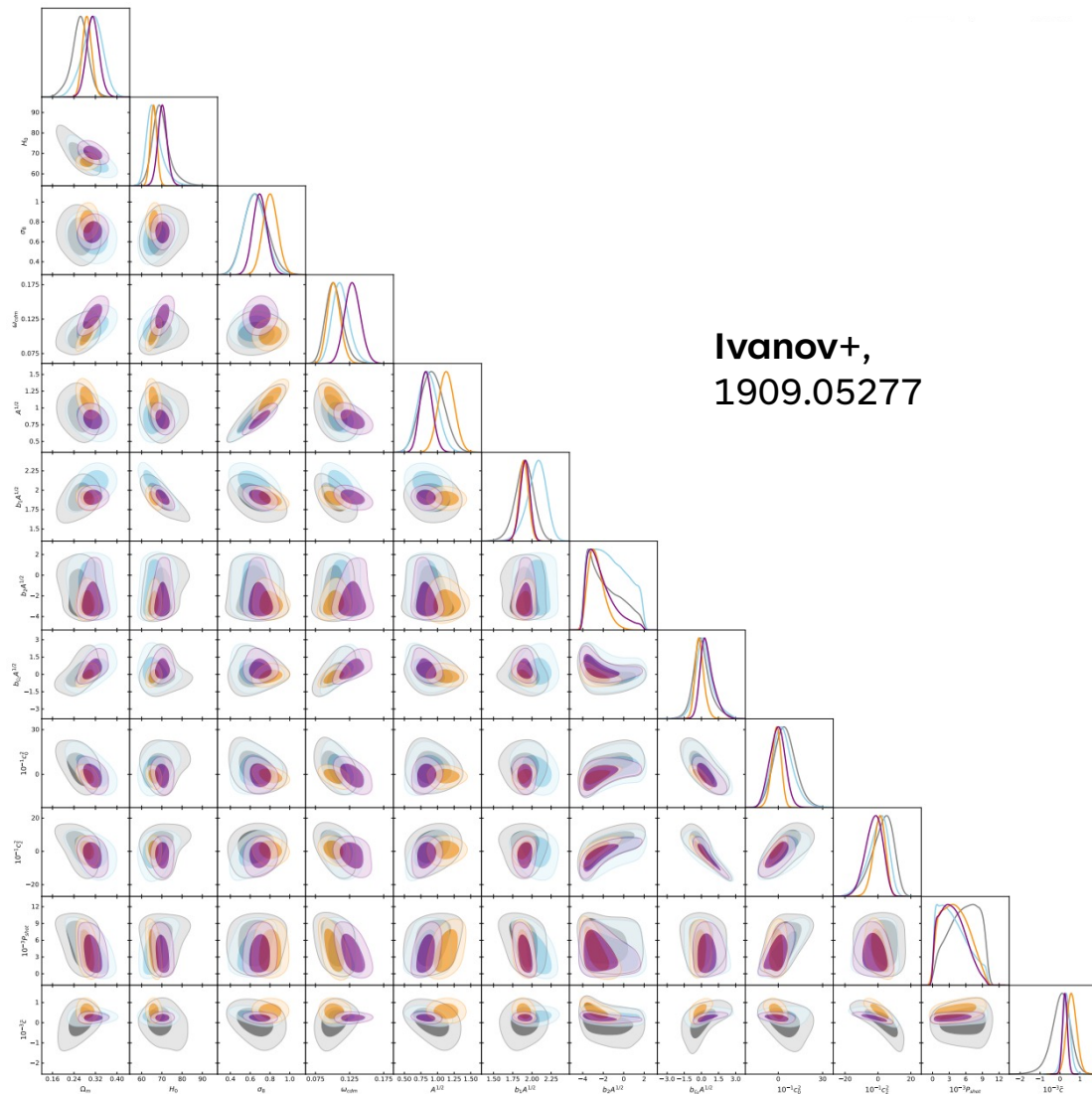
<https://camels.readthedocs.io/>



Picture credit:
Legin+, 2304.03788



How to analyse LSS data?



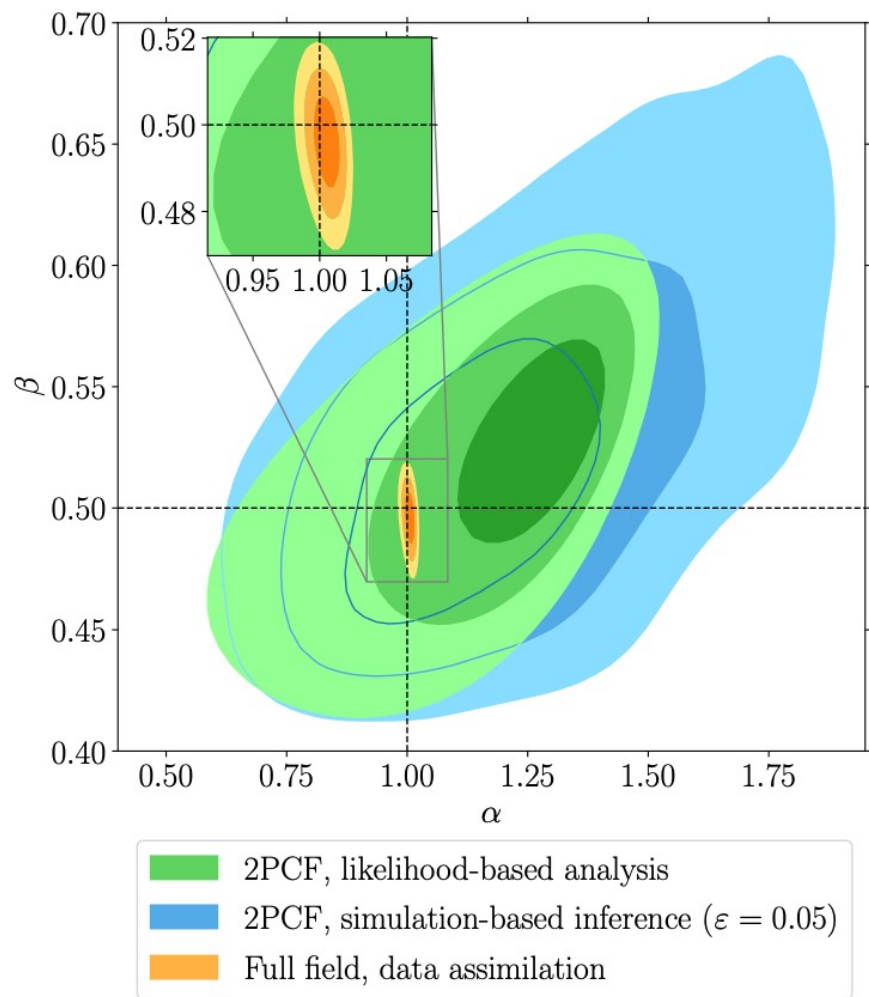
Classic approach:

- Come up with some summary statistic $s(\text{data})$
- 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.

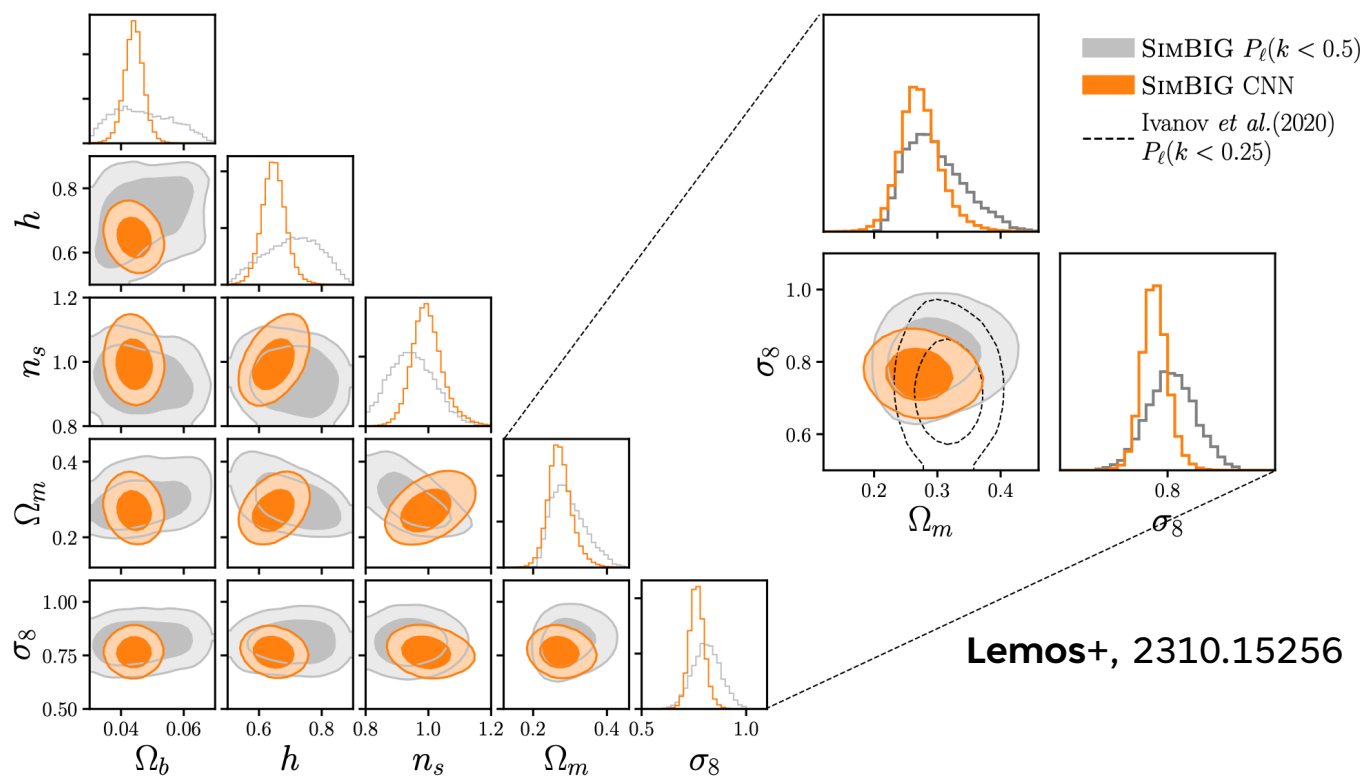
Field-level inference



Leclercq+, 2103.04158

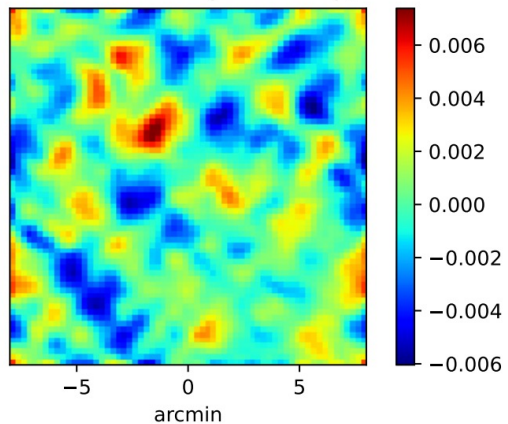
The whole field contains much more information than some summary like a power spectrum!

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



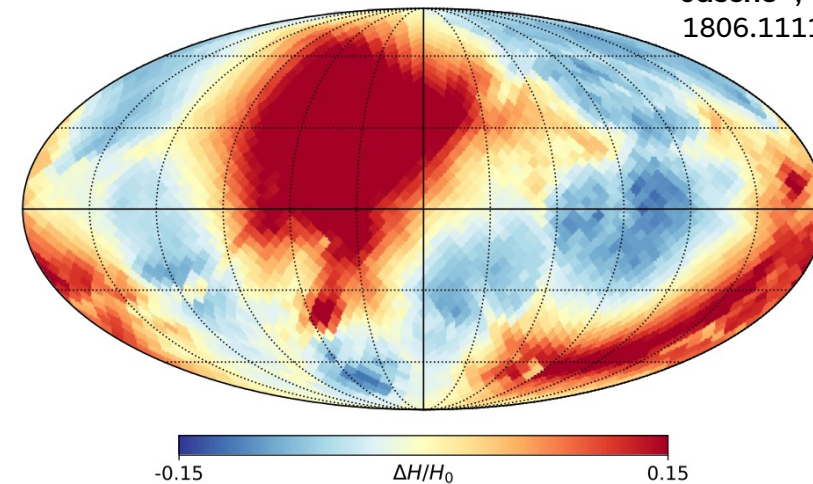
Reconstructing fields

Posterior mean predicted γ_2



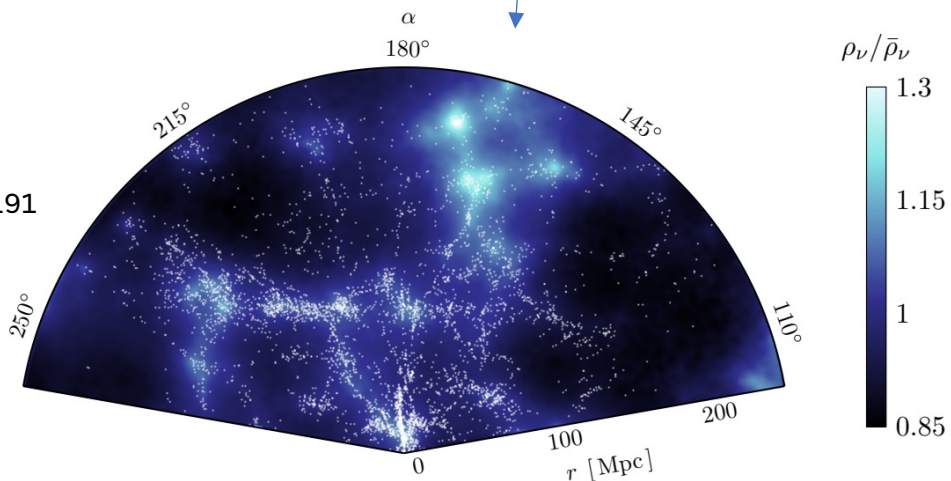
Shear

Mass estimation and evolution
history, velocity fields, BAO,
 f_{NL} , 21 cm, weak lensing...

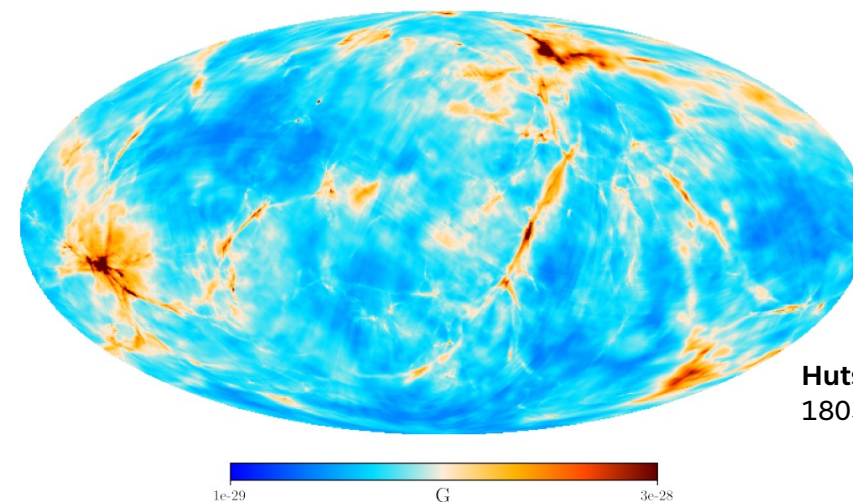


Hubble parameter
uncertainties

Elbers+,
2307.03191



Neutrino field



Magnetic field

Hutschenreuter+,
1803.02629

Initial conditions reconstruction

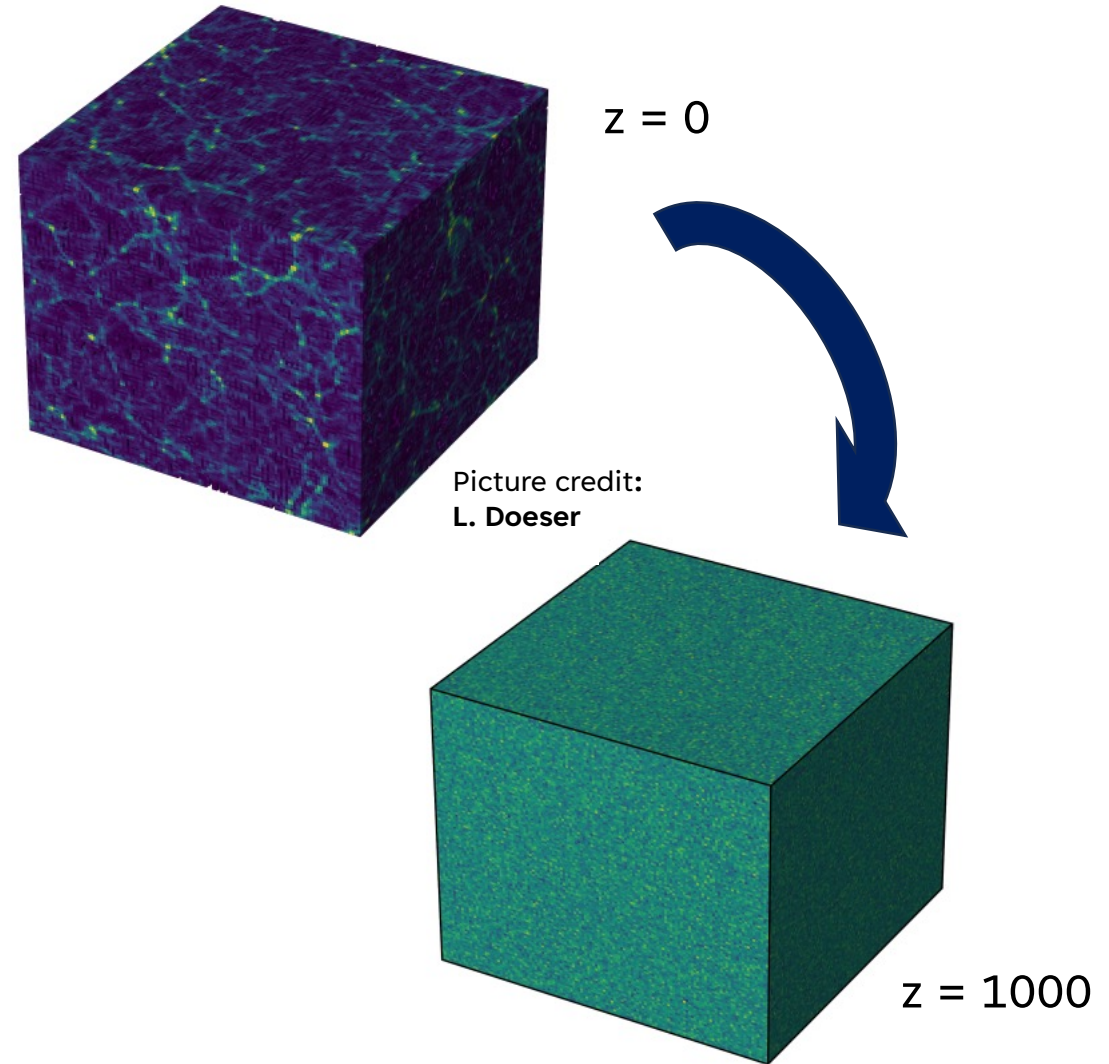
Very early universe had very simple properties!

→ 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

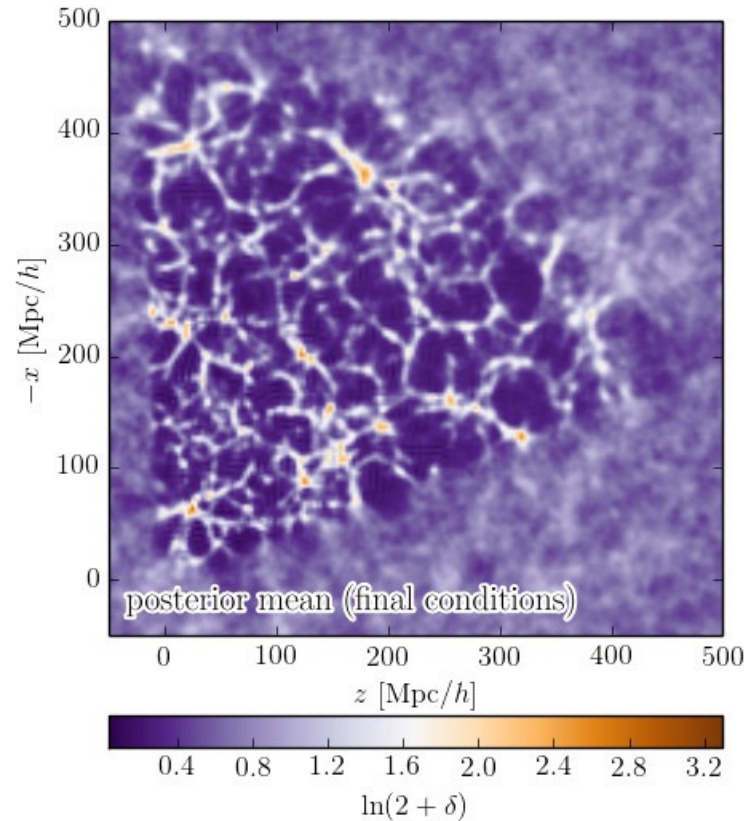
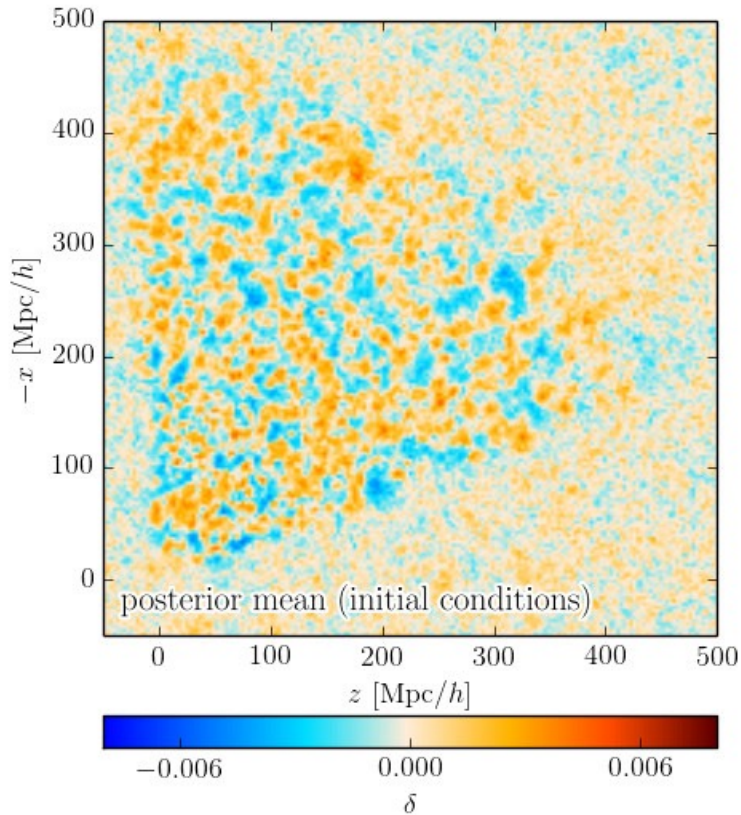
$\{N_i^g\}$ - galaxy catalog data



Gaussian initial conditions

Bayesian Origin Reconstruction from Galaxies (BORG)

Jasche+, 1203.3639, 1806.11117

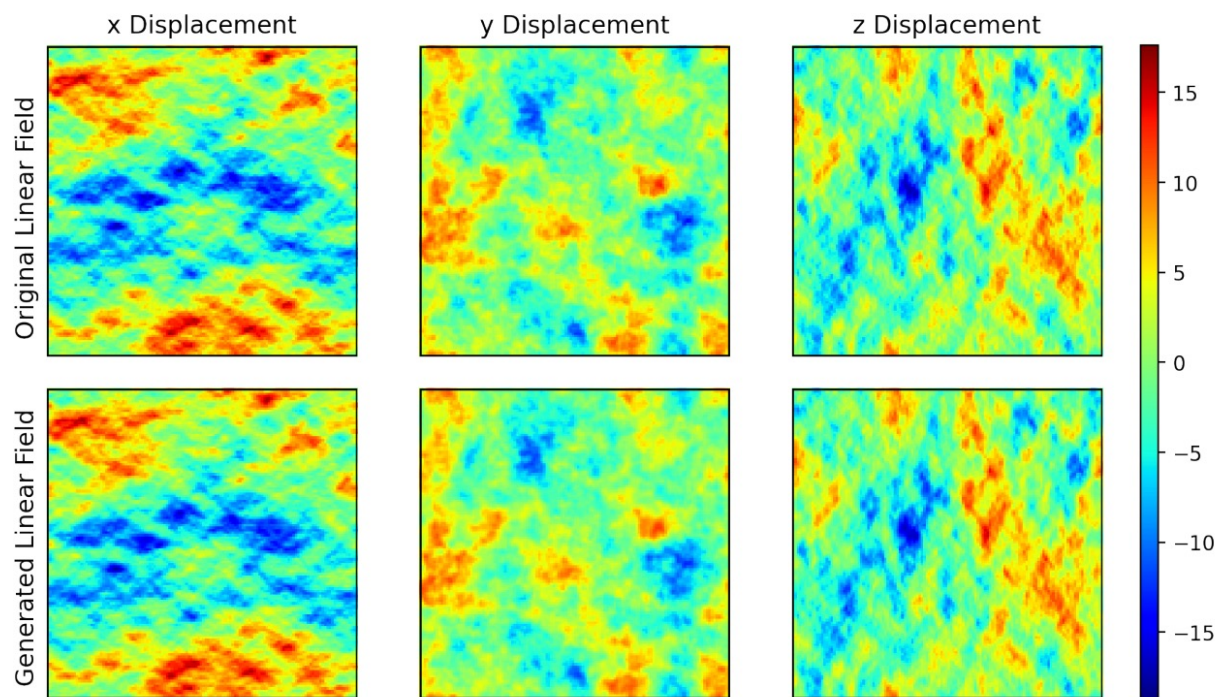


- Most developed method so far
- Explicit likelihood
- Requires gradients from the simulator
- Takes hours to produce a single sample
- Not amortized

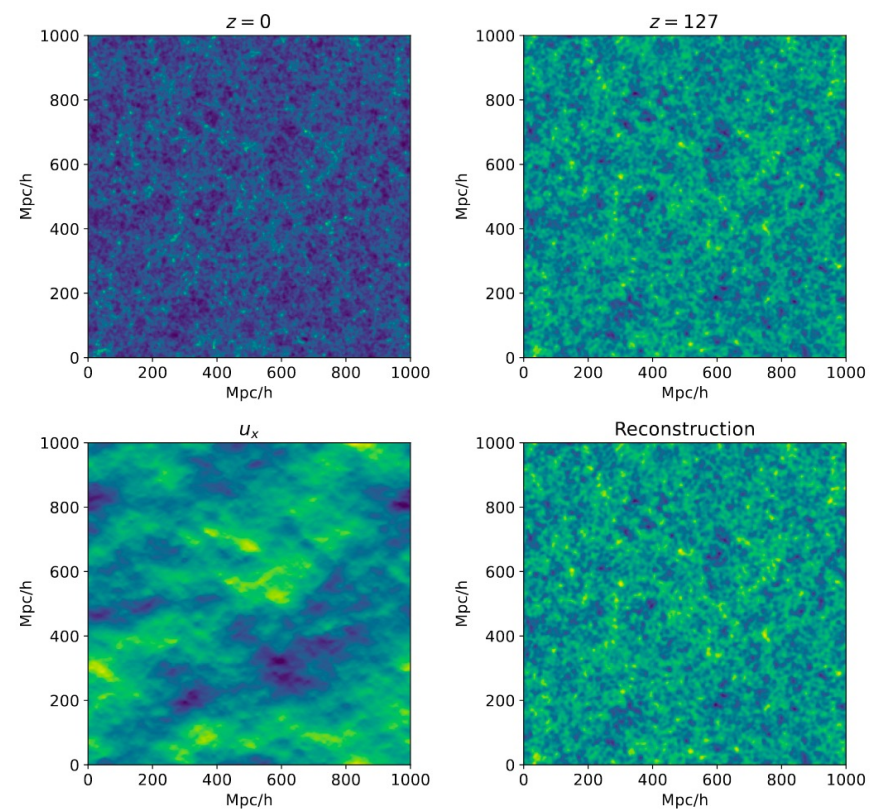
Samples produced with Hamiltonian Monte Carlo

ML approaches

- **Point estimates:** train a neural net to give single deterministic prediction (MAP estimation)



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

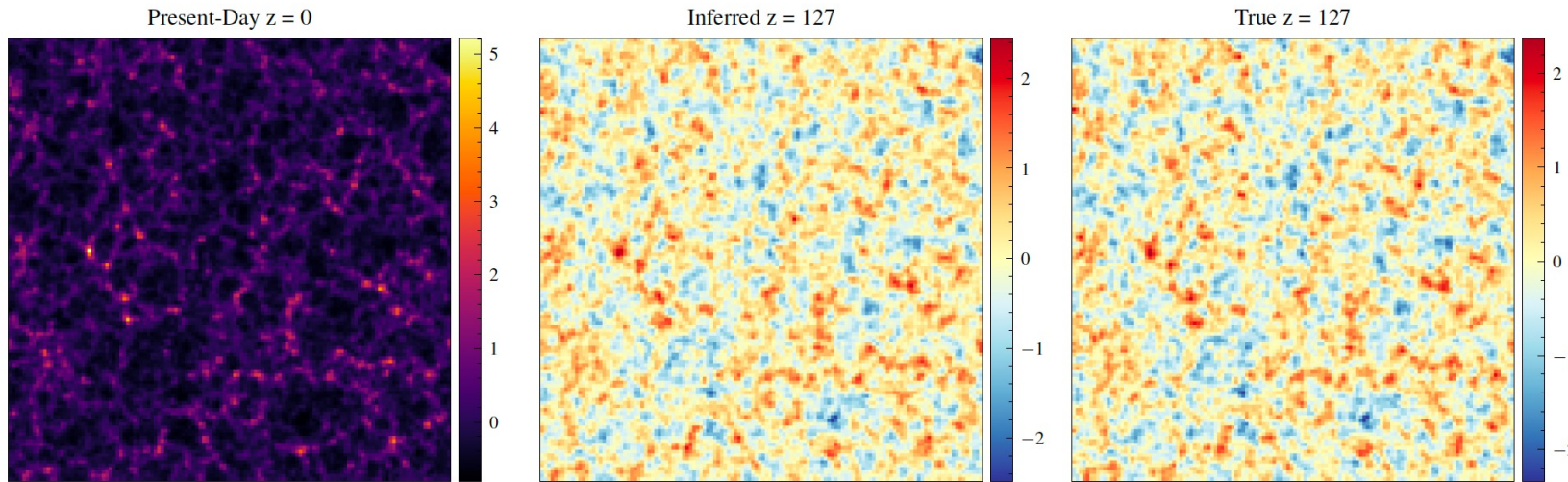
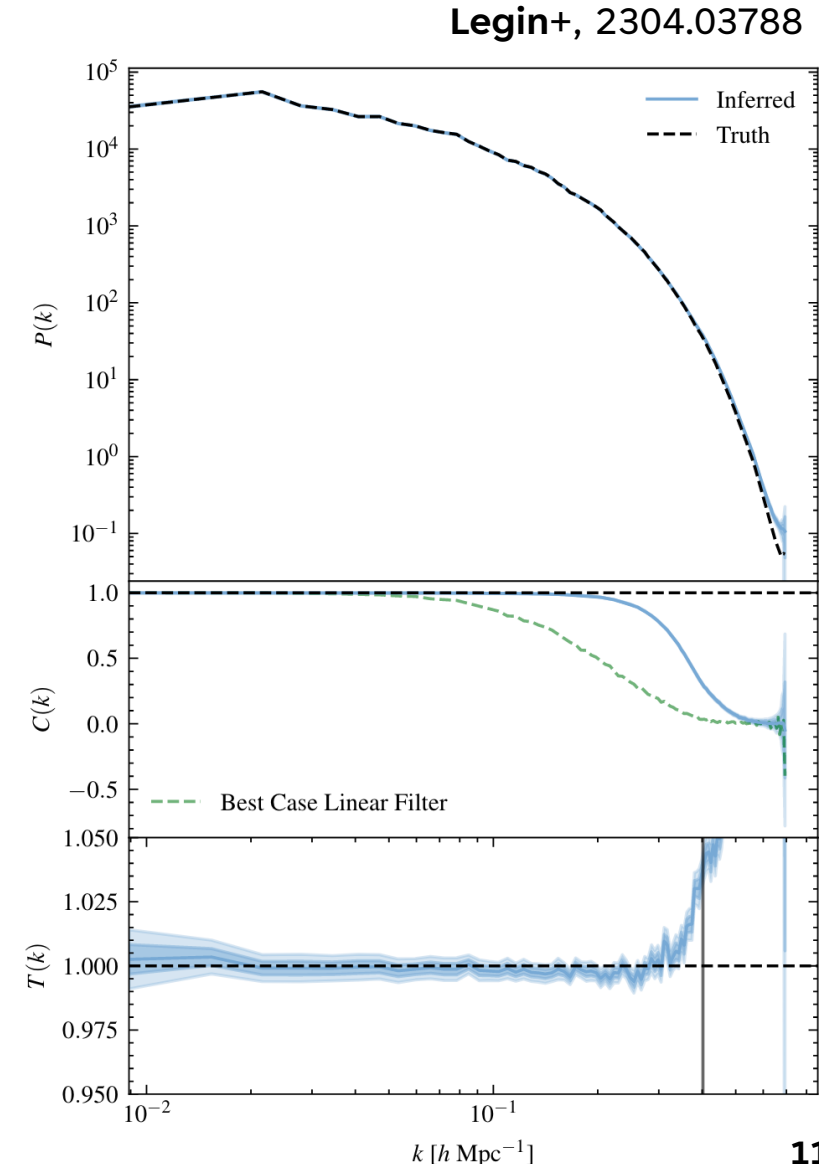
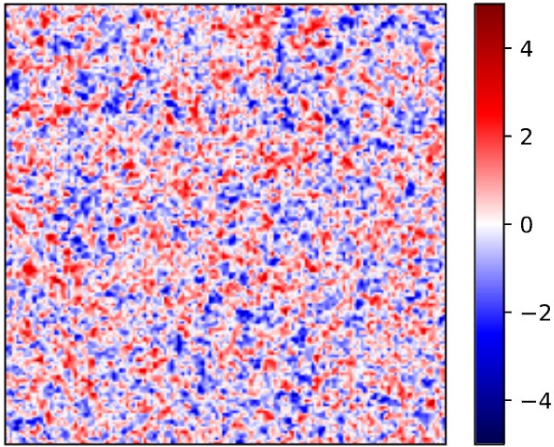


Figure 2. Left: The density field at redshift $z = 0$ for the fiducial Planck cosmology. Center: Initial conditions sampled from the posterior $p(\mathbf{x}|\mathbf{y})$. Right: The true initial conditions. All three density fields span a $1000 \times 1000 \times 125 (h^{-1} \text{Mpc})^3$ region averaged 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.

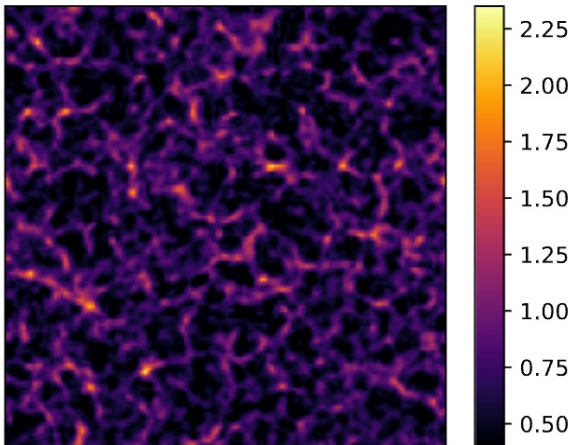


Our setting

$\delta(x)$ slice at $z=1000$



$\log(2 + \delta(x))$ slice at $z=0$



Training data: **2000** 128^3
(1 Gpc/h)³ Simbelmyne or
Quijote simulations

128^3 resolution: \sim million-dimensional parameter space!

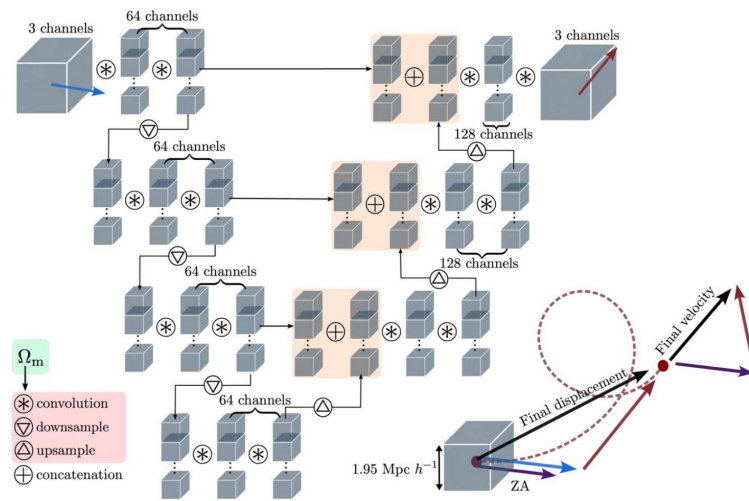
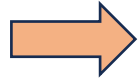
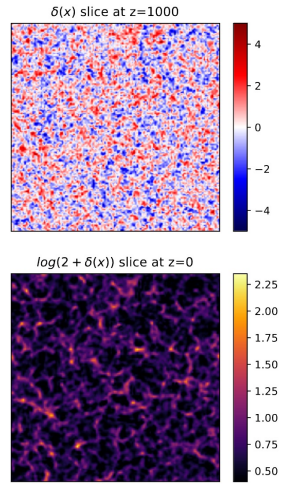
- Want to explore the full posterior, not only get point estimates
- Want to keep things as simple as possible
- Want to be able to do things sequentially: need to estimate the likelihood:

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

- Tried different approaches: autoregressive modelling, sliced score matching...

Our approach

Apply a U-Net to estimate μ_{MLE}

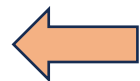


Train the model with a simple MAP loss function

$$\log \mathcal{L}(\vec{\mu}, \mathbf{Q}) = \mathbb{E}_{\vec{x} \sim p(\vec{x})} \left[\frac{1}{2} (\vec{x}_{(i)} - \vec{\mu})^T \mathbf{Q} (\vec{x}_{(i)} - \vec{\mu}) - \frac{1}{2} \text{tr}(\log \mathbf{Q}) \right]$$

\mathbf{Q} matrix diagonal in Fourier space, depends only on $|\mathbf{k}|$

Super-fast sampling from a Gaussian

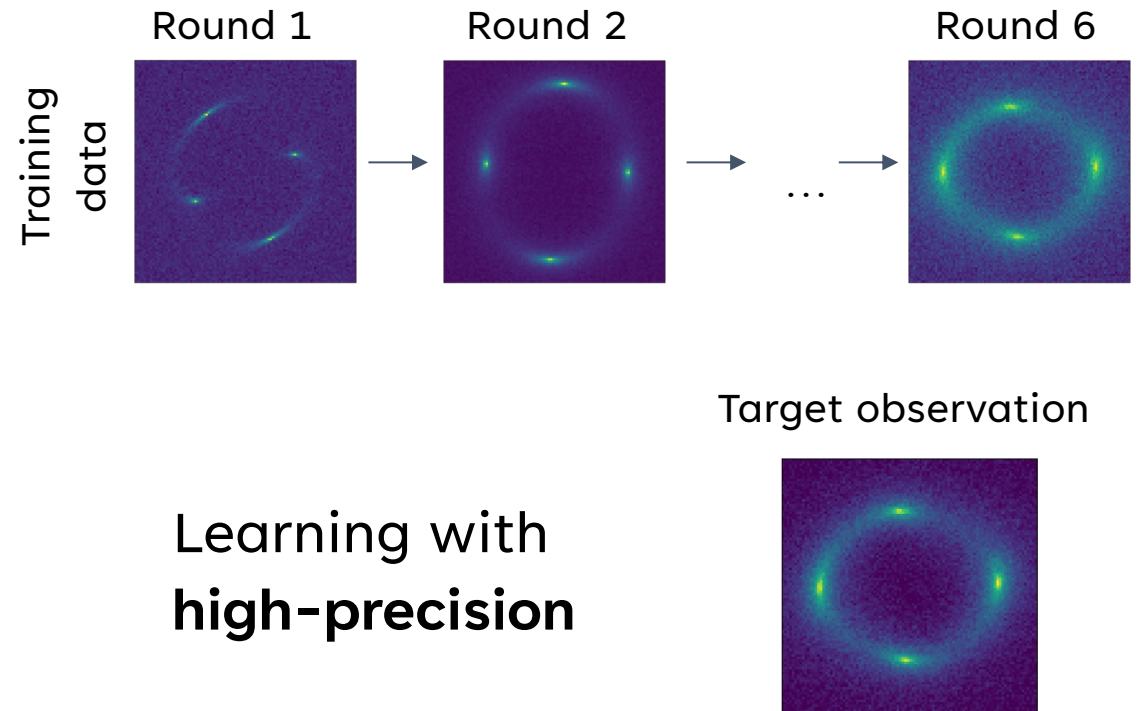
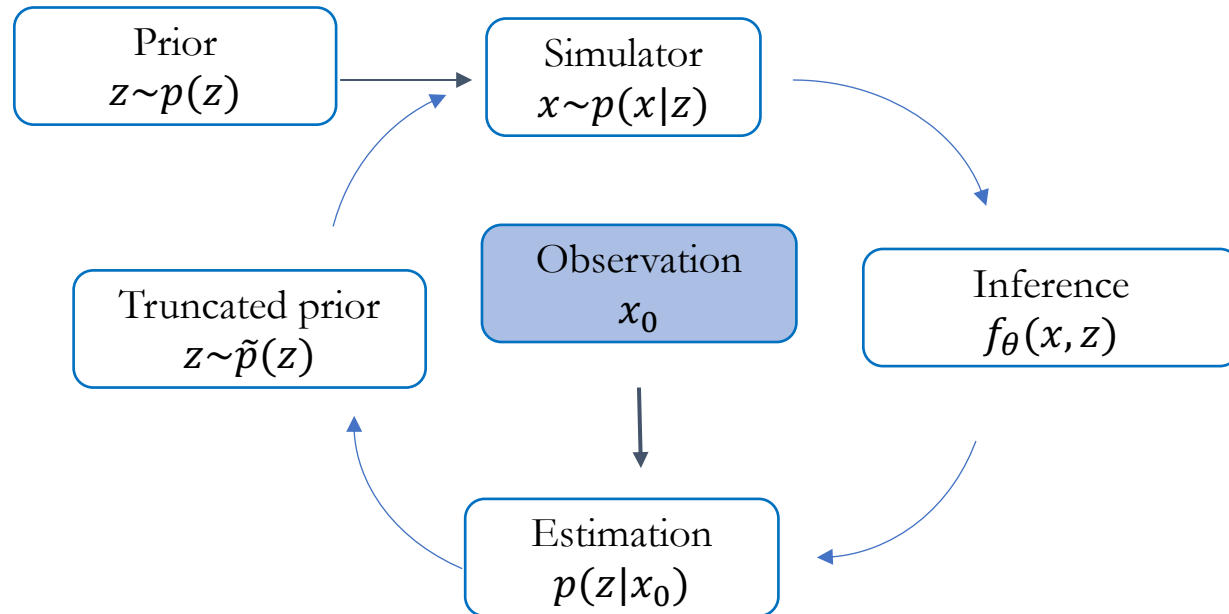


Learn correlation matrix and embedding network simultaneously



Sequential inference

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



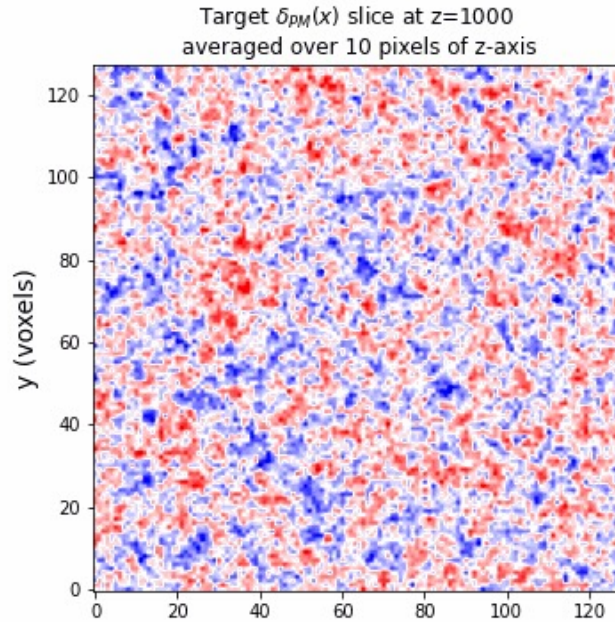
Results

1 hour of training

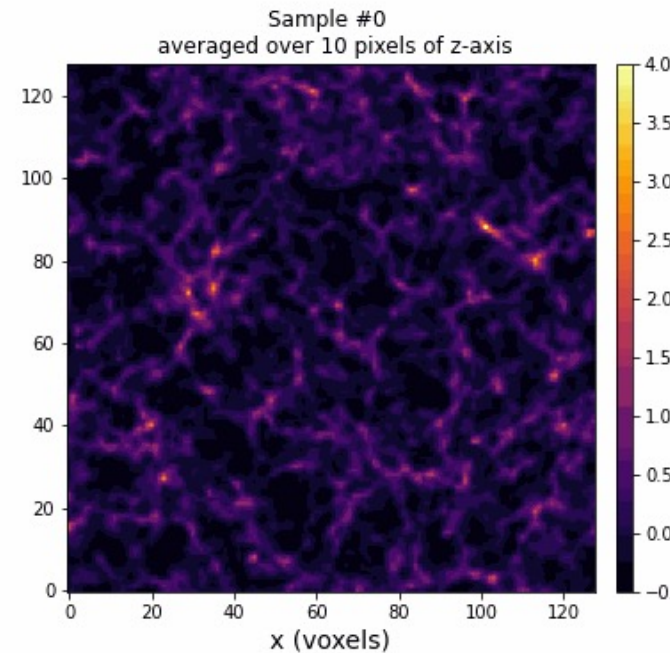
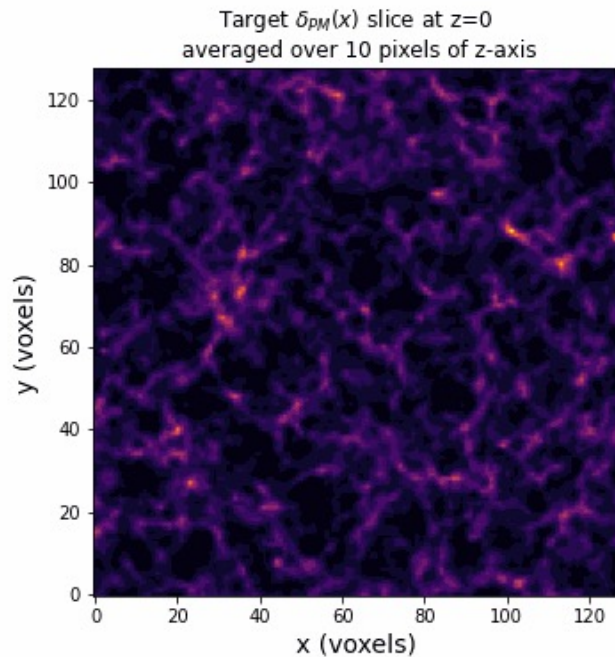
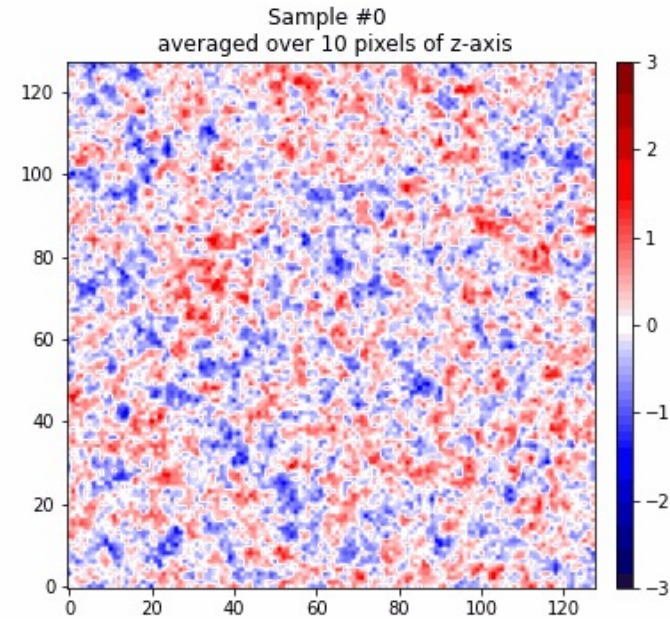
1 GPU

NVIDIA 40GB A100

Target image



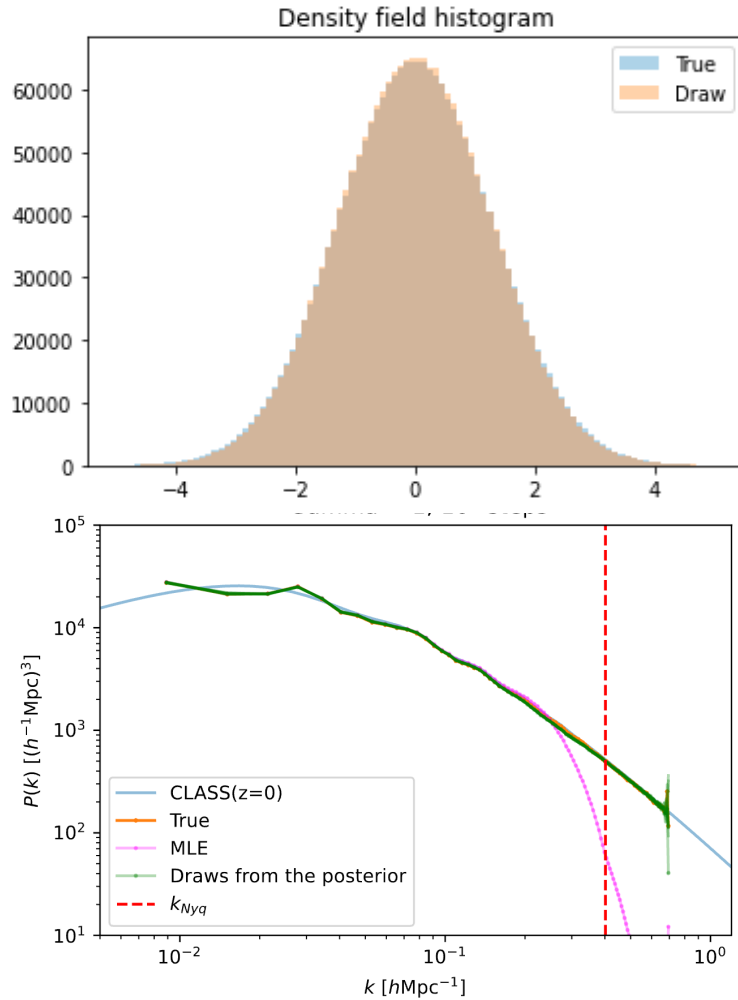
Posterior draws



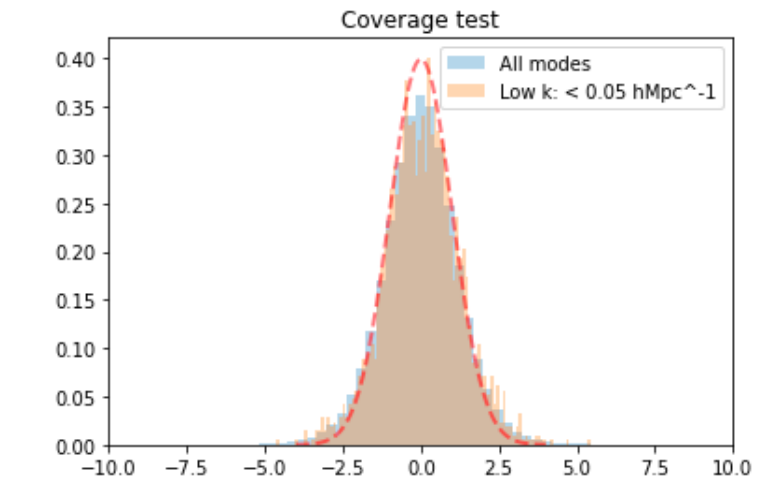
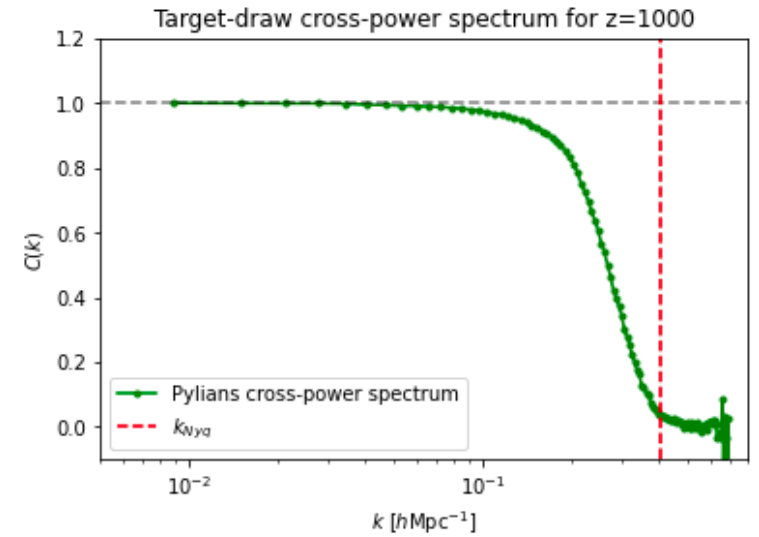
**$(1 \text{ Gpc}/h)^3$
volume**

**128^3
voxels**

Summary statistics comparison



1-2% agreement in the power spectrum



Coverage test shows that samples follow the correct distribution

Summary

- Reconstruction of cosmological initial conditions is an important problem that allows to analyse LSS data in the fullest way
- Methods like BORG achieve the goal, but are very slow, require gradients, run for months on supercomputers...
- Most ML approaches do point estimates, not explore the full posterior
- Our approach does this full exploration, is easy to train and produces samples in a super-fast way; it allows to turn any point estimator into a sampler
- The method is flexible and allows to do inference in an effective sequential way, with many applications in the future!