# 100,000 Universes: Simulation-Based Inference from Galaxy Clustering with `muchísimocks`

**Kate Storey-Fisher**

*with Raúl Angulo, Marcos Pellejero-Ibañez, ++*

Stanford / KIPAC | October 1, 2025

New Physics from Galaxy Clustering | Galileo Galilei Institute

*these slides at tinyurl.com/ksf-ggi-2025*

TNG-100-1, stars, z=0

# 100,000 Universes: Simulation-Based Inference from Galaxy Clustering with `muchísimocks`

*"So many sims, like a lot of sims"*

**Kate Storey-Fisher**

*with Raúl Angulo, Marcos Pellejero-Ibañez, ++*
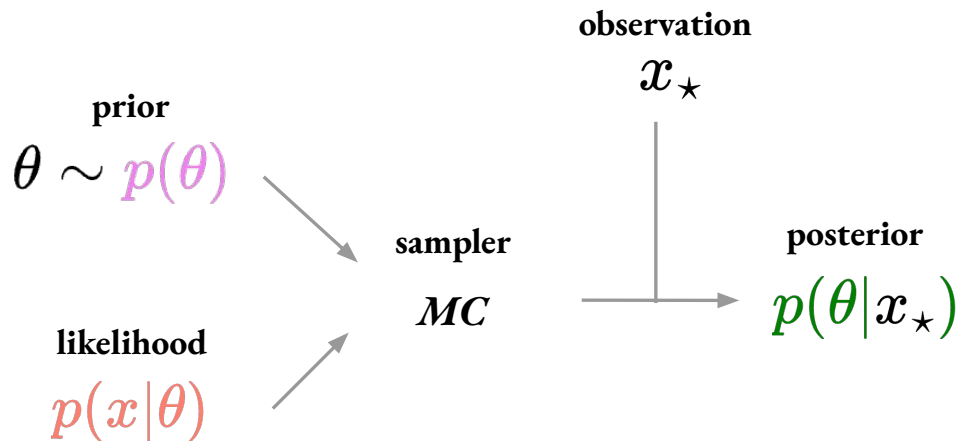
Stanford / KIPAC | October 1, 2025

New Physics from Galaxy Clustering | Galileo Galilei Institute

*these slides at tinyurl.com/ksf-ggi-2025*
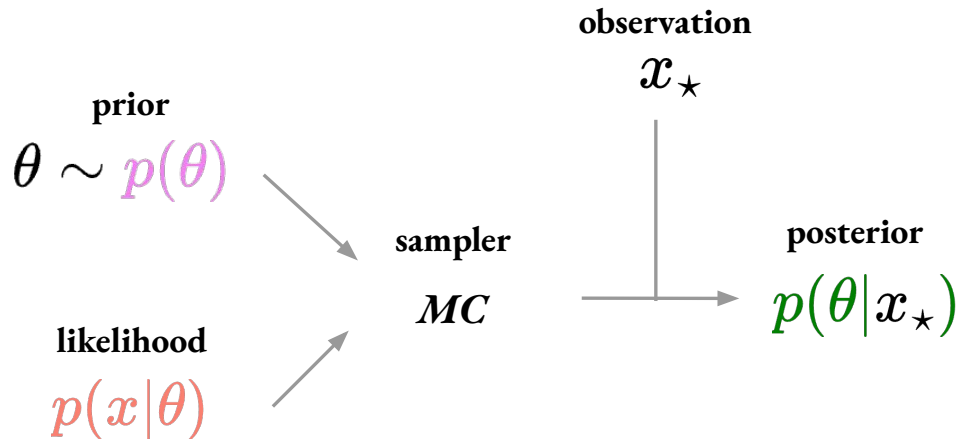
TNG-100-1, stars, z=0

# Likelihood-based inference

- Requires theoretical modeling of statistics
- Hard to model high-order statistics
- Hard to model impact of systematics
- Typically assumes Gaussian likelihood

**prior**

$$\theta \sim p(\theta)$$

**likelihood**

$$p(x|\theta)$$

**sampler**

$$MC$$

**observation**

$$x_\star$$

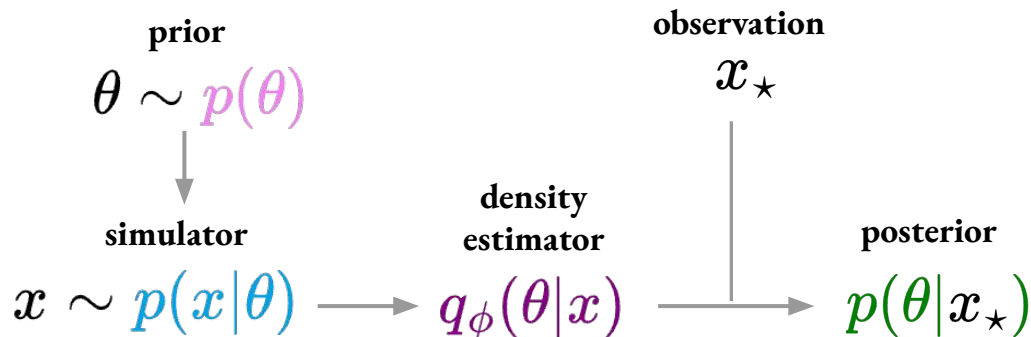**posterior**

$$p(\theta|x_\star)$$

# Likelihood-based inference

- Requires theoretical modeling of statistics
- Hard to model high-order statistics
- Hard to model impact of systematics
- Typically assumes Gaussian likelihood

**prior**
$$\theta \sim p(\theta)$$

**likelihood**
$$p(x|\theta)$$

**sampler**
$$MC$$

**observation**
$$x_\star$$

**posterior**
$$p(\theta|x_\star)$$

# Simulation-based inference

- Only requires forward model of data
- Naturally good for high-order statistics
- Can forward-model systematics
- Does not assume Gaussian likelihood

SBI for galaxy clustering limited by
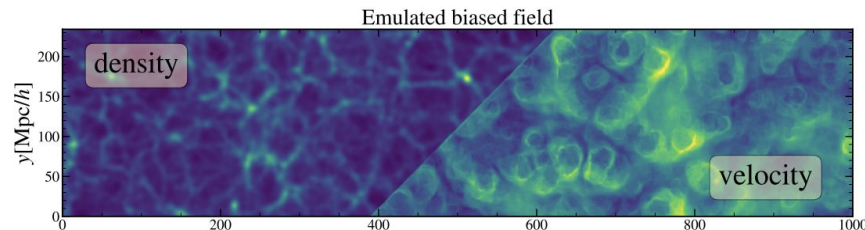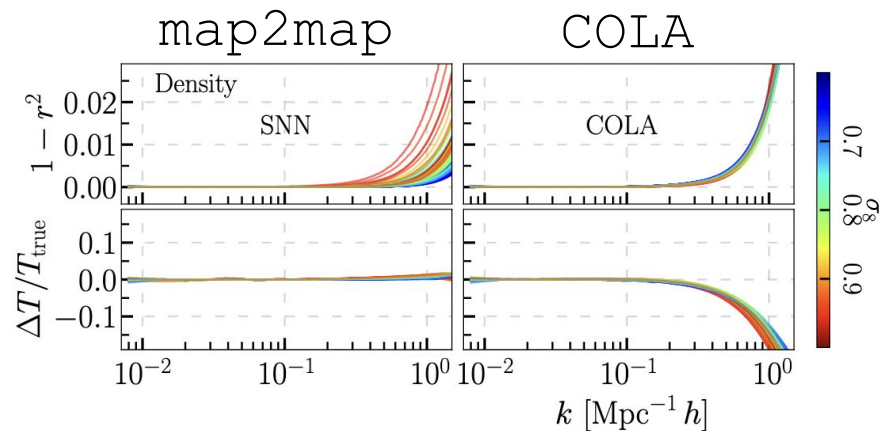cost of training simulations

**prior**
$$\theta \sim p(\theta)$$

**simulator**
$$x \sim p(x|\theta)$$

**density estimator**
$$q_\phi(\theta|x)$$

**observation**
$$x_\star$$

**posterior**
$$p(\theta|x_\star)$$

# The `muchísimocks` forward model: the matter distribution

We use the **map2map full-field emulator** (Jamieson+2022): Predicts the *N*-body displacement and velocity fields from an initial ZA approximation.

- ○ Trained on the `Quijote` simulations: 1000 simulations, 5 cosmological parameters, 1 (Gpc/h)$^3$, 512$^3$ particles

- ○ %-level accuracy on *P(k)* down to *k*=1 Mpc/*h*

- ○ Takes ~4 minutes per emulated field





Emulated biased field

# The `muchísimocks` forward model: galaxy bias

We use the **hybrid Lagrangian bias expansion** to model how galaxies populate the matter distribution:

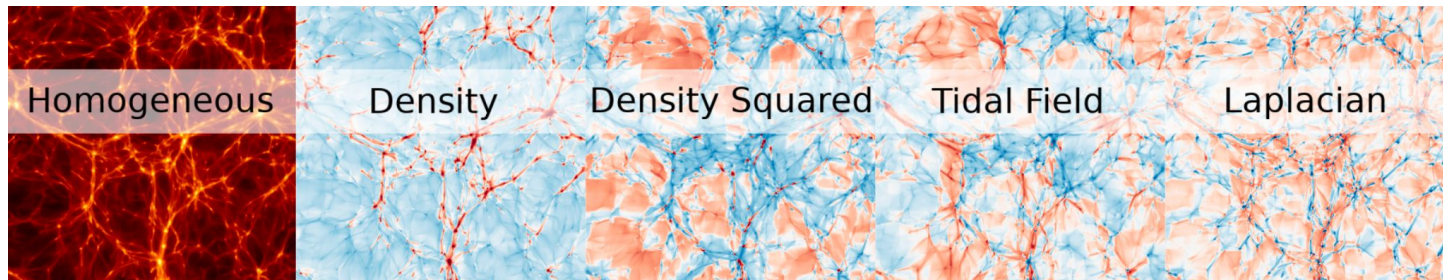Model galaxy density as function of matter density, expanding to $2^{nd}$ order:

$$\delta_g(\boldsymbol{q}) = 1 + b_1\,\delta_L(\boldsymbol{q}) + b_2\left[\delta_L^2(\boldsymbol{q}) - \langle\delta^2\rangle\right] + b_{s^2}\left[s^2(\boldsymbol{q}) - \langle s^2\rangle\right] + b_\Delta\Delta^2\delta(\boldsymbol{q})$$

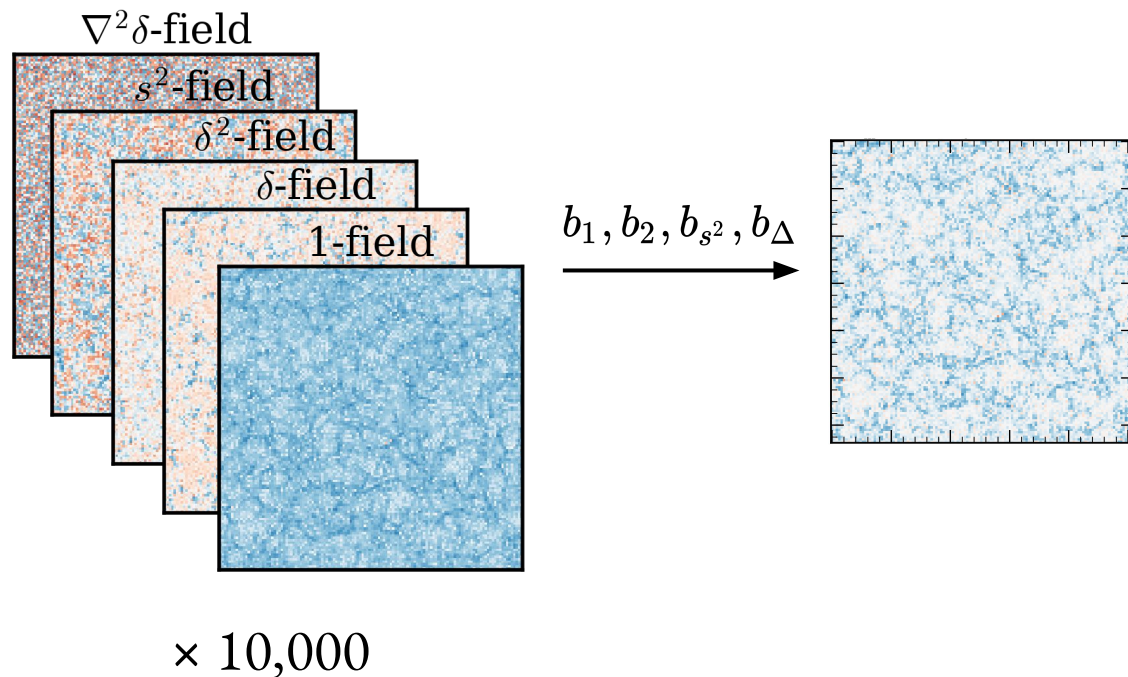Advect Lagrangian to Eulerian space, using simulations to compute the displacement field $\Psi(\boldsymbol{q})$:
$$\boldsymbol{x} = \boldsymbol{q} + \Psi(\boldsymbol{q})$$

Our model parameters are the bias coefficients,
$$b_1\,,\,b_2\,,\,b_{s^2}\,,\,b_\Delta$$



Homogeneous    Density    Density Squared    Tidal Field    Laplacian

# `muchísimocks`: a library of 3D biased tracer fields

$\nabla^2 \delta$-field

$s^2$-field

$\delta^2$-field

$\delta$-field

1-field

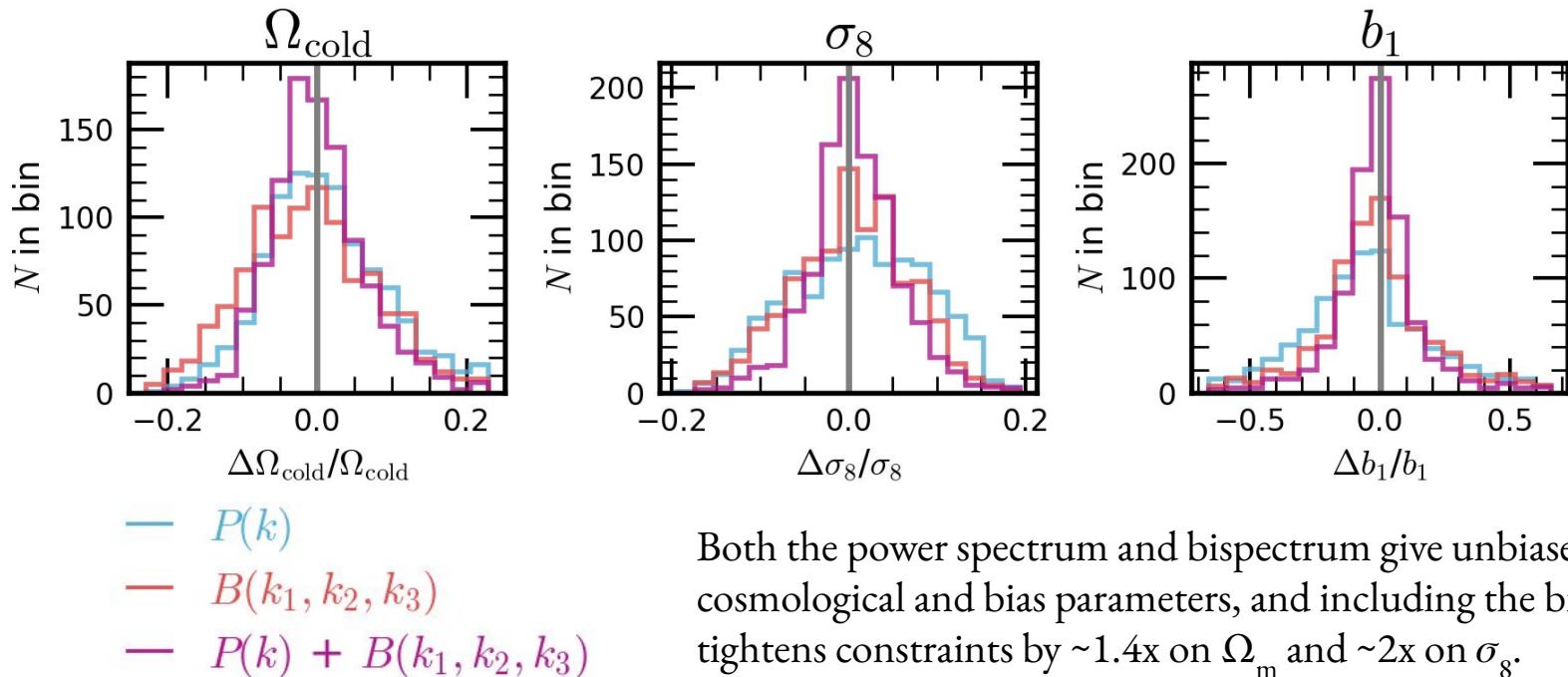$\times$ 10,000

$$b_1, b_2, b_{s^2}, b_\Delta$$

○ 10,000 sets of bias fields, with cosmologies spanning a Latin hypercube of $\sigma_8$, $\Omega_m$, $h$, $\Omega_b$

○ Can combine bias fields with choice of 2$^{nd}$ order bias parameters to construct density field

○ Each $(1\ h^{-1}\text{Gpc})^3$, $512^3$ resolution → cut small-$k$ modes to get $128^3$

○ ~1 TB of mocks (storage-limited); will be publicly available

# Simulation-based inference recovery test with *P(k)* & bispectrum
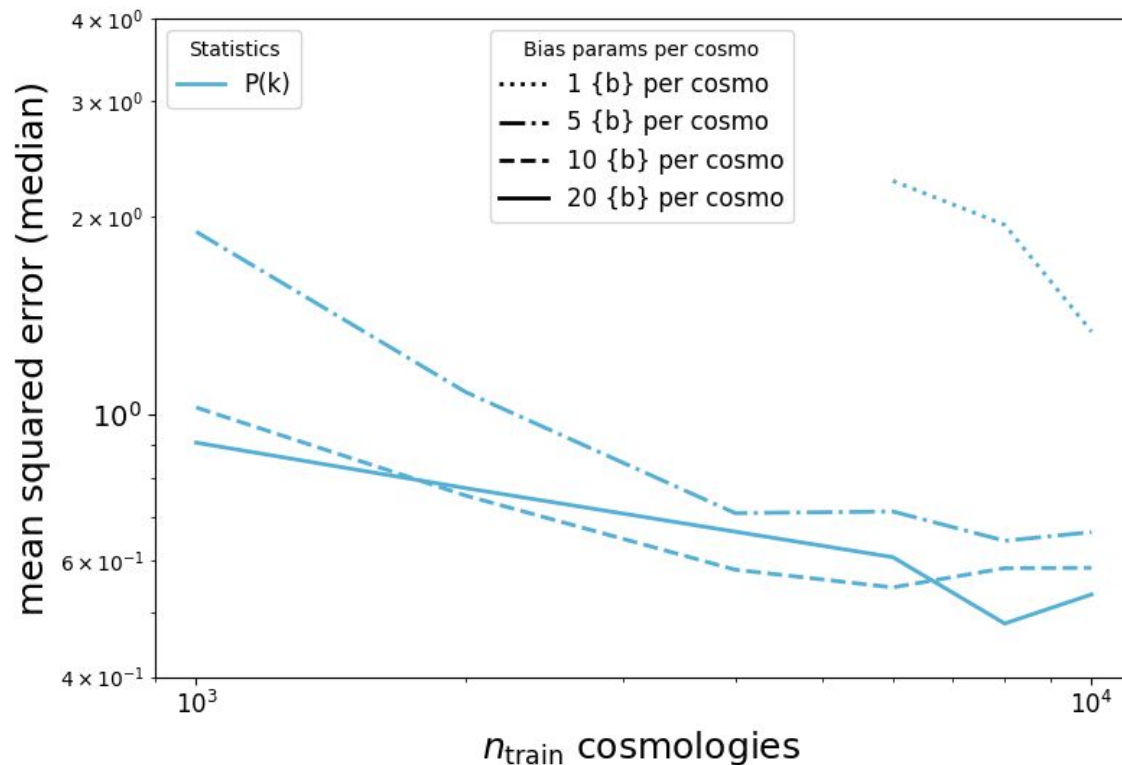
**training:** 10,000 cosmologies, 20x {$b$}
**testing:** 1,000 mocks spanning training set



Both the power spectrum and bispectrum give unbiased cosmological and bias parameters, and including the bispectrum tightens constraints by ~1.4x on $\Omega_m$ and ~2x on $\sigma_8$.
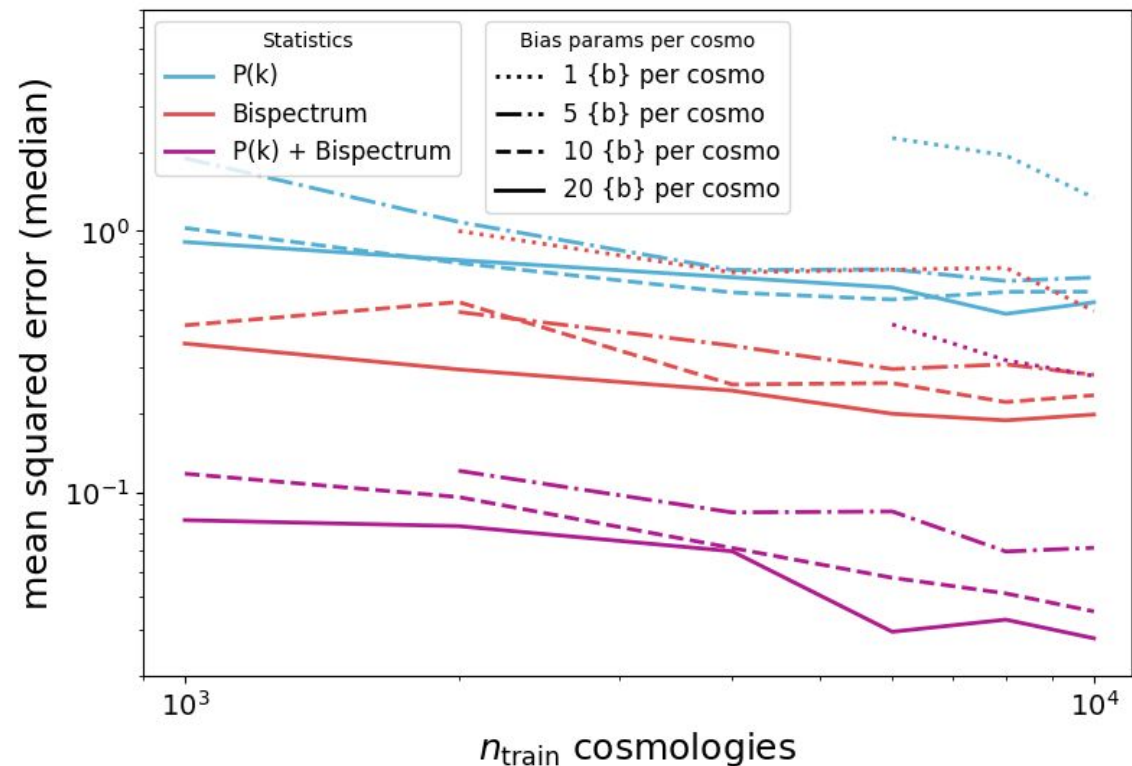
# Inference results: Convergence with size of training data



The accuracy & precision improves significantly from 1→20 bias parameter sets per cosmology, and from 1k→10k cosmologies, with ~convergence around 6k cosmologies.
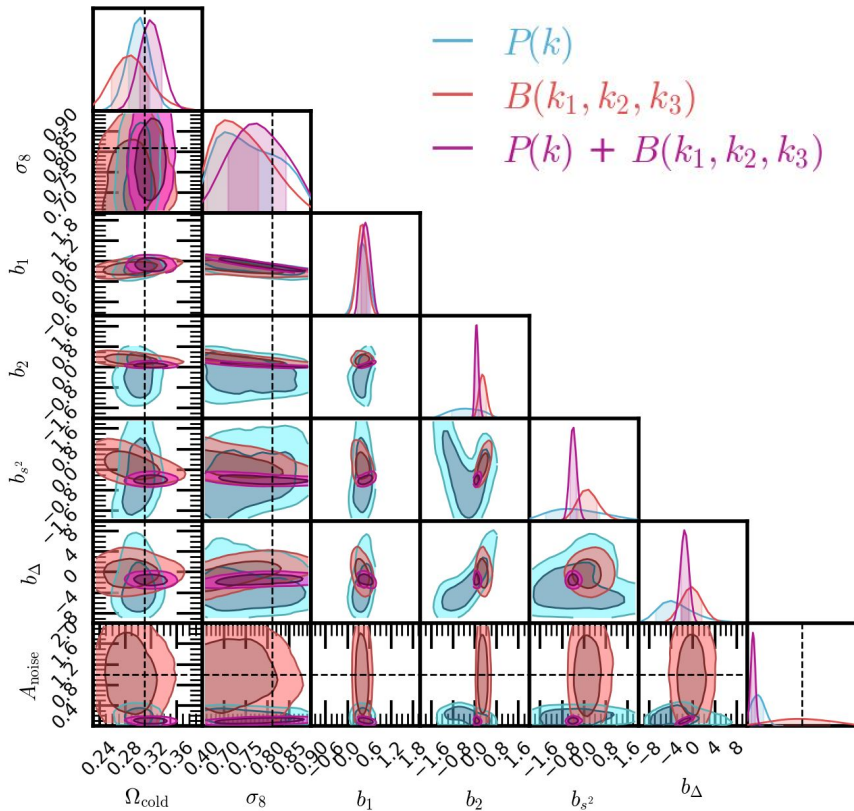
# Inference results: Convergence with size of training data



The accuracy & precision improves significantly from 1→20 bias parameter sets per cosmology, and from 1k→10k cosmologies, with ~convergence around 6k cosmologies.

# Out-of-distribution test: BACCO high-res $N$-body + SHAMe galaxies

training: 10,000 cosmologies, 20x $\{b\}$



- $P(k)$
- $B(k_1, k_2, k_3)$
- $P(k) + B(k_1, k_2, k_3)$

The model is robust to out-of-distribution test data, recovering unbiased cosmological parameters (with the noise parameter likely absorbing some model mismatch).

# Summary & Outlook

○ We present the **`muchísimocks`** library of 10,000 density fields: uses a **field-level emulator** of $N$-body simulations to span cosmological parameter space, and the **hybrid Lagrangian bias expansion** for galaxy bias space; the approach is a path forward for simulation-limited galaxy clustering analyses.

○ We trained a **simulation-based inference** pipeline on the power spectrum and bispectrum, finding good convergence with the number of cosmologies and bias parameter sets (200,000 total training sims). Tests on an out-of-distribution SHAMe mock demonstrate robustness.

○ `muchísimocks` can be utilized for a range of analyses, such as full-field inference, testing beyond-standard statistics, velocity reconstruction, etc *[reach out with ideas!]* - publicly available soon!

○ *Upcoming/in-progress:* Redshift-space, scale dependence, bias expansion checks, bispectrum `map2map` checks, ++

kstoreyf
@stanford.edu
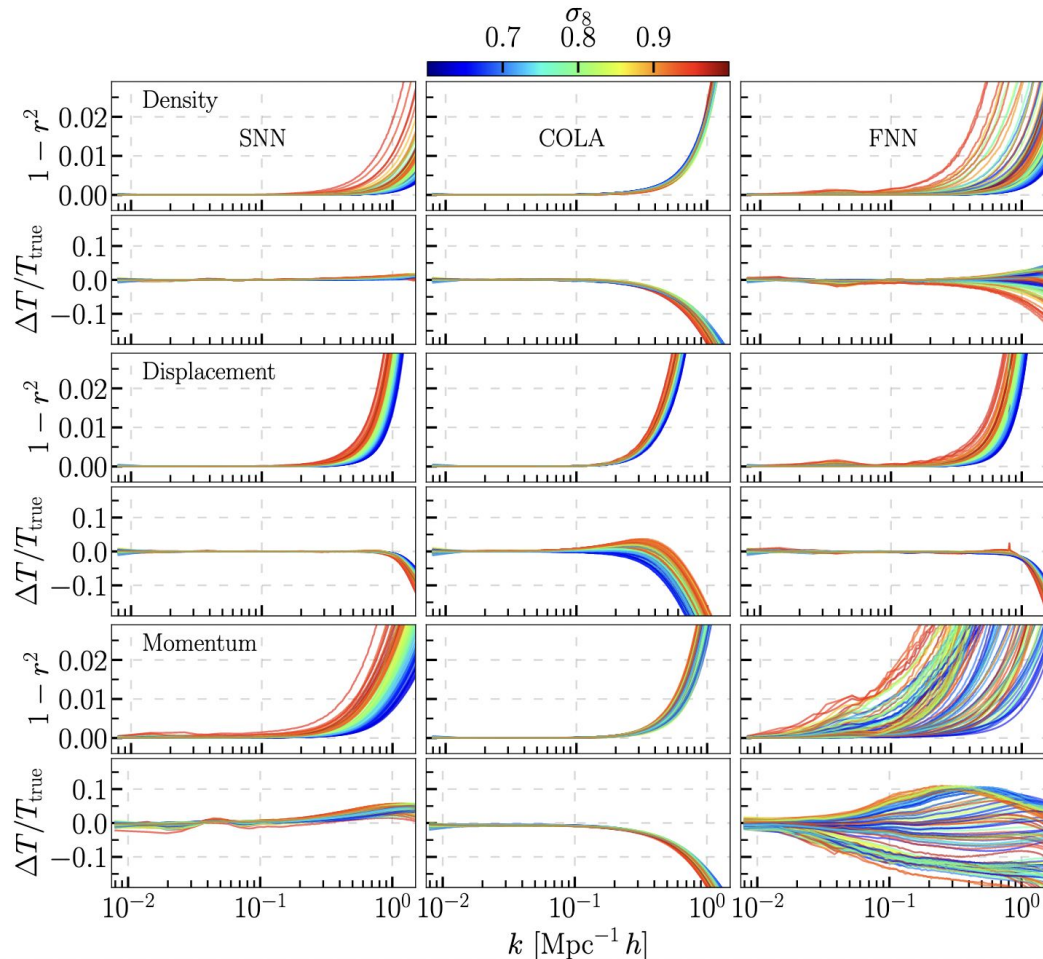
@kstoreyf

cosmo.nyu.edu/ksf

# Extra Slides

TNG-100-1, stars, z=0

# map2map accuracy



Jamieson+2022
([2206.04594](https://arxiv.org/abs/2206.04594))
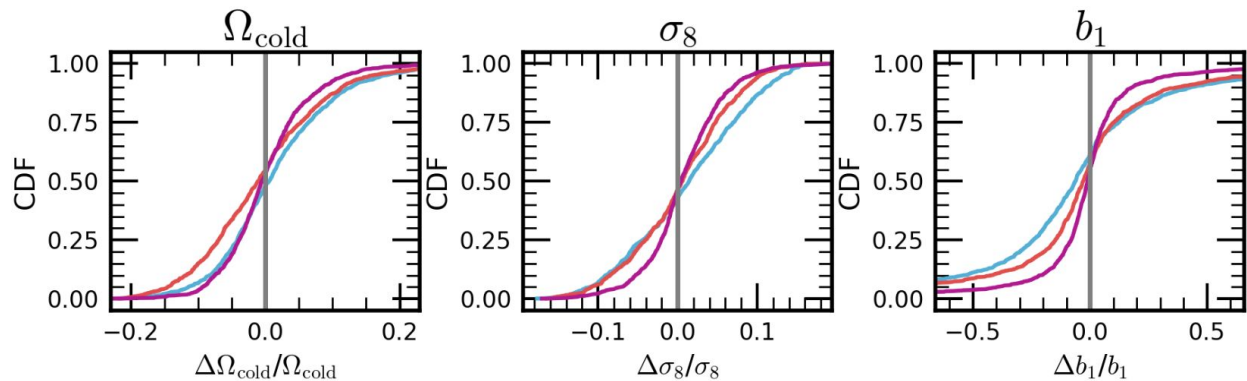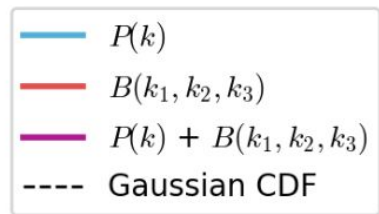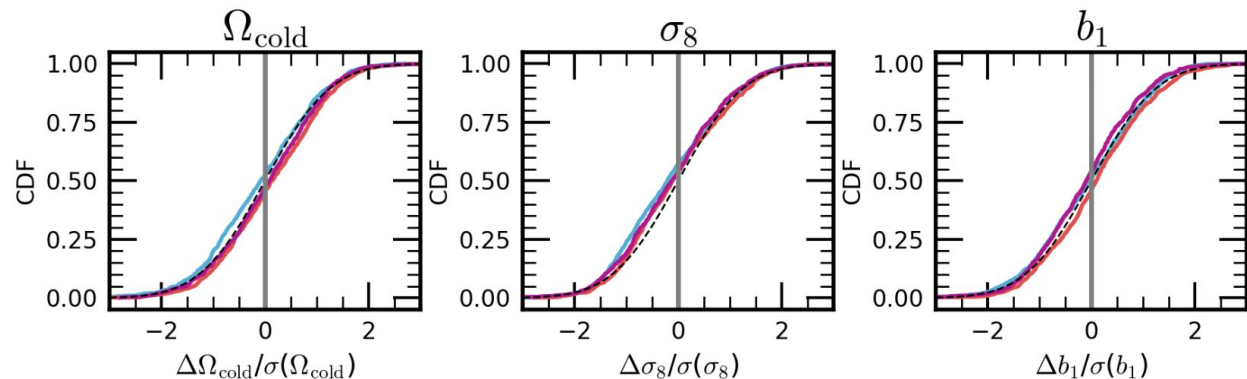
**Fig. 1.** Power spectra errors for the styled neural network (left), COLA (middle), and the fiducial neural network (right). Each pair of rows shows the stochasticities and transfer function errors for the Eulerian density field power spectra (top two), the Lagrangian displacement field power spectra (middle two), and the Eulerian momentum field (bottom two) power spectra. The color of each curve corresponds to its value of $\sigma_8$ according to the color bar at the top.

# Simulation-based inference results: *P(k)* & bispectrum, 20x{*b*}

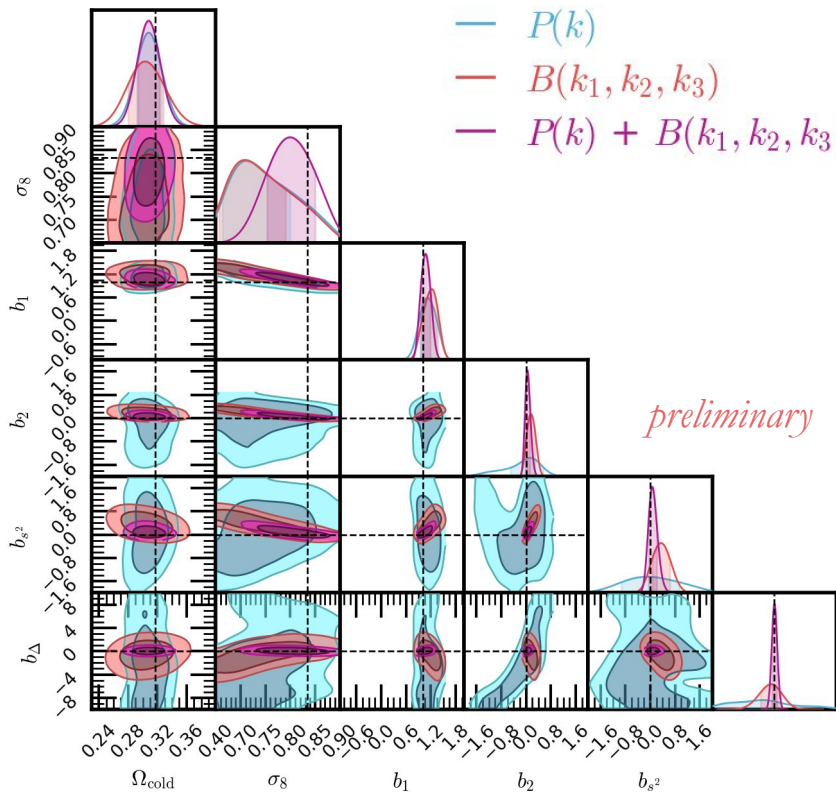coverage test: test set spanning training set space (cosmo & bias varied)



*more step-function histogram = more precise & accurate*

*more gaussian (like dashed line) = more correct uncertainties*

# Inference recovery test results

recovery test: mean of 1000 fixed cosmo mocks   (20x {b})

# Simulation-based inference: neural posterior estimation

Train a neural network to *approximate* the posterior by minimizing the Kullback-Leibler divergence between it and the true posterior, with trainable parameters $\phi$:

$$\min_{\phi} D_{\mathrm{KL}} \left( p(\theta|x)p(x) \| q_{\phi}(\theta|x)p(x) \right) = \min_{\phi} \int p(\theta, x) \log \frac{p(\theta|x)}{q_{\phi}(\theta|x)} d\theta dx$$

*integral→sum because we have discrete samples*

$$\approx \min_{\phi} \sum_i \log p(\theta_i|x_i) - \log q_{\phi}(\theta_i|x_i)$$

*the true posterior is independent of NN parameters*

$$\approx \min_{\phi} \sum_i - \log q_{\phi}(\theta_i|x_i)$$

*minimizing the negative log-likelihood is the same as maximizing the log-likelihood*

$$\approx \max_{\phi} \sum_i \log q_{\phi}(\theta_i|x_i)$$

# Simulation-based inference: density estimation

We use a *normalizing flow* as our model for approximating the posterior, specifically a *masked autoregressive flow*.

a base variable $\mathbf{z}$ is sampled from a Gaussian

the $f$'s are conditioned on the observable $x$

*the flow thus learns to transform noise into samples from the posterior*

$$q_\phi(\theta|x) = \mathcal{N}(\mathbf{z}_0|\mathbf{0}, \mathbf{I}) \prod_{t=1}^{T} \left| \det\left( \frac{\partial f_t}{\partial \mathbf{z}_{t-1}} \right) \right|^{-1}$$

the $\mathbf{z}$'s are transformed via invertible functions $f$ (NN layers) to resemble the parameters $\theta$

$$\mathbf{z}_0 \xrightarrow{f_1} \mathbf{z}_1 \xrightarrow{f_2} \mathbf{z}_2 \xrightarrow{\cdots} \mathbf{z}_T = \boldsymbol{\theta}$$