

Benchmarking field-level inference from galaxy surveys

Hugo Simonfroy,

Ph.D. student supervised by *Arnaud de Mattia, François Lanusse*



The universe recipe (so far)

Cosmological principle + Einstein equation

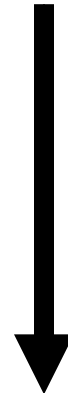
+ Inflation



$$\frac{H}{H_0} = \sqrt{\Omega_r + \Omega_b + \Omega_c + \Omega_\kappa + \Omega_\Lambda}$$

energy content

instantaneous
expansion rate



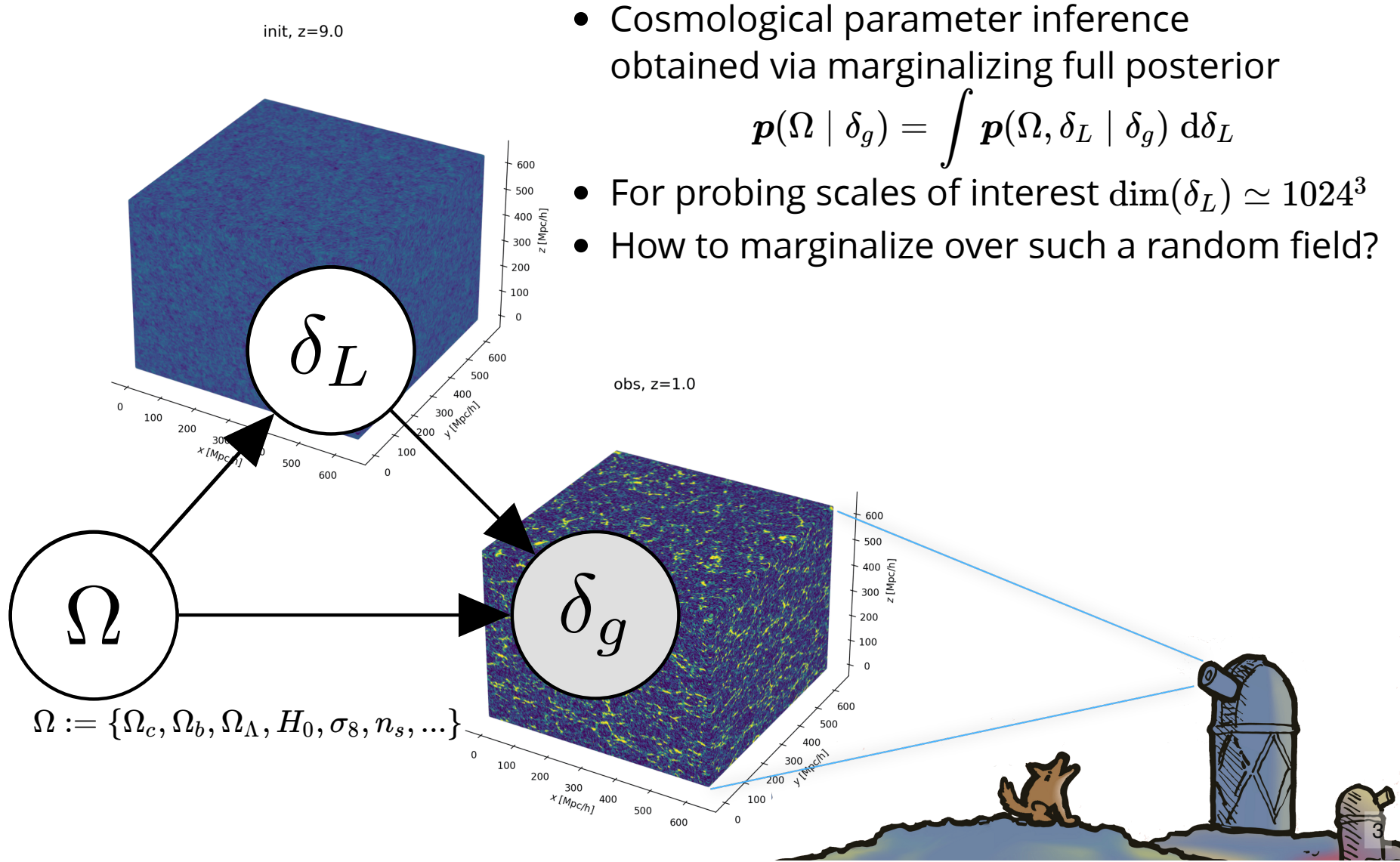
$$\delta_L \sim \mathcal{G}(0, \mathcal{P})$$

initial field primordial
power spectrum

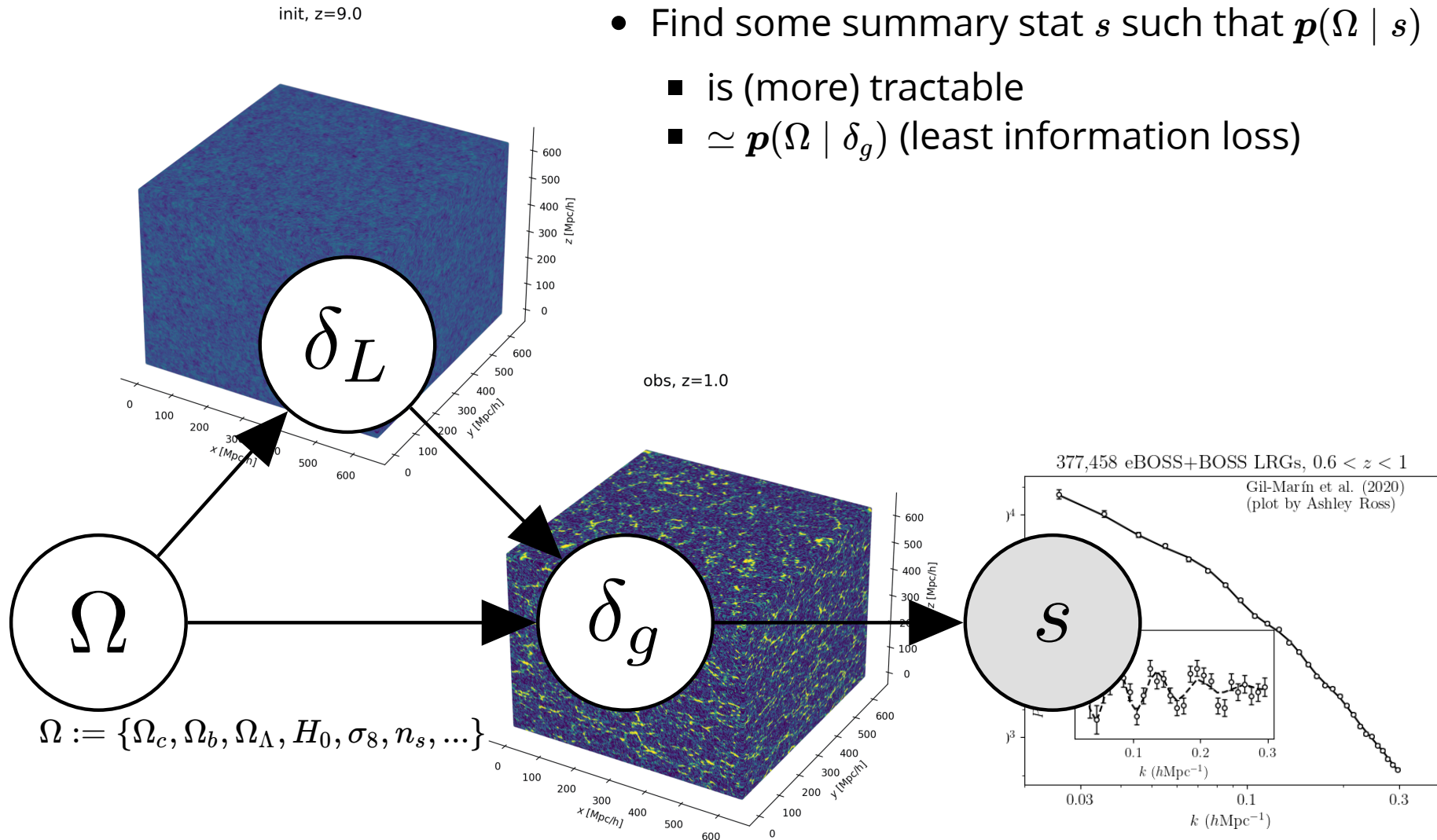
$$\sigma_8 := \sigma(\delta_L * \Pi_8)$$

std. of fluctuations smoothed at 8 Mpc/h

A high-dimensional inference problem

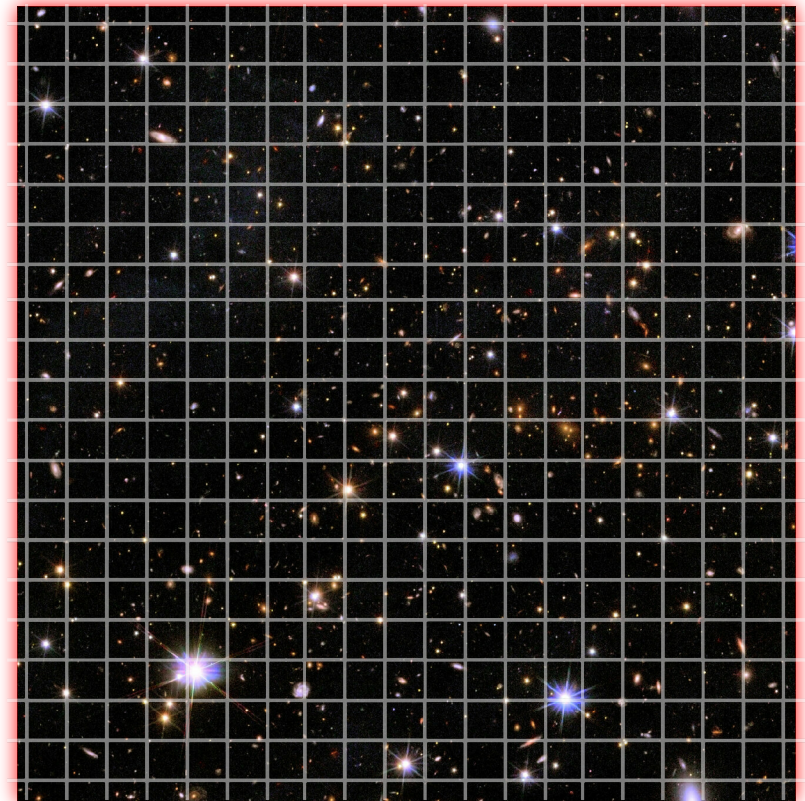
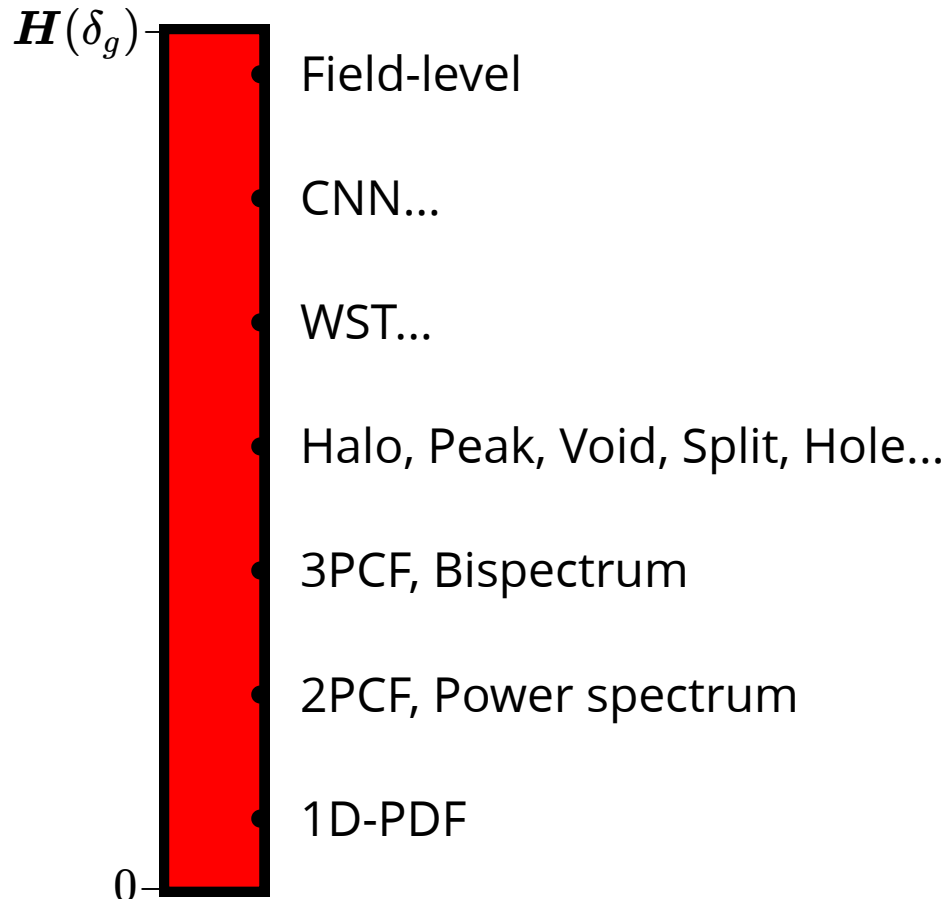


Use summary statistics



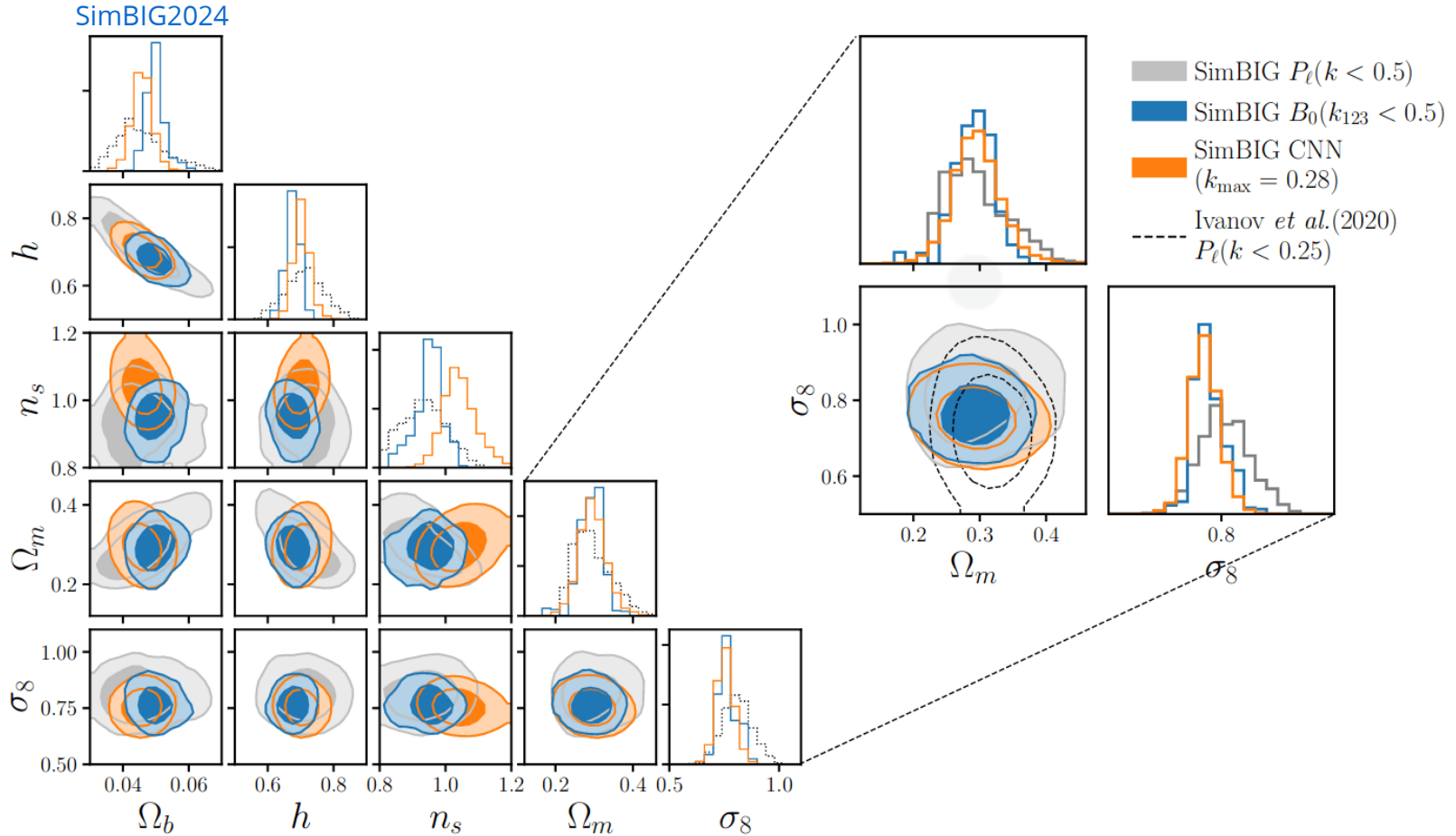
We gotta pump this information up

- At large scales, matter density field almost Gaussian so power spectrum is almost lossless compression.
- To prospect smaller non-Gaussian scales, let's use:



Euclid's view of Perseus

More information is better, but how much better?



Let's push one step further and use the whole field:

field-level inference

Some useful programming tools

- **JAX**

- GPU acceleration
- Just-In-Time (JIT) compilation acceleration
- Automatic vectorization/parallelization
- Automatic differentiation



- **NumPyro**

- Probabilistic Programming Language (PPL)
- Powered by JAX
- Integrated samplers



So JAX in practice?

- GPU accelerate

```
1 import jax.numpy as np
2 # then enjoy
```

- JIT compile

```
1 function = jax.jit(function)
2 # function is so fast now!
```

- Vectorize/Parallelize

```
1 vfunction = jax.vmap(function)
2 pfunction = jax.pmap(function)
3 # for-loops are for-loosers
```

- Auto-diff

```
1 gradient = jax.grad(function)
2 # too bad if you love chain ruling by hand
```



Now let's build a cosmological model

1. Prior on

- Cosmology Ω
- Initial field δ_L
- Dark matter-galaxy connection (Lagrangian galaxy biases) b

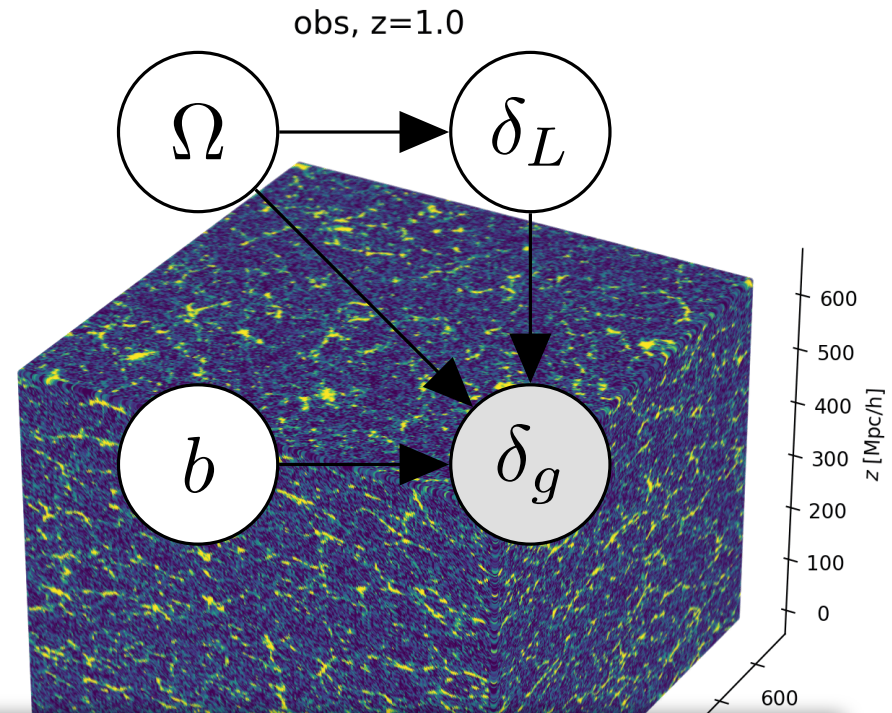
2. Initialize matter particles

3. LSS formation (LPT+PM)

4. Populate matter field with galaxies

5. Galaxy peculiar velocities (RSD)

6. Observational noise

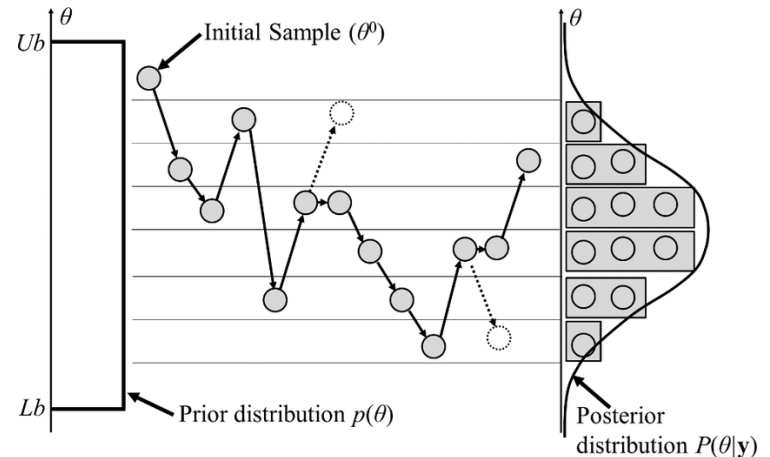


- **Fast** and **differentiable** model
- $\simeq 1024^3$ parameters is huge!
- Need inference methods that **scale to high dimensions**
- Some proposed by [Lavaux+2018](#), [Bayer+2023](#)

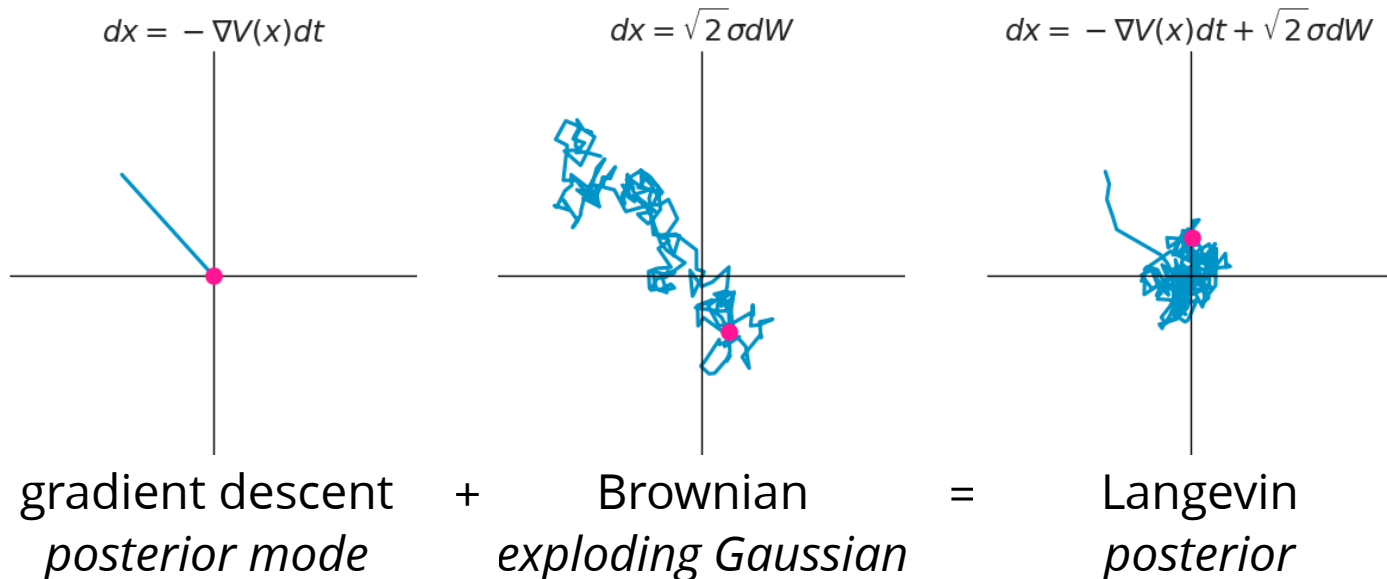
Why care about differentiable model?

- **Classical MCMCs**

- agnostic random moves
+ MH acceptance step
= **blinded natural selection**
- small moves yield **correlated samples**.



- **s.o.t.a. MCMCs** rely on the **gradient of the model log-proba**, to drive dynamic towards highest density regions.

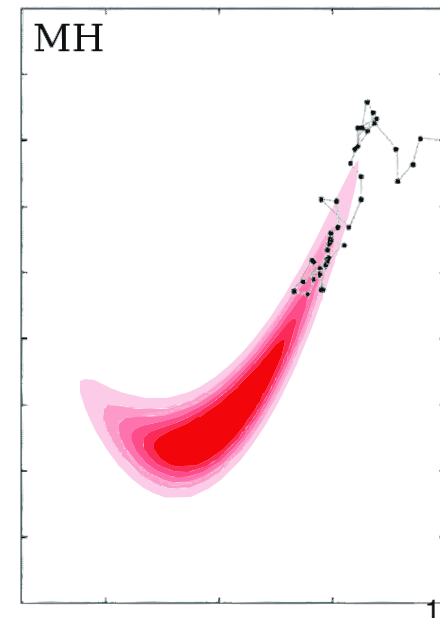
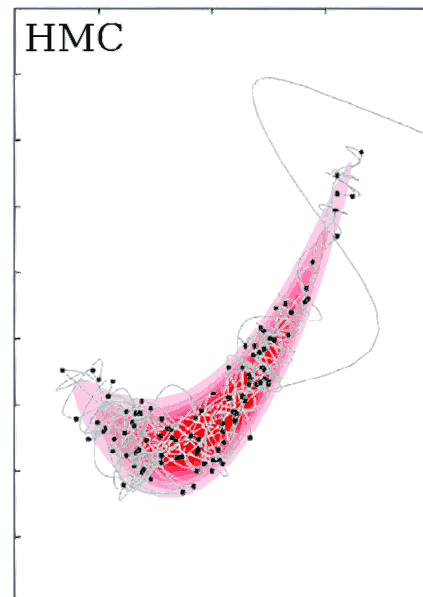


Hamiltonian Monte Carlo (HMC)

- To travel farther, add inertia.
 - sample particle at **position** q now have **momentum** p and **mass matrix** M
 - target $\mathbf{p}(q)$ becomes $\mathbf{p}(q, p) := e^{-\mathcal{H}(q, p)}$, with **Hamiltonian**
$$\mathcal{H}(q, p) := -\log \mathbf{p}(q) + \frac{1}{2} p^\top M^{-1} p$$
 - at each step, **resample momentum** $p \sim \mathcal{N}(0, M)$
 - let (q, p) follow the **Hamiltonian dynamic** during time length L , then arrival becomes new **MH proposal**.

Variations around HMC

- **No U-Turn Sampler (NUTS)**
 - trajectory length L auto-tuned
 - samples drawn along trajectory
- **NUTSGibbs** i.e. alternating sampling over parameter subsets.



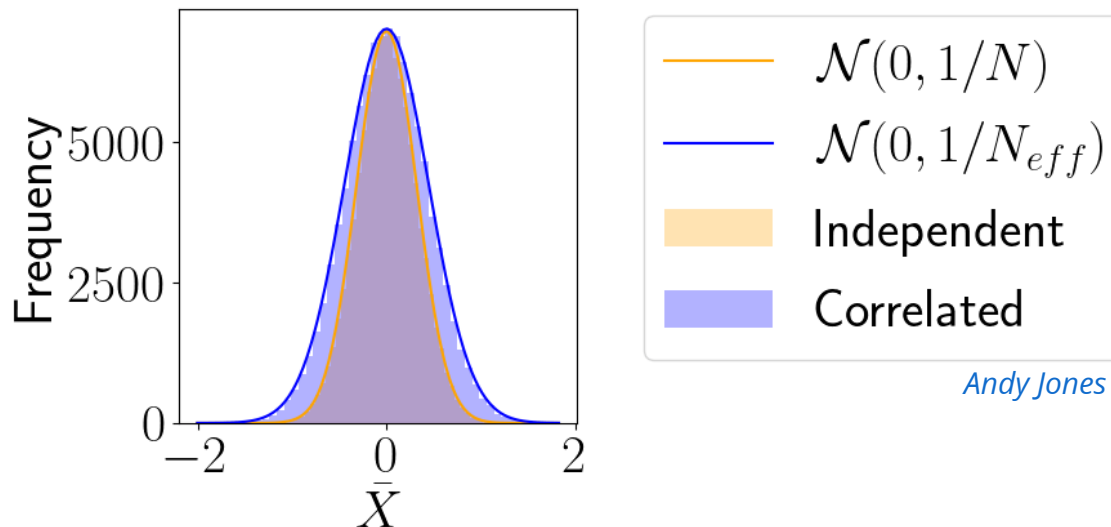
How to compare samplers?

- **Effective Sample Size (ESS)**

- number of i.i.d. samples that yield same statistical power.
- For sample sequence of size N and autocorrelation ρ

$$N_{\text{eff}} = \frac{N}{1 + 2 \sum_{t=1}^{+\infty} \rho_t}$$

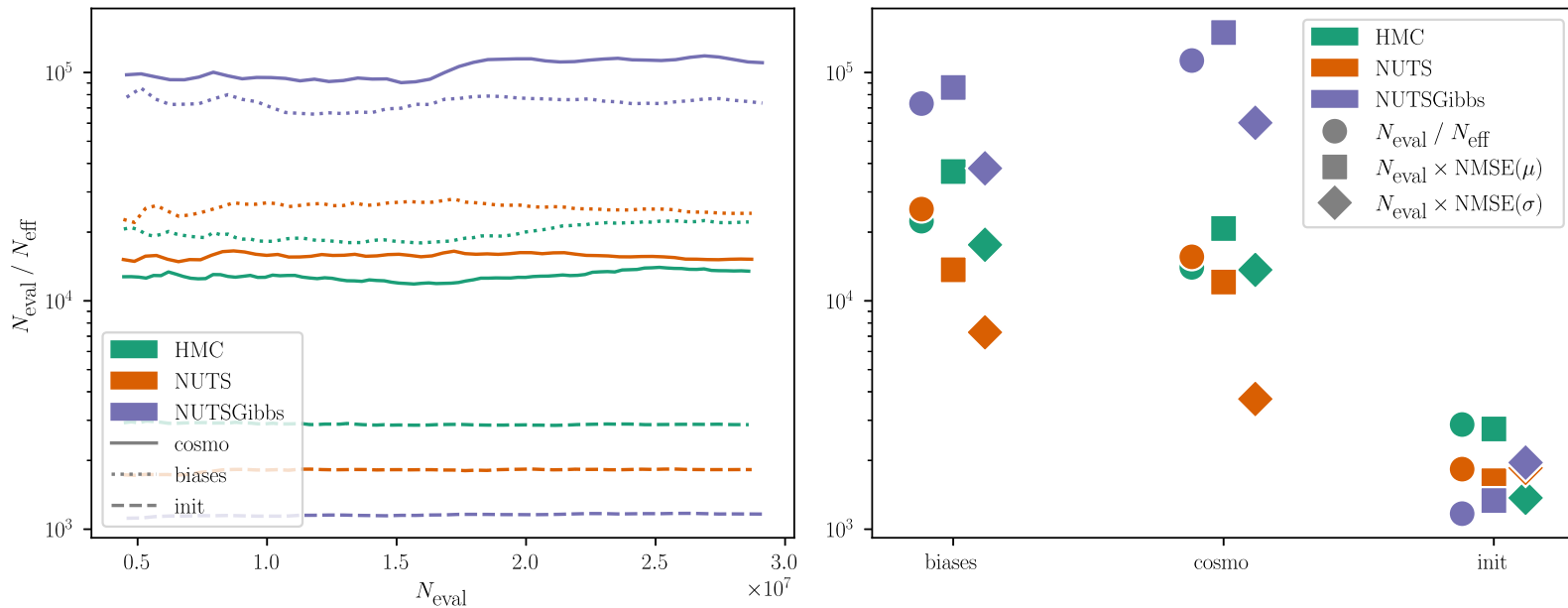
so aim for as less correlated sample as possible.



- Main limiting computational factor is **model evaluation** (e.g. N-body), so characterize MCMC efficiency by $N_{\text{eval}}/N_{\text{eff}}$

Benchmarking

- **model setting:** 64^3 mesh, $(640 \text{ Mpc}/h)^3$ box, 1LPT, second order Lagrangian bias expansion, RSD and Gaussian observational noise.
- **parameter space:** initial field δ_L , cosmology $\Omega = \{\Omega_m, \sigma_8\}$, and galaxy biases $b = \{b_1, b_2, b_{s^2}, b_{\nabla^2}\}$. Total of $64^3 + 2 + 4$ parameters.
- For NUTSGibbs: split sampling between δ_L and the rest (common in lit.)



- Results suggest no particular advantage to splitting sampling between initial field and rest.

Recap...

- **Field-level inference** may be relevant to fully capture cosmological information in data.
- Leverage modern computational tools to build **fast** and **differentiable** cosmological model.
- Performing field-level inference becomes tractable. Many proposed methods in literature.
- Standardized benchmark for **comparing s.o.t.a. MCMC samplers** on field-level inference tasks, selecting proposed methods for **Stage-IV galaxy surveys**.

...and what's next

- Include more proposed samplers.
- Compare to SBI approaches.
- Move towards real applications on DESI data.