



# Implicit Likelihood Inference in Cosmology while efficiently checking for survey systematics



National Action Dark Energy Colloquium 2023  
LAPP, Annecy

Tristan Hoellinger

[hoellin.github.io](https://hoellin.github.io)

Institut d'Astrophysique de Paris  
CNRS & Sorbonne Université

In collaboration with:  
Florent Leclercq (IAP – CNRS & Sorbonne Université)

and the Aquila Consortium

[www.aquila-consortium.org](http://www.aquila-consortium.org)

**5 November 2023**



# Challenges

- Systematic effects

Large surveys (Euclid, LSST) will be dominated by **systematic** rather than statistical uncertainty.

- How to integrate prior knowledge?

- Prior on cosmology  $\omega = (h, \Omega_b, \Omega_m, n_S, \sigma_8)$
- Prior on non measurable physical quantities
  - The initial matter power spectrum  $\theta$  should be close to  $\theta_{\text{Planck},2018}$ : can we benefit from this?
  - Can we benefit from theoretical insights on  $\theta$ ?

[Jasche & Lavaux 2017, 1706.08971](#)

[Laureijs et al., 2011, 1110.3193](#)

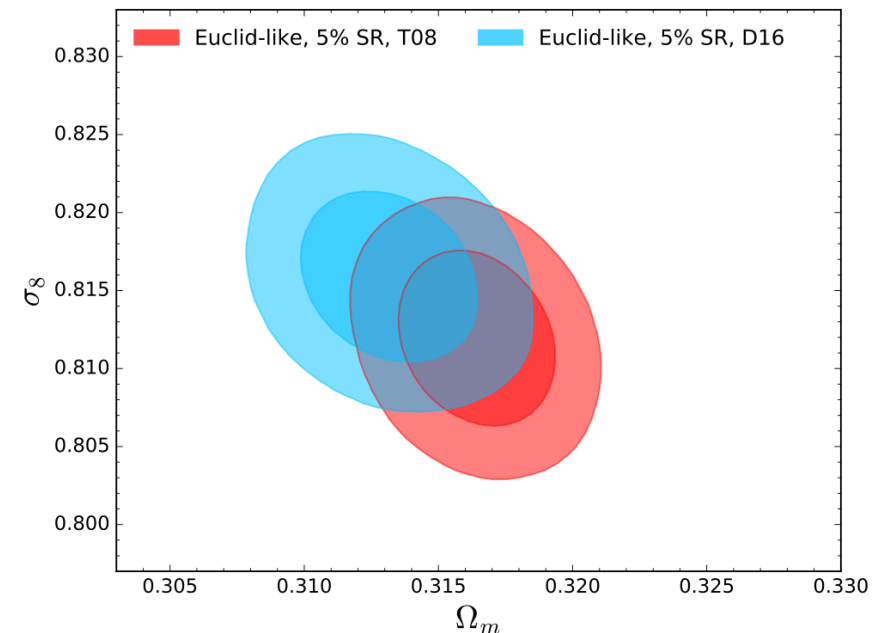


Tristan Hoellinger

- Model misspecification in Bayesian & simulation-based inference (SBI)

- When model differ from data-generating process
- Biased or overly concentrated posteriors

[Müller 2013, 10.3982/ECTA9097](#)



68% and 95% marginalized posterior via MCMC on simulated galaxy catalogues w/ two mass functions To8 and D16.

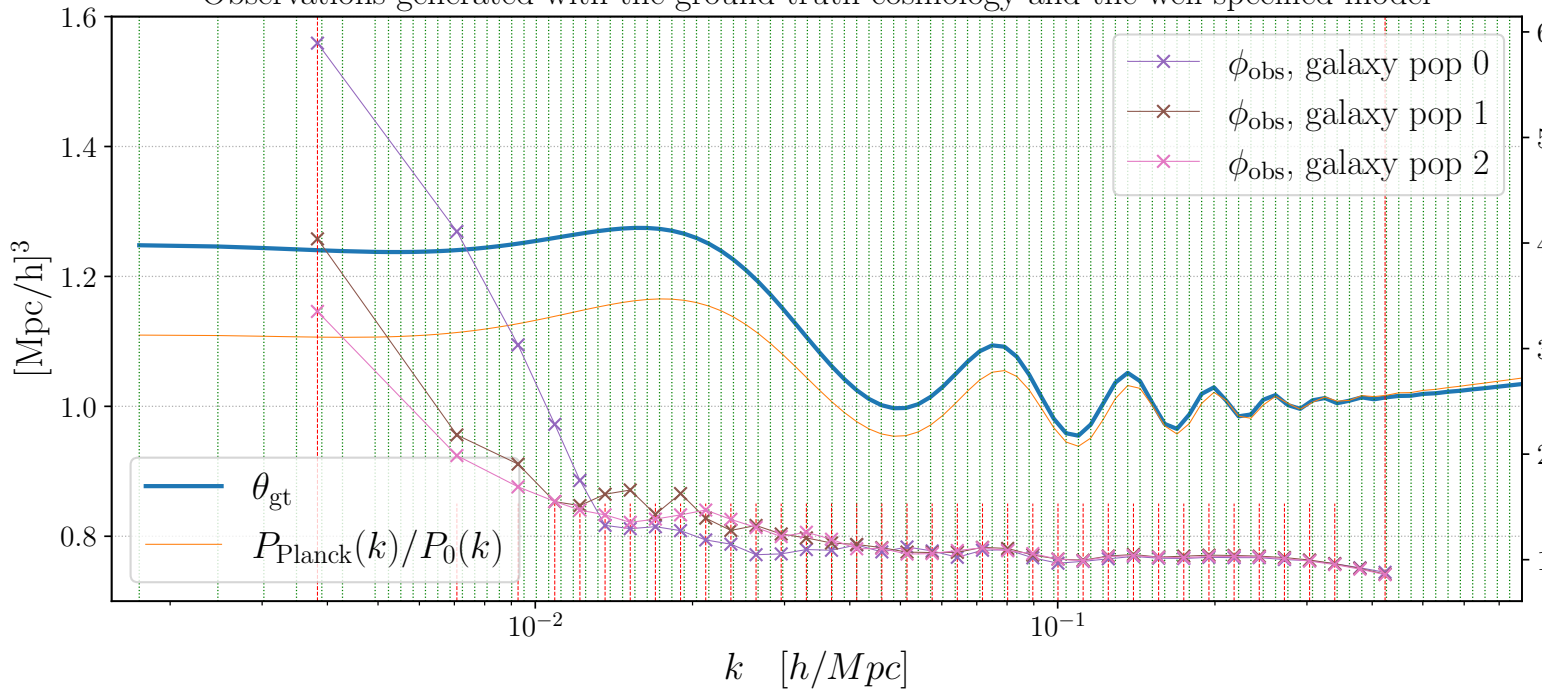
[Salvati, Douspis & Aghanim 2020, 2005.10204](#)

# A general class of Bayesian Hierarchical Model of interest in cosmology

Prior on target parameters  
e.g. Planck 2018

Cheap deterministic simulator  
(Boltzmann solver: CLASS, CAMB)

Observations generated with the ground truth cosmology and the well specified model



$\mathcal{P}(\omega)$

$\omega$

$\mathcal{T}$

$\theta$

$\mathcal{P}(\Phi|\theta)$

$\Phi$

$\tilde{\mathcal{C}}$

$\tilde{\omega}$

← Target parameters  
 $\omega = (h, \Omega_b, \Omega_m, n_S, \sigma_8)$

← Latent function  
(initial matter power spectrum)

← Summary statistic  
from galaxy overdensity field  $\delta^f$

← Compressed data



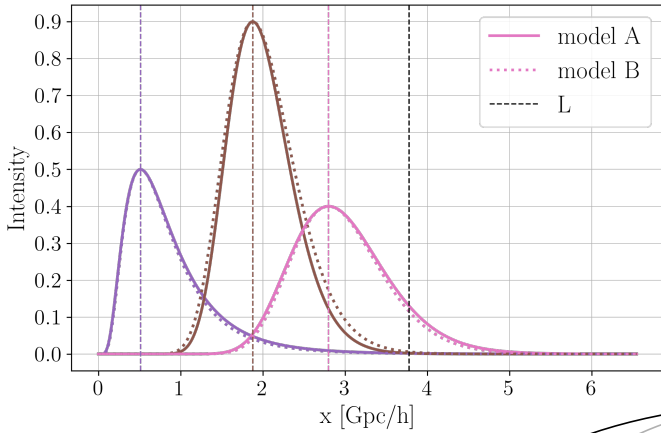
# A general class of Bayesian Hierarchical Model of interest in cosmology

Complex probabilistic observational process

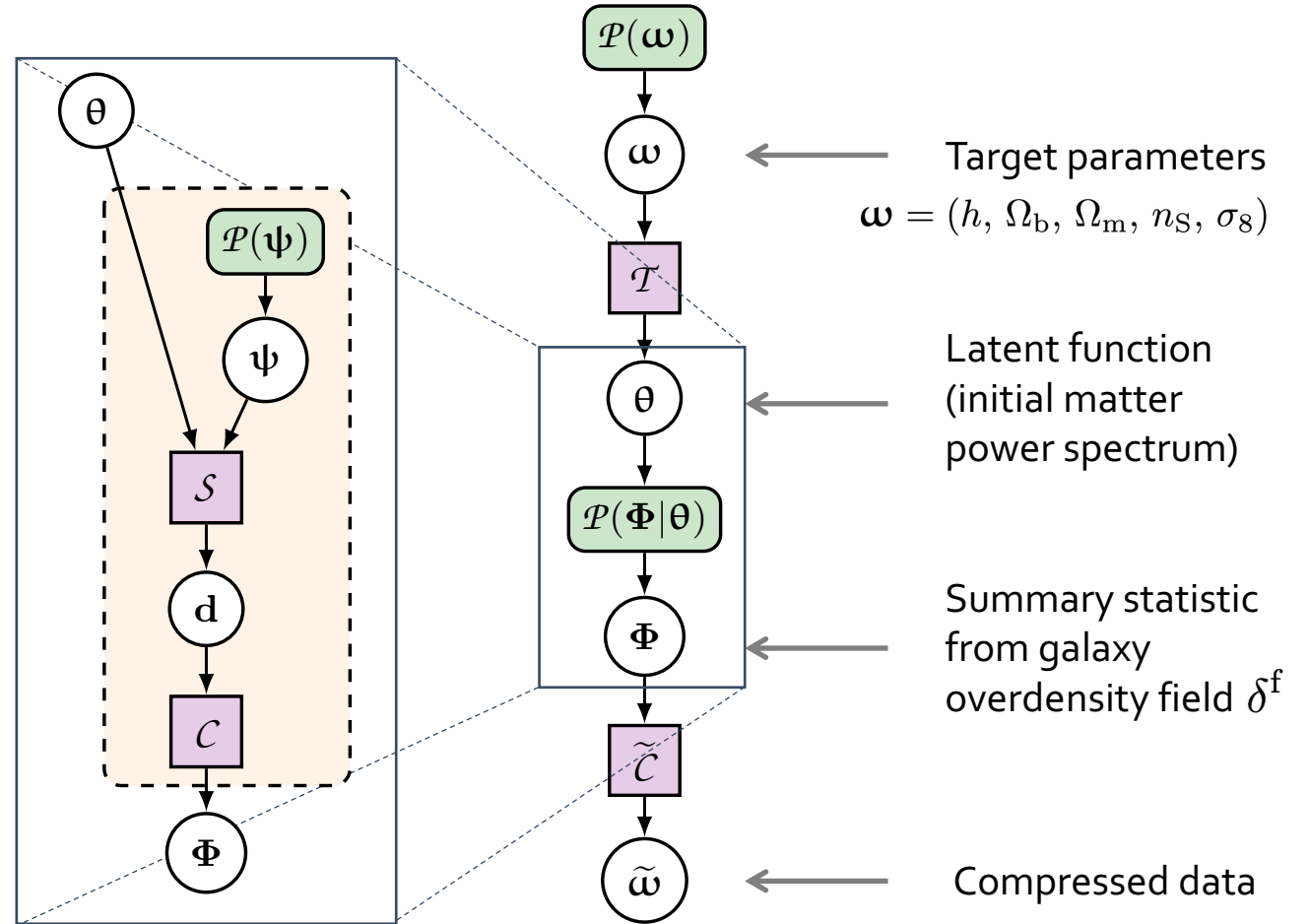
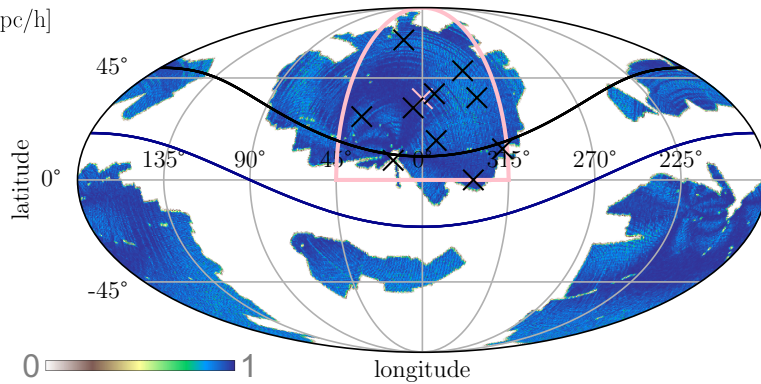
2-LPT, N-body, galaxy formation

Nuisance parameters  $\psi$  (galaxy biases, ...)

Selection functions



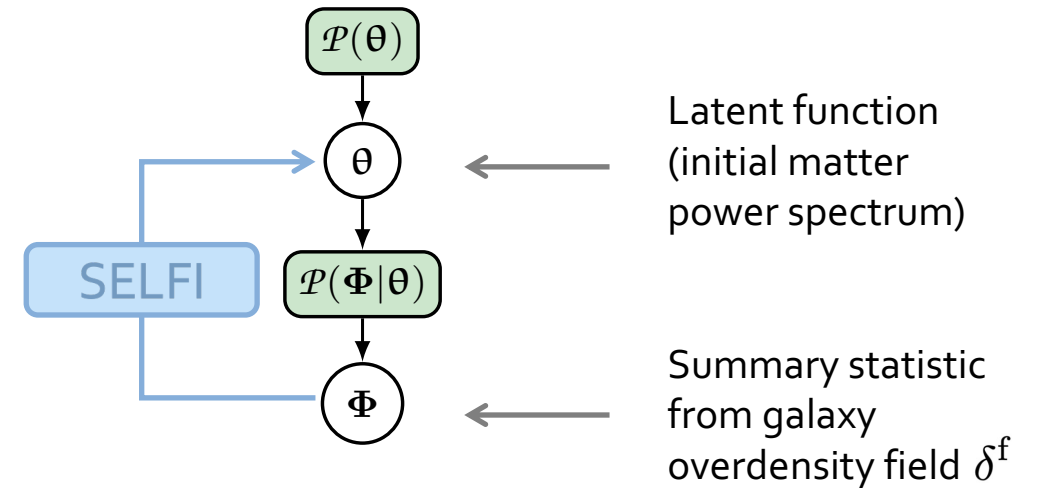
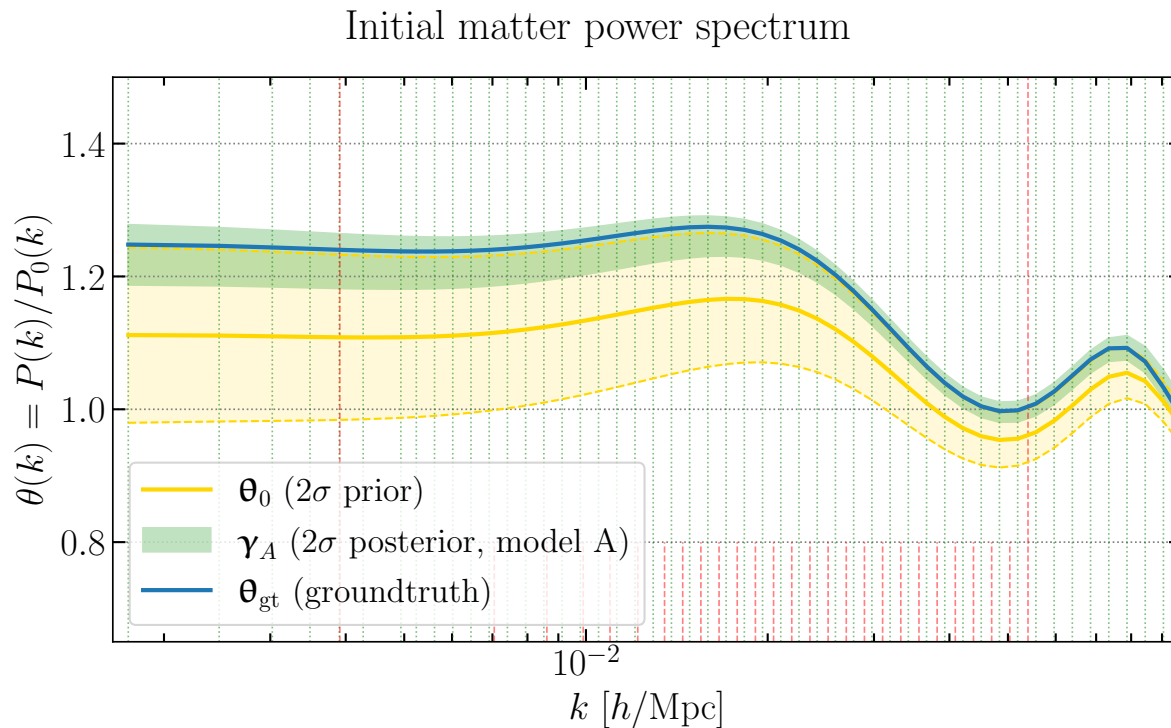
Instrumental effects :  
radial selection, mask...



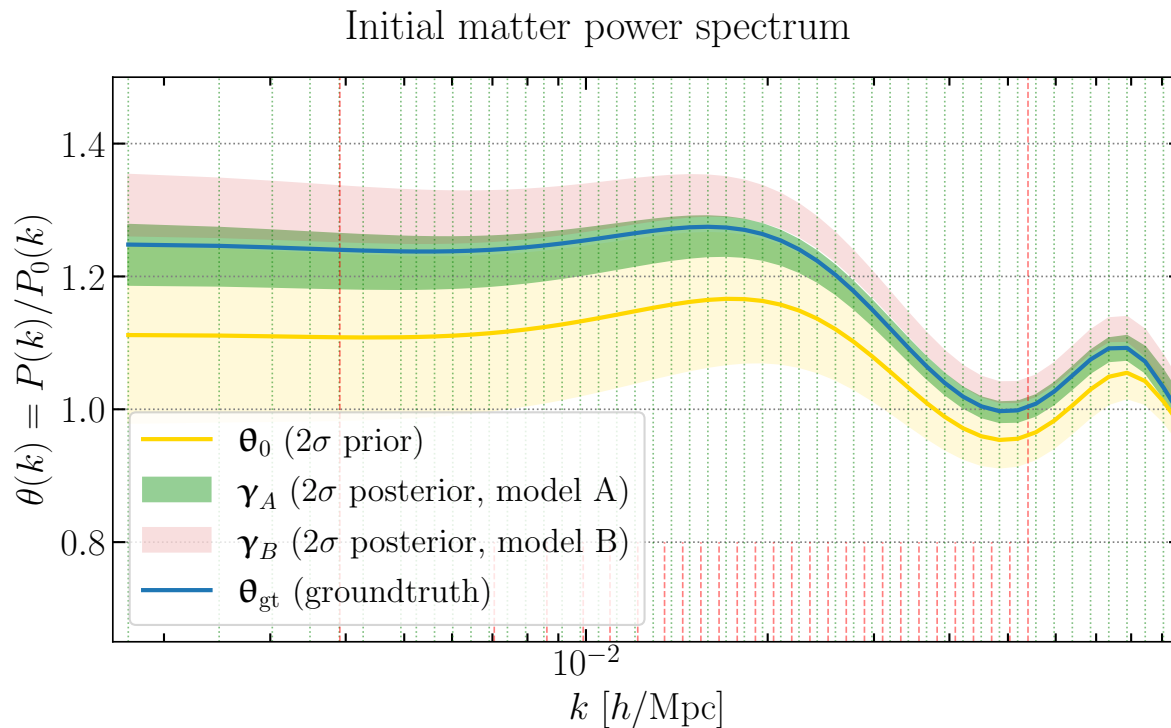


1. Infer a latent function, in this work the initial matter power spectrum  $\theta$

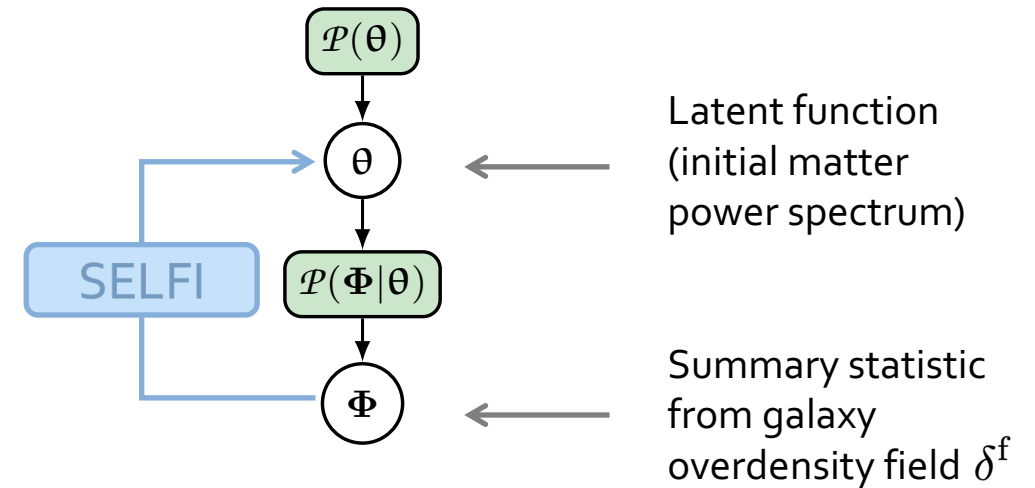
$\theta$ : initial matter powerspectrum normalized by BBKS spectrum



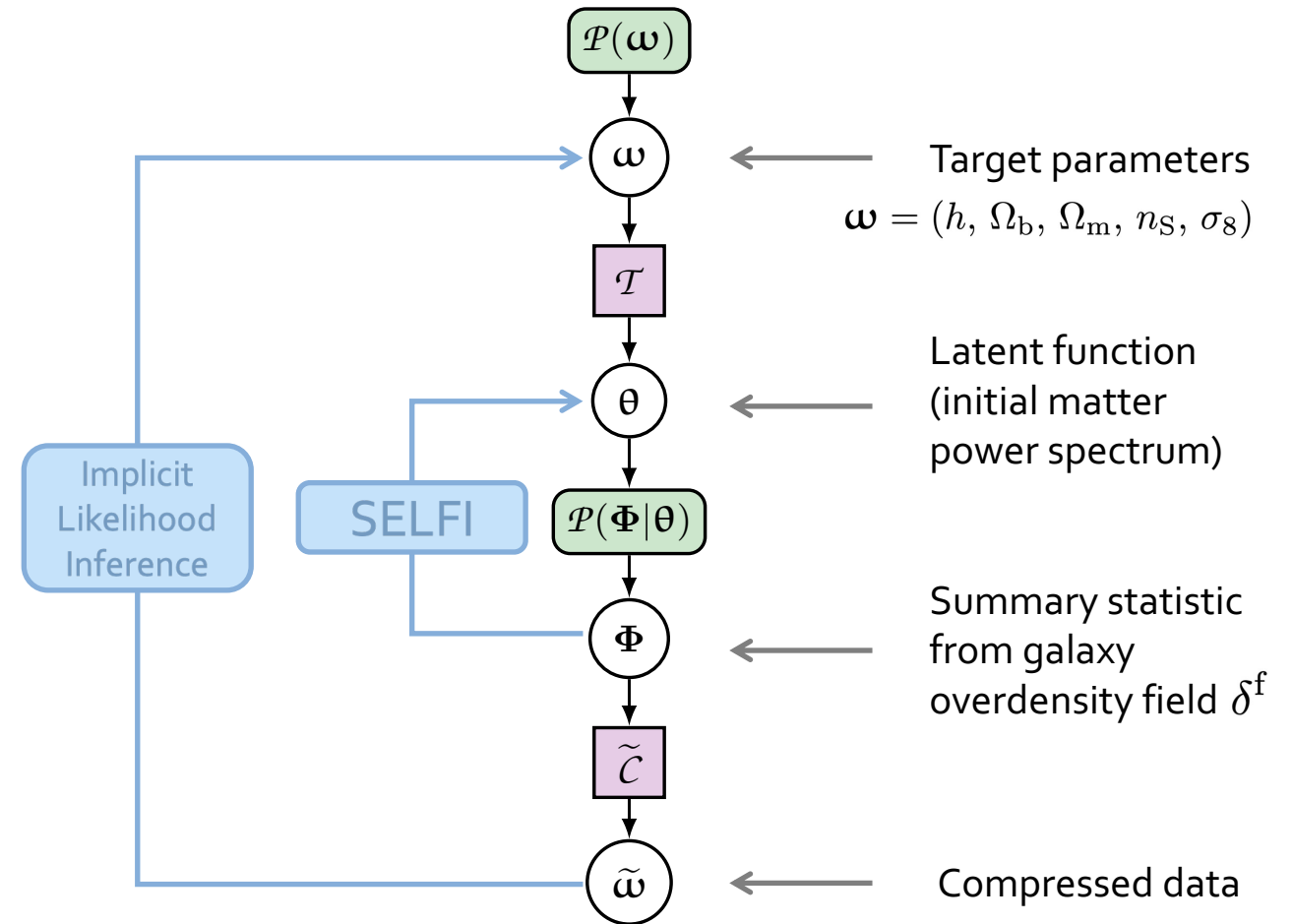
1. Infer a latent function, in this work the initial matter power spectrum  $\theta$
2. Utilize the posterior on  $\theta$  to check for model misspecification

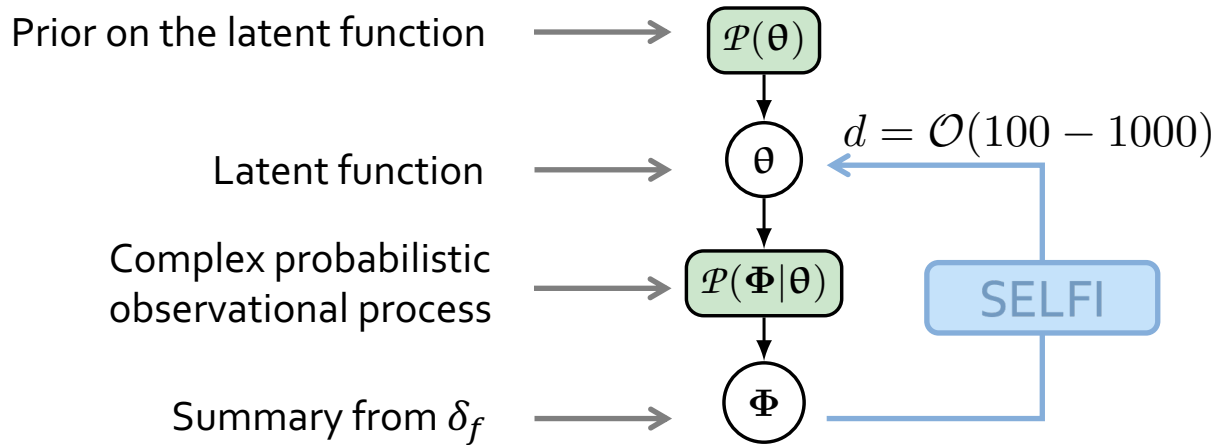


$\theta$ : initial matter powerspectrum normalized by BBKS spectrum



1. Infer a latent function, in this work the initial matter power spectrum  $\theta$
2. Utilize the posterior on  $\theta$  to check for model misspecification
3. Infer the top-level cosmology  $\omega$ 
  - Recycle simulations from step 1. for optimal data compression
  - Use implicit likelihood inference (ABC, Bayesian optimization: BOLFI, other methods)





## Effective posterior:

expansion point  $\theta_0$       observed summaries  $\Phi_O$

Mean:  $\gamma \equiv \theta_0 + \mathbf{\Gamma} (\nabla \mathbf{f}_0)^\top \mathbf{C}_0^{-1} (\Phi_O - \mathbf{f}_0)$

Covariance:  $\mathbf{\Gamma} \equiv [(\nabla \mathbf{f}_0)^\top \mathbf{C}_0^{-1} \nabla \mathbf{f}_0 + \mathbf{S}^{-1}]^{-1}$

covariance of summaries  $\mathbf{C}_0^{-1}$       gradient of the black-box  $\nabla \mathbf{f}_0$       prior covariance  $\mathbf{S}^{-1}$

Assumptions

- Linearization of the black-box data model around an expansion point  $\theta_0$

$$\hat{\Phi}_\theta \approx \mathbf{f}_0 + \nabla \mathbf{f}_0 \cdot (\theta - \theta_0) \equiv \mathbf{f}(\theta)$$

- For step 1. & data compression only, assume:
  - Gaussian prior
  - Gaussian effective likelihood

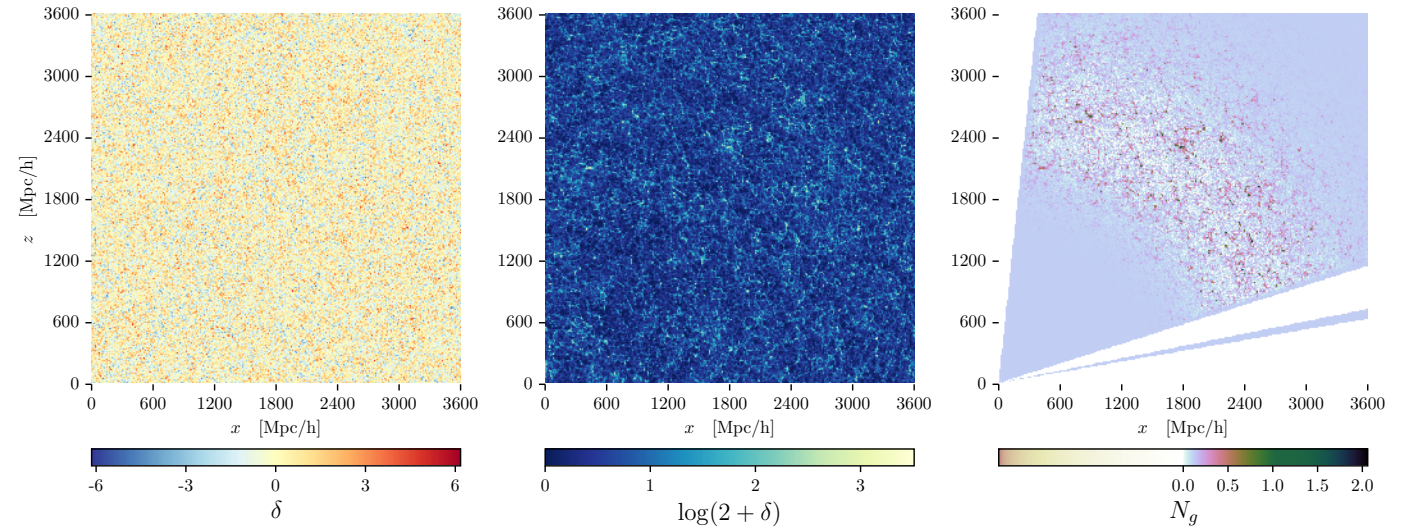
- $\mathbf{f}_0$ ,  $\mathbf{C}_0$  and  $\nabla \mathbf{f}_0$  evaluated through simulations
- The number of simulations is fixed *a priori* (contrary to MCMC)



# Investigating the impact of systematics with SELFI

- Grid specifications:
  - $\theta$  defined on  $S = 100$  support wavenumbers
  - $512^3$  grid, comoving  $L = 3.6$  Gpc/h
- Gravitational evolution with Simbelmynë:
  - [Leclercq, Jasche & Wandelt 2015, 1502.02690](#)
  - Flat  $\Lambda$ -CDM, initial  $\delta^i$  with CLASS
  - $512^3$  DM particles, 2LPT up to  $z = 19$
  - PM grid of  $1024^3$  voxels, COLA to  $z = 0$

From initial matter overdensity field to observed galaxy counts



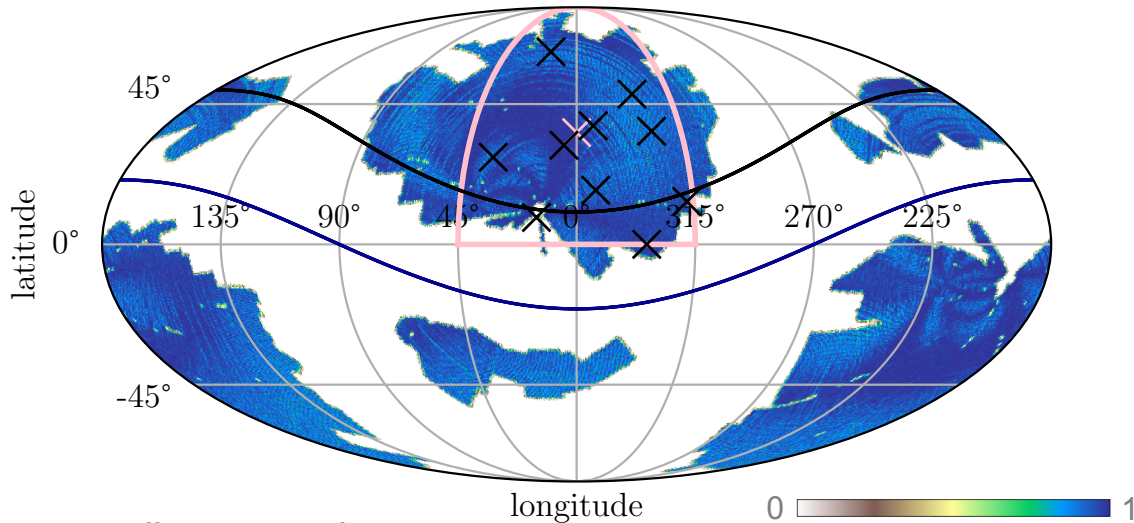
- The observer is at the corner of a cubic box covering 1 octant of the sky, with Euclid-like mask

**Model A**

10 additional masked areas, extinction near galactic plane

**Model B**

no such effects, lower resolution



Hoellinger & Leclercq, in prep.

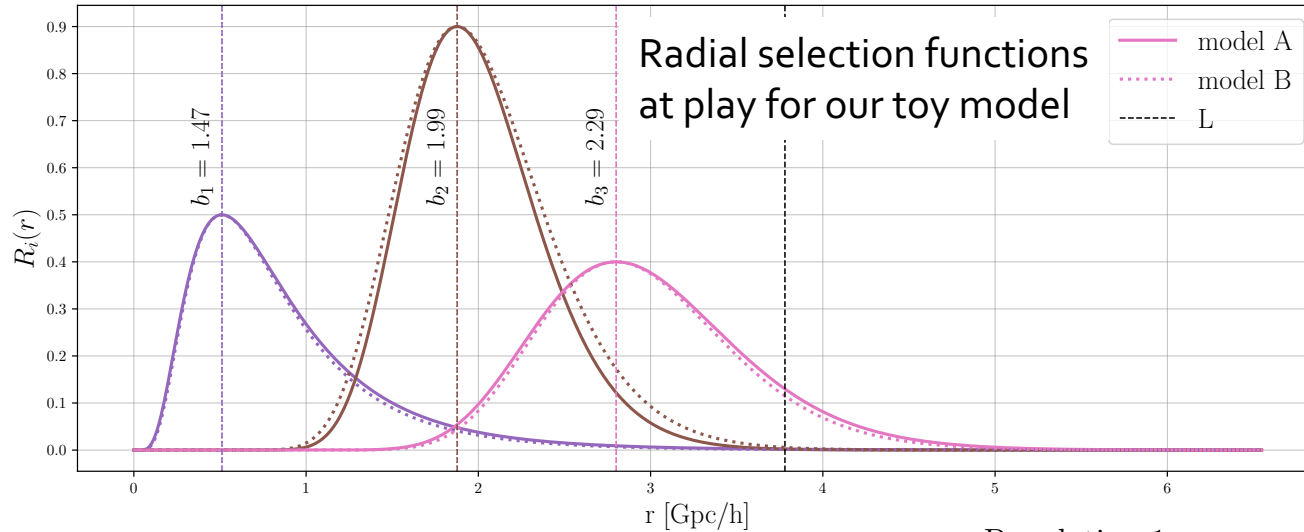
Tristan Hoellinger

Implicit Likelihood Inference in Cosmology while efficiently checking for systematics

2023/11/05



# Investigating the impact of systematics with SELFI



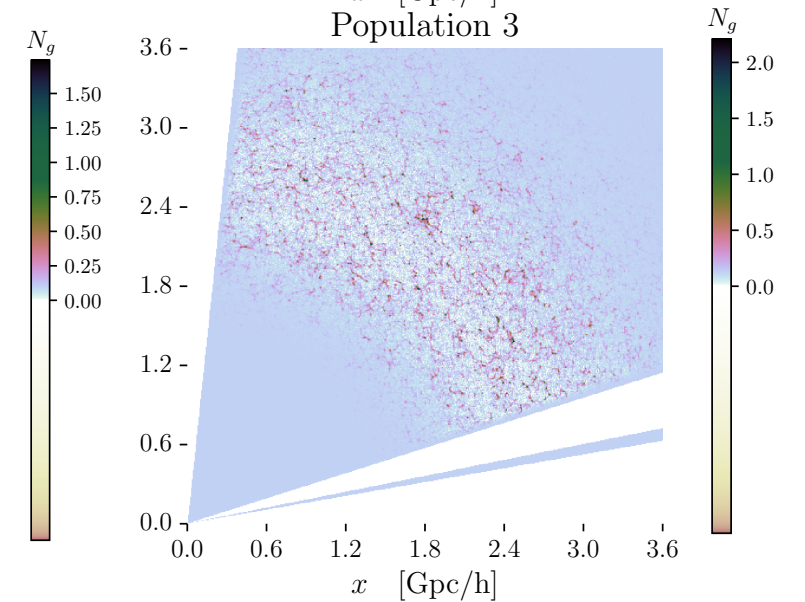
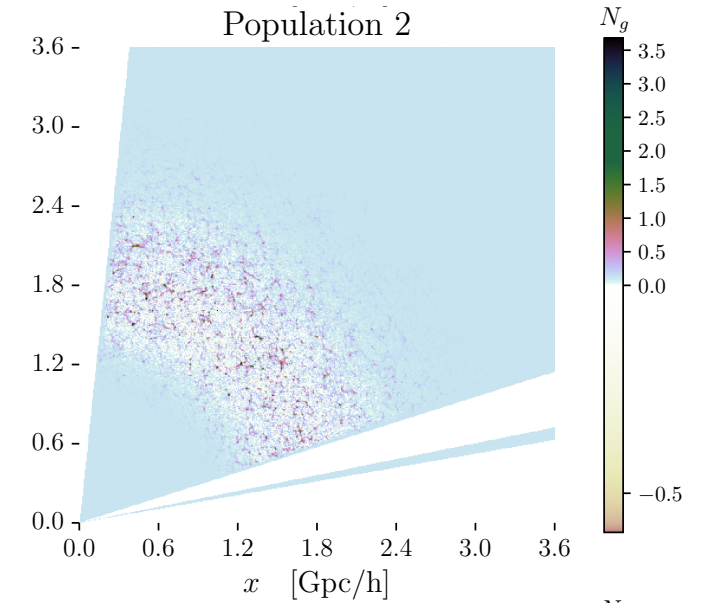
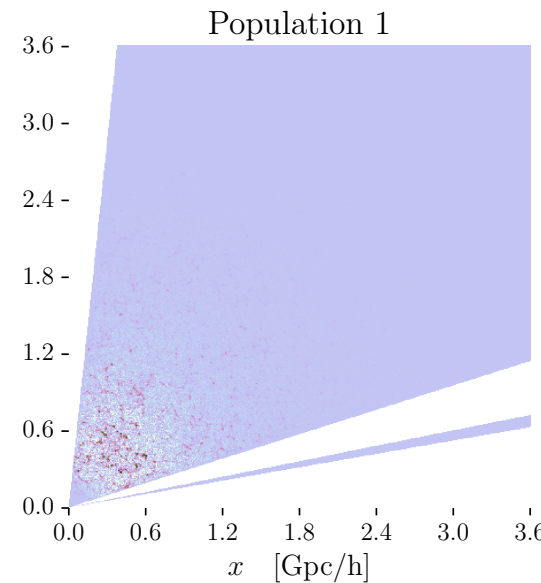
## Model A

- Lognormal selection functions
- Luminosity-dependent galaxy biases

## Model B

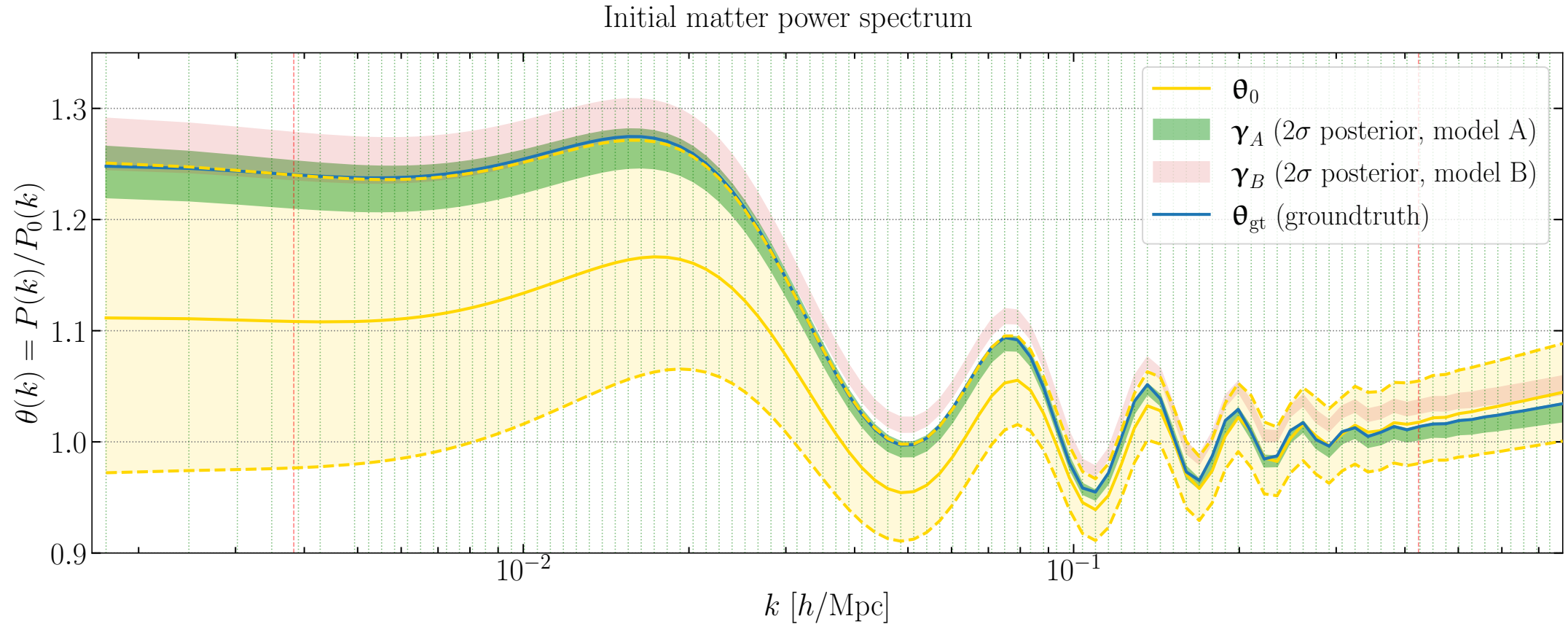
- Misspecified selections and biases
- Effects of order  $\mathcal{O}(1\%)$

$N_e + S \times N_s = 10,200$  simulations  
with the full forward data model





# Investigating the impact of systematics with SELFI



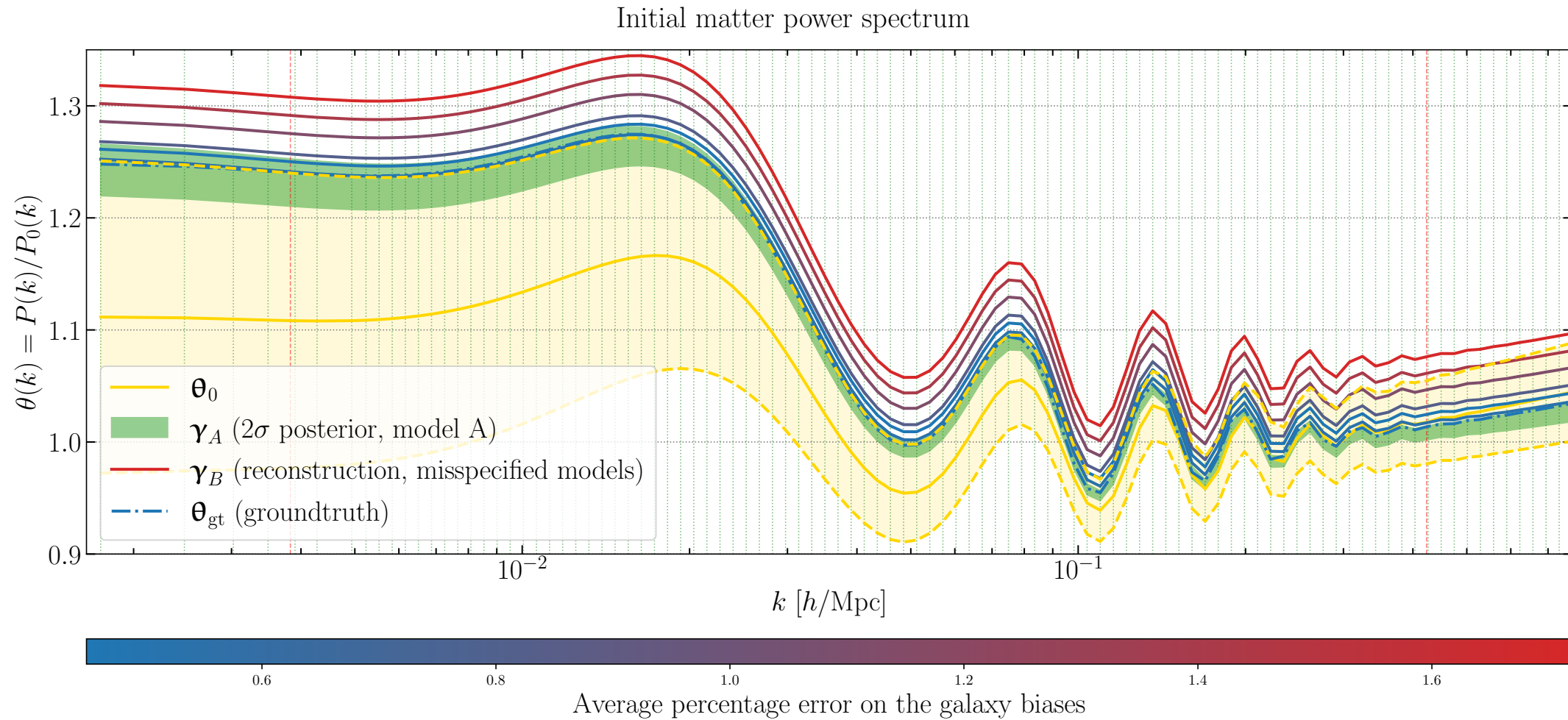
Mahalanobis distances to prior:

**Model A**  $\rightarrow$  1.96

**Model B**  $\rightarrow$  2.91



# Investigating the impact of systematics with SELFI

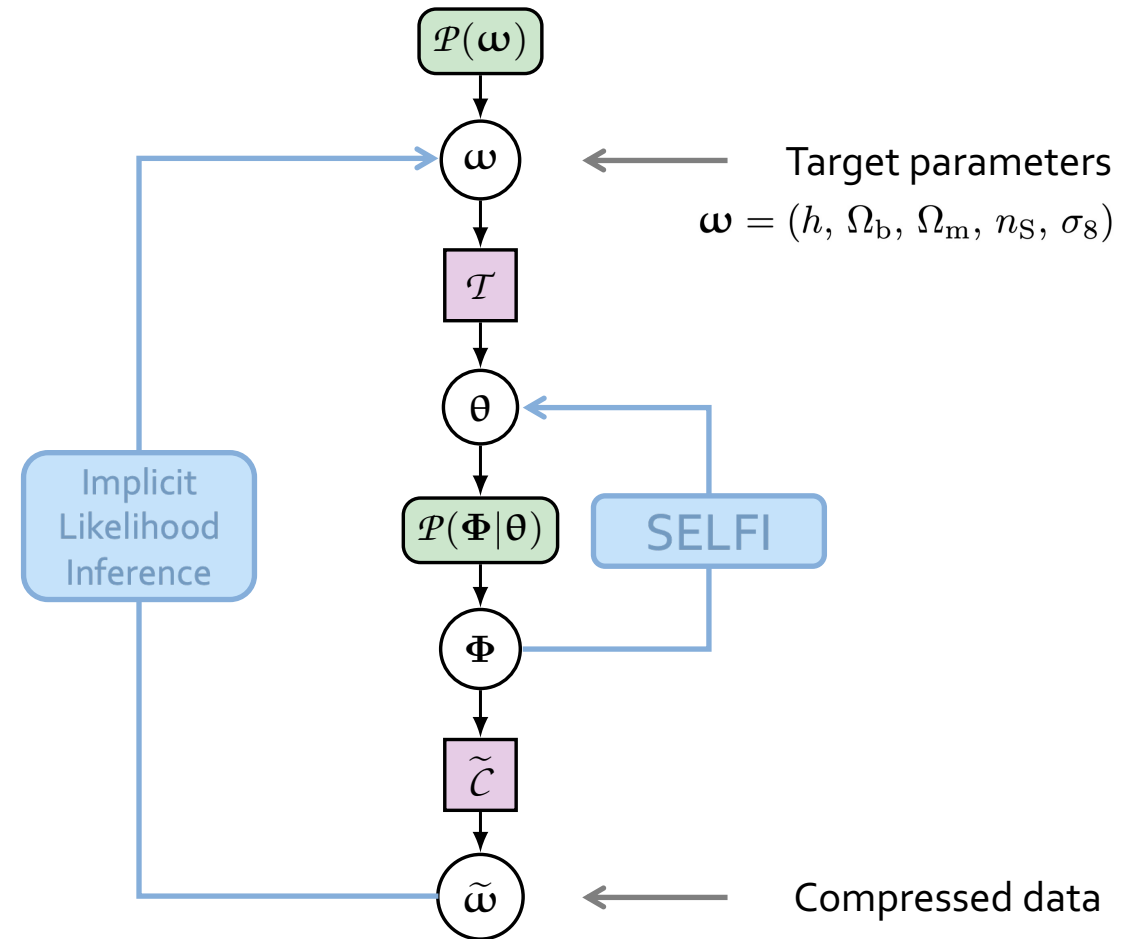


# Optimal data compression

## 3. Infer top-level cosmology $\omega$

We rely on **score compression** to compress the summaries from  $\dim(\Phi) = 111$  to  $\dim(\tilde{\omega}) = \dim(\omega) = 5$

- Score function  $\nabla_{\omega} \hat{\ell}_{\omega_0} \rightarrow$  steepness of  $\hat{\ell}$
- It is a **sufficient statistic** for  $\omega$  for the the linearized log-likelihood, hence it is a natural way to compress the data



# Optimal data compression

$$\mathcal{C}(\Phi) = \tilde{\omega} \equiv \omega_0 + \mathbf{F}_0^{-1} [(\nabla_{\omega} \mathbf{f}_0)^{\top} \mathbf{C}_0^{-1} (\Phi - \mathbf{f}_0)]$$

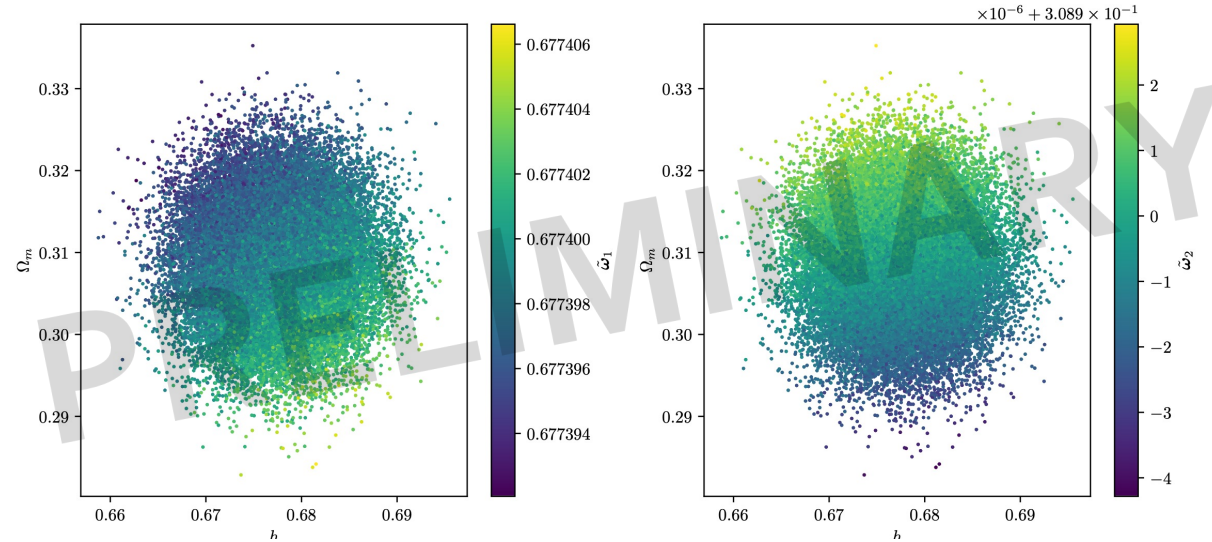
Fisher matrix:  $\mathbf{F}_0 = (\nabla_{\omega} \mathbf{f}_0)^{\top} \mathbf{C}_0^{-1} \nabla_{\omega} \mathbf{f}_0$

$$\nabla_{\omega} \mathbf{f}_0 = \nabla \mathbf{f}_0 \cdot \nabla_{\omega} \mathcal{T}_0$$

Already computed for SELFI      Cheap via finite differences

- The compression is optimal in the sense that it preserves the Fisher content of the data. Hypothesis:
  - the covariance matrix does not vary close to the expansion point  $\nabla_{\omega} \mathbf{C} = 0$
  - The likelihood is gaussian or the following holds:  $\nabla \mathbb{E}_{\theta} [\nabla^T \mathcal{L}] = \nabla \mathbb{E}_{\theta} [\nabla \nabla^T \mathcal{L}]$

- Example for  $\Omega_m$  and  $h$  (simplified data model):



Leclercq 2022, 2209.11057

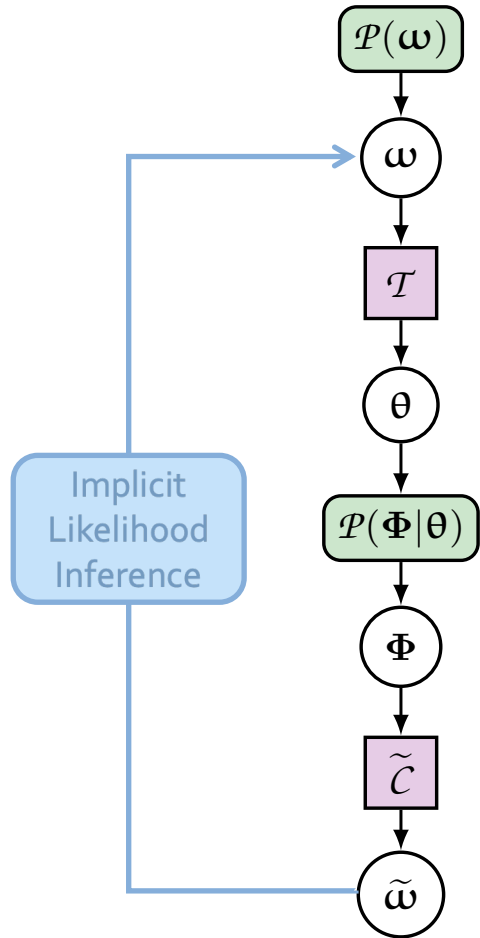
Alsing & Wandelt 2018, 1712.00012

Hoellinger & Leclercq, in prep

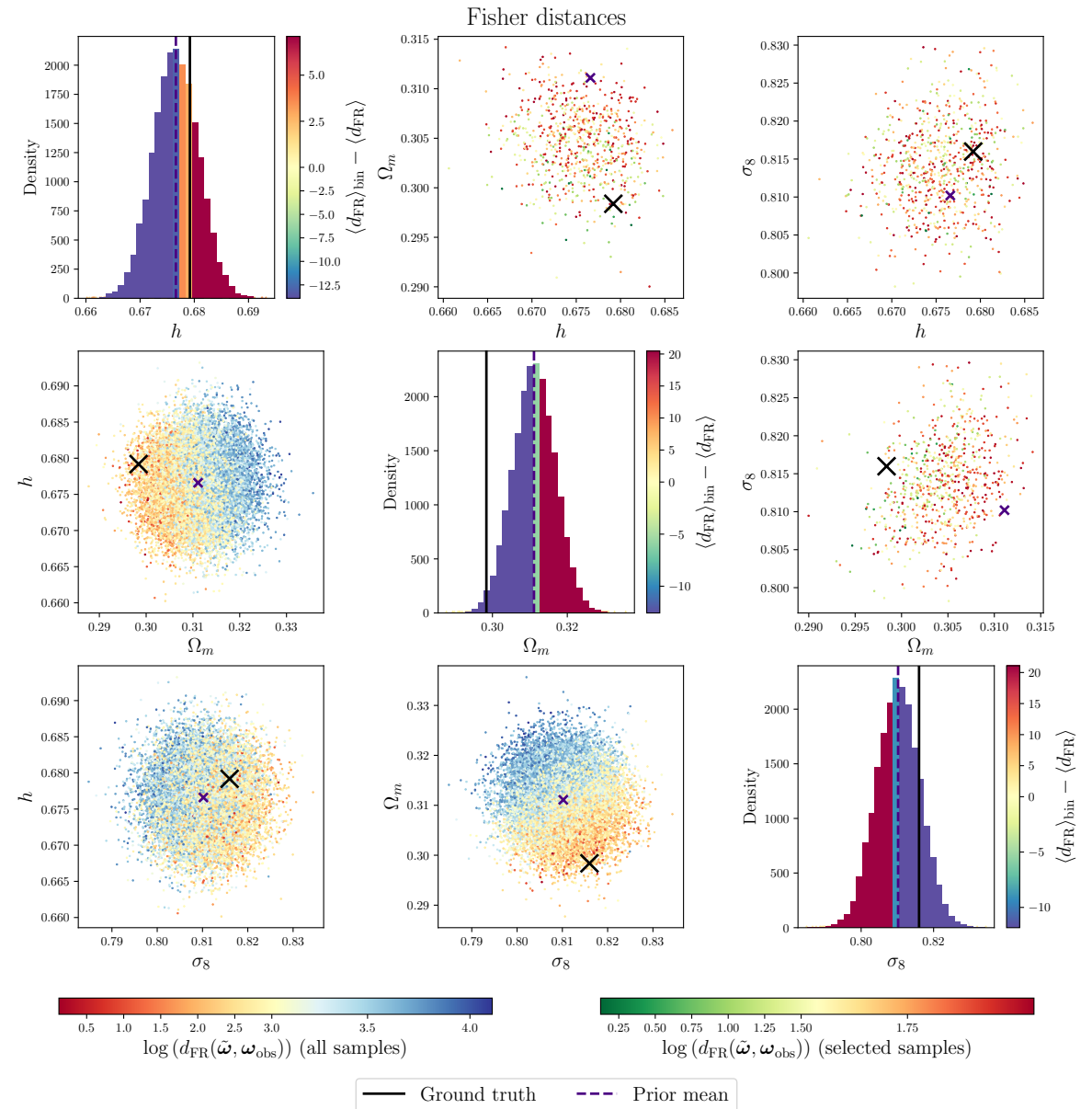


# Optimal data compression

## 3. Infer top-level cosmology $\omega$



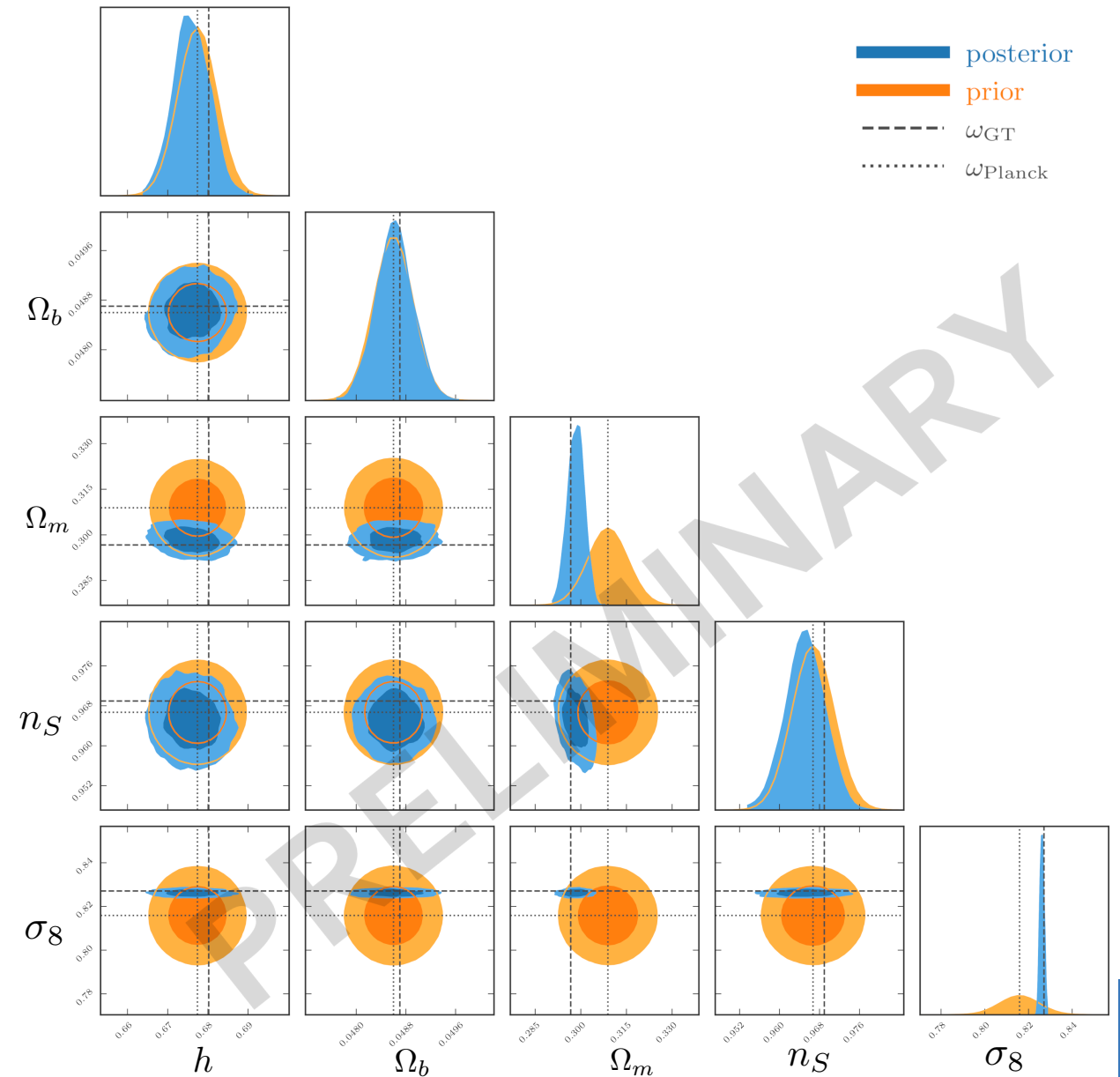
- Assumptions made for steps 1. - 2. do not impact final inference:
    - the Gaussian effective likelihood is not required to infer the cosmology except through data compression
    - lossy data compression usually leaves posteriors unbiased
  - Any algorithm can be used to obtain the posterior  $\mathcal{P}(\omega|\tilde{\omega}_O)$ 
    - Non-parametric approaches can use the Fisher-Rao distance
- $$d_{\text{FR}}(\tilde{\omega}, \tilde{\omega}_O) \equiv \sqrt{(\tilde{\omega} - \tilde{\omega}_O)^\top \mathbf{F}_0(\tilde{\omega} - \tilde{\omega}_O)}$$
- Example with ABC (same data model but smaller dim. for  $\omega$  and  $64^3$  grid)



## Final posterior

- We ran the full SELFI pipeline with a simplified forward data model
  - same instrumental response as before,
  - no gravitational evolution
  - baseline ABC for step 3.
- We got unbiased posteriors for the top-level cosmological parameters
$$\omega = (h, \Omega_b, \Omega_m, n_S, \sigma_8)$$

2313 selected samples out of 700K





- A novel **two-step simulation-based Bayesian** approach, combining SELFIE and SBI, to tackle the issue of model misspecification for a large class of BHM.
- Advantages related to the first step (SELFIE):
  - No need to incorporate any knowledge of the data-generating process in the analysis, even with high dimensional complex black-box simulators.
  - Number of simulations fixed a priori.
  - The computational workload is perfectly parallel.
- Advantages related to the second step (SBI):
  - The **score compressor comes for free**: we recycle simulations from step **1**.
  - General advantages of SBI with respect to likelihood-based methods are preserved.
  - No simplification nor inner knowledge of the forward data model required

- Incorporate parametrization of deviations from  $\Lambda$ -CDM in the top-level parameters:
  - equation of state of dark energy  $w(a)$
  - total neutrino masses  $m_\nu$
  - deviations from gaussianity  $f_{NL}$
- Jointly infer all top-level parameters and nuisance parameters such as galaxy biases
  - This means implicit likelihood inference in dimension  $\mathcal{O}(10 - 20)$
  - We need to investigate and extend advanced Bayesian optimisation strategies to explore the parameter space to avoid curse of dimensionality

- Thanks for listening!

- Main references:

[Alsing & Wandelt 2018, 1712.00012](#)

[Leclercq et al. 2019, 1902.10149](#)

[Leclercq 2022, 2209.11057](#)

[Leclercq, Jasche & Wandelt 2015, 1502.02690](#)

- Code and data availability: wait for SELF<sub>l</sub>2 release in 2024!



[hoellin.github.io](https://hoellin.github.io)

[www.aquila-consortium.org](https://www.aquila-consortium.org)

