



Likelihood-free large-scale structure inference with robustness to model misspecification



Cosmology & Gravitation seminar,
Stockholm University

Florent Leclercq

www.florent-leclercq.eu

Institut d'Astrophysique de Paris
CNRS & Sorbonne Université

In collaboration with:
Wolfgang Enzi (MPA), Alan Heavens (Imperial College),
Tristan Hoellinger (IAP), Jens Jasche (Stockholm
University)

and the Aquila Consortium

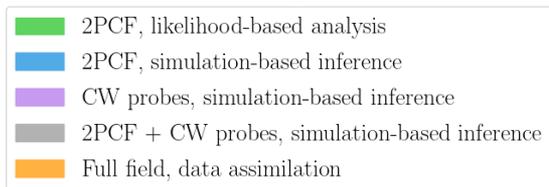
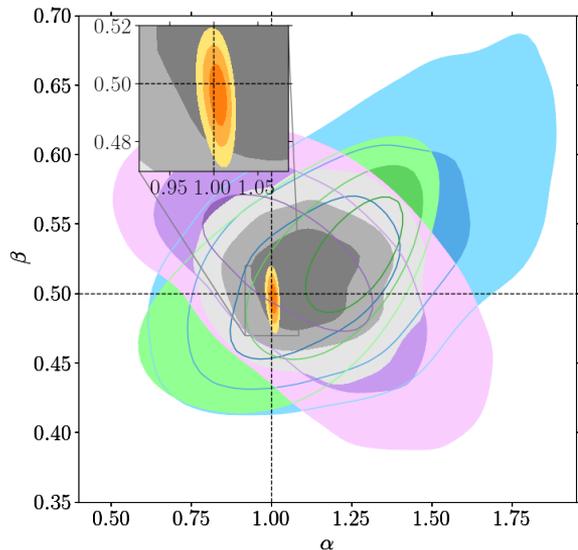
www.aquila-consortium.org

15 June 2023

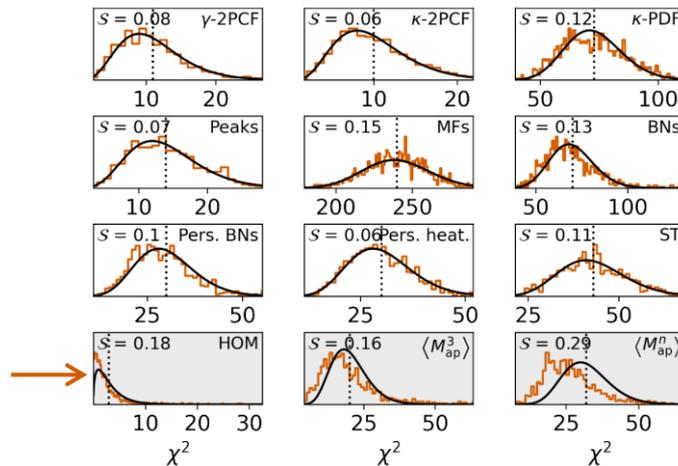
Why I decided to go “likelihood-free” for galaxy clustering additional probes

Note: likelihood-free inference \approx simulation-based inference \approx implicit likelihood inference

- A question of **accuracy**: first, avoid biases.



- Some WL additional probes also have a non-Gaussian distribution.



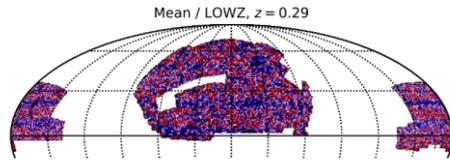
- A question of **precision**: can numerical forward models be used to push further than $k \gtrsim 0.15 h/\text{Mpc}$? The full field contains much more information.



The issue of model misspecification in Bayesian inference and in simulation-based inference (SBI)

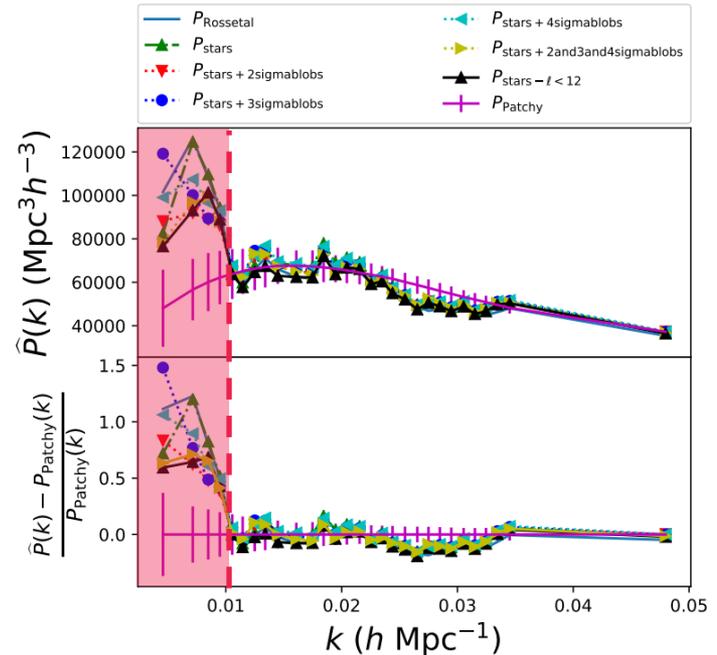
- [Model misspecification](#) arises when model differs from actual data-generating process.
- Field-based inference techniques have a successful track record at handling model misspecification, e.g. automatically reporting unknown data contaminations.

Map of unknown foreground contaminant



- Model misspecification: a major challenge particularly for approaches that marginalise over latent variables, such as [simulation-based inference](#) (SBI).

- Typical cosmological example: the galaxy *power spectrum* at large scales.

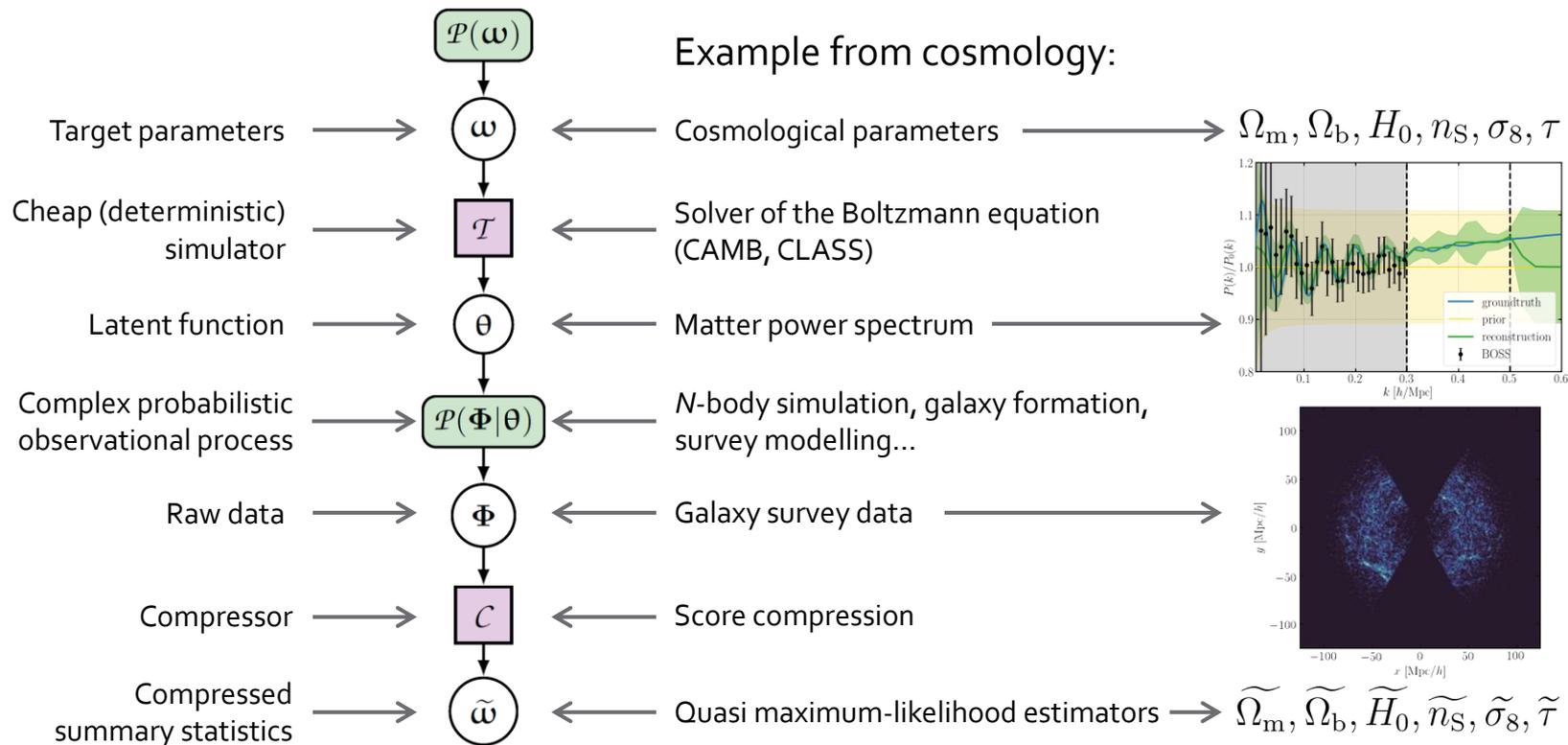


Porqueres, Ramanah, Jasche & Lavaux, 1812.05113
Lavaux, Jasche & FL, 1909.06396

Kalus, Percival et al., 1806.02789

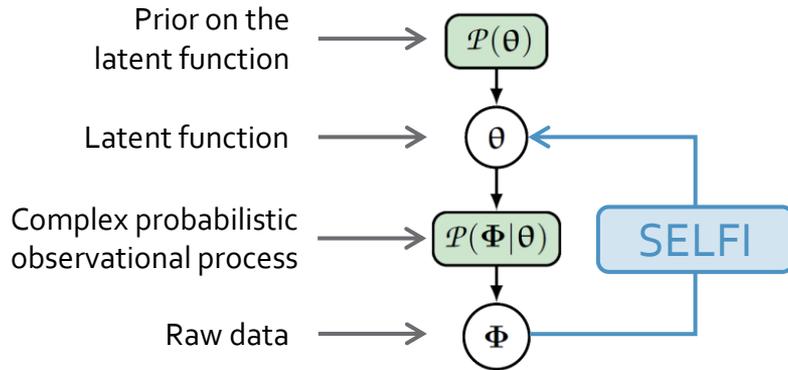


A general class of Bayesian hierarchical models (BHM): Complex observations of a latent function controlled by top-level parameters

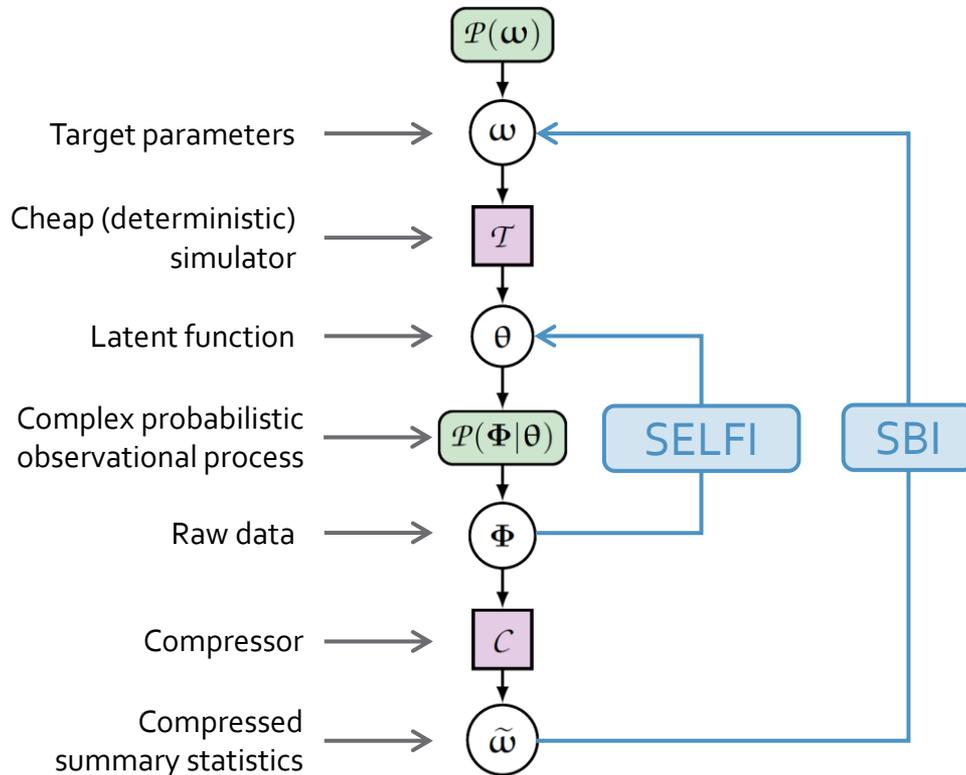


Key idea: a two-step SBI process that recycles simulations

1. Inference of the latent function θ , to check for model misspecification:
 - SELFI algorithm



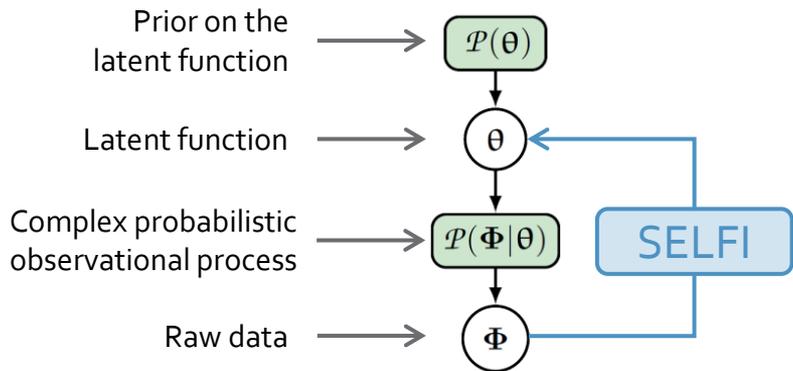
Key idea: a two-step SBI process that recycles simulations



1. Inference of the latent function θ , to check for model misspecification:
 - SELF algorithm
2. Simulation-based inference of ω :
 - Approximate Bayesian Computation (ABC), Likelihood-Free Rejection Sampling
 - Density/ratio estimation (DELFI / NRE)
 - Bayesian optimisation (BOLFI)
 - others...

Important: the simulations necessary for step 1. are recycled for data compression, which is required for step 2.

Latent function inference: the SELFI approach (*Simulator Expansion for Likelihood-Free Inference*)

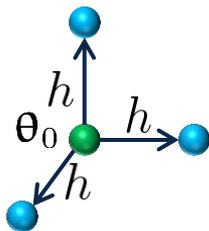


- We aim at inferring the latent function θ , which usually contains most/all of the information on ω .
(initial power spectrum in cosmology, prey/predator population functions in ecology)
- This requires doing SBI in $d = \mathcal{O}(100) - \mathcal{O}(1,000)$
- If we trust the results of earlier experiments, we can Taylor-expand the black-box around an expansion point θ_0 :

$$\hat{\Phi}_\theta \approx \mathbf{f}_0 + \nabla \mathbf{f}_0 \cdot (\theta - \theta_0) + \frac{1}{2} (\theta - \theta_0)^\top \cdot \mathbf{H} \cdot (\theta - \theta_0) + \dots$$

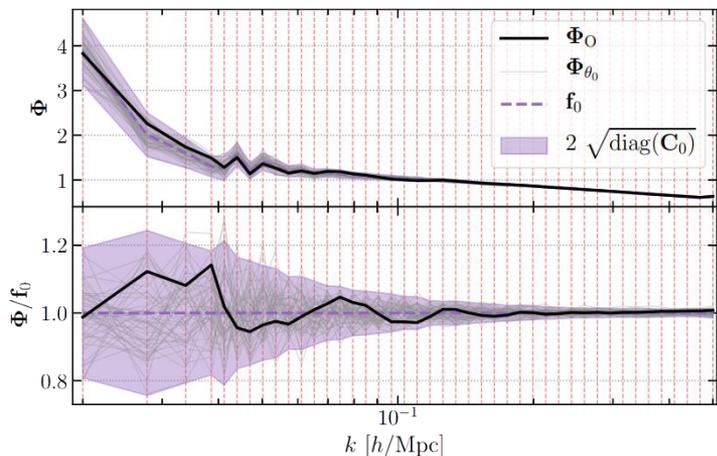
SELFI-2 (second order): coming soon!

- Gradients, Hessian matrix, etc. of the black-box can be evaluated via finite differences in parameter space.

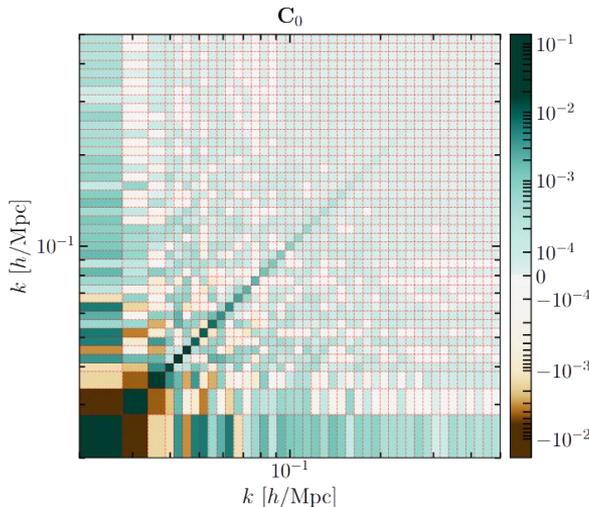


Galaxy clustering additional probes pipeline: diagnostics of the linearised black-box data model

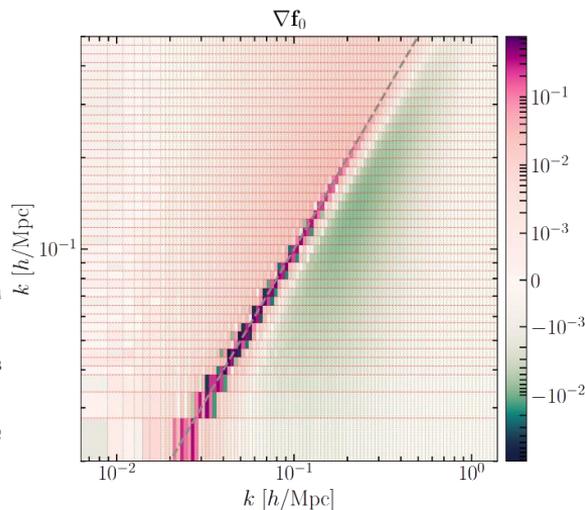
Black-box mean



Black-box covariance



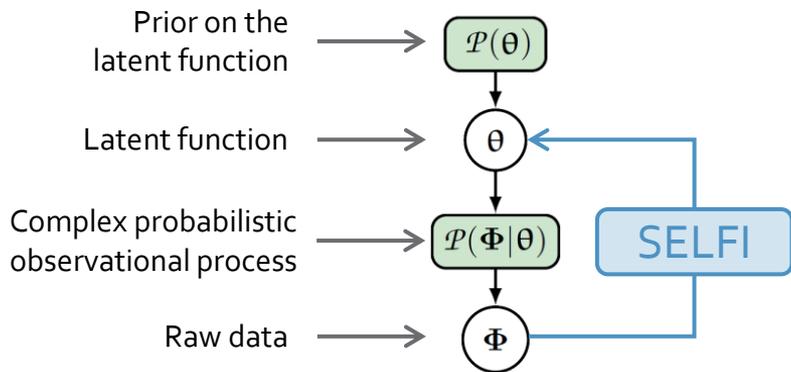
Black-box gradient



- Using only here the (final, non-linearly evolved) power as summary statistics.
- Any additional probe can go in the data vector, since we need the simulations anyway!



Latent function inference: the SELFI approach (*Simulator Expansion for Likelihood-Free Inference*)



- Linearisation of the black-box data model:

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

- Further assume:

- Gaussian prior: $\mathcal{P}(\theta) = \mathcal{G}(\theta_0, \mathbf{S})$
- Gaussian effective likelihood: $\mathcal{P}(\Phi|\theta) = \mathcal{G}[\mathbf{f}(\theta), \mathbf{C}_0]$

- The posterior is Gaussian and analogous to a Wiener filter:

expansion point observed summaries

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

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

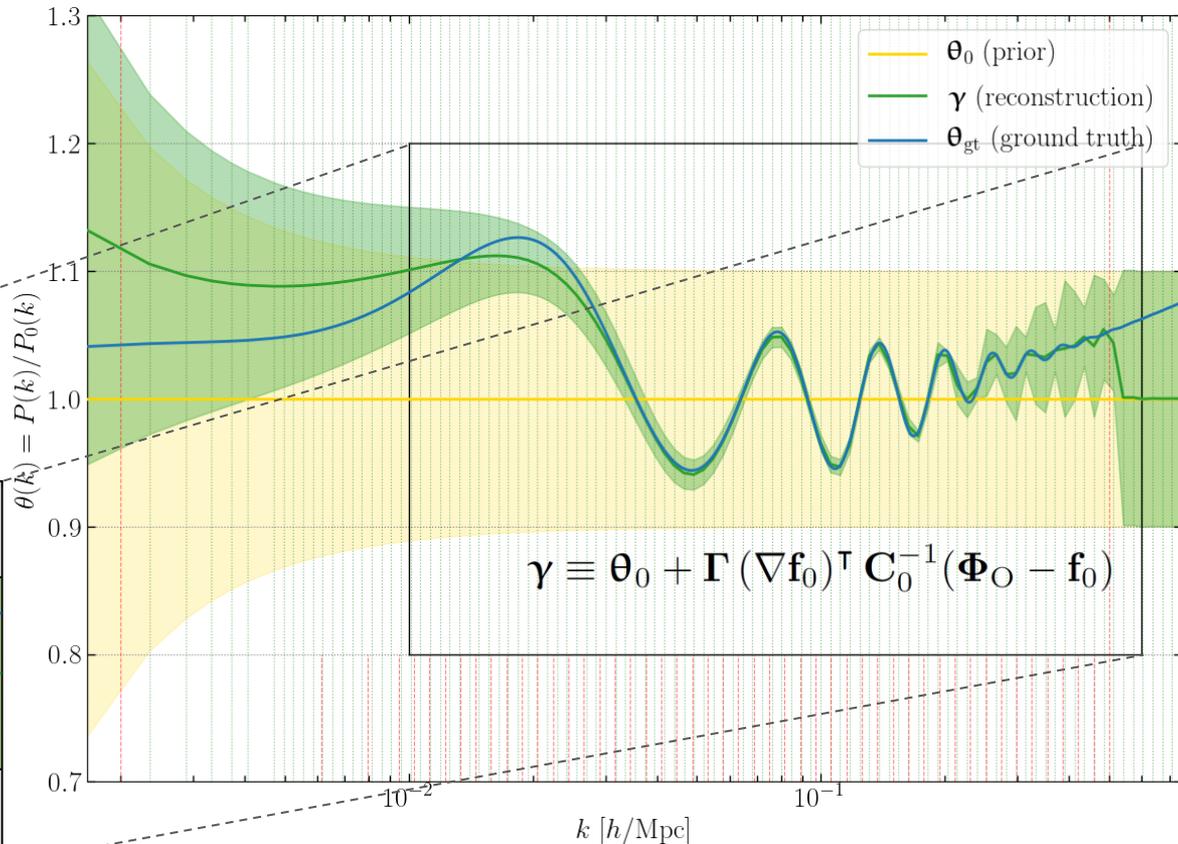
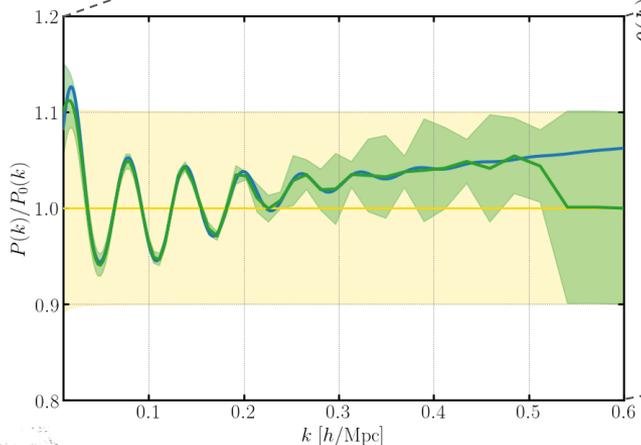
covariance of summaries gradient of the black-box prior covariance

- $\mathbf{f}_0, \mathbf{C}_0$ and $\nabla \mathbf{f}_0$ can be evaluated through simulations only.
- The number of required simulations is fixed *a priori* (contrary to MCMC).
- The workload is perfectly parallel.



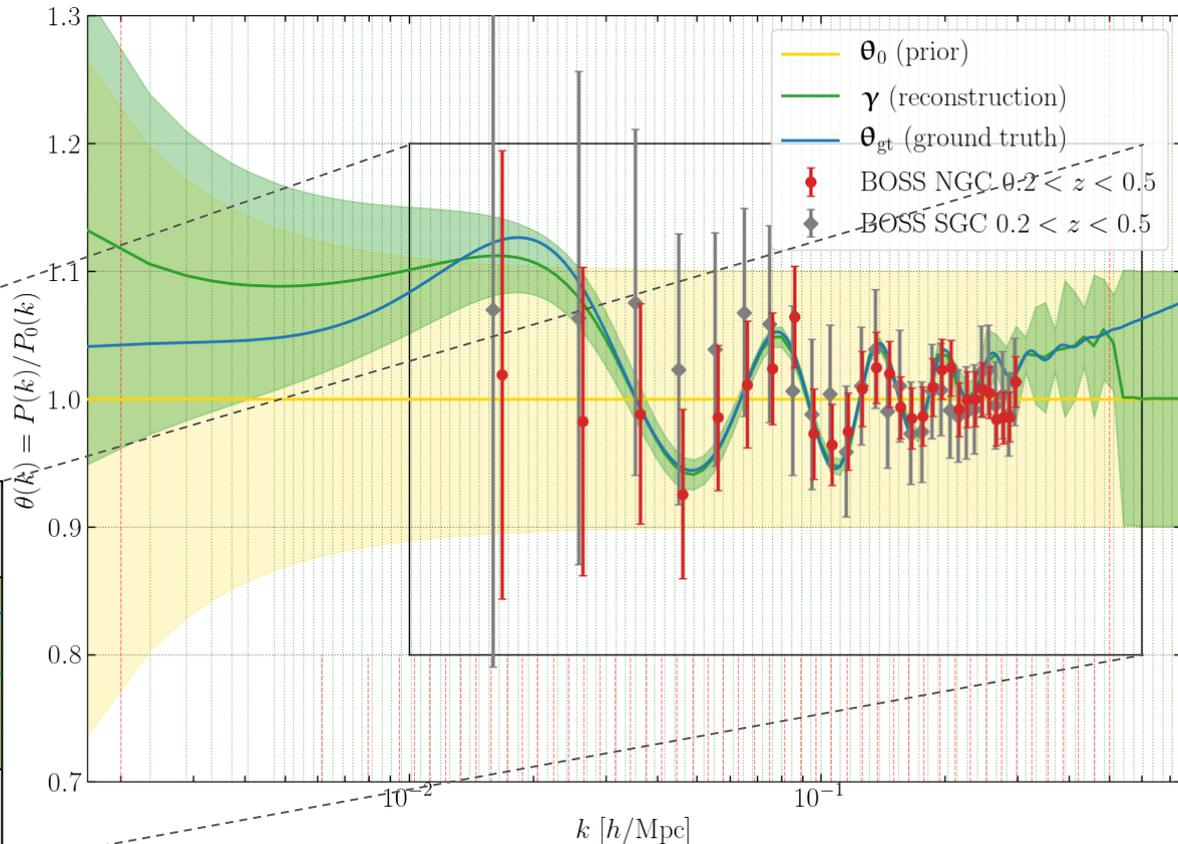
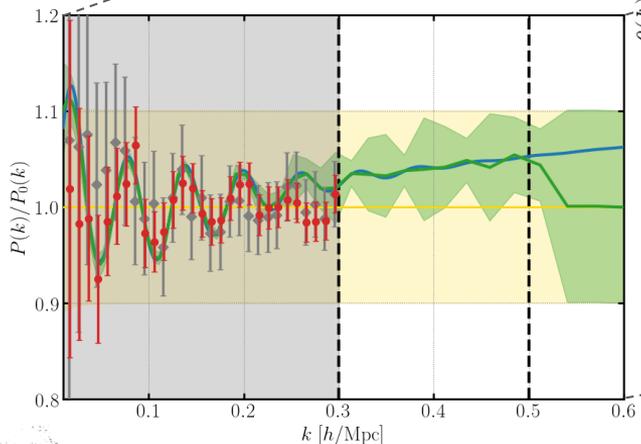
SELI-1 Euclid forecast (cosmic variance limit)

- $V = (3780 \text{ Mpc}/h)^3$
(volume of the Euclid flagship simulation)
- Gaussian random field data model; 6,060 simulations
- 100 parameters are simultaneously inferred



SELI-1 Euclid vs BOSS

- Numerical data models allow using the galaxy power spectrum as summary statistics up to at least $k \gtrsim 0.5 \text{ h/Mpc}$ safely
- $N_{\text{modes}} \propto k^3$: **5 times more modes** are used in the analysis.



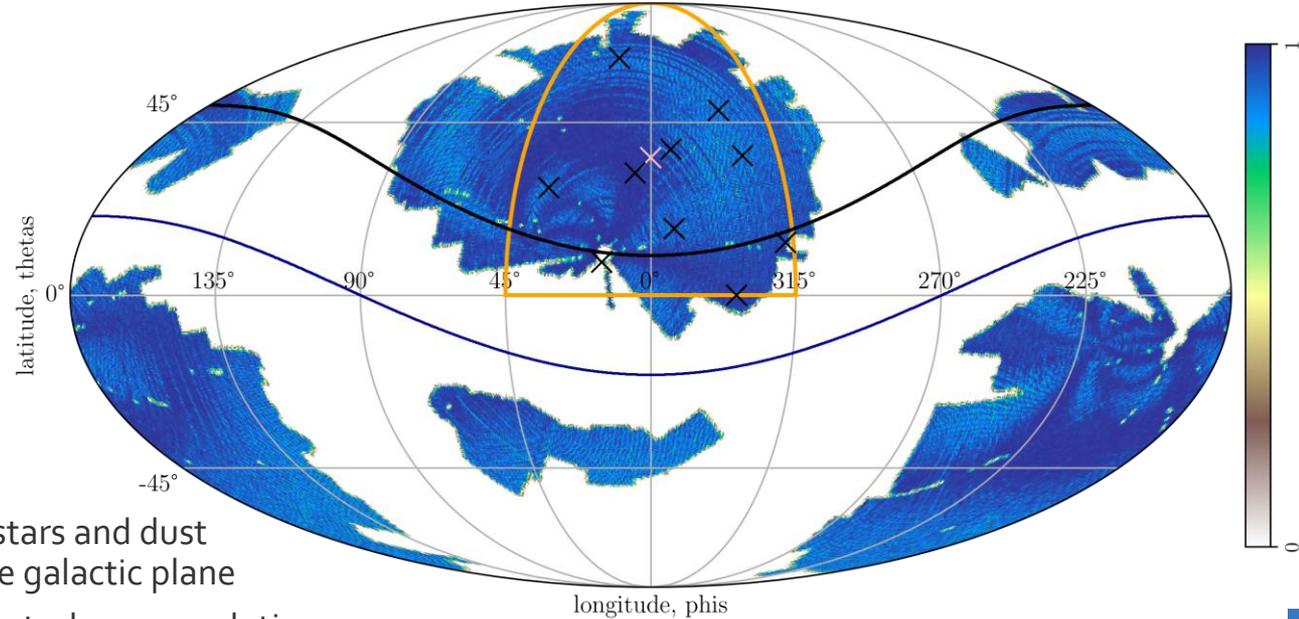
Data points from Beutler et al., 1607.03149



Systematic effects at large scales: mask and selection functions

- $V = (3780 \text{ Mpc}/h)^3$ cubic box, covering one octant of the sky

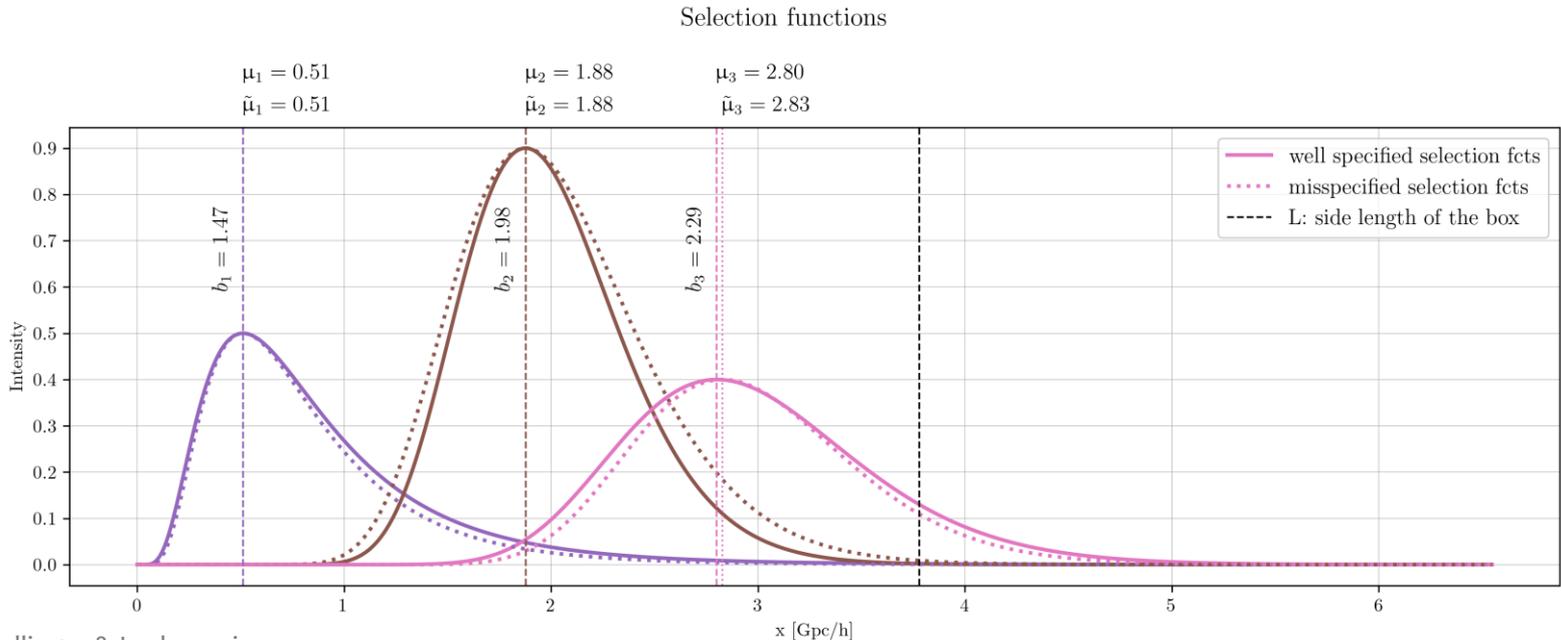
The extracted region of the mask for the observed octant is delimited by the orange triangle



- Two models:
 - Model A: 10 masked stars and dust extinction close to the galactic plane
 - Model B: no such effects, lower resolution

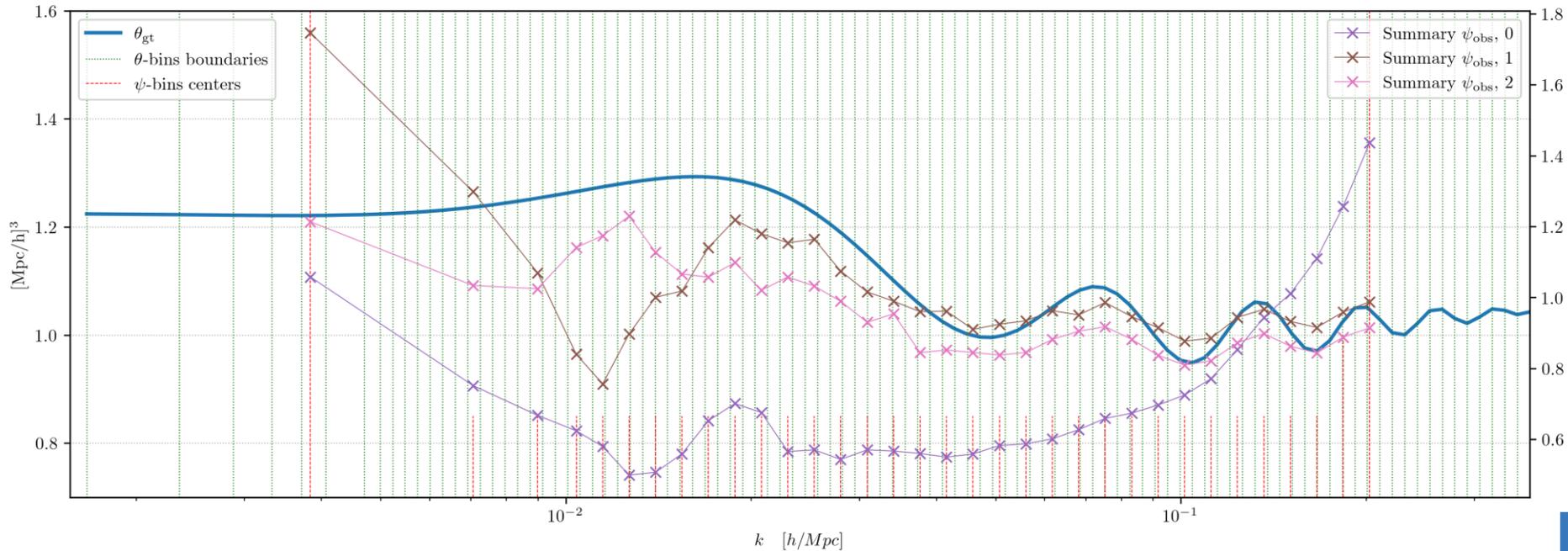
Systematic effects at large scales: mask and selection functions

- Two models:
 - Model A: lognormal selection functions, luminosity-dependent galaxy bias
 - Model B: misspecified selection functions and galaxy biases



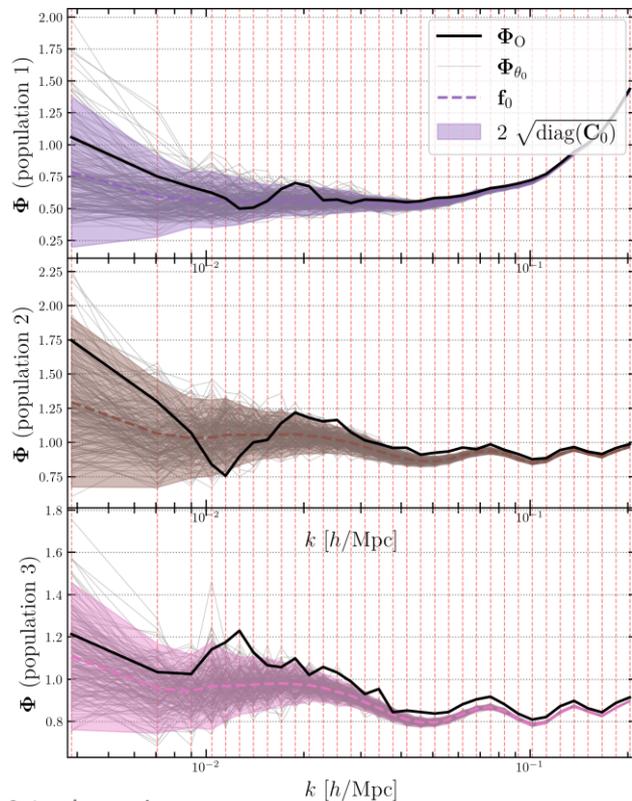
Synthetic observations

- Generated using model A:

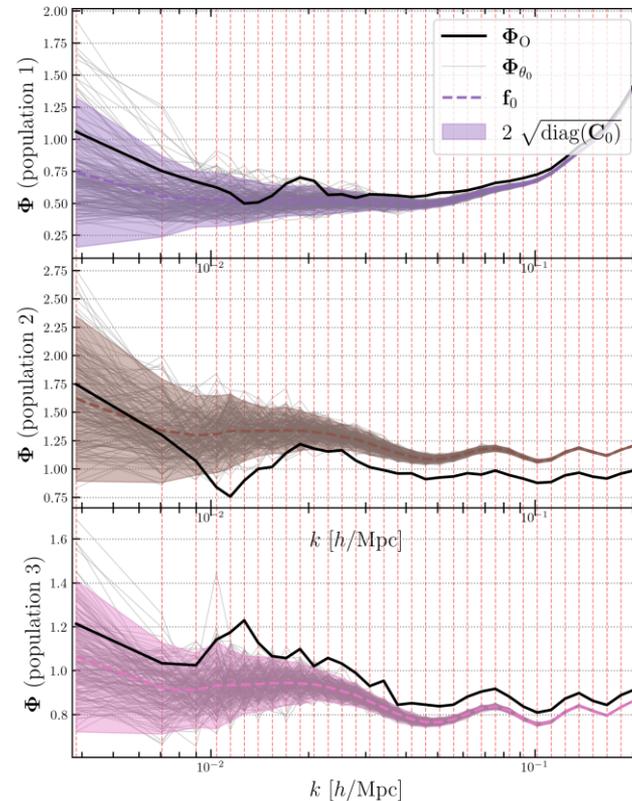


Synthetic observations versus simulations and linearised data models

Model A

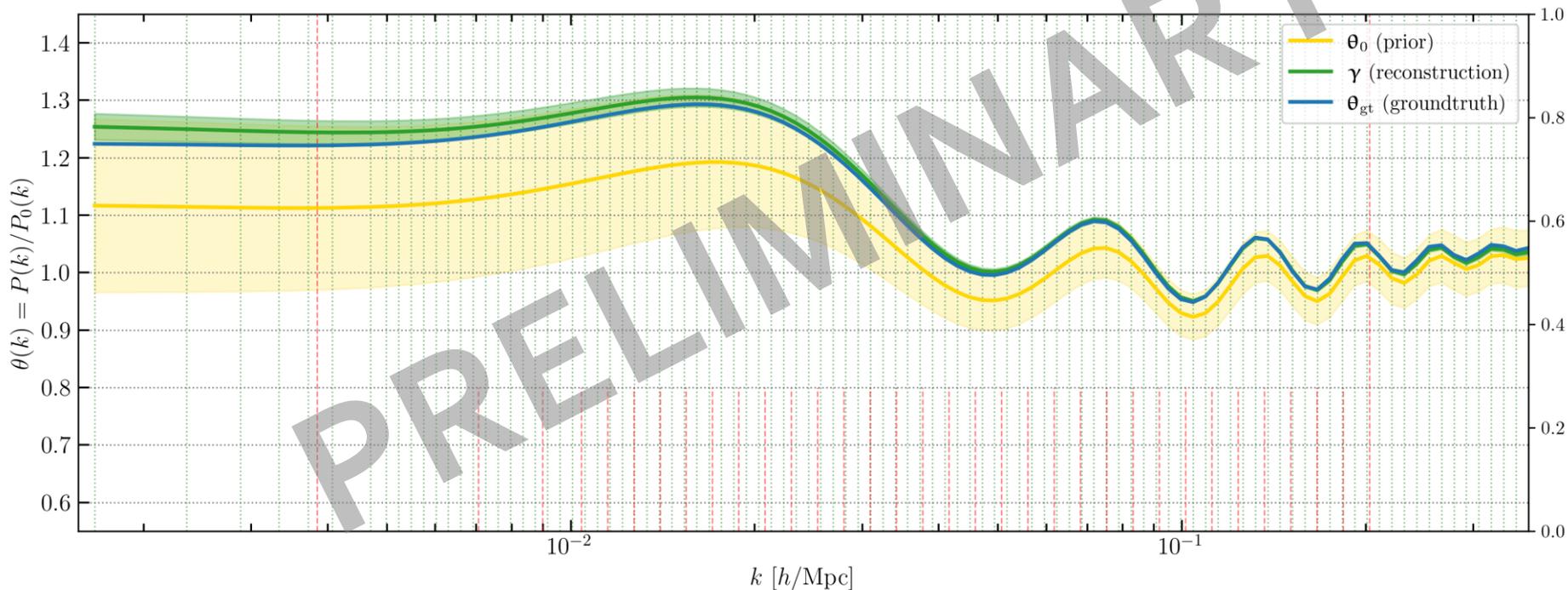


Model B



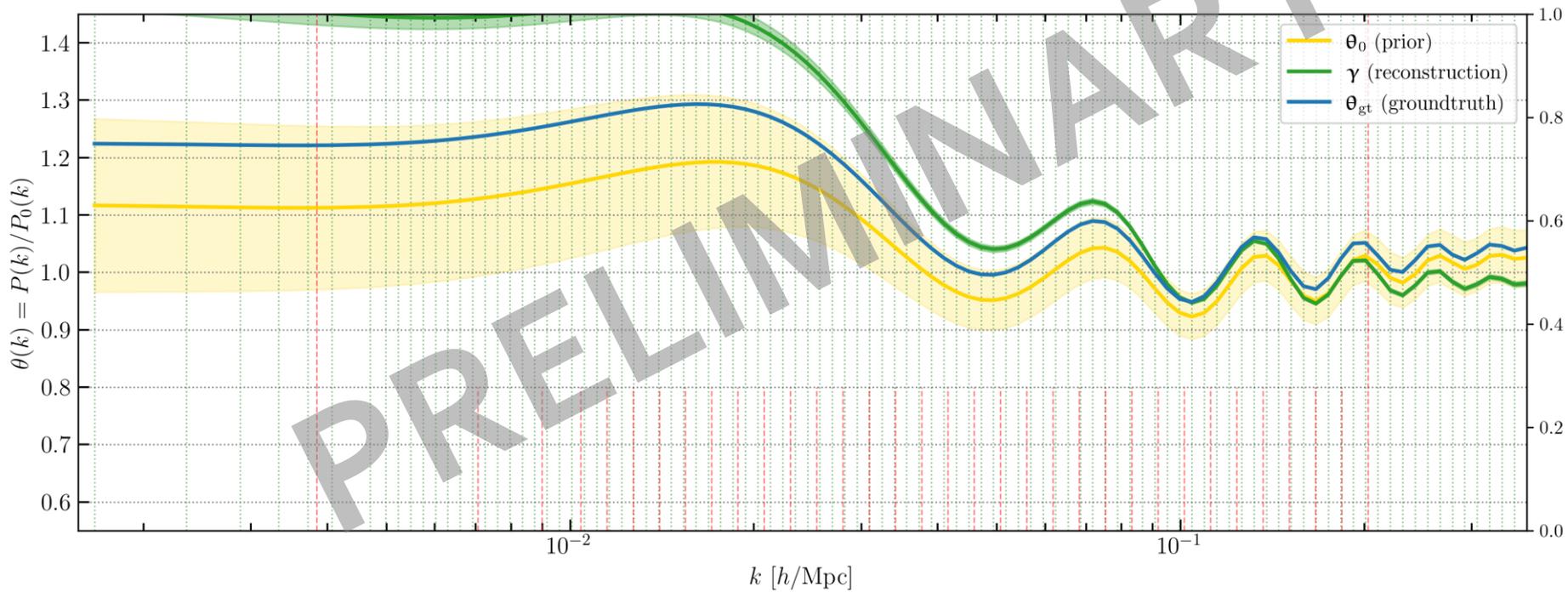
Model A

Reconstruction for the well specified model



Model B

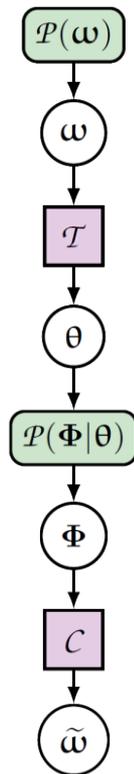
Misspecified model



Check for model misspecification and data compression for SBI

- Qualitatively: the shape of the reconstructed θ is useful as a [check for model misspecification](#) (independent theoretical understanding).
- Quantitatively: we can use the Mahalanobis distance between the reconstruction γ and the prior distribution $\mathcal{P}(\theta)$:

$$d_M(\gamma, \theta_0 | \mathbf{S}) \equiv \sqrt{(\gamma - \theta_0)^\top \mathbf{S}^{-1} (\gamma - \theta_0)}$$



- The score function $\nabla_{\omega} \hat{\ell}_{\omega_0}$ is the gradient of the log-likelihood at fiducial point ω_0 in parameter space.
- A quasi maximum-likelihood estimator for the parameters is

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

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

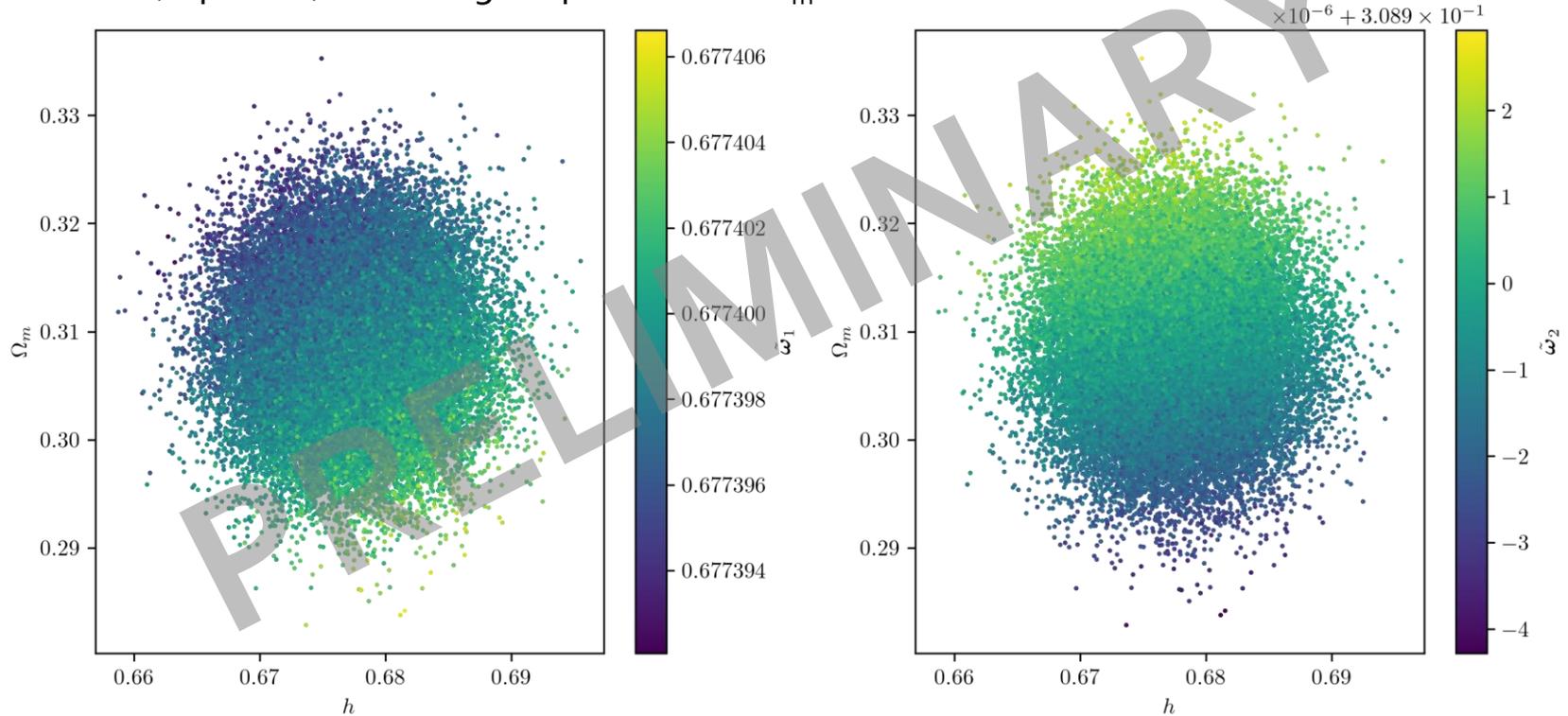
$$\nabla_{\omega} \mathbf{f}_0 = \nabla \mathbf{f}_0 \cdot \nabla_{\omega} \mathcal{T}_0$$

Already computed for SELF1 Cheap via finite differences

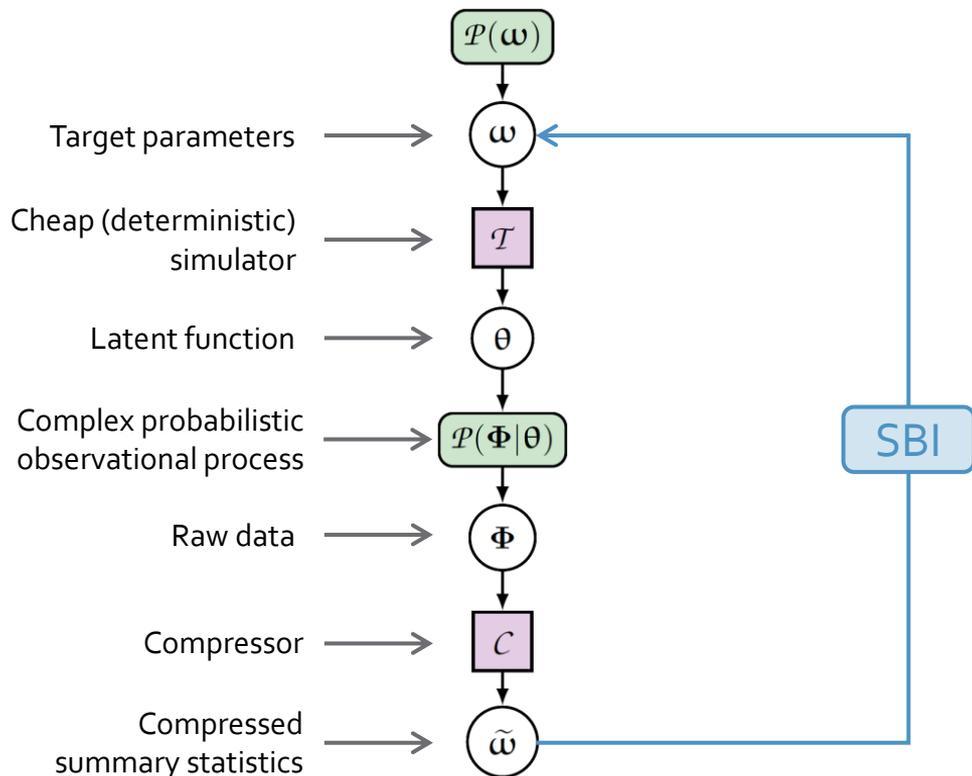
- Score compression is optimal in the sense that it [preserves the Fisher information content](#) of the data.

Score compression of the observed and simulated statistical summaries

- For two (top-level) cosmological parameters Ω_m and h :



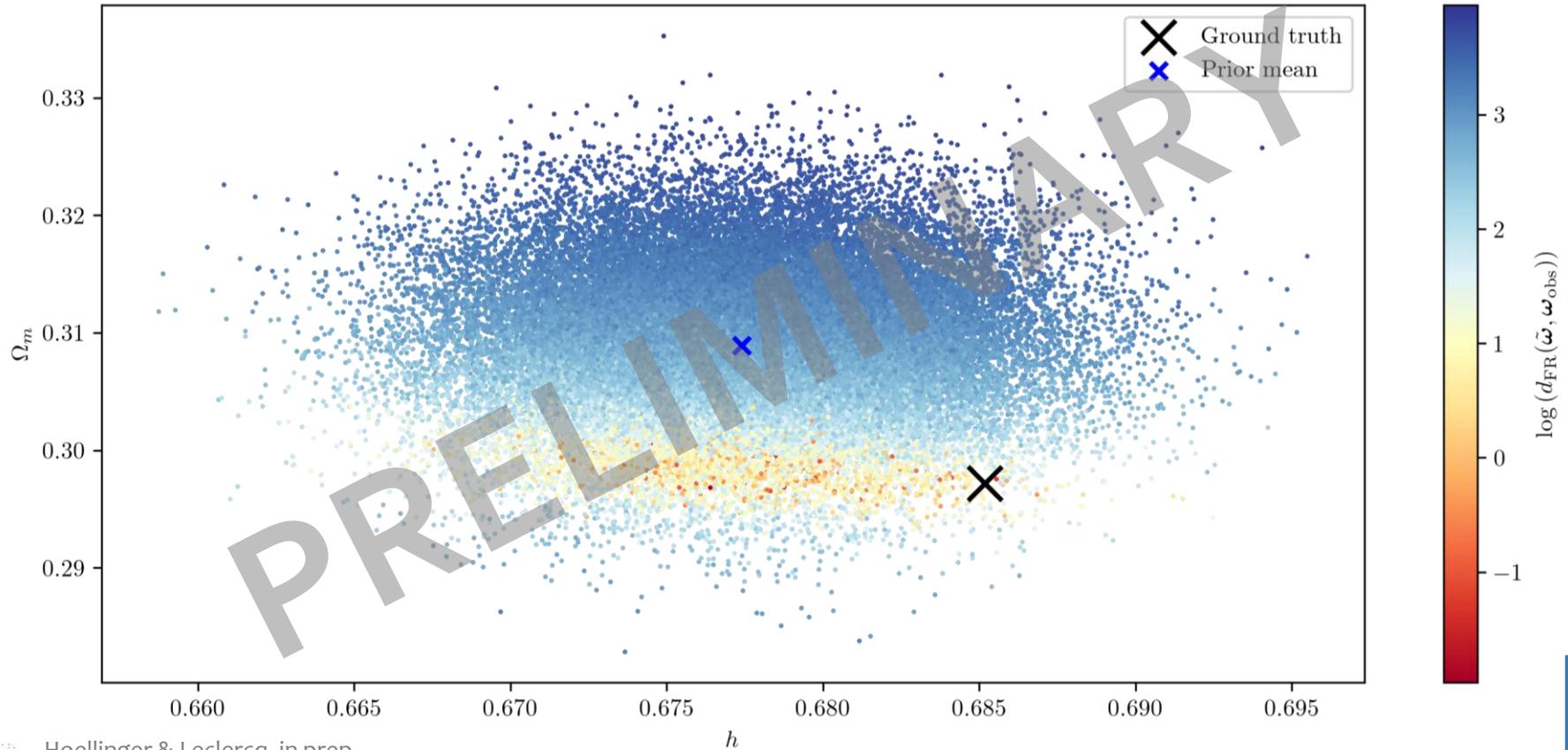
Simulation-based inference of top-level target parameters



- Any SBI algorithm can be used to obtain the posterior $\mathcal{P}(\omega|\tilde{\omega}_O)$.
- Final inference:
 - does not depend on the assumptions made to check for model misspecification,
 - is unbiased (only more conservative) in case data compression is lossy.
- Non-parametric approaches can use the **Fisher-Rao distance** between simulated summaries $\tilde{\omega}$ and observed summaries $\tilde{\omega}_O$:

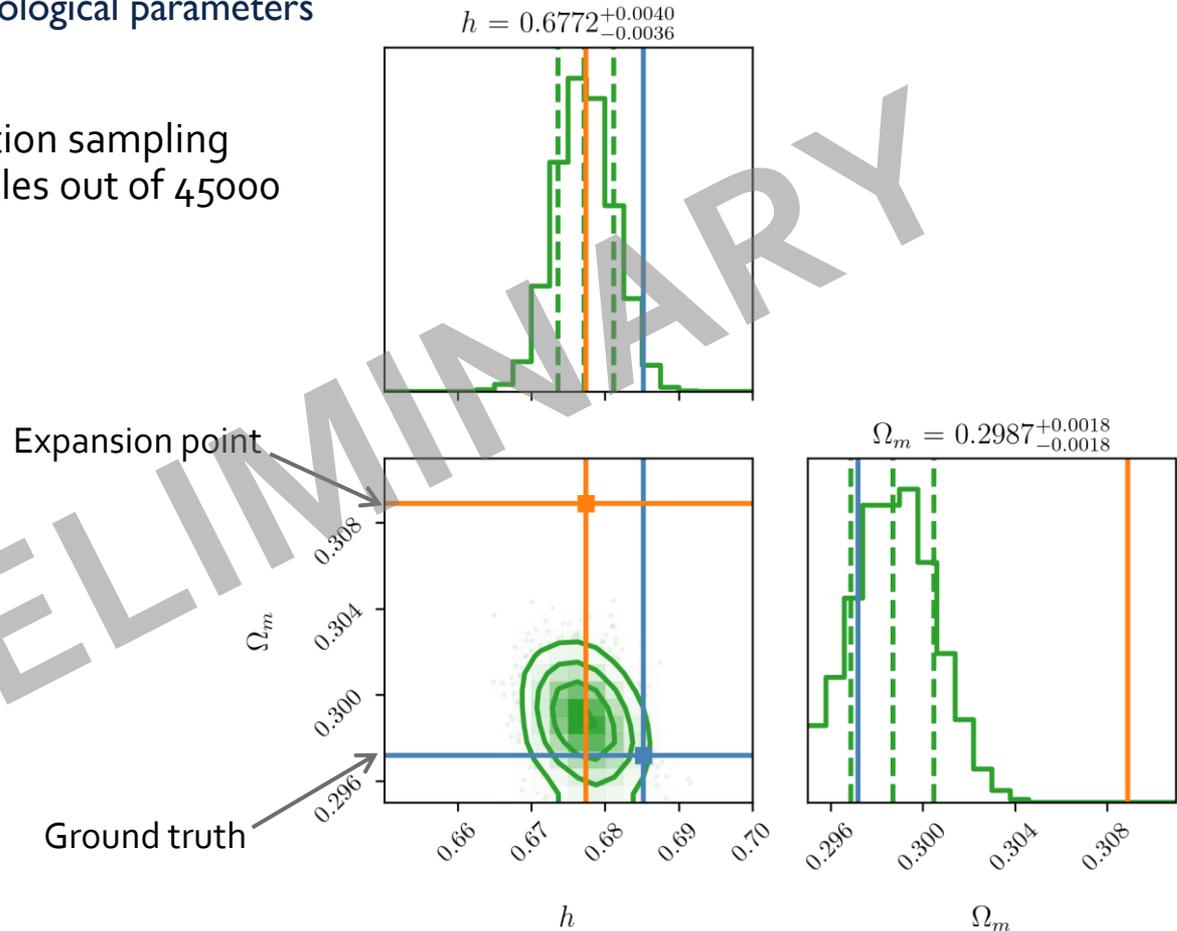
$$d_{\text{FR}}(\tilde{\omega}, \tilde{\omega}_O) \equiv \sqrt{(\tilde{\omega} - \tilde{\omega}_O)^\top \mathbf{F}_0(\tilde{\omega} - \tilde{\omega}_O)}$$

Fisher-Rao distance between simulated and observed summaries



Final SBI posterior on cosmological parameters

- After likelihood-free rejection sampling (ABC): 1687 selected samples out of 45000 simulations:



Conclusion: a science-ready statistical framework for the galaxy clustering additional probes

- A novel [two-step simulation based Bayesian approach](#), combining SELFI and SBI, to tackle the issue of model misspecification for a large class of BHM.
 - Advantages of the first step (SELFI):
 - Even if the inference is in high dimension, the simulator remains a black-box.
 - The number of simulations is fixed *a priori* by the user.
 - The computational workload is perfectly parallel.
 - The linearised data model is trained once and for all independently of the data vector (amortisation).
 - Advantages of the second step (SBI):
 - SELFI quantities provide a score compressor for free.
 - General advantages of SBI with respect to likelihood-based methods are preserved.
 - Inference does not depend on the assumptions made to check for model misspecification.
- A computationally efficient and easily applicable framework to perform [SBI of BHMs while checking for model misspecification](#).

pySELFI is publicly available at <https://pyselfi.florent-leclercq.eu>.



- Thanks for listening!

- References:

[FL, Enzi, Jasche & Heavens, 1902.10149](#)

[FL, 2209.11057](#)

Hoellinger & Leclercq, in prep.



www.florent-leclercq.eu

www.aquila-consortium.org

