# Implicit likelihood inference from galaxy survey data with robustness to model misspecification

## 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), Guilhem Lavaux (IAP)
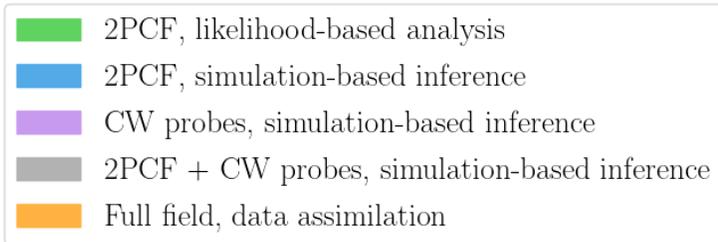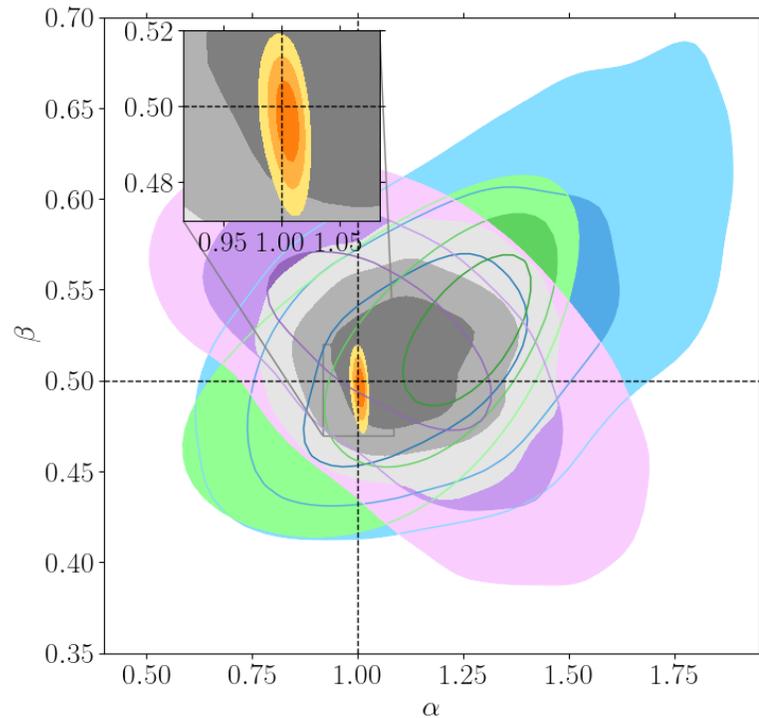

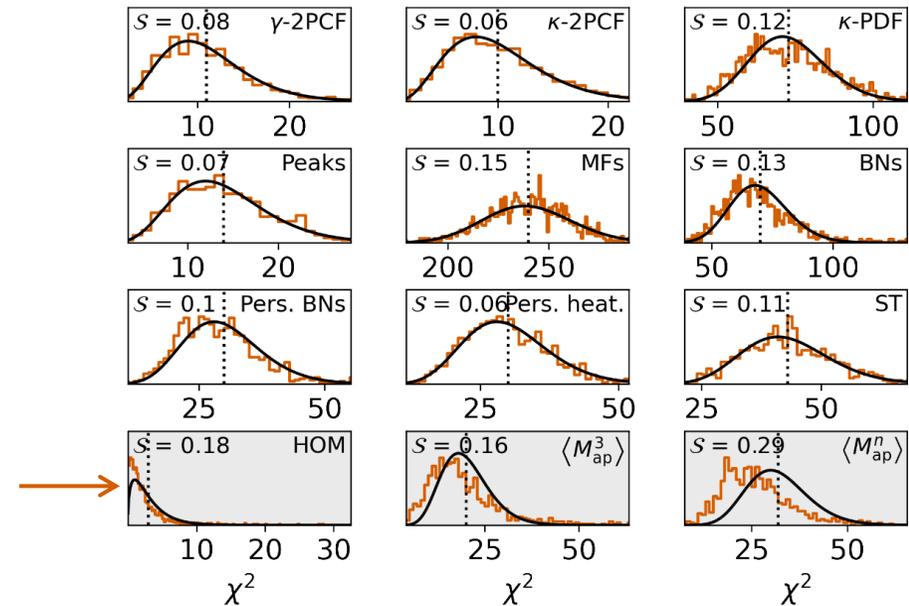and the Aquila Consortium

www.aquila-consortium.org

**12 January 2024**

# Why I decided to go "implicit" for galaxy clustering additional probes

Note: likelihood-free inference (LFI) ≈ simulation-based inference (SBI) ≈ implicit likelihood inference (ILI)

- A question of accuracy: first, avoid biases.



2PCF, likelihood-based analysis
2PCF, simulation-based inference
CW probes, simulation-based inference
2PCF + CW probes, simulation-based inference
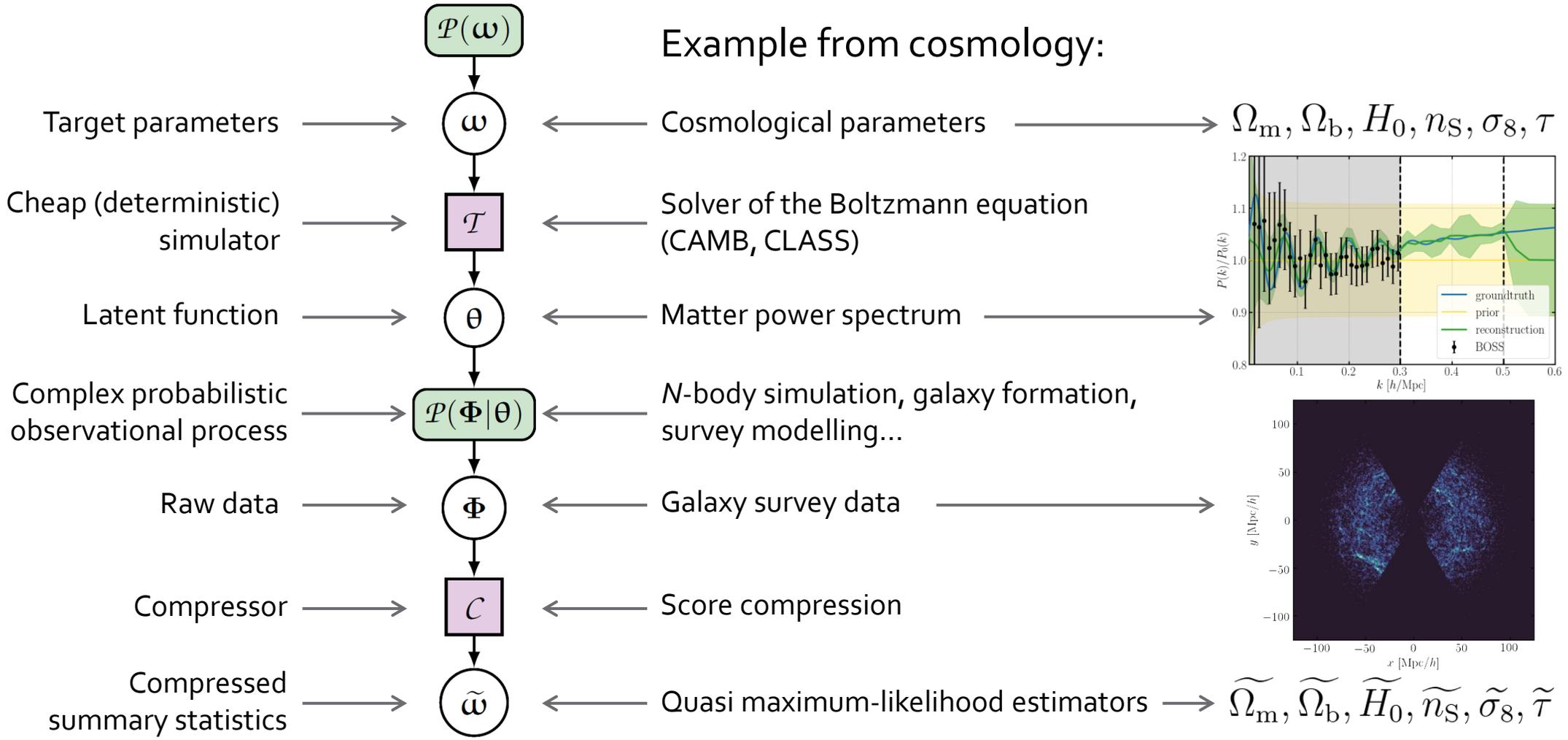Full field, data assimilation

- Some weak lensing additional probes also have a non-Gaussian distribution.



- A question of precision: can numerical forward models be used to push further than $k \gtrsim 0.15\ h/\mathrm{Mpc}$? The full field contains much more information.
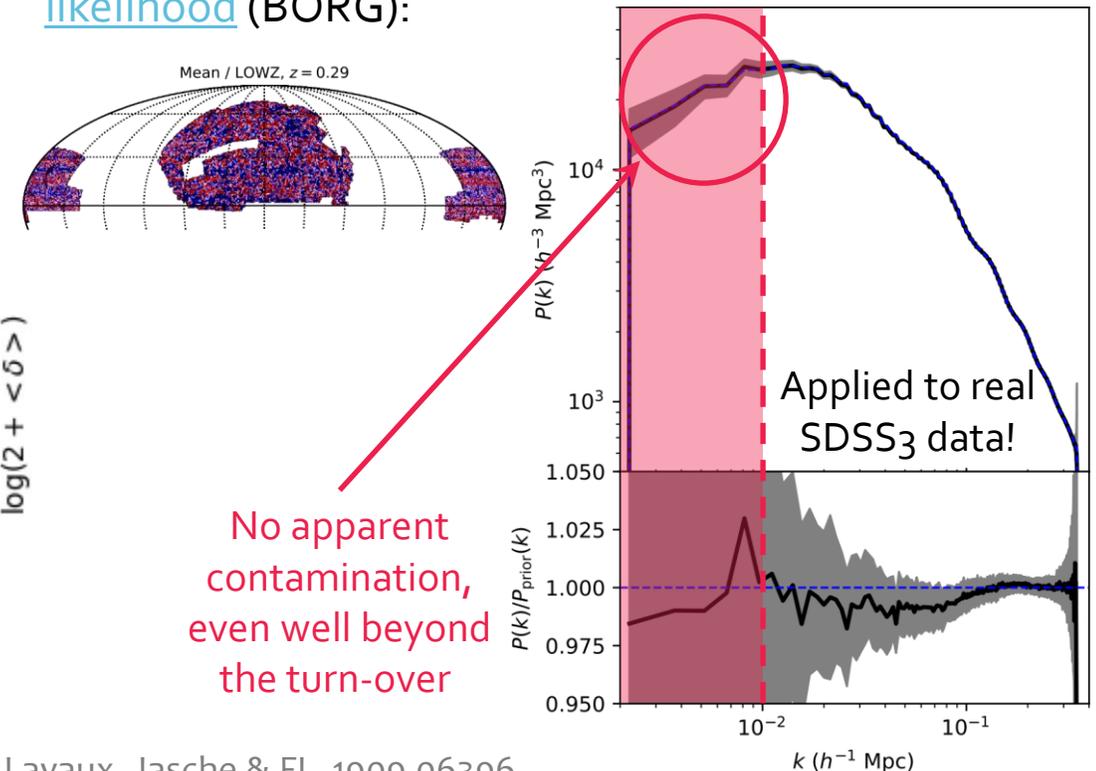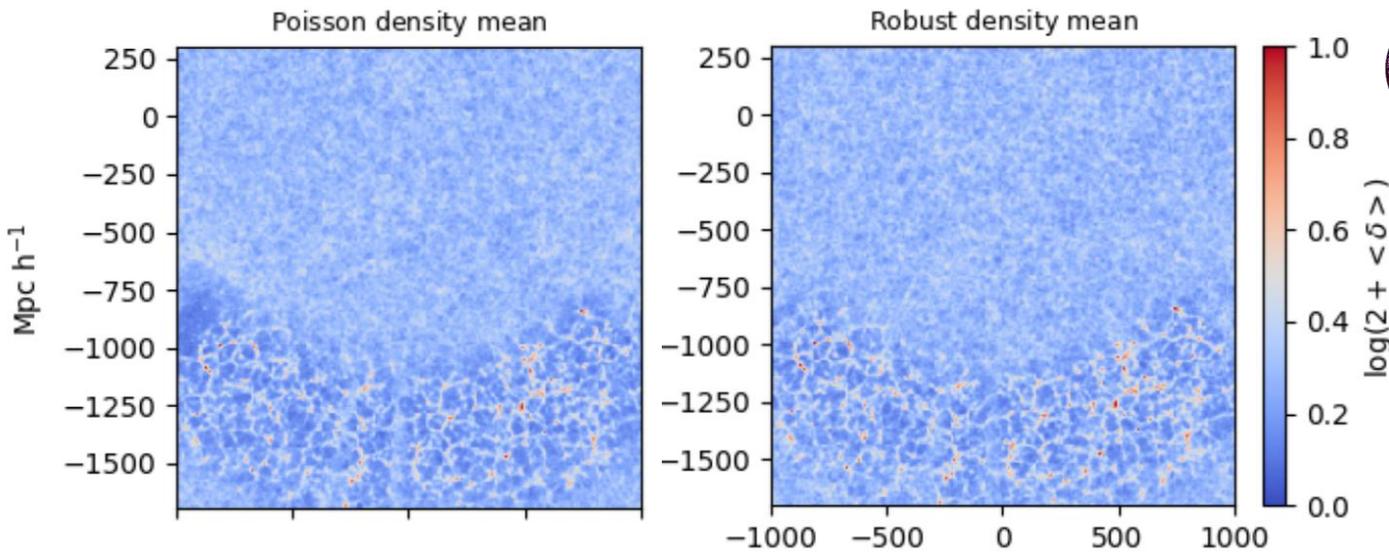
FL & Heavens, 2103.04158

Euclid HOWLS-KP paper 1, Ajani *et al.*, 2301.12890

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

## Example from cosmology:

$\mathcal{P}(\boldsymbol{\omega})$

Target parameters → $\boldsymbol{\omega}$ ← Cosmological parameters → $\Omega_{\mathrm{m}}, \Omega_{\mathrm{b}}, H_0, n_{\mathrm{S}}, \sigma_8, \tau$

Cheap (deterministic) simulator → $\mathcal{T}$ ← Solver of the Boltzmann equation (CAMB, CLASS)

Latent function → $\theta$ ← Matter power spectrum →

Complex probabilistic observational process → $\mathcal{P}(\boldsymbol{\Phi}|\theta)$ ← $N$-body simulation, galaxy formation, survey modelling...

Raw data → $\boldsymbol{\Phi}$ ← Galaxy survey data →

Compressor → $\mathcal{C}$ ← Score compression

Compressed summary statistics → $\widetilde{\omega}$ ← Quasi maximum-likelihood estimators → $\widetilde{\Omega}_{\mathrm{m}}, \widetilde{\Omega}_{\mathrm{b}}, \widetilde{H}_0, \widetilde{n}_{\mathrm{S}}, \widetilde{\sigma}_8, \widetilde{\tau}$

# Model misspecification and unknown systematics with an explicit field-level likelihood

- **Model misspecification** is a long-standing problem for Bayesian inference: when the model differs from the actual data-generating process, posteriors tend to be biased and/or overly concentrated.

- This issue is particularly critical for cosmological data analysis in the presence of systematic effects.

- In cosmology, we are sometimes unable to formulate **any** model that fits the data in some regimes.

- Machine-aided report of unknown systematic effects is possible with an explicit field-level likelihood (BORG):



Mean / LOWZ, $z = 0.29$

Applied to real SDSS$_3$ data!

No apparent contamination, even well beyond the turn-over

Poisson density mean

Robust density mean

$\log(2 + <\delta>)$

Porqueres, Ramanah, Jasche & Lavaux, 1812.05113

Lavaux, Jasche & FL, 1909.06396

# Key idea: a two-step ILI process that recycles simulations
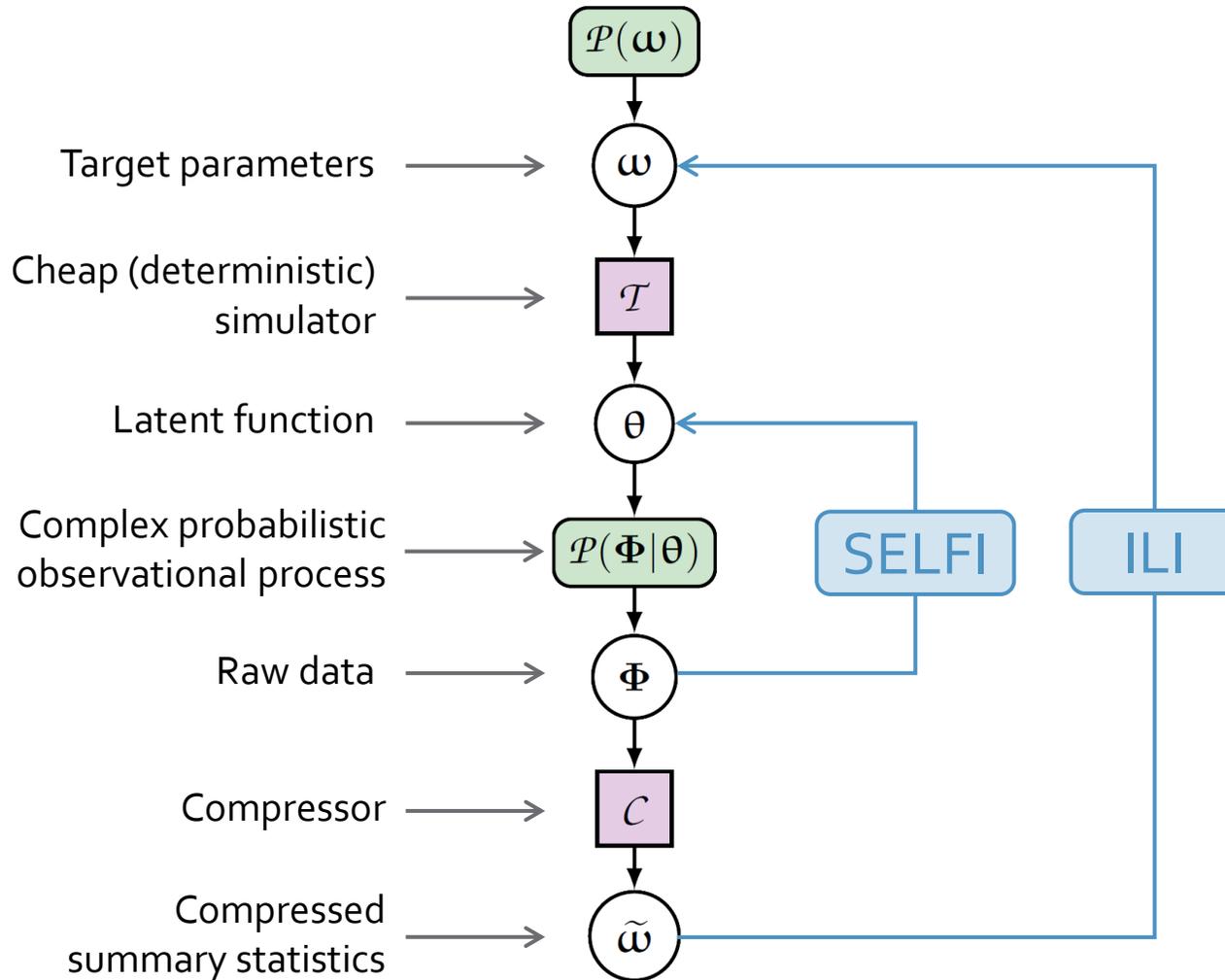
1. Inference of the latent function **θ**, to check for model misspecification:
   - SELFI algorithm

Prior on the latent function $\longrightarrow$ $\mathcal{P}(\theta)$

Latent function $\longrightarrow$ $\theta$

Complex probabilistic observational process $\longrightarrow$ $\mathcal{P}(\Phi|\theta)$

SELFI

Raw data $\longrightarrow$ $\Phi$

# Key idea: a two-step ILI process that recycles simulations
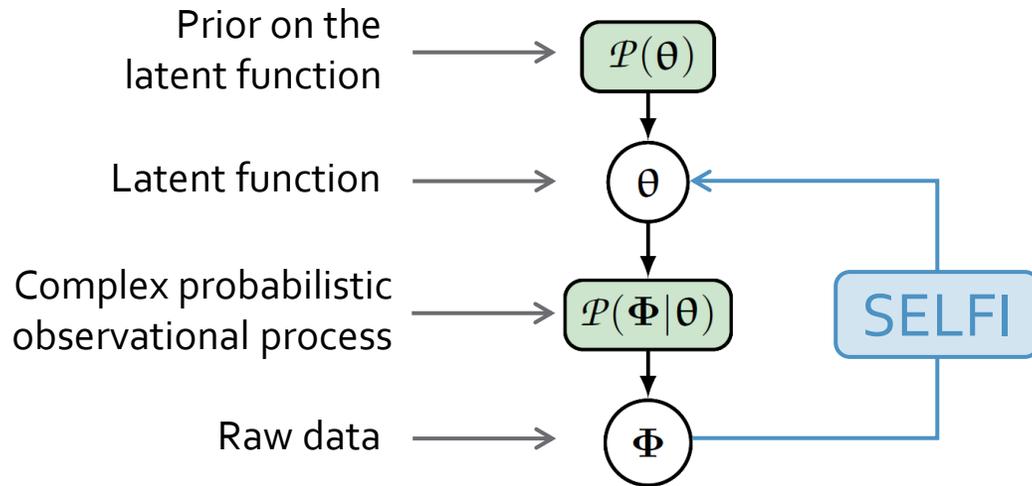


1. Inference of the latent function $\theta$, to check for model misspecification:
   - SELFI algorithm

2. Implicit likelihood inference of $\omega$ :
   - 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.

- Linearisation of the black-box data model:

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

- Further assume:
  - Gaussian prior: $\mathcal{P}(\boldsymbol{\theta}) = \mathcal{G}(\boldsymbol{\theta}_0, \mathbf{S})$
  - Gaussian effective likelihood:

$$\mathcal{P}(\mathbf{\Phi}|\boldsymbol{\theta}) = \mathcal{G}\left[\mathbf{f}(\boldsymbol{\theta}), \mathbf{C}_0\right]$$

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

> expansion point       observed summaries
>
> mean: $\boldsymbol{\gamma} \equiv \boldsymbol{\theta}_0 + \boldsymbol{\Gamma}\,(\nabla\mathbf{f}_0)^{\mathsf{T}}\,\mathbf{C}_0^{-1}(\mathbf{\Phi}_O - \mathbf{f}_0)$
>
> covariance: $\boldsymbol{\Gamma} \equiv \left[(\nabla\mathbf{f}_0)^{\mathsf{T}}\,\mathbf{C}_0^{-1}\nabla\mathbf{f}_0 + \mathbf{S}^{-1}\right]^{-1}$
>
> covariance of summaries    prior covariance
> gradient of the black-box

- $\mathbf{f}_0, \mathbf{C}_0$ 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.

- Numerical data models allow using the galaxy power spectrum as summary statistics up to at least $k \gtrsim 0.5 \, h/\mathrm{Mpc}$ safely

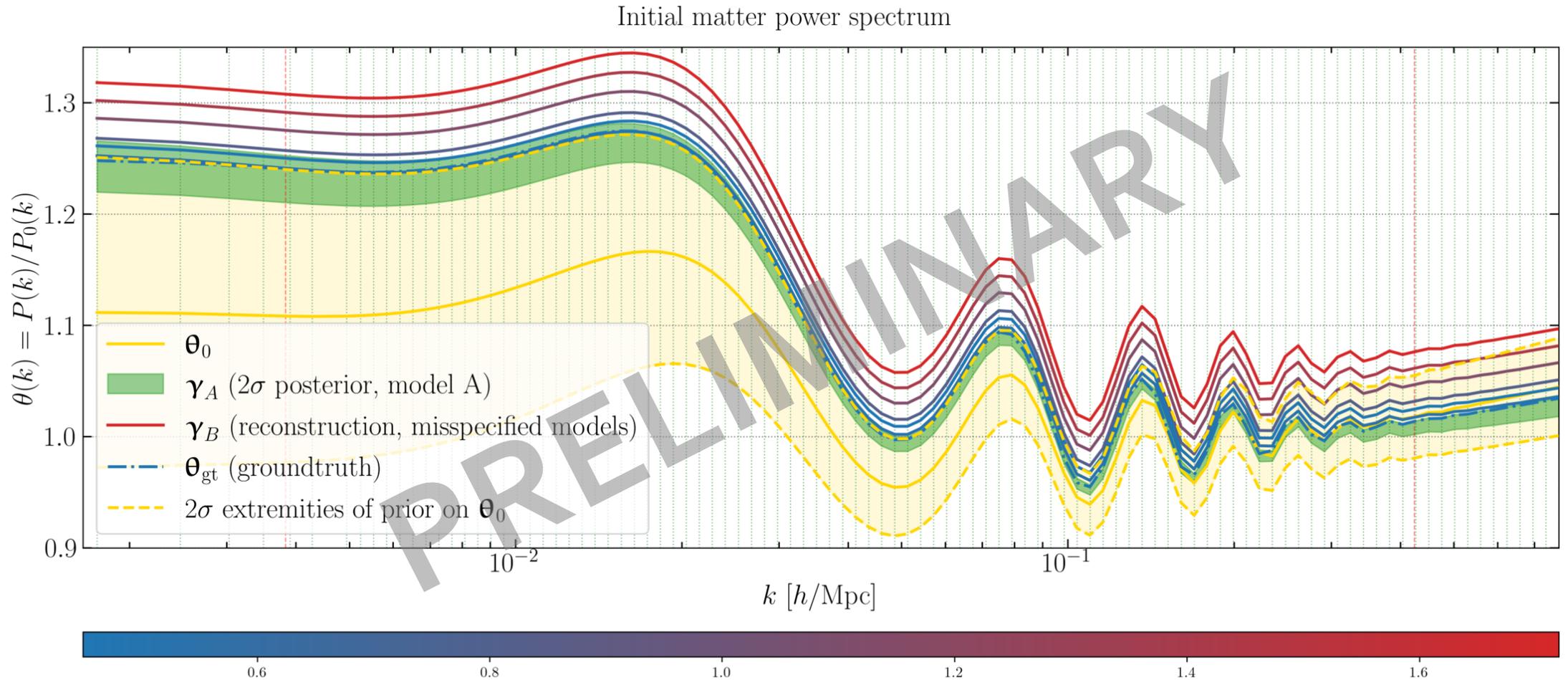- $N_{\mathrm{modes}} \propto k^3$: 5 times more modes are used in the analysis.



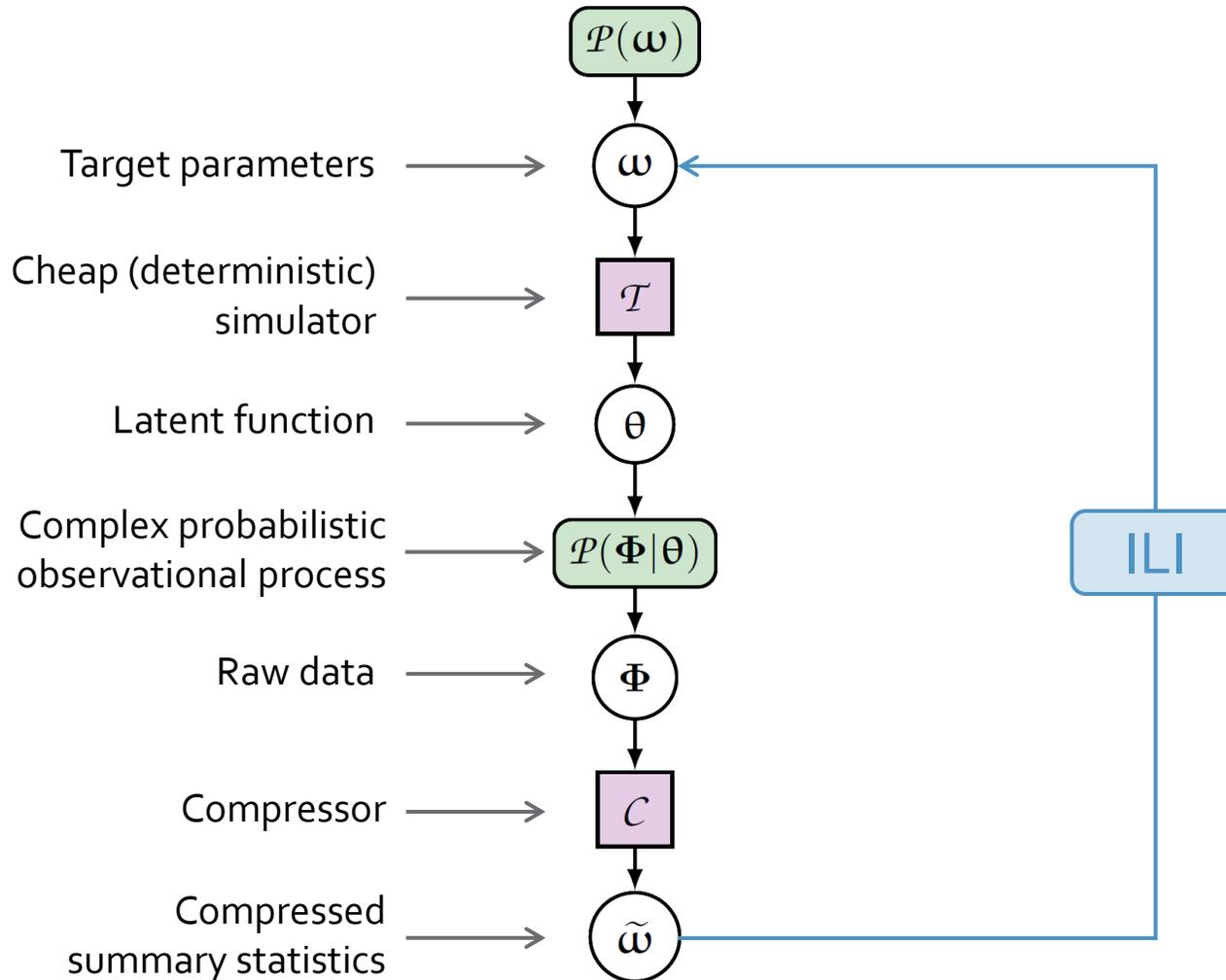Data points from Beutler *et al.*, 1607.03149

- One can utilise the initial matter power spectrum to check for systematics.



Initial matter power spectrum

Legend:
- $\boldsymbol{\theta}_0$
- $\boldsymbol{\gamma}_A$ (2$\sigma$ posterior, model A)
- $\boldsymbol{\gamma}_B$ (reconstruction, misspecified models)
- $\boldsymbol{\theta}_{gt}$ (groundtruth)
- 2$\sigma$ extremities of prior on $\boldsymbol{\theta}_0$

Vertical axis: $\theta(k) = P(k)/P_0(k)$

Horizontal axis: $k$ [$h$/Mpc]

Colorbar: Average percentage error on the galaxy biases

# Step 2: implicit likelihood inference of top-level target parameters



Target parameters → $\omega$

Cheap (deterministic) simulator → $\mathcal{T}$

Latent function → $\theta$

Complex probabilistic observational process → $\mathcal{P}(\boldsymbol{\Phi}|\boldsymbol{\theta})$

Raw data → $\boldsymbol{\Phi}$

Compressor → $\mathcal{C}$

Compressed summary statistics → $\widetilde{\omega}$

$\mathcal{P}(\boldsymbol{\omega})$

ILI
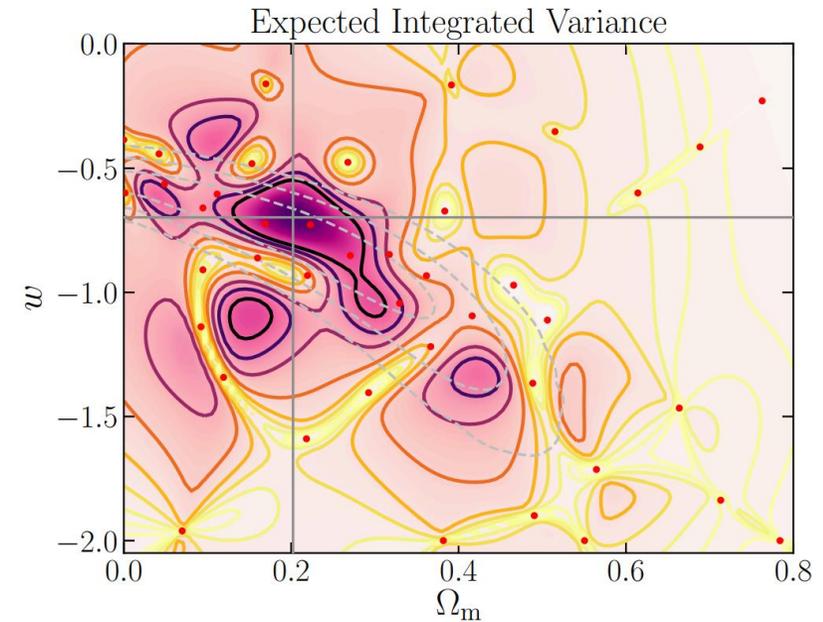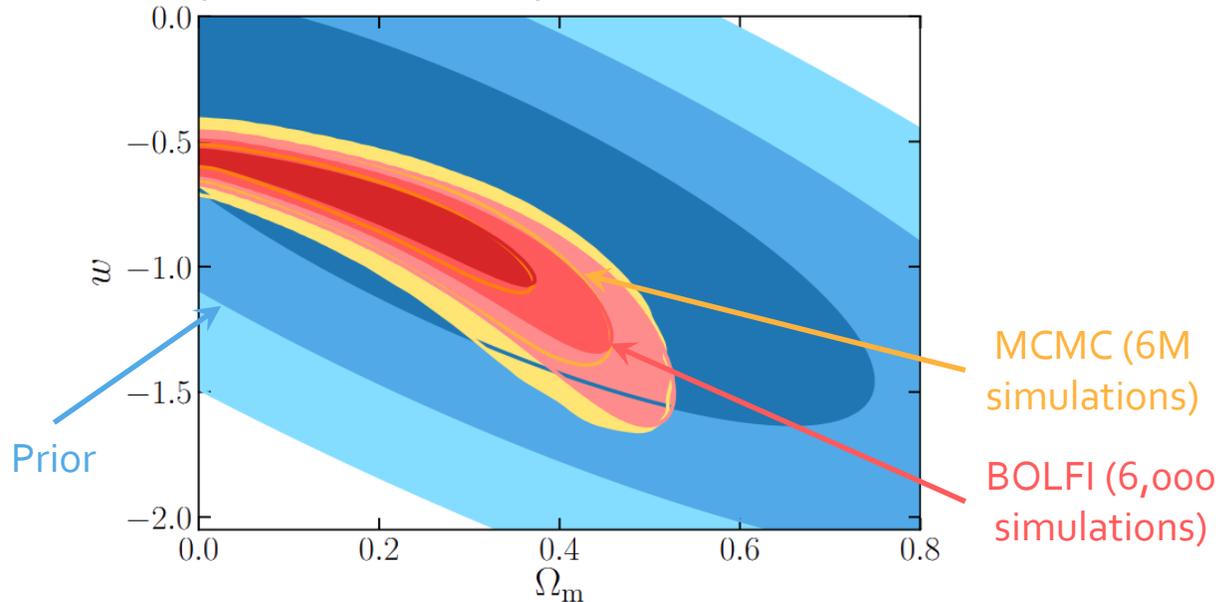
- The simulations used for step 1 can be recycled to write a free score compressor for step 2.

- Any ILI algorithm can be used to obtain the posterior $\mathcal{P}(\boldsymbol{\omega}|\widetilde{\boldsymbol{\omega}}_{\mathrm{O}})$.

- 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 $\widetilde{\omega}$ and observed summaries $\widetilde{\omega}_{\mathrm{O}}$:

$$d_{\mathrm{FR}}(\widetilde{\boldsymbol{\omega}}, \widetilde{\boldsymbol{\omega}}_{\mathrm{O}}) \equiv \sqrt{(\widetilde{\boldsymbol{\omega}} - \widetilde{\boldsymbol{\omega}}_{\mathrm{O}})^{\mathsf{T}} \mathbf{F}_0 (\widetilde{\boldsymbol{\omega}} - \widetilde{\boldsymbol{\omega}}_{\mathrm{O}})}$$

# Dealing with expensive simulators in ILI problems:
## The BOLFI algorithm (*Bayesian Optimisation for Likelihood-Free Inference*)

- The simulator will typically be extremely expensive (*N*-body simulation, halo finding, complex observational effects). We can typically afford O(10,000) evaluations.

- Emulation of the data model is not the only option.

- BOLFI (*Bayesian Optimisation for Likelihood-Free Inference*) uses an acquisition function to place expensive simulations in the parameter space.

- The optimal acquisition function for implicit inference can be derived: the Expected Integrated Variance.
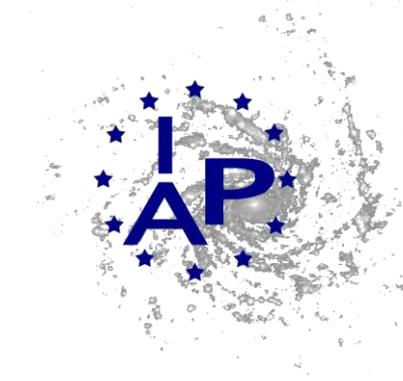
Re-analysis of the JLA supernovae data:

- A novel two-step implicit likelihood inference approach, combining SELFI and BOLFI, to tackle the issue of model misspecification for a large class of BHMs.

- 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 (ILI/BOLFI):
  - SELFI quantities provide a score compressor for free.
  - General advantages of ILI with respect to likelihood-based methods are preserved.
  - Inference does not depend on the assumptions made to check for model misspecification.
  - BOLFI uses active acquisition to deal with expensive simulators.

- ➤ A computationally efficient and easily applicable framework to perform ILI of BHMs while checking for model misspecification.

References:

- Leclercq 2018, 1805.07152, *Bayesian optimisation for likelihood-free cosmological inference*

- Leclercq *et al.* 2019, 1902.10149, *Primordial power spectrum and cosmology from black-box galaxy surveys*

- Leclercq 2022, 2209.11057, *Simulation-based inference of Bayesian hierarchical models while checking for model misspecification*

- Hoellinger & Leclercq, in prep.

https://pyselfi.florent-leclercq.eu: publicly available implementation of SELFI
https://aquila-consortium.org