# Implicit Likelihood Inference while efficiently checking for survey systematics

Euclid Galaxy Clustering meeting, Marseille 2024

## Florent Leclercq

www.florent-leclercq.eu

Institut d'Astrophysique de Paris
CNRS & Sorbonne Université

In collaboration with:
Tristan Hoellinger (IAP)

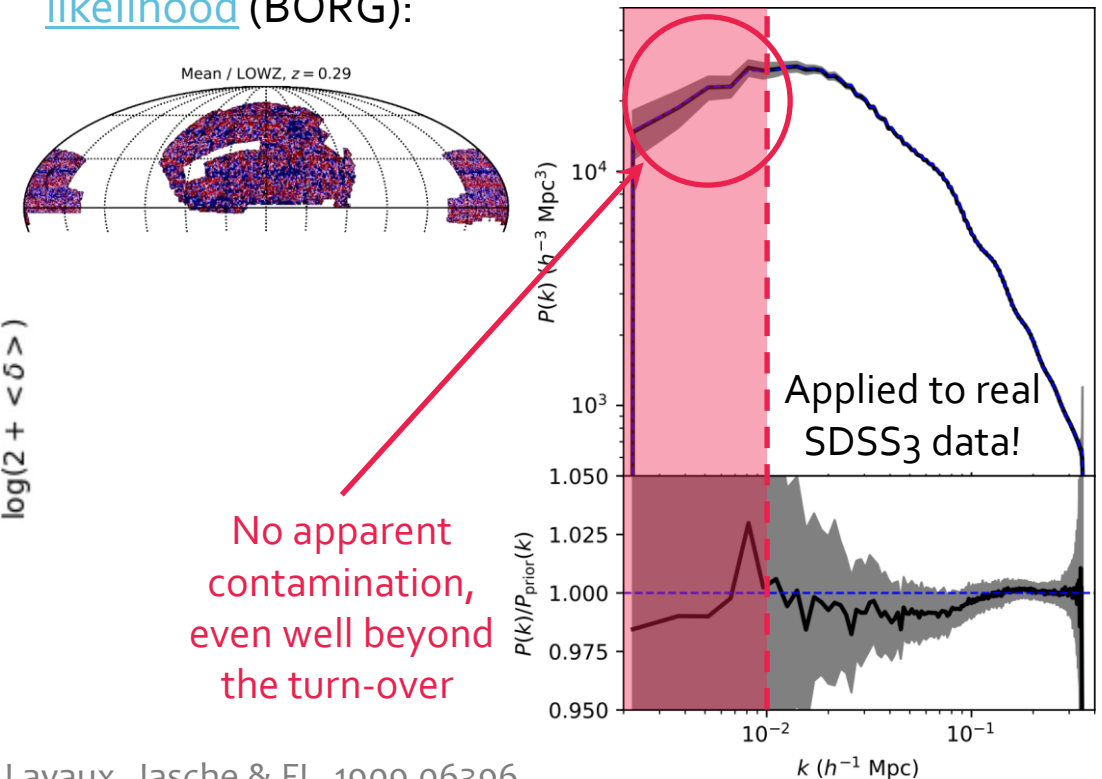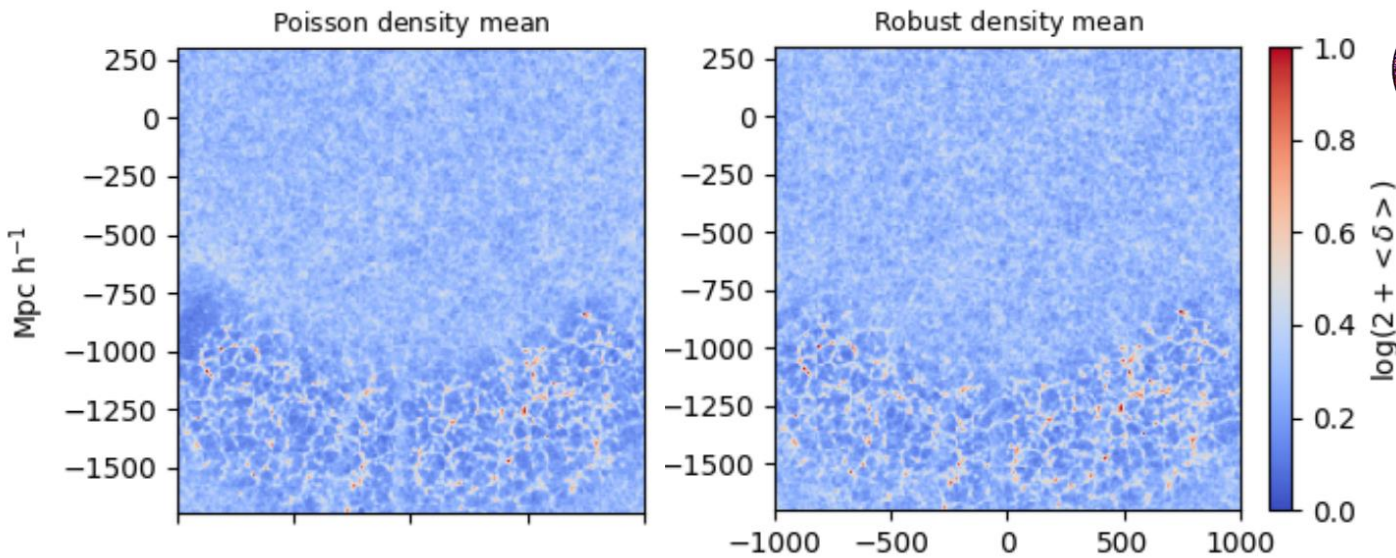and the Aquila Consortium

www.aquila-consortium.org

**1 February 2024**

# 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):



Poisson density mean

Robust density mean

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

Mean / LOWZ, $z = 0.29$

Applied to real SDSS$_3$ data!

No apparent contamination, even well beyond the turn-over

$P(k)$ ($h^{-3}$ Mpc$^3$)
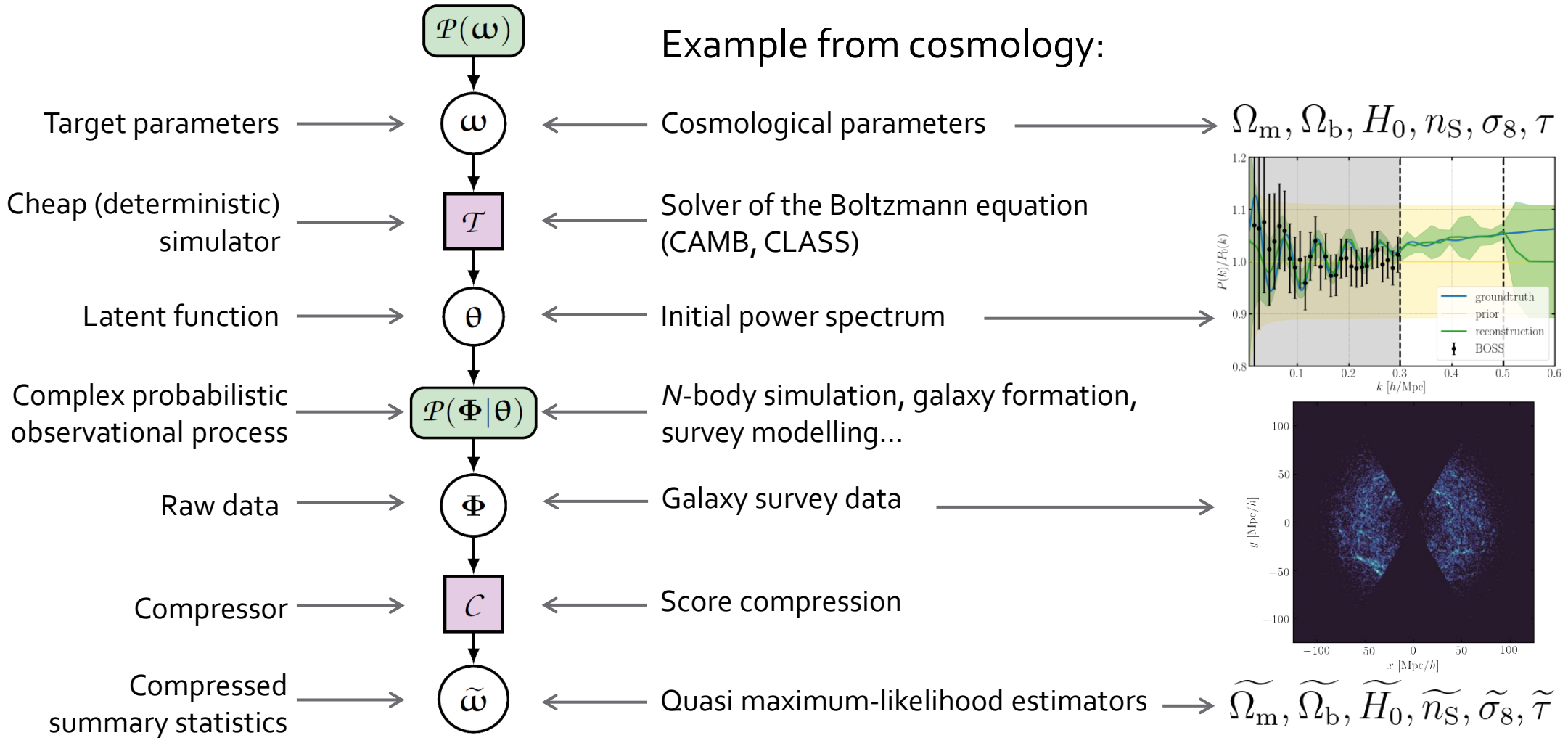
$P(k)/P_{prior}(k)$

$k$ ($h^{-1}$ Mpc)

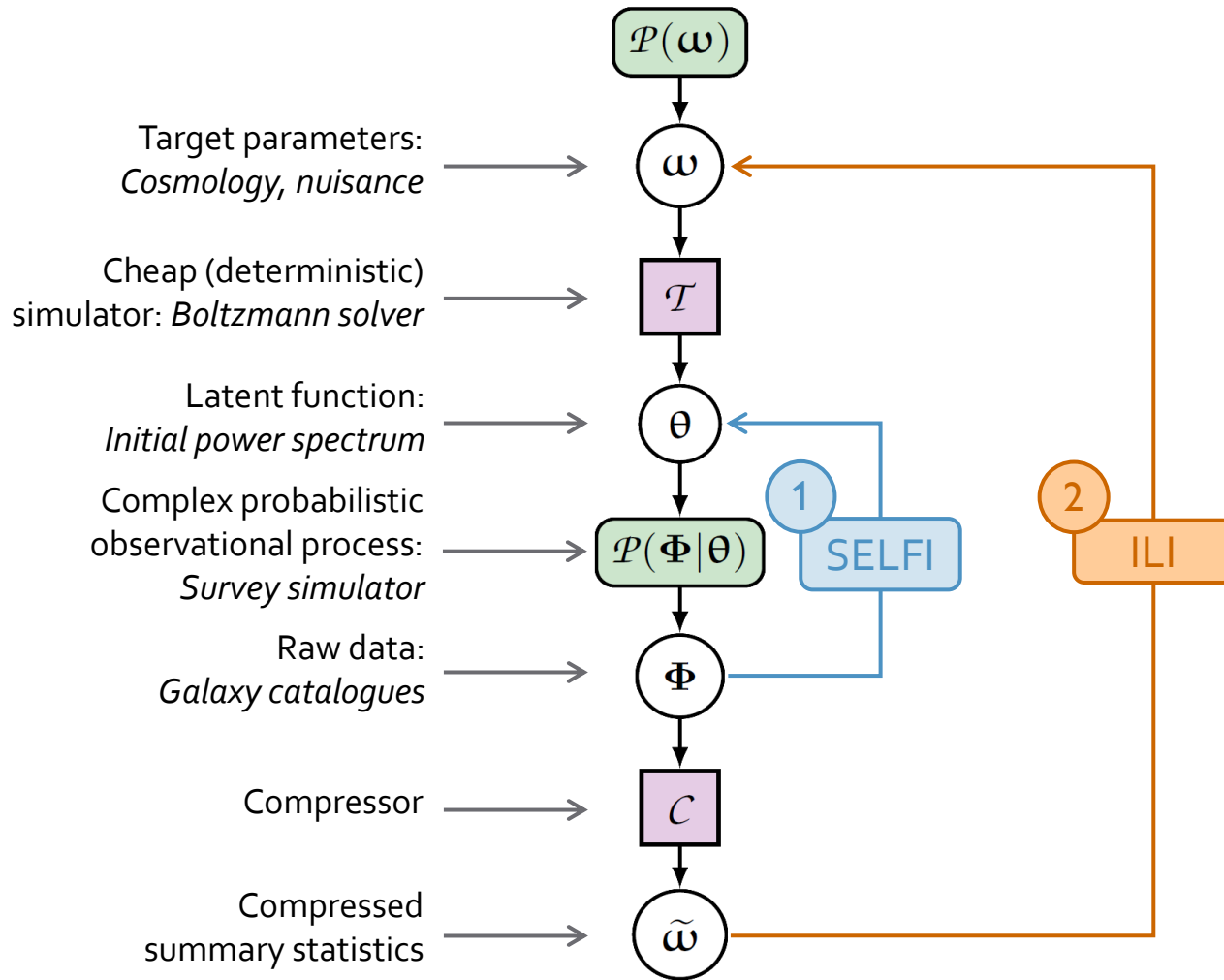Porqueres, Ramanah, Jasche & Lavaux, 1812.05113

Lavaux, Jasche & FL, 1909.06396

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



Example from cosmology:

| | | |
|---|---|---|
| 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 | $\boldsymbol{\theta}$ | Initial power spectrum |
| Complex probabilistic observational process | $\mathcal{P}(\boldsymbol{\Phi}\|\boldsymbol{\theta})$ | $N$-body simulation, galaxy formation, survey modelling... |
| Raw data | $\boldsymbol{\Phi}$ | Galaxy survey data |
| Compressor | $\mathcal{C}$ | Score compression |
| Compressed summary statistics | $\widetilde{\boldsymbol{\omega}}$ | Quasi maximum-likelihood estimators → $\widetilde{\Omega_{\mathrm{m}}}, \widetilde{\Omega_{\mathrm{b}}}, \widetilde{H_0}, \widetilde{n_{\mathrm{S}}}, \widetilde{\sigma_8}, \widetilde{\tau}$ |

# Key idea: a two-step implicit likelihood inference (ILI) process that recycles simulations



Target parameters:
*Cosmology, nuisance*

Cheap (deterministic)
simulator: *Boltzmann solver*

Latent function:
*Initial power spectrum*

Complex probabilistic
observational process:
*Survey simulator*

Raw data:
*Galaxy catalogues*

Compressor

Compressed
summary statistics

**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**

# Initial power spectrum inference:
# the SELFI approach (*Simulator Expansion for Likelihood-Free Inference*)

Prior on the *initial matter power spectrum* →

Latent function: *Initial power spectrum* →

Complex probabilistic observational process: *Survey simulator* →

Raw data: *Galaxy catalogues* →

$\mathcal{P}(\theta)$

$\theta$

1 SELFI

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

$\Phi$

- Linearisation of the black-box:

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

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

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

  expansion point            observed summaries

  mean: $\boldsymbol{\gamma} \equiv \theta_0 + \boldsymbol{\Gamma}\, (\nabla\mathbf{f}_0)^\mathsf{T}\, \mathbf{C}_0^{-1} (\boldsymbol{\Phi}_\mathrm{O} - \mathbf{f}_0)$

  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

FL, Enzi, Jasche & Heavens, 1902.1014; FL, 2209.11057; Hoellinger & Leclercq, in prep.

# Check for model misspecification

- Qualitatively: the shape of the reconstructed initial power spectrum $\theta$ is useful as a check for unknown systematics / model misspecification (using our independent theoretical understanding).

- Quantitatively: we can use the Mahalanobis distance between the reconstruction $\gamma$ and the prior distribution $\mathcal{P}(\theta)$:
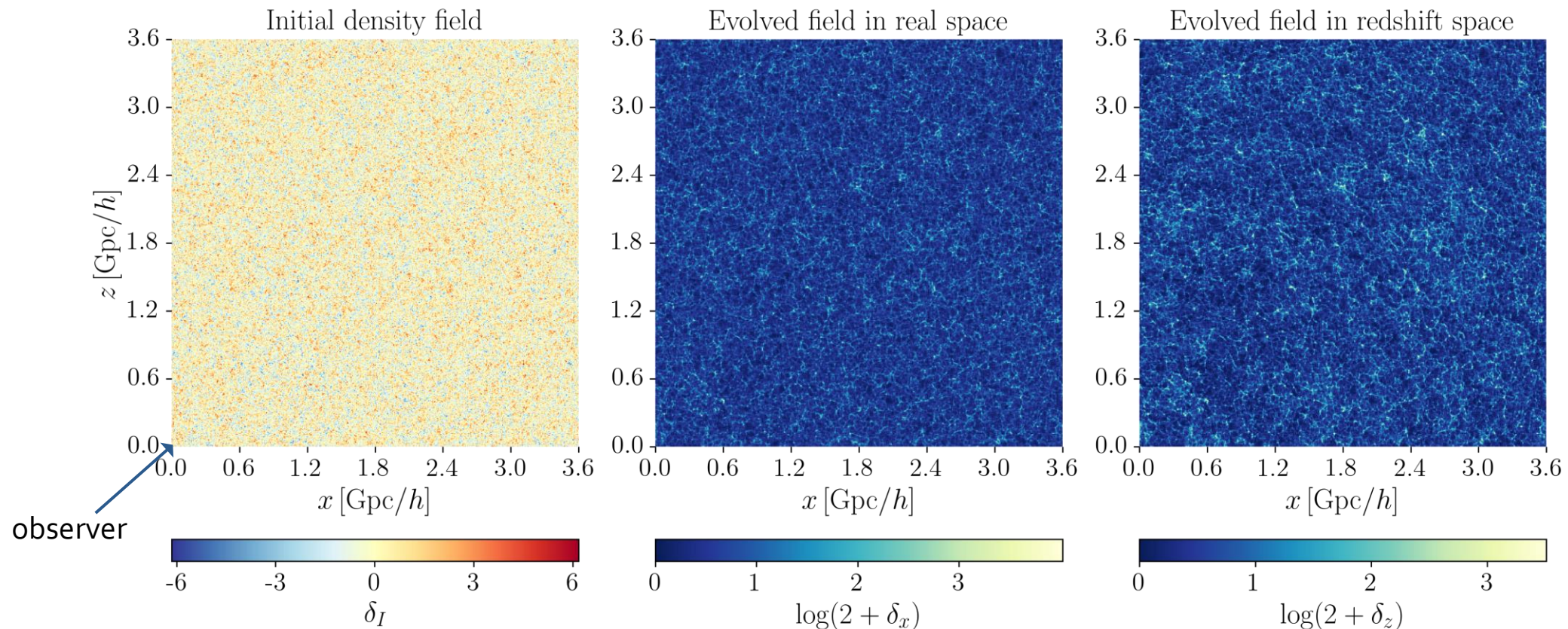
$$d_{\mathrm{M}}(\gamma, \theta_0 | \mathbf{S}) \equiv \sqrt{(\gamma - \theta_0)^{\mathsf{T}}\, \mathbf{S}^{-1} (\gamma - \theta_0)}$$

# Simulator-based data model of galaxy surveys

- θ defined on S = 100 support wavenumbers

- Flat ΛCDM assumed

- Gravitational evolution (*N*-body) using Simbelmynë
  Leclercq, Jasche & Wandelt, 1502.02690; http://simbelmyne.florent-leclercq.eu

  - $512^3$ dark matter particles, 2LPT up to $z = 19$
  - Particle-mesh grid of $1024^3$ voxels, COLA to $z = 0$
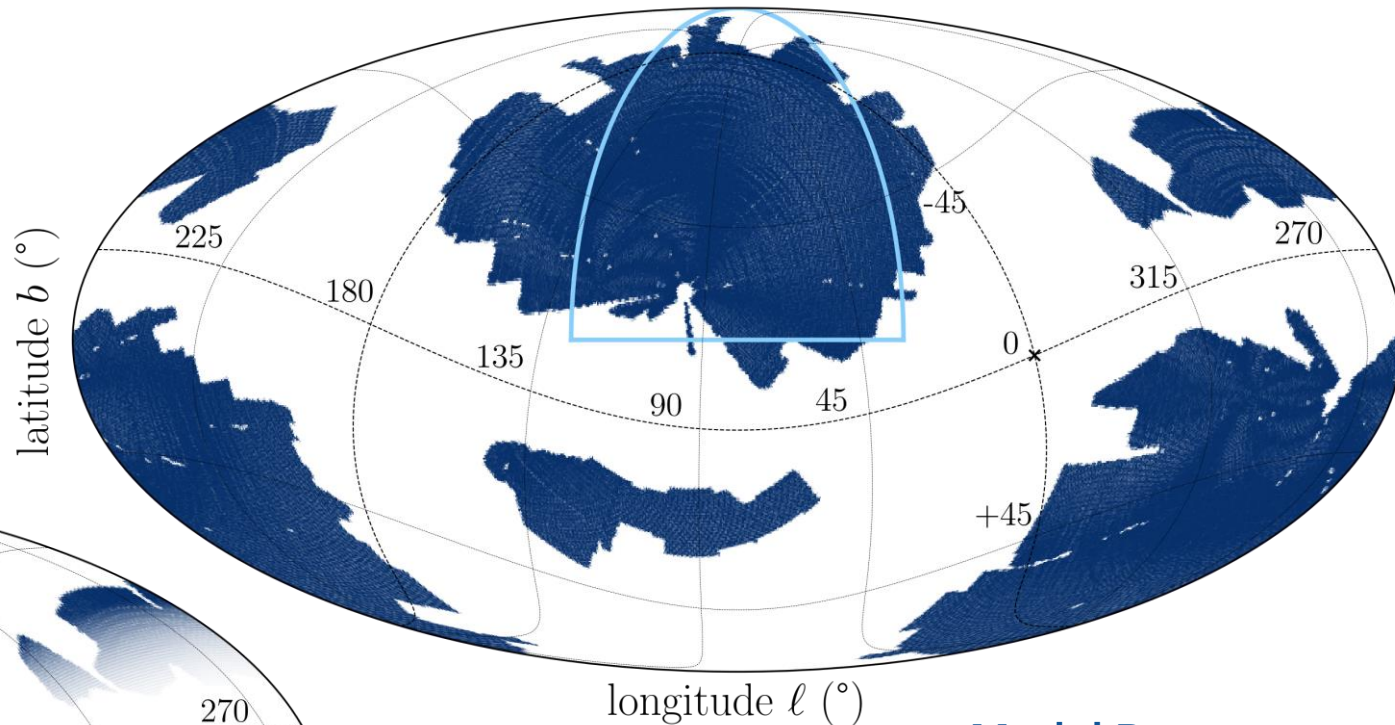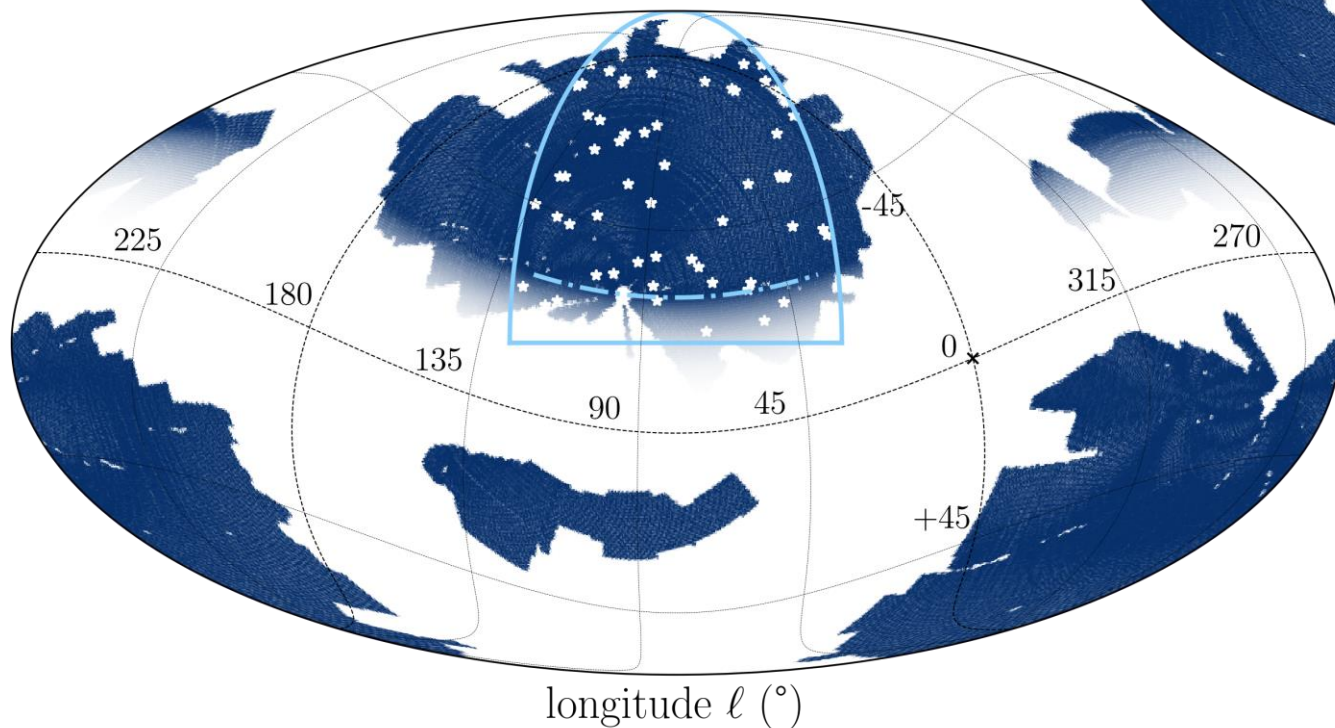


Hoellinger & Leclercq, in prep.

# Systematic effect n°1: survey mask

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



**Model B**
no such effects

**Model A**
80 additional holes
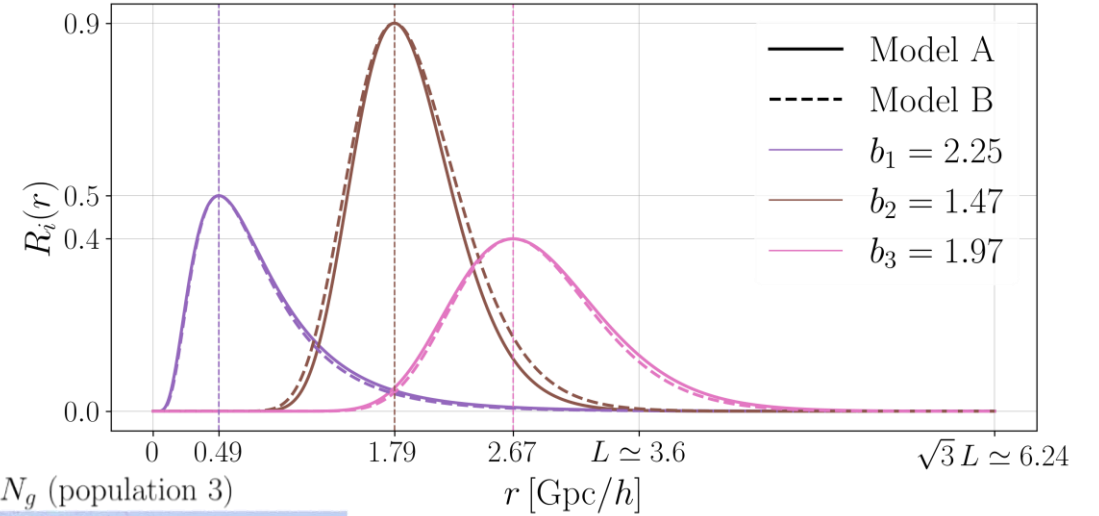extinction from -30° to 0° latitude (galactic)

## Model A

3 simulated populations of galaxies (1 nearby + 2 LRGs) with

- Log-normal selection functions
- Luminosity-dependent galaxy biases



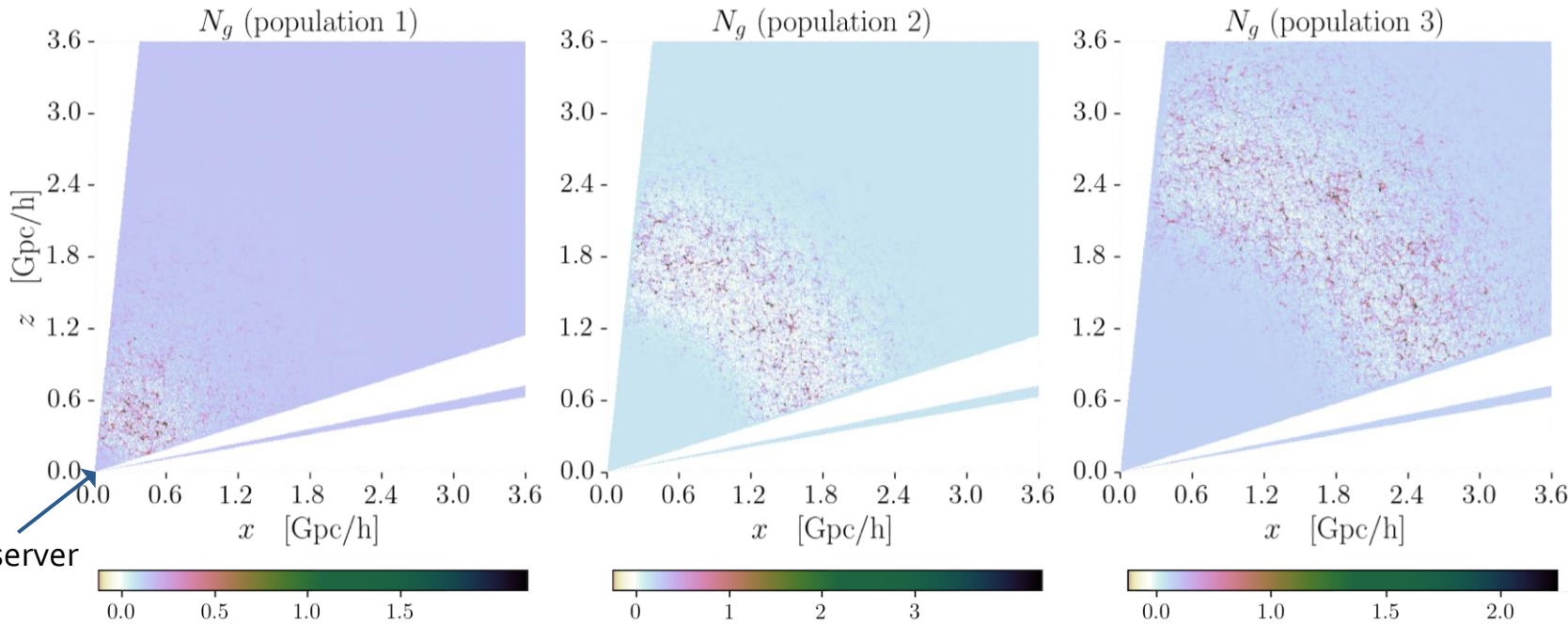Biases based on:
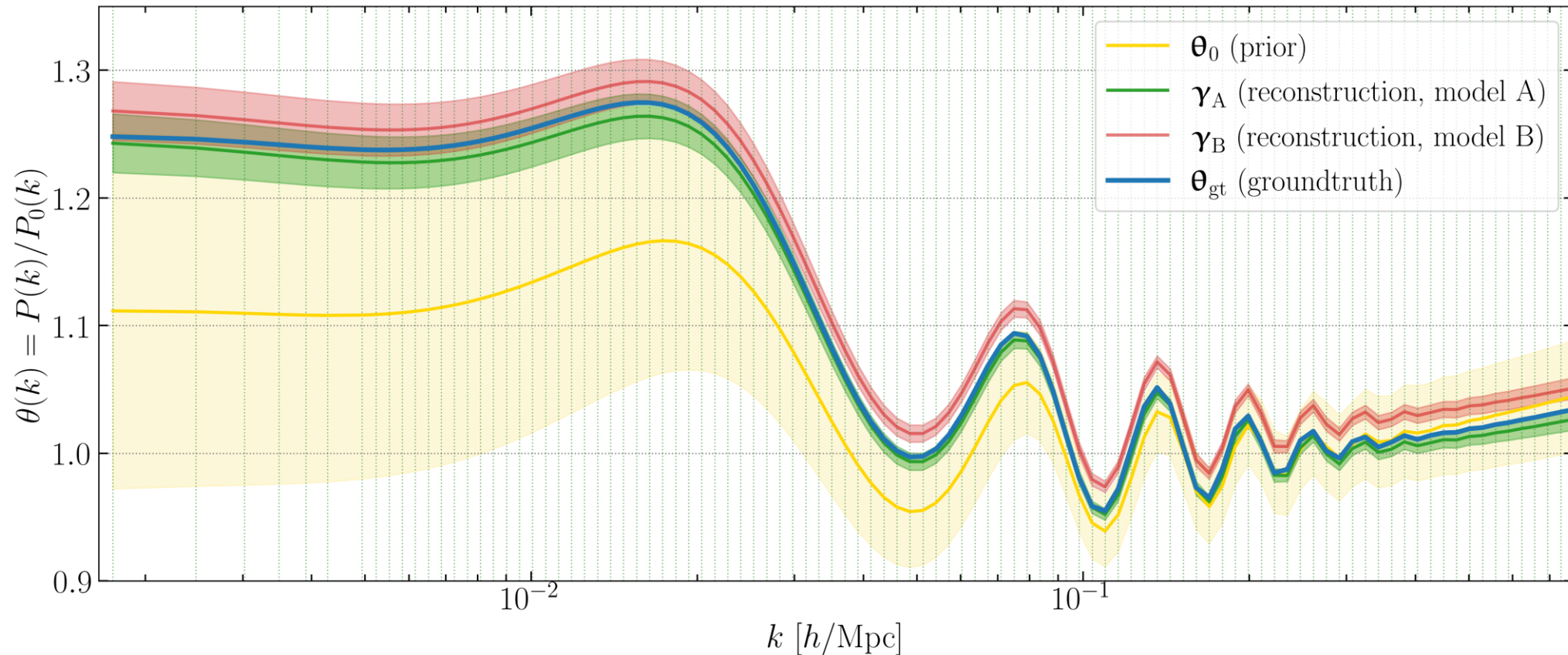Howlett *et al.*, 1409.3238
Gil-Marín *et al.*, 1407.5668

### Model B

- Misspecified selections functions
- Misspecified biases
- Effect sizes $\mathcal{O}(1\%)$



observer

Hoellinger & Leclercq, in prep.

# Check for model misspecification using the SELFI posterior



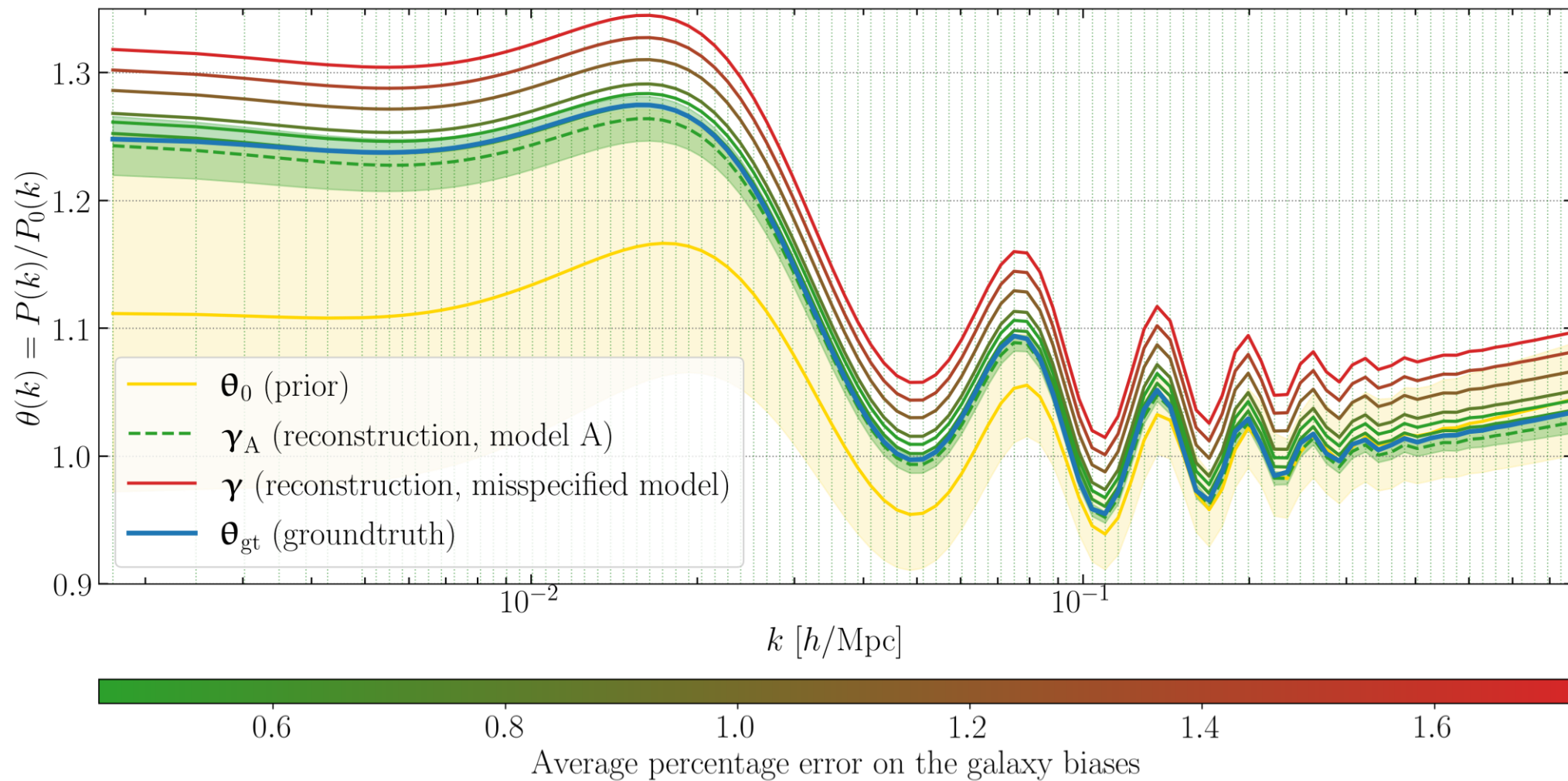Mahalanobis distances to prior:

$$d_{\mathrm{M}}(\boldsymbol{\gamma}, \boldsymbol{\theta}_0 | \mathbf{S}) \equiv \sqrt{(\boldsymbol{\gamma} - \boldsymbol{\theta}_0)^{\mathsf{T}} \mathbf{S}^{-1} (\boldsymbol{\gamma} - \boldsymbol{\theta}_0)}$$

**Model A:** 1.96          **Model B:** 2.91

Hoellinger & Leclercq, in prep.

# Impact of galaxy biases on the posterior initial power spectrum



Figure legend:
- $\boldsymbol{\theta}_0$ (prior)
- $\boldsymbol{\gamma}_A$ (reconstruction, model A)
- $\boldsymbol{\gamma}$ (reconstruction, misspecified model)
- $\boldsymbol{\theta}_{gt}$ (groundtruth)

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

X-axis: $k$ [$h$/Mpc]

Colorbar: Average percentage error on the galaxy biases

Hoellinger & Leclercq, in prep.

- Qualitatively: the shape of the reconstructed initial power spectrum $\theta$ is useful as a [check for unknown systematics](#) / [model misspecification](#) (using our independent theoretical understanding).

- Quantitatively: we can use the Mahalanobis distance between the reconstruction $\gamma$ and the prior distribution $\mathcal{P}(\theta)$:

$$d_{\mathrm{M}}(\gamma, \theta_0 | \mathbf{S}) \equiv \sqrt{(\gamma - \theta_0)^{\mathsf{T}} \mathbf{S}^{-1} (\gamma - \theta_0)}$$

$\mathcal{P}(\omega)$

$\omega$

$\mathcal{T}$

$\theta$

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

$\Phi$

$\mathcal{C}$

$\widetilde{\omega}$

- The score function $\nabla_{\omega}\hat{\ell}_{\omega 0}$ is the gradient of the log-likelihood at fiducial point $\omega_0$ in parameter space.

- A quasi maximum-likelihood estimator for the parameters is

$$\mathcal{C}(\Phi) = \widetilde{\omega} \equiv \omega_0 + \mathbf{F}_0^{-1} \left[ (\nabla_{\omega}\mathbf{f}_0)^{\mathsf{T}} \mathbf{C}_0^{-1} (\Phi - \mathbf{f}_0) \right]$$

Fisher matrix:
$$\mathbf{F}_0 = (\nabla_{\omega}\mathbf{f}_0)^{\mathsf{T}} \mathbf{C}_0^{-1} \nabla_{\omega}\mathbf{f}_0$$
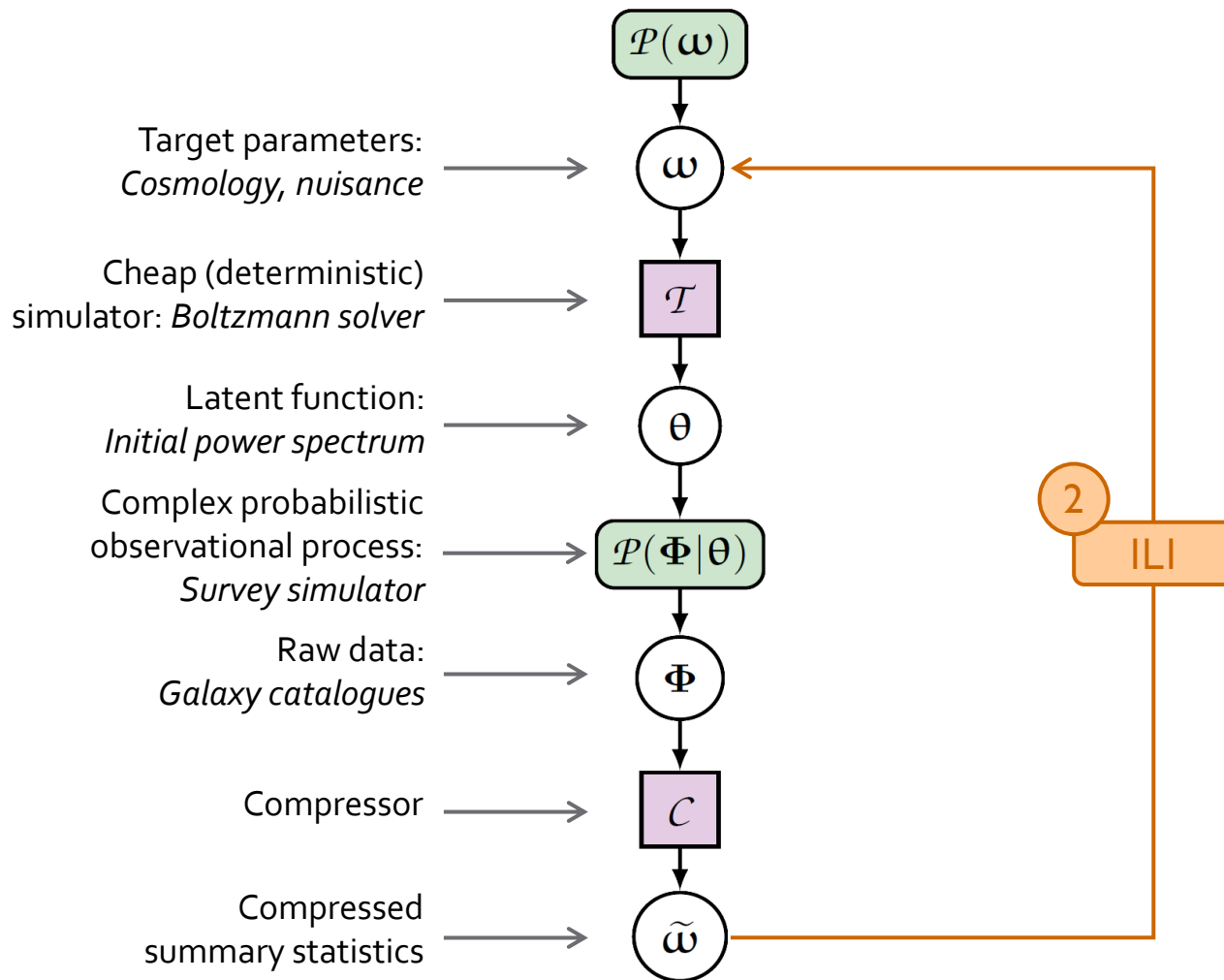$$\nabla_{\omega}\mathbf{f}_0 = \nabla\mathbf{f}_0 \cdot \nabla_{\omega}\mathcal{T}_0$$

Already computed for SELFI   Cheap via finite differences

- Score compression is optimal in the sense that it [preserves the Fisher information content](#) of the data.
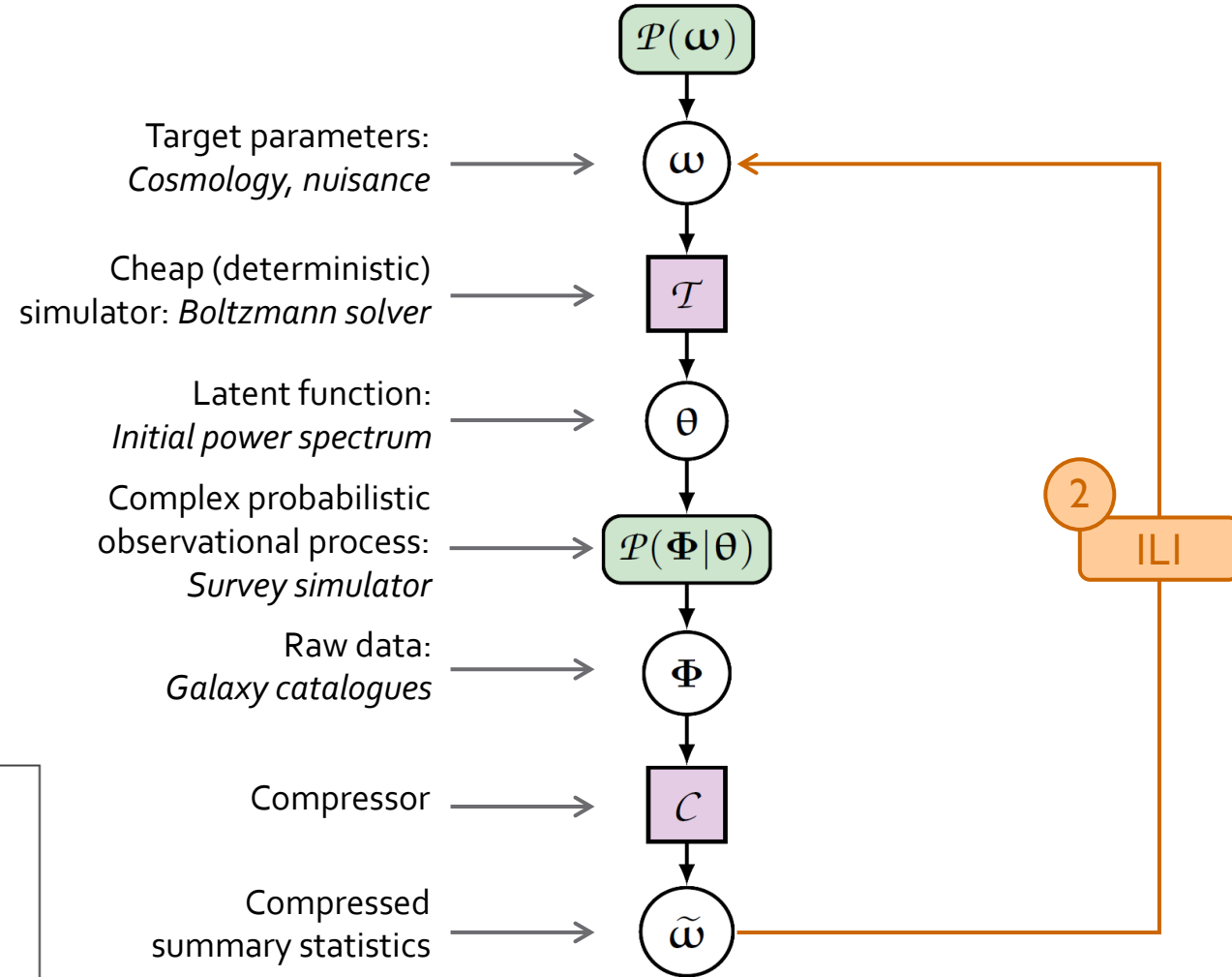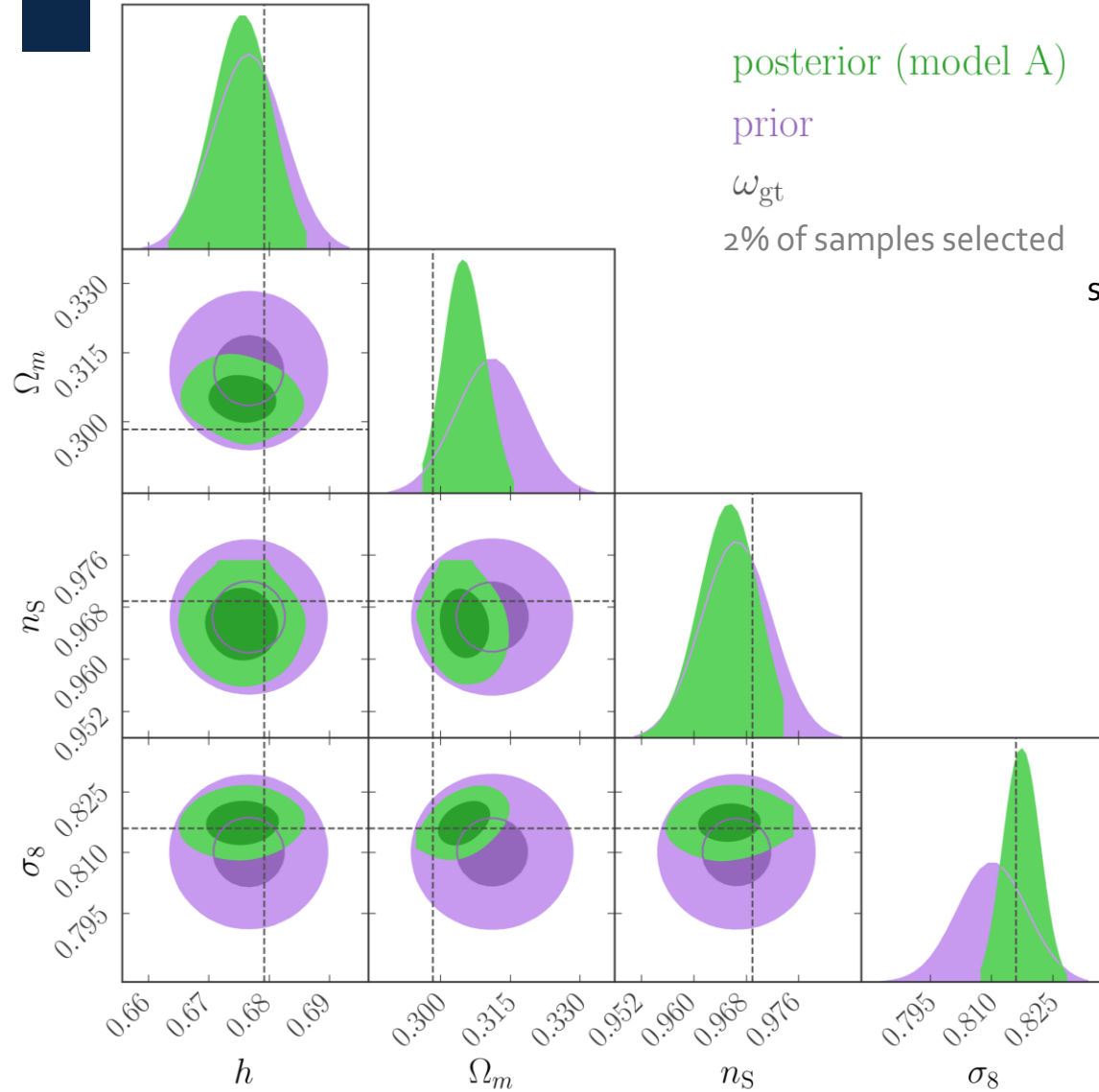
Alsing & Wandelt, 1712.00012

# Implicit likelihood inference of top-level cosmological parameters



Target parameters:
*Cosmology, nuisance*

Cheap (deterministic)
simulator: *Boltzmann solver*

Latent function:
*Initial power spectrum*

Complex probabilistic
observational process:
*Survey simulator*

Raw data:
*Galaxy catalogues*

Compressor

Compressed
summary statistics

- 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 $\widetilde{\omega}_{\mathrm{O}}$:

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

# Posterior on cosmological parameters



posterior (model A)

prior

$\omega_{\mathrm{gt}}$

2% of samples selected

Target parameters:
*Cosmology, nuisance*

Cheap (deterministic)
simulator: *Boltzmann solver*

Latent function:
*Initial power spectrum*

Complex probabilistic
observational process:
*Survey simulator*

Raw data:
*Galaxy catalogues*

Compressor

Compressed
summary statistics

$\mathcal{P}(\omega)$

$\omega$

$\mathcal{T}$

$\theta$

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

$\Phi$

$\mathcal{C}$

$\widetilde{\omega}$

2

ILI

Hoellinger & Leclercq, in prep.

# Posterior on cosmological parameters



posterior (model B)
prior
$\omega_{\mathrm{gt}}$
2% of samples selected

Target parameters:
*Cosmology, nuisance*

Cheap (deterministic)
simulator: *Boltzmann solver*

Latent function:
*Initial power spectrum*

Complex probabilistic
observational process:
*Survey simulator*

Raw data:
*Galaxy catalogues*

Compressor

Compressed
summary statistics

$\mathcal{P}(\omega)$
$\omega$
$\mathcal{T}$
$\theta$
$\mathcal{P}(\Phi|\theta)$
$\Phi$
$\mathcal{C}$
$\widetilde{\omega}$

2
ILI

Hoellinger & Leclercq, in prep.

# Conclusion: the statistical framework is in place for the GC:AP pipeline

- A novel two-step simulation based Bayesian approach, combining SELFI and ILI, 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):
  - 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.

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

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