

Galaxy Clustering with Likelihood-Free Inference

Prospects and Forecasts for Euclid



Florent Leclercq

www.florent-leclercq.eu

Imperial Centre for Inference and Cosmology
Imperial College London

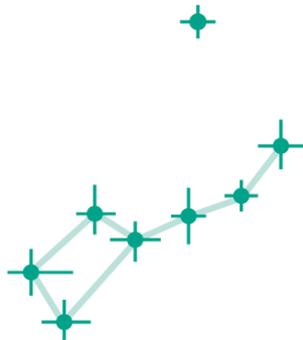
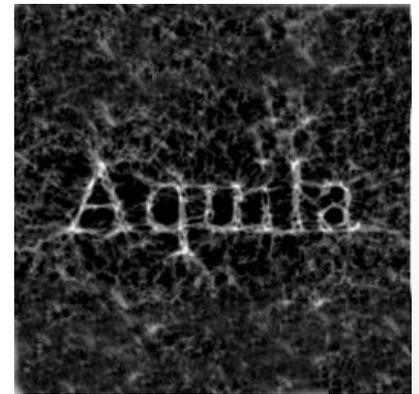


Wolfgang Enzi, Baptiste Faure, Alan Heavens,
Andrew Jaffe, Jens Jasche, Guilhem Lavaux,
Will Percival, Benjamin Wandelt

and the Aquila Consortium

www.aquila-consortium.org

20 April 2020



ICIC

Imperial Centre
for Inference & Cosmology

**Imperial College
London**

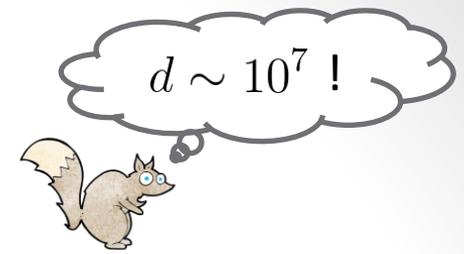
Forward modelling within the Euclid Consortium

- Galaxy Clustering Science Working Group (**GC SWG**)
 - WP Likelihood + IST:L:
 - “Methods to speed up the GC likelihood computation” (includes BOLFI)
 - WP Additional Probes:
 - “Density reconstruction via Bayesian large-scale structure inference” (ARES, HADES, BORG)
 - “Primordial power spectrum from black-box galaxy surveys” (SELF)
- Weak Lensing Science Working Group (**WL SWG**):
 - WP Forward-modelling (DELFI, BOLFI, SELF)
- Theory Science Working Group (**TH SWG**):
 - WP Initial conditions (BORG with f_{NL})
- Cosmological Simulations Science Working Group (**SIM SWG**)
 - WP Machine Learning (emulators, neural networks)

-
- Galaxy clustering with likelihood-free inference (GCLFI) will be proposed as a SP.

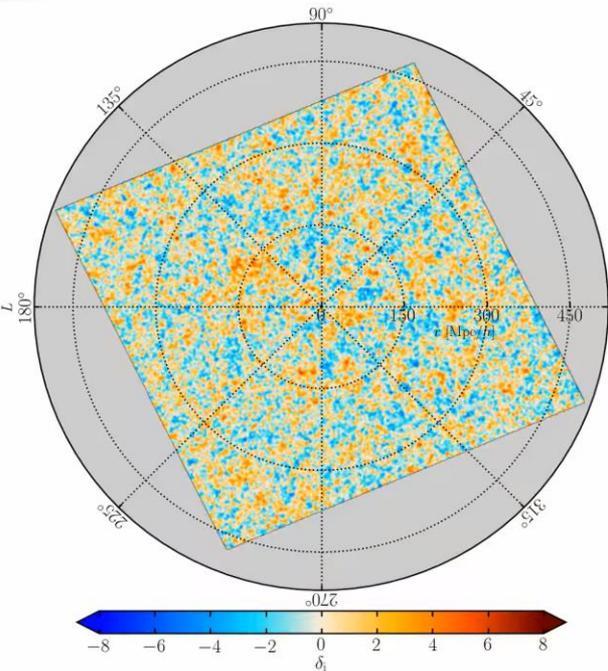
Vocabulary consideration:

What is the likelihood?

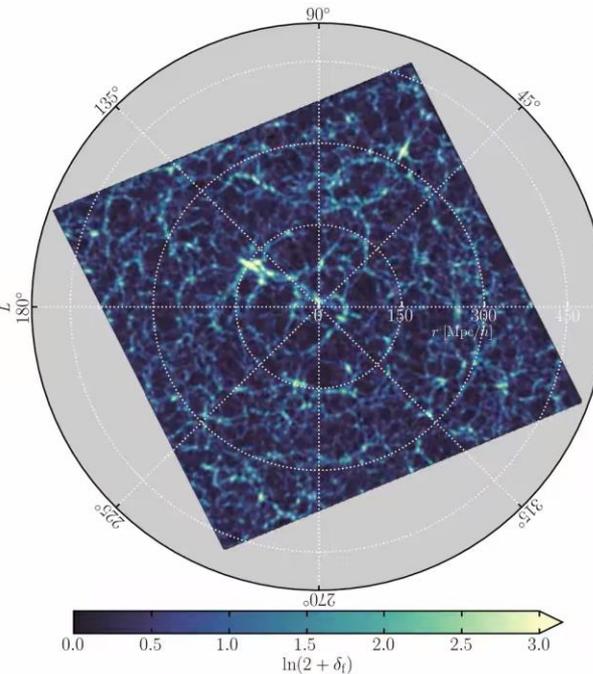


In cosmology, the (true?) likelihood should live at the level of the **map** of the CMB or LSS. e.g. Wiener filtering for the CMB, BORG for the LSS (a 256^3 -dimensional Poisson likelihood):

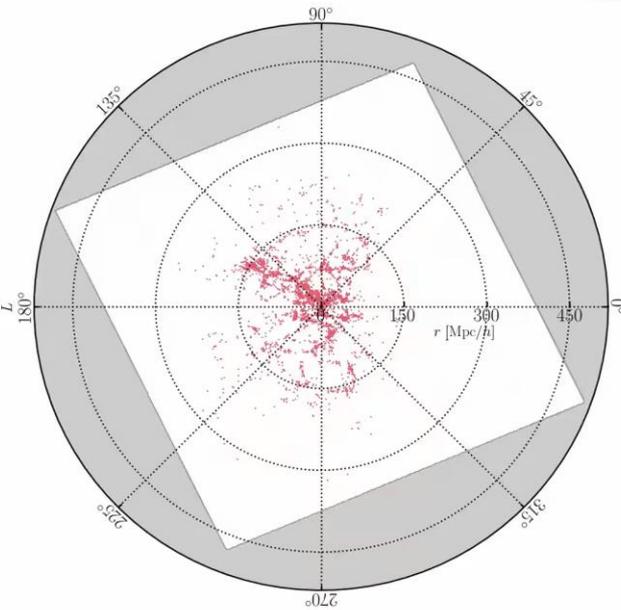
Initial conditions



Final conditions



Observations

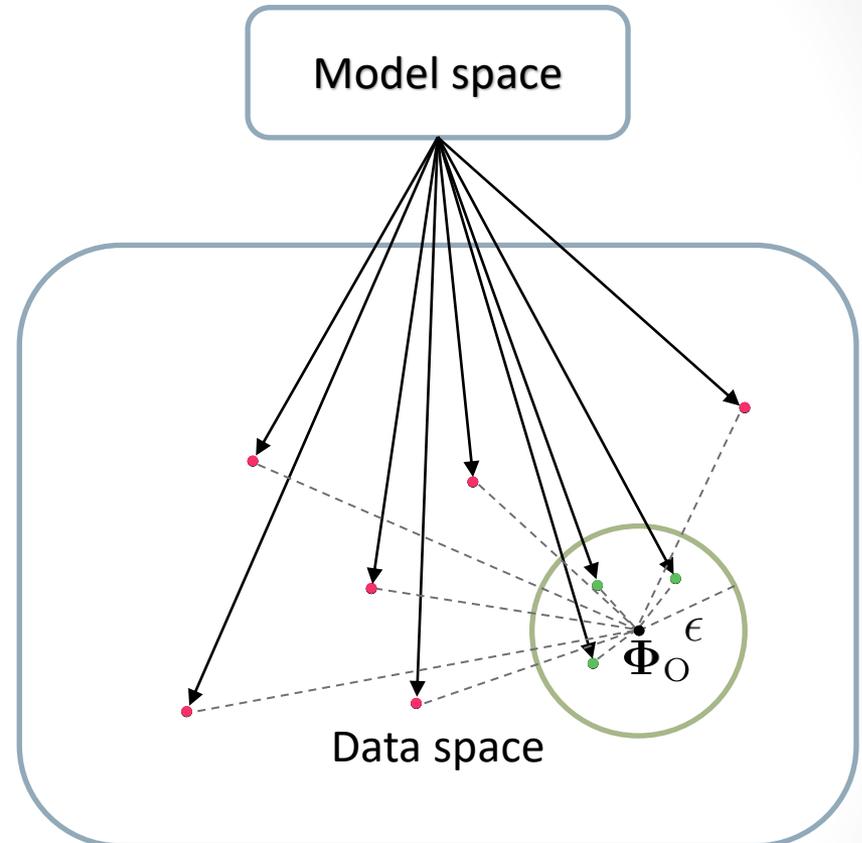


Jasche & Lavaux 2019, 1806.11117 – FL, Lavaux & Jasche, in prep.

Expert knowledge of the likelihood is needed to beat the curse of dimensionality: conditionals/gradients of the likelihood are required by the samplers (Gibbs/Hamiltonian).

Likelihood-free rejection sampling (LFRS)

- Iterate many times:
 - Sample θ from a proposal distribution $q(\theta)$
 - Simulate Φ_θ using the black-box
 - Compute the distance $\Delta(\Phi_\theta, \Phi_O)$ between simulated and observed data
 - Retain θ if $\Delta(\Phi_\theta, \Phi_O) \leq \epsilon$, otherwise reject



ϵ can be adaptively reduced
(Population Monte Carlo)

Beyond LFRS: two scenarios

The “number of simulations” route:

- Specific cosmological models ($d \lesssim 10$), general exploration of parameter space
- Density Estimation for Likelihood-Free Inference (DELFI)
Papamakarios & Murray 2016, 1605.06376
Alsing, Feeney & Wandelt 2018, 1801.01497
- Bayesian Optimisation for Likelihood-Free Inference (BOLFI)
Gutmann & Corander 2016, 1501.03291
FL 2018, 1805.07152

The “number of parameters” route:

- Model-independent theoretical parametrisation ($d \gtrsim 100$), strong existing constraints in parameter space
- Simulator Expansion for Likelihood-Free Inference (SELI)
FL, Enzi, Jasche & Heavens 2019, 1902.10149

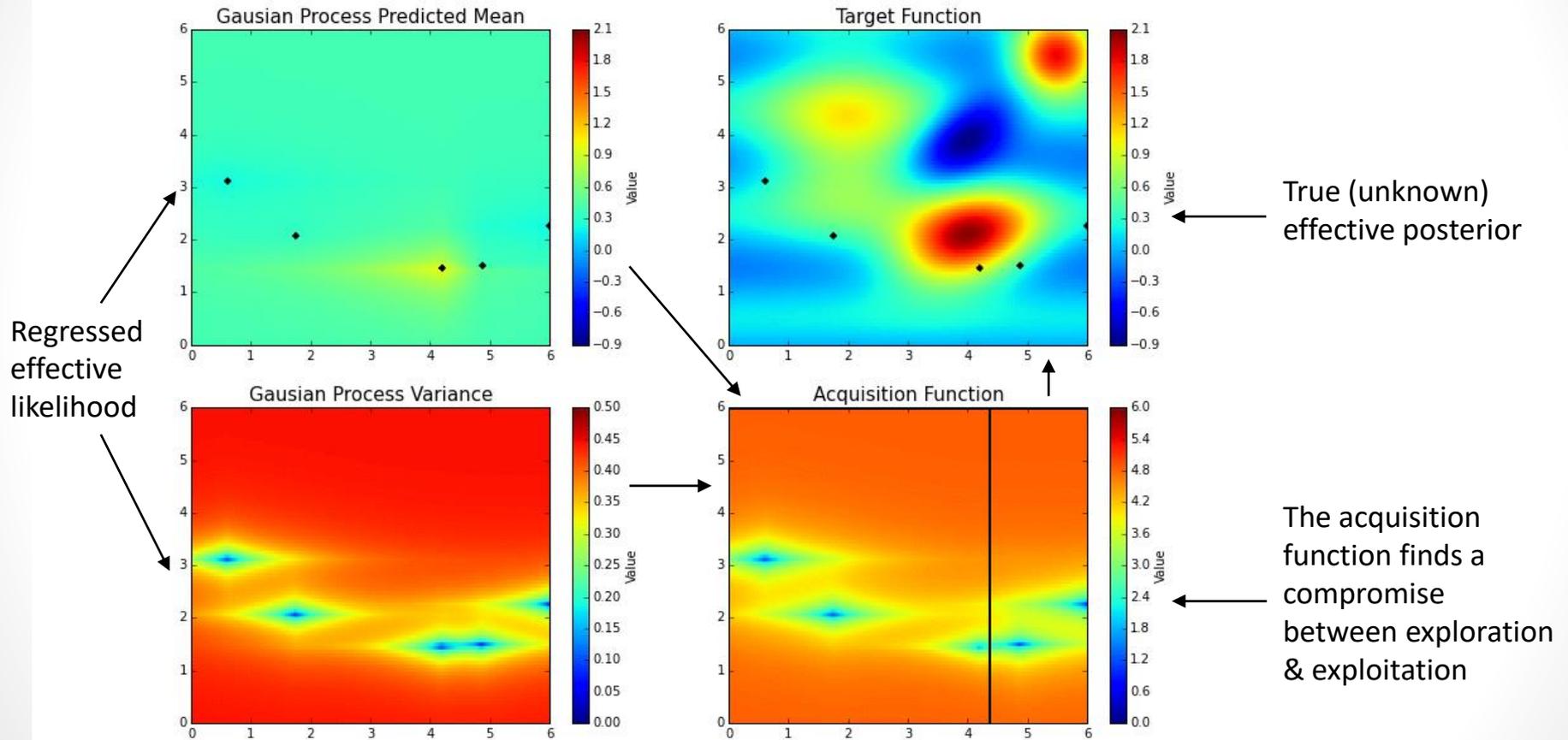
The “number of simulations” route: BOLFI

Bayesian Optimisation for Likelihood-Free Inference

BOLFI: Data acquisition

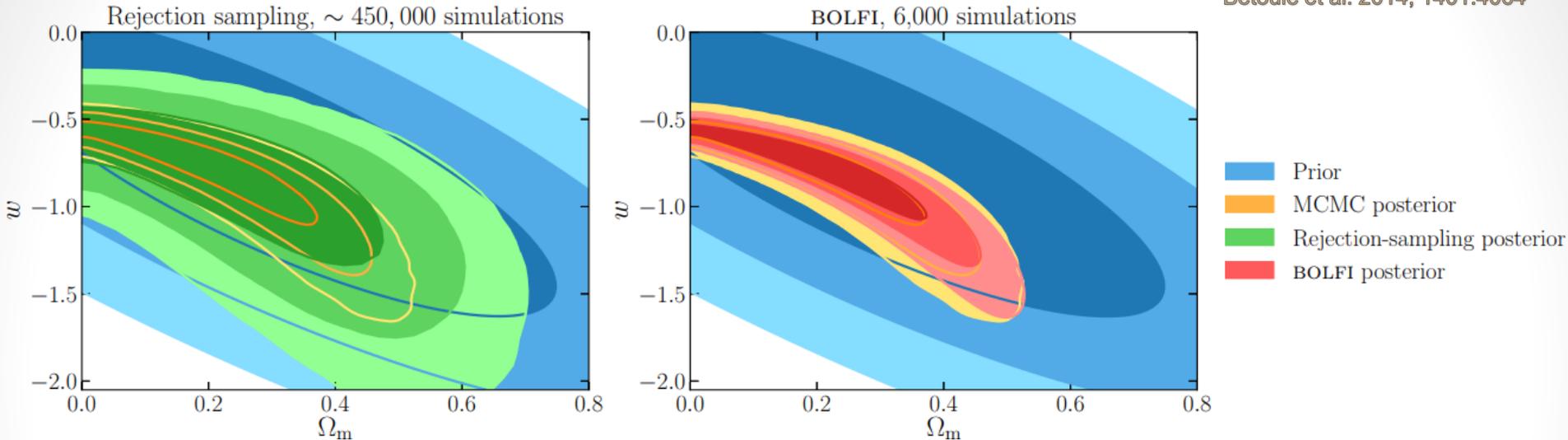
Simulations are obtained from sampling an **adaptively-constructed proposal distribution**, using the regressed effective likelihood

Bayesian Optimization in Action



BOLFI: Re-analysis of the JLA supernova sample

Betoule et al. 2014, 1401.4064



- The number of required simulations is reduced by:
 - 2 orders of magnitude with respect to likelihood-free rejection sampling (for a much better approximation of the posterior)
 - 3 orders of magnitude with respect to exact Markov Chain Monte Carlo sampling

FL 2018, 1805.07152

- Bayesian optimisation can also be applied to the “true” likelihood (if known) or to build an emulator of the data model: see derivative work by

Rogers et al. 2019, 1812.04631 – Takhtaganov et al. 2019, 1905.07410 – Pellejero-Ibañez et al. 2019, 1912.08806 and in the WP:Lik

Standard acquisition functions are suboptimal (at best)

- Goal for Bayesian optimisation: finding the optimum (assumed unique) of a function (\neq exploring a parameter space)
- Examples of acquisition functions :

- The **Expected Improvement**

Gaussian cdf Gaussian pdf

$$\text{EI}(\boldsymbol{\theta}_*) \equiv \underbrace{\sigma(\boldsymbol{\theta}_*)}_{\text{Exploration}} \underbrace{[z\Phi(z) + \phi(z)]}_{\text{Exploitation}} \quad z \equiv \frac{\min(\mathbf{f}) - \mu(\boldsymbol{\theta}_*)}{\sigma(\boldsymbol{\theta}_*)}$$

e.g. Brochu, Cora & de Freitas 2010, 1012.2599 – Takhtaganov *et al.* 2019, 1905.07410

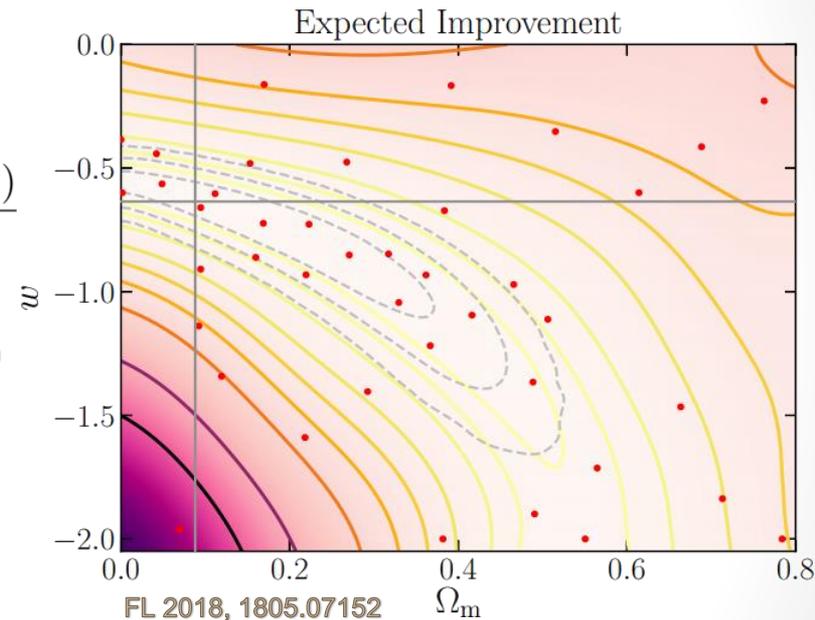
- The **Upper Confidence Bound**

$$\text{UCB}(\boldsymbol{\theta}_*) = \mu(\boldsymbol{\theta}_*) + \alpha\sigma(\boldsymbol{\theta}_*)$$

Rogers *et al.* 2019, 1812.04631 – Pellejero-Ibañez *et al.* 2019, 1912.08806

- Drawbacks of these acquisition rules:
 - Do not take into account prior information
 - Local evaluation rules
 - Too greedy for parameter inference

➔ In most cases these standard rules will lead to a suboptimal exploration of parameter space. But in some situations (e.g. posteriors with non-trivial degeneracies or multiple modes) they can even yield a biased inference.



The optimal acquisition function for parameter inference

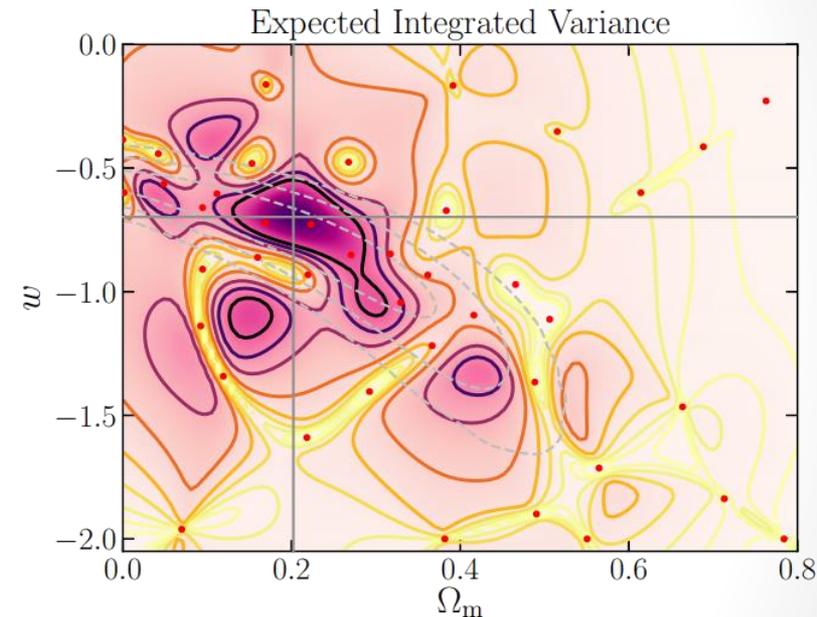
- Goal for parameter inference: minimise the expected uncertainty in the estimate of the (approximate) posterior over the future evaluation of the simulator
- The optimal acquisition function : the **Expected Integrated Variance**

$$\text{EIV}(\theta_*) = \int \frac{\mathcal{P}(\theta)^2}{4} \exp[-\mu(\theta)] [\sigma^2(\theta) - \tau^2(\theta, \theta_*)] d\theta$$

Integral
Prior
Exploitation
Exploration

$$\tau^2(\theta, \theta_*) \equiv \frac{\text{cov}^2(\theta, \theta_*)}{\sigma^2(\theta_*)}$$

- Advantages:
 - Takes into account the prior
 - Non-local (integral over parameter space): more expensive... but much more informative
 - Exploration of the posterior tails is favoured when necessary
 - Analytic gradient



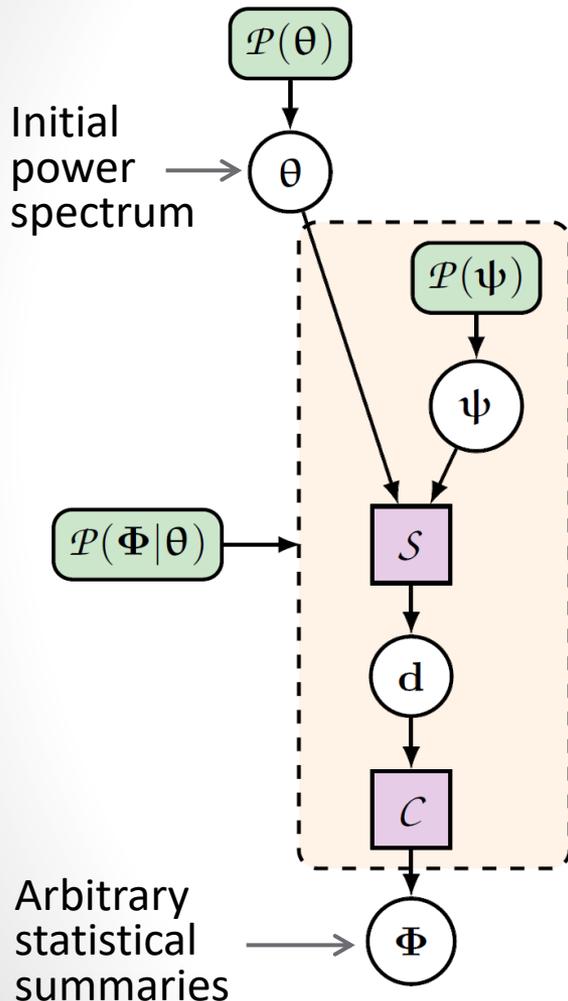
Järvenpää et al. 2017, 1704.00520 (expression of the EIV in the non-parametric approach)

FL 2018, 1805.07152 (expression of the EIV in the parametric approach)

The “number of parameters” route: SELFI

Simulator Expansion for Likelihood-Free Inference

SELEFI: expansion of black-box data models

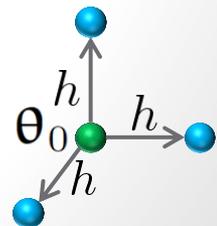


- We aim at inferring the initial power spectrum, which contains (almost?) all of the information
- This requires doing LFI in $d = \mathcal{O}(100) - \mathcal{O}(1,000)$
- If we trust the results of earlier experiments, we can Taylor-expand the black-box around an expansion point θ_0 :

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

SELEFI-2 (second-order): coming soon!

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



SEIFI-1: linearization of the black-box

- Linearization of the black-box:

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

- Gaussian prior + Gaussian effective likelihood

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

expansion point

observed summaries

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

covariance of summaries

gradient of the black-box

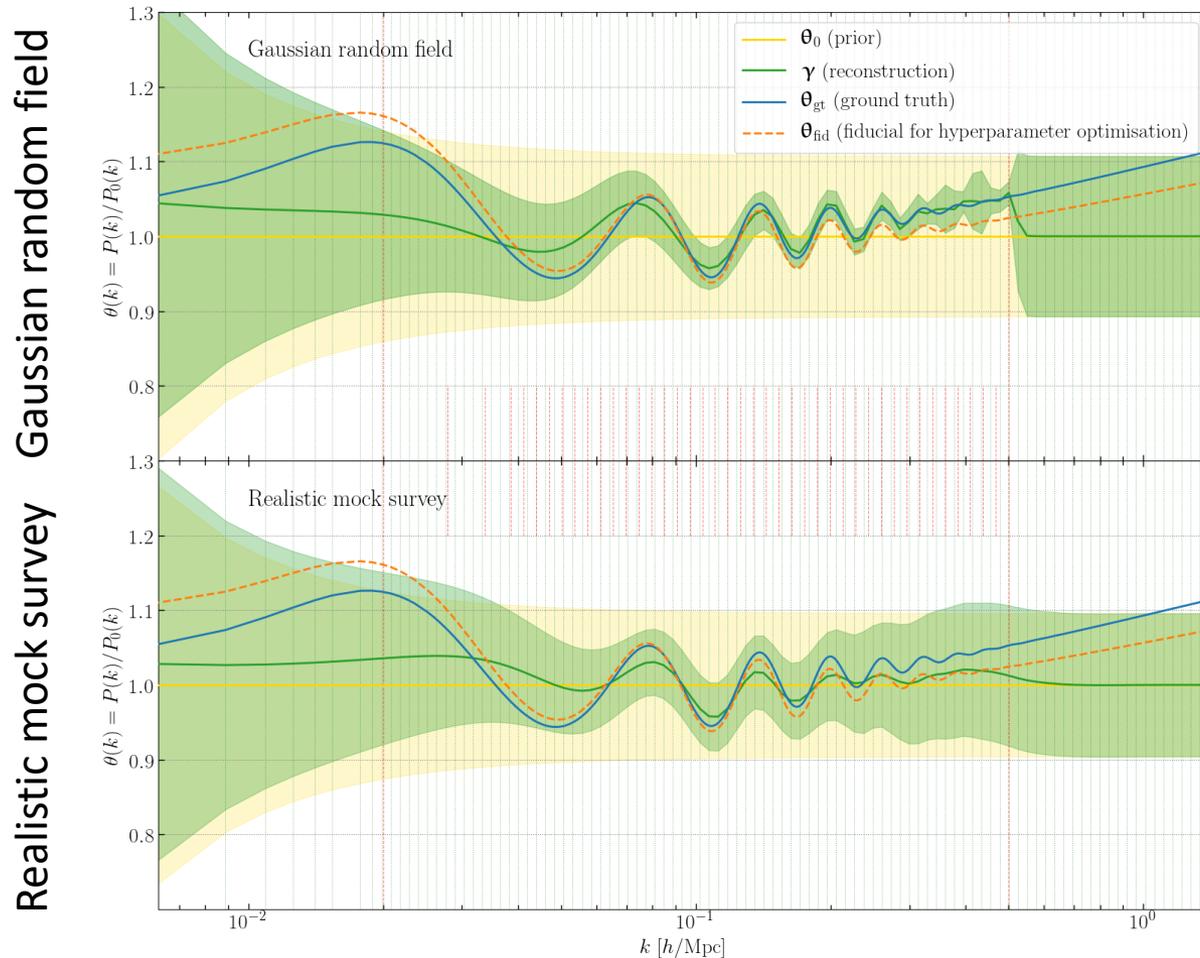
prior covariance

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

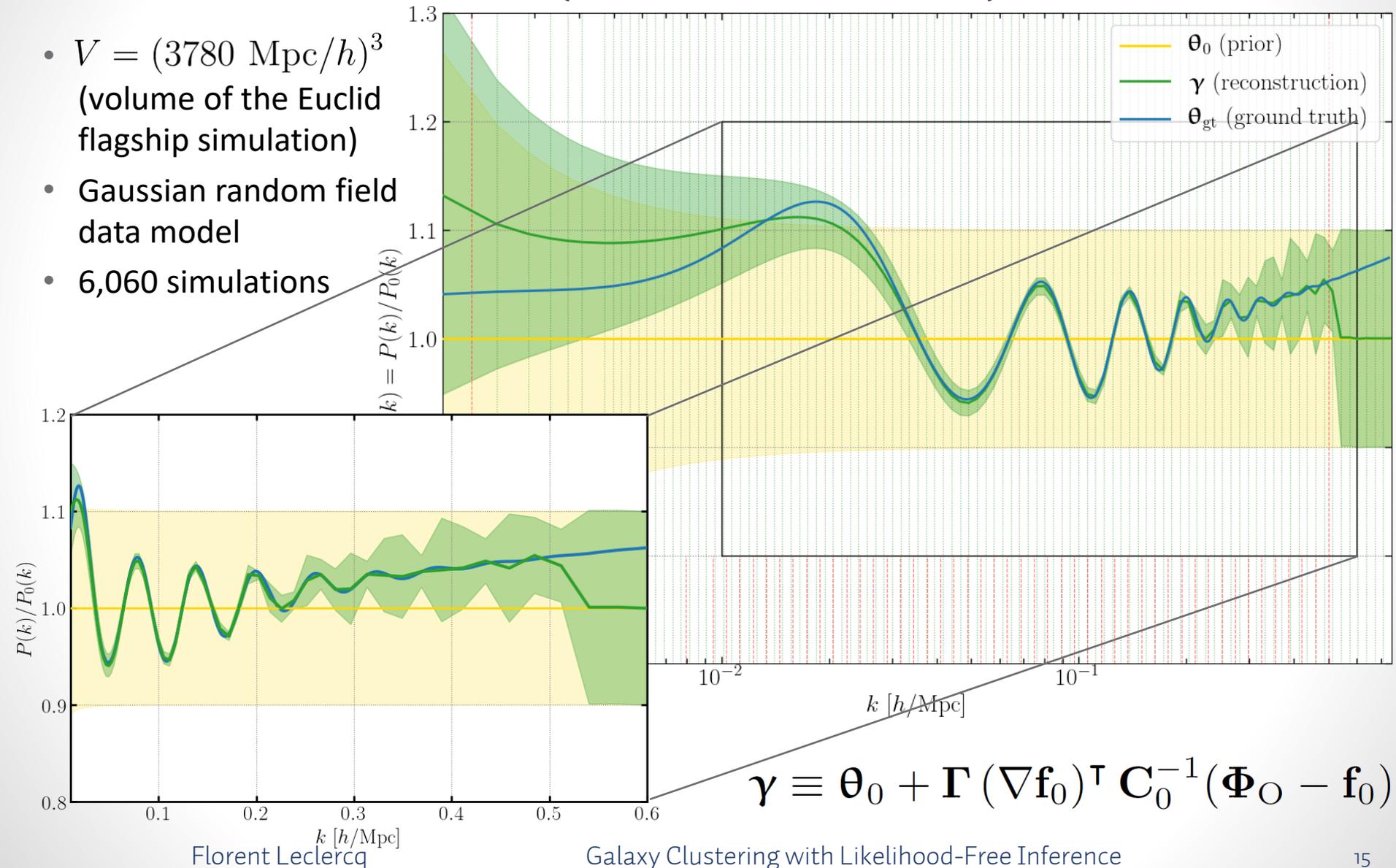
SELFI + numerical model: Proof-of-concept



100 parameters are simultaneously inferred from a black-box data model
1 (Gpc/h)³ only! Much more potential for upcoming data...

SEIFI-1 Euclid forecast (cosmic variance limit)

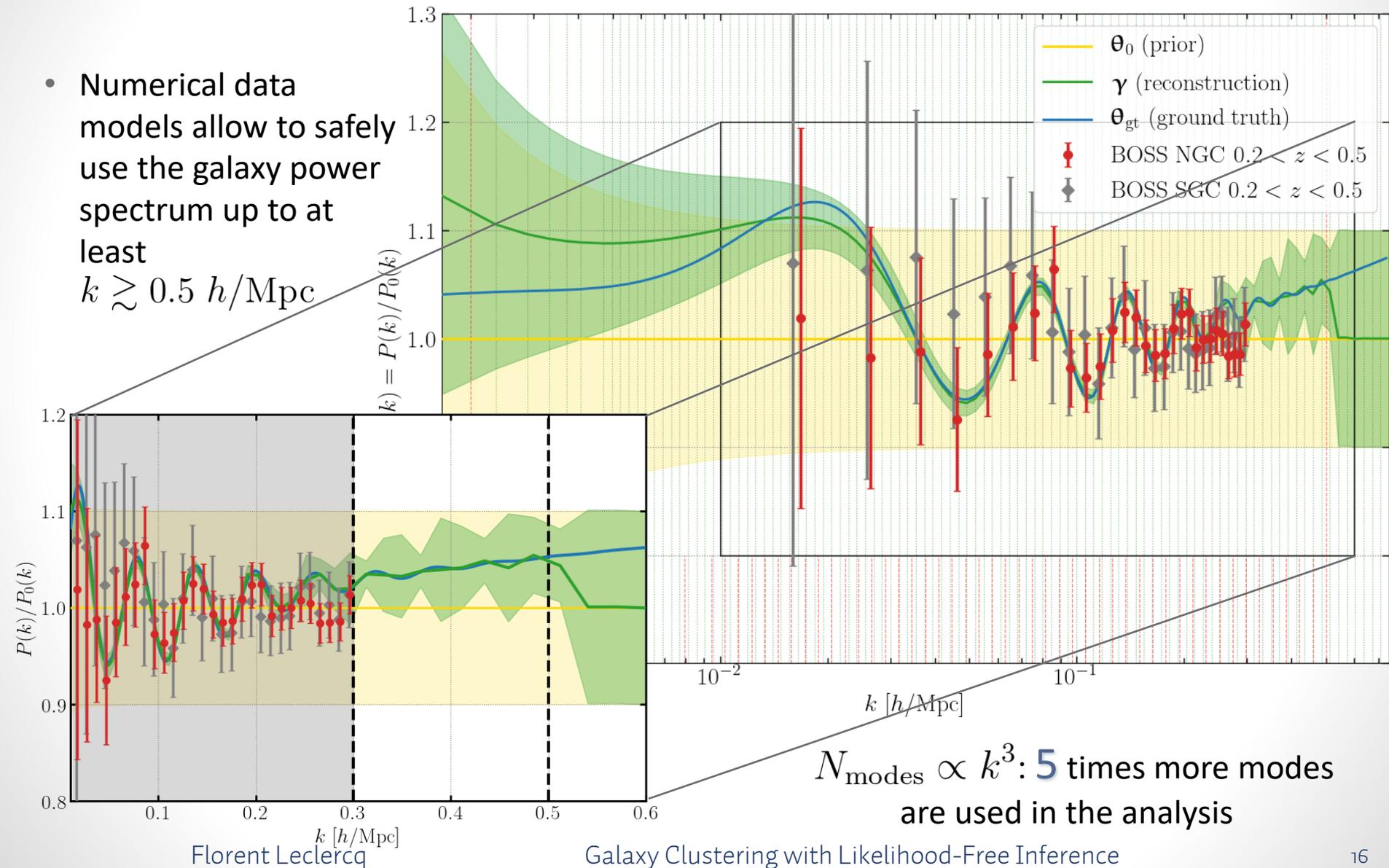
- $V = (3780 \text{ Mpc}/h)^3$
(volume of the Euclid flagship simulation)
- Gaussian random field data model
- 6,060 simulations



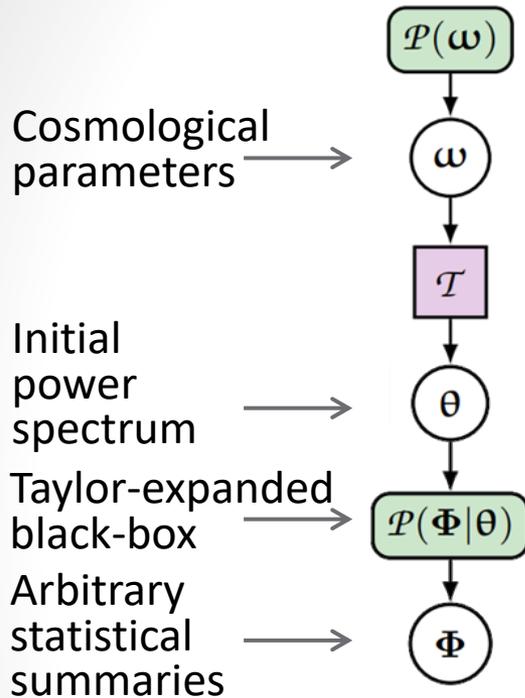
SEIFI-1 Euclid versus BOSS

Data points from
Beutler *et al.* 2016, 1607.03149

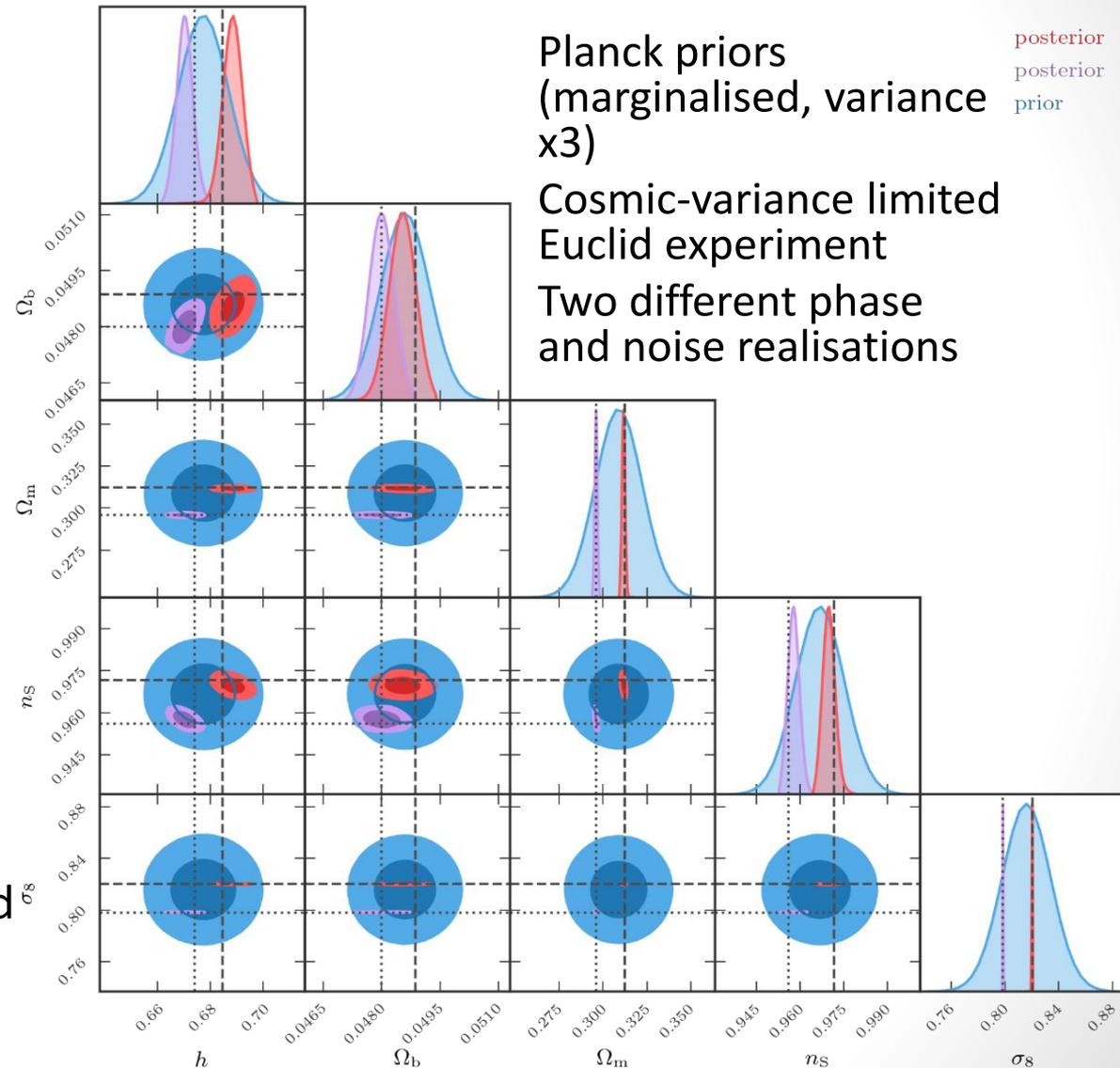
- Numerical data models allow to safely use the galaxy power spectrum up to at least $k \gtrsim 0.5 h/\text{Mpc}$



From initial power spectrum to cosmology

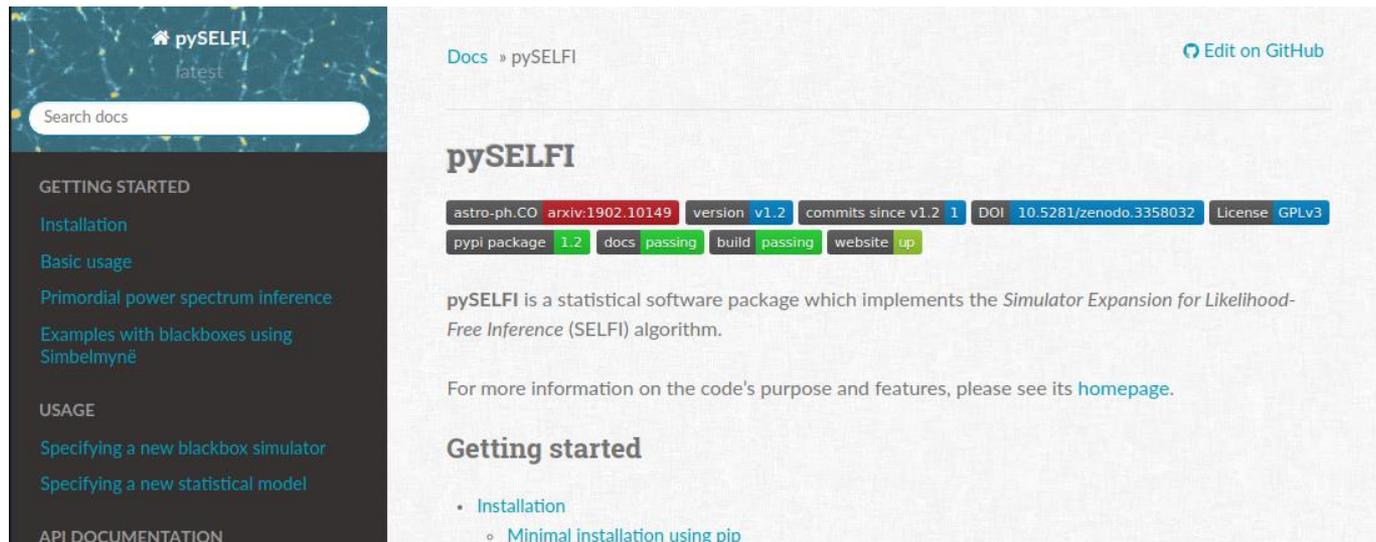


- Robust inference of cosmological parameters can be easily performed *a posteriori* once the linearized data model is learnt



pySELFi is publicly available

- Code homepage: <http://pyselfi.florent-leclercq.eu/>
- Source on GitHub: <https://github.com/florent-leclercq/pyselfi/>
- Documentation on ReadtheDocs: <https://pyselfi.readthedocs.io/en/latest/>
(with templates to use your own black-box)



The screenshot shows the ReadTheDocs interface for pySELFi. On the left is a dark sidebar with a search bar and navigation links under 'GETTING STARTED' (Installation, Basic usage, Primordial power spectrum inference, Examples with blackboxes using Simbelmyne) and 'USAGE' (Specifying a new blackbox simulator, Specifying a new statistical model). The main content area has a light background with a 'Docs » pySELFi' breadcrumb and an 'Edit on GitHub' link. The title 'pySELFi' is followed by a row of badges: astro-ph.CO, arxiv:1902.10149, version v1.2, commits since v1.2 1, DOI 10.5281/zenodo.3358032, License GPLv3, pypi package 1.2, docs passing, build passing, and website up. Below this is a paragraph describing pySELFi as a statistical software package implementing the Simulator Expansion for Likelihood-Free Inference (SELFi) algorithm, with a link to its homepage. A 'Getting started' section follows, containing a link to 'Installation' and a sub-link for 'Minimal installation using pip'.

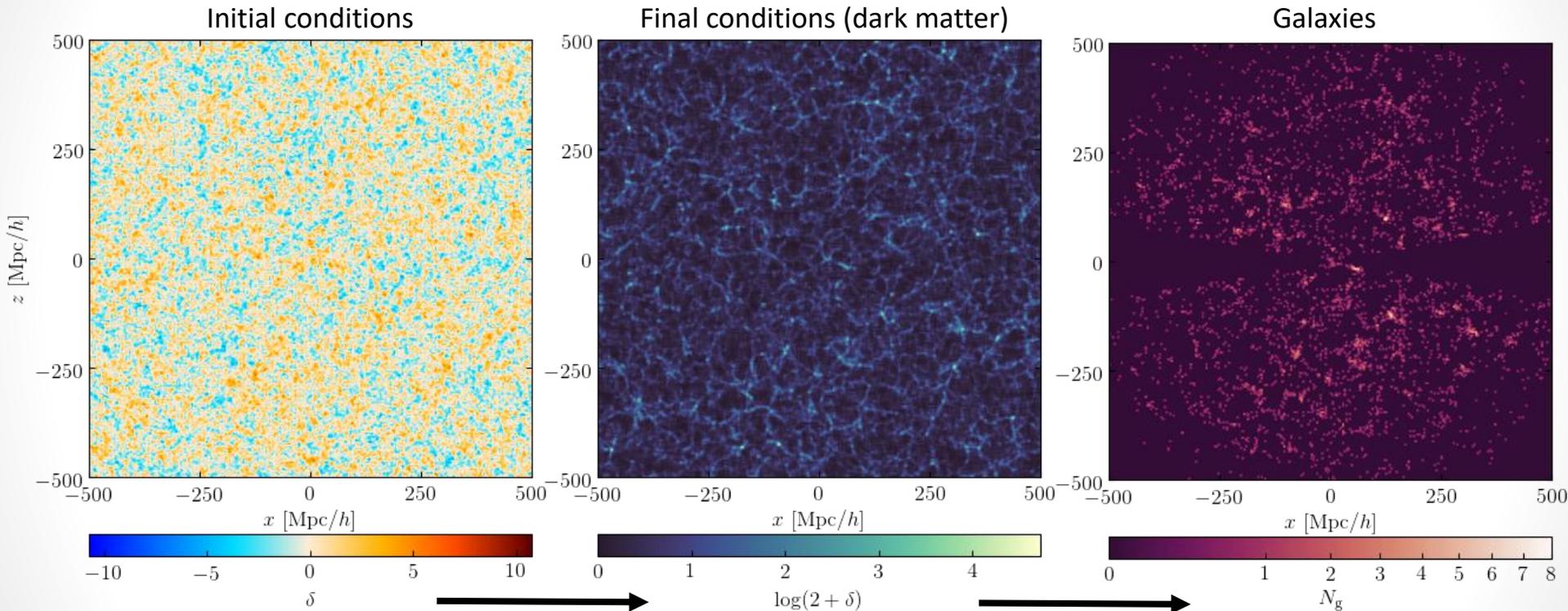
```
pip install pyselfi
```

Numerical data models: Galaxy Clustering and beyond

A black-box: Simbelmynë

Publicly available code:

<https://bitbucket.org/florent-leclercq/simbelmyne/>



Dark matter simulation
with PM/tCOLA/sCOLA

Tassev, Zaldarriaga & Eisenstein 2013, 1301.0322

Tassev, Eisenstein, Wandelt & Zaldarriaga 2015, 1502.07751

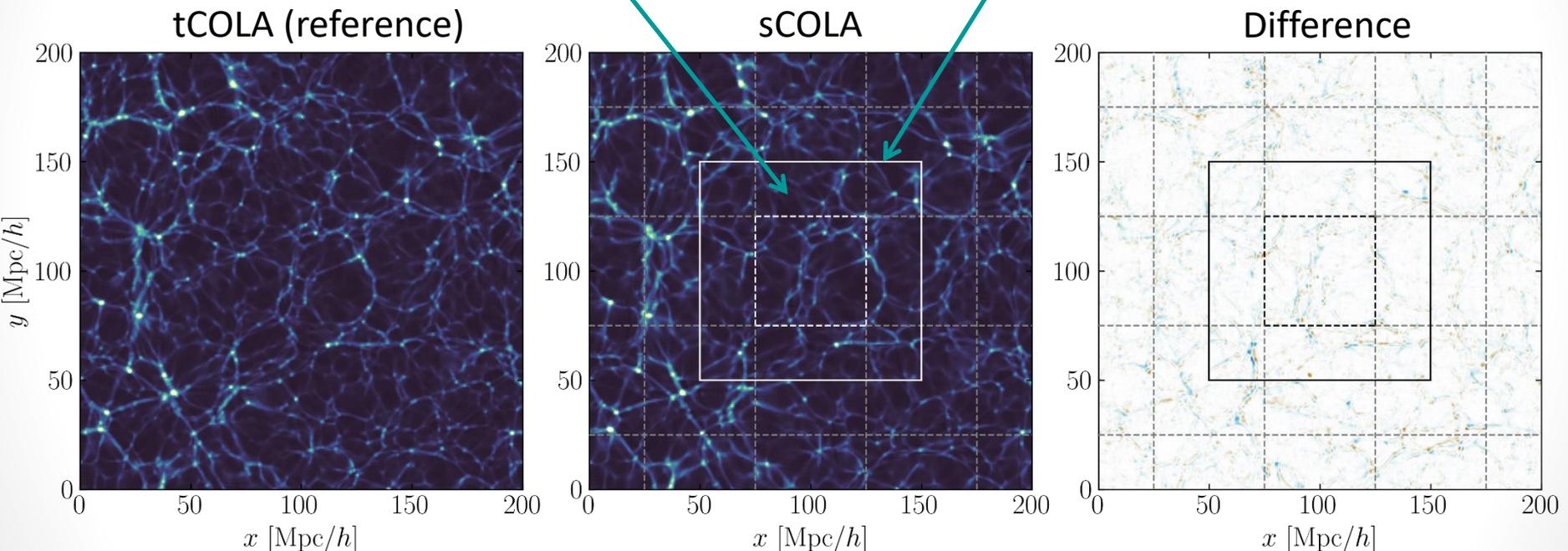
Survey simulation:
Redshift-space distortions, galaxy
bias, selection effects, survey
geometry, instrumental noise

Perfectly parallel cosmological simulations using sCOLA

Can we decouple sub-volumes by using the large-scale analytical solution?

1. A buffer region around each tile

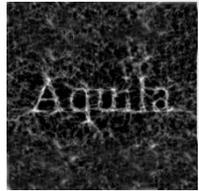
2. Appropriate Dirichlet boundary conditions for the potential



Our new perfectly parallel algorithm unlocks **profoundly new possibilities of computing larger and higher-resolution cosmological simulations**, taking advantage of a variety of hardware architectures (e.g. Cosmology@Home).

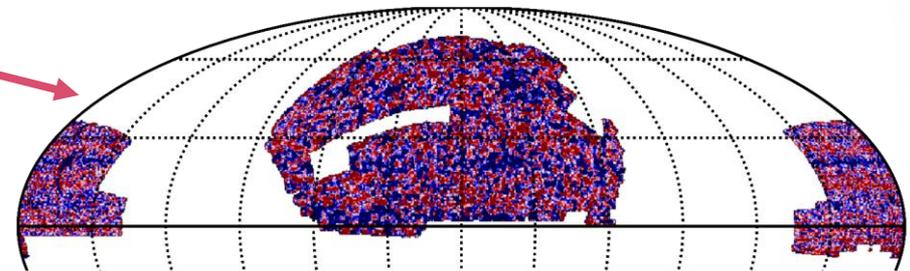
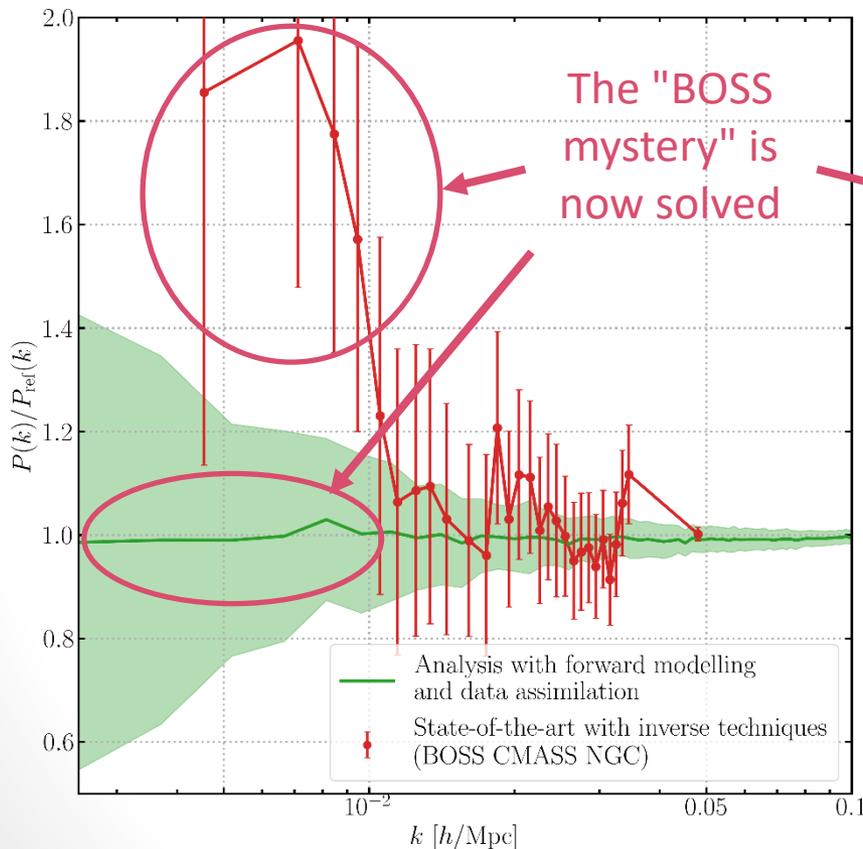
FL, Faure, Lavaux, Wandelt, Jaffe, Heavens, Percival & Noûs, 2023.04925

Forward modelling of known and unknown systematics



- Effects currently accounted for in our numerical data models:

Non-local galaxy biases, redshift-space distortions, light-cone effects, survey geometry, selection effects, foregrounds (including stars, dust, atmosphere, and unknown foregrounds)



Residual unexplained foreground in BOSS data

Summary and concluding thoughts

- **Goal:** developing and using algorithms for targeted science questions, allowing the use of simulators including **all relevant physical and observational effects**.
- Bayesian analyses of galaxy surveys with fully **non-linear numerical black-box models** is not an impossible task!
- BOLFI allows inference within **specific cosmological models** with a very limited simulation budget. The **optimal acquisition function** shall be used.
- SELFI allows inference of the **initial power spectrum** and **cosmological parameters**.
- Our numerical data models are being refined and optimised to prepare for upcoming data, including Euclid.