

Bayesian inference with black-box cosmological models



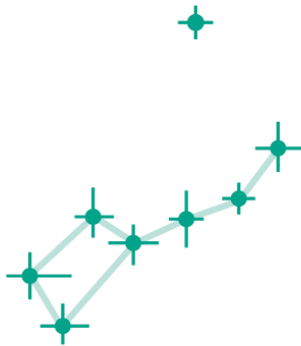
Florent Leclercq

www.florent-leclercq.eu

Imperial Centre for Inference and Cosmology
Imperial College London

Wolfgang Enzi, Jens Jasche, Alan Heavens,
and the Aquila Consortium
www.aquila-consortium.org

September 20th, 2019



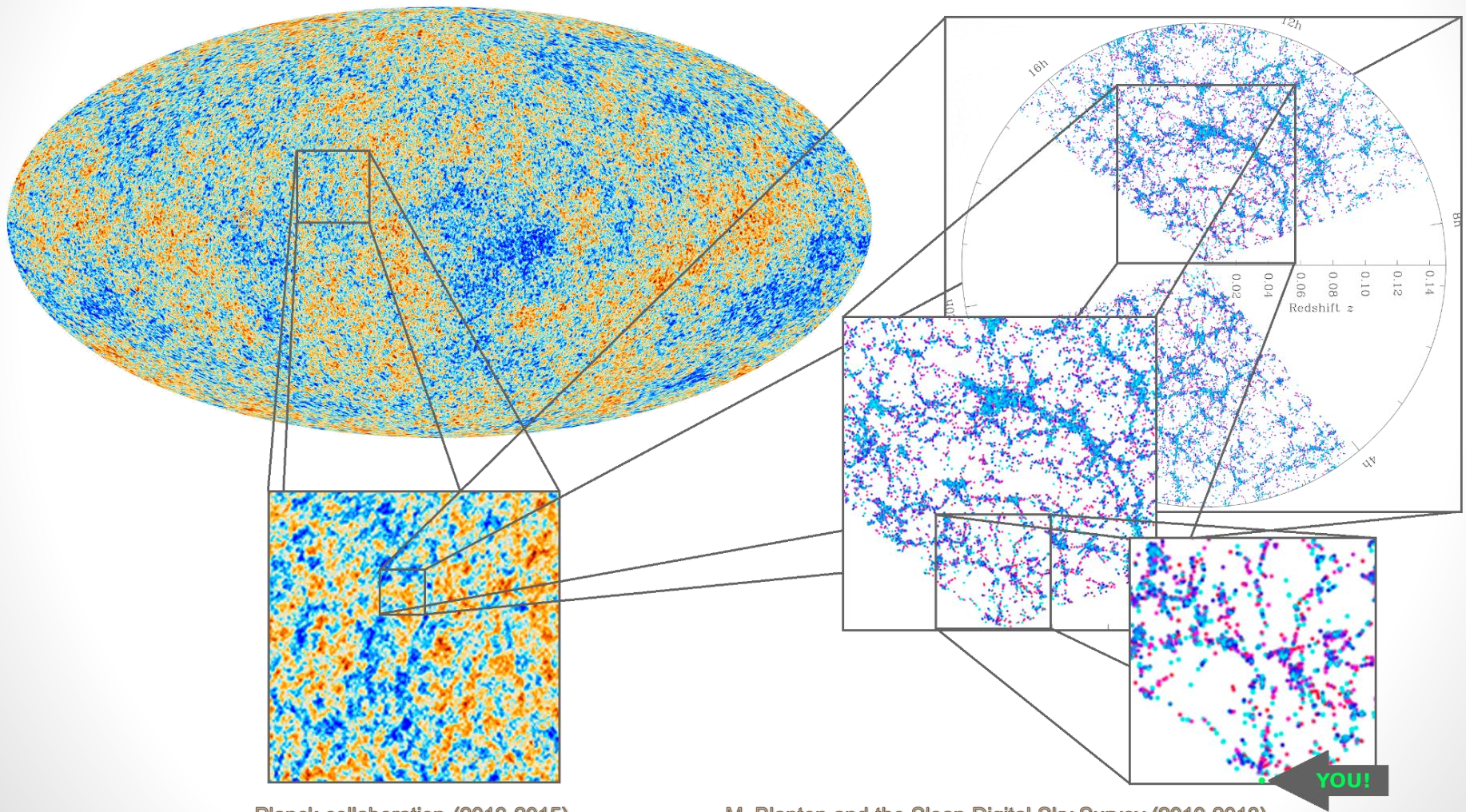
ICIC

Imperial Centre
for Inference & Cosmology

**Imperial College
London**

The big picture: the Universe is highly structured

You are here. Make the best of it...



Planck collaboration (2013-2015)

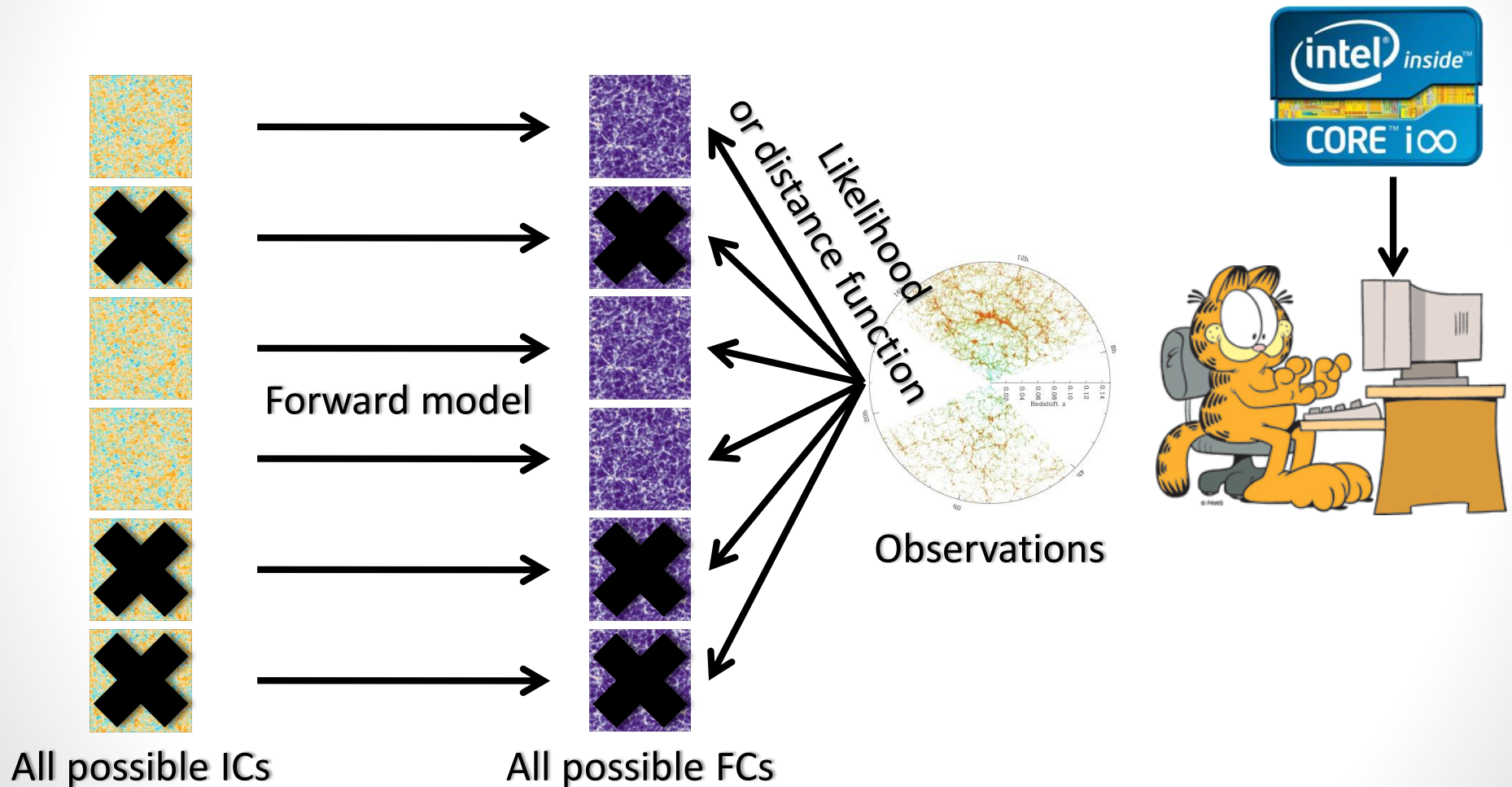
M. Blanton and the Sloan Digital Sky Survey (2010-2013)

What we want to know from the large-scale structure

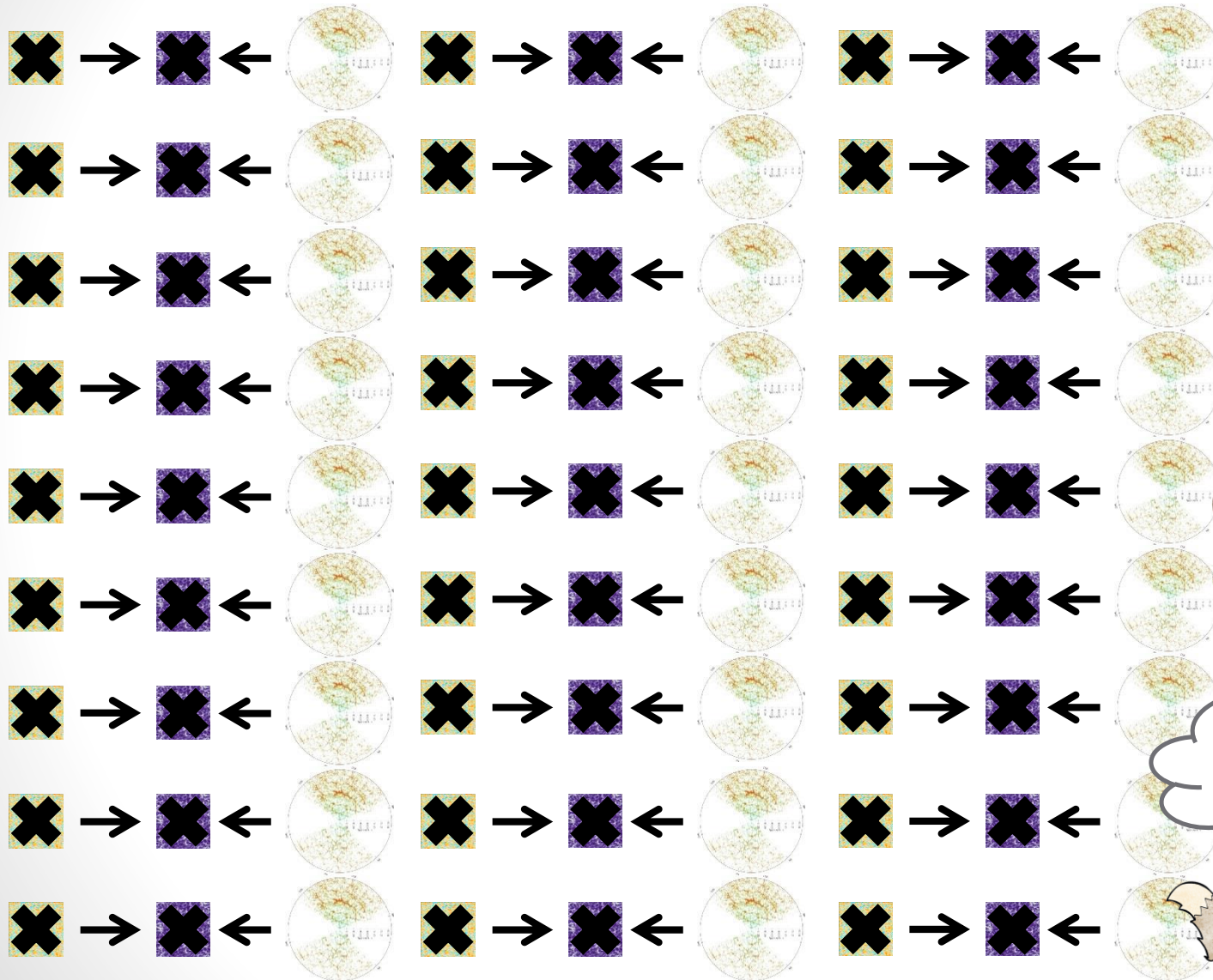
The LSS is a vast source of knowledge:

- **Cosmology:**
 - Λ CDM : cosmological parameters and tests against alternatives,
 - Physical nature of the dark components,
 - Neutrinos : number and masses,
 - Geometry of the Universe,
 - Tests of General Relativity,
 - Initial conditions and link to high energy physics
- **Astrophysics:** galaxy formation and evolution as a function of their environment
 - Galaxy properties (colours, chemical composition, shapes),
 - Intrinsic alignments, intrinsic size-magnitude correlations

Bayesian forward modelling: the ideal scenario



Bayesian forward modelling: the challenge



The (true) likelihood lives in $d \approx 10^7$!



Likelihood-based solution: BORG

Bayesian Origin Reconstruction from Galaxies

Likelihood-based solution:

Exact statistical analysis
Approximate data model

Data assimilation



Likelihood-based solution: BORG at work

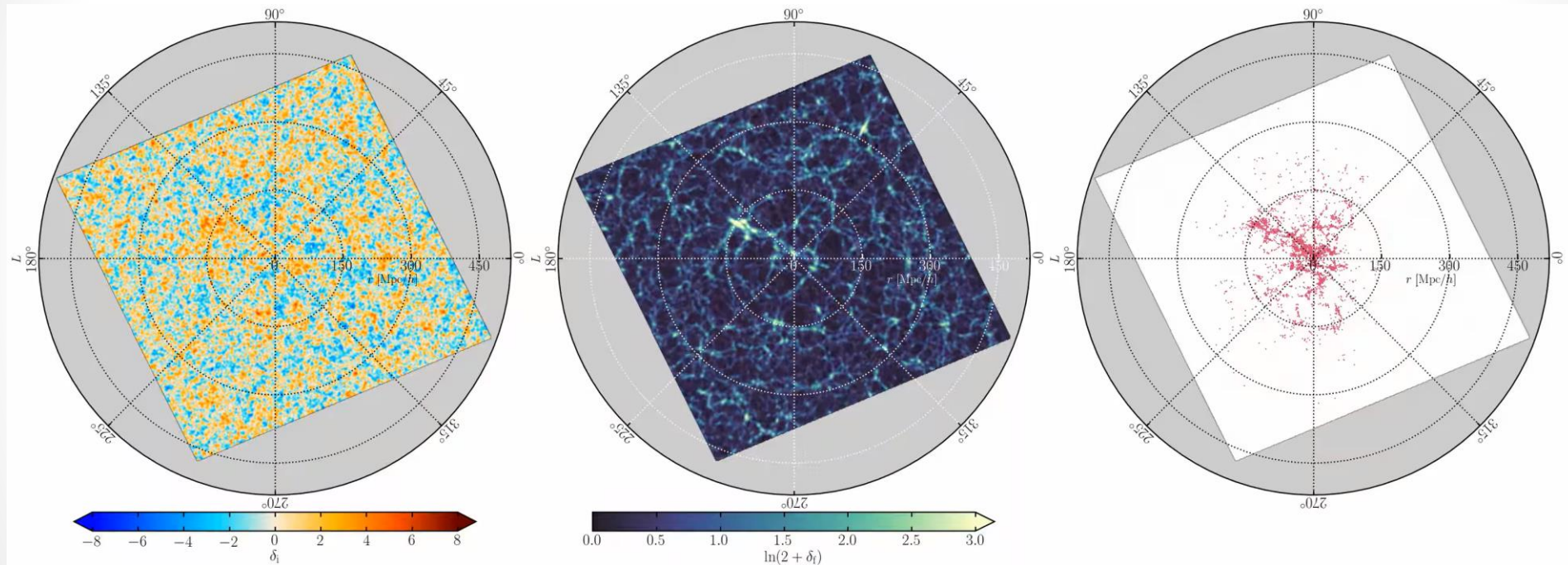


www.aquila-consortium.org/

Initial conditions

Final conditions

Observations



Supergalactic plane

67,224 galaxies, \approx 17 million parameters, 5 TB of primary data products, 10,000 samples, \approx 500,000 forward and adjoint gradient data model evaluations, 1.5 million CPU-hours

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

Likelihood-free solution: BOLFI & SELFI

Bayesian Optimisation for Likelihood-Free Inference
Simulator Expansion for Likelihood-Free Inference

Likelihood-based solution:

Exact statistical analysis
Approximate data model

Data assimilation



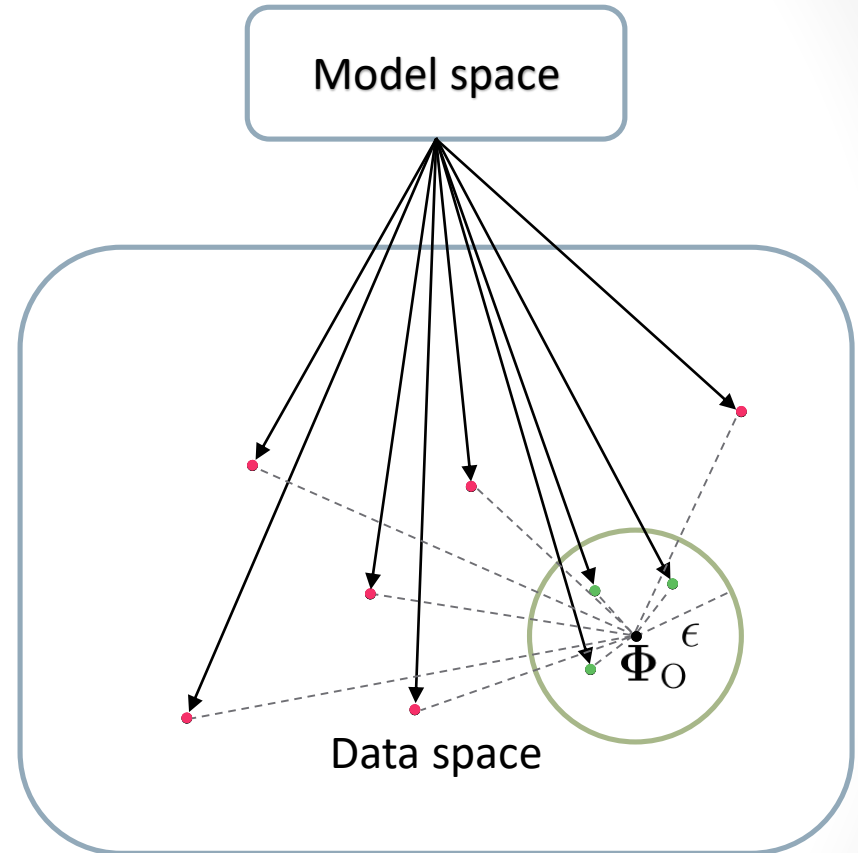
Likelihood-free solution:

Approximate statistical analysis
Arbitrary data model

Generative inference

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)

Effective likelihood approximation:

$$L(\theta) \approx \frac{1}{N} \sum_{i=1}^N \mathbb{I} \left(d(\tilde{d}(\theta), d) \leq \epsilon \right)$$

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

Alsing, Charnock, Feeney & Wandelt 2019, 1903.00007

- 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 (SELFIE)

FL, Enzi, Jasche & Heavens 2019, 1902.10149

I thought of the name **after** developing the method!

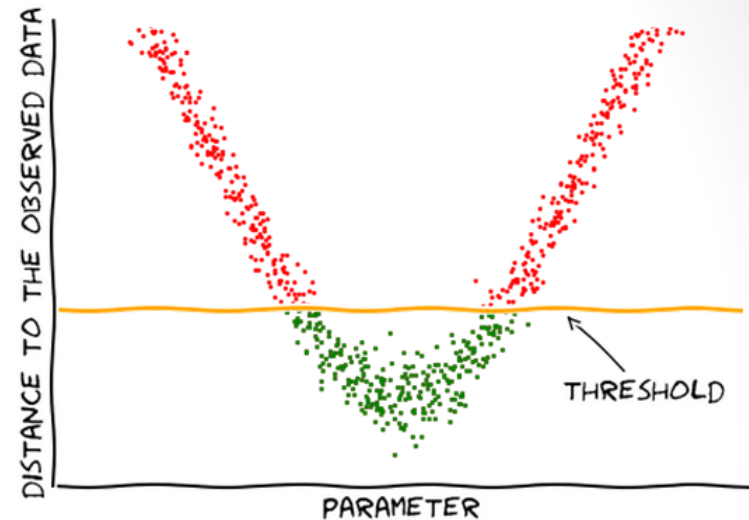


The “number of simulations” route: BOLFI

Bayesian Optimisation for Likelihood-Free Inference

Why is LFRS so expensive?

1. It rejects most samples when ϵ is small
2. It does not make assumptions about the shape of $L(\theta)$
3. It uses only a fixed proposal distribution, not all information available
4. It aims at equal accuracy for all regions in parameter space



$$L(\theta) \approx \frac{1}{N} \sum_{i=1}^N \mathbb{I} \left(d(\tilde{d}(\theta), d) \leq \epsilon \right)$$

Proposed solution: Regression + Active data acquisition

1. It rejects most samples when ϵ is small

➡ Don't reject samples: learn from them!

2. It does not make assumptions about the shape of $L(\theta)$

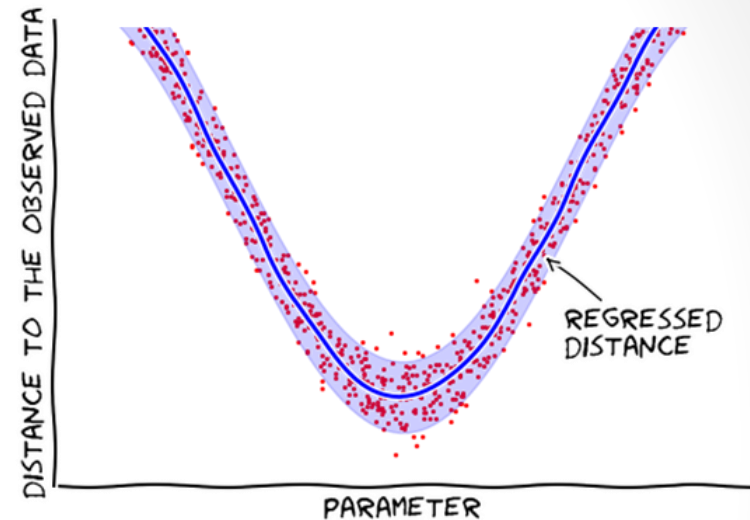
➡ Model the distances, assuming the average distance is smooth

3. It uses only a fixed proposal distribution, not all information available

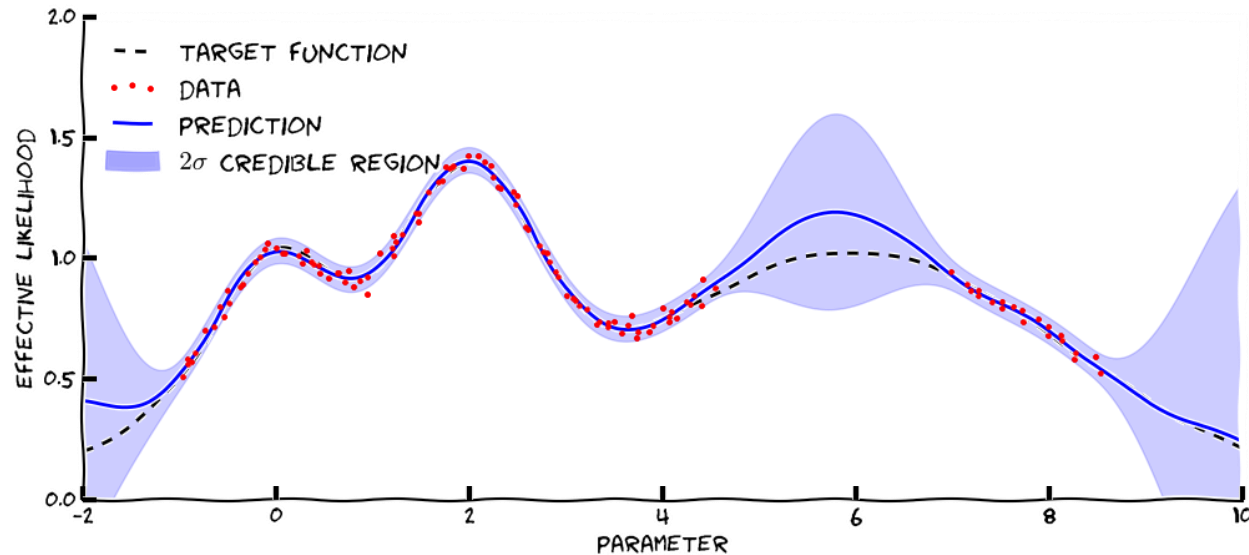
➡ Use Bayes' theorem to update the proposal of new points

4. It aims at equal accuracy for all regions in parameter space

➡ Prioritise "interesting" regions in parameter space

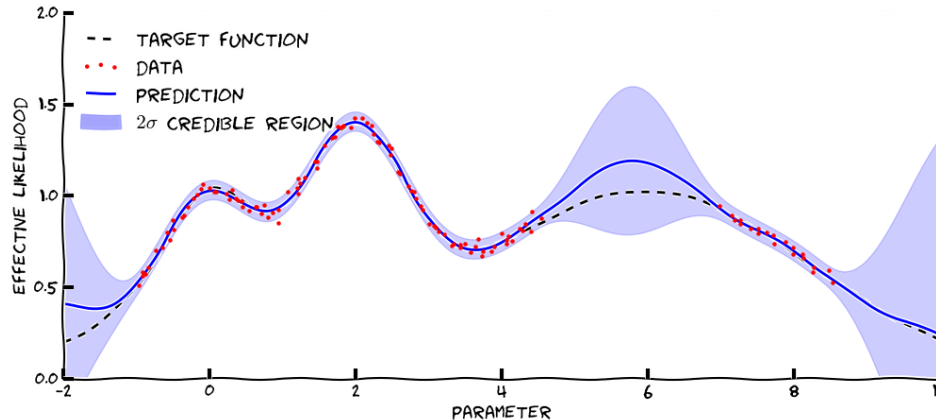


BOLFI: Regression of the effective likelihood



1. “LFRS rejects most samples when ϵ is small”
 - Keep all values (θ_i, d_i) $d_i = d(\tilde{d}(\theta_i), d)$
2. “LFRS does not make assumptions about the shape of $L(\theta)$ ”
 - Model the conditional distribution of distances given this training set

Gaussian process regression (a.k.a. kriging)



- Why?

- It is a **general purpose regressor**: it will be able to deal with a large variety of complex/non-linear features of likelihood functions.
- It provides not only a prediction, but also the **uncertainty of the regression**.
- It allows to **extrapolate** in regions where we have no data points.

$$p(\mathbf{f}|\mathbf{X}) \propto \exp \left[-\frac{1}{2} \sum_{mn} (f(\mathbf{x}_m) - \mu(\mathbf{x}_m))^\top K(\mathbf{x}_m, \mathbf{x}_n) (f(\mathbf{x}_n) - \mu(\mathbf{x}_n)) \right]$$

$$K(\mathbf{x}_m, \mathbf{x}_n) = \underbrace{C_1}_{K_C(C_1)} \times \underbrace{\exp \left[-\frac{1}{2} \left(\frac{\mathbf{x}_m - \mathbf{x}_n}{C_2} \right)^2 \right]}_{K_{\text{RBF}}(C_2)} + \underbrace{C_3 \delta_K^{mn}}_{K_{\text{GN}}(C_3)}$$

The prediction and uncertainty for a new point is:

$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{f}) \propto \exp \left[-\frac{1}{2} \left(\frac{f_* - \alpha(\mathbf{x}_*)}{\sigma(\mathbf{x}_*)} \right)^2 \right]$$

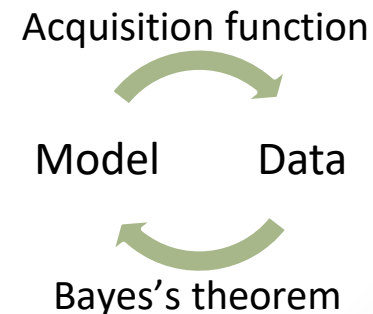
$$\alpha(\mathbf{x}_*) = \mu(\mathbf{x}_*) + K(\mathbf{x}_*, \mathbf{x}_m)^\top K^{-1}(\mathbf{x}_m, \mathbf{x}_n) (\mathbf{f} - \mu(\mathbf{X}))_n$$

$$\sigma(\mathbf{x}_*)^2 = K(\mathbf{x}_*, \mathbf{x}_*) - K(\mathbf{x}_*, \mathbf{x}_m)^\top K^{-1}(\mathbf{x}_m, \mathbf{x}_n) K(\mathbf{x}_*, \mathbf{x}_n)$$

Hyperparameters C_1, C_2, C_3 are automatically adjusted during the regression.

BOLFI: Data acquisition

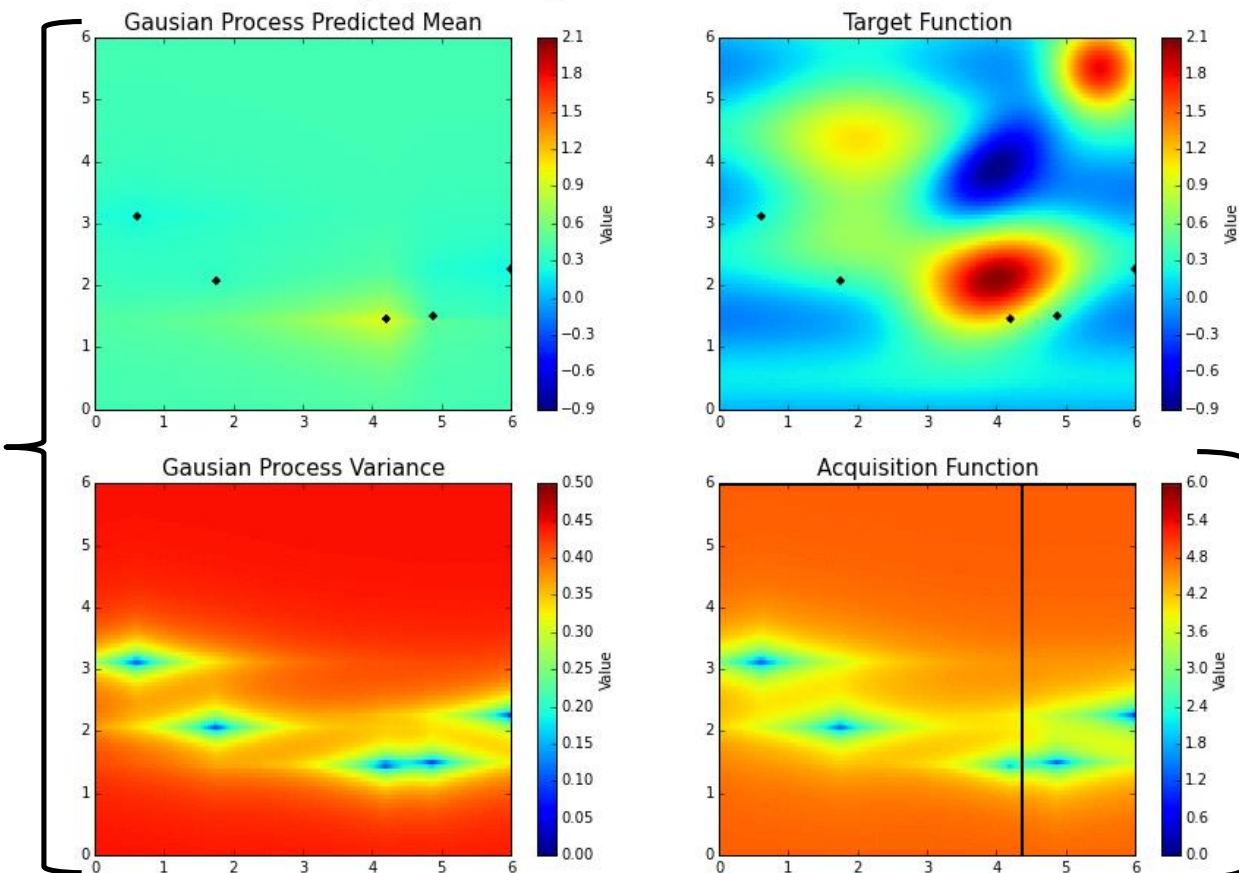
3. “LFRS uses only a fixed proposal distribution, not all information available”
 - Samples are obtained from sampling an **adaptively-constructed proposal distribution**, using the regressed effective likelihood
4. “LFRS aims at equal accuracy for all regions in parameter space”
 - The **acquisition function** finds a compromise between **exploration** (trying to find new high-likelihood regions) & **exploitation** (giving priority to regions where the distance to the observed data is already known to be small)
 - **Bayesian optimisation** (decision making under uncertainty) can then be used



BOLFI: Data acquisition

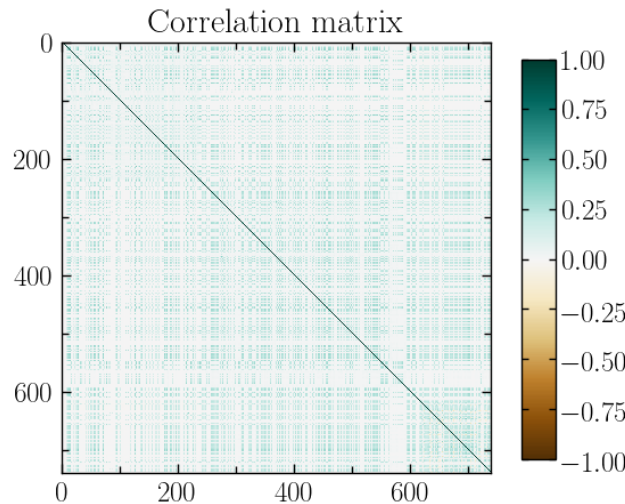
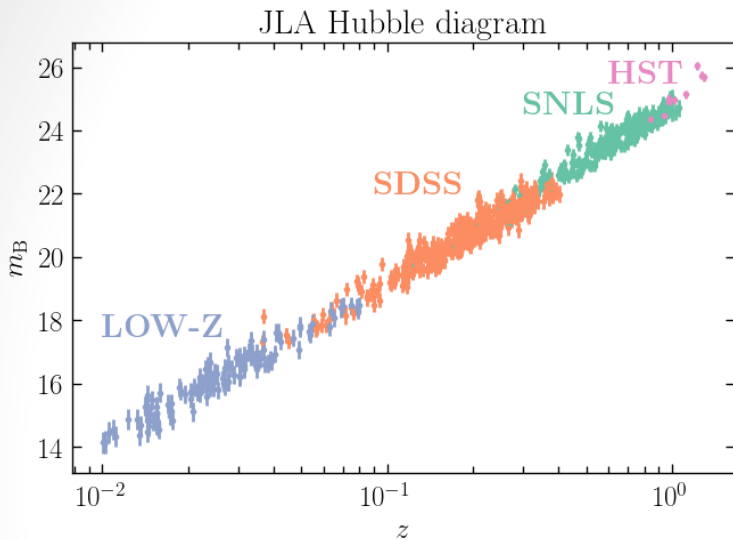
Bayesian Optimization in Action

Regression of the distance between observed and simulated data



Expected utility of the next simulation in parameter space

BOLFI: Re-analysis of the JLA supernova sample



Betoule et al. 2014, arXiv:1401.4064

- 6-parameter model:

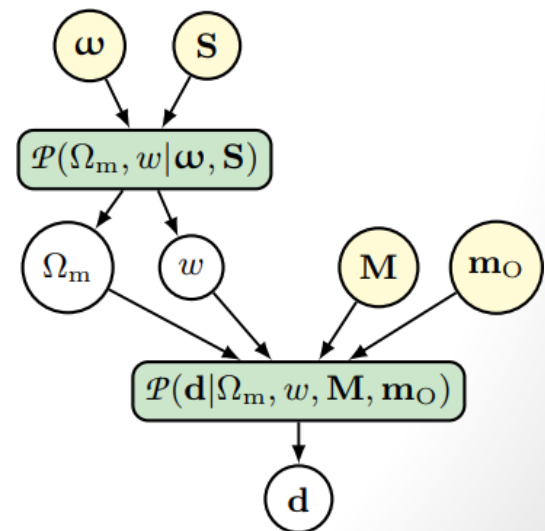
2 cosmological parameters + 4 nuisance parameters

$$m_B = 5 \log_{10} \left[\frac{D_L(z)}{10 \text{ pc}} \right] + \tilde{M}_B(M_{\text{stellar}}, M_B, \delta M) - \alpha X_1 + \beta C$$

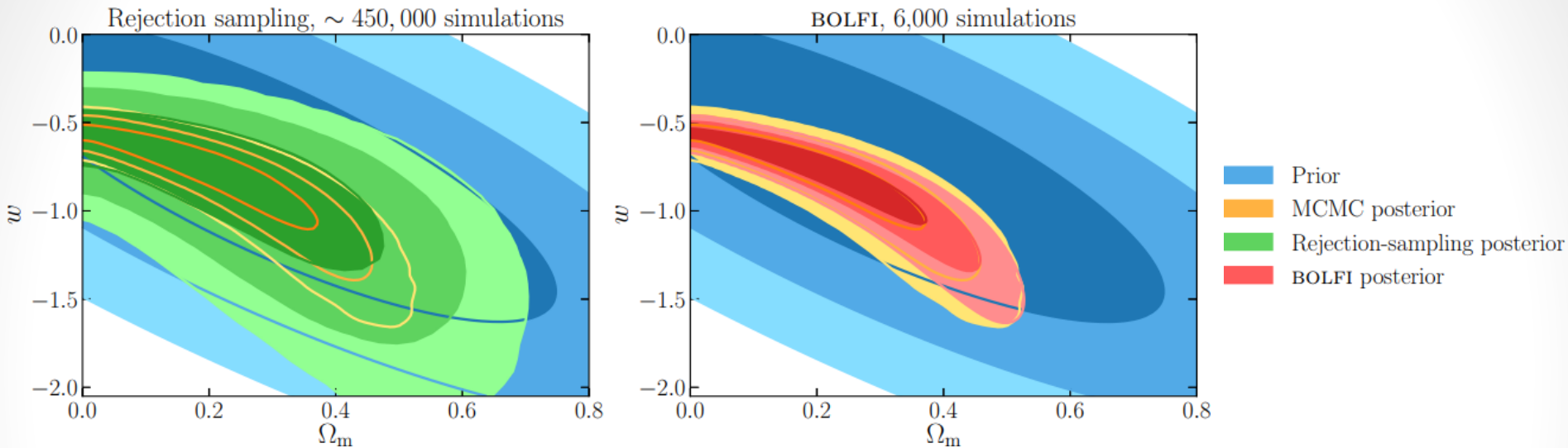
$$\tilde{M}_B(M_{\text{stellar}}, M_B, \delta M) = \bar{M}_B + \delta M \Theta (M_{\text{stellar}} - 10^{10} M_{\odot})$$

$$D_L(z) = \frac{(1+z)c}{H_0} \int_0^z \frac{dz'}{E(z')}$$

$$E(z) \equiv \sqrt{\Omega_m (1+z)^3 + (1 - \Omega_m)(1+z)^{3(w+1)}}$$



BOLFI: Re-analysis of the JLA supernova sample



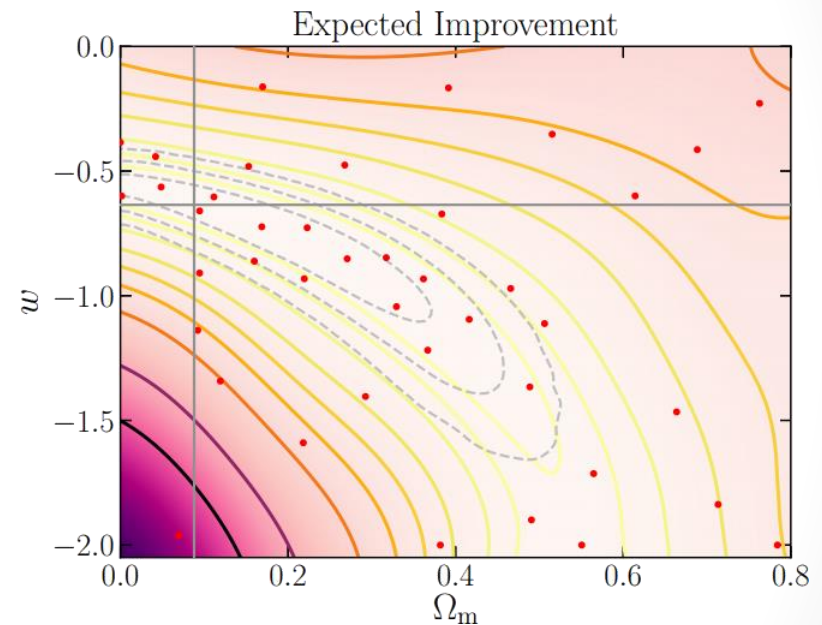
- The **number of required simulations is reduced** by:
 - 2 orders of magnitude with respect to likelihood-free rejection sampling (for a much better approximation of the posterior)
 - 3 orders of magnitude with respect to exact Markov Chain Monte Carlo sampling

Standard acquisition functions are suboptimal

- Goal for Bayesian optimisation: find the optimum (assumed unique) of a function
- Example of acquisition function : the **Expected Improvement**

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

Gaussian cdf Gaussian pdf



- Drawbacks:
 - Do not take into account prior information
 - Local evaluation rules
 - Too greedy for ABC

e.g. Brochu, Cora & de Freitas 2010, arXiv:1012.2599

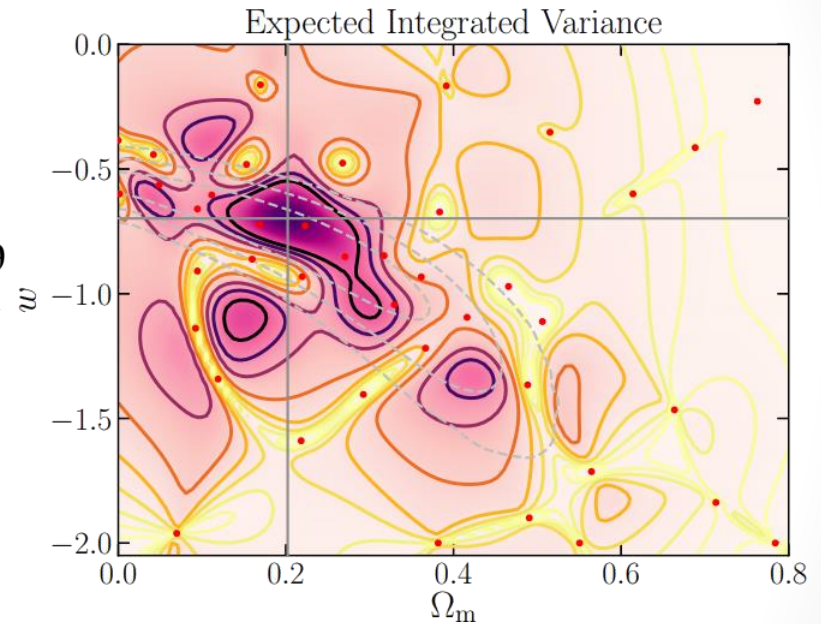
FL 2018, arXiv:1805.07152

The optimal acquisition function for ABC

- Goal for ABC: 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**

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

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



- 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

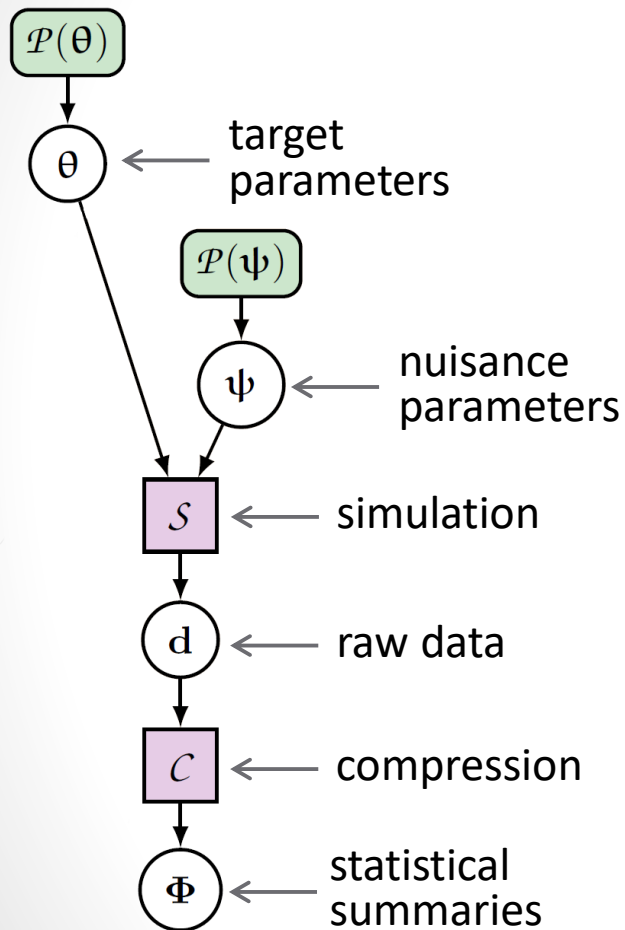
Järvenpää et al. 2017, [arXiv:1704.00520](https://arxiv.org/abs/1704.00520) (expression of the EIV in the non-parametric approach)

FL 2018, [arXiv:1805.07152](https://arxiv.org/abs/1805.07152) (expression of the EIV in the parametric approach)

The “number of parameters” route: SELFI

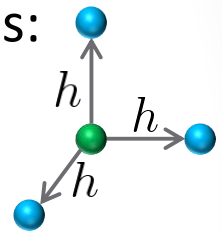
Simulator Expansion for Likelihood-Free Inference

SELFIE: Method



- Gaussian prior + Gaussian effective likelihood
- Linearisation of the black-box around an expansion point + finite differences:

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



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

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

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

expansion point θ_0 observed summaries Φ_O
 covariance of summaries \mathbf{C}_0 gradient of the black-box $\nabla \mathbf{f}_0$ prior covariance \mathbf{S}

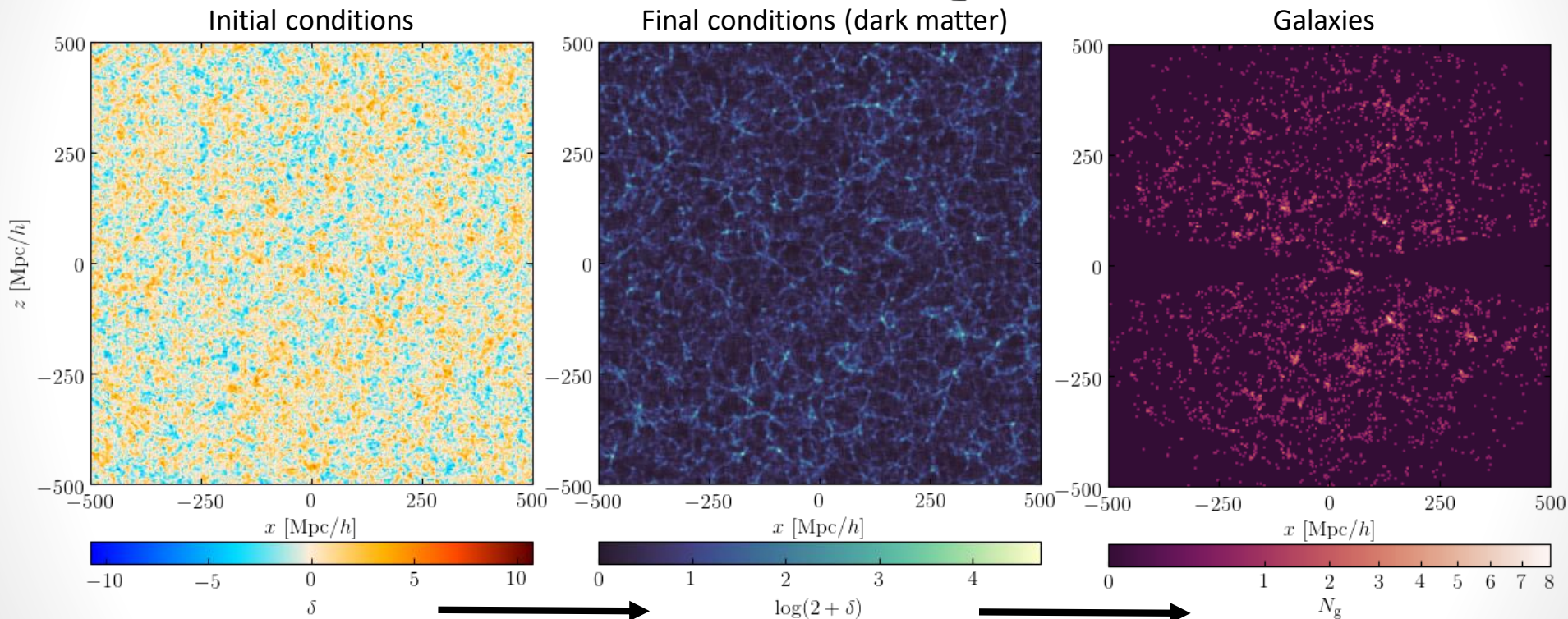
A black-box: Simbelmynë

I'm happy to explain the name during the coffee break...



Publicly available code:

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

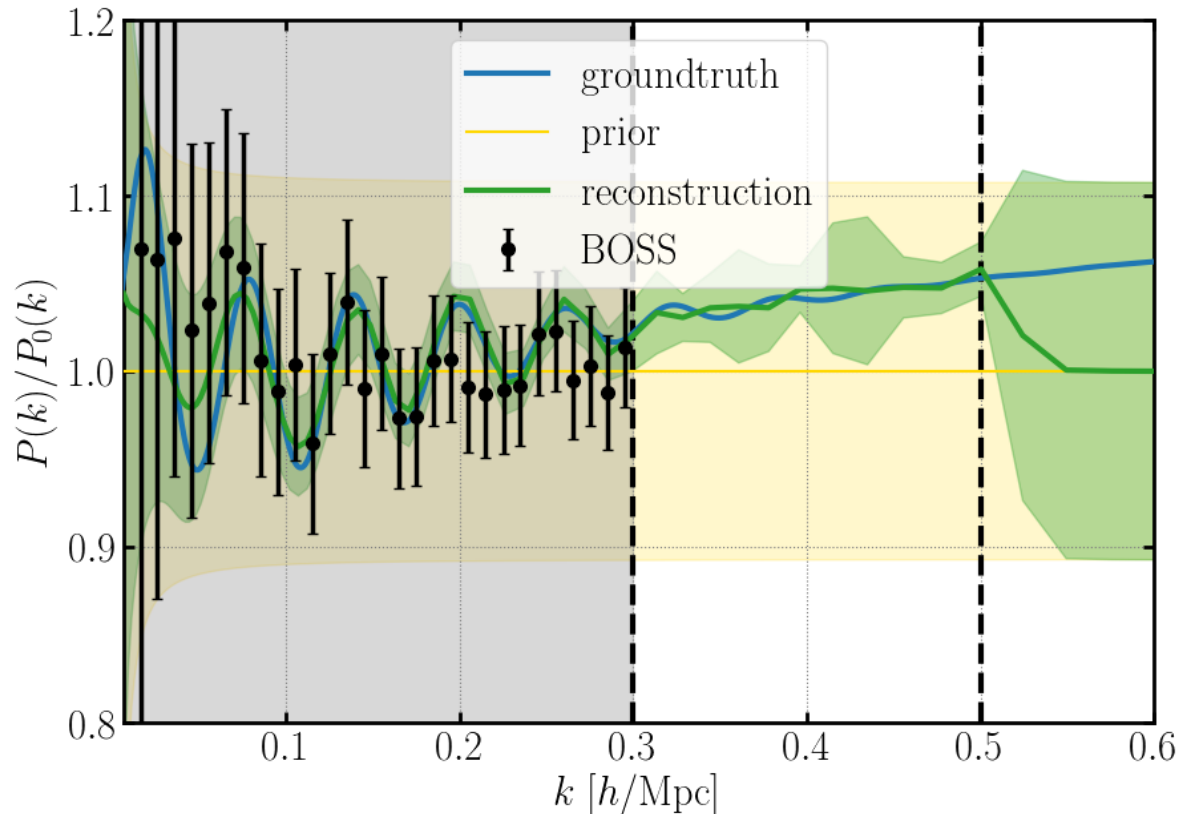


Dark matter simulation
with COLA

Tassev, Zaldarriaga & Eisenstein 2013, arXiv:1301.0322

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

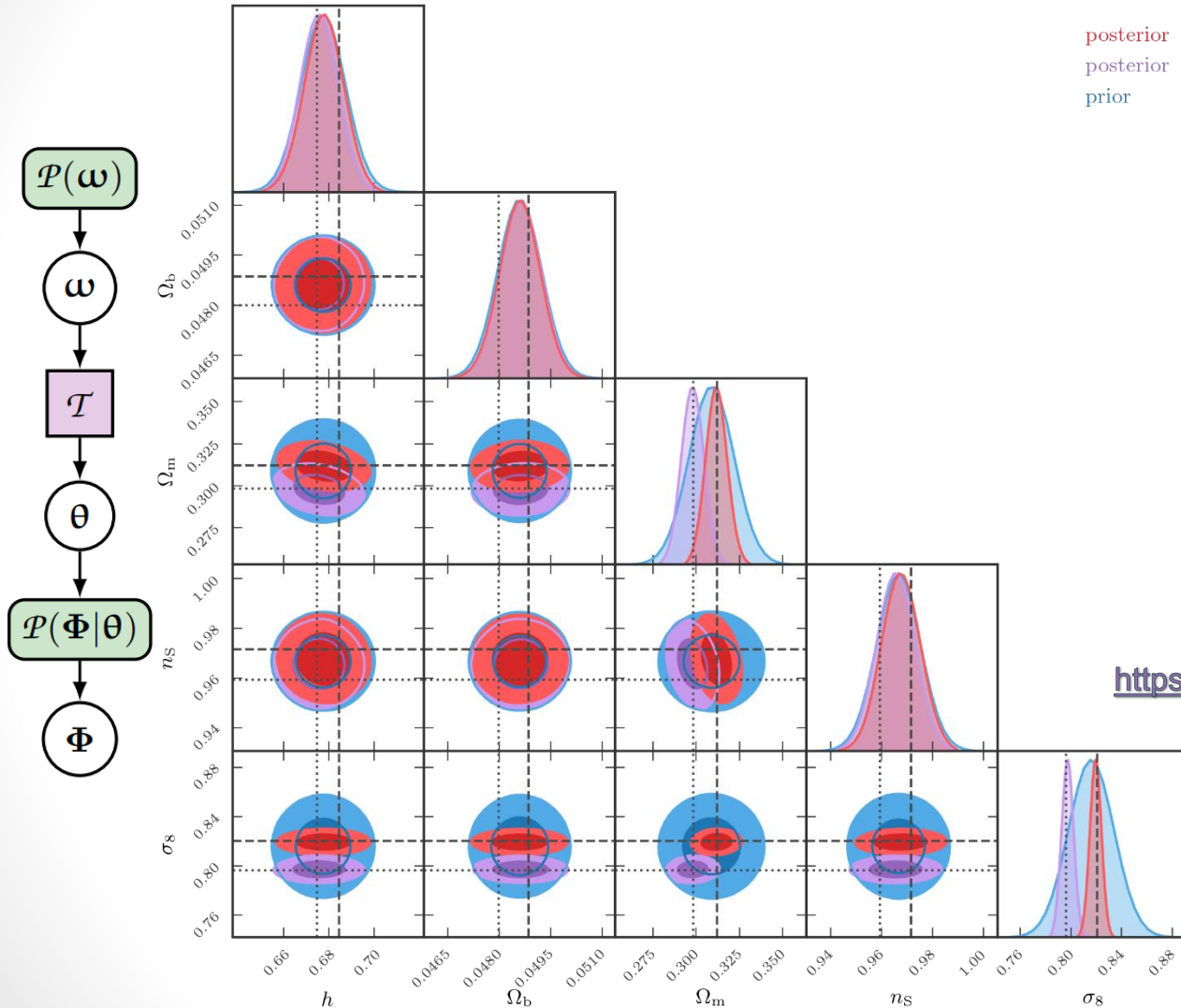
SEIFI + Simbelmynë: Proof-of-concept



100 parameters are simultaneously inferred from a black-box data model

$N_{\text{modes}} \propto k^3$: **5** times more modes are used in the analysis

SELFIE + Simbelmynë: Proof-of-concept



- Robust inference of cosmological parameters can be easily performed *a posteriori* once the linearised data model is learned
- **pyselfi** is publicly available:

<https://github.com/florent-leclercq/pyselfi/>

Concluding thoughts

- **Goal:** developing and using algorithms for targeted 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!
- The “**number of simulations route**” (BOLFI):
 - The optimal acquisition function can be derived: the Expected Integrated Variance.
 - The number of simulations is reduced by several orders of magnitude.
- The “**number of parameters route**” (SELF):
 - High-dimensional likelihood-free problems can be addressed.
 - The computational workload is fixed *a priori* and perfectly parallel.