# Primordial power spectrum and cosmology from black-box galaxy surveys

## Prospects for Euclid
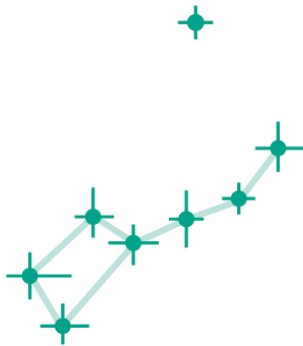
### Florent Leclercq
www.florent-leclercq.eu

Imperial Centre for Inference and Cosmology
Imperial College London

Wolfgang Enzi, Alan Heavens,
Jens Jasche, Guilhem Lavaux,

and the Aquila Consortium
www.aquila-consortium.org

December 11th, 2019
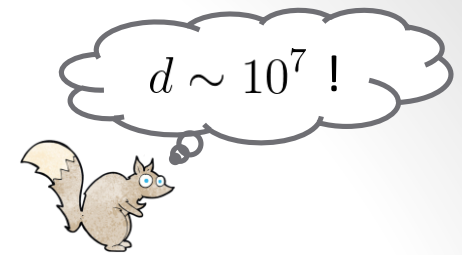
ICIC
Imperial Centre
for Inference & Cosmology

Imperial College
London

# Vocabulary considerations I:
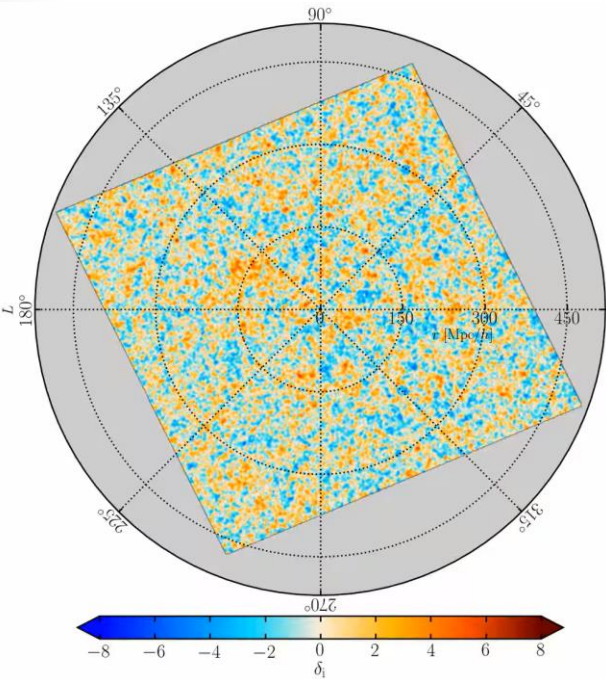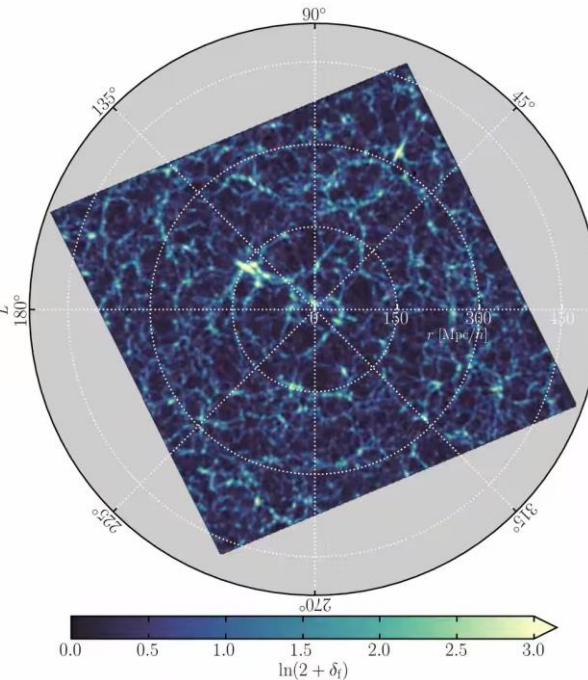*What is the likelihood?*

$d \sim 10^7$ !

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**



Supergalactic plane

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).

# Vocabulary considerations II:

*You may already be an LFI specialist!*

- Likelihood-free inference (LFI) techniques bypass the need for a map-level likelihood, by relying instead only on a "black-box".

- The likelihood is replaced by a measure of the distance/discrepancy $\Delta$ between simulated and observed statistical summaries of the data.

**diagram labels:**

$\mathcal{P}(\theta)$

$\theta$ ← target parameters

$\mathcal{P}(\psi)$

$\psi$ ← nuisance parameters

$\mathcal{S}$ ← simulation

$\mathbf{d}$ ← raw data

$\mathcal{C}$ ← compression

$\Phi$ ← statistical summaries

e.g. $\mathbf{d} =$ full galaxy survey data
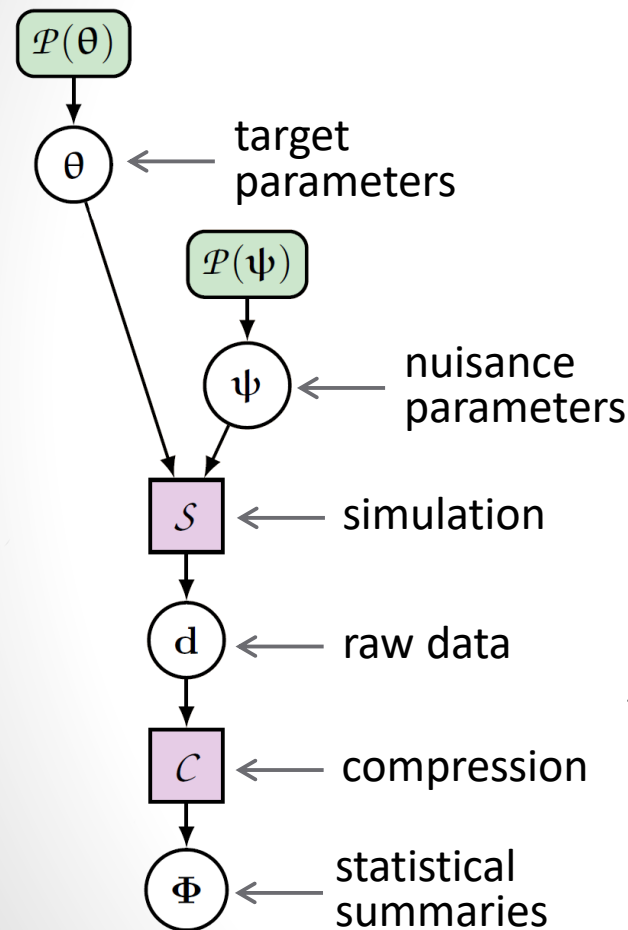$\Phi = \{\widehat{P}(k)\}$ estimated power spectrum
$\Delta =$ Mahalanobis distance with covariance matrix $\Sigma$

$$\Delta(\Phi_{\theta}, \Phi_{O}) = \sqrt{\sum_{k,k'} \left[ \widehat{P}_{\theta}(k) - \widehat{P}_{O}(k) \right]^{\mathsf{T}} \Sigma^{-1}_{k,k'} \left[ \widehat{P}_{\theta}(k') - \widehat{P}_{O}(k') \right]}$$

Note that this is what many people would call…
(square root of -2 times) the log-likelihood!
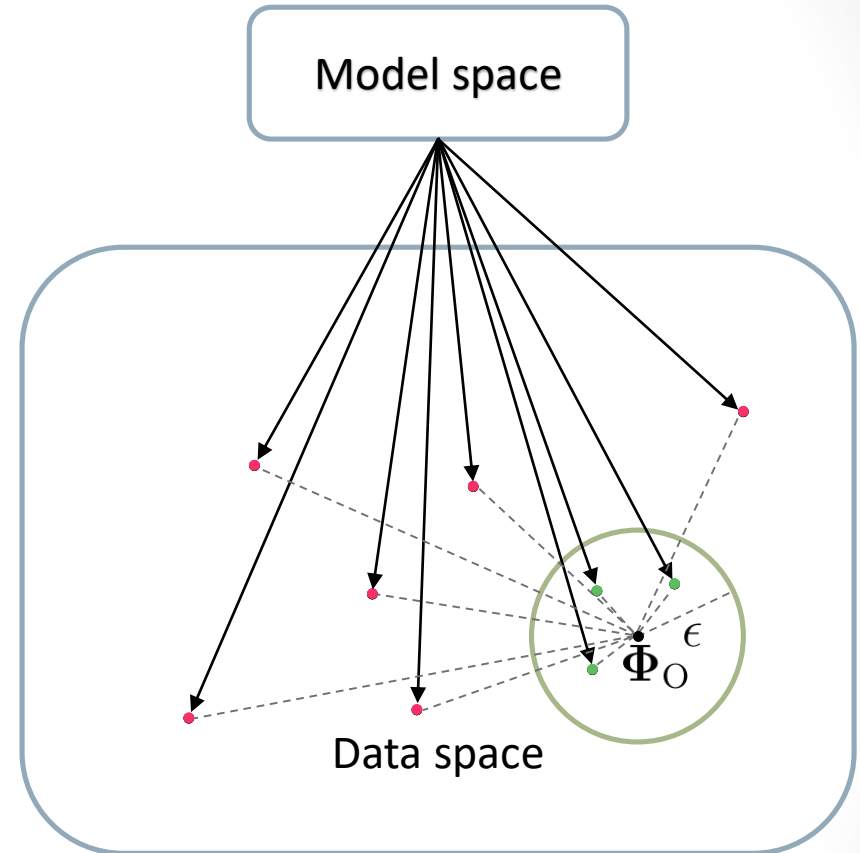
- What is "primordial" depends mostly on your ambition…

# Cosmological synergies with LFI

- Advantages of likelihood-free inference:

  - No expert knowledge (conditionals/gradients of the likelihood) is required

  - Summary statistics need not be physically modelled and can be chosen robustly to model misspecification, e.g.:
    - Microwave sky: cross-power spectra between different frequency maps
    - Imaging surveys: cross-correlation between different bands

  - Joint and self-consistent analyses of correlated data sets is straightforward

- Drawbacks of likelihood-free inference:

  - No inference of the map

  - Relies on (lossy) data compression and statistical approximations
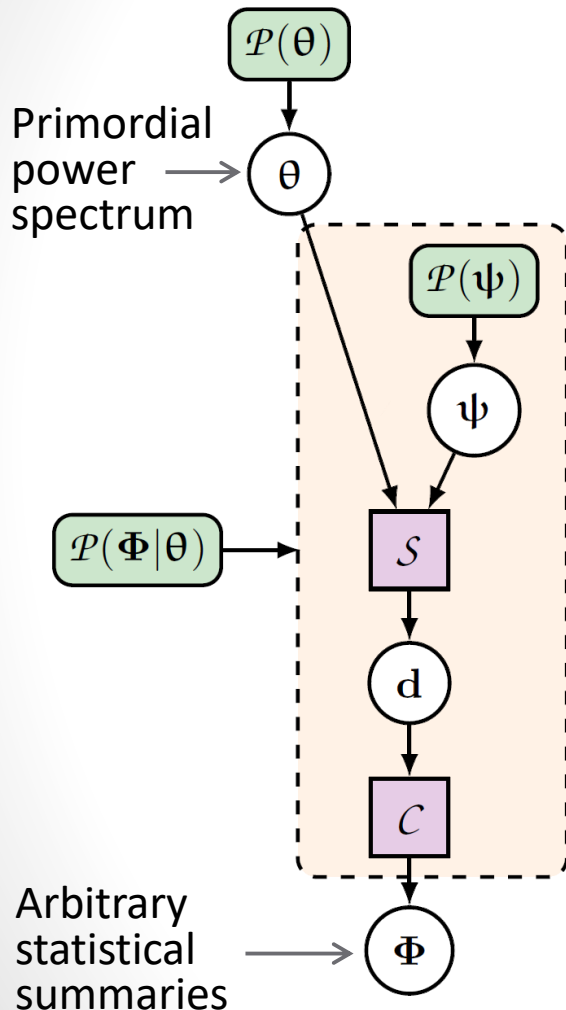
# Likelihood-free rejection sampling (LFRS)

- Iterate many times:
  - Sample $\theta$ from a proposal distribution $q(\theta)$
  - Simulate $\mathbf{\Phi}_\theta$ using the black-box
  - Compute the distance $\Delta(\mathbf{\Phi}_\theta, \mathbf{\Phi}_O)$ between simulated and observed data
  - Retain $\theta$ if $\Delta(\mathbf{\Phi}_\theta, \mathbf{\Phi}_O) \leq \epsilon$, otherwise reject

$\epsilon$ can be adaptively reduced (Population Monte Carlo)



Model space

$\epsilon$
$\mathbf{\Phi}_O$

Data space

# Beyond LFRS: the SELFI approach

Primordial power spectrum →

$\mathcal{P}(\theta)$

$\theta$

$\mathcal{P}(\psi)$

$\psi$

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

$\mathcal{S}$

$\mathbf{d}$

$\mathcal{C}$

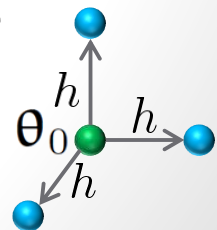Arbitrary statistical summaries →

$\Phi$

- We aim at inferring the primordial 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 $\theta_0$:

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

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

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

$\theta_0$ $h$ $h$ $h$

# SELFI-1: linearization of the black-box

- Linearization of the black-box:

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

- Gaussian prior + Gaussian effective likelihood

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

expansion point                observed summaries

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

$$\boldsymbol{\Gamma} \equiv \left[ (\nabla \mathbf{f}_0)^{\mathsf{T}} \, \mathbf{C}_0^{-1} \nabla \mathbf{f}_0 + \mathbf{S}^{-1} \right]^{-1}$$

prior covariance

covariance of summaries        gradient of the black-box

$\mathbf{f}_0, \mathbf{C}_0$ and $\nabla \mathbf{f}_0$ can be evaluated through simulations only.

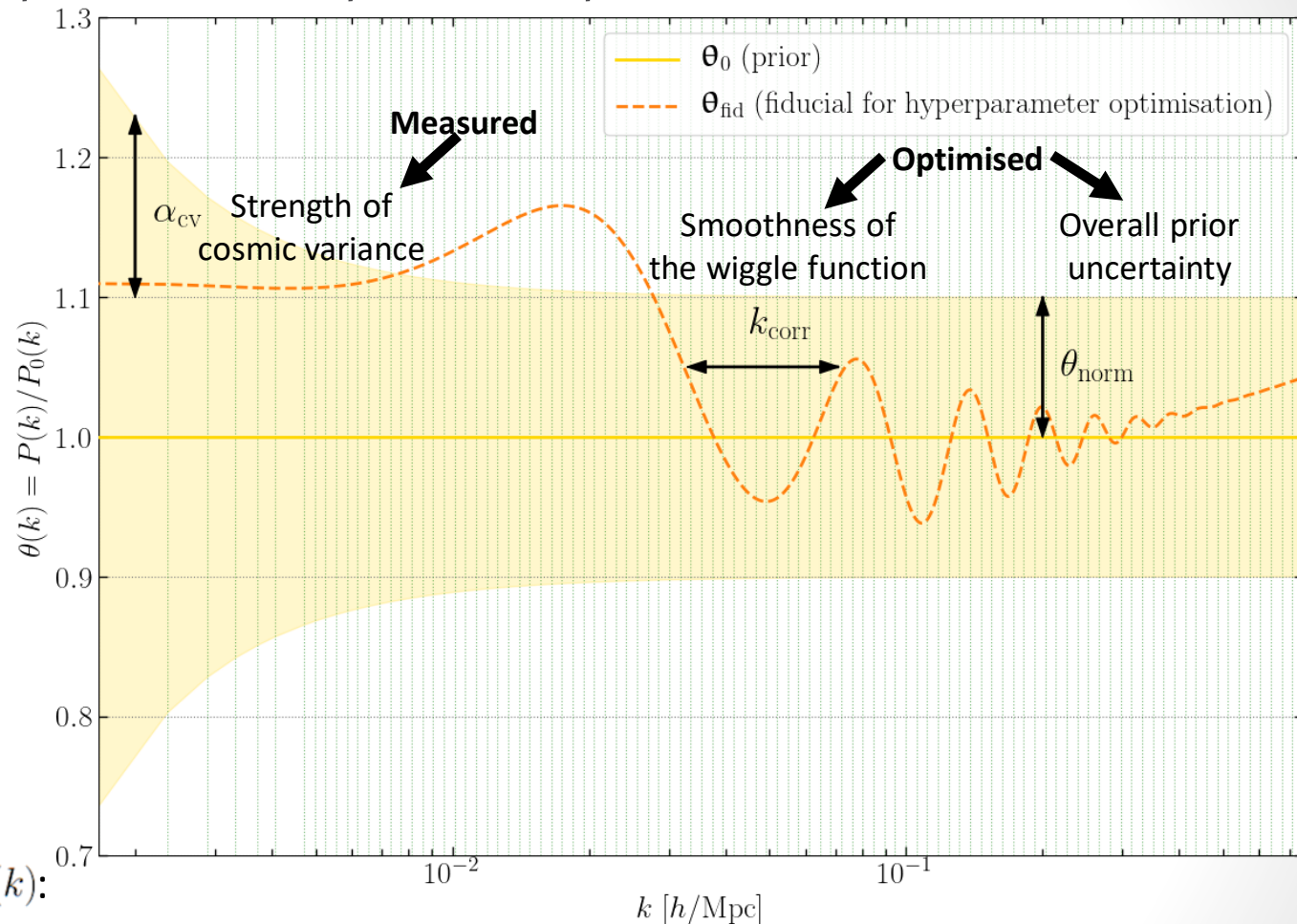The number of required simulations is fixed *a priori* (contrary to MCMC).

The workload is perfectly parallel.

# A prior for the primordial power spectrum

**Assumptions:**

1. the power spectrum is Gaussian-distributed

2. it is strongly constrained to live close to $P_0$,

3. it is a smooth function of wavenumber,

4. and the power spectrum $P_0$ is subject to cosmic variance
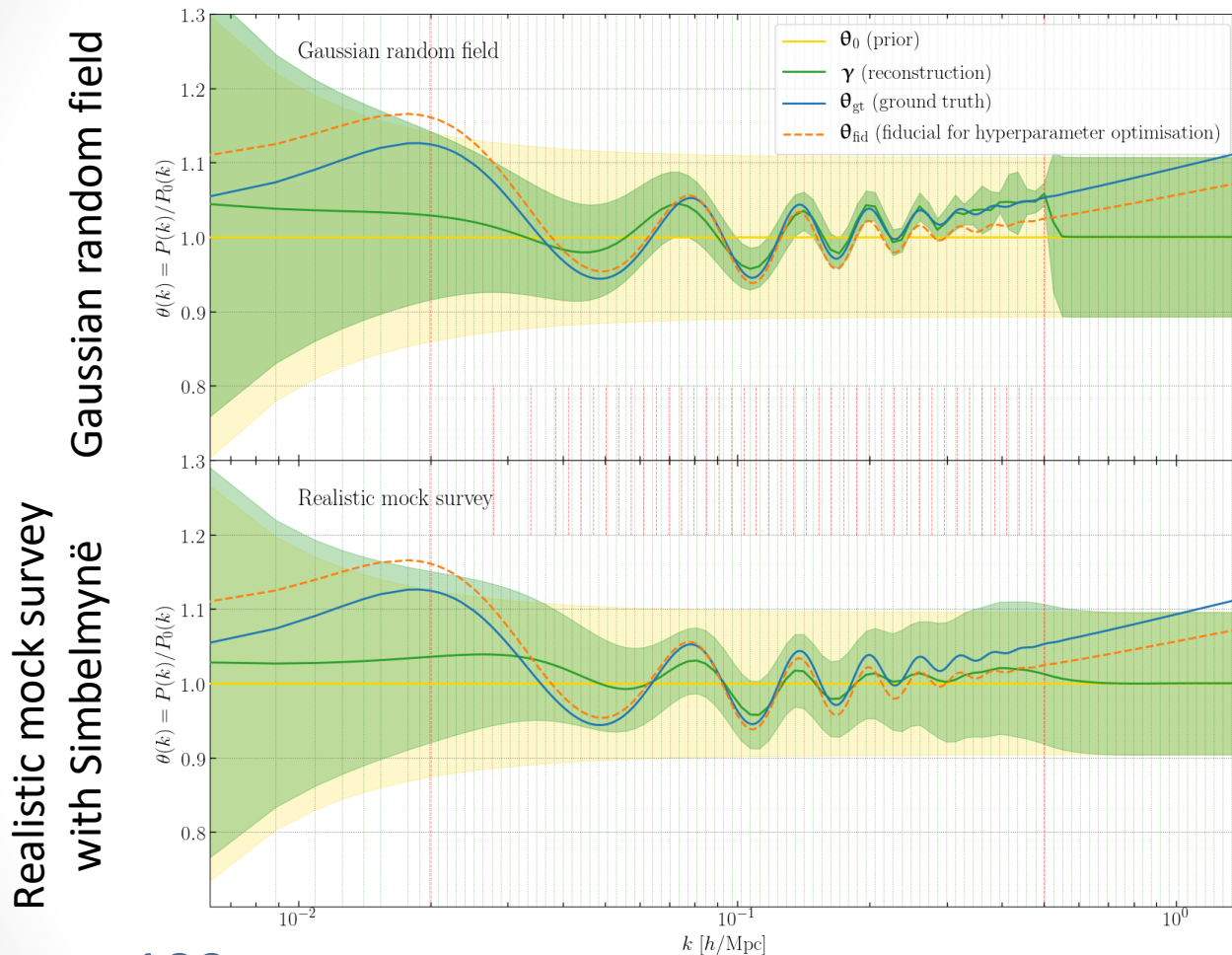


**Prior for $\theta(k) = P(k)/P_0(k)$:**

Mean: $\quad \theta_0 = \mathbf{1}_{\mathbb{R}^S}$ (without baryon acoustic oscillations wiggles)

Covariance: $\quad \mathbf{S} \equiv \theta_{\mathrm{norm}}^2 \, \mathbf{u}\mathbf{u}^\mathsf{T} \circ \mathbf{K} \qquad (\mathbf{K})_{ss'} \equiv \exp\left[-\frac{1}{2}\left(\frac{k_s - k_{s'}}{k_{\mathrm{corr}}}\right)^2\right] \qquad (\mathbf{u})_s \equiv 1 + \sigma_s = 1 + \frac{\alpha_{\mathrm{cv}}}{k_s^{3/2}}$
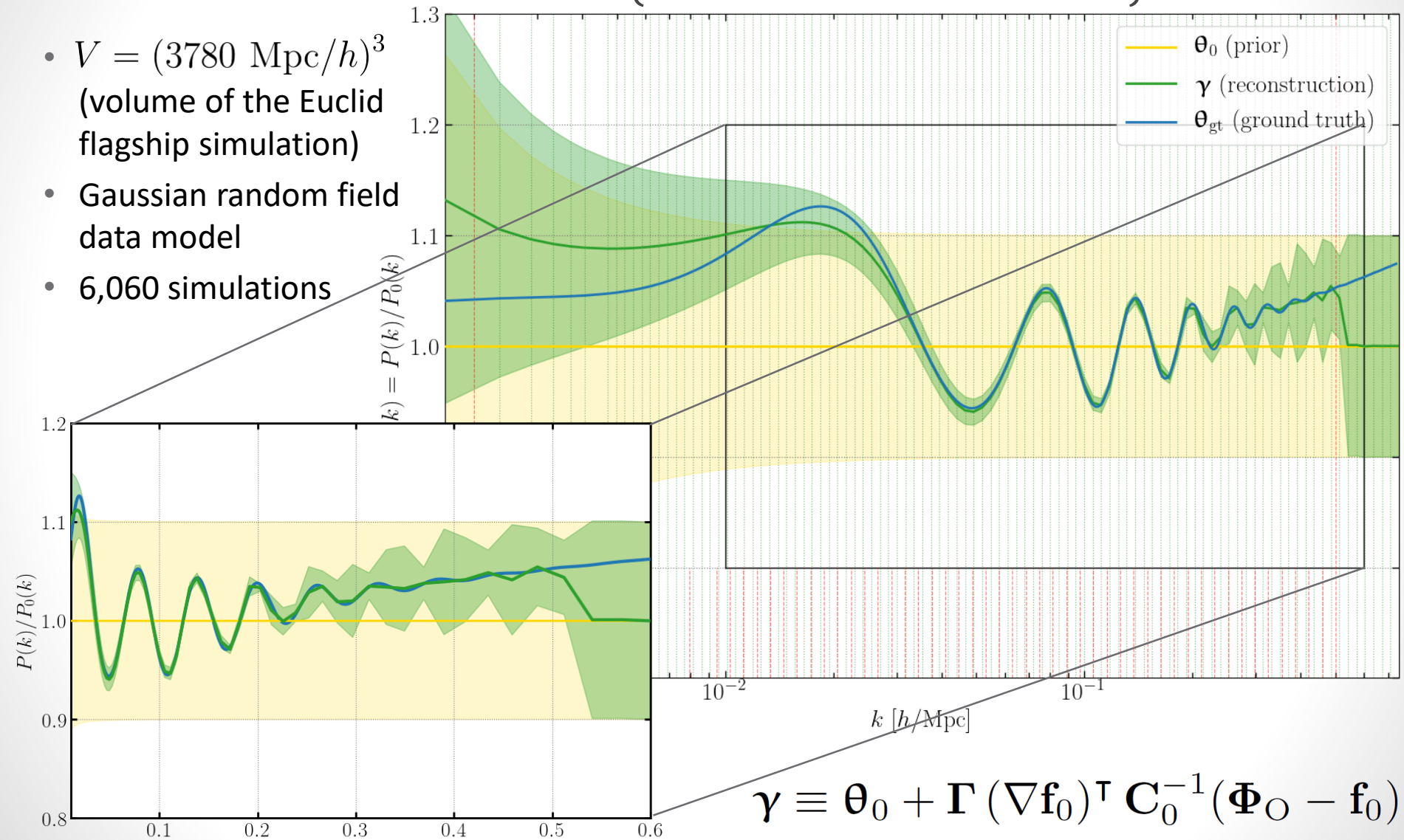
FL, Enzi, Jasche & Heavens 2019, 1902.10149

# 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…

FL, Enzi, Jasche & Heavens 2019, 1902.10149
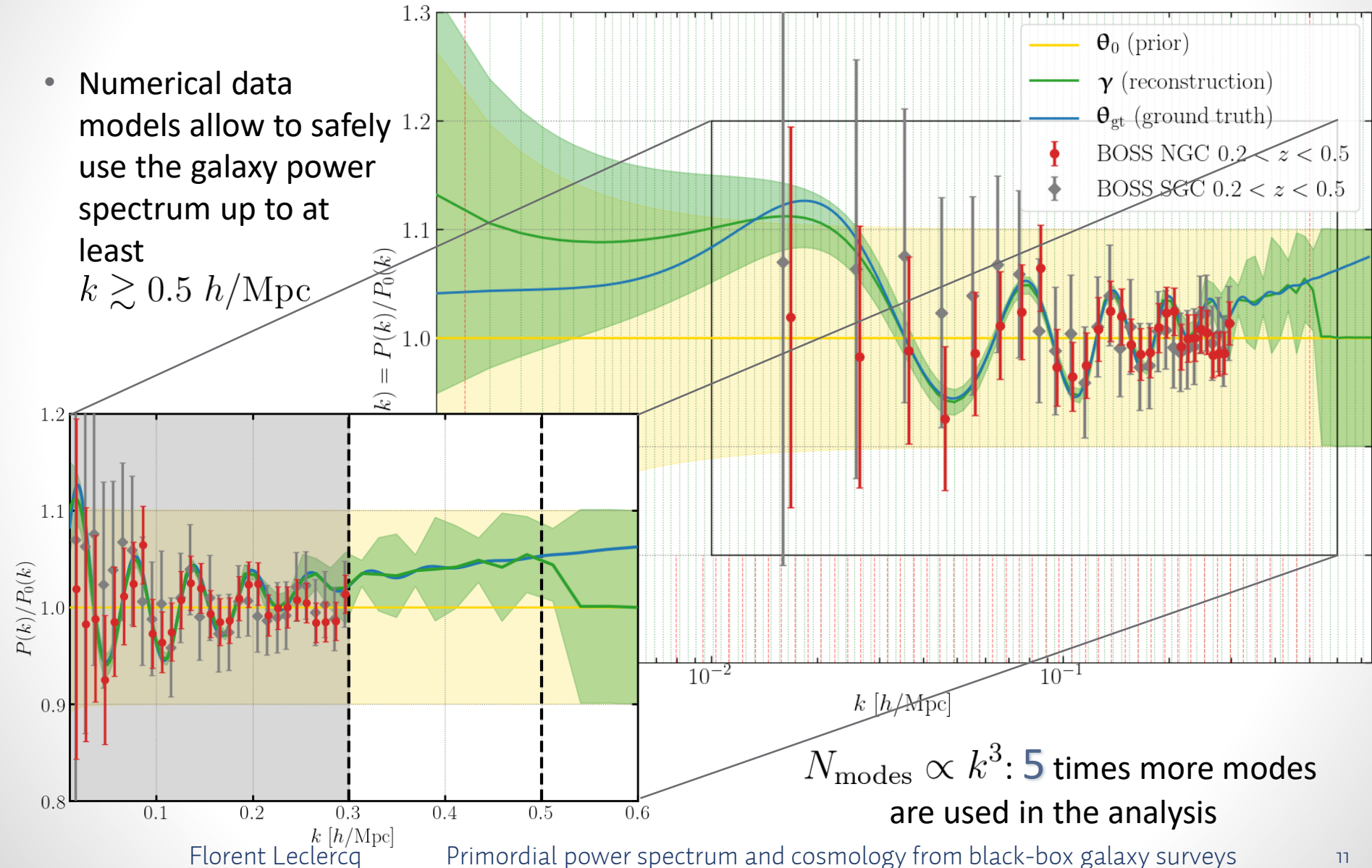
# SELFI-1 Euclid forecast (cosmic variance limit)

- $V = (3780 \text{ Mpc}/h)^3$
  (volume of the Euclid flagship simulation)

- Gaussian random field data model

- 6,060 simulations



$$\boldsymbol{\gamma} \equiv \boldsymbol{\theta}_0 + \boldsymbol{\Gamma} \left(\nabla \mathbf{f}_0\right)^{\mathsf{T}} \mathbf{C}_0^{-1} (\boldsymbol{\Phi}_O - \mathbf{f}_0)$$

# SELFI-1 Euclid versus BOSS

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



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

Legend:
- $\boldsymbol{\theta}_0$ (prior)
- $\boldsymbol{\gamma}$ (reconstruction)
- $\boldsymbol{\theta}_{\mathrm{gt}}$ (ground truth)
- BOSS NGC $0.2 < z < 0.5$
- BOSS SGC $0.2 < z < 0.5$

# Uncertainty quantification

$$\boldsymbol{\Gamma} \equiv \left[ (\nabla \mathbf{f}_0)^{\mathsf{T}} \mathbf{C}_0^{-1} \nabla \mathbf{f}_0 + \mathbf{S}^{-1} \right]^{-1}$$

Prior covariance matrix

Posterior covariance matrix

# From primordial power spectrum to cosmology



Cosmological parameters

Primordial power spectrum

Taylor-expanded black-box

Arbitrary statistical summaries

Planck priors (marginalised, variance x3)

Cosmic-variance limited Euclid experiment

Two different phase and noise realisations

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

FL, Enzi, Jasche & Heavens 2019, 1902.10149

# 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 on black-box)



```
pip install pyselfi
```

# A black-box: Simbelmynë

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

**Publicly available code:**

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



Initial conditions — Final conditions (dark matter) — Galaxies

**Dark matter simulation with COLA**

Tassev, Zaldarriaga & Eisenstein 2013, 1301.0322

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

# tCOLA: Comoving Lagrangian Acceleration (temporal domain)

- Write the displacement vector as: $\mathbf{s} = \mathbf{s}_{\text{LPT}} + \mathbf{s}_{\text{MC}}$

Tassev & Zaldarriaga 2012, arXiv:1203.5785

- Time-stepping (omitted constants and Hubble expansion):

**Standard**:                                    **Modified**:
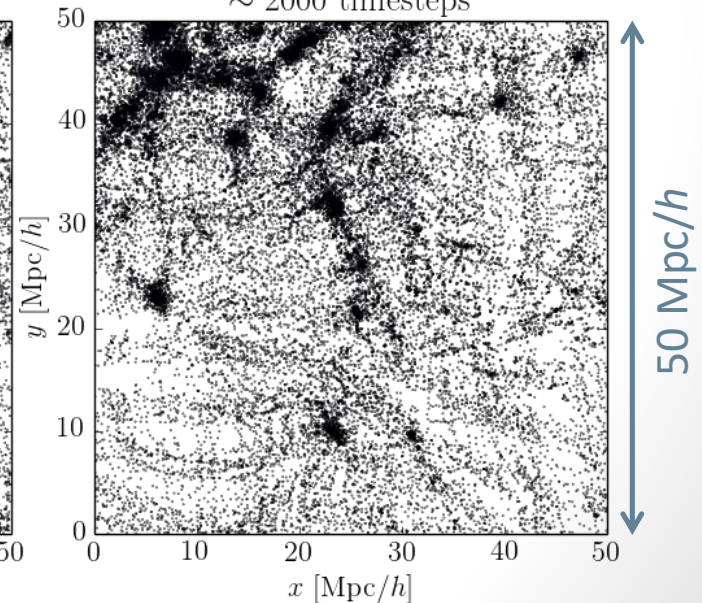
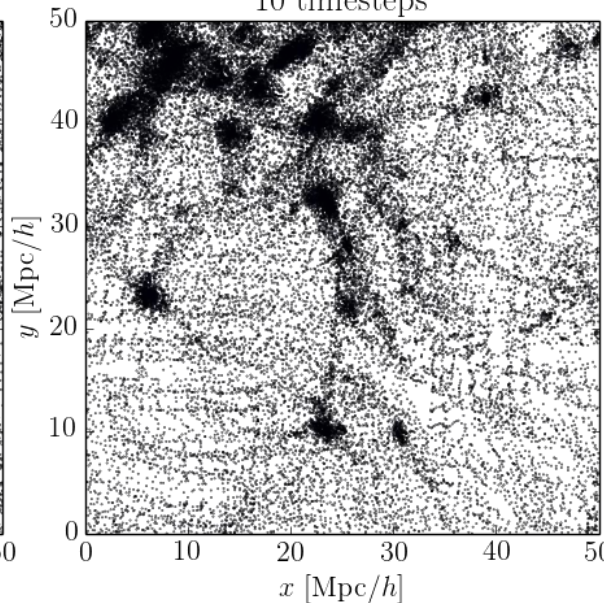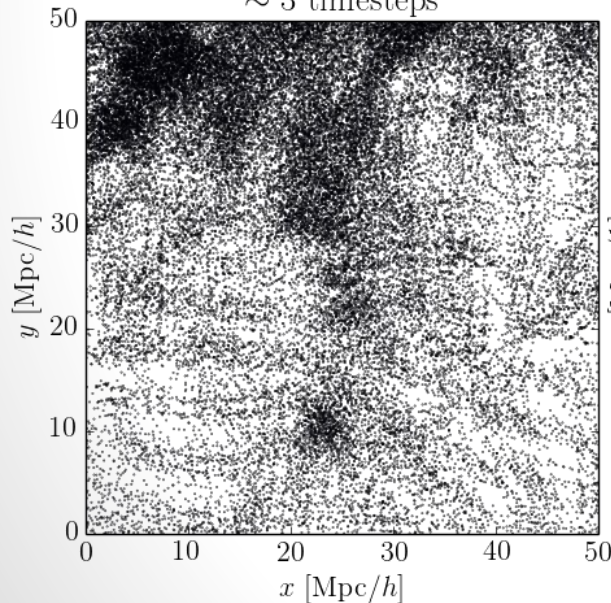$$\partial_\tau^2 \mathbf{s} = -\nabla\Phi \implies \partial_\tau^2 \mathbf{s}_{\text{MC}} = \partial_\tau^2(\mathbf{s} - \mathbf{s}_{\text{LPT}}) = -\nabla\Phi - \partial_\tau^2 \mathbf{s}_{\text{LPT}}$$
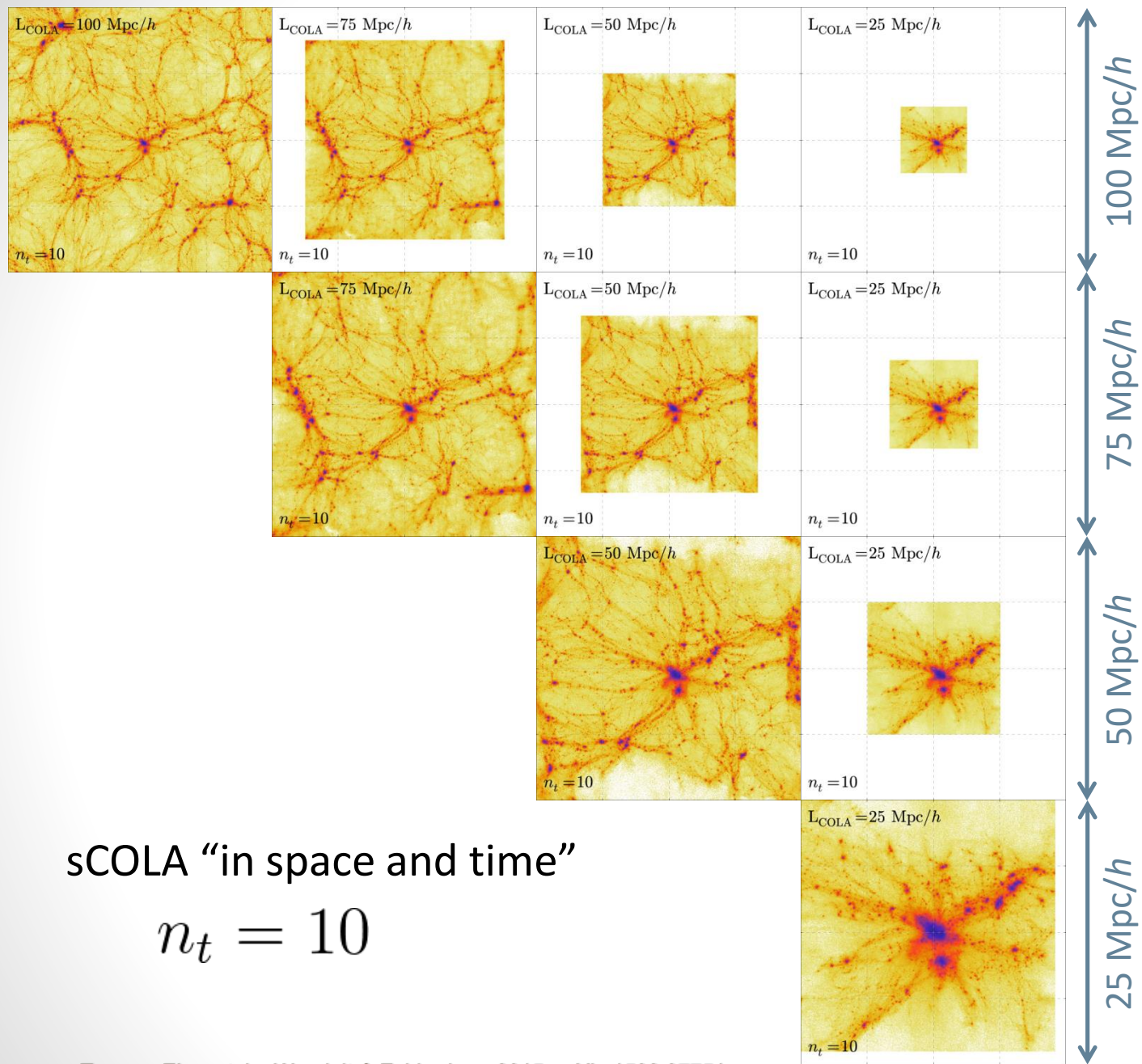


Tassev, Zaldarriaga & Einsenstein 2013, arXiv:1301.0322

sCOLA:
Extension to
the spatial
domain

100 Mpc/h

75 Mpc/h

50 Mpc/h

25 Mpc/h

sCOLA "in space and time"

$$n_t = 10$$

# Perfectly parallel simulations with sCOLA tiling

# Concluding thoughts

- **Goal:** developing an algorithm 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!

- Likelihood-free inference is an easy way to account for cosmological synergies.

- The "number of parameters route" beyond likelihood-free rejection sampling (SELFI):

  - High-dimensional likelihood-free problems can be addressed.
  - The computational workload is fixed *a priori* and perfectly parallel.

- SELFI allows inference of the primordial power spectrum and cosmological parameters.