

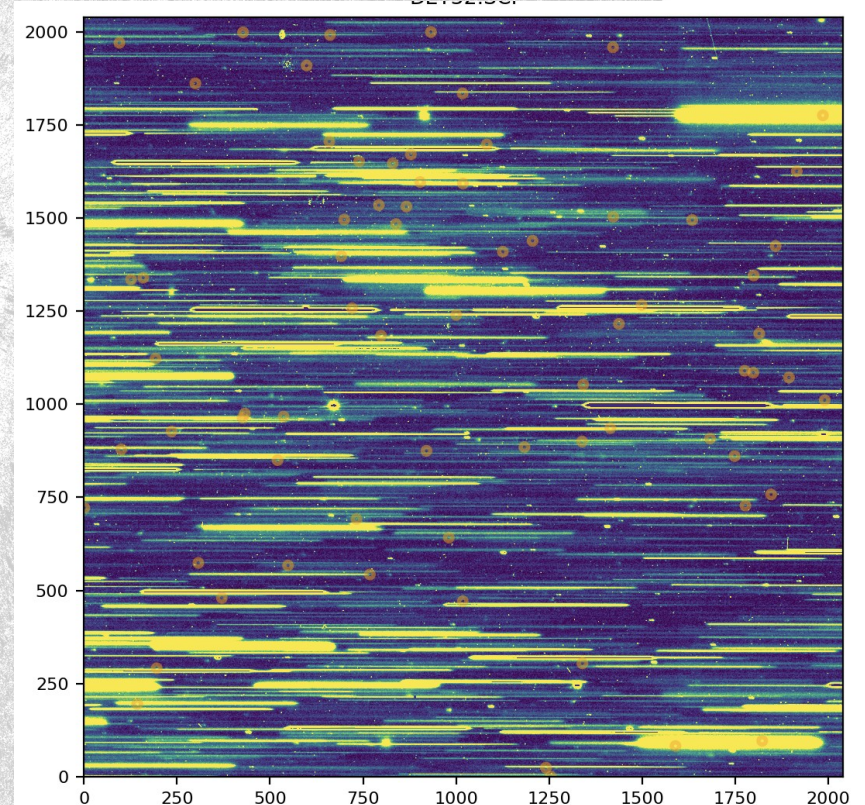
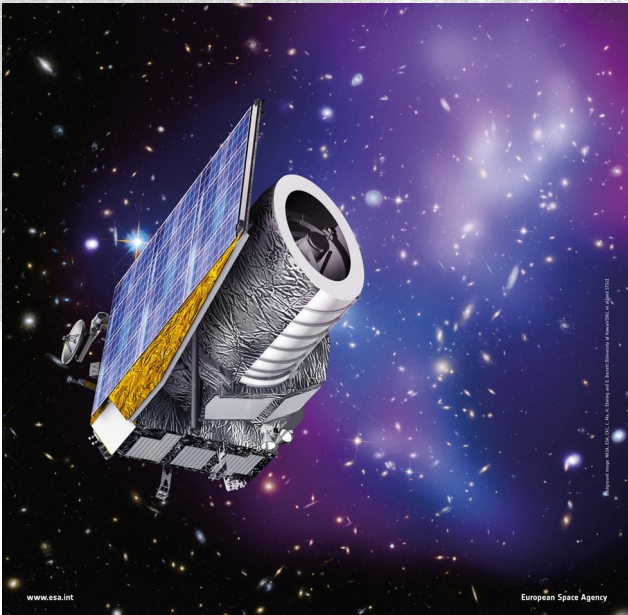
# Mocking the Euclid spectroscopic survey

P. Monaco, Università di Trieste, INAF-OATs, INFN, IFPU

credit: Ben Granett + OU-SIM

Role in Euclid:

co-lead of Observational Systematics WP, Galaxy Clustering SWG  
co-lead of Simulations for Galaxy Clustering WP, CosmoSim SWG  
lead of covariance Processing Functions for OU-LE3





# Mocking the Euclid spectroscopic survey

P. Monaco, Università di Trieste, INAF-OATs, INFN, IFPU

Role in Euclid:

co-lead of Observational Systematics WP, Galaxy Clustering SWG  
co-lead of Simulations for Galaxy Clustering WP, CosmoSim SWG  
lead of covariance Processing Functions for OU-LE3

## Requirements for a simulation:

Sampling BAO scales with subpercent precision:

**box size > 1 Gpc/h**

Past light cone to  $z \sim 2$  with no (few) replications:

**box size  $\gtrsim 4$  Gpc/h**

Sampling the largest scales (PNGs, SSC etc.):

**box size  $\gg 1$  Gpc/h**

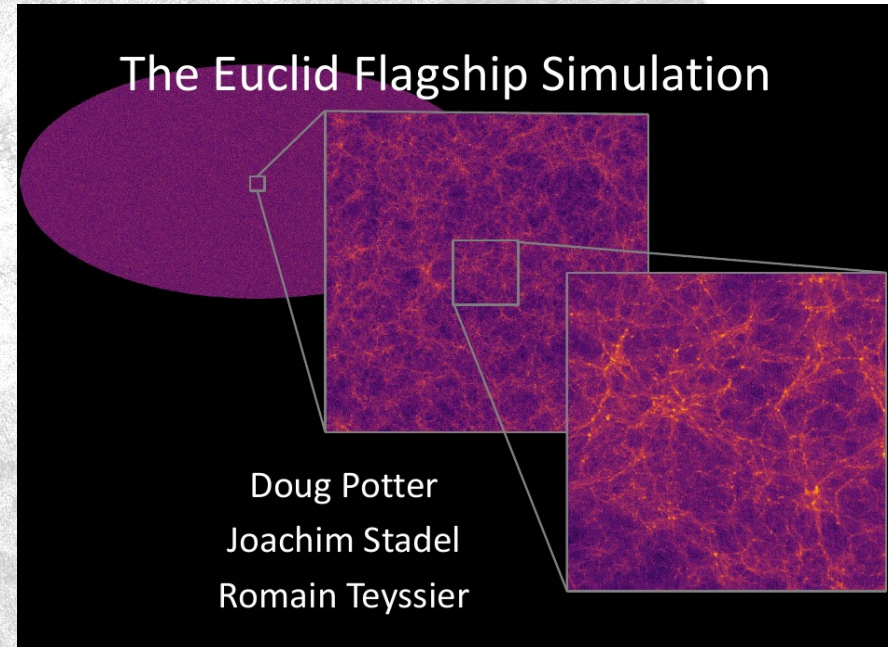
Sampling halos for Euclid spectroscopic sample:

**$M_{\text{part}} \lesssim 3e10 M_{\text{sun}}/h$**

Sampling halos for Euclid photometric sample:

**$M_{\text{part}} \lesssim 3e9 M_{\text{sun}}/h$**

-> at least  $\sim 6000^3$  particles realizations





# The number of needed mocks for the covariance matrix

Noise in the CM due to sampling with a finite number of mocks propagates to the parameter covariance (errorbars) as a function of:

- $N_b$  length of the data vector (number of bins),
- $N_p$  number of parameters
- $N_s$  number of simulations (mocks)

$$\langle \Delta\theta_i \Delta\theta_j \rangle = f \langle \Delta\theta_i \Delta\theta_j \rangle_{\text{ideal}}$$

$$f = 1 + B (N_b - N_p)$$

$$B = \frac{(N_s - N_b - 2)}{(N_s - N_b - 1)(N_s - N_b - 4)}$$

$N_b$  = Number of bins in  $\mathbf{D}$

$N_p$  = number of parameters

(credits: A. Sanchez, L. Blot)

In Euclid we aim at **3500 simulations**, good for two cases:

- $N_b = 600$  with required 10% precision in the errorbar,
- $N_b = 300$  with required 5% precision in the errorbar.

Advanced techniques (tapering, eigenvalue decomposition, shrinkage, fitting a model for the CM, denoising with ML...) will weaken this constraint, higher orders or cross correlation will increase the data vector.



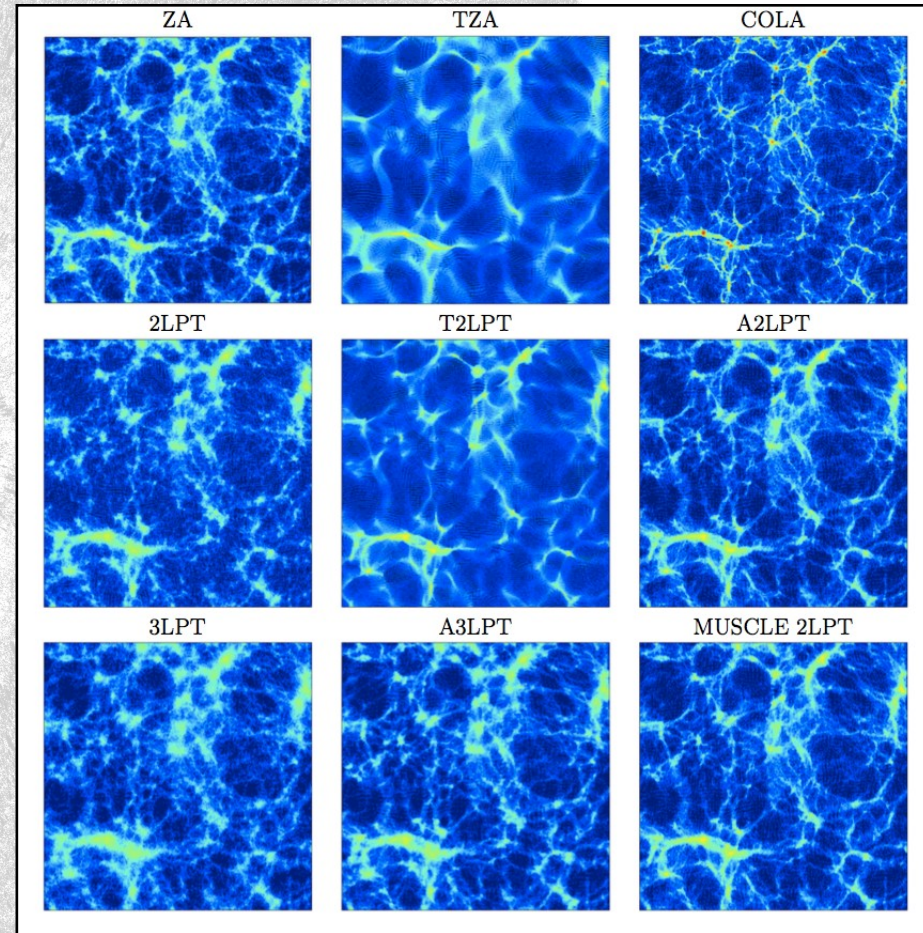
# Approximate methods

Several methods have been proposed to generate catalogs of DM halos with approximate methods (or fast simulations). These can be divided into two broad classes:

- **Lagrangian**, or predictive (peak-patch, pinocchio, cola, fast-pm...)
- **Bias-based**, or calibrated (log-normal, patchy etc., ez-mocks, **BAM**...)

## Criterion for validation:

the relative change of parameter errorbars induced by the use of approximate methods should be  $<10\%$





# PINOCCHIO

PINpointing Orbit Crossing-Collapsed Hierarchical Objects (Monaco+ 2002 ... Munari+2017)

## Computing the “collapse” (OC) time of a fluid element

- Taylor expansion of the gravitational potential:

$$\phi(\vec{q}_0) \simeq \cancel{\phi_0} + \underbrace{\phi_{,i}(\vec{q}_0)(\vec{q} - \vec{q}_0)_i}_{\text{bulk flow}} + \underbrace{\phi_{,ij}(\vec{q}_0)(\vec{q} - \vec{q}_0)_{ij}}_{\text{second-order term}}$$

- => evolution of a homogeneous ellipsoid
- numerical solution, or
- solution with 3LPT up to **ORBIT CROSSING**
- + correction for quasi-spherical cases

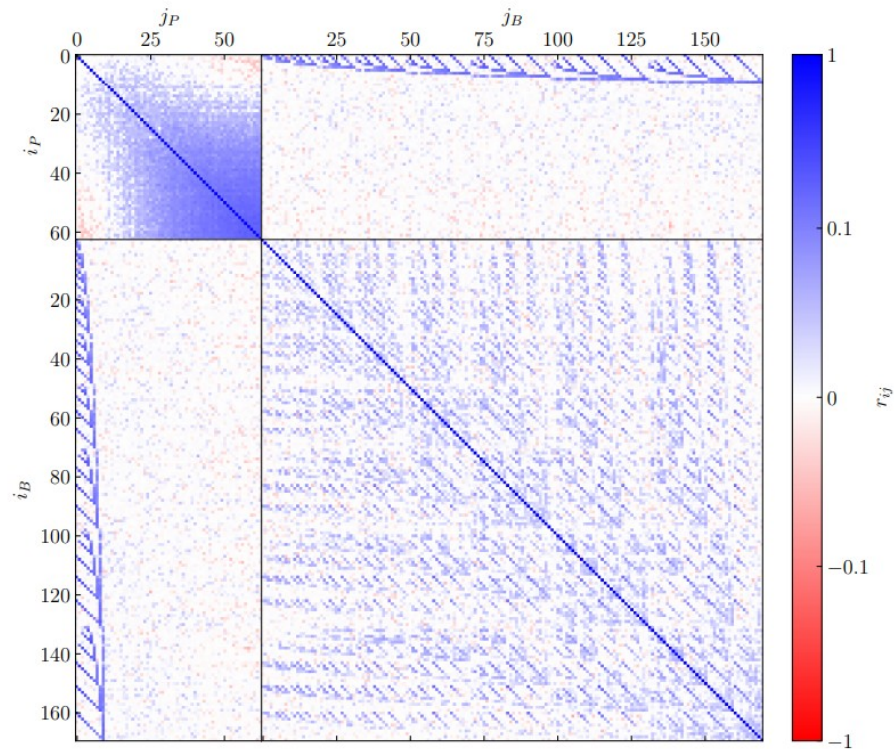
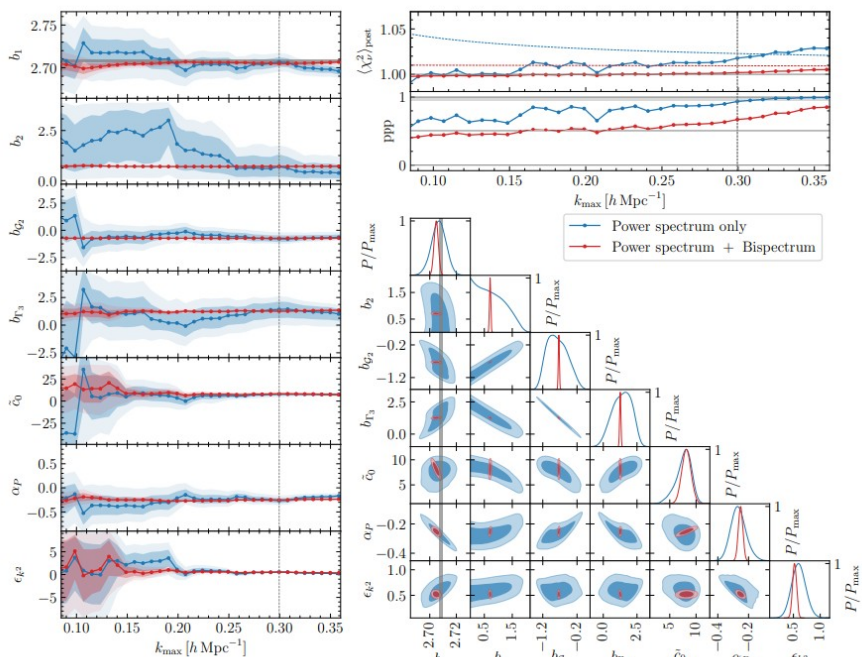
(P.M. 1995, 1997a)





# Covariances for a $P(k) + B(k)$ analysis

Oddo+2021 used 10000 Pinocchio simulations in periodic boxes to construct a CM for fitting  $P(k)$  and  $B(k)$  in real space, demonstrating the power of the joint analysis to constrain bias parameters.



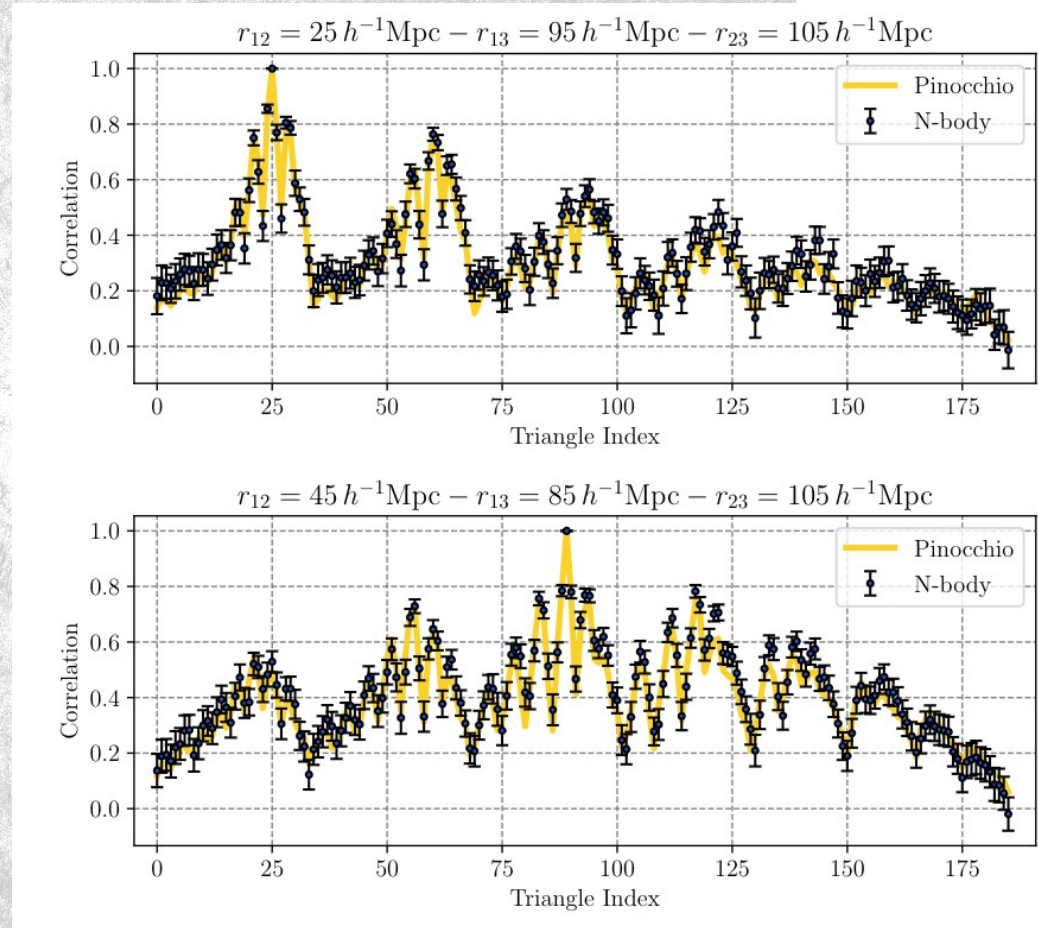


# Covariances for a 3PCF analysis

Veropalumbo+2022 (in prep.) how the CM constructed from 298 N-body simulations and 298 pinocchio simulations run on the same initial conditions compare. Here the sampling noise is exactly the same for the two sets, so the comparison is expected to be better than the (jackknife) errorbars.

## Reason of success:

- 3LPT displacements (that represent the 4-point function) applied to sets of particles in Lagrangian space that approximate well halo particles,
- i.e. it's better than 1-loop PT!



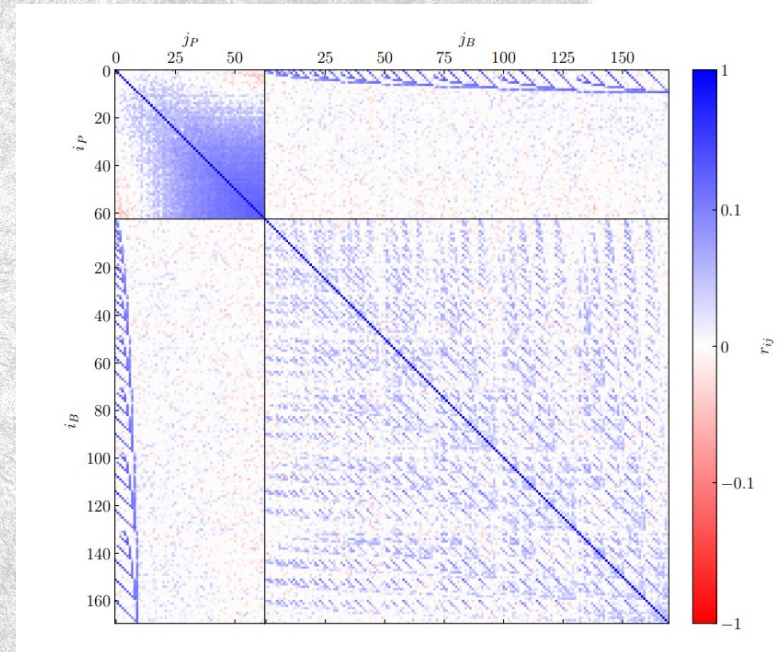


# Available mocks

## Minerva-like (ThirdEuclidMocks)

Simulation performed with PINOCCHIO v4.1.3:

- WMAP-like cosmology
- 1500 Mpc/h box sampled with  $1000^3$  particles
- $M_p = 3.8e11 M_{\text{sun}}/h$ , smallest halo  $M_h = 1.15e13 M_{\text{sun}}/h$  (30 part.)
- **10000 realizations available**
- no lightcones
- outputs at  $z=2, 1.8, 1.35, 1, 0.9, 0.57, 0.32$  and 0
- storage: 8.3 TB
- used in many papers



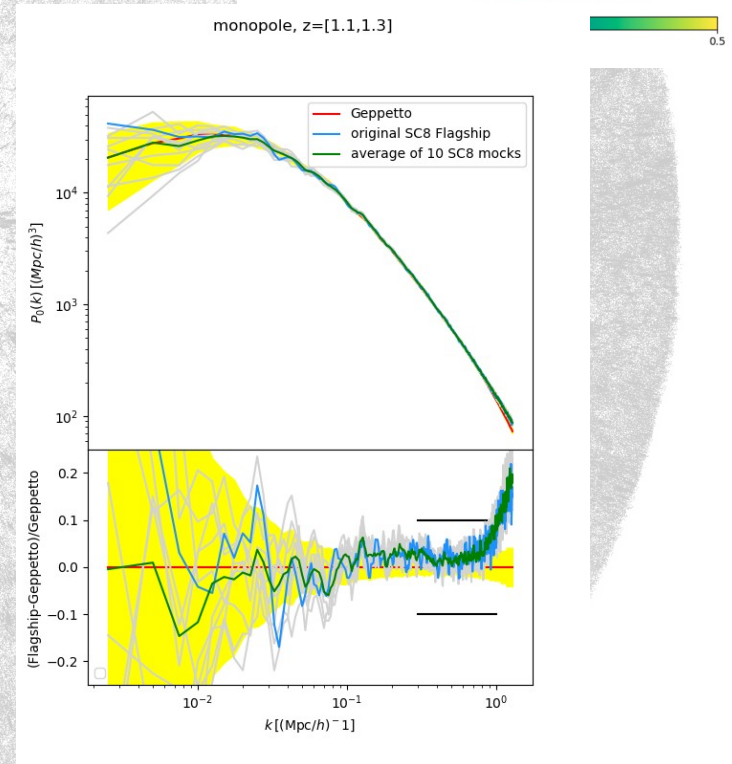
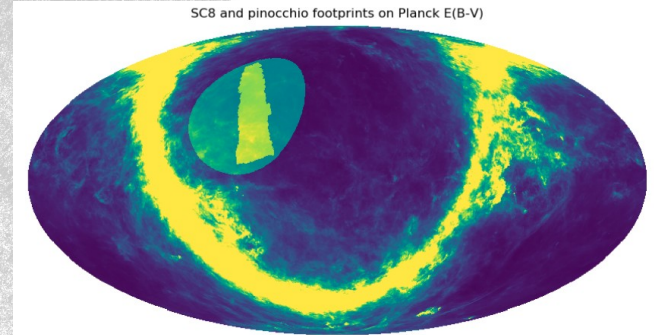


# Available mocks

## GeppettoFC

Simulation performed with PINOCCHIO v4.1.3:

- LambdaCDM cosmology similar to Flagship 1
- 1200 Mpc/h box sampled with  $2160^3$  particles
- $M_p = 1.5e10 M_{\text{sun}}/h$ , smallest halo  $M_h = 1.5e11 M_{\text{sun}}/h$  (10 part.)
- **600 realizations available, 3500 planned**
- light-cone covering a circle of radius 30 deg, 2763 sq deg, starting at  $z=2$  (containing SC8 wide field)
- outputs at  $z=2, 1.8, 1.35, 1, 0.9, 0.5$  and 0 + lightcone + histories
- storage: 37 Gb each (17 Gb for the lightcone)
- used for SC8 mocks





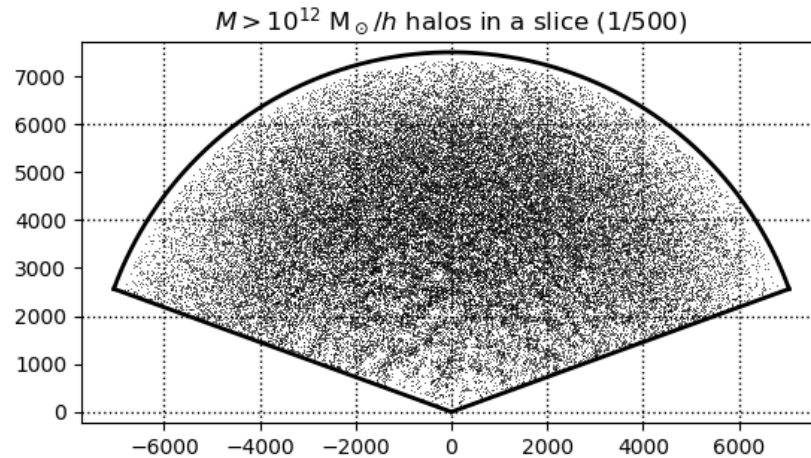
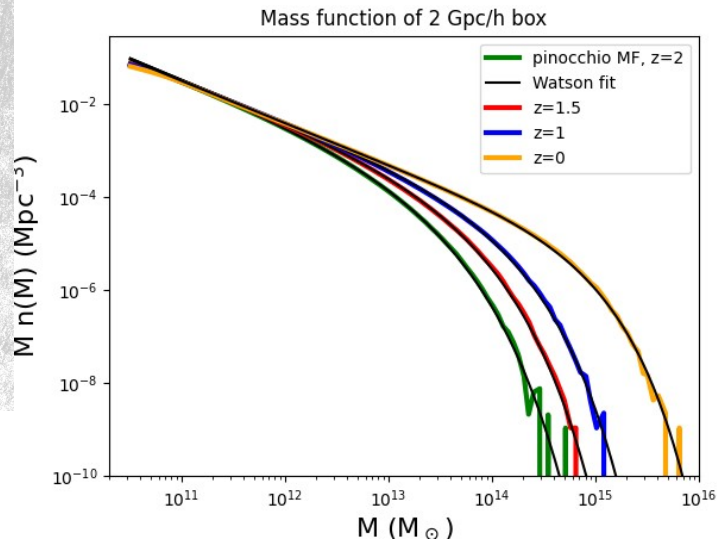
# Available mocks

## EuclidMocks2000

Simulation performed with PINOCCHIO v5.0:

- LambdaCDM cosmology similar to Flagship 1
- 2000 Mpc/h box sampled with  $6144^3$  particles
- $M_p = 3e9 M_{\text{sun}}/h$ , smallest halo  $M_h = 3e10 M_{\text{sun}}/h$  (10 part.), same as the Flagship 1 simulation
- **50 realizations, target: 500**
- light-cone covering  $\sim$ half sky starting at  $z=4$
- central l.o.s. aligned with the  $z$  axis
- 2 outputs ( $z=1$  and  $z=0$ ) + lightcone + histories,
- storage: 573 Gb (450 Gb for the lightcone)
- issue with replications?

computer time provided by [INFN-Euclid](#)



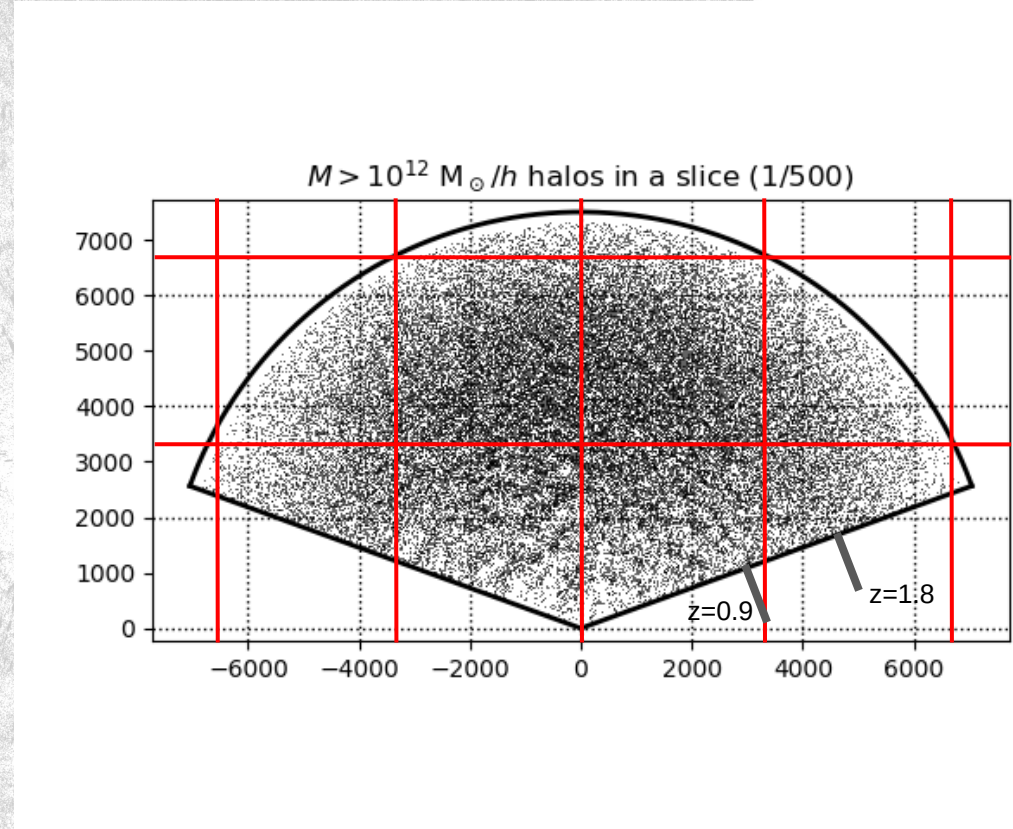


# Available mocks

## EuclidMocks4000

Simulation performed with PINOCCHIO v5.0:

- LambdaCDM cosmology similar to Flagship 1
- 3380 Mpc/h box sampled with  $6144^3$  particles
- $M_p = 1.5e10 M_{\text{sun}}/h$ , smallest halo  $M_h = 1.5e11 M_{\text{sun}}/h$  (10 part.)
- **217 realizations, target: 3500**
- light-cone covering  $\sim$ half sky starting at  $z=4$
- central l.o.s. aligned with the  $z$  axis **or** with the main diagonal **or** random
- 2 outputs ( $z=1$  and  $z=0$ ) + lightcone + histories,
- storage: 220 Gb (50 Gb for the lightcone)
- 52 TB total



computer time provided by [INFN-Euclid](#)



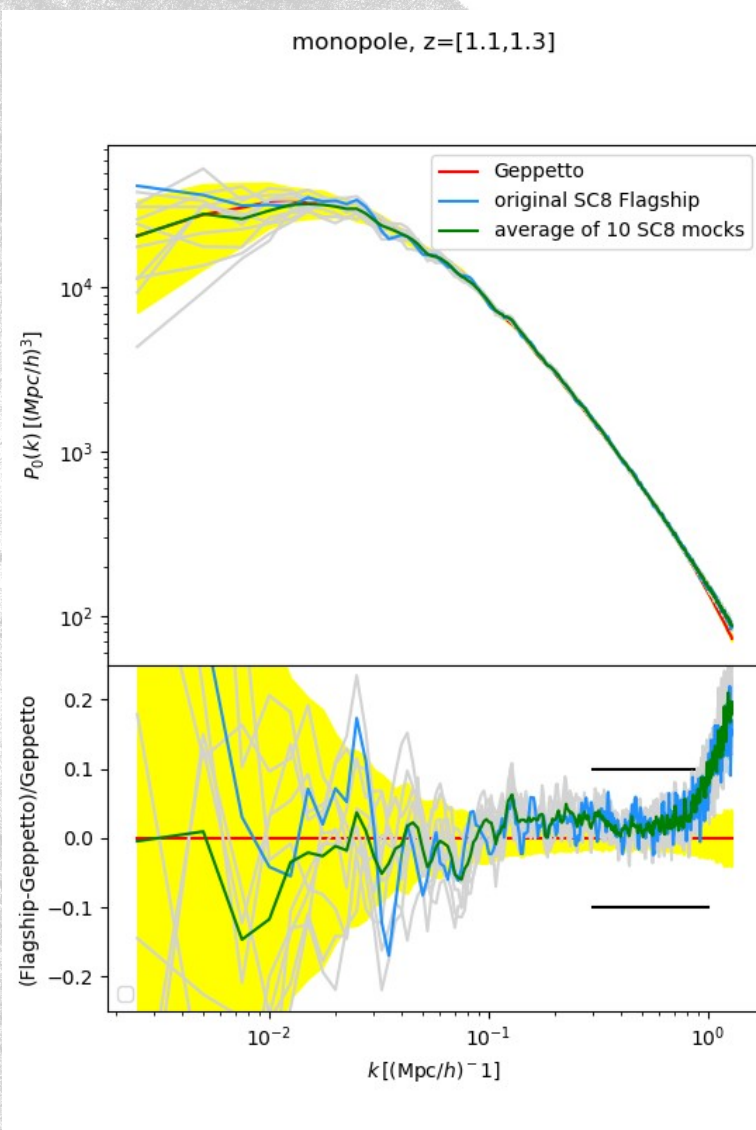
# From halos to galaxies

## Flagship mock galaxy catalog:

- match DM halos to r-band galaxies using an HOD,
- assign SFR, B/T ratio, bulge and disk sizes,
- sample galaxy SEDs according to their properties,
- add emission lines, including internal extinction,
- “predict” observed magnitudes and line fluxes and calibrate them against available observations

## Mocks for the covariance:

- only produce position, redshift, line flux
- “measure” HOD from Flagship to populate halos
- process the Flagship mock with image simulations (**bypasses**)
- compute the probability of detecting a galaxy, already marginalized over galaxy properties, given image noise and exposure time
- use this probability to select galaxies





# From halos to galaxies

## Flagship mock galaxy catalog:

- match DM halos to r-band galaxies using an HOD,
- assign SFR, B/T ratio, bulge and disk sizes,
- sample galaxy SEDs according to their properties,
- add emission lines, including internal extinction,
- “predict” observed magnitudes and line fluxes and calibrate them against available observations

## Mocks for the covariance:

- only produce position, redshift, line flux
- “measure” HOD from Flagship to populate halos
- process the Flagship mock with image simulations (**bypasses**)
- compute the probability of detecting a galaxy, already marginalized over galaxy properties, given image noise and exposure time
- use this probability to select galaxies

## Parallel effort by Galaxy Evolution SWG (De Lucia):

- collect galaxy catalogs from hydro simulations and SAMs,
- project them on the past lightcone (limited areas)
- produce a consistent set of observed magnitudes



# Galaxy Clustering systematics

Assumptions:

+ we restrict to linear galaxy bias

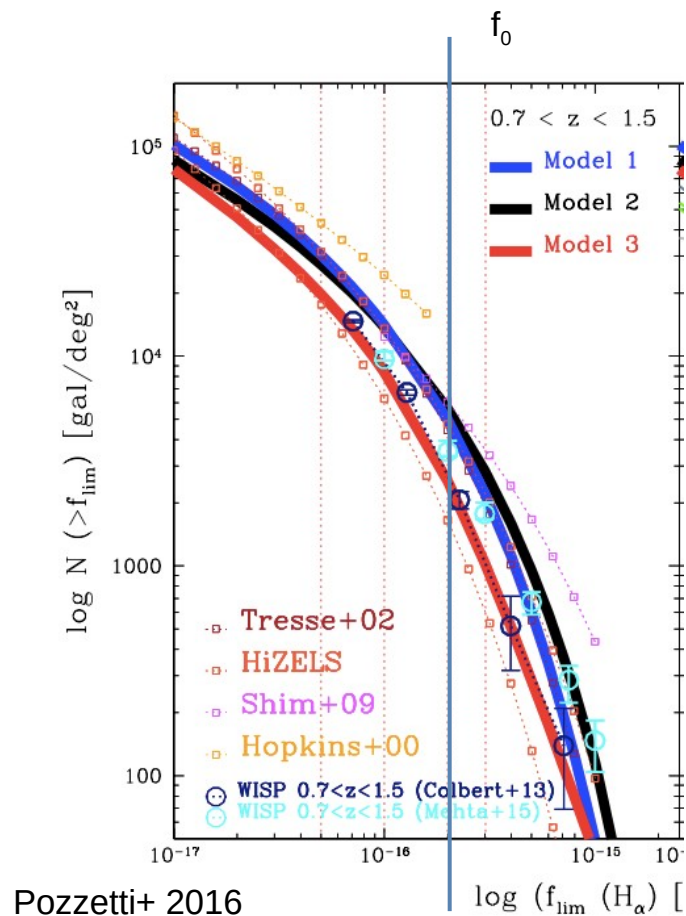
+ galaxy bias **does not depend on luminosity**

then the shape of the luminosity function is universal

We aim at measuring the galaxy number density:

$$n_g = \int_{f_0}^{\infty} [1 + \delta_g(\mathbf{x})] \Phi(f|z) df$$

$f_0$  is a fiducial Halpha flux limit (e.g.  $f_0 = 2e-16$  erg/s/cm<sup>2</sup>).





# Galaxy Clustering systematics

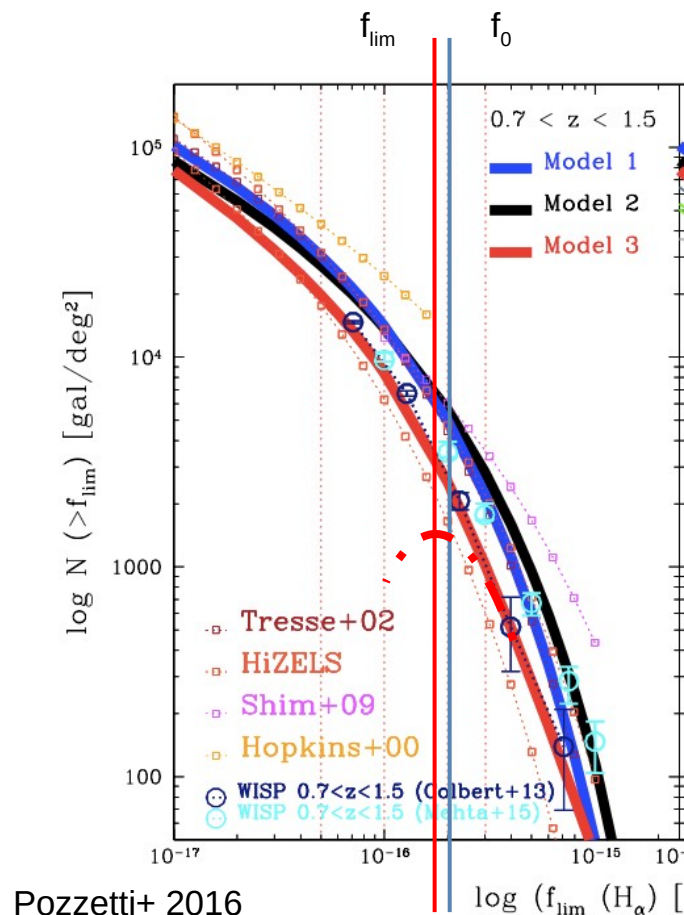
Our measurement involves a completeness function:

$$n_o = \int_0^\infty [1 + \delta_g(\mathbf{x})] \Phi(f|z) \mathcal{C}(f, z, \{N_i\}) df$$

Completeness can be recast in terms of a flux limit  $f_{\text{lim}}$ :

$$n_o = \int_{f_{\text{lim}}(\boldsymbol{\vartheta}, z, \{N_i\})}^\infty [1 + \delta_g(\mathbf{x})] \Phi(f|z) df$$

Here  $\{N_i\}$  represent a generic set of noise terms that determine the probability of a galaxy being detected.



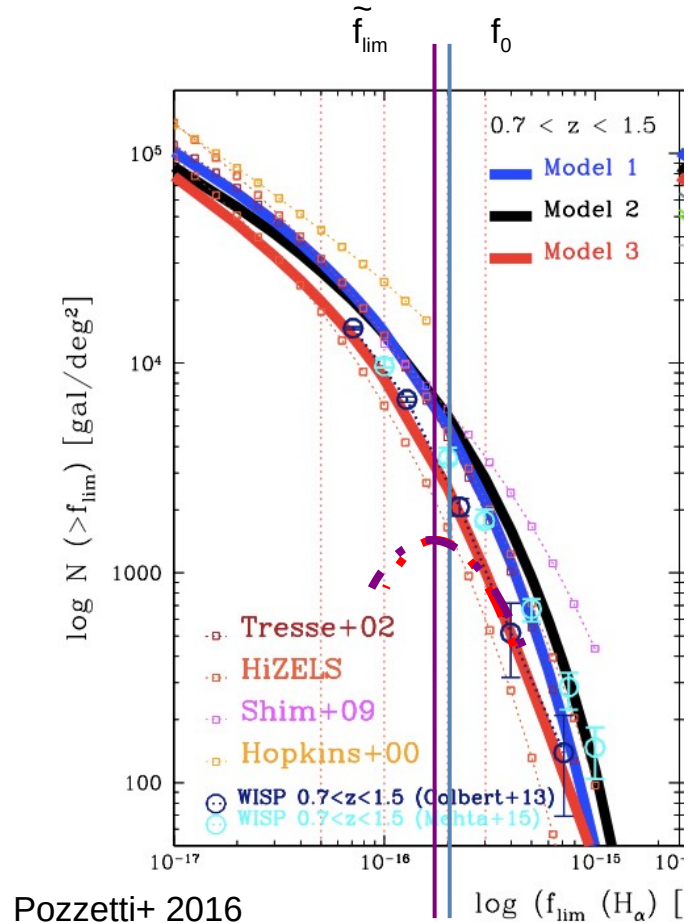


# Galaxy Clustering systematics

Galaxy density is compared with that of a random catalog, created on the basis of an approximation of the completeness function (=visibility mask):

$$n_r = N_r \int_0^\infty \Phi(f|z) \tilde{C}(f, z, \{N_i\}) df$$

$$= \int_{\tilde{f}_{\text{lim}}(\boldsymbol{\vartheta}, z, \{N_i\})}^\infty \Phi(f|z) df$$





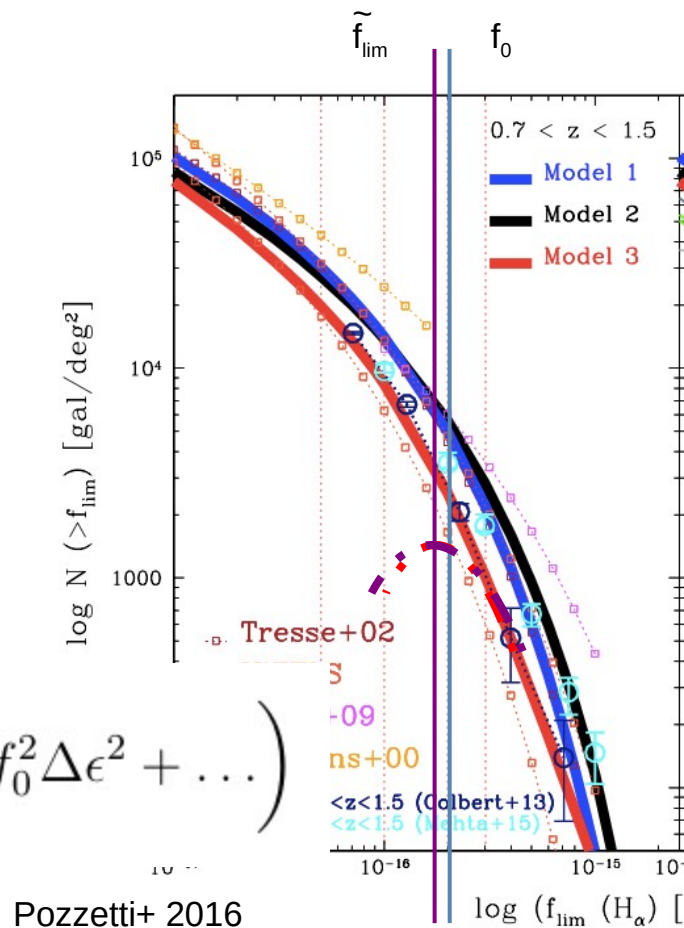
# Galaxy Clustering systematics

The observed density contrast can then be expressed in terms of the galaxy (true) density contrast and the uncertainty in the flux limit ( $\alpha = 1/N_r$ ):

$$f_{\text{lim}}(\boldsymbol{\vartheta}, z, \{N_i\}) = f_0 [1 + \epsilon(\boldsymbol{\vartheta}, z, \{N_i\})]$$

$$\Delta\epsilon = \tilde{\epsilon} - \epsilon$$

$$n_o - \alpha n_r = \delta_g \alpha n_r + (1 + \delta_g) \left( 1 + \Phi f_0 \Delta\epsilon + \frac{1}{2} \frac{d\Phi}{df} f_0^2 \Delta\epsilon^2 + \dots \right)$$





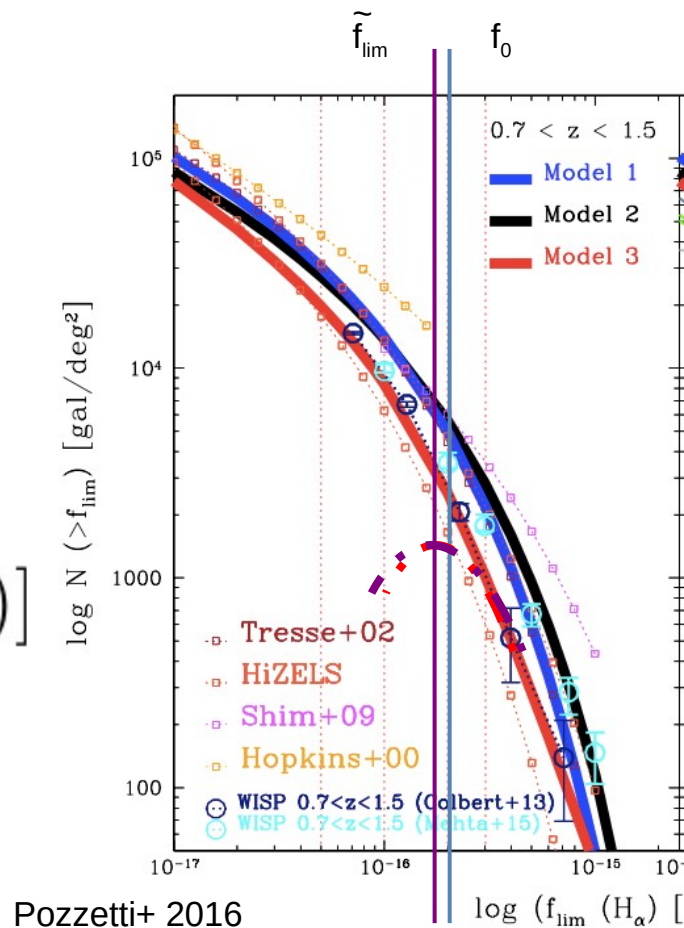
# Galaxy Clustering systematics

This means that the measured signal **is not** a sum of cosmic signal + noise:

$$\delta_o \neq \delta_g + A$$

$$1 + \delta_o(\mathbf{x}) = [1 + \delta_g(\mathbf{x})] [1 + A(\boldsymbol{\vartheta}, z)]$$

$$A = \Phi f_0 \Delta\epsilon + \frac{1}{2} \frac{d\Phi}{df} f_0^2 \Delta\epsilon^2 + \dots$$





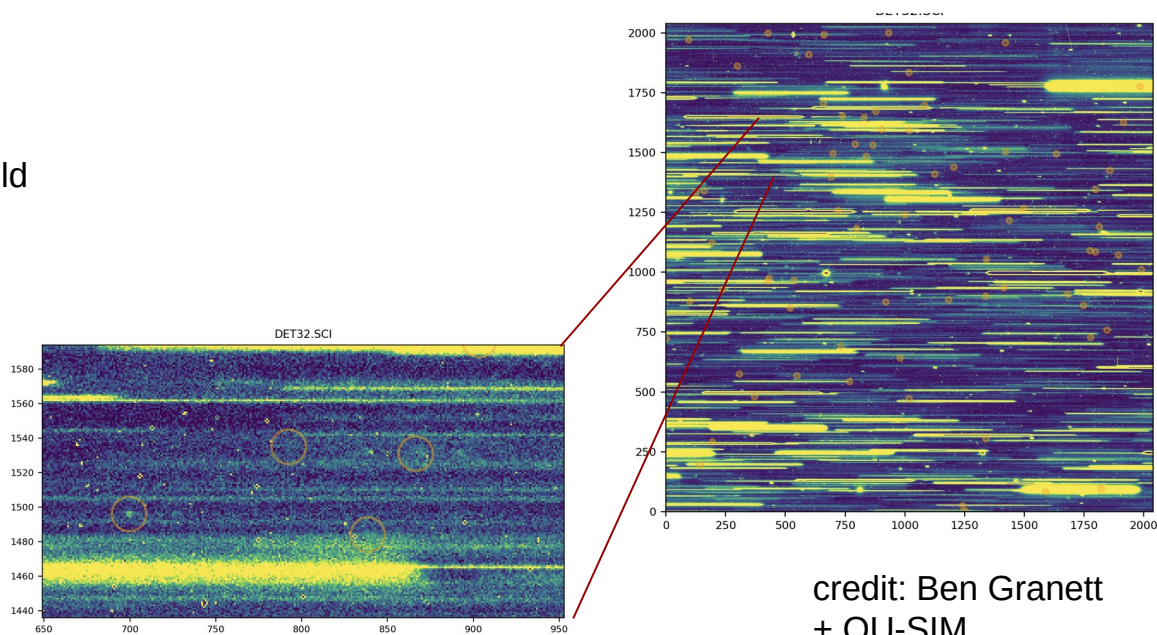
# Mitigation through the random catalog

## Construction of the random catalog

- place a galaxy in a random sky position
- measure the noise level at that position
- add galaxy properties drawn from the deep field
- use a bypass (pypelid) to estimate its SNR
- map SNR to detection probability
- decide if the galaxy gets into the random

Degrees of freedom:

- random error on noise
- calibration error
- assigning properties to the galaxy
- from detection probability to getting or not in the random



credit: Ben Granett  
+ OU-SIM



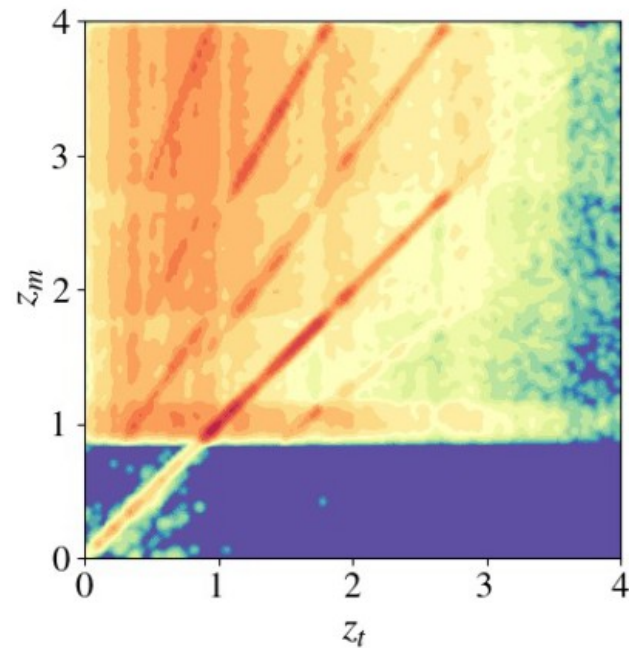


# Redshift errors

observed galaxies =

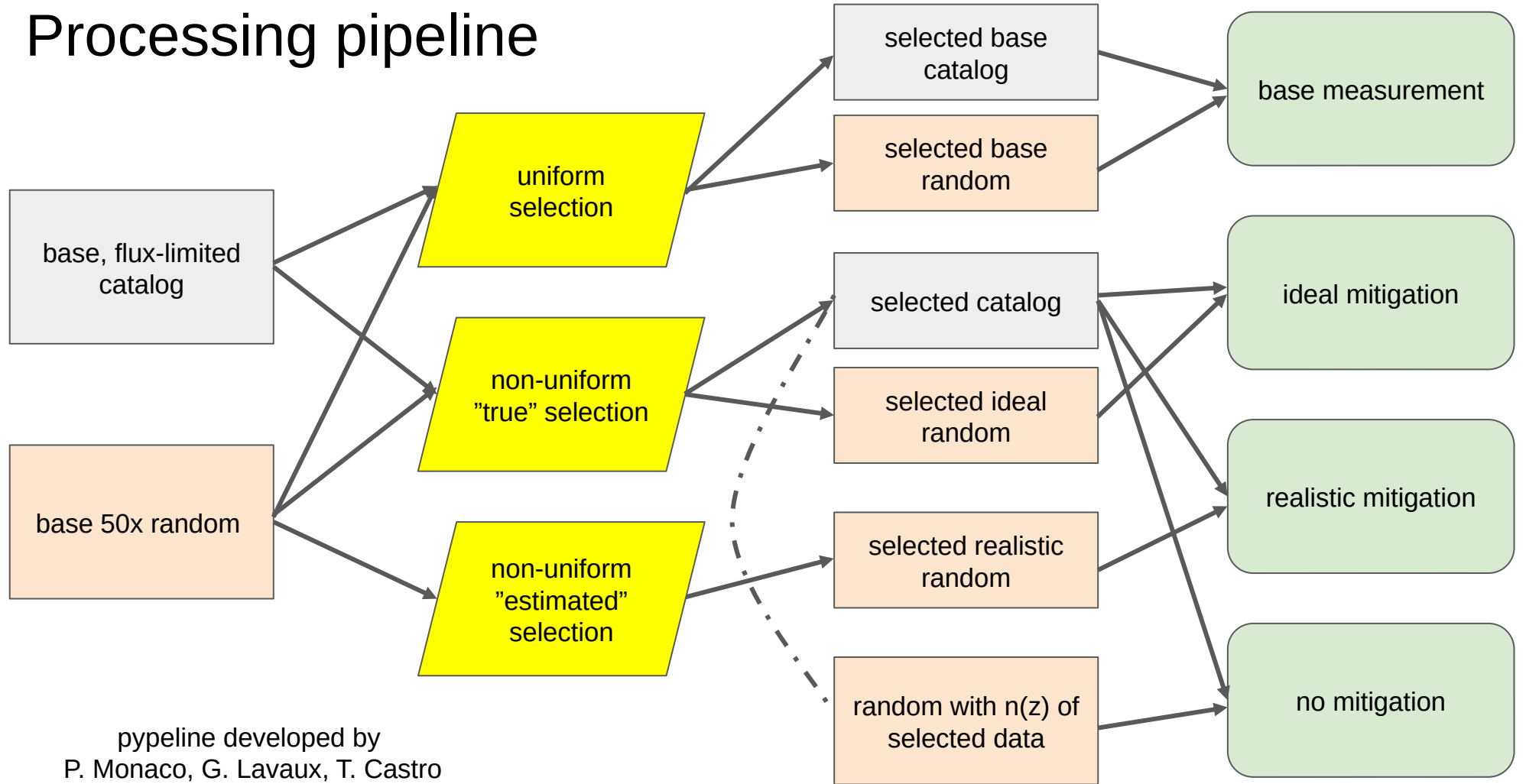
- good galaxies
- + line misidentifications -> rescaled correlation
- + noise spikes -> **uncorrelated?**
- + galactic objects -> correlated like MW

Line misidentifications can be treated at the likelihood level.  
Noise spikes and galactic objects must be represented in the random



credit: Sylvain De La Torre  
+ OU-SPE

# Processing pipeline





# Galaxy Clustering covariance

This propagates to the two-point CF (one more term) in and covariance, with **74 more terms** to add:

$$\xi_o = \xi_g(1 + 2\langle A \rangle) + \xi_A + \xi_g \xi_A$$

$$\begin{aligned}
 C_o &= \langle (\xi_{o12}^{(r)} - \xi_{o12})(\xi_{o34}^{(r)} - \xi_{o34}) \rangle \\
 &= C_g + C_A \\
 &\quad + \xi_{g13}\xi_{A34} + \xi_{g14}\xi_{A23} + \xi_{g23}\xi_{A14} + \xi_{g24}\xi_{A13} \\
 &\quad + \langle \delta_{g1}\delta_{g2}\delta_{g3} \rangle [\langle A_4 \rangle + \xi_{A14} + \langle A_1 A_2 A_4 \rangle] + \text{perm.} && 28 \text{ terms} \\
 &\quad + \langle \delta_{g1}\delta_{g2}\delta_{g3}\delta_{g4} \rangle [\langle A_1 \rangle + \xi_{A12} + \langle A_1 A_2 A_3 \rangle] + \langle A_1 A_2 A_3 A_4 \rangle + \text{perm.} && 19 \text{ terms} \\
 &\quad + \xi_{g12} [\langle A_1 A_3 A_4 \rangle + \langle A_1 A_2 A_3 A_4 \rangle] + \text{perm.} && 18 \text{ terms} \\
 &\quad - \xi_{g12}\xi_{g34}\xi_{A12} - \xi_{g12}\xi_{g34}\xi_{A34} - \xi_{g12}\xi_{A12}\xi_{A34} - \xi_{g34}\xi_{A12}\xi_{A34} \\
 &\quad - \xi_{g12}\xi_{g34}\xi_{A12}\xi_{A34}
 \end{aligned}$$

# Galaxy Clustering covariance

We worked out the expression for the power spectrum in Colavincenzo et al. (2017):

$$C_{ij} \equiv \text{cov}[\hat{P}(k_i), \hat{P}(k_j)] = \langle \delta \hat{P}(k_i) \delta \hat{P}(k_j) \rangle$$

$$C_{ij}^{\text{obs}} = C_{ij}^{\text{cosm}} + C_{ij}^{\text{mask}} + C_{ij}^{\text{mixed}}$$

$$\begin{aligned} C_{ij}^{\text{mixed}} = & \langle \hat{P}_{\text{conv}}(k_i) \hat{P}_{\text{conv}}(k_j) \rangle - \langle \hat{P}_{\text{conv}}(k_i) \rangle \langle \hat{P}_{\text{conv}}(k_j) \rangle + \\ & \langle \hat{P}_{\text{cosmo}}(k_i) \hat{P}_{\text{conv}}(k_j) \rangle - \langle \hat{P}_{\text{cosmo}}(k_i) \rangle \langle \hat{P}_{\text{conv}}(k_j) \rangle + \\ & \langle \hat{P}_{\text{mask}}(k_i) \hat{P}_{\text{conv}}(k_j) \rangle - \langle \hat{P}_{\text{mask}}(k_i) \rangle \langle \hat{P}_{\text{conv}}(k_j) \rangle + \\ & \langle \hat{P}_{\text{cosmo}}(k_i) \hat{G}(k_j) \rangle + \langle \hat{P}_{\text{mask}}(k_i) \hat{G}(k_j) \rangle + \\ & \langle \hat{P}_{\text{conv}}(k_i) \hat{G}(k_j) \rangle + \langle \hat{G}(k_i) \hat{G}(k_j) \rangle \end{aligned}$$

$$\hat{G} = 2\delta_{\mathbf{q}}\delta_{\text{mask},\mathbf{q}} - \delta_{\mathbf{q}}\delta_{\text{conv},\mathbf{q}} + \delta_{\text{mask},\mathbf{q}}\delta_{\text{conv},\mathbf{q}}$$

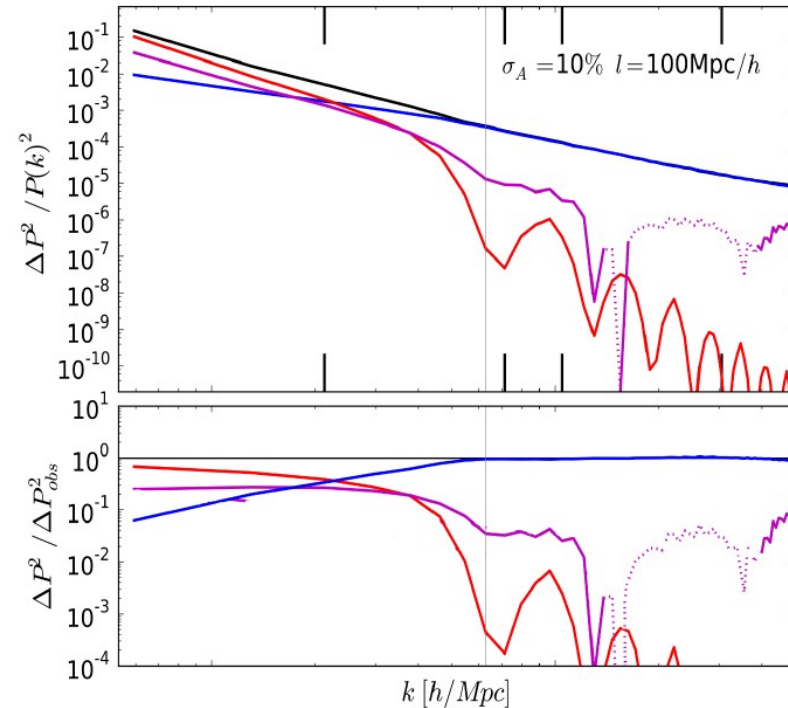
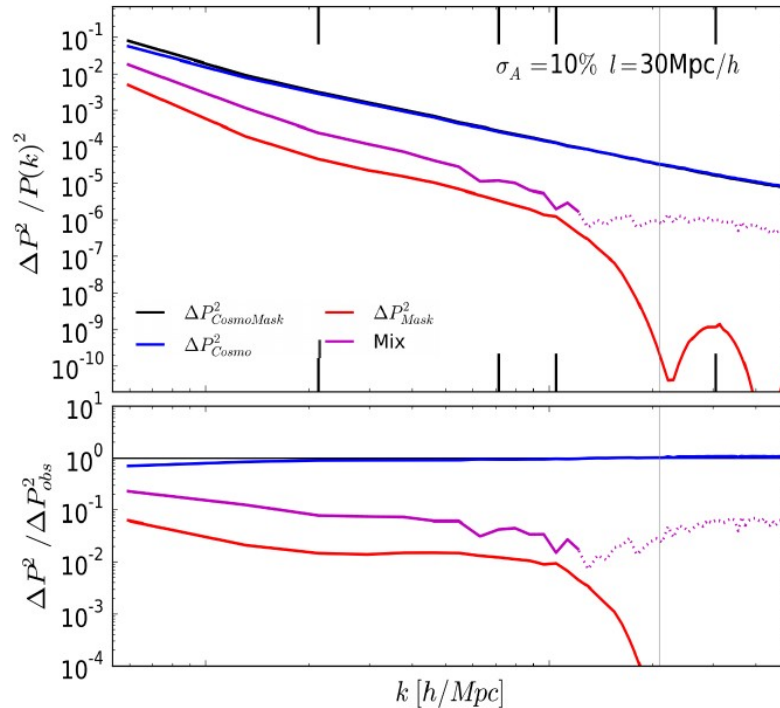


# Galaxy Clustering covariance

## Results: Variance

$$C_{ij}^{\text{obs}} = C_{ij}^{\text{cosm}} + C_{ij}^{\text{mask}} + C_{ij}^{\text{mixed}}$$

10 000 Pinocchio realizations

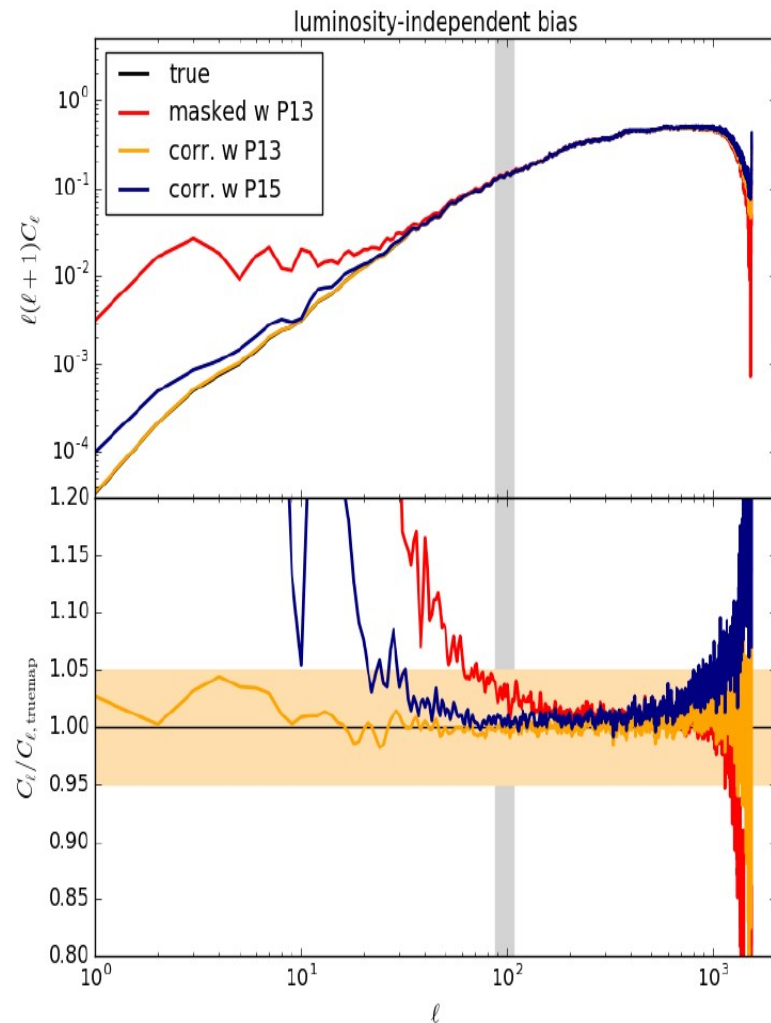


# What is missing

Systematics may have error correlations that are not caught by this scheme.  
Example: MW extinction.

$$1 + \delta_o(\mathbf{x}) = [1 + \delta_g(\mathbf{x})] [1 + A(\boldsymbol{\vartheta}, z)]$$

We need to identify and mitigate “unknown unknowns”.



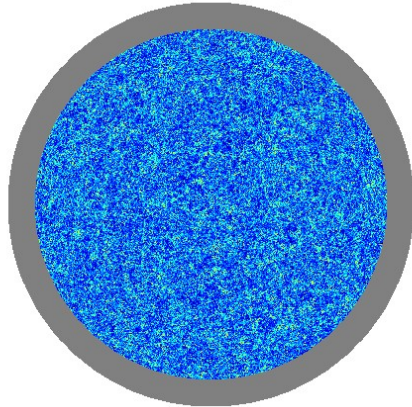


# A blind method to recover the mask of a deep galaxy survey

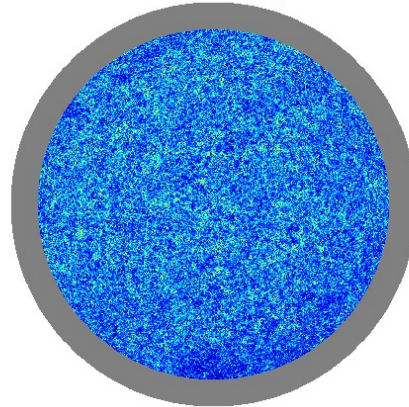
Pierluigi Monaco,<sup>a,b,c</sup> Enea Di Dio<sup>d,e</sup> and Emiliano Sefusatti<sup>b,c</sup>

Journal of Cosmology and Astroparticle Physics, 04, 023 (2019).

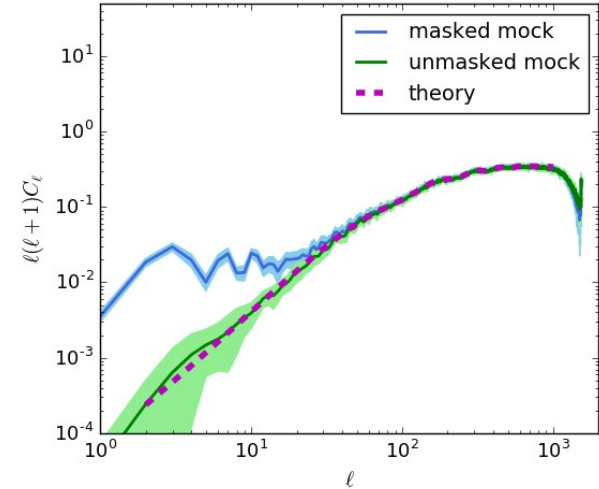
$\delta_o$  of unmasked catalog 10,  $z=1$



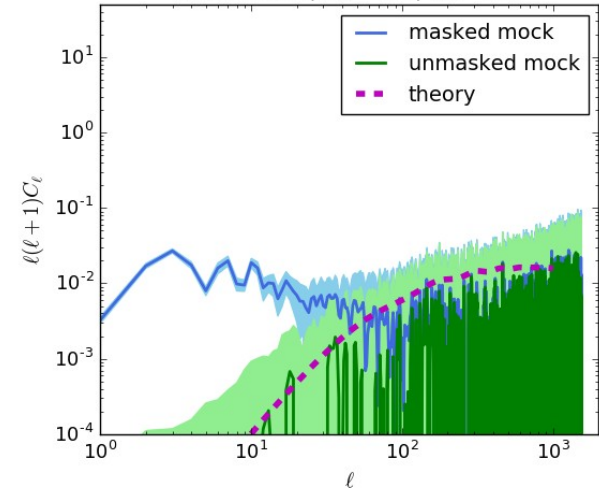
$\delta_o$  of masked catalog 10,  $z=1$



auto correlation,  $z=1.0$



cross corr.,  $z_1=1.0, z_2=2.0$



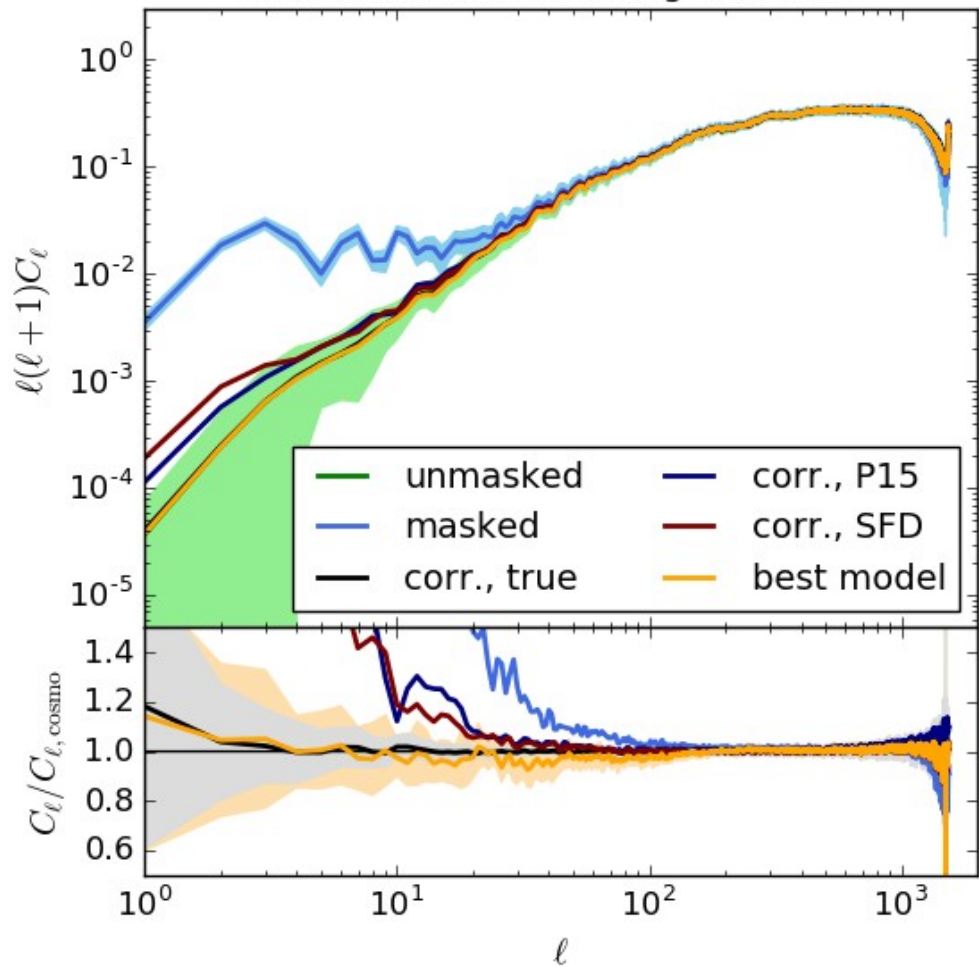
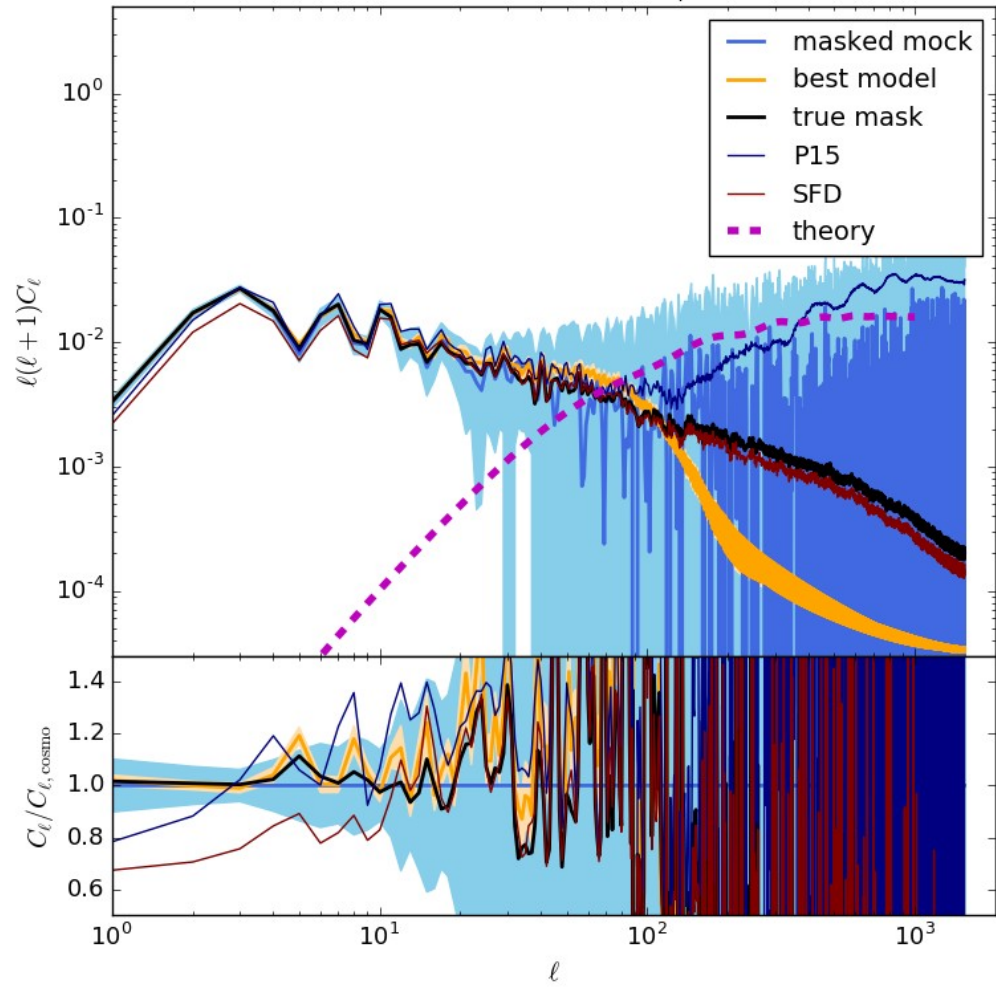
One can construct estimators of the visibility mask by averaging the density contrast over the line of sight, at the 1-point or 2-point level

$$E_{\text{av}}(\boldsymbol{\theta}) \equiv \frac{1}{N_z} \sum_i \delta_r([z_i, \boldsymbol{\theta}]) \simeq -\frac{1}{N_z} \sum_i \frac{S_{Ci}}{S_{Ai}} \mathcal{M} - \frac{1}{N_z} \sum_i \frac{S_{Bi} S_{Ci}^2}{S_{Ai}} \mathcal{M}^2 + O(\mathcal{M}^3). \quad (4.4)$$

$$E_{\text{sq}}(\boldsymbol{\theta}) \equiv \frac{1}{N_p} \sum_i \sum_{j>i} \delta_r([z_i, \boldsymbol{\theta}]) \delta_r([z_j, \boldsymbol{\theta}]) \quad (4.5)$$

$$\simeq \frac{1}{N_p} \sum_i \sum_{j>i} \frac{S_{Ci} S_{Cj}}{S_{Ai} S_{Aj}} \mathcal{M}^2 + \frac{1}{N_p} \sum_i \sum_{j>i} \frac{S_{Ci} S_{Cj} (S_{Bi} S_{Ci} + S_{Bj} S_{Cj})}{S_{Ai} S_{Aj}} \mathcal{M}^3 + O(\mathcal{M}^4).$$



Reconstructed clustering at  $z=1.0$ cross correlation at  $z_1=1.0, z_2=2.0$ 

# Possible synergies

...will we share the mocks?

**but, what about mass resolution?**

would it be convenient to add small halos with a BAM-like model (see De La Torre & Peacock 2012)?