

# Dark Quest. I. Fast and Accurate Emulation of Halo Clustering Statistics and Its Application to Galaxy Clustering

based on arXiv:1811.09504; ApJ accepted + updates

**TAKAHIRO NISHIMICHI (YITP, KYOTO)**

In collaboration with M. Takada, R. Takahashi, K. Osato,  
M. Shirasaki, T. Oogi, H. Miyatake, M. Oguri, R. Murata,  
Y. Kobayashi, N. Yoshida (HSC weak lensing WG)

# THE GAME

- ▶ Efficiently cover the space
- ▶ Controlled error

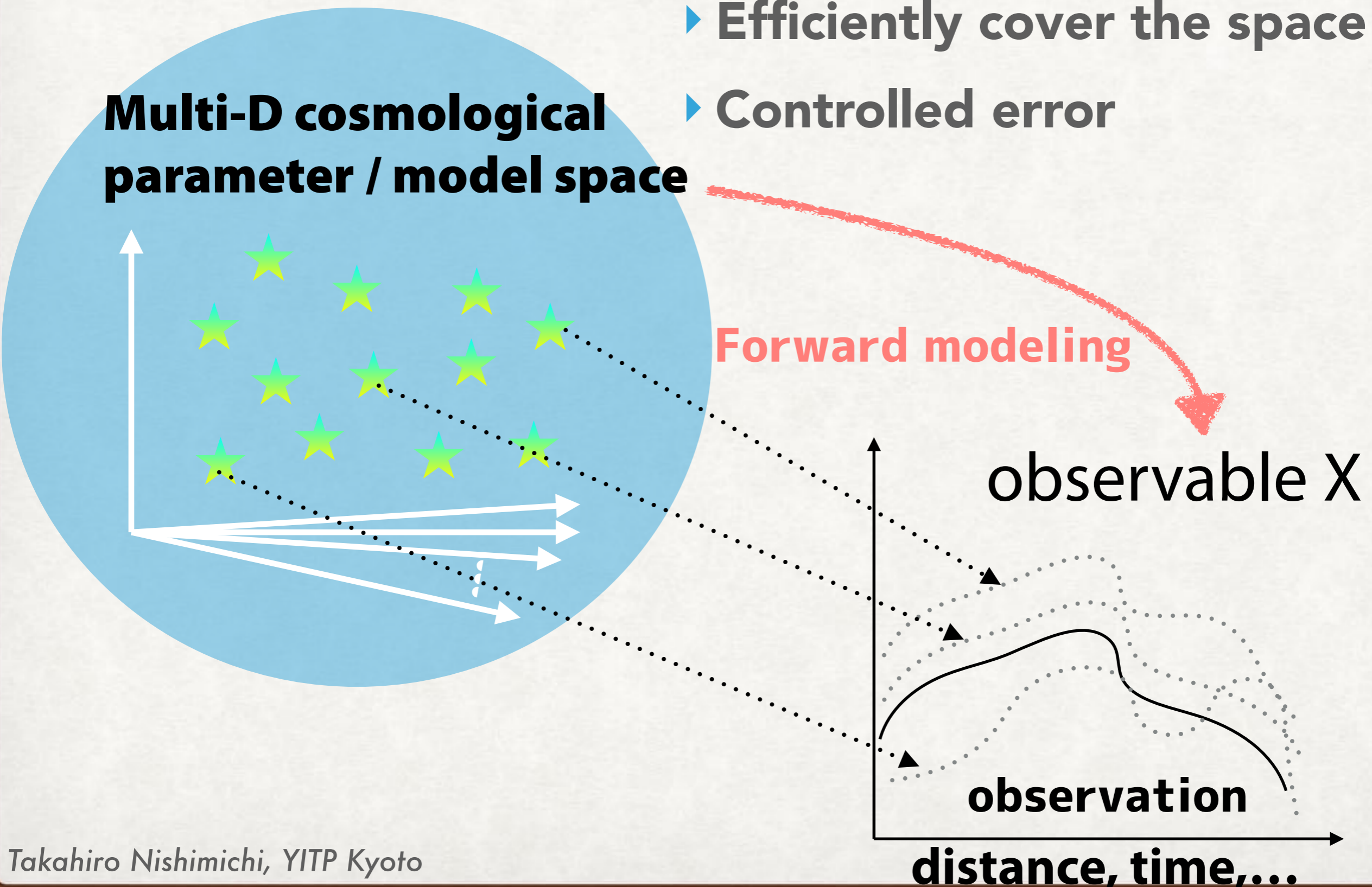
**Multi-D cosmological  
parameter / model space**

**Forward modeling**

**observable X**

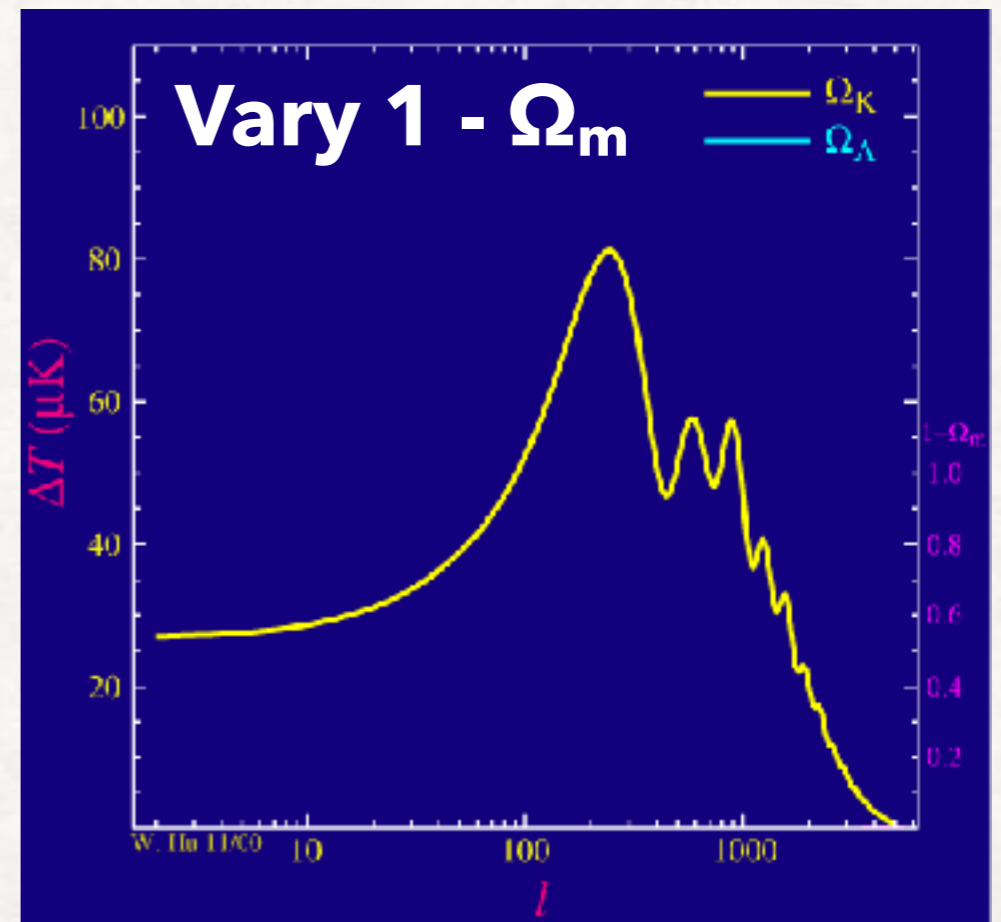
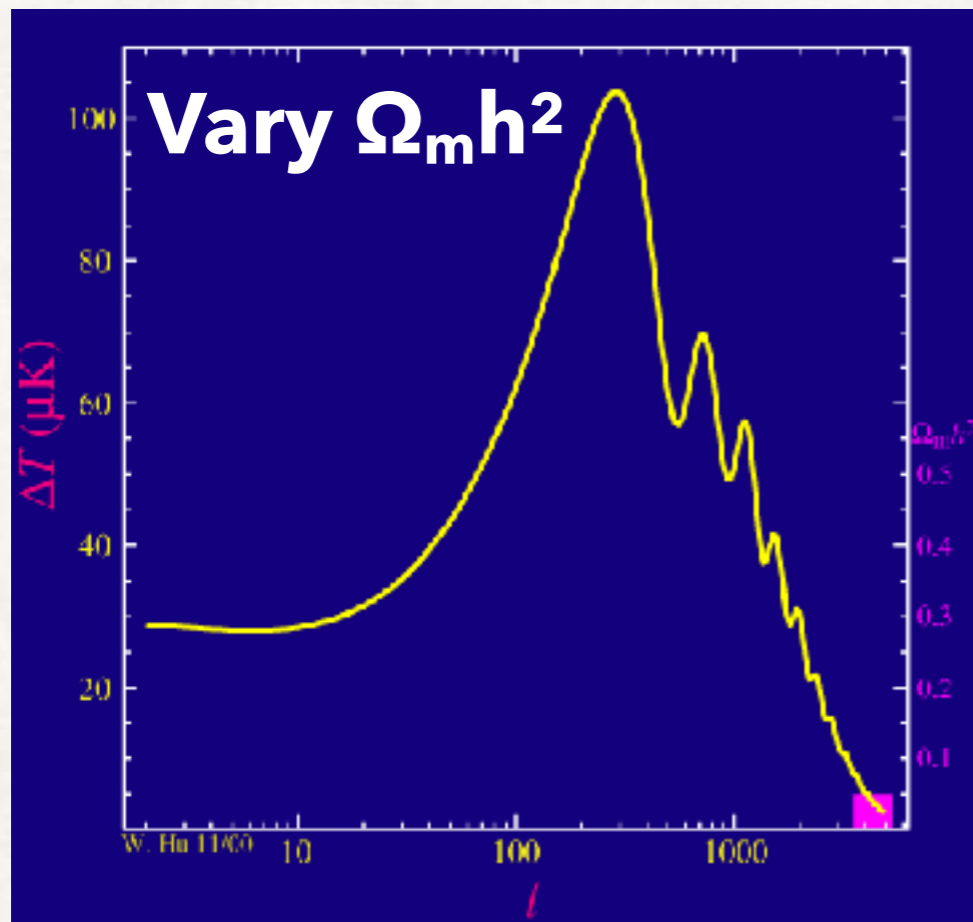
**observation**

**distance, time,...**



# LSS JUST LIKE CMB...?

- This is quite well established in case of CMB analyses.
  - Thanks to **the smallness of fluctuations** -> linear Boltzmann
  - Typically **Markov-Chain Monte Carlo** works comfortably
- Can we do the same for large scale structure?

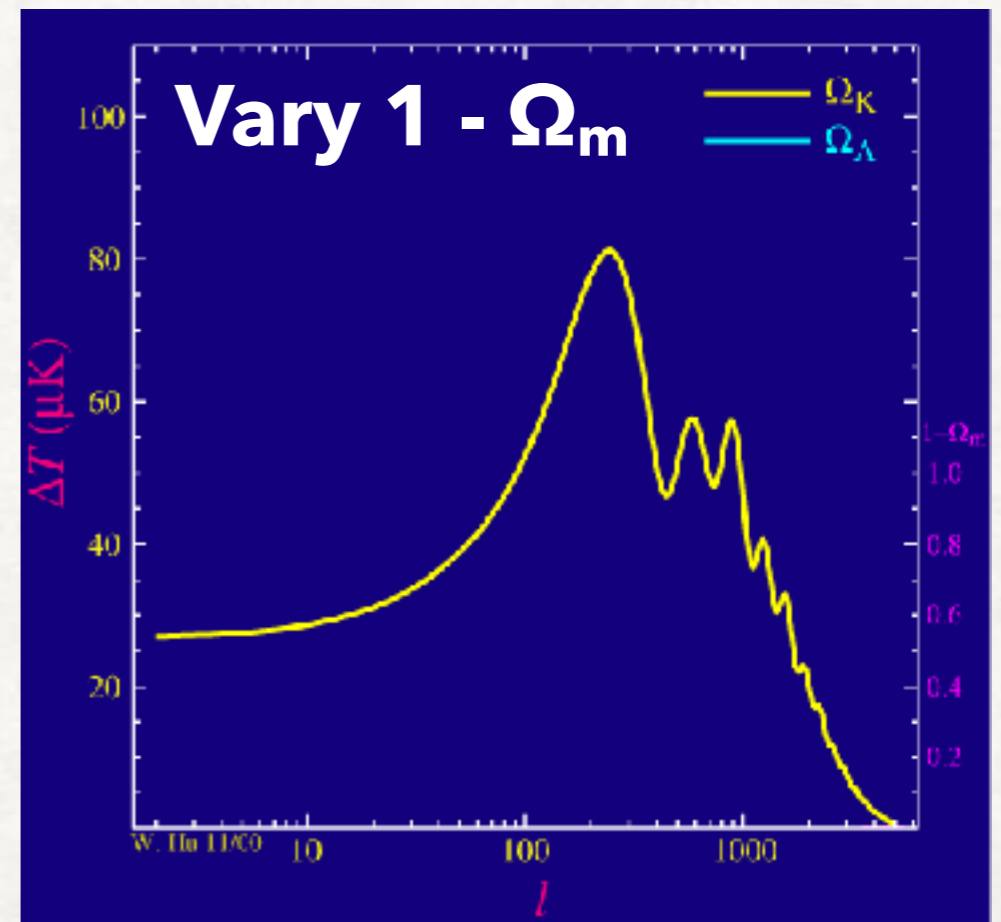
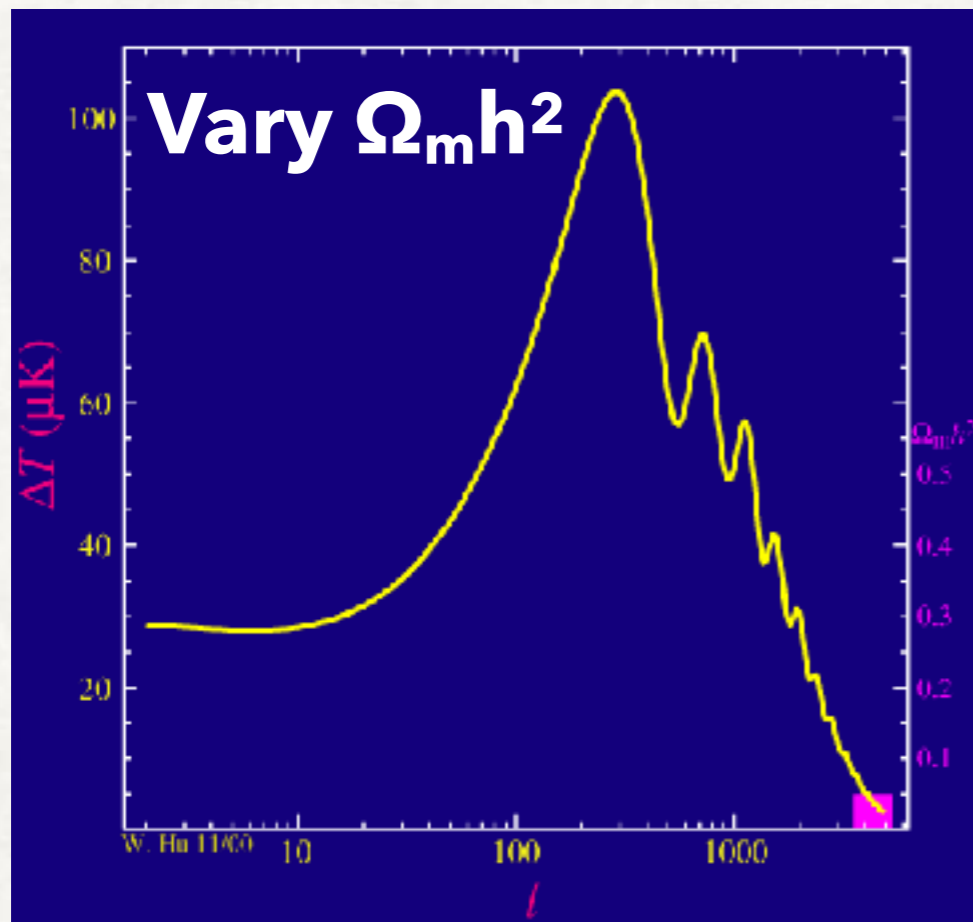


<http://background.uchicago.edu/~whu> © Wayne Hu



# LSS JUST LIKE CMB...?

- This is quite well established in case of CMB analyses.
  - Thanks to **the smallness of fluctuations** -> linear Boltzmann
  - Typically **Markov-Chain Monte Carlo** works comfortably
- Can we do the same for large scale structure?

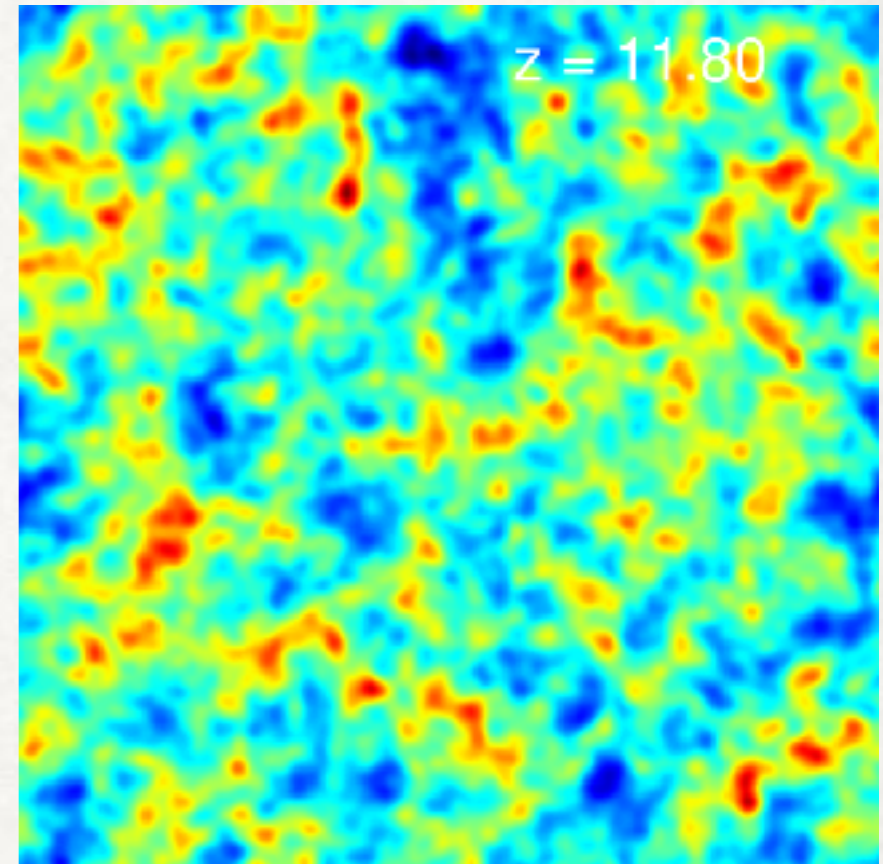


<http://background.uchicago.edu/~whu> © Wayne Hu

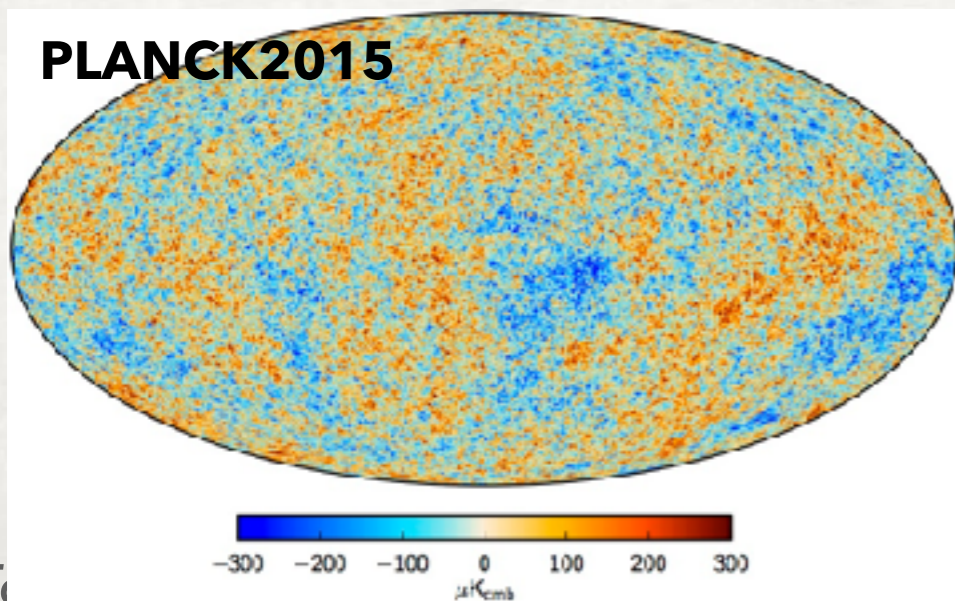


# LARGE SCALE STRUCTURE

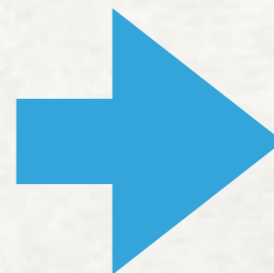
- Complementarities to CMB
  - **Dark energy** dominates the nearby universe
    - Equation of state?
  - **Gravity** is the driver of structure growth
    - Test of GR?
- Understanding nonlinearity is the key



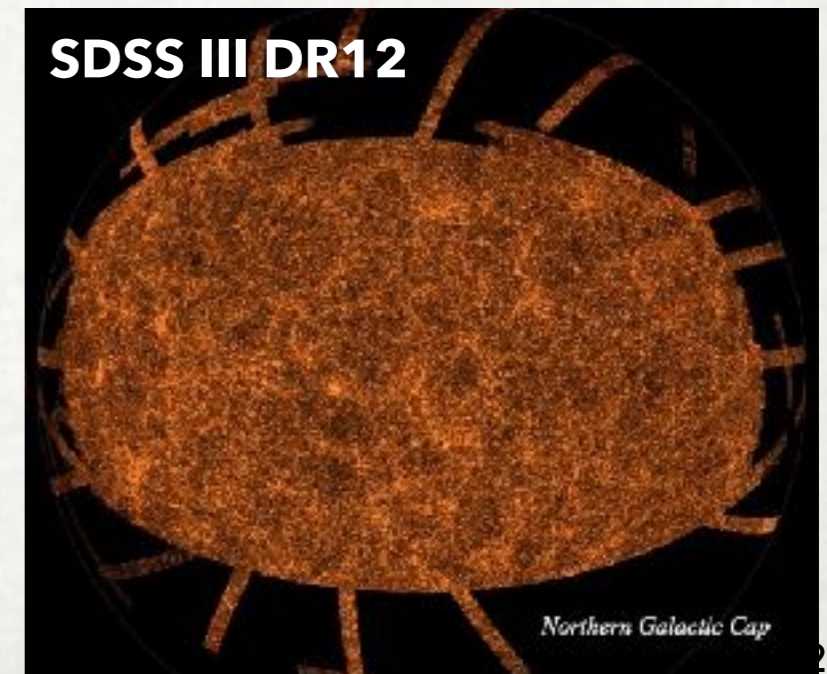
PLANCK2015



gravity !



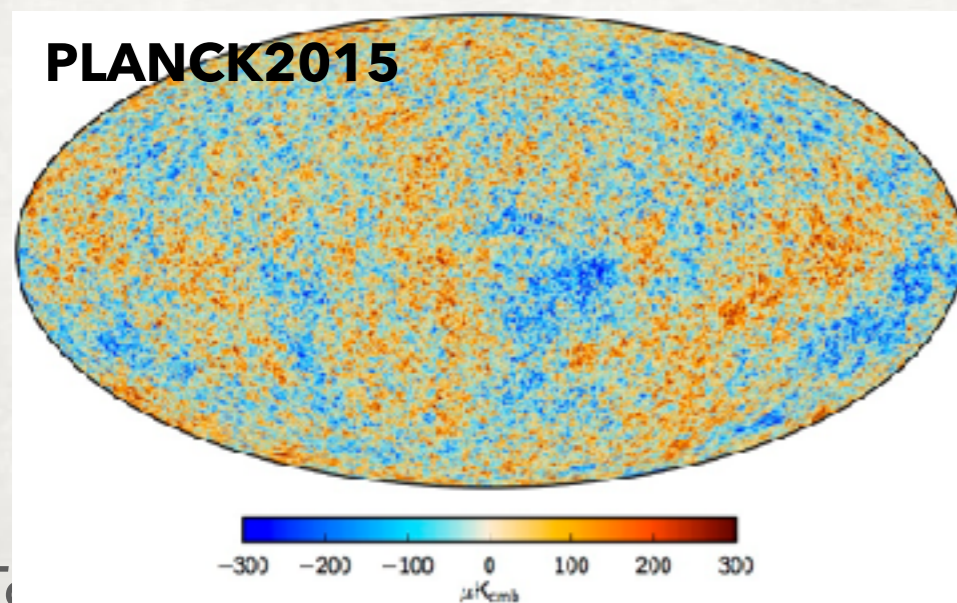
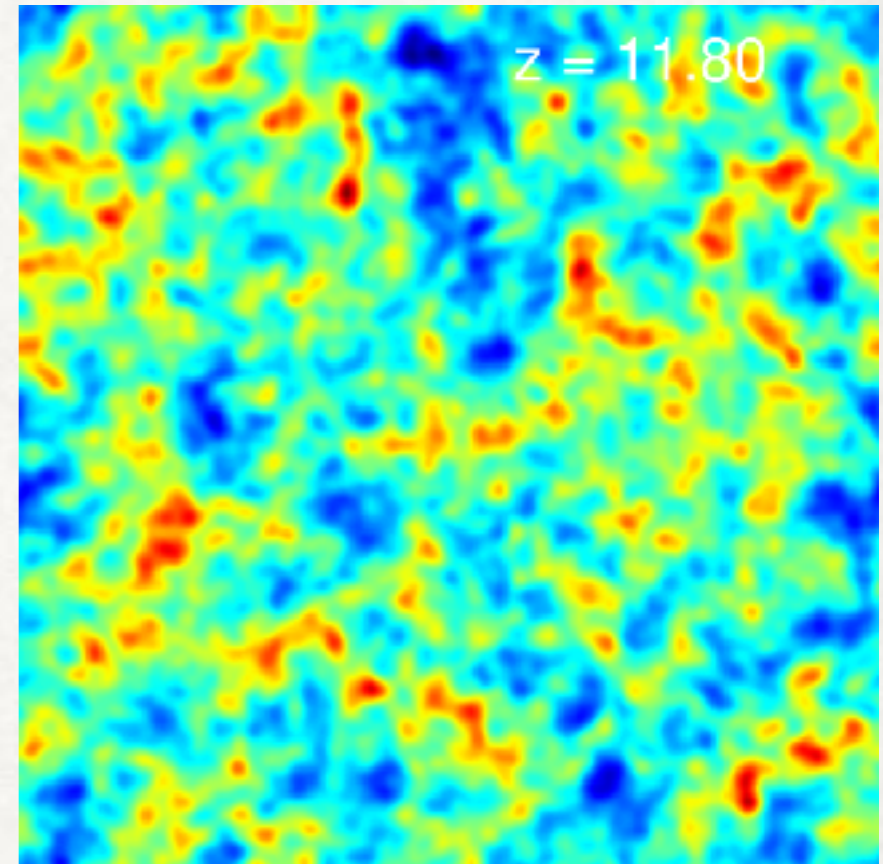
SDSS III DR12



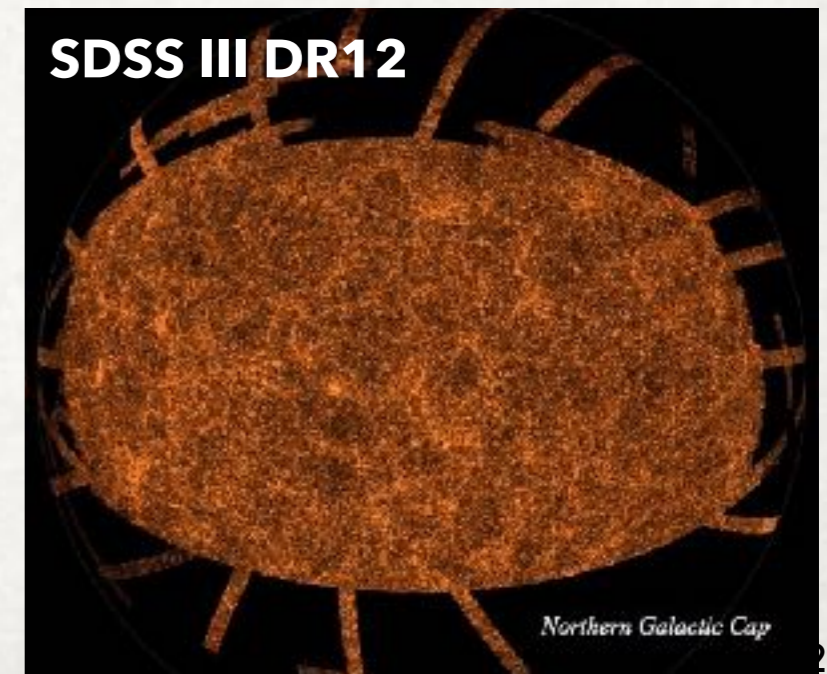
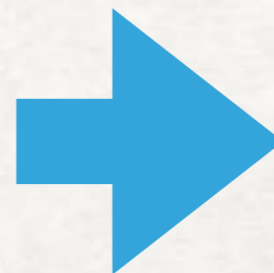


# LARGE SCALE STRUCTURE

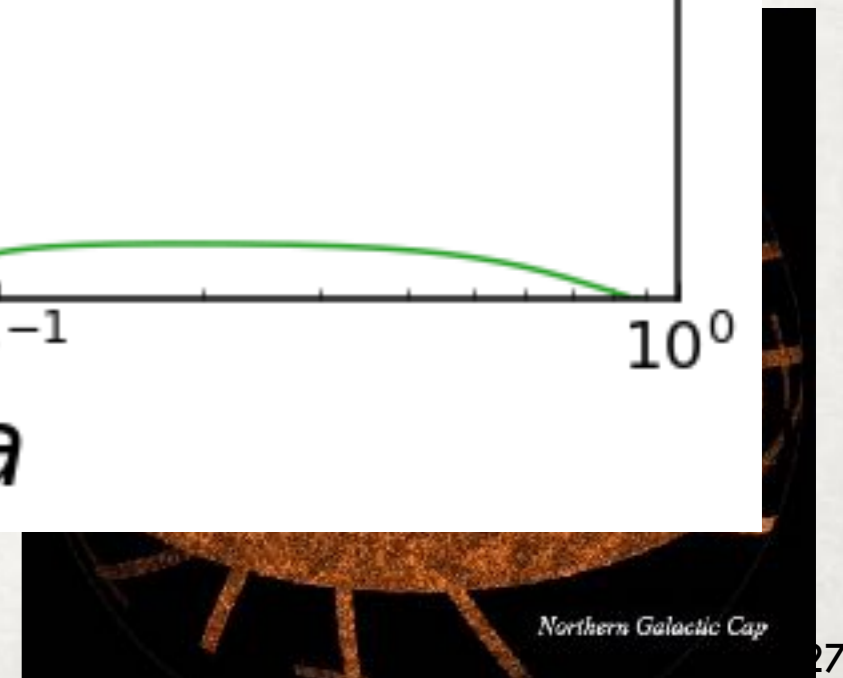
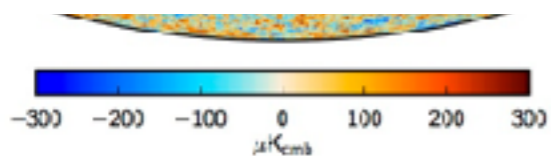
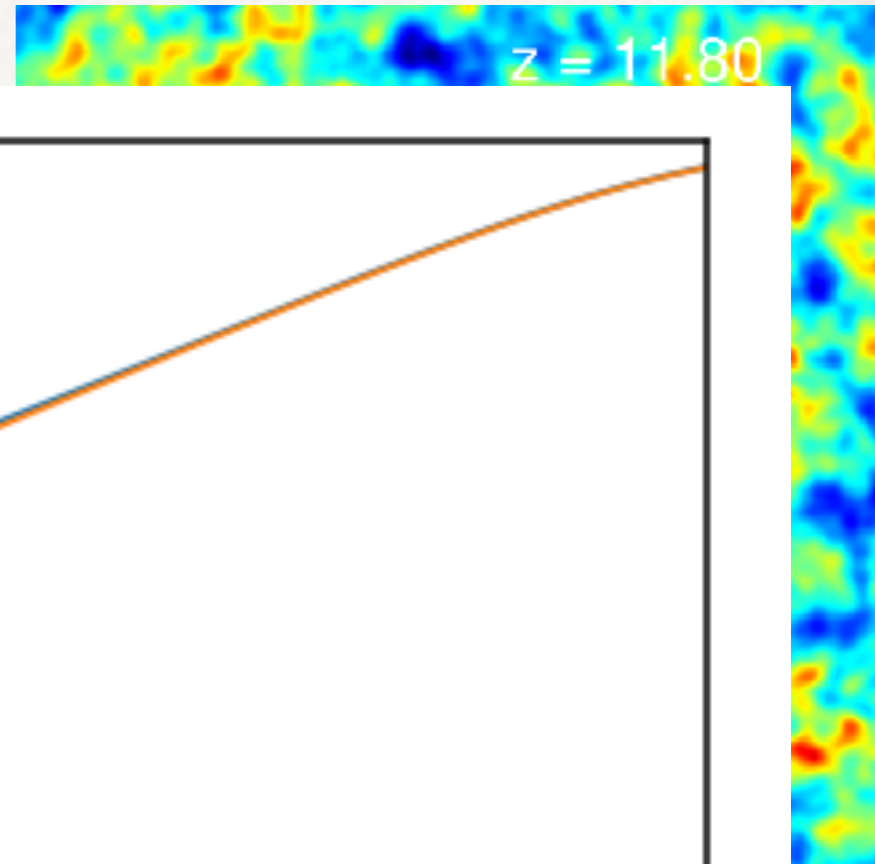
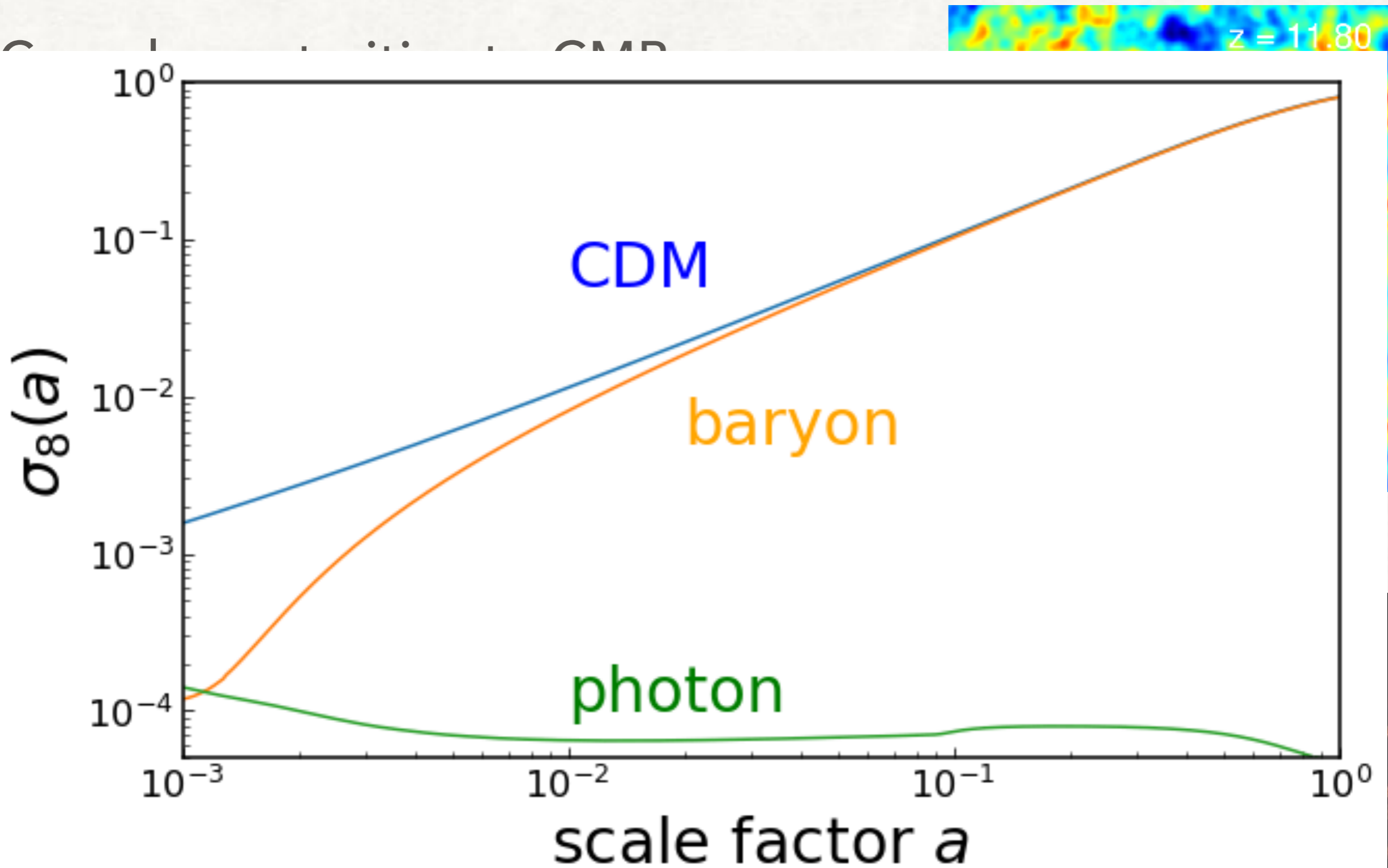
- Complementarities to CMB
  - **Dark energy** dominates the nearby universe
    - Equation of state?
  - **Gravity** is the driver of structure growth
    - Test of GR?
- Understanding nonlinearity is the key



gravity !



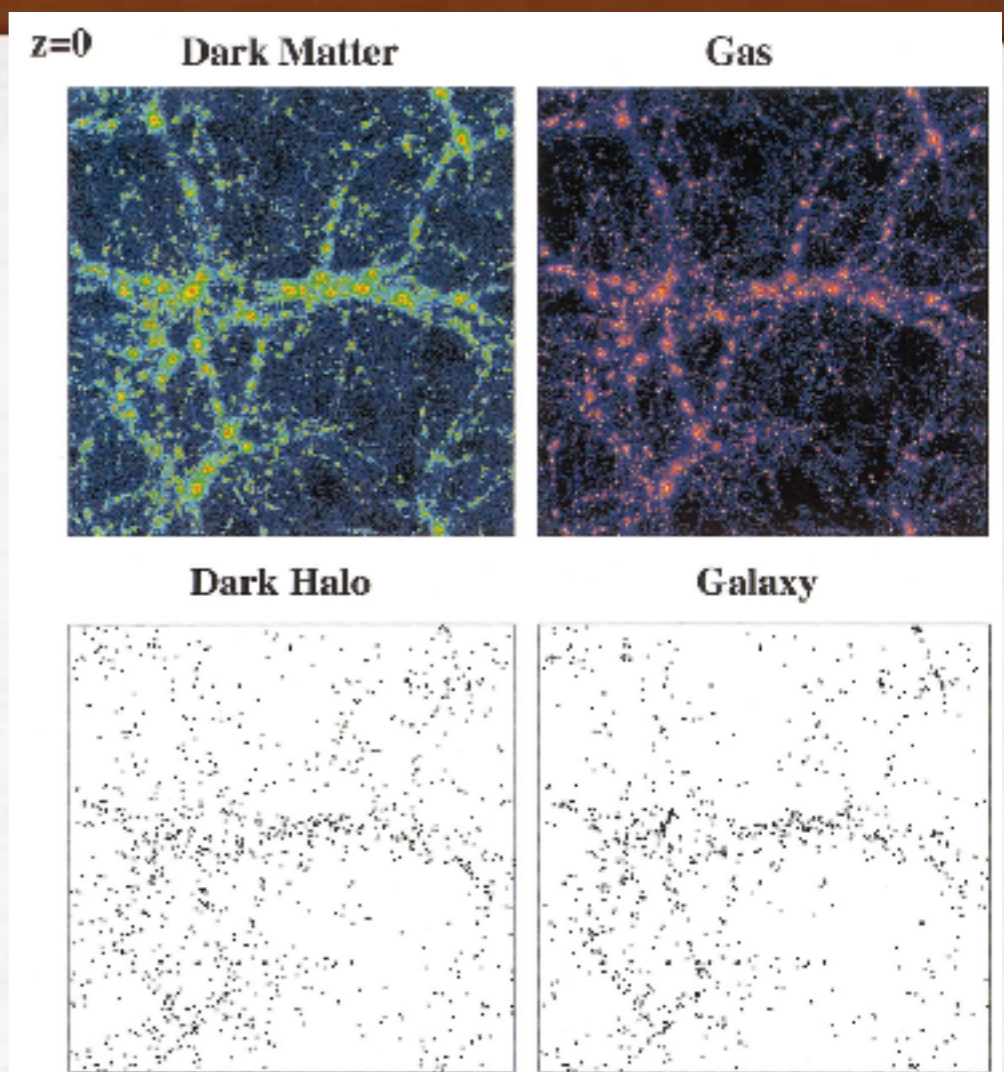
# LARGE SCALE STRUCTURE





# BIAS IS ANNOYING

- Luminous things are “**tracers**” of the underlying matter field (Kaiser '84)
- No first-principle analytical approach available (but hydro sims)
- Have to introduce many (really many!!) nuisance parameters?



Yoshikawa+'01

## The Galaxy Power Spectrum and Bispectrum in Redshift Space

Vincent Desjacques, Donghui Jeong, Fabian Schmidt

(Submitted on 11 Jun 2018)

We present the complete expression for the next-to-leading (1-loop) order galaxy power spectrum and the leading-order galaxy bispectrum in redshift space in the general bias expansion, or equivalently the effective field theory of biased tracers. We consistently include selection effects. These are degenerate with many, but not all, of the redshift-space distortion contributions, and have not been included before. Moreover, we show that, in the framework of effective field theory, a consistent bias expansion in redshift space includes selection contributions. Physical arguments about the tracer sample considered and its observational selection have to be used to identify the relevant contributions. In summary, the next-to-leading order galaxy power spectrum and leading-order galaxy bispectrum is described by 22 parameters, which reduces to 11 parameters if selection effects can be neglected. All contributions to the power spectrum and bispectrum are expressed in terms of 28 independent loop integrals.



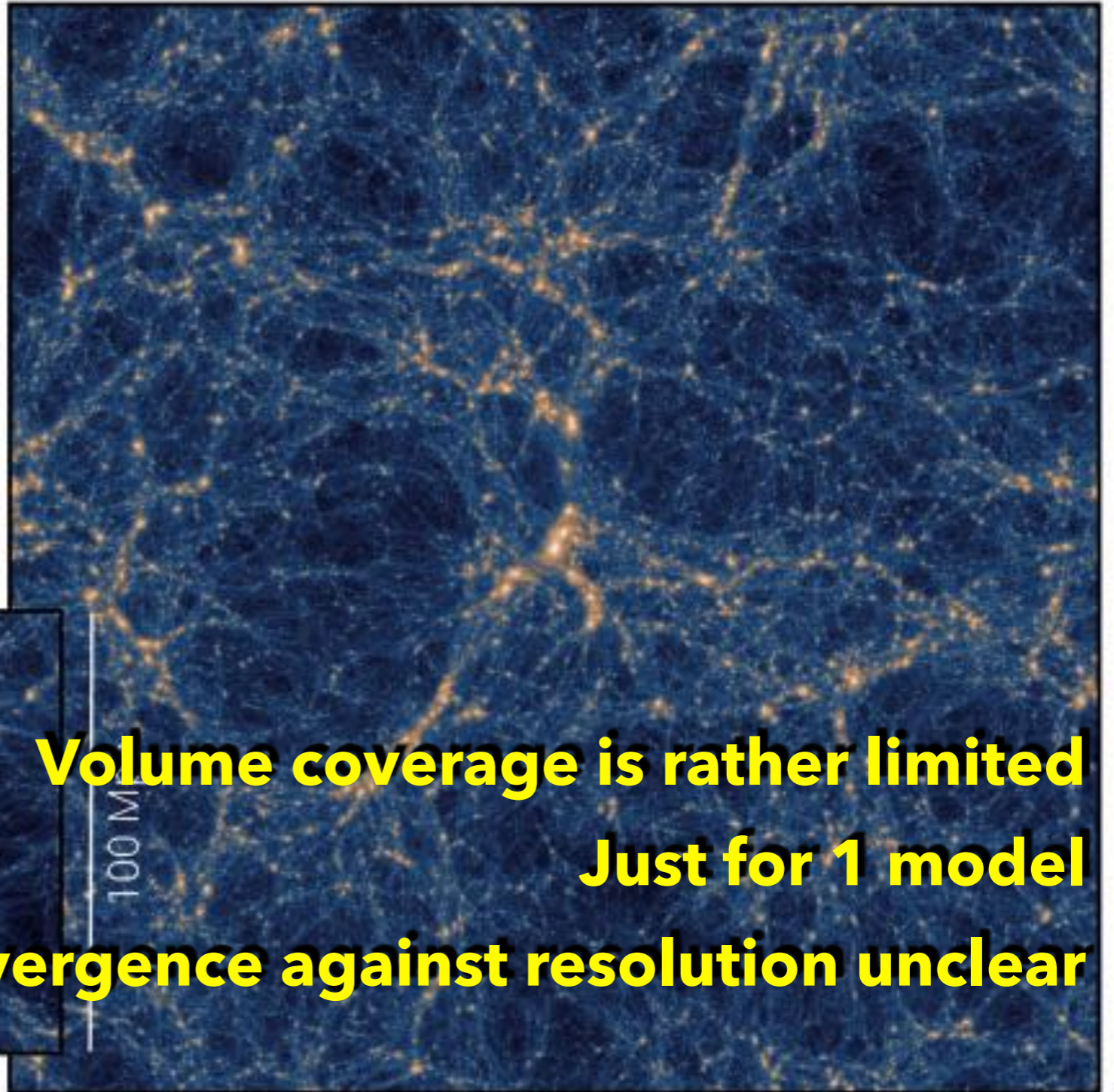
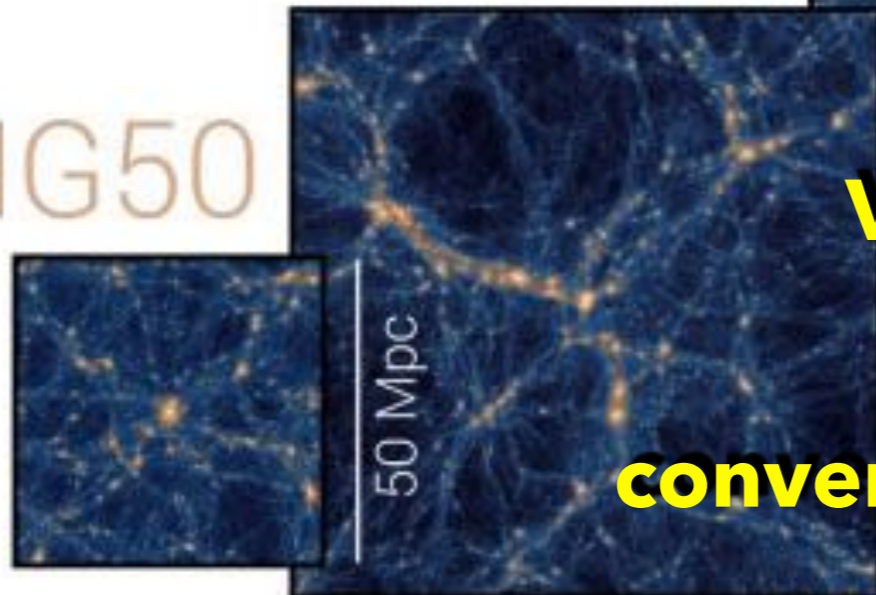
# BIAS IS DIFFICULT

Illustris TNG project  
(PI: V Springel)

TNG300

TNG100

TNG50



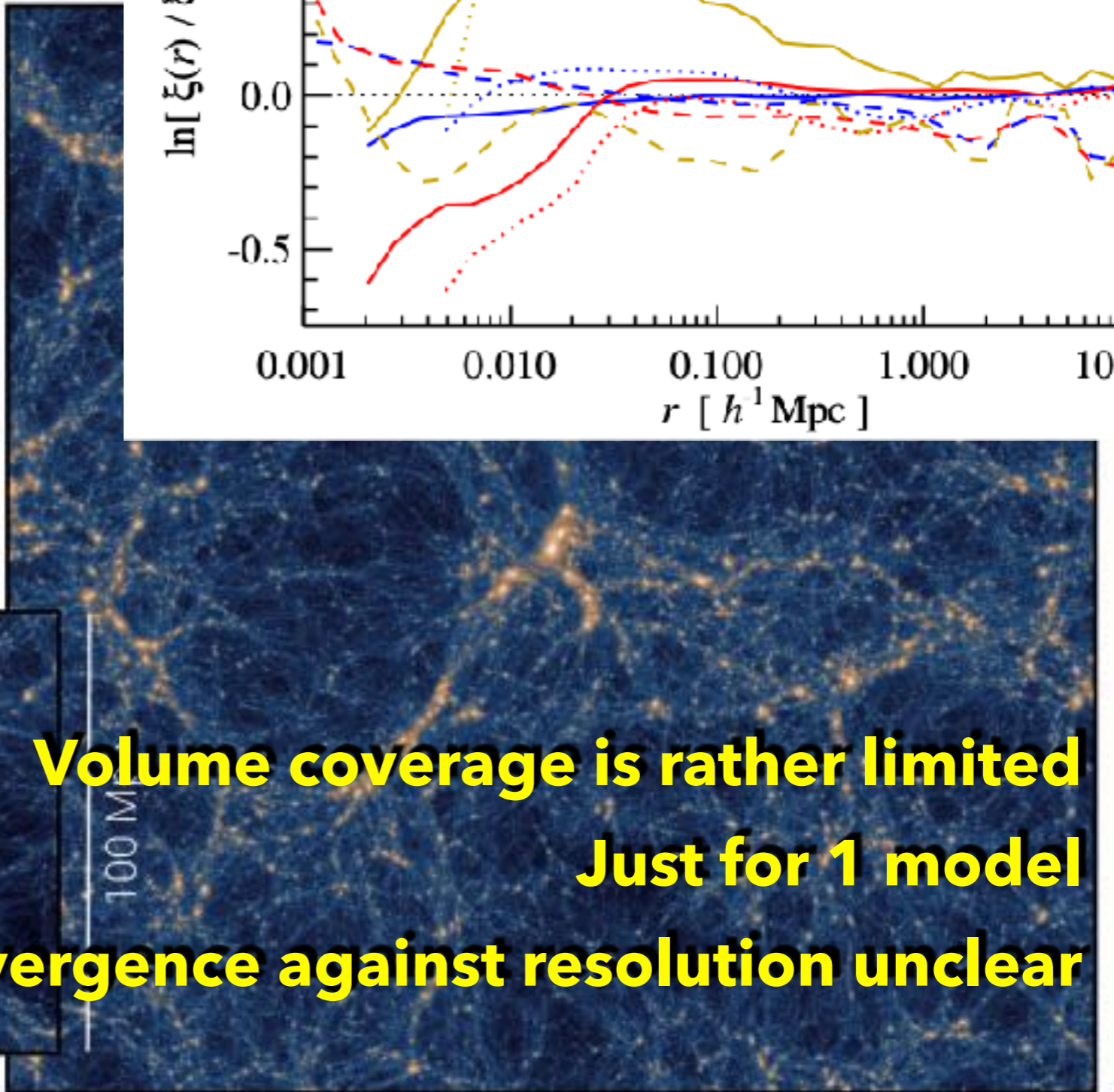
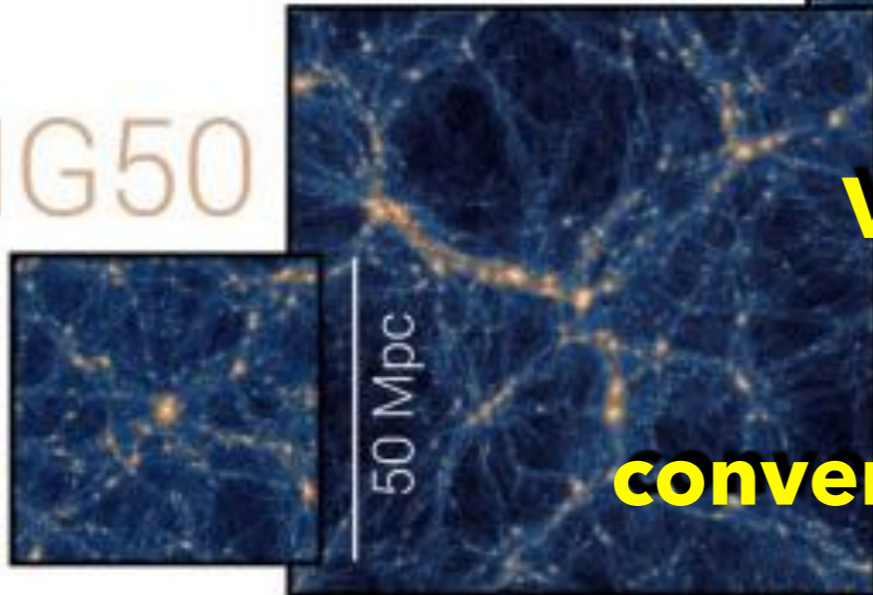
**Volume coverage is rather limited  
Just for 1 model  
convergence against resolution unclear**



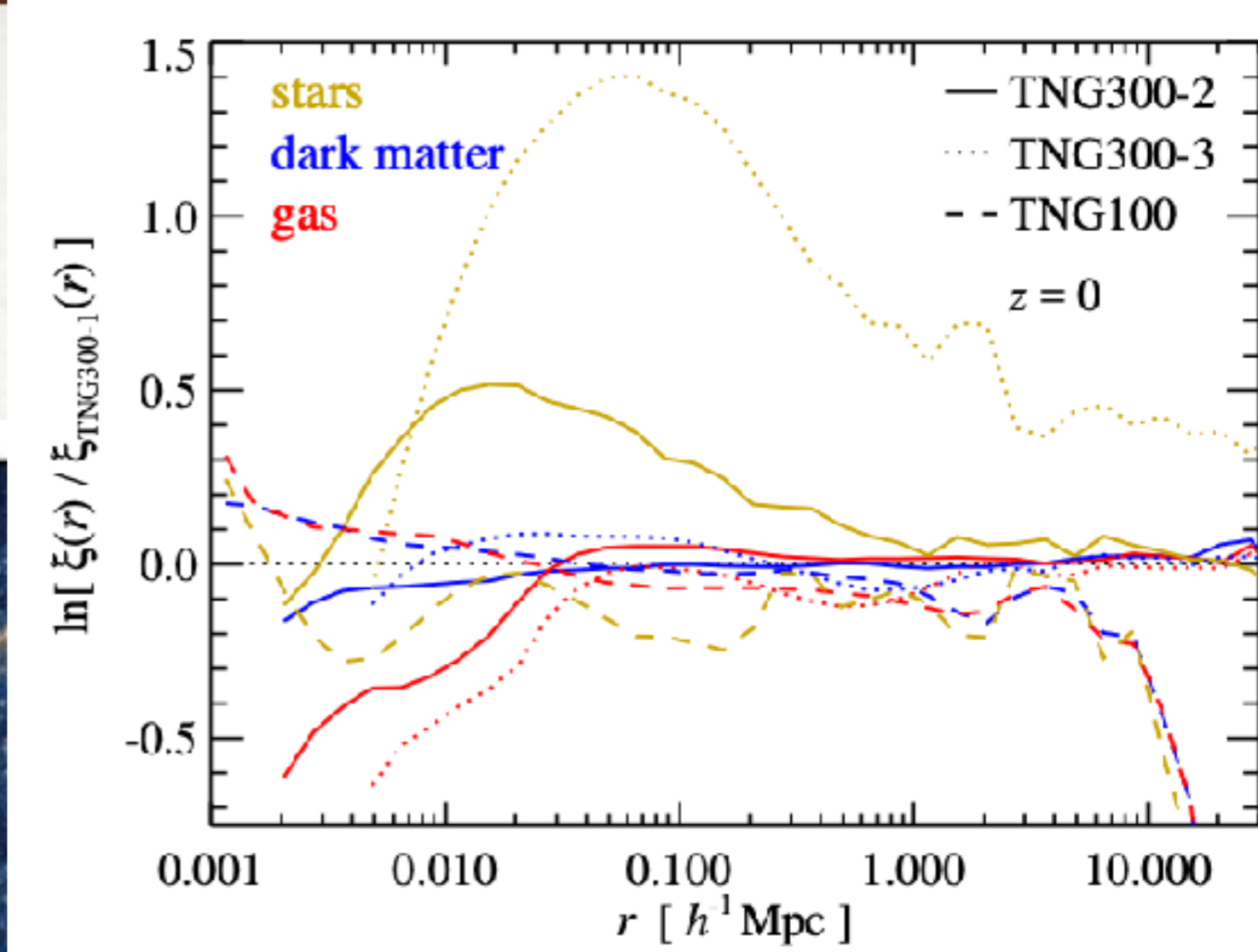
# BIAS IS DIFFICULT

Illustris TNG project  
(PI: V Springel)

TNG300  
TNG100  
TNG50



300 Mpc

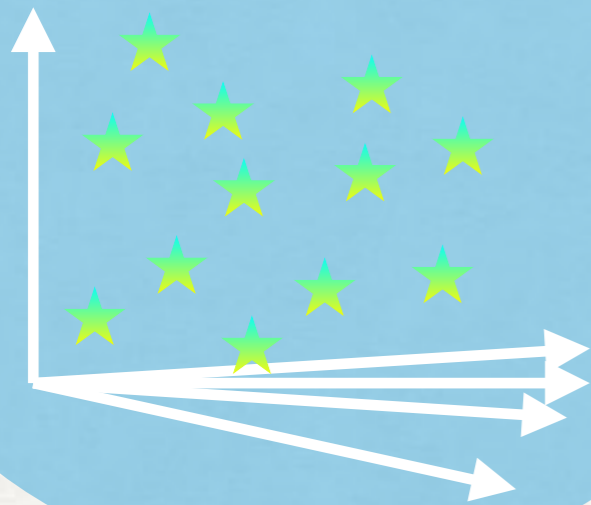


**Volume coverage is rather limited**  
**Just for 1 model**  
**convergence against resolution unclear**



# HOW TO DO IT?

## Cosmological parameters



$\Omega_m, \Omega_b, H_0, A_s, n_s, w, \dots$

**Halos**

Statistics of halos based on N-body and Machine Learning

Analytical calculation

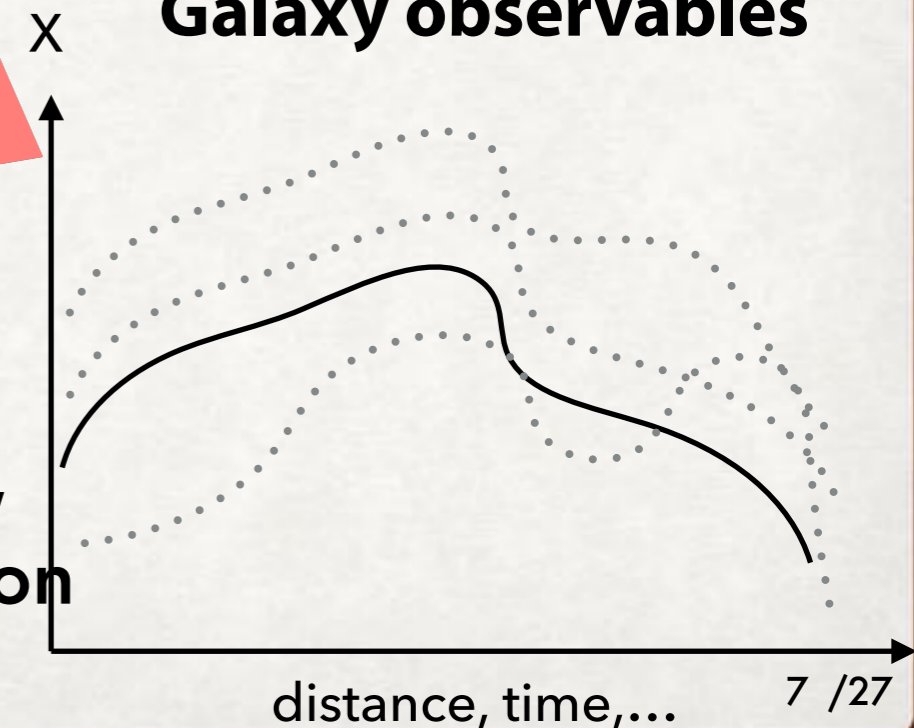
Can be updated if you want, e.g., off-centering, incompleteness, or more complex parameterization

If you are a big fan of halo approach ...

## HOD parameters

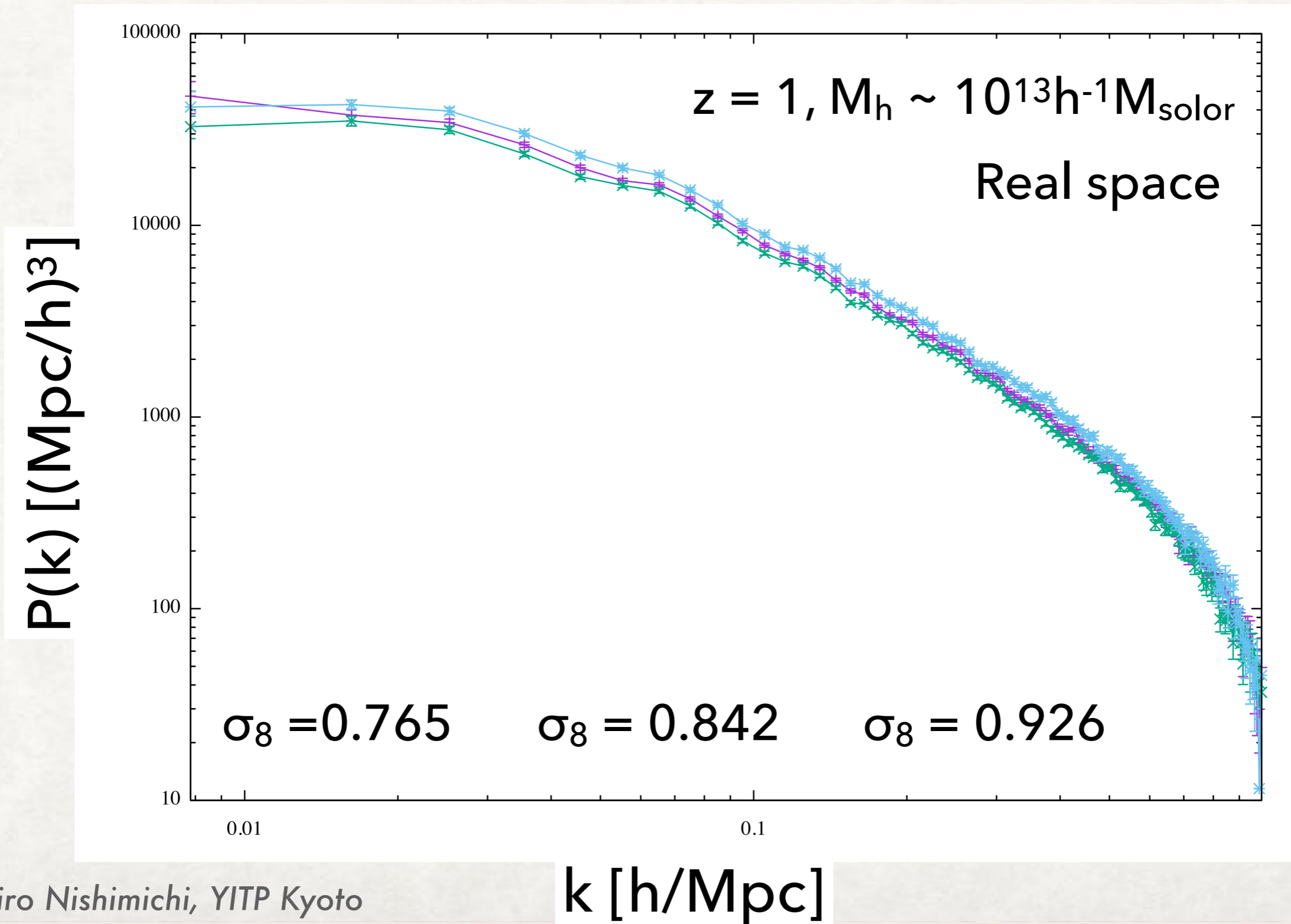
$N(M): M_{\min}, M_1, \alpha, \dots$

## Galaxy observables



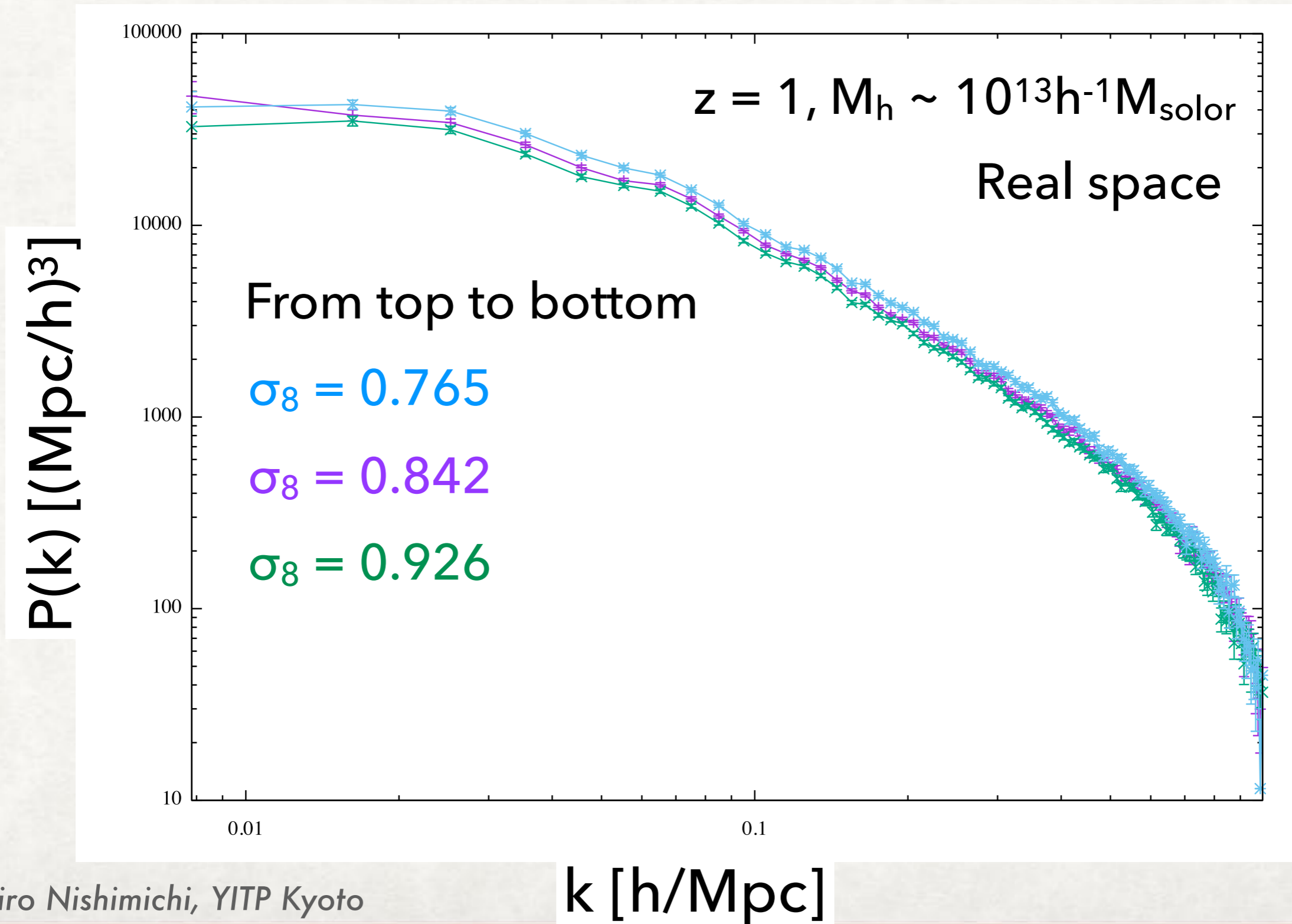
# BIAS IS STILL TRICKY

(EVEN AT THE LEVEL OF HALOS AND AT LINEAR SCALE)



# BIAS IS STILL TRICKY

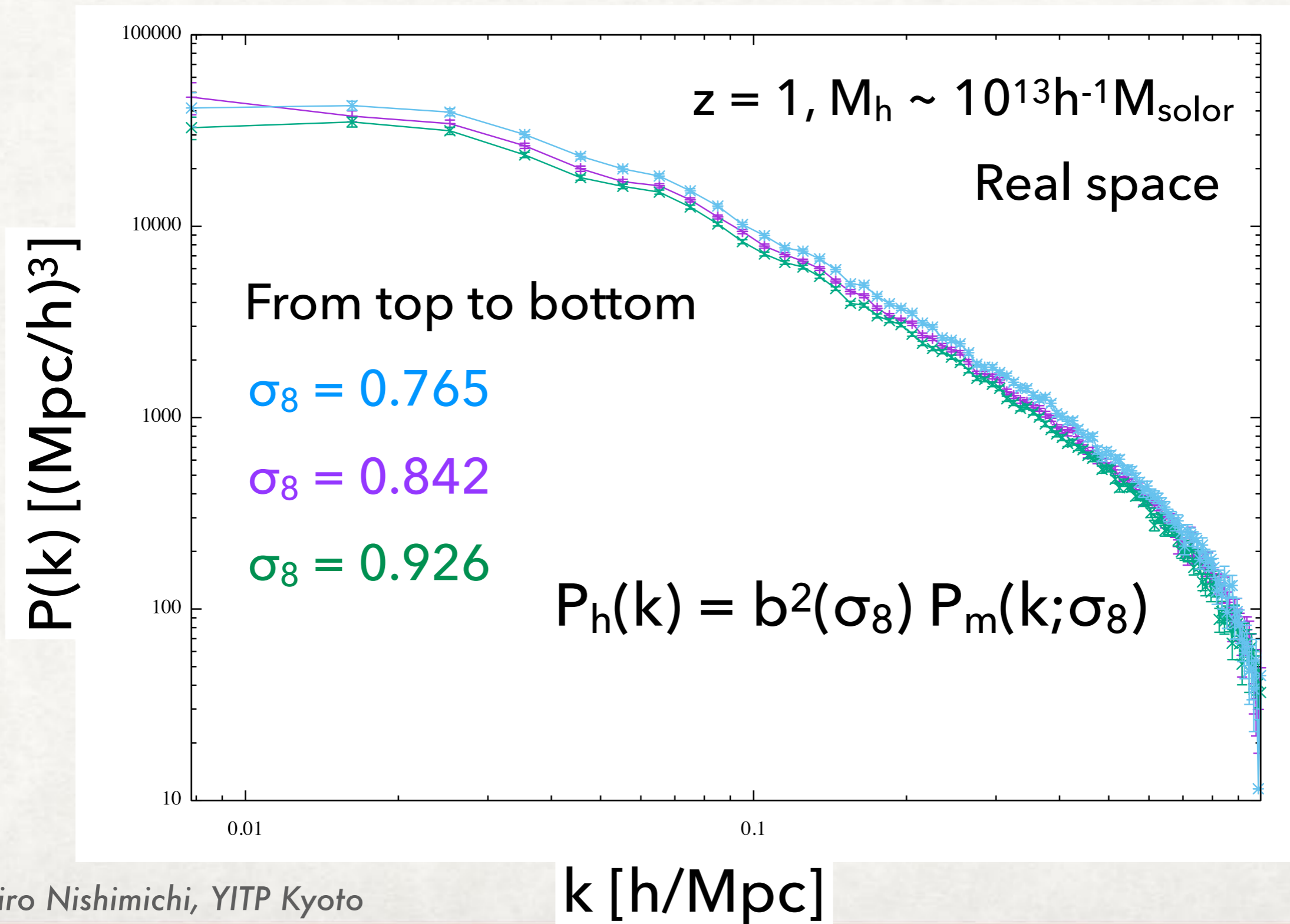
(EVEN AT THE LEVEL OF HALOS AND AT LINEAR SCALE)





# BIAS IS STILL TRICKY

(EVEN AT THE LEVEL OF HALOS AND AT LINEAR SCALE)



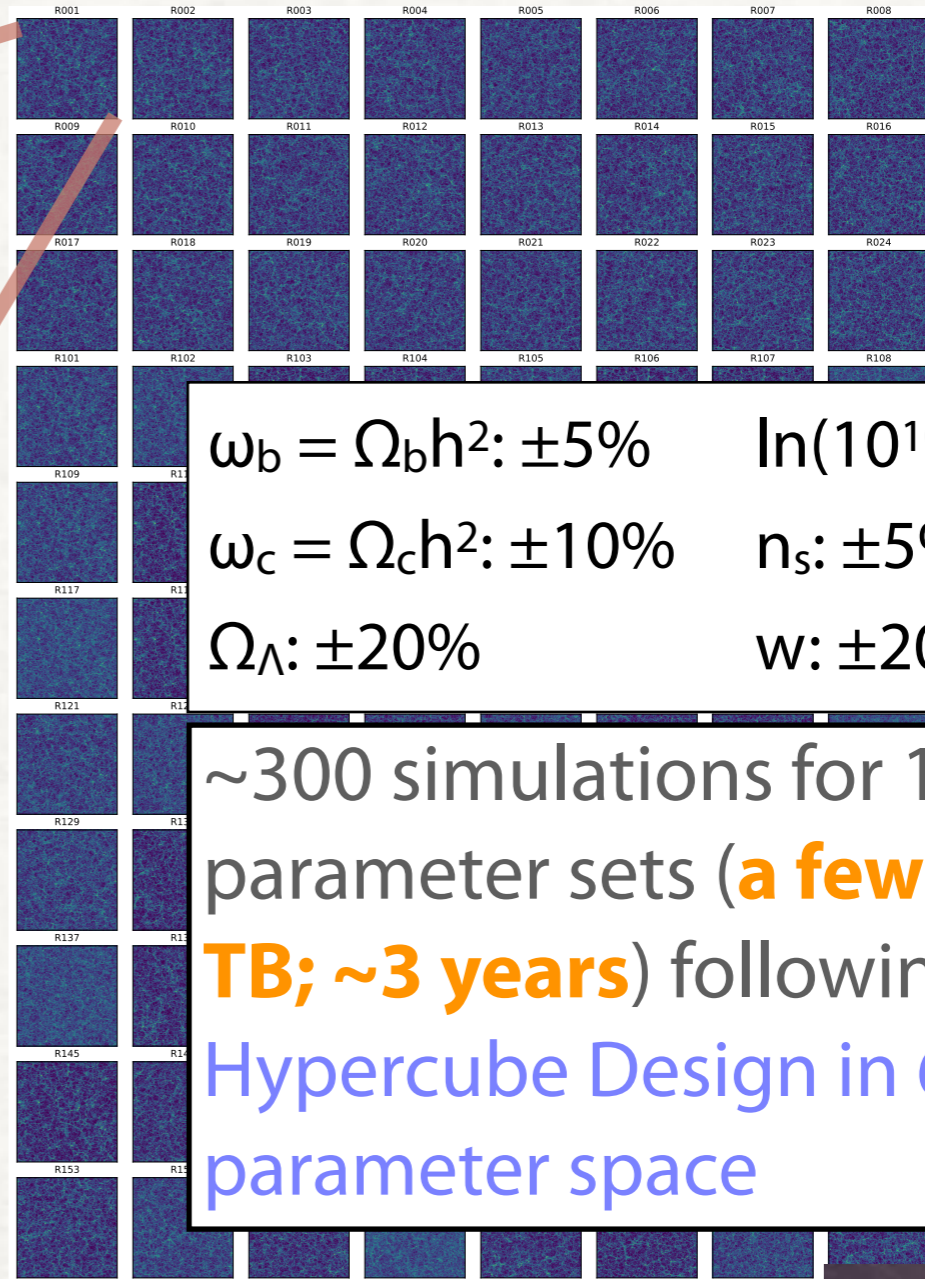
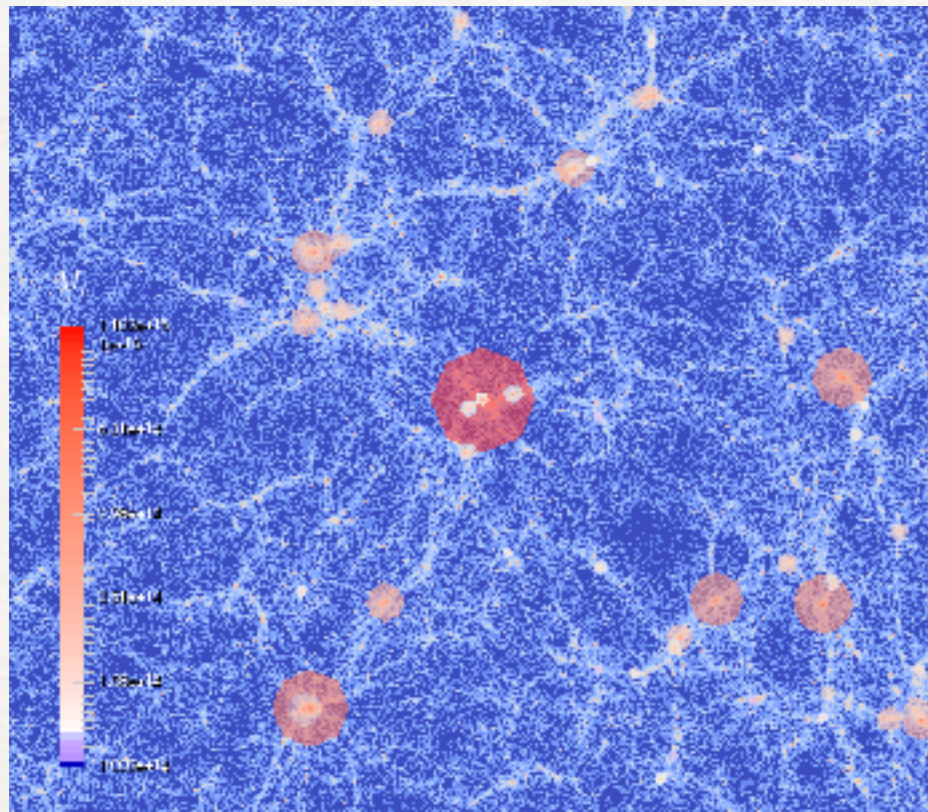


# DARK QUEST PROJECT

arXiv:1811.09504

## COSMIC STRUCTURE SIM ENSEMBLE IN 6D PARAM SPACE

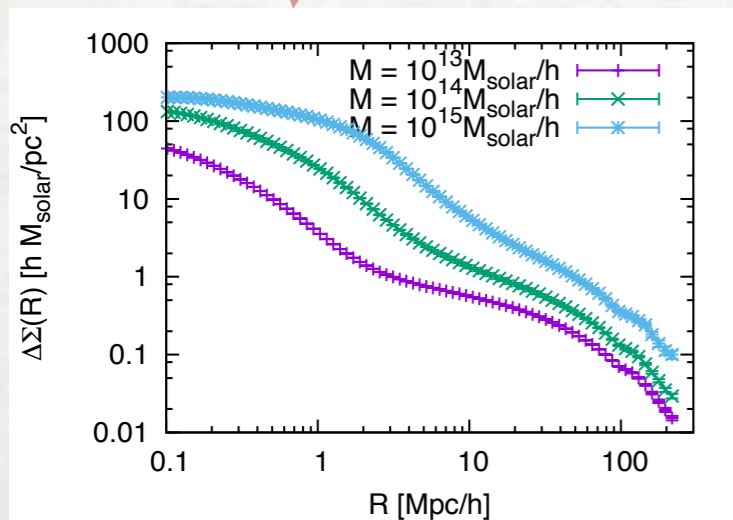
**Zoom of a 2048<sup>3</sup>-body simu**



$\omega_b = \Omega_b h^2: \pm 5\%$        $\ln(10^{10} A_s): \pm 20\%$   
 $\omega_c = \Omega_c h^2: \pm 10\%$        $n_s: \pm 5\%$   
 $\Omega_\Lambda: \pm 20\%$                        $w: \pm 20\%$

~300 simulations for 100 parameter sets (**a few hundred TB; ~3 years**) following a Latin Hypercube Design in 6D parameter space

**Extract halo statistics**



**Machine learning**

**Prediction**

**~ 100 mili sec / one model**  
**with accuracy (~3%)**



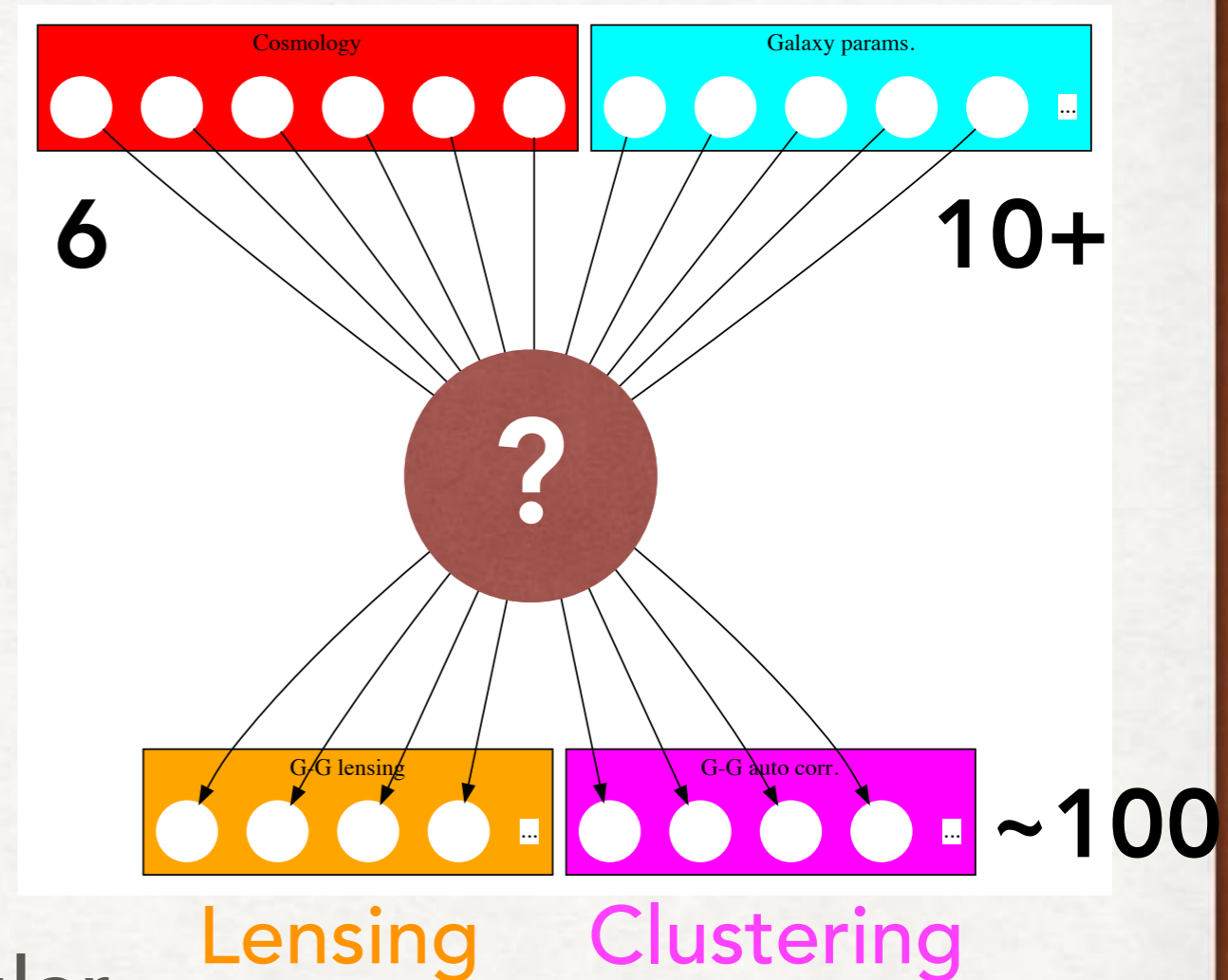
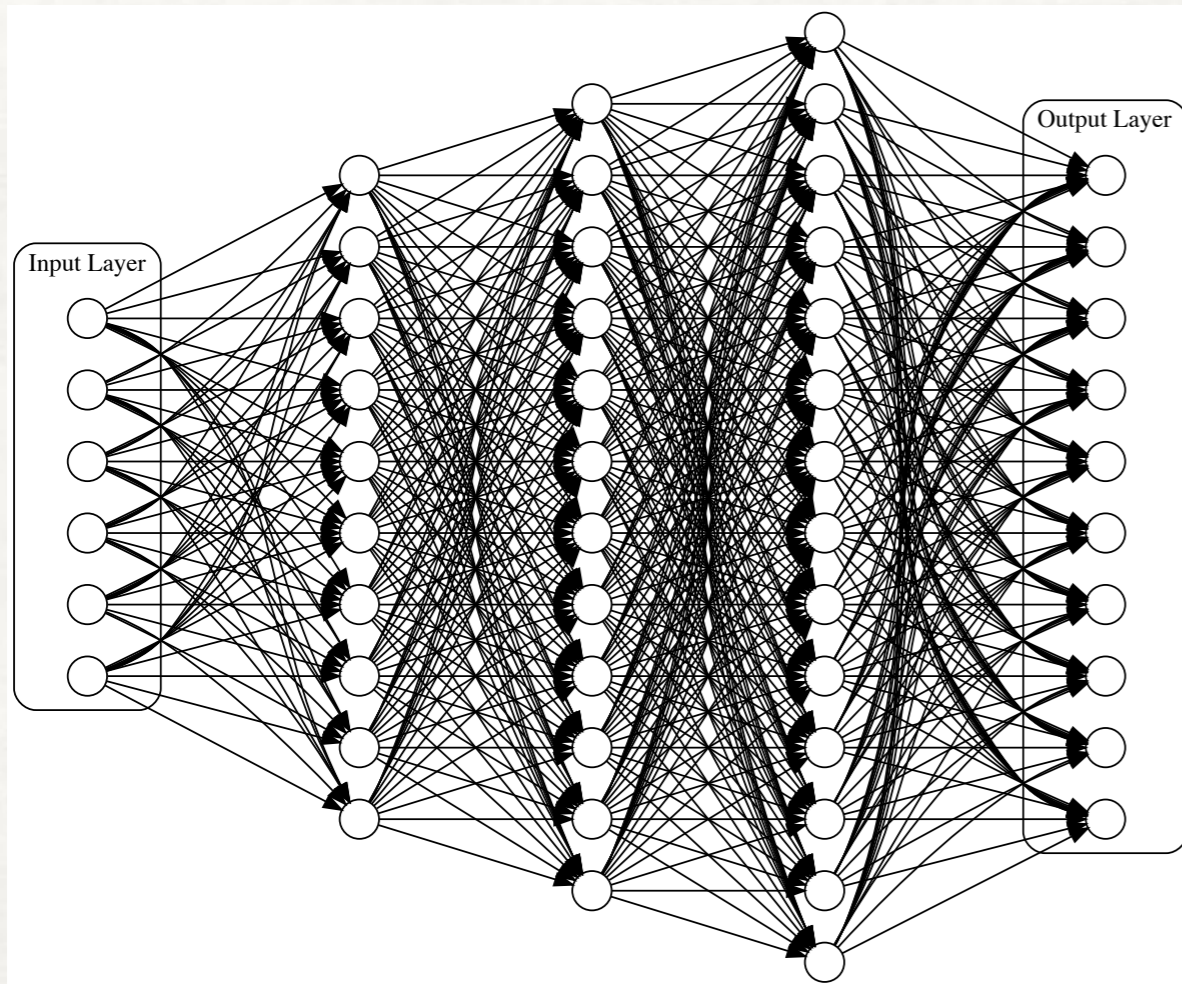


# HOW DO WE SET UP THE ML ARCHITECTURE?

Input      Hidden layers      Output

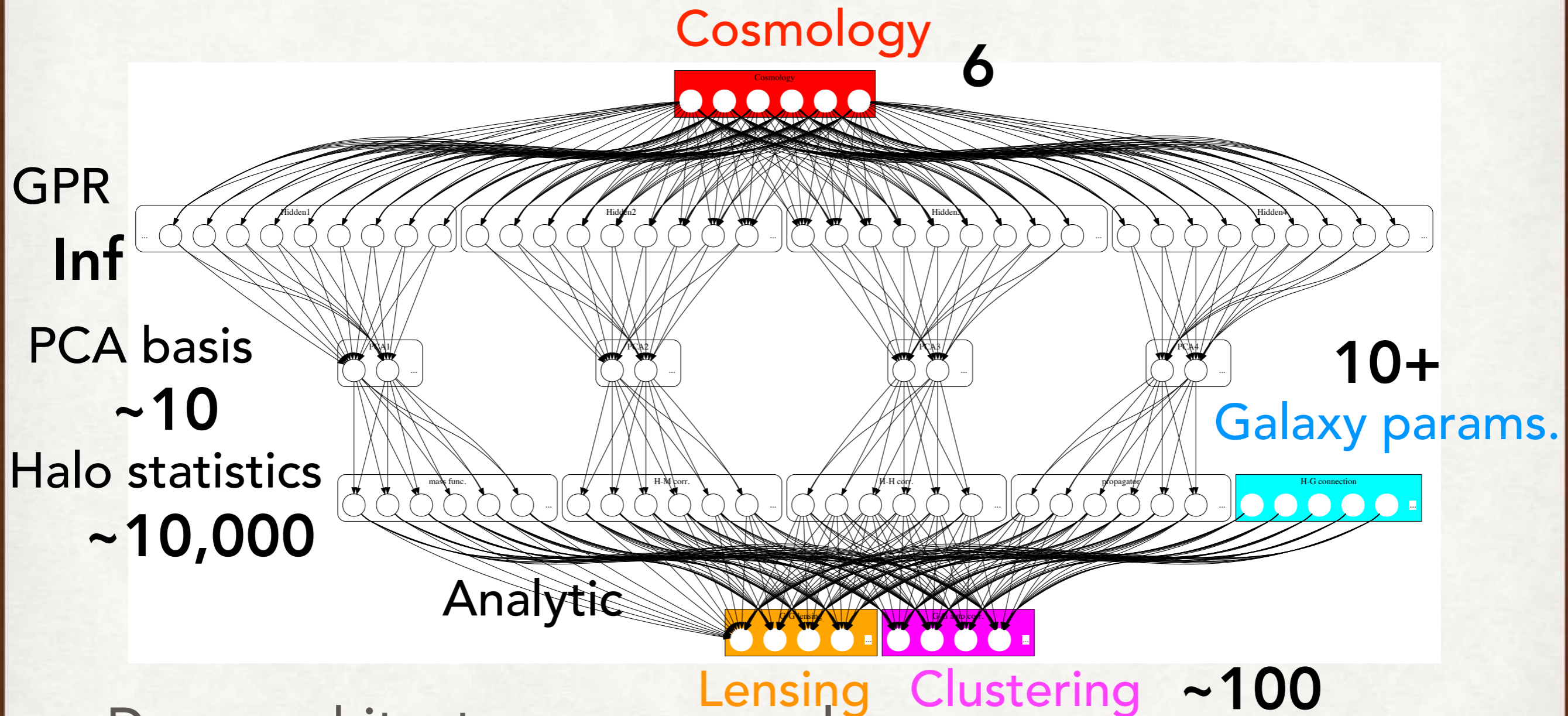
Cosmology

Galaxy params.



- Deep architectures are popular
  - Availability of large ensemble of data
  - Flexibility of the function
  - Tricks to avoid overfitting, i.e., regularization

# HOW DO WE SET UP THE ML ARCHITECTURE?

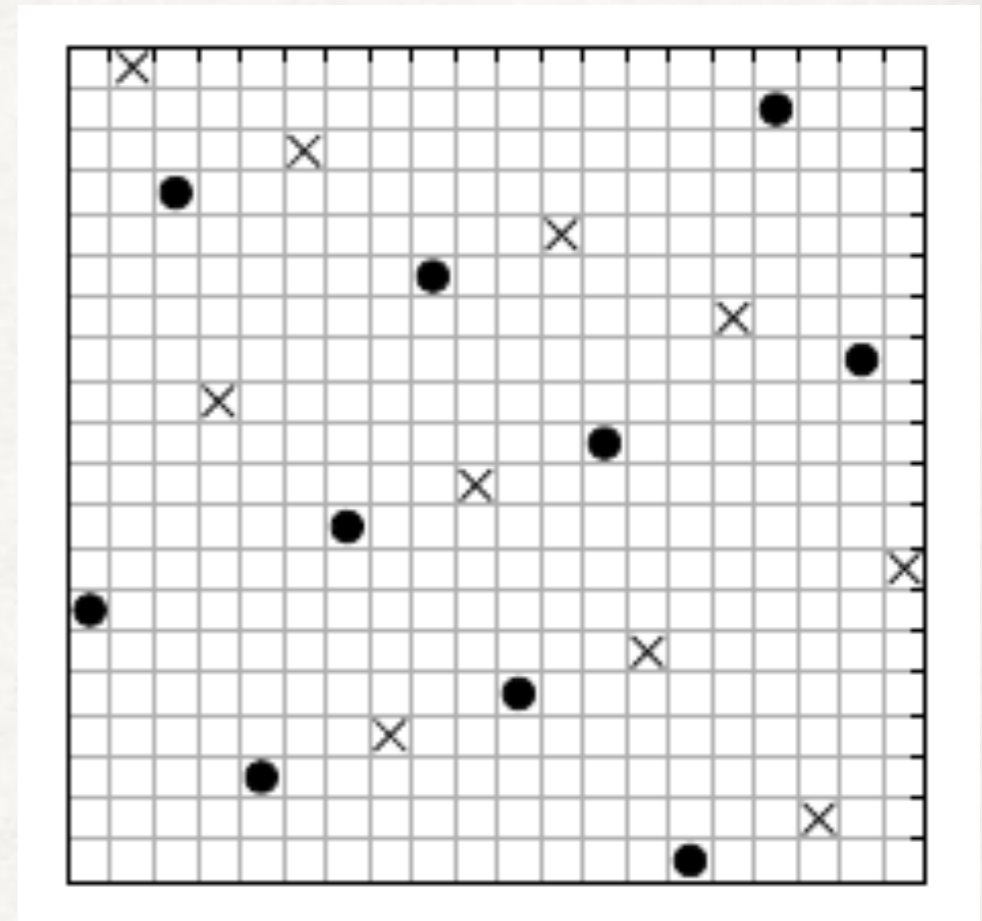


- Deep architectures are popular
  - Availability of large ensemble of data -> **not enough!**
  - Flexibility of the function
  - Tricks to avoid overfitting, i.e., **regularization**



# SLICED LATIN HYPERCUBE DESIGN

- **Latin Hypercube Designs**
  - one and only one in every row and column
  - Good **projection properties**
- **Hierarchical LHD** (Ba, Myers, Brenneman '15)
  - Each symbol forms an LHD
  - All the points together form an LHD
  - **Space filling property** is ensured by



minimizing

$$\Phi(\mathbf{X}_N) = \frac{1}{2} \left( \phi_{\text{all}} + \frac{1}{m} \sum_{t=1}^m \phi_t \right),$$

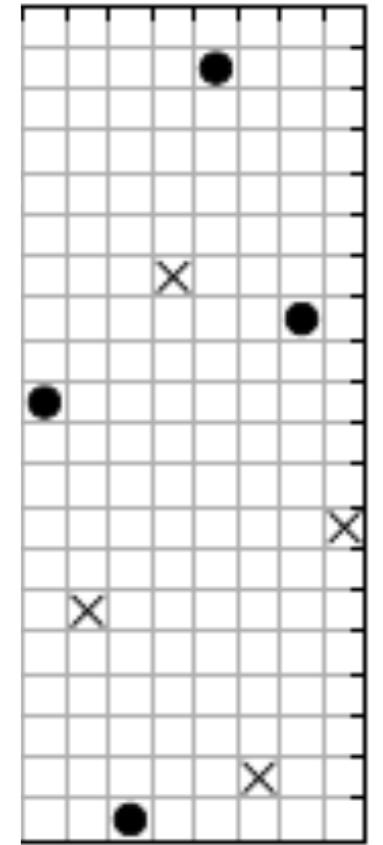
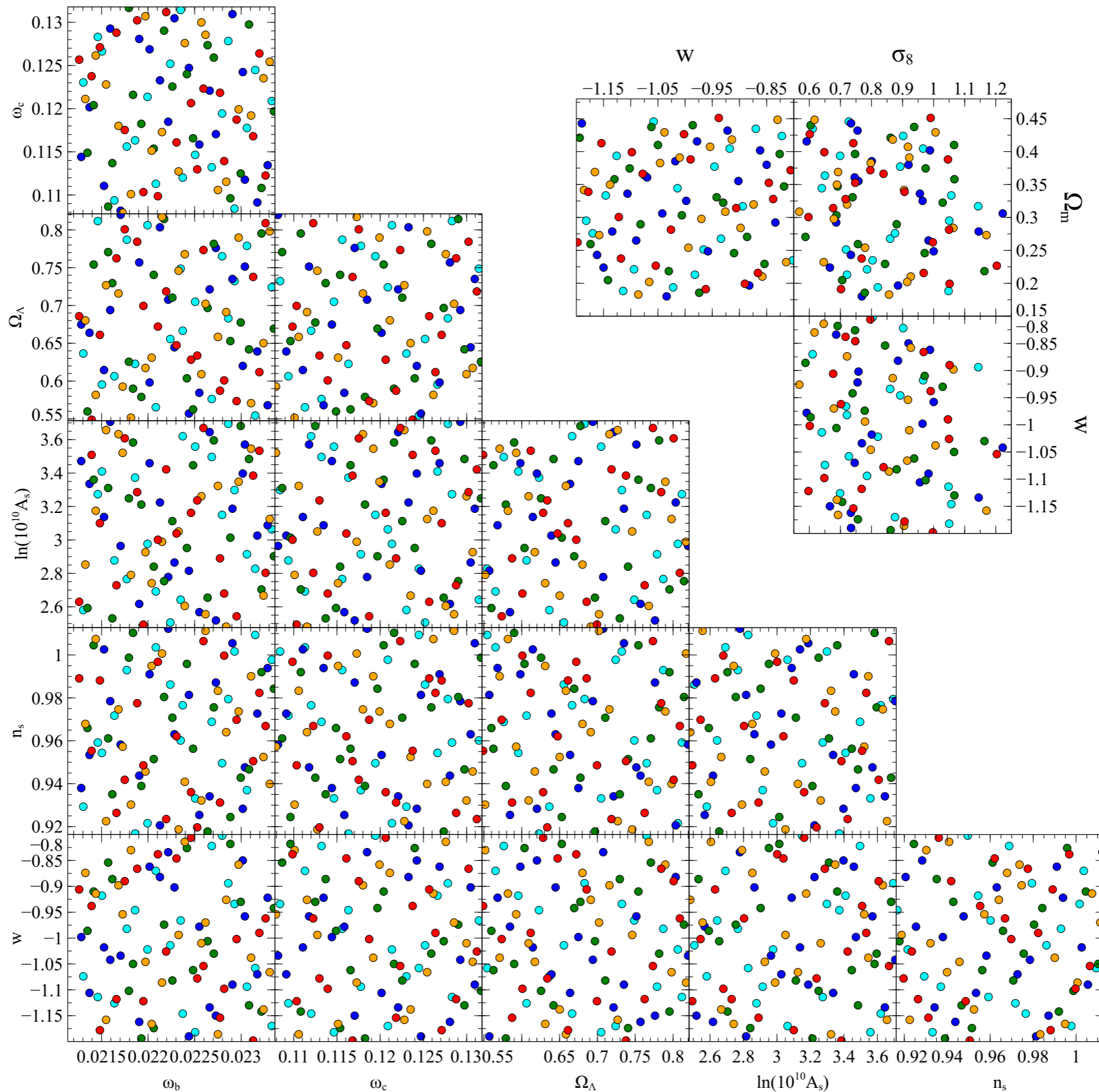
$$\phi(\mathbf{X}_N) = \left( \frac{2}{N(N-1)} \sum_{i \neq j} \frac{1}{d^r(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})} \right)^{1/r}$$

Cost function for **all the samples** within the **t-th slice** with some large  $r$

- Useful
  - for a **stringent cross-validation test**
  - by splitting the sample into **training** and **validation** sets

# SLIC

- Latin
- one
- Go
- Hiera
- Eac
- All
- Spa
- min
- Cost fun
- Use
- fo
- b



$$\left( \frac{1}{(i, \mathbf{x}(j))} \right)^{1/r}$$

me large r



# DIMENSION REDUCTION

- Mass function

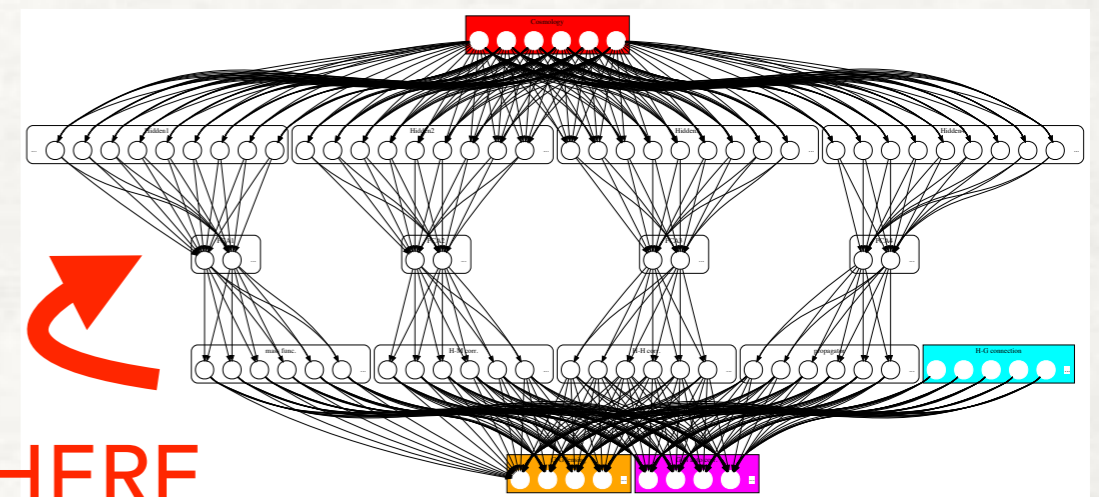
fitting

$$f(\sigma) = A[\sigma^{-a} + b] \exp\left[-\frac{c}{\sigma^2}\right]$$

- Sheth-Tormen type functional form
- b and c from Tinker
- (A, a) at 21 redshifts = **42 component vector -> 6 PCs**

- Halo-matter cross correlation

- (r-bin, n-bin, z-bin) = (66, 13, 21)
- **18,018 components -> 5 PCs**



HERE

- Halo auto correlation

- (r-bin, n<sub>1</sub>-bin, n<sub>2</sub>-bin, z-bin) = (21, 8, 8, 21)
- **28,224 components -> 8 PCs**

- Propagator

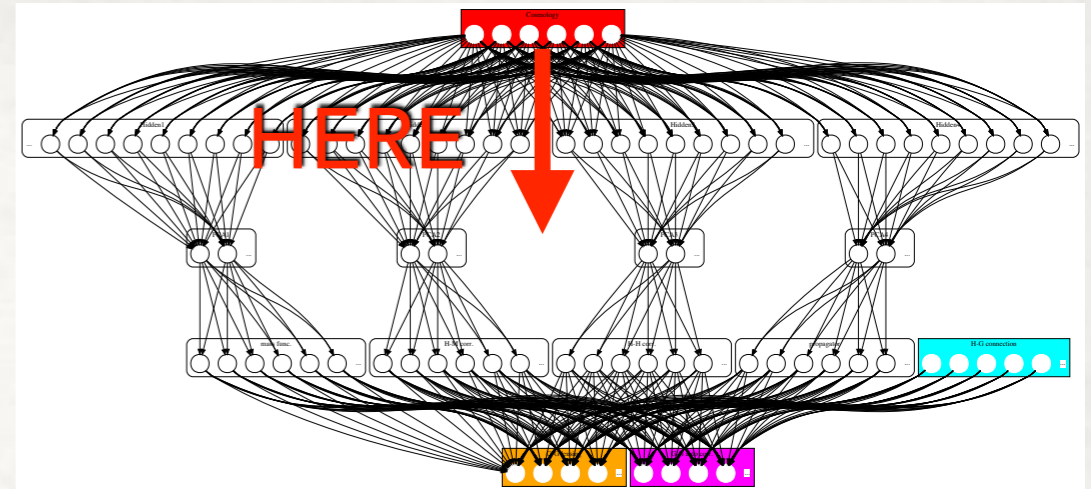
fitting

$$G(k) = [c_0 + c_2 k^2 + c_4 k^4] \exp\left[-k^2 \sigma_d^2 / 2\right]$$

- 3 parameters at 21 redshifts = **63 components -> 3 PCs**

# GAUSSIAN PROCESS REGRESSION

- Take the limit of **infinite # of nodes**
- Regularize the coefficients (i.e., L2 norm) s.t. the outputs follow i.i.d. -> **multivariate Gaussian (c.f., central limit theorem)**
  - Specified by the **correlation (kernel) function**
  - works as a non-parametric regressor or classifier

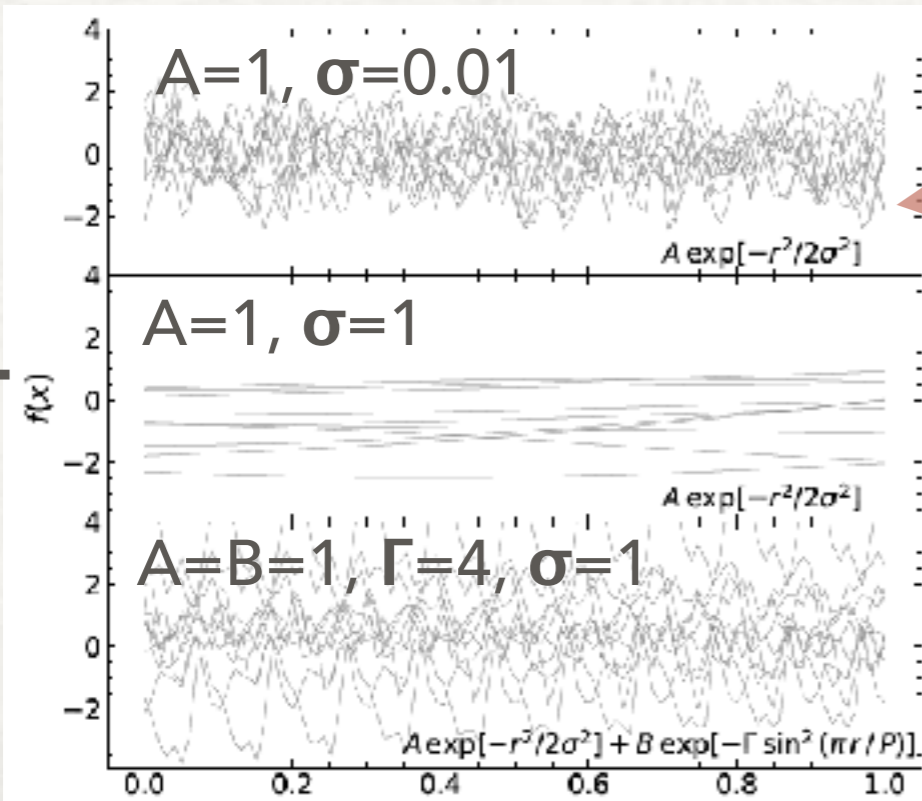


Prior

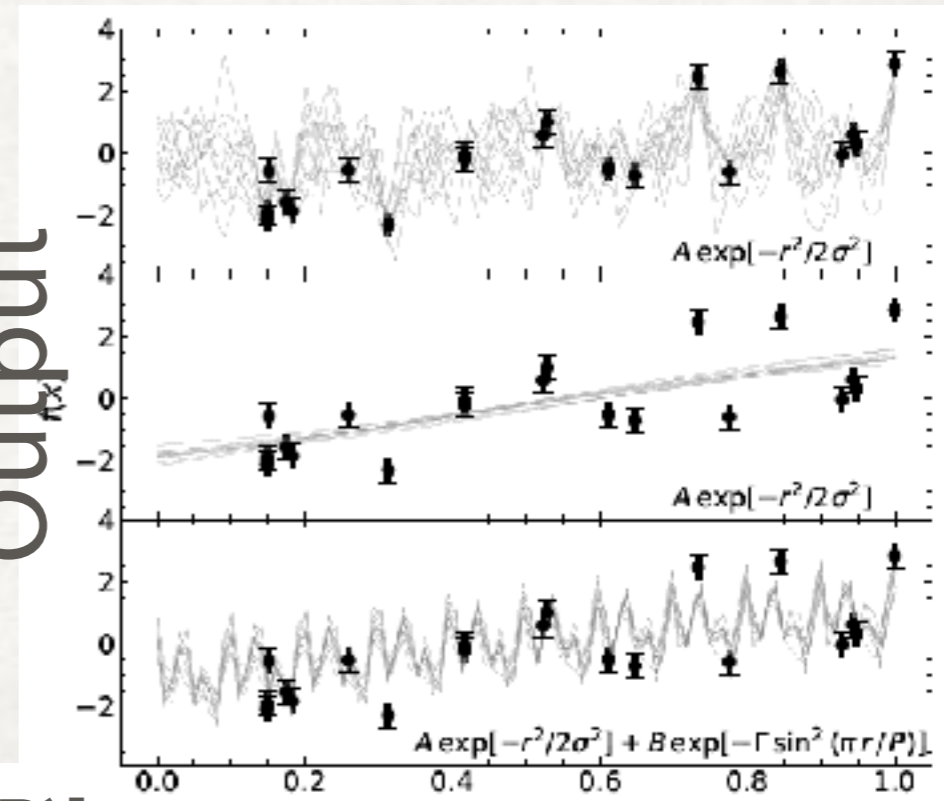
$$A \exp[-r^2/2\sigma^2]$$

Posterior

Output



Output



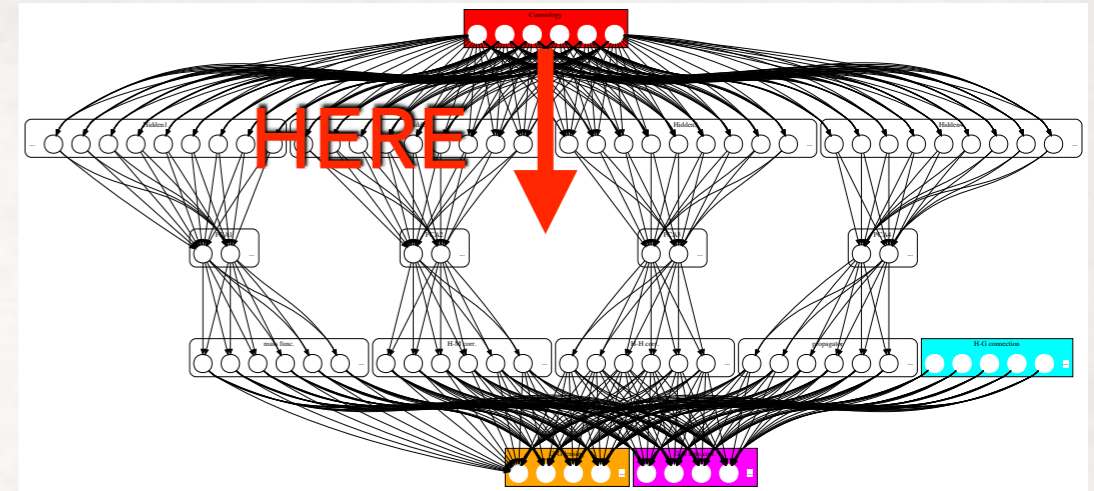
Input

$$+ B \exp[-\Gamma \sin^2(\pi r/P)]$$

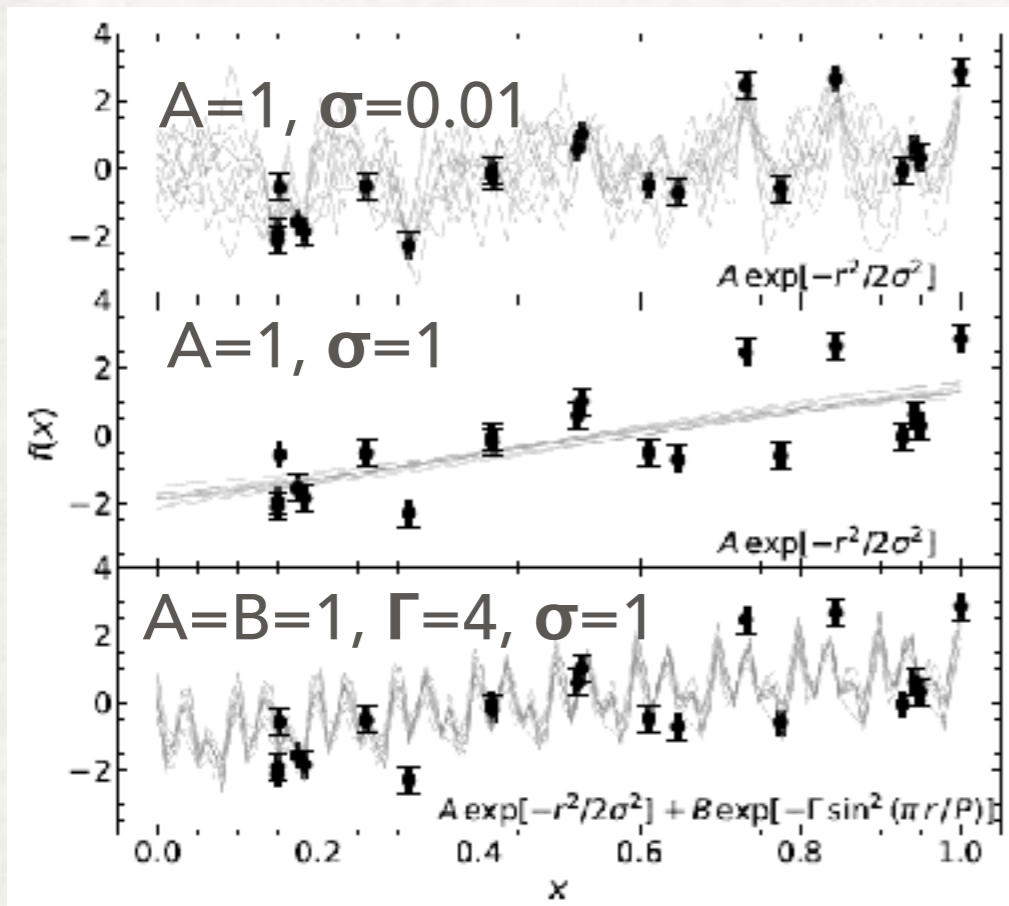
Input



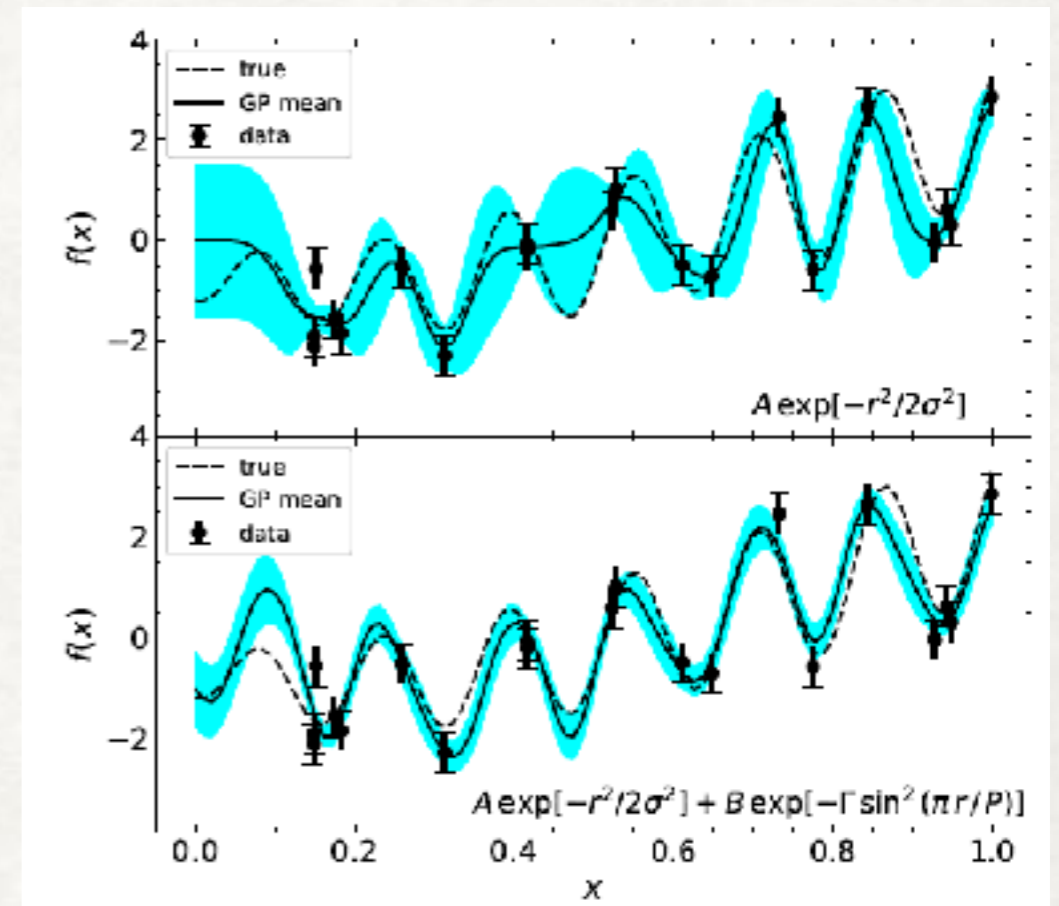
# GAUSSIAN PROCESS REGRESSION



Posterior (before optimization)



Posterior (after optimization)

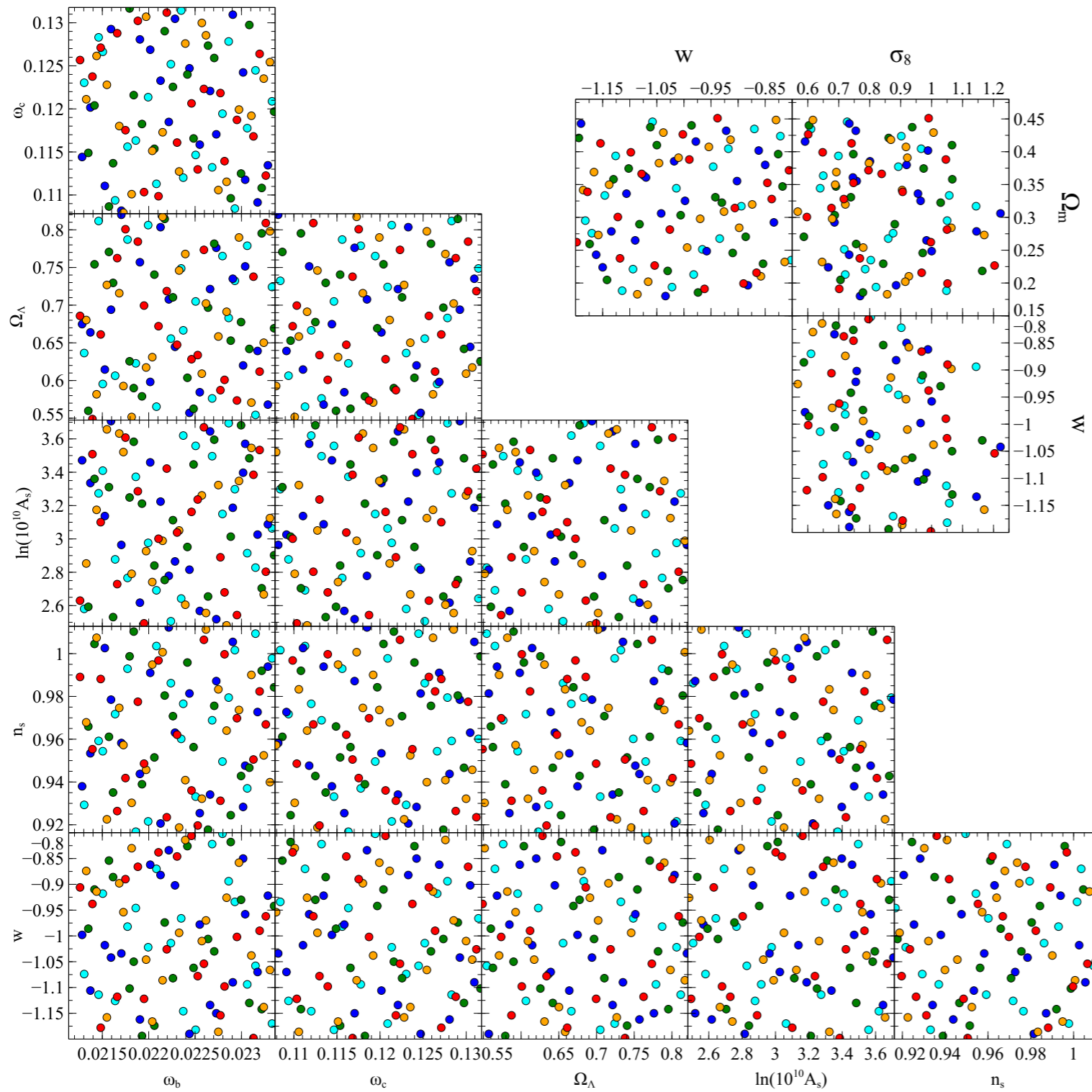


- All the quantities are analytic thanks to Gaussianity
- "Hyperparameters" can be "trained"

$$P(t_{N+1} | \mathbf{t}_N) \propto \exp \left[ -\frac{1}{2} \begin{bmatrix} \mathbf{t}_N & t_{N+1} \end{bmatrix} \mathbf{C}_{N+1}^{-1} \begin{bmatrix} \mathbf{t}_N \\ t_{N+1} \end{bmatrix} \right]$$

$$\begin{aligned} \hat{t}_{N+1} &= \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N \\ \sigma_{\hat{t}_{N+1}}^2 &= \kappa - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k} \end{aligned}$$

# CROSS VALIDATION STUDIES

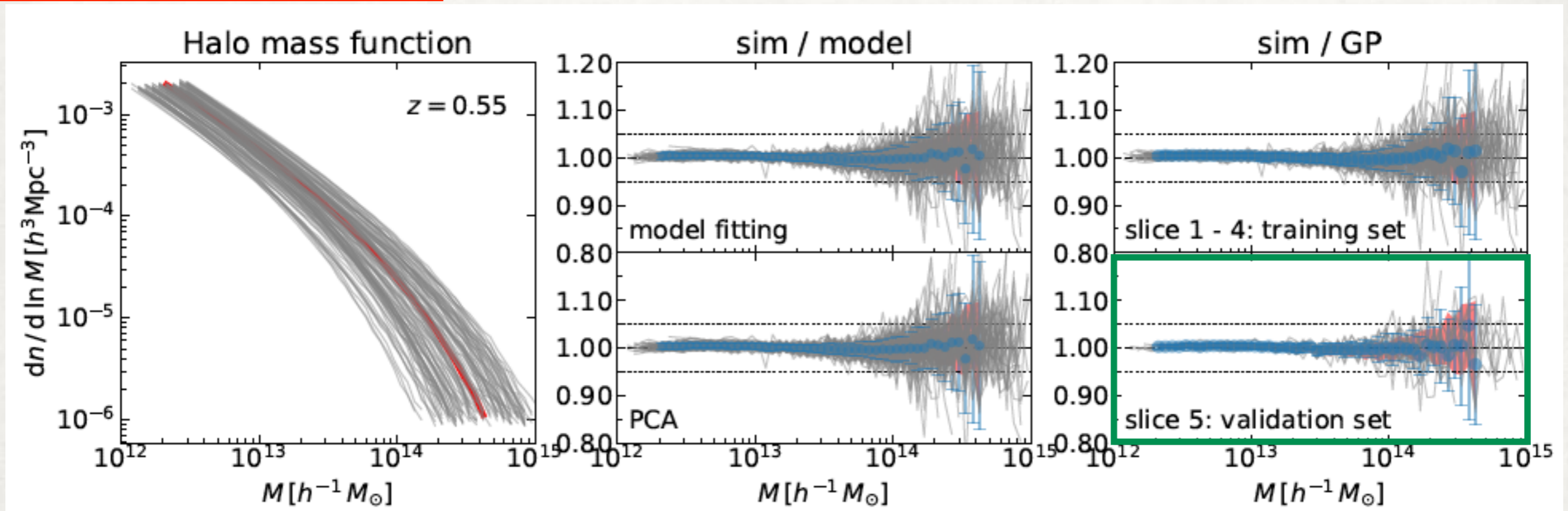


- Train the network using 80 points in 4 different colors.
- Validate the prediction of the network using the remaining 20 points in cyan.



# ACCURACY: HALO MASS FUNCTION

Example plot at  $z = 0.55$



Spread in HMF among the 100 models

Gaussian Process Regression

Upper: Model fitting w/ Sheth-Tormen type function (2 free parameters)

Lower: Compress the 42 (=2 x 21 redshifts) coefficients into 6 PCs

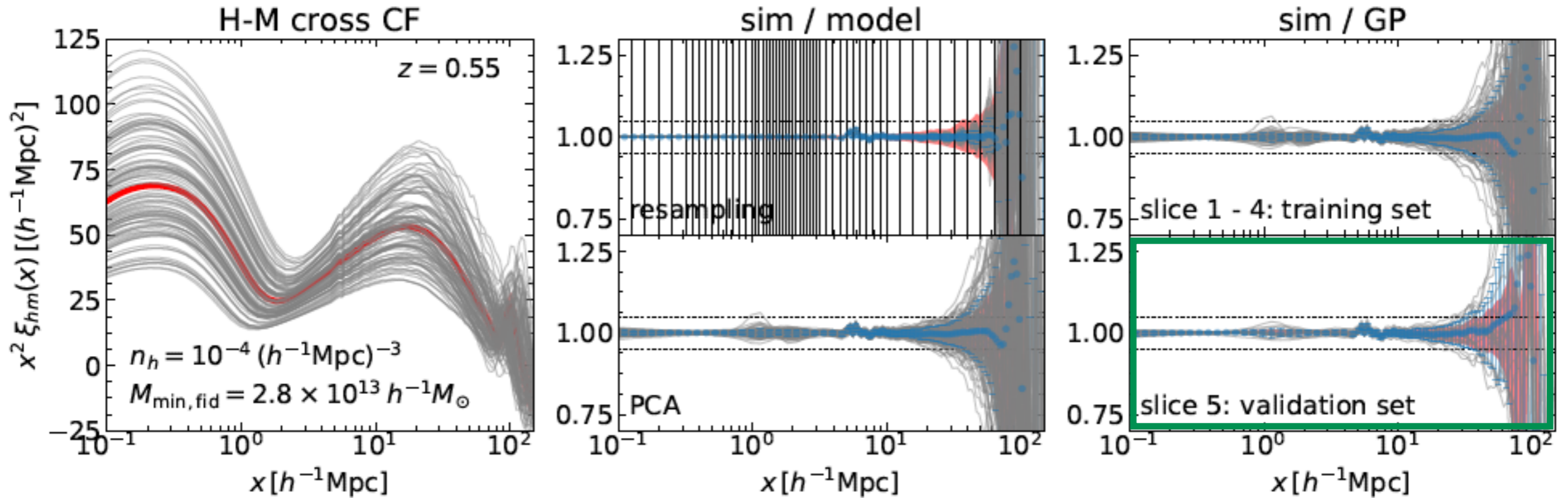
Training set

Validation set

Red shades: scatter of 28 fiducial runs

# ACCURACY: CROSS CORRELATION

Example plot at  $z = 0.55$  for a halo sample with  $10^{-4} (h^{-1}\text{Mpc})^{-3}$



Spread in  $\xi_{hm}$  among the 100 models

Gaussian Process Regression

Upper: Sample fewer number of data points using cubic spline interpolation

Training set

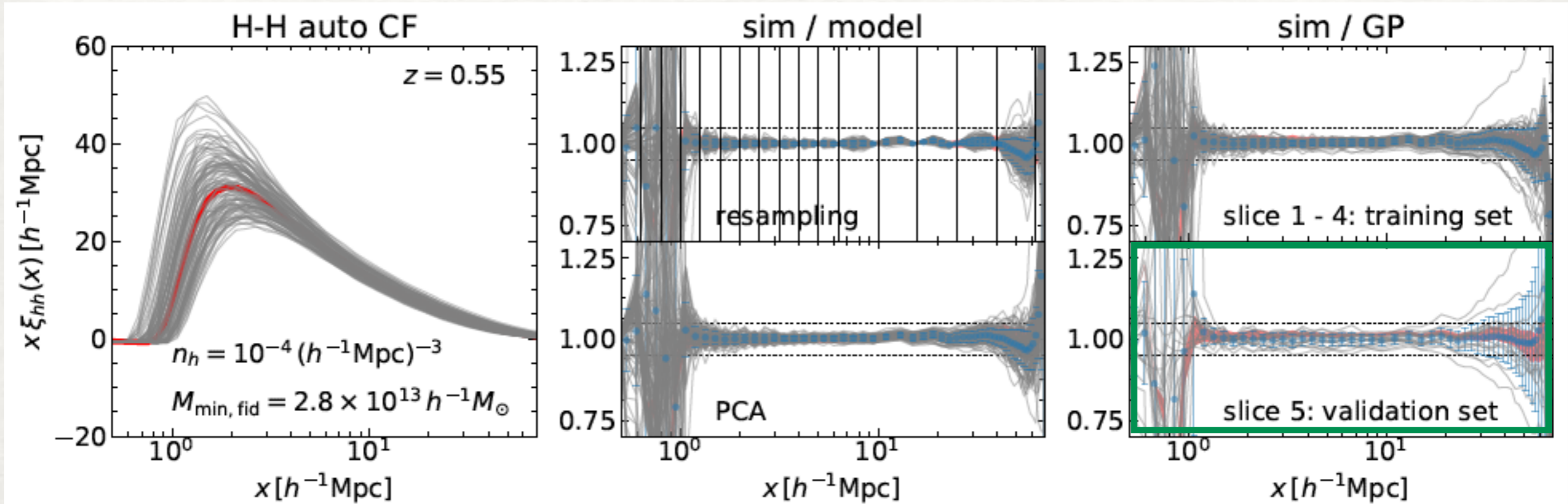
Lower: Compress the 18,018 (=  $66 \times 13 \times 21$  radial x number density x redshift bins) coefficients into 5 PCs

Validation set



# ACCURACY: AUTO CORRELATION

Example plot at  $z = 0.55$  for a halo sample with  $10^{-4} (h^{-1}\text{Mpc})^{-3}$



Spread in  $\xi_{hh}$  among the 100 models

Gaussian Process Regression

Upper: Sample fewer number of data points using cubic spline interpolation

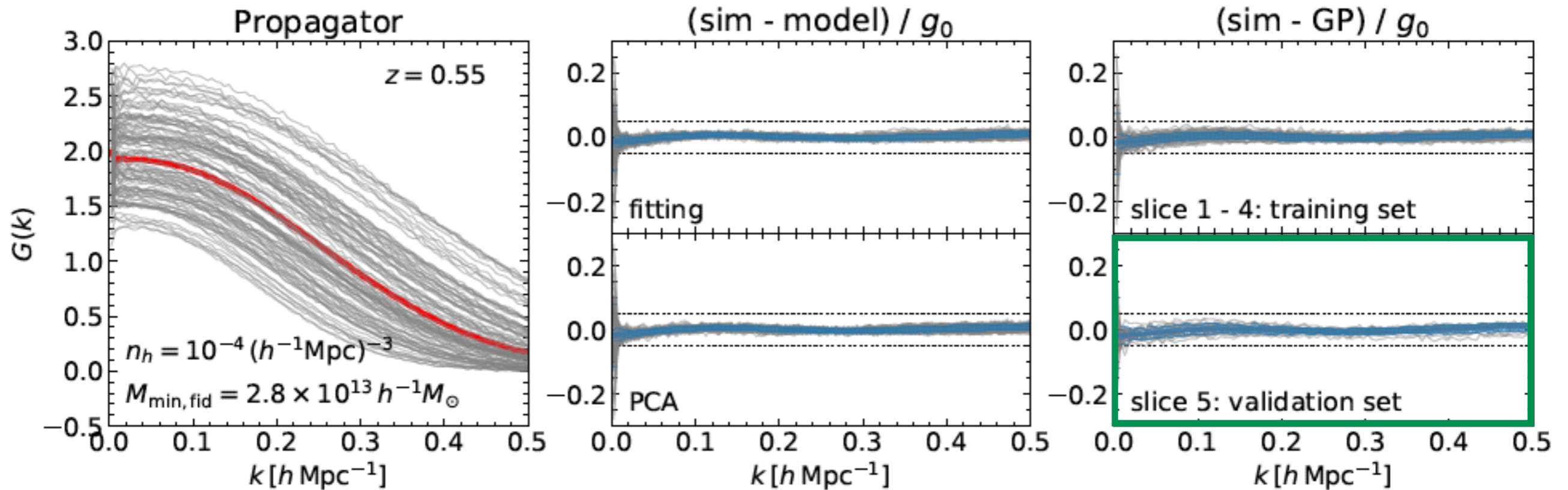
Training set

Lower: Compress the 28,224 ( $= 21 \times 8 \times 8 \times 21$  radial x number density 1 x number density 2 x redshift bins) coefficients into 8 PCs

Validation set

# ACCURACY: PROPAGATOR

Example plot at  $z = 0.55$  for a halo sample with  $10^{-4} (h^{-1}\text{Mpc})^{-3}$



Spread in  $G(k)$  among the 100 models

Gaussian Process Regression

Upper: Model fitting with Gaussian + corrections (3 free params)

Training set

Lower: Compress the 63 (=3 x 21 redshifts) coefficients into PCs

Validation set

Red shades: scatter of 14 fiducial runs



# DARK EMULATOR: WHAT IT CAN DO

## OVERVIEW

```
In [18]: import darkemu
```

```
In [19]: emu = darkemu.base_class()
```

initialize cosmo\_class  
 Initialize xlin emulator  
 initialize xnl emulator  
 Initialize pklin emulator  
 initialize propagator emulator  
 Initialize sigma\_d emulator  
 initialize cross-correlation emulator  
 initialize auto-correlation emulator  
 Initialize hmf emulator  
 Initialize sigmaM emulator

$(\omega_b, \omega_c, \Omega_{de}, \ln(10^{10} A_s), n_s, w)$

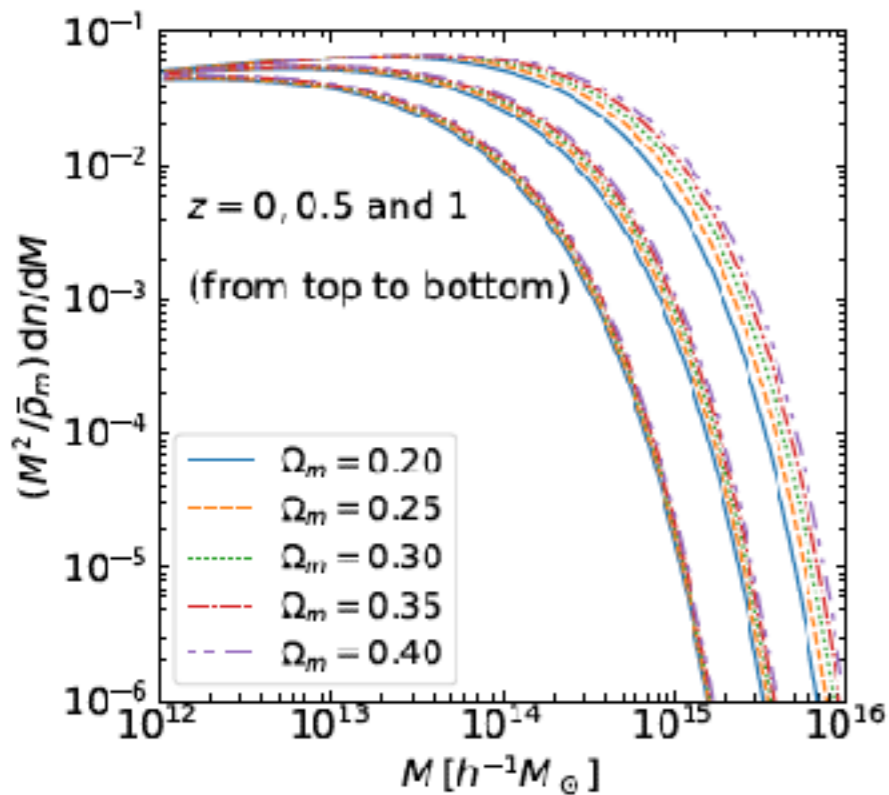
```
In [14]: cparam = np.array([0.02225, 0.1198, 0.6844, 3.094, 0.9645, -1.])
          emu.set_cosmology(cparam)
```

```
emu.get_nhalo(massbins[ii], massbins[ii+1], 1., z)
```

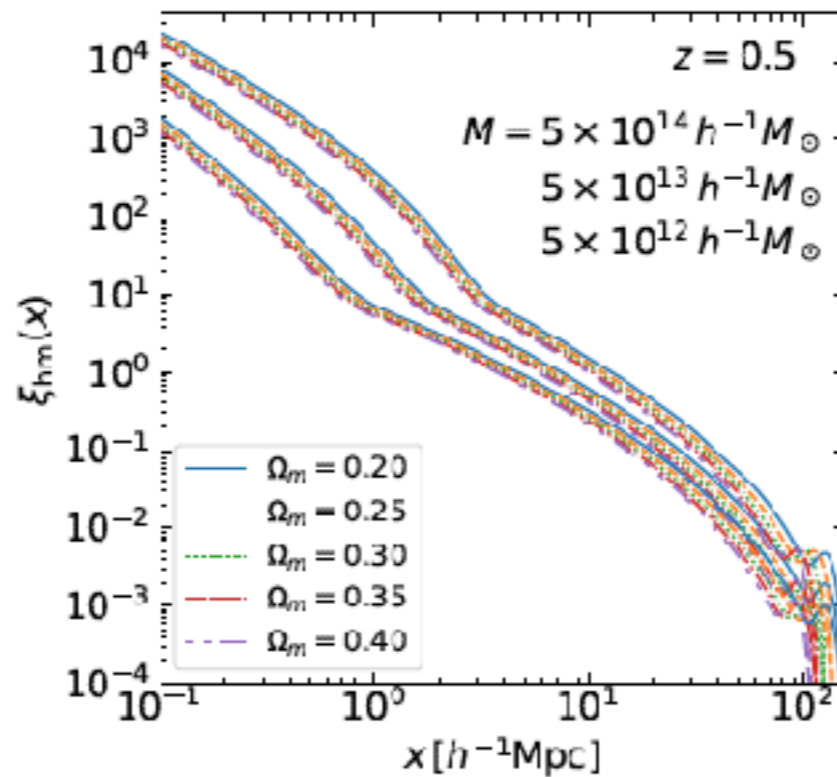
```
emu.get_xicross_mass(rs, Mh, z)
```

```
emu.get_xiauto_mass(rs, Mh, Mh, z)
```

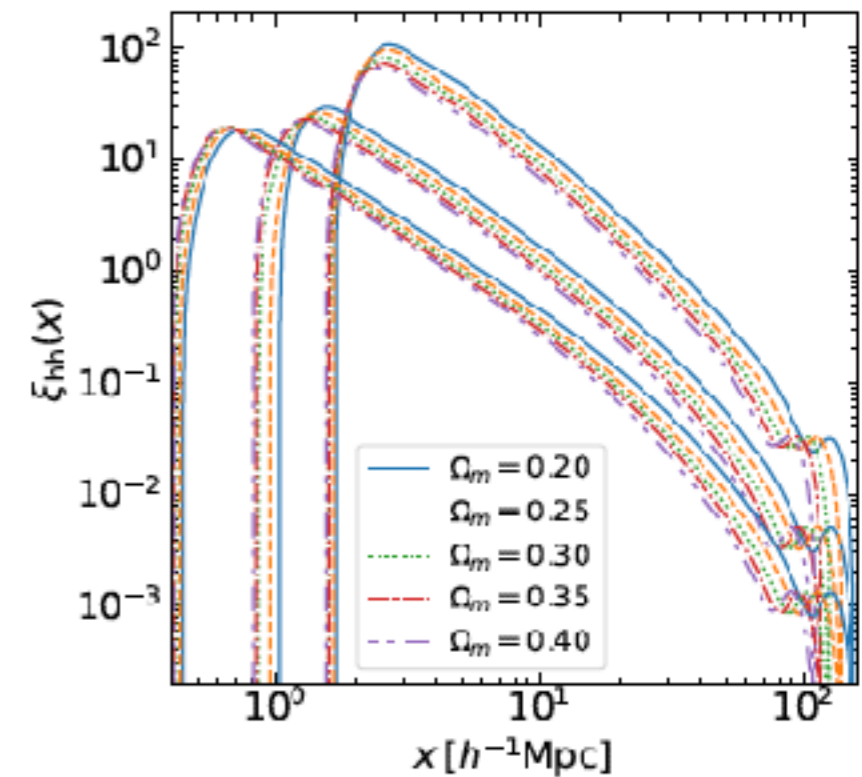
### Halo mass function



### Halo-Matter Cross CF

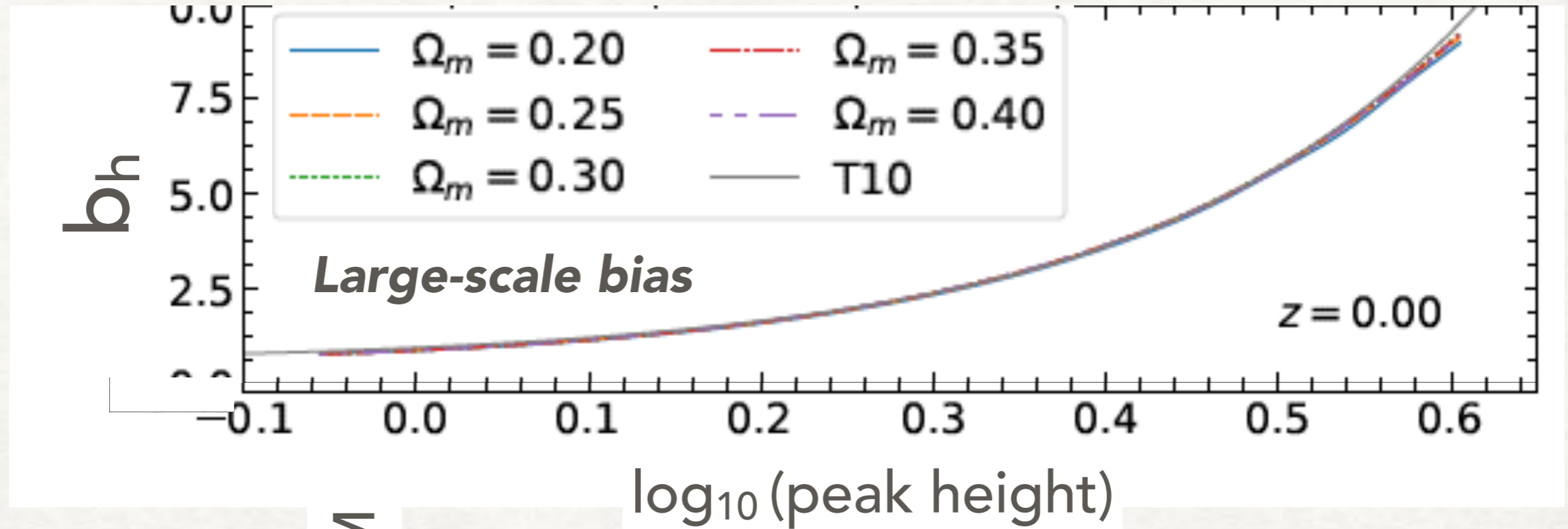


### Halo-Halo Auto CF



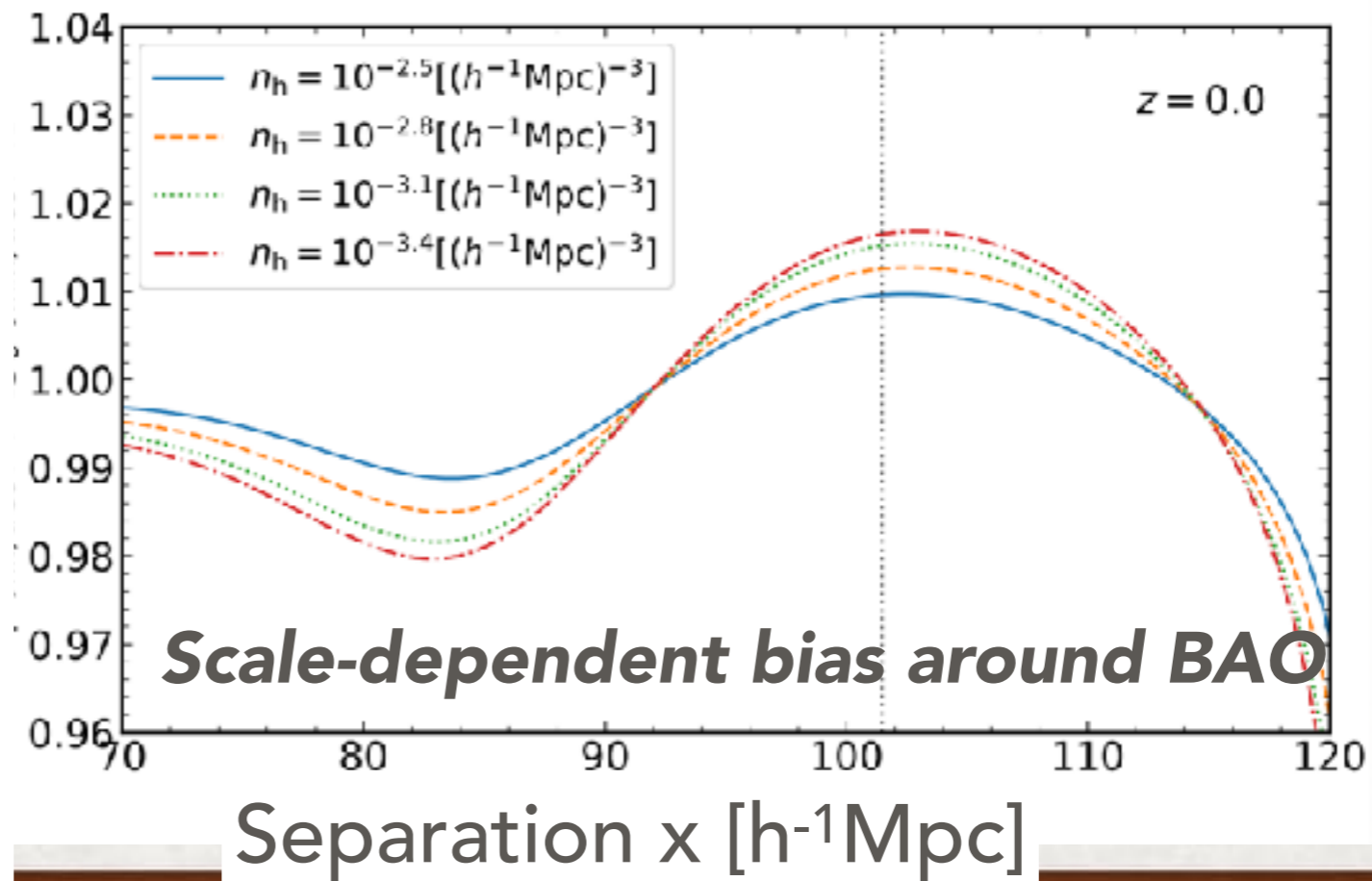
# DARK EMULATOR: WHAT IT CAN DO

## LARGE SCALES



$$\frac{1}{b_h} \sqrt{\frac{\xi_{hh}(x)}{\xi_{mm}(x)}}$$

BAO bump relative to DM

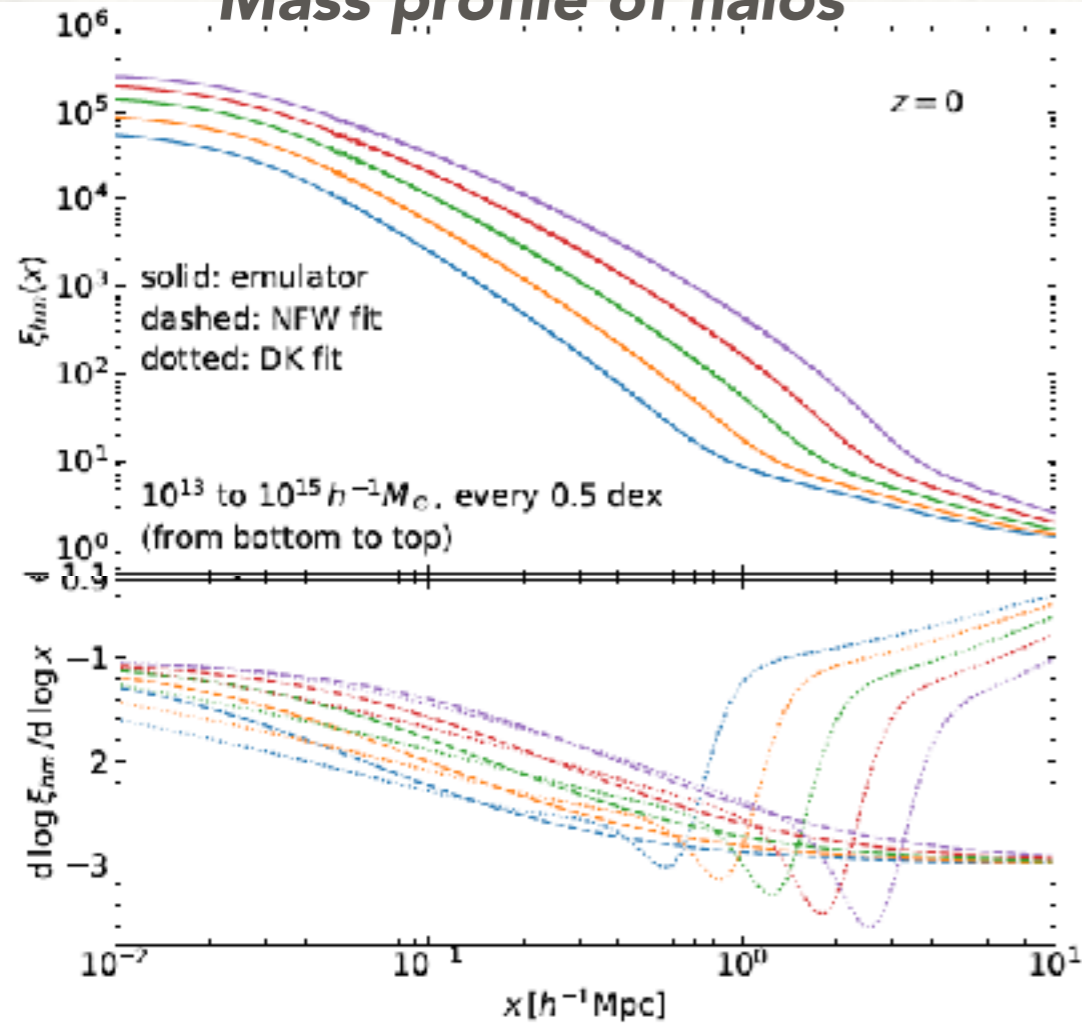




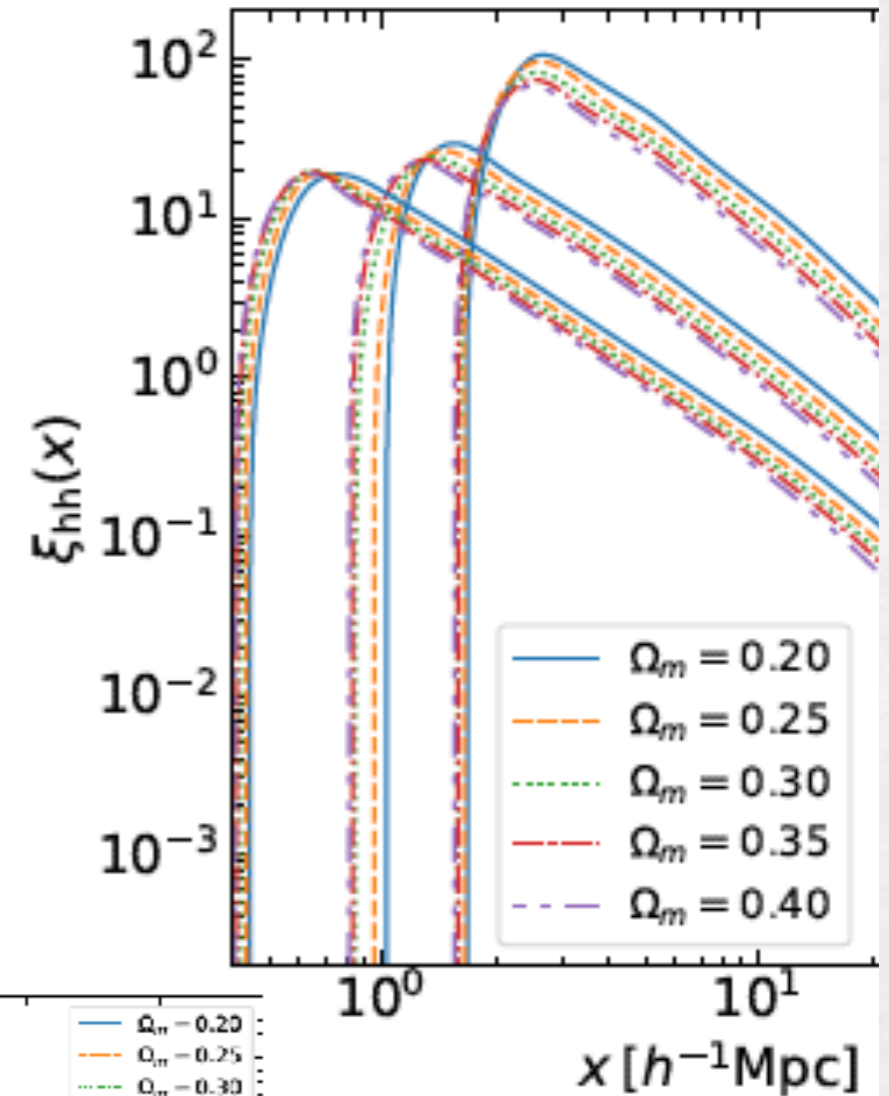
# DARK EMULATOR: WHAT IT CAN DO

## SMALL SCALES

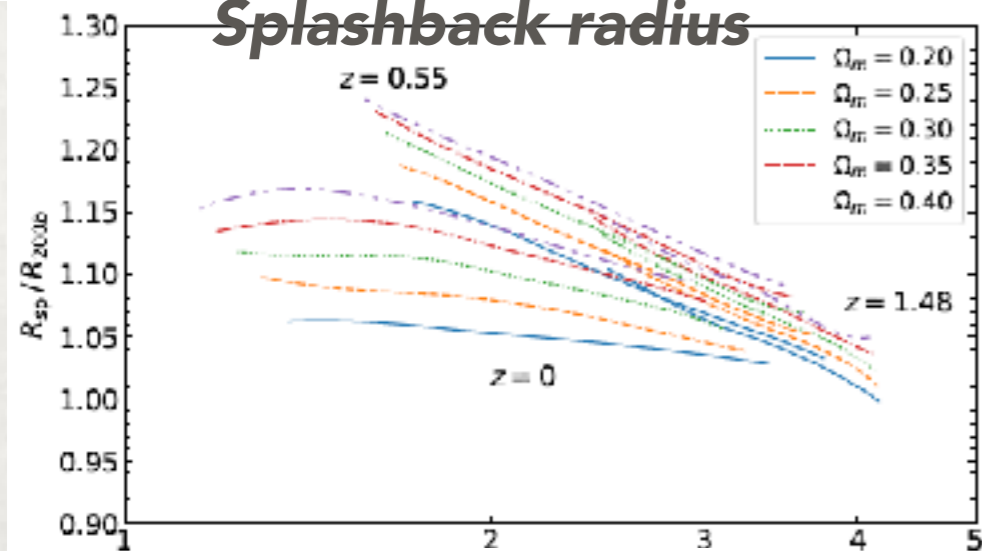
**Mass profile of halos**



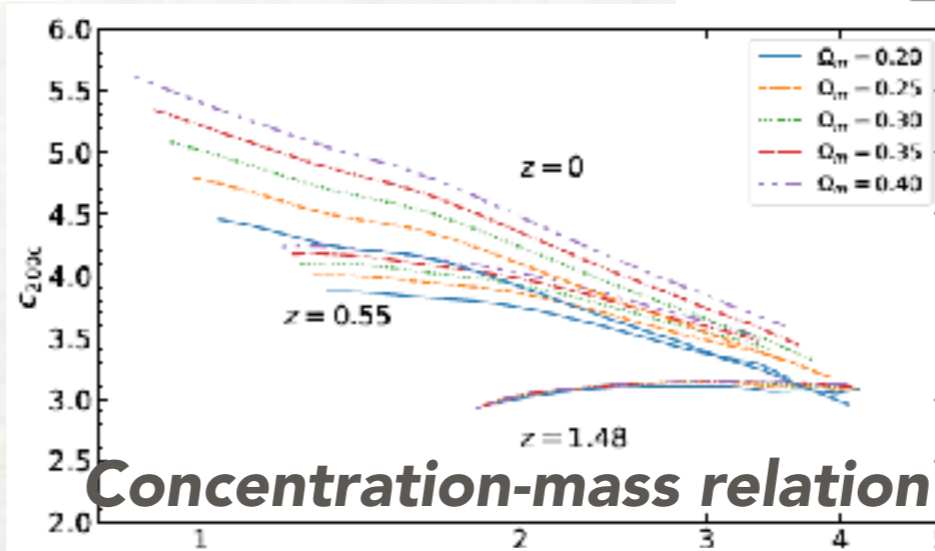
**Halo exclusion effect**



**Splashback radius**



**Concentration-mass relation**



peak height

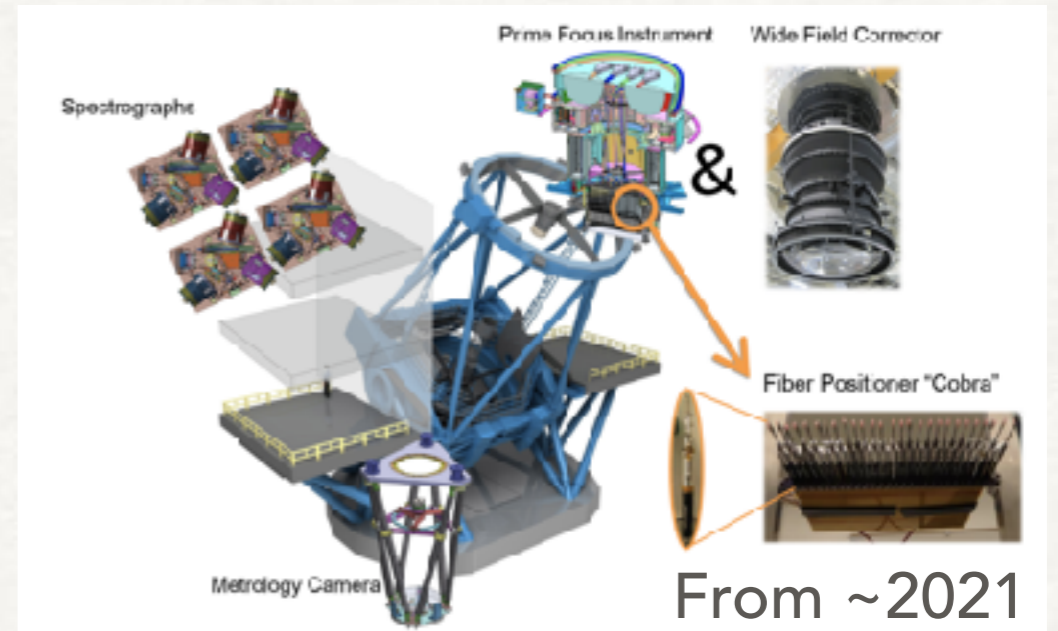
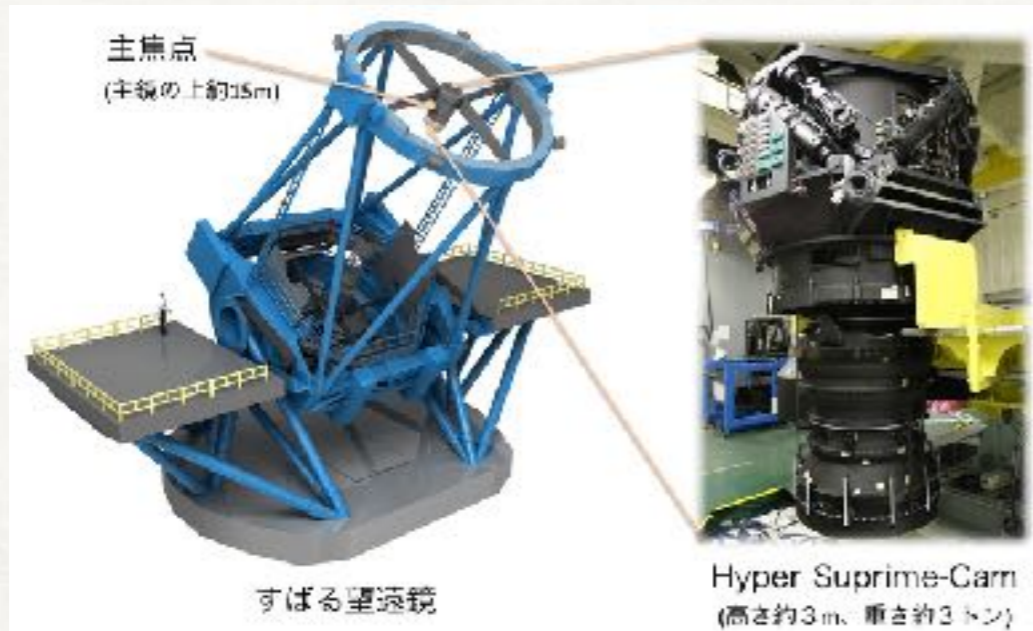
peak height

# BREAK THE DEGENERACY

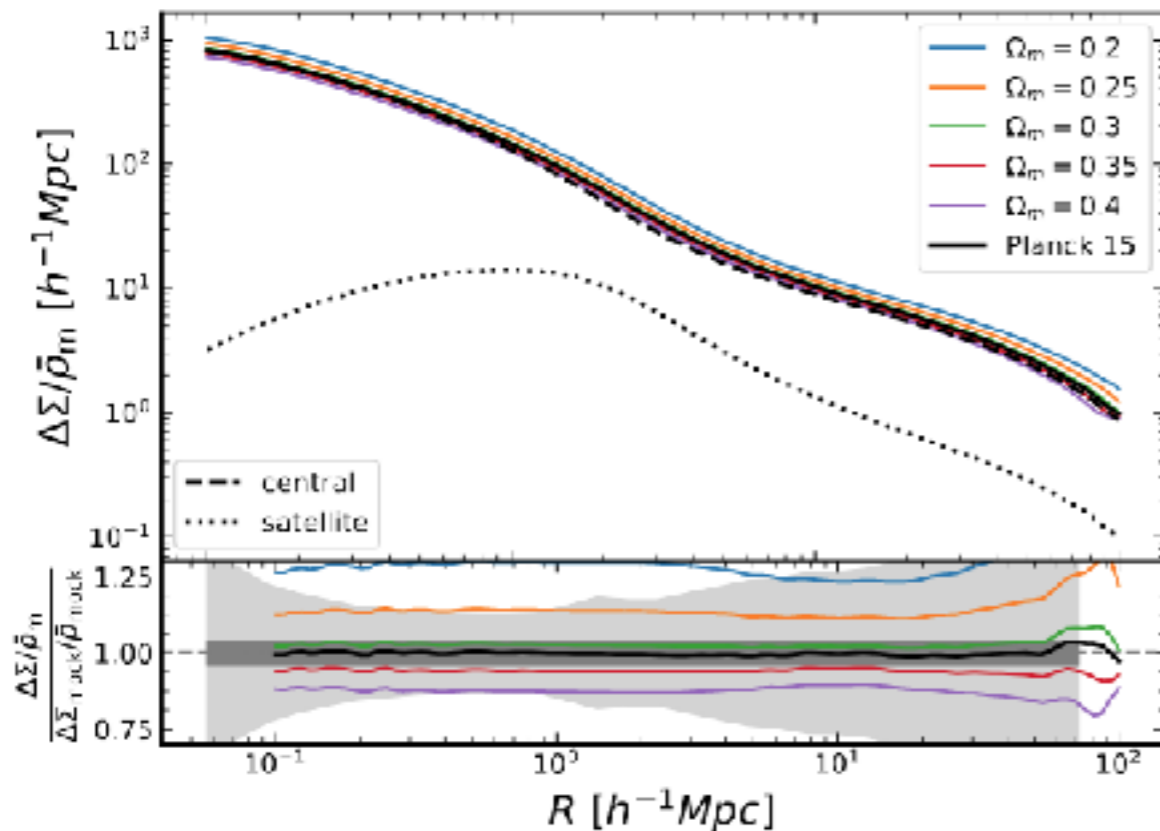
Subaru Measurement of Images and Redshifts (SuMIRe) Project

Hyper Suprime Cam (HSC)

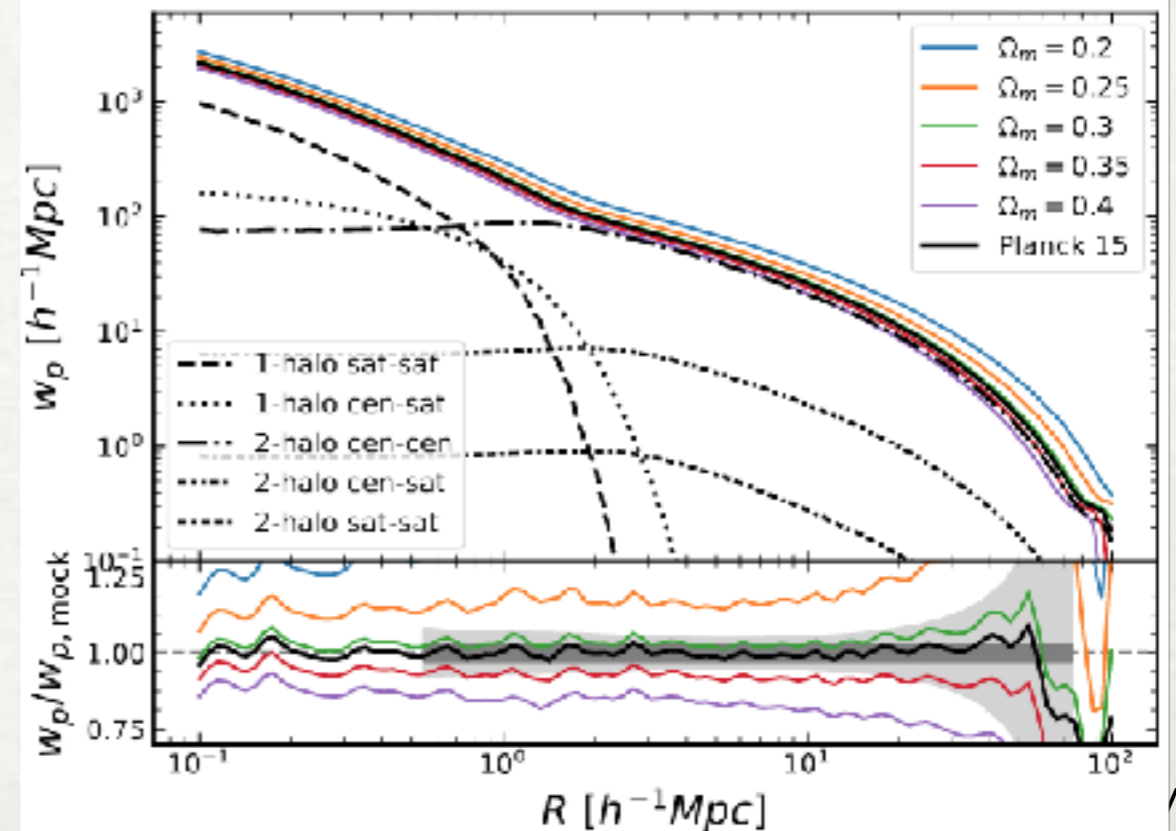
Prime Focus Spectrograph (PFS)



Galaxy-galaxy lensing



Galaxy (projected) clustering



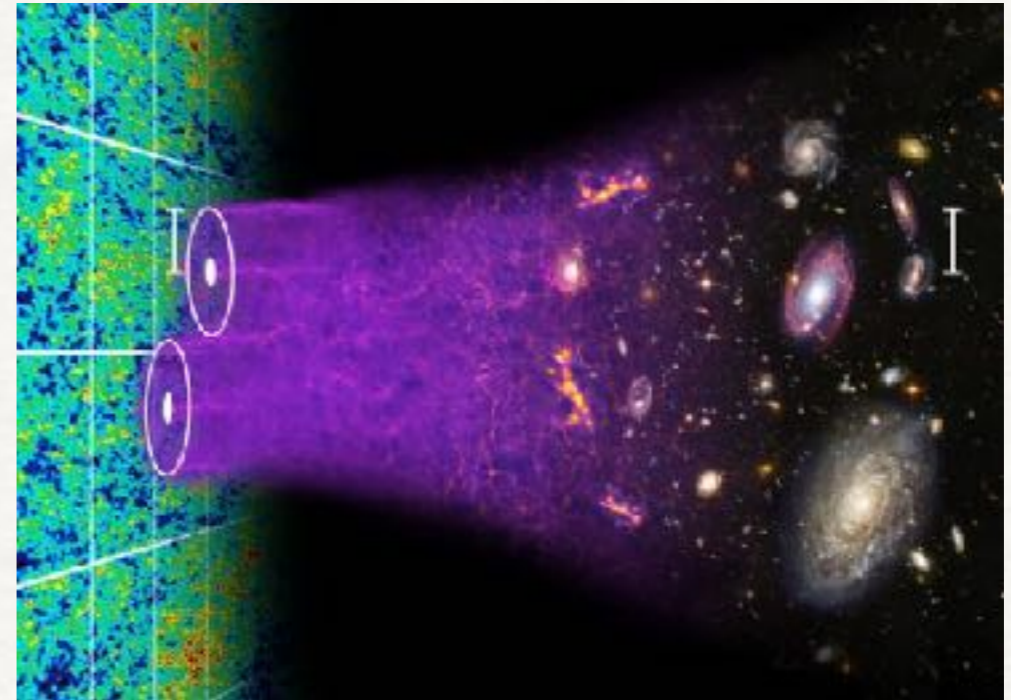
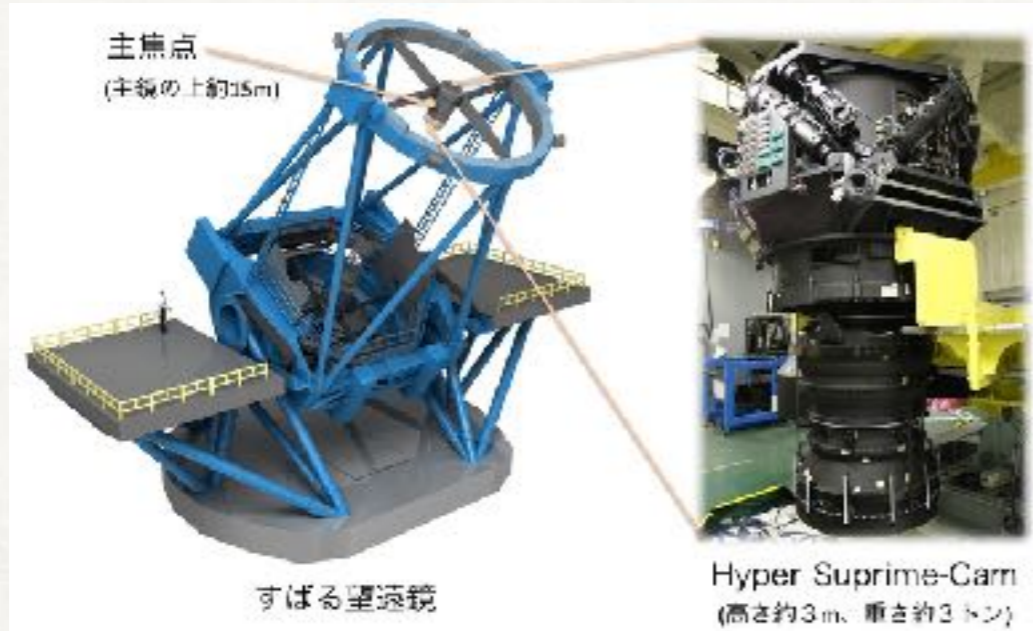


# BREAK THE DEGENERACY

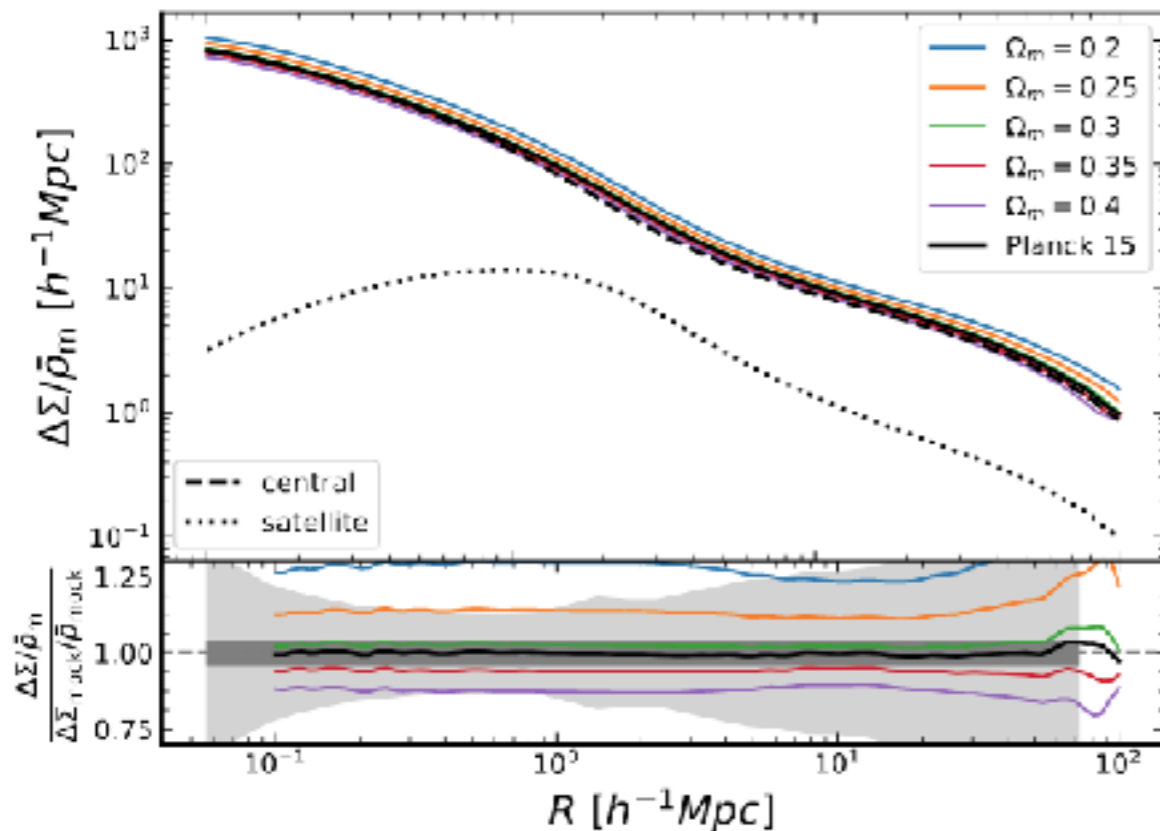
Subaru Measurement of Images and Redshifts (SuMIRe) Project

Hyper Suprime Cam (HSC)

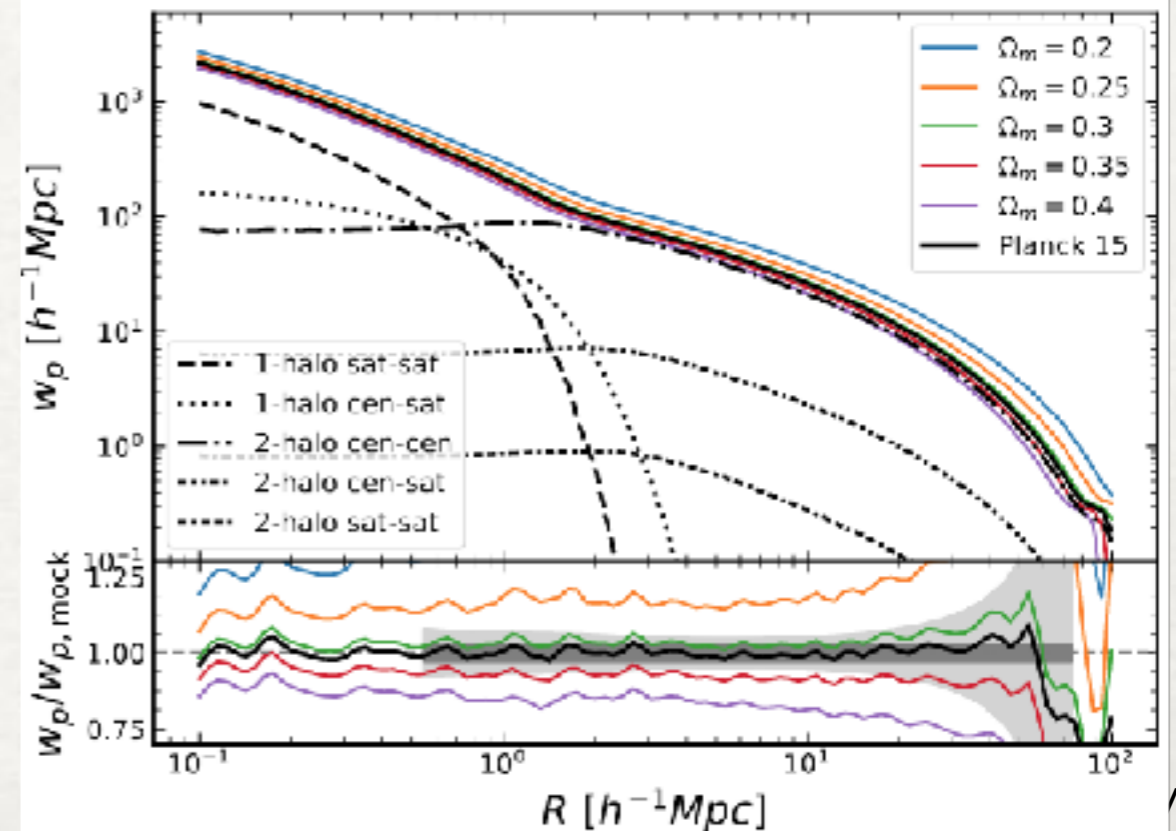
SDSS BOSS (CMASS + LOWZ)



Galaxy-galaxy lensing



Galaxy (projected) clustering



# BREAK THE DEGENERACY

Subaru Measurement of Images and Redshifts (SuMIRe) Project

Hyper Suprime Cam (HSC)

SDSS BOSS (CMASS + LOWZ)

主焦点

$$\Delta \Sigma_g(R) = \bar{\Sigma}_g(< R) - \Sigma_g(R),$$

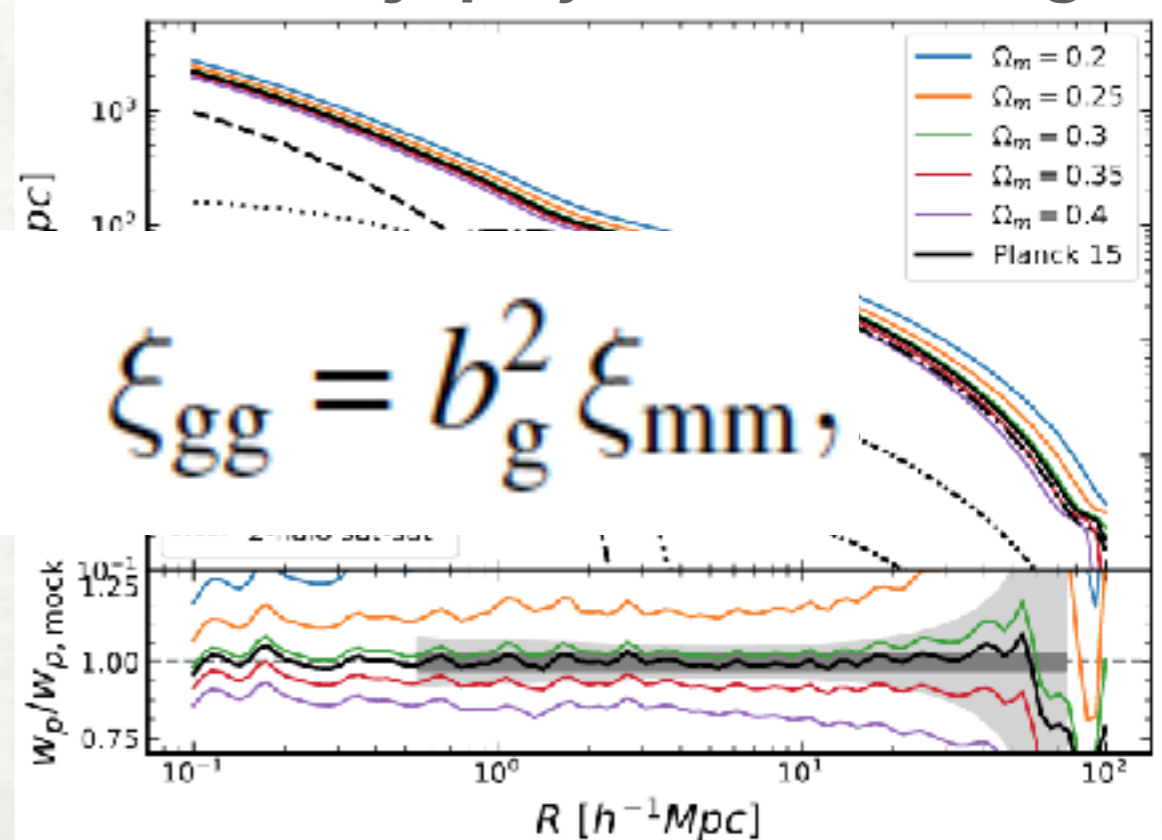
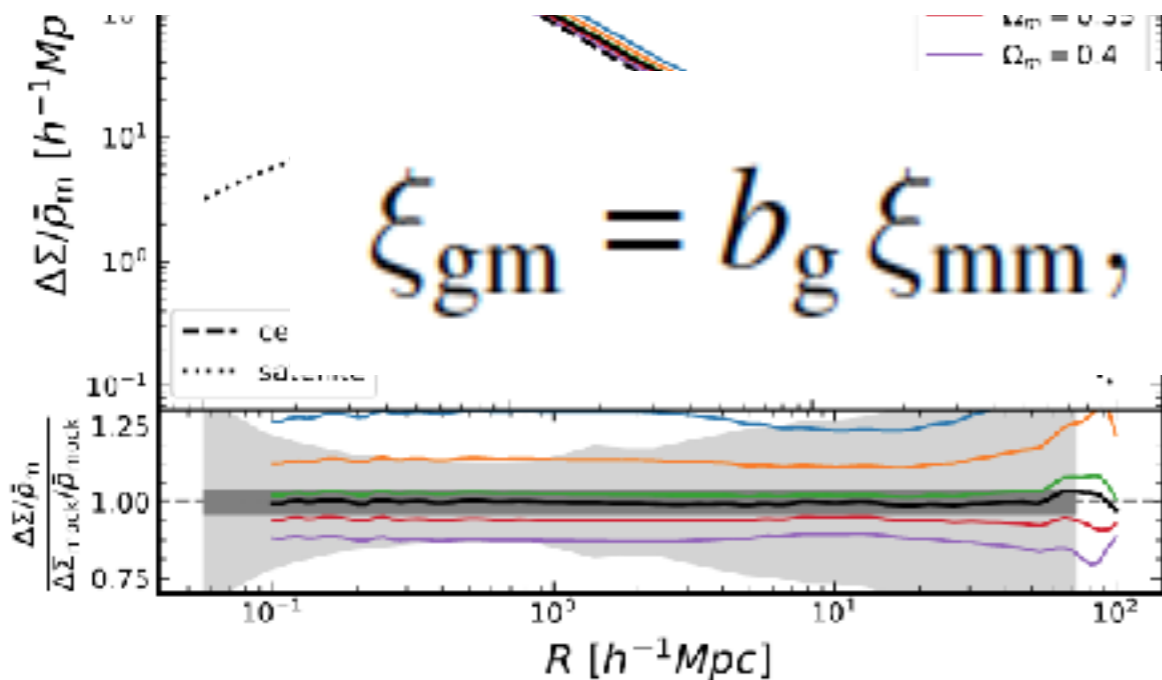
$$\bar{\Sigma}_g(< R) = \frac{2}{R^2} \int_0^R \Sigma_g(y) dy,$$

$$\Sigma_g(R) = \bar{\rho}_{m0} \int \xi_{gm}(R, \pi) d\pi,$$

$$w_{gg}(R) = 2 \int_0^{\pi_{\max}} \xi_{gg}(R, \pi) d\pi,$$

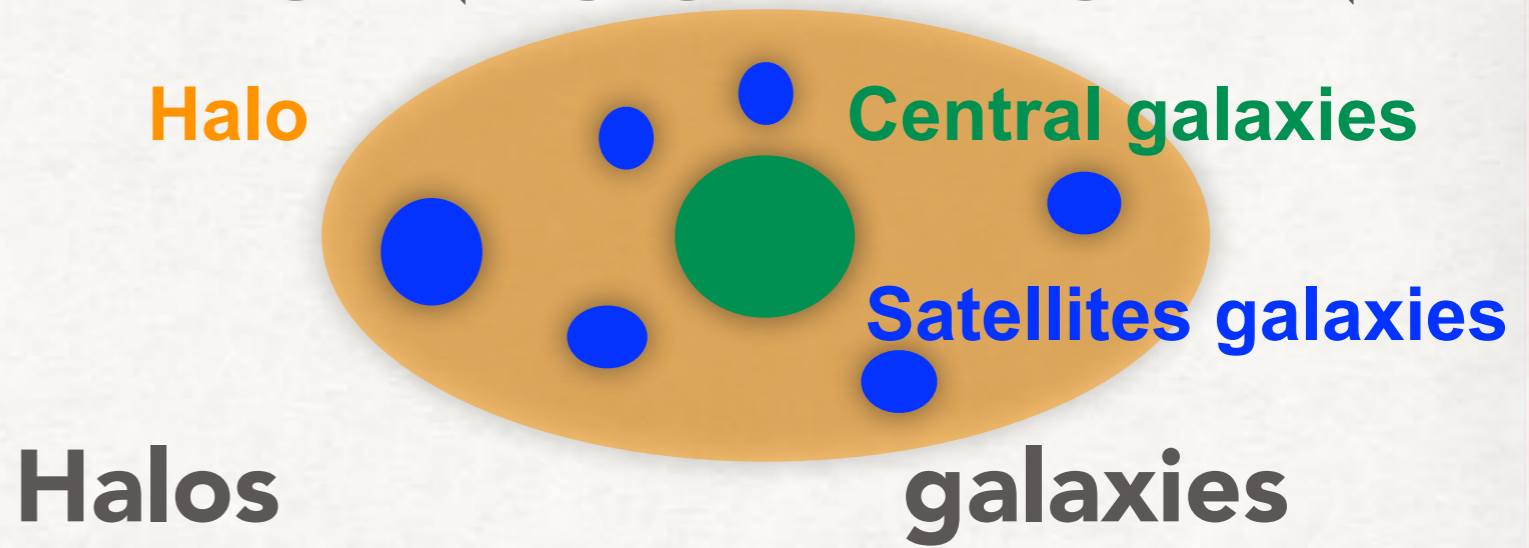
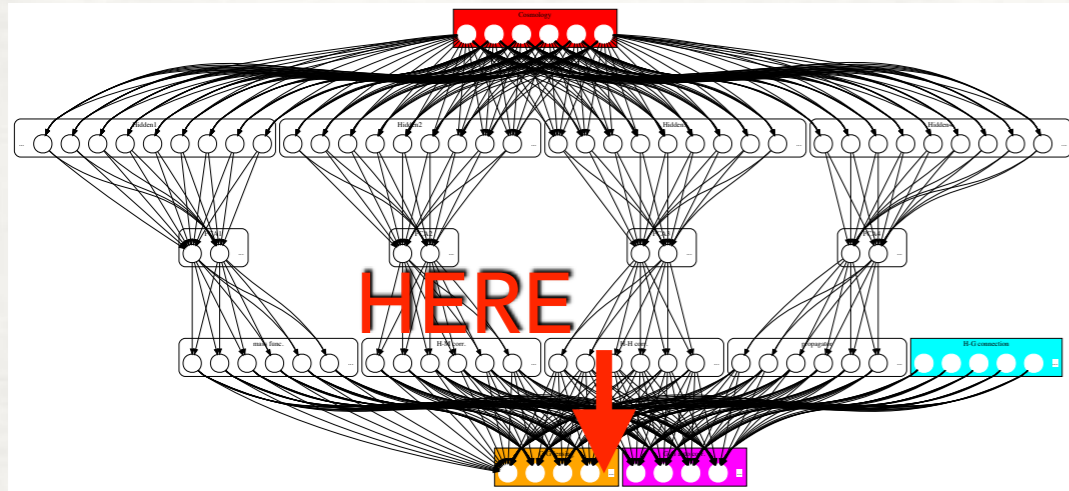


Galaxy (projected) clustering





# CROSS-CORRELATION COEFFICIENT

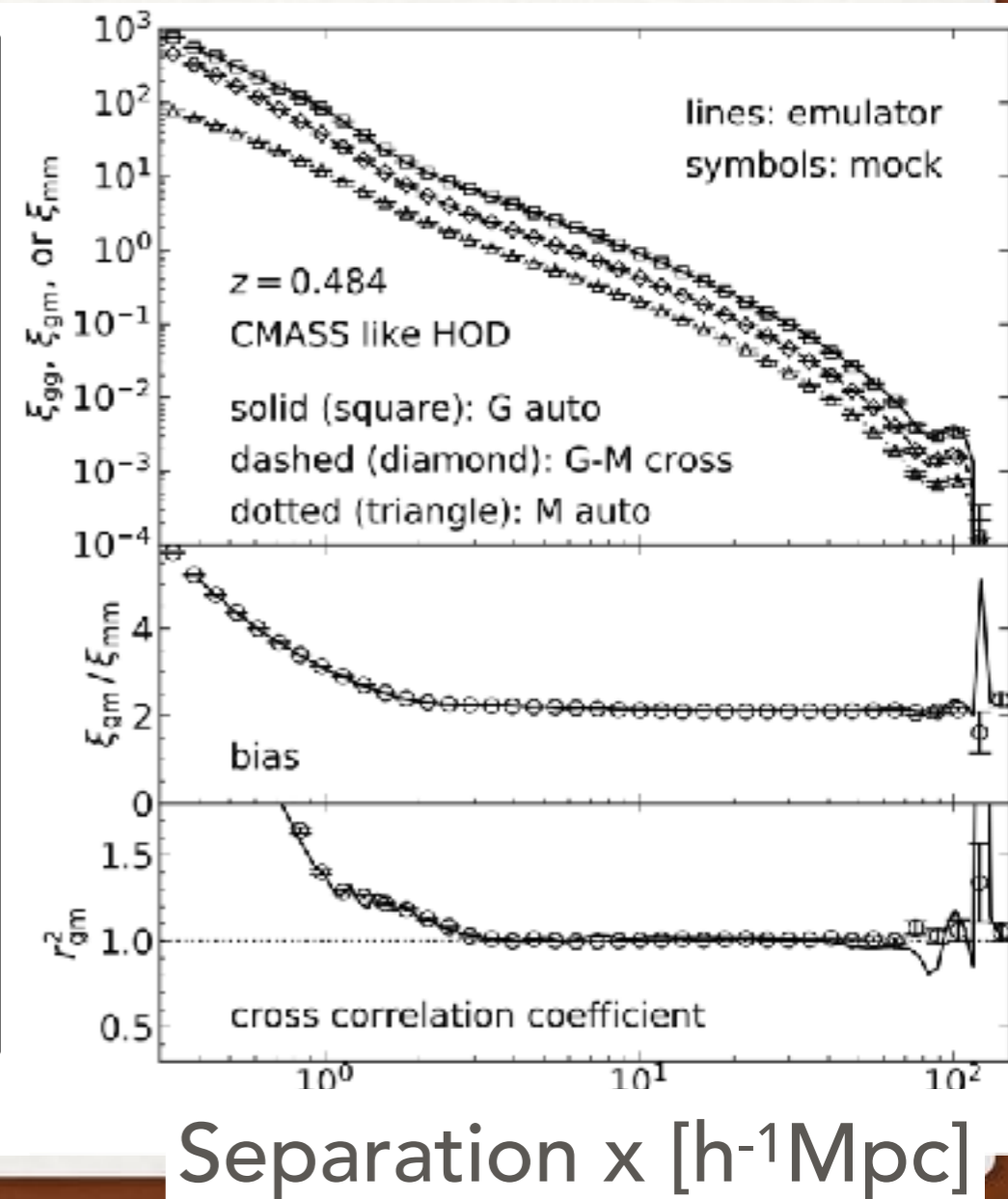
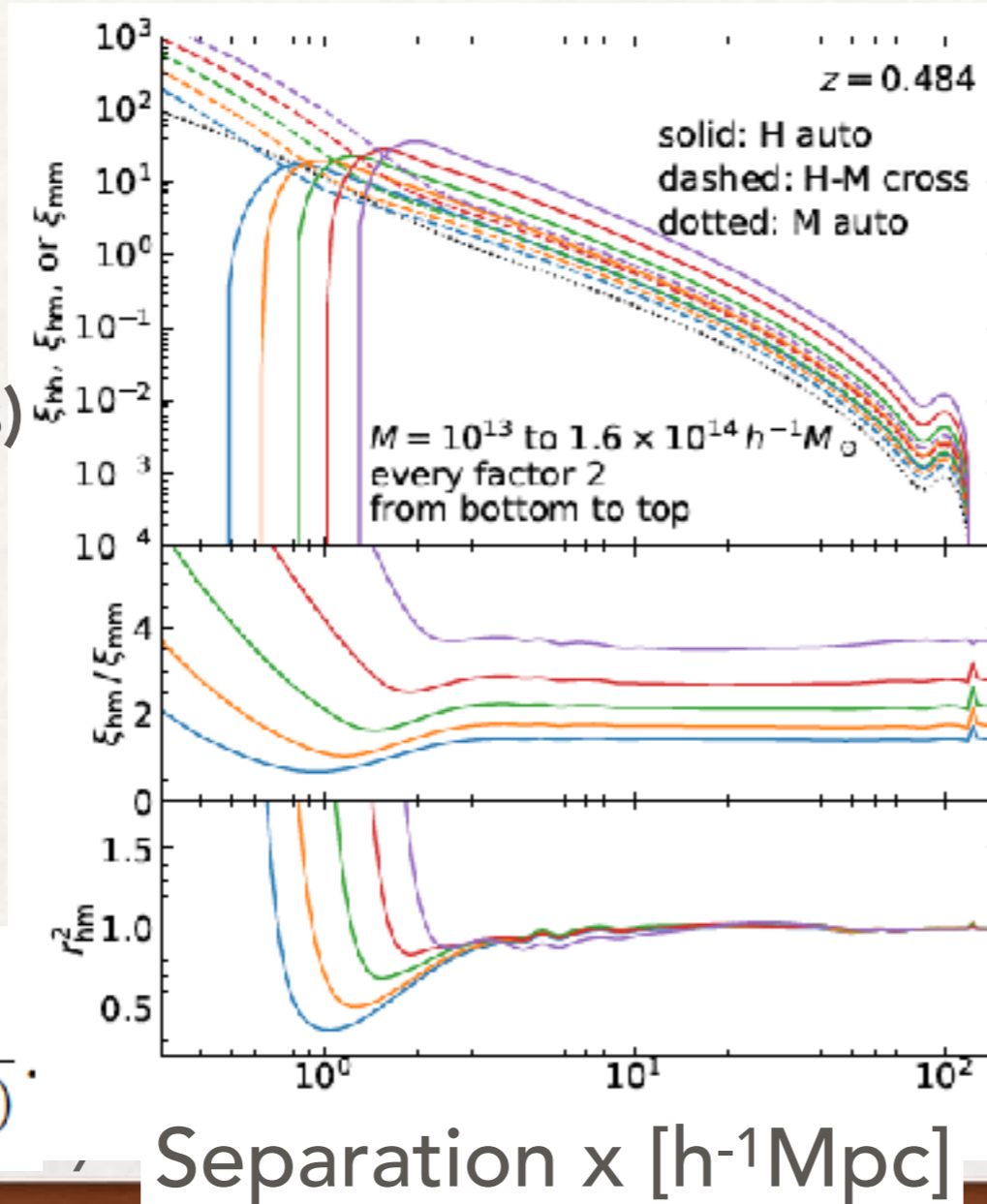


Correlation functions  
(matter & tracers)

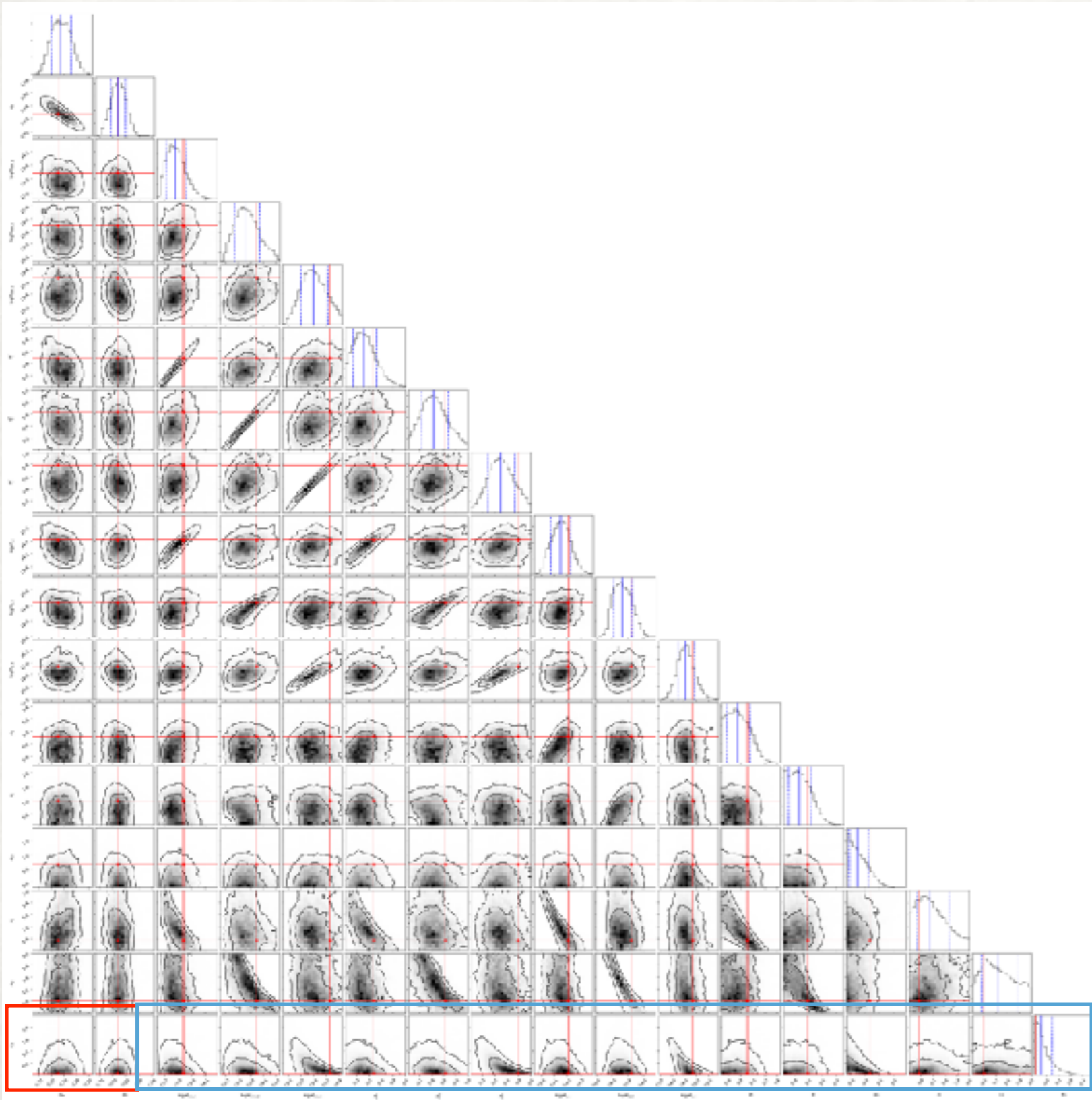
Bias

Correl. coeff.

$$r_{gm}^2(r) \equiv \frac{[\xi_{gm}(r)]^2}{\xi_{gg}(r)\xi_{mm}(r)}$$



# MOCK UNIVERSE VS EMULATOR

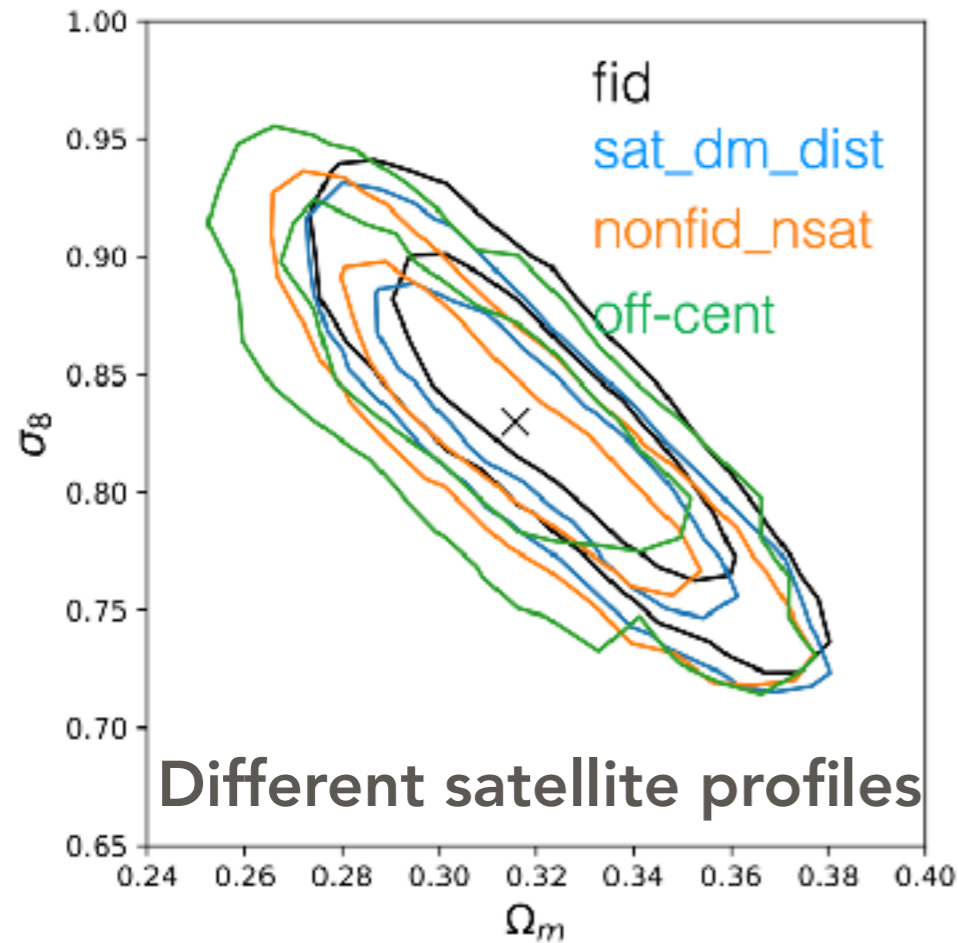


**COSMO**

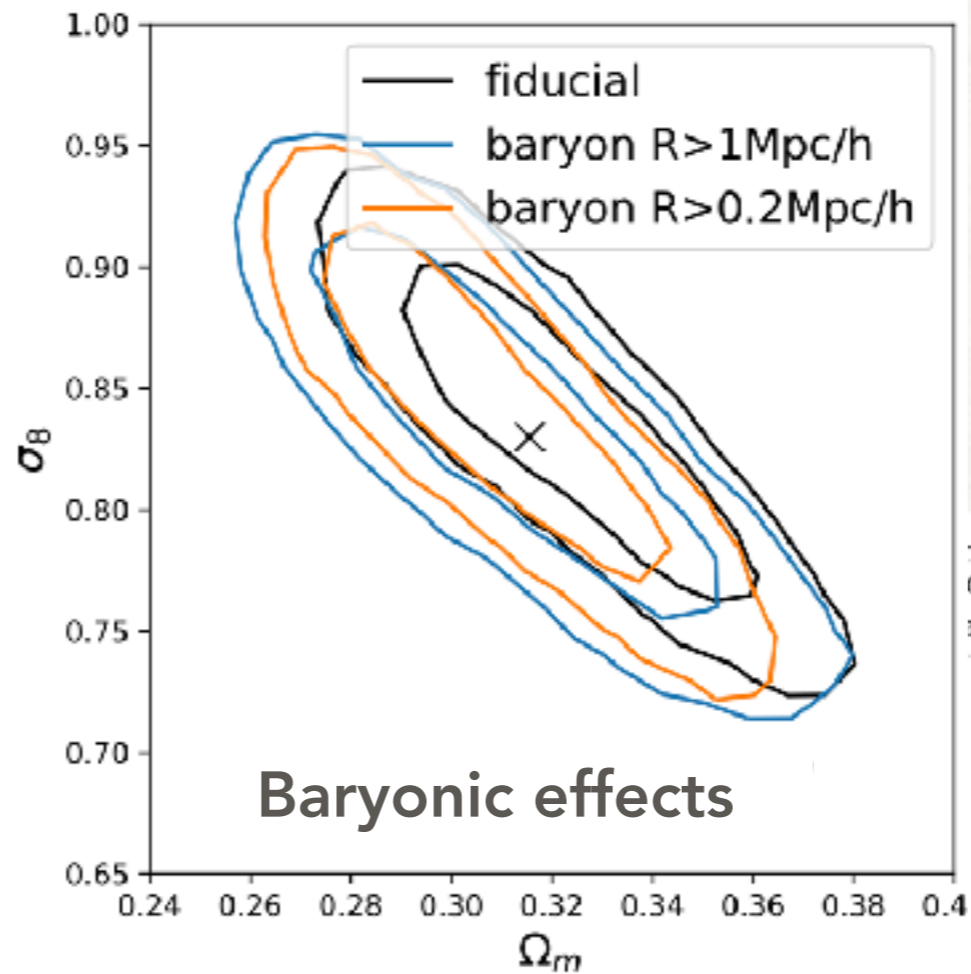
**GALAXY (HOD + extra)**



# MOCK UNIVERSE VS EMULATOR



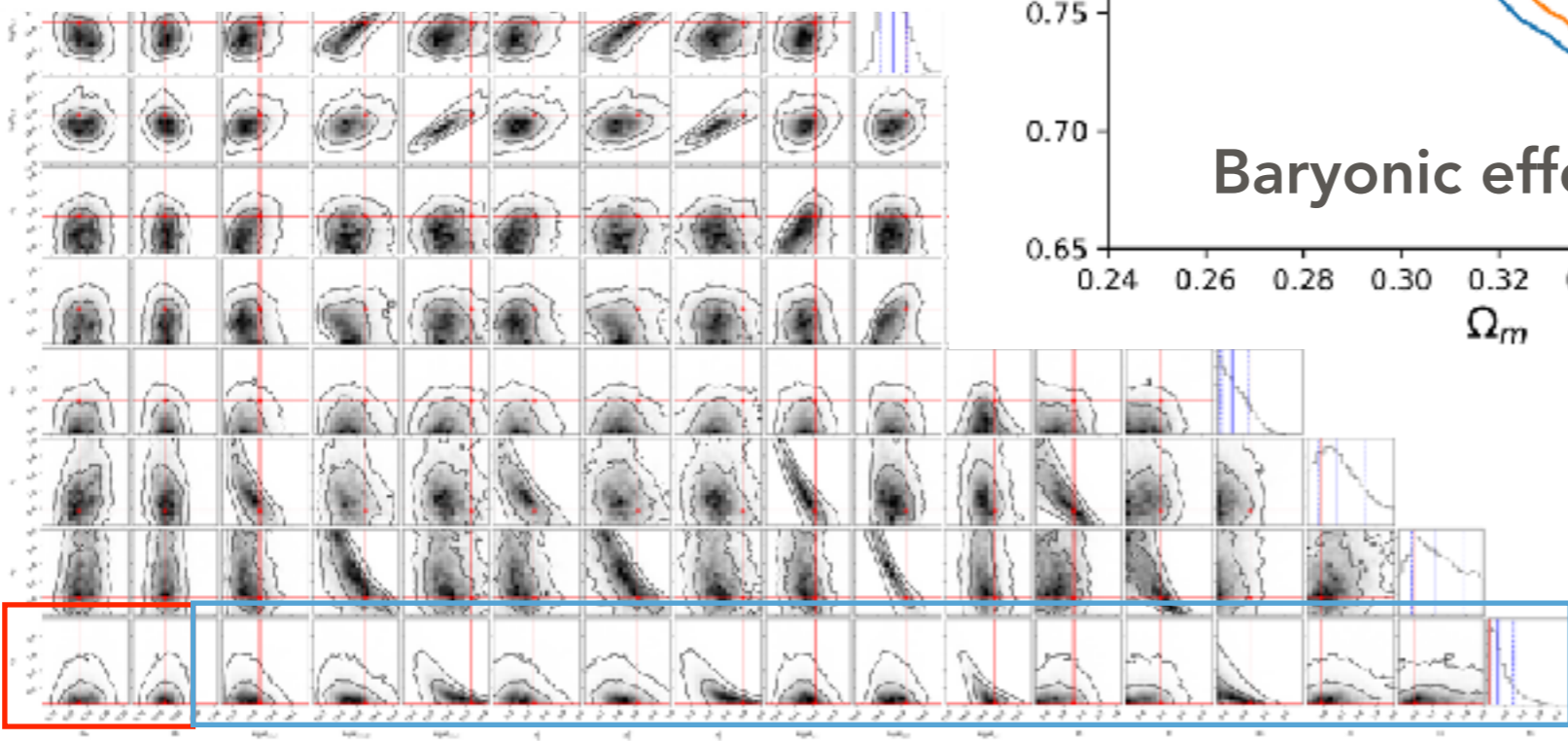
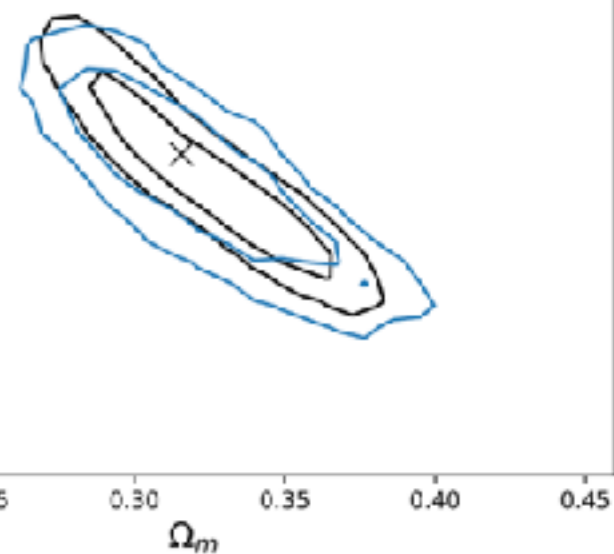
Different satellite profiles



Baryonic effects

fiducial:  $r_{\text{lens}} = [0.2, 30] \text{ Mpc}/h$ ,  $r_{\text{cl}} = [0.2, 30] \text{ Mpc}/h$   
 at\_subhalo:  $r_{\text{lens}} = [2, 30] \text{ Mpc}/h$ ,  $r_{\text{cl}} = [1, 30] \text{ Mpc}/h$

Substructures instead of HOD



COSMO

GALAXY (HOD + extra)

# SUMMARY

- Modeling halo clustering signal based on simulations
  - **Dark Quest** simulation suite with **Latin Hypercube Design**
  - Nonparametric regression based on **Gaussian Process**
  - currently **2~3% accuracy** and hopefully this gets better
  - **“Dark emulator” public after the HSC g-g lensing analysis**
- HSC (g-g lensing) + BOSS (w\_p) analysis
  - **Break the degeneracy!**
  - **“Galaxy parameters” on small scales (g-g lensing)**
  - **Cosmology on large scale (g-g clustering)**
  - **Mock challenge shows good results**
- To follow
  - Redshift-space distortions
  - Expand the input parameter space