# Machine learning to improve galaxy cluster mass estimation

(Jay) Digvijay Wadekar
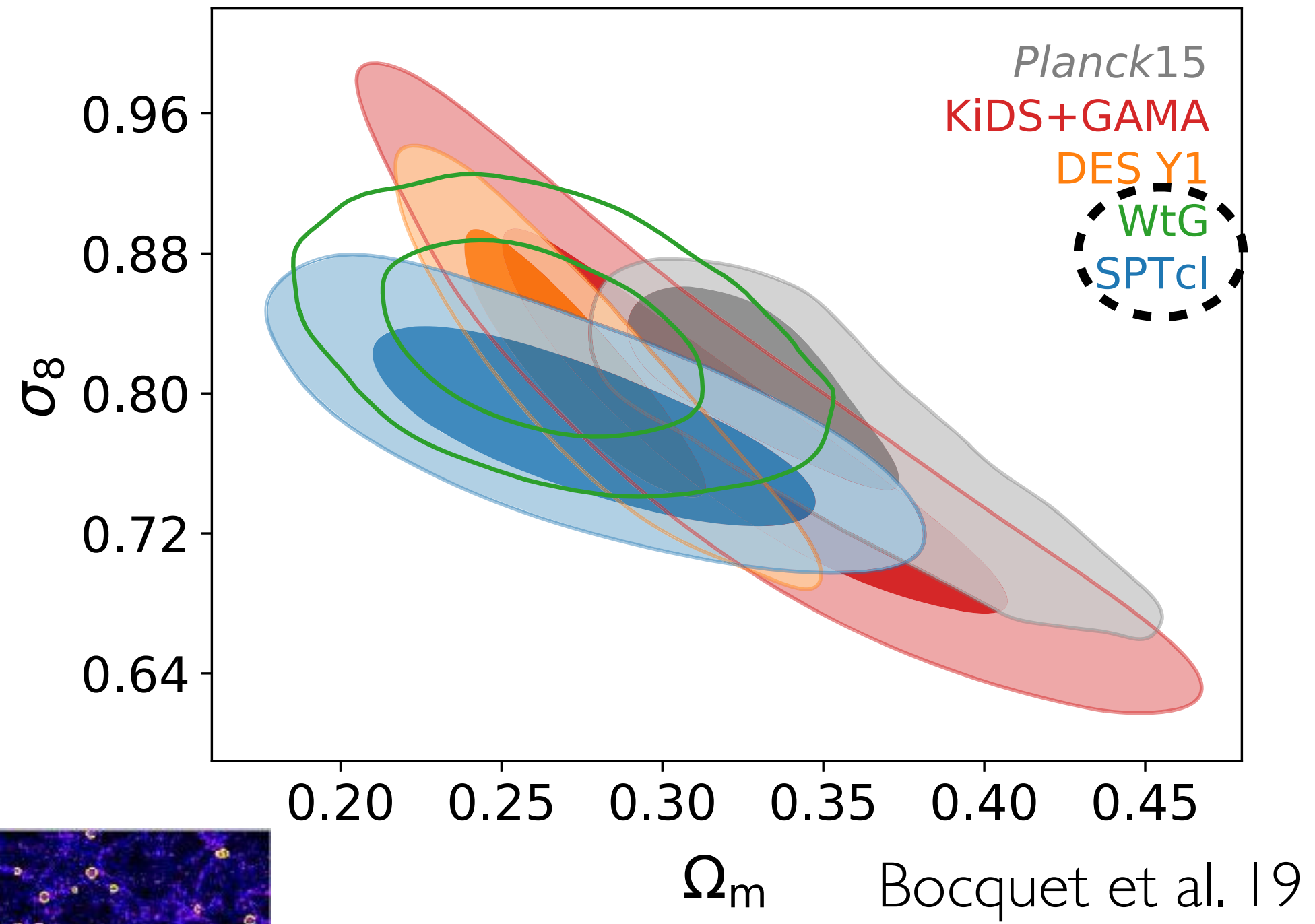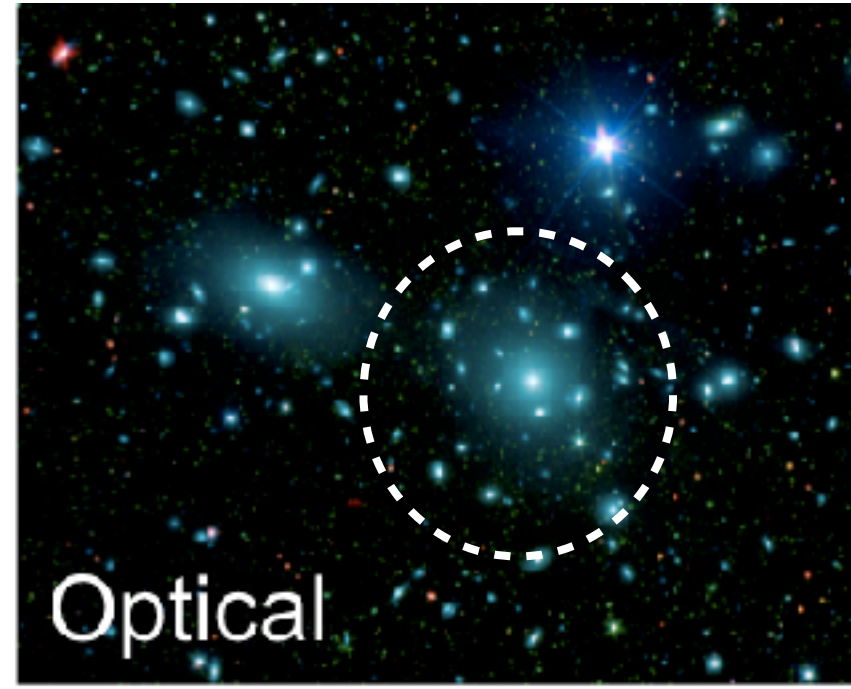
IAS

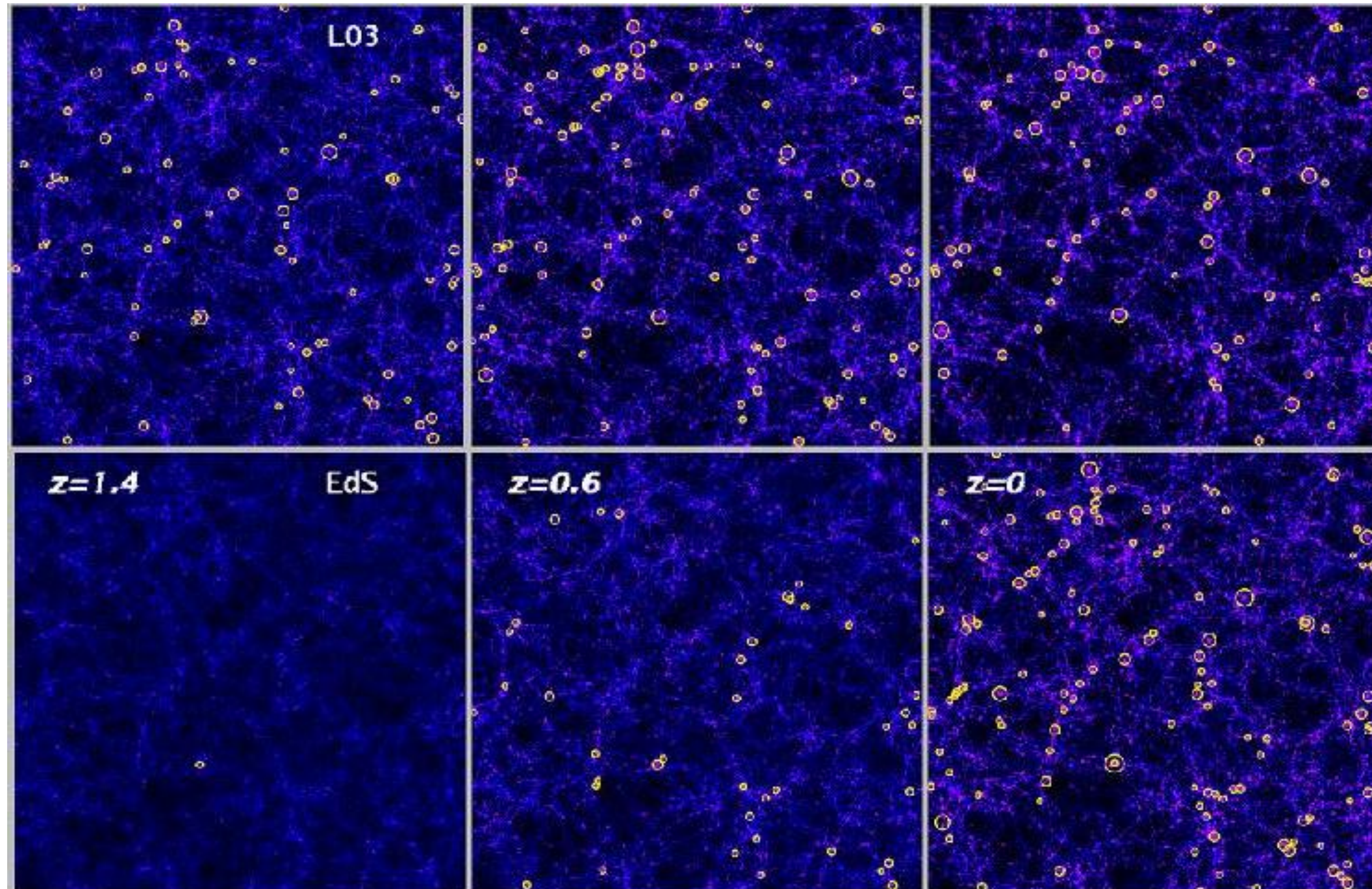with

*L. Thiele*,  F. Villaescusa-Navarro, C. Hill, D. Spergel, *M. Cranmer*,
N. Battaglia, S. Ho, D. Angles-Alcazar, L. Hernquist

# Cluster mass estimation is important for cosmology


Optical



Planck15
KiDS+GAMA
DES Y1
WtG
SPTcl

$\sigma_8$

$\Omega_m$    Bocquet et al. 19



ΛCDM
$\Omega_M = 0.3$
$\Lambda = 0.7$

SCDM
$\Omega_M = 1$

L03

z=1.4    EdS    z=0.6    z=0

normalized to present density

Borgani & Guzzo 01

al. (2005)
et al. (2005)
Allen et al. (2002)

2

"dunkle materie"

# Traditional approaches for cluster mass estimation



Optical | Radio | Microwave | X-ray | Lensing



(Velocity dispersion)

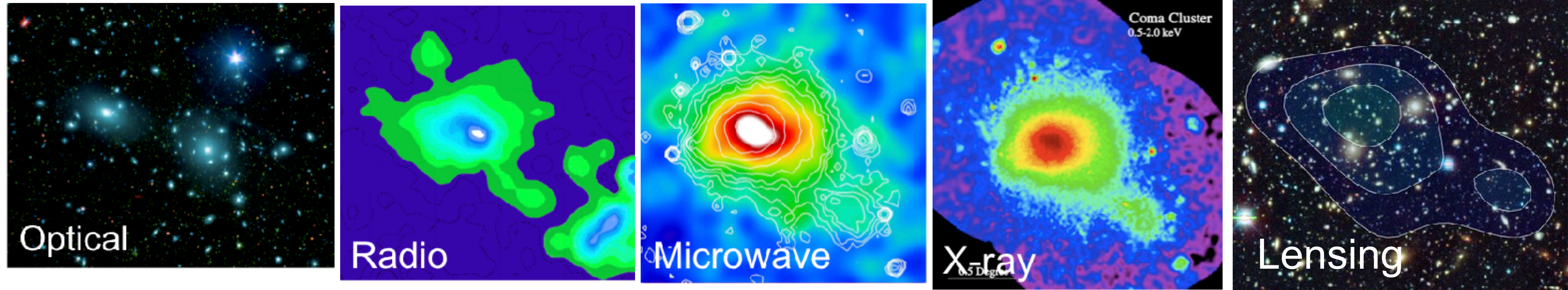$\log_{10}[\sigma_v \ (\mathrm{km \ s^{-1}})]$ vs $\log_{10}[M_{200c} \ (h^{-1}\mathrm{M_\odot})]$

- data
- fit
- $\pm 1\sigma$



$M_{200c} \ (h^{-1}M_\odot)$ vs $Y_{200c}$ (scaled)

- IllustrisTNG
- Scatter mean
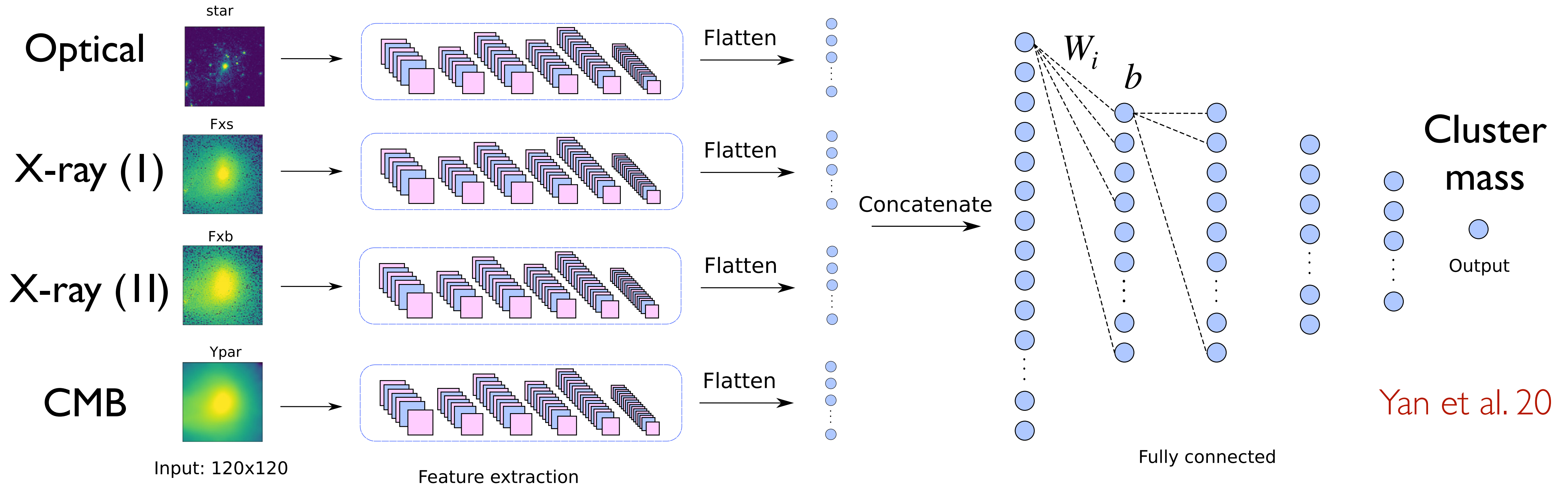- $M \propto Y^{3/5}$ (virial theorem)

$$Y \propto \int_0^{R_{200c}} P_e(r)dV$$

$$\sim M_{\mathrm{gas}}T_{\mathrm{gas}}$$

(thermal energy of gas)

# Machine learning (ML) is a potential alternative

Optical

X-ray (I)

X-ray (II)

CMB

star

Fxs

Fxb

Ypar

Input: 120x120

Feature extraction

Flatten

Flatten

Flatten

Flatten

Concatenate

$W_i$

$b$

Fully connected

Cluster mass

Output

Yan et al. 20
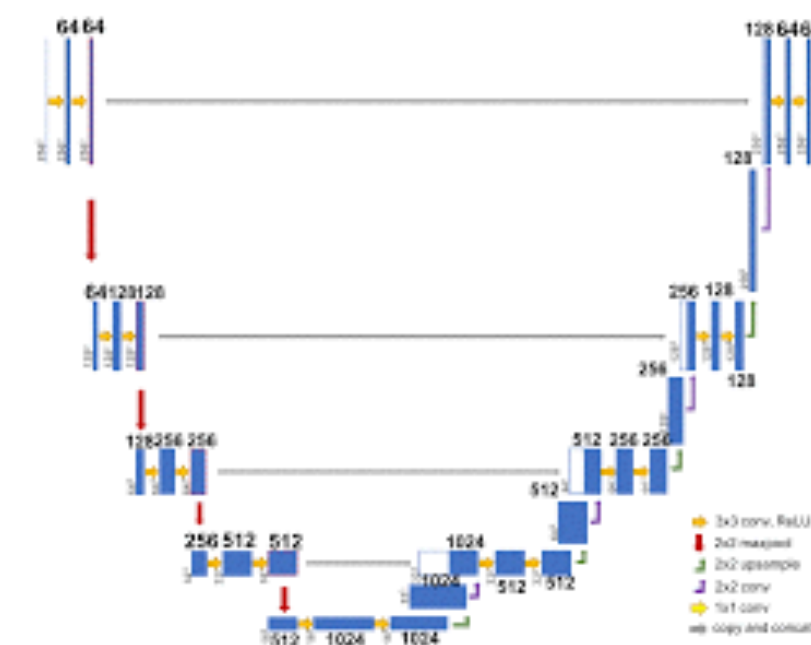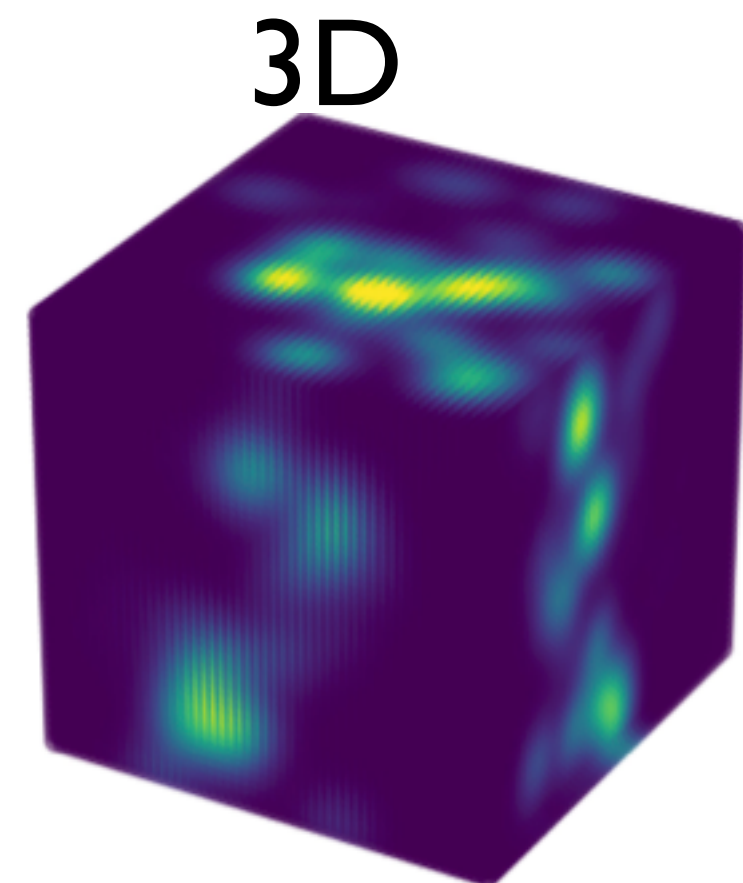
$$x^{(2)} = \text{ReLU}(\sum W_i x_i^{(1)} + b^{(2)})$$

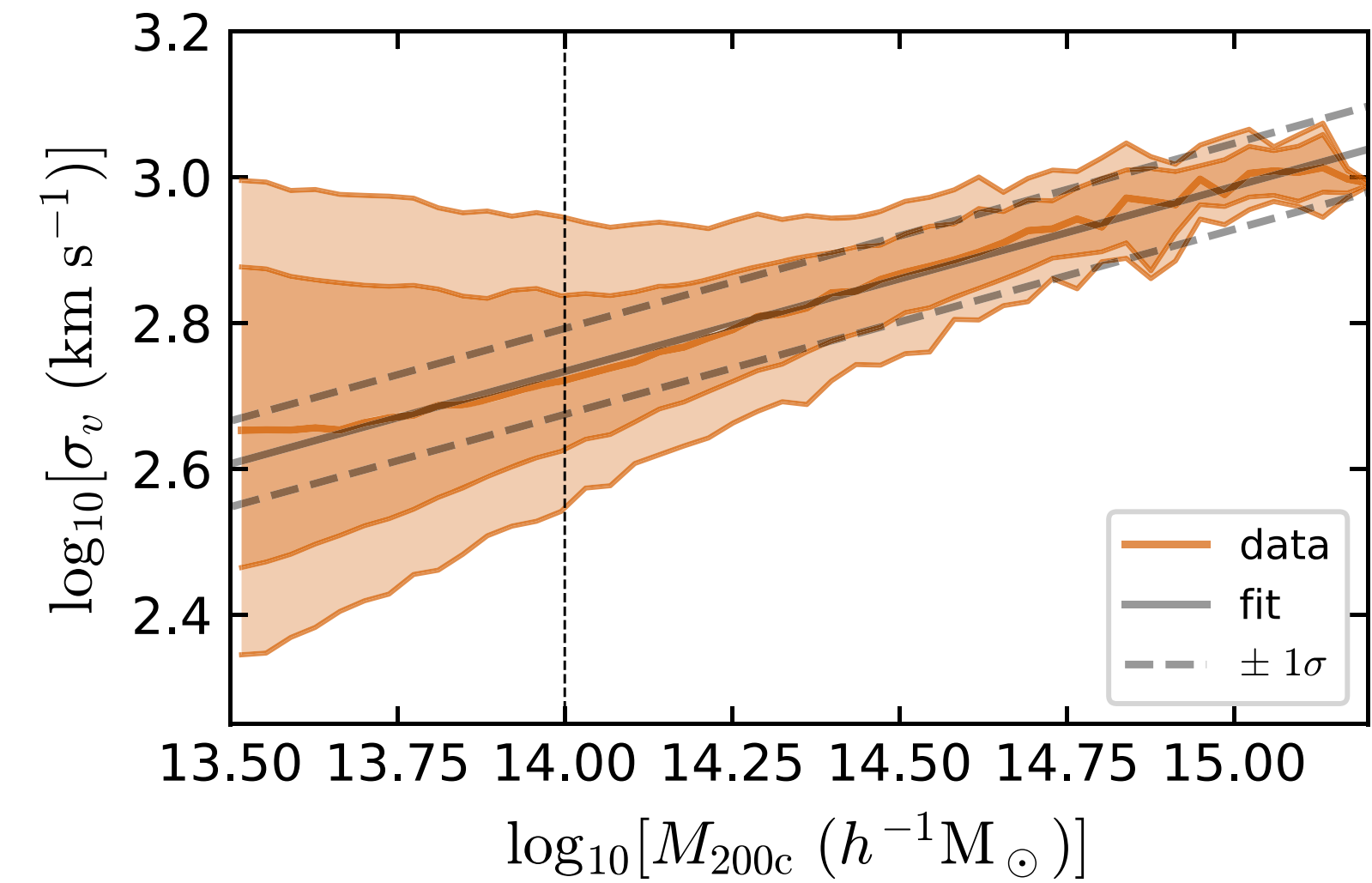# Machine learning (ML) is a potential alternative



- Utilize full dataset instead of just the first order moment ($\sigma_{\text{velocity}}$)
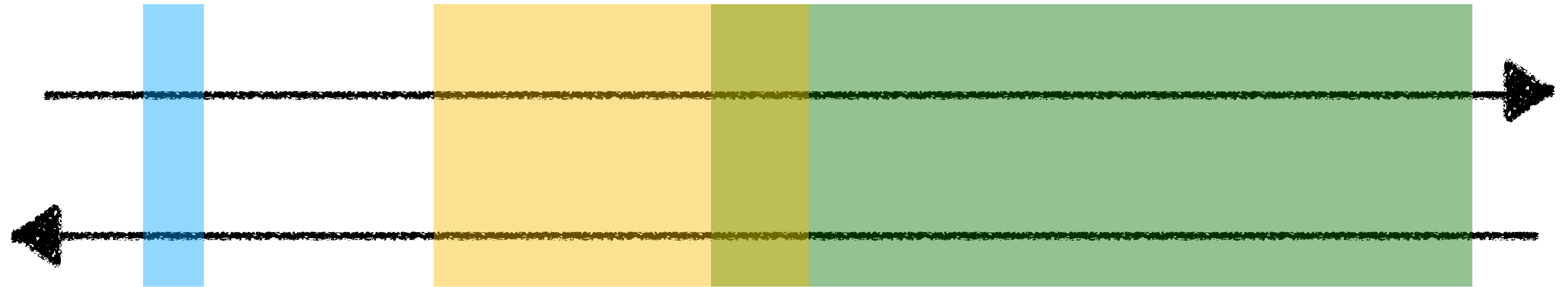
ID

3D

or

Cluster mass

- Ntampaka et al. 2015 (SDM)
- Ho et al. 2020, 2021 (CNN, Bayesian NN)
- Ramanah et al. 2020 (Neural flows)

# Comparison of ML approaches

**Power**
(input dimensionality/
data size)

**Generalizability/
Interpretability**

**Symbolic regression**

**Decision-tree approaches**
(e.g., random forests)

**Deep neural networks**

~10 parameters
~10,000 data points

# Comparison of ML approaches

Power
(input dimensionality/
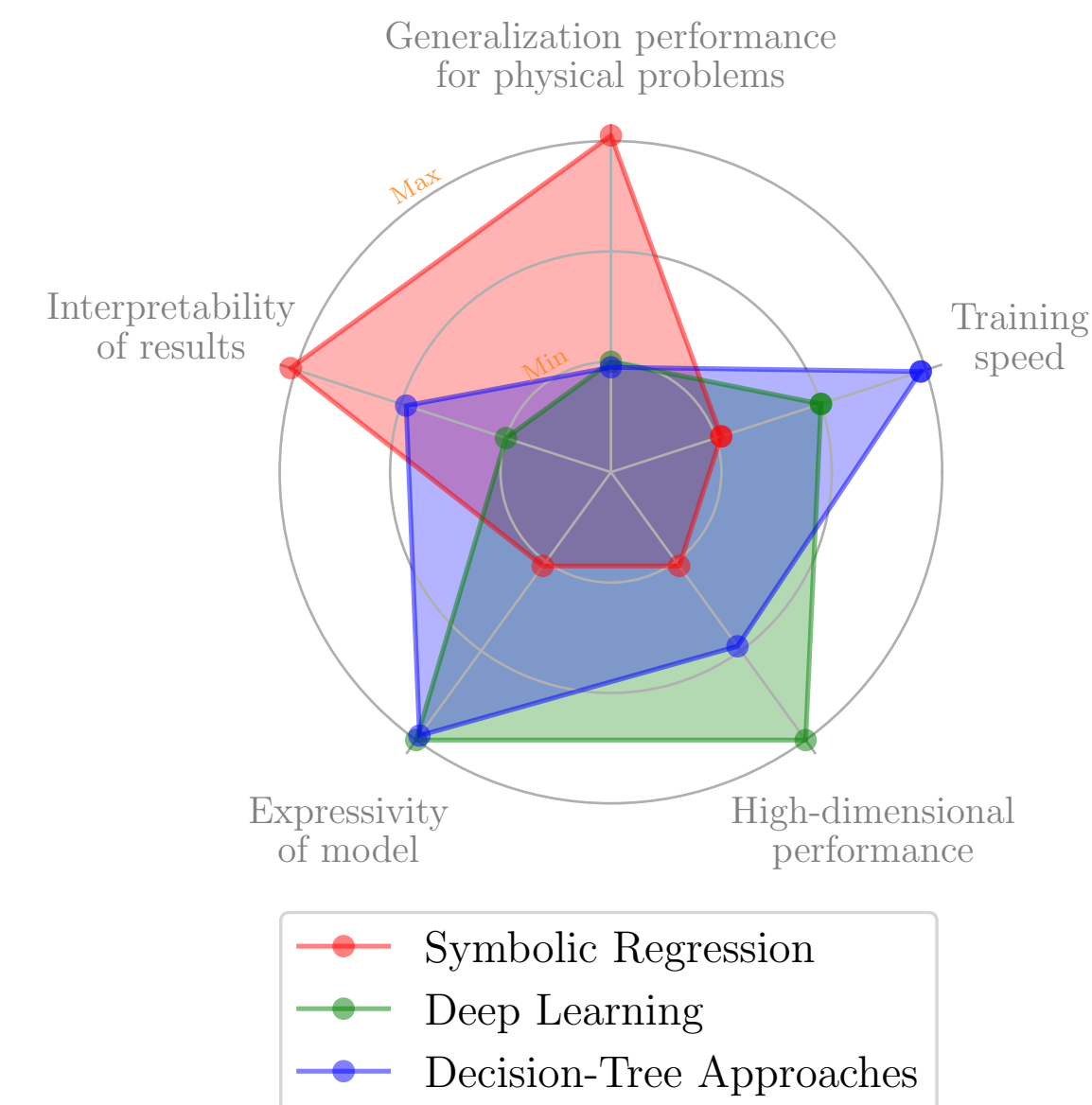data size)

Generalizability/
Interpretability

Symbolic
regression

Decision-tree
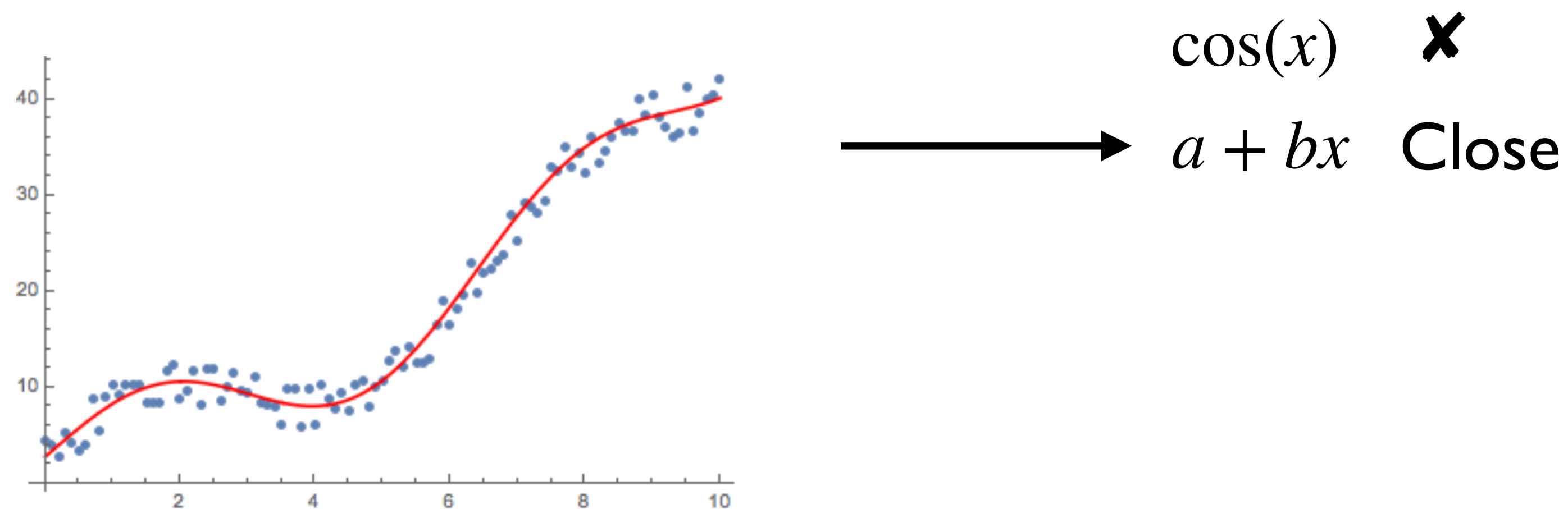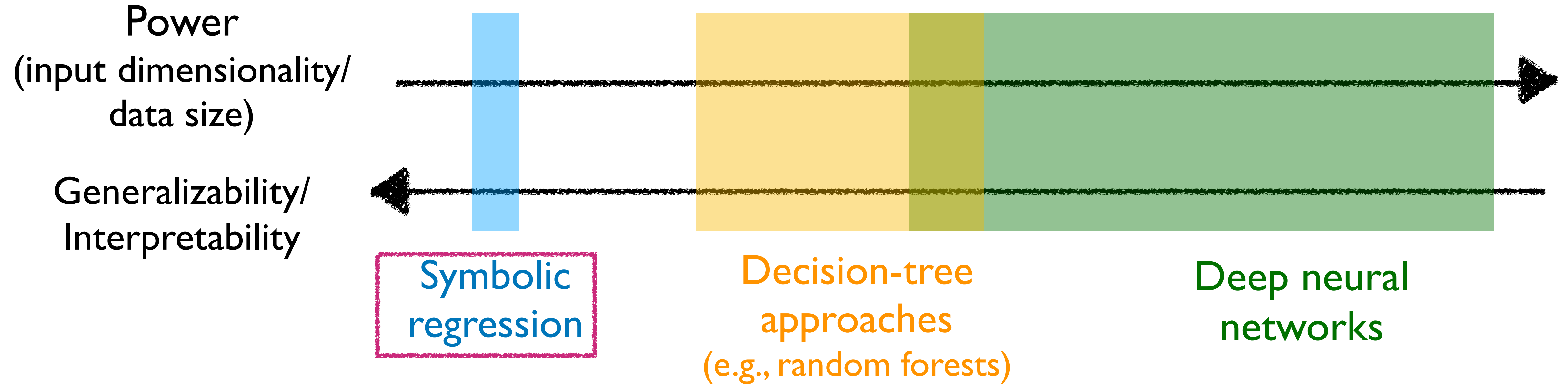approaches
(e.g., random forests)

Deep neural
networks

$\cos(x)$ ✘

$a + bx$ Close

# Comparison of ML approaches

Power
(input dimensionality/
data size)

Generalizability/
Interpretability

**Symbolic regression**

**Decision-tree approaches**
(e.g., random forests)

**Deep neural networks**

$\cos(x)$ ✘

$a + bx$ (Close)

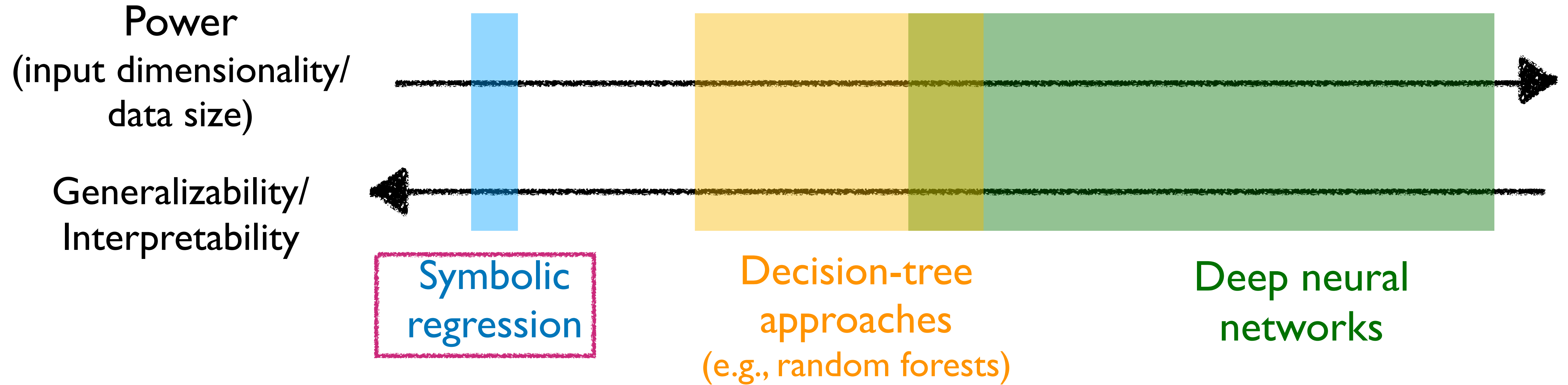$a + bx + d\sin(x)$ (Closer!)

$a + bx + cx^2 + d\sin(x)$ ✔

*PySR package:*
*https://github.com/MilesCranmer/PySR*

# Our approach:
## Symbolic regression + Random Forest

$$M_{\text{cluster}} = f \left( Y_{\text{CMB}}^{3/5} \; , \; \textcolor{red}{\text{observables from other surveys?}} \right)$$
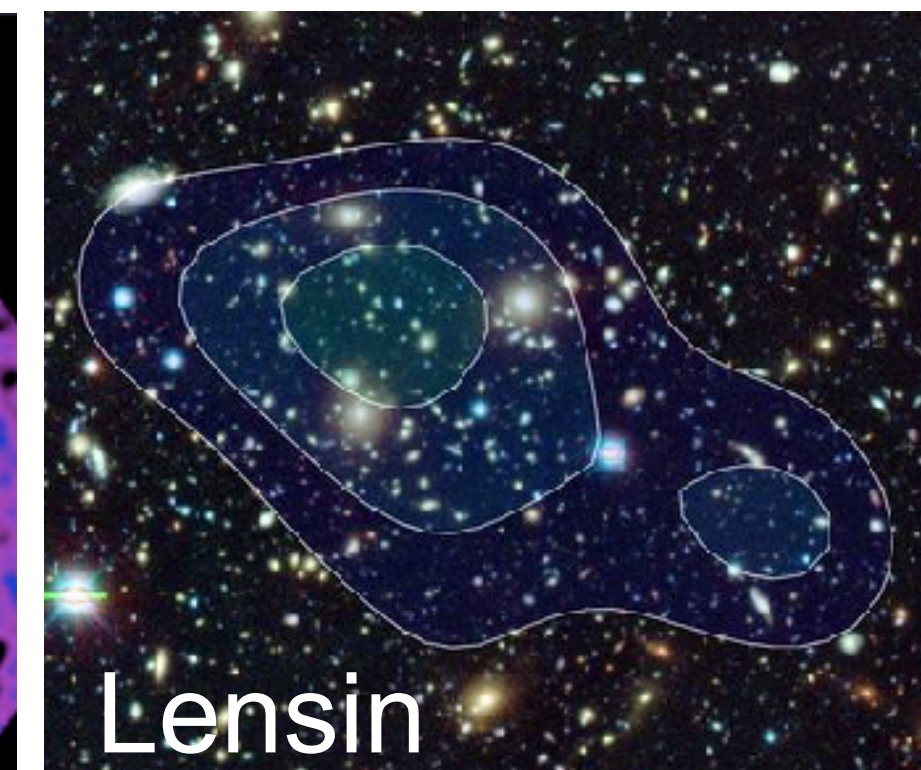


- **X-ray surveys**
  - Gas mass profile
  - Luminosity profile
  - Spectral temperature
  - Gas ellipticity
  - ……

- **Galaxy surveys**
  - Richness
  - Galaxy colors
    (e.g. fraction of red galaxies)
  - Stellar mass
  - ……



Optical   Radio   Microwave   X-ray   Lensin

# Our approach:
# Symbolic regression + Random Forest

$$M_{\text{cluster}} = f\,(\ Y_{\text{CMB}}^{3/5}\ , \text{observables from other surveys? }\,)$$



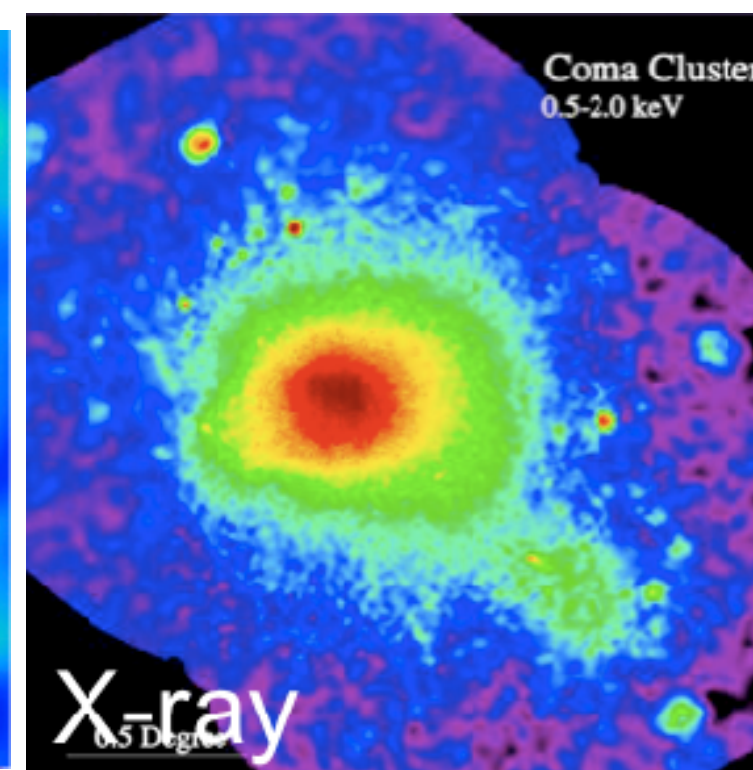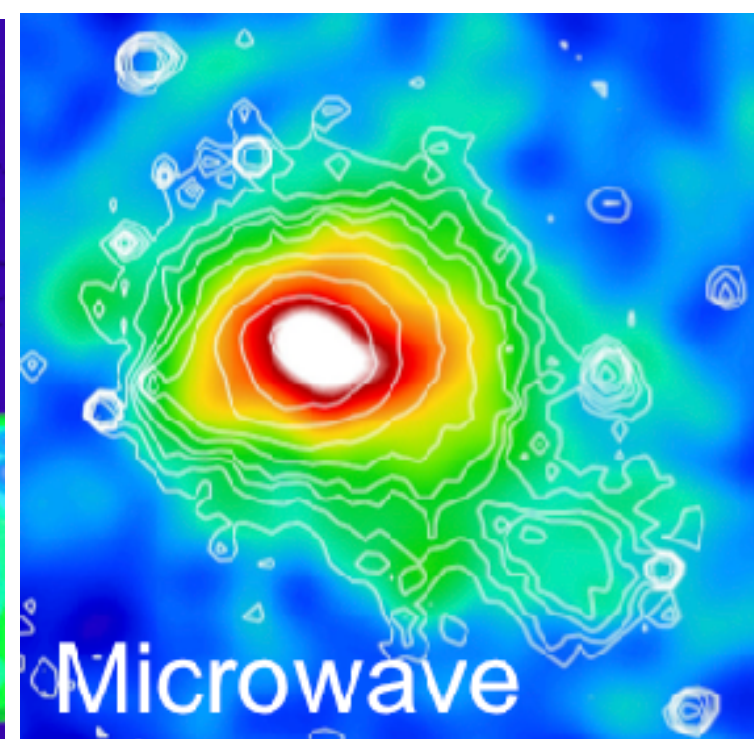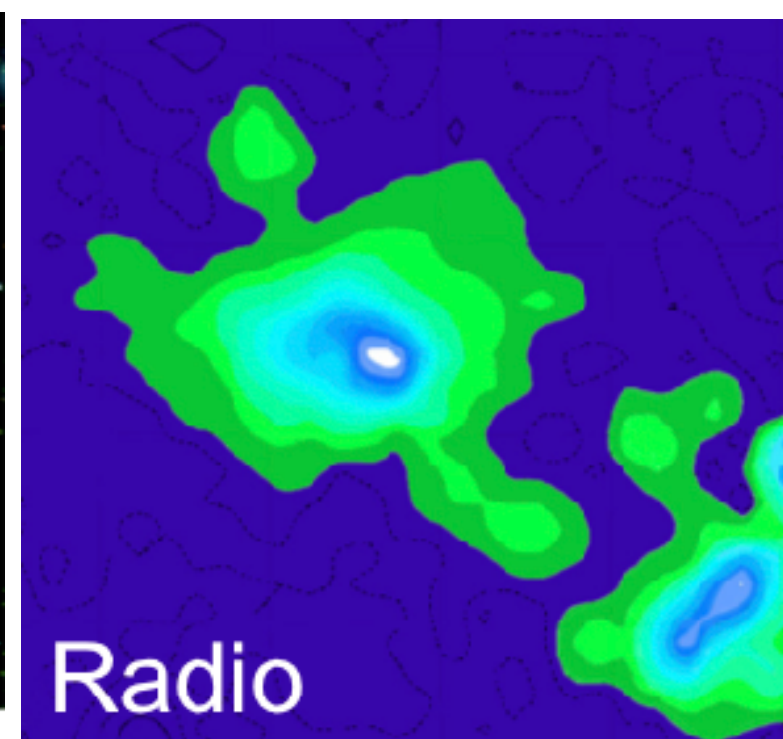- **X-ray surveys**

  - Gas mass profile
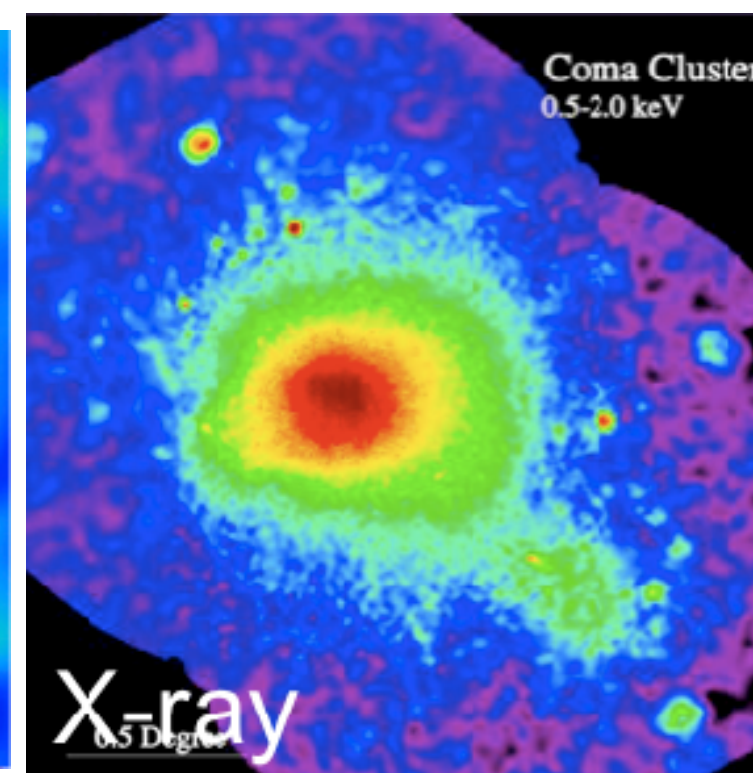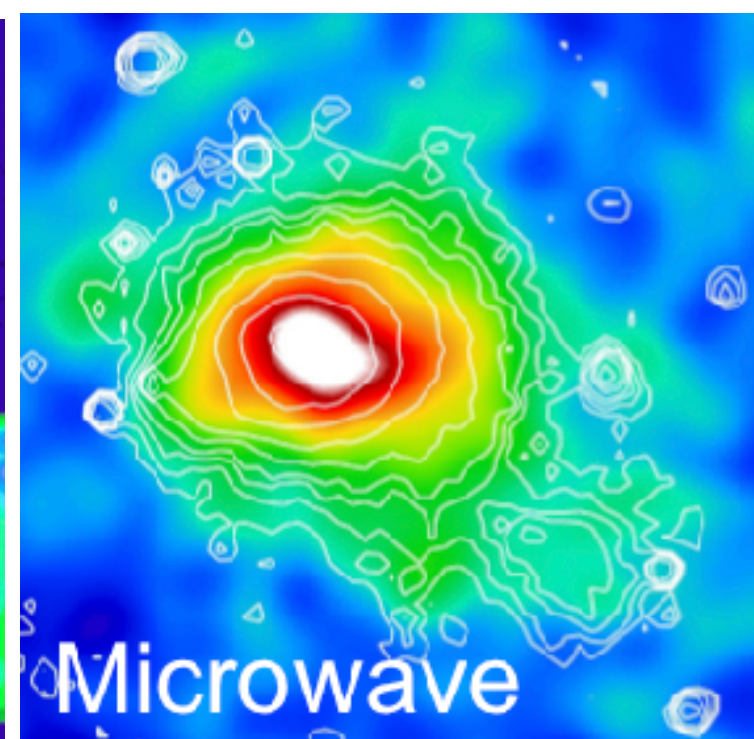  - Luminosity profile
  - Spectral temperature
  - Gas ellipticity
  - ……

- **Galaxy surveys**

  - Richness
  - Galaxy colors
    (e.g. fraction of red galaxies)
  - Stellar mass
  - ……



Optical · Radio · Microwave · X-ray · Lensin

# Results for IllustrisTNG

$$M_{\rm pred}^{(1)} \propto Y^{3/5}$$

$$M_{\rm pred}^{(2)} \propto Y^{3/5} \, {\color{red}(1 - A \, c_{\rm gas})} \qquad c_{\rm gas} \equiv \frac{M_{\rm gas}(r < R_{200c}/2)}{M_{\rm gas}(r < R_{200c})}$$

$$M_{\rm pred}^{(3)} \propto Y^{3/5} {\color{red}\left(\frac{B}{c_{\rm NFW}}\right)^{M_*/M_{\rm gas}}}$$

11

# Results for IllustrisTNG



$$M_{\text{pred}}^{(1)} \propto Y^{3/5}$$

$$M_{\text{pred}}^{(2)} \propto Y^{3/5} \, (1 - A \, c_{\text{gas}})$$

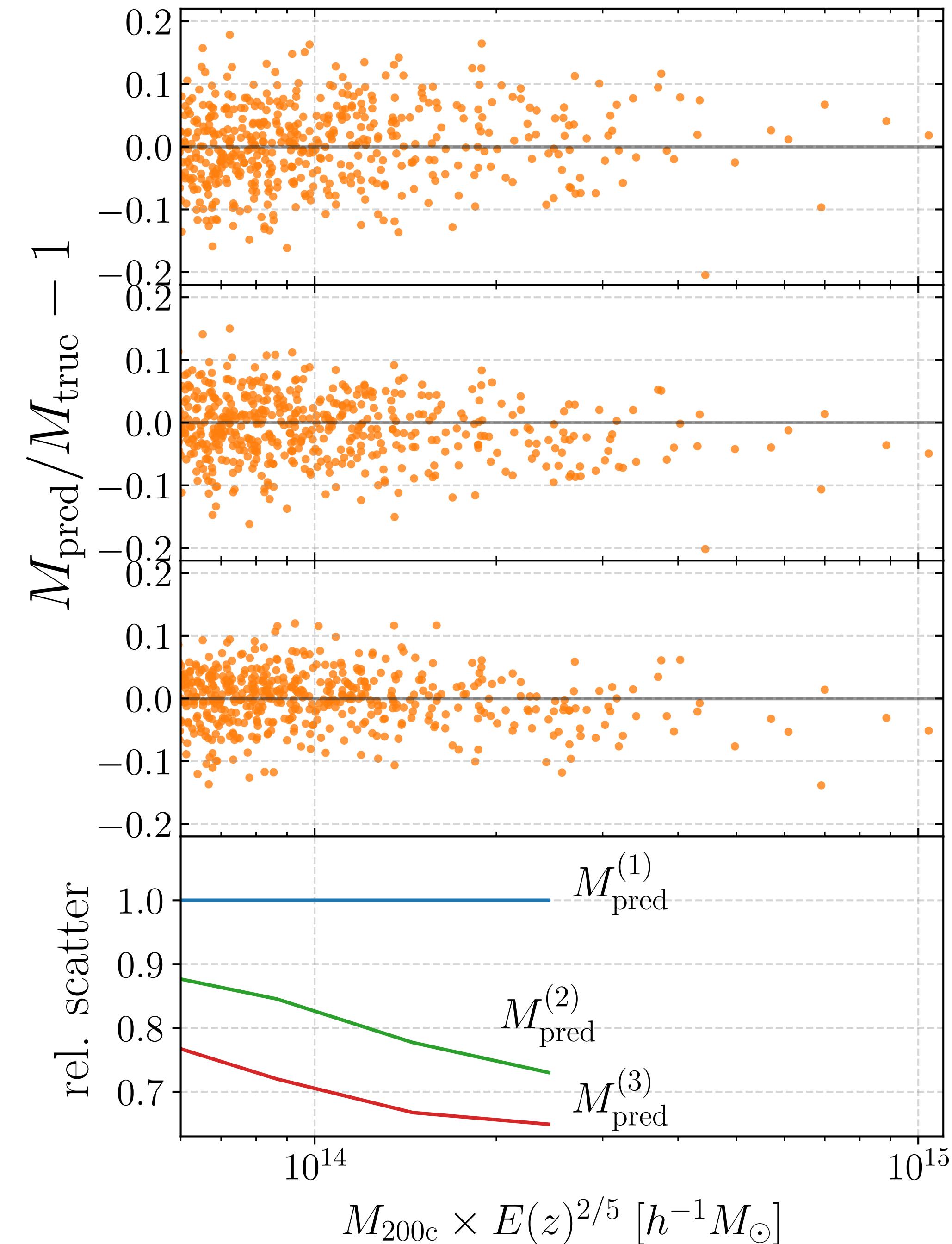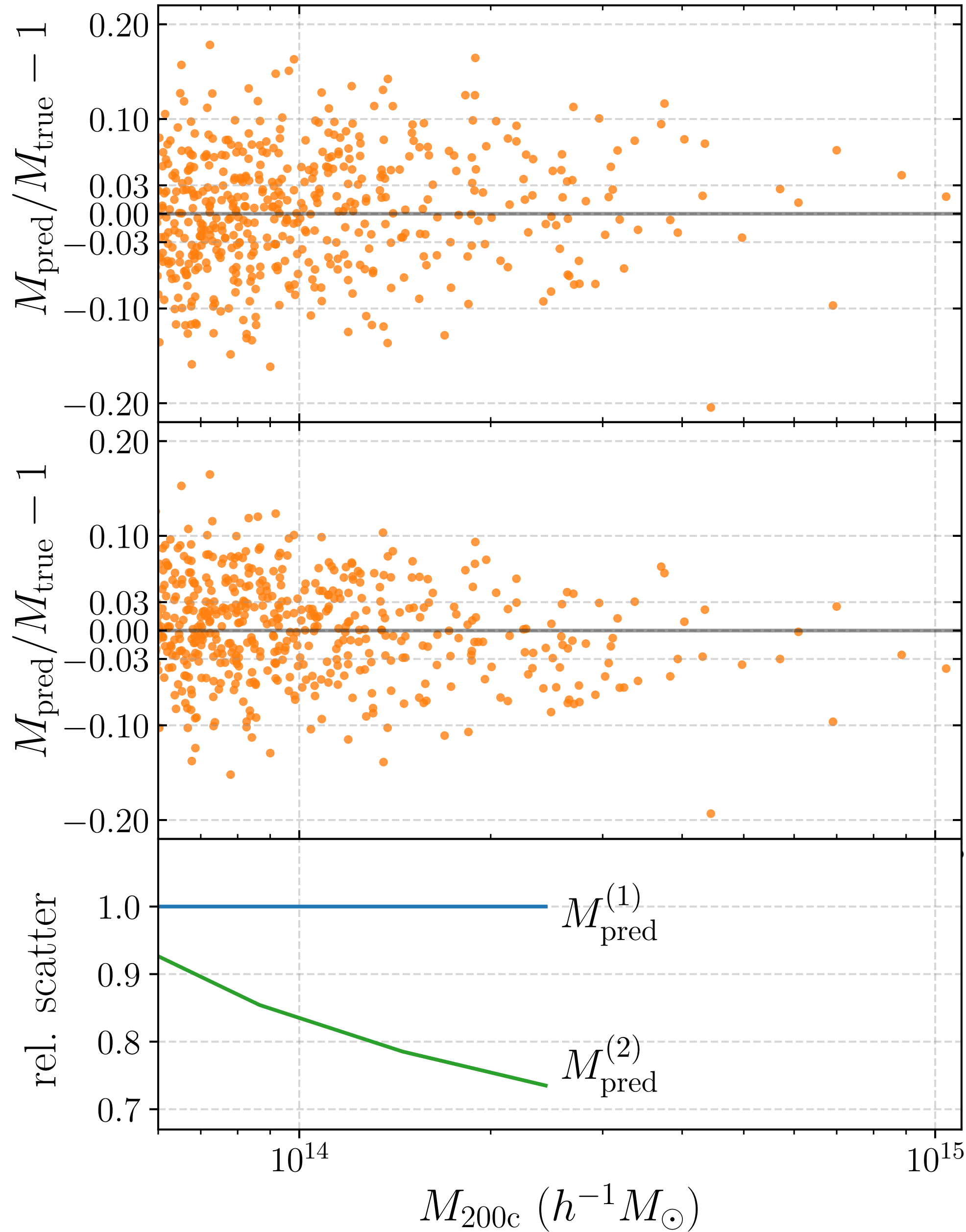$$c_{\text{gas}} \equiv \frac{M_{\text{gas}}(r < R_{200c}/2)}{M_{\text{gas}}(r < R_{200c})}$$

$$M_{\text{pred}}^{(3)} \propto Y^{3/5} \left( \frac{B}{c_{\text{NFW}}} \right)^{M_*/M_{\text{gas}}}$$

## Reasons for dependence:

1. Central regions of clusters are noisier
   (conc. can be used to down-weight central regions)

2. Conversion of gas to stars reduces Y

Kravtsov et al. 06,
Arnaud et al. 10

12

# Results for IllustrisTNG



$$M_{\mathrm{pred}}^{(1)} \propto Y^{3/5}$$

$$M_{\mathrm{pred}}^{(2)} \propto Y^{3/5}\,(1 - A\,c_{\mathrm{gas}}) \qquad c_{\mathrm{gas}} \equiv \frac{M_{\mathrm{gas}}(r < R_{200c}/2)}{M_{\mathrm{gas}}(r < R_{200c})}$$

| X-ray | SZ |
|---|---|
| - High resolution | - Low resolution |
| - Indirect probe | - Direct probe |

13

# Cross-checks

## Excising inner cluster regions



$M_{\rm pred}^{(1)} \propto Y^{3/5}$

$M_{\rm pred}^{(2)} \propto [Y_{\rm 200c} - Y(r < R_{\rm 200c}/4)]^{3/5}$

$M_{\rm pred}^{(1)}$

$M_{\rm pred}^{(2)}$

## Radial dependence of scatter



$M_{200} \in [1-2] \times 10^{14}\ h^{-1} M_\odot$

But IllustrisTNG has only one configuration
of baryonic feedback and initial conditions?

Do the results hold in a more general setting?
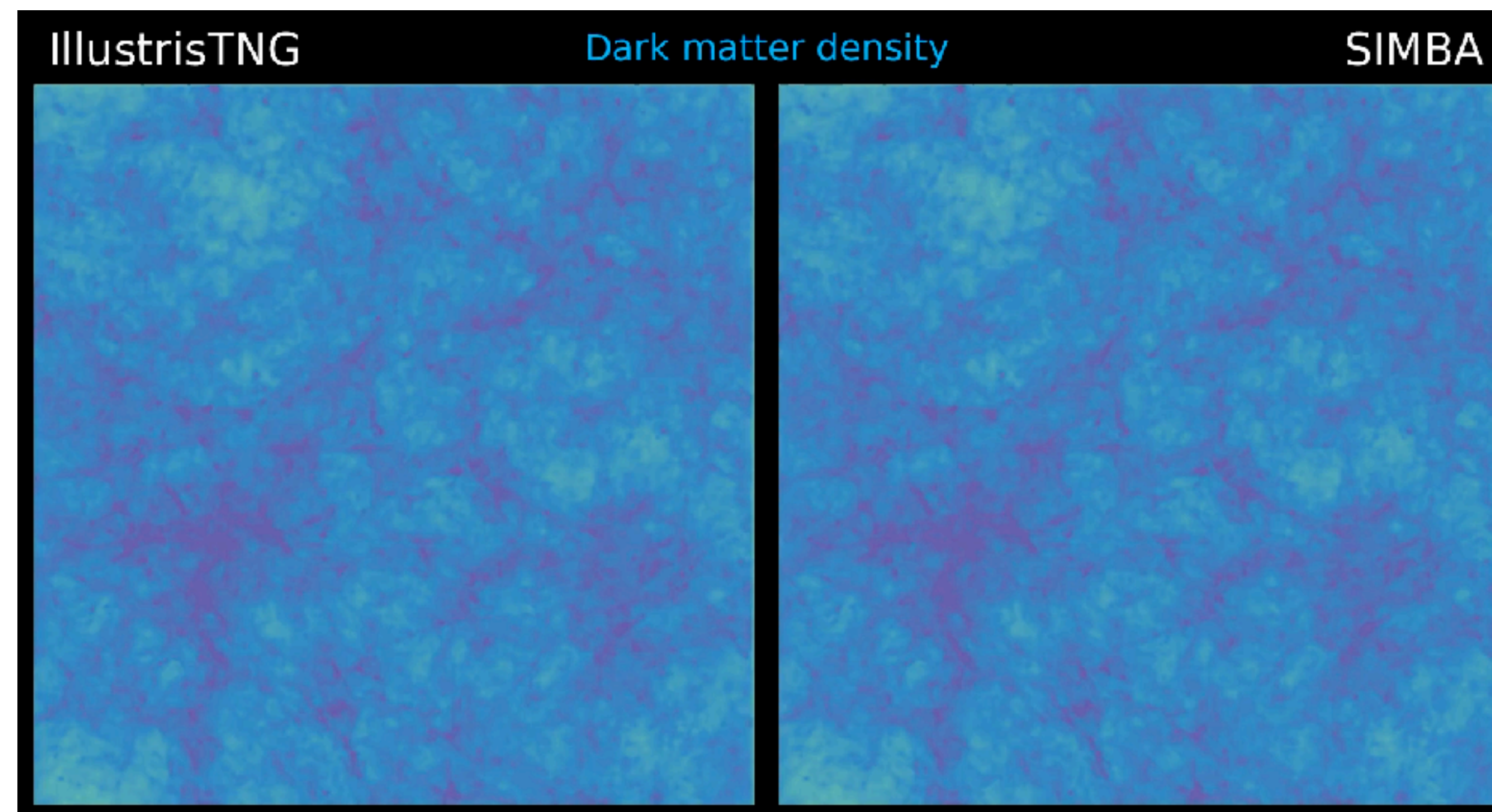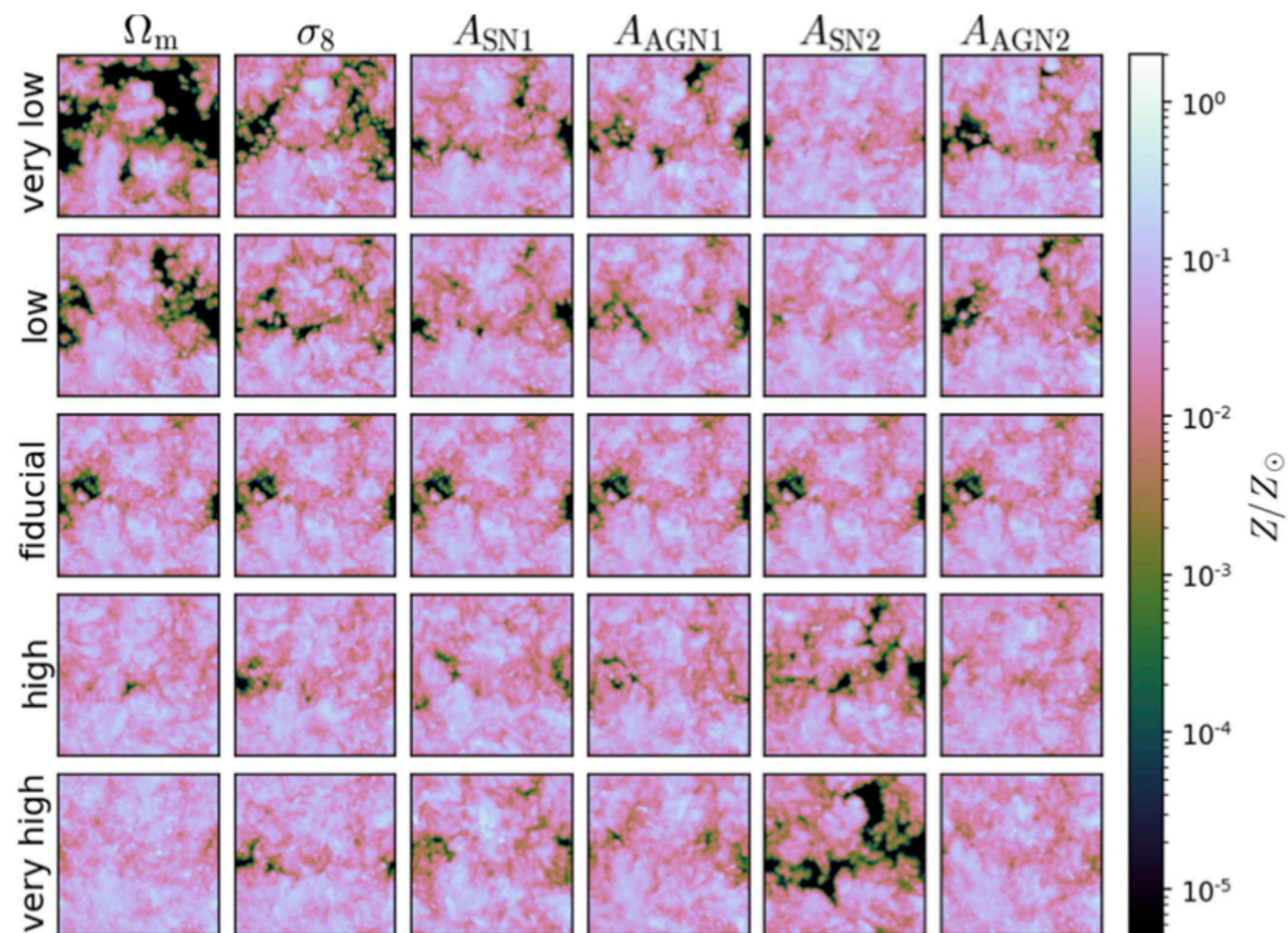
# But IllustrisTNG has only one configuration of baryonic feedback and initial conditions?

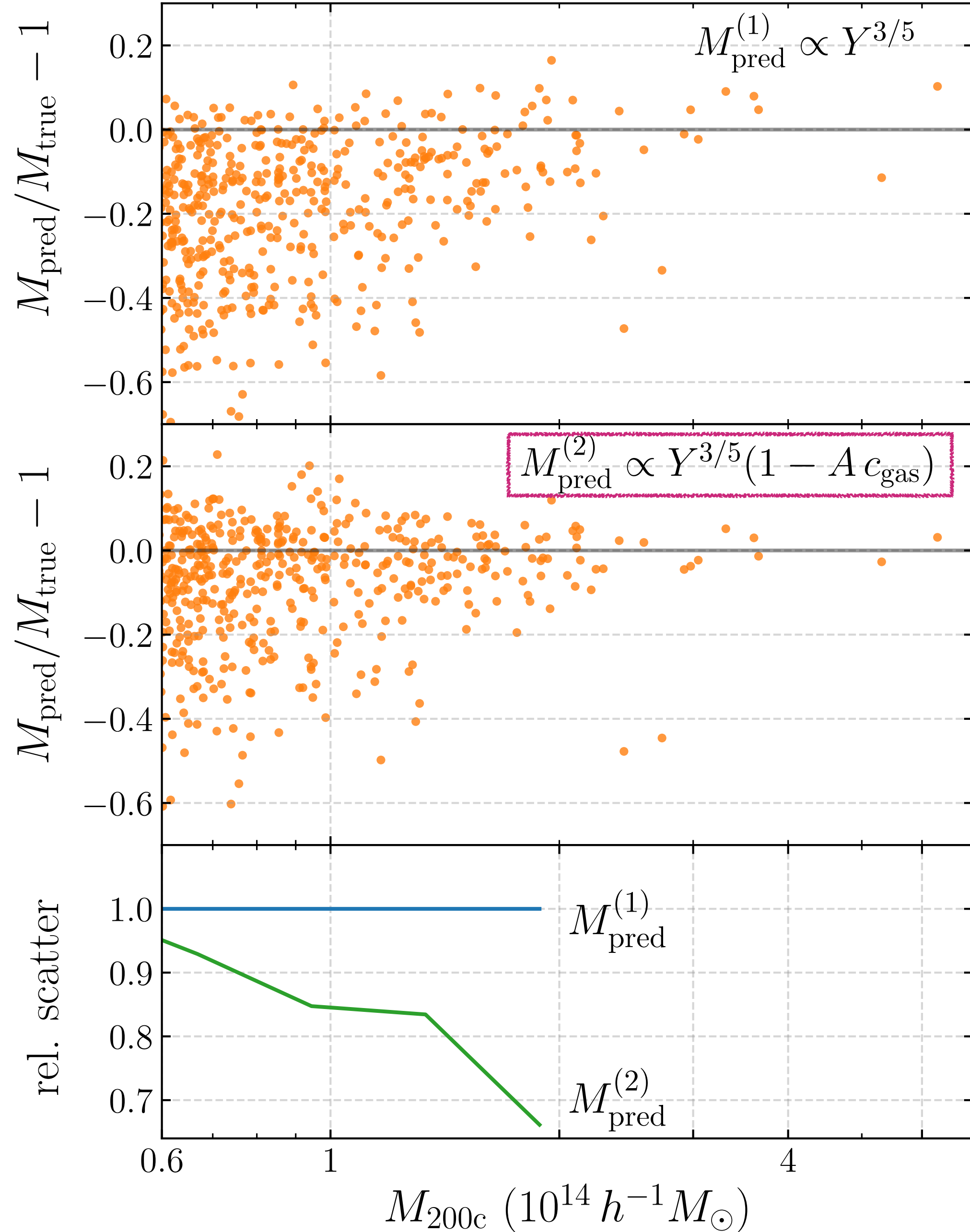## Do the results hold in a more general setting?



CAMELS simulations
Villaescusa-Navarro et al. 21
*https://camels.readthedocs.io/*

Results for
CAMELS

CAMELS - SIMBA

CAMELS - TNG

$M_{\mathrm{pred}}^{(1)} \propto Y^{3/5}$

$M_{\mathrm{pred}}^{(2)} \propto Y^{3/5}(1 - A\,c_{\mathrm{gas}})$

$M_{\mathrm{pred}}^{(2)} \propto Y^{3/5}(1 - B\,c_{\mathrm{gas}})$

$M_{\mathrm{pred}}/M_{\mathrm{true}} - 1$

rel. scatter

$M_{\mathrm{pred}}^{(1)}$

$M_{\mathrm{pred}}^{(2)}$

$M_{200\mathrm{c}}\ (10^{14}\,h^{-1}M_\odot)$

17
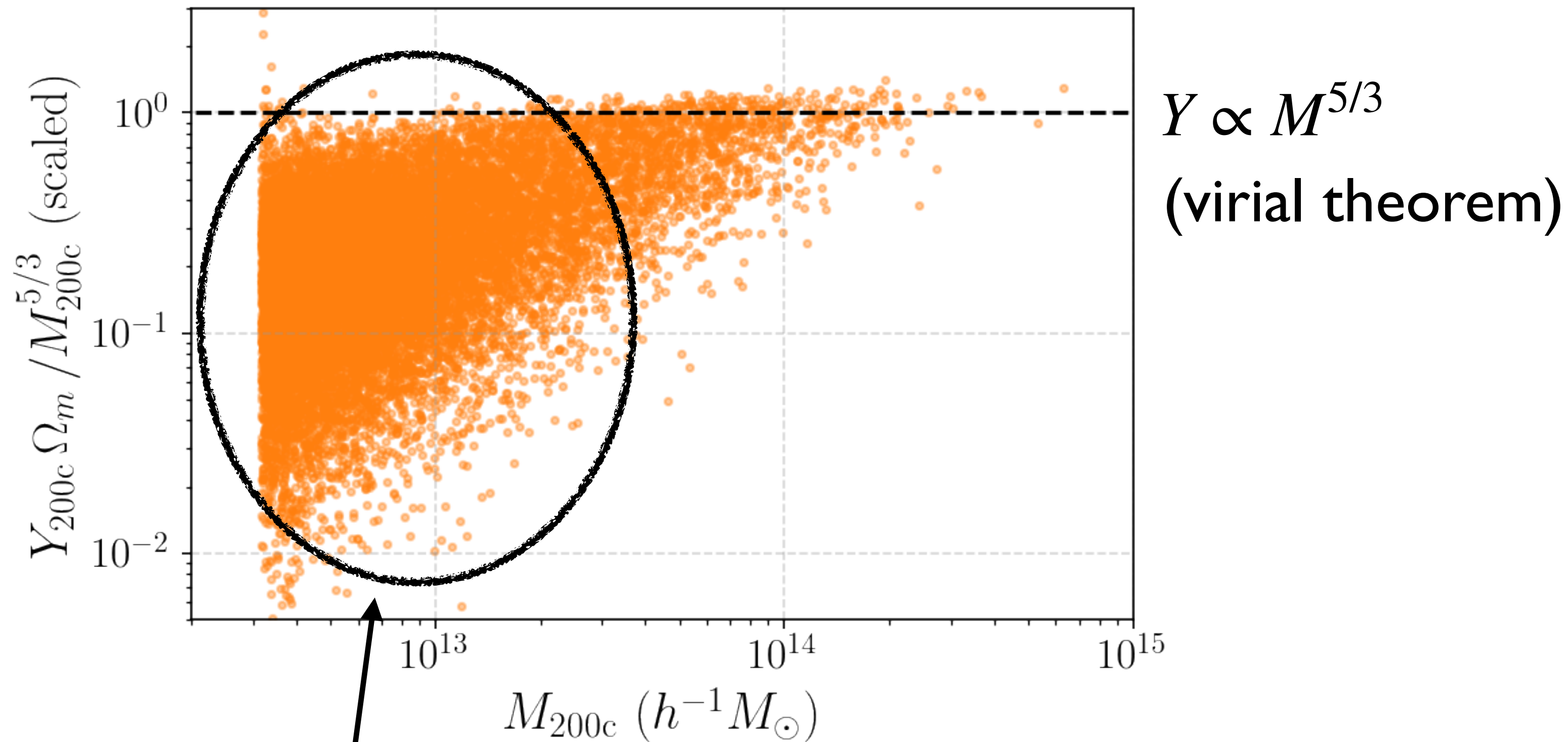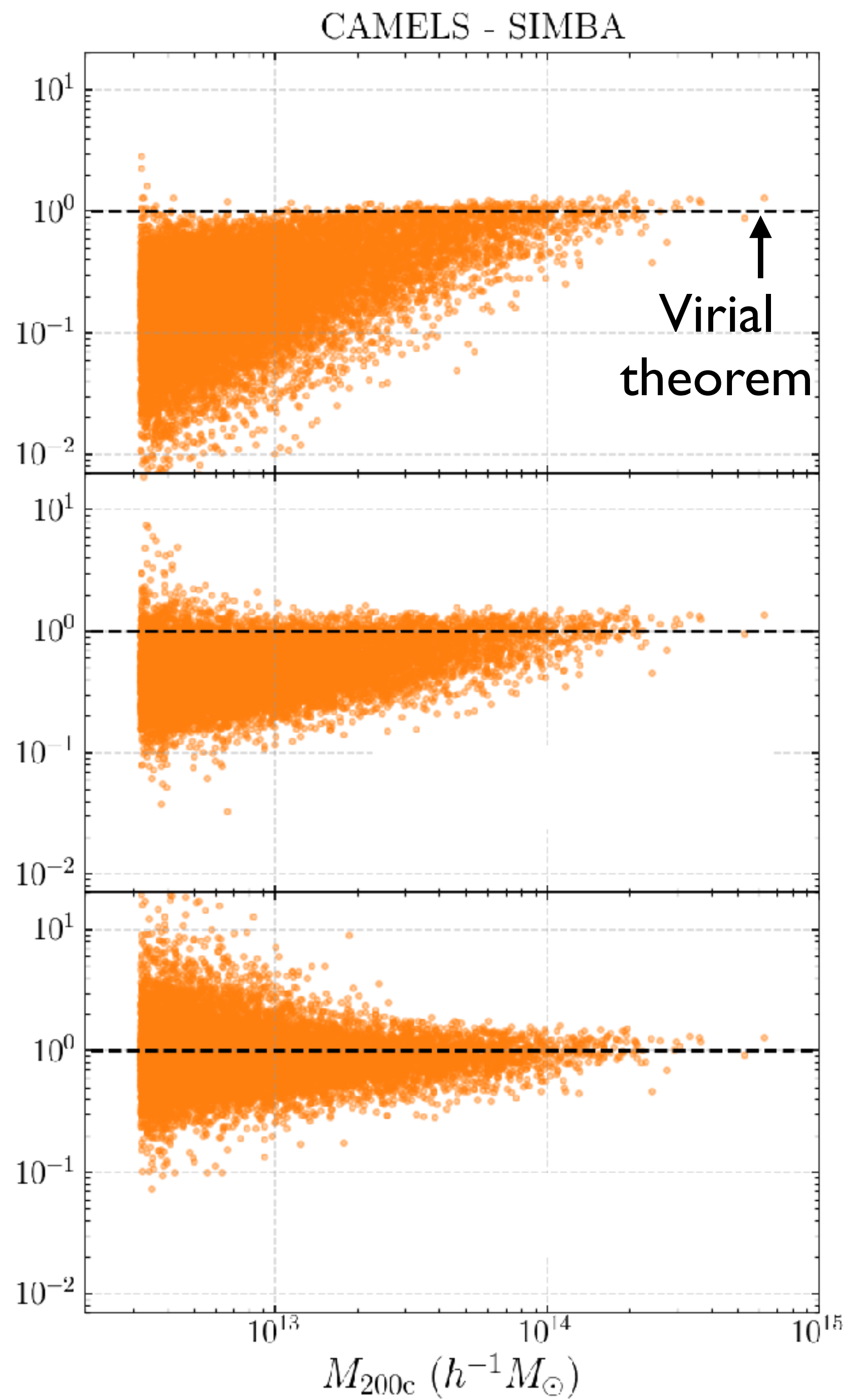
# Reducing deviation from self-similarity (pow. law)



CAMELS - SIMBA

$Y \propto M^{5/3}$

(virial theorem)

Due to ejection of gas from clusters from AGN/SN feedback

CAMELS - SIMBA

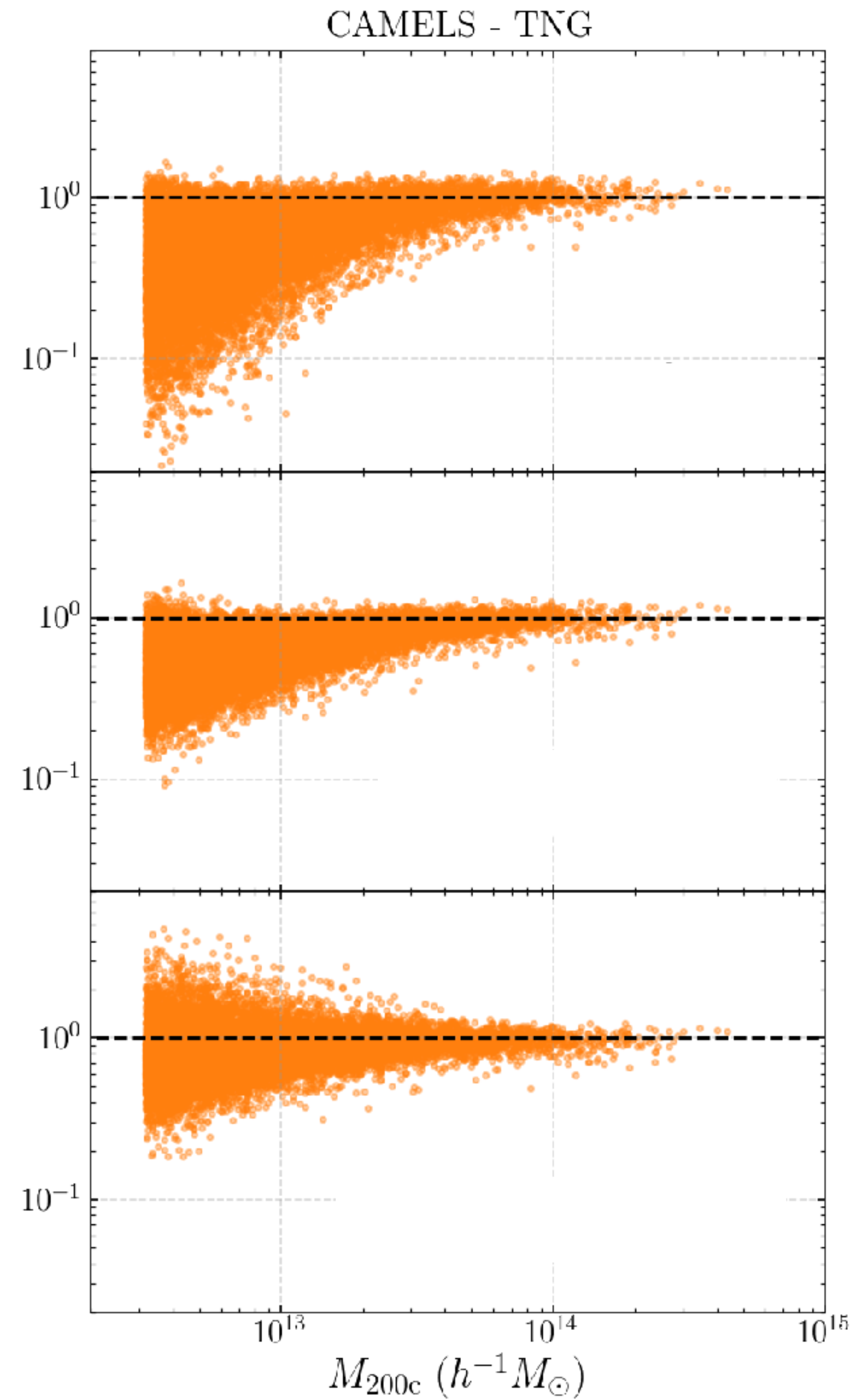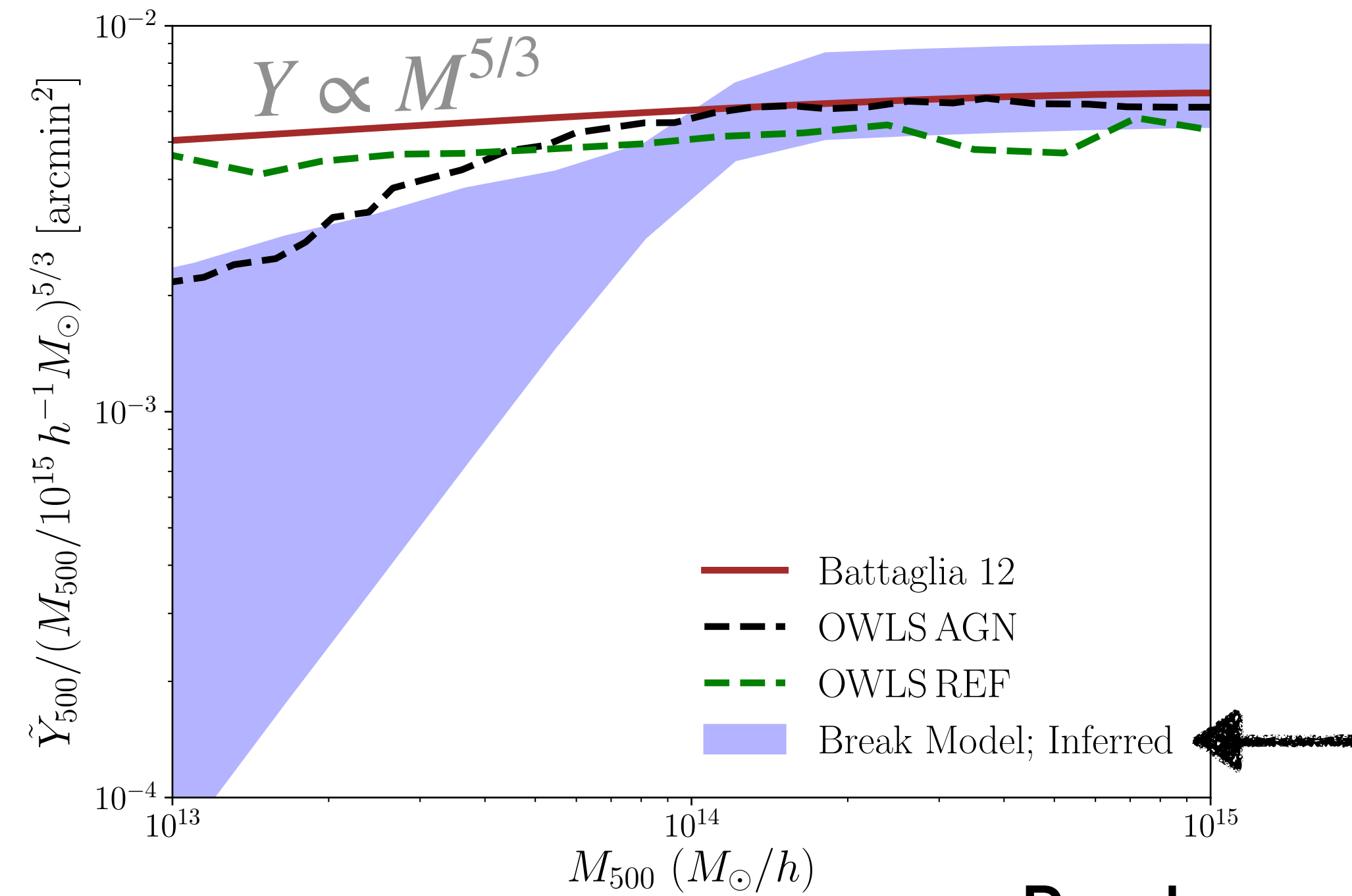CAMELS - TNG

Preliminary results

$Y$

Virial theorem

$Y \left( 1 + \dfrac{M_*(r < R)}{M_{\text{gas}}(r < R)} \right)$

$Y \left[ 1 + \dfrac{M_*(r < R/2)}{M_{\text{gas}}(r < R/2)} \right]$

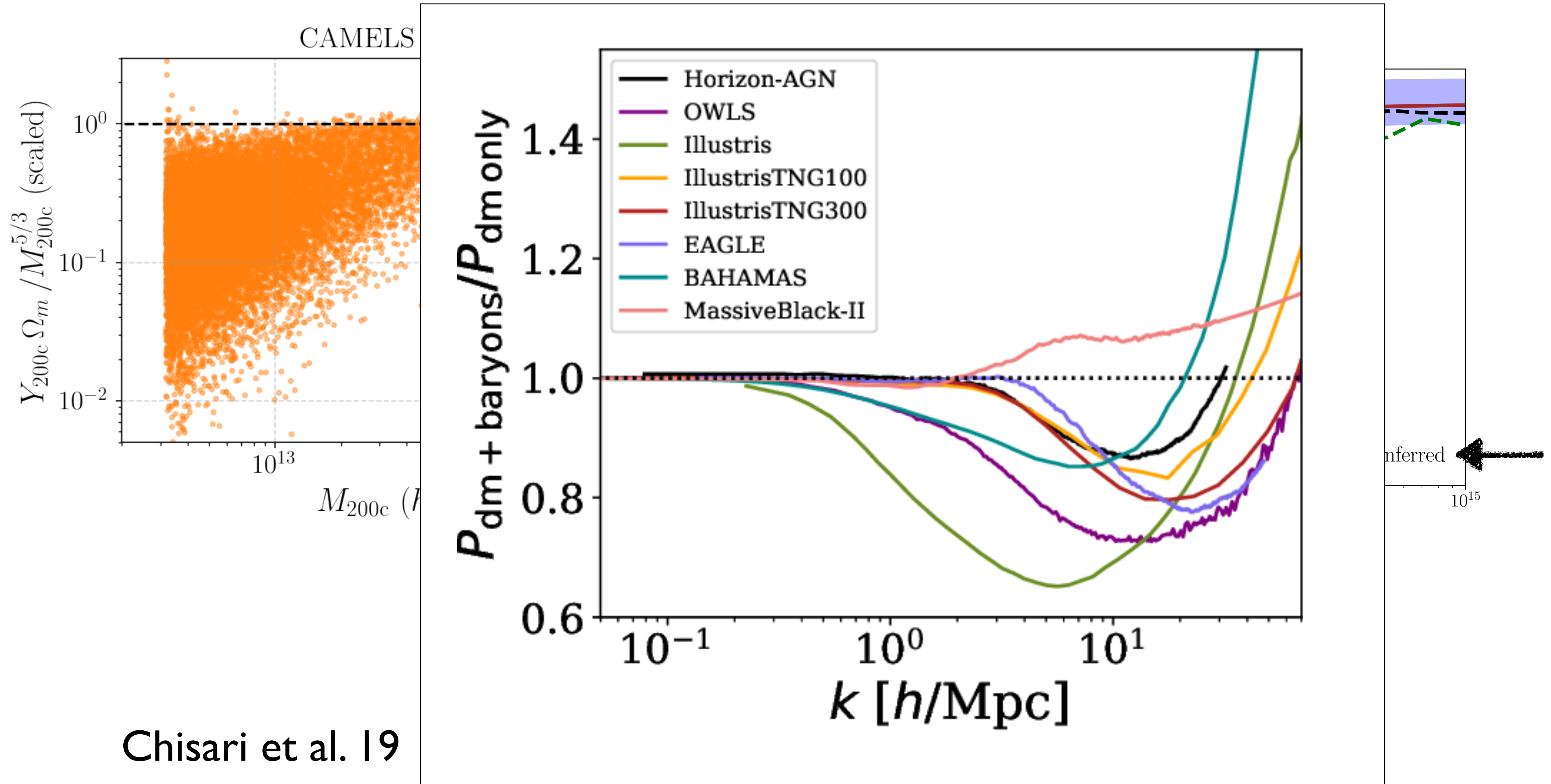$M_{200c} \ (h^{-1} M_\odot)$

$M_{200c} \ (h^{-1} M_\odot)$

# Using the Y-M measurements to constrain baryonic feedback



CAMELS - SIMBA

$Y \propto M^{5/3}$

Battaglia 12
OWLS AGN
OWLS REF
Break Model; Inferred
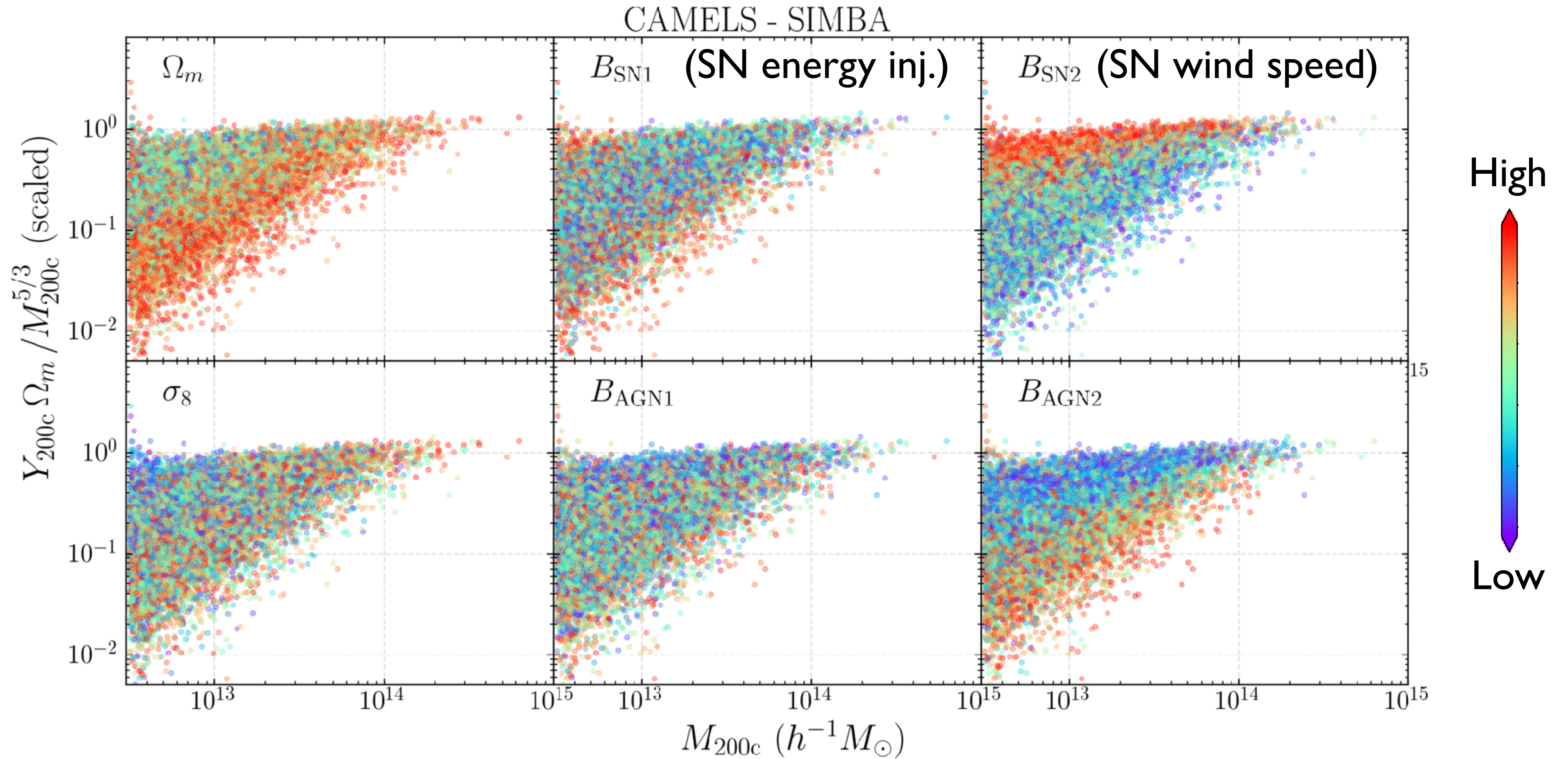
Pandey et al. 21
(ACT x DES)

# Using the Y-M measurements to constrain baryonic feedback



CAMELS

$Y_{200c}\,\Omega_m\,/M_{200c}^{5/3}$ (scaled)

$M_{200c}$

Chisari et al. 19

**Legend:**
- Horizon-AGN
- OWLS
- Illustris
- IllustrisTNG100
- IllustrisTNG300
- EAGLE
- BAHAMAS
- MassiveBlack-II

$P_{\text{dm + baryons}}/P_{\text{dm only}}$

$k\ [h/\text{Mpc}]$

inferred

CAMELS - SIMBA

$\Omega_m$

$B_{\mathrm{SN1}}$ (SN energy inj.)

$B_{\mathrm{SN2}}$ (SN wind speed)

$\sigma_8$

$B_{\mathrm{AGN1}}$

$B_{\mathrm{AGN2}}$

$Y_{200c}\,\Omega_m\,/M_{200c}^{5/3}$ (scaled)

$M_{200c}\,(h^{-1}M_{\odot})$

High
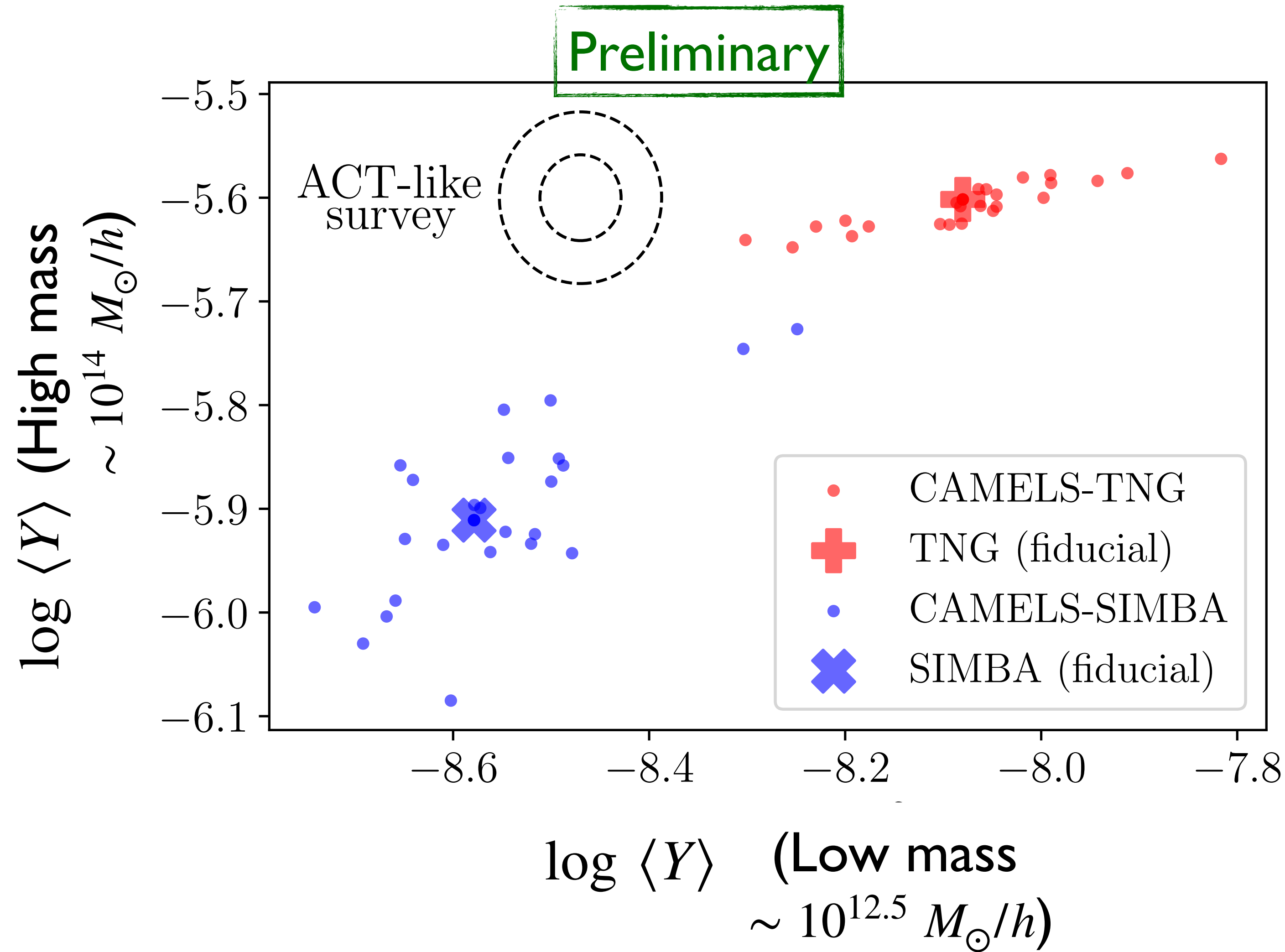
Low
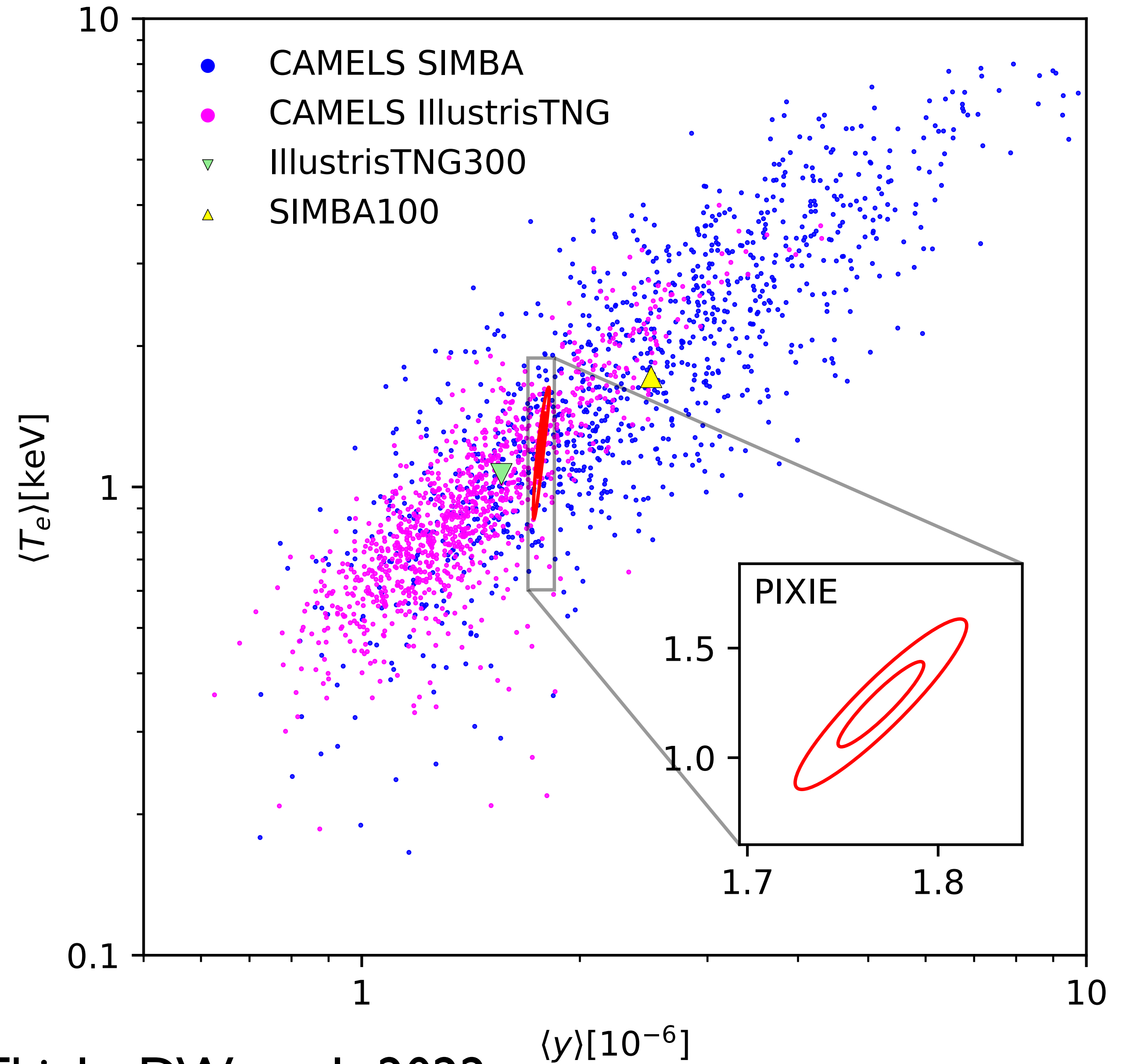
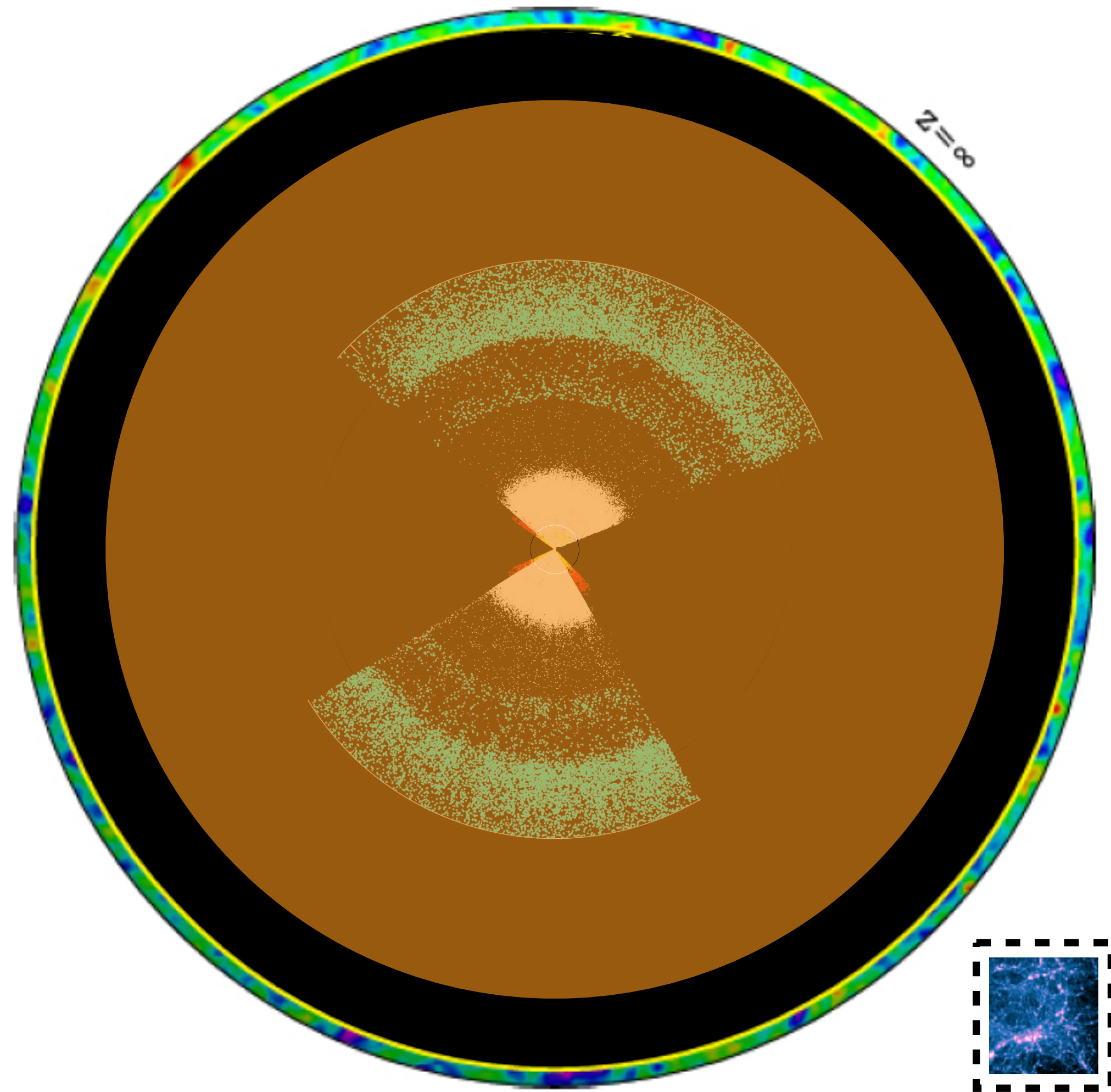# Constraints on sub-grid models

Similarly, CMB spectral distortions
can also constrain baryonic feedback
(% level constraints)
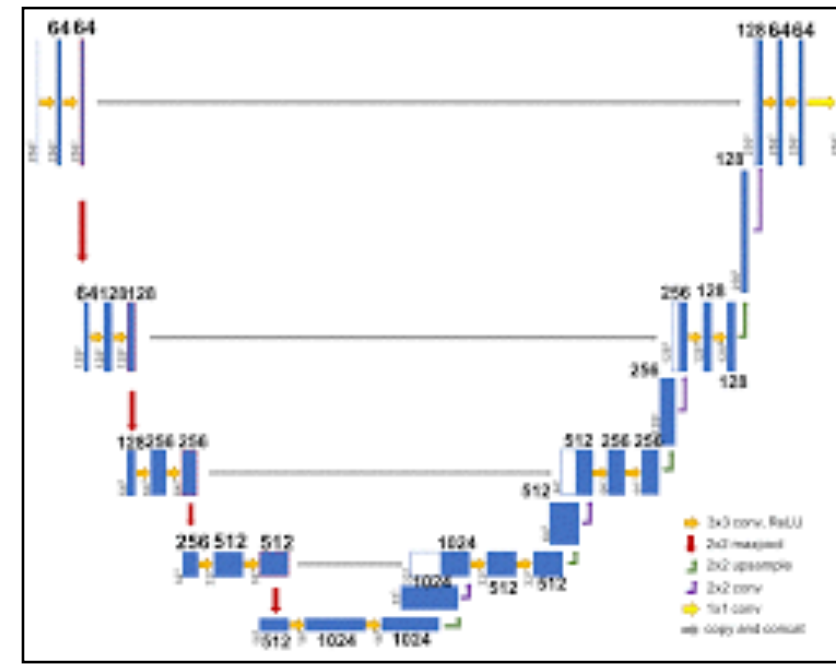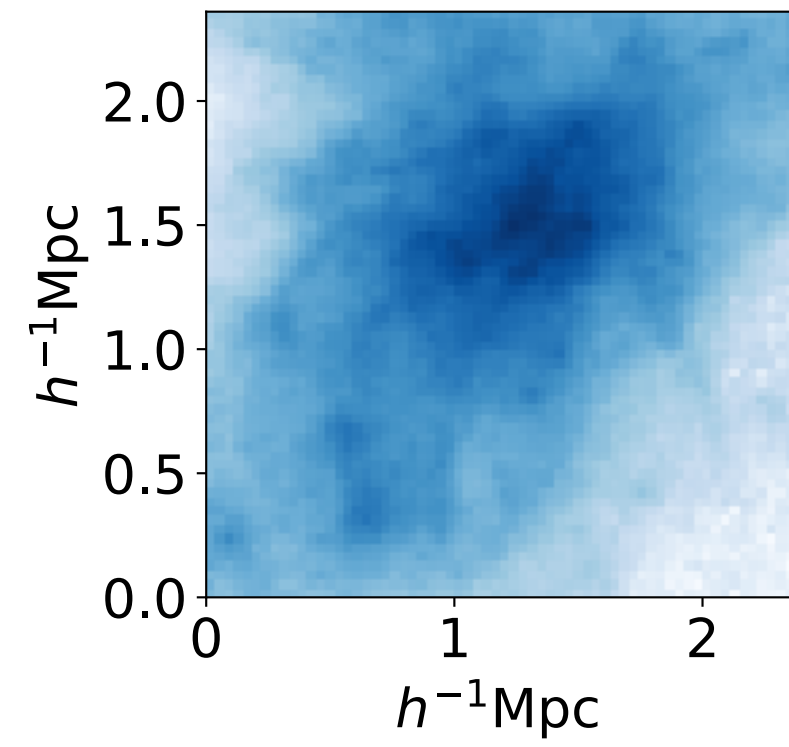


L. Thiele, DW, et al., 2022

# ML for emulation of hydro simulations for future surveys
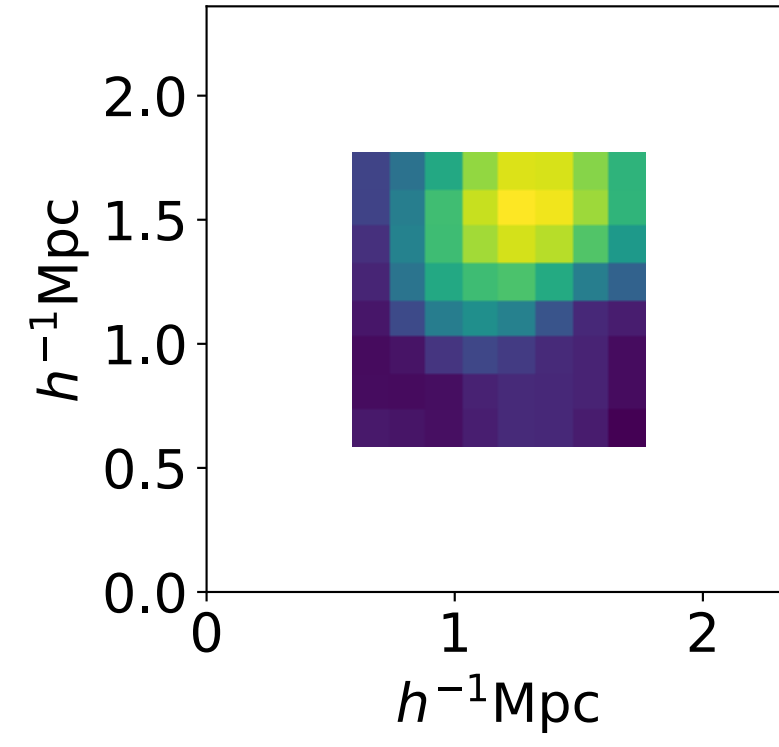


- Volume of upcoming surveys like DESI: $\sim \mathcal{O}$ (10-100 Gpc$^3$)

- Hydro sims are expensive: $\sim$10 million CPU hours for (0.001 Gpc$^3$)

- Needed to study non-linear scales where baryonic effects dominate

# ML for emulation of hydro simulations for future surveys

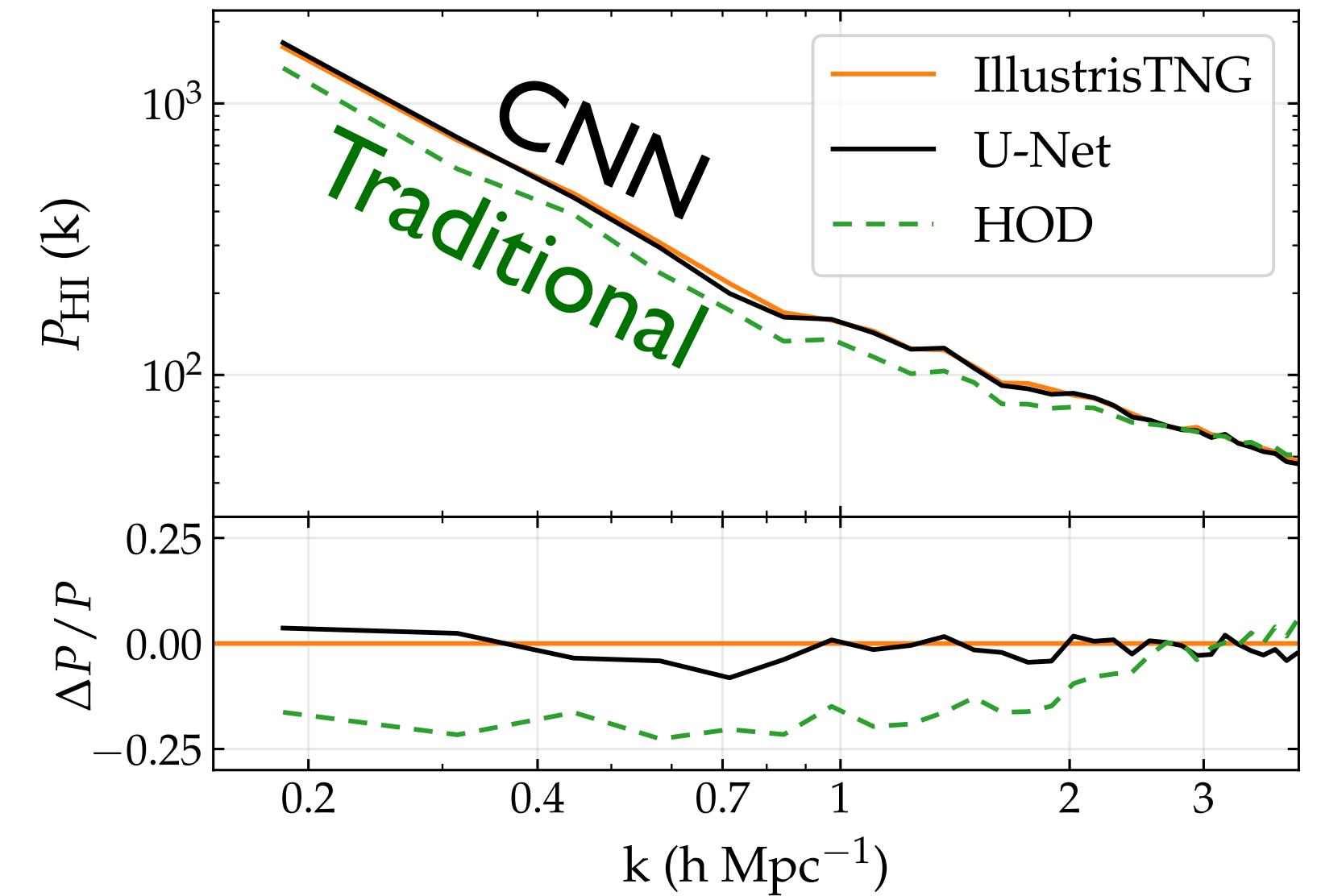**DM (cheap)**



**Neut. hydrogen HI (expensive)**



## Power spec. (2 pt)

# ML for emulation of hydro simulations for future surveys

DM (cheap)



Neut. hydrogen
HI (expensive)



## Power spec. (2 pt)



CNN

Traditional

- IllustrisTNG
- U-Net
- HOD

N-body $\longrightarrow$ Galaxies

N-body $\longrightarrow$ N-body + Neutrinos

ZA (theoretical) $\longrightarrow$ N-body

Sims/Data $\longrightarrow$ Cosmo. parameters

Low res. N-body $\longrightarrow$ High res. N-body

# ML for emulation of hydro simulations for future surveys

**DM (cheap)**



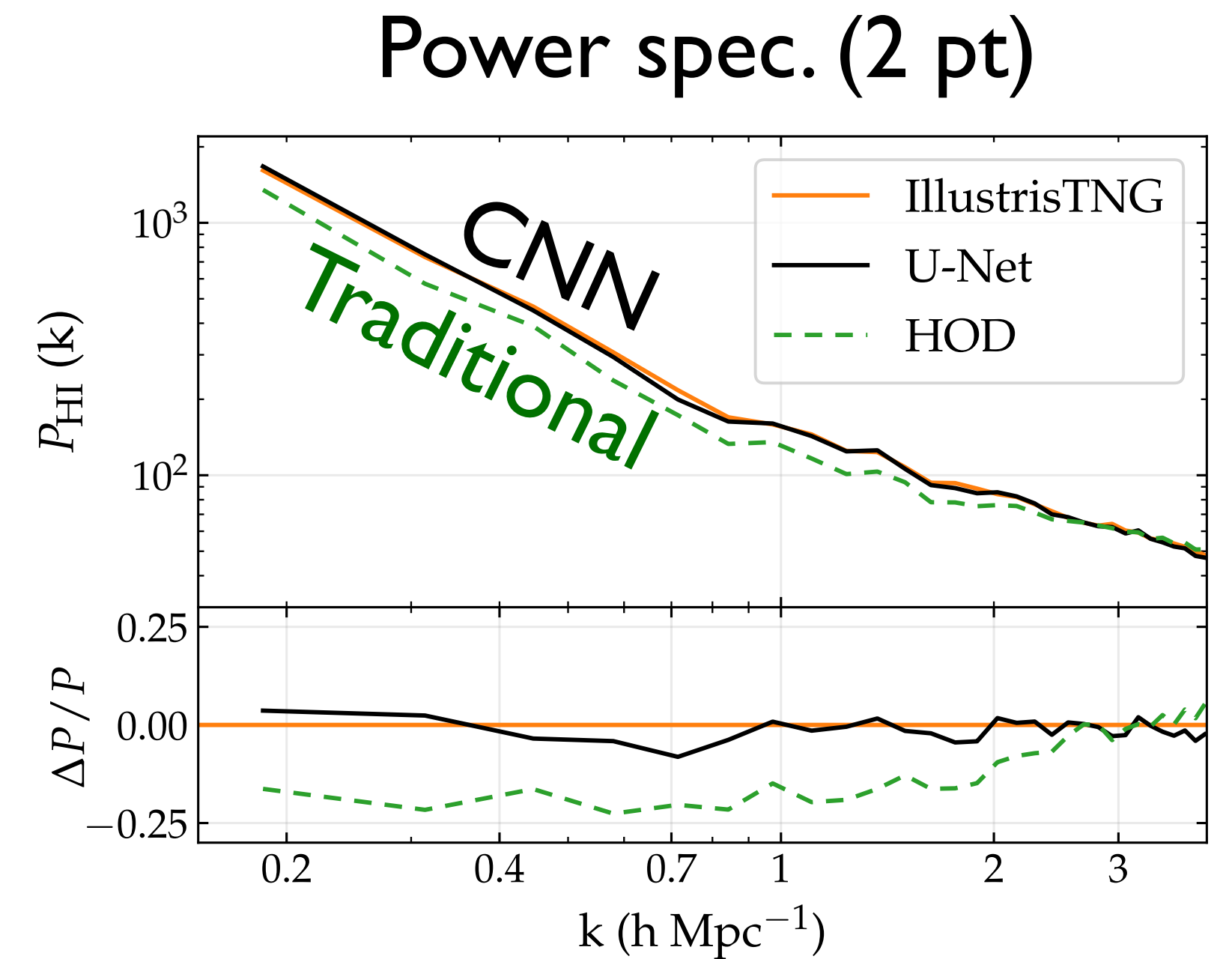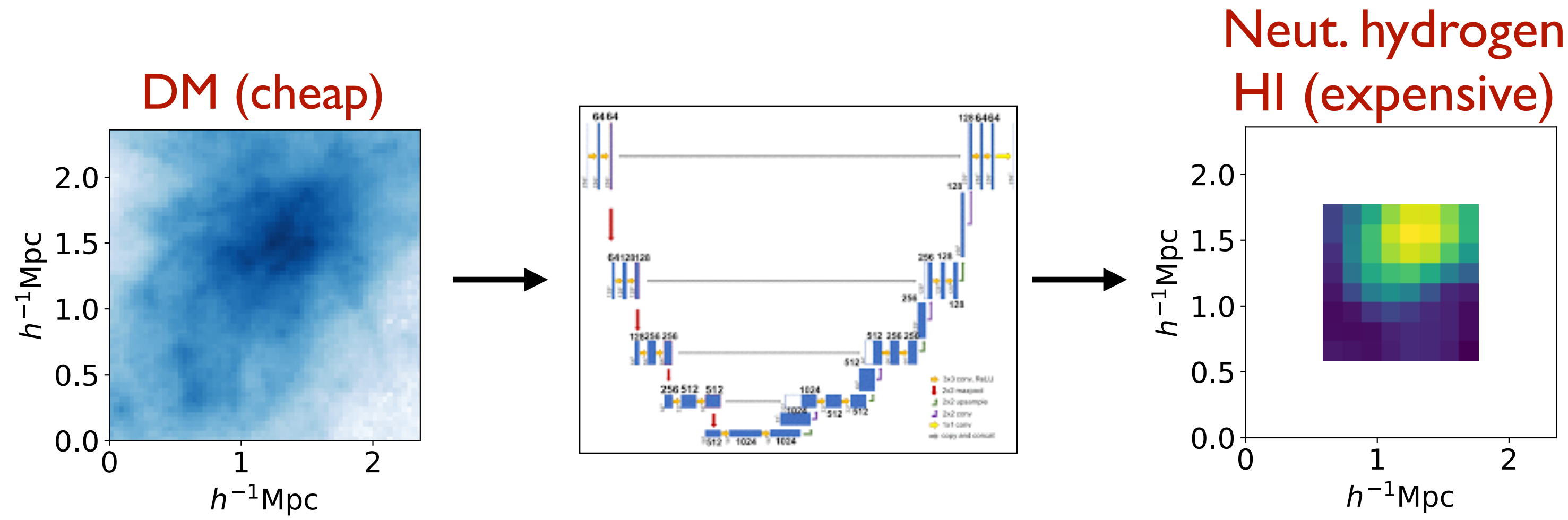**Neut. hydrogen HI (expensive)**



**Power spec. (2 pt)**



N-body $\longrightarrow$ Galaxies

N-body $\longrightarrow$ N-body + Neutrinos

ZA (theoretical) $\longrightarrow$ N-body

Sims/Data $\longrightarrow$ Cosmo. parameters

Low res. N-body $\longrightarrow$ High res. N-body

## Challenges:

1. Robustness to feedback prescriptions
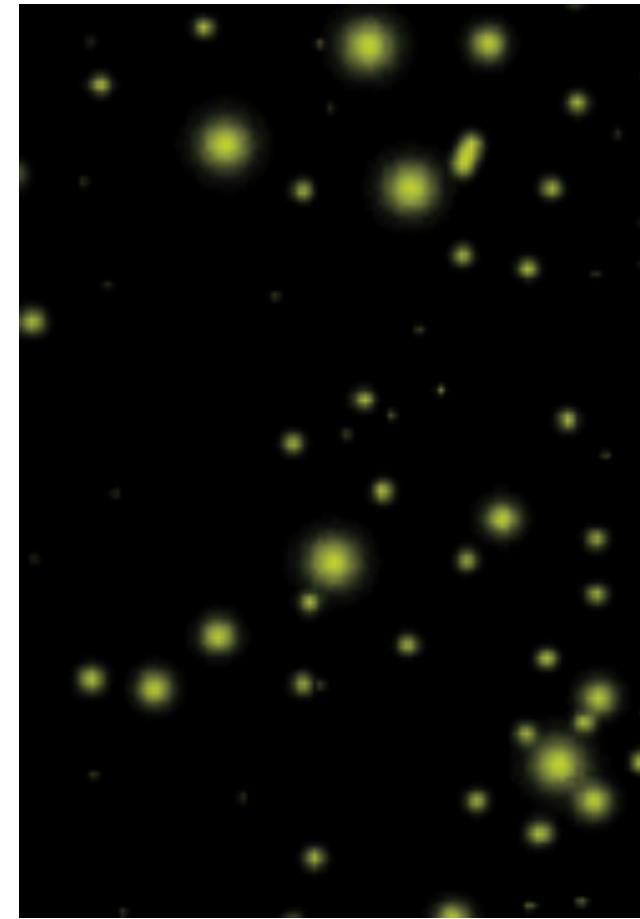2. Robustness to sim resolution
3. Robustness to observational systematics

# ML to model assembly/secondary bias

DM (dark matter)          HI (neutral hydrogen)
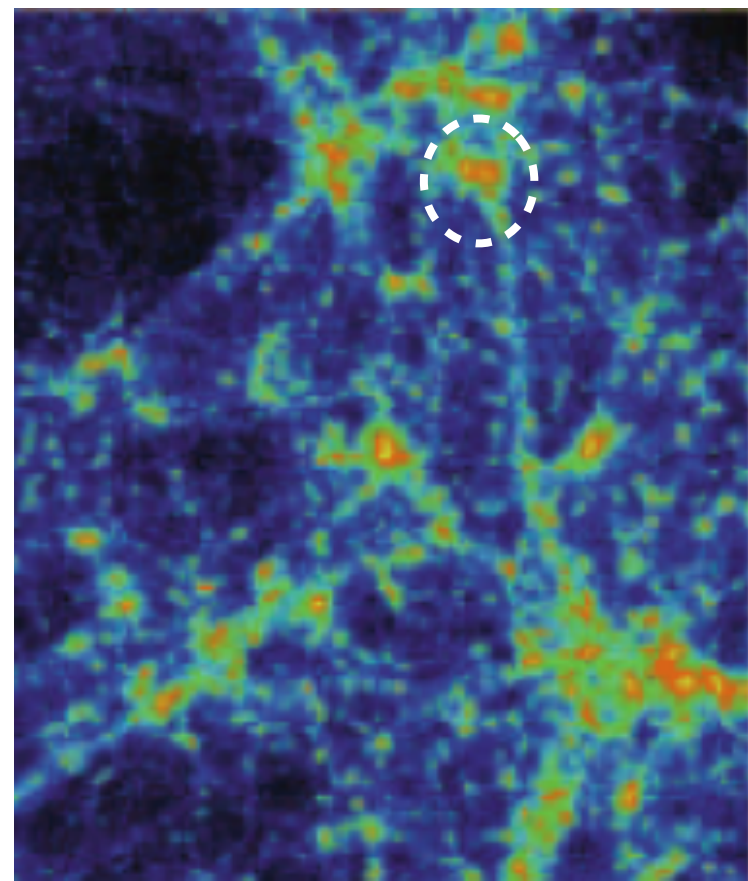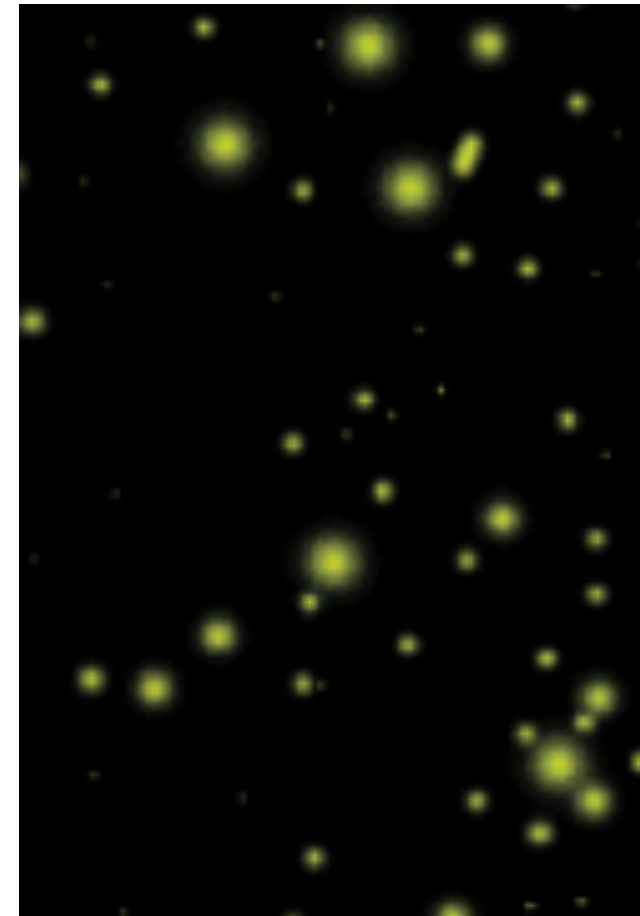
HI mass of halo  =  $f$ (Halo mass only)

(No. of galaxies in halo)

# ML to model assembly/secondary bias

DM (dark matter)          HI (neutral hydrogen)



HI mass of halo  =  $f$ (Halo mass,  *secondary props.* ?)

(No. of galaxies in halo)                    *{local env.,*
                                              *conc.,*
                                              *shear,....}*

# ML to model assembly/secondary bias

**No. of galaxies in a halo**   $= f$ (Halo mass, *environmental shear and overdensity* )

$$N_{\text{sat}}(M_h) = N_{\text{sat}}^{\text{HOD}}(M_h) \times (q' + A)$$

A. Delgado, DW, et al. 21

$$N_{\text{cen}}(M_h) = N_{\text{cen}}^{\text{HOD}}(M_h) \times$$

$$\left[ 1 + B(\delta'_{\text{env}} - \overline{\delta'_{\text{env}}})(1 - N_{\text{cen}}^{\text{HOD}}) \right]$$

**Neutral hydrogen content of halo**   $= f$ (Halo mass, *environmental shear and overdensity*)
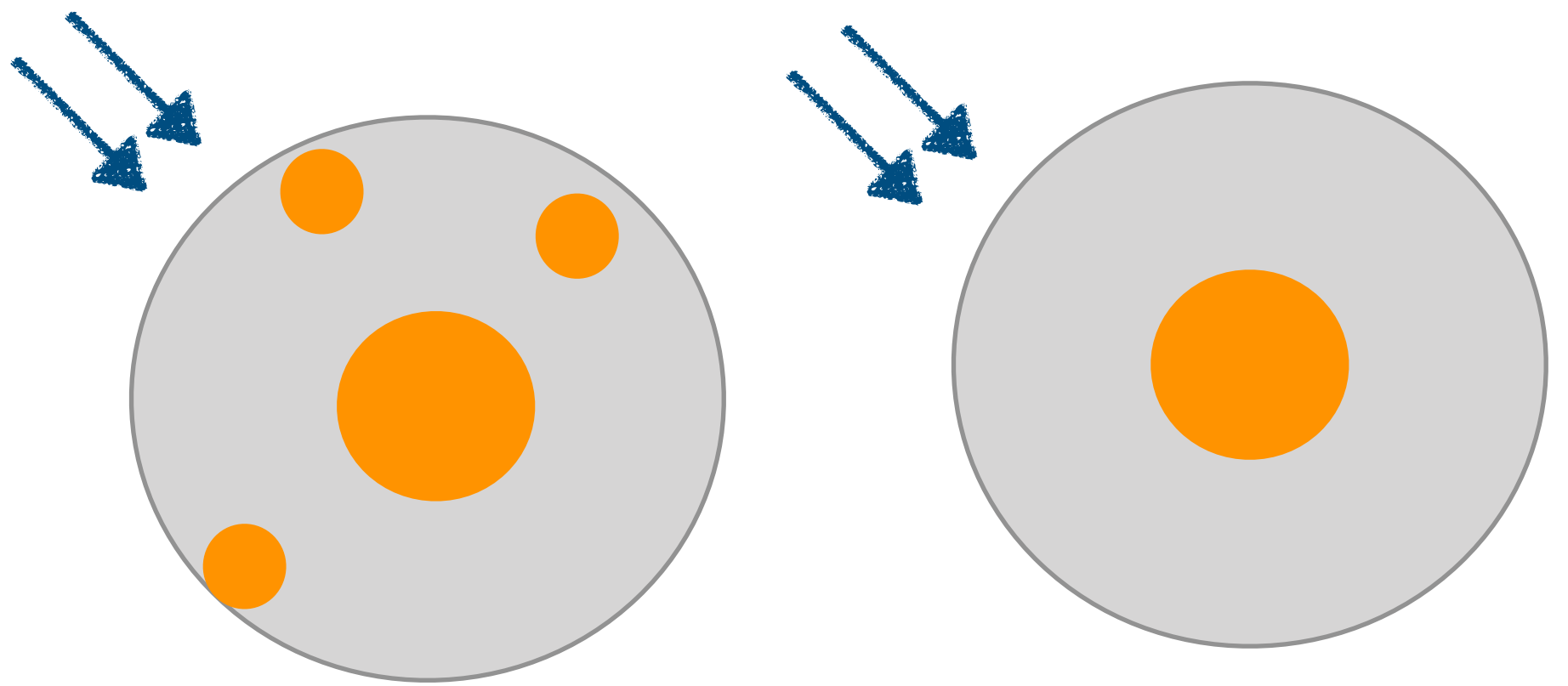
DW et al. 20

$$\frac{M_{\text{HI}}}{M_{\text{HOD}}} = 0.81 + 1.44 \, \alpha'_{0.5} \, m_{10}$$

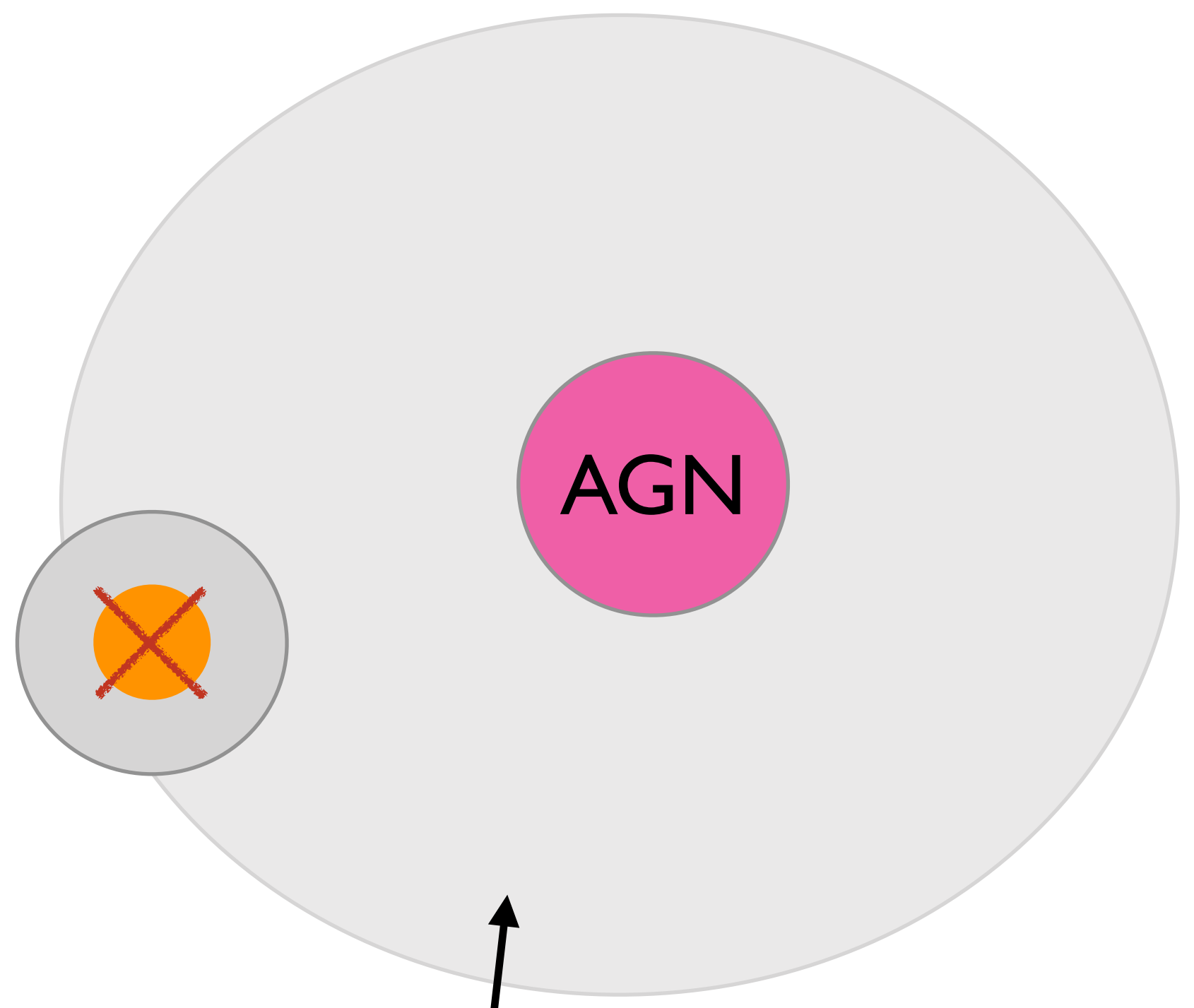$$- 0.57 \, (\alpha'^2_{0.5} \, m_{10}^2 + \alpha'_{0.5} \, \delta'_5)$$

# Why does HI content have env. dependence?

$$\frac{M_{\mathrm{HI}}}{M_{\mathrm{HOD}}} = 0.8 + 1.4\,\alpha'_{0.5}m_{10} - 0.6\,(\alpha'^2_{0.5}m^2_{10} + \alpha'_{0.5}\delta'_5)$$
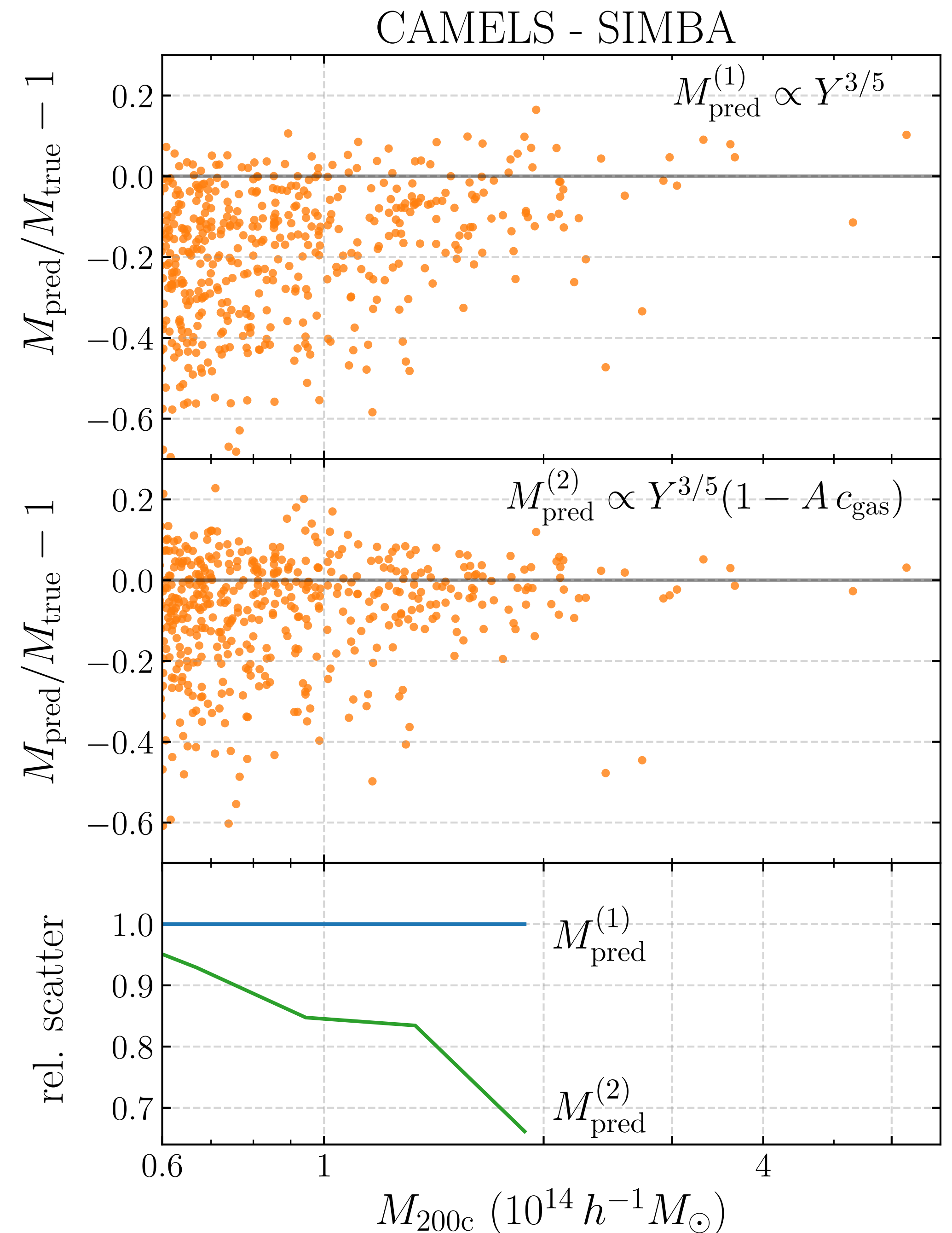
UV + X-ray
background

AGN

Denser env.
→ more mergers

Denser env.: Ionized medium
(ram pressure stripping)

# Summary

★ Symbolic regression can be used to *augment astrophysical scaling relations* and increase their precision

- Using gas conc. reduces scatter in SZ mass estimates by 20-30% for large clusters

- Including stellar to gas mass ratio reduces deviation from self-similarity by factor >2

➡️Suggestions for other scaling relations?



CAMELS - SIMBA

$M_{\rm pred}^{(1)} \propto Y^{3/5}$

$M_{\rm pred}^{(2)} \propto Y^{3/5}(1 - A\,c_{\rm gas})$

$M_{\rm pred}^{(1)}$

$M_{\rm pred}^{(2)}$

$M_{200\rm c}\ (10^{14}\ h^{-1}M_\odot)$
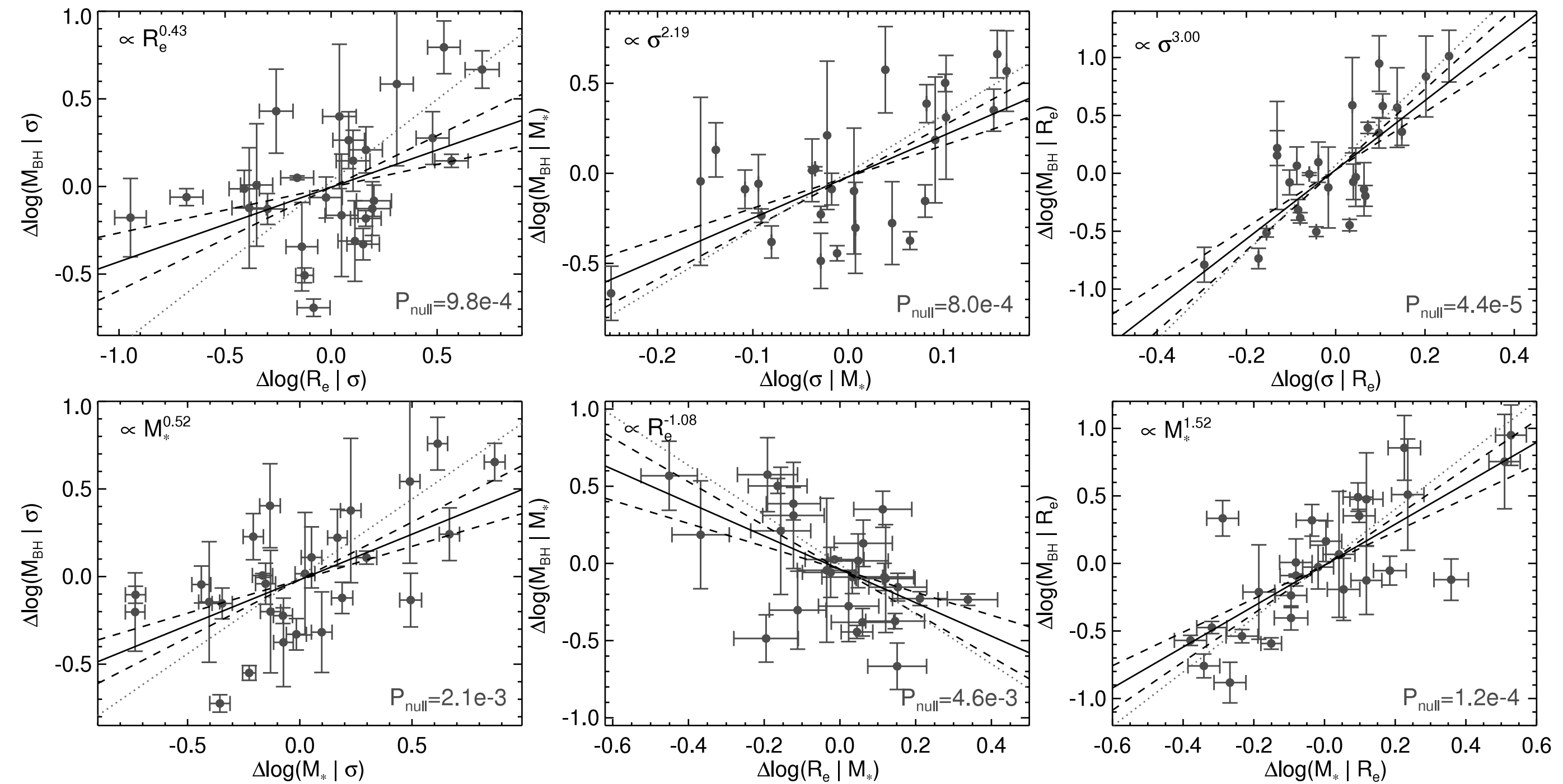
# Application to other scaling relations?

- Philips relation for supernovae

$$M_{\max}(B) = -21.726 + 2.698\,\Delta m_{15}(B)$$

- Cepheid P-L relation

$$M_v = A(\log_{10} P - 1) - B$$

- Tully fisher relation
- Black hole-bulge mass relation
- Fundamental plane relation
- ….



Hopkins et al. 07