

# Current progress and challenges from the **CAMELS** project

(**C**osmology and **A**strophysics with **M**achin**E** **L**earning **S**imulations)

Core team:

Francisco Villaescusa-Navarro

Daniel Anglés-Alcázar

Shy Genel

**Daniel Anglés-Alcázar**

Department of Physics, University of Connecticut

Debating the potential of machine learning in astronomical surveys #2 - ML-IAP/CCA-2023

Flatiron institute / IAP, November 27th - December 1st, 2023

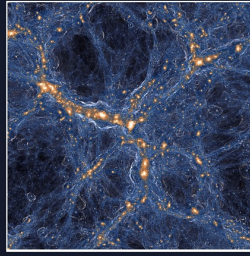
# Cosmology and Astrophysics with Machine Learning Simulations

Core team:

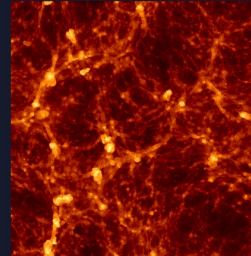
Francisco Villaescusa-Navarro

Daniel Anglés-Alcázar

Shy Genel



IllustrisTNG team



SIMBA team

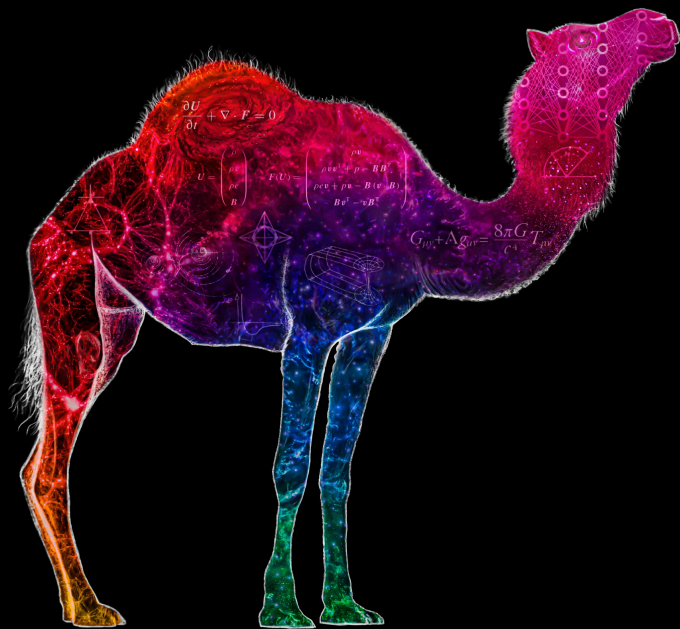


Astrid team



SMAUG collaboration

+ SWIFT-Eagle  
+ Magneticum  
+ Ramses  
+ Enzo  
+ ...



# CAMELS

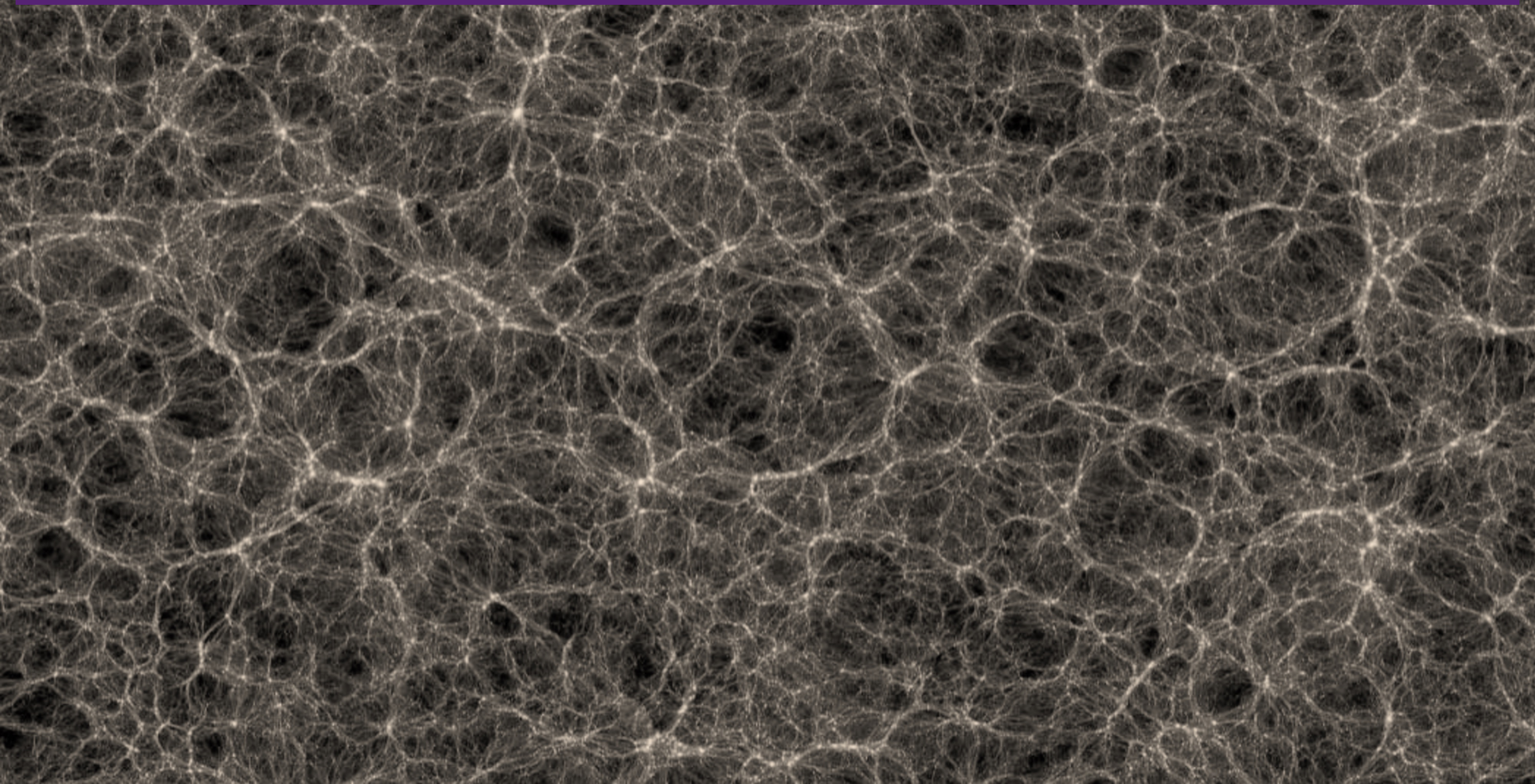
Adrian Bayer  
Alex Barreira  
Ana Maria Delgado  
Andrina Nicola  
Alice Pisani  
Benjamin Oppenheimer  
Benjamin Wandelt  
Blakesley Burkhart  
ChangHoon Hahn  
Colin Hill  
Core Francisco Park  
Daisuke Nagai  
Desika Narayanan  
David Spergel  
Emily Moser  
Erwin T. Lau

Faizan Mohammad  
Gabrielle Paribelli  
Greg Bryan  
Gabiella Contardo  
Helen Shao  
Jay Wadekar  
Jingjing Shi  
Joyce Caliendo  
Lucia Perez  
Lars Hernquist  
Leander Thiele  
Luis F. Machado Poletti  
Matteo Viel  
Matthew Gebhardt  
Megan Tillman  
Michael Eickenberg

Neerav Kaushal  
Nicholas Battaglia  
Oliver Philcox  
Pablo Villanueva-Domingo  
Rachel Somerville  
Romeel Dave  
Stephanie Tonnensen  
Sultan Hassan  
Romain Teyssier  
Ulrich Steinwandel  
Valentina La Torre  
Vid Irsic  
William Coulton  
Yin Li  
Yongseok Jo  
Yueying Ni



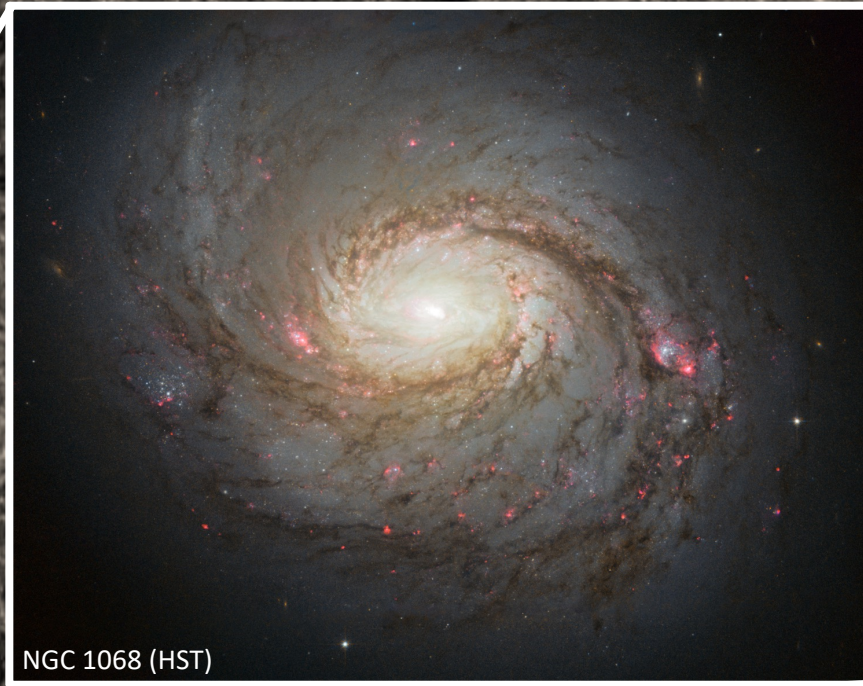
Large scale structure → ...galaxy formation physics... → cosmology?





Large scale structure → ...galaxy formation physics... → cosmology?

We only see the tip of the iceberg!



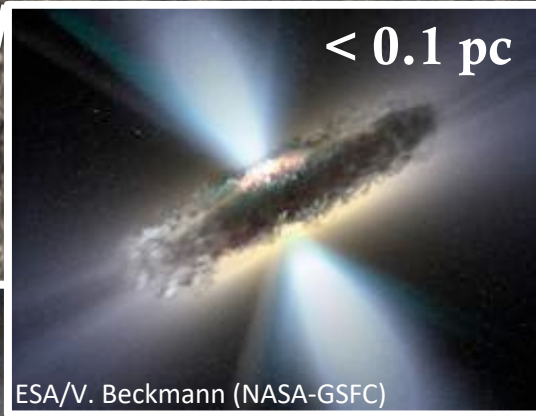
NGC 1068 (HST)

Galaxies form at the centers of dark matter halos

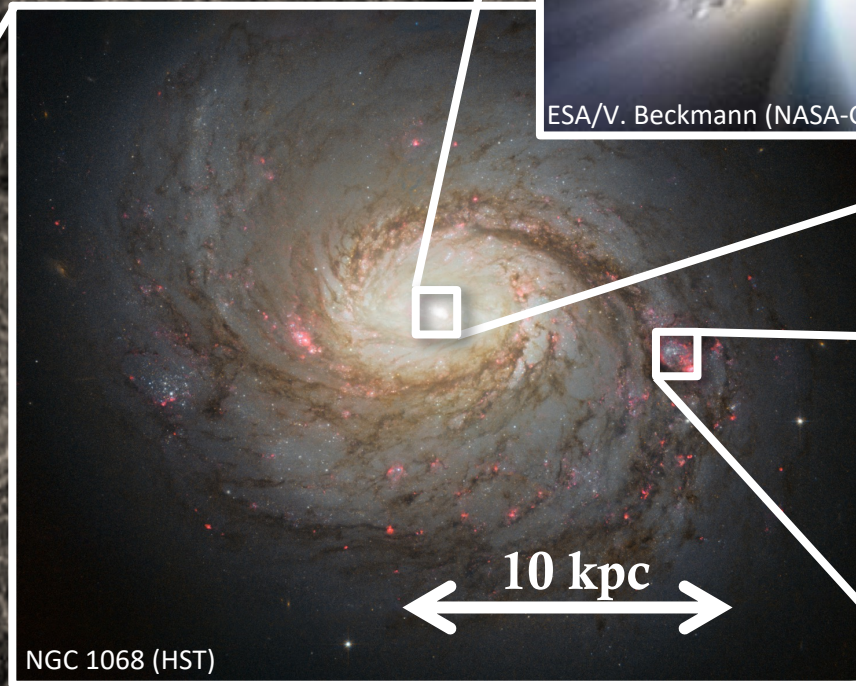


# Large scale structure → ...galaxy formation physics... → cosmology?

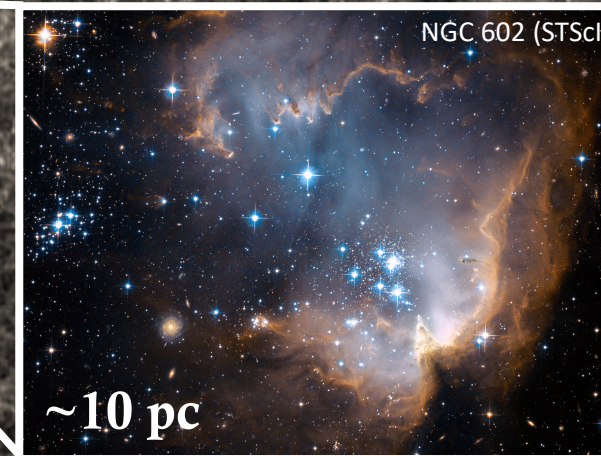
**We only see the tip of the iceberg!**



Supermassive black holes (SMBH) grow at the centers of galaxies and likely affect their evolution via radiation, winds, jets...



Massive stars affect their surrounding interstellar medium through supernovae, radiation, winds...



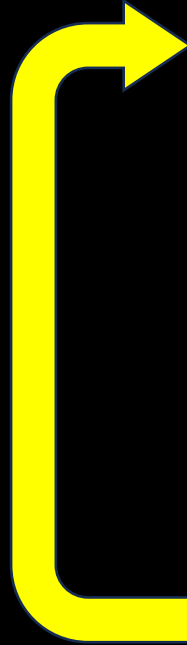
Galaxies form at the centers of dark matter halos

25 Mpc



**Dream goal in galaxy formation simulations:**

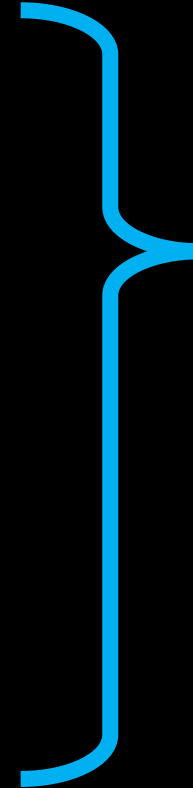
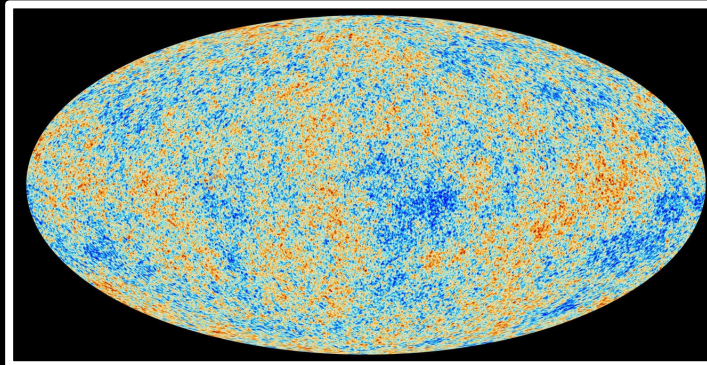
predict detailed properties of millions of galaxies starting from cosmological initial conditions using 'ab-initio' physics



Galaxy surveys



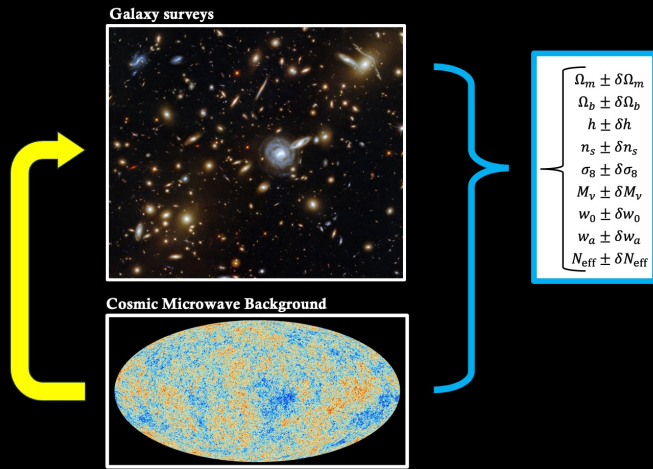
Cosmic Microwave Background



- $\Omega_m \pm \delta\Omega_m$
- $\Omega_b \pm \delta\Omega_b$
- $h \pm \delta h$
- $n_s \pm \delta n_s$
- $\sigma_8 \pm \delta\sigma_8$
- $M_\nu \pm \delta M_\nu$
- $w_0 \pm \delta w_0$
- $w_a \pm \delta w_a$
- $N_{\text{eff}} \pm \delta N_{\text{eff}}$

**Dream goal in Cosmology:**  
Infer accurately and without bias the cosmological parameters using the full amount of information in the observable Universe





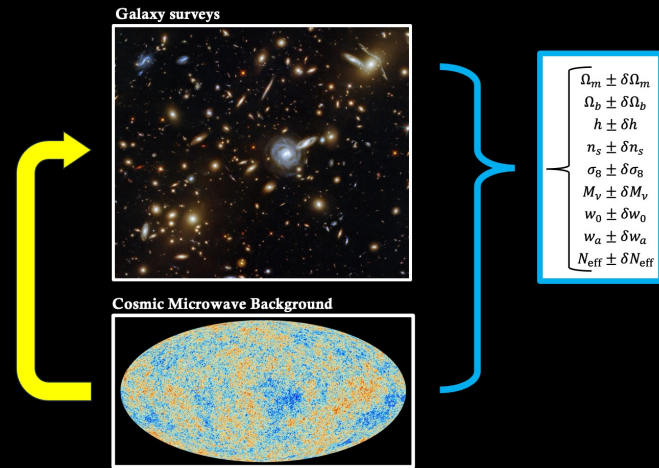
**PROBLEM 1:** Galaxy formation simulations rely on “sub-grid” models for unresolved processes that are still poorly understood

**PROBLEM 2:** Lots of cosmological information on small scales inaccessible due to impact of uncertain astrophysical processes

**PROBLEM 3:** The optimal summary statistic to extract cosmological information is unknown

**PROBLEM 4:** Need to speed up simulations to predict cosmological observables for large cosmological volumes





**PROBLEM 1:** Galaxy formation simulations rely on “sub-grid” models for unresolved processes that are still poorly understood

**PROBLEM 2:** Lots of cosmological information on small scales inaccessible due to impact of uncertain astrophysical processes

**PROBLEM 3:** The optimal summary statistic to extract cosmological information is unknown

**PROBLEM 4:** Need to speed up simulations to predict cosmological observables for large cosmological volumes

## The CAMELS approach

Run thousands of simulations spanning the full range of uncertainty in galaxy formation physics and train machine learning algorithms to extract the maximum amount of cosmological information at the field level while marginalizing over uncertainties in baryonic effects



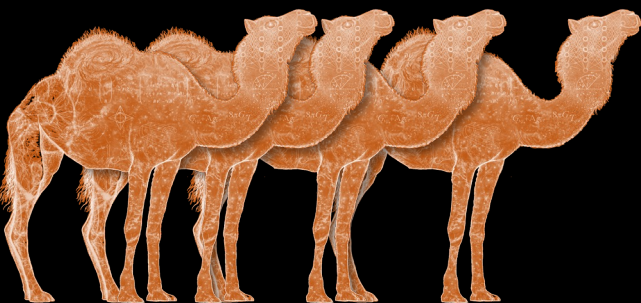


# The CAMELS suites

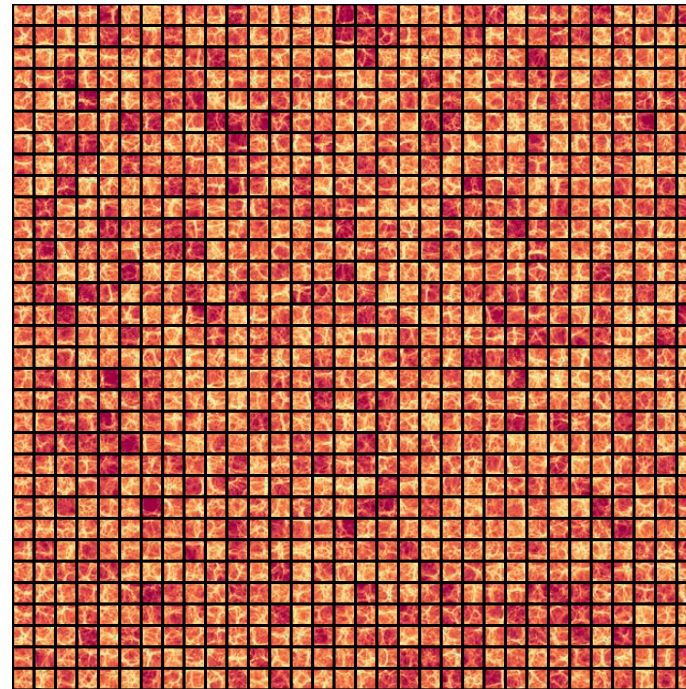
- >10,000 cosmological boxes of  $(25\text{Mpc}/h)^3$
  - >5,000 variations of TNG, SIMBA, and ASTRID
    - cosmological params ( $\Omega_m, \sigma_8, \dots$ )
    - astrophysical params (feedback)
  - >5,000 corresponding DM-only simulations
- Additional simulation sets:
- Same ICs, varying one parameter
  - Fiducial model, varying the ICs

## CAMELS-SAM (Lucía Perez)

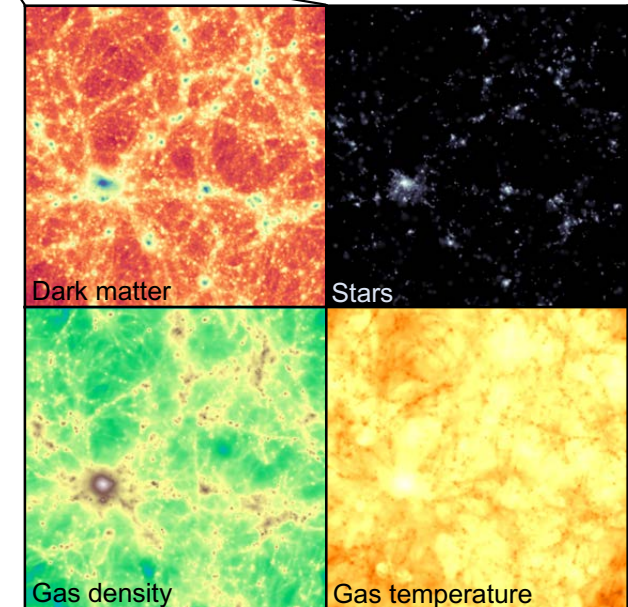
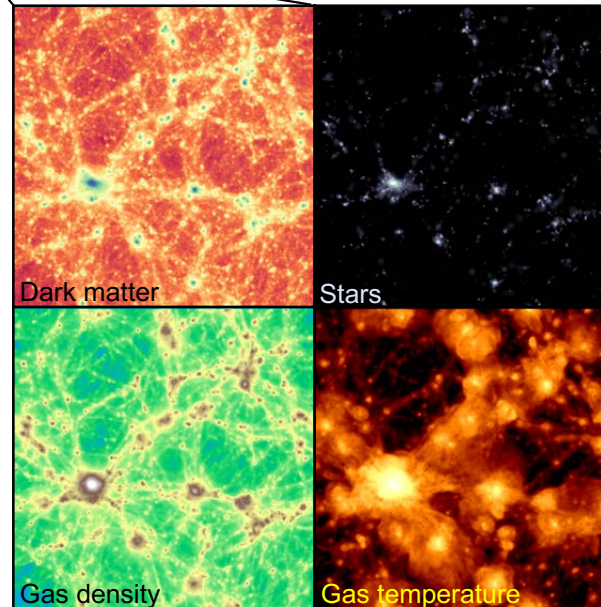
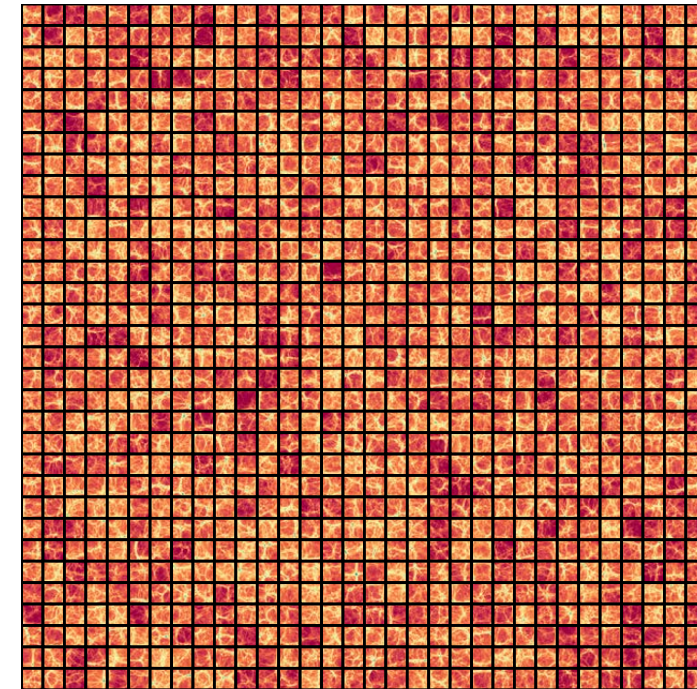
- 1,000 DM-only simulations of  $(100\text{Mpc}/h)^3$
- 2 cosmological params ( $\Omega_m, \sigma_8$ )
  - Santa-Cruz SAM parameter variations



IllustrisTNG



SIMBA





# CAMELS

SAM



Lucia A. Perez

Princeton Future Faculty in  
the Physical Sciences Fellow  
& CCA Flatiron Research  
Fellow

## New large-volume simulation 'hump' of CAMELS project

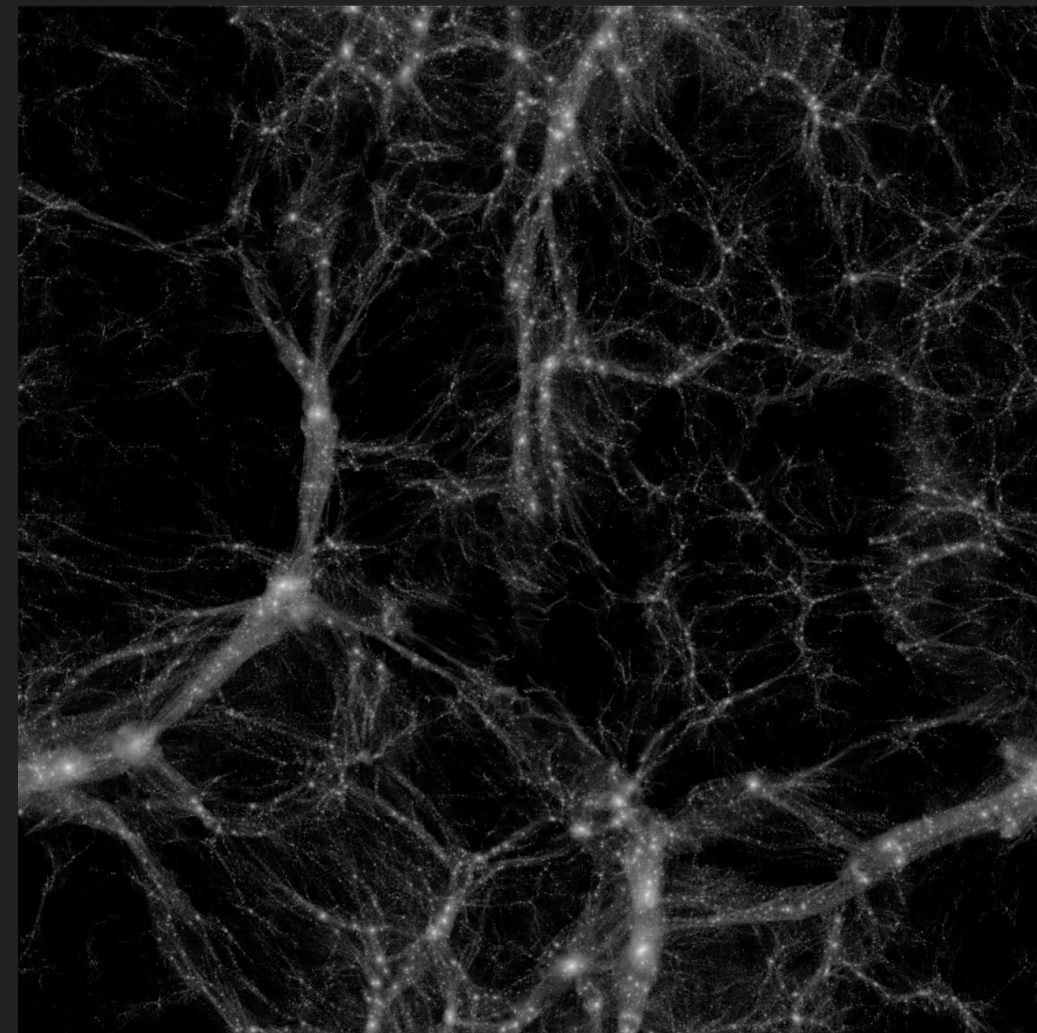
- **CAMELS (Cosmology & Astrophysics with MachinE Learning Simulations):** machine learning data sets to create predictions for observations, marginalize over astrophysics to learn cosmology, and identify useful summary statistics and analyses
- **1000+ N-body simulations:**  $(100 h^{-1} \text{ Mpc})^3$  large ;  $N=640^3$  particles of  $\sim 1-6 \times 10^8 h^{-1} M_{\text{sol}}$  ; 100 snapshots between  $0 < z < 27$
- Cosmological parameter space:  $\Omega_m$  (fraction of energy density in DM+baryons) &  $\sigma_8$  ( $\sim$ amplitude of density fluctuations)
- Run through the **Santa Cruz Semi-Analytic Model:**

“ $A_{\text{SN}}$ ”: mass outflow + reheating rates of cold gas due to SNe + stars

“ $A_{\text{AGN}}$ ”: AGN feedback, how much mass ejected in radio jets?

Data is public! [camels-sam.readthedocs.io](https://camels-sam.readthedocs.io) | [arxiv.org/abs/2204.02408](https://arxiv.org/abs/2204.02408)

Proof-of-concept in Perez+2022: constraining power of galaxy clustering statistics  
(3D two-point correlation function, count-in-cells, Void Probability Function)



LH\_643:  $\Omega_m = 0.131$  ;  $\sigma_8 = 0.986$



# CAMELS public data repository

<https://camels.readthedocs.io>

## The CAMELS project: public data release

FRANCISCO VILLAESCUSA-NAVARRO<sup>1,2</sup> SHY GENEL,<sup>1,3</sup> DANIEL ANGLÉS-ALCÁZAR<sup>4,1</sup> LUCIA A. PEREZ,<sup>5</sup>  
PABLO VILLANUEVA-DOMINGO<sup>6</sup> DIGVIJAY WADEKAR,<sup>7,8</sup> HELEN SHAO,<sup>2</sup> FAIZAN G. MOHAMMAD<sup>9,10</sup>  
SULTAN HASSAN,<sup>1,11</sup> EMILY MOSER<sup>12</sup> ERWIN T. LAU,<sup>13</sup> LUIS FERNANDO MACHADO POLETTI VALLE<sup>14</sup>  
ANDRINA NICOLA,<sup>2</sup> LEANDER THIELE<sup>15</sup> YONGSEOK JO,<sup>16</sup> OLIVER H. E. PHILCOX,<sup>2,8</sup> BENJAMIN D. OPPENHEIMER,<sup>17,13</sup>  
MEGAN TILLMAN<sup>18</sup> CHANGHOON HAHN<sup>19</sup> NEERAV KAUSHAL<sup>19</sup> ALICE PISANI<sup>1,20,2</sup> MATTHEW GEBHARDT,<sup>4</sup>  
ANA MARIA DELGADO,<sup>13</sup> JOYCE CALIENDO,<sup>4,21</sup> CHRISTINA KREISCH,<sup>2</sup> KAZE W.K. WONG,<sup>1</sup> WILLIAM R. COULTON,<sup>1</sup>  
MICHAEL EICKENBERG,<sup>22</sup> GABRIELE PARIMBELLI<sup>23,24,25,26,27</sup> YUEYING NI,<sup>28</sup> ULRICH P. STEINWANDEL<sup>1</sup>  
VALENTINA LA TORRE,<sup>29</sup> ROMEEL DAVE,<sup>30,11,31</sup> NICHOLAS BATTAGLIA,<sup>12</sup> DAISUKE NAGAI,<sup>32</sup> DAVID N. SPERGEL,<sup>1,2</sup>  
LARS HERNQUIST,<sup>13</sup> BLAKESLEY BURKHART,<sup>18,1</sup> DESIKA NARAYANAN,<sup>33,34</sup> BENJAMIN WANDEL,<sup>35,1</sup>  
RACHEL S. SOMERVILLE,<sup>1</sup> GREG L. BRYAN,<sup>36,1</sup> MATTEO VIEL<sup>25,27,26,37</sup> YIN LI<sup>1,22</sup> VID IRSIC,<sup>38,39</sup>  
KATARINA KRALJIC,<sup>40</sup> AND MARK VOGELSBERGER<sup>41</sup>

- Large number of labeled data products in the form of 1D, 2D, and 3D arrays
- Full documentation and metadata available
- Designed to enable a broad range of creative AI applications
- Public access to full data, (limited) local computing, and tutorials

arXiv:2201.01300

# CAMELS

Search docs

## CAMELS

News

Scientific goals

Publications

Data Access

Citation

## SIMULATIONS

General description

Suites and sets

Simulations codes

Simulations chart

Simulation parameters

Redshifts

## DATA PRODUCTS

Data organization

Simulations

SUBFIND catalogues

SubLink catalogues

Rockstar catalogues

AHF catalogues

CAESAR catalogues

Power spectra

Bispectra

Probability distribution functions

VIDE Voids

Lyman-alpha spectra

X-Rays

CAMELS CGM Profiles

CAMELS Multifield Dataset

CAMELS-SAM

Home / CAMELS

Edit on GitHub

## CAMELS

CAMELS stands for Cosmology and Astrophysics with MachinE Learning Simulations, and it is a project that aims at building bridges between cosmology and astrophysics through numerical simulations and machine learning. CAMELS contains 10,680 cosmological simulations – 5,164 N-body and 5,516 state-of-the-art (magneto-)hydrodynamic – and more than 700 Terabytes of data. CAMELS is the largest set of cosmological hydrodynamic simulations ever run.

Type	Code	Subgrid model	Simulations
Hydrodynamic	Arepo	IllustrisTNG	2,143
	Gizmo	SIMBA	1,092
	MP-Gadget	Astrid	2,116
	OpenGadget	Magneticum	77
	SWIFT	EAGLE	77
	Ramses		5
	Enzo		6
N-body	Gadget-III	–	5,164

Introductory video to the CAMELS project:



The video below shows an example of a CAMELS hydrodynamic simulation run with the Ramses code. Gas density and gas temperature are shown in blue and red, respectively as a function of time. CAMELS contains thousands of simulations like this one.

# First CAMELS results very encouraging!



## Cosmology inference... from 2D maps: (Paco, Yueying Ni, Jonah Rose)

- 2D projected maps of 27 fields (dark matter, gas, stars) with 100 kpc/pixel resolution
- It works! 2-3% error in  $\Omega_m$  and  $\sigma_8$  with all fields combined (3-4% with HI only)
- Extracting information down to the smallest scales (1 pixel), marginalizing over baryonic effects
- But... only the total mass field is robust to differences in galaxy formation model (TNG vs SIMBA)

## ... from summary statistics (Andrina Nicola, Lucía Perez, Ana María Delgado)

## ... from galaxy positions/velocities with GNN (Natalí de Santi, Helen Shao)

## ... and from a single galaxy! (Paco, Nicolas Echeverri-Rojas, Chaitanya Chawak)

## Constraining feedback...

## ... with SZ (Emily Moser, Pandey, Shivam), spectral distortions (Leander Thiele), Ly $\alpha$ (Megan Tillman, Blakesley Burkhart)

## Predicting galaxy/halo properties:

- Finding universal Relations in (sub)halo properties (Helen Shao)
- Inferring halo masses from galaxy properties with GNN (Pablo Villanueva-Domingo)
- Halos mass and CGM properties from X-ray and HI maps (Naomi Gluck)
- Reducing the scatter in the SZ flux-mass relation (Digvijay Wadekar)

Not an exhaustive list!



## Emulation (Sultan Hassan, Chris Lovell, Yongseok Jo, Max Lee, Matt Gebhardt), Inpainting (Faizan Mohammad)

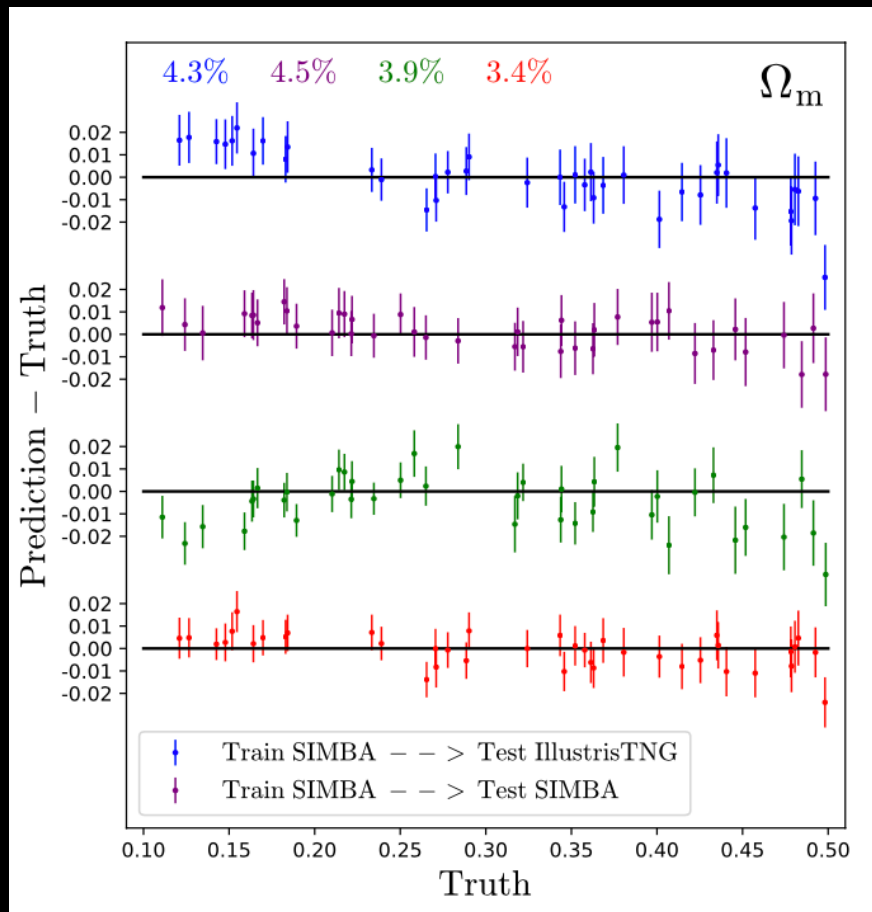


# Cosmological inference at the field level

Villaescusa-Navarro, Anglés-Alcázar, Genel, et al. (2021a,b,c)

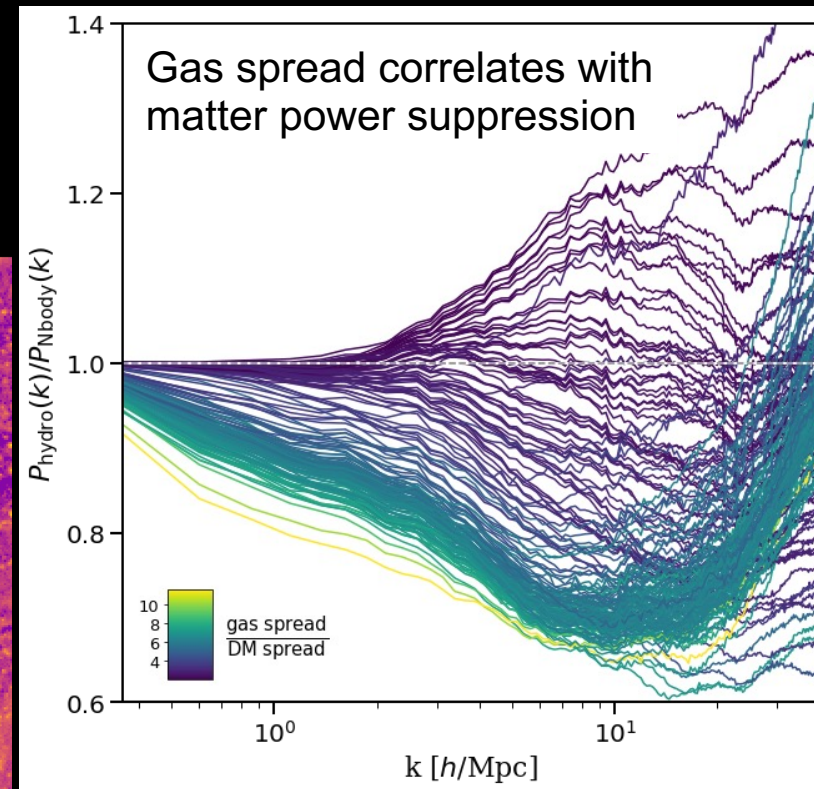
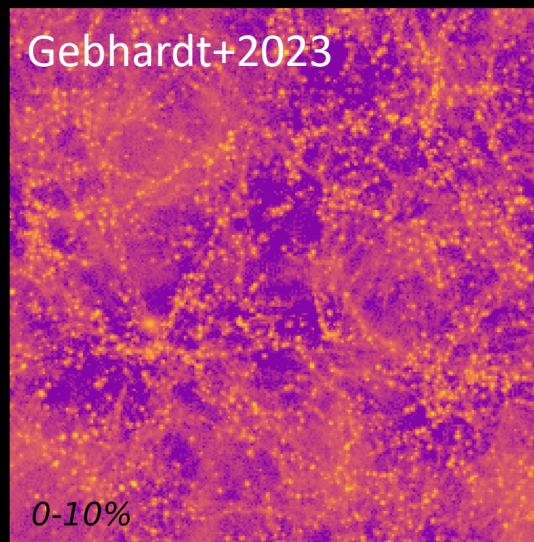
## The 2D total mass field is a robust predictor

- Extracting more information than power spectra
- Down to smallest scale (100 kpc/pixel)
- Marginalizing over baryonic effects



Despite the large impact of baryons  
on the matter power spectrum

Delgado+2023, Gebhardt+2023, Pandey+2023

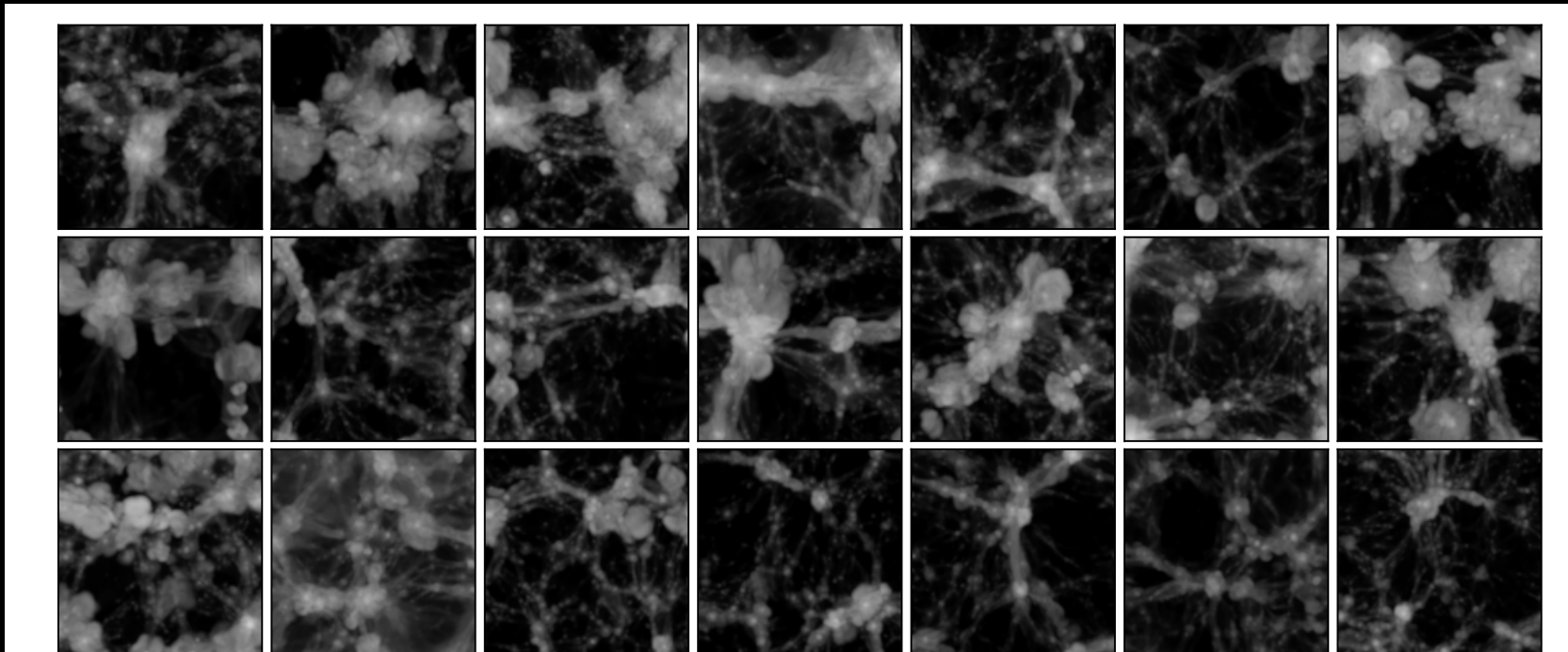
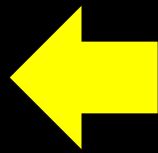
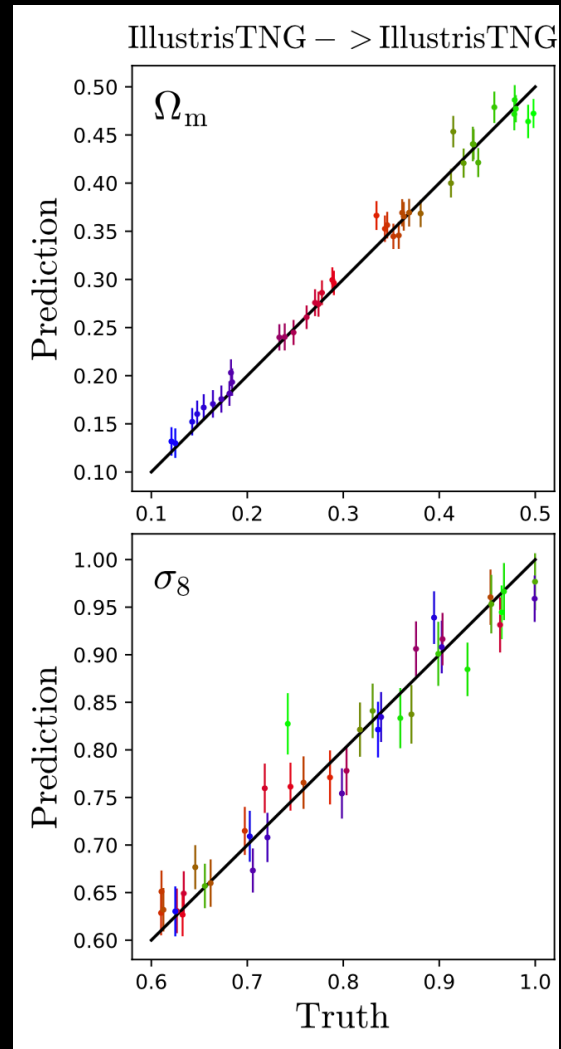




# Cosmological inference at the field level

Villaescusa-Navarro, Anglés-Alcázar, Genel, et al. (2021a,b,c)

→ Train neural network on **temperature maps** to predict input cosmological parameters while marginalizing over sub-grid physics



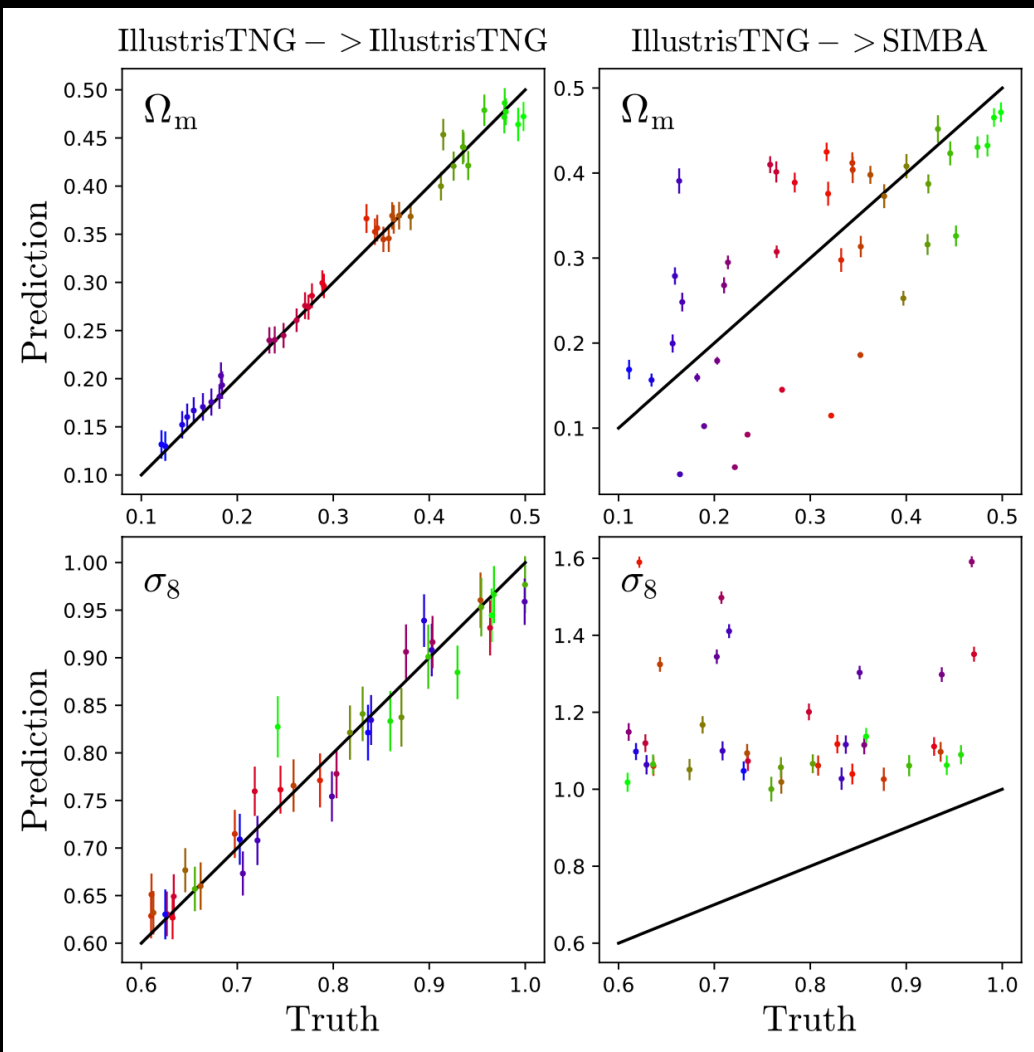
Every map has  $256 \times 256$  pixels, covers an area of  $25 \times 25 (h^{-1} \text{Mpc})^2$ , and has a different cosmology & astrophysics. 15,000 images in total.



# Cosmological inference at the field level

Villaescusa-Navarro, Anglés-Alcázar, Genel, et al. (2021a,b,c)

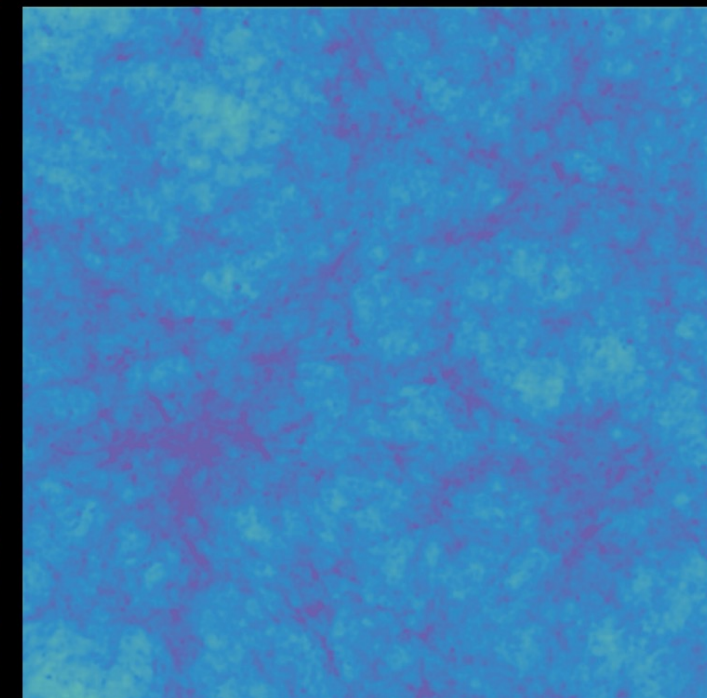
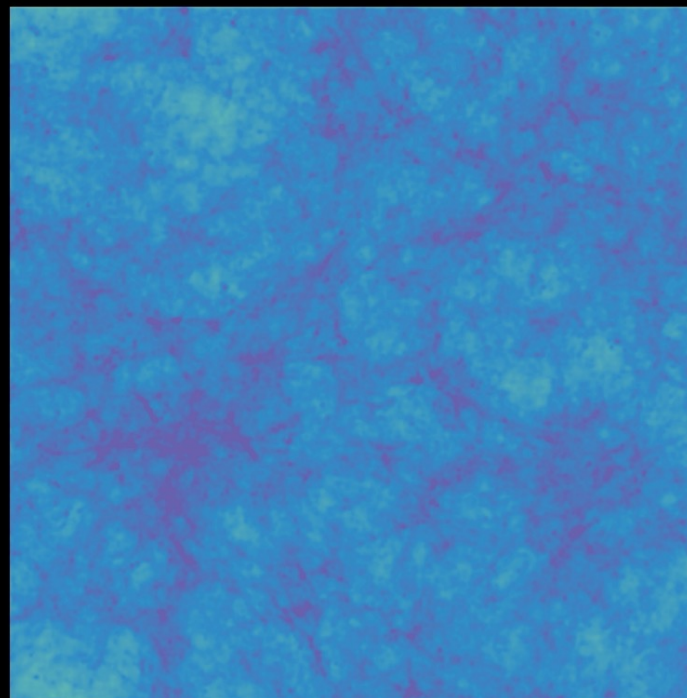
→ Inference from 2D Temperature maps is not robust to galaxy formation physics implementation



IllustrisTNG

Dark matter density

SIMBA



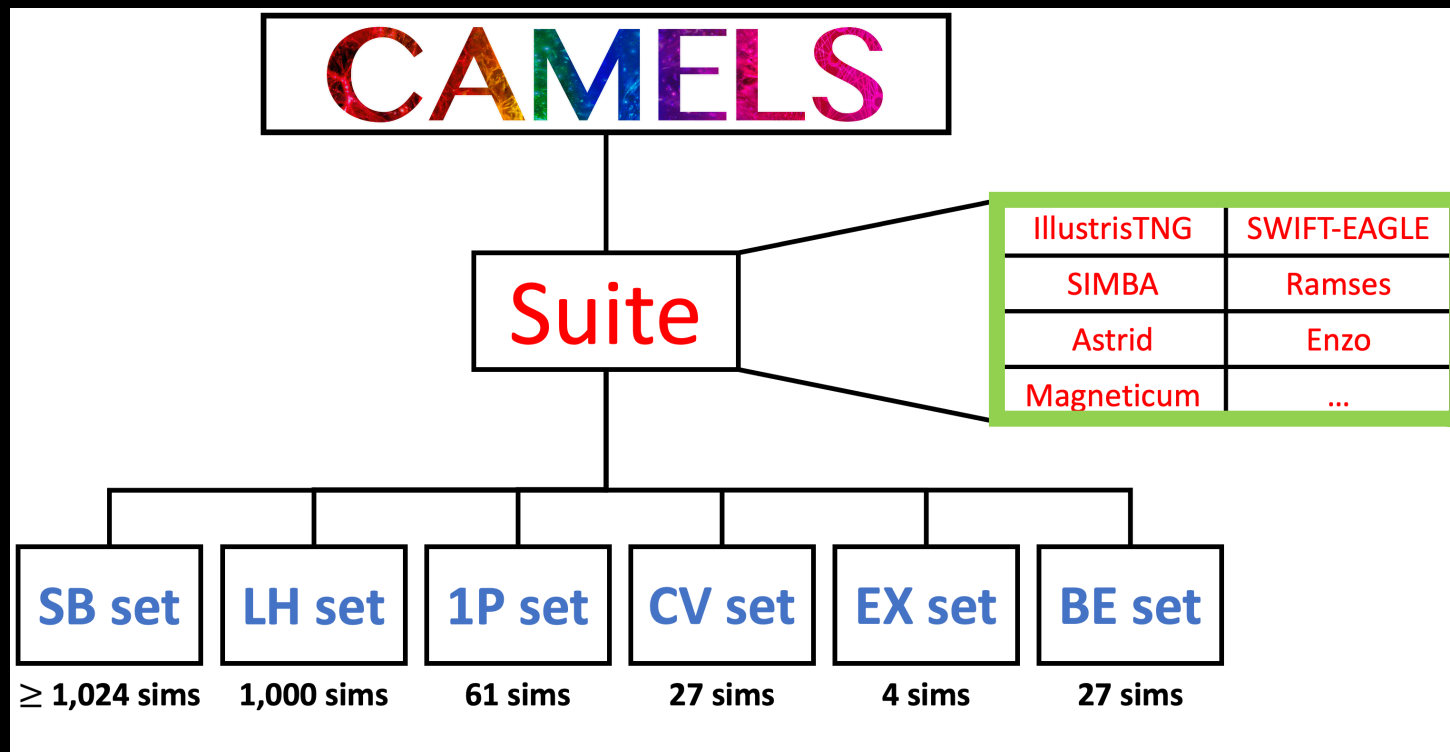
Differences between sub-grid models limit learning across them



# Crucial to expand the range of models in the training set

The CAMELS project: Expanding the galaxy formation model space  
with new ASTRID and 28-parameter TNG and SIMBA suites

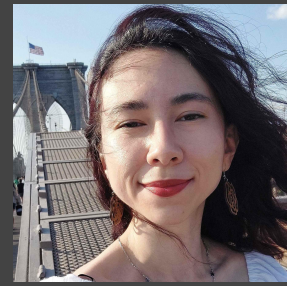
YUEYING NI,<sup>1,2</sup> SHY GENEL,<sup>3,4</sup> DANIEL ANGLÉS-ALCÁZAR,<sup>5,3</sup> FRANCISCO VILLAESCUSA-NAVARRO,<sup>3,6</sup> YONGSEOK JO,<sup>3</sup>  
SIMEON BIRD,<sup>7</sup> TIZIANA DI MATTEO,<sup>2,8</sup> RUPERT CROFT,<sup>2,8</sup> NIANYI CHEN,<sup>2</sup> NATALÍ S. M. DE SANTI,<sup>3,9</sup>  
MATTHEW GEBHARDT,<sup>5</sup> HELEN SHAO,<sup>6</sup> SHIVAM PANDEY,<sup>10,11</sup> LARS HERNQUIST,<sup>1</sup> AND ROMEEL DAVE<sup>12</sup>



“Building robust machine-learning models favors training and testing on the largest possible diversity of galaxy formation models”

Reach out if you would like to bring your model into the public CAMELS dataset!

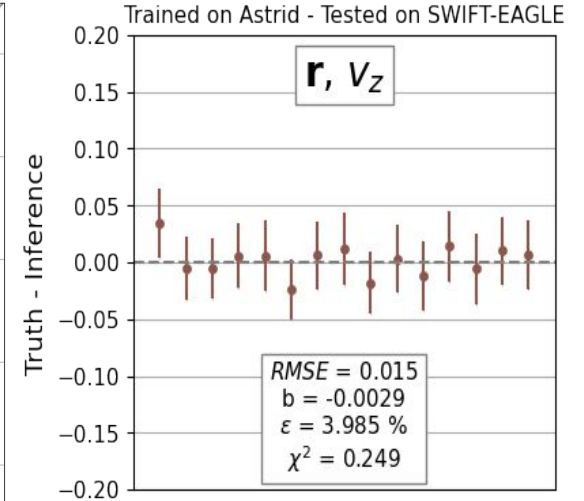
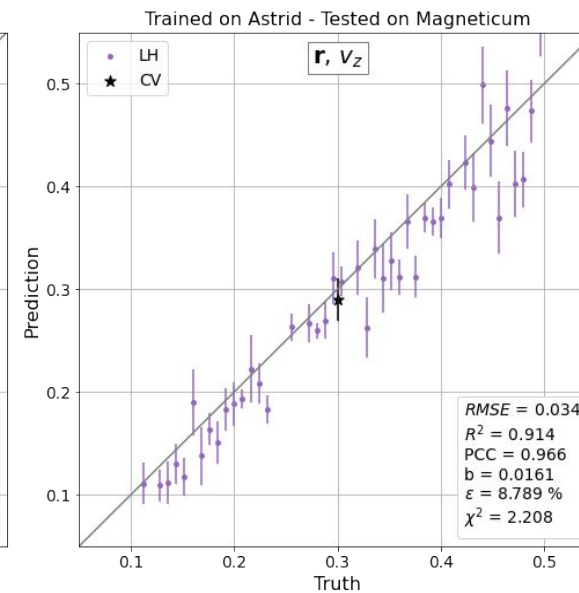
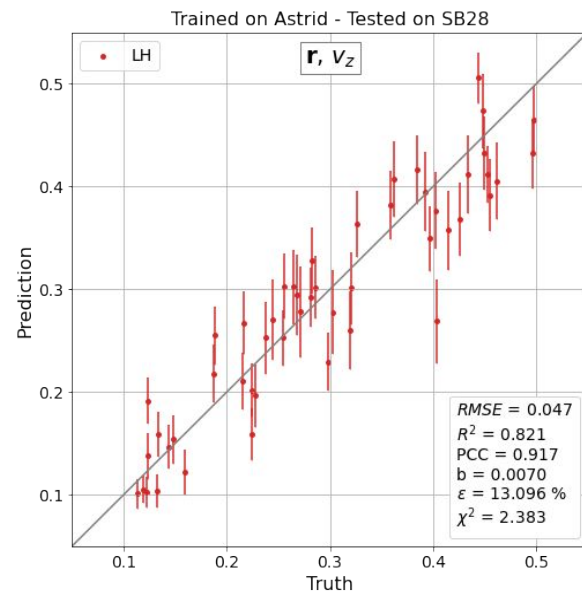
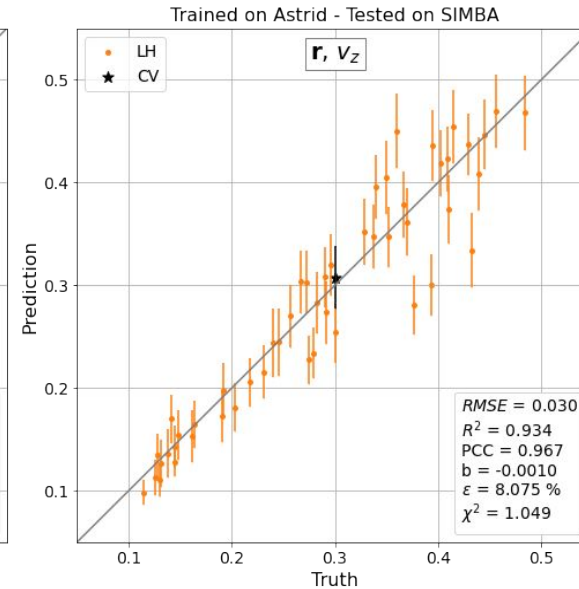
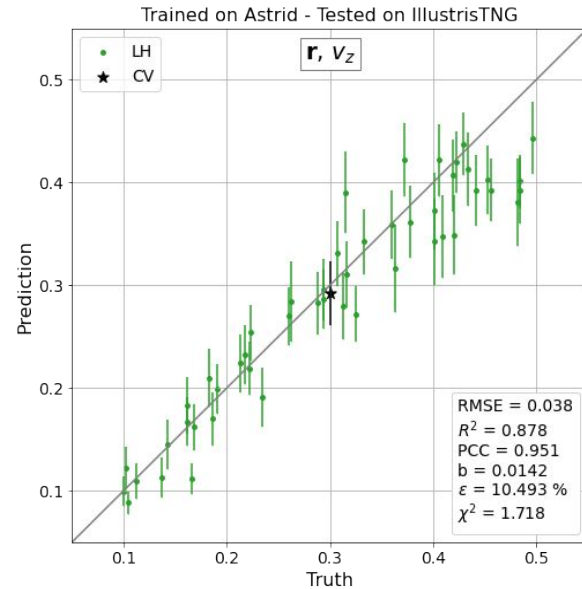
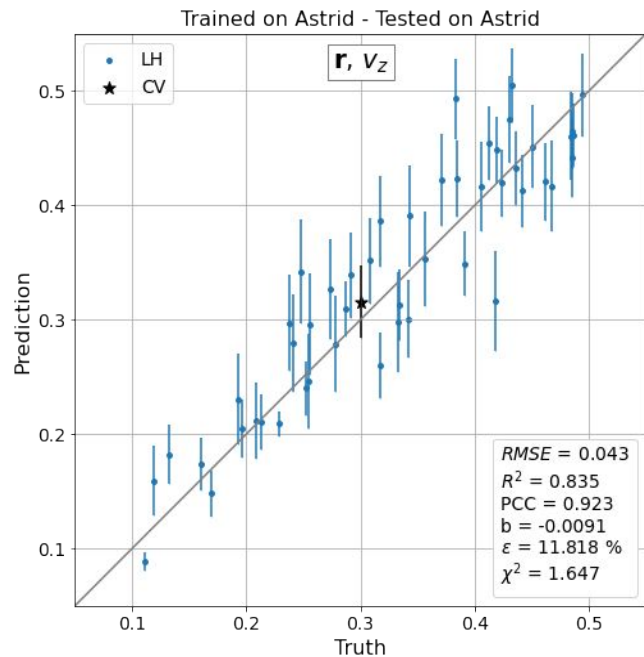




Natalí de Santi  
 Flatiron Institute  
 University of São Paulo  
 natalidesanti@gmail.com

# Robust field-level likelihood-free inference with galaxies





**Dataset:** Galaxies from Astrid  
**Machine Learning Method:** Graph Neural Networks  
**Objective:**  $\Omega_m$  inference



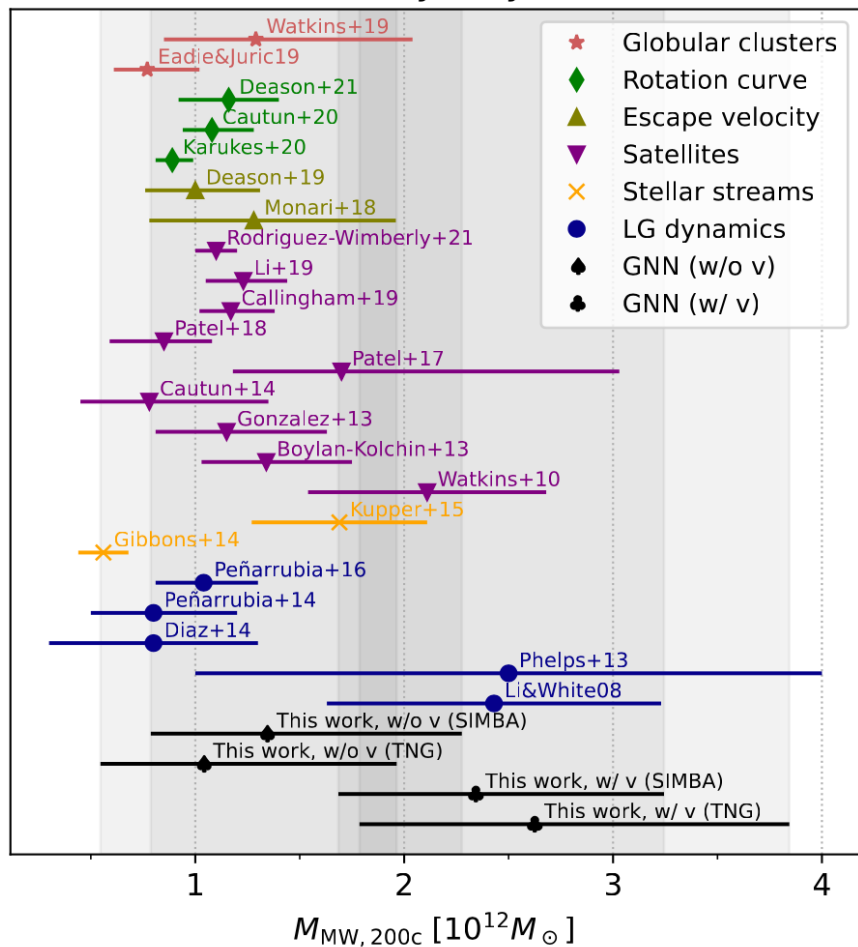
[arXiv: 2302.14101](https://arxiv.org/abs/2302.14101)

- Information came from galaxy positions and velocities;
- The broader variation in Astrid allowed a robust model across 5 different sub-grid physics sets;
- First steps to apply this machinery on real data.

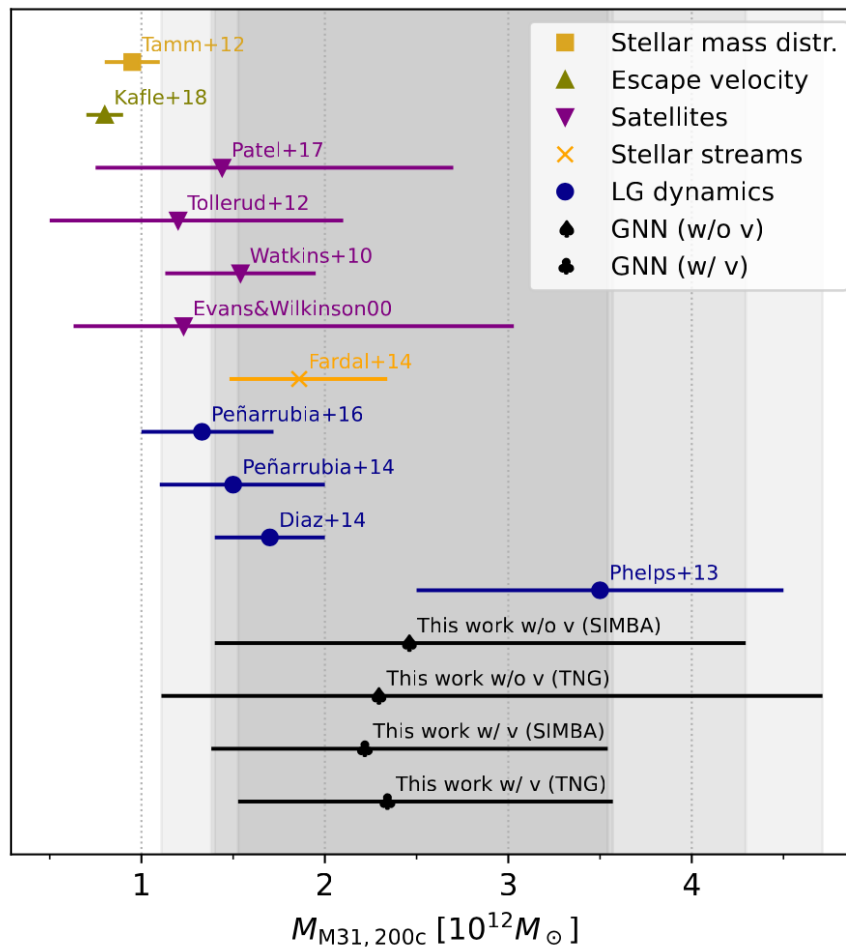
# Weighing the Milky Way and Andromeda with Artificial Intelligence

Pablo Villanueva-Domingo <sup>1,\*</sup> Francisco Villaescusa-Navarro <sup>2,3,†</sup> Shy Genel <sup>2,4</sup> Daniel Anglés-Alcázar <sup>5,2</sup>  
 Lars Hernquist,<sup>6</sup> Federico Marinacci,<sup>7</sup> David N. Spergel,<sup>2,3</sup> Mark Vogelsberger,<sup>8</sup> and Desika Narayanan<sup>9,10</sup>

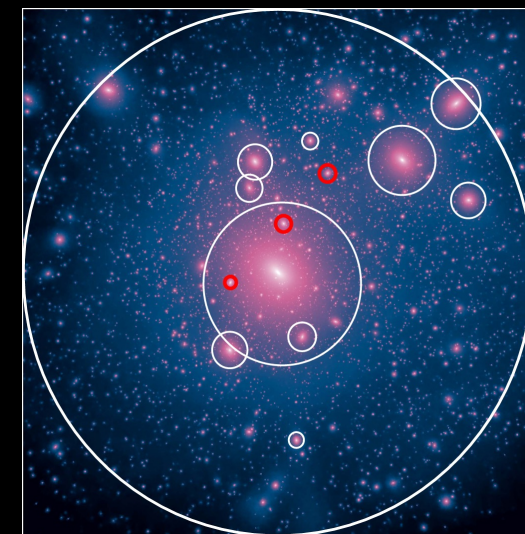
## Milky Way



## Andromeda



Graph Neural Networks trained on positions, velocities, and stellar masses of galaxies to predict halo mass





# THE NEW YORKER

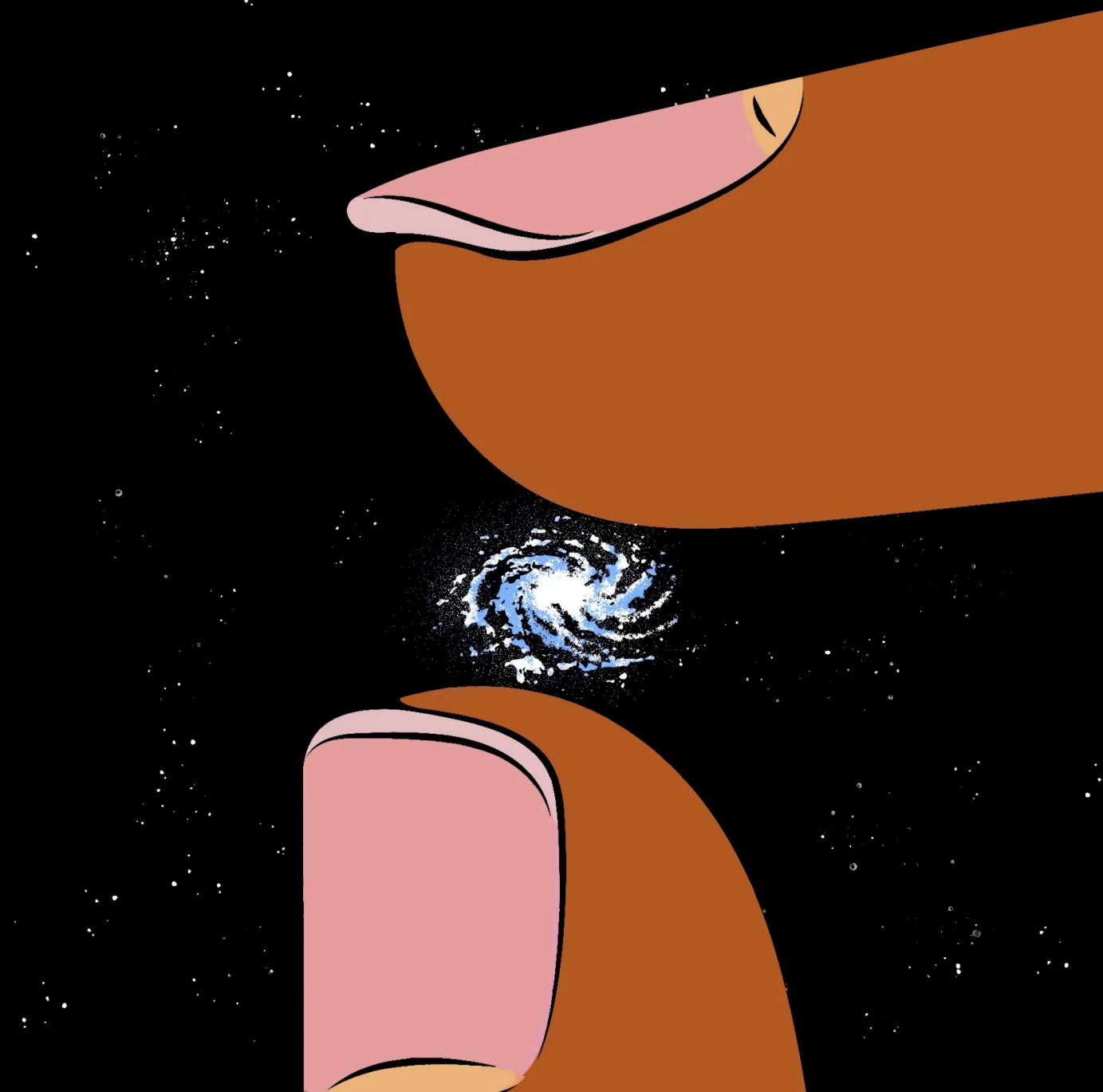
---

ELEMENTS

## WHAT CAN WE LEARN ABOUT THE UNIVERSE FROM JUST ONE GALAXY?

*In new research, begun by an undergraduate, William Blake's phrase "to see a world in a grain of sand" is suddenly relevant to astrophysics.*

By Rivka Galchen  
March 23, 2022

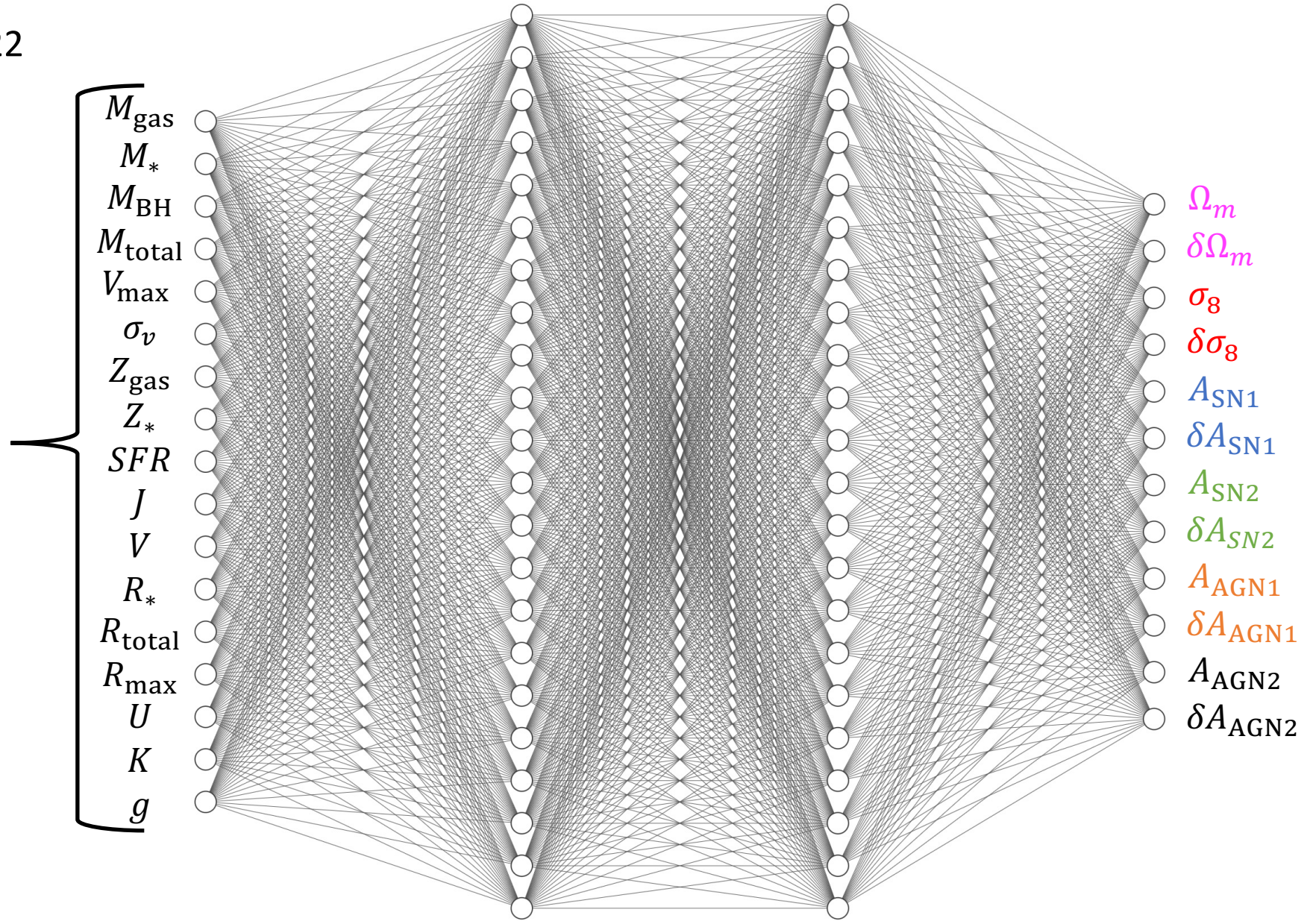


# CAMELS enables testing new ideas: **Cosmology with a single galaxy?**

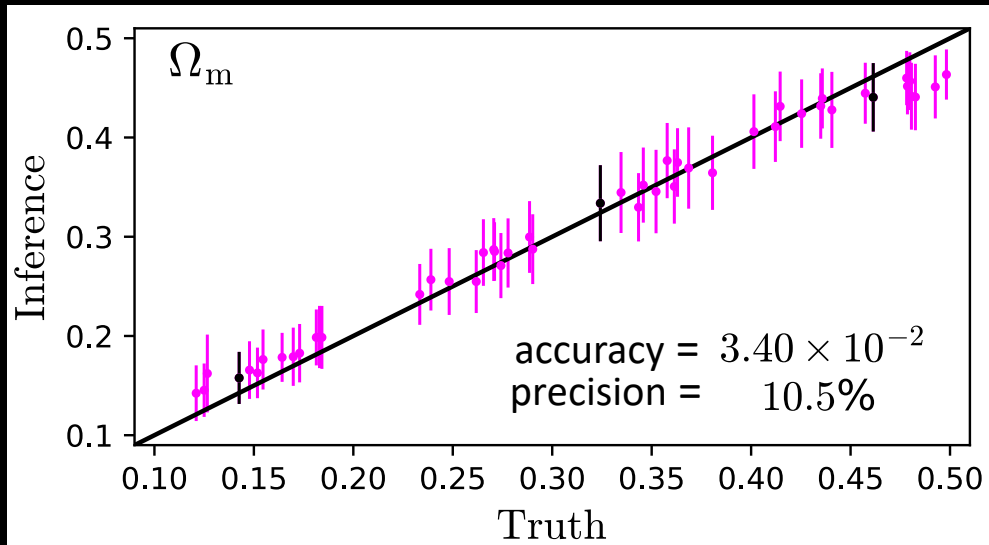
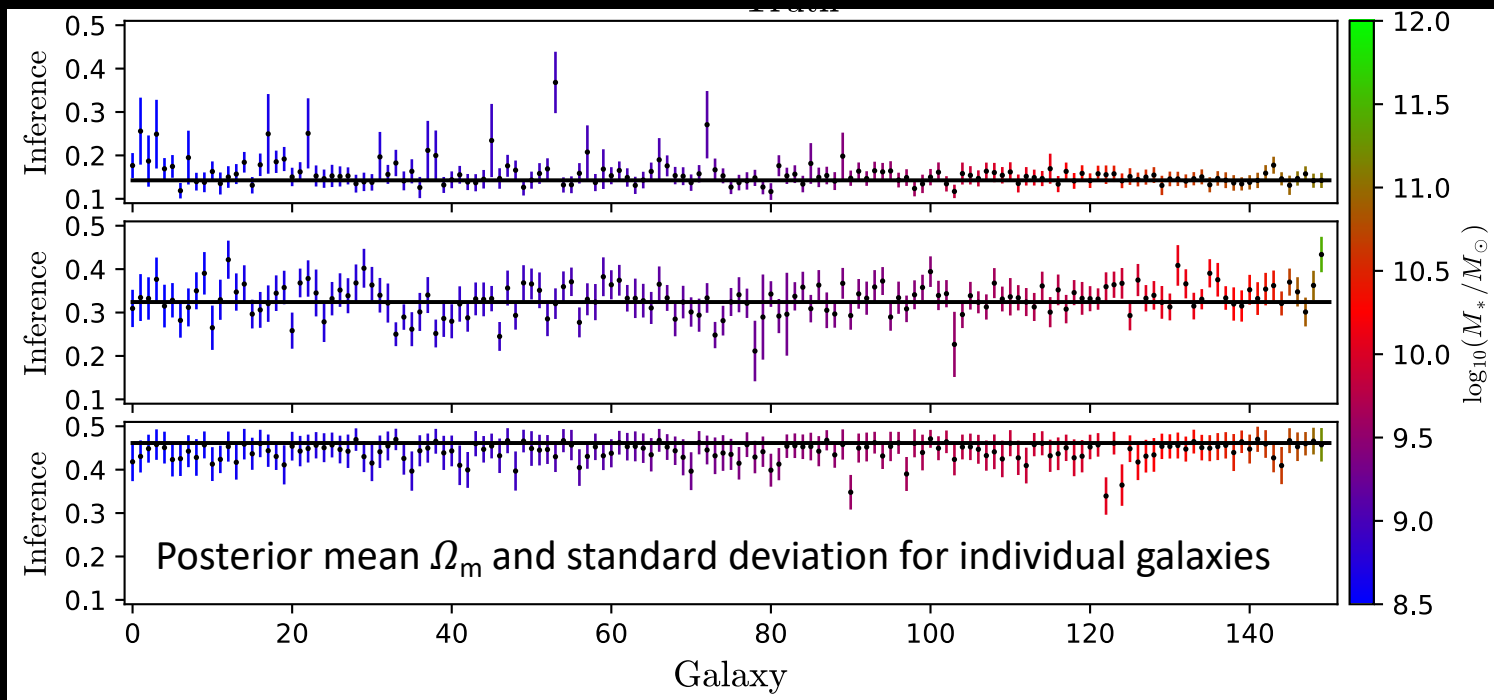
Villaescusa-Navarro+2022

Echeverri-Rojas+2023

Chawak+2023



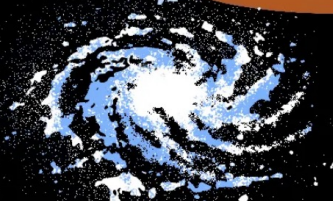
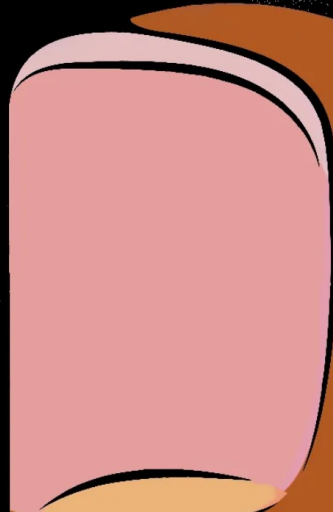




Villaescusa-Navarro+2022  
Echeverri-Rojas+2023  
Chawak+2023.



Most important to measure  $\Omega_m$   
maximum circular velocity  
stellar mass  
stellar metallicity  
Stellar effective radius



# The CAMELS project:

## Cosmology and Astrophysics with Machine Learning Simulations

Villaescusa-Navarro, Anglés-Alcázar, Genel, et al. (2021,2022,2023)

Perez+2023 (arXiv:2204.02408), Ni+2023 (arXiv:2304.02096)

- Largest suite of cosmological hydrodynamic simulations with thousands of model variations designed for machine learning applications
- Encouraging results extracting cosmological information at the field level down to small scales even where astrophysical effects are significant
- Many possible applications in galaxy formation and cosmology (inference, emulating/accelerating simulations, learning physics with AI,...)
- Full dataset publicly available: <https://camels.readthedocs.io>
- Challenges: larger-volume simulations, extending parameter space, interpolation between models, robustness, synthetic observations...

