

DRAFT VERSION NOVEMBER 17, 2021
Typeset using L^AT_EX twocolumn style in AASTeX631

Inferring halo masses with Graph Neural Networks

PABLO VILLANUEVA-DOMINGO ¹, FRANCISCO VILLAESCUSA-NAVARRO ^{2,3}, DANIEL ANGLÉS-ALCÁZAR ^{4,2},
SHY GENEL ^{2,5}, FEDERICO MARINACCI,⁶ DAVID N. SPERGEL,^{2,3} LARS HERNQUIST,⁷ MARK VOGELSBERGER,⁸
ROMEEL DAVE,^{9,10,11} AND DESIKA NARAYANAN^{12,13}

¹*Instituto de Física Corpuscular (IFIC), CSIC-Universitat de València, E-46980, Paterna, Spain*

²*Center for Computational Astrophysics, Flatiron Institute, 162 5th Avenue, New York, NY, 10010, USA*

³*Department of Astrophysical Sciences, Princeton University, Peyton Hall, Princeton NJ 08544, USA*

⁴*Department of Physics, University of Connecticut, 196 Auditorium Road, U-3046, Storrs, CT 06269-3046, USA*

⁵*Columbia Astrophysics Laboratory, Columbia University, New York, NY, 10027, USA*

⁶*Dipartimento di Fisica e Astronomia ‘Augusto Righi’, Università di Bologna, via Gobetti 93/2, 40129, Bologna, Italy*

⁷*Center for Astrophysics — Harvard & Smithsonian, 60 Garden St, Cambridge, MA 02138, USA*

⁸*Kavli Institute for Astrophysics and Space Research, Department of Physics, MIT, Cambridge, MA 02139, USA*

⁹*Institute for Astronomy, University of Edinburgh, Royal Observatory, Edinburgh EH9 3HJ, UK*


¹⁰*Department of Physics & Astronomy, University of the Western Cape, Cape Town 7535, South Africa*

¹¹*South African Astronomical Observatories, Observatory, Cape Town 7925, South Africa*

¹²*Department of Astronomy, University of Florida, Gainesville, FL, USA*

¹³*University of Florida Informatics Institute, 432 Newell Drive, CISE Bldg E251, Gainesville, FL, USA*

ABSTRACT

Understanding the halo-galaxy connection is fundamental in order to improve our knowledge on the nature and properties of dark matter. In this work we build a model that infers the mass of a halo given the positions, velocities, stellar masses, and radii of the galaxies it hosts. In order to capture information from correlations among galaxy properties and their phase-space, we use Graph Neural Networks (GNNs), that are designed to work with irregular and sparse data. We train our models on galaxies from more than 2,000 state-of-the-art simulations from the Cosmology and Astrophysics with Machine Learning Simulations (CAMELS) project. Our model, that accounts for cosmological and astrophysical uncertainties, is able to constrain the masses of the halos with a ~ 0.2 dex accuracy. Furthermore, a GNN trained on a suite of simulations is able to preserve part of its accuracy when tested on simulations run with a different code that utilizes a distinct subgrid physics model, showing the robustness of our method. The PyTorch Geometric implementation of the GNN is publicly available on [GitHub](#) .

1. INTRODUCTION

In 1933 Fritz Zwicky found out that the mass of the Coma cluster should be much larger than the one from its luminous component (Zwicky 1933). That finding pointed out to the existence of an unknown type of non-luminous matter: dark matter. The requirement of unseen matter was later supported by the observation of rotation curves of galaxies (Rubin et al. 1978; Bosma 1978). Nowadays, there is overwhelming evidence for

the existence of dark matter, although we are still ignorant of its fundamental properties (Young 2017).

We now believe that dark matter is the backbone of the distribution of matter in the Universe: it concentrates in high-density regions called halos, that are connected by thin filaments with intermediate density and surrounded by gigantic regions with low density (voids). It is within halos that gas can cluster, cool, and form stars and galaxies (Somerville & Davé 2015). Dark matter halos are therefore the environment where galaxies reside.

Understanding the halo-galaxy connection is fundamental to improve our knowledge on the nature and properties of dark matter. There are two possible directions to take in the halo-galaxy connection. On one

pablo.villanueva.domingo@gmail.com

fvillaescusa@flatironinstitute.edu

hand, given a dark matter halo and its properties and environment, predict the number and distribution of galaxies it hosts. This task is fundamental in order to create galaxy mocks with the correct clustering on all scales needed for forward modeling approaches (Wechsler & Tinker 2018). On the other hand, given a set of galaxies, it may be useful to determine some properties of their host halo such as its mass, spin, and concentration. This task is fundamental to derive cosmological constraints from the abundance of dark matter halos.

There has been an extensive program of weighing the masses of halos and clusters, using a wealth of techniques including gravitational lensing (Mandelbaum 2015; Huang et al. 2020, 2021), rotation curves of galaxies (Sofue & Rubin 2001; Sofue 2015), abundance matching (Behroozi et al. 2010), Sunyaev-Zeldovich effect (Grego et al. 2001), X-rays observations (Landry et al. 2013), velocity dispersion (Saro et al. 2013) and kinematics of satellite galaxies (Wojtak & Mamon 2013; Seo et al. 2020) among others; see, e.g. Old et al. (2015) for a comparison of different techniques for galaxy clusters.

All the above techniques do not make use of all available information. For instance, the abundance matching technique only considers the total stellar mass in the system, disregarding information about its clustering state. In this work we attempt to build a model that can use all available information from observations and/or simulations, e.g. phase-space information, stellar masses, galaxy sizes, to infer the mass of the halo hosting the galaxies. For this, we made use of neural networks and their capacity as universal function approximators.

Different machine learning (ML) algorithms have been already used to perform this task. For instance, Convolutional Neural Networks (CNNs) have been already applied in order to predict dynamical galaxy cluster masses (Ntampaka et al. 2019; Ho et al. 2019; Kodi Ramanah et al. 2020, 2021; Yan et al. 2020; Gupta & Reichardt 2020; de Andres et al. 2021), as well as other ML techniques (Ntampaka et al. 2015, 2016; Green et al. 2019; Armitage et al. 2019; Haider Abbas 2019). Subhalo masses can be inferred from different subhalo properties via multilayer perceptrons (Shao et al. 2021) or different ML algorithms (von Marttens et al. 2021). Other works have been employed different ML algorithms to predict the mass of a halo from the properties of the halo and galactic group (Man et al. 2019; Calderon & Berlind 2019; Lucie-Smith et al. 2020). These ML-based approaches have been shown to outperform other traditional techniques to infer halo masses. However, some of these works make use of N-body computations plus semianalytical galaxy formation models, rather than more accurate hydrodynamical simulations. Moreover,

several of the features employed may not be easily observable, which could complicate their applicability to real data.

Although these approaches make use of global properties of a halo as well as individual features of the galaxies, they do not incorporate explicitly the relationship between galaxies or subhalos, either in the form of clustering in configuration space and/or distribution in phase-space. In this article, we aim at predicting halo masses exploiting the halo-galaxy connection, using a novel method based on *Graph Neural Networks* (GNNs). This type of neural network shares the typical training procedure of other deep learning techniques, but it is applied to data structured in the form of mathematical graphs. To understand their significance, it is useful to compare them to other deep learning frameworks. CNNs are mostly employed with regular data (grids), such as images and 3D grids; CNNs automatically account for translational invariance. Recurrent neural networks, on their side, are designed to treat sequential data, such as chains of characters in natural language or time series. However, GNNs can be applied when dealing with irregular data, where data points may have arbitrary relations.¹ That is the case of point clouds, as a galaxy catalog can be regarded. GNNs have been successfully employed in different fields such as chemistry, computer vision, natural language processing, social networks or particle physics (Wu et al. 2019; Bronstein et al. 2021). There are already some applications of GNNs in cosmology, for instance in order to perform symbolic regression (Cranmer et al. 2019; Cranmer et al. 2020), to predict the redshift of galaxies (Beck & Sadowski 2019), or to allocate resources in an unsupervised way in order to select galaxies (Cranmer et al. 2021). GNNs have also been applied in other physics fields, such as particle physics (Shlomi et al. 2020), but still represent a novel, promising and mostly unexplored way to extract information from irregular data.

We train our GNNs using galaxies from simulations of the Cosmology and Astrophysics with Machine Learning Simulations (CAMELS) project (Villaescusa-Navarro et al. 2021a) to extract information from galaxy properties and their phase-space distribution. Since CAMELS contains thousands of state-of-the-art (magneto-)hydrodynamic simulations with different values of the cosmological and astrophysical parameters, our method accounts for uncertainties in cosmology and astrophysics. Furthermore, since CAMELS contains

¹ See Battaglia et al. (2018) for a comparison of the different deep learning components and their relational inductive biases.

two different suites of hydrodynamic simulations run with two different codes that employ different subgrid physics, we can quantify the robustness of our results to astrophysical uncertainties. Ultimately, we would like to build a tool that can infer the mass of galaxy systems, like our own Milky Way and Andromeda, just from the observed properties of those galaxies and their satellites. Knowing the mass of the host dark matter halos of those system will allow us to make consistency checks within the Λ CDM model.

The article is structured as follows. We start by reviewing the basics on graphs and GNNs in Sec. 2. In Sec. 3 we describe the data we use to train our model together with an outline of the training procedure. The main results of this work are presented in Sec. 4. Some aspects of the interpretability of the GNNs are examined in Sec. 5, followed by a discussion of the main conclusions in Sec. 6.

2. GRAPH NEURAL NETWORKS

In this section we review the basics of graph neural networks, firstly summarizing the fundamentals of graphs, subsequently detailing how to build graphs from the galaxies of halos, and finally introducing how to build a neural network on a graph based on the message passing scheme. We refer the reader to Bronstein et al. (2021); Battaglia et al. (2018); Hamilton (2020) for comprehensive references on graph neural networks and geometric deep learning.

2.1. General concepts on graphs

We start by discussing some generic concepts of graphs and standard definitions, since some astrophysicists, cosmologists and, ML practitioners may not be familiar with the terminology. A *graph* can be defined as a tuple $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} denotes the set of *nodes*, and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ the set of *edges*. For two nodes of the graph $i, j \in \mathcal{V}$, there is an edge connecting them if $(i, j) \in \mathcal{E}$. Nodes i, j are thus coined as *neighbors*. The connectivity can be described by the *adjacency matrix* A_{ij} , which takes value of 1 if the pair (i, j) is connected by an edge, being 0 otherwise. A graph is *undirected* when, if $(i, j) \in \mathcal{E}$, then $(j, i) \in \mathcal{E}$ too. A *loop* is an edge which connects a node to itself. We restrict our discussion to *simple* graphs, i.e., undirected graphs without loops, since these are enough for our purposes (although self-loops can be implicitly employed in some GNN architectures).

Given a node i , we denote its *feature vector* as $\mathbf{x}_i \in \mathbb{R}^{n_{\text{in}}}$, encoding the relevant physical information about the node, with n_{in} the number features. The *feature matrix* $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_{|\mathcal{V}|})^T \in \mathbb{R}^{|\mathcal{V}| \times n_{\text{in}}}$, where $|\mathcal{V}|$ is the

Symbol	Feature
\mathbf{p}	3D comoving position
v	Modulus of relative velocity
M_*	Stellar mass
R_*	Stellar half-mass radius

Table 1. Summary of galactic properties employed to train the GNNs.

number of nodes in the graph, comprises the feature vectors of all nodes. The *neighborhood* of the node $i \in \mathcal{V}$, denoted by \mathcal{N}_i , includes every node which shares and edge with i , and can thus be written as

$$\mathcal{N}_i = \{j \mid A_{ij} = 1\}. \quad (1)$$

Note that \mathcal{N}_i does not include the node i , since we dismiss loops. Analogously, we can define the set of feature vectors of the neighbors of the node i as

$$\mathcal{X}_i = \{\mathbf{x}_j \mid j \in \mathcal{N}_i\}. \quad (2)$$

A graph is *complete* if every pair of nodes is connected by an edge, having thus $\binom{|\mathcal{V}|}{2} = |\mathcal{V}|(|\mathcal{V}| - 1)/2$ edges in total.² The edges set can hence be written as $\mathcal{E} = \mathcal{V} \times \mathcal{V} \setminus \{(i, i) \in \mathcal{V}\}$ (excluding loops). As will be shown later, most of the halos employed here lead to complete graphs.³

2.2. Halos as graphs

The general terminology discussed above allows us to set up our problem. Consider a halo h with mass M_h which hosts several galaxies. We build a graph, denoted by \mathcal{G}_h , by considering all the (central and satellite) galaxies of the halo as the nodes of the graph. The feature vector of a node i , \mathbf{x}_i , will include information about the corresponding galaxy, namely the 3D comoving position, \mathbf{p} , the stellar mass, M_* , the modulus of the velocity, v , and the stellar half-mass radius (i.e., the comoving radius containing half of the stellar mass in the galaxy), R_* . These galactic properties are listed on Table 1. Positions are expressed with respect to the center of the halo (chosen as the point with minimum

² This can be reasoned as the number of pairs that can be chosen from a total of $|\mathcal{V}|$ elements.

³ Complete graphs can be very relevant in ML applications, since the popular architecture Transformers (Vaswani et al. 2017) can be regarded as a particular case of a GNN, a Graph Attention Network (Velićković et al. 2018), where the graph is complete (see Bronstein et al. 2021 for more details).

where \mathbf{h}_i is the output feature vector, \oplus denote a differentiable, permutation invariant *aggregation function*, such as the maximum, the mean or the sum, while ψ is the *message function*, a differentiable function, like a Multi Layer Perceptron (MLP), which depends upon the feature vector of the node, \mathbf{x}_i , as well as on those from its neighbors, \mathbf{x}_j , for $j \in \mathcal{N}_i$. The distinct ways to build ψ lead to different graph layers.⁶

The function ψ maps $\psi : \mathbb{R}^{n_{\text{in}}} \rightarrow \mathbb{R}^{n_{\text{out}}}$, with n_{in} and n_{out} the number of initial and output features, respectively. As shown in Sec. 2.4, n_{in} may not correspond to the number of features considered of each node. A common advantage of GNNs is that they exploit locality of data, since closer nodes may present stronger relations than those farther away. Here, locality is enforced by assuming that the functions ψ are the same for all nodes and graphs.⁷

The final GNN incorporates a last step to infer the global graph target, and thus can be of the form

$$\mathbf{y} = \phi \left(\bigoplus_{i \in \mathcal{G}_h} \mathbf{h}_i \right), \quad (4)$$

where ϕ is a differentiable function, chosen as a MLP with three hidden layers with 300 channels separated by ReLU activation functions, and \mathbf{y} is the output of the network. In our case, since we perform likelihood-free inference to estimate the mean and standard deviation of the posterior of the mass of the halo, the target is given by a vector of two components $\mathbf{y} = [y, \sigma]$, containing the mean prediction y and its expected standard deviation, σ . The targeted quantity is the logarithm of the mass of the halo, $y = \log_{10} [M_h / (M_\odot / h)]$. The likelihood-free approach to estimate the posterior mean and standard deviation is detailed in Sec. 3.2. The equation above can be easily modified to include some global graph quantities u , which may help to train the network, such as the number of nodes or the total stellar mass, writing $\mathbf{y} = \phi ([\bigoplus_{i \in \mathcal{G}_h} \mathbf{h}_i, u])$. We include these global features in our graphs, although we have checked that removing them does not leave a significant impact in the results. Figure 1 shows a sketch of this basic architecture, illustrating the message passing scheme.

Permutation equivariance and invariance are key concepts in GNNs. We say that a function acting on the

node feature matrix $f(\mathbf{X})$ is permutation invariant if for every permutation matrix \mathbf{P} , one has $f(\mathbf{P}\mathbf{X}) = f(\mathbf{X})$, leaving thus unchanged the output. On the other hand, f is permutation equivariant if it transforms as $f(\mathbf{P}\mathbf{X}) = \mathbf{P}f(\mathbf{X})$.⁸ A message passing layer $\text{GL}(\mathbf{x}_i, \mathcal{X}_i)$ should be permutation equivariant, since a reordering of the input nodes produces the same permutation in the outputs, although the output feature space can be different. However, any global quantity of the graph such as the final output of the network, the halo mass, must be permutation invariant, since its value should not depend on the ordering of the nodes. The final aggregation step, $\bigoplus_{i \in \mathcal{G}_h} \text{GL}(\mathbf{x}_i, \mathcal{X}_i)$, can be regarded as a function $f(\mathbf{X})$ which fulfills by construction permutation invariance, $f : \mathbb{R}^{|\mathcal{V}| \times n_{\text{out}}} \rightarrow \mathbb{R}^{n_{\text{out}}}$. In this case, the aggregation \bigoplus may be termed as a global pooling layer, since it reduces dimensionality. Note that the aggregation method here does not have to be the same as in the message passing layer. Furthermore, it is possible to employ different aggregation operators, such as maximum, sum, and mean, and concatenate them, which is the approach used in this work. The last MLP ϕ makes use of n_{out} global features of the graph in order to extract a global property vector \mathbf{y} , the halo mass and its expected standard deviation in our case, $\phi : \mathbb{R}^{n_{\text{out}}} \rightarrow \mathbb{R}^2$. Finally, we can apply multiple successive GNN layers, $\text{GL}(\text{GL}(\mathbf{x}_i, \mathcal{X}_i), \mathcal{H}_i)$, where $\mathcal{H}_i = \{\mathbf{h}_j \mid j \in \mathcal{N}_i\}$ denotes the updated feature vector of the neighbors. In such a case, the node after k updates will encode information from the nodes of its k -hop neighborhood. We take the number of message passing layers as an optimizable hyperparameter, although a single GNN layer is sufficient to provide the best results, as shown in Sec. 3.2.

2.4. Graph layer architecture

So far we have discussed the general shape of the GNN, but the explicit design of the message passing layer remains undetermined. There are many possible ways to construct the specific architecture of the GNN layer, which is defined by the message function ψ and the aggregation scheme \bigoplus . For instance, one could not take into account neighbors to update each node, and employ only the information of the node \mathbf{x}_i to extract the hidden information \mathbf{h}_i . These types of architectures, a subset of GNNs, are known as DeepSets (Zaheer et al. 2017). No aggregation function is thus needed to update the node features, and the hidden layer can just be written as $\mathbf{h}_i = \psi(\mathbf{x}_i)$.

⁶ One could also apply an additional MLP ξ to \mathbf{h}_i and its original feature vector, $\mathbf{h}'_i = \xi([\mathbf{x}_i, \mathbf{h}_i])$, where the square brackets indicate array concatenation, although we dismiss this step for the sake of simplicity, and since we apply several successive graph layers, which leads to a similar effect.

⁷ Similarly, CNNs in, e.g., computer vision tasks guarantee locality by employing the same kernels over all different images.

⁸ See Bronstein et al. (2021) for more details on invariances on graphs.

However, to fully exploit the relations among the nodes, it is useful to incorporate neighborhood and edge information actually performing message passing. Here we assume an architecture based on the edge convolutional layer, coined as EdgeNet (Wang et al. 2019), where for each neighbor j of the node i , the relative vectors $\mathbf{x}_i - \mathbf{x}_j$ are concatenated to the feature vector \mathbf{x}_j . The aggregation function is chosen as the maximum, writing thus the hidden layer as

$$\mathbf{h}_i = \max_{j \in \mathcal{N}_i} \psi([\mathbf{x}_i, \mathbf{x}_i - \mathbf{x}_j]), \quad (5)$$

where the square brackets indicate array concatenation. In this case, the initial input number is given by $n_{\text{in}} = 2n_{\text{feat}}$, where n_{feat} is the number of features. The differentiable function ψ is taken as a MLP with three hidden layers, with 300, 300 and 100 hidden channels, separated by ReLu activation functions. We have checked that other choices of the number of channels do not improve the predictions of the net.

There are other popular architectures, such as PointNet (Qi et al. 2017a,b), which employs only the relative spatial positions for the message passing, $\mathbf{p}_i - \mathbf{p}_j$, rather than the full feature vector, or the Graph Convolutional Network (Kipf & Welling 2017), where the neighbor information is simply incorporated by taking a linear combination of the features and summing over all the neighbors, imitating a convolution operation. We have checked that the EdgeNet (Eq. 5) outperforms the other architectures mentioned above via a hyperparameter optimization procedure, as commented in Sec. 3.2. See, e.g., Wu et al. (2019); Zhou et al. (2018) for a discussion of other different GNN architectures.⁹

3. METHODS

In this section, we specify the details regarding the data employed and the training of the network.

3.1. The CAMELS simulations

The CAMELS project (Villaescusa-Navarro et al. 2021a) comprises a set of state-of-the-art hydrodynamic and N-body simulations, specially suited and designed to train and test ML algorithms. They include thousands of realizations varying two cosmological parameters, namely the matter density parameter Ω_m , and the variance of the linear field on 8 Mpc/ h at $z = 0$, σ_8 , plus four astrophysical parameters controlling the efficiency of supernovae and active galactic nuclei (AGN) feedback. The rest of cosmological parameters are kept

fixed to the values: $\Omega_b = 0.049$, $h = 0.6711$, $n_s = 0.9624$, $w = -1$, $M_\nu = 0$ eV. Each simulation follows the evolution from $z = 127$ to $z = 0$ of 256^3 DM particles and 256^3 gas resolution elements within a box of size periodic volume of size 25 Mpc/ h . The DM particle mass resolution is $\sim 10^8 M_\odot/h$. A collection of 2D maps and 3D grids from the CAMELS simulations have been recently released publicly available as the CAMELS Multifield Dataset (CMD)¹⁰, intended to be a standard cosmological dataset for ML applications (Villaescusa-Navarro et al. 2021b).

The CAMELS project includes two suites of simulations, with different astrophysics modeling and subgrid physics. On the one hand, a suite of magneto-hydrodynamic simulations with the subgrid physics models of IllustrisTNG (Weinberger et al. 2017; Pillepich et al. 2018; Nelson et al. 2019), performed with the code Arepo¹¹ (Weinberger et al. 2020, see also Springel et al. 2018; Marinacci et al. 2018; Naiman et al. 2018 for more details). Its galaxy formation model is based on the predecessor Illustris (Vogelsberger et al. 2013, 2014). On the other hand, another suite employs the SIMBA subgrid physics model (Davé et al. 2019), making use of the code GIZMO¹² (Hopkins 2015). SIMBA is built on its precursor MUFASA (Davé et al. 2016) with the addition of supermassive black hole growth and feedback (Anglés-Alcázar et al. 2017). The hydrodynamics simulations are accompanied by N-body counterparts, run with the GADGET-III code (Springel 2005).¹³

Within each simulation suite, there are different sets, according to the configuration of astrophysical and cosmological parameters. The CV simulations (standing for Cosmic Variance) include 27 simulations sharing their astrophysical and cosmological parameters, fixed to standard values ($\Omega_m = 0.3$ and $\sigma_8 = 0.8$), and varying only the random seed to generate the initial conditions. The LH set (standing for Latin-Hypercube) consists of 1,000 simulations, varying all the astrophysical and cosmological parameters (together with the random seed) along a latin-hypercube.¹⁴ We make use of both sets for training our networks, in order to check whether the architectures considered are robust enough for different cosmologies. The halos and subhalos are identified us-

¹⁰ <https://camels-multifield-dataset.readthedocs.io>

¹¹ <https://arepo-code.org/>

¹² <http://www.tapir.caltech.edu/~phopkins/Site/GIZMO.html>

¹³ See, e.g., Vogelsberger et al. (2020) for a comparison of different models of galaxy formation in cosmological simulations.

¹⁴ Besides CV and LH, other simulation sets are also included in CAMELS, but not employed in this work.

⁹ See also <https://github.com/thunlp/GNNPapers> for a comprehensive list of GNN references.

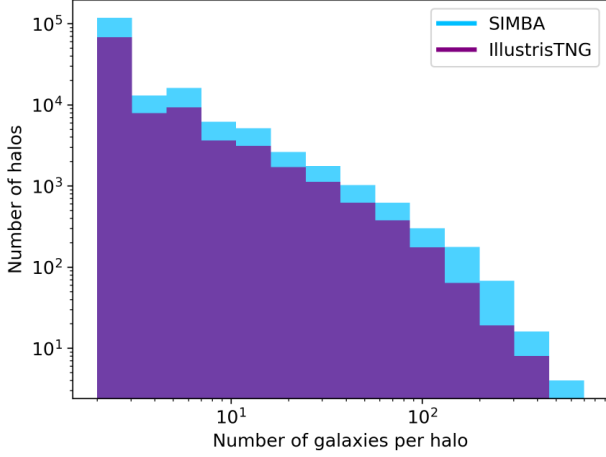


Figure 2. Number of halos sorted by the number of galaxies that they host in the IllustrisTNG (blue) and SIMBA (purple) suites, for the LH sets. On average, the halos of the SIMBA simulations contain more galaxies than the halos of the IllustrisTNG simulations, reflecting the large differences between the two simulation suites given their distinct subgrid physics models.

ing the SUBFIND algorithm (Springel et al. 2001). We define galaxies as subhalos that contain more than 20 star particles. Although the CAMELS suite contains data for several redshifts, only simulations at $z = 0$ are considered in this work. We refer the reader to the CAMELS webpage¹⁵ and Villaescusa-Navarro et al. (2021a) for further details.

As an example of the differences between IllustrisTNG and SIMBA, Fig. 2 shows the number of halos as a function of the number of galaxies the halos contain for the simulations of the LH set. Note that for a fixed halo mass, the SIMBA simulations contain more galaxies than the IllustrisTNG simulations. This reflects the inherent differences between the two codes/subgrids models. While most of the halos in the simulations only host a few galaxies, there are a number of them that contain hundreds of satellites.

3.2. Training procedure

We construct our dataset, the collection of graphs, based on the procedure outlined in Sec. 2.2. We restrict our analysis to halos with more than one galaxy, as mentioned above. The models are trained on simulations of the CV and LH sets separately, where in each case we split the dataset into training (70%), validation (15%) and testing (15%) sets. We employ an Adam optimizer

(Kingma & Ba 2014) and L2 regularization. The batch size is set to 128 and the number of epochs limited to 150. The GNN architecture is implemented following the prescription outlined in Sec. 2.3, concatenating one or several message passing layers and appending at the end a global pooling and MLP to infer the halo mass.

We perform a hyperparameter optimization in order to get the best values for the hyperparameters, following a bayesian model-based optimization procedure with the Tree Parzen Estimator (TPE, Bergstra et al. 2011), making use of the Python package Optuna¹⁶ (Akiba et al. 2019). The hyperparameters considered for this optimization are the learning rate, the weight decay, the number of message passing layers, the distance to define neighborhoods, and the specific architecture (among those defined in Sec. 2.4). We perform at least 75 trials for each suite and set, where each trial is a specific choice of the values of the hyperparameters. The optimization procedure leads to different values for each simulation suite and set, with learning rates ranging between 10^{-5} and 6×10^{-4} , weight decay values between 10^{-8} and 10^{-7} , and the neighborhood distance between 2 and 20 Mpc/h. One unique GNN layer and the EdgeNet architecture are the optimal choices for all cases. The large values of the neighborhood radii obtained imply that most of the graphs created are complete, around 98% for both the CV and LH sets. That means that most of the galaxies are connected with each other within the halo, and most nodes can access to information about each other companion of the graph at one hop.

We follow a likelihood-free bayesian inference approach to sample the expected standard deviation of the outputs. This procedure allows us to reproduce some properties of the posterior (its second centered moment in this case) without the need of a likelihood (Jeffrey et al. 2021). We follow Jeffrey & Wandelt (2020) and design our model to output two quantities: the mean and standard deviation of the halo mass posterior. The loss function needed to achieve that is given by $\mathcal{L} = \mathcal{L}_1 + \mathcal{L}_2$, where

$$\mathcal{L}_1 = \log \left[\sum_{i \in \text{batch}} (y_{\text{truth},i} - y_{\text{infer},i})^2 \right] \quad (6)$$

and

$$\mathcal{L}_2 = \log \left[\sum_{i \in \text{batch}} \left((y_{\text{truth},i} - y_{\text{infer},i})^2 - \sigma_i^2 \right)^2 \right] \quad (7)$$

Note that we have included *log* functions to make sure that the two contributions are similar. For instance,

¹⁵ <https://camels.readthedocs.io>

¹⁶ <https://optuna.readthedocs.io>

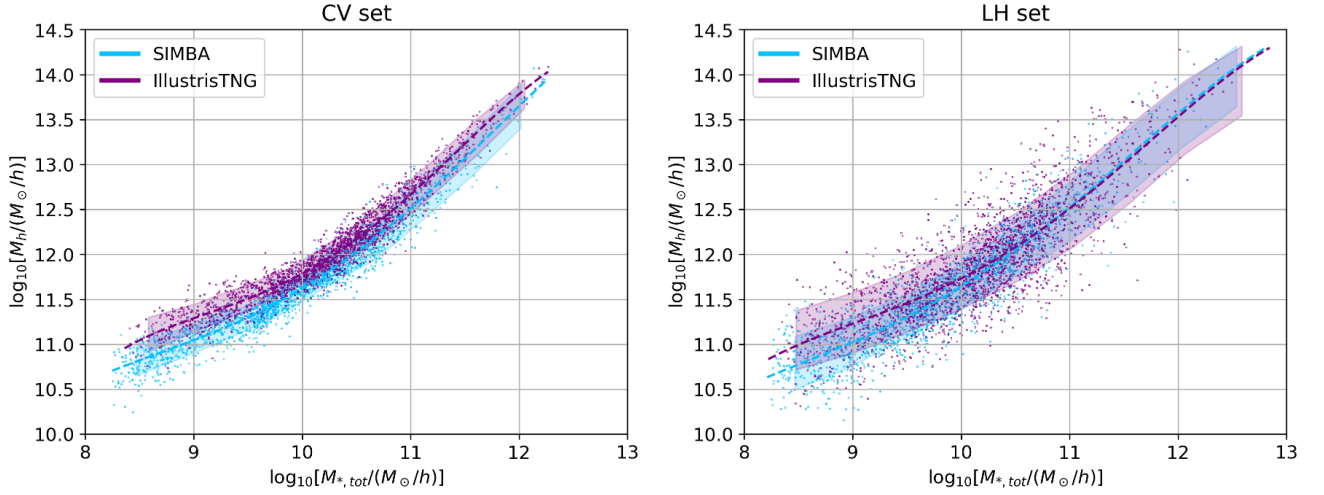


Figure 3. Halo masses with respect to the total stellar mass within, both for SIMBA and IllustrisTNG suites, for samples of halos from the CV (left) and LH (right) sets. Shaded areas denote standard deviation of points while dashed lines correspond to polynomial fits. The LH set covers a large astrophysical and cosmological parameter space, leading to a broader scatter in the masses with respect to the CV set, where parameters are fixed. While such fits can lead to accurate predictions of the halo mass for the CV set, it worsens in the LH set, given the larger dispersion.

if the errors are much smaller than the mean, the loss will be dominated by the mean and the errors may not be accurately computed (see Villaescusa-Navarro et al. 2021b, for further details).

It is worth emphasizing the symmetries fulfilled by the GNNs. By construction, GNNs are permutation invariant, guaranteed by neighborhood aggregation as discussed in Sec. 2.3. Furthermore, given that graphs are written in the center of mass rest frame, our framework is also translational invariant, since any displacement of the galaxy group would not alter the relative coordinates with respect to the center, and hence the graph would be the same. Moreover, for our special case of point clouds, our GNNs should be rotationally invariant, since an arbitrary rotation of all galaxies around the center of the halo should not change the global halo properties. We try to enforce this symmetry by performing random rotations on each graph at every training epoch, as a data augmentation procedure. This practice also helps in alleviating overfitting, which is absent in our training with the proper selection of the hyperparameters.

We write and train the models making use of PyTorch Geometric¹⁷ (Fey & Lenssen 2019). Our implementation of the GNNs, HaloGraphNet, is publicly available on GitHub¹⁸ (Villanueva-Domingo 2021).

4. RESULTS

In this section, the main results of the work are discussed. We first introduce a benchmark model to estimate halo masses from stellar masses, and next we detail the results of training and testing the GNNs in the different CAMELS simulation sets and suites considered, examining also their robustness over different astrophysical modeling.

4.1. Predictions from stellar masses

Before starting to discuss the accuracy of the GNN predictions, it may be useful to set a benchmark model to predict halo masses, to test the worthiness and degree of improvement of using GNNs. A traditional approach exploits the relation between stellar and halo mass, based on abundance matching (Behroozi et al. 2010; Wechsler & Tinker 2018). Figure 3 shows the mass of a halo M_h as a function of its total stellar mass $M_{*,tot} = \sum_i M_{*,i}$ (summing over all stellar particles in the central and satellite galaxies within the halo) in the IllustrisTNG and SIMBA suites, for the CV (left) and the LH (right) sets. One can notice a clear correlation between stellar and halo masses, although with a large scatter. This is specially noteworthy in the LH set, where multiple values of the cosmological and astrophysical parameters are considered, leading to completely different outcomes. Shaded areas denote the standard deviation of the points, which is around ~ 0.2 dex in the CV case, but grows up to $\sim 0.3 - 0.4$ dex in the LH set.

Taking advantage of the expressed correlation, we can build a naive estimator of the halo mass based only on the total stellar mass of galaxies. A simple 4th degree

¹⁷ <https://pytorch-geometric.readthedocs.io>

¹⁸ <https://github.com/PabloVD/HaloGraphNet>

polynomial fit is shown in dashed lines in Fig. 3. This can be regarded as a benchmark model for comparing with our forthcoming models based on GNNs, in order to further evaluate their strength. To check its accuracy for predicting $y = \log_{10}[M_h/(M_\odot/h)]$, we employ the mean relative error ϵ ,

$$\epsilon = \frac{1}{N} \sum_i^N \frac{|y_{\text{truth},i} - y_{\text{infer},i}|}{y_{\text{truth},i}}, \quad (8)$$

with N the number of test halos, as well as the correlation coefficient (or coefficient of determination) R^2 , defined in the usual way as

$$R^2 = 1 - \frac{\sum_i^N (y_{\text{truth},i} - y_{\text{infer},i})^2}{\sum_i^N (y_{\text{truth},i} - \bar{y}_{\text{truth}})^2}, \quad (9)$$

with \bar{y}_{truth} the mean of true values. This naive estimator gives fairly accurate results in the CV set, with relative errors around 1% and a linear correlation coefficient of $R^2 \simeq 0.94$. However, when an analogous fit is attempted in the LH case, the relative error worsens down to $\sim 1.7\%$ ($\sim 2.4\%$) for SIMBA (IllustrisTNG), with $R^2 \simeq 0.84$ ($R^2 \simeq 0.67$). Note that these errors are for the logarithm of the mass, y , rather than for the halo mass itself, which correspond to relative errors in the mass between $\sim 50 - 120\%$. This fact illustrates the non-trivial dependence of the halo and stellar masses on the different astrophysical and cosmological scenarios. Thus, a prediction of the halo mass from only the stellar mass when a broad range of astrophysical models are considered seems to be quite inaccurate. In the following, we shall show how a GNN is able to overcome this difficulty by considering further features and taking advantage of the graph structure of halos.

4.2. Inferring halo masses with GNNs

Here we first discuss the results of training a GNN to predict halo masses using galaxies from simulations of the CV set. As stated in Sec. 3.1, this set contains 27 simulations with fixed fiducial values for the cosmological and astrophysical parameters, only varying the random seed. The left panels of Fig. 4 show the accuracy of the network in the CAMELS CV sets, where the top panel stands for the IllustrisTNG subgrid model and the bottom one for the SIMBA suite. The vertical axis is the difference between the predicted and true logarithms of the halo mass, $\log_{10}(M_h/(M_\odot/h))$, with respect to the true value in the horizontal axis. Error bars have been estimated via likelihood-free inference, sampling the standard deviation of the posterior, as outlined in Sec. 3.2. One can see that the performance is fairly good in both cases, as the linear correlation coefficient of $R^2 = 0.96 - 0.97$ confirms, with a relative

error lower than $\sim 1\%$ in y ($\sim 25 - 40\%$ in the mass itself). The dashed lines show the mean of test points, while the shaded region depicts their actual standard deviation, which extends up to ~ 0.14 dex. Thus, most of the test predictions lie within this region, although there are some outliers. The neural network is thus able to accurately infer the halo mass given some features of its galaxies.

Nevertheless, the previous results only show that the GNN is capable of predicting the mass of a halo when cosmology and astrophysics is known, i.e., when the parameters are fixed. However, the specific values of the relevant parameters for the real Universe are not well known yet, specially the astrophysical ones. In order to marginalize over uncertainties in cosmology and astrophysics, we have trained our network in the LH simulation set, which includes one thousand simulations varying cosmological and astrophysical parameters, together with the random seed (to also incorporate effects of sample variance), as noted in Sec. 3.1. The right panels of Fig. 4 show the GNN predictions for training the GNN in IllustrisTNG (top) and SIMBA (bottom). One can see that the results only slightly worsen from the equivalent case in the CV simulation set, with $R^2 = 0.90 - 0.92$ and relative errors of $\sim 1\%$. It should be noted that a small bias appears at low masses, below $\lesssim 10^{11} M_\odot/h$, which deviates up to 0.4 dex (also present in the SIMBA CV case). Uncertainties are also larger than in the CV set, roughly by a factor of 2, meaning that the model trained in the LH set offers lower precision. This worsening is expected since the network needs to marginalize over the astrophysical and cosmological parameters, differently from the network trained on the CV set. In any case, these results indicate that our model is able to learn the mapping between astrophysics/cosmology and the halo mass, and thus marginalize over the value of the parameters present in the simulations to make an accurate prediction.

To further evaluate the predictive power of GNNs, it is worth comparing these results to those obtained from the naive fit based only on stellar mass outlined in Sec. 4.1. For both the CV and LH sets, the GNN predictions outperform those from the polynomial fit, presenting larger R^2 coefficients and lower relative errors. The improvement is especially clear in the LH case, where the R^2 is significantly better than the benchmark method. Moreover, the scatter around the true values spans ~ 0.14 and $\lesssim 0.2$ dex, compared to the larger standard deviations from the naive fit, which are 0.2 and 0.3 – 0.4 dex respectively. Note that 0.2 dex is a factor up to ~ 1.6 in the mass (rather than in the logarithm), while 0.3 dex is a factor ~ 2 , meaning a $\sim 100\%$

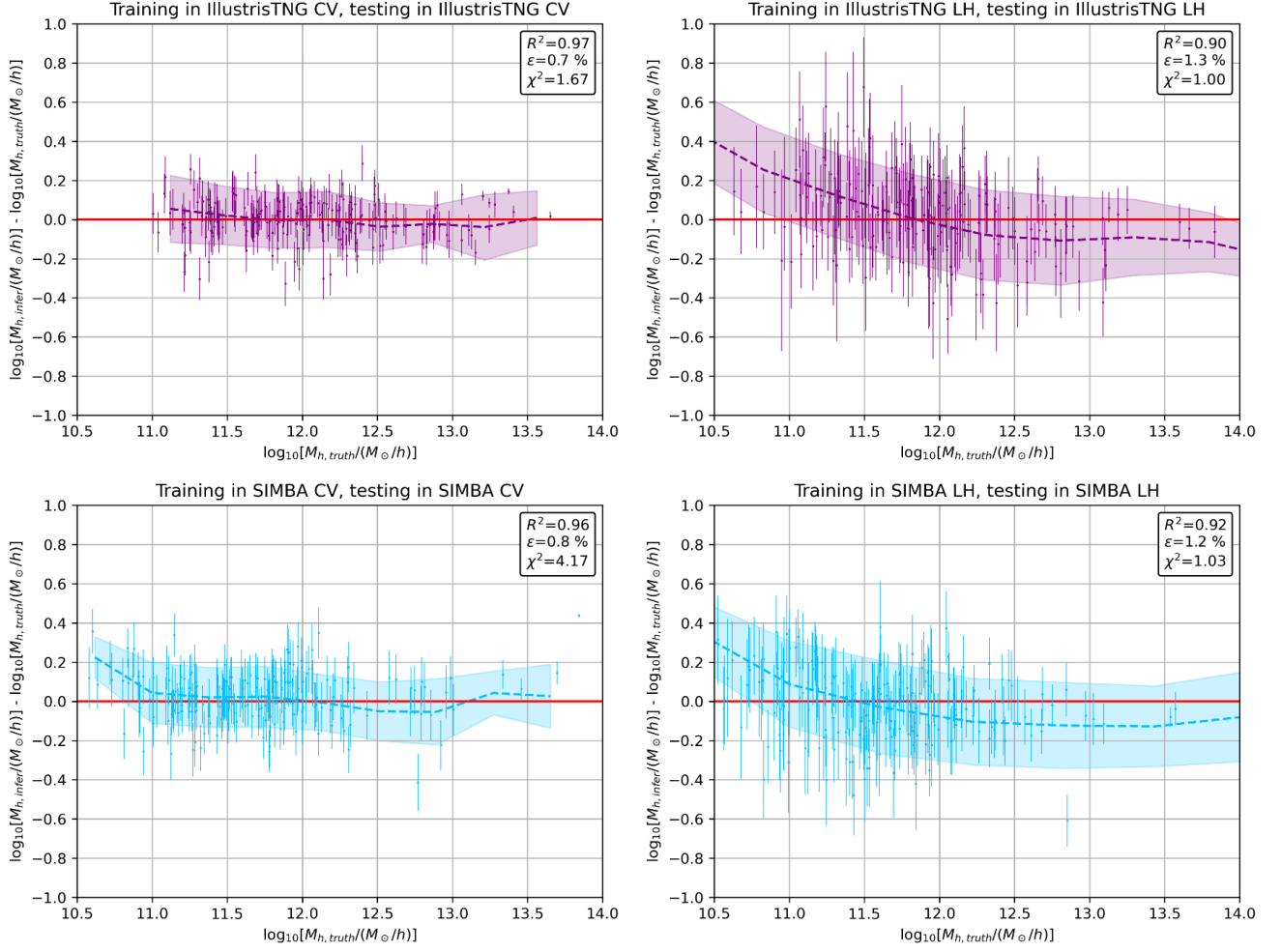


Figure 4. Predicted versus true logarithm of halo masses $\log_{10} [M_h / (M_\odot / h)]$ for the CV (left) and LH sets (right), training in the IllustrisTNG suite (top) and in SIMBA (bottom). A sample of 200 halos in the test dataset is shown in each case. Shaded regions and dashed lines correspond to real standard deviation and mean of test points, respectively. While in the CV set, astrophysical and cosmological parameters are fixed to fiducial values, the LH set comprises a broad range of astrophysical and cosmological scenarios. Even so, the GNN is still able to learn the halo/galaxy relation and predict masses in the LH case, only slightly worsening the prediction with respect to the CV case.

error in the mass. These results demonstrate how taking advantage of the graph structure and further galaxy features, it is possible to attain richer correlations and better results.

Nevertheless, it has to be emphasized that the linear correlation coefficient R^2 and the relative error cannot constitute a complete statistical summary for testing the accuracy of the GNN. It is because our network is also predicting the standard deviation of the target, for which an additional component to the loss has been included, as discussed in Sec. 3.2. Thus, neither R^2 nor the relative error quantify the sampling accuracy of the standard deviation. To test whether the uncertainties are reasonably predicted, it is useful to compute the χ^2 ,

defined as

$$\chi^2 = \frac{1}{N} \sum_i \frac{(y_{\text{truth},i} - y_{\text{infer},i})^2}{\sigma_i^2}. \quad (10)$$

Note that minimizing the loss function contribution from Eq. 7 tends to drive the χ^2 towards unity. This is actually the case in the LH set, where $\chi^2 = 1.00$ for IllustrisTNG and $\chi^2 = 1.03$ for SIMBA, indicating that uncertainties are accurately predicted. In the CV cases, however, while the mean predictions are more accurate, some errors are underestimated, leading to larger χ^2 values, around ~ 4 in SIMBA. Moreover, since the test dataset is smaller, a few outliers with too small standard deviations greatly impact the value of the χ^2 .

One can also compute how many points in the test dataset present an accurate uncertainty by counting the fraction which fulfill the conditions $|y_{\text{truth},i} - y_{\text{infer},i}| \leq \sigma_i$ and $|y_{\text{truth},i} - y_{\text{infer},i}| \leq 2\sigma_i$, i.e., how many points lie within one and two times the standard deviation of the posterior. We find that for the LH cases, the fraction of points fulfilling the above conditions is 69% and 95% respectively for both suites. For the CV sets, these fractions only slightly deviate, 62 and 90% for IllustrisTNG, and 72 and 96% for SIMBA, respectively. For Gaussian distributed errors, these fractions should be around 68 and 95% respectively. Note however that our calculation of the posterior mean and standard deviation does not make any assumption about the form of the posterior. Therefore, the numbers quoted above should be taken with caution and comparison with the Gaussian case should be done in a careful manner.

There is another way to figure out whether uncertainties are correctly sampled. The shaded regions in Fig. 4 represent the actual standard deviation of the test points computed within several mass bins. Therefore, if the predicted uncertainties are accurately sampled, their mean value $\bar{\sigma}$ should correspond to those shaded regions. For the CV case, despite the GNN providing more accurate models, mean uncertainties $\bar{\sigma}$ are underpredicted with respect to actual scatter by $\sim 40\%$ for IllustrisTNG and $\sim 20\%$ for SIMBA. However, in the LH case, uncertainties are better sampled, only deviating $\sim 3\%$ and $\sim 7\%$ for IllustrisTNG and SIMBA respectively. This implies that models trained in LH, despite predicting slightly less accurate results than those from the CV set, provide a better sample of the posterior uncertainties.

4.3. Robustness over different subgrid physics models

Subgrid physics, i.e. the models used to simulate unresolved astrophysical processes such as the feedback from supernovae and black holes, can only be implemented in a phenomenological way and therefore there is not a unique subgrid model that best represents reality. Thus, having a ML model robust over different subgrid scenarios would be needed in order to obtain predictions that do not depend on the particular type of simulation used to train the networks.

To check whether our GNN fulfills this requirement, we have taken the model previously trained on simulations from the IllustrisTNG suite, and we have tested it on galaxies from the SIMBA simulations, which make use of a completely different subgrid physics model for AGN and SN feedback. The top left panel of Fig. 5 shows the predictions for the halo mass in the CV set, i.e., trained in IllustrisTNG CV and tested in SIMBA

CV. We see that the performance becomes worse, and actually a bias appears. This offset may arise from the fact that the IllustrisTNG and SIMBA models make different predictions for the halo-galaxy connection, as shown in Fig. 3. This can be related with the fact that astrophysical parameters are not completely correlated and calibrated between both cases. While the CV set assumes fiducial values for the astrophysical and cosmological parameters, the default values in the IllustrisTNG suite do not correspond to the ones in SIMBA, since they refer to different quantities and physics. The absence of a one-to-one relation between both suites in the CV set may explain why the network fails to make a robust prediction.

The right panel of Fig. 5 depicts the results of testing the network trained on galaxies of the IllustrisTNG LH set on galaxies of the SIMBA LH set. In this case, the bias present in the CV case disappears, and only a broad scatter holds. The absence of the offset can be attributed to the intrinsic marginalization of astrophysical effects carried out by the network. Moreover, the linear correlation coefficient only slightly decreases with respect to testing in IllustrisTNG (top right panel of Fig. 4), and is actually better than in the CV counterpart (top left panel of Fig. 5), in spite of dealing with a much broader astrophysical parameter space. The χ^2 also remains closer to unity than in the CV case, meaning that uncertainties are better sampled in the LH set. The fraction of points fulfilling the conditions $|y_{\text{truth},i} - y_{\text{infer},i}| \leq \sigma_i$ and $|y_{\text{truth},i} - y_{\text{infer},i}| \leq 2\sigma_i$ is 65% and 93% respectively, while much lower for the CV case (20 and 46% respectively). These facts imply that models trained in the LH set generalize better than those trained in the CV one.

An analogous experiment has been carried out, employing the model trained in the SIMBA suite and testing it in IllustrisTNG. The results are shown in the bottom panels of Fig. 5 for the CV set (left) and LH set (right). Note the offset in the CV case, which is similar but opposite (underpredicting the mass) to the one appearing in the top left panel of Fig. 5. This is reasonable, since given that the IllustrisTNG model overpredicts the mass in SIMBA, the contrary should be expected when a SIMBA model is tested in IllustrisTNG. In the LH scenario, predictions are slightly worse, decreasing the truth-prediction correlation down to $R^2 = 0.8$, and also poorer than in the CV case. Both sets present larger χ^2 values and a deviation from the Gaussian counterpart, given that only a fraction of $\sim 50\%$ and $\sim 80\%$ points lie within one and two times the posterior standard deviation, respectively.

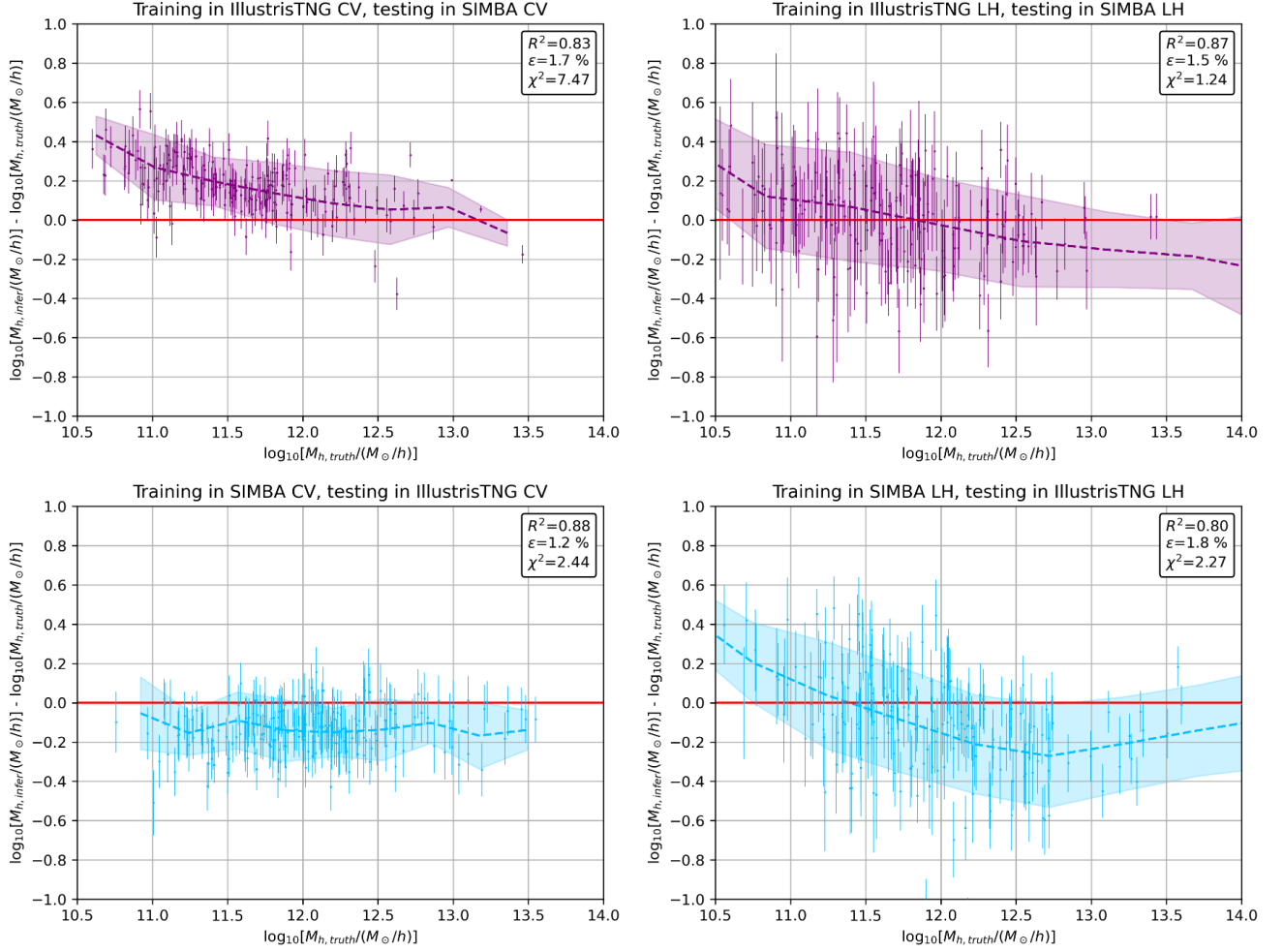


Figure 5. Same as Fig. 4 but either using a model trained with the IllustrisTNG suite and tested with the SIMBA simulations (top) or trained in SIMBA and tested in the IllustrisTNG suite (bottom), for CV (left) and LH (right) sets. A model trained in a given suite worsens its behavior when tested in the other one, appearing biased in the CV case. However, in the LH set, it is possible to find a mapping between the parameter space of both subgrid physics models, alleviating such biases.

It is noteworthy that, for the LH set, a model trained in IllustrisTNG generalizes better in SIMBA than the opposite case. An interpretation of this outcome is problematic. SIMBA usually covers more extreme astrophysical scenarios than IllustrisTNG, presenting also more galaxies per halo (see Fig. 2), facts which could have an impact in this different cross testing. SIMBA simulations usually show a more stochastic behavior, presenting more scatter in some properties such as galaxy scaling relations. This can be due to the small box size compared to the large scale effect of the SIMBA AGN jets (Davé et al. 2019; Villaescusa-Navarro et al. 2021a). Nevertheless, although we observe a higher variance and some outliers, most of the predictions only deviate up to 0.4 dex from the ground truth. In all cases, the mean relative error in the logarithm of the mass lies below

2%. The GNNs are able to recover the true mass in the majority of the cases within the standard deviation uncertainty. This shows that the models are relatively accurate even when they are applied to simulations with a different subgrid physics modeling, manifesting their robustness.

5. INTERPRETABILITY OF THE GNN

Neural networks represent superb tools to deal with multiple problems, but the factors that determine their behavior and performance are difficult to understand. It is always desirable to gain some interpretability of our ML models, in order to determine which features and inputs are the most relevant for predicting the output.

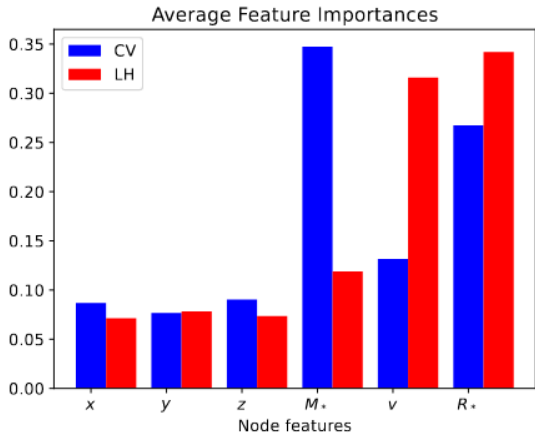


Figure 6. Saliency values of the galaxy features employed to train the GNN. Larger values imply that the network prediction is more influenced by variations of that feature, indicating that such properties are more relevant for the result.

To do so, we make use of the Python library Captum¹⁹ (Kokhlikyan et al. 2020), a package designed for model interpretability and attribute selection. Specifically, we employ the saliency method for computing the gradients of the outputs with respect to the inputs and features employed (Simonyan et al. 2013). In that way, larger gradients, and hence larger saliency values, imply that the prediction is more sensitive to variations of that variable.

We have computed the gradients with respect to the features employed in the training of the GNN, averaged over all galaxies and taken the absolute value. This approach gives us a saliency value for each property, which is shown in Fig. 6 for the IllustrisTNG suite in the CV (blue) and LH (red) sets. For both sets, we can notice that, as could be expected, each coordinate of the position presents a very similar importance, due to implicit rotational invariance present in the problem. Regarding the CV case, the most relevant features to determine the output appear to be the stellar mass, M_* , and subsequently, the stellar half-mass radius, R_* . However, in the LH case, the network focuses more on R_* and v rather than the stellar mass. One can interpret this change of learning behavior from the larger scatter in stellar mass shown in Fig. 3 in LH with respect to the CV set; the network is marginalizing over the scatter due to astrophysics and cosmology. In other words, the large variety of astrophysical and cosmological scenarios in the LH set may require the network to base its predic-

tions on features that do not exhibit such large scatter. While the predictor mostly rely on the M_* - M_h correlation when the scatter is small, velocities and galaxy size become better tracers otherwise. Nevertheless, one has to be cautious with these interpretations, since the saliency values may not give a complete picture. For instance, spatial position importance may be underestimated since its weight could be split into the three coordinates.

It is common in computer vision tasks to calculate the saliency map, which indicates those pixels in an image that are most relevant for the final output (see, e.g., Villanueva-Domingo & Villaescusa-Navarro 2021 for an application in cosmology). In our case, we are dealing with graphs rather than images. Therefore, for a given halo it is possible to compute the *saliency graph*, which shows the nodes whose features are more relevant for predicting the halo mass. Saliency values can be computed using the same procedure as above, but taking the absolute value and averaging over all features at each node. Examples of such saliency graphs for different halos are depicted in Fig. 7, where the color indicates the saliency value, and the size is proportional to the stellar mass. Neighbors are connected by lines. Chosen samples present relative errors lower than 1.5 % to ensure that their saliency graphs are meaningful. Top row stands for the CV set and bottom for LH, both in IllustrisTNG. In the CV set, as one would naively expect, the central galaxies provide the most relevant nodes for the output. These galaxies are also those with larger stellar masses. However, given that the stellar mass is an informative property, as seen in Fig. 6, halos with relatively massive satellites can also show other relevant nodes besides the central one. On the other hand, since in the LH set, stellar mass is less informative, as seen in Fig. 6, the network may focus more on some low mass satellites rather than in central galaxies, which become less important.

It is thus pertinent to ask ourselves which satellites leave a greater impact on the halo mass. Fig. 8 depicts the saliency value of each satellite galaxy as a function of distance to the center, excluding central galaxies. Point size is proportional to the stellar mass of the node. One can notice a trend where closer galaxies become more relevant than those farther away. Moreover, this tendency seems to be relatively independent on their stellar mass. This plot employs the IllustrisTNG CV dataset, although similar qualitative conclusions can be extracted from the other cases.

These tests provide us with enlightening information about how GNNs learn and predict their outputs, as well as which are the most relevant components to un-

¹⁹ <https://captum.ai/>

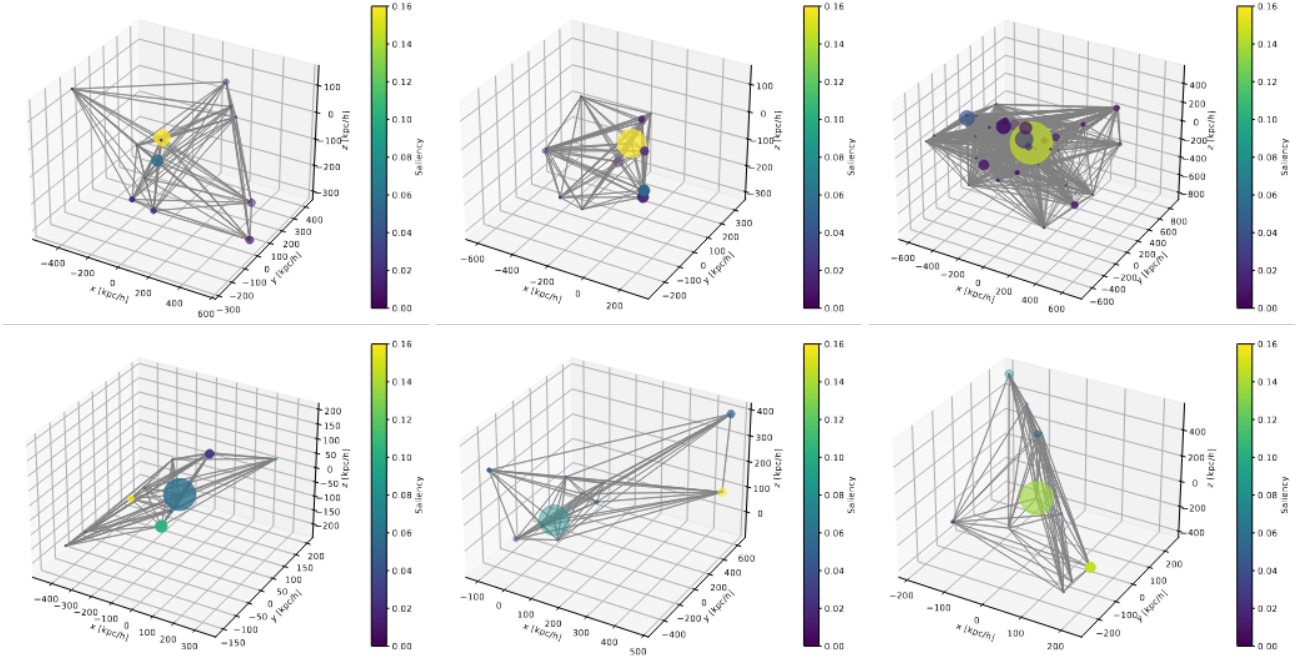


Figure 7. Saliency graphs for six different halos from the IllustrisTNG CV (top) and LH (bottom) datasets. Colors denote the saliency attributes, meaning that galaxies with larger values are more relevant for the prediction. Sizes of the nodes are proportional to the stellar mass of the galaxies. In the CV set, central galaxies usually have the greatest impact on the prediction, although if there are massive satellites, they can also contribute significantly (e.g., right panel). In the LH case, however, more attention is often focused on smaller satellites rather than on central galaxies.

derstand the halo/galaxy connection. Further development is however required in order to obtain a better understanding of the training procedure and the emergent properties of halos from galaxies.

6. DISCUSSION

Constraining the total mass of dark matter halos from galaxies still presents a challenge from both theoretical and observational perspectives, given the large contribution of the dark matter component. In this work we have presented a new method based on artificial intelligence designed to infer the total mass of a halo from the properties of the galaxies it hosts. The point cloud arrangement of halo-galaxy catalogs has been exploited in order to structure halos as mathematical graphs, where galaxies constitute the nodes, connected by proximity. This organization of data makes it possible to employ GNNs, naturally suited to operate with graphs, to extract global permutation invariant quantities, as is the case of the halo mass. The models are fed with different observable galactic features, such as the position, velocity, stellar mass, and stellar half-mass radius.

We have made use of the large collection of state-of-the-art hydrodynamic simulations from the CAMELS project to train our networks, which include thousands of simulations covering different astrophysical scenar-

ios. Training GNNs over this dataset allows us to achieve precise models capable of predicting the mass of a halo with remarkable accuracy. The networks successfully marginalize over a broad astrophysical parameter space, learning the connection between the halo mass and the properties of the galactic components. Our results strongly rely on the velocity and size of satellite galaxies, illustrating the importance of features beyond the stellar mass. Moreover, we have proven that the trained networks in a simulation suite still provide relatively precise predictions when tested in simulations with a different subgrid physics model, illustrating their robustness with respect to the astrophysical modeling. Hyperparameter optimization has been carried out to maximize the performance of the networks. We have performed likelihood-free bayesian inference, providing additionally an estimate for the standard deviation without knowing the actual likelihood.

It is not straightforward to compare our results to ML methods estimating the halo mass in previous literature, given that the datasets employed, problem setups and features considered can be notably different. [Calderon & Berlind \(2019\)](#) apply several ML techniques, such as standard MLPs, to predict the halo mass from galaxy groups data. They train their models in semianalytical

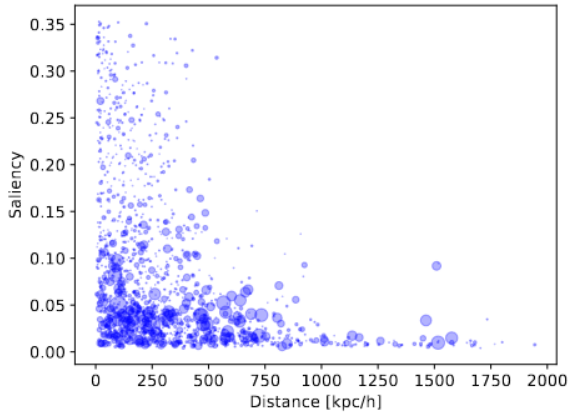


Figure 8. Saliency for each satellite galaxy depending on its distance to the halo center, excluding the central ones. Point diameter indicates the stellar mass of the nodes. There is a general tendency where closer galaxies to the center present larger saliency values, and therefore, are more relevant for predicting the output.

galaxy catalogues calibrated to SDSS data, outperforming traditional methods to infer the dynamical mass. Their 1σ scatter region spans $\sim 0.4 - 0.6$ dex (a fractional difference of $\sim 3-5\%$), while for the same masses, our GNN reduces down to ~ 0.14 dex for the CV set and ~ 0.2 dex for the LH one. In terms of the relative error in the mass (rather than in the logarithm), their scatter gets up to $\sim 250 - 400\%$, while our models reduce it down to 40 % and 60 % for CV and LH respectively. Furthermore, our models make use of fewer galaxy features (6 rather than their 9) and the LH simulations cover a large volume of the astrophysical and cosmological parameter space. Man et al. (2019) face a similar problem of predicting halo masses from group galaxy properties, splitting their datasets in red and blue groups, according to the color of the central galaxies. They train a random forest estimator, obtaining scatter comparable to ours (note the $\sqrt{2}$ factor of difference in their definitions), but employing 9 features with a semianalytical galaxy formation model with fixed parameters. von Marttens et al. (2021) predict galactic properties, such as the total mass, training different ML models in the TNG100 simulation. As in our case, the galaxies considered present stellar masses above $\gtrsim 10^8 M_\odot/h$, although they predict the subhalo mass rather than the halo one. They employ up to 15 features, including photometric, kinematic, and structural properties. When all features are included, they are able to get accurate results, with correlation of $R^2 = 0.92$. This case can be compared to our IllustrisTNG CV set, since it presents fixed astro-

physical and cosmological parameters,²⁰ where we obtain a better correlation of $R^2 = 0.97$ even using only 6 features. Moreover, we are also capable of obtaining a similar accuracy in the LH set than their models in TNG100. Lucie-Smith et al. (2020) train CNNs on the density field of dark matter halos to infer their masses, with an accuracy comparable to the one of our network train on the CV set. However, note that this is a very different task, since that method relies on the 3D dark matter density field rather than on observable features, as is the case of our GNNs. Moreover, the authors employ N-body simulations with fixed cosmology rather than full hydrodynamic simulations. In any case, one has to be cautious at comparing these different approaches, due to the distinct datasets, assumptions and features considered, being appropriate only for illustrating the potential of GNN models.

While the ML method presented here shows reasonable accuracy, it however has some caveats that have to be emphasized. For instance, it may not be obvious whether one galaxy belongs to a halo or not, in cases where two galactic groups are close together. This can be exacerbated in real observations, which take place in redshift-space. This may cause the appearance of interlopers in a halo which could be counted as its own satellites, distorting the results. The effect of the presence of other halos in the environment when building the graph is also disregarded, which could have an impact for some close halos gravitationally bound. A way to deal with the influence of surrounding groups may be to include some global feature regarding the amount of galaxies or halos within a certain distance from the halo which is evaluated. Moreover, we have only trained with halos at $z = 0$, but it would be also convenient to derive models capable of inferring halo masses at earlier times.

The results presented here are obtained by restricting ourselves to a reduced set of several observable quantities. However, it is possible in principle that using further variables, the results could improve. Additional correlations with the halo mass could arise if supplementary features are considered to train the GNNs. For instance, previous works like von Marttens et al. (2021) have shown the importance of photometric variables, such as luminosities at different wavelengths, for inferring the total halo mass. Additionally, considering, e.g., radii enclosing 20% and 80% of the mass (or light) may provide further information about the concentration and morphology of the halo. Among other appealing galactic

²⁰ Although TNG100 has a box size of ~ 100 Mpc/h, larger than the 25 Mpc/h CAMELS box size, the CV set includes 27 simulations varying the random seed.

properties that could be contemplated are the gas mass, the HI mass, metallicities, velocity dispersion, etc. Furthermore, it would be desirable to explore more complex GNN architectures to find whether more accurate results can be achieved, reducing the scatter and the predicted uncertainties. The training of GNNs on galactic properties combined with other observables such as lensing could also enhance its predictive power.


The framework developed in this article has been proposed to infer the total mass of a given halo. However, since it is based on extracting global quantities from galaxy features, it could be generally applied to predict any other global quantity of the halo, such as its concentration, spin, characteristic age, etc. Depending on the specific global quantity considered, different galaxy or subhalo features may be required in order to maximize accuracy, which in principle could differ from those employed here. Moreover, one could apply GNN architectures to other problems, such as edge prediction. For instance, given a set of galaxies, a GNN could be trained to identify those which conform separate halos, as a ML alternative to friend-of-friends algorithms.

It would be desirable to synthesize the predictive power of the GNNs into an analytical formula, via symbolic regression, as has been done already in other cosmological contexts (Cranmer et al. 2019; Cranmer et al. 2020; Shao et al. 2021; Wadekar et al. 2020). The analytical formulae derived in this way would depend upon the messages exchanged between the nodes. Nonetheless, obtaining a model susceptible to be replicated by a concise formula presents further difficulties. Performing symbolic regression requires interpretable models. One way to achieve this is to enforce sparse weights, obtained, e.g., by applying L1 regularization, which can affect the accuracy of the network. Sometimes, increasing the accuracy of a network reduces its extrapolation properties, and the other way around (Shao et al. 2021). Hence, deriving simple analytical formulae via symbolic regression from these GNNs remains as a desirable but challenging step for future work.

The description of halos as graph-structured data makes this kind of network suitable for other types of astrophysical systems, which are characterized by point clouds. Examples of these may include galaxy clusters, globular clusters or even planetary systems. Any point distribution could take advantage of the graph structure presented here for halos, in order to derive global quantities of such systems. An unexplored window remains open to apply all the power of GNNs to astrophysics.

Given that halo masses can be accurately predicted from numerical simulations by using the method shown here, the natural step forward would be applying this kind of model to real data. This task is carried out in a companion paper in preparation (Villanueva-Domingo, P., Villaescusa-Navarro, F., Genel, S., et al. 2021), where GNNs trained with CAMELS simulations are employed to predict the masses of the halos containing the Milky Way and Andromeda. In that article, independent predictions for both halos are presented, which are consistent with other standard methods for estimating the dynamical masses of our galaxy and its companion. This represents a success of applying ML models trained with numerical simulations on real data, illustrating the power of artificial intelligence to enhance our knowledge of the Universe.

DATA AVAILABILITY

The models and implementation of GNNs in PyTorch Geometric, **HaloGraphNet**, are available on GitHub ²¹ (Villanueva-Domingo 2021). Details on the CAMELS simulations can be found in <https://www.camel-simulations.org>.

ACKNOWLEDGEMENTS

We thank Miles Cranmer for enlightening discussions at early stages of this work. The work of PVD is supported by CIDEAGENT/2018/019, CPI-21-108. DAA was supported in part by NSF grants AST-2009687 and AST-2108944, and by the Flatiron Institute, which is supported by the Simons Foundation. The training of the GNNs has been carried out using GPUs from the Tiger cluster at the Princeton University.

REFERENCES

- Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. 2019, arXiv e-prints, arXiv:1907.10902.
<https://arxiv.org/abs/1907.10902>
- ²¹ <https://github.com/PabloVD/HaloGraphNet>
- Anglés-Alcázar, D., Davé, R., Faucher-Giguère, C.-A., Özel, F., & Hopkins, P. F. 2017, MNRAS, 464, 2840,
 doi: [10.1093/mnras/stw2565](https://doi.org/10.1093/mnras/stw2565)
- Armitage, T. J., Kay, S. T., & Barnes, D. J. 2019, MNRAS, 484, 1526, doi: [10.1093/mnras/stz039](https://doi.org/10.1093/mnras/stz039)

- Battaglia, P. W., Hamrick, J. B., Bapst, V., et al. 2018, arXiv e-prints, arXiv:1806.01261. <https://arxiv.org/abs/1806.01261>
- Beck, R., & Sadowski, P. 2019, Refined Redshift Regression in Cosmology with Graph Convolution Networks
- Behroozi, P. S., Conroy, C., & Wechsler, R. H. 2010, ApJ, 717, 379, doi: [10.1088/0004-637X/717/1/379](https://doi.org/10.1088/0004-637X/717/1/379)
- Bergstra, J., Bardenet, R., Bengio, Y., & Kégl, B. 2011, in Advances in Neural Information Processing Systems, ed. J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, & K. Q. Weinberger, Vol. 24 (Curran Associates, Inc.). <https://proceedings.neurips.cc/paper/2011/file/86e8f7ab32cfd12577bc2619bc635690-Paper.pdf>
- Bosma, A. 1978, PhD thesis, -
- Bronstein, M. M., Bruna, J., Cohen, T., & Veličković, P. 2021, arXiv e-prints, arXiv:2104.13478. <https://arxiv.org/abs/2104.13478>
- Calderon, V. F., & Berlind, A. A. 2019, MNRAS, 490, 2367, doi: [10.1093/mnras/stz2775](https://doi.org/10.1093/mnras/stz2775)
- Cranmer, M., Melchior, P., & Nord, B. 2021. <https://arxiv.org/abs/2106.09761>
- Cranmer, M., Sanchez-Gonzalez, A., Battaglia, P., et al. 2020. <https://arxiv.org/abs/2006.11287>
- Cranmer, M. D., Xu, R., Battaglia, P., & Ho, S. 2019, arXiv e-prints, arXiv:1909.05862. <https://arxiv.org/abs/1909.05862>
- Davé, R., Anglés-Alcázar, D., Narayanan, D., et al. 2019, Mon. Not. Roy. Astron. Soc., 486, 2827, doi: [10.1093/mnras/stz937](https://doi.org/10.1093/mnras/stz937)
- Davé, R., Thompson, R., & Hopkins, P. F. 2016, MNRAS, 462, 3265, doi: [10.1093/mnras/stw1862](https://doi.org/10.1093/mnras/stw1862)
- de Andres, D., Cui, W., Ruppín, F., et al. 2021. <https://arxiv.org/abs/2111.01933>
- Fey, M., & Lenssen, J. E. 2019, Fast Graph Representation Learning with PyTorch Geometric, 2.0.2. https://github.com/pyg-team/pytorch_geometric
- Green, S. B., Ntampaka, M., Nagai, D., et al. 2019, Astrophys. J., 884, 33, doi: [10.3847/1538-4357/ab426f](https://doi.org/10.3847/1538-4357/ab426f)
- Grego, L., Carlstrom, J. E., Reese, E. D., et al. 2001, ApJ, 552, 2, doi: [10.1086/320443](https://doi.org/10.1086/320443)
- Gupta, N., & Reichardt, C. L. 2020, ApJ, 900, 110, doi: [10.3847/1538-4357/aba694](https://doi.org/10.3847/1538-4357/aba694)
- Haider Abbas, M. 2019, arXiv e-prints, arXiv:1912.05316. <https://arxiv.org/abs/1912.05316>
- Hamilton, W. L. 2020, Synthesis Lectures on Artificial Intelligence and Machine Learning, 14, 1
- Ho, M., Rau, M. M., Ntampaka, M., et al. 2019, Astrophys. J., 887, 25, doi: [10.3847/1538-4357/ab4f82](https://doi.org/10.3847/1538-4357/ab4f82)
- Hopkins, P. F. 2015, Mon. Not. Roy. Astron. Soc., 450, 53, doi: [10.1093/mnras/stv195](https://doi.org/10.1093/mnras/stv195)
- Huang, S., Leauthaud, A., Hearin, A., et al. 2020, MNRAS, 492, 3685, doi: [10.1093/mnras/stz3314](https://doi.org/10.1093/mnras/stz3314)
- Huang, S., Leauthaud, A., Bradshaw, C., et al. 2021, arXiv e-prints, arXiv:2109.02646. <https://arxiv.org/abs/2109.02646>
- Jeffrey, N., Alsing, J., & Lanusse, F. 2021, Mon. Not. Roy. Astron. Soc., 501, 954, doi: [10.1093/mnras/staa3594](https://doi.org/10.1093/mnras/staa3594)
- Jeffrey, N., & Wandelt, B. D. 2020, arXiv e-prints, arXiv:2011.05991. <https://arxiv.org/abs/2011.05991>
- Kingma, D. P., & Ba, J. 2014. <https://arxiv.org/abs/1502.03167>
- Kipf, T. N., & Welling, M. 2017, Semi-Supervised Classification with Graph Convolutional Networks. <https://arxiv.org/abs/1609.02907>
- Kodi Ramanah, D., Wojtak, R., Ansari, Z., Gall, C., & Hjorth, J. 2020, Mon. Not. Roy. Astron. Soc., 499, 1985, doi: [10.1093/mnras/staa2886](https://doi.org/10.1093/mnras/staa2886)
- Kodi Ramanah, D., Wojtak, R., & Arendse, N. 2021, Mon. Not. Roy. Astron. Soc., 501, 4080, doi: [10.1093/mnras/staa3922](https://doi.org/10.1093/mnras/staa3922)
- Kokhlikyan, N., Miglani, V., Martin, M., et al. 2020, Captum: A unified and generic model interpretability library for PyTorch. <https://arxiv.org/abs/2009.07896>
- Landry, D., Bonamente, M., Giles, P., et al. 2013, MNRAS, 433, 2790, doi: [10.1093/mnras/stt901](https://doi.org/10.1093/mnras/stt901)
- Lucie-Smith, L., Peiris, H. V., Pontzen, A., Nord, B., & Thiagalingam, J. 2020, arXiv e-prints, arXiv:2011.10577. <https://arxiv.org/abs/2011.10577>
- Man, Z.-Y., Peng, Y.-J., Shi, J.-J., et al. 2019, ApJ, 881, 74, doi: [10.3847/1538-4357/ab2ece](https://doi.org/10.3847/1538-4357/ab2ece)
- Mandelbaum, R. 2015, in Galaxy Masses as Constraints of Formation Models, ed. M. Cappellari & S. Courteau, Vol. 311, 86–95, doi: [10.1017/S1743921315003452](https://doi.org/10.1017/S1743921315003452)
- Marinacci, F., Vogelsberger, M., Pakmor, R., et al. 2018, MNRAS, 480, 5113, doi: [10.1093/mnras/sty2206](https://doi.org/10.1093/mnras/sty2206)
- Naiman, J. P., Pillepich, A., Springel, V., et al. 2018, MNRAS, 477, 1206, doi: [10.1093/mnras/sty618](https://doi.org/10.1093/mnras/sty618)
- Nelson, D., Springel, V., Pillepich, A., et al. 2019, Computational Astrophysics and Cosmology, 6, 2, doi: [10.1186/s40668-019-0028-x](https://doi.org/10.1186/s40668-019-0028-x)
- Ntampaka, M., Trac, H., Sutherland, D. J., et al. 2015, Astrophys. J., 803, 50, doi: [10.1088/0004-637X/803/2/50](https://doi.org/10.1088/0004-637X/803/2/50)
- . 2016, Astrophys. J., 831, 135, doi: [10.3847/0004-637X/831/2/135](https://doi.org/10.3847/0004-637X/831/2/135)
- Ntampaka, M., et al. 2019, Astrophys. J., 876, 82, doi: [10.3847/1538-4357/ab14eb](https://doi.org/10.3847/1538-4357/ab14eb)
- Old, L., Wojtak, R., Mamon, G. A., et al. 2015, MNRAS, 449, 1897, doi: [10.1093/mnras/stv421](https://doi.org/10.1093/mnras/stv421)
- Pillepich, A., Springel, V., Nelson, D., et al. 2018, MNRAS, 473, 4077, doi: [10.1093/mnras/stx2656](https://doi.org/10.1093/mnras/stx2656)

- Qi, C. R., Su, H., Mo, K., & Guibas, L. J. 2017a, PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. <https://arxiv.org/abs/1612.00593>
- Qi, C. R., Yi, L., Su, H., & Guibas, L. J. 2017b, PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space. <https://arxiv.org/abs/1706.02413>
- Rubin, V. C., Ford, W. K., J., & Thonnard, N. 1978, ApJ, 225, L107, doi: [10.1086/182804](https://doi.org/10.1086/182804)
- Saro, A., Mohr, J. J., Bazin, G., & Dolag, K. 2013, ApJ, 772, 47, doi: [10.1088/0004-637X/772/1/47](https://doi.org/10.1088/0004-637X/772/1/47)
- Seo, G., Sohn, J., & Lee, M. G. 2020, ApJ, 903, 130, doi: [10.3847/1538-4357/abbd92](https://doi.org/10.3847/1538-4357/abbd92)
- Shao, H., Villaescusa-Navarro, F., Genel, S., et al. 2021. <https://arxiv.org/abs/2109.04484>
- Shlomi, J., Battaglia, P., & Vlimant, J.-R. 2020, doi: [10.1088/2632-2153/abbf9a](https://doi.org/10.1088/2632-2153/abbf9a)
- Simonyan, K., Vedaldi, A., & Zisserman, A. 2013, arXiv e-prints, arXiv:1312.6034. <https://arxiv.org/abs/1312.6034>
- Sofue, Y. 2015, PASJ, 67, 75, doi: [10.1093/pasj/psv042](https://doi.org/10.1093/pasj/psv042)
- Sofue, Y., & Rubin, V. 2001, ARA&A, 39, 137, doi: [10.1146/annurev.astro.39.1.137](https://doi.org/10.1146/annurev.astro.39.1.137)
- Somerville, R. S., & Davé, R. 2015, ARA&A, 53, 51, doi: [10.1146/annurev-astro-082812-140951](https://doi.org/10.1146/annurev-astro-082812-140951)
- Springel, V. 2005, Mon. Not. Roy. Astron. Soc., 364, 1105, doi: [10.1111/j.1365-2966.2005.09655.x](https://doi.org/10.1111/j.1365-2966.2005.09655.x)
- Springel, V., White, S. D. M., Tormen, G., & Kauffmann, G. 2001, MNRAS, 328, 726, doi: [10.1046/j.1365-8711.2001.04912.x](https://doi.org/10.1046/j.1365-8711.2001.04912.x)
- Springel, V., Pakmor, R., Pillepich, A., et al. 2018, MNRAS, 475, 676, doi: [10.1093/mnras/stx3304](https://doi.org/10.1093/mnras/stx3304)
- Vaswani, A., Shazeer, N., Parmar, N., et al. 2017, arXiv e-prints, arXiv:1706.03762. <https://arxiv.org/abs/1706.03762>
- Veličković, P., Cucurull, G., Casanova, A., et al. 2018. <https://arxiv.org/abs/1710.10903>
- Villaescusa-Navarro, F., Anglés-Alcázar, D., Genel, S., et al. 2021a, ApJ, 915, 71, doi: [10.3847/1538-4357/abf7ba](https://doi.org/10.3847/1538-4357/abf7ba)
- Villaescusa-Navarro, F., Genel, S., Angles-Alcazar, D., et al. 2021b, arXiv e-prints, arXiv:2109.10915. <https://arxiv.org/abs/2109.10915>
- Villanueva-Domingo, P. 2021, HaloGraphNet, v1.0, Zenodo, doi: [10.5281/zenodo.5676528](https://doi.org/10.5281/zenodo.5676528)
- Villanueva-Domingo, P., & Villaescusa-Navarro, F. 2021, ApJ, 907, 44, doi: [10.3847/1538-4357/abd245](https://doi.org/10.3847/1538-4357/abd245)
- Villanueva-Domingo, P., Villaescusa-Navarro, F., Genel, S., et al. 2021, Weighing the Milky Way and Andromeda with Artificial Intelligence. <https://arxiv.org/abs/2111.XXXXX>
- Vogelsberger, M., Genel, S., Sijacki, D., et al. 2013, MNRAS, 436, 3031, doi: [10.1093/mnras/stt1789](https://doi.org/10.1093/mnras/stt1789)
- Vogelsberger, M., Marinacci, F., Torrey, P., & Puchwein, E. 2020, Nature Reviews Physics, 2, 42, doi: [10.1038/s42254-019-0127-2](https://doi.org/10.1038/s42254-019-0127-2)
- Vogelsberger, M., Genel, S., Springel, V., et al. 2014, Nature, 509, 177, doi: [10.1038/nature13316](https://doi.org/10.1038/nature13316)
- von Marttens, R., Casarini, L., Napolitano, N. R., et al. 2021, arXiv e-prints, arXiv:2111.01185. <https://arxiv.org/abs/2111.01185>
- Wadekar, D., Villaescusa-Navarro, F., Ho, S., & Perreault-Levasseur, L. 2020, arXiv e-prints, arXiv:2012.00111. <https://arxiv.org/abs/2012.00111>
- Wang, Y., Sun, Y., Liu, Z., et al. 2019, Dynamic Graph CNN for Learning on Point Clouds. <https://arxiv.org/abs/1801.07829>
- Wechsler, R. H., & Tinker, J. L. 2018, ARA&A, 56, 435, doi: [10.1146/annurev-astro-081817-051756](https://doi.org/10.1146/annurev-astro-081817-051756)
- Weinberger, R., Springel, V., & Pakmor, R. 2020, Astrophys. J. Suppl., 248, 32, doi: [10.3847/1538-4365/ab908c](https://doi.org/10.3847/1538-4365/ab908c)
- Weinberger, R., Springel, V., Hernquist, L., et al. 2017, MNRAS, 465, 3291, doi: [10.1093/mnras/stw2944](https://doi.org/10.1093/mnras/stw2944)
- Wojtak, R., & Mamon, G. A. 2013, MNRAS, 428, 2407, doi: [10.1093/mnras/sts203](https://doi.org/10.1093/mnras/sts203)
- Wu, Z., Pan, S., Chen, F., et al. 2019, arXiv e-prints, arXiv:1901.00596. <https://arxiv.org/abs/1901.00596>
- Yan, Z., Mead, A. J., Van Waerbeke, L., Hinshaw, G., & McCarthy, I. G. 2020, MNRAS, 499, 3445, doi: [10.1093/mnras/staa3030](https://doi.org/10.1093/mnras/staa3030)
- Young, B.-L. 2017, Frontiers of Physics, 12, 121201, doi: [10.1007/s11467-016-0583-4](https://doi.org/10.1007/s11467-016-0583-4)
- Zaheer, M., Kottur, S., Ravanbakhsh, S., et al. 2017, arXiv e-prints, arXiv:1703.06114. <https://arxiv.org/abs/1703.06114>
- Zhou, J., Cui, G., Hu, S., et al. 2018, arXiv e-prints, arXiv:1812.08434. <https://arxiv.org/abs/1812.08434>
- Zwicky, F. 1933, Helvetica Physica Acta, 6, 110