Machine-assisted Semi-Simulation Model (MSSM): Estimating Galactic Baryonic Properties from Their Dark Matter Using A Machine Trained on Hydrodynamic Simulations

Yongseok Jo¹, Ji-hoon Kim¹*

¹Center for Theoretical Physics, Department of Physics and Astronomy, Seoul National University, Seoul 08826, Korea

Last updated 2019 Feburary 1; in original form 2019 February 1

ABSTRACT

We present a pipeline to estimate baryonic properties of a galaxy inside a dark matter (DM) halo in DM-only simulations using a machine trained on high-resolution hydrodynamic simulations. As an example, we use the ILLUSTRISTNG hydrodynamic simulation of a $(75 h^{-1} Mpc)^3$ volume to train our machine to predict e.g., stellar mass and star formation rate in a galaxysized halo based purely on its DM content. An extremely randomized tree (ERT) algorithm is used together with multiple novel improvements we introduce here such as a refined error function in machine training and two-stage learning. Aided by these improvements, our model demonstrates a significantly increased accuracy in predicting baryonic properties compared to prior attempts — in other words, the machine better mimics ILLUSTRISTNG's galaxy-halo correlation. By applying our machine to the MULTIDARK-PLANCK DM-only simulation of a large $(1 h^{-1} \text{Gpc})^3$ volume, we then validate the pipeline that rapidly generates a galaxy catalogue from a DM halo catalogue using the correlations the machine found in ILLUSTRISTNG. We also compare our galaxy catalogue with the ones produced by popular semi-analytic models (SAMs). Our so-called machine-assisted semi-simulation model (MSSM) is shown to be largely compatible with SAMs, and may become a promising method to transplant the baryon physics of galaxy-scale hydrodynamic calculations onto a larger-volume DM-only run. We discuss the benefits that machine-based approaches like this entail, as well as suggestions to raise the scientific potential of such approaches.

Key words: galaxies: formation – galaxies: evolution – galaxies: statistics – cosmology: theory – cosmology:dark matter – cosmology:large-scale structure of Universe – methods: numerical – methods: analytical

1 INTRODUCTION

Years of work have been devoted by numerous researchers to the gravitational *N*-body simulations which contains only dark matter (DM) in order to describe the evolution of large scale structures (LSS) in the Universe (e.g., Boylan-Kolchin et al. 2009; Klypin, Trujillo-Gomez & Primack 2011; Angulo et al. 2012; Riebe et al. 2013; Watson et al. 2013; Skillman et al. 2014; Heitmann et al. 2015; Ishiyama et al. 2015). DM-only simulations also provide valuable insights into the spatial and velocity correlations (e.g., White et al. 1987a, 1987b; Jenkins et al. 1998), density profiles of individual halos (e.g., Navarro et al. 1997; Bullock et al. 2001a; Prada et al. 2006, 2012; Klypin et al. 2016), angular momentum profiles and shapes (e.g., Cole et al. 1996; Lemson et al. 1999; Bullock et al. 2001b; Bett et al. 2007) and halo substructures (e.g., Moore et al. 1999; Klypin et al. 2008; Madau et al. 2008).

However, gravitational dynamics alone is clearly not sufficient for understanding our Universe. Baryon physics must be taken into account via one of the two popular methods: hydrodynamic simulations, or semi-analytic models (SAMs). On the one hand, with the advent of high-performance computing units with a large amount of memories, fully hydrodynamics, high-resolution cosmological simulations have become one of the major tools in studying baryonic contributions in the Universe's evolution. Hydrodynamic simulations that treat baryon physics such as individual galaxy formation from ~Mpc scales down to ≤100 pc scales have emerged in recent years despite the expensive computational costs. Prominent examples includes ILLUSTRIS (Vogelsberger et al. 2014a, 2014b; Genel et al. 2014), ILLUSTRISTNG (Pillepich et al. 2018; Springel et al. 2018; Nelson et al. 2018a), HORIZON-AGN (Dubois et al. 2014), EAGLE (Schaye et al. 2015), ROMULUS (Tremmel et al. 2017), MU-FASA (Davé et al. 2016) and SIMBA (Davé et al. 2019). On the other hand, in SAMs and empirical models, halos from DM-only simulations are "colored" with baryons based on relatively simple physical

^{*} Co: me@jihoonkim.org

recipes (e.g., Baugh et al. 2006; Benson 2010; Croton et al. 2016; Rodriguez-Puebla et al. 2017; Cora et al. 2018; Moster et al. 2018; Behroozi et al. 2018). While SAMs inevitably require a set of tunable parameters, the computational cost of typical SAMs is much less than that of high-resolution hydrodynamic simulations. In addition, SAMs make it easy to test and appreciate the importance of physical interactions and parameters in play (Silk & Mamon 2012).

Even with the cutting-edge computing technologies that have allowed us to simulate individual galaxies with high fidelity, the contemporary computation power is insufficient to describe a larger volume of the Universe (i.e., ~Gpc scale) with detailed baryon physics resolved at ≤ 100 pc resolution. To obtain "observable" baryonic signatures populating such a large volume, combining DM-only simulations with a SAM has traditionally been the only strategy that is computationally feasible. But, now with the arrival of machine learning technology, preliminary studies have been carried out to combine DM-only simulations with machine learning algorithms such as random forest (RF) to produce galaxy catalogues (Kamdar et al. 2016b; Agarwal et al. 2018; see also Kamdar et al. 2016a).

Here, in what we call a machine-assisted semi-simulation model (MSSM), a machine - suitable for big data regression is trained to first establish correlations between DM and baryonic properties in fully hydrodynamic simulations (e.g., DM mass and stellar mass in a halo). The machine is then tested and used to estimate various baryonic properties of a DM halo (either in hydrodynamic simulations or in DM-only simulations) based purely on its DM content. A well-constructed machine can generate an extensive galaxy catalogue out of a DM-only simulation of a large volume, within a fraction of time needed for a high-resolution hydrodynamic simulation. Furthermore, this method can be one of the most promising ways to accurately transplant the baryon physics of galaxy-scale hydrodynamic calculations (e.g., ILLUSTRISTNG in a $(75 h^{-1} Mpc)^3$ volume) onto a larger-volume DM-only simulation (e.g., MULTIDARK-PLANCK in a $(1 h^{-1} \text{Gpc})^3$ volume; Klypin et al. 2016). Training the machine with a RF-type algorithm, we could also grasp the degree of contribution or "feature importance" by each of the input features (e.g., halo mass vs. halo angular momentum) in estimating a particular property (e.g., stellar mass). From the intuition gained by feature importances and by comparing the resulting catalogues with SAMs', we will be able to provide insights to improve the SAMs as well.

In this article, we first focus on improving the machine training for MSSM, and compare our machine's accuracy with a simpler baseline model's (Sections 2 and 3.1). Major improvements include: a refined error function in machine training, using historical and environmental factors of a halo as inputs, and the two-stage learning with some predicted baryonic properties as an intermediary (Sections 2.5 and 3.2). Among these, the logarithmic scaling in the error function alleviates the inaccuracy in the lower end of the predicted outputs. A scheme that "links" two machines is introduced; it uses a predicted output from one machine as an input to the next, and is found to be one of the most effective ways to enhance the MSSM's accuracy. Tested with the ILLUSTRISTNG dataset, our pipeline demonstrates a significantly increased accuracy in estimating baryonic properties than previous attempts do (Section 3). Our machine learning and application pipeline, MSSM, is shown to be largely compatible with popular SAMs when generating a galaxy catalogue using the DM-only simulation database MULTIDARK-PLANCK (Section 4).

The remainder of this paper is organized as follows. In Section 2, we explain our methodology focusing on the pipeline of our machinery and the machine learning algorithm. The pre-processing scheme of input datasets is detailed, too. In Section 3, we elaborate on how and how much our MSSM pipeline is improved when trained with the ILLUSTRISTNG dataset. Then in Section 4, we apply our machine to the MULTIDARK-PLANCK dataset, and compare our resulting galaxy catalogue with popular SAM catalogues. In Section 5, we briefly point out a few technical issues of our model, and discuss how its scientific potential could be raised. Finally we summarize and conclude the paper in Section 6.

2 METHODOLOGY

In this section, we describe the pipeline of our model and how we build and train our machine. In particular, we focus on the machine learning algorithm, and how we pre-process the input dataset to improve the machine's accuracy.

2.1 Machine Learning Overview

In our so-called MSSM, we exploit the results of fully hydrodynamic, high-resolution simulations to establish correlations or mappings — not analytic prescriptions — between DM and baryonic properties. Machine learning means training a machine for a task that typically deals with a large amount of data. If we assign two sets of data as "input" and "output", the machine by itself searches for a model and model parameters to take in the input and produce the output. In general, the more amount of data one gives, the more accurate the model becomes. The large datasets from modern cosmological simulations are thus ideal to exploit the novelty of machine learning.

In the *supervised* learning phase of our work,¹ we first divide the halo-galaxy catalogue from a large hydrodynamic simulation into a "training set" and a "test set" (see Section 2.3.1). The machine learns a structure that maps an input to an output based on example input–output pairs, i.e., the training set (e.g., DM mass and stellar mass). The machine looks for an optimized mapping by constantly evaluating the current mapping with an "error function" (or "cost function"; e.g., a widely used metric in public packages is mean square error or MSE, see Section 2.2). Based on this evaluation, the machine returns positive or negative feedback to itself. When the training is completed, one can "score" how well the machine can match the *actual* features in the simulation using the test set (see Section 2.3.2). Based on this score, one may choose to update the learning algorithm or replace it with a different method.

2.2 Chosen Machine Learning Algorithm: Extremely Randomized Tree

The public machine learning package SCIKIT-LEARN offers an easy-to-use python interface and various hyper-parameters to adjust for a chosen regressor (Pedregosa et al. 2011). We use the ExtraTreeRegressor in SCIKIT-LEARN, an extremely randomized tree algorithm (ERT; Geurts et al. 2006).² ERT is a branch of random forest (RF) algorithms which itself is a type of ensemble learning.

¹ Machine learning algorithms are divided into several categories based on the amount and type of supervision in training: supervised learning, unsupervised learning, semi-supervised learning, and reinforcement learning.
² ERT is chosen over other algorithms since the "tacked" multi-expert

² ERT is chosen over other algorithms since the "stacked" multi-expert meta-regressor finds that ERT is almost always the most successful regressor.



Figure 1. Flowchart of our machine-assisted semi-simulation model (MSSM). In the learning phase (*top panel*), we train our machine with a fully hydrodynamic simulation database that contains both dark matter (DM) and baryon data (e.g., ILLUSTRISTNG) to predict the baryonic properties ("output") based on the DM properties ("input"). In the application phase (*bottom panel*), by feeding a DM-only *N*-body simulation (e.g., MULTIDARK-PLANCK) to the trained machine, we produce a catalogue of baryonic predictions. See Section 2.3 for more information about our MSSM pipeline.

We introduce the regressor's basic concept and inner workings here to later explain the improvements we made in the machine.

At the heart of an ERT lies a "decision tree" that is constructed top-down from a root node. The tree partitions the data into subsets which contain instances of similar values; a (leaf) node generally has more than one instance. A "forest" refers to an ensemble of decision trees — i.e., a collection of trees makes a forest. Compared to a plain RF, ERT's additional randomization step arises as the tree nodes are split (i.e., the points of split are randomly chosen), which makes an ERT perform mostly faster than a plain RF.

To best split the nodes, different statistical techniques can be adopted, but a common choice is to use an error function (see Section 2.1). Often in the form of MSE, the error function helps determine the accuracy of an ERT model at each node as

$$MSE_{node} = \frac{1}{N_{node}} \sum_{i \in node} \left(y^i - y_{node} \right)^2,$$
(1)

where $y_{\text{node}} = \frac{1}{N_{\text{node}}} \sum_{i \in \text{node}} y^i$, and N_{node} is the number of instances at the node. It is important to employ an appropriate error function based on the data structure in use. MSE is the most common and widely used error function, but we note that in the reported study we choose a different metric to best serve our cosmological datasets. This will be discussed in detail in Section 3.2.1.³

One of the most salient advantages of ERT is that it is less prone to overfitting, a critical issue in machine learning. If we overtrain the machine with a dataset of often a relatively small size, the machine could end up being skewed towards the particular input– output pairs. In other words, the machine may perform well on that particular datasets with high accuracy, but may not show similar accuracies when fed with different datasets. To mitigate overfitting, ERT uses subsets and boostrap aggregating ("bagging"; see Geurts et al. 2006 for more information), and randomly splits nodes rather than looking for the least "biased" split points.⁴ This way we could reduce the "generalization error" (as opposed to a "sampling error") when the machine is applied to previously unseen data.

2.3 Flowchart of Machine-assisted Semi-Simulation Model (MSSM)

The flowchart of our MSSM, the machine learning and application pipeline, is illustrated in Figure 1. Our goal is to construct a machine to produce a galaxy catalogue by combining a DM-only N-body simulation and a machine learning technique, that is on a par with or better than catalogues made with popular SAMs. Our pipeline is divided into two main parts — (1) the learning phase: train a machine to estimate baryonic data out of DM data using a fully hydrodynamic simulation, and (2) the application phase: apply the trained machine to a DM-only simulation to produce catalogues of galactic baryonic properties.

2.3.1 Learning Phase

In the learning phase, we use only the DM-related features extracted from the ILLUSTRISTNG hydrodynamic simulation of a $(75 \ h^{-1} \text{Mpc})^3$ volume ("TNG100" in Nelson et al. 2018a; see Section 2.4.1 for more information) as input data. We take these DM features such as DM halo mass and halo velocity dispersion as inputs, and the baryonic features such as stellar mass and gas mass

³ In addition to the error function, other hyper-parameters in ERT include: maximum depth of a tree, minimum samples split, maximum number of nodes, etc. The "depth" of a decision tree refers to the distance from a root node to a farthest leaf node. The "size" of a tree is the number of all nodes.

⁴ Therefore, when compared to RF, ERT decreases the "variance" of the model, but increase its "bias". This is so-called bias–variance tradeoff. High "variance" means that the machine is *overfitted* to random noises in a particular training set. High "bias" means that the machine is *underfitted* that it only finds poor mappings between input and output data.

of the halo as desired outputs. These input-output pairs - a "training set" — is used to train the machine via ExtraTreeRegressor (Section 2.2). Note that several historical and environmental characteristics of each halo not included in the native catalogue are computed in the pre-processing step (see Section 2.5 and Table 1 for more information). During the training process, 20% of the IL-LUSTRISTNG data is spared for a test — a "test set" — to score the accuracy of the machine afterwards. Fed with the test set, the resulting machine makes a set of predicted output data (e.g., stellar masses predicted from DM masses); and, by comparing it with the actual values in the simulation (e.g., the actual stellar masses in IL-LUSTRISTNG) we "score" the machine. Common metrics for scoring the linear regression are MSE and Pearson correlation coefficient (PCC); but, in the reported study different measures are also used to evaluate the machine accuracy. We will discuss this in detail in Section 3.1.

It is also worth to mention that ERT in our MSSM not only builds a map connecting inputs and outputs, but also provides the "feature importance" that shows which input feature contributes how much to predict a particular output (e.g., which input feature is more important to predict stellar mass, halo mass or halo angular momentum?). From the feature importance we may update the set of input parameters to increase the machine's accuracy, understand the underlying physics, and potentially provide insights to improve SAMs (see also Sections 1 and 5.1).

2.3.2 Application Phase

In the application phase, the machine from the learning phase is fed with a DM-only simulation data. Here, the MULTIDARK-PLANCK DM-only simulation of a large $(1 \ h^{-1}\text{Gpc})^3$ volume is used as an input ("MDPL2" in Knebe et al. 2018; see Section 2.4.2 for more information). Needless to say, this input data needs to be preprocessed so that it is exactly in the same format and structure as the input used in the learning phase (Section 2.5). A well-optimized machine can swiftly generate a galaxy catalogue once the DM-only simulation dataset is pre-processed. In our study, the machine is able to "paint" baryonic features on ~ 10⁶ halos in a large cosmological volume in just a few tens of minutes. This is a miniscule amount of time when contrasted with what is typically needed for a highresolution hydrodynamic simulation that resolves each galaxy-size halo with proper baryon physics. In Section 6 we will discuss more on how to utilize MSSM for science.

2.4 Simulation Datasets for Machine Inputs

As noted in Section 2.3 and Figure 1, two types of simulations are considered in our MSSM pipeline — (1) in the learning phase: a fully hydrodynamic simulation is used to train a machine, and (2) in the application phase: the trained machine is applied to a DM-only simulation to produce galaxy catalogues.

2.4.1 Hydrodynamic Simulation for The Learning Phase: ILLUSTRISTNG

ILLUSTRIS (Vogelsberger et al. 2014a, 2014b) and ILLUSTRISTNG (Pillepich et al. 2018; Nelson et al. 2018a) are gravito-hydrodynamic simulations performed with a moving-mesh code AREPO (Springel 2010). Both simulations include all relevant galaxy-scale physics to follow the evolution of dark matter, cosmic gas, stars and super massive black holes (SMBHs) from z = 127 to 0, such as radiative

gas cooling (Katz et al. 1996; Wiersma et al. 2009a), star formation (Springel & Hernquist 2003b; Schaye & Dalla Vecchia 2008), stellar evolution and chemical enrichment based on stellar synthesis models (Wiersma et al. 2009b), stellar feedback (Springel & Hernquist 2003a) and SMBH and Active Galactic Nuclei (AGN) feedback (Springel et al. 2005a, 2005b). The more recent ILLUTRISTNG (The Next Generation) updates ILLUSTRIS by including magneto-hydrodynamics (Pakmor et al. 2011; Pakmor & Springel 2013), metal advection (Naiman et al. 2018), updated SMBH physics (Wienberger et al. 2017; Weinberger et al. 2018), various computational improvements (detailed in Pillepich et al. 2018), as well as updated cosmology consistent with Planck Collaboration (2016): $\Omega_{\rm m,0} = 0.3089, \Omega_{\Lambda,0} = 0.6911, \Omega_{\rm b,0} = 0.0486, \sigma_8 = 0.8159, n_{\rm s} = 0.9667, and <math display="inline">h = 0.6774.$

ILLUSTRISTNG is one of the most successful hydrodynamic calculations to date resolving individual galaxies with sophisticated baryon physics in a large enough volume. For this reason, we employ ILLUTRISTNG in the learning phase of our MSSM pipeline (Section 2.3.1). In particular, among the three different box sizes the ILLUTRISTNG database offers, the "TNG100" simulation of a $(75 h^{-1} Mpc)^3$ volume is adopted ("TNG100" dataset as designated in Nelson et al. 2018a), where 100 denotes the simulation's approximate box size in Mpc). The TNG100 run was performed at three different resolutions: TNG100-1, -2 and -3 with TNG100-1 being the highest resolution run. At z = 127, the TNG100-1 data consists of 1820³ DM particles with $m_{\text{DM}} = 7.5 \times 10^6 \,\text{M}_{\odot}$, and 1820³ hydrodynamic cells with $m_{gas} = 1.4 \times 10^6 \text{ M}_{\odot}$. At z = 0 the simulation box holds 4371211 (sub)halos identified with the friends-of-friends halo finder (FOF; Davis et al. 1985) and the SUBFIND subhalo finder (Springel et al. 2001). The publicly available halo catalogue also includes the merger trees built with the SUBLINK code (Rodriguez-Gomez et al. 2015).⁵

2.4.2 DM-only Simulation for The Application Phase: MULTIDARK-PLANCK

MULTIDARK-PLANCK (Riebe et al. 2013; Klypin et al. 2016; Rodríguez-Puebla et al. 2016) is a DM-only gravitational dynamics simulation using L-GADGET-2, a version of GADGET-2 optimized for a run with large number of particles (Springel 2005). The cosmological model adopted is consistent with Planck Collaboration (2014): $\Omega_{m,0} = 0.3071$, $\Omega_{\Lambda,0} = 0.6929$, $\Omega_{b,0} = 0.0482$, $\sigma_8 = 0.8228$, $n_8 =$ 0.96, and h = 0.6777.

In the application phase of our MSSM (Section 2.3.2), the later version of MULTIDARK-PLANCK is used as an input ("MDPL2" dataset as designated in Knebe et al. 2018). Run on a volume of $(1 \ h^{-1} \text{Gpc})^3$ that is large enough to match observational surveys, MDPL2 depicts the large-scale evolution of a significant chunk of the Universe from z = 65 to 0 using 3840³ DM particles with $m_{\text{DM}} = 1.5 \times 10^9 h^{-1} \text{M}_{\odot}$ each. The MDPL2 database publicly provides a halo catalogue for each redshift snapshot identified with the ROCKSTAR code, along with the merger trees built with the CONSISTENT TREES code (Behroozi et al. 2013).⁶

⁵ The ILLUSTRISTNG data is available at http://www.tng-project.org/.

⁶ The MULTIDARK-PLANCK data can be found in the COSMOSIM database at http://www.cosmosim.org/.

		Input Parameter	Definition	on Graphical Description		
This Work	Baseline	DM mass of a halo	Total mass of all DM particles bound to a halo	Spin		
		Velocity dispersion of a halo	Dispersion of all member particles' velocities			
		Maximum velocity of a halo	Maximum of spherically-averaged circular velocity	Mhato, Vdiso, Vmax		
		Angular momentum of a halo	Halo spin parameter			
		Number of all mergers	Number of all mergers throughout the halo's entire history			
		Number of all major mergers	Number of all mergers in which the mass ratios of the participating halos are less than 3:1			
		Last major merger mass ratio	The mass ratio of the most recent major merger along the merger tree	M m + M/m		
		Local density	The local density, $(\sum M_i)/V_{\text{box}}$, estimated for all local halos within a $(2 \text{ Mpc})^3$ volume			
		Number of local halos	Number of all local halos whose mass is larger than 80% of the target halo's mass	Target halo		
		Sum of mass over distance	Sum of mass over distance, $\sum M_i/R_i$, of all local halos within a $(2 \text{ Mpc})^3$ volume	Mi Ri Target halo		
		Maximum mass over distance	Mass over distance, $M_{\text{max}}/R_{\text{max}}$, for the most massive halo in the local volume	R _{mm} heaviest Target halo		

Table 1. All DM-related input parameters utilized to predict baryonic properties of a halo in our machine. See Section 2.5 for more information.

2.5 Pre-processing The Simulation Datasets

Data pre-processing is a pivotal step in machine learning. As noted in Figure 1, we transform the raw database — the ILLUSTRISTNG data for the learning phase, and the MULTIDARK-PLANCK data for the application phase — into a desired input format for the machine.

2.5.1 Pruning The Input Datasets

Becuase the resolutions of MULTIDARK-PLANCK data and ILLUS-TRISTNG data are different, to reconcile it we need to trim input datasets accordingly. MDPL2 dataset resolves dark matter halos down to ~ $2.23 \times 10^9 M_{\odot}$. TNG100-1 dataset resolves dark matter halos down to $7.5 \times 10^6 M_{\odot}$ while resolving baryon down to $1.4 \times 10^6 M_{\odot}$. Therefore, we exclude the halos of masses below $10^9 M_{\odot}$ in TNG100-1 to be compatible with MDPL2. In addition, since halos which do not contain star or gas are not our targets of interest, we have excluded halos whose stellar or gas mass is zero. With these cuts, the actual training set for the learning phase is reduced to ~ 3% of the original TNG100-1 halo catalogue. In Section 5.2 we demonstrate that this training set is still sufficiently large for our learning process.

6 Yongseok Jo & Ji-hoon Kim

2.5.2 Extracting Historical and Environmental Factors

The "baseline" input features to predict baryonic properties include: DM mass, velocity dispersion, and maximum circular velocity of a halo (see Table 1). This set of parameters — straight from public halo catalogues — is largely what prior attempts have used (e.g., Kamdar et al. 2016b). In addition to the baseline parameters, in the present study we aim to capture what we refer to as "historical" and "environmental" factors, and add them to the input dataset. The new features for each halo are extracted (1) from the halo's merger history, and (2) from the halo's local volume.

First, from the halo's merger tree, the following three features are obtained (Table 1): the number of all mergers, the number of all major mergers, and the mass ratio of the last major merger. These characteristics are chosen to *explicitly* quantify the evolution history of a halo imprinted in the merger tree (unlike Agarwal et al. 2018 where the merger-related parameters are implicit). Here, the mass ratio of participating halos must be less than 3:1 to be considered as major merger. Analogous to Rodriguez-Gomez et al. (2015), the mass ratio is calculated when the secondary progenitor reaches its maximum halo mass, t_{max} , before the two halos merge into one in the tree. We take this point t_{max} as the moment of merger.

Second, from the target halo's local volume of $(2 \text{ Mpc})^3$, the following four features are extracted (Table 1): the local density, the number of local halos whose masses are greater than 80% of the target halo's mass, the sum of mass over distance ("semi-potential") of all local halos $\Phi_s = \sum M_i/R_i$, and the mass over distance for the most massive local halo. These parameters aim to characterize the target halo's local environment which has likely affected how the halo has evolved. Extracting these features from the raw dataset leads to the nearest neighbor search and range search problem. It requires us to construct a *k*-d tree that partitions the space into tree structure so that neighboring halos are efficiently located.

Indeed, the *value-added* input datasets containing the additional input features improve the MSSM's accuracy for several feature predictions. This will be discussed in detail in Section 3.2.2.

3 RESULTS 1: IMPROVING A MACHINE THAT PREDICTS BARYONIC PROPERTIES

In Sections 3 and 4, we present the results of our study focusing on the learning phase and the application phase of the MSSM pipeline (Figure 1), respectively. For the rest of the paper, unless the redshift of the data is specified, we only discuss the z = 0 result. We also note that we will focus on the halos of DM masses in the range of approximately $[10^{10},10^{13.5}]\,M_{\odot}$ when presenting our results in e.g., Figures 2 - 6 (but *not* necessarily when training the machine; see Section 2.5.1). It is because (1) for halos of DM masses below $10^{10} M_{\odot}$, the resolutions of ILLUSTRISTNG (Section 3) and MULTIDARK-PLANCK simulations (Section 4) are too coarse for the machine to extract reliable mappings between DM and baryonic features, and (2) ILLUSTRISTNG contains insufficient number of halos of DM masses above $10^{13.5}$ M $_{\odot}$ due to a small simulation box size. It should be noted that the limitation here is not about our model but about the available simulations; $[10^{10},10^{13.5}]\ M_{\odot}$ is indeed also the range for which the SAMs are best optimized.

3.1 How Accurate Is The Machine's Prediction?

We first discuss how well our machine from the learning phase can predict halos' baryonic properties based purely on their DM

features. Shown in Figure 2 are normalized two-dimensional histograms comparing the predicted stellar masses ("predicted output") and the actual stellar masses in a simulation ("desired output" or "answer"), when a test set from the ILLUSTRISTNG run is used. First, shown on left is the "baseline" model that considers only mass, velocity dispersion, and maximum circular velocity of a DM halo as inputs (similar to previous studies such as Kamdar et al. 2016b; see Section 2.5.2). Shown on right is our model that improves the baseline one in various ways to be discussed in Section 3.2, including: a refined error function in machine training (Section 3.2.1), using historical and environmental factors of a halo as inputs (Sections 3.2.2 and 2.5.2), and the two-stage learning with some predicted baryonic properties as an intermediary (Section 3.2.3). We test both models to predict the following baryonic properties: gas mass, stellar mass, central black hole mass, star formation rate (SFR), metallicity, and stellar magnitudes.

Both histograms in Figure 2 are around the ideal prediction line (*black dotted line*), but in the bottom panels, the distribution is markedly tighter in our improved model resulting in the emergence of more concentrated region (*red region*) around the ideal prediction line. To quantify the machine's accuracy we first score each model with two common measures — (1) mean square error (MSE),

$$MSE = \frac{1}{N_{tot}} \sum_{i}^{N_{tot}} \left(y_{pred}^{i} - y_{TNG}^{i} \right), \qquad (2)$$

and (2) Pearson correlation coefficient (PCC),

$$PCC = \frac{cov (y_{pred}, y_{TNG})}{\sigma_{y_{pred}} \sigma_{y_{TNG}}},$$
(3)

where cov() is the covariance of two variables and σ is the standard deviation. In both equations, y_{pred}^i is the predicted *logged* output, and y_{TNG}^i is the desired *logged* output in the simulation. Note that we take the logarithm of the output data because of the similar reason described in Section 3.2.1 — except for stellar magnitudes where y_{pred}^i and y_{TNG}^i are simply the raw data (i.e., not logged).⁸ We find that both measures are significantly improved in our model: MSE decreased from 2.0×10^{-2} to 1.9×10^{-4} , and PCC increased from 0.971 to 0.987.

We have also tried — and eventually adopted — another metric to measure the machine accuracy.⁹ To compute what we call the "mean binned error" (MBE), first, the predicted and desired output pairs, $\left(y_{\text{pred}}^{i}, y_{\text{TNG}}^{i}\right)$, are binned into $\mathcal{N}_{\text{bins}}$ bins according to the y_{TNG}^{i} values. Then, in each bin, the *normalized* mean error is

$$\Gamma_j = \frac{1}{N_j} \sum_{i}^{N_j} \frac{\left| y_{\text{pred}}^i - y_{\text{TNG}}^i \right|}{y_{\text{TNG}}^i}, \qquad (4)$$

where N_i is the number of data in the *j*-th bin. Finally, by averaging

⁷ Stellar magnitudes are the luminosities of all star particles in eight photometric bands — *U*, *B*, *V*, *K*, *g*, *r*, *i*, *z* — as defined in Nelson et al. (2018b). ⁸ Unlike other baryonic properties we consider, the stellar magnitudes are already logged and lie in the range of [-25, -13]. Therefore, the improvement for MSE or PCC suggested here in Section 3.1, or the proposed improvement in Section 3.2.1 is irrelevant for stellar magnitudes.

⁹ This is inspired by the case in which MSE or PCC does not aptly represent the entire $y_{pred}^{i} - y_{TNG}^{i}$ distribution — i.e., PCC can be high even when the datapoints are widely spread out around the $y_{pred}^{i} = y_{TNG}^{i}$ line in Figure 2.



Figure 2. Normalized two-dimensional histogram comparing the *actual* stellar masses of halos in the ILLUSTRISTNG test set, $M_{\star,\text{TNG}}$, and the stellar masses predicted from input DM features of the test set, $M_{\star,\text{pred}}$. Colors indicate the normalized frequency, $n_{\text{bin}} = N_{\text{bin}}/N_{\text{tot}}$, where N_{tot} is the total number of halos and N_{bin} is the number of halos in each two-dimensional bin. Results from two machine learning models are shown: the "baseline" model similar to previous studies (*left panel*; Section 2.5.2) and our model improved for its performance (*right panel*; see Sections 2.5, 3.2 and Table 1 for more information about their differences). The black dotted line indicates an ideal prediction, $M_{\star,\text{pred}} = M_{\star,\text{TNG}}$. The marginal charts at the top and at the right show the distribution of $M_{\star,\text{TNG}}$ and $M_{\star,\text{pred}}$, respectively. See Section 3.1 for more discussion about this figure.

 Γ_i 's across all bins we obtain the MBE as

$$MBE = \frac{1}{\mathcal{N}_{\text{bins}}} \sum_{j}^{\mathcal{N}_{\text{bins}}} \Gamma_j \,. \tag{5}$$

If we replace the mean error in each bin, Γ_j , with the standard deviation in each bin, s_j , then we acquire another accuracy measure "mean binned standard deviation" (MBSD),

$$MBSD = \frac{1}{N_{bins}} \sum_{j}^{N_{bins}} s_j .$$
 (6)

We find that, in general, MBE captures the accuracy of a trained machine better than other metrics do. When predicting stellar masses, our model improves the MBE score from the baseline model's 0.0018 to 0.0013, and MBSD from 0.017 to 0.010. We will extensively use MBE and MBSD in Section 3.2 and in Table 2.

In addition to reducing the machine accuracy down to a numeric score, we also inspect the machine's performance across the output's entire value range. In Figure 3, for six baryonic properties we predict, we compare the probability distribution functions (PDFs) of the two machine learning models, and the *actual* data in the simulation.¹⁰ Again, both the baseline (*blue dot dashed lines*) and our model (*red solid lines*) predict the baryonic properties well, but in general our improved model's PDFs better match the *actual* PDFs in ILLUSTRISTNG — as demonstrated by the residual plots.

3.2 Factors That Improved Our Model

Having overviewed our machine's overall accuracy by comparing it with the *actual* data and with the baseline model, we now focus on each of the factors that improved our model. In the following subsections we explain each of three major improvements we made to our MSSM pipeline (Sections 3.2.1 - 3.2.3), followed by how we identify the best combination of these improvements that exhibits the highest accuracy (Section 3.2.4).

3.2.1 Using A Refined Error Function with Logarithmic Scaling

One of the most common choices for an error function in the machine learning algorithm — including our choice, ERT — is the MSE (see Section 2.2),

$$MSE_{node} = \frac{1}{N_{node}} \sum_{i \in node} \left(y^i - y_{node} \right)^2.$$
(1)

However, a severe problem may arise when our prediction target property has a large dynamic range (e.g., halo gas masses ranging from $10^8 M_{\odot}$ to $10^{12} M_{\odot}$). A simple mathematical argument tells that when naively used with raw *y* values, MSE could be disproportionately more sensitive to larger *y* values. For example, a small fractional error in the $10^{12} M_{\odot}$ range may completely dominate over even a very large fractional error in the $10^8 M_{\odot}$ range. This has caused the naive baseline model (Section 2.5.2) to perform poorly in the lower value range (see e.g., the left panel of Figure 2).

black dashed line in Figure 3 is slightly different from that of Figure 4, the actual halo number density in the ILLUSTRISTNG volume (TNG100-1).

¹⁰ To make the PDF in Figure 3, we sum up the test results of 5 (= 1/0.2) trials of machine learning and testing, where 0.2 is the fractional size of the ILLUSTRISTNG test set (Section 2.3.1). Then, the fractional halo numbers in each bin match the number density in the real Universe. For this reason, the



Figure 3. Probability distribution functions Φ (PDFs) of six baryonic properties — gas mass, stellar mass, central black hole mass, star formation rate (SFR), metallicity, and stellar magnitude (*g* band) — predicted from input DM features in the ILLUSTRISTNG test set. We use two machine learning models to make predictions: the baseline model (*blue dot dashed lines*; see Section 2.5.2) and our improved model (*red solid lines*; see Sections 2.5, 3.2 and Table 1 for more information about their differences). The baryon data in ILLUSTRISTNG used for the training is also shown (*black dashed lines*).¹⁰ The residuals between the predicted and the *actual* PDF, $\log \Phi_{pred} - \log \Phi_{TNG}$, are displayed in the bottom chart of each panel. Overall, our model shows improved accuracy when predicting most baryonic properties of halos. See Section 3 for more discussion about this figure.

Table 2. The mean binned error (MBE), Eq. (5), quantifying how well the machine predicts each of the six baryonic properties — gas mass, stellar mass, central black hole mass, SFR, metallicity, and stellar magnitude (g band) — based on DM features in the ILLUSTRISTNG test set. Each row indicates the MBE score within the respective x-range in Figure 3 when the machine is improved by a single improvement — except the "Best combination" row for which we identified the combination of improvements that yields the best scores for each prediction.¹² Numbers in the parentheses are mean binned standard deviation (MBSD), Eq. (6). See the referenced section in each row for details, and Section 3.2 for more discussion about this table in general.

	Gas mass	Stellar mass	BH mass	SFR	Metallicity	Stellar mag. (g)
$\mathbf{B}_{\mathrm{ascaling}}\left(\left\{ 2,5,2\right\} \right)$	0.0015	0.0018	0.0047	1.71	0.022	0.0012
Baseline (§2.3.2)	(0.023)	(0.017)	(0.020)	(36.10)	(0.099)	(0.0121)
Using an arror function with logarithmic scaling (82.2.1)	0.0010	0.0045	0.0126	1.70	0.010	_8
Using an error function with logarithmic scaling (§5.2.1)	(0.021)	(0.017)	(0.025)	(30.42)	(0.076)	(-)
Using historical and anyiranmental factors (82.2.2, 82.5.2)	0.0014	0.0014	0.0042	1.5	0.018	0.0010
Using instorical and environmental factors (§3.2.2, §2.3.2)	(0.023)	(0.016)	(0.018)	(28.27)	(0.093)	(0.0100)
Two stage learning (82.2.2)	0.0014	0.0016	0.0036	1.11	0.013	0.0005^{13}
Two-stage learning (§5.2.5)	(0.021)	(0.011)	(0.017)	(20.15)	(0.078)	(0.0064)
Best combination $(83.2.4)$	0.0010	0.0013	0.0034	1.00	0.010	0.0005
	(0.020)	(0.010)	(0.016)	(20.23)	(0.070)	(0.0053)

To amend the problem, in the learning phase, we apply *logarith*mic scaling to desired outputs of the training set (i.e., actual baryonic properties in ILLUSTRISTNG — except stellar magnitudes).⁸ Or equivalently, the y variables in the MSE error function, Eq. (1), now mean *logged* outputs, brining y values to the range of O(1). As a result, the equation is no longer heavily biased towards larger y values. Hence, our fix alleviates the inaccuracy in the lower end of the predicted outputs (see e.g., the right panel of Figure 2).¹¹ As seen in the 2nd row of Table 2 where we assemble the scores by each of the improvements, predictions such as gas mass, SFR, and metallicity benefit from the refined error function (e.g., MBE for gas mass prediction for stellar and central black hole masses do not benefit as much from the refined error function alone.

3.2.2 Using Historical and Environmental Factors

As discussed in Section 2.5.2, we extract and add "historical" and "environmental" factors to the input features when we pre-process the data for our MSSM pipeline. The newly added features are extracted (1) from the halo's merger history, and (2) from the halo's local volume, aimed at directly and indirectly capturing the halo's evolution history. The resulting *value-added* dataset includes seven additional input features such as: number of all mergers, number of all major mergers, mass ratio of the last major merger, local density, number of local halos whose masses are greater than 80% of the target halo's mass, etc. (see Section 2.5.2 for details). It improves our model's MBE and MBSD scores when predicting features like stellar mass, central black hole mass, and SFR (see the 3rd row of Table 2). For other features, including these extra factors is not as effective by itself.

3.2.3 Two-stage Learning With Stellar Magnitudes As An Intermediary

Broadly speaking, the accuracy of the ERT machine learning algorithm improves as the number of decision trees or the "size" of each tree increases (Section 2.2).³ Since the increased tree size requires exponentially more computing resources, we often need to limit the "depth" of a tree, and/or prune the nodes that are not functional. In practice, however, it is difficult to grow a large tree and prune them into an efficient shape.

Here we introduce a scheme that "links" two machines, by using a predicted output from one machine as an input to the next. The "two-stage learning" scheme works as follows. Imagine building a machine to predict SFR based only on DM features (e.g., DM mass or velocity dispersion). To increase the machine accuracy the tree must be both deep and large, requiring copious computing resources. A training set may not be informative enough for a machine to establish a meaningful direct mapping between the DM properties and SFR within a practical time limit. Instead, here we first build a machine estimating stellar magnitudes based on DM properties, then use the predicted stellar magnitudes as part of inputs to another machine estimating SFR.^{7,13} By supervising what to estimate first (stellar magnitudes) in order to predict the eventual output (SFR), we effectively "guide" the machine to build one combined, large yet efficient - ERT. Readers should note that we select stellar magnitudes as an "intermediary" because (1) the stellar magnitudes are relatively accurately predicted only from DM features and (2) the stellar magnitudes and SFR are highly correlated in the simulation data.¹⁴ Thus, we argue that in the two-stage machine training, new astrophysical information is provided to the machine by a human supervisor that the stellar magnitudes are a good intermediary between DM properties and SFR. For more discussion on how stellar magnitudes and the two-stage learning can improve the performance of MSSM, see Appendix A.

We find that the two-stage learning technique described here is one of the best strategies to construct a large and efficient ERT, and is also arguably the most effective way to improve the MSSM's

¹¹ An alternative to the logarithmic scaling could be to normalize the raw *y* values. However, the normalized variables lose their physical meanings, so the physically meaningful quantities must be carefully recovered afterwards. In contrast, logarithmic scaling does not lead to the loss of physical meaning. ¹² Each of the MBE/MBSD scores in the table is an average over 200 trials. A machine built in each trial is different due to the randomness in building an ERT, and in choosing a training set (80% of the ILLUSTRISTNG data).

¹³ To predict the *g* band, the other seven bands are used to link the machines. ¹⁴ For example, SFR is more strongly correlated with the stellar magnitudes (e.g., *g* band) than with any other DM features like DM mass. In other words, when predicting SFR, stellar magnitudes' "feature importances" dominate (> 50%; see Section 5.1) over other DM features'.

accuracy. As an example, for the SFR prediction, the two-stage learning scheme improves both MBE and MBSD scores the most when compared with any other improvements (e.g., MBE for SFR prediction decreased from 1.71 to 1.11, and MBSD from 36.10 to 20.15; see the *4th row* of Table 2).

3.2.4 Combining Improvements To Construct The Best Model

Finally, we combine all three improvements discussed above. Rather than using all the improvements at once, we have carefully tested various combinations of improvements per each of baryonic properties. This is because, when combined, one improvement may hurt the other and lead to an unexpected decrease in machine accuracy. The MBE scores for the identified best combinations are shown in the last row of Table 2. The best combinations identified here have been referred as our "improved model", and are used to produce Figures 2 - 6.

In Table 2, readers may notice that the score of a best combination is sometimes the same as that of a single improvement. For example, the MBE for a stellar magnitude prediction is 0.0005 for the best combination, but also for the two-stage learning alone. This means that the two-stage learning technique is the most important and dominant factor in improving the accuracy of stellar magnitude prediction.

4 RESULTS 2: PREDICTING BARYONIC PROPERTIES IN DARK MATTER-ONLY SIMULATIONS

We now turn to the application phase of our MSSM pipeline (Figure 1), and use the machine to estimate baryonic properties for halos in a DM-only *N*-body simulation data. The machine from Section 3 trained with the ILLUSTRISTNG data in the learning phase, is fed with the MULTIDARK-PLANCK DM-only simulation (MDPL2; see Section 2.4.2).¹⁵ The machine is asked to generate a galaxy catalogue with multiple baryonic properties — gas mass, stellar mass, central black hole mass, SFR, metallicity, and stellar magnitudes — filling the entire MULTIDARK-PLANCK volume of $(1 \ h^{-1}\text{Gpc})^{316}$.

4.1 Is The Machine-assisted Semi-Simulation Model (MSSM) Compatible With Semi-Analytic Models (SAMs)?

In Figure 4, for six baryonic properties we estimate, we compare the PDFs of our machine learning model (*red solid lines*), and of a SAM (*green dotted lines*). For a representative SAM, we utilize the MDPL2-SAG catalogue (Cora et al. 2018), one of the three SAMgenerated galaxy catalogues in the MULTIDARK-GALAXIES database (Knebe et al. 2018).¹⁷ We also add the *actual* baryon data in the ILLUSTRISTNG for comparison (TNG100-1; *black dashed lines*). Overall, we find that our MSSM and the SAM (SAG) exhibit largely compatible distribution functions. For certain properties like black hole masses, star formation rate, and stellar magnitudes, there is a sign that the MSSM mimics the distribution of ILLUSTRISTNG more closely — which is what MSSM is specifically designed to do. Yet, there are some clear mismatches due in large part to the small number statistics. For example, in the gas mass distribution, at $M_{gas} \leq 10^{9.5} M_{\odot}$, the MSSM's prediction deviates from ILLUSTRISTNG leading to a sizable gap at the lowest mass end (*1st row, left panel*). The MSSM's prediction for metallicity drops drastically at log(Metallicity) ≥ -1.8 , too (*3rd row, left panel*).

We then consider the relation between the predicted stellar mass and the halo mass, $M_{\star} - M_{\text{halo}}$, in Figure 5. This plot shows how the two halo properties are correlated on a two-dimensional plane (two-dimensional PDF). Since stellar mass is one of the properties the machine can estimate well, our MSSM prediction (red-blue contours in the upper right panel) replicates the actual M_{\star} – M_{halo} relation in the ILLUSTRISTNG run well (top left panel). As a reference, the prediction of three popular SAMs -SAG (Cora et al. 2018), SAGE (Croton et al. 2016), and GALACTICUS (Benson 2010) — are shown here as gray contours demarcating $\rho_{\text{bin, cutoff}} = 0.01 / (\log_{10} M_{\odot})^2$ (see Figure 8 of Knebe et al. 2018). Also as a reference, added to Figure 5 is the result of the baseline model (bottom right panel; see Section 2.5.2 and Table 1). Because of various improvements, our MSSM tends to perform better in the lower mass range (say, $M_{\star} < 10^{9.5} \,\mathrm{M_{\odot}}$) than the baseline model does.

As illustrated in Figures 4 and 5, we find that the MSSM pipeline can be a promising way to transplant the baryon physics of a high-resolution galaxy-scale hydrodynamic simulation (e.g., ILLUSTRISTNG) onto a larger-volume DM-only simulation (e.g., MULTIDARK-PLANCK). It is also worth noting that our machine can "paint" galaxies and their baryonic properties on a large $(1 \ h^{-1}\text{Gpc})^3$ DM-only run, within a fraction of time required for a high-resolution hydrodynamic calculation — a few tens of minutes (at most) versus a few weeks (at least).

4.2 Where The MSSM Can Be Improved

In Figure 6, we plot the probability distribution of halos on the plane of predicted stellar masses and predicted specific star formation rates (sSFR). Shown in each panel is the MDPL-SAG catalogue (black dotted contours; the outermost contour marks $\rho_{bin,\,cutoff}$ = $0.05\,/(log_{10}\,M_{\odot}log_{10}\,yr^{-1}))$ which best matches the observational data (black circles; Elbaz et al. 2011) among SAMs; see Figure 3 of Knebe et al. 2018. Notice that the ILLUSTRISTNG data itself (red-blue contours in the left panel) slightly underpredicts the Elbaz et al. 2011 data at a given stellar mass when compared with MDPL-SAG, but better matches the GALEX-SDSS-WISE Legacy Catalog (black diamonds; Salim et al. 2016). The MSSM prediction behaves in a similar way (red-blue contours in the right panel), which is again exactly what the MSSM is trained to do. However, the two-dimensional distribution of halos is narrower in machine predictions than in the original ILLUSTRISTNG data, as is indicated by the smaller error bars for the binned averages (yellow squares in the right panel). A similar tendency is found in Figure 5 as well,

¹⁵ We note that the DM halos in DM-only simulations and hydrodynamic simulations have experienced different physical processes so are inevitably different. But we also note that the so-called baryonic back-reaction effect is relatively small, justifying our use of a machine trained with hydrodynamic simulations in a different domain. For more discussion, see Appendix B. ¹⁶ The halo catalog of our Machine-assisted Semi-Simulation Model (MSSM) is available at https://sites.google.com/view/yongseok/data-access.

¹⁷ The MULTIDARK-GALAXIES data can be found in the COSMOSIM database at http://www.cosmosim.org/.



Figure 4. Probability distribution functions Φ (PDFs) of six baryonic properties predicted using a DM halo catalogue from the MULTIDARK-PLANCK database. Our improved machine trained with ILLUSTRISTNG is applied to a MULTIDARK-PLANCK dataset to make predictions (*red solid lines*; see Figure 1 and Sections 2.5, 3.2 about our improved model). We compare our prediction with a catalogue by a semi-analytic model (SAM) code SAG (*green dotted lines*; see Section 4). The *actual* baryon data in the ILLUSTRISTNG itself is also shown (TNG100-1; *black dashed lines*).¹⁰ The residuals between the predicted PDF and the simulation's PDF (ILLUSTRISTNG's), $\log \Phi_{pred} - \log \Phi_{TNG}$, are displayed in the bottom chart of each panel. Our machine-assisted semi-simulation model (MSSM) and the SAM show compatible results overall when assigning baryonic properties to halos. See Section 4.1 for more discussion about this figure.



Figure 5. Two-dimensional probability distribution of DM halo masses, M_{halo} , and predicted stellar masses, M_{\star} at z = 0. Colors indicate $\rho_{bin} = N_{bin}/(N_{tot}S_{bin})$, where N_{tot} is the total number of halos, N_{bin} is the number of halos in each two-dimensional bin, and S_{bin} is the bin area. Machines trained with ILLUSTRISTNG are applied to a MULTIDARK-PLANCK dataset to make the PDF predictions: the baseline model (*bottom right panel*) and our improved model (*top right panel*; see Sections 2.5, 3.2 and Table 1 for their differences). Yellow squares represent binned averages. The *actual* baryon data in the ILLUSTRISTNG itself is also presented (*top left panel*). Shown in each panel as gray contours are results by three popular SAMs: SAG, SAGE, and GALACTICUS. See Section 4.1 for more discussion about this figure.

where the halos are distributed in a narrower strip in MSSM predictions but not as much. When only one axis is of a predicted property (e.g., Figure 5), the two-dimensional distribution seems broader than when both x- and y-axis are of predicted properties (e.g., Figure 6).

The narrower distribution of halos likely implies reduced diversity of galaxies of same stellar masses. We suspect that when the machine is asked to predict baryonic features from DM-related features only, it may have been underfitted due to the inherently limited number of available input features. That is, there are only a very few *important* input features that decides the output, so the diversity of resulting outputs is highly restricted (more discussion in Section 5.1). This is the area where our MSSM pipeline should and can be improved in future studies (see Section 6.2).

5 DISCUSSION

In this section, we discuss two topics we find useful to appreciate our MSSM pipeline and its scientific usages.

5.1 Relative Importance of Input Features

Since our machine is built with ERT, a RF-type learning algorithm, we can easily find which input feature contributes more than the others (e.g., halo mass vs. halo angular momentum) in estimating a particular halo property (e.g., stellar mass). The degree of contribution by each of the input features is termed "feature importance". Feature importance is a relative metric among all input features adopted, and lie in the range of [0, 1]. For example, the feature importances of input parameters P_1 , P_2 , P_3 could be 0.6, 0.3, 0.1, respectively, which add up to 1.

Figure 7 shows how relative importances of input features in the baseline model (see Section 2.5.2 and Table 1) change over time when predicting two baryonic properties. At high *z*, the maximum circular velocity is the most responsible in constructing the mappings from inputs to outputs — to both stellar mass (*left panel*) and central black hole mass (*right panel*). However, at lower *z*, the halo mass and velocity dispersion take over and become more dominant. The trends robustly appear across 15 redshift snapshots from z = 7 to 0 we tested, and are highly similar for both mass predictions. At



Figure 6. Two-dimensional probability distribution of predicted stellar masses, M_{\star} , and predicted specific SFRs at z = 0. Colors indicate $\rho_{\text{bin}} = N_{\text{bin}}/(N_{\text{tot}}S_{\text{bin}})$, where N_{tot} is the total number of halos, N_{bin} is the number of halos in each two-dimensional bin, and S_{bin} is the bin area. Yellow squares represent binned averages, black diamonds represent *GALEX*-SDSS-*WISE* Legacy Catalog (GSWLC) from Salim et al. 2016 at $z \sim 0$, and black circles represent a compilation of observations from Elbaz et al. 2011 at $z \sim 0$. Our machine trained with ILLUSTRISTNG is applied to a MULTIDARK-PLANCK dataset to predict the PDF (*right panel*). The *actual* baryon data in the ILLUSTRISTNG itself is also presented (*left panel*). Shown in each panel as black dotted contours is the result by a SAM code, SAG. See Section 4.2 for more discussion about this figure.



Figure 7. Relative importances of input features — halo mass, velocity dispersion, maximum circular velocity — when the machine predicts stellar masses (*left panel*) and central black hole masses (*right panel*) based only on the three DM features of halos in ILLUSTRISTNG (i.e., baseline model; see Section 2.5.2 and Table 1). The evolution of the feature importances are plotted as functions of time. See Section 5.1 for more discussion about this figure.

z = 0, the halo mass is the most important feature in estimating both properties with features importances $\gtrsim 0.4$.

From feature importances we expect to extract physical insights about how cosmological structures have formed and evolved. We may also use features importances to evaluate how effective a new input feature is compared to preexisting ones. For example, a similar test with our improved MSSM reveals that the three input features shown in Figure 7 are still more important than most other newly introduced features in Table 1 (or see Section 2.5.2) most of the time. To raise the scientific potential of MSSM, our next goal would be to develop a set of new inputs whose feature importances are comparable to the three existing ones'.

5.2 Required Training Set Size To Build MSSM

Generally speaking, the size of a training set is one of the deciding factors in the quality of supervised learning. To check whether our TNG100-1 training set (Section 2.4.1) is sufficiently large, we measured the machine accuracy with PCC, Eq. (3), as we increase the size of the training set. In Figure 8, we see the effect of the training set size on the accuracy of the baseline model (that uses just three input features — halo mass, velocity dispersion, maximum circular velocity; see Section 2.5.2 and Table 1). Readers may notice that for all six baryonic properties we estimate, the "learning curves" reach their maximum accuracies with only a surprisingly small number of halos in the training set. For example, for stellar mass and gas mass predictions, ~ 10^3 halos are enough to yield reasonably

good estimates. For stellar magnitudes (g band) and metallicities, $\sim 10^2$ halos seem sufficient for the machine to reach its maximum potential. From the shapes of learning curves one may argue, for example, that the stellar magnitudes are highly correlated with the three input features (steep ascent to PCC ~ 1 only with ~ 10^2 halos), or that SFR is relatively hard to predict no matter how many halos are used in training (steep ascent but only to PCC ~ 0.5).

The baseline model can be well-trained up to its full potential with just $\leq 10^3$ halos, at least for the presented machine learning algorithm. Because the z = 0 training set from TNG100-1 even after aggressive data pruning (Section 2.5.1) still holds ~ 4×10^4 halos, the machine trained with TNG100-1 can be considered to have reached its maximum accuracy.¹⁸ We suspect that if the machine is built with more *important* input features (i.e., not just three features in the baseline model; see Section 5.1), a bigger training set would be needed to converge to the maximum accuracies in the learning curves. Combined with what we see in Sections 4.2 and 5.1, our experiments suggest that the machine's accuracy is limited not necessarily by the data size available for training, but more likely by the number of *important* input features. We will discuss more on potential ways to improve the machine in Section 6.2.

6 CONCLUSION

6.1 Summary

Using machine learning techniques, we have developed a pipeline to estimate baryonic properties of a galaxy based purely on DM-related features of its host halo (Section 2). Our MSSM pipeline was trained with the ILLUSTRISTNG high-resolution hydrodynamic simulation of a $(75 \ h^{-1} \text{Mpc})^3$ volume, so it can establish correlations between DM and baryonic properties (Figure 1). Compared to a simpler baseline model similar to prior studies, our machine's accuracy has been significantly improved by several improvements - such as a refined error function with logarithmic scaling in machine training, considering historical and environmental factors of a halo as inputs, and the two-stage learning with stellar magnitudes as an intermediary. Machine accuracies by each and combinations of these improvements were extensively discussed (Sections 2.5 and 3). For example, the logarithmic scaling in the error function alleviates the inaccuracy in the lower end of the predicted gas masses. The two-stage learning in which predicted stellar magnitudes from one machine is used as an input in the next, is found to be very effective in increasing the prediction accuracy for SFRs.

Once a well-trained machine is in place, in just a few tens of minutes we can rapidly populate a DM-only simulation volume that is large enough to address topics like baryonic acoustic oscillations, with galaxies having basic properties. With our MSSM mimicking ILLUSTRISTNG's galaxy-halo correlation better than previous models, we painted baryonic properties on DM halos in a $(1 \ h^{-1}\text{Gpc})^3$ volume of the MULTIDARK-PLANCK DM-only simulation (Section 4). The resulting MSSM galaxy catalogue¹⁶ is largely compatible with popular SAM catalogues. Furthermore, our MSSM has multiple scientific advantages:



Figure 8. Effect of a training set size on the machine accuracy, Pearson correlation coefficient (PCC), Eq. (3), when the machine predicts various baryonic properties (each of *six panels*) based on three DM features of halos in ILLUSTRISTNG (i.e., baseline model; see Section 2.5.2 and Table 1). The "learning curves" reach their maximum accuracies with only $\leq 10^3$ halos in the training set. See Section 5.2 for more discussion about this figure.

• (1) Within a fraction of time needed for a hydrodynamic simulation, one can efficiently transplant the baryon physics of galaxyscale hydrodynamic calculations onto a much larger volume. Readers should note that, unlike SAMs, this process does not require any recipes with fine-tuned parameters or human bias.

• (2) The ERT algorithm naturally assesses the relative importances of input features in estimating each baryonic properties (Section 5.1). The feature importance enables us to select important input features easily, and refine the machine with newly added input features with higher importance scores.

• (3) From feature importances, and by comparing the MSSM catalogue¹⁶ with SAMs', one can expect to discover physical insights in structure formation and improve the physics models in SAMs.

Despite the many improvements we made over the baseline

¹⁸ To doubly ensure that our z = 0 training set is sufficiently large, we trained a machine with all nine halo catalogues within z < 0.1. Using a ~9 times bigger training set did not significantly improve the machine accuracy, as expected by the saturated learning curves in Figure 8.

model, clearly there is room for further improvements for our MSSM framework. In the present paper, we have assumed that dark matter properties of ILLUSTRISTNG and MULTIDARK-PLANCK are largely similar that we can ignore the baryonic back reaction. But it may introduce inaccuracy in baryon-rich halos (see Appendix B for more discussion). Additionally, the lack of diversity discussed in Section 4.2 needs to be addressed by, for example, finding new input features that are better correlated with a desired output (see Section 6.2 for more discussion).

6.2 Future Work

A well-constructed machine that finds correlations between DM and baryonic contents could open up a new window to understand how our Universe has evolved. Despite important progresses we have made, immediate future projects as well as areas of improvements still remain.

• (1) The analysis in Section 4 compares only the scalar properties of galaxies generated by MSSM and SAMs. In the subsequent study, we compare the spatial distribution of MSSM galaxies with the MULTIDARK-GALAXIES data (Knebe et al. 2018).

• (2) The hyper-parameter space of ERT has not been fully explored.³ We may need to develop a more sophisticated error function for ERT to capture the diverse nature of correlations — not simply linear but complex in a multi-dimensional way — between inputs and outputs. By exploring and tuning the hyper-parameters, we may resolve the underfitting issue described in Section 4.2.

• (3) As noted in Sections 4.2 and 5, the accuracy of the proposed machine-based approach is likely limited by the small number of *important* input features. To raise the scientific potential of MSSM, we will need to find new *important* input features. For example, a set of features characterizing the merging event can be useful — not just the mass ratio, but e.g., collisional orbit parameters, infall rates, etc. These new input features will need to be extracted not from the halo catalogue or merger trees, but from a sequence of simulation snapshots finely spaced in time. One may apply a convolutional neural network to the simulation sequence itself to learn and predict baryonic properties (somewhat similar to Zhang et al. 2019).

ACKNOWLEDGMENTS

The authors thank Yun-Young Choi, Harshil Kamdar, Juhan Kim, Joel Primack, and the anonymous referee for insightful discussion and feedback on our research. Ji-hoon Kim acknowledges support by Research Start-up Fund for the new faculty of Seoul National University (SNU), and by Creative-Pioneering Researchers Program through SNU. This work was also supported by the National Institute of Supercomputing and Network/Korea Institute of Science and Technology Information with supercomputing resources including technical support, grants KSC-2018-S1-0016 and KSC-2018-CRE-0052. The CosmoSIM database used in this paper is a service by the Leibniz-Institute for Astrophysics Potsdam (AIP). The MULTIDARK database was developed in cooperation with the Spanish MultiDark Consolider Project CSD2009-00064.

References

Agarwal S., Davé R., Bassett B. A., 2018, MNRAS, 478, 3410

Angulo R. E., Springel V., White S. D. M., Jenkins A., Baugh C. M., Frenk C. S., 2012, MNRAS, 426, 2046

- Baugh C. M., 2006, Reports on Progress in Physics, 69, 3101
- Behroozi P. S., Wechsler R. H., Conroy C., 2013, ApJ, 770, 57
- Behroozi P., Wechsler R., Hearin A., Conroy C., 2018, arXiv eprints, p. arXiv:1806.07893
- Benson A. J., 2010, Phys. Rep., 495, 33
- Bett P., Eke V., Frenk C. S., Jenkins A., Helly J., Navarro J., 2007, MNRAS, 376, 215
- Boylan-Kolchin M., Springel V., White S. D. M., Jenkins A., Lemson G., 2009, MNRAS, 398, 1150
- Bruzual G., Charlot S., 2003, MNRAS, 344, 1000
- Bullock J. S., Kolatt T. S., Sigad Y., Somerville R. S., Kravtsov A. V., Klypin A. A., Primack J. R., Dekel A., 2001a, MNRAS, 321, 559
- Bullock J. S., Dekel A., Kolatt T. S., Kravtsov A. V., Klypin A. A., Porciani C., Primack J. R., 2001b, ApJ, 555, 240
- Chua K. T. E., Pillepich A., Vogelsberger M., Hernquist L., 2019, MNRAS, 484, 476
- Cole S., Lacey C., 1996, MNRAS, 281, 716
- Cora S. A., et al., 2018, MNRAS, 479, 2
- Croton D. J., et al., 2016, ApJS, 222, 22
- Cui W., Zhang Y., 2017, The Impact of Baryons on the Large-Scale Structure of the Universe. p. 7, doi:10.5772/68116
- Cui W., Borgani S., Dolag K., Murante G., Tornatore L., 2012, MNRAS, 423, 2279
- Davé R., Thompson R., Hopkins P. F., 2016, MNRAS, 462, 3265 Davé R., Anglés-Alcázar D., Narayanan D., Li Q., Rafieferantsoa
- M. H., Appleby S., 2019, arXiv e-prints, p. arXiv:1901.10203
- Davis M., Efstathiou G., Frenk C. S., White S. D. M., 1985, ApJ, 292, 371
- Dubois Y., et al., 2014, MNRAS, 444, 1453
- Duffy A. R., Schaye J., Kay S. T., Dalla Vecchia C., Battye R. A., Booth C. M., 2010, MNRAS, 405, 2161
- Elbaz D., et al., 2011, A&A, 533, A119
- Genel S., et al., 2014, MNRAS, 445, 175
- Geurts P., Ernst D., Wehenkel L., 2006, Machine Learning, 63, 3
- Heitmann K., et al., 2015, ApJS, 219, 34
- Henson M. A., Barnes D. J., Kay S. T., McCarthy I. G., Schaye J., 2017, MNRAS, 465, 3361
- Ishiyama T., Enoki M., Kobayashi M. A. R., Makiya R., Nagashima M., Oogi T., 2015, PASJ, 67, 61
- Jenkins A., et al., 1998, ApJ, 499, 20
- Kamdar H. M., Turk M. J., Brunner R. J., 2016a, MNRAS, 455, 642
- Kamdar H. M., Turk M. J., Brunner R. J., 2016b, MNRAS, 457, 1162
- Katz N., Weinberg D. H., Hernquist L., 1996, ApJS, 105, 19
- Klypin A., Gottlöber S., Kravtsov A. V., Khokhlov A. M., 1999, ApJ, 516, 530
- Klypin A. A., Trujillo-Gomez S., Primack J., 2011, ApJ, 740, 102 Klypin A., Yepes G., Gottlöber S., Prada F., Heß S., 2016, MN-RAS, 457, 4340
- Knebe A., et al., 2018, MNRAS, 474, 5206
- Lemson G., Kauffmann G., 1999, MNRAS, 302, 111
- Madau P., Diemand J., Kuhlen M., 2008, ApJ, 679, 1260
- Martizzi D., Teyssier R., Moore B., Wentz T., 2012, MNRAS, 422, 3081
- Moore B., Ghigna S., Governato F., Lake G., Quinn T., Stadel J., Tozzi P., 1999, ApJ, 524, L19
- Moster B. P., Naab T., White S. D. M., 2018, MNRAS, 477, 1822
- Naiman J. P., et al., 2018, MNRAS, 477, 1206

- Navarro J. F., Frenk C. S., White S. D. M., 1997, ApJ, 490, 493
- Nelson D., et al., 2018a, arXiv e-prints, p. arXiv:1812.05609 Nelson D., et al., 2018b, MNRAS, 475, 624
- Pakmor R., Springel V., 2013, MNRAS, 432, 176
- Pakmor R., Bauer A., Springel V., 2011, MNRAS, 418, 1392
- Pedregosa F., et al., 2011, J. Mach. Learn. Res., 12, 2825
- Pillepich A., et al., 2018, MNRAS, 473, 4077
- Planck Collaboration et al., 2014, A&A, 571, A16
- Planck Collaboration et al., 2016, A&A, 594, A13
- Prada F., Klypin A. A., Simonneau E., Betancort-Rijo J., Patiri S., Gottlöber S., Sanchez-Conde M. A., 2006, ApJ, 645, 1001
- Prada F., Klypin A. A., Cuesta A. J., Betancort-Rijo J. E., Primack J., 2012, MNRAS, 423, 3018
- Riebe K., et al., 2013, Astronomische Nachrichten, 334, 691
- Rodriguez-Gomez V., et al., 2015, MNRAS, 449, 49
- Rodríguez-Puebla A., Behroozi P., Primack J., Klypin A., Lee C., Hellinger D., 2016, MNRAS, 462, 893
- Rodríguez-Puebla A., Primack J. R., Avila-Reese V., Faber S. M., 2017, MNRAS, 470, 651
- Salim S., et al., 2016, ApJS, 227, 2
- Sawala T., Frenk C. S., Crain R. A., Jenkins A., Schaye J., Theuns T., Zavala J., 2013, MNRAS, 431, 1366
- Schaye J., Dalla Vecchia C., 2008, MNRAS, 383, 1210
- Schaye J., et al., 2015, MNRAS, 446, 521
- Silk J., Mamon G. A., 2012, Research in Astronomy and Astrophysics, 12, 917
- Skillman S. W., Warren M. S., Turk M. J., Wechsler R. H., Holz D. E., Sutter P. M., 2014, arXiv e-prints, p. arXiv:1407.2600
- Springel V., 2005, MNRAS, 364, 1105
- Springel V., 2010, MNRAS, 401, 791
- Springel V., Hernquist L., 2003a, MNRAS, 339, 289
- Springel V., Hernquist L., 2003b, MNRAS, 339, 312
- Springel V., White S. D. M., Tormen G., Kauffmann G., 2001, MNRAS, 328, 726
- Springel V., Di Matteo T., Hernquist L., 2005a, MNRAS, 361, 776
- Springel V., et al., 2005b, Nature, 435, 629
- Springel V., et al., 2008, MNRAS, 391, 1685
- Springel V., et al., 2018, in High Performance Computing in Science and Engineering '17. Springer International Publishing, pp 21–36
- Tremmel M., Karcher M., Governato F., Volonteri M., Quinn T. R., Pontzen A., Anderson L., Bellovary J., 2017, MNRAS, 470, 1121
- Vogelsberger M., et al., 2014a, MNRAS, 444, 1518
- Vogelsberger M., et al., 2014b, Nature, 509, 177
- Watson W. A., Iliev I. T., D'Aloisio A., Knebe A., Shapiro P. R., Yepes G., 2013, MNRAS, 433, 1230
- Weinberger R., et al., 2017, MNRAS, 465, 3291
- Weinberger R., et al., 2018, MNRAS, 479, 4056
- White S. D. M., Frenk C. S., Davis M., Efstathiou G., 1987a, ApJ, 313, 505
- White S. D. M., Davis M., Efstathiou G., Frenk C. S., 1987b, Nature, 330, 451
- Wiersma R. P. C., Schaye J., Smith B. D., 2009a, MNRAS, 393, 99
- Wiersma R. P. C., Schaye J., Theuns T., Dalla Vecchia C., Tornatore L., 2009b, MNRAS, 399, 574
- Zhang X., Wang Y., Zhang W., Sun Y., He S., Contardo G., Villaescusa-Navarro F., Ho S., 2019, arXiv e-prints, p. arXiv:1902.05965



Figure A1. Mean binned error (MBE), Eq. (5), of stellar mass prediction as a function of how mock band stellar magnitudes are used as an intermediary in the two-stage learning for our MSSM (see Section 3.2.3).¹² Shown on the *x*-axis are various combinations of mock band magnitudes (e.g., "zKgB" means *z*, *K*, *g*, *B* bands are used as an intermediary in machine training). The *blue dashed line* is for the sequence of combinations shown in the *upper axis*, *B* to *UirVzKgB*. The *red solid line* is for the sequence of combinations shown in the *lower axis*, *U* to *UirVzKgB*. This plot demonstrates that *U* band magnitude is the most dominant feature in predicting stellar mass. The MBE scores are for the entire stellar mass range, not for a smaller range as in Table 2. See Appendix A for more discussion about this figure.

APPENDIX A: VERIFYING STELLAR MAGNITUDES AS INFORMATION CONTAINERS

Stellar magnitudes play an important role in the two-stage learning (Section 3.2.3). As discussed, stellar magnitudes are found to be a good intermediary between e.g., DM halo mass and SFR. Typically, star particles in the simulation are convolved with a stellar population synthesis model (e.g., Bruzual & Charlot 2003) and a photometric filter to produce mock band stellar magnitudes. Therefore, one may argue that additional astrophysical information is provided to the machine as we utilize stellar magnitudes as an intermediary.

To better understand how stellar magnitudes and the two-stage learning help our MSSM to achieve better accuracy, here we evaluate if stellar magnitudes in different bands contain different information. In other words, we check if including more photometric bands improves the MSSM's accuracy. On the x-axis of Figure A1, we list combinations of mock band magnitudes used as an intermediary in the two-stage learning when predicting stellar masses. For the blue dashed line, we start with just one band, B, and add one more band at a time in the order of g, K, z, V, r, i, U (from left to right on the upper axis). This is the ascending order of feature importance among the eight band magnitudes. One can see that as we add more bands, the machine error, MBE, decreases. On the other hand, the red solid line is for the reversed order of combinations starting with U (from left to right on the *lower axis*). Since the U band magnitude has the highest feature importance, the MBE is already near its minimum only with the U band. Adding more bands does not significantly improve the machine accuracy.

Our tests reveal that for stellar mass predictions the U band is dominant; for metallicity predictions, the *i* band is. Because different band magnitudes carry different information about baryonic physics in a galaxy, we expect that including stellar magnitudes in more photometric bands would improve the MSSM's accuracy.



Figure B1. Probability distribution functions Φ (PDFs) of DM halo masses for ILLUSTRISTNG (*red solid line*) and ILLUSTRISTNG-DARK simulations (*blue dotted line*). The shift between the two lines is only less than 1%, and can be safely ignored for our purpose when applying our machines. See Appendix B for more discussion about this figure.

APPENDIX B: EFFECT OF BARYONS ON DARK MATTER HALOS (NEWLY ADDED SECTION)

In Section 4, we feed a DM-only simulation data to the machine trained with a hydrodynamic simulation data to generate a galaxy catalogue. For this to work, an implicit assumption is that DM halos from DM-only simulations and the ones from hydrodynamic simulations starting from an identical IC should have an 1-to-1 match. In hydrodynamic simulation, however, the so-called baryon backreaction may have an effect on the internal properties of a DM halo such as its shape, profile, and circular velocity (e.g., Duffy et al. 2010; Cui et al. 2012; Martizzi et al. 2012; Sawala et al. 2013; Henson et al. 2017; Chua et al. 2019) and possibly some largescale properties (e.g., Cui et al. 2017). Internal structure of DM halo can also be affected by sophisticated baryonic physics such as AGN feedback. In this study, however, we consider only the bulk properties of DM halos such as the ones in Table 1. For our MSSM to work, one of the crucial indicators to inspect would be the DM mass function of halos, not the individual internal structures. Studies have shown that the DM halo mass function of a hydrodynamic simulation including AGN feedback matches well that of a DM-only simulation (e.g., Duffy et al. 2010; Martizzi et al. 2012). Our own comparison of DM halo mass functions from ILLUSTRISTNG and ILLUSTRISTNG-DARK (DM-only run of ILLUSTRISTNG) in Figure B1 reveals high resemblance with only a slight shift (<1%). For these reasons, we have assumed that DM halos from a DM-only simulation can be used as inputs for a machine trained with a hydrodynamic simulation. Further correction and investigation on this issue remains as future work.

This paper has been typeset from a TEX/LATEX file prepared by the author.