

---

# Inferring accreted stellar mass fractions of central galaxies using random forest

---

Arshia Anjum and Sibabrata Biswal  
Group 7  
Machine Learning Course 2023

## Abstract

1 The formation and evolution of the universe, and the galaxies in it is a much  
2 researched topic, with too less labour to look deep into it. Using Machine Learning  
3 to understand this evolution through simulations is the next best option to solve  
4 the mysteries within a satisfactory error bar. One of the main topics for the same  
5 questions the amount of mass formed inside the parent halo/galaxy, compared to  
6 the amount of mass accreted by the same through its satellite halos, and it gains  
7 significance as it provides insights in the assembly history of galaxies, the merger  
8 history as well as the interactions on which galaxy evolution depends. In this report,  
9 we have used random forest (RF) as a means of studying the same and predicting  
10 the ex-situ mass fractions ( $f_{acc}$ ) of various galaxies using the TNG simulation.

## 11 1 Introduction

12 The origin and evolution of galaxies are two of the most active fields of astrophysical research. The  
13 Lambda Cold Dark Matter (or  $\Lambda$ CDM) hypothesis is the most recent manifestation of our knowledge  
14 of the origins of the Universe. It advances the big bang hypothesis by positing that most of the  
15 physical substances in the Universe is made up of a material known as dark matter. Galaxies arise  
16 in the  $\Lambda$ CDM structure creation paradigm by the cooling and condensation of gas at the centre of  
17 dark matter halos. According to the theory, galaxy formation occurs in two stages: an early rapid  
18 production of in-situ stars by gas cooling, followed by a later period of mass increase of ex-situ stars  
19 via accretion of smaller satellite galaxies. These satellite galaxies were earlier considered as the  
20 central galaxies of smaller halos. Satellite galaxies, or the subhalos as we call them after they fall  
21 into the larger halos, loose stellar mass through tidal stripping.

22 One of the reasons we can differentiate between the in-situ and ex-situ mass is that, the accreted  
23 stellar mass makes the outer regions of the parent halo, and are metal poor as compared to the in-situ  
24 mass. The next question which arises is, what is the importance of finding the ex-situ mass fraction  
25 (which will be referred as  $f_{acc}$  throughout the report). The ex-situ mass fraction of a galaxy is the  
26 fraction of its total mass that comes from accreted material, which includes gas and stars that were not  
27 originally formed within the galaxy itself, but were instead acquired through mergers or interactions  
28 with other galaxies. The ex-situ mass fraction derives its significance from the fact that it gives  
29 information on the assembly history of galaxies, the merger history as well as the interactions on  
30 which galaxy evolution depends. Also, it comments on the galaxy properties, surroundings as well  
31 as the age of the galaxy.

32 As mentioned, we use RF to study the process of stellar assembly through the TNG simulations.  
33 IllustrisTNG is a suite of large volume, cosmological, gravo-magnetohydrodynamical simulations run  
34 with the moving-mesh code AREPO. The simulation solves coupled evolution of dark matter, cosmic  
35 gas, stars, supermassive black holes, starting with the highest redshift of 127 to 0, i.e. the present day.

## 36 2 Data

### 37 2.1 The Chosen Data

38 As mentioned in the project proposal, the data used for the model is from the Illustris-TNG simulation.  
39 However, the data we were interested in was TNG-100, which was around 2TB in size, hence more  
40 rigorous to work with. Hence, we decided to write and test the code using the TNG50-4 simulation  
41 data, which is a low resolution simulation, but has the same data format as the high resolution  
42 TNG-100. Once the code is complete, we shall download the TNG-100 files, and run it using those  
43 files to get our model.

### 44 2.2 Data Extraction and Construction

45 Using an authentic API key, the data was extracted from the official website of the Illustris-TNG  
46 simulations. The way to go ahead with the procedure would be, to track the particles in the simulation  
47 at each redshift (point in time) and maintain a label for them. At each redshift, a friends-of-friends  
48 and subfind algorithm shall also be used to identify the halos and subhalos to which the particles  
49 belong. AS we keep track of them, in the present day data file, we would know the assembly history  
50 of the particle, hence be able to judge if it contributes to the in-situ mass or ex-situ mass of the galaxy.  
51 Since this would have been a more cumbersome method, we used an already existing catalogue which  
52 does all this, and gives us the final labels of the particles. Therefore, the data we are working with are:  
53 the snapshots, the group catalogues, the offsets and the stellar assembly catalogues of the simulation.

54 To construct our data, we first work with the stellar assembly catalogue. After reading the file and  
55 sorting for redshift =0, we apply our first constraint, i.e. the mass of the central galaxies of the halos  
56 should be greater than  $10^{10.16} M_{Sun}$ . This constrain exists, because (i) The resolution limit for the  
57 TNG data is around  $7.46 \times 10^8 M_{Sun}$  and (ii) the accuracy of morphology, rotation, and shape of the  
58 galaxies of interest deters below  $10^9 M_{Sun}$ . Hence, making us choose the galaxies whose mass is  
59 not less than  $10^{10} M_{Sun}/h$ . The second constraint had to be to check if the simulation gives faithful  
60 mock images for the chosen subhalos. In the context of TNG simulation, a faithful mock image refers  
61 to a computer-generated image or simulation that accurately represents a real-world phenomenon or  
62 system. There will be a need of SKIRT imaging data for the same. However, this was not possible to  
63 do here, as the data was low resolution. The third constraint was to check if the central galaxy of the  
64 halo is at least 0.5 magnitudes brighter than the satellite galaxies. However, this can be skipped, as  
65 this constraint does not have any manor impact on the model.

66 As we sort the stellar assembly data, we store the index of the subhalos (the subhalo ids). These  
67 ids are then used to get the galaxy features from the snapshot, group catalogue and offset files. The  
68 columns are then concatenated, and used for the model training.

## 69 3 Halo and Galaxy features used

70 We briefly describe the halo and galaxy features present in the data.

- 71 1. SubhaloBHMass: the estimated mass of a black hole that resides within a subhalo.
- 72 2. SubhaloGasMetalFractions: fraction of metals present in the gas component of a subhalo.  
73 (For carbon, nitrogen, oxygen, neon, magnesium, silicon, sulfur, calcium, iron, and nickel)
- 74 3. SubhaloGasMetalFractionsHalfRad: The fraction of metals present in the gas component  
75 of a subhalo within half of the subhalo's maximum circular velocity radius. (For carbon,  
76 nitrogen, oxygen, neon, magnesium, silicon, sulfur, calcium, iron, and nickel)
- 77 4. SubhaloGasMetalFractionsSfr: The fraction of metals present in the gas component of a  
78 subhalo that is actively forming stars. (For carbon, nitrogen, oxygen, neon, magnesium,  
79 silicon, sulfur, calcium, iron, and nickel)
- 80 5. SubhaloGasMetallicity: The metallicity of the gas component of a subhalo, defined as the  
81 fraction of the gas mass that is composed of heavy elements.
- 82 6. SubhaloGasMetallicityHalfRad: The metallicity of the gas component of a subhalo within  
83 half of the subhalo's maximum circular velocity radius.
- 84 7. SubhaloLen: The number of particles used to represent a subhalo in the simulation.

- 85 8. SubhaloMass: The total mass of a subhalo, including all components such as gas, stars, and  
86 dark matter.
- 87 9. SubhaloMassInHalfRad: The total mass of a subhalo within half of the subhalo’s maximum  
88 circular velocity radius.
- 89 10. SubhaloMassInRad: The total mass of a subhalo within the subhalo’s maximum circular  
90 velocity radius.
- 91 11. SubhaloSFRinHalfRad: The star formation rate within half of the subhalo’s maximum  
92 circular velocity radius.
- 93 12. SubhaloSFRinRad: The star formation rate within the subhalo’s maximum circular velocity  
94 radius.
- 95 13. SubhaloSpin: The angular momentum of a subhalo, which can affect its morphology and  
96 evolution.
- 97 14. SubhaloStarMetalFractions: The fraction of metals present in the star component of a  
98 subhalo. (For carbon, nitrogen, oxygen, neon, magnesium, silicon, sulfur, calcium, iron, and  
99 nickel)
- 100 15. SubhaloStarMetalFractionsHalfRad: The fraction of metals present in the star component  
101 of a subhalo within half of the subhalo’s maximum circular velocity radius. (For carbon,  
102 nitrogen, oxygen, neon, magnesium, silicon, sulfur, calcium, iron, and nickel)
- 103 16. SubhaloStarMetallicity: The metallicity of the star component of a subhalo, defined as the  
104 fraction of the star mass that is composed of heavy elements.
- 105 17. SubhaloStarMetallicityHalfRad: The metallicity of the star component of a subhalo within  
106 half of the subhalo’s maximum circular velocity radius.
- 107 18. SubhaloStellarPhotometrics: The properties of the stellar population in a subhalo, including  
108 luminosities and colors. (For U, B, V, K, g, r, i, z)
- 109 19. SubhaloVelDisp: The velocity dispersion of the stars in a subhalo.
- 110 20. SubhaloWindMass: The mass of gas that has been ejected from a subhalo due to feedback  
111 from star formation or black hole activity.
- 112 21. SubhaloHalfmassRad: The radius within which half of the total mass of a subhalo is  
113 contained.
- 114 22. SubhaloSFR: The star formation rate of a subhalo, measured in units of solar masses per  
115 year. This is the rate at which gas is being converted into stars within the subhalo.

116 The features related to satellite galaxies are not used, as the satellite galaxies mass limit is below the  
117 resolution of the simulation, and hence will not be a good testing criterion. Also, most of the galaxies  
118 studied do not have satellite galaxies above the range of  $10^8 M_{Sun}$ .

## 119 4 Machine Learning Methodology

### 120 4.1 Decision Trees

121 Decision trees, a type of machine learning algorithm, is used for classification and regression tasks.  
122 They work by recursively splitting the data-set into subsets based on the most informative feature,  
123 thus creating a tree-like structure. Decision trees can handle both categorical and numerical data,  
124 however, since we need numerical answers, we focus on the regression type trees. They are easy to  
125 interpret and visualize, and can handle noisy data.

### 126 4.2 Random Forest

127 Random forest, that utilizes decision trees for classification and regression, is an ensemble learning  
128 method. It works by constructing multiple decision trees using random subsets of the training data  
129 and features, and combining their predictions through averaging or voting. This helps in reducing  
130 over-fitting, increasing accuracy, and provides measures of feature importance. Random forest has  
131 several hyper-parameters such as the number of trees, the size of the subsets, and the depth of the  
132 trees. These can be tuned using cross-validation to find the optimal combination for the specific  
133 problem.

134 **4.3 Description of the Model**

135 We used RandomForestRegressor from the scikit-learn machine learning library to model the relation-  
136 ship between our input variables and the target variable. The random forest algorithm is an ensemble  
137 learning method that fits multiple decision tree models on randomly selected subsets of the data, and  
138 aggregates the predictions of each individual tree to improve overall predictive accuracy.

139 We utilized several hyperparameters to fine-tune the performance of the random forest regressor.  
140 One key hyperparameter is the number of decision trees in the forest, which we set to 100. Another  
141 important hyperparameter is the "bootstrap" setting, which determines whether each tree is fit on a  
142 bootstrapped sample of the data (with replacement) or the entire dataset. In our analysis, we set the  
143 bootstrap parameter to "True", which enables bootstrapping.

144 We also utilized the "out-of-bag" (OOB) score as a metric to evaluate the performance of our model.  
145 The OOB score measures the predictive accuracy of the model on data points that were not included  
146 in the training set for each individual tree. This provides an estimate of how well the model is likely  
147 to generalize to new, unseen data.

148 **5 Provisional Results**

149 As seen in the TNG50-4 data, the following are the provisional results:

- 150 1. The subhalos that can be used in the dataset are just 265 out of 22869 total subhalos. This is  
151 due to the low resolution data used for the training.
- 152 2. Checking the luminosity of the central galaxies compared to the satellite galaxies as a  
153 constraint, was not helpful, as it should have been.
- 154 3. The major features that we got from the model were: SubhaloMass, SubhaloMassInHalfRad,  
155 SubhaloBHMass, SubhaloStarMetalFractions, SubhaloStarMetallicity

156 Some of the relations observed in the model are expressed as a plot in Fig 1.

157 **6 Future Plans**

158 As regards our future plans in the project, we are interested in carrying out the same procedure in  
159 TNG300 and TNG50 high resolution data, as mentioned during proposal. We would focus on the  
160 specific properties in both the data, one has better statistical properties, while the other has better  
161 structural properties. Also, we will apply mass limits in the data set to split it into 2, which we could  
162 not apply in this due to low data. We will also be applying the SKIRT Imaging data constraints, and  
163 additionally focussing on the observable features alone. Lastly, our highly ambitious goal of trying a  
164 similar procedure for black hole systems also remains.

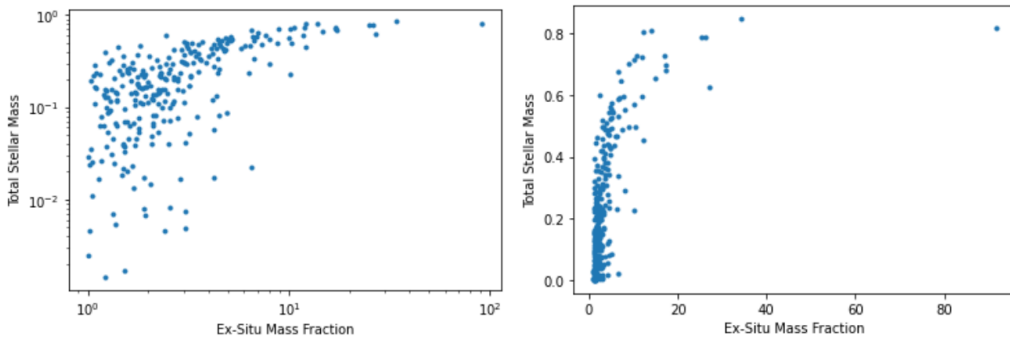


Figure 1: The Ex-Situ Stellar mass fraction's relationship with the Total Stellar mass of the subhalo (i) Log-Log Function (ii) Without Log

165 **References**

- 166 [1] R. Shi, W. Wang, Z. Li, J. Han, J. Shi, V. Rodriguez-Gomez, Y. Peng and Q. Li, *Mon. Not.*  
167 *Roy. Astron. Soc.* **515** (2022) no.3, 3938-3955 doi:10.1093/mnras/stac1541 [arXiv:2112.07203  
168 [astro-ph.GA]].
- 169 [2] V. Rodriguez-Gomez, A. Pillepich, L. V. Sales, S. Genel, M. Vogelsberger, Q. Zhu, S. Wellons,  
170 D. Nelson, P. Torrey and V. Springel, *et al.* *Mon. Not. Roy. Astron. Soc.* **458** (2016) no.3,  
171 2371-2390 doi:10.1093/mnras/stw456 [arXiv:1511.08804 [astro-ph.GA]].
- 172 [3] A. Pillepich, M. Vogelsberger, A. Deason, V. Rodriguez-Gomez, S. Genel, D. Nelson, P. Torrey,  
173 L. V. Sales, F. Marinacci and V. Springel, *et al.* *Mon. Not. Roy. Astron. Soc.* **444** (2014) no.1,  
174 237-249 doi:10.1093/mnras/stu1408 [arXiv:1406.1174 [astro-ph.GA]].
- 175 [4] D. Montenegro-Taborda, V. Rodriguez-Gomez, A. Pillepich, V. Avila-Reese, L. V. Sales,  
176 A. Rodríguez-Puebla and L. Hernquist, doi:10.1093/mnras/stad586 [arXiv:2302.10943 [astro-  
177 ph.GA]].
- 178 [5] S. Tacchella, B. Diemer, L. Hernquist, S. Genel, F. Marinacci, D. Nelson, A. Pillepich,  
179 V. Rodriguez-Gomez, L. V. Sales and V. Springel, *et al.* *Mon. Not. Roy. Astron. Soc.* **487**  
180 (2019) no.4, 5416-5440 doi:10.1093/mnras/stz1657 [arXiv:1904.12860 [astro-ph.GA]].
- 181 [6] S. Wellons, P. Torrey, C. P. Ma, V. Rodriguez-Gomez, M. Vogelsberger, M. Kriek, P. van  
182 Dokkum, E. Nelson, S. Genel and A. Pillepich, *et al.* *Mon. Not. Roy. Astron. Soc.* **449** (2015)  
183 no.1, 361-372 doi:10.1093/mnras/stv303 [arXiv:1411.0667 [astro-ph.GA]].