# Probabilistic cosmic web classification using fast-generated training data

Brandon Buncher[1]⋆ and Matias Carrasco Kind[2,3]

[1]*Department of Physics, University of Illinois, Champaign, IL 61820 USA*
[2]*Department of Astronomy, University of Illinois, Urbana, IL 61801 USA*
[3]*National Center for Supercomputing Applications, Urbana, IL 61801 USA*

10 September 2020

**ABSTRACT**

We present a novel method of robust probabilistic cosmic web particle classification in three dimensions using a supervised machine learning algorithm. Training data was generated using a simplified $\Lambda$CDM toy model with pre-determined algorithms for generating halos, filaments, and voids. While this framework is not constrained by physical modeling, it can be generated substantially more quickly than an N-body simulation without loss in classification accuracy. For each particle in this dataset, measurements were taken of the local density field magnitude and directionality. These measurements were used to train a random forest algorithm, which was used to assign class probabilities to each particle in a $\Lambda$CDM, dark matter-only N-body simulation with $256^3$ particles, as well as on another toy model data set. By comparing the trends in the ROC curves and other statistical metrics of the classes assigned to particles in each dataset using different feature sets, we demonstrate that the combination of measurements of the local density field magnitude and directionality enables accurate and consistent classification of halo, filament, and void particles in varied environments. We also show that this combination of training features ensures that the construction of our toy model does not affect classification. The use of a fully supervised algorithm allows greater control over the information deemed important for classification, preventing issues arising from hyperparameters and mode collapse in deep learning models. Due to the speed of training data generation, our method is highly scalable, making it particularly suited for classifying large datasets, including observed data.

**Key words:** cosmology: large-scale structure of the Universe – dark matter – galaxies: fundamental parameters – halos – methods: data analysis – statistical

## 1 INTRODUCTION

Large-scale structure (LSS) describes the largest scale inhomogeneities in the universe. LSS is comprised of clusters, filaments, and voids. Clusters are small, compact groups of tens to tens of thousands of particles with radii on the scale of Mpc or tens of Mpc. Typically, galaxies form around dark matter halos, large dark matter overdensities with masses on the scale of $10^{11} - 10^{15}$ $M_\odot$ (White & Rees 1978). Galaxy filaments are long, strand-like visible/dark matter overdensities that connect clusters, with lengths between 50 and 80 Mpc $h^{-1}$ (Bharadwaj et al. 2004). While filament properties are not well understood, recent data presented by Pereyra et al. (2019) and Galàrraga-Espinosa et al. (2020) has provided constraints on a variety of properties, including the density profile and total mass. Voids are visible/dark matter underdensities that fill the space between filaments, with a typical density of around $2 \times 10^{-2}$ Mpc$^{-3}$ $h^3$ (Hamaus et al. 2014). Voids are roughly

triaxial ellipsoidal, with volumes on the order of $10^4$ - $10^5$ Mpc$^{-3}$ $h^3$ (Foster & Nelson 2009).

The halo (filament) mass fraction refers to the ratio of the mass of halo (filament) particles relative to the total mass of particles. For the purposes of this paper, the equivalent metric for void particles, the fraction of particles in underdense regions, will hereafter be referred to as the void mass fraction. The relative mass fractions of halos, filaments, and voids have been studied extensively; however, there exists substantial disagreements on their values depending on the methodology used. Using $\Lambda$CDM N-body simulations, the cluster mass fraction has been estimated to be between 9 - 41% (Forero-Romero et al. 2009; Shandarin et al. 2012); the filament mass fraction has been found to be between 18 - 50% (Shandarin et al. 2012; Cautun et al. 2014); and the void mass fraction has been estimated to be between 13 - 27% (Hoffman et al. 2012; Aragon-Calvo et al. 2010a). In this work, all particles have the same mass, so the mass fraction is equivalent to the fraction of particles that are a member of a given class.

For the purposes of this work, we model the underlying dark

---

⋆ buncher2@illinois.edu

matter distribution for halos, filaments, and voids, which provides an excellent representation of galaxies; hence, particles here refer to dark matter particles.

The formation and evolution of galaxies is controlled by a variety of local properties, including the local density of dark and visible matter (Mo et al. 2010). The beginning of galaxy formation is primarily guided by the mass and density profile of the proto-galaxy's dark matter halo (Green & van den Bosch 2019). The collapse of an overdense region of dark and baryonic matter leaves a dark matter halo, a triaxial ellipsoidal dark matter overdensity. The baryonic matter that remains in the gravitational well of the dark matter halo cools, it begins to collapse into star-forming regions. While star formation depends weakly on the local density (Mo et al. 2010), the dynamics of the halo's gravitational collapse, which dictates the proto-galaxy's size and density profile, is governed primarily by this property. The local density also governs the abundance of dark matter halos and, hence, interactions between local galaxies, such as tidal stripping of matter (Mo et al. 2010; Green & van den Bosch 2019). As a result, it is crucial to understand the local dark matter density to understand how galaxies form and evolve.

A galaxy's LSS class provides substantial information about the local density. Thus, it is important to create an efficient, reliable method for determining an individual galaxy's morphological LSS class to understand problems such as these. However, a universal, deterministic algorithm is too complex to construct explicitly (Mo et al. 2010). Various classification algorithms, some of which are non-deterministic and/or non-universal, have been created; several of these are summarized in Table 1 (replicated from Libeskind et al. 2018), which we discuss in greater detail below.

Classifiers that do not use machine learning typically exploit physical or geometric properties of the structures they attempt to classify. These may be further divided into those that classify individual particles and those that determine the location and extent of individual structures. To classify particles, cluster finding algorithms frequently utilize connectedness among particles (Davis et al. 1985; Alpaslan et al. 2013) and/or rely on local geometric information such as density (Kitaura & Angulo 2012). Filament finding algorithms, however, must include information on the local and global density field, as well as some additional information that differentiates them from halos. Filament finding methods are typically geometrical (Tempel et al. 2016; Aragon-Calvo et al. 2004) or topological in nature (Aragon-Calvo et al. 2010b; Sousbie 2011; Bonnaire et al. 2019), though some graph-based methods exist (Alpaslan et al. 2013). Geometric algorithms (Cautun et al. 2012; Kitaura & Angulo 2012) have been used to classify both halos and filaments. However, these algorithms do not assign morphological class values, let alone probabilities, to individual particles. This substantially hinders the effectiveness of classification: due to difficulties in predicting fundamental properties of LSS analytically, the extent of these structures is highly dependent on arbitrary parameters. For example, many algorithms differentiate between structures using arbitrary density/scale cutoffs, either implicitly or explicitly (Libeskind et al. 2018; Tsizh et al. 2019; Falck et al. 2012) Differences between these cutoff values lead to substantial inconsistencies between the classes assigned these algorithms, such as structural mass/volume fractions and the halo mass function (HMF); a discussion of these differences can be found at Cautun et al. (2012).

Lagrangian field classifiers have been used for particle-based classification of halos, filaments, sheets, and voids (Falck et al. 2012; Lavaux & Wandelt 2010; Leclercq et al. 2017). One major drawback to Lagrangian field classifiers is that they require in-

formation about the location and velocity fields (or, equivalently, the initial and final positions of particles). However, these challenges may be circumvented using Bayesian inference methods such as BORG (Jasche & Wandelt 2013; Jasche & Kitaura 2010), a Bayesian inference method that can reconstruct the probabilistic history of the mass and velocity fields of a sample of galaxies. However, differences in the phase-space definitions for each of the structures leads to dramatic difference in class assignments, particularly for voids (Leclercq et al. 2017).

Some machine learning-based methods use deep convolutional neural networks; an example of this method may be found in (Aragon-Calvo 2018). Alternative techniques utilize supervised learning from a variety of time snapshots over the evolution of an N-body simulation (Lucie-Smith et al. 2018; Hui et al. 2018). While ML-based methods are generally more efficient than statistical classifiers, they do have substantial drawbacks. Due to the inability to trace the internal classification methodology of deep learning algorithms, these algorithms provide little information on the hallmark features of structural classes. In addition, deep learning methods are generally highly sensitive to initial hyperparameters, further inhibiting our understanding of why a particular class was chosen for a particular region/particle and potentially introducing arbitrarily-selected biases. This sensitivity may inhibit the generation of a widely-applicable algorithm with reproducible results, as small changes in the test data set may require alteration of these hyperparameters. Methods to avoid these biases have been implemented in other scenarios, such as in the fast generation of cosmic web images (Rodriguez et al. 2018); however, these methods require knowledge of the expected output, for which there is little consensus in the context of LSS classification. While supervised techniques allow greater understanding of the features utilized to determine a particle or region's class, known algorithms require training data extracted from multiple N-body simulation snapshots, which are computationally expensive. In addition, these methods are primarily designed to understand the time-evolution of halo mass distributions, which is not the goal of this project. While it may be possible to adapt these methods to classify halo particles, it would still require multiple time snapshots.

A recent topological cosmic web classifier was presented by Tsizh et al. (2019), in which the authors classified LSS particles by treating the cosmic web as a complex network. Halos were found using a friends-of-friends (FOF) algorithm and used as nodes when constructing the network, and various metrics based on particle position and velocity were used in classification. Unfortunately, classification was not successful, as demonstrated by an average confusion matrix score of 70%. A major contribution to the poor performance of this method stems from the difficulty in classifying halo particles found in large voids as this method, along with many of the topological models described in Table 1 perform classification using a relatively small range of length scales. For example, in Tsizh et al. (2019), the linking lengths used when constructing the network ranged from 1.6 - 2.4 Mpc $h^{-1}$, which fails to cover the radii of even medium-sized voids (Foster & Nelson 2009). As is noted in (Libeskind et al. 2018), local density is highly scale-dependent, indicating that density magnitude alone is ineffective when classifying halo particles (Tsizh et al. 2019; Libeskind et al. 2018). As such, a robust cosmic web classifier must take into account information beyond the local density, and must also ensure that a strict density magnitude cutoff is implicitly used when distinguishing between structures.

To simplify these classification routines, in this paper, we present an efficient ML-based classification routine that does not

fit any of the categories summarized in Table 1. Our algorithm requires substantially less information than others; in particular, we demonstrate that training using only information derived from particle positions in a single toy model generates enough information to classify a particles in a substantially larger N-body simulation. We generate training data using a toy model constructed from pre-determined structural creation algorithms which are distributed pseudo-randomly throughout a particle field. After performing measurements of the local density magnitude and density field directionality for each particle (each of which retains a "true" class inherited from its creation algorithm), we train a random forest ML algorithm to classify particles in an N-body simulation. The classes assigned are then statistically verified through measurements of known features, such as the HMF. As these measurements require only information based on the particles' current positions, we avoid requiring calculations of particle velocity, reducing the potential for error propagation due to multiple sources. Though the toy model lacks many of the physical characteristics of N-body simulation, it requires substantially less time to generate: a $256^3$-particle N-body simulation required thousands of node-hours on a multi-node cluster, while generating a toy model of the same size requires only a small fraction of this amount of time, and substantially less computational resources. We aim to show that the robustness of ML will "fill in" the information missing from the simulation, allowing accurate classification at a fraction of the cost.

To avoid the issues presented by Libeskind et al. (2018) and Tsizh et al. (2019), we propose the inclusion of local density field directionality measurements in addition to measurements of the local density magnitude. These measurements take the place of connectivity measurements (Codis et al. 2018) used in many topological and network-based classifiers. The connectivity of the filamentary skeleton of a network defines a natural filament length scale, a property that allows the prediction of local filament properties (Bardeen et al. 1986). As such, connectivity is especially effective for classification of filament particles (Codis et al. 2018; Kraljic et al. 2019). The local density field directionality provides an alternative to connectivity and other classification methods based on a particle's proximity to a ridge in a network skeleton, with the benefit that it takes into account only the properties of particles local to a given filament particle. In addition, as the directionality value is inherently normalized by the local density magnitude, it depends only on the local directionality and density field contrast. This naturally removes any implicit density magnitude cutoff for classification purposes, making it a robust method for filament classification that takes into account local environmental variations.

An additional benefit provided by supervised models is the ability to assign probabilities to class values. The classes of particles on the border of a given structure are ambiguous, with the class assignment often being due to arbitrary density cutoffs. Probabilistic classification, which is most easily achieved using a supervised model, enables us to quantify this ambiguity, providing additional information that may be correlated with other particle properties to more deeply understand how structure class is tied with other properties.

A final major benefit of this classifier is that it could be applied to observed data. Different regions of a dataset may exhibit differing field depths. To account for this, a toy model simulation could be created for each region that matches the parameters of that region, preventing the bias that would result from deep learning-based algorithms, which cannot distinguish between local field depth and density. Also, the only information required for our classifier is particle position (in particular, particle velocities are not required), making it easier to apply to observed data.

The remainder of this paper is structured as follows: in Section 2, we discuss how we generate the toy model and create training data. In Section 3, we present our class assignments for particles in an N-body simulation and a toy model. In Section 4, we demonstrate the robustness our classifier through an analysis of our results. Finally, Section 5 provides a summary of our conclusions.

Supplementary figures and demonstrations of the methodology can be found on GitHub[1].

## 2 METHODS

Unless otherwise stated, all parameters found in this section are listed in Table 2.

We aim to classify individual LSS particles using a random forest (Breiman 2001) ML algorithm trained using a fast generated data. We developed a toy model that simulated a particle field comprised of halos, filaments, and voids. Measurements of the local, global, and isotropic densities and direction fields were taken for each particle and used to train the ML algorithm. The trained algorithm was used to assign LSS class values for each particle in an N-body simulation we ran, hereafter referred to as "SIM".

SIM is a $\Lambda$CDM model simulation consisting of $N_{\text{tot, SIM}}$ collisionless dark matter particles with particle number density $n_{\text{tot}}$, each of which has a mass of $M_p$. The parameters for this simulation were taken from the WMAP+BAO+$H_0$ results found in (Komatsu et al. 2011). The cosmological parameters used were $\Omega_{m,0} = 0.272$, $\Omega_{\Lambda,0} = 0.728$, and $h = 0.704$, where the Hubble parameter $H_0 = 100\ h$ km s$^{-1}$ Mpc$^{-1}$. Initial conditions were generated using second-order Lagrangian perturbation theory (2LPT) instead of the standard Zeldovic approximation (see Crocce et al. (2006) and Scoccimarro (1998) for an explanation of this code). The primordial linear power spectrum was generated using CAMB. For this cosmology, the power spectrum was normalized using $\sigma_0 = 0.810$ and spectral index $n_s = 0.967$. As the simulation included only dark matter particles, we evolved them using the parallel tree N-body/smoothed particle hydrodynamics (SPH) code GADGET-2 (Springel 2005). Only the tree code was used for this simulation. The simulation started at redshift $z = 50$ (corresponding with scale factor $a = 0.0196$) and evolved until the scale factor $a$ reached 1. For the purposes of this work, we used a snapshot of SIM at $z = 0$.

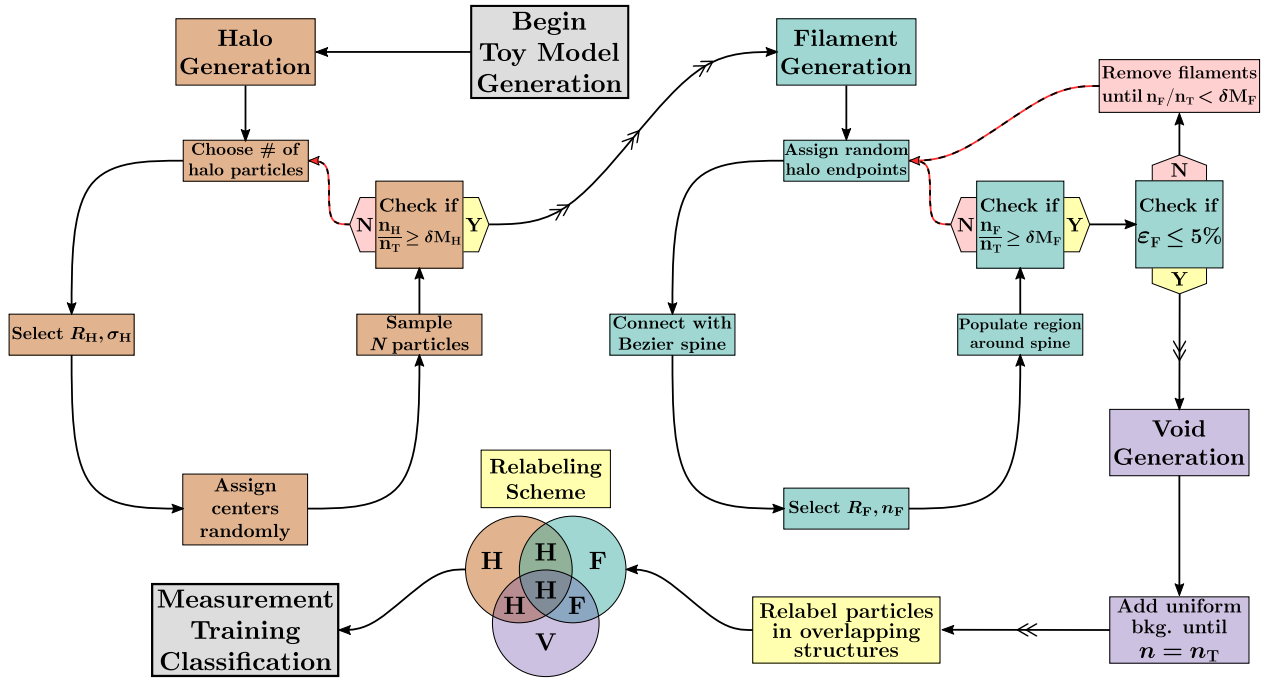### 2.1 Toy Model Simulation Generation

The method for generating a toy model dataset consisted of creation algorithms for halos, filaments, and voids. Note that the toy model only reproduces general structural features: rather than simulating the time evolution of matter due to gravity from the beginning of the universe, each structure is produced without regard to a physical creation process. While this toy model lacks the physical processes seen in N-body simulations, the generation process is substantially faster and more computationally efficient.

In the toy model, each particle's mass is defined as $M_\text{p}$, and the density of the universe $\bar{\rho} = M_\text{p}$ Mpc$^{-3}$ $h^3$, corresponding to a particle number density $n_\text{p} = 1$ Mpc$^{-3}$ $h^3$.

---

[1] https://github.com/bmbuncher/Prob-CWeb

| Method | Web types | Input | Type | Main References |
|---|---|---|---|---|
| MST | filaments | halos | Graph & Percolation | Alpaslan et al. (2013) |
| Bisous; FINE | filaments | halos | Stochastic | Tempel et al. (2016); González & Padilla (2010) |
| T-web; V-web; CLASSIC | all | particles | Hessian | Forero-Romero et al. (2009); Hoffman et al. (2012) Kitaura & Angulo (2012) |
| NEXUS+ | all | particles | Scale-Space, Hessian | Cautun et al. (2012) |
| MMF-2 | all except halos | particles | Scale-Space, Hessian | Aragon-Calvo et al. (2004); Aragon-Calvo (2014) |
| Spineweb; DisPerSE | all except halos | particles | Topological | Aragon-Calvo et al. (2010b); Sousbie (2011) |
| T-Rex | filaments | density field | Topological | Bonnaire et al. (2019) |
| ORIGAMI; MSWA; LICH | all | particles | Phase-Space | Falck et al. (2012); Falck & Neyrinck (2015) Ramachandra & Shandarin (2015); Leclercq et al. (2017) |
| DIVA | voids | particles | Phase-Space | Lavaux & Wandelt (2010) |
| BORG* | all | particles | Phase-Space | Jasche & Wandelt (2013) |
| HADES* | voids | particles | Hamiltonian Monte Carlo | Jasche & Kitaura (2010) |
| This work | all except sheets | particles | Stochastic Geometric | |

**Table 1.** An expanded overview of the methods compared in Libeskind et al. (2018); "all" indicates that the algorithm classifies particles as members of halos, filaments, voids, or sheets/walls. Algorithms marked with an asterisk indicate that these methodologies do not perform classification independently; rather, they provide a probabilistic reconstruction of the density and velocity field which, when paired with a classifier, may be used to assign class probabilities.



**Figure 1.** A diagram showing the steps used to generate a toy model dataset. Values for the listed parameters can be found in Table 2.

A diagram of the toy model creation process can be found in Figure 1

### 2.1.1   Halo Generation

The toy model simulation algorithm begins with halo generation. Halo masses were sampled from a halo mass function (Warren et al. 2006) with minimum (maximum) halo sizes $N_{min}$ and $N_{max}$ parti-

cles; this HMF model was also used to generate the test data. To select varied, visually realistic halo densities, we found an empirical probability density function (defined in Eqns. (1), (2)) that related a halo's radius to its mass. It was assumed that, for halos of a constant density $\rho_H = 200\bar{\rho}$, $R_{200} \sim M_H^{1/3}$. (Hansen et al. 2005) demonstrate through analysis of observation and simulation data that $M_H \sim N_{200}^\alpha$, where $\alpha$ is close to unity, and thus that $R_{200} \sim N_{200}^{1/3}$. We performed a similar fit on SIM using the halo mass $M$ and radii

**Datasets and Measurements**

| | |
|---|---|
| TSIM | Toy model simulation, test dataset |
| SIM | N-body simulation, test dataset |
| *VOR* | Voronoi cell volumes (density magnitude) |
| *CMD* | Distance between a particle and the center of mass of particles within a radius $R_{\rm CME}$ (density magnitude) |
| *MI* | Moment of inertia of particles within a radius of $R_{\rm CME}$ (density magnitude) |
| *ENC* | Number of particles within a radius of $R_{\rm CME}$ (density magnitude) |
| $- R_{\rm CME} \in \{0.2, 0.5, 0.8, 1.0, 1.5, 2.0, 3.0, 5.0, 7.5, 10.0, 12.0\}$ Mpc $h^{-1}$ | |
| *KNN* | Distance to the $k^{\rm th}$-nearest neighbor (density magnitude) |
| $- k \in \{1, 2, 3, 5, 8, 10, 15, 20, 25, 30, 35, 40\}$ | |
| *PCA* | Difference between the maximum and minimum explained variance ratio from a *PCA* decomposition of particles within a radius $R_{\rm PCA}$ (density field directionality) |
| $- R_{\rm PCA} \in \{1.5, 1.6, 1.7, 1.8, 1.9, 2.0\}$ Mpc $h^{-1}$ | |
| $- \sigma_{\rm PCA} = 0.45$ (resampling parameter, see Section 2.2.4) | |
| $- \delta n_{\rm PCA} = 5.5$ (resampling parameter, see Section 2.2.4) | |

**General**

| | |
|---|---|
| $M_p = 7.55 \times 10^{10}\ M_\odot$ | Particle mass in toy model and SIM |
| $n_{\rm tot} = 1.0$ Mpc$^{-3}$ | Particle number density in toy model and SIM |
| $N_{\rm tot,\ SIM} = 256^3$ | Number of particles in SIM |
| $N_{\rm tot,\ TSIM} = 85^3$ | Number of particles in TSIM |
| $L_{\rm SIM} = 256$ Mpc $h^{-1}$ | Side length of SIM particle field |
| $L_{\rm Toy} = 85$ Mpc $h^{-1}$ | Side length of training and TSIM particle fields |
| $\delta M_{\rm H} = 0.42$ | Halo mass fraction in toy model |
| $\delta M_{\rm F} = 0.38$ | Filament mass fraction in toy model |
| $\delta M_{\rm V} = 0.42$ | Void mass fraction in toy model |

**Section 2.1.1 (Halo Generation)**

| | |
|---|---|
| $N_{\rm min} = 8$ | Minimum number of particles in toy model halos |
| $N_{\rm max} = 13245$ | Maximum number of particles in toy model halos (coorresponds with a mass of $10^{11}\ M_\odot$) |
| $R_0 = 0.12$ Mpc $h^{-1}$ | Parameter in Eqn. (1) |
| $\alpha = 0.38$ | Parameter in Eqn. (1) |
| $\sigma_0 = 0.12$ Mpc $h^{-1}$ | Parameter in Eqn. (2) |
| $\beta = 0.16$ | Parameter in Eqn. (2) |

**Section 2.1.2 (Filament Generation)**

| | |
|---|---|
| $R_{\rm F,\ min} = 0.3$ Mpc $h^{-1}$ | Minimimum filament radius |
| $R_{\rm F,\ max} = 0.6$ Mpc $h^{-1}$ | Maximum filament radius |
| $B_{\rm min} = 0.65$ Mpc$^{-4}\ h^4$ | Parameter in Eqn. (3) |
| $B_{\rm max} = 1.15$ Mpc$^{-4}\ h^4$ | Parameter in Eqn. (3) |
| $\lambda_0 = 2.85$ Mpc$^{-1}\ h$ | Parameter in Eqn. (3); represents the filament number density for filaments with radius $R_{\rm F} = R_{\rm F,\ max}$ |
| $\lambda_{\rm F,\ min} \pm \delta\lambda_{\rm F,\ min} = 3.75 \pm 0.25$ Mpc$^{-1}\ h$ | The range of filament number densities for filaments with radius $R_{\rm F} = R_{\rm F,\ min}$ |

**Table 2.** A glossary of acronyms, measurement parameters, and numerical values used throughout this paper

$R(M)$ calculated by a friend-of-friend cluster finding algorithm in Ester et al. (1996) and Turk et al. (2011), determining that

$$\langle R(M)\rangle_{\rm H} = R_0 \left(\frac{M}{M_{\rm gal}}\right)^\alpha, \tag{1}$$

where $\langle R(M)\rangle_{\rm H}$ is the expected halo radius for a given mass $M$; the empirically-found value for $\alpha$ ($\alpha = 0.38$) agrees well with these prior results.

Based on empirical calculations, we assumed that the probability density function for the radii $R(M)$ for halos of a given mass $M$ followed a log-normal distribution; this assumption was based on observation of the radius histogram for halos of a given mass.

For each $M = M_0$ with at least 100 particles, we fit $R(M = M_0)$ to a log-normal histogram where the mean $\mu = \langle R(M = M_0)\rangle$. We then empirically found that the standard deviation $\sigma(M)$ also exhibited a power-law dependence on the mass:

$$\sigma(M) = \sigma_0 \left(\frac{M}{M_{\rm gal}}\right)^\beta \tag{2}$$

To create a halo, the halo mass was sampled from the HMF described by Warren et al. using algorithms from Turk et al. (2011) and Murray et al. (2013), and the halo radius was sampled from the corresponding log-normal distribution. Once a halo's mass and radius were determined, particles were generated by sampling their

radial distance from the halo's center from a truncated spherical normal spatial distribution with standard deviation $\sigma_{\mathrm{H}}(R_{\mathrm{H}}) = \frac{R_{\mathrm{H}}}{5}$.

Halo masses were sampled from the HMF until the halo mass fraction reached the desired value, i.e. $\frac{M_{\mathrm{h, tot}}}{M_{\mathrm{tot}}} \geqslant \delta_{\mathrm{h}}$. Halo centers were pseudo-randomly placed throughout the particle field of side length $L_{\mathrm{Toy}}$, then populated via the process described above.

### 2.1.2  Filament Generation

Filaments were constructed by first creating a spine, then populating the surrounding volume with particles. The spine was created by selecting two halo centers as endpoints, then creating a Bezier curve between them of degree 2 (Hermes 2017). The Bezier nodes were perturbed from the axis connecting the endpoints by $\Delta r$, were $0 \leqslant \Delta r \leqslant L_{\mathrm{F}}$, where $L_{\mathrm{F}}$ is the distance between the two endpoints. Particles were populated within a cylinder of radius $R_{\mathrm{F, min}} \leqslant R_{\mathrm{F}} \leqslant R_{\mathrm{F, max}}$ around the spine. The radial number density was calculated by sampling from a uniform distribution with maximum and minimum values

$$\lambda_{\min}(R_{\mathrm{F}}) = B_{\min}\left(\frac{R_{\mathrm{F, max}} - R_{\mathrm{F}}}{R_{\mathrm{F, max}} - R_{\mathrm{F, min}}}\right)^3 + \lambda_0 \qquad (3)$$

$$\lambda_{\max}(R_{\mathrm{F}}) = B_{\max}\left(\frac{R_{\mathrm{F, max}} - R_{\mathrm{F}}}{R_{\mathrm{F, max}} - R_{\mathrm{F, min}}}\right)^3 + \lambda_0$$

for $R_{\mathrm{F, min}} \leqslant R_{\mathrm{F}} \leqslant R_{\mathrm{F, max}}$.

This corresponded with fixing the filament radial number density $n_{\mathrm{F}}(R_{\mathrm{F, max}}) = n_0$ and allowing the density for the filaments with the minimum radius to vary as $n_{\mathrm{F}}(R_{\mathrm{F, min}}) = n_{\mathrm{F, min}} \pm \delta n_{\mathrm{F, min}}$. These values were chosen to visually match filaments seen in SIM: we approximated the minimum and maximum filament radii and number densities using several prominent filaments in SIM, then applied a bridging function (Eqn. 3) to ensure that smaller radii correlated with higher densities. We will demonstrate that these selections do not strongly affect class label assignment, so this process may be easily replicated for another test dataset.

To create a filament, 1) two endpoint halos were selected (excluding halo pairs that already have a filament generated between one another); 2) a density was selected; 3) the spine was generated; and 4) the cylindrical region around the spine was populated. These particles were placed pseudorandomly along the spine, then perturbed orthogonally from the spine using a truncated normal distribution with standard deviation $\sigma_{\mathrm{F}}(R_{\mathrm{F}}) = \frac{5}{4}R_{\mathrm{F}}$. This value was chosen so that the density at the edge of the filament was roughly $\frac{3}{4}$ the density of the center, ensuring that the filament boundary corresponded closely with the edge of the particle overdensity.

Filaments were created until the filament number density $n_{\mathrm{F}}$ exceeded the desired density $n_{\mathrm{tot}}M_{\mathrm{F}}$. To ensure that the filament number density was close to the desired density, filaments were iteratively destroyed and recreated until

$$\frac{|n_{\mathrm{F}} - n_{\mathrm{tot}}\,\delta M_{\mathrm{F}}|}{n_{\mathrm{tot}}\,\delta M_{\mathrm{F}}} \leqslant 0.05 \qquad (4)$$

Once this condition was satisfied, the filament mass fraction's deviation from $\delta M_{\mathrm{F}}$ (the desired filament mass fraction) was deemed small enough to begin background generation.

### 2.1.3  Background generation

Background (void) particles were sampled from a uniform distribution so that the number of void particles $N_{\mathrm{V, tot}} = N_{\mathrm{tot}} - N_{\mathrm{H, tot}} -$

$N_{\mathrm{F, tot}} \approx N_{\mathrm{tot}}\delta M_{\mathrm{V}}$. Note that, due to the fact that a truncated normal distribution was used to populate both halos and filaments, a sharp cutoff exists at the boundary of each halo and filament. This was a done to simplify the simulation and provide more control over its parameters; we will show that it did not affect our results.

### 2.1.4  LSS Labels

Each particle inherited an LSS class label (halo, filament, or void) from its creation algorithm; however, to prevent contamination of measurement results from these segments, particles were relabelled according to a hierarchy. Particles within the boundaries of a halo were relabelled as halo particles; any remaining particles within the boundaries of a filament were relabelled as filament particles; and the rest remained void particles.

## 2.2  Measurements

Next, measurements of the local, global, and isotropic density and direction field were taken to use as training data. We used five separate measurements of the density magnitude and one measurement of directionality. While each of them may measure similar properties, each carries different information, so a combination can improve classification accuracy and robustness. Throughout the remainder of this paper, we discuss which measurements proved most effective. All measurements were normalized such that all values lay between 0 and 1, and are described below.

### 2.2.1  Voronoi Cell Volume (VOR)

A Voronoi diagram is a method of partitioning of some multidimensional space. For each particle, there is a corresponding Voronoi cell, a region bounded by a convex polytope representing the set of all points that are closer to that point (using a Euclidean distance metric) than to any other point. We created a 3D Voronoi diagram and recorded the Voronoi cell's volume for each particle (Virtanen et al. 2019; Bradford Barber et al. 1996; Gillies et al. 2007). As a Voronoi cell's volume is closely related to the number of nearby particles, we expect that the volume of a particle's corresponding Voronoi cell will act as an effective measure of local density; in particular, we expect it to effectively classify halo particles.

### 2.2.2  Number of Particles Enclosed (ENC), Center of Mass Distance (CMD), and Moment of Inertia (MI)

Using a KD tree, the coordinates for particles within a radius $R_{\mathrm{CME}}$ of each particle were found; these values were used as a metric we will refer to as *ENC*. We found the center of mass for particles in this region, then used the distance between the center of mass and the particle of interest as a training feature. Using the same sets of particles, the moment of inertia was calculated.

For small $R_{\mathrm{CME}}$, these algorithms measure the local density, while for large $R_{\mathrm{CME}}$, they measure the global density. We expect that the information from *ENC*, *CMD*, and *MI* will be most valuable for halo classification

### 2.2.3  Distance to the k-Nearest Neighbor (KNN)

A ball tree was used to find the distance to the *k*-nearest neighbors for each particle, where the *k*-values used can be found in Table 2. For small *k*, this algorithms measures the local density, while for

large $k$, this algorithm measures the global density. We expect this algorithm to primarily influence halo classification due to its close similarity to FOF algorithms (Davis et al. 1985).

*ENC*, *CMD*, and *MI* take into account the properties of all particles within a fixed radius. As a result, they may fail to account for the spatial extent of the structure a particle is a member of. On the other hand, *KNN* measures only the properties of the environment of the closest particles. By training with very small $k$-values, we can obtain information about not only the density near a particular particle, but also of the natural length scale of the structure that particle is a member of, as the $k^{th}$ nearest neighbor for small $k$-values will likely contain only particles that are a member of that structure. As a result, we expect *KNN* measurements to provide information that cannot be obtained with the other density magnitude measurements.

### 2.2.4 *Principal Component Analysis of Local Particles* (PCA)

Principal component analysis (PCA) provides a method for determining the principal component axes, an uncorrelated orthogonal basis set such that the first component takes on the highest possible variance. Using the explained variances, this provided a method for determining the directionality of the data for use in differentiating between filaments and halos.

Prior to performing PCA analysis, the particle field was resampled to ensure an adequate number of particles were contained within each sphere of radius $R_{PCA}$ surrounding a given particle. First, a Gaussian filter with standard deviation $\sigma_{PCA}$ was applied to the particles within a given sphere to smooth the density distribution. After binning the coordinates within a given sphere, the density distribution was resampled and particles placed such that the total number density increased by a factor of $\delta n_{PCA}$. An additional uniform background was added with density $1.0$ Mpc$^{-3}$ $h^3$ to prevent the effects of background particles from being washed out.

After resampling, a PCA decomposition (Pedregosa et al. 2011; Tipping & Bishop 1998) was performed on all particles within a radius $R_{PCA}$ of each particle, and the explained variance ratio for each axis were found. The variance of particles within this region may be described by a covariance matrix $c$, with total variance $\sum_{i,j} C_{i,j} = \sigma_0^2$. After performing the PCA decomposition, $C$ undergoes the transformation $C \rightarrow C'$ such that $C'$ is diagonal and $\text{Tr}(C') = \sigma_0^2$. The explained variance of principal component axis $i$ is $C'_{ii}$, and the explained variance ratio is $\overline{\sigma}_i = \frac{C'_{ii}}{\sigma_0^2}$. After PCA decomposition is performed, a data set that may initially be correlated is transformed to a data set that exhibits no cross correlation. The explained variance ratio describes the proportion of the total variance $\sigma_0^2$ that may be attributed to the variance of particles with respect to a given axis. Intuitively, as the principal component axes correspond to the principal axes when calculating the moment of inertia, this provides information about the spread of particles about an axis such that the mass distribution around that axis is uniform.

After calculating the explained variance ratio $\overline{\sigma}_i$ for each axis for each particle, the directionality value $V$ (the value used for the *PCA* training metric) was found, where

$$V = -\ln\left(\Delta\sigma^2\right), \tag{5}$$
$$\Delta\overline{\sigma}^2 = \overline{\sigma}_{max}^2 - \overline{\sigma}_{min}^2.$$

$\Delta\sigma^2$ describes the difference between the minimum and maximum explained variance ratio. The natural logarithm of this difference was taken to accentuate the differences between the filaments

and halos so that, by using $V$ as a training metric, fewer data points would be required to perform accurate classification.

For a filament, it would be expected that the variance about the spine axis would be much smaller than the variance around the other axes due to the density field being preferentially aligned along this axis, producing a small value for $V_F$. In contrast, as the halo density field tends to exhibit very little directionality preference, it would be expected that the explained variances should vary little between the different axes, meaning that $V_H$ would be large. It may be expected that the explained variance ratios for void particles would exhibit similar properties to those of halos, and hence have a large $V$; however, due to the low density of voids, nearby structures would heavily influence the directionality values of void particles. As a result, it is expected that void particles should exhibit a small value for $V_V$, though a larger spread that of filaments.

### 2.3 Training and Class Assignment

A random forest algorithm (Breiman 2001) is a supervised learning method constructed from several decision trees. Each tree classifies a given particle using a randomized subset of features, and the class assigned to that particle is the class selected by the plurality of trees. As we aim to simplify the classifier as much as possible by minimizing the number of features required for classification, a random forest algorithm ensures that our classification is affected less by statistical fluctuations resulting from a small number of features. In addition, we may assign each particle a class probability based on the number of independent trees that assigned that particle a given class.
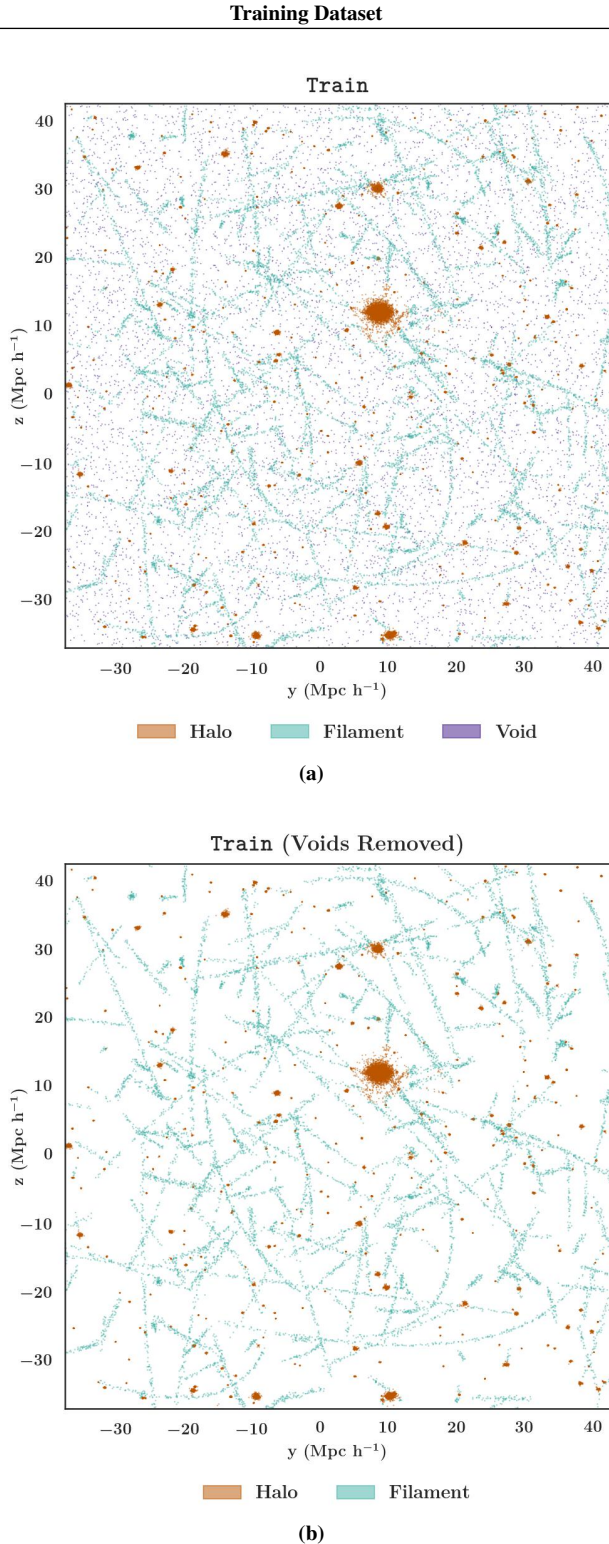
A random forest algorithm (Pedregosa et al. 2011; Breiman 2001) was trained using the measurements from one simulation. While it is typical to use multiple datasets to train a classifier, we found that additional training datasets did not influence the results substantially (likely due to the statistical robustness provided by a random forest algorithm), so only one was used. Using 200 trees, class probabilities were generated for each particle in SIM. Each particle was then assigned the class selected by the plurality of trees; in cases where multiple classes were assigned a plurality, halos were prioritized over filaments, which were prioritized over voids. An FOF clustering algorithm (Davis et al. 1985) was applied to halo particles to differentiate halos from one another, which were used to create an HMF.

## 3 CLASSIFICATION

It is expected that the primary way to differentiate between structures would be through local density magnitude calculations; however, to account for the arbitrary densities used for filaments, density field directionality measurements were included to differentiate between filaments and halos/voids. Of the methods used, only *PCA* provided directionality measurements.

### 3.1 Toy-to-SIM

First, classes were assigned to particles in SIM using a 660,000-particle toy model. A 6 Mpc thick slice of the toy model is shown in Figure 2.

**Training Dataset**



**(a)**



**(b)**

**Figure 2. (a)** A 6 Mpc thick slice of the 3D training dataset used throughout this work. This figure includes void particles (blue). **(b)** The same training data set, but with void particles removed. Removing the void particles provides a better view of the filaments. Both of these figures include the largest halo.

### 3.1.1 *Measurement Histograms*

To provide an initial guess as to which density magnitude methods would provide the most information, histograms were created using measurements of the training dataset. After performing all measurements on the training data set, each measurement was normalized so that all values fell between 0 and 1, ensuring that all measurement methods would be treated equally when training. Measurements that provide the most information for use in classification should show little overlap between the measurements on each structure and exhibit large peaks distinct from one another. Examples of some measurement histograms are displayed in Figure 3.

Figure 3 provides clues as to which measurements will be most effective for classification. The histograms for 3a *CMD* and 3b *VOR* both exhibit a large spike at the origin, indicating that, for each class, the distance from nearly every particle to the center of mass was very small, save for several very large outliers. Due to the lack of differentiation, it appears that *CMD* will provide little information that may be used to differentiate between LSS classes. Similarly, we expect *VOR* to be ineffective at differentiating halos from filaments. Though not shown, the measurement histograms for *MI* and *ENC* appear similar to those of *CMD*.

However, the measurements for 3c *KNN* exhibited much more differentiation between each of the different structure classes. As expected based on local density, halos exhibited the lowest distance to the $8^{th}$-nearest neighbor, followed by filaments, and finally voids. The larger spread on filaments and voids reflects a large chance of contamination by nearby structures due to their low density. As the classes are strongly differentiated from one another, we expect *KNN* to be an effective proxy for local density, and, hence, an effective metric for differentiating between all structures.

While not shown, larger $k$-values led to less differentiation between the classes due to each class exhibiting a greater spread.
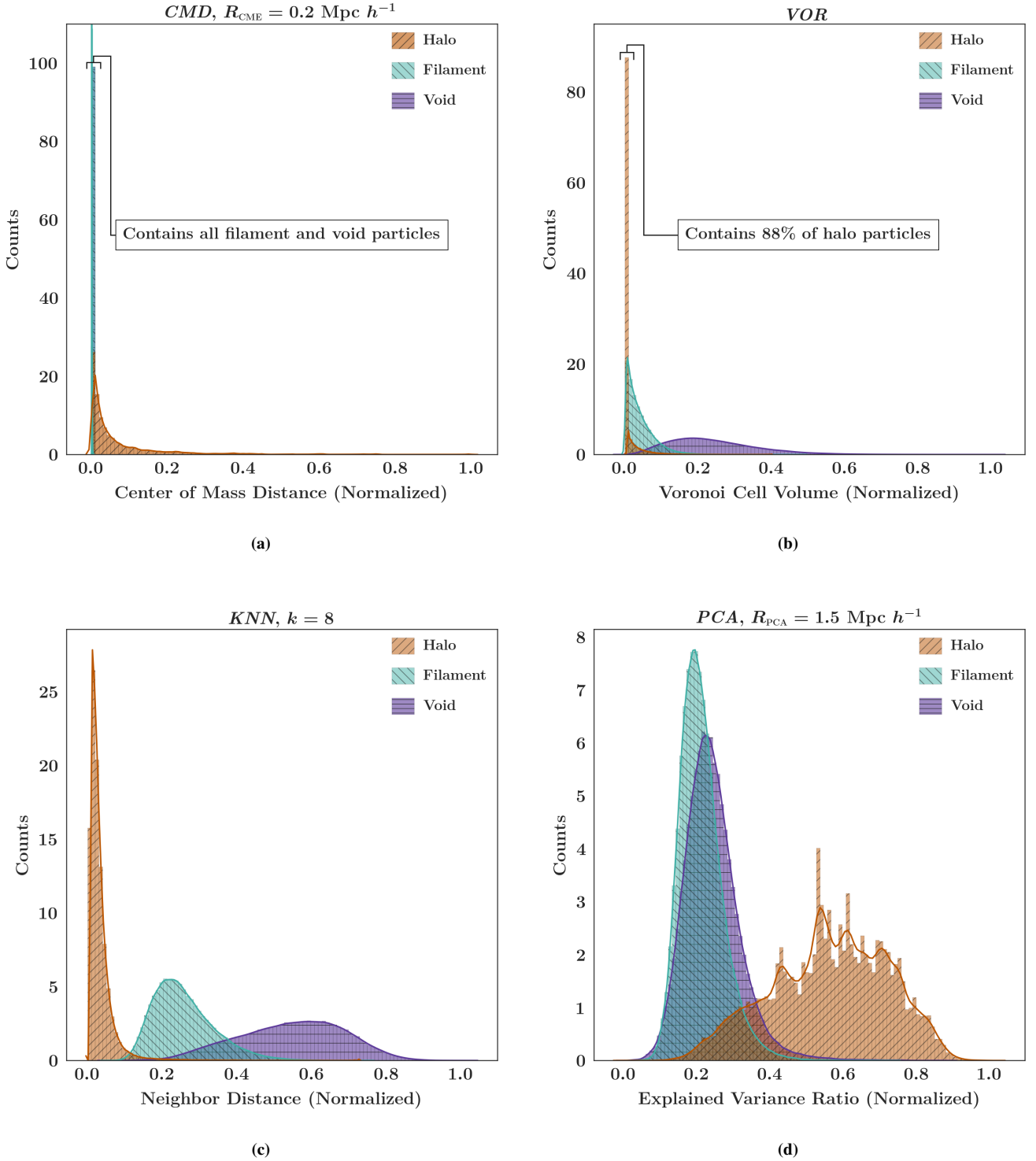
The measurement histograms for 3d *PCA* show strong peaks for filaments and voids, yet a multi-peaked halo distribution with large spread. The strong peak for filaments, especially compared to the poorly-defined halo distribution, indicates that these calculations are effectively measuring the density field directionality, as the roughly spherical halos are not expected to exhibit substantial directionality. The strong peak for voids is likely due to contamination by nearby structures: as voids have very low density, any particles from adjacent halos or filaments would lead to a strong directionality. These results bode well for the use of *PCA* in tandem with *KNN* to create a robust classifier, as *PCA* provides a natural way to differentiate between halos and filaments independent of the local density magnitude. While it may seem that these calculations may cause difficulty when distinguishing between filaments and voids, the strong differentiation between these structures from *KNN* calculations is expected to prevent this issue.

### 3.1.2 SIM *Class Assignment*

Next, class labels were assigned to particles in SIM (described in Section 2). In order to achieve probabilistic classification, each particle was classified using 200 trees, each of which independently assigned that particle a class. The classification probability for a given class was the fraction of estimators that assigned a particular particle that class over all estimators. To compare the class assignments for particles labeled as halos or filaments, define the probability contrast $\overline{\Delta P_i}$ for a halo or filament particle as
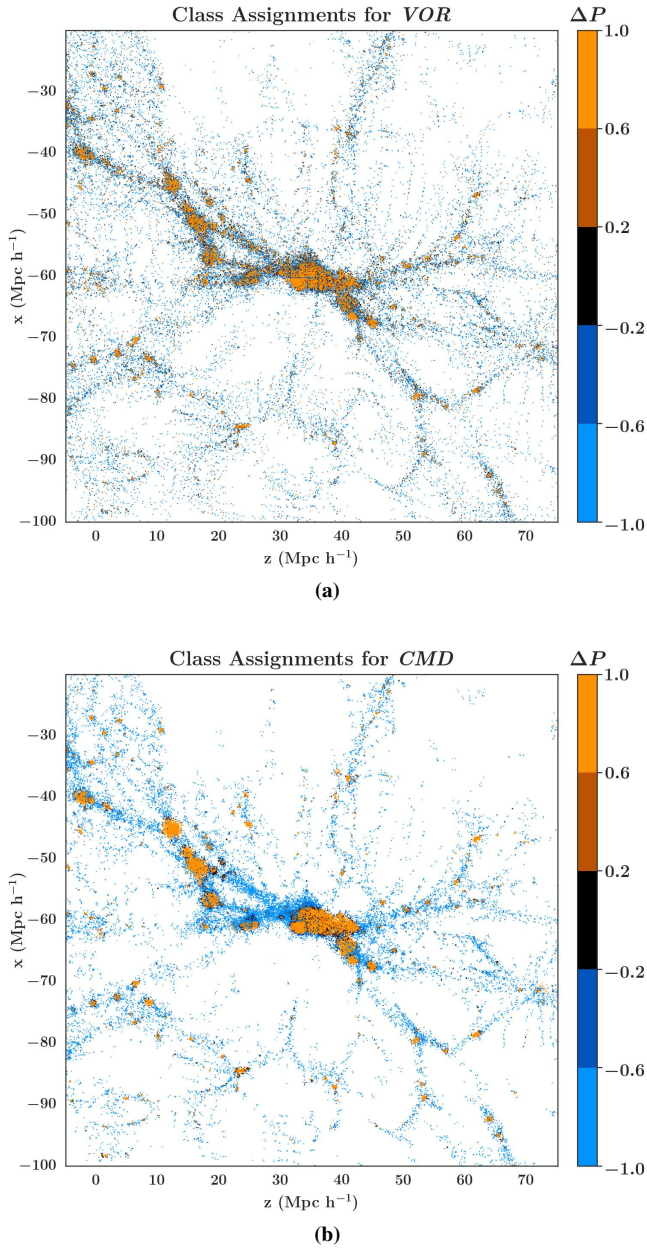
**Measurement Histograms from Training Data (Toy Model)**



**Figure 3.** Measurement histograms from the true values in the toy model for **(a)** *CMD*, **(b)** *VOR*, **(c)** *KNN*, and **(d)** *PCA*. **(a)** and **(b)** show features with less distinguishable distributions, while **(c)** and **(d)** are more distinguishable.

**Toy-to-SIM Probability Contrasts for *CMD* and *VOR***



**(a)**



**(b)**

**Figure 4.** Classification results for SIM; for comparison with TSIM, the full particle field, which spans a box with corners at $(x, y, z) = (\pm 128, \pm 128, \pm 128)$ Mpc $h^{-1}$ has been restricted to a region centered on the largest halo with side lengths 80 Mpc $h^{-1}$ and a depth of 7 Mpc $h^{-1}$. When trained using all features, this halo has $N = 38798$ particles and is centered at $(x, y, z) = (-60.6, 88.4, 35.0)$. **(a)** and **(b)** show the halo and filament class assignments made using *VOR* and *CMD* measurements, respectively, colored based on each particle's probability contrast.

$$\overline{\Delta P}_i = \frac{P_i(H) - P_i(F)}{P_i(H) + P_i(F)}, \tag{6}$$

where $P_i(H)$ ($P_i(F)$) is the probability that particle $i$ is a halo (filament).

Plots of the halo and filament particle distribution, colored by $\overline{\Delta P}$, can be found in Figures 4 and 5. A probability contrast of 1

indicates that every estimator classified that particle as a member of a halo, while a probability contrast of -1 indicates that every estimator classified that particle as a member of a filament. A particle colored black indicates that the estimator was unable to precisely differentiate between the particle's class, indicating that there is ambiguity as to whether it is a halo or filament member. It would be expected that, near the high-density center of halos, the probability contrast would be close to 1, while near the edges, especially where the halo connected to a filament, the probability contrast would be closer to 0. Note that void particles are not displayed in these plots as void probability assignments were generally close to unity.

All particles in SIM were classified by each metric set as detailed below. Figures 4 and 5 show the halo and filament class probability contrasts assigned to particles within a 7 Mpc thick slice of this 3D simulation. While the classes assigned by each feature set shown are visually realistic, we discuss their differences below.

Figures 4a and 4b show the class probabilities assigned by a classifier trained using only *VOR* and *CMD*, respectively. Relative to the assignments by *CMD*, *VOR* overestimated the number of halo and filament particles, indicating that it was not sensitive to the low density void regions. In addition, the class probabilities assigned were generally lower, indicating that *VOR* alone did not provide enough information to distinguish between the classes easily. On the other hand, *CMD* generally assigned very high class probabilities to each particle except for several small, concentrated regions on the border between a halo and filament (appearing as black clumps in Figure 4). These regions had a density magnitude between that of the high-probability halo regions ($P \approx 1$) and the high-probability filament regions ($P \approx -1$). The small width of these regions implies that using only *CMD* calculations imposed a strict density field magnitude cutoff when determining halo and filament membership.
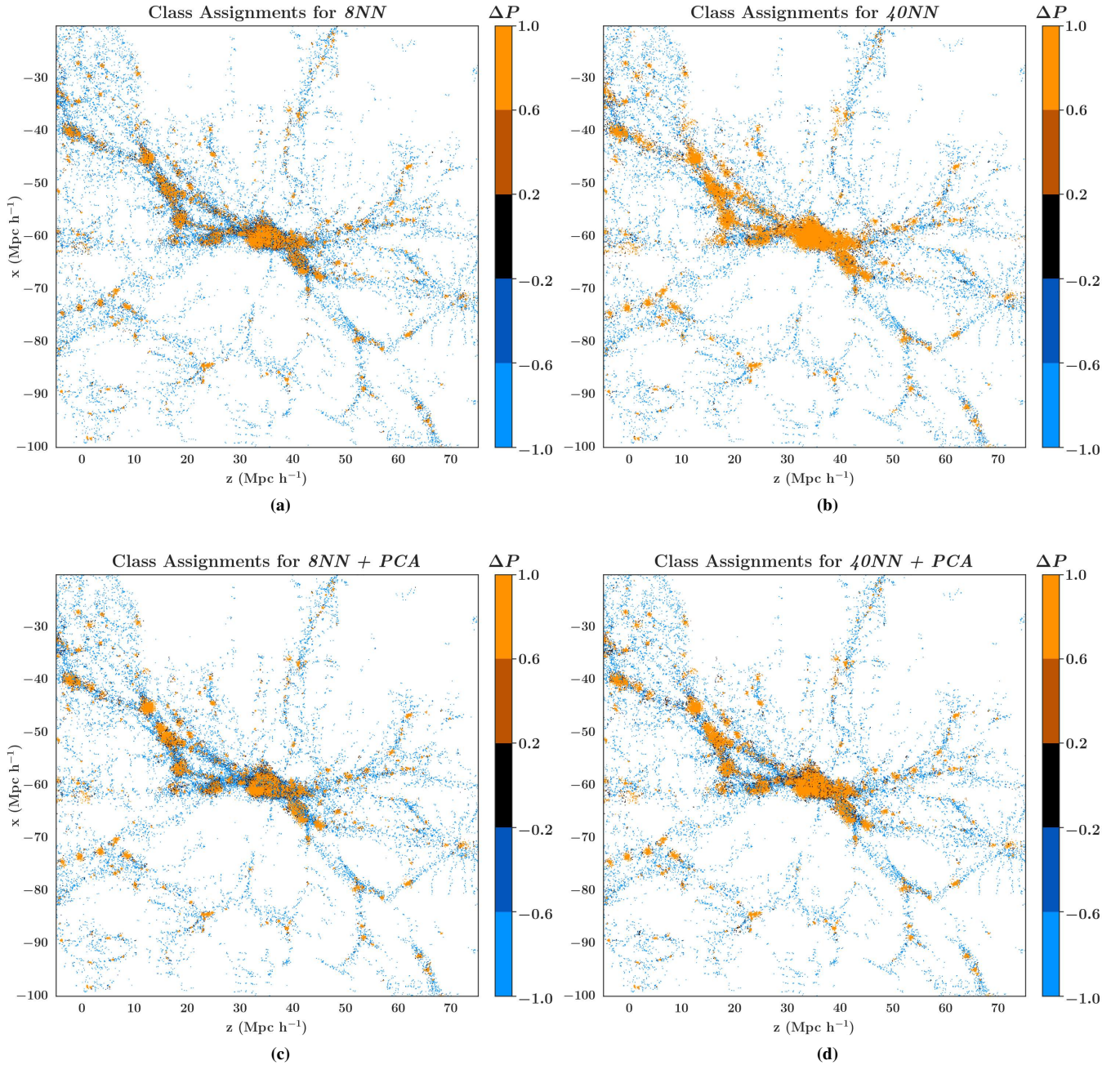
While not shown, the class probability contrast plots for *MI* and *ENC* were both very similar to *CMD*.

Figures 5a and 5b show the class probabilities assigned by a classifier trained using *KNN* for $k \leqslant 8$, 40, respectively. The assignments made by these classifiers are generally similar; however, classification using a classifier $k \leqslant 40$ generally produced larger halos, and particles assigned to those halos had a larger $P(H)$. In addition, both classifiers, particularly the classifier with $k \leqslant 40$, produced elongated halos, demonstrating difficulty in distinguishing between dense filaments and low-density halos.

The inclusion of *PCA* calculations (Figure 5c and 5d) helped eliminate this issue by providing information emphasizing the directionality component of filaments. Including larger $k$ values for *KNN* calculations led to many filament particles being classified as halo particles, likely due to the fact that measurements using large $k$ would often include information from a variety of structure classes, blurring their distinction. However, even though the classifier in Figure 5d used $k$-values much larger than those in 5c, the class probabilities are very similar, especially when compared to Figures 5a and 5b. By including information about the local density field directionality, class assignments were less affected by contamination from distant structures.

Figure 6 shows the HMFs for the halos identified in SIM. For all classifiers, it is clear that classification was least accurate for small halos; elsewhere, the HMFs for *KNN* (+ *PCA*) corresponded very closely with Warren et al. (2006); however, the HMFs for *VOR* and *CMD* were not as accurate. This provides evidence that *KNN* (+ *PCA*) provides important information for differentiating halos from filaments and voids, further supporting the conclusions drawn from Figure 3. Though not shown, the HMFs for *MI* and *ENC* were

**Toy-to-SIM Probability Contrasts for *KNN* (+ *PCA*)**



**Figure 5.** Class assignment results for `SIM`. The halo-filament particles, colored based on their probability contrast (see Eqn. (6)) are shown in **(a)** (*KNN*, $k \leqslant 8$), **(b)** (*KNN*, $k \leqslant 40$), **(c)** (*KNN + PCA*, $k \leqslant 8$), and **(d)** (*KNN + PCA*, $k = 40$). The displayed particles are from the same 7 Mpc thick slice of the 3D N-body simulation `SIM` centered on the halo with the greatest mass as in Figure 4.

nearly identical to those of *VOR* and *CMD*. The inaccuracy for small halos had the same source as for FOF calculations (Springel et al. 2005).

### 3.2 Toy-to-`TSIM`

The true class values for `SIM` are not known, and there is not a known verifiable method that can assign these true classes. However, demonstrating that the statistics of our class assignments are similar for target datasets with very different structural properties demonstrates the robustness and lack of bias of our classifier. We will do this by comparing the statistics of `SIM` particle assignments to classes assigned to a toy model `TSIM`. A plot of `TSIM` (not shown) appears nearly identical to that of the training dataset seen in Figure 2.

Classes were assigned to `TSIM` particles using the same training dataset as `SIM`; these results may be seen in Figures 7 and 8. The large-scale properties of these class assignments were largely the
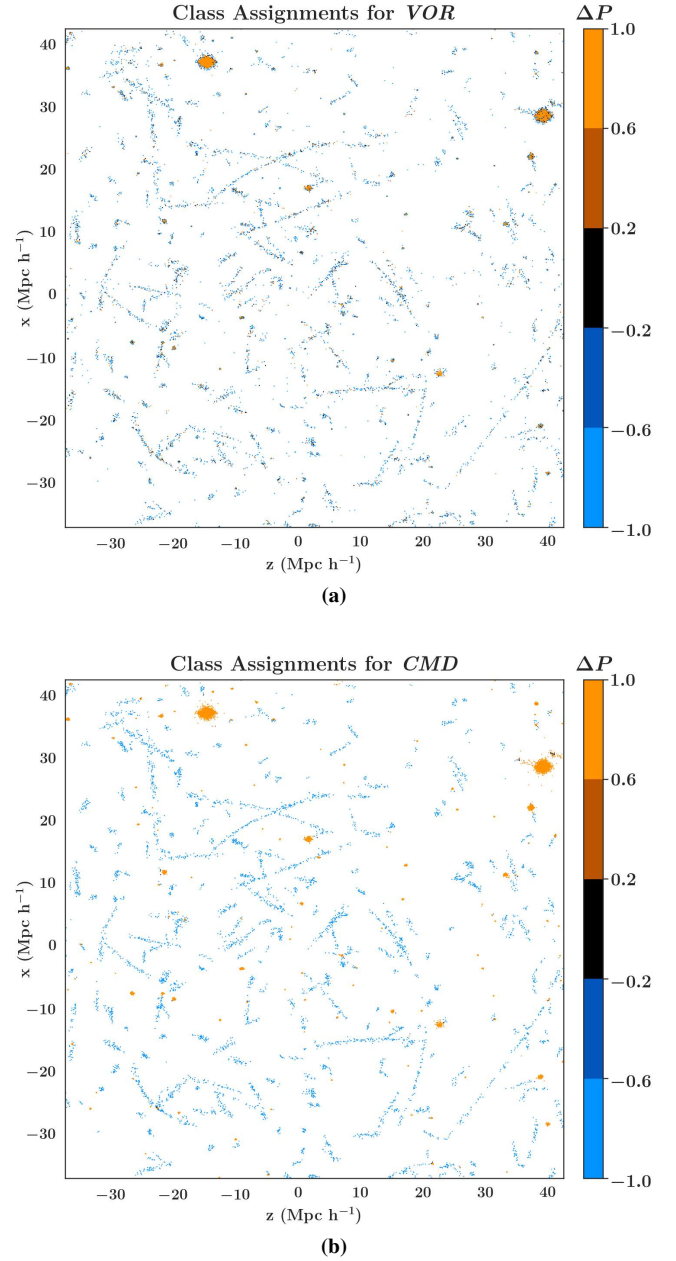
**Halo Mass Functions for** SIM



**(a)**



**(b)**

**Figure 6.** Halo mass functions for SIM created using classifiers trained using **(a)** *VOR* and *CMD* and **(b)** *KNN* (+ *PCA*) for $k \leqslant 8$, 40. Note that not all lines are easily visible; this is because the HMFs in each figure were deviated little from one another. The poor correspondence between the Warren et al. (2006) HMF and that of *VOR* and *CMD*, particularly at small *M*, indicates that *VOR* and *CMD* are ineffective proxies for local density magnitude. In contrast, the HMFs **(b)** deviated little from the Warren et al. (2006) HMF. This may be due to the similarities between *KNN* measurements and FOF algorithms.

same as those of in SIM seen in Figure 4 and 5: *VOR* and *CMD* (7a and 7b, respectively) generally overestimated halo abundance, with assignments made by *VOR* having lower probability assignments than *CMD*; *KNN*-trained classifiers with only small *k* (Figure 8a) assigned lower probabilities to all particles and produced smaller halos than a classifier with large *k* (Figure 8b); and the addition of *PCA* calculations (Figure 8c and 8d) helped remove the differences in class label assignments associated with different ranges in *k*-values.

One notable exception, however, is the classifier trained with only *CMD* (7b). Unlike in SIM (Figure 7b), this classifier did not produce clumps of ambiguous particles near particularly dense halo-filament boundaries in TSIM (Figure 7). This is likely due to
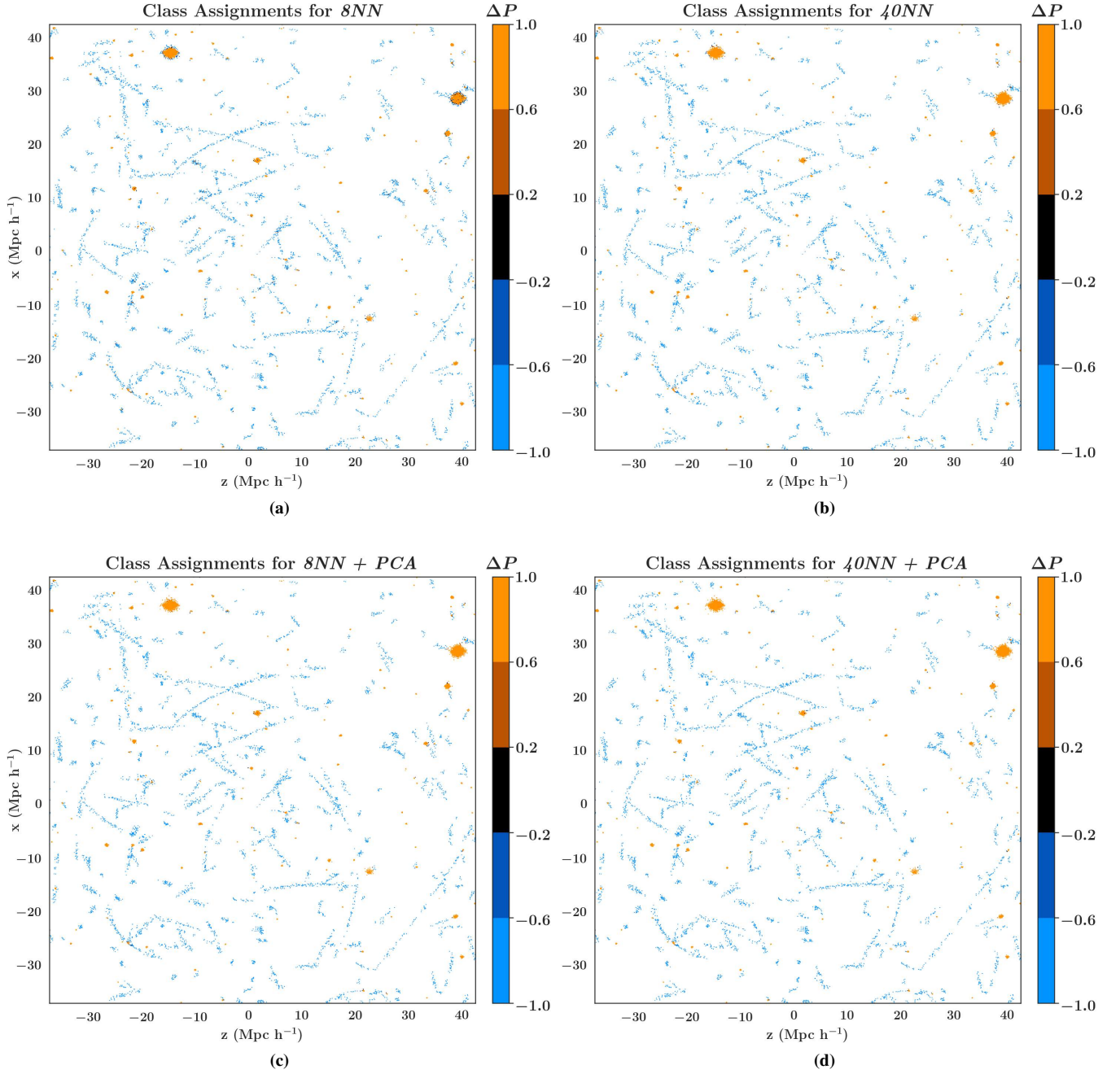
**Toy-to-TSIM Probability Contrasts for** *CMD* **and** *VOR*



**(a)**



**(b)**

**Figure 7.** Class assignments for TSIM. **(a)** and **(b)** show the halo and filament class assignments made using *VOR* and *CMD* measurements, respectively, colored based on each particle's probability contrast. The particles seen are from a 7 Mpc thick slice of the 3D N-body simulation TSIM. This slice includes the largest halo, which has $N = 5482$ particles and is centered at $(x, y, z) = (28.5, -7.3, 39.2)$ Mpc $h^{-1}$.

the fact that the training data used an identical generation procedure to TSIM. This highlights the robustness of *KNN* (and *PCA*) calculations. The assignments made by *CMD* calculations are heavily tied to the generation algorithm, indicating that it is not suited for training using a simplified dataset such as the toy model we developed. On the other hand, the class labels assigned by a classifier trained with *KNN* calculations, especially when paired with *PCA* calculations, are less affected by the exact structure of the toy model,
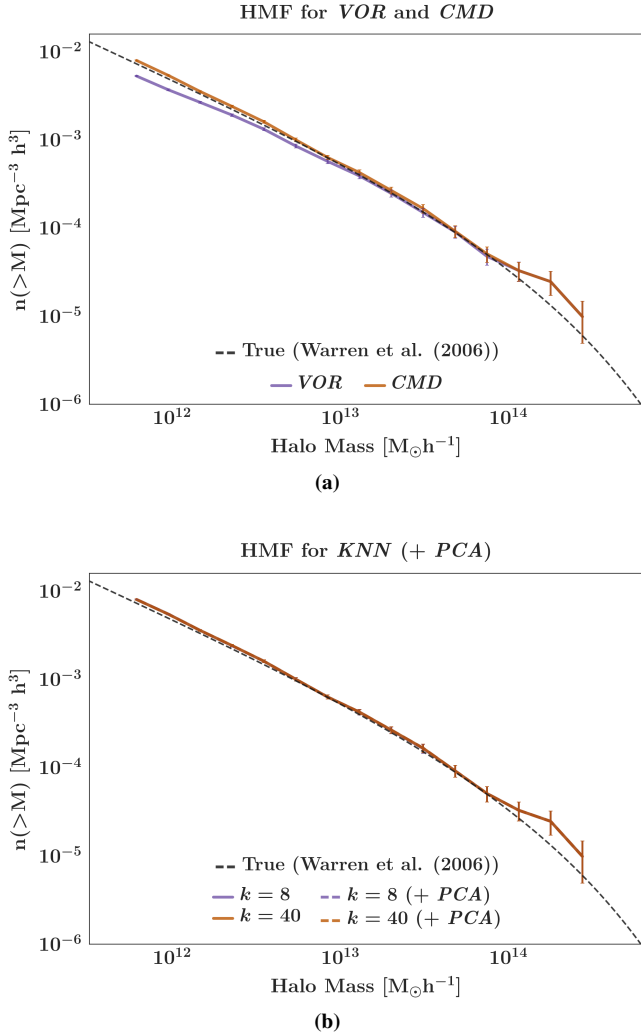
**Figure 8.** Class assignments for `TSIM`. The halo-filament assignments, colored based on the probability contrast, are shown in **(a)** (*KNN*, $k \leqslant 8$), **(b)** (*KNN*, $k \leqslant 40$), **(c)** (*KNN + PCA*, $k \leqslant 8$), and **(d)** (*KNN + PCA*, $k = 40$). The particles seen are from the same 7 Mpc thick slice of the 3D N-body simulation `TSIM` as 7.

enabling their use even when training was performed using a simplified toy model. This indicates that *KNN* calculations establish a natural length scale for halos and filaments together when performing classification, and the addition of *PCA* calculations help distinguish between halos and filaments by establishing distinct length scales for these structures individually. Note that *PCA* calculations are most effective for structures with a length scale that is not much larger than $R_{\mathrm{PCA}}$. Additional discussion of these properties may be found in Section 4.4.

Figure 9 shows the HMFs derived from `TSIM`. For all figures, it is clear that halo classification was least accurate at the extreme ends of the mass range. The high-mass deviation was likely due to statistical fluctuations resulting from there being very few large-mass halos in `TSIM`. The overprediction of low-mass halos by all classifiers other than *VOR* may arise from the generation parameters. In the toy model, the radii of low-mass halos was close to the radii of most filaments. As a result, some filaments in `TSIM` were incorrectly classified as small halos, leading to the discrepancy. The
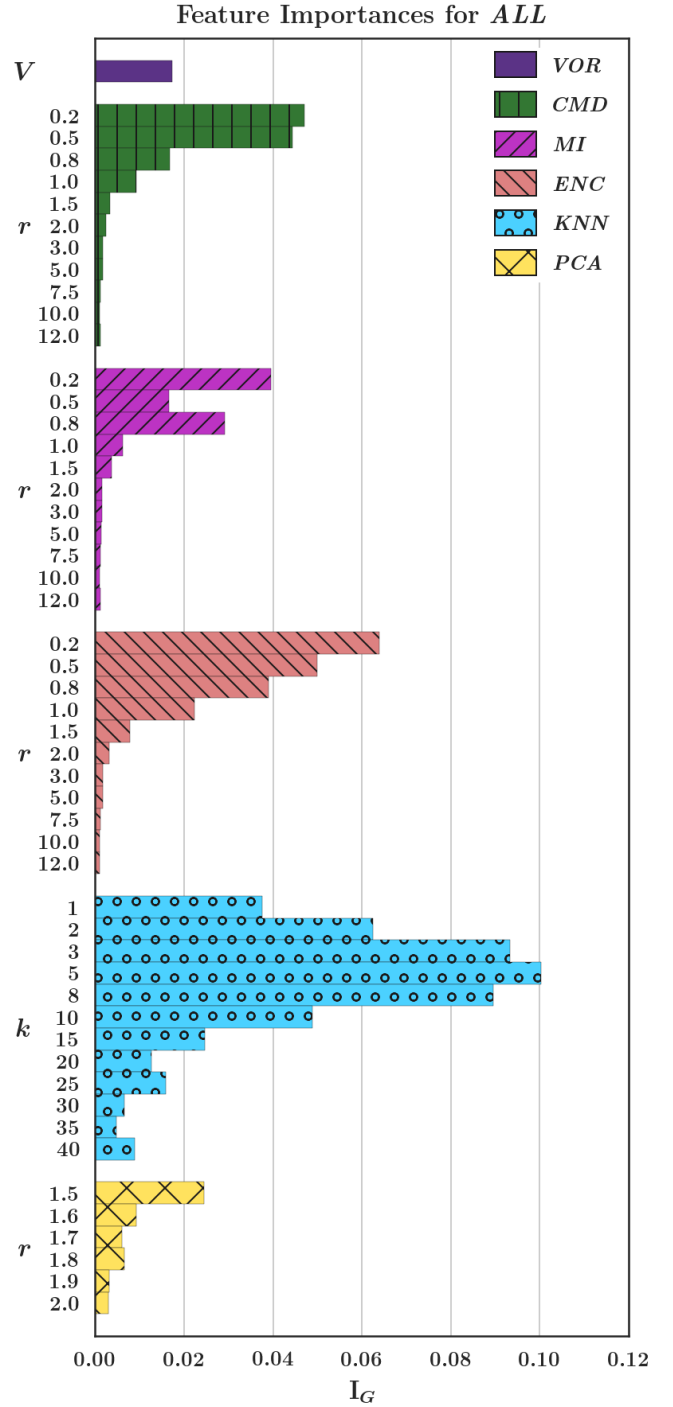
**Halo Mass Functions from TSIM Class Labels**



**Figure 9.** Halo mass functions for TSIM using **(a)** *VOR* and *CMD* and **(b)** *KNN* (+ *PCA*) for $k \leqslant 8, 40$.

*VOR* HMF is inaccurate for low-mass halos because *VOR* is a poor differentiator between LSS classes (see Figures 4a and 7a).

## 4    ANALYSIS

By correlating results from SIM with those of TSIM, we can demonstrate that the toy model is effective as a training data set: if the the properties of class assignments in SIM are statistically and/or visually similar to those in TSIM, which has markedly different properties, we have demonstrated that the methodology is robust enough to be applied to observed data. The goal of this section is not only to demonstrate the validity of our classifier, but also to establish the importance of utilizing measurements of both local density magnitude and directionality to ensure that our class assignments are not strongly influenced by the somewhat arbitrary parameters used to generate structures, particularly filaments, in the toy model.
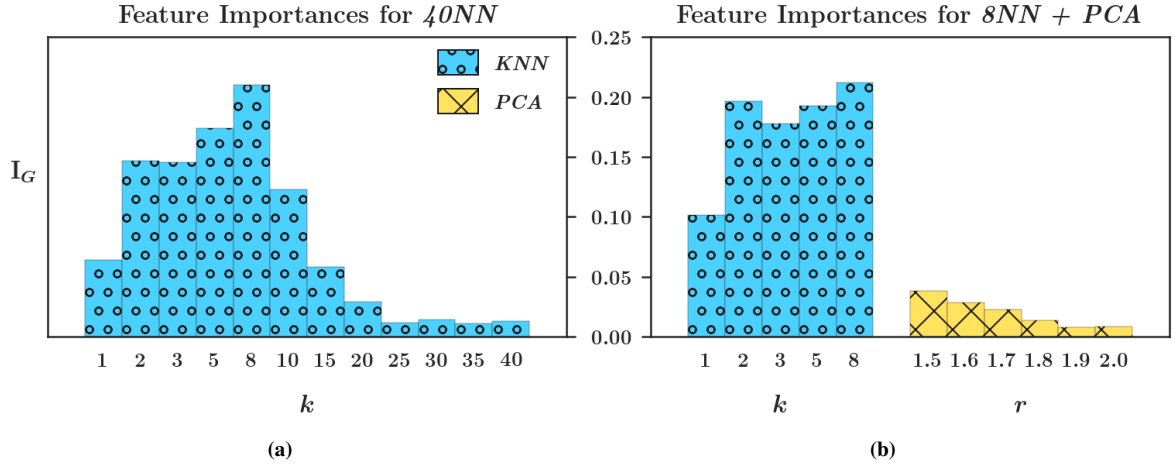
**Feature Importances for Toy-to-SIM**



**Figure 10.** The feature importances $\mathbf{I}_G$ for the Toy-to-SIM calculations using all metrics.

### 4.1    Feature Importances

First, we aimed to decrease the number of measurements required to minimize computation time and understand the role of each metric. We utilized the classifier's feature importances, which describes how relevant each feature was when performing class assignment; this was determined by the frequency a metric was used to choose

**Figure 11.** The feature importances $\mathbf{I}_G$ for the Toy-to-`SIM` calculations using **(a)** *KNN* and **(b)** *KNN + PCA*.

a branch as the classifier descended a tree. Figure 11 shows the feature importances for a variety of different metrics. From here, we see that *KNN* is weighted the most heavily, indicating that it may provide valuable information for the classifier. This expectation correlates with the measurement histogram seen in Figure 3c, where it may be seen that the distance to the eighth-nearest neighbor separates each of the classes distinctly from one another.

Figures 10 and 11 shows the feature importances for several different feature sets. In general, small radii/$k$ were deemed most important. From 11, the feature importances for all features, density magnitude calculations were weighted more heavily than *PCA*, and of the density magnitude calculations, *KNN* was weighted most heavily, reflecting the lack of differentiation between structures in the measurement histograms for the other density magnitude calculations. In both Figures 11 and 11a, the feature importances for *KNN*, the most important measurements were those with $k \leqslant 8$, reaching a peak at $k = 8$. This is likely due to the fact that the smallest halos had 8 particles, and as small halos dominated the halo mass function, they will be utilized most by the training algorithm to determine a feature's importance. These phenomena were also reflected in 11b *KNN + PCA*.

## 4.2 ROC AUC

The receiver operating characteristic (ROC) curve is a way of demonstrating a classifier's ability to accurately discriminate between classes. It consists of a plot of the classifier's true positive rate as a function of its false positive rate. A classifier that cannot discriminate between classes would have a 50% probability of assigning the correct class to a given data point, and would have equal true positive and false positive rates; hence, its ROC curve would appear as a line with unit slope. On the other hand, an effective classifier would have a much larger true positive rate than false positive rate. The area under the curve (AUC) is a measure for a classifier's effectiveness: a classifier with a true positive rate much larger than its false positive rate would have an AUC value near unity, while one with poor discriminatory ability would have AUC = 0.5.
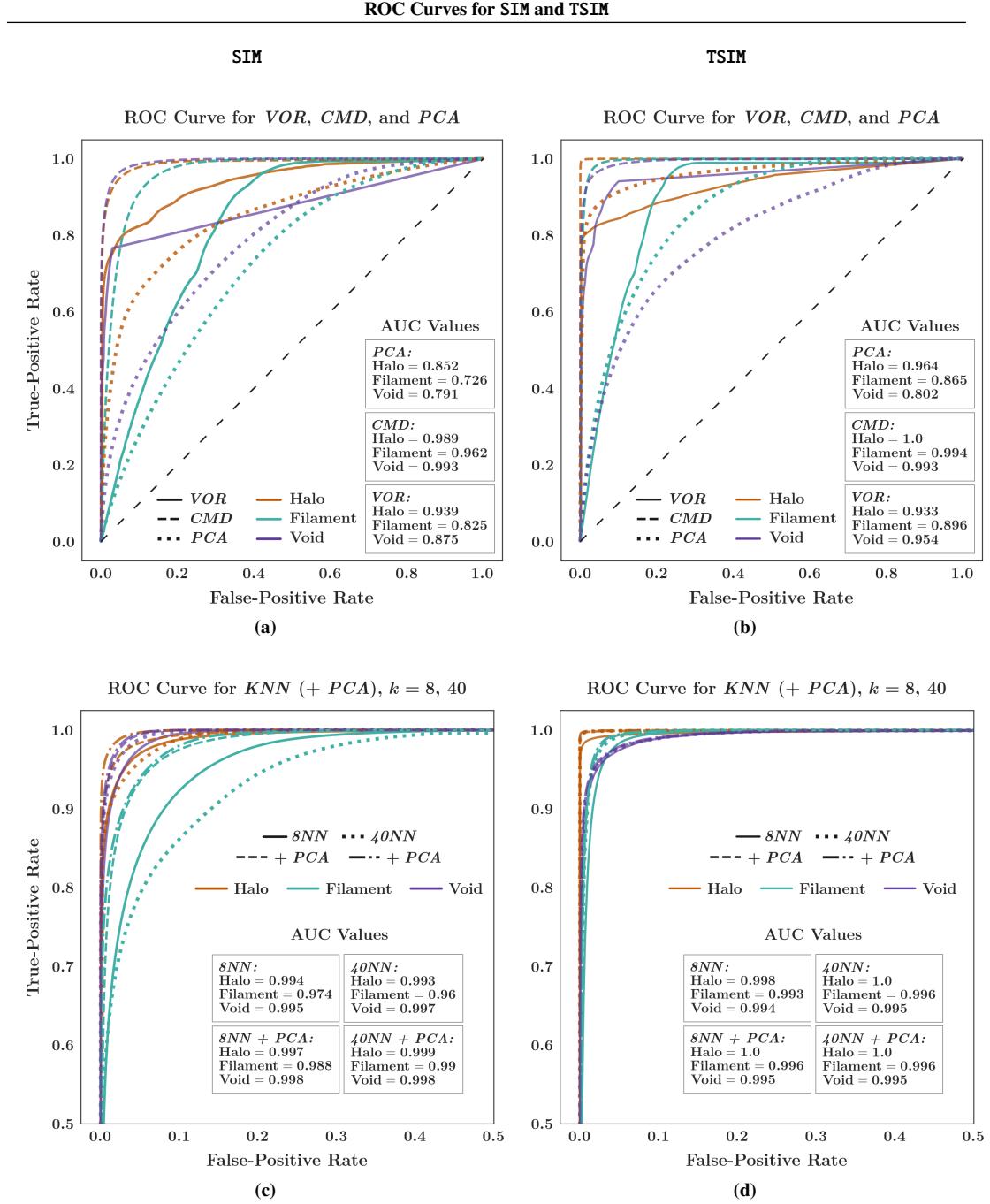
The first method we used to verify our results was through calculating the ROC AUC independently for each class; ROC plots (Pedregosa et al. 2011; Fawcett 2006; Katzma 1989) for some met-

ric combinations can be found in Figure 12. For `TSIM`, the assigned class labels were compared with the true class values in `TSIM`; however, these values did not exist in `SIM`. As a result, we chose the set of class labels assigned by a classifier trained using all features to use as a fiducial comparison dataset. This provides the most generality, as it effectively allows the comparison of each classifier to all others simultaneously. As we found that most metric sets produced similar results, and that the AUC for `TSIM` was maximized for all classes when using a classifier trained using all features, we believe that these class assignments will provide a sufficient approximation to the true class values to use as a fiducial comparison dataset.

In general, the ROC curves in Figure 12 demonstrate that filament classification was the most difficult. In addition, classification on `TSIM` was more accurate than on `SIM`, as evidenced by the shape of the curves and the AUC values.

Figure 12a and 12b show the ROC curves for classifiers trained using only *VOR*, *CMD*, and *PCA*. From this, it can be seen that classifiers trained on *VOR* or *PCA* alone were not effective when classifying particles in `SIM` and on `TSIM`. This further demonstrates that *VOR* does not suffice as a proxy for local density magnitude. The poor performance of *PCA* alone may be attributed to the lack of information provided to the classifier about local density magnitude.

Figure 12c and 12d show the ROC curves for classifiers trained using *KNN* (+ *PCA*) measurements for $k \leqslant 8$, 40. Notably, the addition of *PCA* calculations improved classifier performance for all LSS classes, particularly for filaments. As seen in the `SIM` classes assigned by a classifier trained by *KNN + PCA* (see Figure 5), the addition of *PCA* calculations diminished the dependence of halo classification on the values of $k$ used. Figure 12c demonstrates that this stabilization applies to all classes, especially filaments. The benefits of this stabilization are immense: using large $k$-values allows the classifier to take into account the global environment when performing classification, improving classification of large halos, but lessens its sensitivity to properties of the local environment. Including *PCA* calculations enables the global environment to be used in training without contaminating information about small-scale properties. The combination of small-scale and large-scale information in training enables classification of halos

**Figure 12.** ROC plots for SIM and TSIM. **(a)** and **(b)** show the ROC curves for a classifier trained using *VOR*, *CMD*, an *PCA*, and **(c)** and **(d)** show the same for *KNN (+ PCA)*, $k \leqslant 8$, 40. The dashed diagonal line with unit slope in **(a)** and **(b)** corresponds with a classifier that assigns labels by selecting a class at random; curves above this line indicate that the classifier using those features has predictive capabilities. Note that this line is not visible in **(c)** and **(d)**; this is because the bounds on these figures were altered for illustration purposes.

isolated in large void volumes (a major issue discussed in Tsizh et al. (2019), 2019 and Libeskind et al., 2018).
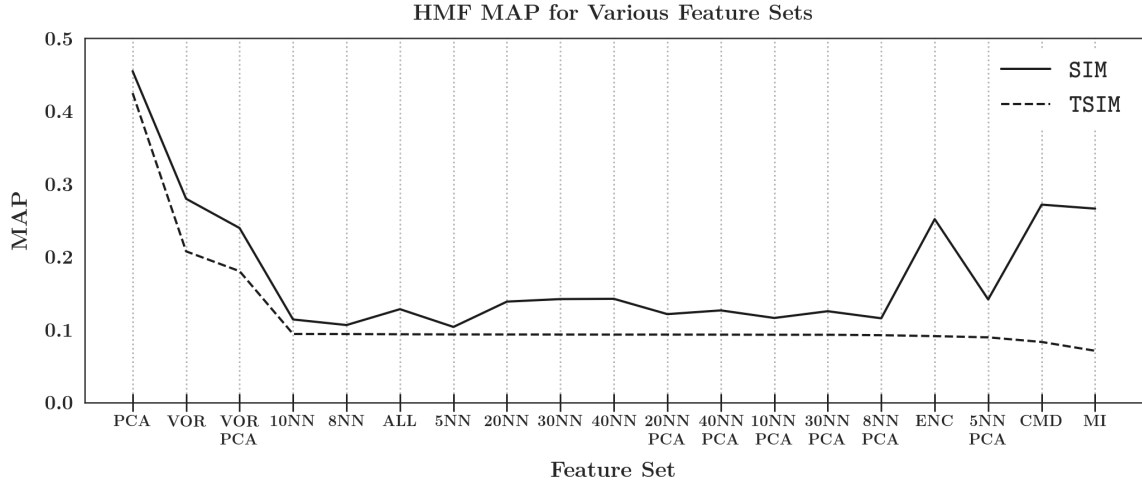
### 4.3 HMF Comparison: Mean Absolute Proportion

We define the mean absolute proportion (MAP) as

$$\text{MAP} = \frac{\mu\left(|n_{\text{Pred}} - n_{\text{Warren}}|\right)}{\mu\left(|n_{\text{Warren}}|\right)}, \qquad (7)$$

where $\mu(n)$ is the mean of $n$ over all $M$ and $n_{\text{Pred}}$ and $n_{\text{Warren}}$ were the empirical and Warren et al. (2006) HMFs, respectively. In Figure 13, we show the MAP for each of the metrics for both SIM and TSIM. The metrics are ordered based on the MAP value for TSIM. From this plot, it can be seen that, as before, classifiers trained using only *VOR* or *PCA* performed substantially worse than all other feature combinations when classifying TSIM particles. In addition, though not previously discussed, a classifier trained using *VOR + PCA* performed poorly when classifying TSIM particles,

HMF MAP for Various Feature Sets



**Figure 13.** A plot of the HMF MAP; the metric sets are ordered such that the `TSIM` MAP decreased from left to right.

emphasizing the importance of a robust density magnitude metric. For `SIM`, these three also performed very poorly; however, unlike in `TSIM` classification, classifiers trained using *CMD*, *MI*, and *ENC* also exhibited a large MAP. This further supports the conclusion that these methods are not effective due to their strong dependence on the training model generation algorithm. Classifiers trained using *KNN* (+ *PCA*) generally performed better for `SIM`, possibly due to the similarity between *KNN* measurements with FOF algorithms.

The analysis performed in Sections 3.1.1 and 4.1 both indicated that *KNN* calculations would likely be the most important, and the results from Sections 4.2 and 4.3 supported this by demonstrating that *KNN* calculations were robust and generated the most accurate results for both `SIM` and `TSIM`. As a result, in the following section, we will use *KNN* alone as the proxy for local density magnitude. In addition, Sections 3.1.1, 4.2, and 4.3 suggest that the inclusion of *PCA* calculations may also be of benefit, particularly for filament classification. In the next section, we present additional measurements to emphasize the importance of *PCA* calculations when creating a robust classifier.

### 4.4 Robustness of *PCA* Calculations

As discussed previously, *PCA* calculations provide substantial benefit when paired with *KNN* calculations; however, as noted in Section 4.3, these benefits are most apparent for structures with length scales no larger than the maximum radius used in *PCA* calculations ($R_{\mathrm{PCA}}$ = 2.0 Mpc $h^{-1}$) Halo vs. filament probability plots for a region where this is true can be seen in Figure 14.

Figure 14a and 14b show the classification results for classifiers trained with *KNN*, $k \leqslant 8$, 40, respectively. As discussed previously, introducing *PCA* calculations (Figure 14c and 14d) curbed the classifier's dependence on the maximum value of $k$ and improved its ability to distinguish between halos and filaments. Regions **I** and **II** in Figure 14 demonstrate this clearly: in 14a and 14b, the classes assigned to particles varied greatly, with the classifier with $k \leqslant 40$ substantially overestimating te number of halo particles. This produced halos that are elongated along one axis, a hallmark trait of filaments. However, *PCA* calculations prevented this issue, as exemplified by the similarity between these regions in 14c and 14d. These strong similarities also indicate that *PCA* calculations help enable the use of global density measurements for classification. While it is clear that the halo radii in regions **I** and **II**
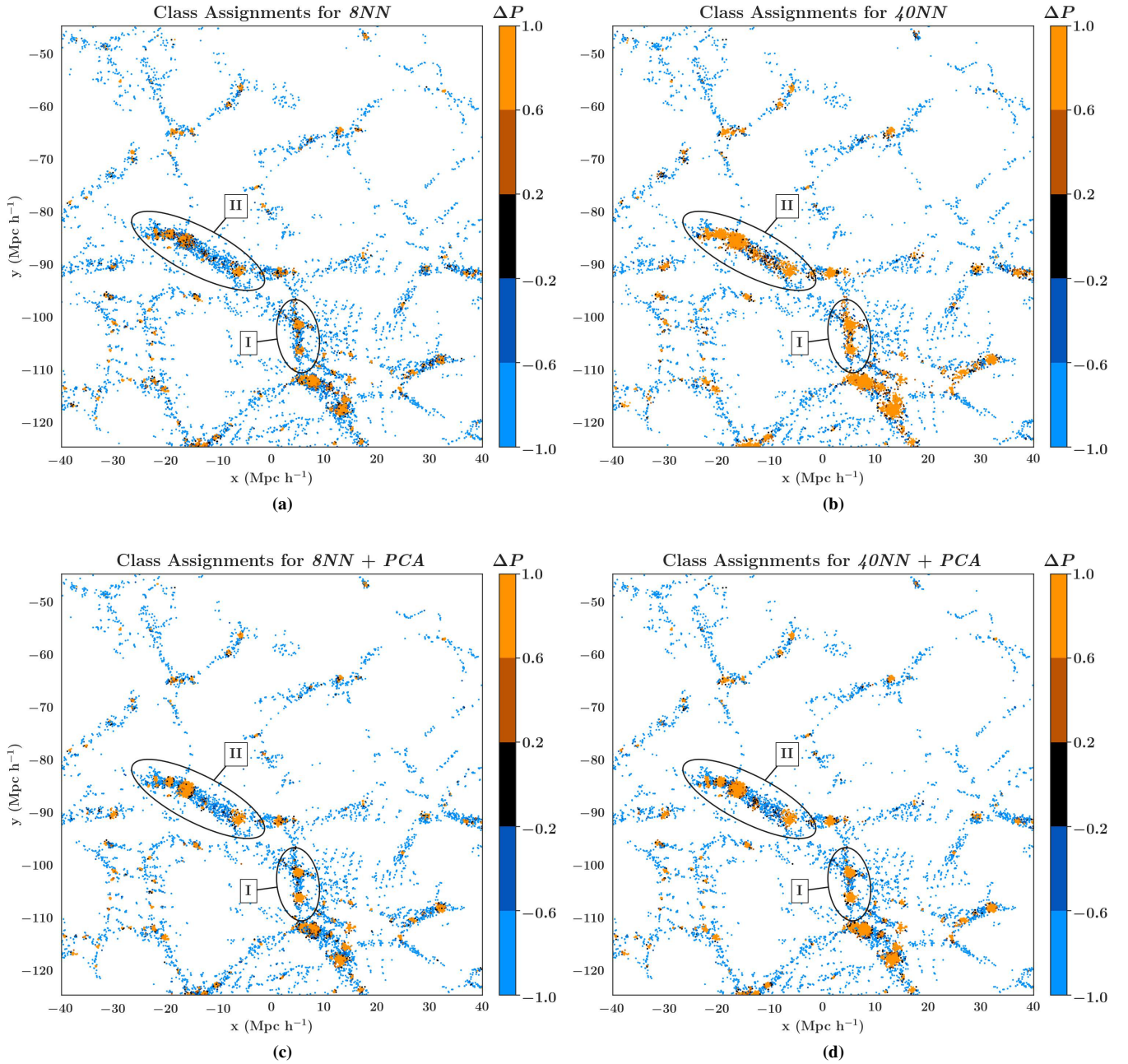
vary substantially between Figures 14a and 14b, these values in 14c and 14d are much closer. Measurements using large $k$ provide information about distant structures. This information is very valuable; however, it can impede classification based on density magnitude by blurring the distinction between the characteristics of the morphological classes. *PCA* calculations help decontaminate the class properties, allowing the inclusion of information on the global density magnitude without losing information about small-scale structural properties.

In addition, due to the ambiguous class of particles on halo-filament boundaries, we expect the particles in these regions would have a probability contrast near zero (these particles are colored black in Figure 14). However, as the density contrast between halos and voids is very large, we expect the class assignments on halo-void boundaries to have higher probabilities. In Figure 14a, there are an abundance of ambiguous particles on halo-filament and halo-void boundaries, as well as in the interior of halos; this is most visible in the halo in the upper-left of Region **II**. In contrast, Figure 14b predominantly lacks ambiguous border particles. The addition of *PCA* calculations in 14c and 14d stabilize these border regions, removing most ambiguous particles inside halo interiors and on halo-void boundaries, and clarify the halo-filament boundary with a thin layer of ambiguous particles, consistently making particle classes more physically realistic.

Though most visible in Regions **I** and **II**, these conclusions apply to many halos shown in Figure 14. As these halos lie in environments with widely varying densities, it is clear that directionality information helped relax the classifier's dependence on density magnitude, instead favoring the local density contrast. Constructing a classifier that is consistent in varied environments has proven difficult (Tsizh et al. 2019; Libeskind et al. 2018), but information on both the local density field magnitude and directionality provides a way to avoid this issue.

#### 4.4.1 Structural Mass Fractions

In the previous sections, we demonstrate that *KNN* and *PCA* calculations together create a robust classifier. Here, we aim to show that the exact construction of our toy model (particularly that of filaments, which was based on visual appearance) did not substantially bias our results. Using this information, we may isolate the gain in classifier effectiveness provided by training using *PCA*.
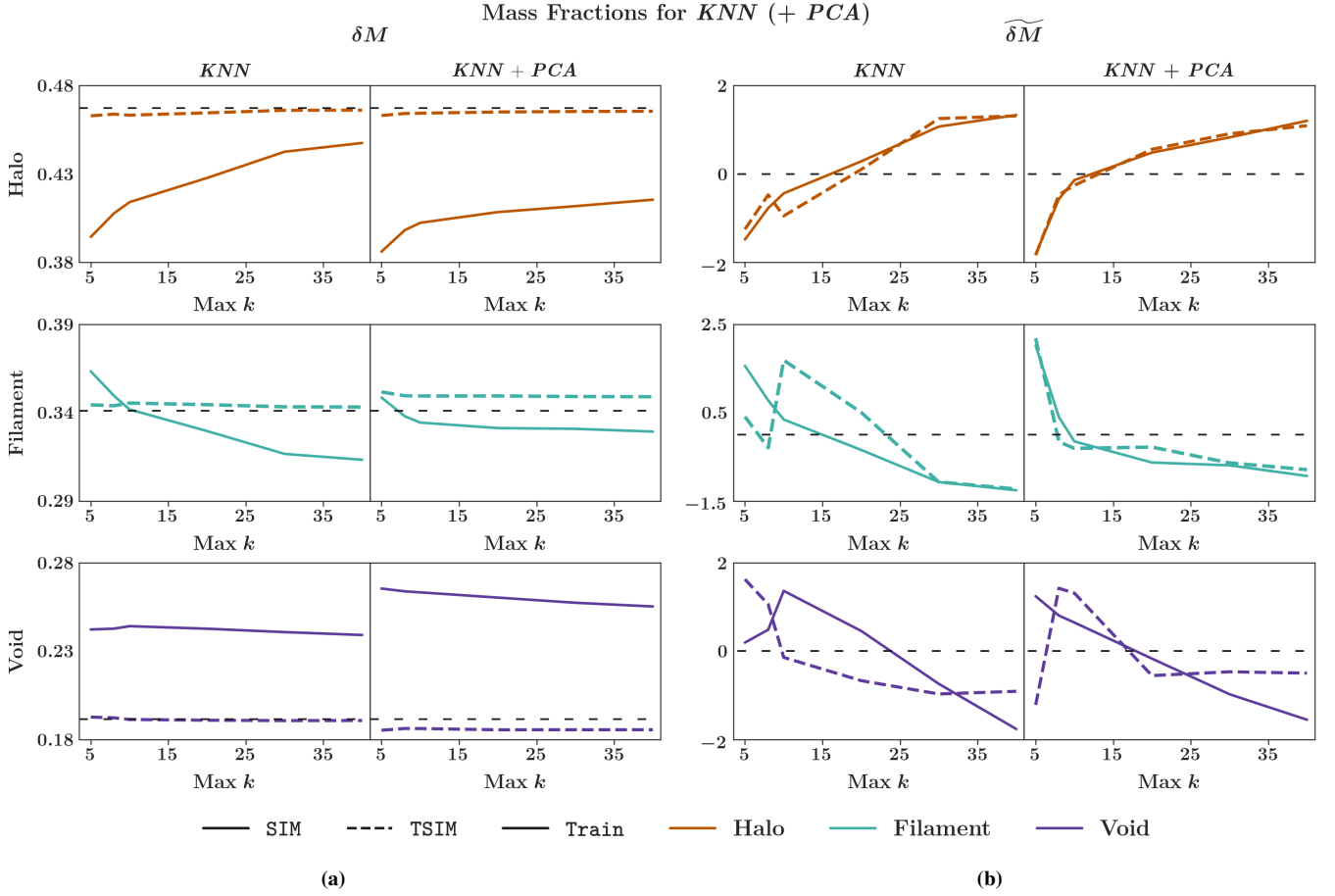
**Figure 14.** Halo vs. filament probability contrast scatterplots for *KNN* ((**a**) and (**b**)) and *KNN + PCA* ((**c**) and (**d**)) for $k \leqslant 8$, 40. Particles are from a slice of width 2.0 Mpc $h^{-1}$. In this region, most halos had radii that were not substantially larger than the maximum radius used for *PCA* calculations ($r = 2.0$ Mpc).

Figure 15 (a) shows mass fractions for halos, filaments, and voids for SIM and TSIM. For all structures, with and without *PCA*, it is clear that the mass fractions for SIM are different from those of TSIM. However, while the mass fractions in TSIM vary little as $k$ changes, those of SIM do (for *KNN* only). With the addition of *PCA* calculations, the SIM mass fractions tend to stabilize to a constant value as $k$ changes. The addition of *PCA* calculations affects minimally affects the mass fraction values of TSIM. The black dashed line shows the true mass fractions for the training dataset (which corresponds very closely to the true mass fractions for TSIM). The fact that the mass fractions for SIM were much further from the

true mass fraction from the training dataset indicates that the exact parameters of the training dataset did not affect the results substantially. This demonstrates that the filament chosen for the toy model did not bias the classes assigned to particles in SIM.

Next, we aim to demonstrate the robustness provided by the addition of *PCA* calculations by determining the dependence of the mass fraction on $k$. To determine whether the classifier's results were dependent on $k$, we assumed the null hypothesis: the assignments made by given classifier are independent of the set of features used for training, so the assignments made by a classifier trained using a particular feature set is taken to be a single trial in a set of

**Figure 15.** (a) The mass fractions $\delta M$ for halos, filaments, and voids for `SIM` and `TSIM` as a function of $k$. (b) The normalized mass fraction $\widetilde{\delta M}$ (defined in Eqn. 8) of particles classified as halos, filaments, and voids as a function of $k$. *KNN* classifiers were trained using only *KNN* measurements for $k \leqslant k_{\max}$; each classifier trained using *PCA* included measurements for all $r \in R_{\mathrm{PCA}}$.

predictions made by a single classifier. Under this assumption, we may normalize the mass fractions using their mean and standard deviation, which were calculated using the mass fractions derived from the test data labelled by each classifier.

Consider a classifier $C_f$ trained using a set of features $f \subseteq F$. Let $\delta M_{f,\mathrm{C}}$ be the fraction of particles assigned to a morphological class C by $C_f$. Then define the normalized mass fraction $\widetilde{\delta M}_{f,\mathrm{C}}$ by

$$\widetilde{\delta M}_{f,\mathrm{C}} = \frac{\delta M_{f,\mathrm{C}} - \overline{\delta M}_{\mathrm{C}}}{\sigma_{\mathrm{C}}}. \tag{8}$$

Here, $\overline{\delta M}_{\mathrm{C}}$ is the mean of the mass fraction over all $f \subseteq F$ for a class C and $\sigma_{\mathrm{C}}$ is the corresponding standard deviation. For `SIM` and `TSIM`, we will compare $\widetilde{\delta M}_{f,\mathrm{C}}$ for the set of classifiers trained using *KNN* to those trained using *KNN* and *PCA*.

Plots of $\widetilde{\delta M}_{f,\mathrm{C}}$ are shown in Figure 15 (b). For a given classifier trained using all $k \leqslant k_{\max}$, the value of $\widetilde{\delta M}_{k_{\max}}$ corresponds with the number of standard deviations between $\delta M_{k_{\max}}$ and the average of $\delta M_{k_{\max}\mathrm{C}}$ over all $k_{\max}$.

While measured mass fractions for `SIM` and `TSIM` were different, the normalized mass fractions each have a mean of 0 and a standard deviation of 1, allowing `SIM` and `TSIM` to be compared to one another directly.

A plot of the normalized mass fractions can be seen in Figure 15 (b), where several trends may be seen. For both `TSIM` and `SIM`, as

$k$ increased, the normalized mass fraction for halos increased, while that of filaments decreased, regardless of whether or not *PCA* calculations were used. In addition, for both `SIM` and `TSIM`, *PCA* calculations generally affected the values of the normalized mass fractions for halos and filaments minimally. Without *PCA* calculations, the `TSIM` normalized mass fraction for both halos and filaments varied widely for small values of $k$; for filaments in particular, this large difference was seen for $k < 30$. In addition, for $k \leqslant 10$, the `TSIM` normalized mass fractions for halos and filaments did not exhibit a particular trend with increasing $k$.

However, the addition of *PCA* calculations substantially improved the classifier's robustness by making the normalized mass fraction values and trends more consistent between `SIM` and `TSIM` for halos and filaments, particularly for $k \leqslant 10$. This is likely because halos are only found as nodes on filaments in the toy model, so for very small halos, using only density calculations blurred the lines between halos and filaments through contamination of filament point measurements by nearby halos. The incorporation of *PCA* calculations helped differentiate filaments from halos, improving the robustness. This claim is further supported by the very large filament normalized mass fraction for $k = 5$: as the smallest halos in the toy model had 8 particles, measurements for $k \leqslant 5$ would not be able to include all particles in a halo. As a result, many halos were classified as filaments due to the fact that directional-

ity effects would dominate the local density magnitude for halos, causing overrepresentation of filaments.

Voids showed no recognizable trend for TSIM, regardless of whether or not *PCA* calculations were included. For SIM, the normalized mass fraction showed a general downward trend, though as before, the normalized mass fraction without *PCA* calculations for $k \leqslant 10$ showed an inconsistent trend.

## 5 CONCLUSIONS

We have presented a novel method for cosmic web classification, demonstrating that supervised machine learning using a simplified toy model as training data provides a potential avenue for robust and efficient cosmic web classification. The simplicity of our toy model indicates that the amount of information required for cosmic web classification is relatively low: assuming appropriate metrics are used, accurate classification can be achieved using even minimally realistic training data. While most known methods require the measurement or inference of the velocity field and/or knowledge of structural properties predicted by analytical methods, we demonstrate a methodology that requires only information about each galaxy's position. The use of a random forest algorithm in particular provides a method for achieving probabilistic classification. In addition, we found that the use of density field directionality measurements in tandem with local density magnitude measurements are crucial for distinguishing between halos, filaments, and voids. In particular, we have provided a method that can classify isolated halos inside large voids, an outstanding problem discussed in Tsizh et al. (2019) and Libeskind et al. (2018). Through calculating and comparing statistical data about our classifications, we found a method to verify our calculations, demonstrating that our algorithm is robust and is not biased by our training data creation algorithm.

The key advantage of our method is the speed and efficiency of toy model generation. While N-body simulations require substantial computational expense and lack true class values, a new toy model can be generated much more efficiently, and this model provides enough information to accurately classify substantially more complicated N-body simulations and potentially observed data. This makes the method especially suitable for large datasets: using a single training dataset, we were able to assign class values to an N-body simulation substantially larger than the toy model. Due to the speed of generation, our method is extremely scalable, as generating additional training datasets would allow us to assign class values to very large N-body simulations at no cost.

In addition, the use of a toy model is particularly suited for cosmic web classification of observed data. Observed datasets contain masked regions and areas with non-uniform depth. The use of a toy model helps account for these issues: for each field, an individual toy model can be generated that matches the density, mask, and depth of that field. By classifying each field using its corresponding training dataset, class assignments would be generated consistently for each field. The use of periodic boundary conditions or padding (as we used here) could avoid issues associated with masked regions.

The ability to assign probabilistic classes to individual galaxies opens to door to a variety of novel data analysis techniques. Observables such as density and physical composition are known to be linked to LSS class membership, so correlations with class probability would enable novel methods for understanding these relationships. For example, spectral analysis may be used to understand chemical composition of galaxies in different environments. Correlating R/G-band magnitudes with cluster-filament probability contrasts could help establish not only the differences in the chemical compositions of clusters and filaments, but also how that composition changes as the cluster-filament boundary is crossed.

As *PCA* calculations clarify halo-filament boundaries, the application of a directionality metric can be used to differentiate halos and filaments in general, as well as study the fundamental properties of LSS. Capturing snapshots of an N-body simulation over large time scales and tracking filament halos as they cross a halo-filament boundary could provide a deeper understanding of the matter inflow from filaments to halos, separating its role in halo formation and collapse from other processes.

While our algorithm is a robust classifier for halos, filaments, and voids, our feature selection is not ideal; in particular, we found that *PCA* calculations did not perform well for filaments and halos with radii substantially larger than the radius used in the *PCA* decomposition calculations. As increasing this radius leads to substantial cross-contamination, future work should focus on identifying and implementing a directionality metric that can more effectively capture properties of large filaments and halos. The information content (Leclercq et al. 2016, 2015) could be used to determine the utility of new metrics.

In addition, we chose not to differentiate between sheets/walls and filaments in our classifier to the complexity of creating a simple algorithm for this purpose. Future work could be devoted to expanding this algorithm to allow sheet/wall classification.

Though untested, we expect our classifier to be just as effective in cosmology models other than ΛCDM. The length scales for LSS are much larger than the scales at which deviations from ΛCDM are detectable. As our training data only reproduces properties of LSS at these larger scales, we expect classification to be independent of the cosmology of the target data set, so the same training data sets may be used to classify fields with equivalent geometric parameters (e.g. average density) but different cosmologies.

## 6 SOFTWARES USED

The Python packages Bezier (Hermes 2017), DBSCAN (Ester et al. 1996), HMF (Murray et al. 2013), matplotlib (Hunter 2007), numpy (Oliphant 2006), scikit-learn (Pedregosa et al. 2011), SciPy (Virtanen et al. 2019), Shapely (Gillies et al. 2007), and yt (Turk et al. 2011) were used extensively in this work.

## 7 ACKNOWLEDGEMENTS

# 8  DATA AVAILABILITY

The data underlying this article will be shared on reasonable request to the corresponding author. Additional figures and basic software demonstrations may be found at https://github.com/bmbuncher/Prob-CWeb.

# REFERENCES

Alpaslan M., et al., 2013, Monthly Notices of the Royal Astronomical Society, 438, 177
Aragon-Calvo M. A., 2014, Mon. Not. Roy. Astron. Soc., 440, 46
Aragon-Calvo M. A., 2018, preprint
Aragon-Calvo M. A., Jones B. J. T., van de Weygaert R., van der Hulst M. J., 2004, Astron. Astrophys., 474, 315
Aragon-Calvo M. A., van de Weygaert R., Jones B. J. T., 2010a, Monthly Notices of the Royal Astronomical Society, 408, 2163
Aragon-Calvo M. A., van de Weygaert R., Jones B. J. T., 2010b, Monthly Notices of the Royal Astronomical Society, 408, 2163
Bardeen J. M., Bond J. R., Kaiser N., Szalay A. S., 1986, ApJ, 304, 15
Bharadwaj S., Bhavsar S. P., Sheth J. V., 2004, ApJ, 606, 25
Bonnaire T., Aghanim N., Decelle A., Douspis M., 2019, arXiv e-prints
Bradford Barber C., Dobkin D. P., Huhdanpaa H., 1996, ACM TRANSACTIONS ON MATHEMATICAL SOFTWARE, 22, 469
Breiman L., 2001, Machine Learning, 45, 5
Cautun M., van de Weygaert R., Jones B. J. T., 2012, MNRAS, 429, 1286
Cautun M., van de Weygaert R., Jones B. J. T., Frenk C. S., 2014, Mon. Not. Roy. Astron. Soc., 441, 2923
Codis S., Pogosyan D., Pichon C., 2018, Mon. Not. Roy. Astron. Soc., 479, 973
Crocce M., Pueblas S., Scoccimarro R., 2006, MNRAS, 373, 369
Davis M., Efstathiou G., Frenk C. S., White S. D. M., 1985, ApJ, 292, 371
Ester M., Kriegel H., Sander J., Xu X., 1996, in Proceedings of the Second International Conference on Knowledge Discovery and Data Mining. AAAI Press, pp 226–231
Falck B., Neyrinck M. C., 2015, Monthly Notices of the Royal Astronomical Society, 450, 3239
Falck B. L., Neyrinck M. C., Szalay A. S., 2012, ApJ, 754, 126
Fawcett T., 2006, Pattern Recogn. Lett., 27, 861
Forero-Romero J. E., Hoffman Y., Gottloeber S., Klypin A., Yepes G., 2009, Mon. Not. Roy. Astron. Soc., 396, 1815
Foster C., Nelson L. A., 2009, ApJ, 699, 1252
Galàrraga-Espinosa D., Aghanim N., Langer M., Gouin C., Malavasi N., 2020, arXiv e-prints
Gillies S., et al., 2007, Shapely: manipulation and analysis of geometric objects
González R. E., Padilla N. D., 2010, MNRAS, 407, 1449
Green S. B., van den Bosch F. C., 2019, Mon. Not. Roy. Astron. Soc., 490, 2091
Hamaus N., Sutter P. M., Wandelt B. D., 2014, Phys. Rev. Lett., 112, 251302
Hansen S. M., McKay T. A., Wechsler R. H., Annis J., Sheldon E. S., Kimball A., 2005, Astrophys. J., 633, 122
Hermes D., 2017, The Journal of Open Source Software, 2, 267
Hoffman Y., Metuki O., Yepes G., Gottlöber S., Forero-Romero J. E., Libeskind N. I., Knebe A., 2012, MNRAS, 425, 2049
Hui J., Aragon M., Cui X., Flegal J. M., 2018, MNRAS, 475, 4494
Hunter J. D., 2007, Computing in Science & Engineering, 9, 90
Jasche J., Kitaura F. S., 2010, Monthly Notices of the Royal Astronomical Society, 407, 29â49842
Jasche J., Wandelt B. D., 2013, Monthly Notices of the Royal Astronomical Society, 432, 894â49913
Katzma McClish D., 1989, Medical Decision Making, 9, 190
Kitaura F., Angulo R. E., 2012, Monthly Notices of the Royal Astronomical Society, 425, 2443
Komatsu E., et al., 2011, ApJ, 192, 18
Kraljic K., et al., 2019, preprint

Lavaux G., Wandelt B. D., 2010, Monthly Notices of the Royal Astronomical Society, 403, 1392â491408
Leclercq F., Jasche J., Wandelt B., 2015, Astron. Astrophys., 576, L17
Leclercq F., Lavaux G., Jasche J., Wandelt B., 2016, Journal of Cosmology and Astroparticle Physics, 2016, 027â49027
Leclercq F., Jasche J., Lavaux G., Wandelt B., Percival W., 2017, Journal of Cosmology and Astroparticle Physics, 2017, 049â49049
Libeskind N. I., et al., 2018, Mon. Not. Roy. Astron. Soc., 473, 1195
Lucie-Smith L., Peiris H. V., Pontzen A., Lochner M., 2018, Mon. Not. Roy. Astron. Soc., 479, 3405
Mo H., van den Bosch F. C., White S., 2010, Galaxy Formation and Evolution
Murray S., Power C., Robotham A., 2013, preprint
Oliphant T., 2006, NumPy: A guide to NumPy, USA: Trelgol Publishing
Pedregosa F., et al., 2011, Journal of Machine Learning Research, 12, 2825
Pereyra L. A., Sgró M. A., Merchán M. E., Stasyszyn F. A., Paz D. J., 2019, arXiv e-prints, p. arXiv:1911.06768
Ramachandra N. S., Shandarin S. F., 2015, Monthly Notices of the Royal Astronomical Society, 452, 1643
Rodriguez A. C., Kacprzak T., Lucchi A., Amara A., Sgier R., Fluri J., Hofmann T., Rèfrègier A., 2018, Comput. Astrophys. Cosmol., 5, 4
Scoccimarro R., 1998, MNRAS, 299, 1097
Shandarin S., Habib S., Heitmann K., 2012, Phys. Rev. D, 85, 083005
Sousbie T., 2011, Monthly Notices of the Royal Astronomical Society, 414, 350
Springel V., 2005, MNRAS, 364, 1105
Springel V., et al., 2005, Nature, 435, 629
Tempel E., Stoica R. S., Kipper R., Saar E., 2016, Astron. Comput., 16, 17
Tipping M. E., Bishop C. M., 1998, Mixtures of Probabilistic Principal Component Analysers
Tsizh M., Novosyadlyj B., Holovatch Y., Libeskind N. I., 2019, preprint
Turk M. J., Smith B. D., Oishi J. S., Skory S., Skillman S. W., Abel T., Norman M. L., 2011, The Astrophysical Journal Supplement Series, 192, 9
Virtanen P., et al., 2019, arXiv e-prints, p. arXiv:1907.10121
Warren M. S., Abazajian K., Holz D. E., Teodoro L., 2006, ApJ, 646, 881
White S. D. M., Rees M. J., 1978, MNRAS, 183, 341

This paper has been typeset from a TeX/LaTeX file prepared by the author.