Benchmarks and Explanations for Deep Learning Estimates of X-ray Galaxy Cluster Masses

Matthew Ho,^{1★} John Soltis,² Arya Farahi,³ Daisuke Nagai,⁴ August Evrard,⁵ and Michelle Ntampaka^{2,6}

- ¹CNRS & Sorbonne Université, Institut d'Astrophysique de Paris (IAP), UMR 7095, 98 bis bd Arago, F-75014 Paris, France
- ²Department of Physics & Astronomy, Johns Hopkins University, Baltimore, MD 21218, USA
- ³Departments of Statistics and Data Science, University of Texas at Austin, Austin, TX 78705, USA
- ⁴Department of Physics, Yale University, New Haven, CT 06520, USA
- ⁵Departments of Physics and Astronomy and Leinweber Center for Theoretical Physics, University of Michigan, Ann Arbor, MI 48109, USA
- ⁶Data Science Mission Office, Space Telescope Science Institute, Baltimore, MD, 21218, USA

Accepted XXX. Received YYY; in original form ZZZ

ABSTRACT

We evaluate the effectiveness of deep learning (DL) models for reconstructing the masses of galaxy clusters using X-ray photometry data from next-generation surveys. We establish these constraints using a catalogue of realistic mock eROSITA X-ray observations which use hydrodynamical simulations to model realistic cluster morphology, background emission, telescope response, and AGN sources. Using bolometric X-ray photon maps as input, DL models achieve a predictive mass scatter of $\sigma_{\ln M_{500c}} = 17.8\%$, a factor of two improvements on scalar observables such as richness $N_{\rm gal}$, 1D velocity dispersion $\sigma_{\rm v,1D}$, and photon count $N_{\rm phot}$ as well as a 32% improvement upon idealised, volume-integrated measurements of the bolometric X-ray luminosity $L_{\rm X}$. We then show that extending this model to handle multichannel X-ray photon maps, separated in low, medium, and high energy bands, further reduces the mass scatter to 16.2%. We also tested a multimodal DL model incorporating both dynamical and X-ray cluster probes and achieved marginal gains at a mass scatter of 15.9%. Finally, we conduct a quantitative interpretability study of our DL models and find that they greatly down-weight the importance of pixels in the centres of clusters and at the location of AGN sources, validating previous claims of DL modelling improvements and suggesting practical and theoretical benefits for using DL in X-ray mass inference.

Key words: methods: data analysis – cosmology: large-scale structure of Universe – galaxies: nuclei –galaxies: clusters: general – galaxies: clusters: intracluster medium – X-rays: galaxies: clusters.

1 INTRODUCTION

Galaxy clusters are large collections of dark matter, gas, and galaxies dynamically bound through gravitational attraction. They are the most massive systems in the Universe, and, as such, act as observable tracers of peak density regions in the cosmic web. The spatial and temporal distribution of galaxy clusters encodes information regarding the growth of large-scale structures throughout the Universe's evolution (Allen et al. 2011, for a review). For example, the abundance of clusters as a function of mass is a widely used probe of the Universal matter density Ω_m and the amplitude of primordial density fluctuations σ_8 (e.g., Vikhlinin et al. 2009; Mantz et al. 2010; Ade et al. 2016; Abbott et al. 2020). A plethora of multiwavelength cluster surveys are underway to constrain physical models of gravity, dark matter, and dark energy (e.g., Pillepich et al. 2018; Raghunathan et al. 2022).

Accurate and precise estimates of galaxy cluster masses are of paramount importance to these analyses, as a cluster's mass sets the scale for processes both within the system and in its interactions with the cosmic web (Pratt et al. 2019, for a recent review). However, 90% of a cluster's mass is hidden away in its host dark matter halo and

* E-mail: matthew.annam.ho@gmail.com

must be inferred indirectly from the spatial and energy distributions of photons observed on the sky. A wide diversity of techniques has been developed to constrain cluster masses using data from X-ray (e.g., Kravtsov et al. 2006; Pratt et al. 2009; Giles et al. 2017; Mantz et al. 2016), microwave (e.g., Nagai 2006; Kay et al. 2012), and optical surveys (e.g., Saro et al. 2013; Wojtak et al. 2018). Generally, these methods establish a physical connection between observable signals and cluster masses and then use simulations (e.g., Nagai et al. 2007a; Planelles et al. 2014; Biffi et al. 2016; Pop et al. 2022a,b) or complementary observations (e.g., Ade et al. 2011; Reichardt et al. 2013; Schellenberger et al. 2015; Mulroy et al. 2019; McClintock et al. 2019) to calibrate their mass inference.

The recent popularity of this problem has led to considerable work to reduce the level of mass scatter with novel data analysis methods. These approaches benefit from their ability to utilise complex signals in cluster observables for mass inference, which are otherwise difficult to model analytically. Deep neural networks trained on mock observations from hydrodynamical simulations have shown a strong ability to improve the accuracy of cluster mass estimates using X-ray (Ntampaka et al. 2019; Green et al. 2019; Yan et al. 2020; Krippendorf et al. 2023), microwave (Cohn & Battaglia 2020; Wadekar et al. 2023a,b; de Andres et al. 2022), and optical data (Ntampaka et al. 2015; Ho et al. 2019, 2021, 2022; Kodi Ramanah et al. 2020). These

techniques, while promising, require extensive validation on mock data before they can be reliably extended to observational samples (Ntampaka et al. 2021).

Today, the development of techniques for precise mass inference on X-ray images is of particular significance due to the recent launch of the extended ROentgen Survey with an Imaging Telescope Array (eROSITA Merloni et al. 2012). eROSITA is designed to perform a deep survey of the sky through the X-ray energy band (0.5 - 10 keV). It is projected to detect approximately ~ 100,000 clusters over its four-year operating timeline (Pillepich et al. 2018). eROSITA will observe these clusters at a lower angular resolution than its predecessor, the Chandra X-ray Observatory, but will reach a much larger sample at a well-modelled selection function, making it an ideal instrument for cosmological inference. In addition, the clusters detected by eROSITA will be subject to follow-up with spectroscopic, optical, and microwave instruments of related surveys, such as SDSS-V, DESI, Euclid, and Rubin, opening the door toward large-scale multiwavelength studies of cluster physics. The use of accurate, precise, and efficient techniques for inference is essential to capitalise on the wealth of eROSITA data.

In this paper, we forecast and explain the mass estimation performance of neural network analysis of X-ray observations for the recently launched eROSITA telescope. We utilize a catalogue of mock observations of 3,285 distinct clusters in the Magneticum hydrodynamical simulation (Dolag et al. 2016), each including realistic noise, instrument response, and simulation-driven realisations of cluster morphology, core physics, and AGN contamination. We train and test modern neural network models on this catalogue and benchmark their mass predictions against the common observable mass proxies. We implement methods for further reducing scatter using multiwavelength probes, specifically by applying neural networks on multi-band X-ray images and joint X-ray and spectroscopic data. We then conduct a quantitative interpretability study of our ML models, investigating and evaluating the improved behaviour of neural networks in the context of cluster mass prediction.

The structure of this paper is as follows. Section 2 presents the simulation data and the procedure for generating mock observations. Section 3 describes our comparative baseline, a covariance analysis of scalar mass proxies commonly used in X-ray analysis. Section 4 introduces our CNN modelling approach for single-band X-ray images and reports their predictive performance against our baselines. Section 5 describes the performance improvements of CNN models, which incorporate multiband X-ray and spectroscopic dynamical information. Section 6 explains the results of our interpretability study of the aforementioned CNN models. Lastly, we present our conclusions and suggestions for future work in Section 7.

2 DATA

Successful implementation of data-driven inference methods, such as the neural networks described in Section 4 and Section 5, requires training catalogues that accurately reflect the realism and distribution of observational data. The mock catalogues curated in this analysis are among the most realistic X-ray cluster mocks currently applied to deep learning models and are the first to incorporate realistic modelling of AGN sources (See Ntampaka et al. 2019; Green et al. 2019; Yan et al. 2020, for a comparison). Our dataset is built on mock eROSITA catalogues generated in Soltis et al. (2022). In the following sections, we describe the primary properties of the mock catalogue, as well as adjustments made to it for the purposes of this work. For details on the original catalogue, see Soltis et al. (2022).

Property	Min.	Median	Max.
$\log_{10} \left[M_{500c} \left(h^{-1} \mathrm{M}_{\odot} \right) \right]$	13.50	13.94	15.07
z	0.07	0.21	0.47
$N_{ m phot}$	17	1072	2.37×10^{5}
$N_{ m ICM}$	9	762	2.00×10^{5}
$N_{ m AGN}$	0	225	2.03×10^{5}
$N_{ m gal}$	28	112	520

Table 1. Summary properties of the mock X-ray (Section 2.1) and dynamical catalogues (Section 2.2). Each row shows the minimum, median, and maximum values of all clusters in our training and test catalogues. The quantities, in order, are: the logarithmic cluster mass, redshift, total observed X-ray photons, total observed photons from ICM sources, total observed photons from AGN sources, and the number of galaxies within the dynamical selection cut (Section 3, for details). The photon counts reported here are the values prior to the imposed redshift normalisation with Equation 1.

The base simulation for our mock catalogue is the Magneticum Box 2/hr hydrodynamical simulation (Hirschmann et al. 2014), provided by the Cosmological Web Portal (Ragagnin et al. 2017). The Magneticum simulation (Dolag et al. 2016) is a cosmological TreePM-SPH simulation code containing physical models for radiative cooling, star formation, stellar populations, and black hole and AGN feedback. Hydrodynamical simulations such as Magneticum are a key cornerstone of our understanding of galaxy clusters and have been used extensively to study X-ray properties of clusters (e.g., Kravtsov et al. 2006; Nagai et al. 2007b; Biffi et al. 2018; Green et al. 2019). The Box2/hr run is simulated in a cube of side length 352 h^{-1} Mpc, an initial particle number of 2×1564^3 , a dark matter particle mass of $M_{\rm DM} = 6.9 \times 10^{18}~M_{\odot}$, and a gas mass of $M_{\rm gas} = 1.4 \times 10^8 \ M_{\odot}$. It assumes a WMAP7 cosmology of $\Omega_m = 0.272, \Omega_b = 0.0456, \Omega_{\Lambda} = 0.728, h = 0.704, n_s = 0.963, and$ $\sigma_8 = 0.809$ (Komatsu et al. 2009). The halos and subhalos in this volume are identified via the Subfind Friends-of-Friends halo finder (Springel et al. 2001; Dolag et al. 2009). Throughout this work, halo mass is defined as the spherical overdensity mass at a density contrast of 500 times the critical density of the universe, henceforth referred to as M_{500c} in the cosmology invariant units of $h^{-1}M_{\odot}$.

We selected a sample of 3285 galaxy clusters from the Magneticum halo catalogue, chosen to have a roughly uniform distribution across the redshift range $0.07 \le z \le 0.47$ and across the mass range $3.16 \times 10^{13} M_{\odot} \le M_{500c} \le 1.17 \times 10^{15} M_{\odot}$. This was done to train our inference model to have a uniform prior over mass and redshift, a favourable property for cosmological analyses. The distribution of the cluster properties of our sample is characterised by Table 1. Each simulated cluster is viewed from a single projected perspective along the z box axis, and is therefore included in our mock catalogue only once. In the following subsections, we describe the X-ray and spectroscopic observables calculated for each mock cluster from this projected perspective.

2.1 Mock eROSITA Observations

The mock X-ray observations used in this work are designed to match the observing conditions described in the eROSITA instrument specifications (Merloni et al. 2012). Specifically, our observations project a field of view spanning approximately 1°, a pixel size of 9.6", and an observation time of 2 ks. The centre of each cluster is assumed to be located at a sky position of 10' right ascension and 10' declination. We have simulated intra-cluster medium (ICM) emission using the PHOX algorithm (Biffi et al. 2012, 2013) incorporating an observation depth of 10 Mpc. Light from AGN sources simulated in the Mag-

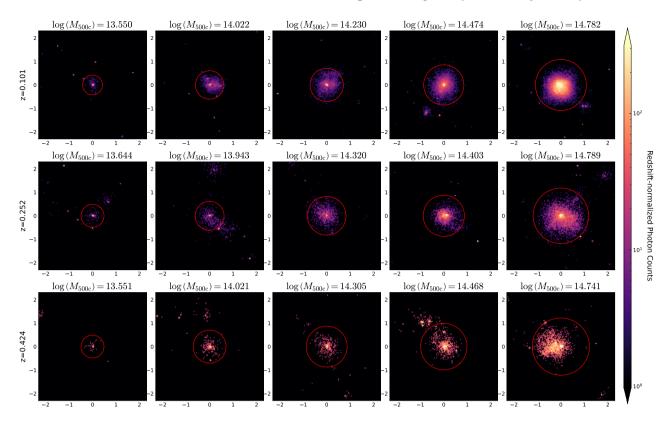


Figure 1. Mock eROSITA photon maps for twelve example clusters in the Magneticum simulation. Each row selects clusters from a different Magneticum snapshot, and selected clusters are sorted in increasing mass from left to right. Photon counts in each sky-projected bin are normalised in scale to account for the system redshift (Equation 1). Each subplot also shows a red circle indicating the R_{500c} of each cluster. M_{500c} masses are presented in units of $h^{-1}\text{Mpc}$.

neticum volume is projected into our observations using the method described in Biffi et al. (2018). These sources trace the positions and properties of Stellar Mass Black Holes (SMBH) in the Magneticum simulation and thus follow the distribution of large-scale structure contaminants expected in real observations. Additionally, we have generated realisations of background emission, instrument response, and point spread function that are consistent with the eROSITA telescope design and exposure using the SIXTE algorithm (Dauser et al. 2019). For more information on modelling X-ray contamination in the eROSITA telescope with SIXTE, see Clerc et al. (2018).

The SIXTE software produces a list of photons expected to be observed by eROSITA for each simulated cluster. Table 1 tabulates the minimum, median, and maximum number of photons that we receive in these images, from ICM and AGN sources. Photons observed within each mock SIXTE pixel are assumed to sit at their pixel center. We then map these photon sky coordinates to flat-sky 2D images. First, we convert the photon list in projected sky coordinates to comoving coordinates using the Magneticum WMAP7 cosmology (Komatsu et al. 2009). We then measure a projected 2D histogram of photon counts across a square aperture of length 2.3 h^{-1} Mpc at a pixel resolution of 128 × 128. This has the effect of 'zoomingin' high-redshift systems and 'zooming-out' low-redshift systems, resulting in a set of 2D photon maps which are scale-invariant to redshift.. Since this data is derived from the same angular resolution and instrumental effects as the SIXTE configuration, this operation is repeatable on real observations. We also normalise the photon counts in each pixel to luminosity distance by assuming a log-linear relationship between redshift and total measured photons. Explicitly, the luminosity normalisation follows

$$n_{ij} = m_{ij} \times \frac{N_0}{10^{\gamma z_{\rm clu} + \eta}},\tag{1}$$

where m_{ij} and n_{ij} are the raw and redshift-normalised per-pixel photon counts, respectively, N_0 is a pivot value fixed at an approximate mean total photon count 4000, $z_{\rm clu}$ is the redshift of the cluster, and m and b are linear regression slopes learned from fitting $\log \sum_{ij} m_{ij} = \gamma z_{\rm clu} + \eta$ for all mock observations in our training set. These normalising transformations produce a set of images whose angular scale and absolute magnitude are invariant to redshift, apart from the shot noise and AGN contamination associated with their evolved environments. We note here that these operations are possible with knowledge of the cluster redshift, which will be attainable for eROSITA clusters through the SPIDERS spectroscopic followup program (Clerc et al. 2020).

From this step, we derive two versions of the X-ray images. We first produce single-band images integrated over the full energy range of the eROSITA instrument (0.5-10~keV). These are to be used in our CNN modelling of single-band photon maps, as presented in Section 4. Examples of these images are shown in Figure 1. We also produce multi-band images in the three eROSITA energy bands, namely the soft (0.5-1.2~keV), medium (1.2-2.0~keV), and hard bands (2.0-7.0~keV). These are to be used in Section 5 as an improved X-ray probe of the cluster mass. The sole difference between the multi-band and single-band images is in their energy separation, i.e., the total photons received by each pixel across all energies in the multi-band are equal to those of the single band. An example of a multi-band image is shown in Figure 2.

Lastly, to reduce the dynamic range of the pixel values, we apply a logarithmic scaling to the images before they are passed as input to the neural network. The scaling we use is,

$$x_{ij} = \log_{10} \left[n_{ij} + 1 \right], \tag{2}$$

where n_{ij} is the redshift-normalised number of photons detected at the j-th pixel of the i-th row, and x_{ij} is the corresponding scaled pixel value used as input of the neural network.

2.2 Mock Spectroscopic Follow-up

In addition to the X-ray images, we generated a matching catalogue of spectroscopic follow-up observations to evaluate the multiwavelength models in Section 5. These catalogues follow the design of those presented in Ho et al. (2019, 2021). We assume that an observer at redshift z = 0 measures the exact spectra of galaxies around each Magneticum cluster from the same line-of-sight perspective as our X-ray observations. The galaxy positions and velocities are seeded from Magneticum subhalos found with the Subfind algorithm. We then convert these measurements to comoving positions, velocities, and stellar mass estimates. We place cuts on these measured properties to assign cluster membership to galaxies with stellar mass $M_{\rm star} \geq 10^{9.5}~h^{-1}{\rm M}_{\odot}$ and within a dynamical cylinder of projected radius $R_{\rm proj} \leq 2.3h^{-1}{\rm Mpc}$ and line-of-sight velocity $|v_{los}| \le 3785 \text{ km s}^{-1}$, relative to the cluster center. With this large dynamical cut, we incorporate the various systematics that impact observational studies of cluster dynamics, including simulated physics of dynamical substructure (Old et al. 2018), cluster mergers (Evrard et al. 2008), halo environment (White et al. 2010), and interloping galaxies (Wojtak et al. 2018). Table 1 shows the minimum, median, and maximum number of galaxies in each observation after selection cuts.

To turn these galaxy catalogues into regular inputs amenable to CNN architectures, we apply the density estimation and sampling technique introduced in Ho et al. (2019) and extended in Kodi Ramanah et al. (2021). For each cluster in our training set, we use kernel density estimators (KDEs; Scott 2015) to estimate the distribution of galaxies in the 3D dynamical phase space, $\{x_{\text{proj}}, y_{\text{proj}}, v_{\text{los}}\|$. The KDEs applied here are Gaussian with a fixed bandwidth scaling factor of 0.15. We then sample this distribution on a regular $64 \times 64 \times 64$ grid that spans our dynamic selection cut, that is, $|v_{\text{los}}| \leq 3785 \text{ km s}^{-1}$, $|x_{\text{proj}}| \leq 2.3h^{-1}\text{Mpc}$, and $|y_{\text{proj}}| \leq 2.3h^{-1}\text{Mpc}$. This operation has the effect of smoothing our discrete list of galaxy positions and velocities into a fixed-size input. The resulting 3D input box is crucial to utilising convolutional filters in our neural network. Figure 2 shows an example of the distribution of galaxies around a Magneticum cluster in a cube of dynamical phase space.

3 BASELINE SCALAR OBSERVABLES

We first develop a quantitative baseline for predicting the host halo mass using scalar observables. Table 2 details the complete list of observables used in this analysis. A subset of these quantities are integrated directly over a 3D comoving sphere of radius R_{500c} , including temperature T, bolometric X-ray luminosity L_X , the gas mass $M_{\rm gas}$, the stellar mass $M_{\rm star}$, and the 1D projected velocity dispersion $\sigma_{\rm v,1D,true}$. We consider these as idealised observables without systematic errors and expect that they tightly correlate with the host halo mass (Evrard et al. 2008; Mantz et al. 2016; Farahi et al. 2018b; Mulroy et al. 2019), but they are difficult or impossible to know exactly from observation. In contrast, we include several realistic observables

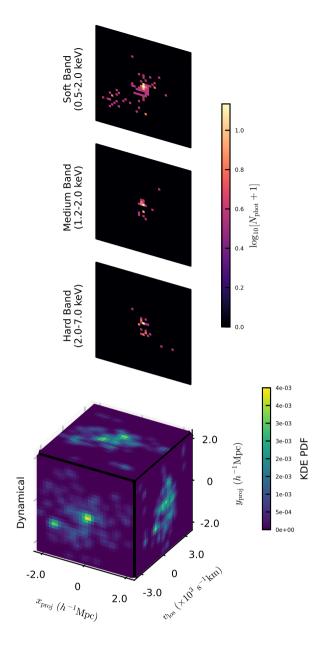


Figure 2. Example multiwavelength observables for a single Magneticum cluster at redshift z=0.47 and mass $M_{500c}=1.7\times10^{14}~h^{-1}{\rm M_\odot}$. The top three images show mock eROSITA photon maps in the soft, medium, and hard bands. The bottom subplot shows the KDE-estimated distribution of galaxies in dynamical phase space $\{x_{\rm proj}, y_{\rm proj}, v_{\rm los}\}$. The X-ray images and the dynamical cube are oriented identically such that the line of sight is oriented to the top-right, as indicated by the red arrow.

calculated directly from the mock catalogues described in Section 2. These are contaminated with observational systematics and more accurately reflect real measurement conditions. These include the total X-ray photon count in both the whole aperture, $N_{\rm phot}$, and the central region $N_{\rm phot}$,500c of our mock eROSITA images, as well as the richness $N_{\rm gal}$ and total 1D projected velocity dispersion $\sigma_{\rm v,1D}$ as measured in our mock spectroscopic catalogue. $N_{\rm phot}$ is calculated over every pixel in our wide aperture of $R_{\rm proj} \leq 2.3~h^{-1}{\rm Mpc}$, while $N_{\rm phot}$,500c only counts pixels within $R_{\rm proj} \leq R_{\rm 500c}$. We construct the $N_{\rm gal}$ and $\sigma_{\rm v,1D}$ probes from the population of galaxies selected

Observable	Description
T L_X $M_{ m gas}$ $M_{ m star}$ $\sigma_{ m v,1D,true}$	Temperature Bolometric X-ray luminosity Gas mass within R_{500c} Stellar mass within R_{500c} Projected 1D velocity dispersion
$N_{ m gal} \ \sigma_{ m v,1D} \ N_{ m phot} \ N_{ m phot,500c}$	Richness Total projected 1D velocity dispersion Total received photon count Photon count within R_{500c}

Table 2. Baseline scalar observable proxies for cluster mass. The top five quantities represent 'idealised' measurements of these observables, while the bottom four are 'realistic' measurements contaminated with observational systematics. Cluster centres, idealised hydrodynamical properties, and R_{500c} 's are produced by the Subfind algorithm (Dolag et al. 2009) and reported in Ragagnin et al. (2017).

within a small dynamical cylinder of aperture $R_{\rm proj} \leq 1.25 h^{-1} {\rm Mpc}$ and velocity cut $|v_{\rm los}| <= 2000~{\rm km~s^{-1}}$ and further refined via 3σ clipping on the $v_{\rm los}$ distribution, as a simplistic form of interloper removal. $\sigma_{\rm v,1D,true}$ and $\sigma_{\rm v,1D}$ are both calculated with the gapper estimator of velocity dispersion.

We evaluate the predictive power of these observables using a power-law model for multi-property cluster statistics as presented in Evrard et al. (2014). Here, we review the framework of this model and refer the reader to the original paper for further details. Consider a set of observable cluster properties **S** and their natural logarithmic counterparts $s_i = \ln S_i$, chosen in some given unit definition. Then consider a desired predictor, in this case, the cluster mass M_{500c} and its associated logarithm-scaled value $\mu = \ln(M_{500c}/M_p)$, normalised to a fixed pivot mass $M_p = 10^{14} \ h^{-1} \rm M_{\odot}$. In this setting, we can build a model in which the vector of features **s** scales linearly with our logarithmic mass μ ,

$$\langle \mathbf{s} | \mu \rangle = \pi + \alpha \mu, \tag{3}$$

where the vectors π and α are the intercepts and slopes, respectively, of the individual scaling laws. These parameters are redshift-dependent, and, in our analysis, we fit distinct π and α for each redshift bin of our mock catalogue. We then assume that the uncertainty about this mass-observable relationship is normally distributed with the mean given in Equation (3). The accuracy of this assumption is tested and verified against simulation (Farahi et al. 2018a; Anbajagane et al. 2020). In this case, the joint distribution of features s is completely described by the covariance matrix,

$$C_{ij} = \langle (s_i - \langle s_i \mid \mu \rangle) (s_j - \langle s_j \mid \mu \rangle) \rangle. \tag{4}$$

We then assume a local, first-order Taylor expansion to the halo mass function (Rozo et al. 2014),

$$n(\mu) = Ae^{-\beta\mu},\tag{5}$$

where A and β are the local amplitude and slope of the mass function evaluated at the pivot $\mu = 0$. This allows us to invert the relationship of Equation (3) to derive the predictive mean and variance of the mass given the information contained in the observables,

$$\langle \mu \mid \mathbf{s} \rangle = \left[\alpha^T \mathbf{C}^{-1} (\mathbf{s} - \boldsymbol{\pi}) - \beta \right] \sigma_{\mu \mid \mathbf{s}}^2,$$
 (6)

$$\sigma_{\mu|\mathbf{s}}^2 = \left(\alpha^T \mathbf{C} \alpha\right)^{-1}. \tag{7}$$

This procedure can determine the minimal scatter on cluster mass of any set of observables S.

Using Equation (7), we measure the predictive variance of each individual observable in our baseline suite and list them in Table 3. We also perform an exhaustive search over all combinations of ideal and realistic observables to find those with the lowest scatter. The combinations with the least scatter for their number of included observables are also shown in Table 3. The improvement in these combinations of observables is due to the strong correlations between their mass-dependent residuals, shown as a heatmap in Figure 3.

The scatter we find on observable mass proxies suggests that they are a reliable baseline for mass estimation. Among the idealised observables, $M_{\rm gas}$ has the lowest scatter at ~ 8.5%, consistent with empirical data in Mulroy et al. (2019). The most direct, yet idealised observable for eROSITA measurements, the bolometric X-ray luminosity L_X , only reaches a mass scatter of $\sim 26\%$. Among the realistic observables, $N_{\rm phot,500c}$ has the lowest scatter at ~ 34% followed by $N_{\rm gal}$ and $N_{\rm phot}$ at ~ 40%. The improvement in using $N_{\rm phot,500c}$ over N_{phot} is due to the removal of AGN contaminates outside of the R_{500c} cut. Another notable point is the poor performance of the total velocity dispersion $\sigma_{v,1D}$ (~ 70% scatter) and its idealised counterpart $\sigma_{v,1D,true}$ (~ 38% scatter). It is evident both from this study and previous works that dynamics measurements from galaxy spectra are high scatter predictors of cluster mass Saro et al. (2013), particularly for mass definitions at low radii. These results suggest that simple linear fits to dynamical measurements could be poor predictors of M_{500c} , and motivate the use of nonlinear fits such as the ML models in Section 5.2.

Figure 3 shows strong correlations between the mass residuals of various observables, suggesting that we can combine their information into better mass estimators (Table 3). The correlation matrix follows that of Farahi et al. (2019), with the only discrepancy being in the anti-correlation of $T - M_{gas}$, which was also found in Kravtsov et al. (2006). These results show that the most informative two-observable probe of M_{500c} is $\{M_{gas}, M_{star}\}$ with a residual scatter of $\sim 4.9\%$, better than the commonly used probe of $Y_X \propto M_{\rm gas}T$ (Kravtsov et al. 2006) which we measure to have a mass scatter of ~ 5.1%. We find that with only three observables, T, $M_{\rm gas}$, and $M_{\rm star}$, the mass scatter can be constrained to the level of $\sim 4\%$. Including all idealised observables, this reduces to ~ 3.9%. However, this idealised scatter is far from what is possible with linear combinations of realistic observables. The combination of all realistic observables (i.e., $\{N_{\rm gal}, \sigma_{\rm v, 1D}, N_{\rm phot}, N_{\rm phot, 500c}\}$) only reaches a minimum scatter of ~ 30%, which is still higher than that of the individual idealised observables apart from $\sigma_{v,1D,true}$. These mass scatters of realistic and idealised proxy observables set the lower and upper bounds, respectively, for what we can expect to be the predictive performance of the ML models presented in subsequent sections.

4 CNNS FOR SINGLE-BAND X-RAY

In this section, we describe, implement, and test CNN models for X-ray mass estimation in the context of eROSITA mock observations generated in Section 2. The models presented in this section take the single-band photon maps as input and produce an estimate of logarithmic M_{500c} as output. This problem framework allows for direct comparison to baseline scalar proxies derived from photon counts, as well as previous works in deep learning X-ray modelling presented by Ntampaka et al. (2019), Green et al. (2019), and Yan et al. (2020).

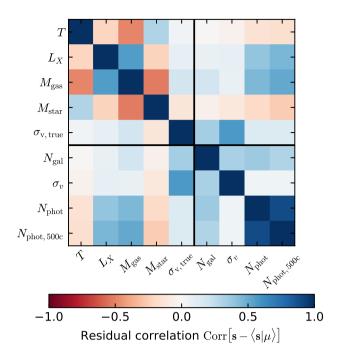


Figure 3. Residual correlation matrix of the scalar observables which form our comparative baseline for cluster mass inference. This covariance is presented under the power law mass dependence of the model described in Section 3. The values shown here are equal to the covariance described in Equation 4 scaled by each observable's individual scatter, i.e., $C_{ij}/(\sigma_{s_i|\mu}\sigma_{s_j|\mu})$. The black vertical and horizontal lines divide the idealized and the realistic observables, as described in Table 2. The correlation here is taken from clusters in our mock catalog at redshift z = 0.17.

4.1 Modelling

The data-driven ML models applied in this work are deep neural networks (Goodfellow et al. 2016), a class of popular data science tools that have been taught to model increasingly complex problems, both in science (Carleo et al. 2019) and beyond (LeCun et al. 2015). Functionally, we can think of deep neural networks as a class of highly non-linear and differentiable functions $f(\mathbf{x}; \boldsymbol{\theta})$, where \mathbf{x} is the model input (in our case, X-ray images) and θ is the model parameterization (i.e. weights and biases). The objective of training a neural network is to find some parameter setting θ^* that minimises a loss function $\mathcal{L}(\theta, \{x, y\})$ given a training sample $\{x, y\}$. Typically, this optimisation is done using a version of gradient descent. The loss function characterises what we want the neural network to do. In our case, we have a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, m_i)\}_{i=1}^N$ of N pairs of X-ray images $\mathbf{x}^{(i)}$ and known logarithmic masses $m = \log_{10} M_{500c}$. Our goal is to learn from this dataset a model parameterization that best predicts the masses given the X-ray images. To do so, we use a mean squared error (MSE) loss averaged over the training set,

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^{N} \left[f\left(\mathbf{x}^{(i)}, \boldsymbol{\theta}\right) - m_i \right]^2.$$
 (8)

Optimisation of this function in θ results in a neural network tuned to predict the mean of the posterior distribution of the logarithmic cluster mass given the X-ray image input (Ho et al. 2021).

The specific neural network model used to process X-ray images for this task is a convolutional neural network (CNN; LeCun et al. 1998). CNNs are widely recognised as the gold standard for com-

puter vision tasks in machine learning, due to their innate ability to learn and understand non-linear features from subregions of images. We use a feed-forward CNN architecture inspired by the model introduced in Ntampaka et al. (2019). Following the 128×128 single-channel input image, the network architecture is as follows:

- 1. 2D convolutional layer with 24 filters of size 5×5
- 2. 2D convolutional layer with 10 filters of size 3×3
- 3. Max-pooling layer of pooling size 4×4
- 4. 2D convolution layer with 24 filters of size 5×5
- 5. 2D convolutional layer with 10 filters of size 3×3
- 6. Max-pooling layer of pooling size 2×2
- 7. Flattening operation
- 8. Dense layer with output size of 128×1
- 9. Dense layer with output size of 128×1
- 10. Dense layer with output size of 64×1
- 11. Dense layer with output size of 1×1

The final layer with one node represents the predictor variable, i.e., the logarithmic M_{500c} . All layers except the last use a rectified linear activation function (ReLU) and an L2 weight regularisation with magnitude 10^{-4} . The model is trained using the Adam optimisation scheme (Kingma & Ba 2015) at a learning rate of 10^{-3} .

We use an eight-fold cross-validation procedure to evaluate the predictive performance of these models on our mock catalogue. To avoid potential train-test data leakage caused by correlated formation environments, the Magneticum clusters are first separated into eight equally-sized quadrants in the simulation volume. Each quadrant is assigned to a unique cross-validation fold. Eight independent models are then each trained on seven folds (~ 85% of the data) and tested on the last eighth fold. The tested fold is not repeated among all of our models, so we end with a set of predictions for all clusters in our catalogue. This process ensures that no clusters (or their local neighbors) can be used for both training and testing at the same time. During training, the training folds are further randomly divided into 90% training data and 10% validation data. We use an early-stopping criterion of no validation loss improvement over 20 epochs. Model training generally converges within ~ 150 epochs. During testing, we measure predictive bias and scatter averaged across the eight cross-validation folds.

4.2 Results

Figure 4 shows the distribution of predicted masses and residual scatter of the clusters in our independent test set relative to the best realistic observable, $N_{\rm phot,500c}$. We characterise the performance of this model in terms of predictive residuals, defined as

$$\epsilon = \log_{10} \left[\frac{M_{\text{pred}}}{M_{500c}} \right]. \tag{9}$$

The variance of this quantity integrated over our test set asymptotically approaches the Bayesian predictive variance of our estimators under a fixed-variance Gaussian likelihood and is complementary to the predictive variances described in Section 3. We note that the $N_{\rm phot,500c}$ scatter shown in Figure 4 was calculated by fitting a power-law regression on the training set and evaluating it on the test set. It is slightly different to the variance estimate calculated with Equation (7) due to the finite size and specific prior of the test set.

The mass predictions of our CNN model demonstrate remarkably low scatter, with an average residual of $\sim 0.1\%$ and a residual scatter of $\sim 18\%$. This scatter is lower than the scatter of $N_{\rm phot,500c}$ by 48% and even outperforms the maximal combination of realistic observables (i.e. $\{N_{\rm gal}, \sigma_{\rm v,1D}, N_{\rm phot}, N_{\rm phot,500c}\}$) by $\sim 35\%$ (see Table 3).

Model	Scatter (dex)	Scatter (%)
T	0.0510	11.75
L_X	0.1144	26.33
$M_{ m gas}$	0.0367	8.46
$M_{ m star}$	0.0545	12.56
$\sigma_{ m v,1D,true}$	0.1631	37.56
$N_{ m gal}$	0.1725	39.72
$\sigma_{ m v,1D}$	0.3057	70.39
$N_{ m phot}$	0.1738	40.02
N _{phot,500c}	0.1479	34.07
$M_{\rm gas}, M_{\rm star}$	0.0214	4.93
$T, M_{\rm gas}, M_{\rm star}$	0.0176	4.06
$T, L_X, M_{\rm gas}, M_{\rm star}$	0.0170	3.90
$T, L_X, M_{\rm gas}, M_{\rm star}, \sigma_{\rm v,1D,true}$	0.0169	3.90
$N_{\rm gal}, N_{\rm phot,500c}$	0.1252	28.82
$N_{\rm gal},\sigma_{\rm v,1D},N_{\rm phot,500c}$	0.1198	27.59
$N_{\rm gal},\sigma_{\rm v,1D},N_{\rm phot},N_{\rm phot,500c}$	0.1187	27.34
Single-band X-ray CNN	0.0773	17.80
Multi-band X-ray CNN	0.0703	16.18
X-ray+Spec-z CNN	0.0691	15.90

Table 3. Residual scatter for predicting M_{500c} with baseline proxy models (Section 3) and machine learning models (Sections 4 and 5). Each row shows the predictive scatter in dex and percent, for reference. Models are divided by their type, in the order: idealised observables, realistic observables, combinations of ideal proxies, combinations of realistic proxies, and machine learning models. We performed an exhaustive search over all combinations of observables, but only listed here the lowest scatter combination for each number of observables. The lowest scatter in each section is highlighted, for reference. The reported scatter for each observable was calculated using Equation (7).

We note that this is also lower than the scatter of real X-ray measurements calibrated with weak lensing, which is on the order of 20% - 50% (Zhang et al. 2008; Mahdavi et al. 2013). Our models are also $\sim 32\%$ more accurate than the idealised L_X proxy model, suggesting that the CNN finds additional information in the X-ray photon maps beyond the bolometric luminosity. An example of such information can be the surface density profile, which is informative of $M_{\rm gas}$, a much lower scatter mass estimator. Another piece of information can be morphology measures, including concentration. Morphology measures are correlated with the formation or accretion that can explain scatter about the mass–X-ray bolometric luminosity relation (Hartley et al. 2008; Parekh et al. 2015; Fujita et al. 2018; Farahi et al. 2020).

Ntampaka et al. (2019) and Yan et al. (2020) used neural network estimators trained on mock X-ray observations built from the IllustrisTNG (Nelson et al. 2019) and BAHAMAS (McCarthy et al. 2017) hydrodynamical simulations, respectively. Both studies built mock catalogues from models of ICM emission and configured observation time and resolution for the Chandra telescope. Relative to this previous work, the mocks used in this study have a lower signal-to-noise ratio, with eROSITA's lower resolution and observation time, and more realistic observation systematics, with the newly added inclusion of instrument response and AGN contamination. As a result, the residual mass scatters reported in these studies, ~ 12% and ~ 16%, respectively, are inevitably lower than the scatter derived from our single-band models (~ 18%). This impact of adding realistic errors to mock data has been seen in other works as well. For example the SZ-based estimators of de Andres et al. (2022) saw a similar, notable increase in scatter between mass predictions on idealised (12% scatter) and Planck-like (20% scatter) mock observations. We argue that the constraints presented here will serve as accurate and robust forecasts for the true bounds of the mass reconstruction of clusters in real eROSITA data.

Green et al. (2019) produced a similar study using ensemble regressors trained on Chandra and eROSITA mocks from the Magneticum simulation. This model extracted morphological parameters such as surface brightness concentration, smoothness, and asymmetry from each X-ray image and used them as features in a random forest regression (Breiman 2001). These models achieved a ~ 16% mass scatter for estimators when trained on either Chandra or eROSITA X-ray maps, a notable result considering the differences in signal-to-noise and telescope design of each sample. This is a ~ 2% improvement on the single-band mass scatter achieved in this work. However, compared to this work, the mocks used in Green et al. (2019) did not include the presence of AGN sources and were calculated over a shallow redshift range $0 \le z \le 0.29$, compared to the range $0 \le z \le 0.47$ of this analysis. In addition, the calculation of some morphological parameters in Green et al. (2019) requires assumed knowledge of the cluster R_{500c} , which in practice is subject to scatter. In contrast, the mass constraints presented in this work implicitly incorporate uncertainties of AGN contamination, redshiftdependent effects, and imperfect knowledge of R_{500c} , and therefore are a high-fidelity forecast of the expected eROSITA mass reconstruction. Nonetheless, we recommend that these methods, and those of Ntampaka et al. (2019) and Yan et al. (2020), be evaluated in equal contexts in a future comparative study.

5 MULTIWAVELENGTH MODELS

One challenge of next-generation survey science will be how to usefully combine nonlinear information in multiwavelength probes to perform increasingly precise astronomical inference. In this section, we study improvements to the single-band CNNs of the previous section by including multi-band and spectroscopic information in our model input.

5.1 Multi-band X-ray

The stratification of X-ray photons into several wide energy bands is useful for constraining cluster mass (Yan et al. 2020). This stratification allows one to utilise the different spectral profiles of ICM, AGN, and background noise sources to better characterise the physical modelling. eROSITA's instrument will scan the sky in the X-ray band, recording both photon counts and pulse height amplitudes (PHAs), a rough proxy for photon frequency. This will then allow us to roughly split our X-ray observations into different photon energy bands, namely the soft (0.5-1.2 keV), medium (1.2-2.0 keV), and hard bands (2.0-7.0 keV).

Here, we apply a neural network on the multi-band X-ray images of Section 2.1 to place tighter constraints on cluster mass. The modelling of multi-band X-ray images is nearly identical to that of single-band images in Section 4.1, except for the fact that we are now dealing with multi-channel inputs. The photon maps, now split into eROSITA's soft, medium, and hard bands, are concatenated together like RGB channels in a photograph. They are then subject to the same CNN architecture as detailed for the single-band images, except that the convolutional filters in the first layer have three channels instead of one. The multi-band models are trained using exactly the same procedure as the single-band models, including the same optimiser, learning rate, and ten-fold cross-validation split.

Figure 4 shows the distribution of mass predictions made by our multi-band model on the independent test set. Remarkably, the use of

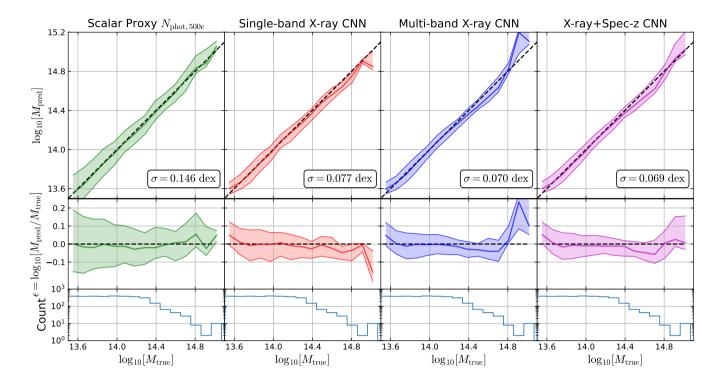


Figure 4. True versus predicted cluster mass for baseline scalar proxies and ML models. Each column displays the predictive performance of a single model in our study. The upper subplots show the median and $[16^{th}-84^{th}]$ percentile confidence intervals of the predicted masses, while the middle subplots show the corresponding residuals (Equation 9). The bottom plots show the histogram of true cluster masses in the test set, which is the same for each model. From left-to-right, these models are the proxy model for the number of photons detected within R_{500c} (Section 3), the single-band X-ray CNN model (Section 4), the multi-band X-ray CNN model (Section 5.1), and the CNN model using both X-ray and dynamical information (Section 5.2), at maximum galaxy sampling. The mass definition used here for M_{true} and M_{pred} is M_{500c} with units $h^{-1}M_{\odot}$.

multi-band images reaches a scatter of $\sim 16.2\%$, a reduction of 9% beyond the single-band scatter. This suggests that the model should utilise the knowledge of band separation to form its predictions better. This interpretation is further explored in Section 6.3. The residual scatter of the multi-band model greatly outperforms both the idealised L_X and realistic observables. Also, like the single-band models, the multi-band models have a low mean residual of $\sim -0.2\%$ but exhibit mean biases at the high-mass end. This is a result of limited training and test data in the high-mass regime and not necessarily indicative of poor modelling, as evident in previous works (Ntampaka et al. 2019; Green et al. 2019).

During the editorial review of this paper, the preprint of Krippendorf et al. (2023) was made public wherein the authors applied a very similar CNN architecture to mock multiband eROSITA images from the eFEDS simulations (Comparat et al. 2019). In this approach, the authors split eROSITA observations into ten energy bands, which served as distinct channels for input into their CNN model. They then fit a regression over cluster mass using a maximum likelihood loss function. Due to methodological differences between Krippendorf et al. (2023) and this work (e.g. the type of base simulation, method of mock generation, cluster selection function, image preprocessing, and training loss), it is not possible to make a thorough, apples-to-apples comparison of the two approaches. Despite these differences, it's worth noting that the conclusions of both manuscripts largely agree on the expected level of mass scatter for eROSITAobserved clusters (16.2% in this work vs 18.8% from Krippendorf et al. (2023)).

5.2 Joint X-ray and Spectroscopic Data

The information encoded in the dynamics of cluster galaxies can be a very potent probe of the system mass. Through the use of spectroscopic follow-up campaigns such as The SPectroscopic IDentification of eROSITA Sources (SPIDERS; Furnell et al. 2018), the eROSITA survey will begin to build wide-field multi-wavelength representations of galaxy clusters including both X-ray and dynamical data. Previous works (Ho et al. 2019, 2021; Kodi Ramanah et al. 2020, 2021) have shown that neural networks can effectively model the connection between dynamical information and cluster mass. However, most of these studies have been confined to predictions of larger mass definitions, such as M_{200c} . Part of the reason for this is that dynamics are a very weak probe of cluster core physics, as is evident from the $\sim 38\%$ scatter of $\sigma_{v,1D,true}$ in our baseline analysis (Section 3). However, it is possible that the inclusion of spectroscopic catalogues will help explain the X-ray systematics, particularly emission from halo environments or projected AGN. So, the question remains as to whether the inclusion of spectroscopic probes will be useful in mass reconstruction for eROSITA clusters or whether X-ray maps capture most of the relevant information.

To test this hypothesis, we construct a neural network architecture to process the X-ray multi-band images from Section 2.1 and dynamical cubes built in Section 2.2. The architecture is designed to learn informative, localised features from both X-ray and dynamical data and combine them to produce a maximally informative mass reconstruction. The model architecture is split into three sections, a 2D convolutional feature extractor for X-ray-only input, a 3D convo-

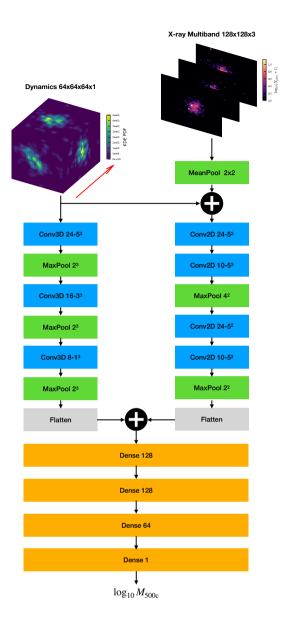


Figure 5. Flowchart neural network architecture for combining X-ray and spectroscopic information as described in Section 5.2. Each box represents a tensor transformation applied to the previous step. Convolutional layers shown in blue include the number of filters and the filter size. Pooling layers in green show the pooling size. Dense layers in orange show the hidden layer width, and are each followed by a dropout layer of probability p = 0.1. 'Plus' sign operators indicate a concatenation operation. All layers use a Rectified Linear Unit (ReLU) activation function.

lutional extractor for X-ray and dynamical input, and a dense network to take the joint features and produce a cluster mass estimate. The motivation behind this choice is that the X-ray data is the primary driver of the M_{500c} information, and the dynamical input is only added to augment or explain the features of the X-ray. Therefore, X-ray inputs have their own path to compression, while the dynamical inputs are inherently tied to the projected X-ray photon counts at their respective sky positions. The X-ray feature extractor has the same design as the convolutional stages of our multi-band model, whereas the dynamical extractor utilises a 3D convolutional architecture following from Kodi Ramanah et al. (2021). The entire neural network

pipeline, including both extractor and dense networks, is trained simultaneously. A representation of the full architecture is shown in Figure 5.

The optimiser, learning rate, and training procedures for this model are identical to those of the single-band model (Section 4). However, because this model utilises 3D convolutions, the forward pass and backpropagation operations are considerably slower to compute. In total, training the 3D convolutional model takes $\sim 15\times$ longer than training the 2D convolutional models on a single Nvidia Ampere A100 GPU.

Figure 4 shows that the inclusion of dynamical information adds little improvement to the performance of the multi-band models. The joint X-ray multi-band and dynamical model reaches a test scatter of $\sim 15.9\%$, which is only slightly smaller than that of the multi-band model at ($\sim 16.2\%$). These results are even assuming full spectroscopic follow-up, i.e., that all galaxies within our selection cut were measured spectroscopically as well. In practice, we only will receive a partial picture of the spectroscopic redshifts of galaxies around each cluster, suggesting that the scatter we achieve here is a lower bound for reality. Although it is possible that the modelling and training applied here are not sufficiently flexible to learn complex multiwavelength relationships, it is rather more likely that most of the information about M_{500c} contained in dynamical probes is equally well covered by X-ray measurements (e.g. Nagai et al. 2007a).

6 INVESTIGATING CLUSTER MASS ESTIMATES

In this section, we analyse the models presented in the previous sections to derive an understanding of the learned behaviour which allows for improved mass estimates. We test several hypotheses presented in the literature to explain improved neural network performance, as well as present new findings from our explainability study.

The primary diagnostic tool used for our interpretation of neural network behaviour is the saliency map. Saliency maps are popular explainers for image-recognition tasks and function by quantifying the importance of individual pixels in an input image for a given modeling task. The specific method we use is a gradient-based saliency map (Simonyan et al. 2014), which measures the pixel-wise sensitivity to model outputs by taking the gradient of the output with respect to the input. Specifically, we can consider a saliency map as the output of applying the operator $\mathbb S$ to a given model output $f(\mathbf x; \boldsymbol{\theta})$ evaluated at a single input image $\mathbf x$, as defined:

$$\mathbb{S}\left[f(\mathbf{x};\boldsymbol{\theta})\right] = \left|\frac{\partial f(\mathbf{x};\boldsymbol{\theta})}{\partial \mathbf{x}}\right|_{\mathbf{x}=\mathbf{x}}.$$
 (10)

The result of this operation is an image of the same dimensionality as \mathbf{x} , wherein each pixel is a measure of the local sensitivity of the model output with respect to the corresponding input pixel. Here, we use the absolute value of the gradient as a sensitivity metric instead of a measure of relative directional change. In general, the function $f(\mathbf{x}; \boldsymbol{\theta})$ can be any differentiable parametric function, but, in this work, we consider f to be estimators of logarithmic mass, either as the scaling models from Section 3 or the neural network models from Sections 4 and 5. For our application, a high value derived from Equation (10) suggests that the current setting of the pixel value is important for mass prediction. Figure 6 shows several examples of single-band X-ray images and their corresponding saliency maps when subject to our CNN mass estimator.

As an example, we derive the expected saliency for a baseline scaling model using the N_{phot} observable (Section 3). Following our proxy formalism (Evrard et al. 2014), a single scalar observable s_a

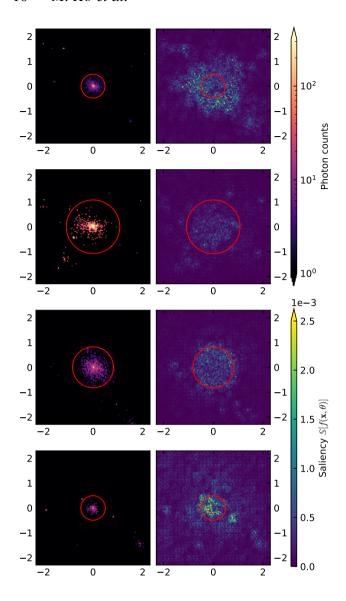


Figure 6. Four example single-band eROSITA X-ray mocks and their corresponding saliency maps when subject to CNN mass estimators. Each row represents a random cluster in our test set. The left column shows the original single-band X-ray image input, while the right column shows the saliency following Equation (10). Each subplot also shows a red circle indicating the R_{500c} of each cluster for scale reference. Subplot x- and y-axes are labelled in units of h^{-1} Mpc.

will predict a mean logarithmic mass $\ln M$ as,

$$\langle \ln M \mid s_a \rangle = (s_a - \pi_a)/\alpha_a - \beta \sigma_{\ln M \mid a}^2,$$
 (11)

where π_a and α_a are the scaling intercept and slope parameters for proxy s_a , β is the local slope of the halo mass function (Equation (5)), and $\sigma_{\ln M|a}^2$ is the proxy predictive variance (Equation (7)). For a photon count based proxy, our proxy is $s_a = \ln N_{\rm phot} = \ln \left[\sum_{ij} n_{ij} \right] = \ln \left[\sum_{ij} (10^{x_{ij}} + 1) \right]$, wherein n_{ij} are the photon counts in each pixel and x_{ij} are their transformed counterparts after normalisation with Equation 2. We note that the estimator for the $N_{\rm phot}$,500c proxy follows this same form, except summed over a smaller radius. Using this definition, we can turn our proxy model of mean logarithmic mass

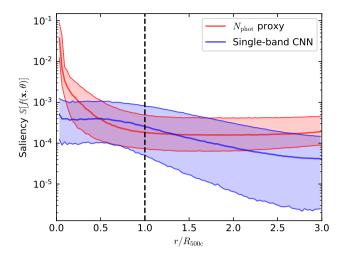


Figure 7. Distribution of input saliencies as a function of radius from the centre of the image for two different mass estimators, a scalar $N_{\rm phot}$ proxy in red and a CNN model in blue. We show the median and 16-84th confidence interval of each distribution, binned over a radius. Saliency values for the scalar proxy were calculated analytically with Equation (13), whereas those for the CNN were calculated via backpropagation with Equation (10).

(Equation 11) into a mass estimator,

$$g(\mathbf{x}; \boldsymbol{\theta}) = \frac{\langle \ln M \mid \ln N_{\text{phot}} \rangle}{\ln 10},$$
(12)

where $\theta = \left\{ \pi_{\ln N_{\text{phot}}}, \alpha_{\ln N_{\text{phot}}}, \beta, \sigma_{\ln M \mid \ln N_{\text{phot}}}^2 \right\}$ and the ln 10 factor scales the $g(\mathbf{x}; \theta)$ output to $\log_{10} M$. Then, we can apply our saliency operator to this N_{phot} proxy mass estimator to infer the following closed-form expression:

$$\mathbb{S}\left[g(\mathbf{x};\boldsymbol{\theta})\right]_{ij} = \frac{10^{x_{ij}}}{\alpha_{\ln N_{\text{phot}}} N_{\text{phot}}} \propto \frac{n_{ij} + 1}{N_{\text{phot}}}.$$
 (13)

Under this model, the saliency would suggest that the more photons a pixel detects, the more important it is to the model. This is an undesirable property of the $N_{\rm phot}$ scaling model, especially as we seek to exclude, for example, contaminating AGN sources in exchange for more informative ICM emission. In subsequent sections, we will use this saliency example as a relative comparison to the CNN models under investigation.

6.1 Cluster cores

Ntampaka et al. (2019), Yan et al. (2020), de Andres et al. (2022), and Ntampaka & Vikhlinin (2022) also apply interpretability methods to explain the decisions of cluster mass estimators. In the former three studies, the authors apply the Google Deep Dream algorithm (Mordvintsev et al. 2018) as their interpretability diagnostic, a gradient-based saliency variant that uses gradient ascent to find the perturbation to the input image which would result in the maximal output shift. In the latter study, traditional gradient saliency is used. All four studies qualitatively find that the neural network significantly down-weights the cluster core region when making a mass prediction. The explanation for this behaviour is that the cluster core is unique with its highly stochastic, non-linear physics and thus is a poor signal from which to derive cluster mass.

Here, we quantitatively test the hypothesis that neural networks

learn to down-weight the central regions of an X-ray image when predicting M_{500c} . Figure 7 shows the distribution of saliency values assigned to pixels within and outside of the cluster core as a function of radius for our CNN model. The CNN saliency is shown relative to that of the $N_{\rm phot}$ proxy, the saliency for which we derived analytically in Equation (13). We compare these two models. In this context, we estimate the relative per-pixel importance given by (i) the $N_{\rm phot}$ proxy model and (ii) the CNN model.

Figure 7 shows that the radial distribution of saliency values for the $N_{\rm phot}$ proxy model and the CNN model differ significantly at the central and outer regions of clusters. We observe that, although the neural network models are not totally disqualifying information in the central regions of the cluster, the information in the core is significantly down-weighted considering the number of detected photons in the region. At around $r = R_{500c}$, the prediction sensitivity for $N_{\rm phot}$ and the CNN is close to the same, suggesting that, in this ICM region, the CNN effectively counts the number of photons to predict mass. However, moving below this radius, the $N_{\rm phot}$ saliency increases exponentially, whereas the CNN saliency remains constant. At a radius of $r = 0.07 R_{500c}$, the median N_{phot} saliency is approximately 10× that of the CNN model. This is caused by the fact that high photon counts in the core of the cluster are extremely crucial to determining the mass prediction with the $N_{\rm phot}$ proxy, but only mildly important to the neural network. Another interesting result from this diagnostic is that beyond $r > R_{500c}$, the CNN starts again to down-weight the importance of pixels. We note that the CNN model does not directly estimate the true R_{500c} , though it may infer it from the mass-luminosity relation or the comoving surface brightness profile. This may be indicative of the fact that CNN recognises that, outside of $r > R_{500c}$, the presence of environmental feedback and AGN contamination makes this region an unreliable probe of central cluster mass and therefore disregards it.

6.2 AGN Attention

Another hypothesis for improving neural networks is the ability of CNN models to recognise and remove contaminating artefacts within subregions of an image (Ho et al. 2019; Kodi Ramanah et al. 2020). In the case of X-ray probes, this behaviour would allow the CNN to identify AGN contaminates and mitigate their contribution to the total photon emission. The separation of sources is impossible to do exactly in real observations, but we can study the CNN sensitivity to various sources in the idealised environment of simulation.

Inspired by Benchmarking Attribution Methods (BAMs; Yang & Kim 2019), we construct a quantitative test to measure the attention paid by the neural network to ICM and AGN sources. We can consider a given input image \mathbf{n} to consist of a contribution of ICM emission \mathbf{n}_{ICM} and AGN emission \mathbf{n}_{AGN} , such that $\mathbf{n} = \mathbf{n}_{\text{ICM}} + \mathbf{n}_{\text{AGN}}$. Then, for each input image $\mathbf{x}^{(i)}$, we can define two binary masks, $\mathbf{M}_{\text{ICM}}^{(i)}$ and $\mathbf{M}_{\text{AGN}}^{(i)}$, which quantify which pixels were dominated by each respective source, i.e., $\mathbf{M}_{\text{ICM}} = \Theta\left[\mathbf{n}_{\text{ICM}} - \mathbf{n}_{\text{AGN}}\right]$ and $\mathbf{M}_{\text{AGN}} = \Theta\left[\mathbf{n}_{\text{AGN}} - \mathbf{n}_{\text{ICM}}\right]$, where Θ is the element-wise Heaviside function. Then, we can define a functional \mathbb{S}_{mask} which quantifies the average saliency for all pixels integrated over each source mask, as follows:

$$\mathbb{S}_{\text{mask}}\left[f(\mathbf{x};\boldsymbol{\theta});\mathbf{M}\right] = \frac{1}{\sum_{ij} M_{ij}} \sum_{ij} M_{ij} \cdot \mathbb{S}\left[f(\mathbf{x};\boldsymbol{\theta})\right]_{ij},\tag{14}$$

where all sums are taken over all matrix elements. Figure 8 shows an example cluster's ICM and AGN components, as well as its binary masks and saliency map. To compare global properties of each source mask, we further derive a summary statistic of Equation 14

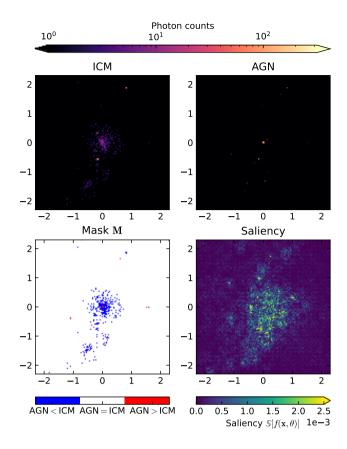


Figure 8. Example ICM and AGN emissions and corresponding saliency for mock X-ray observations of an example cluster. This demonstrates the relative attention paid by our mass estimation model to respective sources in the X-ray image, as discussed in Section 6.2. Top left: ICM emission. Top right: AGN emission. Bottom left: Mask of ICM vs. AGN contributions for use in Equation 14. Bottom right: Saliency map for a trained single-band CNN mass estimator derived from Equation 10.

as $\langle \mathbb{S}_{\text{mask}} [f(\mathbf{x}; \boldsymbol{\theta}); \mathbf{M}] \rangle$, wherein the operation $\langle \cdot \rangle$ is taken as an average over all inputs $\mathbf{x}^{(i)}$ and corresponding masks $\mathbf{M}^{(i)}$ in our test set.

For CNN models, we find that the contribution to the saliency from AGN sources is approximately log-normally distributed with median and [16th-84th] percentile interval of $\mathbb{S}_{\text{mask}}\left[f(\mathbf{x}; \boldsymbol{\theta}); \mathbf{M}_{\text{AGN}}\right] \sim 10^{-3.85^{+0.71}_{-0.74}}$, whereas that of ICM sources is $\mathbb{S}_{\text{mask}}[f(\mathbf{x}; \boldsymbol{\theta}); \mathbf{M}_{\text{ICM}}] \sim 10^{-3.53^{+0.47}_{-0.67}}$. While the range of each distribution is wide, we note that the median value of ICM saliencies is twice that of AGN saliencies. This is particularly notable given that pixels dominated by AGN emission, that is when $\mathbf{n}_{AGN} > \mathbf{n}_{ICM}$, receive on average > 13 times the number of photons as those dominated by ICM emission, a result of the intense radiation emitted from AGN sources. From Equation (13), this would suggest that AGN-dominated pixels would attend for 13 times the saliency of ICM pixels for a typical N_{phot} proxy scaling relation. The fact that the CNN saliency of AGN sources is less than that of ICM sources clearly indicates that CNN models treat AGN sources with reduced importance compared to ICM sources, a definitive departure from the behaviour of scalar proxy-based methods.

However, we note that the AGN attention in our models is nonzero, suggesting that changing AGN photon counts would indeed impact our CNN mass estimates, albeit in a reduced capacity to ICM emis-

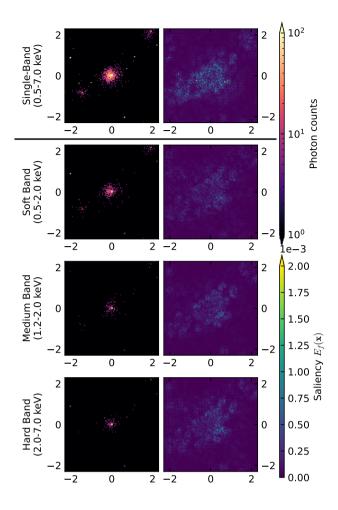


Figure 9. X-ray photon and saliency maps across different X-ray bands for a single example cluster in our test set. The top row shows the input X-ray photon map and resultant saliency map when doing mass inference with a single-band X-ray CNN. The bottom three rows show the photon maps and saliencies for the same mass inference with a multi-band X-ray CNN. The bottom images are separated into soft, medium, and hard bands.

sion. We propose an explanation for this in that the CNN model must learn the corpus of possible AGN profiles in order to recognise their location and remove their photon counts. In this case, any pixel-wise deviation of the established profile will significantly affect the recognition of the AGN source, thereby deviating the result. Qualitative evidence for this behaviour exists in Figure 8 wherein pixels with no detected photons located between the extended AGN source and the central ICM are given high saliency. If these pixels were to be non-zero, the model would recognise evidence of the extended AGN source as being part of the ICM and accordingly increase its mass estimate. We recommend exploring this hypothesis in future studies.

6.3 Multi-band Attention

In Section 5.1, we showed that stratification of X-ray photon maps into soft, medium and hard bands improves the mass prediction performance of the neural network model. The intuition of previous work suggests that peak ICM emission occurs in the soft bands, while AGN emission is more dominant in the hard bands (Biffi et al. 2018). Therefore, the separation of the soft band from other energies

would provide a cleaner sample from which we can study the ICM. We apply the saliency maps to the multi-band models to study the sensitivity of predictions to each energy band.

Figure 9 shows the saliency maps for an example multi-band CNN prediction compared to the saliency map of a single-band prediction. We derive two interpretations on the behaviour of multi-band CNN models: First, the distribution of pixel saliencies in the soft, medium, and hard bands are very nearly identical, with median value and $[16^{\text{th}}-84^{\text{th}}]$ percentile intervals of $10^{-4.15^{+0.56}_{-0.68}}$, $10^{-4.17^{+0.59}_{-0.68}}$, and $10^{-4.15^{+0.60}_{-0.72}}$, respectively. The saliency intensities follow roughly the same spatial distribution for each cluster image. This is suggestive that multi-band information is closely shared across the channels by the CNN architecture but does not allow us to draw any definitive conclusions about the relative importance of measuring one band over another. Second, we clearly notice that AGN emission present in the input X-ray photon maps is not as significant in the saliency of the multi-band models as they are for the single-band models. This is reflected in repeating the analysis of Section 6.2 for the multi-band models, in which we find that the saliency of AGN-dominated pixels is $10^{-3.91^{+0.61}_{-0.74}}$ versus that of ICM-dominated pixels $10^{-3.77^{+0.49}_{-0.63}}$. We can then gather that the AGN-to-ICM saliency ratio for multi-band models is 50% lower than that of the single-band models. We suggest that this is a result of the spectral separation of AGN and ICM sources in multi-band images.

7 CONCLUSION

In this work, we present forecasts for the expected limits of mass reconstruction for eROSITA clusters with deep learning models. To ensure we place reliable constraints, we validate our methodology on a highly realistic catalogue of mock X-ray observations derived from the Magneticum hydrodynamic simulations. The mock catalog used here includes systematics such as AGN contaminants, cluster morphology, background emission, and eROSITA-specific instrument response, making it the most robust mock X-ray catalogue currently applied to deep learning mass estimation. Using this catalogue, we build a quantitative baseline of scalar mass proxies and characterise their intrinsic scatter and mutual information. We then introduce a CNN architecture based on the one used in Ntampaka et al. (2019) designed to learn the mapping between photon maps and logarithmic mass $\log_{10} M_{500c}$. For single-band, bolometric X-ray photon maps, we can constrain M_{500c} masses to within 18% scatter, a factor of two improvements on realistic mass proxies, and a 32% improvement on an idealised bolometric luminosity L_X measurement. This result suggests that neural networks successfully utilise high-order features to reduce predictive uncertainty.

We go on to suggest that these models can be improved with the inclusion of multiwavelength inputs. We show that a CNN model trained on X-ray photon maps separated into soft, medium, and high energy bands further reduces mass scatter to 16.2%, a 9% reduction from the single-band models. We also find that the inclusion of dynamical information from spectroscopic follow-up offers little-to-no scatter improvement for M_{500c} . We argue that M_{500c} information contained in cluster dynamics is mostly or entirely covered by photon maps.

Lastly, we investigate a series of hypotheses as to why CNNs show such substantial improvements over conventional scalar observables. We quantitatively demonstrate that CNNs down-weight the importance of photons emitted by cluster cores by greater than a factor of 10, motivating the belief that these features are too noisy for

reliable mass inference. We also show that the CNNs significantly downweight, but do not ignore, contamination from AGN emission in X-ray images. Lastly, we perform tests of feature attribution for multi-band X-ray models, but do not find definitive evidence that different X-ray bands contribute more or less to the overall mass prediction. However, we find that multi-band models are better at reducing the importance of AGN emission photons. We suggest this results from better source separation in different photometric bins.

In conclusion, DL models can learn a reliable, informative model of X-ray clusters, infer low-scatter estimates of their mass, and be well-calibrated with realistic mock simulations. Their improvements can be attributed to physically-motivated manipulation of information, including core excision and automatic removal of AGN contamination. We recommend these techniques for further application to the upcoming eROSITA cluster catalogue, complementing existing proxy-based mass estimators.

ACKNOWLEDGEMENTS

We thank the anonymous referee for the helpful comments and suggestions used to improve the manuscript during the review process. We acknowledge Klaus Dolag and Antonio Ragagnin for providing the Magneticum 2b simulations used in this work. MH is supported by the Simons Collaboration on "Learning the Universe" and was supported by NSF AI Institute: Physics of the Future, NSF PHY-2020295, and the McWilliams-PSC Seed Grant Program. JS, DN, and MN acknowledge support from the NASA ATP Grant (80NSSC22K0821). DN is also supported by NSF (AST-2206055) and NASA (TM3-24007X) grants. AE is supported by NASA grant (80NSSC22K0476). The computing resources necessary to complete this analysis were provided by the Pittsburgh Supercomputing Center. The X-ray mock data set was collected using computational resources at the Maryland Advanced Research Computing Center (MARCC).

DATA AVAILABILITY

The Magneticum simulations and respective halo catalogues and PHOX X-ray mocks used in this analysis are publicly available through the Cosmology Web Portal ¹. Mock eROSITA photometry can be generated using preset configurations in the SIXTE code² (Dauser et al. 2019). Pre-generated dynamical catalogues for Magneticum halos are available from the authors upon reasonable request.

All code for the data analysis and machine learning models investigated in this analysis has been made available on Github³. Preprocessed training and test catalogues and pre-trained models will be made available upon reasonable request.

REFERENCES

Abbott T., et al., 2020, Physical Review D, 102, 023509
Ade P. A., et al., 2011, Astronomy & Astrophysics, 536, A11
Ade P., et al., 2016, Astronomy & Astrophysics, 594, A24
Allen S. W., Evrard A. E., Mantz A. B., 2011, Annual Review of Astronomy and Astrophysics, 49, 409

Anbajagane D., Evrard A. E., Farahi A., Barnes D. J., Dolag K., McCarthy I. G., Nelson D., Pillepich A., 2020, Monthly Notices of the Royal Astronomical Society, 495, 686

Biffi V., Dolag K., Boehringer H., Lemson G., 2012, Monthly Notices of the Royal Astronomical Society, 420, 3545

Biffi V., Dolag K., Boehringer H., 2013, Monthly Notices of the Royal Astronomical Society, 428, 1395

Biffi V., et al., 2016, The Astrophysical Journal, 827, 112

Biffi V., Dolag K., Merloni A., 2018, Monthly Notices of the Royal Astronomical Society, 481, 2213

Breiman L., 2001, Machine learning, 45, 5

Carleo G., Cirac I., Cranmer K., Daudet L., Schuld M., Tishby N., Vogt-Maranto L., Zdeborová L., 2019, Reviews of Modern Physics, 91, 045002
Clerc N., et al., 2018, Astronomy & Astrophysics, 617, A92

Clerc N., et al., 2020, Monthly Notices of the Royal Astronomical Society, 497, 3976

Cohn J. D., Battaglia N., 2020, MNRAS, 491, 1575

Comparat J., et al., 2019, Monthly Notices of the Royal Astronomical Society, 487, 2005

Dauser T., et al., 2019, Astronomy & Astrophysics, 630, A66

Dolag K., Borgani S., Murante G., Springel V., 2009, Monthly Notices of the Royal Astronomical Society, 399, 497

Dolag K., Komatsu E., Sunyaev R., 2016, Monthly Notices of the Royal Astronomical Society, 463, 1797

Evrard A. E., et al., 2008, ApJ, 672, 122

Evrard A. E., Arnault P., Huterer D., Farahi A., 2014, Monthly Notices of the Royal Astronomical Society, 441, 3562

Farahi A., Evrard A. E., McCarthy I., Barnes D. J., Kay S. T., 2018a, Monthly Notices of the Royal Astronomical Society, 478, 2618

Farahi A., et al., 2018b, A&A, 620, A8

Farahi A., et al., 2019, Nature Communications, 10, 2504

Farahi A., Ho M., Trac H., 2020, MNRAS, 493, 1361

Fujita Y., Umetsu K., Rasia E., Meneghetti M., Donahue M., Medezinski E., Okabe N., Postman M., 2018, ApJ, 857, 118

Furnell K. E., et al., 2018, Monthly Notices of the Royal Astronomical Society, 478, 4952

Giles P. A., et al., 2017, MNRAS, 465, 858

Goodfellow I., Bengio Y., Courville A., 2016, Deep Learning. MIT Press Green S. B., Ntampaka M., Nagai D., Lovisari L., Dolag K., Eckert D.

Green S. B., Ntampaka M., Nagai D., Lovisari L., Dolag K., Eckert D., ZuHone J. A., 2019, ApJ, 884, 33

Hartley W. G., Gazzola L., Pearce F. R., Kay S. T., Thomas P. A., 2008, MNRAS, 386, 2015

Hirschmann M., Dolag K., Saro A., Bachmann L., Borgani S., Burkert A., 2014, Monthly Notices of the Royal Astronomical Society, 442, 2304

Ho M., Rau M. M., Ntampaka M., Farahi A., Trac H., Póczos B., 2019, The Astrophysical Journal, 887, 25

Ho M., Farahi A., Rau M. M., Trac H., 2021, The Astrophysical Journal, 908, 204

Ho M., Ntampaka M., Rau M. M., Chen M., Lansberry A., Ruehle F., Trac H., 2022, Nature Astronomy, 6, 936

Kay S. T., Peel M. W., Short C. J., Thomas P. A., Young O. E., Battye R. A., Liddle A. R., Pearce F. R., 2012, MNRAS, 422, 1999

Kingma D. P., Ba J., 2015, in Bengio Y., LeCun Y., eds, 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings. http://arxiv. org/abs/1412.6980

Kodi Ramanah D., Wojtak R., Ansari Z., Gall C., Hjorth J., 2020, Monthly Notices of the Royal Astronomical Society, 499, 1985

Kodi Ramanah D., Wojtak R., Arendse N., 2021, Monthly Notices of the Royal Astronomical Society, 501, 4080

Komatsu E., et al., 2009, The Astrophysical Journal Supplement Series, 180, 330

Kravtsov A. V., Vikhlinin A., Nagai D., 2006, The Astrophysical Journal, 650, 128

Krippendorf S., et al., 2023, arXiv preprint arXiv:2305.00016

LeCun Y., Bottou L., Bengio Y., Haffner P., 1998, Proceedings of the IEEE, 86, 2278

LeCun Y., Bengio Y., Hinton G., 2015, nature, 521, 436

http://www.magneticum.org/vos.html

https://www.sternwarte.uni-erlangen.de/sixte/

https://github.com/McWilliamsCenter/halo_cnn

Mahdavi A., Hoekstra H., Babul A., Bildfell C., Jeltema T., Henry J. P., 2013, The Astrophysical Journal, 767, 116

Mantz A., Allen S. W., Rapetti D., Ebeling H., 2010, MNRAS, 406, 1759 Mantz A. B., et al., 2016, MNRAS, 463, 3582

McCarthy I. G., Schaye J., Bird S., Le Brun A. M. C., 2017, MNRAS, 465, 2936

McClintock T., et al., 2019, Monthly Notices of the Royal Astronomical Society, 482, 1352

Merloni A., et al., 2012, arXiv preprint arXiv:1209.3114

Mordvintsev A., Pezzotti N., Schubert L., Olah C., 2018, Distill, 3, e12

Mulroy S. L., et al., 2019, MNRAS, 484, 60

Nagai D., 2006, ApJ, 650, 538

Nagai D., Vikhlinin A., Kravtsov A. V., 2007a, The Astrophysical Journal, 655, 98

Nagai D., Kravtsov A. V., Vikhlinin A., 2007b, ApJ, 668, 1

Nelson D., et al., 2019, Computational Astrophysics and Cosmology, 6, 1

Ntampaka M., Vikhlinin A., 2022, The Astrophysical Journal, 926, 45

Ntampaka M., Trac H., Sutherland D. J., Battaglia N., Póczos B., Schneider J., 2015, ApJ, 803, 50

Ntampaka M., et al., 2019, The Astrophysical Journal, 876, 82

Ntampaka M., Ho M., Nord B., 2021, arXiv preprint arXiv:2111.14566

Old L., et al., 2018, Monthly Notices of the Royal Astronomical Society, 475, 853

Parekh V., van der Heyden K., Ferrari C., Angus G., Holwerda B., 2015, A&A, 575, A127

Pillepich A., Reiprich T. H., Porciani C., Borm K., Merloni A., 2018, Monthly Notices of the Royal Astronomical Society, 481, 613

Planelles S., Borgani S., Fabjan D., Killedar M., Murante G., Granato G., Ragone-Figueroa C., Dolag K., 2014, Monthly Notices of the Royal Astronomical Society, 438, 195

Pop A.-R., et al., 2022a, arXiv e-prints, p. arXiv:2205.11528

Pop A.-R., et al., 2022b, arXiv e-prints, p. arXiv:2205.11537

Pratt G. W., Croston J. H., Arnaud M., Böhringer H., 2009, A&A, 498, 361

Pratt G., Arnaud M., Biviano A., Eckert D., Ettori S., Nagai D., Okabe N., Reiprich T., 2019, Space Science Reviews, 215, 1

Ragagnin A., Dolag K., Biffi V., Bel M. C., Hammer N. J., Krukau A., Petkova M., Steinborn D., 2017, Astronomy and Computing, 20, 52

Raghunathan S., et al., 2022, ApJ, 926, 172

Reichardt C., et al., 2013, The Astrophysical Journal, 763, 127

Rozo E., Evrard A. E., Rykoff E. S., Bartlett J. G., 2014, Monthly Notices of the Royal Astronomical Society, 438, 62

Saro A., Mohr J. J., Bazin G., Dolag K., 2013, The Astrophysical Journal, 772, 47

Schellenberger G., Reiprich T., Lovisari L., Nevalainen J., David L., 2015, Astronomy & Astrophysics, 575, A30

Scott D. W., 2015, Multivariate density estimation: theory, practice, and visualization. John Wiley & Sons

Simonyan K., Vedaldi A., Zisserman A., 2014, in Bengio Y., LeCun Y., eds, 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Workshop Track Proceedings. http://arxiv.org/abs/1312.6034

Soltis J., Ntampaka M., Wu J. F., ZuHone J., Evrard A., Farahi A., Ho M., Nagai D., 2022, The Astrophysical Journal, 940, 60

Springel V., White S. D., Tormen G., Kauffmann G., 2001, Monthly Notices of the Royal Astronomical Society, 328, 726

Vikhlinin A., et al., 2009, ApJ, 692, 1060

Wadekar D., et al., 2023a, Proceedings of the National Academy of Sciences, 120, e2202074120

Wadekar D., et al., 2023b, Monthly Notices of the Royal Astronomical Society, 522, 2628

White M., Cohn J., Smit R., 2010, Monthly Notices of the Royal Astronomical Society, 408, 1818

Wojtak R., et al., 2018, Monthly Notices of the Royal Astronomical Society, 481, 324

Yan Z., Mead A., Van Waerbeke L., Hinshaw G., McCarthy I., 2020, Monthly Notices of the Royal Astronomical Society, 499, 3445

Yang M., Kim B., 2019, CoRR, abs/1907.09701

Zhang Y.-Y., Finoguenov A., Böhringer H., Kneib J.-P., Smith G., Kneissl R., Okabe N., Dahle H., 2008, Astronomy & Astrophysics, 482, 451 de Andres D., et al., 2022, Nature Astronomy, 6, 1325

This paper has been typeset from a $T_EX/I = T_EX$ file prepared by the author.