

INVERTIBLE MAPPING BETWEEN FIELDS IN CAMELS

Sambatra Andrianomena^{1,2} Sultan Hassan^{3,2,4} Francisco Villaescusa-Navarro⁴

ABSTRACT

We build a bijective mapping between different physical fields from hydrodynamic CAMELS simulations. We train a CycleGAN on three different setups: translating dark matter to neutral hydrogen (Mcdm-HI), mapping between dark matter and magnetic fields magnitude (Mcdm-B), and finally predicting magnetic fields magnitude from neutral hydrogen (HI-B). We assess the performance of the models using various summary statistics, such as the probability distribution function (PDF) of the pixel values and 2D power spectrum ($P(k)$). Results suggest that in all setups, the model is capable of predicting the target field from the source field and vice versa, and the predicted maps exhibit statistical properties which are consistent with those of the target maps. This is indicated by the fact that the mean and standard deviation of the PDF of maps from the test set is in good agreement with those of the generated maps. The mean and variance of $P(k)$ of the real maps agree well with those of generated ones. The consistency tests on the model suggest that the source field can be recovered reasonably well by a forward mapping (source to target) followed by a backward mapping (target to source). This is demonstrated by the agreement between the statistical properties of the source images and those of the recovered ones.

1 INTRODUCTION

The upcoming generation of surveys (e.g. SKA) will be able to map neutral hydrogen via HI intensity mapping (Santos et al., 2015). This powerful cosmological probe will help us further our understanding of large-scale structure. To extract the relevant information about the matter field from these surveys, we usually resort to summary statistic, such as power spectrum. This is challenging in the non-linear regime due to the contamination of the signal by the baryonic physics and hence requires higher order statistics. The other approach is to carry out the analysis at the field level, e.g. inference or direct mapping between fields. Previous studies demonstrated the feasibility of building a mapping at the field level between baryons and dark matter (Wadekar et al., 2021; Villanueva-Domingo & Villaescusa-Navarro, 2021; Bernardini et al., 2022). In light of those existing works, and by utilizing generative adversarial networks (CycleGAN), we aim at building a bijective map between dark matter and two observables, namely neutral hydrogen and magnetic field magnitude. This is crucial since with a single training, it is possible to either paint the dark matter from simulation with baryons or directly infer its distribution from maps of observables obtained from different surveys in the near future.

2 METHODS

2.1 DATA

In this study, we use the publicly available CAMELS Multifields Dataset (CMD) (Villaescusa-Navarro et al., 2022) which contains thousands of 2D field maps generated from state-of-the-art hydrodynamics simulations (Villaescusa-Navarro et al., 2021). We consider 256×256 pixels 2D

¹South African Radio Astronomy Observatory, Cape Town 7925, andrianomena@gmail.com.

²Department of Physics & Astronomy University of the Western Cape, Cape Town 7535.

³Center for Cosmology and Particle Physics, Department of physics, New York University, 726 Broadway, New York, NY, 10003, shassan@flatironinstitute.org.

⁴Center for Computational Astrophysics, Flatiron Institute, New York, NY 10010, villaescusa.francisco@gmail.com.

maps, which cover an area of $25 \times 25 (h^{-1}\text{Mpc})^2$ at $z = 0$, of dark matter (Mcdm), neutral hydrogen (HI) and magnetic fields magnitude (B) from the IllustrisTNG LH set. The choice of the fields in this simple scenario is based on the aim of inferring the matter distribution from observables. In total, each field corresponds to 15000 2D images, each characterized by a set of 6 parameters: matter density (Ω_8), the amplitude of matter power spectrum (σ_8), the stellar feedbacks ($A_{\text{SN}1}$, $A_{\text{SN}2}$) and AGN feedbacks ($A_{\text{AGN}1}$, $A_{\text{AGN}2}$).

2.2 MODEL AND TRAINING

To build an invertible mapping between two different fields, we make use of CycleGAN (Zhu et al., 2017), an improvement on “pix2pix” method (Isola et al., 2017) which is trained on paired examples to achieve image-to-image translation. The approach consists of building a function $\mathcal{G} : X \rightarrow Y$ that maps a source field X to a target field Y (forward mapping), simultaneously with another function $\mathcal{F} : Y \rightarrow X$ that translates the target to the source field. To this end, two generators G_X (representing \mathcal{F} and producing the source field images) and G_Y (representing \mathcal{G} and producing the target field images) are trained with two adversarial discriminators D_X and D_Y respectively. Following the prescription in Zhu et al. (2017), there are two main components to the loss function for the training. The adversarial loss – for each of the pair (G_X, D_X) and (G_Y, D_Y) – is given by (Zhu et al., 2017)

$$\mathcal{L}_{\text{GAN}}(G_X, D_X, Y, X) = \mathbb{E}_{x \sim p_{\text{data}}(x)}[\log D_X(x)] + \mathbb{E}_{y \sim p_{\text{data}}(y)}[\log(1 - D_X(G_X(y)))] \quad (1)$$

and

$$\mathcal{L}_{\text{GAN}}(G_Y, D_Y, X, Y) = \mathbb{E}_{y \sim p_{\text{data}}(y)}[\log D_Y(y)] + \mathbb{E}_{x \sim p_{\text{data}}(x)}[\log(1 - D_Y(G_Y(x)))] \quad (2)$$

where $x \sim p_{\text{data}}(x)$ and $y \sim p_{\text{data}}(y)$ are the data distributions of the source images ($x \in X$) and target images ($y \in Y$) respectively. The consistency loss ensures that the functions \mathcal{G} and \mathcal{F} are inverse of each other such that $G_X(G_Y(x)) \approx x$ and $G_Y(G_X(y)) \approx y$. In other words, an input x (or y) is recovered by applying G_X on $G_Y(x)$ (or applying G_Y on $G_X(y)$). The consistency loss is

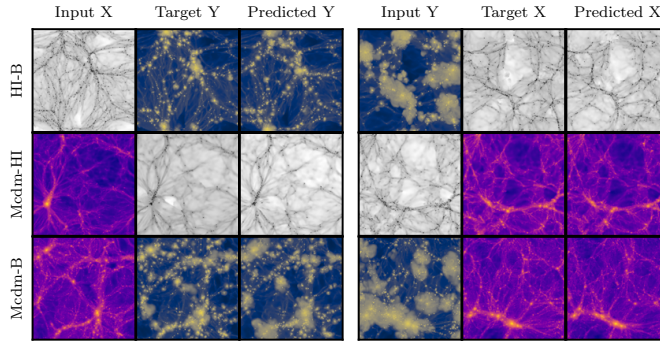


Figure 1: The top, middle, and bottom rows show the results from mapping HI-B, Mcdm-HI, and Mcdm-B respectively. It is worth noting that in a X-Y setting, X and Y designate the source and target fields, respectively.

given by (Zhu et al., 2017)

$$\mathcal{L}_{\text{cycle}} = \mathbb{E}_{x \sim p_{\text{data}}(x)}[\|G_X(G_Y(x)) - x\|_1] + \mathbb{E}_{y \sim p_{\text{data}}(y)}[\|G_Y(G_X(y)) - y\|_1], \quad (3)$$

where $\|\cdot\|_1$ denotes the mean absolute error (L1Loss). To further enforce a unique prediction of a given input such that $G_Y(y) \approx y$ and $G_X(x) \approx x$, an identity loss

$$\mathcal{L}_{\text{id}} = \mathbb{E}_{y \sim p_{\text{data}}(y)}[\|G_Y(y) - y\|_1] + \mathbb{E}_{x \sim p_{\text{data}}(x)}[\|G_X(x) - x\|_1] \quad (4)$$

is used. The total loss is then given by

$$\mathcal{L}_{\text{tot}} = \mathcal{L}_{\text{GAN}}(G_X, D_X, Y, X) + \mathcal{L}_{\text{GAN}}(G_Y, D_Y, X, Y) + \lambda_{\text{cycle}}\mathcal{L}_{\text{cycle}} + \lambda_{\text{id}}\mathcal{L}_{\text{id}}, \quad (5)$$

where both λ_{cycle} and λ_{id} are constants that characterize the contributions of the consistency cycle loss and identity loss respectively. Based on the prescription in Zhu et al. (2017), we have that $\lambda_{\text{cycle}} = 10$ and $\lambda_{\text{id}} = 5$ during training.

The generators G_Y and G_X , which have the same architecture, comprise a stage that downsamples the inputs (similar to encoding), 9 residual layers mimicking a bottleneck in the variational encoder and finally, a stage that upsamples the output from the bottleneck (similar to decoding). The discriminators D_Y and D_X , which are also identical, consist of chaining up 5 convolutional layers where the first three downsample the input by using stride = 2. To have a bit more control on the topology of the generated maps, we condition the input of each component of the model (G_X , D_X , G_Y and D_Y) on the underlying cosmology and astrophysics, i.e. the parameters Ω_8 , σ_8 , A_{SN1} , A_{SN1} , A_{AGN1} and A_{AGN2} . The array of parameters of shape 1×6 is passed through a dense layer with 4096 units whose output is reshaped to 64×64 , upsampled to 256×256 via interpolation, and finally concatenated along the channel to the input image. The model is trained for 100 epochs with 12000 unpaired examples, i.e. shuffling the set of source images such that they don't match the target images, using Adam optimizer with a learning rate of 0.0002. By setting the batch size to 1, following Zhu et al. (2017), each epoch takes about 109 minutes on a NVIDIA GeForce GTX 1080 Ti.

3 RESULTS

We present in Figure 1 some predictions by the generators (G_Y and G_X) using input images from test set. Each column in Figure 1 corresponds to the input, target and prediction in a given setup, e.g.

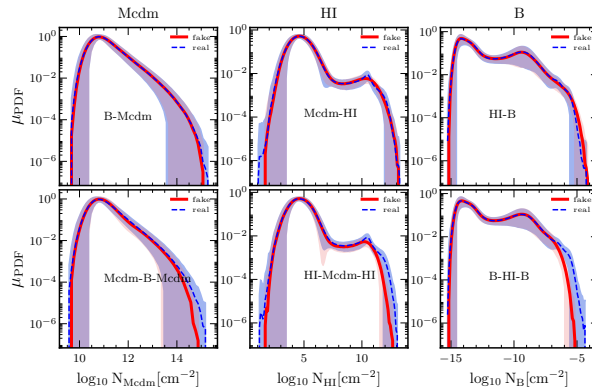


Figure 2: Comparing probability distribution function (PDF) of pixel values of the simulated and generated maps. Results corresponding to each field are shown in each column. Solid red and dashed blue lines denote PDFs of fake and real maps, respectively. Whereas red and blue shaded areas correspond to the standard deviations of the PDFs of fake and real maps, respectively. Each column corresponds to PDF of fields in different setups. The top rows show the comparison between the PDFs of predicted and true maps of a field in each setup. The bottom rows show the consistency test.

HI-B. In the first three columns of each row, we show the result from the forward mapping, i.e. the input is the source field X and the output is the target field Y. The last three columns of each row show the results related to the backward mapping, i.e. the input is the target field Y which is translated to the source field X. Visually, the output of the map by the generators are in good agreement with the inputs and the quality is comparable to that of the data from IllustrisTNG, overall. However it appears that predicting the magnetic field B from dark matter or neutral hydrogen seems to be a bit more challenging, as evidenced by the more noticeable difference in the map features between the ground truth and prediction (see 2nd and 3rd columns of both top and bottom rows).

The first metric we use to assess how well the model performs is the probability distribution function (PDF) of pixel intensities. The test set for each setup comprises of 1000 images unseen by the model during training. We then compute the mean μ_{PDF} and standard deviation σ_{PDF} of PDF of each field in the test set in each setup. In Figure 2, we present the PDFs of each field in some setups in each column. The top rows show the results from either forward or backward mappings, whereas the results of the consistency test – assessing if an input is recovered by applying forward and

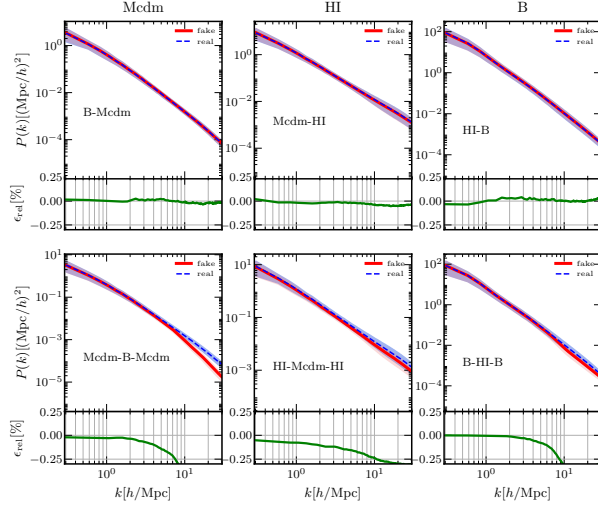


Figure 3: Comparing the resulting power spectra of the simulated and generated maps. Results corresponding to each field are shown in each column. Solid red and dashed blue lines denote the averaged power spectrum ($P(k)$) of fake and real maps, respectively. Whereas red and blue shaded areas correspond to the standard deviations of the $P(k)$ of fake and real maps, respectively. The relative error between the two power spectra for each field in each setup is indicated by the solid green line.

backward mappings sequentially – is presented at the bottom row. Results suggest that the model is capable of predicting the source and target fields in each setup using G_X and G_Y respectively, as evidenced by both μ_{PDF} and σ_{PDF} of the data in good agreement with those of the generated maps of each field. We also find that in general, the bijective mapping is achieved to a good accuracy, i.e. $G_X(G_Y(x)) \approx x$ and $G_Y(G_X(y)) \approx y$, as demonstrated by the μ_{PDF} and σ_{PDF} of the data which agree well with those of the recovered outputs. The relatively small discrepancy at the high end of the distributions can be accounted for by the small number of pixels having those values in the training data, i.e. overdense regions are rare. The other metric used in our investigation is the auto-power spectrum $P(k)$. Similar to the results presented in Figure 2, we show the $P(k)$ of the images produced by the forward/backward mapping and the consistency test at the top and bottom rows of Figure 3 respectively. Each panel shows the mean and standard deviation of the $P(k)$ of both real (dashed blue) and fake maps (solid red), and the solid green line at the bottom of each panel shows the relative difference $\left(\frac{P_{fake}(k)}{P_{real}(k)} - 1\right)$ between the two $P(k)$'s (fake and real maps). It is clear that the forward and backward mappings are able to produce maps with clustering properties in good agreement with those of the data. Moreover, the mean and variance of $P(k)$ of the recovered maps from the consistency check agree reasonably well with those of the maps from IllustrisTNG, regardless of the increasing difference on larger scales $k > 10 h/\text{Mpc}$. Our result for Mcdm-HI (see Figure 3 top middle panel) is comparable to what was obtained by Wadekar et al. (2021) where a standard UNet architecture was used to convert dark matter density to HI maps.

4 CONCLUSION

We have made use of CycleGAN model to build a bijective model that can map different physical fields from 2D maps created from the CAMELS state-of-the-art hydrodynamic simulations. Results show that the predicted maps exhibit statistical properties that agree well with those from the dataset used for training. Moreover, the condition for bijective mapping is met, as demonstrated by the consistency test. By applying the composition function $\mathcal{F} \circ \mathcal{G}(x)$ (or $\mathcal{G} \circ \mathcal{F}(y)$), an input map x (or y) is recovered reasonably while the statistics being preserved. This work represents a step forward towards establishing an efficient *direct* mapping between different observables, and hence maximizing the scientific return of future multi-wavelength surveys.

ACKNOWLEDGMENTS

SA acknowledges financial support from the South African Radio Astronomy Observatory (SARAO). FVN and SH acknowledge support provided by the Simons Foundation. SH also acknowledges support for Program number HST-HF2-51507 provided by NASA through a grant from the Space Telescope Science Institute, which is operated by the Association of Universities for Research in Astronomy, incorporated, under NASA contract NAS5-26555. The CAMELS project is supported by NSF grant AST 2108078.

REFERENCES

- Mauro Bernardini, Robert Feldmann, Daniel Anglés-Alcázar, Mike Boylan-Kolchin, James Bullock, Lucio Mayer, and Joachim Stadel. From ember to fire: predicting high resolution baryon fields from dark matter simulations with deep learning. *Monthly Notices of the Royal Astronomical Society*, 509(1):1323–1341, 2022.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1125–1134, 2017.
- Mario G Santos, Philip Bull, David Alonso, Stefano Camera, Pedro G Ferreira, Gianni Bernardi, Roy Maartens, Matteo Viel, Francisco Villaescusa-Navarro, Filipe B Abdalla, et al. Cosmology with a ska hi intensity mapping survey. *arXiv:1501.03989*, 2015.
- Francisco Villaescusa-Navarro, Daniel Anglés-Alcázar, Shy Genel, David N Spergel, Rachel S Somerville, Romeel Dave, Annalisa Pillepich, Lars Hernquist, Dylan Nelson, Paul Torrey, et al. The camels project: Cosmology and astrophysics with machine-learning simulations. *The Astrophysical Journal*, 915(1):71, 2021.
- Francisco Villaescusa-Navarro, Shy Genel, Daniel Angles-Alcazar, Leander Thiele, Romeel Dave, Desika Narayanan, Andrina Nicola, Yin Li, Pablo Villanueva-Domingo, Benjamin Wandelt, et al. The camels multifield data set: Learning the universe’s fundamental parameters with artificial intelligence. *The Astrophysical Journal Supplement Series*, 259(2):61, 2022.
- Pablo Villanueva-Domingo and Francisco Villaescusa-Navarro. Removing astrophysics in 21 cm maps with neural networks. *The Astrophysical Journal*, 907(1):44, 2021.
- Digvijay Wadekar, Francisco Villaescusa-Navarro, Shirley Ho, and Laurence Perreault-Levasseur. Hinet: Generating neutral hydrogen from dark matter with neural networks. *The Astrophysical Journal*, 916(1):42, 2021.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pp. 2223–2232, 2017.