

DENOISING DIFFUSION PROBABILISTIC MODELS TO PREDICT THE NUMBER DENSITY OF MOLECULAR CLOUDS IN ASTRONOMY

Duo Xu¹, Jonathan C. Tan^{*12}, Chia-Jung Hsu¹, Ye Zhu³⁴

¹Department of Astronomy
University of Virginia
Charlottesville, VA 22904-4235, USA
xuduol17@virginia.edu

²Department of Space, Earth & Environment
Chalmers University of Technology
SE-412 96 Gothenburg, Sweden

³Department of Computer Science
Illinois Institute of Technology
Chicago, IL 60616, USA

⁴Department of Computer Science
Princeton University
Princeton, NJ 08544, USA

ABSTRACT

Denoising Diffusion Probabilistic Models (DDPMs) have become the mainstream generative approach in the Machine Learning and Computer Vision area, achieving the state-of-the-art performance in synthesizing high-quality images, videos, and audio. In this work, we bring the DDPMs out of the data generation tasks, but to a new scientific application field in astronomy for inferring the volume or number density of giant molecular clouds (GMCs) from projected mass surface density maps. Specifically, we adopt magnetohydrodynamic (MHD) simulations with different global magnetic field strengths and large-scale dynamics, *i.e.*, noncolliding and colliding GMCs. We train a DDPM on both mass surface density maps and their corresponding mass-weighted number density maps from different viewing angles for all the simulations. We compare our performance with a more traditional empirical two-component and three-component power-law fitting method and with a more traditional neural network machine learning approach (CAS1-2D). Experiments show that DDPMs achieve an order of magnitude improvement in the accuracy of predicting number density compared to that by other methods, demonstrating the promising potential of applying DDPMs in astrophysics.

1 INTRODUCTION

Giant molecular clouds (GMCs) are the birthplace of stars (Shu et al., 1987). Investigating the physical and chemical conditions of GMCs is a crucial step to understand the evolution of the interstellar medium (ISM) and the formation of stars (McKee & Ostriker, 2007; Heyer & Dame, 2015). Among all the physical quantities of GMCs, the number density or volume density is one of the most fundamental properties that relates to various physical quantity estimations, such as the free-fall time, the magnetic field strength (Chandrasekhar & Fermi, 1953) and chemical reaction rates (Tielens & Hagen, 1982). However, it is difficult to quantify the number density of GMCs from observations. The traditional approach of estimating the number density of GMCs is based on observations of column density and certain assumptions on the geometry of the clouds, for example, a cylindrical geometry for filamentary structures or spherical geometry for dense cores (André et al., 2014). Bisbas et al. (2021; 2023) proposed an empirical power-law to convert the observed column density to the mean number density of GMCs based on the MHD simulations from Wu et al. (2017), which works decently but with noticeable scatter. Another method to constrain the number density of GMCs is utilizing density “probes”, such as cyanoacetylene (HC₃N, Avery et al., 1982; Schloerb et al., 1983; Li & Goldsmith, 2012). The relative intensity of different transitions of HC₃N is sensitive to the number density of the cloud, which makes it possible to constrain the mean number density directly by observing multiple transitions of HC₃N. Li & Goldsmith (2012) successfully observed $J = 2 - 1$

and 10 – 9 transitions of HC_3N in Taurus B213 filament and constrained the molecular gas number density to be around $1.8 \times 10^4 \text{ cm}^{-3}$. Unfortunately, the line ratio of HC_3N can only probe the number density at a relatively narrow range, between 10^4 and 10^6 cm^{-3} (Li & Goldsmith, 2012), which limits its ability to infer the number density of the full range of structures that exist in GMCs. Consequently, a novel method to infer the number density of GMCs under a variety of physical conditions with high precision is in great demand. Machine learning makes it possible to learn from both the morphology of the cloud and their column density to infer the mean number density rather than using a simple average power-law conversion.

Machine learning has gained great popularity among astronomers. For example, Convolutional neural networks have been successfully applied to identify structures (Xu et al., 2020a;b; 2022a) and infer physical quantities, such as protostellar outflow inclination angles and magnetic field directions (Xu et al., 2022b). More recently, Denoising Diffusion Probabilistic Models (DDPMs) Sohl-Dickstein et al. (2015); Ho et al. (2020), originated from the natural thermodynamics problem, simulate a random walk process in the data space, which shares the intrinsic alignment and consistency with most scientific problems. For example, the observed emission or absorption maps of GMCs that derive the column density map are part of the results of the random walk of photons in the ISM, *i.e.*, the scattering process. Inferring the raw mass-weighted number density distribution based on the observed column density inherits the basic concepts of diffusion models. In addition, DDPMs have demonstrated their proficiency and robustness in image and audio generation (Sohl-Dickstein et al., 2015; Ho et al., 2020; Zhu et al., 2023b; Rombach et al., 2022; Zhu et al., 2023a), which are suitable for the prediction task in astronomy. Therefore, our primary objective in this work is to demonstrate the great potential for applying DDPMs in astronomy studies.

2 METHOD

The core formulation of diffusion models is inspired by non-equilibrium thermodynamics (Sohl-Dickstein et al., 2015), which models a stochastic Markov chain of T steps in two directions. The forward direction q , also known as the diffusion process, gradually adds stochastic Gaussian noises to a data sample x_0 following:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}), \quad (1)$$

where $\{\beta_t\}_t^T = 1$ are pre-scheduled variances. The other direction, often referred as reverse direction or generative process p , denoises a noisy sample x_T from a standard Gaussian distribution to a data sample x_0 as:

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{1 - \beta_t}}(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}}\epsilon_t^\theta(\mathbf{x}_t)) + \sigma_t z_t, \quad (2)$$

where ϵ_t^θ is the learnable noise predictor, $z_t \sim \mathcal{N}(0, \mathbf{I})$, and the variance of the reverse process σ_t^2 is set to be $\sigma^2 = \beta$. The actual neural network training process aims to learn the noise predictor to optimize the variational lower bound on negative log likelihood:

$$\mathcal{L} = \mathbb{E}_q[-\log p(\mathbf{x}_T) - \sum_{t \geq 1} \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1})}]. \quad (3)$$

To adapt the diffusion models in our context, we deploy a diffusion model with a similar parameter setup as in Ho et al. (2020). Specifically, the diffusion model has in total $T = 1000$ steps, and is optimized using the variational loss in Equation 3, which makes valid the assumption that the reverse direction converges to the Gaussian stochastic diffusion process. In particular, we provide \mathbf{x}_c as an additional input condition to make the prediction follow each individual observational sample. We train the diffusion model for 400 epochs and evaluate the performance on a sample in the test set.

3 EXPERIMENTS

Data. We carry out ideal MHD simulations based on the set-up of Wu et al. (2020), which were conducted using the adaptive mesh refinement (AMR) code Enzo (Bryan et al., 2014) with a top grid resolved by 256^3 and 4 additional levels of refinement, achieving a resolution of 0.03125 pc^1 within a 128 pc^3 domain. The simulations include heating/cooling based on a photo-dissociation model

¹1 pc = $3.086 \times 10^{16} \text{ m}$

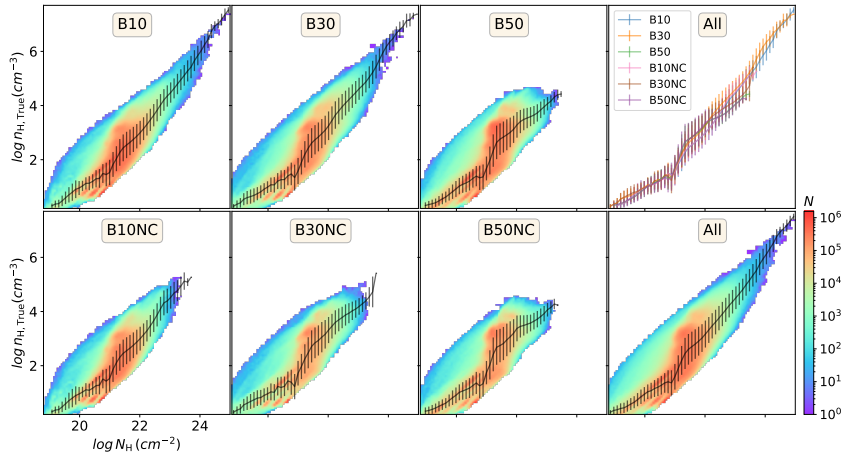


Figure 1: Relation between the LOS mass-weighted number density and the column density for different simulations. The rainbow color background indicates the 2D histogram of the distribution between number density and column density for each simulation. The lines with errorbars represent the mean and their standard deviation of each column density bin.

Table 1: Power-Law Fitting Results. Label, number of power-law components, break points of power-law fittings in log scale, and power indices of each power-law component.

Label	Components	Break Points	Power Indices
2PL	2	20.82	0.61/1.47
3PL	3	20.73/22.69	0.61/1.41/1.49
2PL-i	2	20.19	0.67/1.42
3PL-i	3	21.35/21.79	0.79/4.06/1.26

with FUV radiation field of $G_0 = 4$ Habings and cosmic ray ionization rate $\zeta = 10^{-16} \text{ s}^{-1}$, self-gravity, magnetic fields and supersonic turbulence. The simulations do not include star formation or feedback, so represent the structures that develop in the early phases of collapse up to the onset of star formation. We select two different types of setup of the large-scale dynamics of the GMCs, noncolliding and colliding GMCs. For each setup, we implement three different magnetic field strengths, 10, 30, and 50 μG . In the colliding cases, the clouds have relative velocities of 10 km s^{-1} , and are offset by $0.5R_{\text{GMC}}$. The GMCs have an initial temperature of 15 K, but soon establish a multiphase temperature structure. The simulations are run for 5 Myr. We take two evolutionary stages, 3 and 4 Myr, for each setup run.

To enhance the diversity of the data set, we generate column density maps and their corresponding line-of-sight (LOS) mass-weighted number density across different scales by adopting different AMR levels with different physical resolutions. The image size in pixels is 128×128 , with multiple physical scales, including 32, 16, 8, 4 and 2 pc. In total, we have 7179 images in the data set, in which 70% are used for the training set, and the remaining 30% are a test set. Figure 1 shows the correlation between the LOS mass-weighted number density and the column density for different simulations. Although the column density range is not the same for different simulations, it is obvious that the relation between the mass-weighted number density and column density is similar.

Results. We adopt several different approaches to convert column density to mass-weighted number density on the LOS. We start with power-law fitting on the relation between the LOS mass-weighted number density of H nuclei (n_{H}) and the column density of H nuclei (N_{H}) for all simulations in Figure 1. We adopt two-components and three components power law to fit the $n_{\text{H}} - N_{\text{H}}$ relation, i.e., $n_{\text{H}} = f(N_{\text{H}})$. Meanwhile, we conduct the “inverted” fitting, i.e., adopting two-component and three component power laws to fit the $N_{\text{H}} - n_{\text{H}}$ relation $N_{\text{H}} = f(n_{\text{H}})$ and then derive the inverse function $n_{\text{H}} = f^{-1}(N_{\text{H}})$. We summarize the fitting results in Table 1. We then follow the power-law fitting results to convert the column density to mass-weighted number density. We show the fitting results in Figure 2. In addition, we present the result from machine learning approaches, including CASI-2D (Van Oort et al., 2019) and diffusion model, in Figure 2. It is obvious that there is significant dispersion between the true mass-weighted number density $n_{\text{H, True}}$ and the predicted number density $n_{\text{H, Pred}}$ that is converted by power-law and that predicted by CASI-2D. The pre-

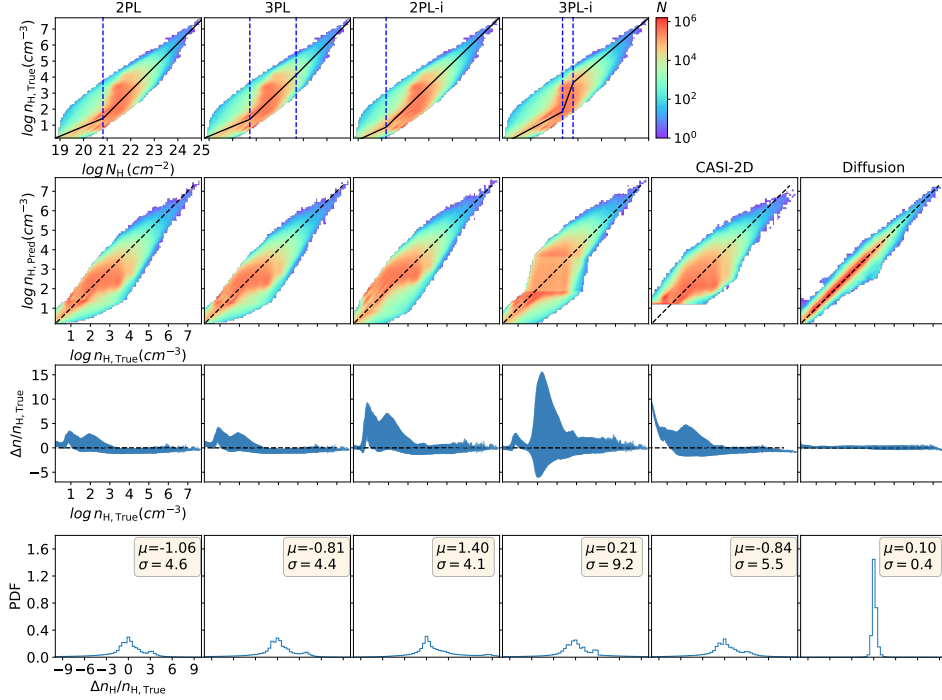


Figure 2: Summary of the performance of different approaches to convert the gas column density to number density on all the data samples.

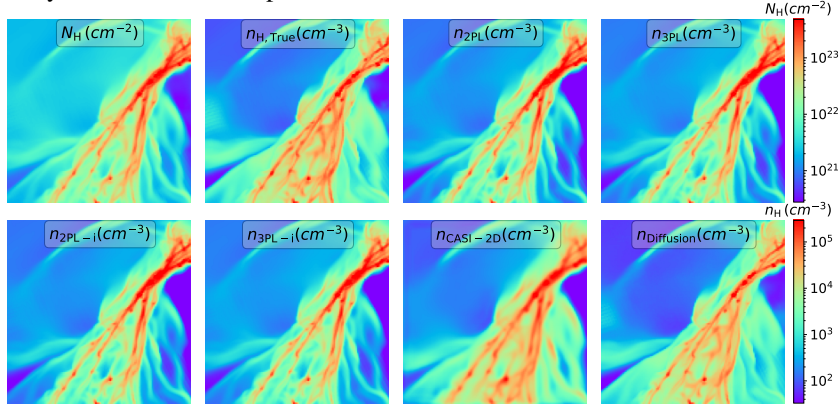


Figure 3: Comparison between different approaches to converting gas column density to number density on a sample in the test set.

dicted number density by the diffusion model has a much smaller dispersion around the true density. We list the dispersion between $n_{H, \text{True}}$ and $n_{H, \text{Pred}}$ in Figure 2. The predicted number density by the diffusion model is an order of magnitude better than that by all the other approaches. We show a sample image from the test set and apply different approaches to obtain the predicted number density in Figure 3. Although the column density map is similar to the true number density map, there is noticeable difference at relatively high column densities, where more structures are highlighted in the true number density map but appear faint in the column density map. The number density map predicted by the diffusion model is able to recreate the actual structures across a wide density range.

4 CONCLUSION

In this work, we deploy the DDPMs to predict the LOS mass-weighted number density of GMCs from column density maps in astronomy. We achieve large improvements in the performance for both synthetic test samples and real observational data, demonstrating the effectiveness and great potential for astronomical applications using diffusion models.

REFERENCES

- P. André, J. Di Francesco, D. Ward-Thompson, S. I. Inutsuka, R. E. Pudritz, and J. E. Pineda. From Filamentary Networks to Dense Cores in Molecular Clouds: Toward a New Paradigm for Star Formation. In Henrik Beuther, Ralf S. Klessen, Cornelis P. Dullemond, and Thomas Henning (eds.), *Protostars and Planets VI*, pp. 27–51, January 2014. doi: 10.2458/azu_uapress_9780816531240-ch002.
- L. W. Avery, J. M. MacLeod, and N. W. Broten. A model of Taurus Molecular Cloud 1 based on HC3N observations. *The Astrophysical Journal*, 254:116–125, March 1982. doi: 10.1086/159713.
- Oliver J. Bartlett, David M. Benoit, Kevin A. Pimblet, Brooke Simmons, and Laura Hunt. Noise reduction on single-shot images using an autoencoder. *Monthly Notices of the Royal Astronomical Society*, March 2023. doi: 10.1093/mnras/stad665.
- Rainer Beck. Magnetic fields in spiral galaxies. , 24:4, December 2015. doi: 10.1007/s00159-015-0084-4.
- Burger Becker, Mattia Vaccari, Matthew Prescott, and Trienko Grobler. CNN architecture comparison for radio galaxy classification. *Monthly Notices of the Royal Astronomical Society*, 503(2): 1828–1846, May 2021. doi: 10.1093/mnras/stab325.
- Thomas G. Bisbas, Jonathan C. Tan, and Kei E. I. Tanaka. Photodissociation region diagnostics across galactic environments. *Monthly Notices of the Royal Astronomical Society*, 502(2):2701–2732, April 2021. doi: 10.1093/mnras/stab121.
- Thomas G. Bisbas, Ewine F. van Dishoeck, Chia-Yu Hu, and Andreas Schrupa. PDFCHEM: A new fast method to determine ISM properties and infer environmental parameters using probability distributions. *Monthly Notices of the Royal Astronomical Society*, 519(1):729–753, February 2023. doi: 10.1093/mnras/stac3487.
- L. Bisigello, C. J. Conselice, M. Baes, M. Bolzonella, M. Brescia, S. Cavuoti, O. Cucciati, A. Humphrey, L. K. Hunt, C. Maraston, L. Pozzetti, C. Tortora, S. E. van Mierlo, N. Aghanim, N. Auricchio, M. Baldi, R. Bender, C. Bodendorf, D. Bonino, E. Branchini, J. Brinchmann, S. Camera, V. Capobianco, C. Carbone, J. Carretero, F. J. Castander, M. Castellano, A. Cimatti, G. Congedo, L. Conversi, Y. Copin, L. Corcione, F. Courbin, M. Cropper, A. Da Silva, H. Degaudenzi, M. Douspis, F. Dubath, C. A. J. Duncan, X. Dupac, S. Dusini, S. Farrens, S. Ferriol, M. Frailis, E. Franceschi, P. Franzetti, M. Fumana, B. Garilli, W. Gillard, B. Gillis, C. Giocoli, A. Grazian, F. Grupp, L. Guzzo, S. V. H. Haugan, W. Holmes, F. Hormuth, A. Hornstrup, K. Jahnke, M. Kümmel, S. Kermiche, A. Kiessling, M. Kilbinger, R. Kohley, M. Kunz, H. Kurki-Suonio, S. Ligorì, P. B. Lilje, I. Lloro, E. Maiorano, O. Mansutti, O. Marggraf, K. Markovic, F. Marulli, R. Massey, S. Maurogordato, E. Medinaceli, M. Meneghetti, E. Merlin, G. Meylan, M. Moresco, L. Moscardini, E. Munari, S. M. Niemi, C. Padilla, S. Paltani, F. Pasian, K. Pedersen, V. Pettorino, G. Polenta, M. Poncet, L. Popa, F. Raison, A. Renzi, J. Rhodes, G. Riccio, H. W. Rix, E. Romelli, M. Roncarelli, C. Rosset, E. Rossetti, R. Saglia, D. Sapone, B. Sartoris, P. Schneider, M. Scodeggio, A. Secroun, G. Seidel, C. Sirignano, G. Sirri, L. Stanco, P. Tallada-Crespí, D. Tavagnacco, A. N. Taylor, I. Tereno, R. Toledo-Moreo, F. Torradeflot, I. Tutusaus, E. A. Valentijn, L. Valenziano, T. Vassallo, Y. Wang, A. Zacchei, G. Zamorani, J. Zoubian, S. Andreon, S. Bardelli, A. Boucaud, C. Colodro-Conde, D. Di Ferdinando, J. Graciá-Carpio, V. Lindholm, D. Maino, S. Mei, V. Scottez, F. Sureau, M. Tenti, E. Zucca, A. S. Borlaff, M. Ballardini, A. Biviano, E. Bozzo, C. Burigana, R. Cabanac, A. Cappi, C. S. Carvalho, S. Casas, G. Castignani, A. Cooray, J. Coupon, H. M. Courtois, J. Cuby, S. Davini, G. De Lucia, G. Desprez, H. Dole, J. A. Escartin, S. Escoffier, M. Farina, S. Fotopoulou, K. Ganga, J. Garcia-Bellido, K. George, F. Giacomini, G. Gozaliasl, H. Hildebrandt, I. Hook, M. Huertas-Company, V. Kansal, E. Keihanen, C. C. Kirkpatrick, A. Loureiro, J. F. Macías-Pérez, M. Magliocchetti, G. Mainetti, S. Marcin, M. Martinelli, N. Martinet, R. B. Metcalf, P. Monaco, G. Morgante, S. Nadathur, A. A. Nucita, L. Patrizzii, A. Peel, D. Potter, A. Pourtsidou, M. Pöntinen, P. Reimberg, A. G. Sánchez, Z. Sakr, M. Schirmer, E. Sefusatti, M. Sereno, J. Stadel, R. Teyssier, C. Valieri, J. Valiviita, M. Viel, and Euclid Collaboration. Euclid preparation: XXIII. Derivation of galaxy physical properties with deep machine learning using mock fluxes and H-band images. *Monthly Notices of the Royal Astronomical Society*, January 2023. doi: 10.1093/mnras/stac3810.

- Greg L. Bryan, Michael L. Norman, Brian W. O’Shea, Tom Abel, John H. Wise, Matthew J. Turk, Daniel R. Reynolds, David C. Collins, Peng Wang, Samuel W. Skillman, Britton Smith, Robert P. Harkness, James Bordner, Ji-hoon Kim, Michael Kuhlen, Hao Xu, Nathan Goldbaum, Cameron Hummels, Alexei G. Kritsuk, Elizabeth Tasker, Stephen Skory, Christine M. Simpson, Oliver Hahn, Jeffrey S. Oishi, Geoffrey C. So, Fen Zhao, Renyue Cen, Yuan Li, and Enzo Collaboration. ENZO: An Adaptive Mesh Refinement Code for Astrophysics. *The Astrophysical Journal Supplement Series*, 211(2):19, April 2014. doi: 10.1088/0067-0049/211/2/19.
- S. Chandrasekhar and E. Fermi. Magnetic Fields in Spiral Arms. *The Astrophysical Journal*, 118: 113, July 1953. doi: 10.1086/145731.
- Leverett Davis. The strength of interstellar magnetic fields. *Phys. Rev.*, 81:890–891, Mar 1951. doi: 10.1103/PhysRev.81.890.2. URL <https://link.aps.org/doi/10.1103/PhysRev.81.890.2>.
- H. Domínguez Sánchez, M. Huertas-Company, M. Bernardi, D. Tuccillo, and J. L. Fischer. Improving galaxy morphologies for SDSS with Deep Learning. *Monthly Notices of the Royal Astronomical Society*, 476(3):3661–3676, February 2018. doi: 10.1093/mnras/sty338.
- C. Gheller and F. Vazza. Convolutional deep denoising autoencoders for radio astronomical images. *Monthly Notices of the Royal Astronomical Society*, 509(1):990–1009, January 2022. doi: 10.1093/mnras/stab3044.
- Munan Gong, Eve C. Ostriker, and Mark G. Wolfire. A Simple and Accurate Network for Hydrogen and Carbon Chemistry in the Interstellar Medium. *The Astrophysical Journal*, 843(1):38, July 2017. doi: 10.3847/1538-4357/aa7561.
- Shoubaneh Hemmati, Eric Huff, Hooshang Nayyeri, Agnès Ferté, Peter Melchior, Bahram Mobasher, Jason Rhodes, Abtin Shahidi, and Harry Teplitz. Deblending Galaxies with Generative Adversarial Networks. *The Astrophysical Journal*, 941(2):141, December 2022. doi: 10.3847/1538-4357/aca1b8.
- Eric Herbst and Ewine F. van Dishoeck. Complex Organic Interstellar Molecules. *Annual Review of Astronomy and Astrophysics*, 47(1):427–480, September 2009. doi: 10.1146/annurev-astro-082708-101654.
- Mark Heyer and T. M. Dame. Molecular Clouds in the Milky Way. *Annual Review of Astronomy and Astrophysics*, 53:583–629, August 2015. doi: 10.1146/annurev-astro-082214-122324.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 6840–6851. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf>.
- Jes K. Jørgensen, Arnaud Belloche, and Robin T. Garrod. Astrochemistry During the Formation of Stars. *Annual Review of Astronomy and Astrophysics*, 58:727–778, August 2020. doi: 10.1146/annurev-astro-032620-021927.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.
- Mark R. Krumholz, Christopher F. McKee, and Joss Bland-Hawthorn. Star Clusters Across Cosmic Time. *Annual Review of Astronomy and Astrophysics*, 57:227–303, August 2019. doi: 10.1146/annurev-astro-091918-104430.
- Adam K. Leroy, Fabian Walter, Elias Brinks, Frank Bigiel, W. J. G. de Blok, Barry Madore, and M. D. Thornley. The Star Formation Efficiency in Nearby Galaxies: Measuring Where Gas Forms Stars Effectively. , 136(6):2782–2845, December 2008. doi: 10.1088/0004-6256/136/6/2782.
- Di Li and Paul F. Goldsmith. Is the Taurus B213 Region a True Filament?: Observations of Multiple Cyanoacetylene Transitions. *The Astrophysical Journal*, 756(1):12, September 2012. doi: 10.1088/0004-637X/756/1/12.

- Hao Liu, Yan Xu, Jiasheng Wang, Ju Jing, Chang Liu, Jason T. L. Wang, and Haimin Wang. Inferring Vector Magnetic Fields from Stokes Profiles of GST/NIRIS Using a Convolutional Neural Network. *The Astrophysical Journal*, 894(1):70, May 2020. doi: 10.3847/1538-4357/ab8818.
- Christopher F. McKee and Eve C. Ostriker. Theory of Star Formation. *Annual Review of Astronomy and Astrophysics*, 45(1):565–687, September 2007. doi: 10.1146/annurev.astro.45.051806.110602.
- Planck Collaboration, P. A. R. Ade, N. Aghanim, M. I. R. Alves, M. Arnaud, D. Arzoumanian, M. Ashdown, J. Aumont, C. Baccigalupi, A. J. Banday, R. B. Barreiro, N. Bartolo, E. Battaner, K. Benabed, A. Benoît, A. Benoit-Lévy, J. P. Bernard, M. Bersanelli, P. Bielewicz, J. J. Bock, L. Bonavera, J. R. Bond, J. Borrill, F. R. Bouchet, F. Boulanger, A. Bracco, C. Burigana, E. Calabrese, J. F. Cardoso, A. Catalano, H. C. Chiang, P. R. Christensen, L. P. L. Colombo, C. Combet, F. Couchot, B. P. Crill, A. Curto, F. Cuttaia, L. Danese, R. D. Davies, R. J. Davis, P. de Bernardis, A. de Rosa, G. de Zotti, J. Delabrouille, C. Dickinson, J. M. Diego, H. Dole, S. Donzelli, O. Doré, M. Douspis, A. Ducout, X. Dupac, G. Efstathiou, F. Elsner, T. A. Enßlin, H. K. Eriksen, D. Falceta-Goncalves, E. Falgarone, K. Ferrière, F. Finelli, O. Forni, M. Frailis, A. A. Fraisse, E. Franceschi, A. Frejsel, S. Galeotta, S. Galli, K. Ganga, T. Ghosh, M. Giard, E. Gjerløw, J. González-Nuevo, K. M. Górski, A. Gregorio, A. Gruppuso, J. E. Gudmundsson, V. Guillet, D. L. Harrison, G. Helou, P. Hennebelle, S. Henrot-Versillé, C. Hernández-Montenegro, D. Herranz, S. R. Hildebrandt, E. Hivon, W. A. Holmes, A. Hornstrup, K. M. Huffenberger, G. Hurier, A. H. Jaffe, T. R. Jaffe, W. C. Jones, M. Juvela, E. Keihänen, R. Kesitalo, T. S. Kisner, J. Knoche, M. Kunz, H. Kurki-Suonio, G. Lagache, J. M. Lamarre, A. Lasenby, M. Lattanzi, C. R. Lawrence, R. Leonardi, F. Levrier, M. Liguori, P. B. Lilje, M. Linden-Vørnle, M. López-Cañiego, P. M. Lubin, J. F. Macías-Pérez, D. Maino, N. Mandolesi, A. Mangilli, M. Maris, P. G. Martin, E. Martínez-González, S. Masi, S. Matarrese, A. Melchiorri, L. Mendes, A. Mennella, M. Migliaccio, M. A. Miville-Deschênes, A. Moneti, L. Montier, G. Morgante, D. Mortlock, D. Munshi, J. A. Murphy, P. Naselsky, F. Nati, C. B. Netterfield, F. Noviello, D. Novikov, I. Novikov, N. Oppermann, C. A. Oxborrow, L. Pagano, F. Pajot, R. Paladini, D. Paoletti, F. Pasian, L. Perotto, V. Pettorino, F. Piacentini, M. Piat, E. Pierpaoli, D. Pietrobon, S. Plaszczynski, E. Pointecouteau, G. Polenta, N. Ponthieu, G. W. Pratt, S. Prunet, J. L. Puget, J. P. Rachen, M. Reinecke, M. Remazeilles, C. Renault, A. Renzi, I. Ristorcelli, G. Rocha, M. Rossetti, G. Roudier, J. A. Rubiño-Martín, B. Rusholme, M. Sandri, D. Santos, M. Savelainen, G. Savini, D. Scott, J. D. Soler, V. Stolyarov, R. Sudiwala, D. Sutton, A. S. Suur-Uski, J. F. Sygnet, J. A. Tauber, L. Terenzi, L. Toffolatti, M. Tomasi, M. Tristram, M. Tucci, G. Umata, L. Valenziano, J. Valiviita, B. Van Tent, P. Vielva, F. Villa, L. A. Wade, B. D. Wandelt, I. K. Wehus, N. Ysard, D. Yvon, and A. Zonca. Planck intermediate results. XXXV. Probing the role of the magnetic field in the formation of structure in molecular clouds. *Astronomy and Astrophysics*, 586:A138, February 2016. doi: 10.1051/0004-6361/201525896.
- R. Rao, R. M. Crutcher, R. L. Plambeck, and M. C. H. Wright. High-Resolution Millimeter-Wave Mapping of Linearly Polarized Dust Emission: Magnetic Field Structure in Orion. , 502(1): L75–L78, July 1998. doi: 10.1086/311485.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- Amélie Saintonge and Barbara Catinella. The Cold Interstellar Medium of Galaxies in the Local Universe. *Annual Review of Astronomy and Astrophysics*, 60:319–361, August 2022. doi: 10.1146/annurev-astro-021022-043545.
- F. P. Schloerb, R. L. Snell, and J. S. Young. Structure of dense molecular gas in TMC 1 from observations of three transitions of HC3N. *The Astrophysical Journal*, 267:163–173, April 1983. doi: 10.1086/160854.
- Frank H. Shu, Fred C. Adams, and Susana Lizano. Star formation in molecular clouds: observation and theory. *Annual Review of Astronomy and Astrophysics*, 25:23–81, January 1987. doi: 10.1146/annurev.aa.25.090187.000323.
- Anton A. Smirnov, Sergey S. Savchenko, Denis M. Poliakov, Alexander A. Marchuk, Aleksandr V. Mosenkov, Vladimir B. Il'in, George A. Gontcharov, Javier Román, and Jonah Seguíne. Prospects

- for future studies using deep imaging: analysis of individual Galactic cirrus filaments. *Monthly Notices of the Royal Astronomical Society*, 519(3):4735–4752, March 2023. doi: 10.1093/mnras/stac3765.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2256–2265, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/sohl-dickstein15.html>.
- Lyman Spitzer. *Physical processes in the interstellar medium*. 1978. doi: 10.1002/9783527617722.
- Sam F. Sweere, Ivan Valtchanov, Maggie Lieu, Antonia Vojtekova, Eva Verdugo, Maria Santos-Lleo, Florian Pacaud, Alexia Briassouli, and Daniel Cámpora Pérez. Deep learning-based super-resolution and de-noising for XMM-newton images. *Monthly Notices of the Royal Astronomical Society*, 517(3):4054–4069, December 2022. doi: 10.1093/mnras/stac2437.
- Linda J. Tacconi, Reinhard Genzel, and Amiel Sternberg. The Evolution of the Star-Forming Interstellar Medium Across Cosmic Time. *Annual Review of Astronomy and Astrophysics*, 58: 157–203, August 2020. doi: 10.1146/annurev-astro-082812-141034.
- A. G. G. M. Tielens and W. Hagen. Model calculations of the molecular composition of interstellar grain mantles. *Astronomy and Astrophysics*, 114(2):245–260, October 1982.
- Colin M. Van Oort, Duo Xu, Stella S. R. Offner, and Robert A. Gutermuth. CASI: A Convolutional Neural Network Approach for Shell Identification. *The Astrophysical Journal*, 880(2):83, August 2019. doi: 10.3847/1538-4357/ab275e.
- V. Wakelam, I. W. M. Smith, E. Herbst, J. Troe, W. Geppert, H. Linnartz, K. Öberg, E. Roueff, M. Agúndez, P. Pernot, H. M. Cuppen, J. C. Loison, and D. Talbi. Reaction Networks for Interstellar Chemical Modelling: Improvements and Challenges. , 156(1-4):13–72, October 2010. doi: 10.1007/s11214-010-9712-5.
- Benjamin Wu, Jonathan C. Tan, Fumitaka Nakamura, Sven Van Loo, Duncan Christie, and David Collins. GMC Collisions as Triggers of Star Formation. II. 3D Turbulent, Magnetized Simulations. *The Astrophysical Journal*, 835(2):137, February 2017. doi: 10.3847/1538-4357/835/2/137.
- Benjamin Wu, Jonathan C. Tan, Duncan Christie, and Fumitaka Nakamura. GMC Collisions as Triggers of Star Formation. VII. The Effect of Magnetic Field Strength on Star Formation. *The Astrophysical Journal*, 891(2):168, March 2020. doi: 10.3847/1538-4357/ab77b5.
- Duo Xu, Stella S. R. Offner, Robert Gutermuth, and Colin Van Oort. Application of Convolutional Neural Networks to Identify Stellar Feedback Bubbles in CO Emission. *The Astrophysical Journal*, 890(1):64, February 2020a. doi: 10.3847/1538-4357/ab6607.
- Duo Xu, Stella S. R. Offner, Robert Gutermuth, and Colin Van Oort. Application of Convolutional Neural Networks to Identify Protostellar Outflows in CO Emission. *The Astrophysical Journal*, 905(2):172, December 2020b. doi: 10.3847/1538-4357/abc7bf.
- Duo Xu, Stella S. R. Offner, Robert Gutermuth, Shuo Kong, and Hector G. Arce. A Census of Protostellar Outflows in Nearby Molecular Clouds. *The Astrophysical Journal*, 926(1):19, February 2022a. doi: 10.3847/1538-4357/ac39a0.
- Duo Xu, Stella S. R. Offner, Robert Gutermuth, and Jonathan C. Tan. A Census of Outflow to Magnetic Field Orientations in Nearby Molecular Clouds. *The Astrophysical Journal*, 941(1):81, December 2022b. doi: 10.3847/1538-4357/aca153.
- Duo Xu, Chi-Yan Law, and Jonathan C. Tan. Application of Convolutional Neural Networks to Predict Magnetic Fields’ Directions in Turbulent Clouds. *The Astrophysical Journal*, 942(2):95, January 2023. doi: 10.3847/1538-4357/aca66c.

Shangjia Zhang, Zhaohuan Zhu, and Mingon Kang. PGNets: planet mass prediction using convolutional neural networks for radio continuum observations of protoplanetary discs. *Monthly Notices of the Royal Astronomical Society*, 510(3):4473–4484, March 2022. doi: 10.1093/mnras/stab3502.

Ye Zhu, Yu Wu, Zhiwei Deng, Olga Russakovsky, and Yan Yan. Boundary guided mixing trajectory for semantic control with diffusion models. *arXiv preprint arXiv:2302.08357*, 2023a.

Ye Zhu, Yu Wu, Kyle Olszewski, Jian Ren, Sergey Tulyakov, and Yan Yan. Discrete contrastive diffusion for cross-modal music and image generation. In *International Conference on Learning Representations*, 2023b.

A BACKGROUND IN ASTROPHYSICS

In this section, we extend the background of the astrophysical problem we target to solve in this work.

Giant molecular clouds (GMCs) are one of the most important components of the interstellar medium (ISM) in galaxies. The ISM is the material that fills the space between stars, consisting of gas and dust, as well as cosmic rays and magnetic fields (Spitzer, 1978). The ISM plays a critical role in the life cycle of galaxies, regulating the rate of star formation (e.g., Leroy et al., 2008; Tacconi et al., 2020). GMCs are particularly important because they contain the majority of the dense gas in the ISM, which is required for the formation of stars (e.g., Shu et al., 1987; Saintonge & Catinella, 2022). The physical conditions within GMCs are highly complex and dynamic, with variations in space and time of density, temperature, velocity, and magnetic field strength and direction. These conditions can lead to a wide range of phenomena, including the formation of protostars and star clusters (e.g., McKee & Ostriker, 2007; Heyer & Dame, 2015; Krumholz et al., 2019), and the formation of complex organic molecules (e.g., Herbst & van Dishoeck, 2009; Jørgensen et al., 2020). Consequently, investigating the physical and chemical conditions of GMCs is a crucial step towards understanding the complex physical processes that occur within the Milky Way, as well as the properties and evolution of galaxies throughout the Universe.

Among all the physical quantities of GMCs, the density (ρ), i.e., mass per unit volume,² is one of the most fundamental properties that relates to various physical quantity estimations, such as the free-fall time, the magnetic field strength (Davis, 1951; Chandrasekhar & Fermi, 1953; Beck, 2015) and chemical reaction rates (Tielens & Hagen, 1982; Wakelam et al., 2010; Gong et al., 2017). The free-fall time depends on density as $t_{\text{ff}} \propto \rho^{-1/2}$. Considering measurement of magnetic fields, one method commonly used is the Davis-Chandrasekhar-Fermi (DCF) method (Davis, 1951; Chandrasekhar & Fermi, 1953; Beck, 2015). This estimates the plane-of-sky (POS) component of the magnetic field using polarized thermal dust emission (Rao et al., 1998; Planck Collaboration et al., 2016). The DCF method is based on the assumption that the magnetic field in the ISM is in a state of equipartition with the turbulent kinetic energy of the gas. This means that the magnetic field strength is proportional to the square root of the gas density and the turbulent velocity dispersion of the gas. Thus, in the DCF method, a good estimation of the gas density is required to obtain an accurate estimation of the magnetic field strength. Similarly, the gas density is a crucial factor affecting the rates of astrochemical reactions (Wakelam et al., 2010; Gong et al., 2017). Thus, a precise estimation of the gas density within GMCs is crucial for accurate prediction of molecular abundances and a better understanding of the chemical evolution in GMCs.

However, it is difficult to quantify the number density of GMCs from observations. The traditional approach of estimating the number density of GMCs is based on observations of column density and certain assumptions on the geometry of the clouds, for example, a cylindrical geometry for filamentary structures or spherical geometry for dense cores (André et al., 2014). Bisbas et al. (2021; 2023) proposed an empirical power-law to convert the observed column density to the mean number density of GMCs based on the MHD simulations from Wu et al. (2017), which works decently but with noticeable scatter. Another method to constrain the number density of GMCs is utilizing density “probes”, such as cyanoacetylene (HC_3N , Avery et al., 1982; Schloerb et al., 1983; Li & Goldsmith, 2012). The relative intensity of different transitions of HC_3N is sensitive to the number density of the cloud, which makes it possible to constrain the mean number density directly by observing multiple transitions of HC_3N . Li & Goldsmith (2012) successfully observed $J = 2 - 1$ and $10 - 9$ transitions of HC_3N in Taurus B213 filament and constrained the number density of H_2 molecules, $n_{\text{H}_2} \sim (1.8 \pm 0.7) \times 10^4 \text{ cm}^{-3}$. Note, $n_{\text{H}} = 2n_{\text{H}_2}$ under the assumption that all H is in the form of H_2 . Unfortunately, the line ratio of HC_3N can only probe the number density at a relatively narrow range, between $\sim 10^4$ and 10^6 cm^{-3} (Li & Goldsmith, 2012), which limits its ability to infer the number density of the full range of structures that exist in GMCs. Consequently, a novel method to infer the number density of GMCs under a variety of physical conditions with high precision is in great demand. Machine learning makes it possible to learn from both the morphology

²We will also use the number density of H nuclei as a metric of density. Under the assumption of an abundance of one He nucleus for every 10 H nuclei in interstellar gas, we have a mass per H nucleus of $\mu_{\text{H}} = 1.4m_{\text{H}} = 2.34 \times 10^{-24} \text{ g}$. Thus $n_{\text{H}} = 1 \text{ cm}^{-3}$ is equivalent to $\rho = 2.34 \times 10^{-24} \text{ g cm}^{-3}$.

of the cloud and their column density to infer the mean number density rather than using a simple average power-law conversion.

Machine learning has gained great popularity among astronomers. For example, convolutional neural networks (CNNs) have been successfully applied to a series of tasks, including galaxy classification (Domínguez Sánchez et al., 2018; Becker et al., 2021), identification of structures like protostellar outflows, stellar wind-driven bubbles and Galactic cirrus filaments (Xu et al., 2020a;b; 2022a; Smirnov et al., 2023) and infer physical quantities based on observations, such as protostellar outflow inclination angles, magnetic field directions, stellar masses, exoplanet masses, galactic redshifts, and galactic star-formation rates (Liu et al., 2020; Zhang et al., 2022; Xu et al., 2022b; 2023; Bisigello et al., 2023). Furthermore, CNNs have also been utilized to mitigate the impact of noise in astronomical observations (Gheller & Vazza, 2022; Bartlett et al., 2023). For instance, Gheller & Vazza (2022) employed CNNs to remove noise and artifacts of radio interferometric images, while Bartlett et al. (2023) employed CNNs to diminish the impact of noise on various observation targets and preserve the morphology of galaxies. Meanwhile, Generative Adversarial Networks (GANs) have been applied to a variety of tasks (Hemmati et al., 2022; Sweere et al., 2022). Sweere et al. (2022) applied GANs to generate super-resolution and de-noised images from the *XMM-Newton* telescope. Hemmati et al. (2022) utilized GANs to effectively deblend galaxies from HST observations. More recently, Denoising Diffusion Probabilistic Models (DDPMs) have demonstrated their proficiency and robustness in image generation and editing (Sohl-Dickstein et al., 2015; Ho et al., 2020; Zhu et al., 2023a), which are suitable for the prediction task in astronomy.

B DISCUSSIONS ON DDPMs OVER OTHER METHODS

While this work seems to be an intuitive and direct application to an existing astrophysical problem, we provide our extended discussions on the reasons for the superiority of DDPMs compared to other existing deep learning based methods for task of number density predication of molecular clouds in astronomy.

The performance of the proposed diffusion model exhibits an improvement of one order of magnitude compared to some existing ML based methods such as CASI-2D (Van Oort et al., 2019) which is a convolutional neural network similar to Variational Autoencoders (VAEs) (Kingma & Welling, 2014). In the context of machine learning, we discuss the reasons for the significant improvement achieved by our proposed diffusion model as follows. The DDPMs are formulated based on the Markov stochastic process and model a random walk in the data space, which aligns with most existing physics problems and is consistent with intrinsic properties of the natural world. In contrast, CNNs and VAEs were originally designed for image classification and generation tasks in computer vision and lack explicit connections to the physical world. For instance, the observed structure of GMCs is likely to be shaped, at least in part, by turbulent motions that involve compressions in a series of quasi-random directions. Then the overall mass surface density map is constructed by summing a series of quasi-independent patches of volume density along the line of sight. Thus, inferring the raw mass-weighted number density distribution based on the observed column density inherits the basic concepts of diffusion models. Furthermore, with respect to information loss, DDPMs maintain the same data dimensionality throughout the entire denoising (i.e., prediction) process, thus better preserving the information conveyed by the original data. In contrast, CNNs and VAEs involve dimension reduction and information compression during training, resulting in inevitable information loss for the prediction objective. In terms of traceability and interpretability, we employ pre-defined Gaussian transition kernels to introduce and remove noise at each diffusion step in DDPMs. This provides us with superior traceability and interpretability for the data transition compared to CNNs and VAEs, whose traceability relies on a relatively vague gradient descent optimization direction.