

Towards characterizing dark matter subhalo perturbations in stellar streams with graph neural networks

PETER XIANGYUAN MA ^{1,2} KEIR K. ROGERS,^{3,4} TING S. LI ⁵ RENÉE HLOŽEK,^{4,5} JEREMY WEBB,^{6,5} RUTH HUANG,⁵
AND JULIAN MEUNIER⁵

¹*Department of Astronomy, UC Berkeley, 501 Campbell Hall, Berkeley, CA, 94720, United States of America*

²*Department of Mathematics, University of Toronto, 40 St. George Street, Toronto, ON, M5S 2E4, Canada*

³*Department of Physics, Imperial College London, Blackett Laboratory, Prince Consort Road, London, SW7 2AZ, United Kingdom*

⁴*Dunlap Institute for Astronomy and Astrophysics, University of Toronto, 50 St. George Street, Toronto, ON, M5S 3H4, Canada*

⁵*David A. Dunlap Department of Astronomy and Astrophysics, University of Toronto,
50 St. George Street, Toronto, ON, M5S 3H4, Canada*

⁶*Department of Science, Technology and Society, Division of Natural Science, York University,
218 Bethune College, Toronto, ON, M3J 1P3, Canada*

ABSTRACT

The phase space of stellar streams is proposed to detect dark substructure in the Milky Way through the perturbations created by passing subhalos — and thus is a powerful test of the cold dark matter paradigm and its alternatives. Using graph convolutional neural network (GCNN) data compression and simulation-based inference (SBI) on a simulated GD-1-like stream, we improve the constraint on the mass of a $[10^8, 10^7, 10^6] M_{\odot}$ perturbing subhalo by factors of $[11, 7, 3]$ with respect to the current state-of-the-art density power spectrum analysis. We find that the GCNN produces posteriors that are more accurate (better calibrated) than the power spectrum. We simulate the positions and velocities of stars in a GD-1-like stream and perturb the stream with subhalos of varying mass and velocity. Leveraging the feature encoding of the GCNN to compress the input phase space data, we then use SBI to estimate the joint posterior of the subhalo mass and velocity. We investigate how our results scale with the size of the GCNN, the coordinate system of the input and the effect of incomplete observations. Our results suggest that a survey with $10\times$ fewer stars (300 stars) with complete 6-D phase space data performs about as well as a deeper survey (3000 stars) with only 3-D data (photometry, spectroscopy). The stronger constraining power and more accurate posterior estimation motivate further development of GCNNs in combining future photometric, spectroscopic and astrometric stream observations.

Keywords: Machine Learning — Simulation Based Inference — Galactic Dynamics — Dark Matter

1. INTRODUCTION

The cold dark matter (CDM) model remains most preferred given astrophysical and cosmological observations (Aghanim et al. 2020), but its fundamental nature is undetermined despite an extensive direct detection program (e.g., Cooley et al. 2022). The hierarchical clustering within the CDM paradigm is successful in matching theory with observations of structure formation from the largest scales (\sim Gpc) down to galactic scales (\sim

Mpc): dark matter (DM) forms halos from galaxy cluster masses ($\sim 10^{15} M_{\odot}$) down to dwarf galaxy masses ($\sim 10^8 M_{\odot}$). The CDM model also predicts that substructure will form on even smaller scales, i.e., in subhalos with masses $< 10^8 M_{\odot}$ that do not host any baryonic matter. However, this regime remains largely untested. It is therefore a focus of current and upcoming photometric and spectroscopic surveys, e.g., the Vera C. Rubin Observatory (Rubin, Drlica-Wagner et al. 2019), *Euclid* (Blanchard et al. 2020), the Dark Energy Spectroscopic Instrument (DESI, Valluri et al. 2022) and the Nancy Grace Roman Space Telescope (Roman, Eifler et al. 2021), to infer robustly the existence of dark substructure within the Milky Way (MW). Doing so would

be a powerful confirmation of the CDM paradigm. On the other hand, proving the absence (or, indeed, enhancement) of substructure would help distinguish between theoretically well-motivated alternatives to CDM.

Although subhalos are not directly observable, there are indirect probes of their presence. Strong gravitational lenses are sensitive to substructure along the line of sight from source to observer, both within and outside the lens (Mandelbaum et al. 2006; Massey et al. 2010; Vegetti et al. 2010, 2023), with current sensitivity down to $\sim (10^7 - 10^8)M_\odot$. The Lyman- α forest, a spectral feature formed in the intergalactic medium (IGM), traces quasi-linear DM fluctuations (in the filaments and voids where the IGM resides) on the smallest scales currently accessible (Croft et al. 1999, 2002; McDonald et al. 2000; Weinberg 2003; Boera et al. 2019; Villasenor et al. 2023). Rogers & Peiris (2021a,b); Rogers et al. (2022) set a limit on the allowed minimum half-mode halo mass of $7.2 \times 10^7 M_\odot$ (95 % c.l.). The luminosity function of MW satellite galaxies can be related to the mass function of the subhalos that host them (Vale & Ostriker 2004; Macciò et al. 2010; Wechsler & Tinker 2018; Nadler et al. 2021). Combined with strong lensing, the minimum half-mode halo mass (for a warm DM transfer function) is limited above $\sim 10^7 M_\odot$ (95 % c.l., Nadler et al. 2021). One of the most powerful probes of substructure in upcoming surveys is the phase space of streams of stars in the Milky Way, e.g., the Legacy Survey of Space and Time (LSST) from *Rubin* is forecast to probe down to $\sim 10^5 M_\odot$ (Drlica-Wagner et al. 2019).

1.1. *Stellar streams as a probe of dark substructure*

Stellar streams are long, thin and dynamically cold trails of stars (Lynden-Bell & Lynden-Bell 1995) formed from tidally disrupted globular clusters and dwarf galaxies (Newberg & Carlin 2016). The formation of a stream in the absence of DM substructure can be modeled (see § 1.4) by accounting for the dynamical tidal disruption (Bovy 2014), the time-varying background potential of the MW (Minchev et al. 2010), merger with the Large Magellanic Cloud and interaction with baryonic (as opposed to DM) substructure like giant molecular clouds (Erkal et al. 2019). Any additional perturbations in the structure of the stream (both star positions and kinematics) will indicate interactions with DM subhalos (Johnston et al. 2002; Ibata et al. 2002; Erkal & Belokurov 2015; Erkal et al. 2016; Bovy 2016). Dozens of streams have been discovered and characterized in the MW halo, thanks to deep- and wide-field photometric and spectroscopic surveys (e.g., Koposov et al. 2014; Shipp et al. 2018; Jethwa et al. 2018; Li et al. 2019), as well as astrometric data provided by *Gaia* (e.g., Ibata

et al. 2019, 2021; Malhan & Rix 2024). Several stellar streams have already been studied for evidence of density variations and potential signatures of dark halo interactions. Notable examples include GD-1 (Grillmair & Dionatos 2006; Price-Whelan & Bonaca 2018; Bonaca et al. 2019), Palomar 5 (Odenkirchen et al. 2001; Carlberg et al. 2012; Erkal et al. 2017) and Atlas-Aliqa Uma (Koposov et al. 2014; Li et al. 2021; Hilmi et al. 2024), among others (e.g., Patrick et al. 2022; Tavangar et al. 2022).

The current state-of-the-art in the full statistical analysis of GD-1 (i.e., beyond modeling specific, identifiable features) is to measure the 1D angular power spectrum of the star density along the arc of the stream, probing subhalo masses from $10^6 M_\odot$ to $10^8 M_\odot$ (Bovy et al. 2017; Hermans et al. 2021; Banik et al. 2021b).

It is manifest that the 1D angular power spectrum does not extract all the information contained in the stream. First, information in the second angular coordinate is lost as it is known that streams contain structures away from the main orbit like spurs of stars, which could be a clear signature of a subhalo perturbation (Erkal et al. 2016). Second, by combining photometry with astrometry and spectroscopy, the full 6-D phase space of the stream can be built up (at least for a subset of stars observed in a given stream). E.g., the additional kinematic data from *Gaia* (Brown et al. 2021) provide two velocity components. *Rubin* LSST is expected to measure proper motions over a decade-long baseline at the current *Gaia* precision but with images that are three magnitudes deeper (Bonaca & Price-Whelan 2024). Spectroscopic measurements, e.g., from the Southern Stellar Stream Spectroscopy Survey (S^5 , Li et al. 2019) or DESI (Cooper et al. 2023), also provide line-of-sight velocities helping to complete the 6-D phase space. It is therefore necessary to prepare future analysis pipelines to leverage the increasingly rich data that we anticipate from deep and wide photometric and spectroscopic surveys.

1.2. *Machine learning and graph neural networks*

The challenge in using streams to constrain DM is to infer the model parameters of a subhalo (or subhalos) interacting with a stream (in particular, the subhalo mass) given the (incomplete) 6-D phase space of hundreds or thousands of observable stars in a stream. As discussed above, the 1D angular power spectrum is a lossy compression of the data. However, field-level inference, where inference occurs directly from the positional data (although some form of neural compression is, in practice, often necessary), is increasingly demonstrated to be a viable and powerful approach, e.g., for the cosmic mi-

crowave background (Caldeira et al. 2019), galaxy clustering (Lemos et al. 2024), the Lyman- α forest (Nayak et al. 2024), astrometric lensing (Mishra-Sharma 2022). In these examples, data are compressed using a convolutional neural network. Nguyen et al. (2024) directly infers from the galaxy field without compression.

Given that we want to combine photometric, astrometric and spectroscopic observations, our data are better described as a point cloud, where each point (star) is labeled with a 6-D phase space vector. For such data, graph convolutional neural networks (GCNNs, Kipf & Welling 2016) are a suitable compression. GCNNs exploit the graph-like geometry that arises in modeling point cloud systems (usually from the spatial arrangement of these points). The graph structure is formed by identifying graph nodes as stars and graph edges as connections between neighbouring stars. The graph can thus capture pairwise interactions present in our data. GCNNs are successfully demonstrated to infer DM density profiles from star data in dwarf galaxies (Nguyen et al. 2023), to infer baryonic properties from dark matter subhalo properties (Wu & Jespersen 2023), to carry out symbolic regression from DM-only simulations (Cranmer et al. 2020). In this work, we investigate using GCNNs to compress the phase space of a stream down to estimators of the perturbing subhalo mass and velocity.

1.3. Simulation-based inference

After this neural compression of the data, standard approaches to parameter inference are typically no longer feasible, as the likelihood function can not now be easily analytically formed. However, since we can simulate the stream, we can therefore forward model the data and learn the posterior or likelihood directly from samples of the joint distribution of data and parameters. This concept encompasses a broad class of algorithms called simulation-based inference (SBI, Diggle & Gratton 1984; Gutmann et al. 2016; Papamakarios et al. 2019). Within this class, machine learning models, e.g., Gaussian mixture models (Duda & Hart 1974) or normalizing flows (Tabak & Vanden-Eijnden 2010; Tabak & Turner 2012; Papamakarios et al. 2019), learn either the posterior (Papamakarios & Murray 2016; Greenberg et al. 2019; Deistler et al. 2022), the likelihood (Papamakarios et al. 2018) or the ratio of likelihood to evidence (Hermans et al. 2019; Miller et al. 2022). SBI is increasingly used in, e.g., cosmology given the need to model complex and multivariate likelihoods (e.g., Leclercq 2018; Alsing et al. 2019; Cole et al. 2022; Chen et al. 2023; Lemos et al. 2023; Lin et al. 2023). SBI is also successfully applied in the analysis of stellar streams.

E.g., Hermans et al. (2021) uses SBI to infer the warm DM particle mass given stream density variations; Alvey et al. (2023) uses SBI given the full phase space of a stream; Nibauer et al. (2022) uses SBI to reconstruct the galactic acceleration field. Here, we use neural posterior estimation with a normalizing flow model (Tejero-Cantero et al. 2020) to learn the posterior distribution, using the GCNN as a data compression.

1.4. Stellar stream simulations

Since the GCNN and SBI models require $\mathcal{O}(10^3 - 10^4)$ training simulations, we must simulate stellar streams with a balance of computational efficiency and physical accuracy that gives robust inference. The most accurate, but most computationally expensive, approach is a hydrodynamical simulation of the Milky Way, e.g., the Feedback in Realistic Environments (FIRE, Hopkins et al. 2018; Shipp et al. 2023) simulations or the Numerical Investigation of Hundred Astrophysical Objects (NIHAO, Wang et al. 2015; Butsky et al. 2016) simulations. These simulations assume a fluid model of the formation of galaxies and their satellites, streams and substructure. They incorporate baryonic feedback (both mechanical and radiative) sources that also shape the DM distribution, e.g., stellar and supernovae feedback, feedback from active galactic nuclei and black holes, star formation quenching. These codes are computationally expensive: a lower-resolution FIRE simulation with 10^6 particles requires > 10 hours (Hopkins et al. 2018). There are faster N -body simulations of dissolving star clusters (e.g., Küpper et al. 2008, 2010; Ngan et al. 2015; Webb & Bovy 2019) that neglect detailed hydrodynamics, but model internal cluster evolution with much higher precision. It is still intractable to run these simulations in sufficient number for the training of our machine learning models.

Best suited for our purposes are particle spray-based methods (e.g., Fardal et al. 2015; Chen et al. 2024) that generate realistic representations of stellar streams without the need for a detailed cluster or galaxy simulation. We use the STREAMSPRAYDF algorithm (Fardal et al. 2015) as implemented in the GALPY Galactic dynamics package (Bovy 2015; Qian et al. 2022). We model a single stream by randomly ejecting stars from a progenitor. In this proof-of-principle work, we consider interactions between a GD-1-like stream and a single subhalo, although a real stream will appreciably interact with $\mathcal{O}(100)$ subhalos (Diemand et al. 2007, 2008; Erkal et al. 2016). This approach is computationally tractable to run in large numbers and is demonstrated to reproduce observed stream properties like their extent and the “streaky” features arising from the disruption of the

progenitor (Fardal et al. 2015). In § 4, we discuss how to incorporate additional physics like multiple subhalos (Erkal et al. 2016), a more physical tidal disruption that reproduces features like the stream cocoon (Carlberg 2018; Malhan et al. 2019; Gialluca et al. 2021; Qian et al. 2022) and a time-dependent background potential (Buist & Helmi 2015; Koppelman & Helmi 2021; Brooks et al. 2024).

In this work, we set out to prove in principle whether the GCNN data compression and SBI inference method described above can reliably estimate the posterior distribution of the mass and velocity of a subhalo that interacted with a stellar stream. We describe the simulations, data compression/encoder and SBI in § 2. In § 3, we present the results of a comparison of our new approach to the current state-of-the-art 1D angular power spectrum analysis. In § 4, we discuss our results and then conclude in § 5.

2. METHODS

We describe the simulation procedure in § 2.1 and the training set in § 2.1.1. We describe the 1D angular power spectrum in § 2.2.1 and the GCNN model in § 2.2.2. After compressing the simulated data, we present in § 2.3 the SBI method that we use to estimate the posterior distribution of the simulation parameters. In § 2.4, we explain the performance metrics that we will use in § 3.

2.1. STREAMSPRAYDF *simulations*

In this proof-of-principle work, we simulate a GD-1-like stream interacting with a single subhalo in a static MW potential. We discuss more sophisticated simulations that we will use in future work in § 4.5. We use the STREAMSPRAYDF algorithm (Fardal et al. 2015) implemented in GALPY (Bovy 2015; Qian et al. 2022) to model GD-1 by randomly ejecting stars from a progenitor according to the tidal disruption model of Fardal et al. (2015). We integrate the orbit of each star within a fixed background gravitational potential, with and without the moving gravitational potential of a subhalo. For each training simulation described in § 2.1.1, we vary the random realization of stars along the GD-1 stream orbit and systematically vary the mass and velocity of the perturbing subhalo in order to span the prior parameter volume. In detail, for each simulation, we follow these steps (see also Fig. 1):

1. We set the state of the GD-1 progenitor (mass, radius, position, velocity as determined in Grillmair & Dionatos (2006); Koposov et al. (2010); Gialluca et al. (2021)) today at time $t = 0$. The progenitor mass is set to a constant $10^4 M_\odot$ and scaled radius of 0.01 kpc.

2. We integrate the progenitor orbit backwards from $t = 0$ in a static MW potential with no substructure until $t = t_i$ in the past. We choose $t_i = -2$ Gyr in order to match the observed length of the GD-1 stream (Gialluca et al. 2021). For the MW potential, we use the MWPOTENTIAL2014 setting in GALPY, which combines a Navarro-Frenk-White (NFW, Navarro et al. 1996) halo potential, a Miyamoto-Nagai disc potential (Miyamoto 1975) and a power-law density spherical bulge potential with an exponential cut-off (Cardone et al. 2005).
3. Using the STREAMSPRAYDF algorithm (Fardal et al. 2015), we evolve the progenitor forward in time, ejecting stars during integration to obtain an evolved stream without a subhalo interaction. We eject $n_{\text{stars}} = 3000$ stars, with each stream arm containing 1500 stars, in order to have a converged estimate of the angular power spectrum.¹
4. We select a target star with which the subhalo will most closely interact. We choose the time of interaction to be -200 Myr.
5. We set the state of the subhalo (mass, radius, position, velocity relative to the stream, impact parameter, angle of approach to the stream) at the time of interaction. We vary the mass and relative velocity in the training set (§ 2.1.1). For the subhalo potential, we use a Hernquist potential (Hernquist 1990).
6. We integrate the subhalo orbit backwards from the time of interaction in the MW potential until t_i in the past.
7. We evolve the progenitor forward in time again from t_i to today, ejecting stars to form the GD-1 stream, now with the background MW potential and the moving subhalo potential to model the subhalo interaction.
8. The final result is the 6D phase space today of the stars in a GD-1-like stream.

In Fig. 2, we visualize examples of the simulated GD-1 streams as they appear today with and without subhalo perturbations. We project the streams into right ascension (RA), declination (DEC) and line-of-sight position

¹ We assess convergence as when statistical fluctuations in the power spectrum as the number of stars increases drops below 10%.

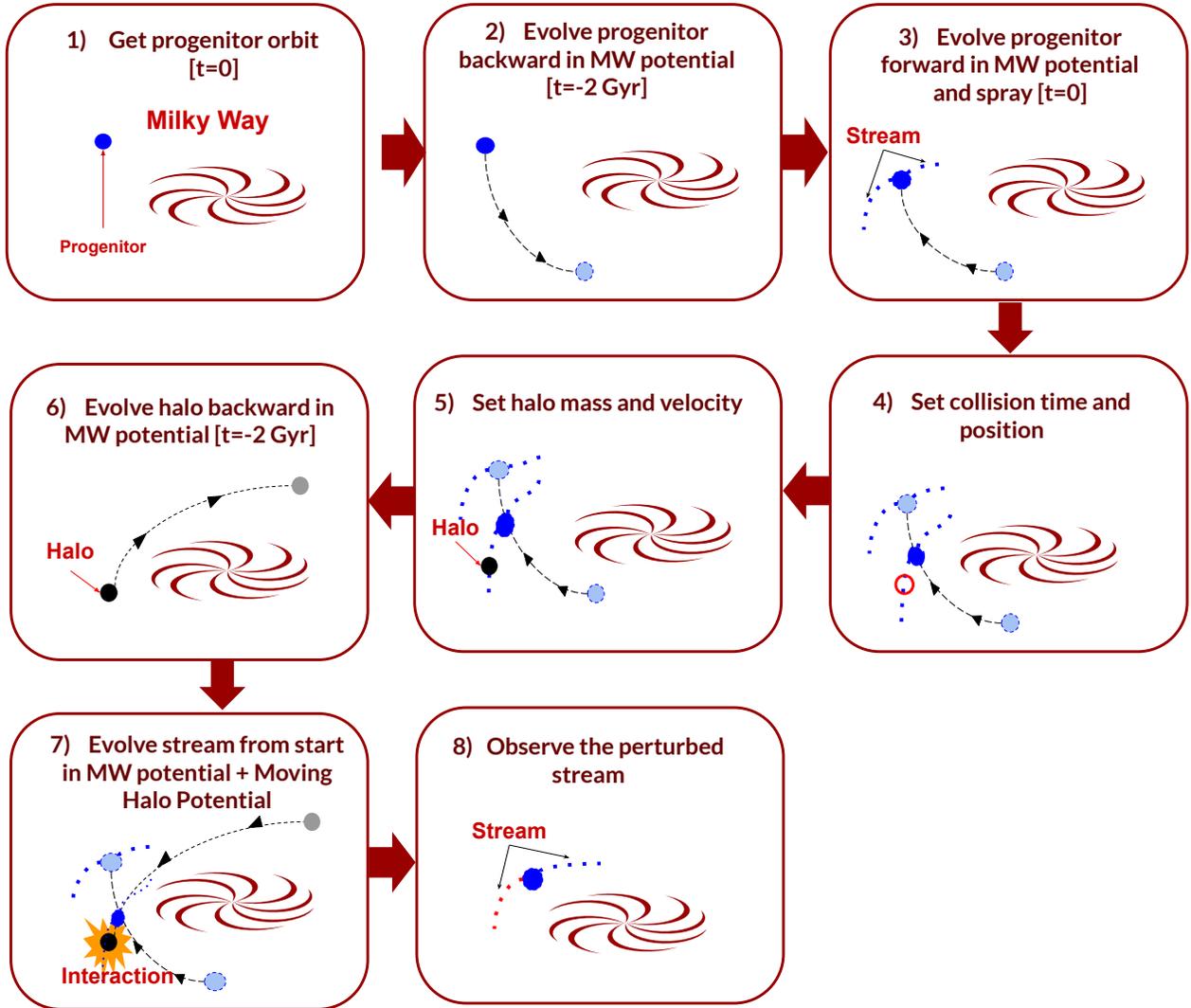


Figure 1 The steps in the stream simulation procedure (see main text for details).

and velocity coordinates, as well as ϕ_1 and ϕ_2 coordinates which are transformations of RA and DEC using the rotation matrix defined in Koposov et al. (2010). ϕ_1 and ϕ_2 are defined such that the stream has maximum extent along ϕ_1 . A perturbing subhalo is qualitatively able to produce observed stream features like gaps, spurs and blobs (Bonaca et al. 2019). The strength of the stream perturbation decreases as the subhalo decreases in mass to the extent that a gap may no longer form, but rather there is a reduction in the number density of stars near the point of interaction. This qualitative change in behavior occurs around $10^7 M_\odot$ consistent with the previous simulations of Bonaca et al. (2019). However, there is a visual degeneracy between the effect of a slower, less massive subhalo and a faster, more massive subhalo. This effect arises because a slower subhalo spends more time closer to the stream and thus

exerts a greater impulse on the stars. Thus, in order to break this degeneracy between subhalo mass and velocity, it will be crucial to extract as much information as possible from the full stream phase space, in particular the velocity distribution as shown in the bottom row of Fig. 2. As expected, a subhalo perturbation produces a tail of high-speed stars in the stream. An animation of an example stream simulation can be found at <https://www.youtube.com/watch?v=d2nWvScdbJQ>.

2.1.1. Generation of training simulations

For the machine learning (ML) in § 2.2 and § 2.3, we generate two training sets, respectively with 40,000 and 80,000 simulations, in order to test how performance scales with the size of the ML model. Each training set is reused for the encoder and SBI models. In each set, we vary subhalo masses m distributed evenly in

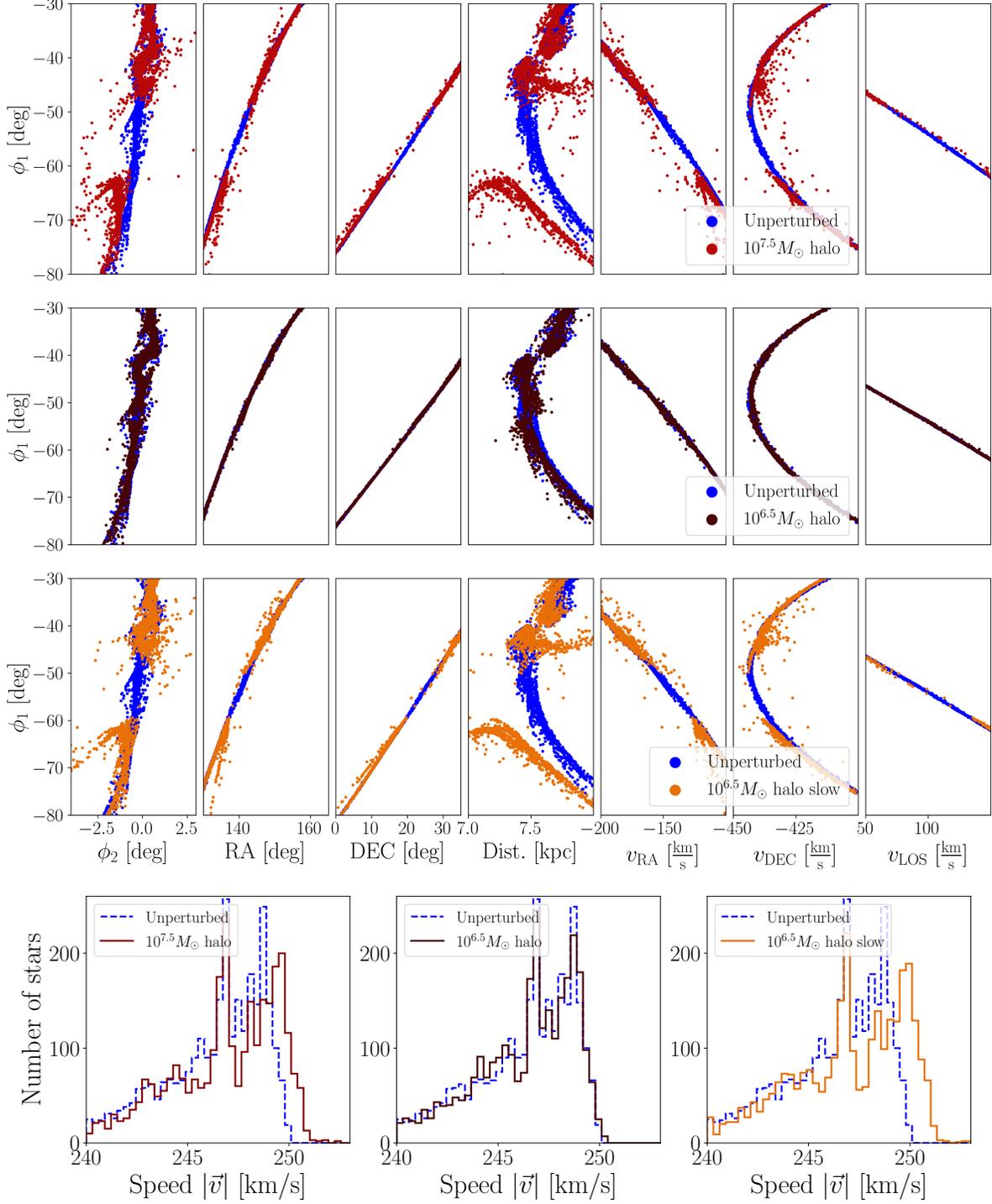
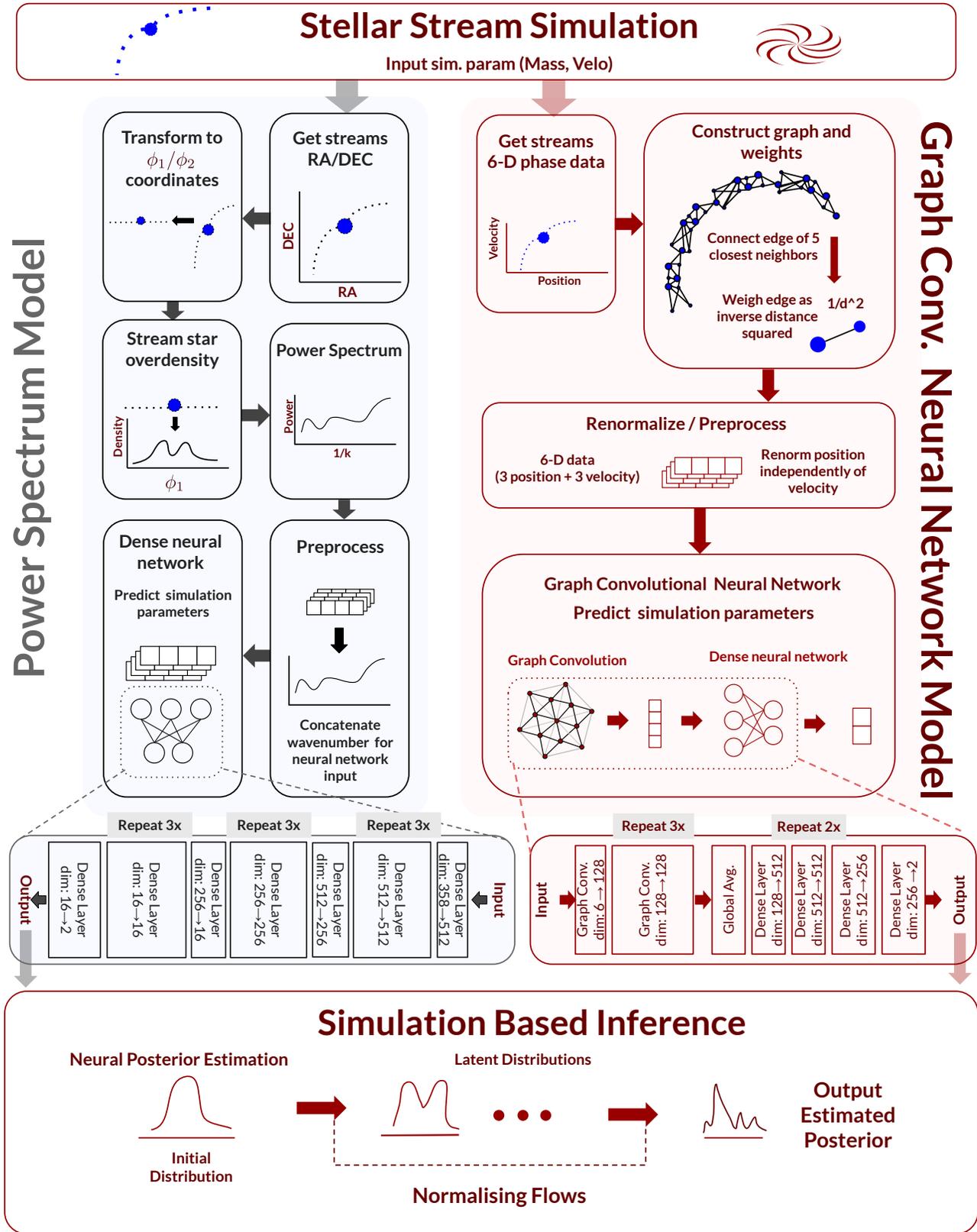


Figure 2 *In the top three rows, the phase space of the GD-1 stellar stream simulations today. We compare the stream modeled without a subhalo (blue) to the perturbed stream. In the top panel, the perturber is a subhalo of mass $10^{7.5} M_\odot$ with a velocity of 220 km/s relative to the stream; in the middle, mass = $10^{6.5} M_\odot$ with the same velocity; in the bottom, mass = $10^{6.5} M_\odot$ with relative velocity = 22 km/s. From left to right, we project the streams into the sky coordinates ϕ_1 , ϕ_2 , right ascension (RA), declination (DEC), heliocentric distance (Dist.), transverse velocity arising from the RA v_{RA} and DEC v_{DEC} components of the proper motion and the line-of-sight velocity v_{LOS} . To decorrelate the position and velocity vectors, we use the transverse velocity defined as the product of the proper motion and the heliocentric distance. In the bottom row, the distribution of the speeds of stars in the same simulations as above.*



Simulation Based Inference

Neural Posterior Estimation

Initial Distribution

Latent Distributions

Normalising Flows

Output Estimated Posterior

Figure 3 The steps in the encoder (§ 2.2) and simulation-based inference (§ 2.3) parts of the pipeline. We consider two possible encoders: (*left*) the 1D angular power spectrum compressed by a dense neural network (§ 2.2.1) and (*right*) a graph convolutional neural network (GCNN) compressing directly from the stream phase space. The GCNN layer structure is specifically shown for the baseline GCNN model variant (Table 1).

the range $5 \leq \log \frac{m}{M_\odot} \leq 10$. We include masses down to the forecast observational limit of the *Rubin* LSST survey (Drlica-Wagner et al. 2019). The upper prior bound excludes catastrophic interactions with massive satellites that would destroy the stream. We simultaneously vary subhalo velocities relative to the stream v in the range $22 \leq \frac{v}{\text{km/s}} \leq 1100$. The upper prior bound arises as anything faster would leave no appreciable effect on the stream, while the lower bound approaches the limit where the subhalo is moving with the stream. In each simulation, we vary the random realization of stars that are ejected from the progenitor, but always with a progenitor orbit that matches GD-1 observations. We sample the prior volume with a Latin hypercube (McKay et al. 1979). The creation of the 40,000 simulation set takes approximately one week of computation on a 128 core machine with 1TB of RAM. The numerical integration in GALPY uses a fourth-order Runge–Kutta method implemented in the C language with temporal discretization into 2000 intervals.

2.2. Encoder

The encoder is the layer of computation within our pipeline that compresses the stream phase space into a set of informative summary statistics. In this work, we choose to compress maximally to estimators of the simulation parameters, i.e., subhalo mass and velocity. We compare two encoders: the current state-of-the-art in data analysis, the 1D angular power spectrum of the star density (then compressed using a neural network to parameter estimators, PS, § 2.2.1); and a graph convolutional neural network in order to compress directly from the phase space (GCNN, § 2.2.2). For the GCNN encoder, we then introduce four model variants: (i) a larger GCNN with twice the number of training simulations (GCNN-L); (ii) a change of the input coordinate system from galactocentric to heliocentric (GCNN helio (6-D)); (iii) the effect of a reduced phase space (GCNN helio (3-D), only RA, DEC, v_{LOS}); and (iv) a reduction in the number of stars from 3000 to 300 (GCNN (6D) 300 stars). We summarize these variants in Table 1. We show in Fig. 3 a schematic summary of the full pipeline from input simulations (§ 2.1) through the two different encoders that we consider to the final simulation-based inference step (§ 2.3).

2.2.1. 1D angular power spectrum

The 1D angular power spectrum [PS; $P(k_{\phi_1})$] of the star number overdensity $\delta(\phi_1)$ quantifies the variance in the number density of stars as a function of angular wavenumber k_{ϕ_1} along the extent of the stream:

$$P(k_{\phi_1}) = \left| \hat{\delta}(k_{\phi_1}) \right|^2 = \left| \int \delta(\phi_1) e^{-2\pi i k_{\phi_1} \phi_1} d\phi_1 \right|^2. \quad (1)$$

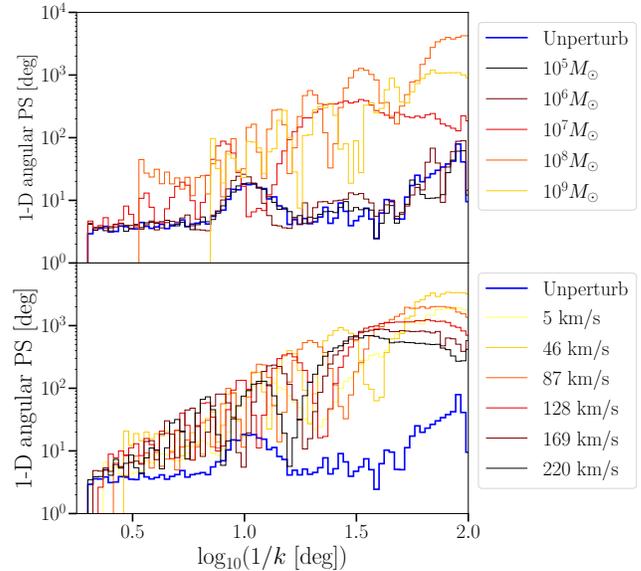


Figure 4 The 1D angular power spectrum $P(k_{\phi_1})$ of the GD-1 star number overdensity today as a function of the inverse of the angular wavenumber k_{ϕ_1} , for an unperturbed stream (*blue*) and streams with different perturbers. *In the top panel*, we vary the mass of the perturbing subhalo with velocity fixed to 110 km/s; *in the bottom panel*, we vary the velocity of the subhalo with mass fixed to $10^{7.5} M_\odot$. Here, we show the power spectrum averaged over 100 random stream realizations in order to mitigate sample variance. A more massive and/or a slower subhalo increases the power on a scale of tens of degrees. A characteristic sinc pattern arises in the Fourier space from the tophat-like gap that forms in the stream density.

Here, $\hat{\delta}(k_{\phi_1})$ is the Fourier transform of $\delta(\phi_1)$ and k_{ϕ_1} is the variable conjugate to ϕ_1 . $\delta(\phi_1) = \frac{n(\phi_1) - \bar{n}(\phi_1)}{\bar{n}(\phi_1)}$, where $n(\phi_1)$ is the number density of stars and $\bar{n}(\phi_1)$ is the average number density in the stream. The PS is sensitive only to Gaussian-distributed star perturbations in the ϕ_1 coordinate.

Figure 4 demonstrates the sensitivity of the power spectrum to stream perturbers of different mass and velocity. We show an increase in power on larger angular scales as the mass of the perturber increases. This trend arises as the perturbation scatters the spatial distribution of stars, in particular on the angular scale of the gap that forms in the stream. Indeed, we see a characteristic sinc pattern in the power spectrum, which is the Fourier transform of the tophat-like gap that forms. We also show an increase in power on larger angular scales as the velocity of the perturber decreases. The less the velocity of the perturber, the more time the subhalo is close to the stream. This effect increases the impulse of

Model	Encoder	Coordinates	# sim.	# stars	# model pars.
PS	Power spectrum	ϕ_1	4×10^4	3000	7×10^5
GCNN	Graph CNN	Galactocentric spherical 6-D	4×10^4	3000	7×10^5
GCNN-L	Graph CNN	Galactocentric spherical 6-D	8×10^4	3000	1.4×10^6
GCNN helio (6-D)	Graph CNN	[RA, DEC, dist., v_{RA} , v_{DEC} , v_{LOS}]	8×10^4	3000	1.4×10^6
GCNN helio (3-D)	Graph CNN	[RA, DEC, v_{LOS}]	8×10^4	3000	1.4×10^6
GCNN (6-D) 300 stars	Graph CNN	[RA, DEC, dist., v_{RA} , v_{DEC} , v_{LOS}]	8×10^4	300	1.4×10^6

Table 1: The encoder model variants that we consider. *From left to right*, we give the model name, the type of encoder, the input coordinate system, the number of training simulations, the number of input stars, the number of encoder neural network hyperparameters that are trained.

the subhalo on the stream and scatters more stars. The same sinc pattern is generated, where a slower subhalo causes a larger gap leading to a more strongly oscillating pattern in Fourier space. There is thus a degeneracy in the power spectrum between increasing the perturber mass and decreasing its velocity.

In order to infer simulation parameters from the power spectrum, we must construct a likelihood function. Following, e.g., Banik et al. (2021b), we do not assume a functional form for the likelihood² and instead use simulation-based inference (§ 2.3), which also allows a more direct comparison to the GCNN approach (§ 2.2.2). In order to use SBI, we must reduce the dimensionality of the simulated data further as the density estimation task in SBI scales with the sum of the dimensions of the data and parameters. We do this compression using a dense neural network to compress the power spectrum to two numbers, which are estimators of the subhalo mass and velocity.

We proceed by normalizing $\log P(k_{\phi_1})$ to the unit interval as input to the neural network. In order also to feed the neural network information on the wavenumber dependence of the power spectrum, we concatenate to the power spectrum vector, a vector of inverse wavenumbers $\frac{1}{k_{\phi_1}}$ at the centers of the bins that we use, also normalized to the unit interval. We use 179 bins distributed evenly in $\log \frac{1}{k_{\phi_1}}$, thus meaning that the first layer of the network has 358 elements.

The dense neural network applies to the input vector a series of composed parameterized functions (layers). In succession, each n^{th} layer \mathbf{l}_n first linearly transforms the previous layer output \mathbf{l}_{n-1} by matrix multiplying an optimizable weight matrix $\mathbf{W}^{(n)}$ and then adding an optimizable bias vector $\mathbf{b}^{(n)}$ before applying a non-linear activation function α (Hornik et al. 1989):

$$\mathbf{l}_n = \alpha \left(\mathbf{W}^{(n)} \mathbf{l}_{n-1} + \mathbf{b}^{(n)} \right). \quad (2)$$

² Banik et al. (2021b) use a form of SBI called approximate Bayesian computation.

The structure of these layers, i.e., the number of elements in each, is shown in Fig. 3, while the details are given in Appendix A. We use Leaky ReLU as the activation function. The weights and biases are optimized by minimizing a loss function $\mathcal{L}(\vec{x}, \vec{\theta}; p_\gamma)$ of the final layer. We use the mean absolute error (MAE) between the output of the network $p_\gamma(\vec{x})$ and the true subhalo parameters $\vec{\theta} = [m, v]$ (normalized to the prior volume defined in § 2.1.1) for a set of training simulations:

$$\mathcal{L}(\vec{x}, \vec{\theta}; p_\gamma) = \frac{1}{2} \sum_{i=1}^2 \left| p_\gamma^{(i)}(\vec{x}) - \theta_i \right|, \quad (3)$$

where $\vec{x} = \left[P(k_{\phi_1}), \frac{1}{k_{\phi_1}} \right]$ is the input vector to the network, $p_\gamma^{(i)}$ represents the neural network transformation of the input and i indexes over the two model parameters (mass and velocity).

We use backpropagation to train the network (Hinton & Williams 1986). Backpropagation is an efficient means of computing the gradients of the loss with respect to the network’s hyperparameters. This method is useful for the gradient-based optimization algorithm, adaptive moment estimation (ADAM, Kingma & Ba 2014), that we use. In building the neural network, we optimize its configuration by performing a neural architecture search. We do this using Bayesian optimization (see Appendix A). We split the input simulation set (see § 2.1.1), with 80% for training and 20% for validation; the test set has 9000 simulations drawn as a Latin hypercube across the prior area. We use the ADAM optimizer with a constant learning rate of 10^{-3} and a batch size of 32. We train the network until saturation (early stopping). We give additional training details in Appendix B.

2.2.2. Graph convolutional neural network

As a less lossy alternative to the power spectrum (§ 2.2.1), we also compress the stream to subhalo parameter estimators using a neural network directly applied to the phase space, i.e., without the power spectrum

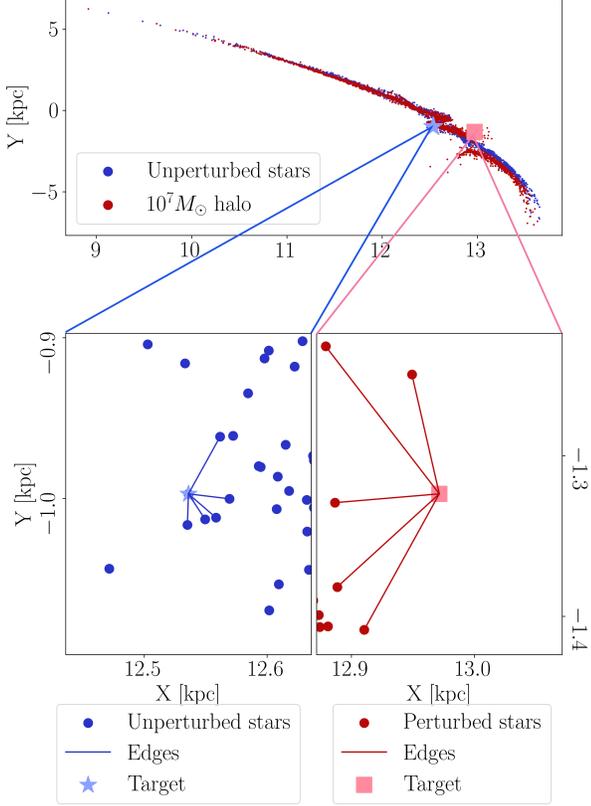


Figure 5 The *top panel* projects GD-1 into galactocentric X and Y coordinates, with *blue* points showing the positions of stars in an unperturbed realization and *red* points showing stars in a realization with a perturbing subhalo of mass $10^7 M_\odot$ (same random seed in each realization). In the *bottom panels*, we show zoom-ins on the same target star in each realization (unperturbed on the *left* and perturbed on the *right*). In the zoom-ins, we show the graph that we construct for the target star, connecting the star with edges to the five nearest neighbors. We illustrate how the graph is sensitive to the geometry of the stream that the perturbation induces.

“pre-compression”. Rather than use a dense neural network as described in § 2.2.1, we use a graph convolutional neural network, at least for the first part of the pipeline. First, we construct a graph from the stream phase space.

In Fig. 5, we show an example of the graphs that we construct (specifically, part of the total stream graph for a single target star). To form a graph, we connect each star in the stream, as a node, to its k -nearest neighbors using an edge, where we set $k = 5$. At each node, we list the features of each star as its 6-D phase space position (or a subset of the full 6-D information depending on the

model variant, see Table 1), i.e., its position and velocity vector in the chosen coordinate system. Each edge e_{ij} is weighted by the square of the inverse distance from star i to neighboring star j . In Fig. 5, we highlight a single star and its connections and see that graphs capture local geometric structures that are highly sensitive to subhalo perturbations. Edges connecting perturbed stars (*bottom right panel*) are longer compared to those connecting unperturbed stars (*bottom left panel*). By connecting only the five closest neighbors, the graph is kept sparse which generates deeper networks and avoids oversmoothing (Zhang et al. 2024a). We observe that smaller values of k lead to a lack of connectivity, which degrades training performance. Therefore, we choose $k = 5$ as a fixed hyperparameter.

We proceed by normalizing the input features (i.e., the phase space vector) at each graph node by the standard deviation of each feature across the training set. Having constructed a graph with normalized features from the stream phase space, we then pass the graph as input to a GCNN. We first apply a series of graph convolutional layers. Each convolutional layer applies three operations: a message passing (where the features at each node are replaced by an aggregation of the features at connected nodes), a linear affine transformation (the application of a weight matrix and bias term) and a non-linear activation (these last two operations are identical to the neural network in Eq. (2)).

In other words, each n^{th} layer is now a matrix \mathbf{l}_n consisting of features (along rows) for each node (stars along columns). First, the previous layer \mathbf{l}_{n-1} is matrix multiplied by a symmetric adjacency matrix \mathbf{A} that encodes the message passing between connected nodes (i.e., the graph convolution). Then, as before, this convolved layer is matrix multiplied by an optimizable weight matrix $\mathbf{W}^{(n)}$ and then added to a bias term, which is now a bias matrix $\mathbf{b}^{(n)}$, i.e., a set of biases (along rows) for each node (stars along columns). Finally, we apply a non-linear activation function σ :

$$\mathbf{l}_n = \sigma \left(\mathbf{A} \mathbf{l}_{n-1} \mathbf{W}^{(n)} + \mathbf{b}^{(n)} \right). \quad (4)$$

The adjacency matrix \mathbf{A} has elements:

$$A_{ij} = \begin{cases} e_{ij}, & \text{if } i^{\text{th}} \text{ node connected to } j^{\text{th}} \text{ node,} \\ 0, & \text{else.} \end{cases} \quad (5)$$

In this way, the adjacency matrix performs the graph convolution by, at each node, replacing each feature by a weighted sum of the equivalent feature at connected nodes. This is analogous to an image convolution, as applied in a standard convolutional neural network, where

each pixel in each layer is an aggregation of nearby pixels weighted by a kernel. Here, the image kernel is replaced by the adjacency matrix.

The structure of the GCNN layers, i.e., the number of features at each node, is shown in Fig. 3, while the details are given in Appendix A. Our choice of edge weighting has the property that $e_{ij} = e_{ji}$ which enforces equivariance symmetry. This bi-directionality means that pairwise force interactions are captured in the compression network.

After the graph convolutional layers, we apply a layer that, for each of the final 128 features, averages the feature across all nodes (stars).³ We then pass this single vector of features through a series of standard layers [Eq. (2)], analogous to the dense neural network described in § 2.2.1. The final layer is again estimators of the two subhalo parameters. The structure of these final layers is also given in Fig. 3, while the details are again given in Appendix A. We otherwise follow the same steps as in § 2.2.1, i.e., we use the same loss function, optimizer, backpropagation, early stopping and training-validation split. We again use Bayesian optimization to optimize the architecture configuration (see Appendix A, meaning that the activation function is instead ReLU); we give additional training details in Appendix B. The training of the baseline GCNN model takes 19 hours on an NVIDIA RTX 6000 Ada Generation graphics processing unit. The slight modifications to the structure of the layers necessary for the GCNN encoder model variants (Table 1) are explained in § 3.3 and 3.4.

2.3. Simulation-based inference

Having compressed the stream phase space by the two alternative encoders presented in § 2.2, we now proceed to estimating the posterior distribution $\mathcal{P}(\vec{\theta}|\vec{\theta})$ of the simulation parameters $\vec{\theta} = [m, v]$ given these compressed summary statistics $\vec{\theta} = p_\gamma(\vec{x})$. Rather than assuming a simple approximation for the posterior or likelihood (e.g., a Gaussian function), we instead use SBI to learn a parameterized form of the posterior from a training set of forward models (simulations).

We use neural posterior estimation, as implemented in the `sbi` package (Tejero-Cantero et al. 2020), which, using neural networks, directly learns the posterior, rather than other approaches that learn the likelihood (Papamakarios et al. 2018) or the ratio of likelihood to evidence (Hermans et al. 2019; Miller et al. 2022). The parameterized distribution that we use to model the

posterior is a normalizing flow $\mathcal{P}_\phi(\vec{\theta}|\vec{\theta})$ (Papamakarios & Murray 2016; Greenberg et al. 2019; Papamakarios et al. 2019; Deistler et al. 2022). Normalizing flows take a simple base distribution (in our case, a Gaussian) and apply a series of parameterized, invertible, volume-preserving (normalized) transformations to match a more complex target distribution. In our case, the target is the true posterior $\mathcal{P}(\vec{\theta}|\vec{\theta})$. Specifically, we implement a masked autoregressive flow (Papamakarios et al. 2017). Here, we factor the posterior distribution into a series of conditional probabilities via the chain rule of probabilities (autoregressive). When training the flow, we mask variables to ensure the sequential variable ordering so we obey the chain rule.

In order to learn the parameters ϕ of this flow model, we use a neural network with simulated estimators $\vec{\theta}$ and true parameters $\vec{\theta}$ as input. We optimize this neural network by minimizing the divergence between the true and target distributions. This minimization is equivalent to minimizing a loss function $\mathcal{L}(\vec{\theta}, \vec{\theta})$ that is the negative log posterior for the training samples $\{\vec{\theta}_i, \vec{\theta}_i\}$:

$$\mathcal{L}(\vec{\theta}, \vec{\theta}) = \sum_i -\ln \mathcal{P}_\phi(\vec{\theta}_i|\vec{\theta}_i), \quad (6)$$

where, here, i indexes over the training simulations.

The training, validation and test simulation split follows the same proportions as for the encoder models. We optimize the SBI model using ADAM with a constant learning rate of 5×10^{-4} . Otherwise, we use the default `sbi` settings. Having learnt a model for the posterior distribution, we draw samples for a given mock data vector in order to generate the posterior intervals in § 3 (Casella et al. 2004).

2.4. Performance metrics

In § 3.1, we compare the performance of the power spectrum (§ 2.2.1) and GCNN (§ 2.2.2) encoders by calculating the fractional mean absolute error on the test set:

$$\text{MAE}_{\text{frac},i} \equiv \frac{|\hat{\theta}_i - \theta_i|}{\theta_i}, \quad (7)$$

where i indexes the two simulation parameters (subhalo mass and velocity) and $\hat{\theta}_i$ is the parameter estimator given either the power spectrum or GCNN compression. The smaller this metric, the more informative is the compressed summary statistic $\hat{\theta}_i$.

In § 3.2, we compare the performance of the SBI (§ 2.3) given power spectrum and GCNN encoders, with regards to both the precision and accuracy of estimated posterior distributions. In order to quantify the im-

³ We visualize the final graph convolutional layer in Appendix C.

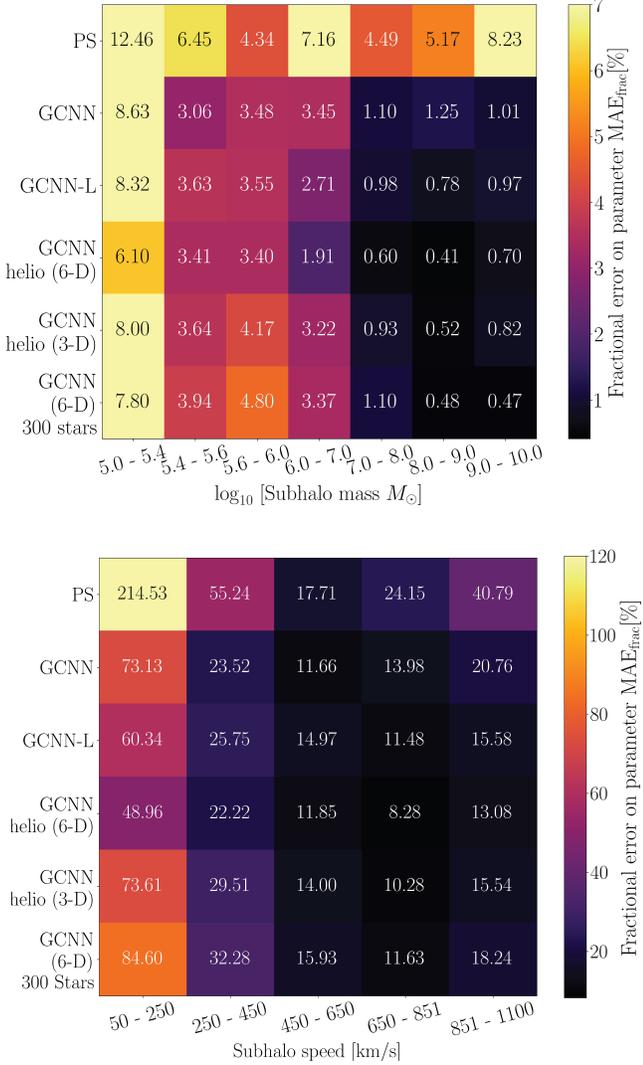


Figure 6 In the top panel, the fractional mean absolute error [as shown in Eq. (7)] on the true subhalo mass for a test set of stream simulations. From left to right in this panel, we show MAE_{frac} averaged over all test samples in the given mass bin. From top to bottom in this panel, we vary the encoder model variant (see Table 1). In the bottom panel, same as the top panel except computed for bins of subhalo speed.

provement in the precision of the subhalo mass constraint, we calculate the ratio of the subhalo mass posterior standard deviation σ relative to the best-performing encoder σ_{best} , which is the GCNN helio (6-D) model variant (Table 1):

$$\sigma_{\text{frac}} \equiv \frac{\sigma}{\sigma_{\text{best}}}. \quad (8)$$

In other words, we give this ratio such that larger values indicate that the marginalized posterior on the subhalo

mass is broader (than the best case) and can conclude that the encoder has extracted less information from the stream. We also sometimes (Fig. 7) calculate the full figure of merit for the 2D posterior:

$$FoM = \frac{1}{\sqrt{\det \mathbf{C}}}, \quad (9)$$

where \mathbf{C} is the posterior covariance between mass and velocity. The figure of merit quantifies the area of the peak of the posterior. By comparing this quantity, we assess the total information gain accounting for both subhalo mass and velocity.

We also assess the accuracy of the estimated posterior distribution, i.e., how correctly has the SBI procedure calculated the true posterior distribution given the input compressed data? First, we report the negative log probability of the true parameters averaged over the test set, i.e., the loss function [Eq. (6)] evaluated at the true parameters of the test set. The lower this metric, the higher probability is the true point (averaged over the test set), which is indicative of a better calibrated set of posteriors (Lueckmann et al. 2021).

Further, we perform a simulation-based calibration test (Cook et al. 2006; Talts et al. 2018). This test checks whether the peaks of posterior distributions generated by the SBI model are consistent with the true simulation parameters with the expected frequency; or, in other words, whether the credible intervals of the posterior have good coverage probabilities. We compute the posterior distribution given each test stream. We then draw samples from the posterior distribution and calculate a rank by counting how many samples fall below the true simulation parameter in terms of the value of the posterior probability. A well-calibrated set of posteriors will have a uniform distribution of true value ranks. If the rank distribution deviates from uniformity by having many ranks in the tails, this indicates that the posterior area is typically underestimated, i.e., over-confident posterior constraints or over-fitting. If the rank distribution has many ranks in the center, this indicates that the posterior area is typically overestimated, i.e., under-confident constraints or under-fitting.

We quantify the consistency of the measured rank distribution with a uniform distribution by the Kolmogorov-Smirnov (KS) test (Kolmogorov 1933; Smirnov 1948; Talts et al. 2018). The KS test returns p values on the null hypothesis that the measured ranks are drawn from a uniform distribution. We reject this null hypothesis if $p < 0.05$; we otherwise report that the posteriors are well calibrated. In practice, we find that posteriors near the edge of the prior can be poorly calibrated (owing to discontinuity in the uniform prior that

Encoder model	$(10^{5.75} - 10^6) M_{\odot}$	$(10^6 - 10^7) M_{\odot}$	$(10^7 - 10^8) M_{\odot}$	$(10^8 - 10^9) M_{\odot}$
PS	1.70 ± 0.02	2.65 ± 0.25	6.73 ± 0.76	10.55 ± 1.00
GCNN	1.15 ± 0.01	1.30 ± 0.01	1.67 ± 0.28	2.02 ± 0.10
GCNN-L	1.15 ± 0.05	1.23 ± 0.01	1.47 ± 0.15	1.83 ± 0.29
GCNN helio (6-D)	1	1	1	1
GCNN helio (3-D)	1.18 ± 0.03	1.37 ± 0.04	1.45 ± 0.03	1.10 ± 0.07
GCNN (6-D) 300 stars	1.18 ± 0.06	1.48 ± 0.14	1.57 ± 0.16	0.97 ± 0.11

Table 2: The ratio [as defined in Eq. (8)] of the inferred subhalo mass uncertainty (marginalized posterior standard deviation for a test set of stream simulations) for the encoder model variant given in the left hand column over the best performing case [GCNN helio (6-D)]. As discussed in § 3.2.1, we report posteriors only for subhalo masses $5.75 \leq \log \frac{m}{M_{\odot}} \leq 9$ and speeds $27 \leq \frac{v}{\text{km/s}} \leq 450$, where the GCNN encoder distributions pass the simulation-based calibration test. The exception to this is the GCNN (6-D) 300 stars variant, where we report for $27 \leq \frac{v}{\text{km/s}} \leq 350$. We give the mean and standard deviation of this ratio for test samples from the subhalo mass range given on the top row (and for speeds in the well-calibrated range given above). Subhalo mass constraints are always stronger given GCNN encoders than the power spectrum approach.

we use) and so we report only a central area for each encoder model where $p \geq 0.05$.

3. RESULTS

3.1. Encoder performance

Figure 6 evaluates the performance metric defined in Eq. (7) (the fractional mean absolute error on the estimated subhalo mass and velocity) given the different encoder model variants defined in Table 1. When averaged over all test samples, the power spectrum approach (§ 2.2.1) has an absolute validation MAE [Eq. (3)] of 0.157. In comparison, the GCNN model (§ 2.2.2) has an absolute validation MAE of 0.064. The lower the MAE, the more informative is the data compression. This conclusion follows since the subhalo mass and velocity completely characterize each simulation, barring the incompressible noise of random star ejections from the stream progenitor given the tidal disruption model that we use (§ 2.1).

Figure 6 breaks down this information gain by varying subhalo mass and velocity. In this section, we consider only the first two rows of each subpanel in Fig. 6, i.e., we compare the power spectrum approach with the baseline GCNN model. We consider the effect of increasing the size of the GCNN model in § 3.3 (GCNN-L) and the effects of coordinate system and incomplete observations in § 3.4 [GCNN helio (6-D), GCNN helio (3-D), GCNN (6-D) 300 stars].

The power spectrum model degrades in performance towards the edges of both the subhalo mass and velocity prior ranges. There are two physical effects that contribute to this behavior. First, when the mass is large and/or the velocity is low, the impulse of the subhalo on the stream is larger and so the stream is far more disrupted. Indeed, this disruption can occur to such an

extent that the perturbed stars are no longer confined to variations in one angular coordinate ϕ_1 , instead spreading out in all spatial dimensions. In this case, the 1D angular power spectrum becomes an increasingly poor compression of the stream phase space.

On the other hand, when the mass is small and/or the velocity is high, the impulse on the stream is smaller and so the stream is far less disrupted. In this case, indeed all encoder model variants perform increasingly poorly simply because there is less and less information in the stream about the weak interaction. A $10^5 M_{\odot}$ subhalo is considered about the lower limit of what can be detected in a single stellar stream using upcoming observatories (Drlica-Wagner et al. 2019, although the details of this forecast depend on many other properties of the stream and its perturbers). Finally, the fractional MAE for the subhalo velocity is much larger than for the subhalo mass (again for all encoder model variants). This means that the stream phase space is more sensitive to variations in subhalo mass than velocity.

The baseline GCNN model performs better than the power spectrum in all twelve subhalo mass and velocity bins in Fig. 6. This improvement reflects that the GCNN produces more informative summary statistics across the prior area, largely driven by having more information as input, i.e., the full 6-D phase space rather than only 1-D for the power spectrum. The same physical effects from varying subhalo mass and velocity, as above, determine variations in the GCNN performance.

There are also two numerical effects from the encoder model training. First, performance tends to degrade at the lower ends of the prior ranges because the absolute MAE [Eq. (3)] is smaller meaning that training and validation samples in this region contribute less to the total model loss function. This effect arises despite

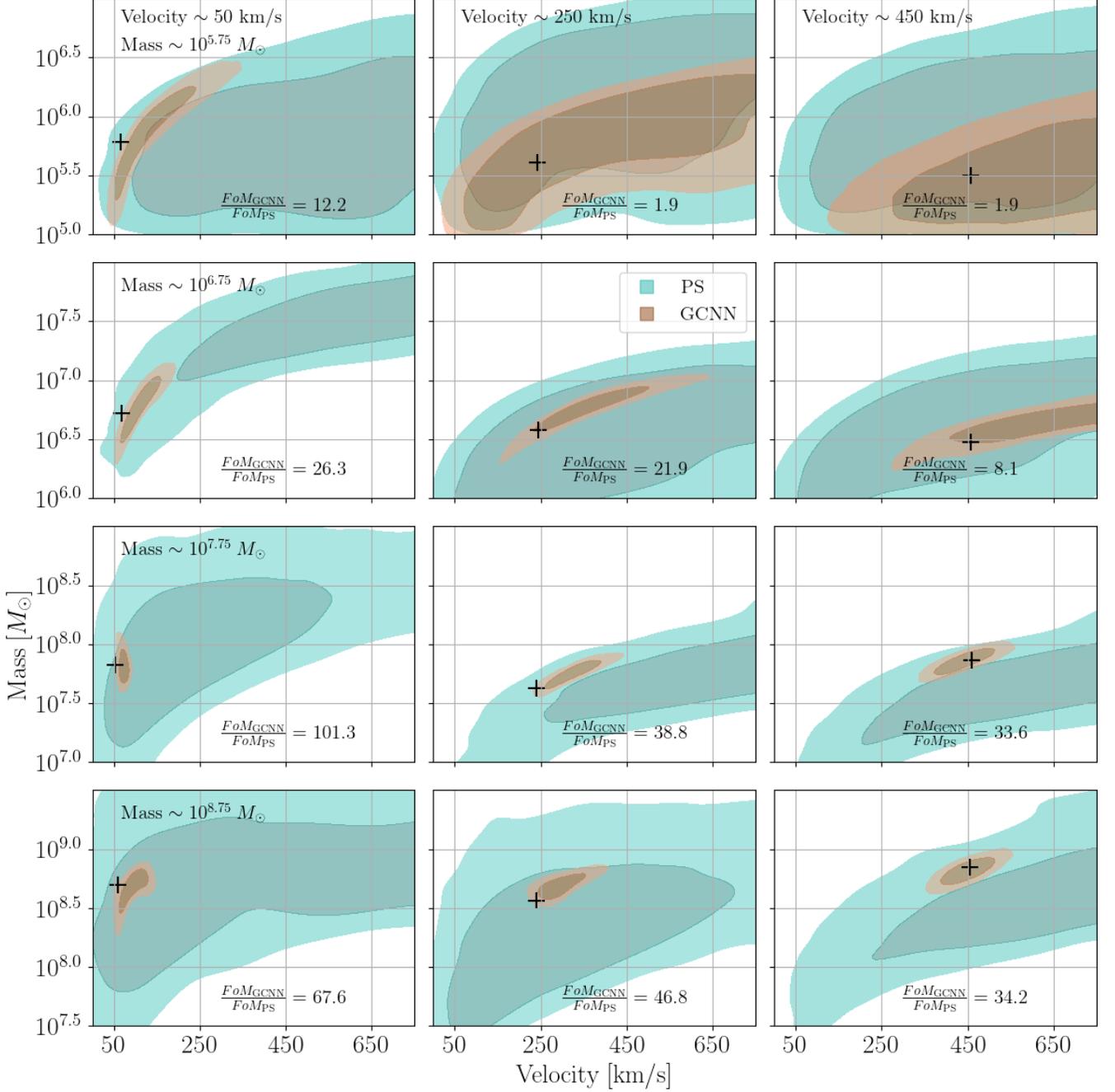


Figure 7 The posterior distribution as estimated by SBI (§ 2.3) given a test set of stream simulations with true parameter values indicated by the black crosses. The *cyan* contours show posteriors given the power spectrum encoder; *brown* contours show posteriors given the baseline GCNN encoder (see Table 1). *From top to bottom*, we vary the subhalo mass of the test simulation (see labels in left-hand column); *from left to right*, we vary the subhalo velocity of the test simulation (see labels in top row). The darker and lighter contours respectively indicate the 68% and 95% credible intervals. *In each panel*, we give the ratio of the figures of merit [Eq. (9)] given the two encoders. The figure of merit (the posterior area) is always improved (smaller) when using the GCNN encoder, with the greatest improvement at higher mass and/or lower velocity. There is some scatter in the figure of merit ratio and the position of the truth relative to the posterior owing to the noise in each test simulation. In § 3.2.1, we confirm that all these posteriors are well calibrated.

the normalization of training data described in § 2.2. Second, performance also tends to degrade at the edges of the prior ranges because the number of nearby training simulations decreases, i.e., there is half the training density at the edge than at the center of the prior. The first numerical effect could be ameliorated by using a different loss function, e.g., the fractional MAE that we use as a performance metric in this section. The second numerical effect could be ameliorated by either varying the density of training samples across the prior range or extending the prior range so that the parameter space of physical interest always remains well sampled. Having demonstrated our primary objective, that the GCNN can consistently extract more information than the power spectrum, we defer to future work these fine tunings of the model training.

3.2. Posterior distributions

Figure 7 shows examples of posterior distributions given test stream data, comparing the use of the power spectrum and baseline GCNN encoders (Table 1). As discussed below, we report posteriors only for $5.75 \leq \log \frac{m}{M_\odot} \leq 9$ and $27 \leq \frac{v}{\text{km/s}} \leq 450$, where the GCNN encoder distributions pass the simulation-based calibration test. Across this well-calibrated area of subhalo mass and velocity, the GCNN encoder leads to an improved figure of merit [Eq. (9)], which we take as the area of the posterior being reduced, i.e., tighter constraints on the subhalo parameters (up to over a factor of a hundred gain in FoM). This result is consistent with the encoder performance that we present in § 3.1, where we demonstrate that the GCNN extracts more information from the stream phase space than the power spectrum.

Both the power spectrum and GCNN encoders lead to some positive degeneracy between subhalo mass and velocity. This result is consistent with the effects of these parameters that we discuss in § 2.1 and 2.2.1, where we demonstrate that the effect of a more massive subhalo can be counteracted by increasing its velocity. For the power spectrum encoder, the velocity is typically much more poorly constrained than the mass (as a fraction of the prior range). This result is consistent with the encoder performance, where we find that the fractional MAE on the velocity is always larger than for the mass.

Although the GCNN always leads to more precise posterior contours, this is particularly true as mass increases and/or velocity decreases. In these regimes, the subhalo perturbations are strongest and there is the biggest information gain from using the full phase space. For lower mass at higher velocity, there is a degree of information saturation (although the figure of merit does improve by about a factor of two and so we have not reached sat-

uration). In this regime, the subhalo perturbations are weakest. Towards the limit of no discernible perturbation, the two encoders will tend towards the same result. The transition between these two regimes is fairly rapid; we expect this behavior from the stream simulations we show in Fig. 2 and the power spectra we show in Fig. 4. In both the full phase space and the power spectra, we see a rapid transition in behavior above and below a mass threshold of $\sim 10^7 M_\odot$ and a velocity threshold of ~ 100 km/s, indicating that this is a physical effect manifesting in the simulations and which is correctly feeding through to the SBI.

With regards to determining fundamental properties of dark matter, it is most powerful to measure the masses of subhalos as these can be directly related to the subhalo mass function. In this respect, parameters of the stream perturbation, like the relative velocity between subhalo and stream, are largely nuisances (although it is conceivable that the velocity distribution of substructure could potentially be a signature of different fundamental physics). Fig. 8 therefore shows examples of posterior distributions of the subhalo mass, having marginalized over subhalo velocity. In this section, we consider only the two encoder model variants on the far left in each subpanel, i.e., we compare the power spectrum approach with the baseline GCNN model (see § 3.3 and 3.4 for the other encoders). The GCNN encoder always leads to stronger subhalo constraints than the power spectrum, again driven by the more informative phase space compression. The relative amount of improvement decreases as the true subhalo mass decreases, reflecting the increasing saturation in the amount of information to be extracted. The improvement is fairly independent of the true subhalo velocity.

Table 2 quantifies the gain in constraining power from using the GCNN encoder by evaluating the performance metric defined in Eq. (8) (the ratio of subhalo mass uncertainty for a given encoder over the best case; again, we consider the GCNN encoder variants in § 3.3 and 3.4). We find that the improvement scales from a factor of five for subhalo masses $> 10^8 M_\odot$ to 1.5 for masses $< 10^6 M_\odot$.

3.2.1. Simulation-based calibration test

Having discussed the information gain from using the GCNN data compression (see also § 3.1), we now assess how robust is the parameter inference by the posterior calibration tests discussed in § 2.4. Fig. 9 shows the cumulative distribution function of the true parameter ranks within sets of posterior samples; or, the coverage probability for each credible limit. In Table 3, we report the KS p value assessing the consistency between these

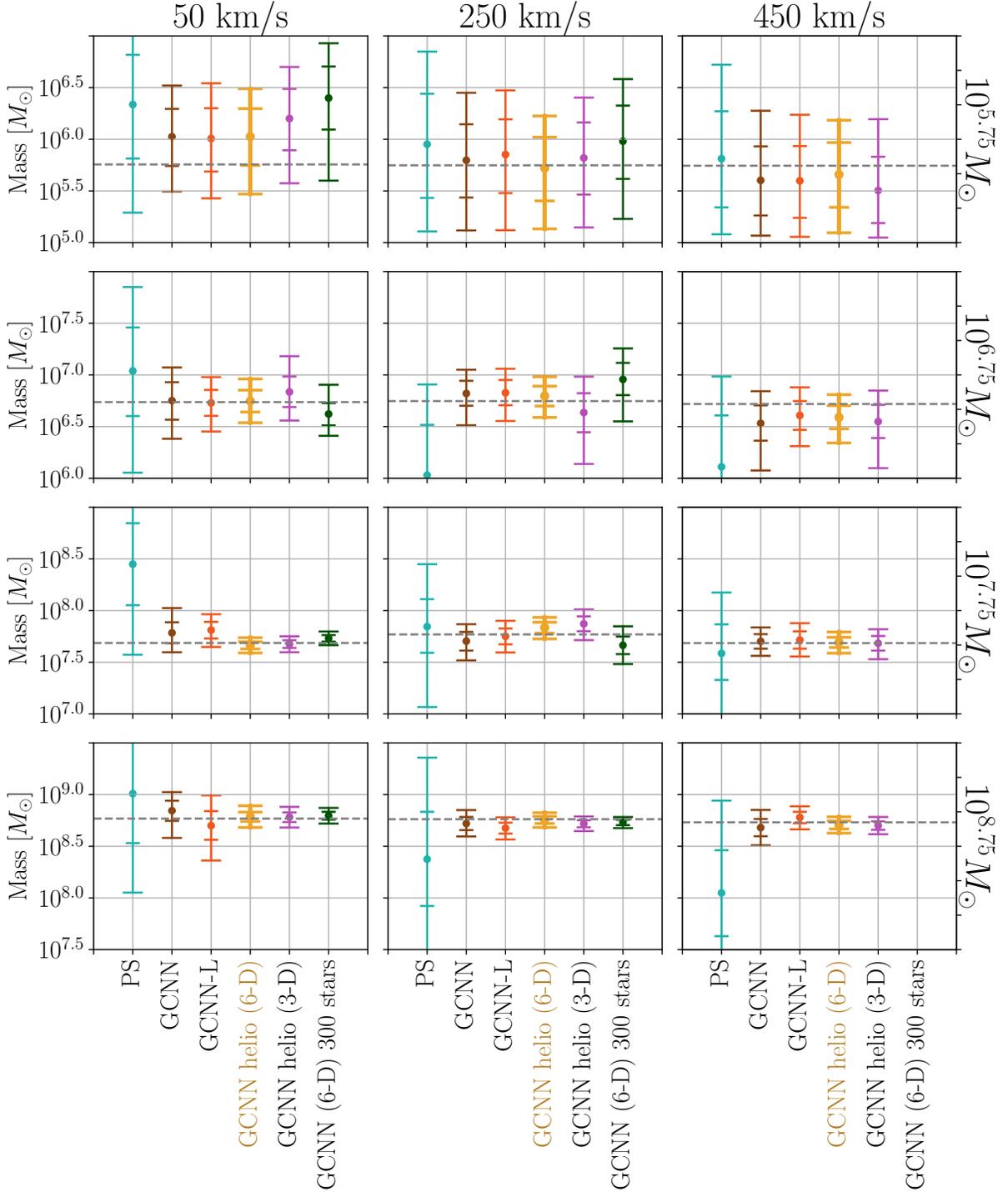


Figure 8 The marginalized posterior distribution of subhalo mass as estimated by SBI (§ 2.3) given a test set of stream simulations with true subhalo masses indicated by the grey dashed lines. *In each panel*, we show subhalo mass constraints given each encoder model variant (see labels at the bottom and Table 1). *From top to bottom*, we vary the subhalo mass of the test simulation with the true mass given on the right hand side; *from left to right*, we vary the subhalo velocity of the test simulation with the true velocity given on the top. The inner and outer ticks respectively indicate the 68% and 95% credible intervals, while the point indicates the mean. We do not report the GCNN (6-D) 300 stars results in the highest velocity bin as this model fails the calibration test in this regime. The GCNN helio (6-D) encoder model variant returns the strongest subhalo mass constraints.

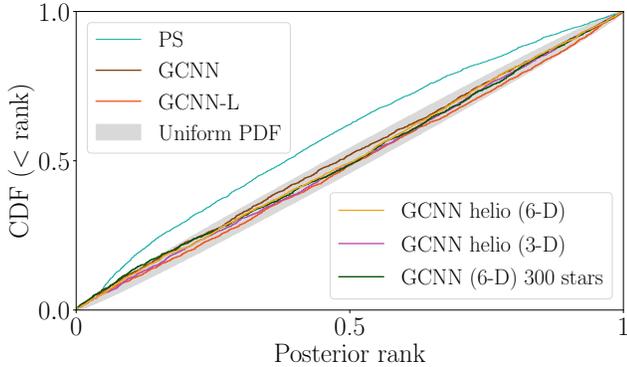


Figure 9 The cumulative distribution function (CDF) of the rank (normalized to the unit interval) of the true subhalo parameters in a set of posterior distribution samples, for a test set of stream simulations. We interpret the y -axis as the coverage probability for a given credible limit on the x -axis. As discussed in § 3.2.1, we report the posterior calibration only for subhalo masses $5.75 \leq \log \frac{m}{M_\odot} \leq 9$ and speeds $27 \leq \frac{v}{\text{km/s}} \leq 450$, where the GCNN encoder distributions pass the simulation-based calibration test. The exception to this is the GCNN (6-D) 300 stars variant, where we report for $27 \leq \frac{v}{\text{km/s}} \leq 350$. In other words, for these ranges only, the test GCNN posteriors are accurately calibrated and the rank distributions are consistent with a uniform distribution, i.e., the gray area which indicates the 99% c.l. of a uniform probability distribution function. We find that the GCNN posteriors are not only more constraining but also more accurate than when using the power spectrum.

rank CDFs and the expected uniform CDF (gray area in Fig. 9). Specifically, by searching through different sub-areas of the prior, we identify a common area where the GCNN encoders are well-calibrated, i.e., $p \geq 0.05$. This well-calibrated area is for true subhalo masses $5.75 \leq \log \frac{m}{M_\odot} \leq 9$ and true speeds $27 \leq \frac{v}{\text{km/s}} \leq 450$, except for the GCNN (6-D) 300 stars variant, where the area is for $27 \leq \frac{v}{\text{km/s}} \leq 350$ (see § 3.4). It is for these areas only that we report results in Figs. 7, 8 and 9 and Tables 2 and 3.

We therefore find that posteriors towards the edge of the prior outside these areas are not well calibrated and therefore unreliable (i.e., when the true parameters lie outside this area). We attribute this performance to a combination of two effects. First, the edge of the prior is where the encoder performance degrades (as discussed in § 3.1). Second, there is a numerical effect where, since the posterior is cut off by the uniform prior, a discontinuity has arisen that breaks a condition for the calibration test. We discuss in § 3.1 how to mitigate the encoder

Encoder model	Neg. log prob.	KS p value
PS	-1.30 ± 0.90	$\ll 0.05$
GCNN	-4.70 ± 1.61	0.08
GCNN-L	-4.90 ± 1.62	0.06
GCNN helio (6-D)	-5.71 ± 1.63	0.27
GCNN helio (3-D)	-5.08 ± 1.97	0.52
GCNN (6-D) 300 stars	-5.40 ± 2.07	0.08

Table 3: Performance metrics (as defined in § 2.4) on the accuracy of posterior estimation given each encoder model variant (*from top to bottom*). As discussed in § 3.2.1, we report the posterior calibration only for subhalo masses $5.75 \leq \log \frac{m}{M_\odot} \leq 9$ and speeds $27 \leq \frac{v}{\text{km/s}} \leq 450$, where the GCNN encoder distributions pass the simulation-based calibration test, i.e., $p \geq 0.05$. The exception to this is the GCNN (6-D) 300 stars variant, where we report for $27 \leq \frac{v}{\text{km/s}} \leq 350$. For the negative log probability, we give the mean and standard deviation for test samples from this well-calibrated range only. For the power spectrum approach, we report only that $p \ll 0.05$ as it is significantly in the tail where numerical effects manifest. Posteriors are always more accurate given GCNN encoders than the power spectrum approach.

edge effect, while the second effect could be mitigated by using a prior with continuous support (e.g., a Gaussian) or instead learning the likelihood rather than the posterior. In any case, we find that our GCNN posteriors are well calibrated throughout the parameter space of most interest and we defer to future work a more thorough investigation of these edge effects. In particular, we find that the GCNN posteriors are always better calibrated than the power spectrum approach, meaning that using the GCNN is both more informative and more accurate. In Table 3, we also compare the negative log probabilities of the well-calibrated test set finding that the GCNN encoder model again performs better than the power spectrum encoder (lower negative probability).

3.3. Performance scaling with number of training simulations

We consider here the performance of the first of the four GCNN encoder model variants that we introduce in § 2.2 and Table 1: the GCNN-L encoder. For the GCNN-L encoder, relative to the baseline GCNN encoder discussed in § 3.1 and 3.2, we double the number of neural network model hyper-parameters from 7×10^5 to 1.4×10^6 and, in turn, double the number of training simulations from 40,000 to 80,000. Increasing the size (two times more parameters) and depth (one additional layer) of a neural network can increase its expressivity,

i.e., its ability to express complex transformations from input to output, and so we investigate if the data compression can be more informative as a consequence.

In Fig. 6, we see little systematic change in the accuracy of the encoder model from GCNN to GCNN-L; in eight out of twelve test subhalo mass and velocity bins, we see a reduced fractional error in parameter estimation. In any case, the size of the gain or loss is always much less than the gain relative to the power spectrum approach. In Fig. 8 and Table 2, we assess the performance of the parameter inference, specifically for subhalo mass constraints. There is a discernible, though marginal, trend with up to a 12% reduction in parameter uncertainties. We do not over-interpret this result, though note that this trend agrees with the change in encoder accuracy for subhalo masses $> 10^6 M_\odot$. Finally, in Fig. 9 and Table 3, we report that the posteriors are well calibrated (KS $p = 0.06$) for the area identified in § 3.2.1. The negative log probability is slightly lower than for the baseline GCNN indicating better performance. Overall, we find no strong scaling with the size of the encoder network in performance throughout the inference pipeline. For the further encoder model variants, we retain the larger architecture of the GCNN-L model.

3.4. Performance in the limit of incomplete observations

We consider here the performance of the final three GCNN encoder model variants that we introduce in § 2.2 and Table 1: the GCNN helio (6-D), GCNN helio (3-D) and GCNN (6-D) 300 stars encoders. For the GCNN helio (6-D) encoder, relative to the baseline GCNN and the GCNN-L encoders, we change the input coordinate system of the stream phase space. We change the coordinates from galactocentric spherical (i.e., the radius, polar angle and azimuthal angle relative to the center of the Galaxy and their velocity equivalents) to a set of heliocentric coordinates (RA, DEC, heliocentric distance, v_{RA} , v_{DEC} , v_{LOS}). In doing this, we simulate more closely to real data, where we observe in angular coordinates on the sky, measure line of sight velocities from stellar spectra, infer heliocentric distances from stellar parallaxes and measure angular velocities indirectly through proper motions. We convert proper motions into angular velocities in order to decorrelate the position and velocity vectors as input features to the encoder. We may expect some sensitivity of performance to the change of coordinate system, despite no change in input information content, as discussed in § 4.4. We quantitatively investigate this effect here.

In Fig. 6, we see that the GCNN helio (6-D) encoder performs more accurate parameter estimation than both the GCNN and GCNN-L models in ten out of twelve test subhalo mass and velocity bins. This result feeds through to stronger subhalo mass constraints after parameter inference (Fig. 8 and Table 2). We find that the improvement (relative to GCNN-L) scales from a factor of 1.15 for subhalo masses $< 10^6 M_\odot$ up to a factor of 1.83 for masses $> 10^8 M_\odot$. The fact that the gain is largest for the most disruptive subhalo interactions suggests that it is the projection of clear features in the stream like gaps and spurs that drives this result. Indeed, Fig. 2 demonstrates that features like spurs are more discernible in certain coordinates (for the simulations in the figure, ϕ_2 and heliocentric distance rather than ϕ_1 , RA and DEC). It follows that changing from galactocentric to heliocentric coordinates will change the prominence of certain stream features. Doing this transformation in turn changes the feature space which is input to the encoder model, thus making an informative data compression an easier or harder task. We discuss in § 4.4 other examples where the feature embedding space affects model performance. Fig. 9 and Table 3 again confirm that the posteriors are well-calibrated in the area we have identified ($p = 0.27$, see § 3.2.1). The negative log probability is even lower than for the GCNN and GCNN-L models, indicating further improved performance.

With the GCNN helio (3-D) and GCNN (6-D) 300 stars encoders, we probe how the performance of the inference pipeline scales in the limit of a reduced set of observations. For the other GCNN encoders, we allow access to the full 6-D phase space of 3000 GD-1 stars. This setting is conceivable through a future combination of deep photometric, spectroscopic and astrometric observations, but is beyond our current capability. For the GCNN helio (3-D) encoder, we input 3000 stars but only with a 3-D phase space, which we take as RA, DEC and v_{LOS} . This setting roughly mocks up a future deep photometric survey with follow-up spectroscopy. The number of data is reduced from 18,000 to 9000. For the GCNN (6-D) 300 stars encoder, we input only 300 stars but with a full 6-D phase space. This setting roughly mocks up observations we will have in the near future, where we will have complete data for a subset of the GD-1 stars. The number of data here is reduced further to 1800. We do not set out to perform detailed forecasts for particular observational settings, but rather we seek to understand the trade-off in performance from having more stars with incomplete data or fewer stars with complete data. We leave for future work a detailed forecast using our new data analysis method.

In Fig. 6, we see that the GCNN encoders with reduced input perform less precise parameter estimation than the GCNN helio (6-D) model in all but one of the twelve test subhalo mass and velocity bins. This result is expected as supplying less information to the model will lead to a less informative summary statistic. However, both the GCNN helio (3-D) model with 9000 input data and the GCNN (6-D) 300 stars model with 1800 input data perform better at parameter estimation than the power spectrum approach (with 3000 input data before any compression, i.e., the ϕ_1 coordinates of 3000 stars). This result underscores the combined benefit of a more informative compression scheme (GCNN over PS) and gathering multi-dimensional data (more than 1-D). Between the two reduced input models, we find that the GCNN helio (3-D) performs more accurate parameter estimation in nine out of twelve parameter bins.

In Fig. 8 and Table 2, we find that the scaling of the constraining power of the parameter inference is consistent with the scaling of the information extracted in the data compression. The 1σ uncertainty on subhalo mass is worse than when using all the input data across the prior range that we consider (apart from one of eight test subhalo mass bins), but always better than when using the power spectrum pre-compression, even if there are fewer input data (1800 for the GCNN (6-D) 300 stars model, 3000 for the PS model). The degradation from using incomplete observations relative to the GCNN helio (6-D) model is strongest for intermediate subhalo masses $\sim 10^7 M_\odot$. We see the effects of information saturation at lower masses, where all models perform increasingly similarly as the subhalo perturbation weakens. For the most massive subhalos $\sim 10^8 M_\odot$, we even see a marginal gain in constraining power from inputting fewer data (the GCNN (6-D) 300 stars model). This result tells us that, when the perturbation is very strong, there are a few stars that dominate the signal. Therefore, adding more unperturbed stars as input can make this signature less clear in the feature embedding space.

Finally, in Fig. 9 and Table 3, for the GCNN helio (3-D) variant, we report well-calibrated posteriors ($p = 0.52$) in the range for subhalo masses $5.75 \leq \log \frac{m}{M_\odot} \leq 9$ and speeds $27 \leq \frac{v}{\text{km/s}} \leq 450$. For the GCNN (6-D) 300 stars variant, we report a slightly reduced well-calibrated range ($p = 0.08$) for $27 \leq \frac{v}{\text{km/s}} \leq 350$, suggesting that the reduced amount of data can make the SBI task harder. The reduced amount of data in both variants leads to worse negative log probabilities, albeit still better than the power spectrum encoder.

4. DISCUSSION

4.1. Improved model parameter estimators

In § 3, we present results demonstrating that the GCNN encoders return more accurate determinations of the perturbing subhalo mass and speed than when using the power spectrum. This information gain persists, irrespective of the size of the GCNN model, the coordinate system of the data or the subset of the data that is input to the GCNN.

We discuss here what drives this gain in information being extracted from the stream data. First, the GCNN ingests a richer dataset since it includes the full 6-D phase space. On the other hand, the power spectrum model uses only the perturbations along the ϕ_1 coordinate, which is a lossy compression. Second, the dataset is inherently a particle-like system. Thus, there are symmetries underlying the system like equivariance of which the GCNN takes full advantage due to the bidirectional message passing architecture. This symmetry is not used in the power spectrum approach since the inputs are compressed by binning. Further, the binning preprocessing for the power spectrum also smoothes out potential signals in the data. Lastly, the GCNN model makes fewer assumptions about the data. Current empirical evidence in the trends of deep learning has seen an increase in abstraction of the features in the data. This means that the less hard coded is the feature extraction, the better the model can perform. The improved flexibility of the GCNN to adapt to a given problem often yields better results at the cost of computational expense and some loss of interpretability, which we observe in our experiments.

4.2. Breaking model parameter degeneracies

As discussed in § 2.1, there is a physical degeneracy in the stream simulations between increasing subhalo mass and speed. As the subhalo becomes more massive, there is more stream perturbation, but this effect can be counteracted by increasing the relative velocity between subhalo and stream so that the impulse is reduced. In § 3.2, we see this degeneracy manifest in the posterior distributions whether using the power spectrum or GCNN encoders. A striking result is that for the power spectrum model, the subhalo velocity is largely unconstrained and independent of the true parameters. This lack of sensitivity is consistent with the smaller dynamic range in power spectra that we see in Fig. 4 as we vary subhalo velocity compared to mass. In Fig. 7, we see that using the GCNN encoder can significantly break this degeneracy in a mass- and velocity-dependent way. As the mass decreases and/or the velocity decreases, the GCNN model returns much stronger constraints on the subhalo velocity, thus (partially) breaking the degeneracy.

acy with subhalo mass. We attribute this behavior to the fact that the GCNN has additional input stream velocity information, which is more sensitive to the subhalo velocity than the power spectrum. This result emphasises the importance of providing stream star velocities in constraining the mass of perturbing subhalos.

4.3. Saturation in GCNN performance

In § 3.3, we find that increasing the size and depth of the GCNN model leads to marginal gain in parameter estimation and the precision and accuracy of posterior estimation. This result indicates that we can likely achieve similar performance with fewer training simulations than previously anticipated. This can be attributed to a known property of over-smoothing (Zhang et al. 2024a) of the GCNN, where large graph models [number of hyperparameters $\sim \mathcal{O}(10^6)$] have plateauing performance as the number of hyperparameters increases with more layers (Liu et al. 2024). Thus, if increasing model capacity cannot yield substantial improvements, it is unsurprising that increasing the amount of training data also does not. For this problem, we suggest that $\mathcal{O}(10^3 - 10^4)$ training simulations thus suffices while remaining computationally efficient. Further, in future work, we will investigate more efficient use of training data through sequential SBI methods that iteratively generate training data in a more optimal fashion.

4.4. Effect of the stream coordinate system

In § 3.4, we find up to a factor of 1.83 improvement in the subhalo mass constraint from changing from galactocentric to heliocentric coordinates. This is initially surprising as an invertible coordinate transformation does not directly yield any additional information about the dynamics of the system, yet performance improves. This phenomenon is well studied in the field of representation learning (Goodfellow et al. 2016). The choice of embedding or representation of data is central in many deep learning models. E.g., positional encoding, which is an invertible map to represent position data using sinusoids, when applied to physics inspired neural networks (Huang et al. 2021; Zhang et al. 2024b), improves the learning outcomes. In natural language processing, tokenization using WORD2VEC (Mikolov et al. 2013) or word embeddings in transformer models (Vaswani et al. 2017) adopt a “basis change” to boost model performance. Although it appears surprising at first that a coordinate change can yield noticeable model improvements, this is a norm rather than an exception. The common argument is that the choice of coordinates can disentangle features. It is arguable that stream perturbations are more aligned along a particular axis in the

heliocentric coordinates than the galactocentric coordinates. Indeed, we see in Fig. 2 that, for these example simulations, the spur is more or less discernible depending on the coordinates in which it is projected. For real data, it will therefore be critical to assess the sensitivity of results to the chosen coordinates.

4.5. Future simulation improvements

While this work demonstrates a proof-of-concept method for characterizing dark matter subhalo perturbations in stellar streams, future iterations will incorporate more sophisticated simulation methods with a view to a robust data analysis. First, we will vary additional simulation parameters such as the impact parameter, the angle of approach of the subhalo to the stream, the time of closest approach and the subhalo density profile and size. It is important to consider these parameters in addition as their values are unknown for any given stream and so we must constrain them using data. Second, we will incorporate multiple subhalo interactions. The host (MW) halo contains a distribution of subhalos (of varying mass, velocity, etc.) and typically tens or hundreds of these will appreciably interact with a stream. Thus, folding in a realistic subhalo distribution and multiple interactions will be critical for a future comparison to data. We also anticipate investigating more realistic particle-spray simulation approaches (Carlberg 2018; Malhan et al. 2019; Gialluca et al. 2021; Qian et al. 2022), modeling time-dependent variations in the background potential (Buist & Helmi 2015; Koppelman & Helmi 2021; Brooks et al. 2024) and including baryonic perturbers like giant molecular clouds (Banik et al. 2021a). Finally, for a reliable comparison to data, we must model the properties of the survey(s) when generating the input simulations. E.g., we can model measurement uncertainties by sampling from a distribution that resembles instrument uncertainties and we can further account for completeness and masking by post-processing the simulations before input to the encoder.

5. CONCLUSIONS

We have demonstrated that the combination of a graph convolutional neural network compression and simulation-based inference can improve constraints on the mass of a perturbing subhalo by factors of three to eleven compared to the current state-of-the-art density power spectrum analysis of a GD-1-like stellar stream. When varying the mass and velocity of the subhalo in our simulations, we find that the GCNN encoder consistently improves estimation of these parameters relative to the power spectrum. The ability of the GCNN to

exploit information in the kinematics of the stream improves inference of the subhalo velocity, which in turn breaks the degeneracy with subhalo mass. We find that posterior distributions given the GCNN compression are significantly better calibrated, meaning that this method is simultaneously more precise and more accurate.

We perform three experiments with the GCNN encoder. First, we find that there is no strong scaling in the performance of the method as we increase the complexity of the neural network and increase the number of training simulations. We conclude that $\mathcal{O}(10^3 - 10^4)$ stream simulations are sufficient in training compression networks for future analyses, thus informing the computational expense of the problem. Second, we find some sensitivity in our results to the coordinate system of the stream as input to the pipeline. This result is consistent with the deep learning literature and we postulate that identifying a projection of the stream that accentuates the features associated with a subhalo interaction can enhance the ability of the encoder to extract information.

Finally, while we do not perform detailed forecasts for the power of this method applied to upcoming data with a full treatment of errors and systematic effects in this proof-of-principle study, we approximate two observational settings achievable in the near future. We compare a setting with 3000 stars with sky coordinates and line-of-sight velocities (a combination of deep photometry and spectroscopy but without associated astrometry) with a setting in which the dataset consists of 300 stars with full 6-D phase space data (photometry, spectroscopy and astrometry). We find an improvement in subhalo mass constraints by factors of two to eleven for masses from $(10^6 - 10^9) M_\odot$ compared to the constraints when using the power spectrum approach. The improvement is about the same for the two set-

tings suggesting that collecting exhaustive information on a few hundred stream stars can be about as powerful as a deeper survey of thousands of members with only RA, DEC and line-of-sight velocities. We conclude that the graph neural network method we have introduced here will be powerful in maximizing the information return about the MW subhalo population from upcoming deep- and wide-field photometric and spectroscopic surveys (e.g., *Rubin*, *Euclid*, DESI, *Roman*). In future work, we will develop this method further to incorporate more sophisticated simulation approaches and the modeling of survey properties with a view to future data analysis.

6. ACKNOWLEDGMENTS

The authors thank Jo Bovy for insightful discussions and support with using GALPY. KKR is supported by an Ernest Rutherford Fellowship from the UKRI Science and Technology Facilities Council (grant number ST/Z510191/1). The Dunlap Institute is funded through an endowment established by the David Dunlap family and the University of Toronto. TSL acknowledges financial support from the Natural Sciences and Engineering Research Council (NSERC) of Canada through grant RGPIN-2022-04794. RH acknowledges support from the Natural Sciences and Engineering Research Council of Canada Discovery Grant Program and the Connaught Fund. The authors at the University of Toronto acknowledge that the land on which the University of Toronto is built is the traditional territory of the Haudenosaunee and, most recently, the territory of the Mississaugas of the New Credit First Nation. They are grateful to have the opportunity to work in the community on this territory. PXM thanks the friends at CERN for the generous vehicular support and fantastic company at the early stages of this work.

APPENDIX

A. NEURAL ARCHITECTURE OPTIMIZATION

We optimize the architecture of the power spectrum and baseline GCNN compression networks using the `scikit-optimize` Bayesian optimization (Snoek et al. 2012) package⁴. We search for the architecture that minimizes the validation loss [Eq. (3)]. Table 4 shows the architecture parameters that we vary and the values that we consider. For this search, we train with only 1000 simulations and the ADAM optimizer for 200 epochs. The bold values in Table 4 are the optimal choice given this procedure.

⁴ <https://scikit-optimize.github.io/stable>.

Power spectrum encoder	
Number of dense layers #1	1, 2, 3
Number of dense layers #2	1, 2, 3
Number of dense layers #3	1, 2, 3
Dense layer #1 width	256, 512 , 1024, 2048
Dense layer #2 width	64, 128, 256
Dense layer #3 width	16 , 32, 64
Activation function	Leaky ReLU , ReLU, Sigmoid, Tanh
Graph convolutional neural network (GCNN) encoder	
Number of graph convolutional layers	1, 2, 3 , 4, 5
Graph convolutional layer width	6, 32, 64, 128 , 512
Number of dense layers	1, 2 , 3, 4, 5
Dense layer width	32, 64, 128, 512 , 1024
Activation function	Leaky ReLU, ReLU , Sigmoid, Tanh

Table 4: Results of the encoder neural architecture search using Bayesian optimization. The *top* part shows results for the power spectrum encoder, while the *bottom* part shows results for the baseline GCNN encoder. Each of the parameters on the *left* is optimized while searching across the values on the *right*. The optimal architectures are indicated in *bold*.

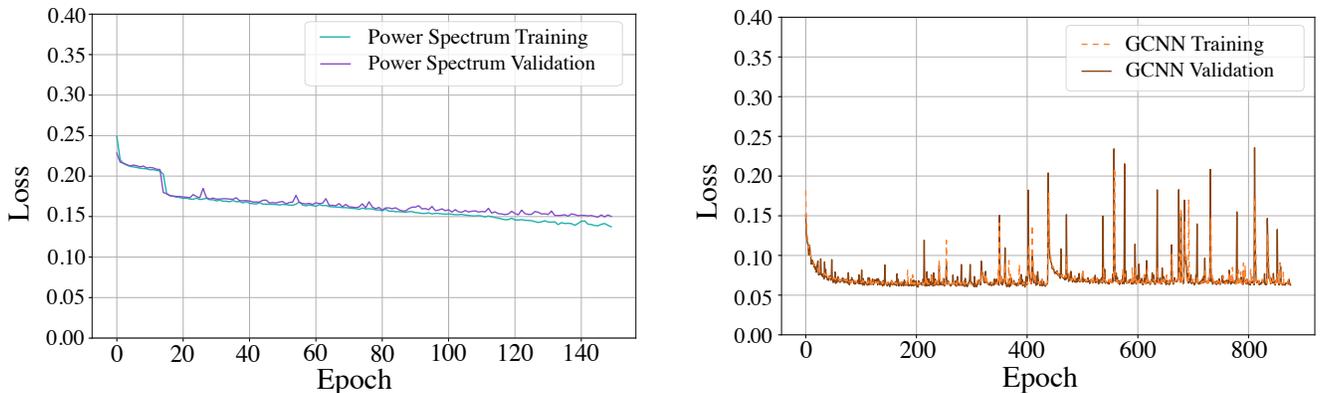


Figure 10 The training and validation loss curves as a function of epoch for the power spectrum (*left*) and baseline GCNN (*right*) encoders. The GCNN training returns a lower final loss, while the training and validation losses do not diverge indicating that overfitting has been avoided.

B. TRAINING AND VALIDATION LOSS CURVES

Figure 10 shows the training and validation losses [Eq. (3)] as a function of epoch for the power spectrum and baseline GCNN encoders. The GCNN encoder finishes with a lower loss indicating a better fit, while the training and validation curves do not diverge indicating that overfitting to the training set has not occurred.

C. VISUALIZATION OF A GRAPH CONVOLUTIONAL LAYER

The GCNN approach to analyzing stellar streams that we have introduced, although more powerful, is necessarily more black box than measuring the density power spectrum. Nonetheless, we can open up the black box in order to increase interpretability and to inform better network construction. Fig. 11 visualizes the final graph convolutional layer in the baseline GCNN model. This is the final layer before the aggregation across nodes. We perform an experiment where we fix the GCNN model and then vary the input to the network by changing only the mass of the perturbing subhalo (random seed also fixed). In this way, we see which features are activated by our parameter of interest. The latent space remains sparse, i.e., only a few features are significantly activated. As the mass increases,

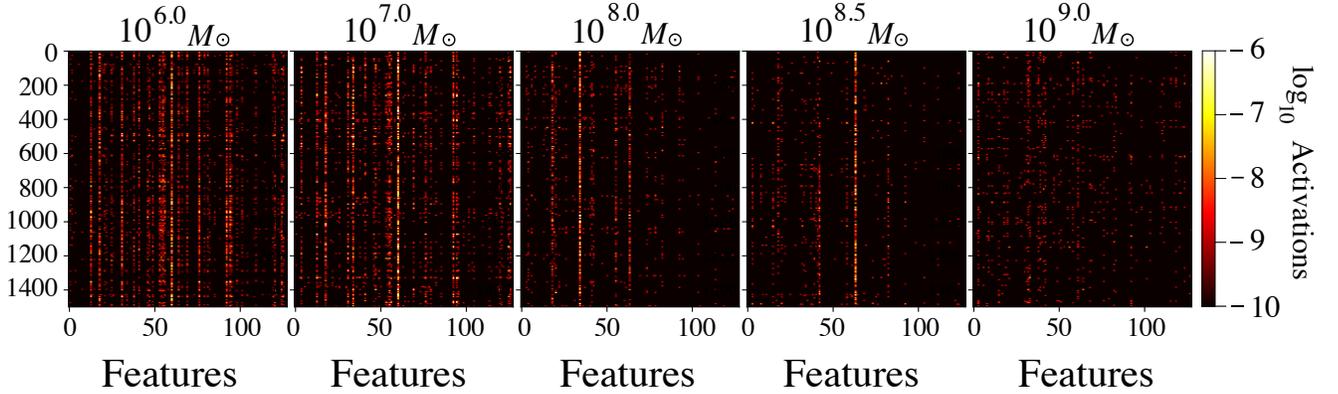


Figure 11 The activation in the final graph convolutional layer of the baseline GCNN model for one arm of the stellar stream (dimensions of 1500 stars (nodes) by 128 features). *From left to right*, as indicated at the top, we vary the mass of the perturbing subhalo in the input simulation but otherwise fix all parameters including the random seed and the velocity of the subhalo relative to the stream to 440 km/s. As the mass varies, different features are activated in the network.

an increasingly small number of nodes (stars) and features are significantly activated supporting the success of the GCNN (6-D) 300 stars variant in this regime.

REFERENCES

- Aghanim, N., Akrami, Y., Ashdown, M., et al. 2020, *Astronomy & Astrophysics*, 641, A6, doi: [10.1051/0004-6361/201833910](https://doi.org/10.1051/0004-6361/201833910)
- Alsing, J., Charnock, T., Feeney, S., & Wandelt, B. 2019, *Monthly Notices of the Royal Astronomical Society*, doi: [10.1093/mnras/stz1960](https://doi.org/10.1093/mnras/stz1960)
- Alvey, J., Gerdes, M., & Weniger, C. 2023, *Monthly Notices of the Royal Astronomical Society*, 525, 3662–3681, doi: [10.1093/mnras/stad2458](https://doi.org/10.1093/mnras/stad2458)
- Banik, N., Bovy, J., Bertone, G., Erkal, D., & de Boer, T. 2021a, *Journal of Cosmology and Astroparticle Physics*, 2021, 043, doi: [10.1088/1475-7516/2021/10/043](https://doi.org/10.1088/1475-7516/2021/10/043)
- Banik, N., Bovy, J., Bertone, G., Erkal, D., & de Boer, T. J. L. 2021b, *Monthly Notices of the Royal Astronomical Society*, 502, 2364–2380, doi: [10.1093/mnras/stab210](https://doi.org/10.1093/mnras/stab210)
- Blanchard, A., Camera, S., Carbone, C., et al. 2020, *Astronomy & Astrophysics*, 642, A191, doi: [10.1051/0004-6361/202038071](https://doi.org/10.1051/0004-6361/202038071)
- Boera, E., Becker, G. D., Bolton, J. S., & Nasir, F. 2019, *The Astrophysical Journal*, 872, 101, doi: [10.3847/1538-4357/aafec4](https://doi.org/10.3847/1538-4357/aafec4)
- Bonaca, A., Hogg, D. W., Price-Whelan, A. M., & Conroy, C. 2019, *The Astrophysical Journal*, 880, 38, doi: [10.3847/1538-4357/ab2873](https://doi.org/10.3847/1538-4357/ab2873)
- Bonaca, A., & Price-Whelan, A. M. 2024, *Stellar Streams in the Gaia Era*, arXiv, doi: [10.48550/ARXIV.2405.19410](https://doi.org/10.48550/ARXIV.2405.19410)
- Bovy, J. 2014, *The Astrophysical Journal*, 795, 95, doi: [10.1088/0004-637x/795/1/95](https://doi.org/10.1088/0004-637x/795/1/95)
- . 2015, *The Astrophysical Journal Supplement Series*, 216, 29, doi: [10.1088/0067-0049/216/2/29](https://doi.org/10.1088/0067-0049/216/2/29)
- . 2016, *Physical Review Letters*, 116, doi: [10.1103/physrevlett.116.121301](https://doi.org/10.1103/physrevlett.116.121301)
- Bovy, J., Erkal, D., & Sanders, J. L. 2017, *MNRAS*, 466, 628, doi: [10.1093/mnras/stw3067](https://doi.org/10.1093/mnras/stw3067)
- Brooks, R. A. N., Sanders, J. L., Lilleengen, S., Petersen, M. S., & Pontzen, A. 2024, *Monthly Notices of the Royal Astronomical Society*, 532, 2657–2673, doi: [10.1093/mnras/stae1565](https://doi.org/10.1093/mnras/stae1565)
- Brown, A. G. A., Vallenari, A., Prusti, T., et al. 2021, *Astronomy & Astrophysics*, 649, A1, doi: [10.1051/0004-6361/202039657](https://doi.org/10.1051/0004-6361/202039657)
- Buist, H. J. T., & Helmi, A. 2015, *Astronomy & Astrophysics*, 584, A120, doi: [10.1051/0004-6361/201526203](https://doi.org/10.1051/0004-6361/201526203)
- Butsky, I., Macciò, A. V., Dutton, A. A., et al. 2016, *Monthly Notices of the Royal Astronomical Society*, 462, 663–680, doi: [10.1093/mnras/stw1688](https://doi.org/10.1093/mnras/stw1688)
- Caldeira, J., Wu, W., Nord, B., et al. 2019, *Astronomy and Computing*, 28, 100307, doi: [10.1016/j.ascom.2019.100307](https://doi.org/10.1016/j.ascom.2019.100307)
- Cardone, V. F., Piedipalumbo, E., & Tortora, C. 2005, *Monthly Notices of the Royal Astronomical Society*, 358, 1325–1336, doi: [10.1111/j.1365-2966.2005.08834.x](https://doi.org/10.1111/j.1365-2966.2005.08834.x)

- Carlberg, R. G. 2018, *The Astrophysical Journal*, 861, 69, doi: [10.3847/1538-4357/aac88a](https://doi.org/10.3847/1538-4357/aac88a)
- Carlberg, R. G., Grillmair, C. J., & Hetherington, N. 2012, *The Astrophysical Journal*, 760, 75, doi: [10.1088/0004-637x/760/1/75](https://doi.org/10.1088/0004-637x/760/1/75)
- Casella, G., Robert, C. P., & Wells, M. T. 2004, *Lecture notes-monograph series*, 342
- Chen, H., Speagle, J., & Rogers, K. K. 2023, in *37th Conference on Neural Information Processing Systems*. <https://arxiv.org/abs/2311.16238>
- Chen, Y., Valluri, M., Gnedin, O. Y., & Ash, N. 2024, *arXiv e-prints*, arXiv:2408.01496, doi: [10.48550/arXiv.2408.01496](https://doi.org/10.48550/arXiv.2408.01496)
- Cole, A., Miller, B. K., Witte, S. J., et al. 2022, *Journal of Cosmology and Astroparticle Physics*, 2022, 004, doi: [10.1088/1475-7516/2022/09/004](https://doi.org/10.1088/1475-7516/2022/09/004)
- Cook, S. R., Gelman, A., & Rubin, D. B. 2006, *Journal of Computational and Graphical Statistics*, 15, 675–692, doi: [10.1198/106186006x136976](https://doi.org/10.1198/106186006x136976)
- Cooley, J., et al. 2022. <https://arxiv.org/abs/2209.07426>
- Cooper, A. P., Kogosov, S. E., Allende Prieto, C., et al. 2023, *The Astrophysical Journal*, 947, 37, doi: [10.3847/1538-4357/acb3c0](https://doi.org/10.3847/1538-4357/acb3c0)
- Cranmer, M., Sanchez-Gonzalez, A., Battaglia, P., et al. 2020, *Discovering Symbolic Models from Deep Learning with Inductive Biases*, arXiv, doi: [10.48550/ARXIV.2006.11287](https://doi.org/10.48550/ARXIV.2006.11287)
- Croft, R. A. C., Weinberg, D. H., Bolte, M., et al. 2002, *The Astrophysical Journal*, 581, 20–52, doi: [10.1086/344099](https://doi.org/10.1086/344099)
- Croft, R. A. C., Weinberg, D. H., Pettini, M., Hernquist, L., & Katz, N. 1999, *The Astrophysical Journal*, 520, 1–23, doi: [10.1086/307438](https://doi.org/10.1086/307438)
- Deistler, M., Goncalves, P. J., & Macke, J. H. 2022, *Truncated proposals for scalable and hassle-free simulation-based inference*, arXiv, doi: [10.48550/ARXIV.2210.04815](https://doi.org/10.48550/ARXIV.2210.04815)
- Diemand, J., Kuhlen, M., & Madau, P. 2007, *The Astrophysical Journal*, 667, 859–877, doi: [10.1086/520573](https://doi.org/10.1086/520573)
- Diemand, J., Kuhlen, M., Madau, P., et al. 2008, *Nature*, 454, 735–738, doi: [10.1038/nature07153](https://doi.org/10.1038/nature07153)
- Diggle, P. J., & Gratton, R. J. 1984, *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 46, 193–212, doi: [10.1111/j.2517-6161.1984.tb01290.x](https://doi.org/10.1111/j.2517-6161.1984.tb01290.x)
- Drlica-Wagner, A., Mao, Y.-Y., Adhikari, S., et al. 2019, *arXiv e-prints*, arXiv:1902.01055, doi: [10.48550/arXiv.1902.01055](https://doi.org/10.48550/arXiv.1902.01055)
- Duda, R. O., & Hart, P. E. 1974, in *A Wiley-Interscience publication*. <https://api.semanticscholar.org/CorpusID:12946615>
- Eifler, T., Miyatake, H., Krause, E., et al. 2021, *Monthly Notices of the Royal Astronomical Society*, 507, 1746–1761, doi: [10.1093/mnras/stab1762](https://doi.org/10.1093/mnras/stab1762)
- Erkal, D., & Belokurov, V. 2015, *Monthly Notices of the Royal Astronomical Society*, 450, 1136–1149, doi: [10.1093/mnras/stv655](https://doi.org/10.1093/mnras/stv655)
- Erkal, D., Belokurov, V., Bovy, J., & Sanders, J. L. 2016, *Monthly Notices of the Royal Astronomical Society*, 463, 102–119, doi: [10.1093/mnras/stw1957](https://doi.org/10.1093/mnras/stw1957)
- Erkal, D., Kogosov, S. E., & Belokurov, V. 2017, *MNRAS*, 470, 60, doi: [10.1093/mnras/stx1208](https://doi.org/10.1093/mnras/stx1208)
- Erkal, D., Belokurov, V., Laporte, C. F. P., et al. 2019, *Monthly Notices of the Royal Astronomical Society*, 487, 2685–2700, doi: [10.1093/mnras/stz1371](https://doi.org/10.1093/mnras/stz1371)
- Fardal, M. A., Huang, S., & Weinberg, M. D. 2015, *Monthly Notices of the Royal Astronomical Society*, 452, 301–319, doi: [10.1093/mnras/stv1198](https://doi.org/10.1093/mnras/stv1198)
- Gialluca, M. T., Naidu, R. P., & Bonaca, A. 2021, *The Astrophysical Journal Letters*, 911, L32, doi: [10.3847/2041-8213/abf491](https://doi.org/10.3847/2041-8213/abf491)
- Goodfellow, I., Bengio, Y., & Courville, A. 2016, *Deep Learning* (MIT Press)
- Greenberg, D. S., Nonnenmacher, M., & Macke, J. H. 2019, *Automatic Posterior Transformation for Likelihood-Free Inference*, arXiv, doi: [10.48550/ARXIV.1905.07488](https://doi.org/10.48550/ARXIV.1905.07488)
- Grillmair, C. J., & Dionatos, O. 2006, *The Astrophysical Journal*, 643, L17–L20, doi: [10.1086/505111](https://doi.org/10.1086/505111)
- Gutmann, M. U., Cor, J., & er. 2016, *Journal of Machine Learning Research*, 17, 1. <http://jmlr.org/papers/v17/15-017.html>
- Hermans, J., Banik, N., Weniger, C., Bertone, G., & Louppe, G. 2021, *Monthly Notices of the Royal Astronomical Society*, 507, 1999–2011, doi: [10.1093/mnras/stab2181](https://doi.org/10.1093/mnras/stab2181)
- Hermans, J., Begy, V., & Louppe, G. 2019, *Likelihood-free MCMC with Amortized Approximate Ratio Estimators*, arXiv, doi: [10.48550/ARXIV.1903.04057](https://doi.org/10.48550/ARXIV.1903.04057)
- Hernquist, L. 1990, *The Astrophysical Journal*, 356, 359, doi: [10.1086/168845](https://doi.org/10.1086/168845)
- Hilmi, T., Erkal, D., Kogosov, S. E., et al. 2024, *arXiv e-prints*, arXiv:2404.02953, doi: [10.48550/arXiv.2404.02953](https://doi.org/10.48550/arXiv.2404.02953)
- Hinton, D. R. G., & Williams, R. 1986, *Nature*. [rumelhart1986learning](https://doi.org/10.1038/30918a)
- Hopkins, P. F., Wetzel, A., Kereš, D., et al. 2018, *Monthly Notices of the Royal Astronomical Society*, 480, 800–863, doi: [10.1093/mnras/sty1690](https://doi.org/10.1093/mnras/sty1690)
- Hornik, K., Stinchcombe, M., & White, H. 1989, *Neural networks*, 2, 359

- Huang, X., Alkhalifah, T., & Song, C. 2021, in First International Meeting for Applied Geoscience & Energy Expanded Abstracts (Society of Exploration Geophysicists), doi: [10.1190/segam2021-3584127.1](https://doi.org/10.1190/segam2021-3584127.1)
- Ibata, R., Malhan, K., Martin, N., et al. 2021, *ApJ*, 914, 123, doi: [10.3847/1538-4357/abfcc2](https://doi.org/10.3847/1538-4357/abfcc2)
- Ibata, R. A., Lewis, G. F., Irwin, M. J., & Quinn, T. 2002, *MNRAS*, 332, 915, doi: [10.1046/j.1365-8711.2002.05358.x](https://doi.org/10.1046/j.1365-8711.2002.05358.x)
- Ibata, R. A., Malhan, K., & Martin, N. F. 2019, *ApJ*, 872, 152, doi: [10.3847/1538-4357/ab0080](https://doi.org/10.3847/1538-4357/ab0080)
- Jethwa, P., Torrealba, G., Navarrete, C., et al. 2018, *Monthly Notices of the Royal Astronomical Society*, 480, 5342–5351, doi: [10.1093/mnras/sty2226](https://doi.org/10.1093/mnras/sty2226)
- Johnston, K. V., Spergel, D. N., & Haydn, C. 2002, *The Astrophysical Journal*, 570, 656–664, doi: [10.1086/339791](https://doi.org/10.1086/339791)
- Kingma, D. P., & Ba, J. 2014, Adam: A Method for Stochastic Optimization, arXiv, doi: [10.48550/ARXIV.1412.6980](https://doi.org/10.48550/ARXIV.1412.6980)
- Kipf, T. N., & Welling, M. 2016, Semi-Supervised Classification with Graph Convolutional Networks, arXiv, doi: [10.48550/ARXIV.1609.02907](https://doi.org/10.48550/ARXIV.1609.02907)
- Kolmogorov, A. 1933, *G. Ist. Ital. Attuari.*
- Koposov, S. E., Irwin, M., Belokurov, V., et al. 2014, *MNRAS*, 442, L85, doi: [10.1093/mnrasl/slu060](https://doi.org/10.1093/mnrasl/slu060)
- Koposov, S. E., Rix, H.-W., & Hogg, D. W. 2010, *ApJ*, 712, 260, doi: [10.1088/0004-637X/712/1/260](https://doi.org/10.1088/0004-637X/712/1/260)
- Koppelman, H. H., & Helmi, A. 2021, *Astronomy & Astrophysics*, 649, A55, doi: [10.1051/0004-6361/202039968](https://doi.org/10.1051/0004-6361/202039968)
- Küpper, A. H. W., Kroupa, P., Baumgardt, H., & Heggie, D. C. 2010, *Monthly Notices of the Royal Astronomical Society*, 401, 105–120, doi: [10.1111/j.1365-2966.2009.15690.x](https://doi.org/10.1111/j.1365-2966.2009.15690.x)
- Küpper, A. H. W., Macleod, A., & Heggie, D. C. 2008, *Monthly Notices of the Royal Astronomical Society*, 387, 1248, doi: [10.1111/j.1365-2966.2008.13323.x](https://doi.org/10.1111/j.1365-2966.2008.13323.x)
- Leclercq, F. 2018, *Physical Review D*, 98, doi: [10.1103/physrevd.98.063511](https://doi.org/10.1103/physrevd.98.063511)
- Lemos, P., Cranmer, M., Abidi, M., et al. 2023, *Machine Learning: Science and Technology*, 4, 01LT01, doi: [10.1088/2632-2153/acbb53](https://doi.org/10.1088/2632-2153/acbb53)
- Lemos, P., Parker, L., Hahn, C., et al. 2024, *Physical Review D*, 109, doi: [10.1103/physrevd.109.083536](https://doi.org/10.1103/physrevd.109.083536)
- Li, T. S., Koposov, S. E., Zucker, D. B., et al. 2019, *Monthly Notices of the Royal Astronomical Society*, 490, 3508–3531, doi: [10.1093/mnras/stz2731](https://doi.org/10.1093/mnras/stz2731)
- Li, T. S., Koposov, S. E., Erkal, D., et al. 2021, *The Astrophysical Journal*, 911, 149, doi: [10.3847/1538-4357/abeb18](https://doi.org/10.3847/1538-4357/abeb18)
- Lin, K., von wiersheim Kramsta, M., Joachimi, B., & Feeney, S. 2023, *Monthly Notices of the Royal Astronomical Society*, 524, 6167–6180, doi: [10.1093/mnras/stad2262](https://doi.org/10.1093/mnras/stad2262)
- Liu, J., Mao, H., Chen, Z., et al. 2024, arXiv preprint arXiv:2402.02054
- Lueckmann, J.-M., Boelts, J., Greenberg, D., Goncalves, P., & Macke, J. 2021, in *Proceedings of Machine Learning Research*, Vol. 130, *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, ed. A. Banerjee & K. Fukumizu (PMLR), 343–351. <https://proceedings.mlr.press/v130/lueckmann21a.html>
- Lynden-Bell, D., & Lynden-Bell, R. M. 1995, *Monthly Notices of the Royal Astronomical Society*, 275, 429–442, doi: [10.1093/mnras/275.2.429](https://doi.org/10.1093/mnras/275.2.429)
- Macciò, A. V., Kang, X., Fontanot, F., et al. 2010, *Monthly Notices of the Royal Astronomical Society*, 402, 1995–2008, doi: [10.1111/j.1365-2966.2009.16031.x](https://doi.org/10.1111/j.1365-2966.2009.16031.x)
- Malhan, K., Ibata, R. A., Carlberg, R. G., Valluri, M., & Freese, K. 2019, *The Astrophysical Journal*, 881, 106, doi: [10.3847/1538-4357/ab2e07](https://doi.org/10.3847/1538-4357/ab2e07)
- Malhan, K., & Rix, H.-W. 2024, *The Astrophysical Journal*, 964, 104, doi: [10.3847/1538-4357/ad1885](https://doi.org/10.3847/1538-4357/ad1885)
- Mandelbaum, R., Seljak, U., Kauffmann, G., Hirata, C. M., & Brinkmann, J. 2006, *Monthly Notices of the Royal Astronomical Society*, 368, 715–731, doi: [10.1111/j.1365-2966.2006.10156.x](https://doi.org/10.1111/j.1365-2966.2006.10156.x)
- Massey, R., Kitching, T., & Richard, J. 2010, *Reports on Progress in Physics*, 73, 086901, doi: [10.1088/0034-4885/73/8/086901](https://doi.org/10.1088/0034-4885/73/8/086901)
- McDonald, P., Miralda-Escude, J., Rauch, M., et al. 2000, *The Astrophysical Journal*, 543, 1–23, doi: [10.1086/317079](https://doi.org/10.1086/317079)
- McKay, M. D., Beckman, R. J., & Conover, W. J. 1979, *Technometrics*, 21, 239, doi: [10.2307/1268522](https://doi.org/10.2307/1268522)
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. 2013, Efficient Estimation of Word Representations in Vector Space, arXiv, doi: [10.48550/ARXIV.1301.3781](https://doi.org/10.48550/ARXIV.1301.3781)
- Miller, B. K., Weniger, C., & Forré, P. 2022, Contrastive Neural Ratio Estimation for Simulation-based Inference, arXiv, doi: [10.48550/ARXIV.2210.06170](https://doi.org/10.48550/ARXIV.2210.06170)
- Minchev, I., Boily, C., Siebert, A., & Bienayme, O. 2010, *Monthly Notices of the Royal Astronomical Society*, 407, 2122–2130, doi: [10.1111/j.1365-2966.2010.17060.x](https://doi.org/10.1111/j.1365-2966.2010.17060.x)
- Mishra-Sharma, S. 2022, *Machine Learning: Science and Technology*, 3, 01LT03, doi: [10.1088/2632-2153/ac494a](https://doi.org/10.1088/2632-2153/ac494a)
- Miyamoto, M. ; Nagai, R. 1975, *Astronomical Society of Japan*, 27, 533

- Nadler, E., Drlica-Wagner, A., Bechtol, K., et al. 2021, *Physical Review Letters*, 126, doi: [10.1103/physrevlett.126.091101](https://doi.org/10.1103/physrevlett.126.091101)
- Nadler, E. O., Birrer, S., Gilman, D., et al. 2021, *ApJ*, 917, 7, doi: [10.3847/1538-4357/abf9a3](https://doi.org/10.3847/1538-4357/abf9a3)
- Navarro, J. F., Frenk, C. S., & White, S. D. M. 1996, *The Astrophysical Journal*, 462, 563, doi: [10.1086/177173](https://doi.org/10.1086/177173)
- Nayak, P., Walther, M., Gruen, D., & Adiraju, S. 2024, *Astronomy & Astrophysics*, 689, A153, doi: [10.1051/0004-6361/202348485](https://doi.org/10.1051/0004-6361/202348485)
- Newberg, H. J., & Carlin, J. L. 2016, *Tidal Streams in the Local Group and Beyond: Observations and Implications* (Springer International Publishing), doi: [10.1007/978-3-319-19336-6](https://doi.org/10.1007/978-3-319-19336-6)
- Ngan, W., Bozek, B., Carlberg, R. G., et al. 2015, *The Astrophysical Journal*, 803, 75, doi: [10.1088/0004-637x/803/2/75](https://doi.org/10.1088/0004-637x/803/2/75)
- Nguyen, N.-M., Schmidt, F., Tucci, B., Reinecke, M., & Kostić, A. 2024, How much information can be extracted from galaxy clustering at the field level?, arXiv, doi: [10.48550/ARXIV.2403.03220](https://doi.org/10.48550/ARXIV.2403.03220)
- Nguyen, T., Mishra-Sharma, S., Williams, R., & Necib, L. 2023, *Physical Review D*, 107, doi: [10.1103/physrevd.107.043015](https://doi.org/10.1103/physrevd.107.043015)
- Nibauer, J., Belokurov, V., Cranmer, M., Goodman, J., & Ho, S. 2022, *The Astrophysical Journal*, 940, 22, doi: [10.3847/1538-4357/ac93ee](https://doi.org/10.3847/1538-4357/ac93ee)
- Odenkirchen, M., Grebel, E. K., Rockosi, C. M., et al. 2001, *ApJL*, 548, L165, doi: [10.1086/319095](https://doi.org/10.1086/319095)
- Papamakarios, G., & Murray, I. 2016, in *Advances in Neural Information Processing Systems*, ed. D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, & R. Garnett, Vol. 29 (Curran Associates, Inc.), https://proceedings.neurips.cc/paper_files/paper/2016/file/6aca97005c68f1206823815f66102863-Paper.pdf
- Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., & Lakshminarayanan, B. 2019, doi: [10.48550/ARXIV.1912.02762](https://doi.org/10.48550/ARXIV.1912.02762)
- Papamakarios, G., Pavlakou, T., & Murray, I. 2017, *Masked Autoregressive Flow for Density Estimation*, arXiv, doi: [10.48550/ARXIV.1705.07057](https://doi.org/10.48550/ARXIV.1705.07057)
- Papamakarios, G., Sterratt, D. C., & Murray, I. 2018, *Sequential Neural Likelihood: Fast Likelihood-free Inference with Autoregressive Flows*, arXiv, doi: [10.48550/ARXIV.1805.07226](https://doi.org/10.48550/ARXIV.1805.07226)
- Patrick, J. M., Kuposov, S. E., & Walker, M. G. 2022, *MNRAS*, 514, 1757, doi: [10.1093/mnras/stac1478](https://doi.org/10.1093/mnras/stac1478)
- Price-Whelan, A. M., & Bonaca, A. 2018, *The Astrophysical Journal Letters*, 863, L20, doi: [10.3847/2041-8213/aad7b5](https://doi.org/10.3847/2041-8213/aad7b5)
- Qian, Y., Arshad, Y., & Bovy, J. 2022, *Monthly Notices of the Royal Astronomical Society*, 511, 2339–2348, doi: [10.1093/mnras/stac238](https://doi.org/10.1093/mnras/stac238)
- Rogers, K. K., Dvorkin, C., & Peiris, H. V. 2022, *Phys. Rev. Lett.*, 128, 171301, doi: [10.1103/PhysRevLett.128.171301](https://doi.org/10.1103/PhysRevLett.128.171301)
- Rogers, K. K., & Peiris, H. V. 2021a, *Physical Review Letters*, 126, doi: [10.1103/physrevlett.126.071302](https://doi.org/10.1103/physrevlett.126.071302)
- . 2021b, *Phys. Rev. D*, 103, 043526, doi: [10.1103/PhysRevD.103.043526](https://doi.org/10.1103/PhysRevD.103.043526)
- Shipp, N., Drlica-Wagner, A., Balbinot, E., et al. 2018, *ApJ*, 862, 114, doi: [10.3847/1538-4357/aacdab](https://doi.org/10.3847/1538-4357/aacdab)
- Shipp, N., Panithanpaisal, N., Necib, L., et al. 2023, *The Astrophysical Journal*, 949, 44, doi: [10.3847/1538-4357/acc582](https://doi.org/10.3847/1538-4357/acc582)
- Smirnov, N. 1948, *Annals of Mathematical Statistics*.
- Snoek, J., Larochelle, H., & Adams, R. P. 2012, *Practical Bayesian Optimization of Machine Learning Algorithms*, arXiv, doi: [10.48550/ARXIV.1206.2944](https://doi.org/10.48550/ARXIV.1206.2944)
- Tabak, E. G., & Turner, C. V. 2012, *Communications on Pure and Applied Mathematics*, 66, 145–164, doi: [10.1002/cpa.21423](https://doi.org/10.1002/cpa.21423)
- Tabak, E. G., & Vanden-Eijnden, E. 2010, *Communications in Mathematical Sciences*, 8, 217–233, doi: [10.4310/cms.2010.v8.n1.a11](https://doi.org/10.4310/cms.2010.v8.n1.a11)
- Talts, S., Betancourt, M., Simpson, D., Vehtari, A., & Gelman, A. 2018, *Validating Bayesian Inference Algorithms with Simulation-Based Calibration*, arXiv, doi: [10.48550/ARXIV.1804.06788](https://doi.org/10.48550/ARXIV.1804.06788)
- Tavangar, K., Ferguson, P., Shipp, N., et al. 2022, *ApJ*, 925, 118, doi: [10.3847/1538-4357/ac399b](https://doi.org/10.3847/1538-4357/ac399b)
- Tejero-Cantero, A., Boelts, J., Deistler, M., et al. 2020, *Journal of Open Source Software*, 5, 2505, doi: [10.21105/joss.02505](https://doi.org/10.21105/joss.02505)
- Vale, A., & Ostriker, J. P. 2004, *Monthly Notices of the Royal Astronomical Society*, 353, 189–200, doi: [10.1111/j.1365-2966.2004.08059.x](https://doi.org/10.1111/j.1365-2966.2004.08059.x)
- Valluri, M., Chabanier, S., Irsic, V., et al. 2022, *Snowmass2021 Cosmic Frontier White Paper: Prospects for obtaining Dark Matter Constraints with DESI*, arXiv, doi: [10.48550/ARXIV.2203.07491](https://doi.org/10.48550/ARXIV.2203.07491)
- Vaswani, A., Shazeer, N., Parmar, N., et al. 2017, *Attention Is All You Need*, arXiv, doi: [10.48550/ARXIV.1706.03762](https://doi.org/10.48550/ARXIV.1706.03762)
- Vegetti, S., Koopmans, L. V. E., Bolton, A., Treu, T., & Gavazzi, R. 2010, *Monthly Notices of the Royal Astronomical Society*, 408, 1969–1981, doi: [10.1111/j.1365-2966.2010.16865.x](https://doi.org/10.1111/j.1365-2966.2010.16865.x)
- Vegetti, S., Birrer, S., Despali, G., et al. 2023, *Strong gravitational lensing as a probe of dark matter*, arXiv, doi: [10.48550/ARXIV.2306.11781](https://doi.org/10.48550/ARXIV.2306.11781)

- Villasenor, B., Robertson, B., Madau, P., & Schneider, E. 2023, *Physical Review D*, 108, doi: [10.1103/physrevd.108.023502](https://doi.org/10.1103/physrevd.108.023502)
- Wang, L., Dutton, A. A., Stinson, G. S., et al. 2015, *Monthly Notices of the Royal Astronomical Society*, 454, 83–94, doi: [10.1093/mnras/stv1937](https://doi.org/10.1093/mnras/stv1937)
- Webb, J. J., & Bovy, J. 2019, *Monthly Notices of the Royal Astronomical Society*, 485, 5929–5938, doi: [10.1093/mnras/stz867](https://doi.org/10.1093/mnras/stz867)
- Wechsler, R. H., & Tinker, J. L. 2018, *Annual Review of Astronomy and Astrophysics*, 56, 435–487, doi: [10.1146/annurev-astro-081817-051756](https://doi.org/10.1146/annurev-astro-081817-051756)
- Weinberg, D. H. 2003, in *AIP Conference Proceedings (AIP)*, doi: [10.1063/1.1581786](https://doi.org/10.1063/1.1581786)
- Wu, J. F., & Jespersen, C. K. 2023, Learning the galaxy-environment connection with graph neural networks, arXiv, doi: [10.48550/ARXIV.2306.12327](https://doi.org/10.48550/ARXIV.2306.12327)
- Zhang, X., Xu, Y., He, W., Guo, W., & Cui, L. 2024a, A Comprehensive Review of the Oversmoothing in Graph Neural Networks (Springer Nature Singapore), 451–465, doi: [10.1007/978-981-99-9637-7_33](https://doi.org/10.1007/978-981-99-9637-7_33)
- Zhang, Z., Lin, C., & Wang, B. 2024b, *Scientific Reports*, 14, doi: [10.1038/s41598-024-57137-4](https://doi.org/10.1038/s41598-024-57137-4)