

A neural simulation-based inference approach for characterizing the Galactic Center γ -ray excess

Siddharth Mishra-Sharma^{1,*} and Kyle Cranmer^{1,2,†}

¹*Center for Cosmology and Particle Physics, Department of Physics,
New York University, New York, NY 10003, USA*

²*Center for Data Science, New York University, 60 Fifth Ave, New York, NY 10011, USA*
(Dated: July 22, 2021)

The nature of the *Fermi* γ -ray Galactic Center Excess (GCE) has remained a persistent mystery for over a decade. Although the excess is broadly compatible with emission expected due to dark matter annihilation, an explanation in terms of a population of unresolved point sources remains viable. The effort to uncover the origin of the GCE is hampered in particular by an incomplete understanding of diffuse emission of Galactic origin, which can lead to spurious features that make it difficult to robustly differentiate smooth emission, as expected for a dark matter origin, from more “clumpy” emission expected from a population of relatively bright, unresolved PSs. We use machine learning-based likelihood-free inference methods, in particular conditional density estimation with normalizing flows, in order to extract more information from *Fermi* γ -ray data of the Galactic Center with the aim of characterizing the contribution of unresolved point sources to the GCE. We find a significantly reduced preference for a point source origin of the GCE compared to analyses using traditional methods, with about half of the GCE flux attributable to a population of dim point sources.

I. INTRODUCTION

Dark matter (DM) represents one of the major unsolved problems in particle physics and cosmology today. The traditional Weakly-Interacting Massive Particle (WIMP) paradigm envisions production of dark matter in the early Universe through freeze-out of dark sector particles weakly coupled to the Standard Model (SM) sector. In this scenario, one of the most promising avenues of detecting a dark matter signal is through an observation of excess γ -ray photons at \sim GeV energies from DM-rich regions of the sky produced through the cascade of SM particles resulting from DM self-annihilation.

The *Fermi* γ -ray Galactic Center Excess (GCE), first identified over a decade ago using data from the *Fermi* Large Area Telescope (LAT) [1], is an excess of photons in the Galactic Center with properties—such as energy spectrum and spatial morphology—broadly compatible with expectation due to annihilating DM [2–16]. The nature of the GCE remains contentious however, with competing explanations in terms of a population of unresolved astrophysical point sources (PSs), in particular millisecond pulsars (MSPs), remaining viable [9, 17–23]. Analyses of the morphology of the excess have shown it to prefer a spatial distribution tracing the stellar bulge in the Galactic Center [15, 24, 25] rather than the expected distribution due to DM annihilation. Furthermore, analyses leveraging the statistics of photons in the Galactic Center have shown the γ -ray data to prefer a point source origin of the excess [26–29]. Recent studies have however pointed out the potential of unknown

systematics—such as the poorly understood morphology of the diffuse foreground emission and unmodeled point source populations—to affect the conclusions of these analyses [30–32]. [SM: Tighten references here.]

Due to the high dimensionality of γ -ray data, a description of the photon map in terms of hand-crafted summary statistics such as the probability distribution of photon counts [26, 33] or a wavelet decomposition of the photon map [27, 34, 35] has traditionally been necessary in order to enable computationally tractable analyses. While effective, this reduced description necessarily involves loss of information compared to the original γ -ray map, increasing susceptibility to mismodeled features. On the other hand, recent developments in machine learning have given rise to analysis techniques that can extract more information from high-dimensional datasets and, consequently, more robustly hedge against unknown systematics in the data compared to traditional analyses based on specific data summaries. Machine learning methods have demonstrated promise for analyzing γ -ray data [36] and in particular for understanding the nature of the *Fermi* GCE [37, 38].

In this paper, we leverage recent developments in the field of simulation-based inference (SBI, also referred to as likelihood-free inference; see, *e.g.*, Ref. [39] for a recent review) in order to weigh in on the nature of the GCE. In particular, we employ conditional density estimation techniques based on normalizing flows [40, 41] in order to characterize the contributions of various modeled components, including “clumpy” PS-like and “smooth” DM-like emission spatially tracing the GCE, to the γ -ray photon sky at \sim GeV energies in the Galactic Center region. We employ graph-based neural network architecture in order to automatically extract summary statistics from γ -ray maps optimized for the downstream task of estimating the distribution of parameters characterizing the contri-

* sm8383@nyu.edu; ORCID: 0000-0001-9088-7845

† kyle.cranmer@nyu.edu; ORCID: 0000-0002-5769-7094

bution of modeled components to the GCE.

This paper is organized as follows. In Sec. II we describe our forward model and the analysis framework based on neural simulation-based inference. In Sec. III we validate our analysis on various mock observations of the *Fermi* GCE. Section IV presents an application of the method to *Fermi* γ -ray data, including a discussion of systematic variations on the analysis. In Sec. V we quantify the susceptibility of the analysis to mismodeling of the signal and background templates using data-driven techniques. We conclude in Sec. VI.

II. METHODOLOGY

We begin by describe the various ingredients of our forward model and the datasets used. We then detail our analysis methodology, going over in turn the general principles behind simulation-based inference, posterior estimation using normalizing flows, and learning representative summary statistics from high-dimensional γ -ray maps with graph neural networks.

A. Datasets and the forward model

We use the datasets and templates from Ref. [42] (packaged with Ref. [43]) to create the simulated maps for forward modeling. The data and templates used correspond to 413 weeks of *Fermi*-LAT Pass 8 data collected between August 4, 2008 and July 7, 2016. The top quarter of photons in the energy range 2–20 GeV by quality of PSF reconstruction (corresponding to the PSF3 event type) in the event class ULTRACLEANVETO are used. The recommended quality cuts are applied, corresponding to zenith angle less than 90° , LAT.CONFIG = 1, and DATA_QUAL > 0.1.¹ The maps are spatially binned using HEALPix [44] with nside=128.

The simulated data maps are a combination of smooth (*i.e.*, Poissonian) and PS contributions. Each PS population is completely specified by its spatial distribution, described by a spatial template, and the distribution of photon counts, specified by a source-count distribution. PS populations spatially correlated with (*i*) the GCE, modeled using a (squared) generalized Navarro-Frenk-White (NFW) [45, 46] profile with inner slope $\gamma = 1.2$ motivated by previous GCE analyses [8, 10, 47],

$$\rho(r) \propto \frac{1}{(r/r_s)^\gamma (1 + r/r_s)^{3-\gamma}} \quad (1)$$

where $r_s = 20$ kpc is the Milky Way scale radius, [SM: Specify distance to GC.] and (*ii*) the Galactic disk, modeled as a doubly-exponential profile motivated by studies

of Galactic millisecond pulsar populations [48, 49],

$$n(R, z) \propto \exp(-R/5 \text{ kpc}) \exp(-|z|/1 \text{ kpc}) \quad (2)$$

where R and z are the radial and vertical Galactic cylindrical coordinates. Photon counts from a generated PS population are put down on a map according to the *Fermi* PSF at 2 GeV, modeled as a King function, using the algorithm implemented in the code package NPTFit-Sim [50]. The source-count distribution (SCD) dN/dS of each PS population, describing the differential number of sources per photon counts, is modeled as a doubly-broken power law,

$$\frac{dN_p}{dS} = A_{\text{PS}} \begin{cases} \left(\frac{S}{S_{b,1}}\right)^{-n_1}, & S \geq S_{b,1} \\ \left(\frac{S}{S_{b,1}}\right)^{-n_2}, & S_{b,1} > S \geq S_{b,2} \\ \left(\frac{S_{b,2}}{S_{b,1}}\right)^{-n_2} \left(\frac{S}{S_{b,2}}\right)^{-n_3}, & S_{b,2} > S \end{cases} \quad (3)$$

specified by the breaks $\{S_{b,1}, S_{b,2}\}$, slopes $\{n_1, n_2, n_3\}$, and an overall normalization A_{PS} . We note that these parameters specify the spatially-averaged properties of the PS population—variation due to non-uniform exposure of the LAT instrument is accounted for in putting down simulated photon counts.

In addition to the PS-like emission, we account for Poissonian astrophysical emission in the simulated maps. These contributions include: (*i*) the Galactic diffuse foreground emission, described in the baseline model by Model O from Ref. [28], [SM: More description of diffuse emission.] (*ii*) DM-like emission following a generalized squared-NFW profile as in Eq. (1) with inner slope $\gamma = 1.2$, (*iii*) spatially isotropic emission, (*iv*) resolved PSs from the *Fermi* 3FGL catalog [51], and (*v*) emission from the *Fermi* bubbles [52]. The relative normalizations of these templates are included as parameters of the forward model. The latter three templates are obtained from Ref. [42]. The final maps are obtained by combining a Poisson-fluctuated realization of the summed astrophysical templates with the PS maps. The inner regions of the Galactic plane are masked at $|b| < 2^\circ$, and a radial cut $r < 25^\circ$ defines the region of interest for the fiducial analysis. We mask resolved PSs from the 3FGL catalog at a containment radius of 0.8° . [SM: Specify motivation for mask size.]

The forward model is thus specified by 18 parameters—6 parameters for the overall normalizations of the Poissonian templates, and 6×2 parameters modeling the source-count distribution associated with GCE-correlated and disk-correlated PS populations. Simulated samples are produced using parameters drawn from priors motivated by a Poissonian fit to the real *Fermi* data in order to improve sample efficiency. Since PS-like and Poissonian components of the model are exactly degenerate in the limit of each PS producing $\lesssim 1$ photon counts in expectation (see Ref. [29] for a detailed discussion of this degeneracy), we place priors on the expected counts contribution per pixel for the

¹ https://fermi.gsfc.nasa.gov/ssc/data/analysis/documentation/Cicerone/Cicerone_Data_Exploration/Data_preparation.html

GCE-correlated PS-like and Poissonian emission in order to mitigate biases caused by an induced prior preferring one model over the other and preventing the expression of the degeneracy. [SM: Specify priors, maybe in a table. Cite compound Poisson generator paper.]

B. Simulation-based inference

Irrespective of domain, of central interest in parameter estimation is often the probability distribution of a set of parameters of interest θ given some data x —the posterior distribution $p(\theta | x)$. Bayes’ theorem can be used to obtain the posterior as $p(\theta | x) = p(\theta) p(x | \theta) / \mathcal{Z}$, where $p(x | \theta)$ is the likelihood and $\mathcal{Z} \equiv \int p(x) d\theta$ is the Bayesian evidence. In practice, parameters other than θ —latent variables z —are often involved in the data-generation process, and computing the likelihood involves marginalizing over the latent space, $p(x | \theta) = \int dz p(x | \theta, z)$. In typical problems of interest, the high dimensionality of the latent space often means that this integral is intractable, necessitating simplifications in statistical treatment as well as theoretical modeling.

Simulation-based inference (SBI) refers to a class of methods for performing inference when the data-generating process does not have a tractable likelihood. In this setting, a model is defined through a simulator as a probabilistic program, often known as a forward model. Samples x from the simulator then implicitly define a likelihood, $x \sim p(x | \theta)$. In the simplest realizations of SBI, samples x' generated from a given prior proposal distribution $p(\theta)$ can be compared to a given dataset of interest x , with the approximate posterior defined by samples that most closely resemble x according to some similarity metric. Such methods—usually grouped under the umbrella of Approximate Bayesian Computation (ABC) [53]—are not uncommon in astrophysics and cosmology. Nevertheless, they suffer from several downsides. The curse of dimensionality usually necessitates reduction of data to representative, hand-crafted lower-dimensional summary statistics $s(x)$, resulting in loss of information. A notion and measure of distance between summaries from the implicit model and those derived from the dataset of interest is necessary, leading to inexact inference. Additionally, the ABC analysis must be performed anew for each new target dataset.

Recent methods [54–70] have leveraged advancements in machine learning, in particular the ability of neural networks to extract useful features from high-dimensional data and to flexibly approximate functions and distributions, in order to address these issues, enabling new ways of performing inference on complex models defined through simulations; see Ref. [39] for a review of recent developments.

C. Conditional density estimation with normalizing flows

In this paper, we approximate the joint posterior $p(\theta | x)$ through a parameterized distribution $\hat{p}_\phi(\theta | s)$ conditioned on *learned* summaries $s = s(x)$ from the simulated samples x . This class of simulation-based inference techniques, known as conditional neural density estimation [67], directly models the posterior distribution given a set of samples drawn from a simulator according to some prior proposal distribution $p(\theta)$.

We employ normalizing flows [40, 41], which provide an efficient way of constructing flexible probability distributions. Normalizing flows model the conditional distribution $\hat{p}_\phi(\theta | s)$ as a series of invertible transformations, denoted f and having a tractable inverse and Jacobian, from a base distribution $\pi_z(z)$, chosen here to be a standard Gaussian $z \sim \mathcal{N}(0, \mathbb{1})$, to the target distribution:

$$\hat{p}(\theta | x) = \pi_z(f^{-1}(\theta)) \left| \det \left(\frac{\partial f^{-1}}{\partial \theta} \right) \right| \quad (4)$$

Specifically, we use Masked Autoregressive Flows (MAFs) [71] for posterior estimation. The MAF is built upon blocks of affine transformations where scaling and shifting factors for each dimension are computed with a Masked Autoencoder for Distribution Estimation (MADE) [72]. For a single block, the transformation from θ to z is expressed as

$$z_i = (\theta_i - \mu_i) \cdot \exp(-\alpha_i) \quad (5)$$

where $\mu_i = f_{\mu_i}(\theta_{1:i-1}; x)$ and $\alpha_i = f_{\alpha_i}(\theta_{1:i-1}; x)$ are scaling and shift factors modeled by a MADE according to the autoregressive condition. This allows for an analytically tractable determinant,

$$\left| \det \left(\frac{\partial f^{-1}}{\partial \theta} \right) \right| = \exp \left(- \sum_i \alpha_i \right) \quad (6)$$

and a forward pass through the flow according to Eq. (5). Multiple transformations can be composed together as $f = f_1 \circ f_2 \circ \dots \circ f_K$ in order to model more expressive posteriors,

$$\hat{p}(\theta | x) = \pi_z(f^{-1}(\theta)) \prod_{i=1}^K \left| \det \left(\frac{\partial f_i^{-1}}{\partial z_{i-1}} \right) \right|. \quad (7)$$

The log-probability of the posterior can then be computed using Eq. (6):

$$\log \hat{p}(\theta | x) = \log [\pi_z(f^{-1}(\theta))] - \sum_{i=1}^K \sum_{j=1}^N \alpha_j^i, \quad (8)$$

which acts as the optimization objective. Here, we use 8 MAF transformations, each made up of a 2-layer MADE with 128 hidden units. Each transformation is conditioned on summaries $s(x)$ extracted from the γ -ray maps (described below) by including these as inputs into each transformation block. [SM: Clean up discussion of normalizing flows.]

D. Learning summary statistics with (graph) neural networks

The curse of dimensionality makes it computationally prohibitive to condition the density estimation task on the raw dataset x *i.e.*, the γ -ray pixel counts map in the region of interest (ROI). Representative summaries $s = s_\varphi(x)$ of the data must therefore be extracted in order to enable a tractable analysis, where φ parameterizes the data-to-summary transformation. Although many choices for data summaries are possible—*e.g.*, a Principal Component Analysis (PCA) decomposition of the photon counts map, an angular power spectrum decomposition of the photon counts map, or simply a histogram of the photon counts—in this paper, we use a neural net to automatically learn low-dimensional summaries that are efficiently suited for the specific downstream task at hand.

Graph construction and architecture

The *DeepSphere* architecture [73, 74], with a configuration similar to that employed in Ref. [37], is used to extract representative summaries and is briefly outlined here. *DeepSphere* is a graph-based convolutional neural network (CNN) architecture tailored to data sampled on a sphere, and in particular is able to leverage the hierarchical structure of data in the HEALPix representation. This makes it well-suited for our purposes.

The HEALPix sphere is represented as a weighted undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, A)$, where \mathcal{V} is the set of $N_{\text{pix}} = |\mathcal{V}|$ vertices, \mathcal{E} is the set of edges, and A is the weighted adjacency matrix. Each pixel i is represented by a vertex $v_i \in \mathcal{V}$. Each vertex v_i is then connected to the 8 vertices v_j which represent the neighboring pixels j of pixel i , forming edges $(v_i, v_j) \in \mathcal{E}$. Given those edges, we define the weights of the adjacency matrix A over neighboring pixels following the weighing scheme given in Ref. [73].

We use the combinatorial graph Laplacian, defined as $L = D - A$, where D is the diagonal degree matrix, and can be used to define a Fourier basis on a graph. By construction symmetric positive semi-definite, the graph Laplacian can be decomposed as $L = U\Lambda U^T$, where U is an orthonormal eigenvector matrix and Λ is a diagonal eigenvalue matrix. The Laplacian eigenvectors then define the graph Fourier basis, with the graph Fourier transform \tilde{x} of a signal x on a graph being its projection $\tilde{x} = U^T x$. Given a convolutional kernel h , graph convolutions can be efficiently performed in the Fourier basis as $h(L)x = Uh(\Lambda)U^T x$.

The *DeepSphere* convolutional kernel h is defined as a linear combination of Chebyshev polynomials, $h_\theta(L) = \sum_{k=0}^K \theta_k T_k(L)$ where T_k are the order- k Chebyshev polynomials and θ_k are the $K + 1$ filter coefficients to be fit.

The graph filtering operation can then be expressed as

$$h_\theta(L)x = U \left(\sum_{k=0}^K \theta_k T_k(\Lambda) \right) U^T x = \sum_{k=0}^K \theta_k T_k(L)x. \quad (9)$$

We set $K = 4$ as the maximum Chebyshev polynomial order, having checked that larger values do not qualitatively affect the results of the study.

Following Ref. [74], the feature extraction architecture is built out of layers which progressively coarsen the pixel representation of the γ -ray maps while increasing the number of filter channels at each step. Starting with HEALPix resolution `nside=128`, each graph convolution operation is followed by a BatchNorm, a ReLU nonlinearity, and a max pooling operation which downsamples the representation into a coarser resolution, starting with `nside=128` until resolution `nside=1` after the final convolutional layer. The number of filter channels is doubled at each convolution until a maximum of 256. The output of the final convolutional layer is augmented with 2 additional auxiliary variables—the log-mean and log-standard deviation of the γ -ray maps within the region of interest—and passed through a fully-connected layer with 1024 hidden units outputting a desired number of summary features, which is 128 in our fiducial configuration. All input features (*i.e.*, individual non-zero pixels as well as auxiliary variables) are individually normalized to zero mean and unit variance. [SM: Tighten discussion, specify how masking is done.]

Optimization, training, and evaluation

The graph-based and normalizing flow networks are trained simultaneously. 10^6 samples are generated using the prior proposal distribution, and models are optimized with batch size 128 using the AdamW [75, 76] optimizer with initial learning rate 10^{-3} and weight decay 10^{-5} , using cosine annealing to decay the learning rate across epochs. Training proceeds for up to 50 epochs with early stopping if the validation loss (evaluated on 15% of held-out samples) has not improved after 8 epochs.

After training, given a new dataset (either real or simulated *Fermi* data in our ROI), the posterior is obtained by sampling the flow within the prior using rejection sampling, conditioning each flow transformation on summaries returned by the graph-based neural network with the new dataset as input.

III. TESTS ON SIMULATED DATA

We begin by validating our pipeline on simulated *Fermi* data. We create simulated datasets by drawing parameter values from ranges motivated by a fit of the model to real *Fermi* data in our fiducial ROI, and test the ability of our model to infer the presence of either DM-like

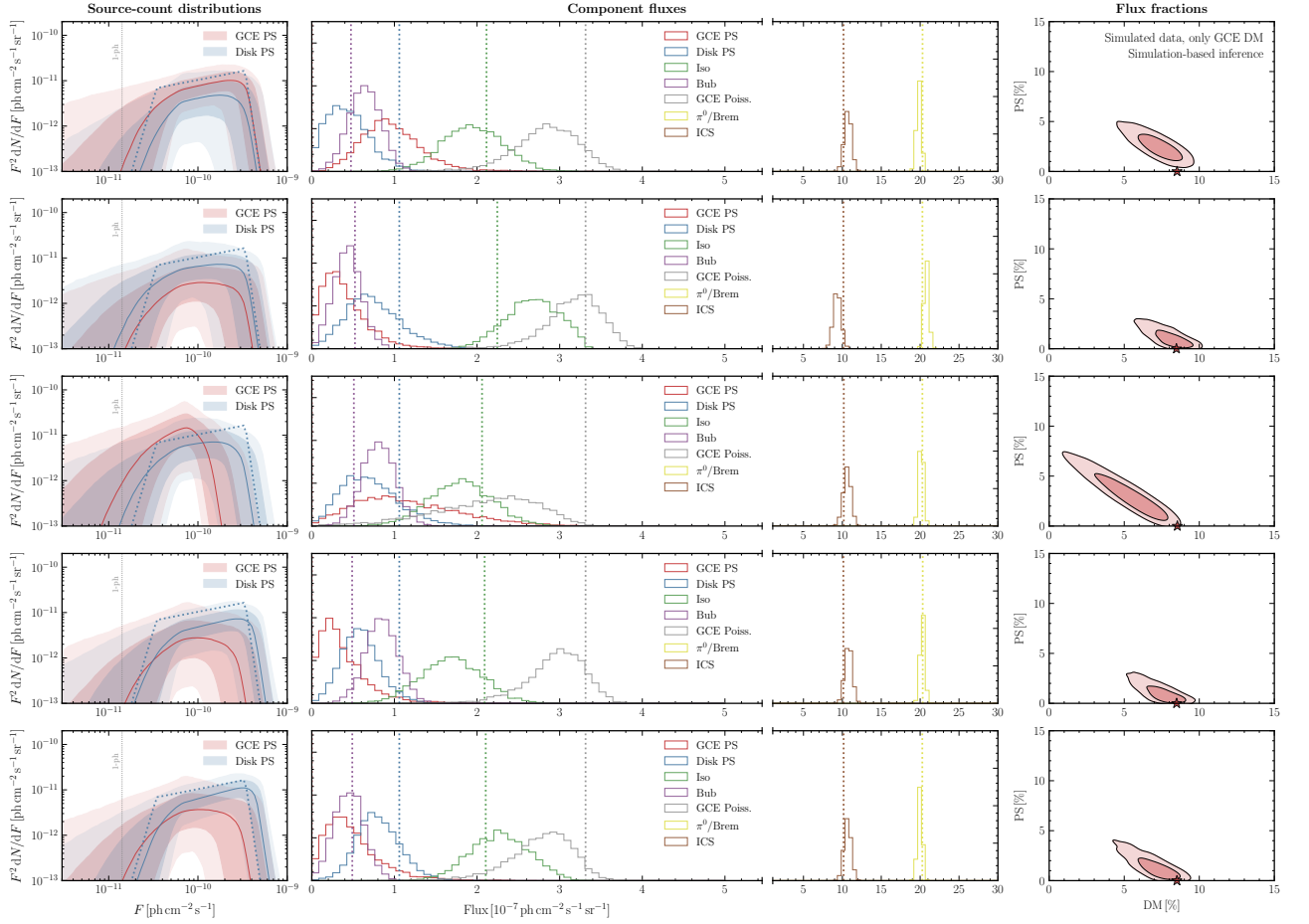


FIG. 1. Results on simulated *Fermi* data where the GCE consists of purely DM-like emission, with different rows corresponding to different simulated realizations. The left column shows the inferred source-count distribution for GCE-correlated (red) and disk-correlated (blue) PS. The middle panel shows the posteriors for the Poissonian templates. The right panel shows the joint posterior on DM-like and PS-like emission. The dotted lines corresponding to the true simulated quantities in the left two panels, and the starred point in the right panel. DM-like emission is inferred in each case, with the other posteriors corresponding well to their true simulated values.

or PS-like signals on top of the modeled astrophysical background.

Figure 1 shows results of the analysis pipeline conditioned on five simulated realizations of maps where the GCE consists of purely DM-like emission. The left column shows the middle-68/95% containment of the point-wise posterior on the source-count distributions of GCE- and disk-correlated point source emission in red and blue, respectively. The middle column shows the posteriors on various modeled emission components, excluding emission from resolved 3FGL PSs as the posterior in that case is largely unconstrained owing to the fact that resolved PSs are masked out in the analysis. The right column shows the fraction of DM- and PS-like emission in proportion to the total inferred flux in the ROI. The true underlying parameter values from which the data was generated are represented by dotted lines in the left and middle columns, and by star markers in the right column. We see that, in all cases shown, the pipeline successfully

recovers the presence of DM-like emission, with little flux attributed to unresolved PSs. Some PS-like emission is inferred in most cases as well however, due to a combination of degeneracy with both disk-correlated and DM-like flux. The overall flux of all components corresponds well to their true underlying values.

Figure 2 shows the corresponding results for simulated data containing PS-like emission correlated with the GCE. Here, simulations were produced such that the highest break of the GCE-correlated PS SCD was contained between 5 and 20 expected photon counts, since we found that the method cannot robustly attribute an SCD that corresponds to a peak dimmer than $\lesssim 5$ photon counts to a PS population. We see that PS-like emission is successfully inferred in each case, while at the same time exemplifying a degeneracy with the Poissonian component. Furthermore, as seen in the left column, the method is able to characterize the contribution of the two modeled PS components through the inferred source-

count distribution. Some degeneracy between GCE- and disk-correlated PSs is seen, although the true SCDs are seen to lie within the 95% containment interval of the inferred point-wise SCD posteriors in each case.

IV. RESULTS ON *FERMI* DATA

We finally apply our formalism to the real *Fermi* dataset in our ROI. As a point of comparison, we also run the NPTF pipeline on the data using the same spatial templates and prior assumptions described in Sec. II A. The NPTF likelihood implemented in `NPTFit` from Ref. [43] is used, and the posterior distribution is constructed through nested sampling using `dynesty` with 500 live points. The results of this analysis are shown in the bottom panel of Fig. 3. Consistent with previous analyses using a similar configuration, a preference for PS-like emission is seen, with the analysis attributing the majority of the GCE to PS-like emission. [SM: More detailed description of NPTF.]

The top panel of Fig. 3 shows results using the neural simulation-based analysis pipeline introduced in this paper. Although posteriors for the astrophysical background templates are seen to be broadly consistent with the NPTF analysis, the preference for PSs is significantly reduced in this case. In fact, about half of the emission is attributed to PS-like and Poissonian emission each. We also note that the inferred GCE-correlated SCD peaks at values lower than those inferred from previous NPTF analyses, which have generally found the bulk of expected emission from PSs to lie just below the 3FGL PS detection threshold [26] at $\sim 2\text{--}3 \times 10^{-10} \text{ ph cm}^{-2} \text{ s}^{-1}$. [SM: Quantify where the SCD peaks and how many PSs it predicts.]

A. Signal injection test on data

A crucial self-consistency test is the ability of the analysis to recover an artificial signal injected onto the real γ -ray data. As shown in Ref. [32], early applications of the 1-point PDF based methods like NPTF to the GCE would generally fail this closure test, with implications for characterizing the nature of PSs in the Galactic Center explored in Refs. [28, 29]. In particular, it was shown that the closure test can help diagnose underlying issues associated with mismodeling of the diffuse foreground emission, which have the potential to bias the characterization of PS populations. We perform a version of this test within our framework, testing the ability of our method to recover different mock signals injection onto the real *Fermi* data.

Figure 4 shows the results of this test, with the different rows corresponding to different signal configuration—purely DM, bright PSs, medium-bright PSs, and dim PSs. Bright, medium-bright, and dim PS configurations are taken to peak at 20, 10, and 5 photon counts respec-

tively. The leftmost columns shows the fiducial analysis on *Fermi* data, with subsequent columns showing signals of progressively larger sizes injected onto the data, up to approximately the size of the original GCE signal. The dotted horizontal and vertical lines show the expected total emission on top of the median fluxes for the PS and DM components of the GCE inferred without any additional injected signal, respectively.

The additional injected signal is seen to be reconstructed correctly within the inferred 95% confidence interval in all four cases. For the DM signal (top row), the brightest tested DM signal is seen to partially reconstruct as PS-like, which could be attributed to the larger magnitude of Poisson fluctuations in this case mimicking the effect of an unresolved PS population. The injected PS signals (rows 2-4) are correctly reconstructed in all cases, with the dimmer PS signals showing a more prominent flat direction with Poissonian emission, as expected.

B. Systematic variations on the analysis

We explore several systematic variations on the fiducial analysis, shown in Fig. 5.

- *Variation of the diffuse foreground model:* In addition to diffuse Model O considered in the fiducial analysis, we consider Models A and F from Ref. [11] to model the diffuse foreground emission. Results for these variations are shown in the left panel of Fig. 5 for Models A (blue) and F (green) compared to the fiducial Model O (red). Although the overall GCE flux is seen to varying by up to a factor of ~ 2 between diffuse models, the overall conclusion regarding the relative amounts of flux attributed to PS-like and Poissonian emission remained unchanged, with similar proportions of the GCE attributed to each component. A slight preference for smooth emission is seen when using Models A and F, and the results between these two models and broadly similar.
- *Variations on the ROI size:* Although the GCE signal is concentrated predominantly in the inner 10° , the use of the larger 25° ROI in this work is motivated by the fact that a larger region may better constrain various spatially-extended modeled emission components. On the other hand, there is also the potential for more susceptibility to mismodeling effects when using a larger ROI. The right panel of Fig. 5 shows analysis results using smaller ROI sizes— 10° , 15° and 20° . These are seen to be completely consistent with the fiducial analysis in the 25° ROI, with a wider posterior in the smaller ROIs as expected since these contain less information than fiducial ROI.

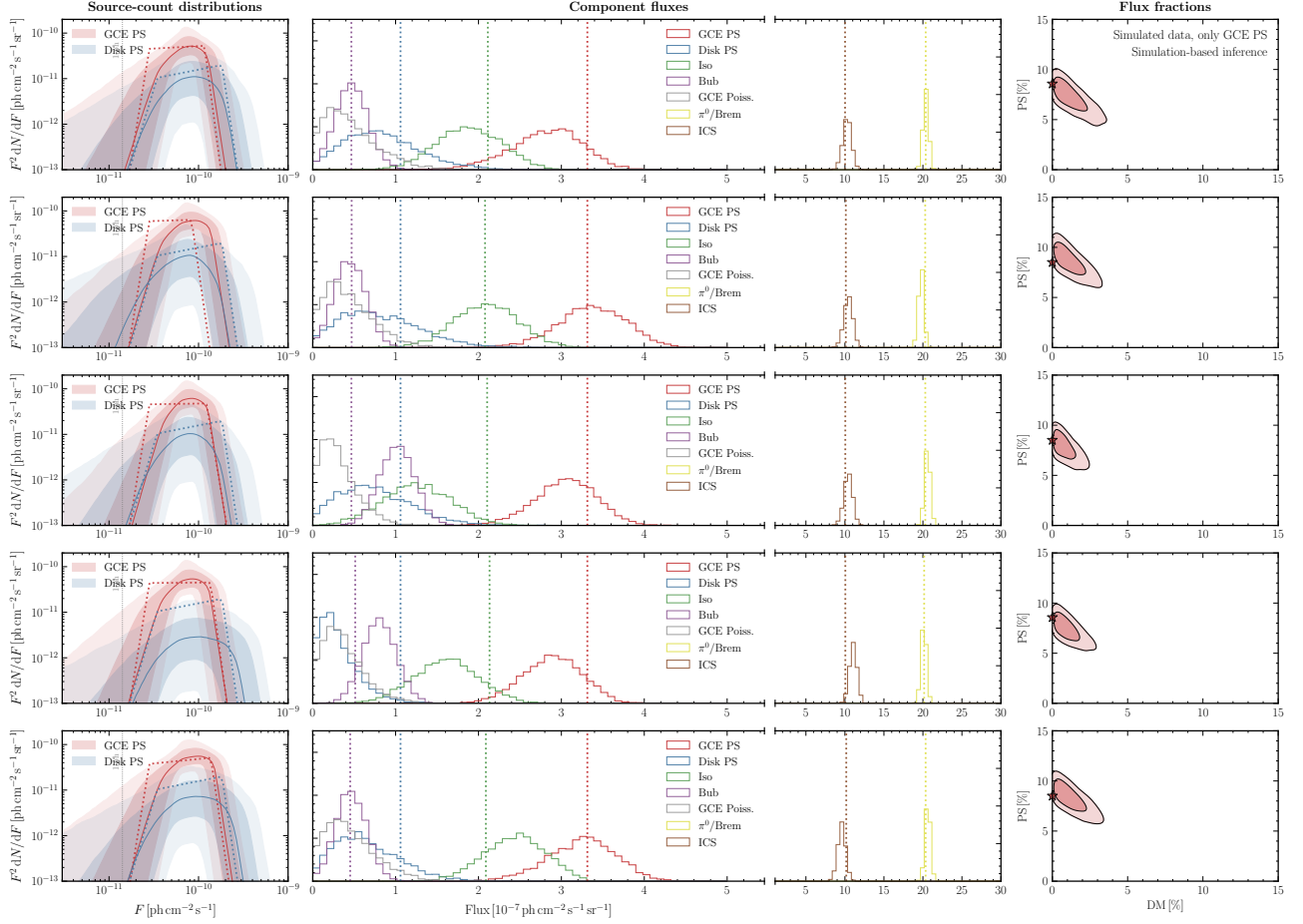


FIG. 2. Same as Fig. 1, but for simulated data where the simulated GCE consists of purely PS-like emission. PS-like emission is inferred in each case, with the other posteriors corresponding well to their true simulated values.

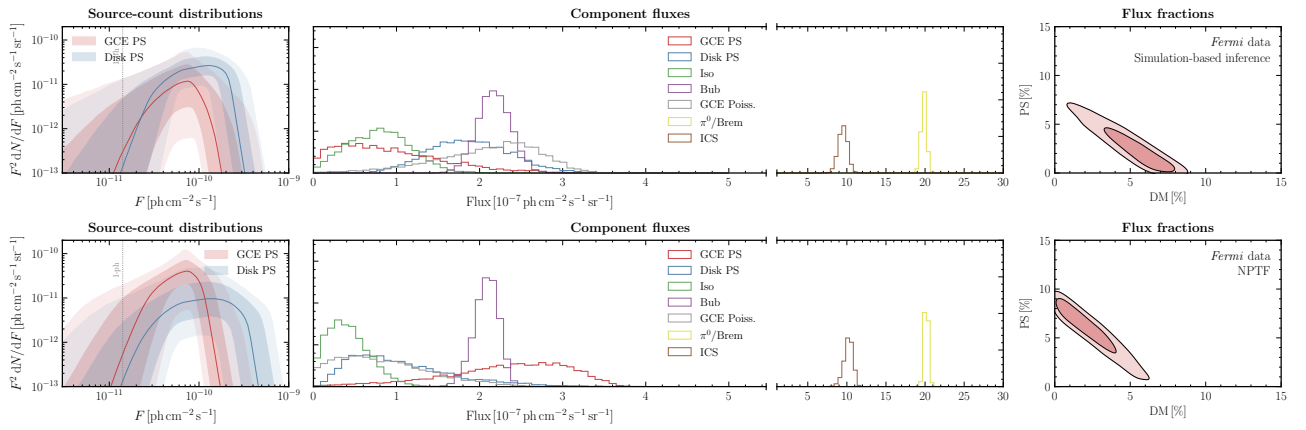


FIG. 3. Results of the fiducial analysis on data. (Top) Analysis using simulation-based inference with normalizing flows, and (bottom) using the 1-point PDF likelihood implemented in the non-Poissonian template fitting (NPTF) framework. While moderate preference for a PS-like origin of the GCE is seen in the case of the NPTF analysis, the simulation-based inference analysis finds a nearly perfect degeneracy between smooth and PS-like emission, with about half of the GCE emission attributed to each component.

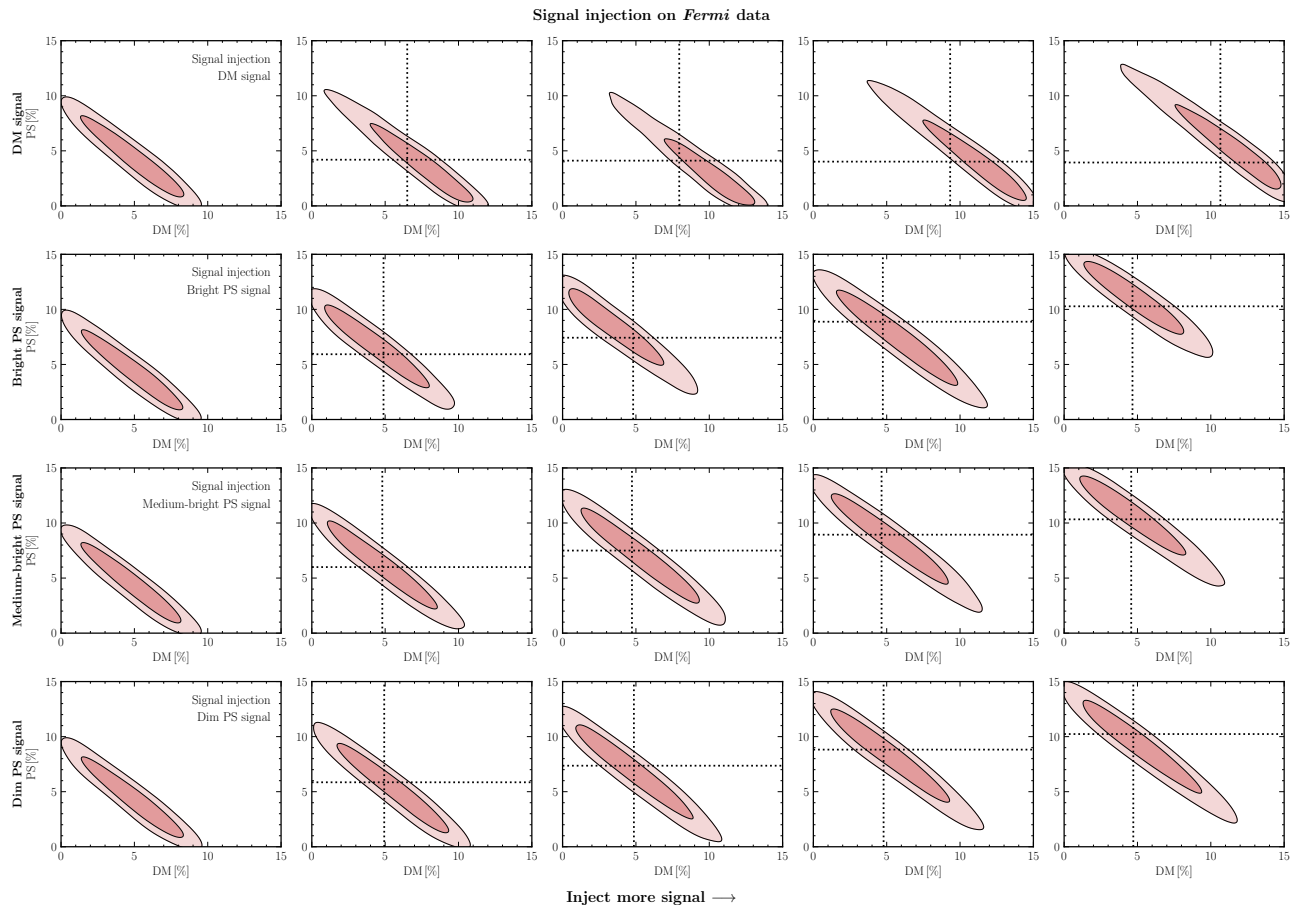


FIG. 4. Joint posterior for flux fraction of PS-like and DM-like emission when an artificial DM signal is injected onto the real data. The different rows correspond to different signal types, from top to bottom, purely DM, dim PSs (peaking at 5 expected counts per PS), moderately-bright PSs (peaking at 10 expected counts per PS), and bright PSs (peaking at 20 expected counts per PS). The leftmost panels show the fiducial analysis on *Fermi* data, with subsequent panels showing results with progressively larger signals injected onto the data. The dotted lines show the expected total emission on top of the median initial inferred flux. The additional DM and PS signals are seen to be reconstructed into the correct components in all cases.

V. SUSCEPTIBILITY TO MISMODELING

A key challenge in γ -ray analyses of the Galactic Center analyses is that associated with effects of mismodeled signal and background templates [26, 28–31]. In this section we assess the susceptibility of our simulation-based inference pipeline to these systematics, exploring the effect of background and signal mismodeling in turn.

Diffuse foreground mismodeling

In order to test the effect of large-scale foreground mismodeling, we construct a data-driven model of such mismodeling and assess the ability of our method to recover either smooth or PS-like emission in the face of such mismodeling. Following Ref. [77], we perform a Poissonian template analysis on the *Fermi* dataset x , modulating the diffuse model template T_{dif} as described by the bremsstrahlung and neutral pion decay components of

diffuse Model O by an (exponentiated) Gaussian process (GP):

$$x \sim \text{Pois} \left(\sum_{i \neq \text{dif}} A_i T_i + \exp(f) A_{\text{dif}} T_{\text{dif}} \right). \quad (10)$$

The other Poissonian templates T_i , including a GCE DM template and the Inverse Compton component of the diffuse foreground model, are treated as before using an overall normalization factor A_i . $f \sim \mathcal{N}(m, K)$ is the GP component with mean m set to zero, and the covariance K described through the Matérn kernel with smoothness parameter $\nu = 5/2$. We refer to Ref. [77] for further details of the analysis, as well a validation of the GP-augmented template fitting pipeline on simulated data.

The median Gaussian process describing the multiplicative mismodeling obtained relative to the real *Fermi* data when using our fiducial diffuse Model O is shown in Fig. 6. The most severe mismodeling is inferred to be concentrated in the central and southern regions of

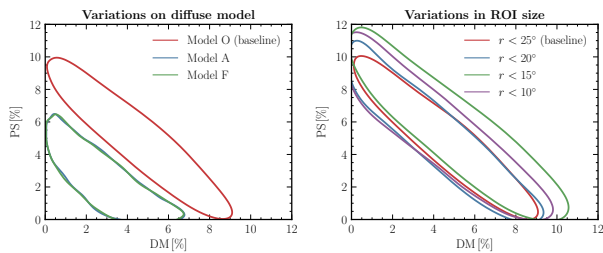


FIG. 5. Joint posterior for flux fraction of PS-like and DM-like emission on real *Fermi* data for different diffuse models (*left*) and ROI sizes (*right*). In varying diffuse models, the fiducial Model O (red) is compared with results obtained using Models A (blue) and F (green). Although the overall GCE flux is seen to varying by up to a factor of ~ 2 between diffuse models, no evidence for PS-like emission is seen. As seen in the right panel, results remain consistent for smaller ROI sizes.

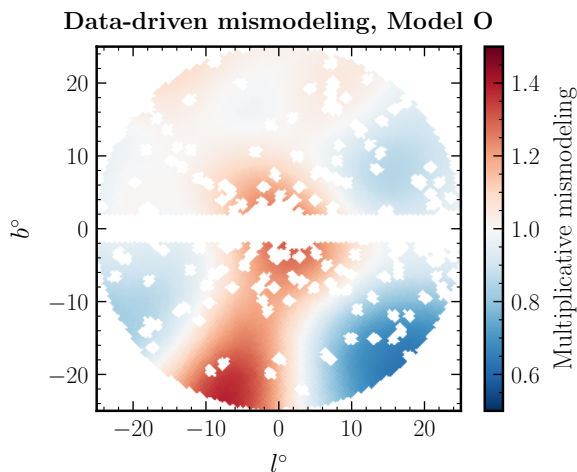


FIG. 6. The median Gaussian process description of multiplicative mismodeling associated with diffuse foreground Model O when applied to the real *Fermi* data.

the fiducial ROI. We modulate the bremsstrahlung and neutral pion decay-tracing components of Model O using samples drawn from the inferred Gaussian process, producing simulated data with the aim of mocking up the scenario of large-scale diffuse mismodeling. These simulated samples are then analyzed with our standard pipeline, using the unmodulated Model O to model the diffuse emission.

The results of this test are shown in Figs. 7 and 8, for simulated samples consisting of purely DM-like and PS-like emission, respectively. It can be seen that while large-scale mismodeling can distort the total flux attributed to either DM or PS-like emission, preference for the true underlying nature of the signal remains robust in either case. The marginalized PS flux as well as the inferred SCD is consistent with the underlying truth in all cases. The DM flux tends to be overestimated in either case however, which may be attributed to the centrally-

concentrated mismodeling as seen in Fig. 6. This is also reflected in the fact that the inferred inverse Compton component flux tends to be underestimated, with the residual flux attributable to the DM template. [SM: Specify what kind of PSs are injected.]

Signal mismodeling

Besides mismodeling of the diffuse foreground emission, another major potential concern when characterizing the GCE PS population is associated with mismodeling of the signal emission itself. In particular, as pointed out in Refs. [30, 31], a North-South asymmetry in a putative dark matter signal, if unaccounted for, could lead to the inference of a spurious PS population associated with the purely smooth, asymmetric signal in the traditional NPTF framework. Refs. [30, 31] found preference for such a scenario in real *Fermi* data, with the signal in the Northern hemisphere a factor of ~ 2 larger than that in the Southern hemisphere if the signal in the two regions is floated separately. [SM: Explain why this happens.]

We test the impact of a North-South-asymmetric dark matter signal within our framework by running our pipeline on simulated datasets where the dark matter-like signal in the Northern hemisphere of the ROI is 3 times larger than that in the Southern hemisphere, larger than the ~ 2 preference in real data found in Refs. [30, 31]. The results of this test on 3 such simulated realizations is shown in Fig. 9. We see that the presence of a substantially asymmetric DM signal has only a marginal impact on the inferred posteriors, and does not lead to a spurious preference for a PS population as was found in Refs. [30, 31] in the NPTF framework. We attribute this to the fact that the *DeepSphere*-based feature extractor can account for pixel-to-pixel correlations in the γ -ray counts map, and can thus be sensitive to *local* PS-like structures. In contrast, the 1-point PDF-based NPTF framework, being agnostic to the ordering of the pixels, can notice spurious PS-like structures in the “residuals” associated with an asymmetric signal when analyzed with a symmetric template. As done in Ref. [28], we emphasize that the presence of an asymmetry in the GCE signal, if not attributed to diffuse mismodeling, would point towards astrophysical explanations of the GCE since a true dark matter signal would not be expected to be significantly asymmetric.

VI. DISCUSSION AND CONCLUSIONS

In this work, we have leveraged recent advances in neural simulation-based inference in order to characterize a putative point source population that may be responsible for the observed *Fermi* Galactic Center Excess. Consistent with Ref [37] which used a Bayesian neural networks and first leveraged the *DeepSphere* graph-based network, our

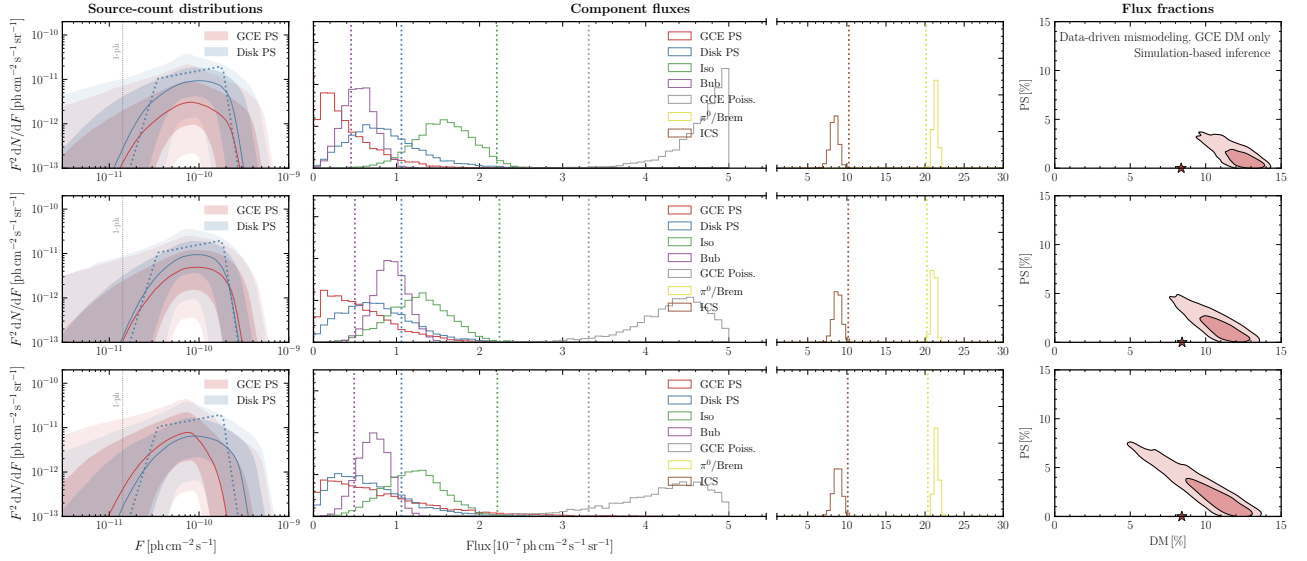


FIG. 7. Same as Fig. 1, but for simulated data where the GCE consists of purely DM-like emission and the diffuse model is modulated by draws from the Gaussian process description of diffuse mismodeling. DM-like emission is inferred in each case, although the magnitude of emission is overestimated as some of the diffuse mismodeling is absorbed into the Poissonian GCE component.

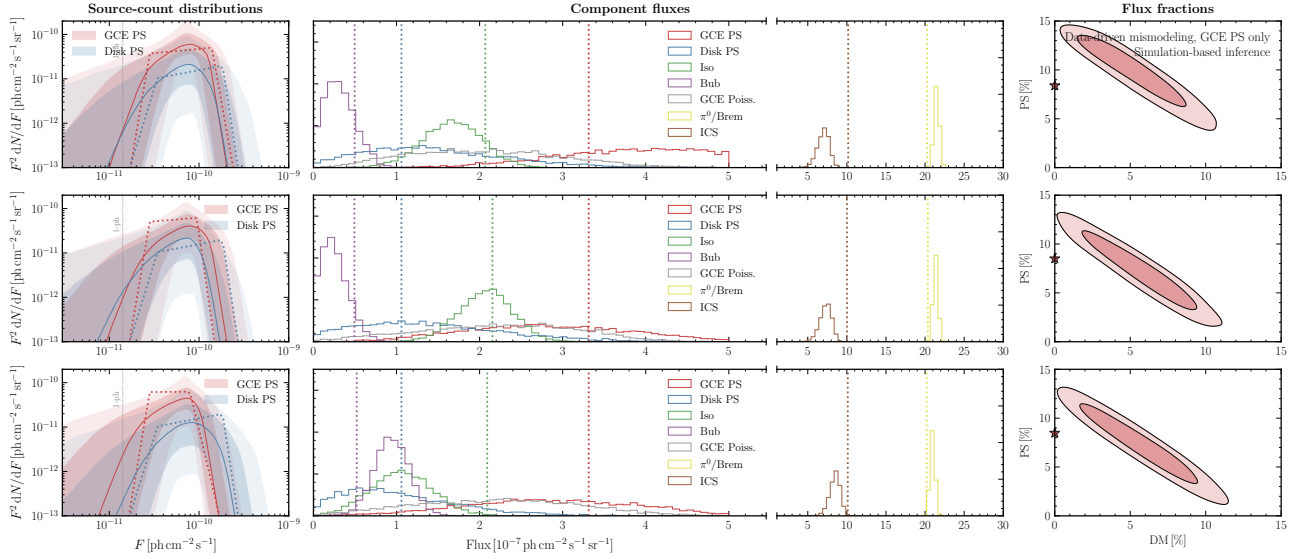


FIG. 8. Same as Fig. 1, but for simulated data where the GCE consists of a PS population peaking at 10 expected photon counts and the diffuse model is modulated by draws from the Gaussian process description of diffuse mismodeling. PS-like emission is inferred in each case, although a non-trivial DM component is inferred as some of the diffuse mismodeling is absorbed into the Poissonian GCE component.

analysis based on conditional posterior density estimation with normalizing flows shows a significantly reduced preference for a γ -ray PS population as the explanation for the GCE compared to previous analyses based on the 1-point PDF or a wavelet decomposition of the Galactic Center photon map. In particular, we find that roughly half of the GCE emission is attributable to a PS population, with the inferred source-count distribution peaking at significantly lower expected photon counts than those

found in previous analyses based on the NPTF framework [26], where the SCD is seen to peak just below the threshold for resolution of individual PSs. We have additionally shown our framework to be robust to large-scale diffuse foreground and signal mismodeling of the kind previously discussed in the literature as potential sources of bias. We have used a novel Gaussian Process-based method to construct a data-driven model of large-scale mismodeling. Our conclusions are also robust to the

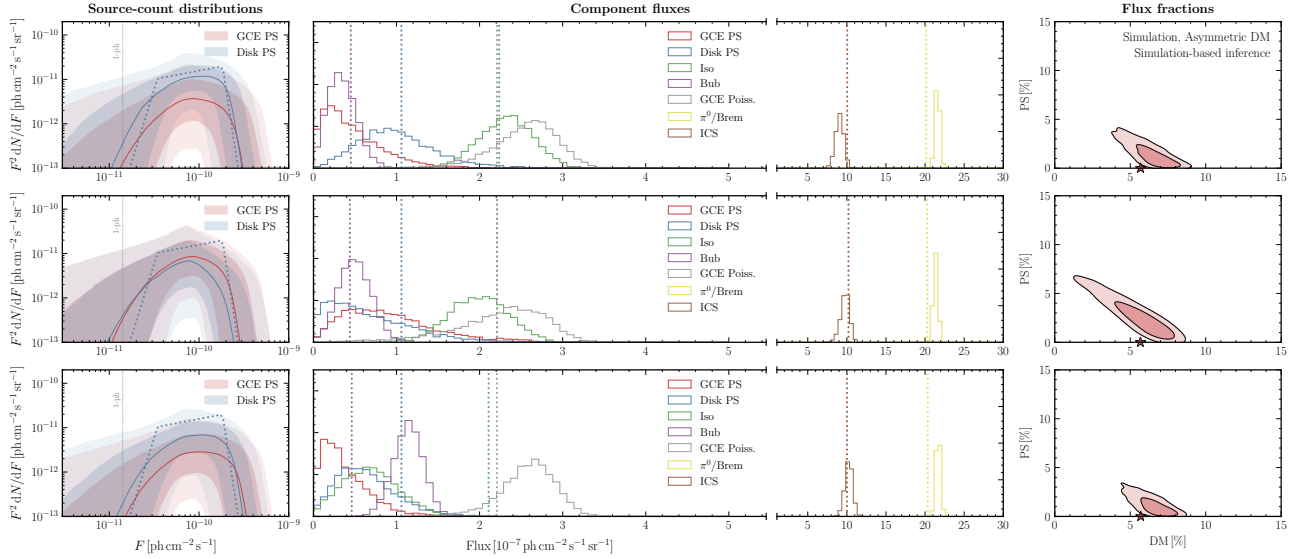


FIG. 9. Same as Fig. 1, but for simulated data where the GCE consists of purely DM-like emission with a North-South asymmetry; the signal in the Northern hemisphere is larger by a factor of 3. The mismodeled signal is seen to have marginal qualitative effect on the recovery of DM-like emission.

systematic variations we have explored, including variations on the the diffuse foreground model and size of the region of interest. As in any Galactic Center γ -ray analysis, given the poorly understood astrophysical emission in this region, we caution of the potential of unknown systematics, such as small-scale mismodeling on the scale of the size of the LAT point-spread function, to bias the results of our analysis.

Several improvements to the framework presented here are possible. The inclusion of energy-binning information in the analysis can be implemented by splitting up the data and template maps into individual bins and inputting these as separate channels in the graph-convolutional feature extraction neural network. The use of more complex feature extraction and flow architectures can additionally improve the robustness of our results. While we have considered a simulated-based inference framework based on posterior density estimation with normalizing flows, alternative frameworks based on likelihood-ratio estimation [55–57, 59, 63] or flow-based likelihood estimation [78, 79] can provide complementary ways to characterize the γ -ray PS population in the Galactic Center. Additionally, the use of sequential active-learning methods [79] and methods that extract additional latent information from the simulator [55–57, 80, 81] can significantly improve the sample efficiency of the analysis and allow for extensions to more complex latent spaces, in particular for an energy-binned analysis. Since diffuse mismodeling is the largest source of uncertainty in any analysis that aims to characterize the GCE, we also note the possibility of using adversarial learning methods [82] to account for systematic differences between the modeled and real *Fermi* data. Alternatively, generative modeling of the diffuse foreground either in a

GP-based data-driven framework or using, *e.g.*, autoencoders trained on an ensemble of plausible diffuse model scenarios, can provide a principled way to account for the large latent space associated with the diffuse foreground. [SM: Cite examples.] These extensions can lead to a more robust characterization of a putative PS population in the GCE, and we leave their study to future work.

The code used to obtain the results in this paper is available at <https://github.com/smsharma/fermi-flows>.

[SM: Beef up attribution to previous papers on different topics.]

ACKNOWLEDGMENTS

We thank Johann Brehmer, Florian List, Nick Rodd, and Tracy Slatyer for helpful conversations. KC is partially supported by NSF awards ACI-1450310, OAC-1836650, and OAC-1841471, the NSF grant PHY-1505463, and the Moore-Sloan Data Science Environment at NYU. SM is supported by the NSF CAREER grant PHY-1554858, NSF grants PHY-1620727 and PHY-1915409, and the Simons Foundation. This work made use of the NYU IT High Performance Computing resources, services, and staff expertise. This research has made use of NASA’s Astrophysics Data System. This research made use of the *astropy* [83, 84], *dynesty* [85], *IPython* [86], *Jupyter* [87], *matplotlib* [88], *MLflow* [89], *NPTFit* [43], *NumPy* [90], *pandas* [91], *Pyro* [92], *PyTorch* [93], *PyTorch Geometric* [94], *PyTorch Lightning* [95], *seaborn* [96], *sbi* [97], *scikit-learn* [98], *SciPy* [99], and *tqdm* [100]

software packages. We acknowledge the use of the code

repository associated with Ref. [37], in particular the associated data products and templates.²

-
- [1] W. B. Atwood *et al.* (Fermi-LAT), *Astrophys. J.* **697**, 1071 (2009), [arXiv:0902.1089 \[astro-ph.IM\]](#).
- [2] L. Goodenough and D. Hooper, (2009), [arXiv:0910.2998 \[hep-ph\]](#).
- [3] D. Hooper and L. Goodenough, *Phys. Lett. B* **697**, 412 (2011), [arXiv:1010.2752 \[hep-ph\]](#).
- [4] A. Boyarsky, D. Malyshev, and O. Ruchayskiy, *Phys. Lett. B* **705**, 165 (2011), [arXiv:1012.5839 \[hep-ph\]](#).
- [5] D. Hooper and T. Linden, *Phys. Rev. D* **84**, 123005 (2011), [arXiv:1110.0006 \[astro-ph.HE\]](#).
- [6] K. N. Abazajian and M. Kaplinghat, *Phys. Rev. D* **86**, 083511 (2012), [Erratum: *Phys.Rev.D* 87, 129902 (2013)], [arXiv:1207.6047 \[astro-ph.HE\]](#).
- [7] D. Hooper and T. R. Slatyer, *Phys. Dark Univ.* **2**, 118 (2013), [arXiv:1302.6589 \[astro-ph.HE\]](#).
- [8] C. Gordon and O. Macias, *Phys. Rev. D* **88**, 083521 (2013), [Erratum: *Phys.Rev.D* 89, 049901 (2014)], [arXiv:1306.5725 \[astro-ph.HE\]](#).
- [9] K. N. Abazajian, N. Canac, S. Horiuchi, and M. Kaplinghat, *Phys. Rev. D* **90**, 023526 (2014), [arXiv:1402.4090 \[astro-ph.HE\]](#).
- [10] T. Daylan, D. P. Finkbeiner, D. Hooper, T. Linden, S. K. N. Portillo, N. L. Rodd, and T. R. Slatyer, *Phys. Dark Univ.* **12**, 1 (2016), [arXiv:1402.6703 \[astro-ph.HE\]](#).
- [11] F. Calore, I. Cholis, and C. Weniger, *JCAP* **03**, 038 (2015), [arXiv:1409.0042 \[astro-ph.CO\]](#).
- [12] K. N. Abazajian, N. Canac, S. Horiuchi, M. Kaplinghat, and A. Kwa, *JCAP* **07**, 013 (2015), [arXiv:1410.6168 \[astro-ph.HE\]](#).
- [13] M. Ajello *et al.* (Fermi-LAT), *Astrophys. J.* **819**, 44 (2016), [arXiv:1511.02938 \[astro-ph.HE\]](#).
- [14] T. Linden, N. L. Rodd, B. R. Safdi, and T. R. Slatyer, *Phys. Rev. D* **94**, 103013 (2016), [arXiv:1604.01026 \[astro-ph.HE\]](#).
- [15] O. Macias, C. Gordon, R. M. Crocker, B. Coleman, D. Paterson, S. Horiuchi, and M. Pohl, *Nature Astron.* **2**, 387 (2018), [arXiv:1611.06644 \[astro-ph.HE\]](#).
- [16] H. A. Clark, P. Scott, R. Trotta, and G. F. Lewis, *JCAP* **07**, 060 (2018), [arXiv:1612.01539 \[astro-ph.HE\]](#).
- [17] K. N. Abazajian, *JCAP* **03**, 010 (2011), [arXiv:1011.4275 \[astro-ph.HE\]](#).
- [18] D. Hooper, I. Cholis, T. Linden, J. Siegal-Gaskins, and T. Slatyer, *Phys. Rev. D* **88**, 083009 (2013), [arXiv:1305.0830 \[astro-ph.HE\]](#).
- [19] F. Calore, M. Di Mauro, and F. Donato, *Astrophys. J.* **796**, 1 (2014), [arXiv:1406.2706 \[astro-ph.HE\]](#).
- [20] I. Cholis, D. Hooper, and T. Linden, *JCAP* **06**, 043 (2015), [arXiv:1407.5625 \[astro-ph.HE\]](#).
- [21] J. Petrović, P. D. Serpico, and G. Zaharijas, *JCAP* **02**, 023 (2015), [arXiv:1411.2980 \[astro-ph.HE\]](#).
- [22] Q. Yuan and K. Ioka, *Astrophys. J.* **802**, 124 (2015), [arXiv:1411.4363 \[astro-ph.HE\]](#).
- [23] T. D. Brandt and B. Kocsis, *Astrophys. J.* **812**, 15 (2015), [arXiv:1507.05616 \[astro-ph.HE\]](#).
- [24] O. Macias, S. Horiuchi, M. Kaplinghat, C. Gordon, R. M. Crocker, and D. M. Nataf, *JCAP* **09**, 042 (2019), [arXiv:1901.03822 \[astro-ph.HE\]](#).
- [25] R. Bartels, E. Storm, C. Weniger, and F. Calore, *Nature Astron.* **2**, 819 (2018), [arXiv:1711.04778 \[astro-ph.HE\]](#).
- [26] S. K. Lee, M. Lisanti, B. R. Safdi, T. R. Slatyer, and W. Xue, *Phys. Rev. Lett.* **116**, 051103 (2016), [arXiv:1506.05124 \[astro-ph.HE\]](#).
- [27] R. Bartels, S. Krishnamurthy, and C. Weniger, *Phys. Rev. Lett.* **116**, 051102 (2016), [arXiv:1506.05104 \[astro-ph.HE\]](#).
- [28] M. Buschmann, N. L. Rodd, B. R. Safdi, L. J. Chang, S. Mishra-Sharma, M. Lisanti, and O. Macias, *Phys. Rev. D* **102**, 023023 (2020), [arXiv:2002.12373 \[astro-ph.HE\]](#).
- [29] L. J. Chang, S. Mishra-Sharma, M. Lisanti, M. Buschmann, N. L. Rodd, and B. R. Safdi, *Phys. Rev. D* **101**, 023014 (2020), [arXiv:1908.10874 \[astro-ph.CO\]](#).
- [30] R. K. Leane and T. R. Slatyer, *Phys. Rev. Lett.* **125**, 121105 (2020), [arXiv:2002.12370 \[astro-ph.HE\]](#).
- [31] R. K. Leane and T. R. Slatyer, *Phys. Rev. D* **102**, 063019 (2020), [arXiv:2002.12371 \[astro-ph.HE\]](#).
- [32] R. K. Leane and T. R. Slatyer, *Phys. Rev. Lett.* **123**, 241101 (2019), [arXiv:1904.08430 \[astro-ph.HE\]](#).
- [33] S. K. Lee, M. Lisanti, and B. R. Safdi, *JCAP* **05**, 056 (2015), [arXiv:1412.6099 \[astro-ph.CO\]](#).
- [34] B. Balaji, I. Cholis, P. J. Fox, and S. D. McDermott, *Phys. Rev. D* **98**, 043009 (2018), [arXiv:1803.01952 \[astro-ph.HE\]](#).
- [35] S. D. McDermott, P. J. Fox, I. Cholis, and S. K. Lee, *JCAP* **07**, 045 (2016), [arXiv:1512.00012 \[astro-ph.HE\]](#).
- [36] S. Caron, K. Dijkstra, C. Eckner, L. Hendriks, G. Jóhannesson, B. Panes, R. Ruiz De Austri, and G. Zaharijas, (2021), [arXiv:2103.11068 \[astro-ph.HE\]](#).
- [37] F. List, N. L. Rodd, G. F. Lewis, and I. Bhat, *Phys. Rev. Lett.* **125**, 241102 (2020), [arXiv:2006.12504 \[astro-ph.HE\]](#).
- [38] S. Caron, G. A. Gómez-Vargas, L. Hendriks, and R. Ruiz de Austri, *JCAP* **05**, 058 (2018), [arXiv:1708.06706 \[astro-ph.HE\]](#).
- [39] K. Cranmer, J. Brehmer, and G. Louppe, *Proceedings of the National Academy of Sciences* **117**, 30055 (2020).
- [40] G. Papamakarios, E. Nalisnick, D. J. Rezende, S. Mohamed, and B. Lakshminarayanan, *arXiv preprint arXiv:1912.02762* (2019).
- [41] D. Rezende and S. Mohamed, in *International Conference on Machine Learning* (PMLR, 2015) pp. 1530–1538.
- [42] S. Mishra-Sharma, N. L. Rodd, and B. R. Safdi, “Supplementary material for NPTFit,” (2016).
- [43] S. Mishra-Sharma, N. L. Rodd, and B. R. Safdi, *Astron. J.* **153**, 253 (2017), [arXiv:1612.03173 \[astro-ph.HE\]](#).
- [44] K. M. Gorski, E. Hivon, A. J. Banday, B. D. Wandelt, F. K. Hansen, M. Reinecke, and M. Bartelman, *Astrophys. J.* **622**, 759 (2005), [arXiv:astro-ph/0409513](#).

² https://github.com/FloList/GCE_NN

- [45] J. F. Navarro, C. S. Frenk, and S. D. M. White, *Astrophys. J.* **462**, 563 (1996), [arXiv:astro-ph/9508025 \[astro-ph\]](#).
- [46] J. F. Navarro, C. S. Frenk, and S. D. White, *Astrophys. J.* **490**, 493 (1997), [arXiv:astro-ph/9611107 \[astro-ph\]](#).
- [47] B. Zhou, Y.-F. Liang, X. Huang, X. Li, Y.-Z. Fan, L. Feng, and J. Chang, *Phys. Rev. D* **91**, 123010 (2015), [arXiv:1406.6948 \[astro-ph.HE\]](#).
- [48] D. R. Lorimer *et al.*, *Mon. Not. Roy. Astron. Soc.* **372**, 777 (2006), [arXiv:astro-ph/0607640 \[astro-ph\]](#).
- [49] R. T. Bartels, T. D. P. Edwards, and C. Weniger, *Mon. Not. Roy. Astron. Soc.* **481**, 3966 (2018), [arXiv:1805.11097 \[astro-ph.HE\]](#).
- [50] N. L. Rodd and M. W. Toomey, *NPTFit-Sim* (2017).
- [51] F. Acero *et al.* (Fermi-LAT), *Astrophys. J. Suppl.* **218**, 23 (2015), [arXiv:1501.02003 \[astro-ph.HE\]](#).
- [52] M. Su, T. R. Slatyer, and D. P. Finkbeiner, *Astrophys. J.* **724**, 1044 (2010), [arXiv:1005.5480 \[astro-ph.HE\]](#).
- [53] D. B. Rubin, *The Annals of Statistics* **12**, 1151 (1984).
- [54] J. Alsing, T. Charnock, S. Feeney, and B. Wandelt, *Monthly Notices of the Royal Astronomical Society*, [stz1960](#) (2019), [arXiv: 1903.00007](#).
- [55] J. Brehmer, K. Cranmer, G. Louppe, and J. Pavez, *Physical Review D* **98**, 052004 (2018), [arXiv: 1805.00020](#).
- [56] J. Brehmer, G. Louppe, J. Pavez, and K. Cranmer, *Proceedings of the National Academy of Sciences* **117**, 5242 (2020), [arXiv: 1805.12244](#).
- [57] J. Brehmer, K. Cranmer, G. Louppe, and J. Pavez, *Physical Review Letters* **121**, 111801 (2018), [arXiv: 1805.00013](#).
- [58] J. Brehmer and K. Cranmer, [arXiv:2010.06439 \[hep-ex, physics:hep-ph, physics:physics, stat\]](#) (2020), [arXiv: 2010.06439](#).
- [59] K. Cranmer, J. Pavez, and G. Louppe, [arXiv:1506.02169 \[physics, stat\]](#) (2016), [arXiv: 1506.02169](#).
- [60] K. Cranmer, J. Brehmer, and G. Louppe, [arXiv:1911.01429 \[cs, stat\]](#) (2020), [arXiv: 1911.01429](#).
- [61] C. Durkan, I. Murray, and G. Papamakarios, [arXiv:2002.03712 \[cs, stat\]](#) (2020), [arXiv: 2002.03712](#).
- [62] D. S. Greenberg, M. Nonnenmacher, and J. H. Macke, [arXiv:1905.07488 \[cs, stat\]](#) (2019), [arXiv: 1905.07488](#).
- [63] J. Hermans, V. Begy, and G. Louppe, [arXiv:1903.04057 \[cs, stat\]](#) (2020), [arXiv: 1903.04057](#).
- [64] J.-M. Lueckmann, J. Boelts, D. S. Greenberg, P. J. Gonçalves, and J. H. Macke, [arXiv:2101.04653 \[cs, stat\]](#) (2021), [arXiv: 2101.04653](#).
- [65] J.-M. Lueckmann, G. Bassetto, T. Karaletsos, and J. H. Macke, [arXiv:1805.09294 \[cs, stat\]](#) (2019), [arXiv: 1805.09294](#).
- [66] L. Pacchiardi and R. Dutta, [arXiv:2104.03889 \[stat\]](#) (2021), [arXiv: 2104.03889](#).
- [67] G. Papamakarios and I. Murray, [arXiv:1605.06376 \[cs, stat\]](#) (2018), [arXiv: 1605.06376](#).
- [68] G. Papamakarios, D. C. Sterratt, and I. Murray, [arXiv:1805.07226 \[cs, stat\]](#) (2019), [arXiv: 1805.07226](#).
- [69] S. Wqvist, J. Frellsen, and U. Picchini, [arXiv:2102.06522 \[cs, stat\]](#) (2021), [arXiv: 2102.06522](#).
- [70] D. Zhao, N. Dalmaso, R. Izbicki, and A. B. Lee, (2021).
- [71] G. Papamakarios, T. Pavlakou, and I. Murray, in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17 (Curran Associates Inc., Red Hook, NY, USA, 2017) pp. 2335–2344.
- [72] M. Germain, K. Gregor, I. Murray, and H. Larochelle, in *International Conference on Machine Learning* (PMLR, 2015) pp. 881–889.
- [73] M. Defferrard, M. Milani, F. Gusset, and N. Perraudin, [arXiv preprint arXiv:2012.15000](#) (2020).
- [74] N. Perraudin, M. Defferrard, T. Kacprzak, and R. Sgier, *Astron. Comput.* **27**, 130 (2019), [arXiv:1810.12186 \[astro-ph.CO\]](#).
- [75] D. P. Kingma and J. Ba, in *ICLR (Poster)* (2015).
- [76] I. Loshchilov and F. Hutter, in *International Conference on Learning Representations* (2019).
- [77] S. Mishra-Sharma and K. Cranmer, in *34th Conference on Neural Information Processing Systems* (2020) [arXiv:2010.10450 \[astro-ph.HE\]](#).
- [78] C. Winkler, D. Worrall, E. Hooeboom, and M. Welling, [arXiv preprint arXiv:1912.00042](#) (2019).
- [79] G. Papamakarios, D. Sterratt, and I. Murray, in *The 22nd International Conference on Artificial Intelligence and Statistics* (PMLR, 2019) pp. 837–848.
- [80] J. Brehmer, F. Kling, I. Espejo, and K. Cranmer, *Comput. Softw. Big Sci.* **4**, 3 (2020), [arXiv:1907.10621 \[hep-ph\]](#).
- [81] M. Stoye, J. Brehmer, G. Louppe, J. Pavez, and K. Cranmer, (2018), [arXiv:1808.00973 \[stat.ML\]](#).
- [82] G. Louppe, M. Kagan, and K. Cranmer, (2016), [arXiv:1611.01046 \[stat.ML\]](#).
- [83] A. M. Price-Whelan *et al.*, *Astron. J.* **156**, 123 (2018), [arXiv:1801.02634](#).
- [84] T. P. Robitaille *et al.* (Astropy), *Astron. Astrophys.* **558**, A33 (2013), [arXiv:1307.6212 \[astro-ph.IM\]](#).
- [85] J. S. Speagle, *Monthly Notices of the Royal Astronomical Society* **493**, 3132 (2020).
- [86] F. Perez and B. E. Granger, *Computing in Science and Engineering* **9**, 21 (2007).
- [87] T. Kluyver *et al.*, in *ELPUB* (2016).
- [88] J. D. Hunter, *Computing In Science & Engineering* **9**, 90 (2007).
- [89] A. Chen, A. Chow, A. Davidson, A. DCunha, A. Ghodsi, S. A. Hong, A. Konwinski, C. Mewald, S. Murching, T. Nykodym, *et al.*, in *Proceedings of the fourth international workshop on data management for end-to-end machine learning* (2020) pp. 1–4.
- [90] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant, *Nature* **585**, 357 (2020).
- [91] W. McKinney, in *Proceedings of the 9th Python in Science Conference*, edited by S. van der Walt and J. Millman (2010) pp. 51 – 56.
- [92] E. Bingham, J. P. Chen, M. Jankowiak, F. Obermeyer, N. Pradhan, T. Karaletsos, R. Singh, P. Szerlip, P. Horsfall, and N. D. Goodman, *Journal of Machine Learning Research* (2018).
- [93] A. Paszke *et al.*, in *Advances in Neural Information Processing Systems 32*, edited by H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Curran Associates, Inc., 2019) pp. 8024–8035.
- [94] M. Fey and J. E. Lenssen, in *ICLR Workshop on Representation Learning on Graphs and Manifolds* (2019).

- [95] W. Falcon *et al.*, “Pytorchlightning/pytorch-lightning: 0.7.6 release,” (2020).
- [96] M. Waskom *et al.*, “mwaskom/seaborn: v0.8.1 (september 2017),” (2017).
- [97] A. Tejero-Cantero *et al.*, *Journal of Open Source Software* **5**, 2505 (2020).
- [98] F. Pedregosa *et al.*, *Journal of Machine Learning Research* **12**, 2825 (2011).
- [99] P. Virtanen *et al.*, *Nature Methods* (2020), <https://doi.org/10.1038/s41592-019-0686-2>.
- [100] C. O. da Costa-Luis, *JOSS* **4**, 1277 (2019).