

SURFing to the Fundamental Limit of Jet Tagging

Ian Pang,^{1,*} Darius A. Faroughy,^{1,†} David Shih,^{1,‡} Ranit Das,^{1,2,§} and Gregor Kasieczka^{3,¶}

¹*NHETC, Dept. of Physics and Astronomy, Rutgers University, Piscataway, NJ 08854, USA*

²*Institute for Theoretical Physics, Universität Heidelberg, Germany*

³*Institut für Experimentalphysik, Universität Hamburg, 22761 Hamburg, Germany*

Beyond the practical goal of improving search and measurement sensitivity through better jet tagging algorithms, there is a deeper question: what are their upper performance limits? Generative surrogate models with learned likelihood functions offer a new approach to this problem, provided the surrogate correctly captures the underlying data distribution. In this work, we introduce the SURrogate ReFerence (SURF) method, a new approach to validating generative models. This framework enables exact Neyman–Pearson tests by training the target model on samples from another tractable surrogate, which is itself trained on real data. We argue that the EPiC-FM generative model is a valid surrogate reference for JETCLASS jets and apply SURF to show that modern jet taggers may already be operating close to the true statistical limit. By contrast, we find that autoregressive GPT models unphysically exaggerate top vs. QCD separation power encoded in the surrogate reference, implying that they are giving a misleading picture of the fundamental limit.

I. INTRODUCTION

Jet tagging is a central task in collider physics. Over the past decade, machine learning has driven major advances in jet tagging, with increasingly sophisticated architectures achieving very high classification performance on simulated datasets [1–11]. This success naturally raises a key question: *have current jet taggers already reached the fundamental limit of jet tagging, or does a gap remain between practical performance and the true statistical optimum?*

The Neyman-Pearson (NP) limit, defined by the likelihood ratio, is the best possible discriminant between two different underlying physics processes – such as top and QCD jets – that any classifier could achieve if it had access to the exact data likelihoods [12]. In practice, however, this limit cannot be evaluated directly because the true likelihood of the data-generating process is unknown. It therefore remains unclear how close existing classifiers are to this ultimate bound.

Recently, Ref. [13] proposed using autoregressive GPT-style generative models to probe this limit for top vs. QCD jets from the JETCLASS dataset [14]. These models operate on discretized, tokenized representations of jet constituents and yield explicit log-likelihoods, enabling the computation of likelihood ratios between jet classes. When applied to top–QCD discrimination, GPT-based likelihoods were found to produce receiver-operator-characteristic (ROC) curves that appear to substantially exceed those of all existing classifiers trained on GPT samples. This was interpreted as evidence that current jet taggers fall well short of the true NP limit.

In this work, we reevaluate this conclusion in several ways. First, we upgrade the GPT model of Refs. [13, 15] with a number of enhancements, including voxel tokenization, positional encoding and weight tying. Despite this, we find that the improved GPT model trained on tops and QCD from JETCLASS still inflates their separation relative to state-of-the-art classifiers by more than an order of magnitude in rejection factor over a wide range of signal of efficiencies. In fact, this is an even larger gap than the original one of Refs. [13, 15]. Along the way, we show that GPT-generated jets are perfectly separable from continuous jets, while being close to the discretized jets that the GPT models were trained on. However, we show this cannot explain the inflated separation of the GPT jets – the discretized top and QCD jets are if anything *less* separable than their continuous counterparts.

Secondly, we assess the fundamental limit of jet tagging using an alternative generative modeling framework that also gives tractable likelihoods: conditional flow matching (CFM) [16–18] with permutation equivariance, in the form of the EPiC-FM architecture [19, 20]. Unlike the GPT framework, EPiC-FM operates directly at the level of the continuous jet constituents and does not require any discretization or tokenization. Using the log-likelihood ratios derived from EPiC-FM trained on JETCLASS QCD and top jets, we find no evidence for an inflated fundamental limit: the ROC curve obtained from the exact NP test is fairly close (within a factor of two) to the state-of-the-art classifier performance on EPiC-FM samples.

These results raise the question – which jet generative model is giving the more accurate picture of the fundamental limit of jet tagging? Is there an enormous gap between the fundamental limit and our best classifiers, as the GPT models would suggest? Or is there essentially no gap between the two, as the EPiC-FM models would suggest?

To clarify the origin of this discrepancy and to provide a general framework for evaluating generative models,

* ian.pang@physics.rutgers.edu

† daruis.faroughy@rutgers.edu

‡ shih@physics.rutgers.edu

§ das@thphys.uni-heidelberg.de

¶ gregor.kasieczka@uni-hamburg.de

we introduce the SURrogate ReFerence (SURF) method. Here the key idea is to use one generative model with tractable likelihoods (EPiC-FM) to evaluate the other generative model with tractable likelihoods (GPT). The first one serves as the *surrogate reference model* for the second. By various ways we validate EPiC-FM as a surrogate reference model, checking that it does not introduce spurious artifacts into the JETCLASS data that are hard to detect.¹ We replace JETCLASS samples with samples drawn from EPiC-FM, train the GPT model on these samples, and repeat all the studies and tests we want to do with the GPT model. But now with tractable likelihoods for both the surrogate reference and the target generative model, ground truth NP-optimal ROC curves can be calculated and compared against one another, giving us more information than when we had trained the GPT samples on JETCLASS directly.

Within this framework, we find that the NP ROC curve derived from the GPT likelihoods far exceeds that of the EPiC-FM surrogate itself. This demonstrates unambiguously that the GPT model artificially inflates the apparent separation between top and QCD jets beyond the physical limit represented by the surrogate reference. Evidently, the GPT model is introducing unphysical artifacts into the EPiC-FM surrogate samples that it is trained on, leading them to be more separable than they should be. We investigate the reasons for this and find evidence that GPT models exhibit signs of overfitting to their training data. We discuss how this form of mismodeling produces likelihood ratios that are inherently difficult for classifiers to reproduce, leading to artificially inflated ROC curves.

In the end, we find no evidence of an inflated “fundamental limit” of jet tagging. Instead, the overall picture, especially the concurrence of top vs. QCD ROC curves from OmniLearn trained on JETCLASS and on EPiC-FM together with the exact likelihoods of EPiC-FM, suggests that the fundamental limit is not far away from state-of-the-art classifiers.

The outline of our paper is as follows. In Sections II and III, we present the model descriptions, baseline evaluations, and top vs. QCD separation results obtained using the GPT model and the EPiC-FM model, respectively. Section IV introduces the SURF method and its corresponding results. Possible reasons for the inflated GPT top vs. QCD curve are examined in Section V. We conclude with an outlook in Section VI, while additional technical background and log-likelihood details are provided in the appendices.

¹ The same tests fail for the GPT model, which is why we do not recommend running the SURF method in the other order, with GPT as the surrogate reference.

II. GPT: JETS AS TOKENIZED SEQUENCES

A. Description of the model

One popular generative modeling framework for jets is to view them as p_T -ordered sequences of tokenized particles. We follow the approach of [13], where the kinematic features of each constituent are discretized into bins and mapped to tokens. The resulting sequences can then be modeled autoregressively with GPT-style autoregressive language models, which learn the joint distribution over jet constituents by factorizing it into a product of conditional probabilities for each tokenized particle given the preceding ones. We note that VQ-VAE-based tokenization is also possible [21, 22]; however, obtaining a well-defined likelihood in physical space is nontrivial because the learned codebook partitions feature space into irregular regions without a known volume element. Establishing a principled mapping from sequence log-probabilities to physical-space log-densities in that setting is left to future work, so we adopt a fixed-bin vocabulary here.

Concretely, within each jet $\mathbf{x} = \{\mathbf{x}_i\}_{i=1}^N$ we order the constituents $\mathbf{x}_i = (\log p_T, \Delta\eta, \Delta\phi)_i$ by descending p_T and discretize the kinematic features into $N_{p_T} = 40$, $N_{\Delta\eta} = 30$, and $N_{\Delta\phi} = 30$ evenly spaced bins. Here $p_{T,i}$ is the transverse momenta of particle i , $\Delta\eta_i \equiv \eta_i - \eta_J$ is the pseudorapidity offset from the jet axis, and $\Delta\phi_i \equiv \phi_i - \phi_J$ is the azimuthal-angle offset. Each three-dimensional voxel in feature space is mapped to a unique token, giving rise to a vocabulary of size $V = N_{p_T} \cdot N_{\Delta\eta} \cdot N_{\Delta\phi} = 36,000$. We restrict each jet to the 40 highest- p_T particles before discretization. The pseudo-rapidity and azimuthal bins are restricted to the range $(-0.8, 0.8)$ with overflows. Each jet \mathbf{x} is thus represented by a sequence $\mathbf{t}_x = (t_1, \dots, t_N)$, where each token t_i represents an individual particle. Sequences are padded to a maximum length of $N_{\max} = 40$ using a dedicated [PAD] token. Additional special tokens, [START] and [END], are appended to the vocabulary to indicate the beginning and end of each jet sequence.

We implement our GPT model with the HuggingFace Transformers library [23]. Model and training details can be found in appendix B. When compared to Ref. [13], our setup differs in three key ways. (i) *Voxel tokenization*: we use a single, joint vocabulary over $(\log p_T, \Delta\eta, \Delta\phi)$ voxels and embed each token directly into $\mathbb{R}^{d_{\text{embd}}}$, whereas Ref. [13] tokenizes each binned feature separately and sums three feature-specific embeddings to form the per-particle representation. (ii) *Positional information*: we add learned positional embeddings (reflecting the p_T ordering) to the token embeddings before the self-attention layers; Ref. [13] omits positional encodings. (iii) *Weight tying*: we tie the last linear layer (language head) to the input token embedding matrix, using the same weights for input embedding and output logits. Weight tying reduces the number of trainable parameters and has been shown to improve generalization in language models [24, 25]. Empirically, we find these choices improve

the generative performance of our GPT model.

Drawing a sample from a trained GPT produces tokens $\mathbf{t}_x = (\mathbf{t}_0 = [\text{START}], \mathbf{t}_1, \dots, \mathbf{t}_N, \mathbf{t}_{N+1} = [\text{END}])$ which correspond to binned, p_T -ordered jet constituents. The log-likelihood of a tokenized jet is computed by summing the next-token probabilities

$$\log p_\theta(\mathbf{t}_x) = \sum_{n=1}^{N+1} \log p_\theta(\mathbf{t}_n \mid \mathbf{t}_0, \dots, \mathbf{t}_{n-1}), \quad (1)$$

along the sequence and stopping at the `[END]` symbol.

The tokenized jets can be mapped back to the physical space of continuous jet constituents by identifying each constituent token with its corresponding voxel in $(\log p_T, \Delta\eta, \Delta\phi)$ space and then smearing it within its voxel according to a uniform density. The corresponding log-likelihood in the physical space is then given by

$$\log p_\theta^{\text{cont}}(\mathbf{x}) = \log p_\theta(\mathbf{t}_x) - N \log dV + \sum_k \log N_k! - \log N! \quad (2)$$

Here the second term is the correction to the log-likelihood from the uniform smearing (dV is the volume of the voxel), and the third and fourth terms are necessary to remove the spurious effects of p_T -ordering in the sequence (N_k is the number of constituents with identical discretized p_T).

Clearly, with the GPT models there are two choices for what we consider as the “original data” – we can use either the original continuous JETCLASS samples or their bin-smeared counterparts as the baseline. Ostensibly only the former case is of interest – we are ultimately interested in generative models for continuous jets. However, the latter is also a useful comparison point, and we will consider both in the following.

B. Baseline evaluation of the model

We first perform a baseline evaluation of GPT models trained directly on JETCLASS data. For brevity, we evaluate performance solely using the AUC of a classifier trained to distinguish generated samples from real data [26, 27], although other evaluation strategies are possible [28–31].

We train an OmniLearn classifier² to separate JETCLASS data from GPT-generated jets. We find that OmniLearn separates both datasets almost perfectly (AUC = 0.9949 for QCD and 0.9734 for tops; Table I, first row). The near-perfect separability indicates that GPT samples lie far from the continuous data manifold.

Comparison	QCD	Top
JETCLASS vs GPT	0.9949 ?	0.9734 ?
JETCLASS-bs vs GPT	0.5564 ?	0.5965 ?
JETCLASS vs JETCLASS-bs	0.9940 ?	0.9714 ?

TABLE I. AUC scores for distinguishing between jets from JETCLASS, bin-smeared JETCLASS, and GPT trained on JETCLASS, using the OmniLearn classifier. Below the diagonal line the space is reserved for the true NP-optimal classifier scores, which cannot be computed because JETCLASS likelihoods are unavailable.

Is this due to an issue with the GPT models, or is it due to the binning of the JETCLASS data itself? We answer this question with two more classifier tests: GPT samples vs. bin-smeared JETCLASS, and continuous JETCLASS vs. bin-smeared JETCLASS. The first is shown in Table I, second row, indicating OmniLearn AUCs of 0.5564 and 0.5965 for QCD and tops respectively. This reproduces the good performance achieved by the original Aachen GPT-style jet generator [15] and in fact improves upon it.³ Meanwhile the latter is shown in the third row of Table I, indicating OmniLearn AUCs of 0.9940 and 0.9714 for QCD and tops respectively. We conclude that the separability of GPT jets from original continuous JETCLASS is driven entirely by the binning of the data – this alone takes the data off-manifold with respect to the original continuous samples. By contrast, GPT jets are quite close to the binned JETCLASS that they were trained on.

C. Top vs. QCD separation with GPT

We next turn to an evaluation of how well the GPT model reproduces the separation of tops vs. QCD. We compare three classifiers in Figure 2 (each classifier is evaluated on the samples from the generative process it’s defined/trained with):

1. “GPT Optimal”: The GPT model’s likelihood ratio discriminator. Note that here it doesn’t matter whether one uses the autoregressive likelihood for

² For all our results, we use OmniLearn from Refs. [8, 9] to represent current state-of-the-art classifiers. We pretrain it on the entire JETCLASS training dataset and fine-tune it on the jets used in our study. In particular, we use 1M jets from each class, with an 80/10/10 split into training, validation, and test sets.

³ Retraining OmniLearn on bin-smeared Aachen GPT-generated jets from Ref. [32] vs. bin-smeared JETCLASS jets yields AUC values of 0.63 for QCD jets and 0.75 for top jets, considering the 40 hardest constituents of each jet.

tokens $\log p_\theta(\mathbf{t}_x)$ or the bin-smeared likelihood in the continuous physical space $\log p_\theta^{cont}(\mathbf{x})$ – the extra correction factors in (2) cancel out.

2. “JetClass OmniLearn”: An OmniLearn classifier trained on JETCLASS jets.
3. “GPT OmniLearn”: An OmniLearn classifier trained on samples drawn from the GPT model.

We see that while the OmniLearn classifiers trained on original JETCLASS and on GPT outputs agree, the NP-optimal classifier derived from the GPT model gives a much higher ROC curve. Evidently, the GPT tops vs. QCD jets are highly separable, and in a way that the OmniLearn classifier cannot detect. This reproduces the observation of Ref. [13], and in our setup the separability is even more pronounced.

However, unlike that work, we point out that there are two possible interpretations of this result: either

- (a) The NP-optimal classifier between GPT tops and QCD is representative of the NP-optimal classifier between JETCLASS tops and QCD, in which case there is a major gap between state-of-the-art jet taggers and the true “fundamental limit” of jet tagging (the conclusion reached in [13]),

or

- (b) the GPT model introduces unphysical artifacts into JETCLASS that make the GPT tops much more separable from GPT QCD than their JETCLASS counterparts, and in a way that the classifier is somehow unable to learn. Then the fact that the state-of-the-art classifiers cannot reproduce the NP-optimal result on GPT-generated jets is unfortunate, but possibly besides the point if this result is due to unphysical artifacts.

We next turn to another popular generative modeling framework for jets – permutation-equivariant conditional flow matching – and explore what it has to say about the separation of top vs. QCD jets.

III. EPIC-FM: JETS AS PARTICLE CLOUDS

A. Description of the model

We represent jets as *unordered* particle clouds with kinematic features $\mathbf{x} = \{\mathbf{x}_i\}_{i=1}^N$, where each constituent carries $\mathbf{x}_i = (p_T, \Delta\eta, \Delta\phi)_i$. We restrict each jet to at most the 40 highest- p_T constituents (i.e. $N_{\max} = 40$), discarding softer particles beyond this cutoff. We use the conditional flow-matching (CFM) paradigm to train a continuous normalizing flow equipped with an time-dependent EPiC architecture [19] that approximates the velocity field \mathbf{u}_t of the ODE trajectories connecting

source to target samples. The setup for our architecture closely follows the one in [20], with the main difference being that the only conditional model inputs that we consider are the time variable $t \in [0, 1]$, the jet-type (QCD and top), and the number of constituents in each jet. For further details on the model and training see appendix A 1.

Once trained, we generate jets of a specified class by integrating the ODE $\dot{\mathbf{x}}_t = \mathbf{u}_t^\theta(\mathbf{x}_t)$ from $t = 0$ to $t = 1$ with a forward Euler integrator with step size $\Delta t = 3.33 \times 10^{-3}$. We estimate the log-density of the generated jets by integrating the dynamics with backward Euler steps in the reverse time direction. Writing the total derivative as $d/dt = \partial_t + \mathbf{u}_t^\theta \cdot \nabla_{\mathbf{x}}$, the log-density evolves according to

$$\frac{d}{dt} \log p_t(\mathbf{x}) = -\text{Tr} \left[\frac{\partial \mathbf{u}_t^\theta}{\partial \mathbf{x}} \right], \quad (3)$$

which integrates to

$$\log p_1(\mathbf{x}_1) = \log p_0(\mathbf{x}_0) - \int_0^1 dt \text{Tr} \left[\frac{\partial \mathbf{u}_t^\theta}{\partial \mathbf{x}} \right]. \quad (4)$$

If the model is well trained, then this solution is expected to approximate the target log-density. In practice, we accumulate the Jacobian-trace term via automatic differentiation; numerical details and validation (fixed-point refinement, step-size scans) are provided in appendix A 2.

B. Baseline evaluation of the model

We again perform a 2-sample classifier test between EPiC-FM generated jets and JETCLASS jets. We find AUCs of 0.6729 for QCD and 0.7703 for tops. These results indicate that EPiC-FM captures the main features of the simulation with moderate levels of mismodeling likely due to finite training data and finite model capacity.

C. Top vs. QCD separation with EPiC-FM

Next we investigate whether EPiC-FM accurately models the separation of top vs. QCD jets. In Figure 2 we compare the ROC curves of three top vs. QCD classifiers (each classifier is evaluated on the samples from the generative process it’s defined/trained with):

1. “EPiC-FM Optimal”: The NP-optimal EPiC-FM likelihood ratio discriminator
2. “JetClass OmniLearn”: An OmniLearn classifier. trained on JETCLASS jets.
3. “EPiC-FM OmniLearn”: An OmniLearn classifier trained on samples from the EPiC-FM model.

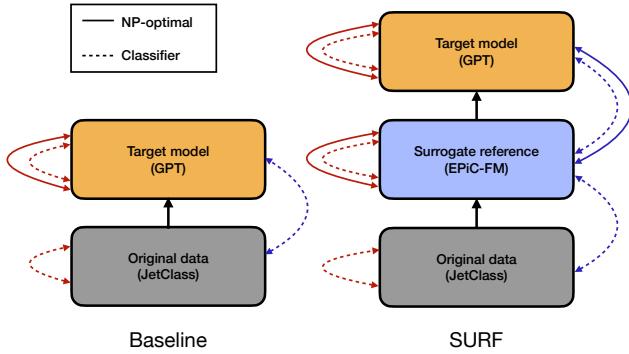


FIG. 1. Illustration of the SURF method. The conventional baseline approach is shown on the left and the SURF approach is shown on the right. The red lines denote interclass (e.g. top vs QCD) tests and the blue lines denote 2-sample tests. The solid lines correspond to NP-optimal tests, while dashed lines correspond to trained classifier tests (e.g. OmniLearn). Here we see the NP-optimal tests enabled by the SURF method.

We see from Figure 2 that all three ROC curves are fairly consistent with one another. While we cannot totally rule out the possibility that JETCLASS has subtle features that neither EPiC-FM nor OmniLearn can detect, the concurrence of JETCLASS OmniLearn with both EPiC-FM ROC curves is highly suggestive that the NP-optimal JETCLASS ROC curve is also close by, and that the SOTA classifiers are close to the fundamental limit.

In any case, the OmniLearn classifier is evidently powerful enough to learn the true separation of the EPiC-FM top and QCD samples, unlike what we saw for the GPT jets.⁴ Together with the 2-sample classifier tests described in the previous subsection, this gives us confidence that EPiC-FM is not introducing unphysical artifacts into JETCLASS and thus constitutes a valid surrogate reference model.

IV. SURF METHOD

We have two generative models that give different possibilities for the “fundamental limit” of jet tagging – GPT suggests a huge gap between the true NP-optimal ROC curve of tops vs. QCD and the performance of state-of-the-art classifiers, while EPiC-FM suggests almost no gap. Which is correct?

Here we propose a new method that provides a plausible answer to this puzzle: the *SUrrogate ReFerence (SURF) method*. As illustrated in Figure 1, the core idea is to use a generative model with tractable likelihoods –

Comparison	QCD	Top
EPiC-FM vs GPT	0.9955 1.000	0.9746 1.000
EPiC-FM-bs vs GPT	0.5505 0.5640	0.5579 0.6406
EPiC-FM vs EPiC-FM-bs	0.9950 0.9980	0.9733 0.9967

TABLE II. Same as Table I but with JETCLASS replaced with the EPiC-FM surrogate reference. In each cell, the upper-left value corresponds to the OmniLearn classifier AUC, and the lower-right value corresponds to the NP-optimal classifier AUC, which is now available thanks to the tractable likelihoods of the surrogate reference.

the *surrogate reference model* – to evaluate another target generative model whose likelihood is likewise tractable. We achieve this by first training the surrogate reference model on the original data, and then training the target model on samples generated by the surrogate reference model.

We have seen that EPiC-FM seems to be a valid surrogate for JETCLASS, in that it does not seem to introduce any subtle artifacts into JETCLASS that OmniLearn cannot detect. So the SURF method in this context is to use samples from EPiC-FM as a new reference dataset to train GPT on. Then we will have access to all the classifiers as before, plus a fourth one: the true NP-optimal ROC curve for the reference dataset. This is what’s lacking for JETCLASS, and is the value added by the SURF method.

As a first application of the SURF method, we perform the 2-sample classifier tests for top and QCD jets. These are summarized in the first and second rows of Table II for continuous and bin-smeared jets respectively. We confirm using the EPiC-FM surrogate reference that the GPT-generated jets are well-separated from continuous jets, and that the OmniLearn classifier is able to pick this up. We also confirm that the GPT-generated jets are quite close to bin-smeared jets, so the OmniLearn classifier was not misleading on this front. As shown in the third row of Table II, we observe the same behavior previously identified for JETCLASS: bin-smeared EPiC-FM jets are nearly perfectly separable from their continuous counterparts.

Next we turn to the main question of interest: whether these GPT models accurately represent the true degree of top vs. QCD separation or not. The relevant ROC curves are shown in Figure 3. Here in addition to the analogues of the three classifiers considered previously – GPT Optimal, Surrogate Reference OmniLearn and GPT Om-

⁴ The difference in R50 values between EPiC-FM Optimal and EPiC-FM OmniLearn is within a factor of 2. This is consistent with recent gains in top tagging from improved architectures and more data [7, 11].

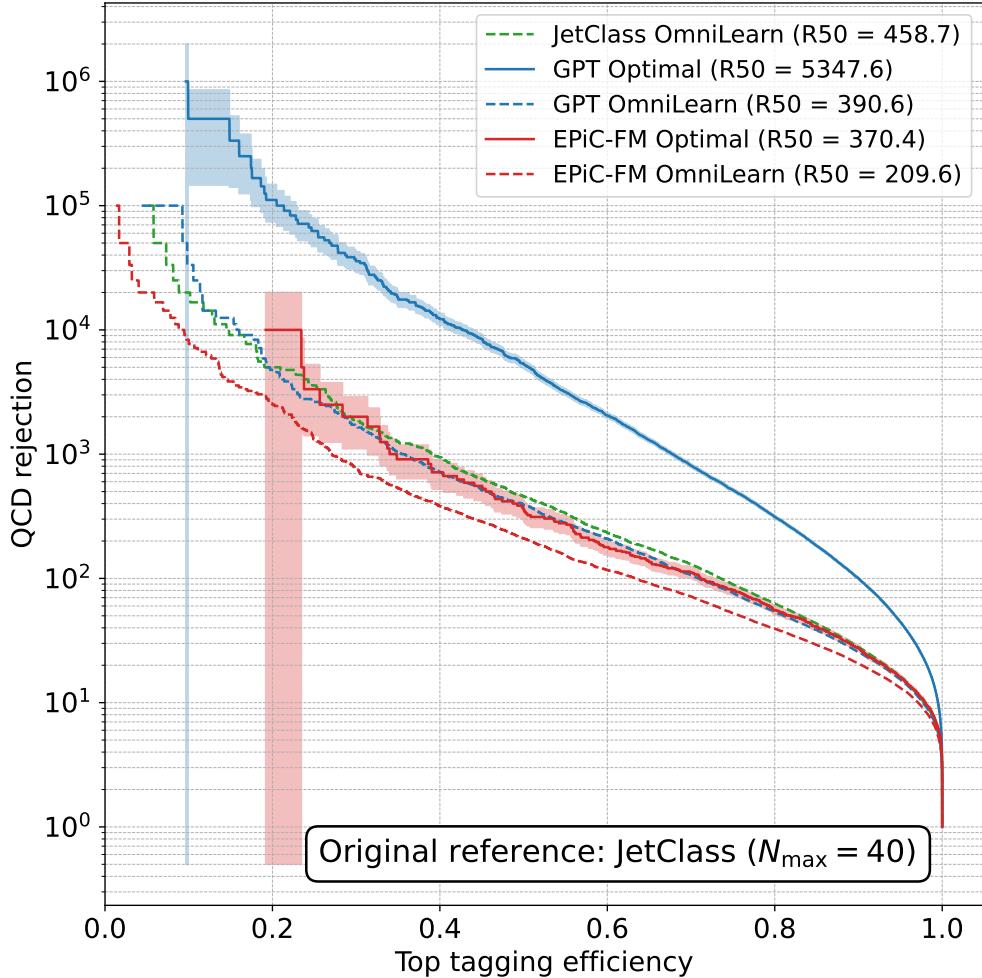


FIG. 2. ROC curves for top vs. QCD jet classification derived from the **JetClass** dataset and the generative models trained on it. Note that, although physical jets can contain more than 40 constituents, in this study each jet is represented only by its 40 hardest constituents (i.e. $N_{\max} = 40$). Solid lines indicate the performance of the NP-optimal classifier obtained directly from the true log-likelihood ratio. Dashed lines correspond to classifiers trained with OmniLearn. The QCD rejection at a top tagging efficiency of 50% (R50) is shown in parentheses for each curve. Shaded bands indicate the statistical uncertainty on the optimal ROC curves, estimated from binomial counting errors on the background sample and propagated to the QCD rejection axis.

niLearn, we are able to add a fourth ROC curve, Surrogate Reference Optimal, thanks to the tractable likelihoods of the EPiC-FM surrogate reference model.

We see that the GPT model vastly inflates the ground-truth ROC curve of the surrogate reference model! By contrast, the OmniLearn classifier curves follow the Surrogate Reference Optimal curve closely, showing that classifier-based discrimination remains consistent with the physically meaningful separation encoded in the EPiC-FM surrogate. This confirms that the inflated NP-optimal behavior originates from artifacts in the GPT likelihood landscape, not from deficiencies in classifier performance as was claimed in Ref. [13].

As expected, we observe the same behavior – both qualitatively and quantitatively – as in the baseline GPT evaluation (Figure 2), where we found the NP-optimal ROC curve for GPT jets to be substantially larger than

the OmniLearn classifiers trained on the reference data and on GPT samples. This cross-validation builds further confidence in our interpretation of the hugely inflated NP-optimal ROC curve from GPT tops and QCD.

V. INVESTIGATING ROC INFLATION IN THE GPT MODEL

Since top and QCD jets are already extremely well separated ($\text{AUC} \geq 0.98$ with state-of-the-art taggers), the ROC curve of any approximate model becomes highly sensitive to even small mismodelings of the jet distributions.

We illustrate this sensitivity using a 10D Gaussian toy model with background given by $\mathcal{N}(\mathbf{0}, I_{10})$ and signal given by $\mathcal{N}(\mathbf{1}, I_{10})$. To emulate mismodeling, we in-

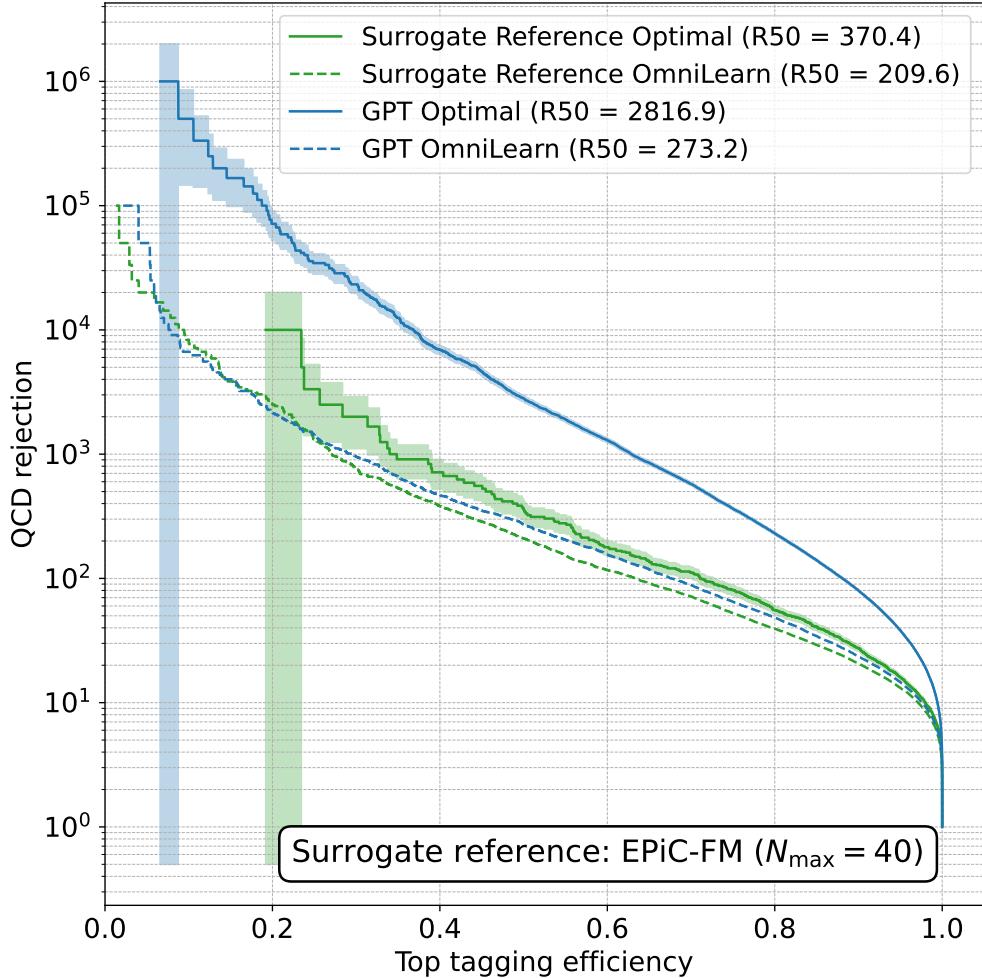


FIG. 3. ROC curves for top vs. QCD jet classification derived from the **EPiC-FM surrogate reference** samples and the GPT models trained on them. Note that, although physical jets can contain more than 40 constituents, in this study each jet is represented only by its 40 hardest constituents (i.e. $N_{\max} = 40$). Solid lines indicate the performance of the NP-optimal classifier obtained directly from the true log-likelihood ratio. Dashed lines correspond to classifiers trained with OmniLearn. The QCD rejection at a top tagging efficiency of 50% (R50) is shown in parentheses for each curve. Shaded bands indicate the statistical uncertainty on the optimal ROC curves, estimated from binomial counting errors on the background sample and propagated to the QCD rejection axis.

introduce shifted versions of the background and signal, described by Gaussians with slightly displaced means: $\mathcal{N}(-0.1 \mathbf{1}_{10}, I_{10})$ and $\mathcal{N}(1.1 \mathbf{1}_{10}, I_{10})$.

While the AUCs between the original and shifted distributions are around 0.6, the ROC curve between the shifted signal and background appears highly inflated, as shown in Figure 4. Because the original signal and background already exhibit minimal overlap, even a small shift in the likelihood ratios can have an outsized impact on the ROC curve. This example highlights how minor deviations in otherwise well-separated classes can artificially enhance apparent separation power.

Next we turn to the question of what kind of mismodeling could be causing the inflated ROC curve in the GPT-generated jets.

A. Bin-smearing?

The fact that the GPT-generated jets are attempting to match bin-smeared jets, and the bin-smeared jets are off-manifold compared to continuous jets, could be a significant source of mismodeling.

However, Figure 5 shows that this is not the issue. Using the EPiC-FM model as the reference, we can compute the true likelihood of both continuous and bin-smeared jets, so we can obtain the true ROC curve between continuous tops vs. QCD as well as bin-smeared tops vs. QCD. We see that if anything, the bin-smeared jets are *less* separable than their original continuous counterparts. For JETCLASS, where the true likelihood is not available and we instead rely on an OmniLearn classifier, we observe the same qualitative behavior that bin-

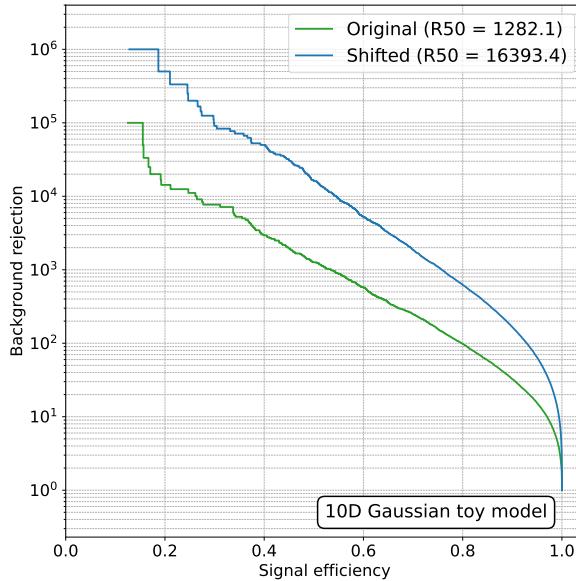


FIG. 4. ROC curves from the 10-dimensional Gaussian toy model. Although the shifted signal and background distributions differ only slightly from their original counterparts, their mutual ROC curve appears highly inflated. Here, “Original” denote the original signal and background, while “Shifted” indicate the slightly displaced versions. The background rejection at a signal efficiency of 50% (R50) is shown in parentheses for each curve. This demonstrates how small mismodelings of well-separated classes can artificially exaggerate separation power.

smearing reduces separability. So the binning of the jets, despite being a major deviation from the continuous jets, cannot be the cause of the inflated top vs. QCD ROC curve.

B. Overfitting

Finally, we consider another potential source of the mismodeling that affects GPT-generated jets but not EPiC-FM-generated jets: overfitting. Shown in Figure 6 are comparisons of the train vs. validation loss curves for GPT and EPiC-FM models. We see that the GPT jets overfit almost immediately (increasing gap between train vs. val losses), despite their val loss continuing to decrease slowly. Overfitting in GPT-style models is a known phenomenon, discussed already in the original GPT-2 paper [33]. Meanwhile the EPiC-FM loss curves are in perfect agreement between train vs. validation. Hence, we do not find any evidence for the EPiC-FM model overfitting.⁵ More broadly, flow-based models such as diffusion

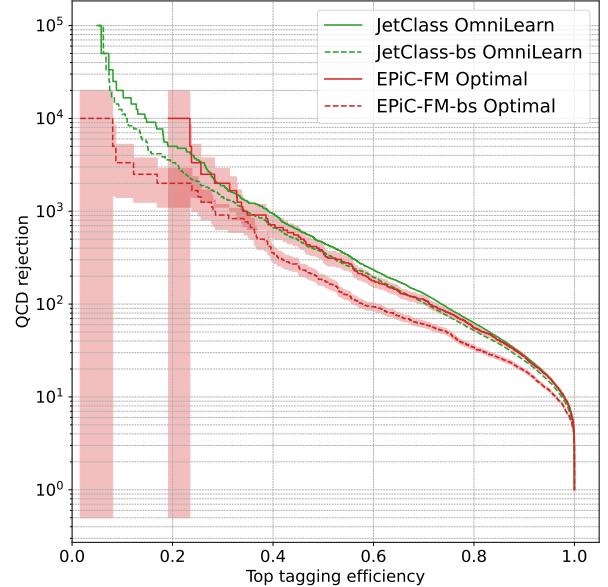


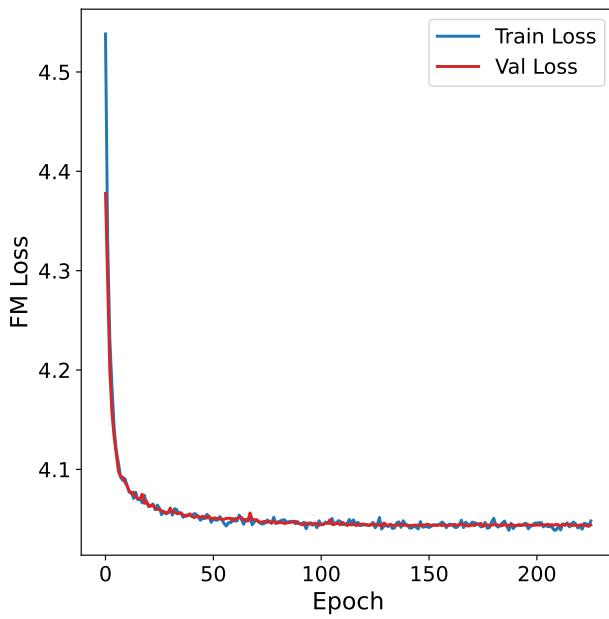
FIG. 5. Comparison of top vs. QCD discrimination for both EPiC-FM surrogate reference jets and JETCLASS jets. The EPiC-FM curves, shown in red (continuous) and red dashed (bin-smeared, labeled “bs”), represent the true NP optimal. Shaded bands indicate the statistical uncertainty on the optimal ROC curves, estimated from binomial counting errors on the background sample and propagated to the QCD rejection axis. For JETCLASS, where the true likelihood is not available, the solid and dashed blue curves show the corresponding OmniLearn classifiers trained on continuous and bin-smeared (“bs”) JETCLASS samples, respectively. Across both EPiC-FM and JETCLASS jets, bin-smearing leads to slightly reduced separability between top and QCD jets.

and flow-matching have often been observed to generalize well in machine learning applications, though the underlying reasons remain an active area of investigation (see e.g., Refs. [34–36]).

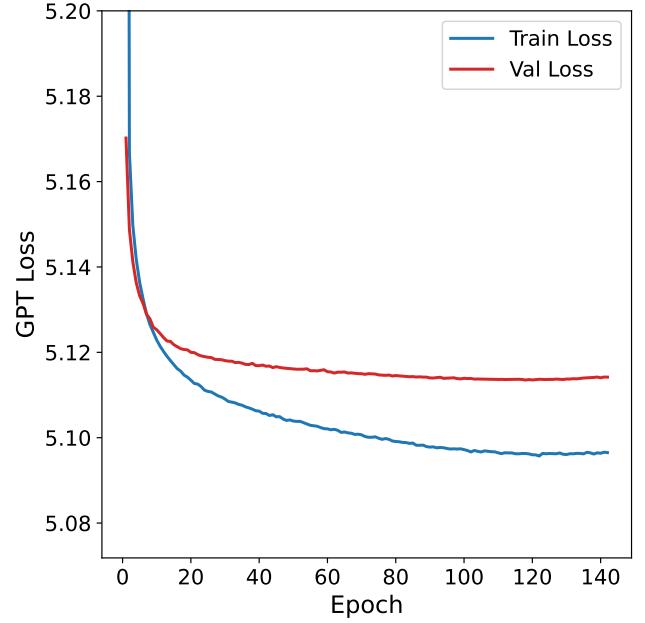
Finally, we directly examine the effects of model overfitting on the optimal likelihood ratio performance. We train separate GPT models at fixed capacity on subsets of JETCLASS ranging from 10⁵ to 10⁷ jets for tops and QCD. As the training set size decreases – particularly when the number of examples becomes small relative to the number of model parameters and in the absence of strong regularization – overfitting becomes more severe. Indeed, as shown in Figure 7, the models trained on smaller datasets exhibit a more inflated ROC curve, indicating that overfitting amplifies the apparent separation.

The overfitting hypothesis could also potentially explain why state-of-the-art classifiers cannot detect the inflated ROC curve for GPT jets. Overfitting in generative models means that the model is too close to the training data, starting to memorize the training data and become a lookup table, i.e. the model is reducing to a sum of (smeared out) delta functions centered around the training data. Compared to the smooth distribution describing the reference data, a sum of smeared out

⁵ Although the EPiC-FM loss is not equivalent to the log-likelihood, we checked that the EPiC-FM model with the best validation loss gives consistent log-likelihoods on the train and validation set.



(a) EPiC-FM train and validation loss curves.



(b) GPT (tops) train and validation loss curves.

FIG. 6. Training vs. validation loss curves for EPiC-FM and GPT. The GPT plot is shown for tops, with qualitatively similar behavior observed for QCD. GPT overfits almost immediately (diverging train/val losses), while EPiC-FM maintains closer alignment.

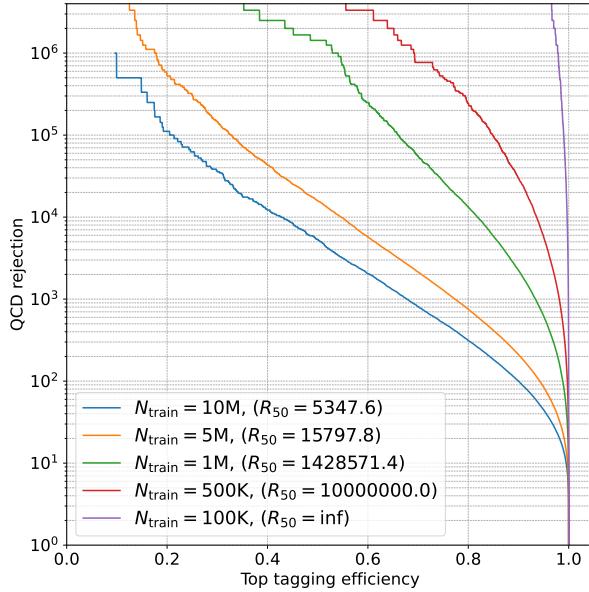


FIG. 7. Top vs. QCD ROC curves for GPT models trained on increasing fractions of the JETCLASS dataset, evaluated with the NP-optimal classifier.

delta functions is a high-frequency deformation. In the GPT case, this shows up as memorization of specific token co-occurrence patterns within the training jets, in the extreme, near-copying of jet (sub)sequences. Such behavior can distort multi-particle (higher-order) corre-

lations, manifesting as ‘high-frequency’ artifacts in the data manifold. Neural networks are well-known to have difficulty learning high frequency modes; see Refs. [37, 38] for discussions and studies of this. This could be the ultimate reason why the state-of-the-art classifiers cannot detect the inflated GPT ROC curves.

VI. CONCLUSION

We revisited recent claims that autoregressive GPT-style jet generators expose a large performance gap between ML-based jet taggers and the NP optimal “fundamental limit” of top vs. QCD discrimination [13]. Using our own independently trained GPT-based jet generator – with an improved architecture from that used in the original study – we reproduced their reported behavior. However, when comparing with the tractable EPiC-FM flow-matching model, which provides explicit likelihoods for unordered particle clouds, we find no evidence for such an inflated limit.

To resolve the apparent tension between the EPiC-FM and GPT results, we introduced the SURF method, which enables exact NP tests even when the real data likelihood is intractable. By training the target GPT model on samples from a tractable EPiC-FM surrogate, SURF allows direct NP comparisons between the surrogate and the target. We find that the GPT likelihood ratios dramatically overstate the true top vs. QCD sepa-

ration encoded in the EPiC-FM surrogate. The “fundamental limit” ROCs obtained from GPT models therefore reflect artifacts of the generative process rather than performance limitations of the jet taggers.

We speculate that the inflated GPT ROC curve could be due to overfitting, and more generally high frequency artifacts introduced by the GPT model. Such high-frequency modes can be exploited by exact NP tests but are difficult for neural-network classifiers to learn, given their well-known low-frequency spectral bias [37, 38]. In future work it would be interesting to explore these high frequency artifacts in more detail. Existing work in the ML literature (e.g., [39]) suggests potential methods for making such artifacts more explicit. It would also be interesting to develop new classifier architectures that can detect these high frequency artifacts. This would lead to more sensitive classifier metric tests.

Furthermore, our study shows that GPT models trained on tokenized jets are faithful only to their discretized references, while being nearly perfectly separable from the original continuous jets. These results indicate that discretization can introduce significant off-manifold effects, which may limit the fidelity of generative models based on binning and tokenization if not carefully mitigated. Alternative discretization schemes could potentially reduce the degree of separability from the continuous reference. For example, Refs. [21, 22] use VQ-VAEs to tokenize jet constituents. Exploring such strategies, and their impact on preserving on-manifold fidelity, would be an interesting direction for future study.

Overall, the consistent picture between JETCLASS and the EPiC-FM surrogate reference, particularly the concurrence of three ROC curves – classifiers trained on samples from JETCLASS and EPiC-FM samples, and the NP-optimal likelihood ratio from EPiC-FM – point to the fundamental limit of top tagging being not far from the performance of state-of-the-art classifiers. However, we cannot rule out the possibility that there remains a gap between the fundamental limit of top tagging and SOTA classifiers – the EPiC-FM surrogate and OmniLearn may both miss subtle features in JETCLASS that make tops and QCD more separable than they seem.

We can think of several avenues for further probing the fundamental limit of jet tagging. For example, improving the surrogate reference should be a high priority. Our EPiC-FM surrogate was far from perfect – our classifier tests showed that they are partially separable from the original continuous jets ($AUC \approx 0.7$). It would be worthwhile to repeat this study with a more expressive surrogate reference (for example replacing EPiC with a transformer) that is closer to the JETCLASS data to confirm the conclusions reached here. It would also be good to repeat this study with a more faithful target generative model that does not inflate the true ROC curve of the surrogate reference. Taken together (a more faithful surrogate and target), this could start to reveal a gap between SOTA classifiers and the fundamental limit.

Recently some studies were initiated into the true the-

oretical limits on top tagging [40]. It would be interesting to bring the empirical, data-driven approach explored here into contact with the theory-driven approach and elucidate further the true fundamental limit of jet tagging.

We focused on just the momentum four-vectors of jet constituents in this work. For many jet tagging tasks, including top tagging, additional features are highly informative, such as track displacements and particle ID. These features are also included in JETCLASS, and it would be interesting to repeat this study for the full set of JETCLASS features to obtain an even fuller picture of the fundamental limit of jet tagging.

Finally, the SURF method provides a general route to perform exact NP tests using a tractable surrogate reference. Extending this approach to other classes of generative models, datasets, and applications – both within and beyond high-energy physics – would enable more systematic analysis of generative models.

CODE AVAILABILITY

The code for the GPT model and the EPiC-FM model can be found at github.com/dfaroughy/JetSequences and github.com/Ian-Pang/jet_cfm, respectively.

ACKNOWLEDGEMENTS

IP, DAF, DS and RD are supported by DOE grant DOE-SC0010008. RD is also supported by Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under grant 396021762 – TRR 257 Particle Physics Phenomenology after the Higgs Discovery. GK is supported by the DFG under the German Excellence Initiative – EXC 2121 Quantum Universe – 390833306. This research used resources of the National Energy Research Scientific Computing Center, a DOE Office of Science User Facility supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231 using NERSC award HEP-ERCAP0027491.

Appendix A: Details of the EPiC-FM model

1. Architecture and training

The EPiC network is composed of EPiC layers [41] – a permutation-equivariant variant of Deep Sets [42], also known in high-energy physics as particle flow networks (PFNs) [43]. These layers couple per-particle embeddings to a learnable global context, enabling the model to capture particle correlations while exactly preserving the permutation symmetry of the constituents [19]. The setup for our architecture closely follows the one

in [20], with the main difference being that the only conditional model inputs that we consider are the time variable $t \in [0, 1]$, the jet type (QCD and tops), and the number of constituents in each jet. Time is embedded into a vector space \mathbb{R}^{16} using a sinusoidal embedding network and the jet-type labels are embedded into \mathbb{R}^8 using a linear lookup table. We stack 12 EPiC-Layers with a local hidden dimension of $h_{\text{loc}} = 300$ and a global hidden dimension $h_{\text{glob}} = 32$. The variable number of particles per jet is handled via binary masking. During training the multiplicity is given by the true number of constituents in each jet, while during generation it is sampled from an empirical distribution computed from the JETCLASS dataset.

During training, we construct source–target pairs by sampling the source from a three-dimensional Gaussian and randomly pairing it with a target jet (QCD or top) from the dataset. We use 10 million top jets and 10 million QCD jets, split 80/20 into train/validation.

Optimization uses Adam [44] with an initial learning rate of 10^{-4} and a reduce-on-plateau schedule triggered by the validation loss with a patience of 10 epochs. We train for a total of 1000 epochs with a patience window of 50 epochs. We select the best checkpoint with the lowest validation loss.

2. Validation of log-likelihood computation

With a trained FM model \mathbf{u}^θ , we can calculate the log-likelihood iteratively by taking the following inference step

$$\log p_{t-\Delta t}(\mathbf{x}_{t-\Delta t}) = \log p_t(\mathbf{x}_t) - \Delta t \operatorname{Tr}(\operatorname{Jac}[\mathbf{u}_{t-\Delta t}^\theta(\mathbf{x}_t)]),$$

at every time step Δt . We compute the Jacobian matrix `Jac[]` of the neural network in pytorch via `torch.func.jacrev`. We verified that using `torch.autograd.functional.jacobian` yields identical results.

A natural concern is whether flow matching models yield *reliable* likelihood *ratios*, especially in light of claims about the precision of likelihood computation in related models [13]. Constant offsets in log-likelihoods cancel in this ratio and therefore do not affect ROC curves. Operationally, we require numerical stability of *differences* in log-likelihoods accumulated along the FM trajectory. We assessed this with the following two tests.

a. Consistency of Forward and Backward Euler Integration

As mentioned in Sec. III A, the forward Euler step is used for generation and the backward step is used for log-likelihood computation. Therefore, it is vital that the backward step is as close as possible to an exact inverse of the forward step. This guarantees consistent trajectories and preserves the theoretical assumption that the underlying flow is perfectly reversible. Any mismatch between

the two steps violates this assumption and might lead to incorrect log-likelihood computation.

Unlike for regular normalizing flows, consistency of forward and backward steps is not guaranteed for continuous normalizing flows. The reason can be seen by the following:

$$\text{Forward step: } \mathbf{x}_{t+\Delta t} = \mathbf{x}_t + \Delta t \mathbf{u}_t^\theta(\mathbf{x}_t), \quad (\text{A1})$$

$$\text{Backward step: } \mathbf{x}'_t = \mathbf{x}_{t+\Delta t} - \Delta t \mathbf{u}_t^\theta(\mathbf{x}_{t+\Delta t}). \quad (\text{A2})$$

A backward update (Eq. A2) generally fails to recover the original point, since

$$\mathbf{x}'_t - \mathbf{x}_t = \Delta t(\mathbf{u}_t^\theta(\mathbf{x}_{t+\Delta t}) - \mathbf{u}_t^\theta(\mathbf{x}_t)) \neq 0$$

This means that the forward and backward updates are not exact inverses of each other; they only approach invertibility in the limit of infinitesimally small step sizes, $\Delta t \rightarrow 0$.

We observe that the backward step would be an exact inverse of the forward step if, in Eq. A2, the vector field \mathbf{u}_t is evaluated at \mathbf{x}_t rather than at $\mathbf{x}_{t+\Delta t}$. However, this is precisely what is unavailable during the computation of the backward step, indeed this is meant to be the *output* of the backward step.

The problem is solved if we view the backward step as an *implicit equation* for \mathbf{x}_t :

$$\mathbf{x}_t = \mathbf{x}_{t+\Delta t} - \Delta t \mathbf{u}_t^\theta(\mathbf{x}_t), \quad (\text{A3})$$

This can be solved for \mathbf{x}_t via fixed-point iteration.

Figure 8a displays the distribution of the component-wise relative error

$$\frac{x'_{t,i} - x_{t,i}}{x_{t,i}},$$

where i indexes all components of the flattened jet (i.e., both constituent and feature). Results are shown for backward updates computed with and without fixed-point refinement over 300 integration time steps. When employing the iterative backward update, we performed five fixed-point refinement iterations. The iterative method yields a sharply peaked distribution around zero and strongly suppresses the heavy tails observed in the naive backward update.

Nevertheless, we show in Figure 8b that the error has minimal impact on approximate NP classifier performance computed based on log-likelihood values obtained with and without the iterative procedure.⁶ This study demonstrates the robustness of the FM-based log-likelihood ratio computation to errors accumulated in the backward update. For the main results presented in the paper, we employed the naive backward update to efficiently generate a large training dataset for the jet taggers at reasonable computational costs.

⁶ In the iterative case, error bands are wider, as computational constraints limited the evaluation sample size used to construct the ROC curve.

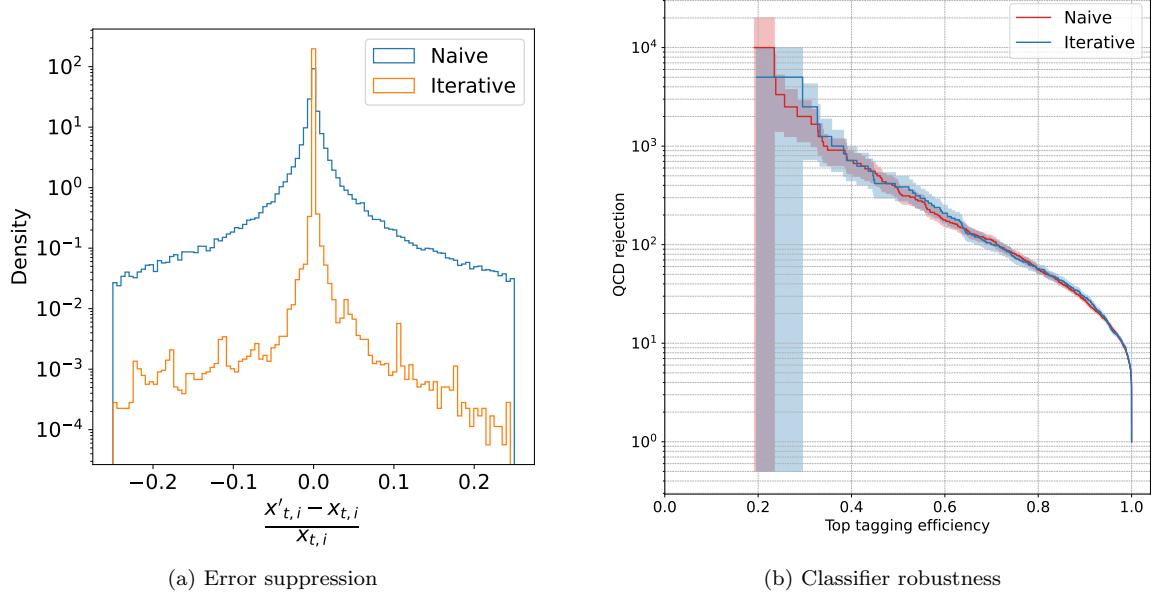


FIG. 8. Effect of fixed-point refinement in the backward update. “Naive” denotes the backward update computed without any fixed-point iterations; “Iterative” denotes the backward update performed with five fixed-point refinement iterations. (a) Component-wise relative error $\frac{x'_{t,i} - x_{t,i}}{x_{t,i}}$ is sharply peaked around zero when using fixed-point iterations. Here we show results for a few top jets; similar behavior is observed across larger jet samples, including QCD jets. (b) The ROC curves built from FM-based log-likelihood ratios with and without refinement nearly coincide, showing negligible impact on classifier performance. Shaded bands indicate the statistical uncertainty on the optimal ROC curves, estimated from binomial counting errors on the background sample and propagated to the QCD rejection axis.

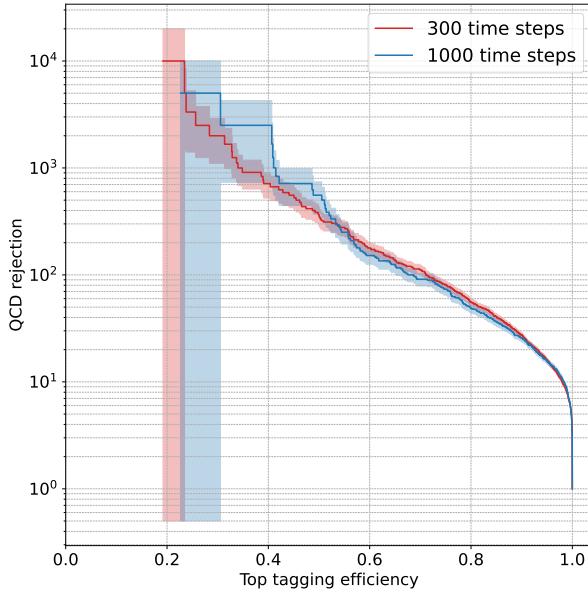


FIG. 9. Comparison of NP-optimal ROC curves with number of inference/sampling timesteps. The blue curve shows the performance when both inference and sample generation use 300 time steps, while the red curve corresponds to 1000 time steps. The shaded bands denote binomial counting errors on the background sample and propagated to the QCD rejection axis.

b. Number of time steps

To assess the dependence of the approximate Neyman-Pearson classifier performance on the number of time steps used for sample generation and inference the log-likelihood values, we compared two configurations – 300 versus 1000 time steps. For each configuration, the same number of time steps was used for sampling and inference. No fixed-point iteration refinement was used in either configuration. We found that using a 1000 time steps for generation results in slightly higher fidelity jets. Nevertheless, as shown in Figure 9, both configurations yield statistically consistent classification performance. The larger error bands for the 1000-step case are from the smaller evaluation sample size – imposed by computational constraints – used to construct the ROC curve. For our main results in the paper, we used the 300-step configuration in order to obtain a large training dataset for the jet taggers while keeping computational costs manageable.

Beyond these internal consistency checks, we also verify correctness in a controlled toy setting where the analytic NP ROC is known. This provides an end-to-end demonstration that FM likelihood ratios reproduce exact NP behavior in a case with closed-form ground truth.

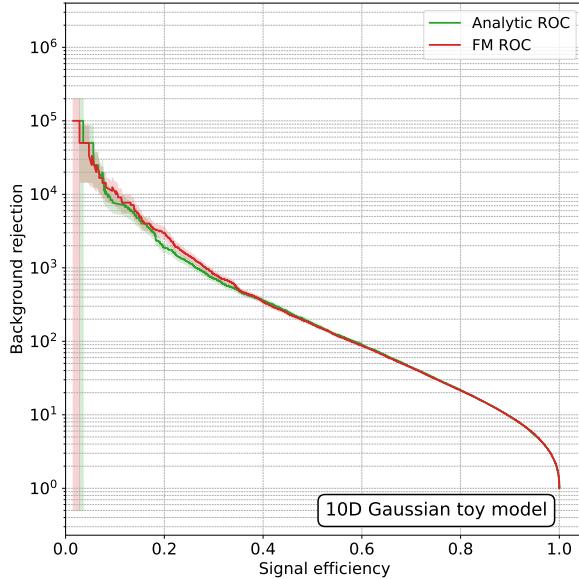


FIG. 10. ROC curves for 10D gaussian toy example used for flow matching log-likelihood validation. The shaded bands denote binomial counting errors on the background sample and propagated to the background rejection axis.

c. Validation on a toy 10D Gaussian model

We construct two ten-dimensional Gaussian distributions with shared unit covariance but shifted means:

$$p_{\text{sig}} = \mathcal{N}(+0.4 \mathbf{1}_{10}, I_{10}), \quad p_{\text{bg}} = \mathcal{N}(-0.4 \mathbf{1}_{10}, I_{10}).$$

In this setting, the NP test has a closed-form solution, and the corresponding ROC curve can be derived analytically. We train separate FM models, implemented as simple MLP-based vector-field networks, on samples from each Gaussian. After training, we generate samples from the FM models and evaluate the NP log-likelihood ratio using the FM log-density estimators on these gen-

erated samples. Figure 10 shows that the ROC curve obtained from FM likelihood ratios coincides quite well with the analytic NP ROC and is mostly within statistical fluctuations.

This toy test confirms that flow matching recovers the correct ROC behavior in a regime where the ground truth NP statistic is exactly known, supporting the validity of our approach before applying it to more complex jet data.

Taken together, these checks validate that FM-based log-likelihood ratios are numerically stable under integrator choices and, in the Gaussian toy model, reproduce the exact NP ROC. They do not claim absolute calibration of log-densities; rather, they establish that the relative log-densities that define our ROC curves are robust for the analyses presented here.

Appendix B: Details of the GPT model

For our model, we use the `GPT2LMHeadModel` implementation of the GPT-2 transformer architecture [45] provided by HuggingFace [23], configured with an initial token embedding to $\mathbb{R}^{d_{\text{embd}}}$, followed by $n_{\text{layer}} = 12$ transformer blocks, each with $n_{\text{head}} = 8$ attention heads. We set $d_{\text{embd}} = 256$ for the embedding/hidden size, $d_{\text{inner}} = 1024$ for the feed-forward hidden size, $p_{\text{drop}} = 0.1$ for all dropout layers and GELU for all non-linear activation functions. This resulting model contains approximately 16 million trainable parameters. In addition to the standard causal mask that enforces autoregressive ordering, we apply an attention mask to ignore [PAD] tokens. The final “language head” then maps the last hidden states to vocabulary logits through a linear layer.

The model is trained by maximizing the autoregressive likelihood (1) with $t_0 \equiv [\text{START}]$, excluding [PAD] tokens. Unlike EPiC-FM, which conditions on jet type, we train separate class-specific GPT models – one for top jets and one for QCD jets. Optimization uses Adam [44] with an effective batch size of 256 jets for up to 150 epochs: a cosine-annealing schedule decays the learning rate from 10^{-3} to 10^{-4} over the first 100 epochs, followed by 50 epochs at a fixed rate of 10^{-4} . The final model is selected by the lowest validation loss.

-
- [1] H. Qu and L. Gouskos, ParticleNet: Jet Tagging via Particle Clouds, *Phys. Rev. D* **101**, 056019 (2020), [arXiv:1902.08570 \[hep-ph\]](#).
 - [2] A. Butter *et al.*, The Machine Learning Landscape of Top Taggers, *SciPost Phys.* **7**, 014 (2019), [arXiv:1902.09914 \[hep-ph\]](#).
 - [3] A. M. Sirunyan *et al.* (CMS), Identification of heavy, energetic, hadronically decaying particles using machine-learning techniques, *JINST* **15** (06), P06005, [arXiv:2004.08262 \[hep-ex\]](#).
 - [4] S. Gong, Q. Meng, J. Zhang, H. Qu, C. Li, S. Qian, W. Du, Z.-M. Ma, and T.-Y. Liu, An Efficient Lorentz Equivariant Graph Neural Network for Jet Tagging, *JHEP* **07**, 030, [arXiv:2201.08187 \[hep-ph\]](#).
 - [5] A. Bogatskiy, T. Hoffman, D. W. Miller, J. T. Offermann, and X. Liu, Explainable Equivariant Neural Networks for Particle Physics: PELICAN (2023), [arXiv:2307.16506 \[hep-ph\]](#).
 - [6] H. Qu, C. Li, and S. Qian, Particle Transformer for Jet Tagging (2022), [arXiv:2202.03772 \[hep-ph\]](#).
 - [7] J. Brehmer, V. Bresó, P. de Haan, T. Plehn, H. Qu, J. Spinner, and J. Thaler, A Lorentz-equivariant transformer for all of the LHC, *SciPost Phys.* **19**, 108 (2025), [arXiv:2411.00446 \[hep-ph\]](#).
 - [8] V. Mikuni and B. Nachman, Solving key challenges in collider physics with foundation models, *Phys. Rev. D* **111**, L051504 (2025), [arXiv:2404.16091 \[hep-ph\]](#).
 - [9] V. Mikuni and B. Nachman, A Method to Simultaneously

- Facilitate All Jet Physics Tasks, *Phys.Rev.D* **111**, 054015 (2025), arXiv:2502.14652 [hep-ph].
- [10] ATLAS Collaboration, Accuracy versus precision in boosted top tagging with the ATLAS detector, *JINST* **19**, P08018, arXiv:2407.20127 [hep-ex].
- [11] W. Bhimji, C. Harris, V. Mikuni, and B. Nachman, OmniLearned: A Foundation Model Framework for All Tasks Involving Jet Physics (2025), arXiv:2510.24066 [hep-ph].
- [12] J. Neyman and E. S. Pearson, On the Problem of the Most Efficient Tests of Statistical Hypotheses, *Phil. Trans. Roy. Soc. Lond. A* **231**, 289 (1933).
- [13] J. Geuskens, N. Gite, M. Krämer, V. Mikuni, A. Mück, B. Nachman, and H. Reyes-González, The Fundamental Limit of Jet Tagging (2024) arXiv:2411.02628 [hep-ph].
- [14] H. Qu, C. Li, and S. Qian, JetClass: A large-scale dataset for deep learning in jet physics, 10.5281/zenodo.6619768 (2022).
- [15] T. Finke, M. Krämer, A. Mück, and J. Tönshoff, Learning the language of QCD jets with transformers, *JHEP* **06**, 184, arXiv:2303.07364 [hep-ph].
- [16] M. S. Albergo and E. Vanden-Eijnden, Building normalizing flows with stochastic interpolants, arXiv preprint arXiv:2209.15571 (2022).
- [17] Y. Lipman, R. T. Chen, H. Ben-Hamu, M. Nickel, and M. Le, Flow matching for generative modeling, arXiv preprint arXiv:2210.02747 (2022).
- [18] A. Tong, K. Fatras, N. Malkin, G. Huguet, Y. Zhang, J. Rector-Brooks, G. Wolf, and Y. Bengio, Improving and generalizing flow-based generative models with minibatch optimal transport, arXiv preprint arXiv:2302.00482 (2023).
- [19] E. Buhmann, C. Ewen, D. A. Faroughy, T. Golling, G. Kasieczka, M. Leigh, G. Quétant, J. A. Raine, D. Sengupta, and D. Shih, EPiC-ly Fast Particle Cloud Generation with Flow-Matching and Diffusion (2023), arXiv:2310.00049 [hep-ph].
- [20] J. Birk, E. Buhmann, C. Ewen, G. Kasieczka, and D. Shih, Flow matching beyond kinematics: Generating jets with particle identification and trajectory displacement information, *Phys. Rev. D* **111**, 052008 (2025), arXiv:2312.00123 [hep-ph].
- [21] T. Golling, L. Heinrich, M. Kagan, S. Klein, M. Leigh, M. Osadchy, and J. A. Raine, Masked particle modeling on sets: towards self-supervised high energy physics foundation models, *Mach. Learn. Sci. Tech.* **5**, 035074 (2024), arXiv:2401.13537 [hep-ph].
- [22] J. Birk, A. Hallin, and G. Kasieczka, OmniJet- α : The first cross-task foundation model for particle physics, *Mach. Learn. Sci. Tech.* **5**, 035031 (2024), arXiv:2403.05618 [hep-ph].
- [23] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. L. Scao, S. Gugger, M. Drame, Q. Lhoest, and A. M. Rush, Transformers: State-of-the-art natural language processing, in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations* (Association for Computational Linguistics, 2020) pp. 38–45.
- [24] O. Press and L. Wolf, Using the output embedding to improve language models, Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics (2017).
- [25] H. Inan, K. Khosravi, and R. Socher, Tying word vectors and word classifiers: A loss framework for language modeling, in *International Conference on Learning Representations (ICLR)* (2017).
- [26] C. Krause and D. Shih, CaloFlow: Fast and Accurate Generation of Calorimeter Showers with Normalizing Flows, *Phys. Rev. D* **107**, 113003 (2021), arXiv:2106.05285 [physics.ins-det].
- [27] O. Amram *et al.*, CaloChallenge 2022: A Community Challenge for Fast Calorimeter Simulation, *Reports on Progress in Physics* 10.1088/1361-6633/ae1304 (2024), arXiv:2410.21611 [physics.ins-det].
- [28] S. H. Lim, K. A. Raman, M. R. Buckley, and D. Shih, GalaxyFlow: upsampling hydrodynamical simulations for realistic mock stellar catalogues, *Mon. Not. Roy. Astron. Soc.* **533**, 143 (2024), arXiv:2211.11765 [astro-ph.GA].
- [29] R. Kansal, A. Li, J. Duarte, N. Chernyavskaya, M. Pierini, B. Orzari, and T. Tomei, Evaluating generative models in high energy physics, *Phys. Rev. D* **107**, 076017 (2023), arXiv:2211.10295 [hep-ex].
- [30] R. Das, L. Favaro, T. Heimel, C. Krause, T. Plehn, and D. Shih, How to Understand Limitations of Generative Networks, *SciPost Phys.* **16**, 031 (2023), arXiv:2305.16774 [hep-ph].
- [31] P. Cappelli, G. Grossi, M. Letizia, H. Reyes-González, and M. Zanetti, Learning to Validate Generative Models: a Goodness-of-Fit Approach (2025), arXiv:2511.09118 [stat.ML].
- [32] H. Reyes-Gonzalez, M. Krämer, A. Mück, B. Nachman, V. Mikuni, J. Geuskens, and N. Gite, Resources for the fundamental limit of jet tagging, 10.5281/zenodo.14023638 (2024).
- [33] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, Language models are unsupervised multi-task learners (2019).
- [34] Z. Kadkhodaie, F. Guth, E. P. Simoncelli, and S. Mallat, Generalization in diffusion models arises from geometry-adaptive harmonic representations, arXiv preprint arXiv:2310.02557 (2023).
- [35] J. J. Vastola, Generalization through variance: how noise shapes inductive biases in diffusion models, arXiv preprint arXiv:2504.12532 (2025).
- [36] Q. Bertrand, A. Gagneux, M. Massias, and R. Emonet, On the closed-form of flow matching: Generalization does not arise from target stochasticity, arXiv preprint arXiv:2506.03719 (2025).
- [37] N. Rahaman, A. Baratin, D. Arpit, F. Draxler, M. Lin, F. Hamprecht, Y. Bengio, and A. Courville, On the spectral bias of neural networks, in *International conference on machine learning* (PMLR, 2019) pp. 5301–5310.
- [38] B. Ronen, D. Jacobs, Y. Kasten, and S. Kritchman, The convergence rate of neural networks for learned functions of different frequencies, *Advances in Neural Information Processing Systems* **32** (2019).
- [39] Y. Zhang and X. Xu, Diffusion noise feature: Accurate and fast generated image detection, arXiv preprint arXiv:2312.02625 (2023).
- [40] A. J. Larkoski, Systematic Interpretability and the Likelihood for Boosted Top Quark Identification (2024), arXiv:2411.00104 [hep-ph].
- [41] E. Buhmann, G. Kasieczka, and J. Thaler, EPiC-GAN: Equivariant Point Cloud Generation for Particle Jets, *SciPost Phys.* **15**, 130 (2023), arXiv:2301.08128 [hep-ph].
- [42] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. R.

- Salakhutdinov, and A. J. Smola, Deep sets, *Advances in neural information processing systems* **30** (2017).
- [43] P. T. Komiske, E. M. Metodiev, and J. Thaler, Energy Flow Networks: Deep Sets for Particle Jets, *JHEP* **01**, 121, [arXiv:1810.05165 \[hep-ph\]](#).
- [44] D. P. Kingma and J. Ba, Adam: A method for stochastic optimization, arXiv preprint [arXiv:1412.6980](#) (2014).
- [45] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, Language models are unsupervised multi-task learners, *OpenAI Blog* **1** (2019).