

LEARNING AT THE EDGE: TAILED-UNIFORM SAMPLING FOR ROBUST SIMULATION-BASED INFERENCE

CHAIPAT TIRAPONGPRASERT^{1*} AND MATTHEW HO^{1*}

¹ Columbia Astrophysics Laboratory, Columbia University, 550 West 120th Street, New York, NY 10027, USA
Version January 15, 2026

ABSTRACT

We introduce the TAILED-UNIFORM proposal distribution for simulation-based inference. Instead of sampling parameters uniformly within bounded regions, we extend the distribution beyond prior boundaries with smooth Gaussian tails. This eliminates sharp discontinuities that cause neural posterior estimators to fail near parameter space boundaries. The method requires minimal hyperparameter tuning, with tail widths of 10–30% of the prior width proving robust across problems. We demonstrate these benefits on a synthetic Gaussian linear task and cosmological parameter inference from the matter power spectrum. We also demonstrate that boundary pathologies are systematic rather than data-deprived: adding more uniform samples provides negligible benefit, while TAILED-UNIFORM outperforms uniform sampling even with $8\times$ fewer simulations. This advantage grows in higher dimensions, where boundaries dominate parameter space volume. All code is publicly available at github.com/chaipattira/tailed-uniform-sbi.

Subject headings: astrophysics, machine learning, deep learning, simulation-based inference, neural posterior estimation, Bayesian methods

1. INTRODUCTION

Today, researchers—among them astronomers—face a growing challenge, as many systems of interest (e.g., exoplanets, black hole binaries, or galaxy clusters) are not amenable to direct probing (Feigelson and Babu 2012; Dodelson and Schmidt 2020). Thus emerges a new class of problems known as “inverse problems,” in which the key goal is to identify the model parameters that generated some observed data. The key challenge is that almost all mathematical models of real-world phenomena are complex in that they demand high-dimensional parameter space. This makes it prohibitively expensive to compute the probability density for some observation (also known as the likelihood in the Bayesian parlance).

Traditionally, sampling from arbitrarily complex distributions has been done using Markov Chain Monte Carlo (MCMC), which proposes candidate parameter values and accepts or rejects them based on their relative likelihood (Metropolis et al. 1953; Hastings 1970; Gelman et al. 2013). By definition, the method requires evaluating the likelihood, which results in the cost of running the full forward simulation accumulating at each step in the chain. MCMC also presupposes that the model can be defined analytically, which is not the case for most complex models (Pretorius 2005; Campanelli et al. 2006; Batalha et al. 2017). Even in cases where evaluating likelihood is possible, MCMC can get “stuck” in some local minima when the posterior has tricky geometry, involving multi-modal distributions, strong parameter degeneracies, or highly anisotropic structure (Roberts et al. 1997; Brooks et al. 2011).

Simulation-based inference (SBI) has emerged as the gold standard framework for tackling these challenges

(Cranmer et al. 2020). Instead of evaluating the explicit likelihood, which has proven to be unknown or intractable, we can now leverage model simulations to learn approximate posteriors straight from simulated training dataset. Needless to say, most exciting SBI methods today rely on neural density estimation (Papamakarios and Murray 2016; Greenberg et al. 2019; Lueckmann et al. 2017), which incorporates some flavors of deep neural networks to learn the posterior distribution. Trained on pairs of parameters and simulations, these neural networks learn the mapping between observations and the corresponding parameters that generated them. The cost of training simulations is also “amortized” over time, allowing the neural network to be trained once and used to repeatedly infer new observations.

Despite these advances, current SBI methodologies face scalability issues for a myriad of reasons. The first and perhaps most prominent is that as the number of parameters increases, so does the computational cost of generating simulations in high-dimensional space. In high dimensions, the volume of parameter space expands exponentially that achieving adequate coverage requires an unacceptably large number of simulations.

The second shortcoming lies in the method by which we generate model realizations. Often, the de facto technique involves a uniform sampling of some parameter combinations within some Latin hypercube (McKay et al. 1979)

$$\theta \sim \mathcal{U}([\theta_{\min,1}, \theta_{\max,1}] \times \cdots \times [\theta_{\min,d}, \theta_{\max,d}]), \quad (1)$$

which has the advantage of covering the whole region of interest. But there is no such thing as a free lunch: the training distribution, sampled with this hypercube, will

exhibit a sharp discontinuity at the boundary, where the density collapses to zero. Driven by stochastic gradient descent, our neural network can interpolate the internal regions at ease but falls short near those boundaries. We believe that the boundary coverage should ideally be improved with infinite computational resources. But due to practical constraints, we must prioritize the center of the parameter space, resulting in undersampling on the outskirts. The absence of adequate support corrodes the posterior estimates. We also believe that such pathology will be exacerbated in higher dimensions by the “curse of dimensionality,” (Bellman 1961) according to which the volume fraction of points close to boundaries scales exponentially.

A third, tightly related limitation arises from the standard workflow in large-scale simulation campaigns, like the CAMELS project (Cosmology and Astrophysics with Machine Learning Simulations), where simulations are generated once with a fixed prior, published as public datasets, and then reused for inference across multiple studies (Villaescusa-Navarro et al. 2021, 2023). Any researcher who wants to perform inference with a prior that stretches beyond the initial simulation bounds will likely obtain lower posterior quality, as the emulator has never seen the training data in that region before. Without access to the original, prohibitively priced simulator, this rigidity imposes a systematic constraint on all future analyses.

To address these problems, we present a novel parameter sampling technique for simulation-based inference. In particular, we propose TAILED-UNIFORM, a hybrid proposal that combines uniform cores with smooth Gaussian tails and show that empirically, neural nets trained with our proposal achieve superior performance near the boundary and maintain consistent posterior quality across the overall parameter space. We also validate TAILED-UNIFORM on both synthetic benchmarks and realistic cosmological inference tasks, revealing that TAILED-UNIFORM enables robust inference with competitive simulation budgets.

The remainder of this paper is organized as follows. Section 2 reviews the fundamentals of simulation-based inference and neural posterior estimation, then introduce the mathematical formulation of the TAILED-UNIFORM proposal. Our method is validated on a controlled Gaussian linear benchmark in Section 3, and ablation studies are conducted to investigate the effects of dimensionality, simulation budget, and tail width. In Section 4, we demonstrate practical utility on cosmological parameter inference from matter power spectra, showing how TAILED-UNIFORM mitigates boundary pathology along degenerate parameter directions. Lastly, we provide a summary of our findings in Section 5.

2. METHODS

This section lays the groundwork for our proposed sampling technique. We propose TAILED-UNIFORM and show that the boundary discontinuities that afflict conventional uniform sampling are ameliorated, at least in theory, by smooth Gaussian tails.

2.1. Simulation-based Inference

For clarity, we review, very briefly first, the notion of simulation-based inference (SBI). We define a simulator

as a black box function \mathcal{M} that aims to mimic a physical, possibly stochastic, generative process; that is, it maps some parameters $\theta \in \Theta \subseteq \mathbb{R}^d$ to observable quantities $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^D$ (Cranmer et al. 2020).

To generate each data-parameter pair (θ_i, \mathbf{x}_i) , we: (1) sample parameters $\theta_i \sim \tilde{\mathcal{P}}(\theta)$ from the proposal prior, (2) run the forward model $\mathbf{x}_i = \mathcal{M}(\theta_i) + \epsilon_i$ where ϵ_i represents some systematic noise, and (3) store the resulting data-parameter pairs. We can then feed these into an inverse model to make inferences on some real data $\mathbf{x} = \mathbf{x}_0$.

Bayesian workflow The goal of SBI is to discern the posterior distribution, which encodes all the “relevant” information (that is, point estimates, error bars, and the shape of the posterior) pertaining to the parameters of interest given an observed dataset \mathbf{x}_0 . In most cases, such posterior is not available in closed form, but because it is contingent upon the provided data, we can invoke Bayes’ theorem to write

$$\mathcal{P}(\theta | \mathbf{x}_0) = \frac{\mathcal{P}(\mathbf{x}_0 | \theta) \mathcal{P}(\theta)}{\mathcal{P}(\mathbf{x}_0)}. \quad (2)$$

Here, $\mathcal{P}(\mathbf{x}_0 | \theta)$ is the likelihood (the probability of the data given certain parameters), $\mathcal{P}(\theta)$ is the prior (the proxy for our initial belief regarding the true distribution of the parameters), and $\mathcal{P}(\mathbf{x}_0) = \int \mathcal{P}(\mathbf{x}_0 | \theta) \mathcal{P}(\theta) d\theta$ is the evidence, which is merely a normalization constant.

Neural posterior estimation (NPE) Our objective then is to construct a neural architecture $q_w(\theta | \mathbf{x})$ with weights w that outputs a probability distribution over θ which approximates mapping from observational data to full posterior distribution $\mathcal{P}(\theta | \mathbf{x})$ (Papamakarios and Murray 2016; Greenberg et al. 2019). To train this network to accurately represent the conditional probability, we maximize the joint likelihood of our training data $\mathcal{D}_{\text{train}}$. This is akin to minimizing the negative log-probability loss

$$\mathcal{L}_{\text{NPE}} := -\mathbb{E}_{\mathcal{D}_{\text{train}}} \log q_w(\theta | \mathbf{x}) \quad (3)$$

$$= -\mathbb{E}_{\mathcal{D}_{\text{train}}} \log \left[\frac{\mathcal{P}(\theta)}{\tilde{\mathcal{P}}(\theta)} q_w(\theta | \mathbf{x}) \right], \quad (4)$$

where the expectation $\mathbb{E}_{\mathcal{D}_{\text{train}}}$ in the loss function is taken over all parameter-observation pairs $\{(\theta_i, \mathbf{x}_i)\}_{i=1}^N$ in the training dataset.

If such training data exhibits sufficient diversity and the neural architecture possesses adequate coverage of the parameter space $\Theta \subset \mathbb{R}^d$, then minimization of \mathcal{L}_{NPE} over the network weights w will drive $q_w(\theta | \mathbf{x})$ toward the true posterior distribution $\mathcal{P}(\theta | \mathbf{x})$ (Hornik et al. 1989). As a result, we can sample at inference time by evaluating the log-density of simulated parameters beforehand during training.

Priors It is crucial to make clear the difference between proposal and assumed priors (Cranmer et al. 2020):

- The proposal prior $\tilde{\mathcal{P}}(\theta)$, by definition, is the empirical distribution of parameters θ found in the training dataset $\mathcal{D}_{\text{train}}$. For one thing, it will depend on our sampling strategy (i.e., a Latin hypercube) used to create training simulations.

- On the other hand, our preconceived notion of the global parameter distribution is encoded in the assumed prior $\mathcal{P}(\theta)$, which, we should note, is an experimental design choice. It reflects scientific understanding, theoretical constraints, or previous empirical knowledge about plausible parameter values. In cosmological inference—for example—we may adopt priors centered on values from landmark observational campaigns like the Planck satellite measurements (Planck Collaboration et al. 2020).

Broadly speaking, we can choose the assumed and proposed priors to be identical. This alignment, however, is not always preferable. Take a scenario where we sample training simulations from a uniform distribution $\tilde{\mathcal{P}}(\theta) = \mathcal{U}([\theta_{\min}, \theta_{\max}])$ to ensure adequate coverage, while our scientific understanding may instead motivate a Gaussian prior $\mathcal{P}(\theta) = \mathcal{N}(\mu, \Sigma)$. In such a case, the neural network can only learn to assign reasonable probability densities to regions where it has observed training examples, that is, $\tilde{\mathcal{P}}(\theta) > 0$. Furthermore, since we do not have any support in regions that lack training data, the neural output will tend to zero $q_w(\theta|\mathbf{x}) \approx 0$, leading to poor inference quality near the boundaries.

2.2. The Tailed-Uniform Proposal

We will now describe the theoretical foundations of the TAILED-UNIFORM, our suggested proposal that offers a smooth, continuously differentiable density across \mathbb{R}^d . We hope that by assigning weights beyond the primary region of interest, we can alleviate the sharp discontinuities at the boundary.

One-dimension derivation To kick off the discussion, we consider a one-dimensional ($d = 1$) case and define the probability density function as

$$\tilde{\mathcal{P}}_{\text{TailedUniform}}(x; a, b, \sigma) = \begin{cases} A \cdot \mathcal{N}(a, \sigma^2), & x \leq a \\ B \cdot \mathcal{U}(a, b), & x \in [a, b] \\ A \cdot \mathcal{N}(b, \sigma^2), & x \geq b \end{cases} \quad (5)$$

where a and b define the boundaries of the uniform core region, and σ controls the width of the Gaussian tails. We also impose continuity at the boundary to guarantee that our distribution is well-defined, which in turn establishes the values of the normalization constants.

Putting things together, we have

$$A = \frac{\sqrt{2\pi\sigma^2}}{\sqrt{2\pi\sigma^2} + (b-a)} \text{ and } B = \frac{b-a}{\sqrt{2\pi\sigma^2} + (b-a)}, \quad (6)$$

which ensures that the distribution integrates to unity across the entire domain, and the Gaussian tail matches the uniform core at the left and right boundaries.

Generalizing to higher dimensions To generalize to a multivariate parameter space $\theta \in \mathbb{R}^d$, we write it as a product of independent uni-variate TAILED-UNIFORM distributions

$$\tilde{\mathcal{P}}_{\text{TailedUniform}}(\theta; \mathbf{a}, \mathbf{b}, \boldsymbol{\sigma}) = \prod_{i=1}^d \tilde{\mathcal{P}}_{\text{TailedUniform}}(\theta_i; a_i, b_i, \sigma_i), \quad (7)$$

where $\mathbf{a} = (a_1, \dots, a_d)$ and $\mathbf{b} = (b_1, \dots, b_d)$ define the hypercube boundaries, and $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_d)$ controls the tail-widths in each dimension.

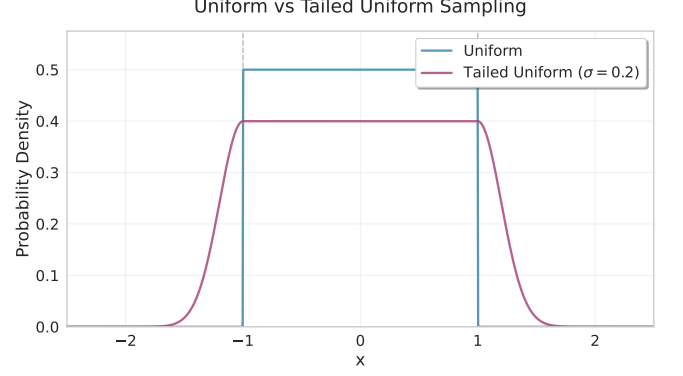


FIG. 1.— The standard uniform distribution (blue) has constant probability density between -1 and 1 with zero probability outside this range. The tailed uniform distribution (magenta) maintains a uniform density in the same central region $[-1, 1]$ but extends beyond these boundaries with Gaussian tails characterized by standard deviation $\sigma = 0.2$.

Here, the key hyperparameter is the tail width σ_i , which controls the balance between boundary smoothness and sampling efficiency. Smaller values of σ_i concentrate more samples within the target region $[a_i, b_i]$ but provide less smoothing at the boundaries. Larger values, on the other hand, provide better boundary smoothness but allocate more resources to regions far from the true posterior support. In practice, we recommend setting σ_i as about 10 – 40 percent of the box width; that is,

$$\sigma_i = \alpha \cdot (b_i - a_i), \quad (8)$$

where $\alpha \in [0.1, 0.4]$ (see Section 3.3 for more detail).

3. TOY PROBLEM

We first demonstrate the advantages of TAILED-UNIFORM sampling on a simple toy problem where analytical ground truth is available. This controlled setting allows us to isolate the effects of boundary pathology without confounding factors from complex simulators, providing clear intuition before tackling realistic scientific inference tasks in Section 4.

3.1. Gaussian Linear Task

We formulate a toy problem using the Gaussian Linear task from the simulation-based inference benchmark (sbibm) (Lueckmann et al. 2021).

Model As a proof-of-concept, we consider a simple, linear, 2-dimensional Gaussian simulator $\mathcal{M} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$, which is defined as

$$\mathcal{M}(\theta) = \mathbf{A}\theta = \mathbf{I}_2\theta = \theta \quad (9)$$

with a standard Gaussian prior $\mathcal{N}(\mathbf{0}, \mathbf{I}_2)$ centered in the bounded parameter space $\theta \in [-1, 1]^2$. We hence have a likelihood of the form

$$\mathcal{P}(\mathbf{x} | \theta) = \mathcal{N}(\mathbf{x}; \theta, \mathbf{I}_2) = \frac{1}{2\pi} \exp\left(-\frac{1}{2}\|\mathbf{x} - \theta\|^2\right), \quad (10)$$

which implies that both the likelihood and prior are Gaussian.

Ground Truth We can leverage the aforementioned conjugate Gaussian-Gaussian structure (Gelman et al. 2013) to obtain a closed-form analytical posterior $\mathcal{P}(\theta | \mathbf{x}) = \mathcal{N}\left(\frac{\mathbf{x}}{2}, \frac{\mathbf{I}_2}{2}\right)$, where the posterior covariance $\frac{\mathbf{I}_2}{2}$

results from $(\mathbf{I}_2^{-1} + \mathbf{I}_2^{-1})^{-1} = \frac{\mathbf{I}_2}{2}$. This gives us reference posterior samples from which to evaluate performance.

Training Our goal is to investigate two distinct strategies for generating the training data:

1. Baseline (*Uniform* proposal): Training data is generated by sampling parameters *uniformly* across the 2-dimensional hypercube $\tilde{\mathcal{P}}_{\text{Uniform}} = \mathcal{U}([-1, 1]^2)$.
2. Proposed (TAILED-UNIFORM proposal): Training data is generated using our hybrid distribution $\tilde{\mathcal{P}}_{\text{TailedUniform}}(\mathbf{a} = [-1, -1], \mathbf{b} = [1, 1], \boldsymbol{\sigma} = [0.2, 0.2])$, where the tail width $\sigma_i = 0.1 \times (b_i - a_i)$ represents 10% of each parameter range.

For both proposals, we generate $N = 4000$ simulation pairs $\{(\boldsymbol{\theta}_i, \mathbf{x}_i)\}_{i=1}^N$ by: (1) sampling $\boldsymbol{\theta}_i$ from the respective proposal distribution, (2) running the forward simulator $\mathbf{x}_i = \mathcal{M}(\boldsymbol{\theta}_i) + \boldsymbol{\epsilon}_i$ with $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_2)$ and (3) training identical neural posterior estimators (NPE) on each dataset. Note again that the assumed prior in both cases is a Gaussian distribution centered at $(0, 0)$. In what follows, we will associate *Uniform* and TAILED-UNIFORM to the NPE trained with these proposals (rather than the proposal itself).

Model Architectures We experiment with ensembles of neural density estimators comprising Masked Autoregressive Flows (MAF) (Papamakarios et al. 2017) and Masked Autoencoder for Distribution Estimation (MADE) (Germain et al. 2015) architectures, each featuring 16 hidden features and 5 coupling layers. Training uses a batch size of 64 and a learning rate of 5×10^{-5} . Our pipeline is built on the LtU-ILI (Learning the Universe—Implicit Likelihood Inference) framework (Ho et al. 2024), which offers standardized data handling and ensemble utilities for simulation-based inference. It is worth pointing out that there is nothing special about the above configuration. The Gaussian Linear task, with its analytically tractable posterior and low dimensionality, is simple enough for any garden-variety neural network to learn the target distribution.

3.2. Validation

Model validation aims to ascertain whether the approximate posteriors generated by our neural density estimators $q_w(\boldsymbol{\theta} | \mathbf{x}_0)$ will be “good enough” when faced with real unlabeled observational data \mathbf{x}_0 . This is no trivial matter, as we do not a priori know the learned weights in the hidden layers, and the internal workings of neural networks are more or less invisible.

The bias-variance tradeoff One useful proxy for gauging the quality of our learned posterior is the degree to which it can concentrate probability mass around the true parameter $\boldsymbol{\theta}_{\text{true}}$. Practically worthless, for instance, will be a model that regurgitates the prior distribution, as it offers no additional insight whatsoever. After all, the goal of inference is to extract as much information as possible from the observed data \mathbf{x}_0 to narrow down the possible parameter values $\boldsymbol{\theta}$. Yet, if one only optimizes for the training set, the resulting model might end up overfitting and fail to generalize when exposed to new data.

Achieving the sweet spot between bias and variance (Gelman and Rubin 1992) is thus the characteristic of

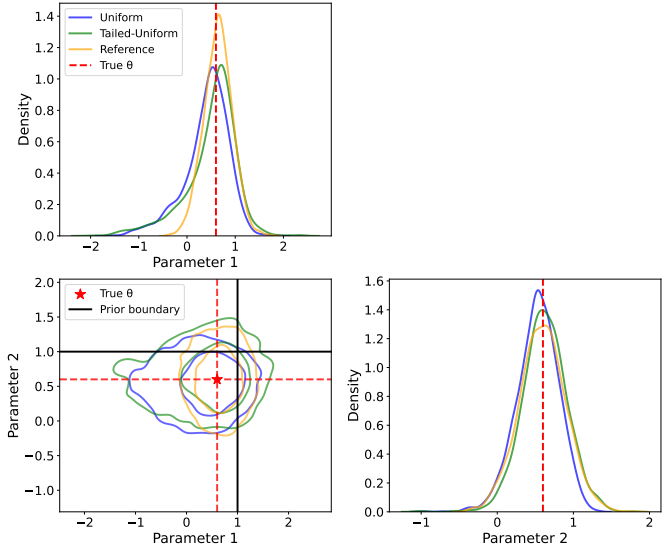


FIG. 2.— Corner plots comparing posterior estimation performance for the boundary test case. The TAILED-UNIFORM (green) demonstrates superior boundary behavior compared to the *Uniform* (blue), closely matching the MCMC reference (yellow). The red dashed line indicates the true parameter value.

a “good” model, as we want a model that is versatile enough to extract genuine information from observations (low bias) while remaining generalizable across different realizations of the data (low variance).

TABLE 1
POSTERIOR PERFORMANCE FOR $\boldsymbol{\theta}_{\text{true}} = (0.6, 0.6)$.

Method	Parameter Estimates		C2ST
	θ_1	θ_2	
Reference	0.446 ± 0.300	0.641 ± 0.308	0.500
TAILED-UNIFORM	0.339 ± 0.547	0.708 ± 0.286	0.458
<i>Uniform</i>	0.259 ± 0.446	0.588 ± 0.272	0.409
C2ST Improvement			+12%

3.2.1. Single-Point Analysis

As a first pass, we select a test point $\boldsymbol{\theta}_{\text{true}} = (0.6, 0.6)$ near the boundary of parameter space—precisely where we anticipate the *Uniform* proposal to struggle most. We generate an observation $\mathbf{x}_0 = \mathcal{M}(\boldsymbol{\theta}_{\text{true}}) + \boldsymbol{\epsilon}$ for this test point, then draw $M = 1000$ posterior samples from each method (MCMC reference, TAILED-UNIFORM, and *Uniform*). The goal is to assess how well each method constrains the parameters by examining the distribution and concentration of samples around the true value.

Quantitative Metrics To quantify the error between our learned posteriors and the true posterior, we employ the Classifier Two-Sample Test (C2ST) (Lopez-Paz and Oquab 2017), which is the most sensitive metric for evaluating distributional fidelity in SBI models Lueckmann et al. (2021). The C2ST trains a logistic regression classifier to differentiate between samples from two distributions. If the two distributions are identical, the classifier can do no better than random guessing, yielding a score of 0.5 (zero distributional error). Conversely, scores that deviate from 0.5 signal that the classifier can

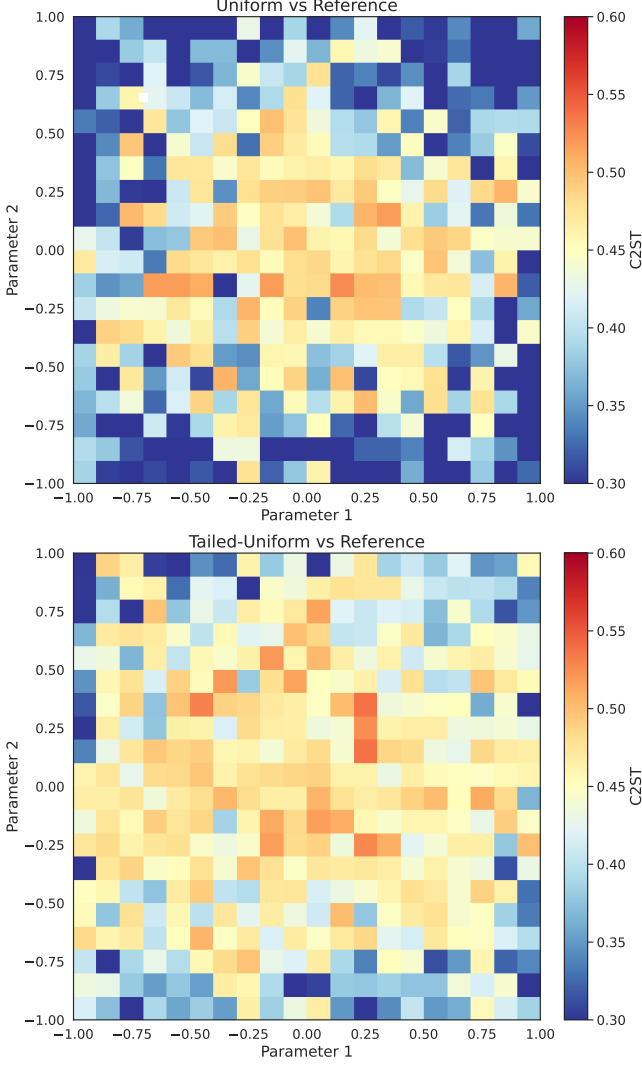


FIG. 3.— Spatial distribution of C2ST performance across the parameter space, with blue regions indicating poor distributional matching ($\text{C2ST} \ll 0.5$) and red/orange regions indicating good performance ($\text{C2ST} \approx 0.5$). **Top:** *Uniform* versus analytical reference, showing systematic boundary degradation with glaring blue regions near parameter space edges. **Bottom:** TAILED-UNIFORM versus reference, demonstrating consistent performance across the entire parameter space.

easily distinguish between the two distributions.

As anticipated in Figure 2, the TAILED-UNIFORM (orange) keeps a tight concentration around the true parameter values, while the *Uniform* (blue) leaks into unusable parameter regions. Table 1 further corroborates TAILED-UNIFORM’s superior performance: the *Uniform* score of 0.409 shows that a classifier can readily distinguish between its posterior samples and the reference, whereas TAILED-UNIFORM’s C2ST score of 0.458 indicates that its learned posterior is almost identical (close to 0.5) to the true posterior.

3.2.2. Spatial Performance Analysis

Next we need a way to assess TAILED-UNIFORM’s performance across the entire parameter space. And so, we discretize a uniform rectangular grid with $n = 20$ test locations along each dimension of the parameter space $[-1, 1]^2$, yielding a total of 400 evaluation points. We

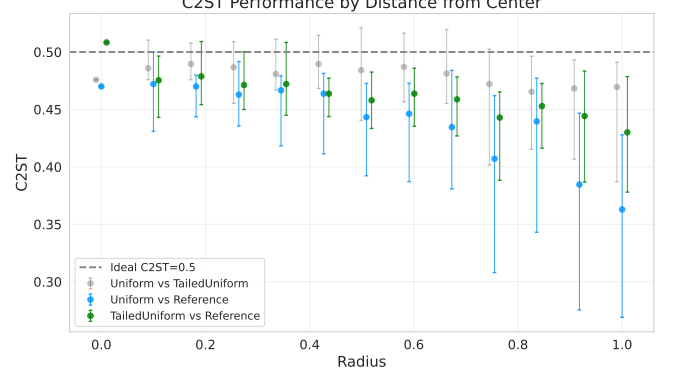


FIG. 4.— C2ST performance degradation as a function of distance from parameter space center. The blue curve reveals systematic deterioration of *Uniform* near boundaries, while the green curve demonstrates that TAILED-UNIFORM maintains consistent performance across all radii. The gray curve quantifies the increasing divergence between methods, with boundary regions showing substantial differences in posterior approximation quality. Values closer to 0.5 indicate better distributional matching.

want to exhaustively sample areas of varying distances, ranging from the center of the priors, where traditional techniques work well, to points near the boundaries. As before, we generate observations for each test point, draw $M = 1000$ posterior samples from all three methods, and calculate pairwise C2ST scores at each test location.

Figure 3 shows how *Uniform*’s performance deteriorates as it drifts away from the central regions, as evidenced by the proliferation of blue patches near the boundary. In contrast, our TAILED-UNIFORM maintains more red and orange coloration (with C2ST scores around 0.44–0.49) across the parameter space.

Furthermore, we observe that the C2ST degradation pattern itself is rotationally invariant; that is, its performance depends only on the radial distance (no angular dependence). This spherical symmetry allows us to visualize the performance trends by binning the test points as a function of their distance from the center, as seen in Figure 4. Evidently, TAILED-UNIFORM exhibits robustness and mitigates the undesired boundary pathology that characterizes the traditional sampling technique.

3.3. Prior Hyper-Optimization

Our subsequent evaluations aim to answer three operational questions: the optimal tail width, whether additional training data can mitigate boundary pathology, and whether TAILED-UNIFORM’s advantages persist in higher dimensions.

3.3.1. Sensitivity to the prior tail widths

The tail width σ represents a bias-variance tradeoff: narrower tails mitigate the boundary pathologies, while wider tails allocate more simulations outside the primary region of interest, thus decreasing constraining power. To assess the importance of the tails, we train one *Uniform* baseline and several TAILED-UNIFORM models, each with different σ values as fractions of prior width: $\alpha \in \{0.01, 0.05, 0.1, 0.2, 0.4\}$.

Figure 5 reveals that performance improves monotonically with tail width. At the prior boundary ($r = 1.0$), C2ST scores improve from 0.331 ($\sigma = 0.01$, equivalent

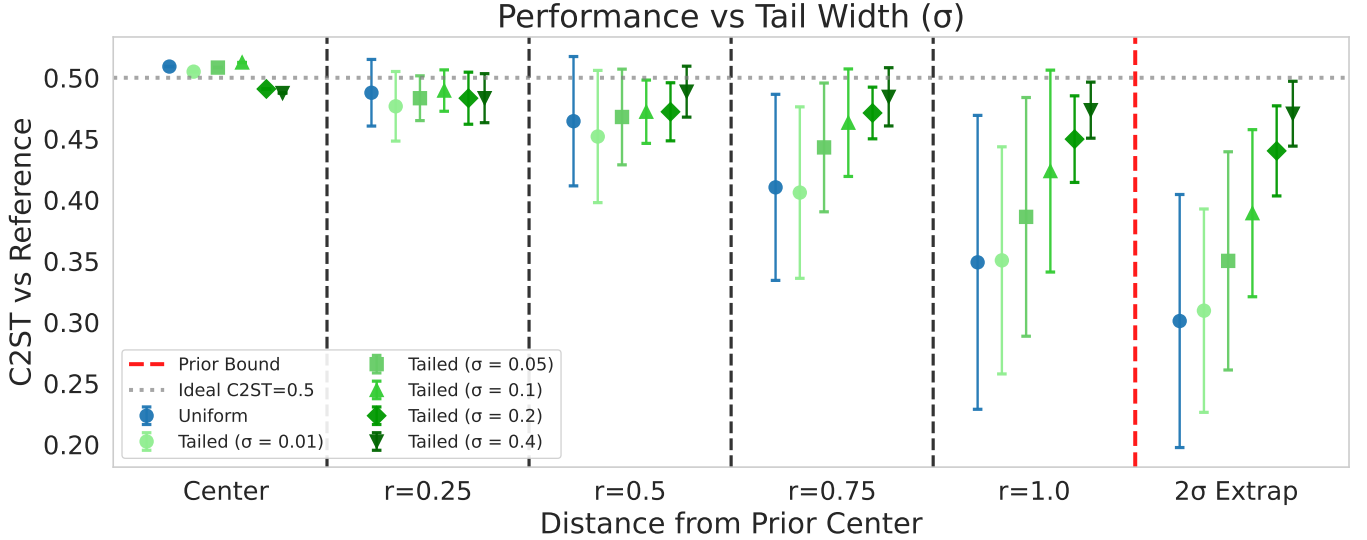


FIG. 5.— C2ST performance versus distance from center, stratified by tail width σ . The *Uniform* baseline (blue) exhibits systematic boundary degradation. The TAILED-UNIFORM’s with varying tail widths (shown in green; lighter shades for smaller σ values) exhibit robustness across a broad range of tail widths. Notably, performance becomes quite good for $\sigma \geq 0.2$

to *Uniform*) to 0.465 ($\sigma = 0.4$, near-ideal), demonstrating the value of wider tails. As might be expected given its higher sample density in the core region, *Uniform* performs better close to the center of the prior. However, it is important to note that while *Uniform* achieves marginally superior performance near the center, this comes at the expense of significant degradation in performance near the boundaries. TAILED-UNIFORM, on the other hand, sacrifices a reasonably small amount of accuracy to maintain robust inference across the entire parameter space. The reassuring robustness indicates that even moderate tail samples suffice for learning smooth density transitions, which will prove insightful in higher dimensions. At $\sigma = 0.01$, TAILED-UNIFORM degenerates as anticipated toward *Uniform* behavior with characteristic boundary degradation, confirming that the tails themselves (not some arbitrary artifacts) are responsible for the observed improvements.

3.3.2. Ablation on the number of simulations

If insufficient data in the near-boundary regions causes degradation, just increasing the training set size could very well solve the problem. After all, neural networks are universal function approximators, so given enough training data, they should be able to learn complex mappings, including sharp transitions at parameter space boundaries. To interrogate this hypothesis, we test an ablation over the training set size, varying $N \in \{2000, 4000, 8000, 16000\}$ while holding all other hyperparameters fixed.

Empirically, we find in Figure 6 that increasing the number of simulations barely helps with the NPE’s performance, suggesting that the boundary degradation stems not from absence of data, but from an inbuilt discontinuous mapping that cannot be learned from samples within the prior alone. Since no amount of supplementary data can smooth a non-smooth function, the hard truncation at prior boundaries poses a fundamental learning barrier. In contrast, we see that every single TAILED-UNIFORM ($N = 2000$) outperforms even the largest *Uniform* ($N = 16000$) with 8 times more training

data, demonstrating the impact of different proposals. One can achieve better posterior quality at a fraction of the computational cost just by sampling in a judicious manner.

3.3.3. Scaling to higher dimensions

While we have shown TAILED-UNIFORM to be effective in two dimensions, the critical question is whether this benefit transfers to higher-dimensional parameter spaces, which are necessary for the majority of astrophysical inference problems. By the curse of dimensionality, high-dimensional spaces become ever more sparse, as the volume of a d -dimensional hypercube grows exponentially while the number of samples remains fixed. Corners, edges, and boundary regions come to dominate the landscape.

We now examine the probability of obtaining samples near boundaries. For a d -dimensional hypercube $[0, 1]^d$, we define the interior region as the set of points at least a distance ε away from any boundary, which occupies a fraction $P_{\text{int}}^{\text{Uniform}}(d, \varepsilon) = (1 - 2\varepsilon)^d$ of the total volume. Complementarily, the probability that a uniformly-sampled point lies in the boundary shell (within ε of any face) grows exponentially as $P_{\text{bdry}}^{\text{Uniform}}(d, \varepsilon) = 1 - (1 - 2\varepsilon)^d$, which approaches unity at high dimensions.

The Tail Allocation Budget For TAILED-UNIFORM, samples extend beyond the hypercube boundaries into tail regions. With independent marginals, the probability that a sample lands inside the d -dimensional hypercube is the product

$$P_{\text{cube}}^{\text{TailedUniform}}(d) = \prod_{i=1}^d B_i = B^d = \left(\frac{1}{1 + \alpha\sqrt{2\pi}} \right)^d, \quad (11)$$

assuming identical α across dimensions. This means the probability of sampling in the tail regions grows as $P_{\text{tail}}^{\text{TailedUniform}}(d) = 1 - B^d$. By the same token, the probability of landing in the interior becomes $P_{\text{int}}^{\text{TailedUniform}}(d, \varepsilon) = B^d \cdot (1 - 2\varepsilon)^d$, while the boundary shell captures $P_{\text{bdry}}^{\text{TailedUniform}}(d, \varepsilon) = B^d [1 - (1 - 2\varepsilon)^d]$.

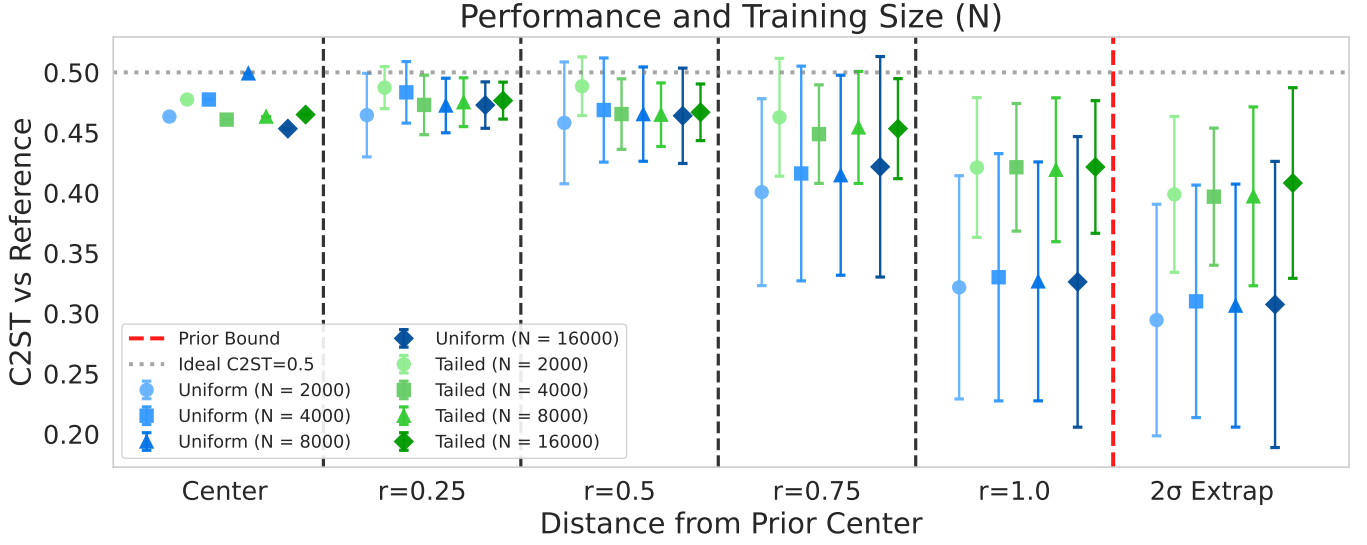


FIG. 6.— C2ST performance versus distance from parameter space center, stratified by training set size N (circles: 2000, squares: 4000, triangles: 8000, diamonds: 16000). Increasing N provides minimal benefit for both proposals.

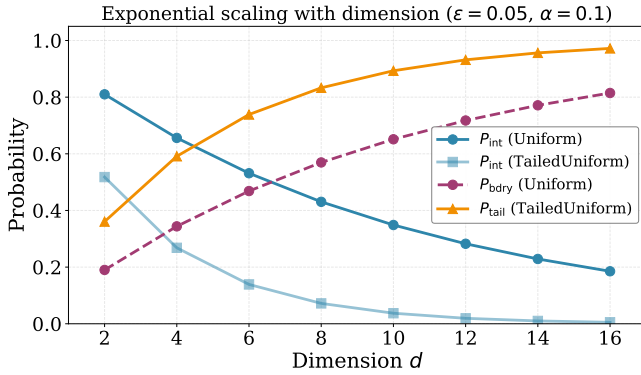


FIG. 7.— Exponential scaling of sample allocation with dimension for *Uniform* and TAILED-UNIFORM proposals ($\epsilon = 0.05$, $\alpha = 0.1$). The probability of sampling in the interior region P_{int} (solid blue) decays exponentially for both methods, while the boundary region P_{bdry} (dashed purple) grows correspondingly for *Uniform*. For TAILED-UNIFORM, samples increasingly concentrate in the tail regions P_{tail} (solid orange), rising from 36% at $d = 2$ to 97% at $d = 16$. The crossing point around $d = 4$ marks where tail allocation exceeds interior allocation. Despite this apparent “waste,” we see that TAILED-UNIFORM maintains superior performance.

Figure 7 summarizes these quantities for representative dimensions.

At first glance, it might seem that by bleeding samples into the tails, too many samples will wind up outside the region of interest, thus diminishing the overall performance of TAILED-UNIFORM. To refute this suspicion, we evaluate NPE trained with *Uniform* and TAILED-UNIFORM proposals on a Gaussian linear inference task across $d \in \{4, 8, 16\}$ dimensions.

Boundary Dominance in Higher Dimensions: Figure 8 shows that while both methods degrade with increasing dimensionality, TAILED-UNIFORM regularly outperforms *Uniform*, with the performance gap widening at higher dimensions. This somewhat paradoxical finding suggests that the ratio of information gained from exterior support to cost of interior sample dilution remains favorable or even improves with dimension. Extra samples from TAILED-UNIFORM flow into the tail regions, surrounding the boundaries and offering regular-

ization. Indeed, at $d = 16$, boundaries are so dominant that the marginal value of additional interior samples is low, as the network has ample interior coverage even when 99.99% of TAILED-UNIFORM samples lie outside the region of interest. Meanwhile, the marginal value of tail samples is high, as they are the only source of information near the boundary, which explains why TAILED-UNIFORM still achieves somewhat better C2ST scores than *Uniform* at the boundary and in the extrapolation regime (2σ beyond the prior boundary). We acknowledge that at very high dimensions, dilution effects will inevitably dominate, but these regimes are beyond the scope of most practical simulation-based inference tasks.

4. SCIENCE EXPERIMENT

4.1. Matter Power Spectrum

We next demonstrate the versatility of the TAILED-UNIFORM proposal in a popular inference benchmark in cosmology: the inference of cosmological parameters from the matter power spectrum (Dodelson and Schmidt 2020). To recreate the problem, we aim to predict the posterior distribution of two cosmological parameters $\theta = (\Omega_m, h)$ using the forward simulator $\mathcal{M} : \mathbb{R}^2 \rightarrow \mathbb{R}^{64}$ that maps a two-dimensional parameter vector to the power spectrum $\mathbf{P} \in \mathbb{R}^{64}$. Here, Ω_m is the matter density parameter that governs the growth rate of perturbations and h is the dimensionless Hubble parameter, which determines the physical scale of the horizon at matter-radiation equality.

In modern cosmology, the power spectrum $\mathbf{P}(k)$ measures the amplitude of density contrast fields $\delta(\mathbf{k})$ in Fourier space and is strictly dependent on the magnitude of the wave vector $k = |\mathbf{k}|$ (Dodelson and Schmidt 2020). This can be thought of as the ensemble average of all universes with the same cosmological parameters but different initial fluctuations. We also assume a flat Λ CDM cosmology with $w_0 = -1$, $w_a = 0$, and negligible neutrino masses $m_\nu = 0$, and fix the remaining cosmological parameters to Planck 2018 (Planck Collaboration et al. 2020) values (the primordial scalar

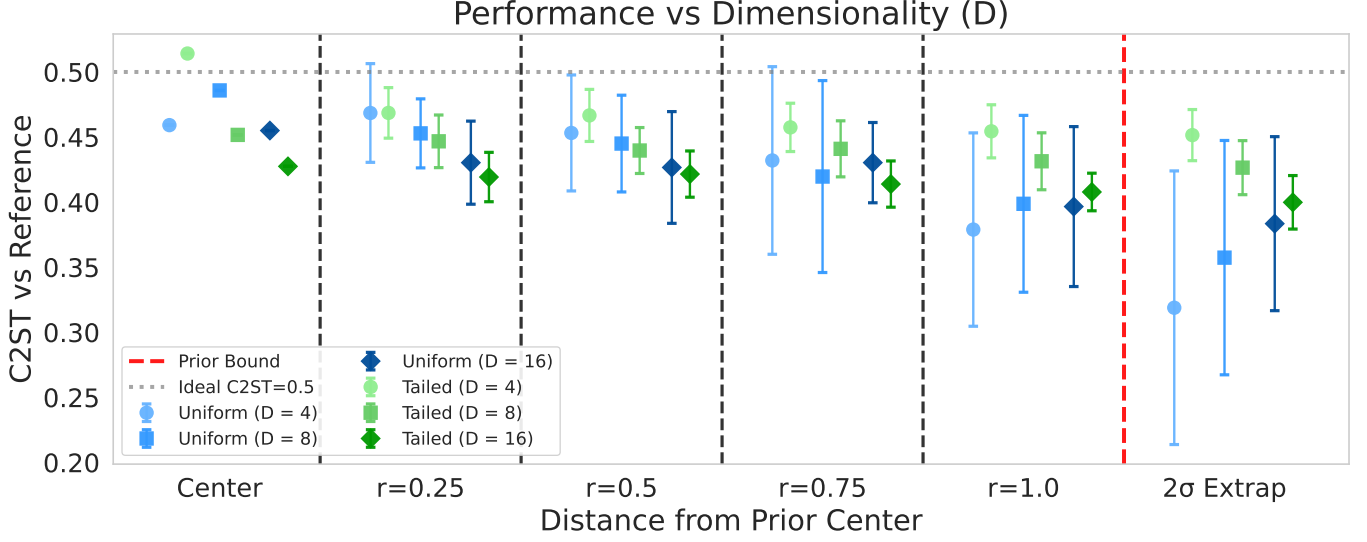


FIG. 8.— C2ST performance versus distance from parameter space center in different representative dimensions. TAILED-UNIFORM maintains superior performance with the advantage growing at higher dimensions, particularly near boundaries ($r \geq 0.75$) and in the extrapolation regime.

amplitude $A_s = 2.105 \times 10^{-9}$, the baryon density parameter $\Omega_b h^2 = 0.02242$, and the scalar spectral index $n_s = 0.9665$).

Model To circumvent the need to solve the perturbation equations numerically, we use the **syren-new** emulator (Sui et al. 2024), which leverages symbolic regression to produce a closed functional form for the theoretical power spectrum $\mathbf{P}_{\text{theory}}$ at negligible compute power. The power spectrum from **syren-new** is not the observed power spectrum per se, as our universe (armed with specific initial conditions) is merely one realization of some stochastic process. If we observe a sufficiently large volume, different regions within that volume will serve as independent samples of the same underlying process (Baumann 2022). But since galaxy surveys cannot cover an infinite volume, the spatial average will not be *exactly* equal to the ensemble average.

For a comoving survey volume of $V = L^3$, the longest wavelength that fits in the box (i.e., the fundamental mode) is $k_f = |\mathbf{k}| = \frac{2\pi}{L}$, and by assuming the universe to be homogeneous and isotropic (i.e., rotationally invariant) on large scales (Dodelson and Schmidt 2020), we then find the number of modes inside a spherical shell at wavenumber k with thickness $\Delta k = k_f$ to be

$$N_k = \frac{4\pi k^2 \Delta k}{k_f^3} = \frac{L^3 k^2 k_f}{2\pi^2}. \quad (12)$$

This is called “cosmic variance,” which arises from the limited number of independent modes available in our finite observation volume (Dodelson and Schmidt 2020; Baumann 2022).

Together, we can construct an observation as

$$\mathbf{x}_i = \mathbf{P}(k | \boldsymbol{\theta}_i) = \underbrace{\mathbf{P}_{\text{theory}}(k | \boldsymbol{\theta}_i)}_{\text{syren-new emulator}} + \underbrace{\boldsymbol{\epsilon}_i}_{\text{cosmic noise}}, \quad (13)$$

where $\boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}, \text{diag}(\sigma_1^2, \dots, \sigma_{N_k}^2))$ represents the heteroskedastic (scale-dependent) cosmic noise with a variance $\sigma_k^2 = \frac{2\mathbf{P}_{\text{theory}}^2(k | \boldsymbol{\theta}_i)}{N_k}$.

Ground Truth As in Section 3, we have the luxury of direct sampling via MCMC (Brooks et al. 2011), as the forward model is explicit and the joint likelihood $\mathcal{P}(\mathbf{x} | \boldsymbol{\theta})$ can be evaluated point-wise. We choose a truncated Gaussian prior centered at the midpoint of $\Omega_m \in [0.24, 0.40]$ and $h \in [0.61, 0.73]$. Assuming independent Gaussian noise in each k -bin, we can write the log-likelihood as

$$\log \mathcal{P}(\mathbf{x} | \boldsymbol{\theta}) = -\frac{1}{2} \sum_{i=1}^N \left(\frac{\mathbf{x}_i - P_{\text{theory},i}(\boldsymbol{\theta})}{\sigma_i(\boldsymbol{\theta})} \right)^2 - \sum_{i=1}^N \log \sigma_i(\boldsymbol{\theta}), \quad (14)$$

where the first term quantifies the chi-squared misfit between theory and observation, and the second term accounts for the uncorrelated nature of the noise. Putting things together, the log posterior we aim to sample is $\log \mathcal{P}(\boldsymbol{\theta} | \mathbf{x} = \mathbf{x}_0) = \log \mathcal{P}(\mathbf{x}_0 | \boldsymbol{\theta}) + \log \mathcal{P}(\boldsymbol{\theta})$.

To explore this posterior distribution, we run four MCMC walkers via the affine-invariant ensemble sampler in the **emcee** package (Foreman-Mackey et al. 2013), each generating 2000 samples after a burn-in period of 500 steps to ensure convergence. In total, we collect $4 \times 2000 = 8000$ posterior samples, which we treat as our ground truth.

Training and Model Architectures In this experiment, we choose a comoving box size of $L = 1000, h^{-1}\text{Mpc}$ and measure the power spectrum at $N = 64$ logarithmically spaced wavenumbers spanning from the fundamental mode $k_f = 2\pi/L$ to the Nyquist frequency $k_{\text{nyq}} = \pi N/L$, evaluated at the present epoch ($a = 1.0$). Recall again that we want to compare the performance of neural networks trained with two distinct proposals $\tilde{\mathcal{P}}_{\text{Uniform}}$ and $\tilde{\mathcal{P}}_{\text{TailedUniform}}$, which by the relative complexity of the problem, will most likely be affected by our architectural choices. One could now proceed by trial and error: training dozens of models with different configurations and choosing the ones with the best log-probability. But in an effort to avoid unnecessary angst, we choose to perform Bayesian optimization via the **Optuna** framework (Akiba et al. 2019) with the

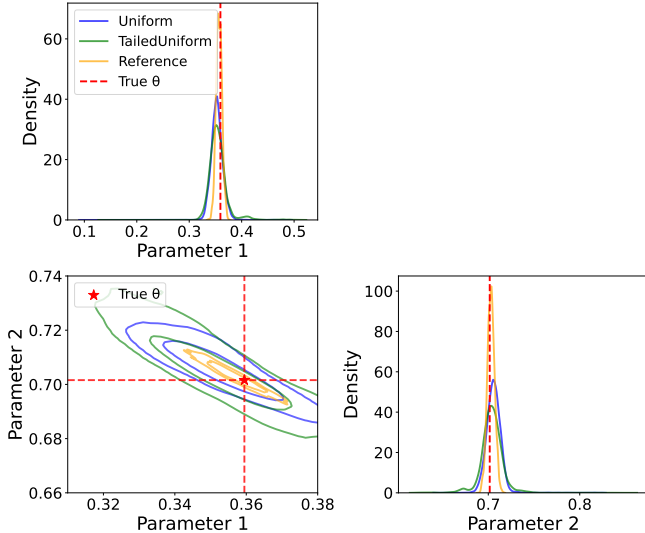


FIG. 9.— Corner plots comparing posterior estimation performance for test point $\theta_{\text{true}} = (0.359, 0.702)$ near the parameter space boundary. The TAILED-UNIFORM (green) closely matches the MCMC reference (yellow), maintaining tight concentration along the banana-shaped posterior. In contrast, the *Uniform* (blue) exhibits excessive spread along the degenerate direction, leaking probability mass into regions inconsistent with the data. The characteristic anti-correlation between Ω_m and h is visible in all methods, but only TAILED-UNIFORM accurately captures the posterior near the boundary.

search space consisting of the flow architecture, number of hidden features, number of transform layers, batch size, and learning rate. Once the hyperparameter search concludes, we rank all trained networks by their validation log-probability and ensemble the top 10 performing models.

4.2. Validation

To assess the performance of our competing proposals in this cosmological inference task, we adopt the same spatial evaluation framework as in Section 3. We should first discuss why Ω_m and h are given special consideration, which led us to select them as the inference parameter. There is a variety of factors, the most notable of which is the strong “degeneracies” between them.

Degeneracies Degeneracy refers to situations in which multiple parameters replicate the same observable. In our case, the matter power spectrum at small and intermediate scales is sensitive to the physical matter density, which happens to scale as $\Omega_m h$, rather than by Ω_m and h independently (Peacock and Dodds 1992). The two parameters are thereby anti-correlated.

As a result, the information embedded in our posterior will be weak along certain directions in the parameter space—that is, where Ω_m trades off against h such that $\Omega_m h$ remains about constant—leading to the slanted banana-shaped contours (as also reflected in the posteriors’ performance in Figure 10). Such geometric degeneracies in the power spectrum from different combinations of Ω_m and h , combined with the likelihood being non-linear, make this a challenging inference problem.

4.2.1. Single-Point Analysis

To study this phenomenon in detail, we first select a test point $\theta = (0.359, 0.702)$ positioned 20% from the

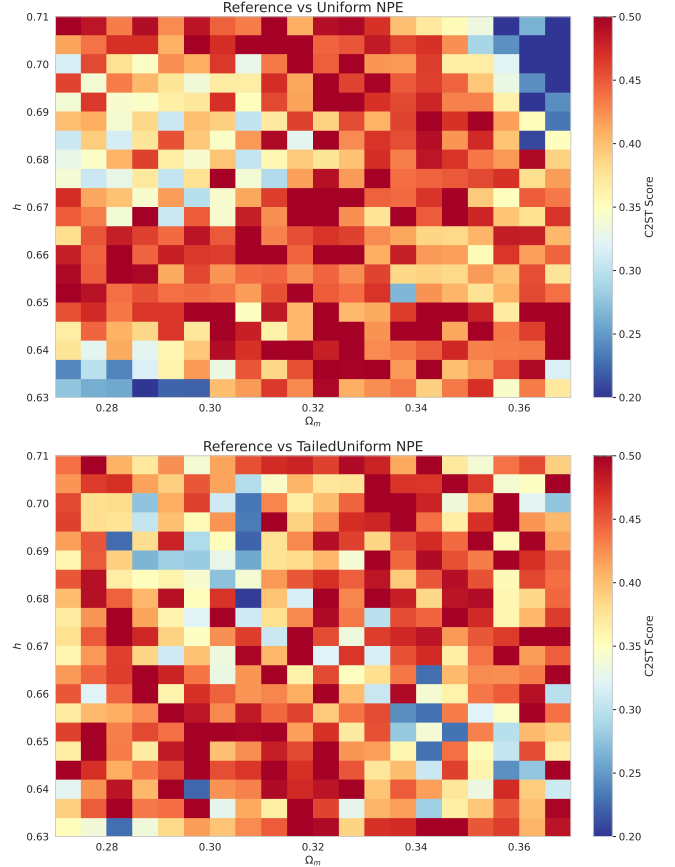


FIG. 10.— Spatial distribution of C2ST performance across the (Ω_m, h) parameter space, with darker blue regions indicating poor distributional matching ($\text{C2ST} \ll 0.5$) and orange-red regions indicating good performance ($\text{C2ST} \approx 0.5$). **Top:** *Uniform* versus MCMC reference, showing systematic degradation with blue patches near parameter space edges and corners. **Bottom:** TAILED-UNIFORM NPE versus reference, demonstrating consistent orange coloration across the parameter space.

top-right corner. We generate observations and draw $M = 1000$ posterior samples from the MCMC reference, *Uniform*, and TAILED-UNIFORM.

TABLE 2
POSTERIOR PERFORMANCE FOR TEST POINT $\theta_{\text{true}} = (0.359, 0.702)$

Method	Parameter Estimates		C2ST
	Ω_m	h	
Reference	0.357 ± 0.006	0.703 ± 0.004	0.500
TAILED-UNIFORM	0.355 ± 0.018	0.704 ± 0.013	0.452
<i>Uniform</i>	0.353 ± 0.012	0.706 ± 0.008	0.395
C2ST Improvement			+14%

Figure 9 shows that *Uniform* has notable dispersion along the degenerate direction. This leakage shows that when observations overlap boundary regions of the training distribution, *Uniform* has learned an inaccurate, spewed-out posterior and is unable to constrain parameters. The TAILED-UNIFORM, on the other hand, reproduces the banana-shaped contours of the MCMC reference, with samples tightly concentrated around the true parameter values. Table 2 quantifies these results. While

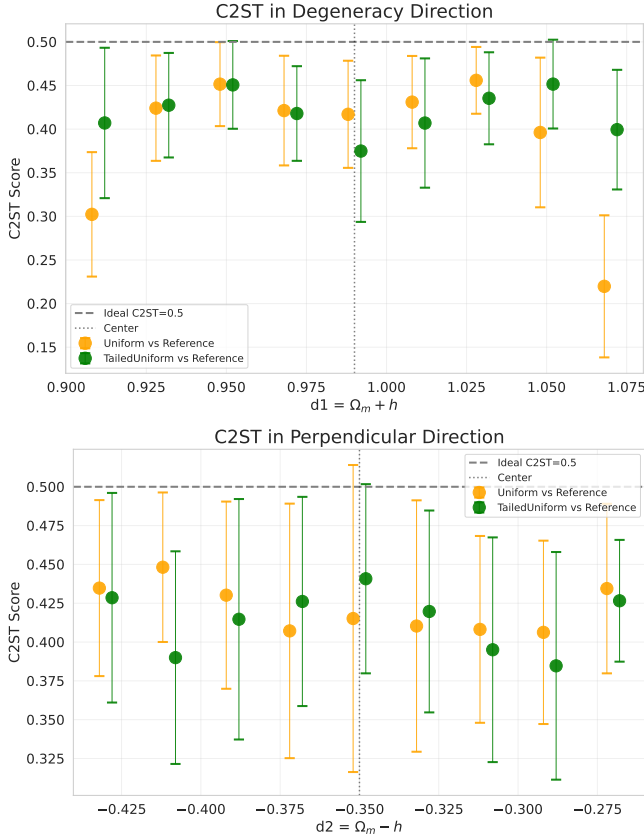


FIG. 11.— C2ST performance as a function of position along two orthogonal directions in parameter space. **Top:** Along the degeneracy direction, the *Uniform* (yellow) exhibits performance degradation at the boundaries. The *TAILED-UNIFORM* (green) maintains robust performance throughout. **Bottom:** Along the perpendicular direction, both methods perform similarly.

all three methods recover similar point estimates for both parameters, the C2ST scores reveal substantial differences in posterior quality. The *Uniform* has a C2ST of only 0.395, while the *TAILED-UNIFORM* scores 0.452, which is very close to the ideal value of 0.5, showing a 14% increase in distributional fidelity. This improvement demonstrates that *TAILED-UNIFORM* can capture the full covariance structure of the banana-shaped posterior, more so along the degenerate direction.

4.2.2. Spatial Performance Analysis

As before, we discretize a uniform rectangular grid with $n = 20$ test locations along each dimension, spanning $\Omega_m \in [0.27, 0.37]$ and $h \in [0.63, 0.71]$. For each of the 400 evaluation points, we generate observations, draw $M = 1000$ posterior samples from all three methods, and calculate pairwise C2ST scores.

Both methods achieve similar C2ST scores in Figure 10, but by eye inspection, *TAILED-UNIFORM* appears to provide reliable inference everywhere, whereas *Uniform* is only productive near the center. To connect the performance degradation trend to the underlying physics, we propose a different binning strategy. We first define the degeneracy direction as $d_1 = \Omega_m + h$, which describes movement along the banana-shaped posterior where the physical matter density $\Omega_m h$ varies slowly. Perpendicular to this, we define $d_2 = \Omega_m - h$, which

captures variation across the degeneracy.

Figure 11 shows that boundary pathology is anisotropic, with the most severe manifestation along degenerate directions where posteriors are elongated into the prior boundaries. As a result, both approaches perform about the same (excluding statistical fluctuations) along the perpendicular direction d_2 . However, across the degeneracy direction d_1 , the *Uniform* performance sinks at both boundaries, with C2ST scores falling to 0.22–0.30. The *Uniform* is unable to learn the posterior tails because the elongated banana-shaped posterior has spread into and beyond areas with sparse training coverage. By extending the simulation distribution beyond the prior boundaries, the *TAILED-UNIFORM* enables robust, unbiased inference across the full parameter space.

5. CONCLUSION

We have identified and addressed a fundamental limitation of neural posterior estimation: the boundary pathology caused by sharp discontinuities in uniform proposal distributions, which systematically degrades posterior quality near parameter space boundaries and is exacerbated both in high-dimensional spaces and when parameter degeneracies extend posteriors toward prior edges. To resolve this limitation, we proposed the *TAILED-UNIFORM* proposal, a hybrid distribution that maintains uniform density within the primary region of interest while extending smooth Gaussian tails beyond the boundaries, providing neural networks with the gradient information necessary to learn continuous density transitions.

In the two-dimensional Gaussian Linear benchmark, *TAILED-UNIFORM* maintains consistent performance (C2ST ≈ 0.44 – 0.49) throughout the domain, while *Uniform* exhibits systematic degradation toward parameter space edges (C2ST scores falling below 0.35). We also validated our method for inferring cosmological parameters (Ω_m, h) from matter power spectra. Empirically, *TAILED-UNIFORM* provides support along the Ω_m – h degeneracy, as elongated posteriors extend toward prior boundaries.

Perhaps most critically, our ablation studies demonstrate that the degradation of the original sampling strategy is structural rather than statistical: increasing training data yields negligible improvement for *Uniform*, whereas *TAILED-UNIFORM* with as few as $N = 2,000$ simulations consistently outperforms *Uniform* with eight times the training budget. This confirms that the problem lies in the discontinuous proposal itself rather than insufficient coverage. Our method also exhibits counterintuitive scaling behavior in higher dimensions: despite allocating an increasingly large fraction of samples to tail regions as dimensionality grows (from 24% at $d = 2$ to 86% at $d = 16$), *TAILED-UNIFORM* maintains superior performance over uniform sampling across all tested dimensions.

In conclusion, *TAILED-UNIFORM* gives neural posterior estimators the support they need to learn accurate posteriors close to boundaries by expanding the training distribution beyond the primary region of interest with smooth Gaussian tails. Our method provides a parsimonious, computationally cheap modification to existing SBI workflows that results in significant improvements in posterior quality across inference tasks.

While we have demonstrated TAILED-UNIFORM’s effectiveness up to $d = 16$ dimensions, the scaling behavior at very high dimensions ($d \gtrsim 32$) remains unexplored. At high dimensionality, the dilution of interior samples may come to dominate over the benefits of boundary smoothing. Additionally, while we focused on neural posterior estimation (NPE), we maintain that TAILED-UNIFORM is architecture-invariant and can be used with

neural likelihood estimation (NLE) and neural ratio estimation (NRE); comparative studies across these methodologies would further establish the universality of our approach.

We appreciate Shy Genel’s insightful comments on a previous draft. This project was developed as part of the Simons Collaboration on “Learning the Universe.”

REFERENCES

- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2623–2631, 2019.
- Natalie M Batalha, Shawn D Domagal-Goldman, Ramses Ramirez, and James F Kasting. Climate and habitability of terrestrial exoplanets. *Handbook of Exoplanets*, pages 2375–2394, 2017.
- Daniel Baumann. *Cosmology*. Cambridge University Press, Cambridge, 2022. ISBN 978-1108838078.
- Richard Bellman. Adaptive control processes. *Princeton University Press*, 1961.
- Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. *Handbook of Markov Chain Monte Carlo*. CRC Press, 2011.
- Manuela Campanelli, Carlos O Lousto, Pedro Marronetti, and Yosef Zlochower. Accurate evolutions of orbiting black-hole binaries without excision. *Physical Review Letters*, 96(11):111101, 2006.
- Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based inference. *Proceedings of the National Academy of Sciences*, 117(48):30055–30062, 2020.
- Scott Dodelson and Fabian Schmidt. *Modern Cosmology*. Academic Press, 2nd edition, 2020.
- Eric D Feigelson and Jogesh G Babu. Modern statistical methods for astronomy: with R applications. *Cambridge University Press*, 2012.
- Daniel Foreman-Mackey, David W Hogg, Dustin Lang, and Jonathan Goodman. emcee: the MCMC hammer. *Publications of the Astronomical Society of the Pacific*, 125(925):306, 2013.
- Andrew Gelman and Donald B Rubin. Inference from iterative simulation using multiple sequences. *Statistical Science*, 7(4):457–472, 1992.
- Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian data analysis*. CRC Press, 3rd edition, 2013.
- Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. MADE: Masked autoencoder for distribution estimation. In *International Conference on Machine Learning*, pages 881–889. PMLR, 2015.
- David Greenberg, Marcel Nonnenmacher, and Jakob Macke. Automatic posterior transformation for likelihood-free inference. In *International Conference on Machine Learning*, pages 2404–2414. PMLR, 2019.
- W Keith Hastings. Monte carlo sampling methods using markov chains and their applications. *Biometrika*, 57(1):97–109, 1970.
- Matthew Ho, Gwendolyn M Farren, Helen Shao, Natalí Anau Montel, Chirag Modi, Bruno Régalo-Saint Blancard, Pablo Lemos, Adam Coogan, Yashar Hezaveh, Laurence Perreault-Levasseur, Benjamin D Wandelt, and Natalí SM de Santi. Ltu-ili: An all-in-one framework for implicit inference in astrophysics and cosmology. *The Open Journal of Astrophysics*, 7, 2024. arXiv:2402.05137.
- Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.
- David Lopez-Paz and Maxime Oquab. Revisiting classifier two-sample tests. In *International Conference on Learning Representations*, 2017. arXiv:1610.06545.
- Jan-Matthis Lueckmann, Pedro J Gonçalves, Giacomo Bassetto, Kaan Öcal, Marcel Nonnenmacher, and Jakob H Macke. Flexible statistical inference for mechanistic models of neural dynamics. *Advances in Neural Information Processing Systems*, 30, 2017.
- Jan-Matthis Lueckmann, Jan Boelts, David Greenberg, Pedro J Gonçalves, and Jakob H Macke. Benchmarking simulation-based inference. In *International Conference on Artificial Intelligence and Statistics*, pages 343–351. PMLR, 2021.
- M D McKay, R J Beckman, and W J Conover. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 21(2):239–245, 1979.
- Nicholas Metropolis, Arianna W Rosenbluth, Marshall N Rosenbluth, Augusta H Teller, and Edward Teller. Equation of state calculations by fast computing machines. *The Journal of Chemical Physics*, 21(6):1087–1092, 1953.
- George Papamakarios and Iain Murray. Fast ϵ -free inference of simulation models with Bayesian conditional density estimation. In *Advances in Neural Information Processing Systems*, volume 29, pages 1028–1036, 2016.
- George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density estimation. *Advances in Neural Information Processing Systems*, 30, 2017.
- J. A. Peacock and S. J. Dodds. Non-linear evolution of cosmological power spectra. *Monthly Notices of the Royal Astronomical Society*, 258(1):1P–11P, sep 1992. doi:10.1093/mnras/258.1.1P.
- Planck Collaboration, N Aghanim, Y Akrami, et al. Planck 2018 results. vi. cosmological parameters. *Astronomy & Astrophysics*, 641:A6, 2020. arXiv:1807.06209.
- Frans Pretorius. Evolution of binary black-hole spacetimes. *Physical Review Letters*, 95(12):121101, 2005.
- Gareth O Roberts, Andrew Gelman, and Walter R Gilks. Weak convergence and optimal scaling of random walk metropolis algorithms. *The Annals of Applied Probability*, 7(1):110–120, 1997.
- Ce Sui, Deaglan J Bartlett, Shivam Pandey, Harry Desmond, Pedro G Ferreira, and Benjamin D Wandelt. Syren-new: Precise formulae for the linear and nonlinear matter power spectra with massive neutrinos and dynamical dark energy. *arXiv preprint arXiv:2410.14623*, 2024.
- Francisco Villaescusa-Navarro, Daniel Anglés-Alcázar, Shy Genel, et al. The camels project: Cosmology and astrophysics with machine learning simulations. *The Astrophysical Journal*, 915(1):71, 2021.
- Francisco Villaescusa-Navarro, Shy Genel, Daniel Anglés-Alcázar, et al. The camels project: Public data release. *The Astrophysical Journal Supplement Series*, 265(2):54, 2023.

provides fast and easy peer review for new papers in the **astro-ph** section of the arXiv, making the reviewing process simpler for authors and referees alike. Learn more at <http://astro.theoj.org>.

This paper was built using the Open Journal of Astrophysics L^AT_EX template. The OJA is a journal which